



HAL
open science

Contributions to a Theory of Pure Exploration in Sequential Statistics

Antoine Barrier

► **To cite this version:**

Antoine Barrier. Contributions to a Theory of Pure Exploration in Sequential Statistics. Machine Learning [cs.LG]. Ecole normale supérieure de lyon - ENS LYON, 2023. English. NNT: 2023ENSL0024 . tel-04192097

HAL Id: tel-04192097

<https://theses.hal.science/tel-04192097>

Submitted on 31 Aug 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



THÈSE

en vue de l'obtention du grade de Docteur, délivré par
l'ÉCOLE NORMALE SUPÉRIEURE DE LYON

École Doctorale N° 512
École Doctorale en Informatique, Mathématiques de Lyon (InfoMaths)

Discipline : Mathématiques

Soutenue publiquement le 20/07/2023, par :

Antoine BARRIER

Contributions à une théorie de l'exploration pure en statistique séquentielle

Devant le jury composé de :

MAILLARD, Odalric-Ambrym	Chargé de recherche	INRIA, Univ. de Lille	Rapporteur
VERNADE, Claire	Professeure	Université de Tübingen	Rapporteuse
PERCHET, Vianney	Professeur des universités	ENSAE (CREST)	Examineur
GARIVIER, Aurélien	Professeur des universités	CNRS, ÉNS de Lyon	Directeur de thèse
STOLTZ, Gilles	Directeur de recherche	CNRS, Univ. Paris-Saclay	Co-directeur

Remerciements

Chers Aurélien et Gilles, merci de m'avoir fait confiance pour ces trois années d'aventures! De par vos différences dans votre approche de la recherche et dans vos habitudes de travail, j'ai appris à vos côtés plus que je ne l'aurais imaginé. Aurélien, je suis toujours impressionné par ta vitesse de raisonnement et par ta maîtrise de tant de domaines des statistiques : une mise en contexte ultra succincte te suffit pour jongler d'un problème à un autre. Tu as su instaurer une ambiance de travail très agréable au sein du petit (mais grossissant) groupe de statisticiens de l'UMPA grâce aux réunions hebdomadaires. Gilles, ton efficacité et ton organisation à toute épreuve m'interrogent sur ton appartenance à l'espèce humaine : 100% de mails répondus dans la demi-heure (peu importe l'heure et le jour d'envoi), une qualité de rédaction inégalable (tout comme ton sens de la formule) et une précision infinie lors de tes relectures (merci pour toutes les typos corrigées dans ce manuscrit et tout au long de cette thèse) ... Tout simplement épatant et inspirant!

Claire et Odalric, merci à tous les deux d'avoir accepté la lourde tâche de rapporter ma thèse, j'en suis honoré. Vos nombreux commentaires ont contribué à améliorer le manuscrit et à aboutir à cette version finale.

Vianney, je suis également très heureux de te compter parmi mon jury. En fait, tu es celui qui m'a initié à l'apprentissage statistique et à l'intelligence artificielle à Cachan, avec tes exposés en L3 et surtout ton cours en M1. C'est donc une belle manière de boucler la boucle!

Mes remerciements ne sauraient se restreindre à ces trois années de thèse. Avant elles, c'est un long parcours qui a mené à l'élaboration du projet, où tant de personnes sont intervenues. Je profite de ces débuts de remerciements pour mettre en avant les professeurs passionnés et investis ---en mathématiques ou non--- que j'ai eu la chance d'avoir depuis mon enfance.

Les années de prépa à Champo ne me remémorent que de bons souvenirs, grâce notamment à de formidables enseignants. Parmi eux, il y a bien sûr JML, qui a plus qu'outrepasé ses fonctions depuis plusieurs années : tes nombreux conseils depuis la fin de la prépa m'ont beaucoup aidé, particulièrement durant cette thèse, et j'ai hâte de continuer à discuter de tes projets de sobriétés environnementale et numérique autour d'un bon burger.

Les études cachannaises ont aussi été mathématiquement très fructueuses, un merci particulier à Frédéric Pascal qui a été notre référent et interlocuteur préférentiel pendant ces quatre années, et à Dieter, mon encadrant de stage de M1, pour sa bienveillance et pour m'avoir guidé vers Aurélien.

Cette thèse n'aurait pas été si agréable sans un environnement de travail adéquat. L'ensemble des membres de l'UMPA ainsi que les étudiants du département de mathématiques (avec bien sûr la meilleure des promos d'agrégatifs cette année) y ont amplement contribué.

Je remercie les membres permanents pour toutes les discussions et tous les conseils, notamment Grégory et Thomas avec lesquels j'ai préparé mes enseignements. Une pensée à Claude, exemple de dévouement et de passion, me vient évidemment.

Merci aux secrétaires Jessica, Laure, Magalie et Virginia pour leur sympathie et leur invincible efficacité.

Merci à tous les membres de l'équipe de statistiques d'Aurélien : Alexandre, Aymen, Élise, Hugues, Meh-rasa, Pierre et Tomáš. Il est toujours intéressant de discuter de bandits, de \mathcal{RL} , ou d'autres problèmes de statistiques avec vous.

Je ne compte plus les moments de vie vécus (au laboratoire ou non) avec les membres non permanents : merci à Alexandre, Aymen, Basile, Benoît, Denis, Cécile, Charlie, Corentin, Héloïse, Hugues, Jules, Juliana,

Léo, Matthieu, Paul, Raphaël (χ2), Riccardo, Ronan, Simon, Thomas, Valentine, Vanessa, Vianney, William et tant d'autres ... Un merci particulier aux incroyables co-chargés de TD Juliana et Léo, et aux co-locataires du bureau 404 (dont certains cherchent encore l'existence) Hugues¹ et Basile, qui ont mis une belle ambiance (parfois peu propice au travail ...) ces derniers temps, mais aussi avant eux Matthieu et Felipe.

Au-delà du laboratoire, j'ai eu la chance de participer à de nombreux séminaires qui furent tous de beaux moments de rencontres avec de jeunes chercheurs : je remercie (la liste est loin d'être exhaustive) Annette, Dorian, El Mehdi, Emmanuel, Eugenio, Joseph, Julien, Juliette, Marc, Mete, Quentin et Thibault. Marc, merci également pour nos tentatives de collaboration cette année et pour tous les bons moments passés au CIRM et à Singapour!

Si ces trois années sont passées si vite, c'est grâce à un environnement mathématique stimulant, mais aussi et surtout grâce aux moments passés entre amis ou en famille qui ont rythmé cette thèse et parsemé tant de souvenirs.

Cette thèse fut l'occasion de multiplier les activités sportives. En tête de liste, la découverte d'innombrables terrains de jeux alpins. Mention spéciale aux sorties (((rando-))) trail, vélo et ski de fond avec la team lyonnaise. En tête Rémy avec ses performances XXL à ski de fond, et son taux de nitro imbattable dans les montées de cols. Mais aussi Mattéo avec ses chaussures de trail dernier cri, ses cargaisons de sel et bien sûr sa force légendaire pour soulever des lits. Ou encore Pierre dont la capacité à soigner ses chutes ---aussi bien à vélo que dans les fameuses Pierrades--- n'a pas d'égal (et ça c'est tip top! Ou idyllique, je ne sais plus ...).

Que de beaux souvenirs me reviennent également des treks dans les Écrins avec les champions Fanny, Marie, Paul et Rémi, le Mercantour avec les boss Rémy et Mattéo, et la Vanoise avec les athlètes Jules, Mattéo, Pomme, Romain et Suzanne. Et à tous les autres qui ont partagé de beaux moments en montagne, à pied ou à vélo : Loïc, Florent, Maël, Hugues, Major, Paul, Quentin, ...

La bonne ambiance avec les pongistes, insufflée à l'AS par Vincent puis Fabrice, mais aussi et surtout au GBTI par tout un club, des jeunes aux moins jeunes (il y aurait tant de monde à citer ...), comptent aussi parmi les super moments sportifs de ces dernières années.

Il y a bien sûr des amis de longue date qui ont été, encore et toujours, particulièrement présents tout au long de cette thèse. Loïc, Louane et Nathan, vous êtes juste formidables ... Il s'en est encore passé des choses : les week-ends de l'ambiance, Saint-Mandrier, Imbours, Europa Park, Center Parcs, LouThan Parc², et tellement d'autres ... J'ai hâte de voir la suite du programme!

Je rajoute à cette liste les parisiens qui m'ont accueilli lors de mes passages dans la capitale. Les cachannais Anatole le Magicien (en tournée dans tout Paris), Doudouche le hipster, Jobic l'indescriptible et Farf à la bonne humeur inébranlable. Les anciens de Champo ne sont pas en reste : Major et le rituel apéro / pâtes au saumon / film à Bourg-la-Reine, Fanny qui répond toujours présente dès qu'il est possible de se voir (mille mercis pour ta curiosité malgré la distance cette année!), les sessions de jeux chez Ben et Cam's, et Stan et son intrépide aventure dans le monde de la finance. Last but not least, mon cousin Simon, toujours si accueillant avec ses bonnes bouteilles et à l'origine de tant de bons moments lors de mes années parisiennes!

Merci aussi aux infatigables Andy et Benoît, toujours prêts à organiser mille et un traquenards.

J'ai aussi eu la chance d'être proche de ma famille pendant cette thèse. Leur soutien a été constant et très important pour moi. Je remercie mes parents Nathalie et Thierry, mon frère Mathis, mais aussi mes grands-parents lyonnais Claudette et René qui ont tous tant œuvrés pour que je ne manque jamais de rien.

Et enfin, Berthine, mille mercis pour ton soutien ces derniers mois. ¡Espero que este sea el comienzo de una gran aventura!

¹ oups, je fais partie des personnes qui oublient malheureusement si souvent les « h » de ton nom et prénom, cher Hugues van Hassel ...

² c'est comme ça que je nommerai désormais votre appartement, tant il est synonyme de divertissement et d'amusement!

Et pour toi, cher lecteur, qui n'ira pas beaucoup plus loin dans la lecture de ce manuscrit, voici un petit pêle-mêle (comme dirait l'autre) de quelques formules anonymes entendues durant ces trois années :

« Un mathématicien c'est pas comme un joueur d'échec, il a le droit de se tromper. »

« Si tu étais à Ninja Warrior, je dirais que tu es au pied du mur. »

« En recherche, on ne fait que jouer au loto. »

« Tes maîtres de thèse sont tes boulets. »

« L'intelligence artificielle, c'est pas des maths. »

Pour terminer, sans doute la question la plus posée ces trois dernières années :

« Qu'est-ce que ça donne pour des bandits Gaussiens à trois bras ? »

Sommaire · Contents

Résumé et liste des publications	4
Abstract and origin of the materials	6
Chapitre 1 Vue d'ensemble des résultats	9
1 Exploration pure dans les problèmes de bandits	10
2 Identification de meilleur bras à confiance fixée	11
3 Identification de meilleur bras à budget fixé	21
Notation	27
Chapter 2 Pure Exploration in Multi-Armed Bandits	31
1 A Sequential Learning Problem	32
2 Best-Arm Identification with a Fixed-Confidence	35
3 Best-Arm Identification with a Fixed-Budget	58
Chapter 3 About the Fixed-Confidence Sample Complexity Optimization Problem for Gaussian Variables	71
1 Introduction	72
2 Solving the Optimization Problem	74
3 Bounds and Computation of the Problem Characteristics	77
4 Monotonicity of the $\max\text{-min}$ Problem	80
5 Regularity Properties	86
6 Conclusion	91
Chapter 4 A Fixed-Confidence Strategy with Non-Asymptotic Guarantees	93
1 Introduction	94
2 The Exploration-Biased-Sampling Strategy	97
3 Theoretical Results	103
4 Numerical Experiments	104
5 Proof of the Non-Asymptotic Bounds of Theorem 4.5	110
6 Technical Results	116
7 Conclusion	121
Chapter 5 A Non-Parametric Theory of Fixed-Budget Best-Arm Identification	123
1 Introduction	124
2 Overview of the Results and more Extended Literature Review	126
3 Upper Bound for the Successive-Rejects Strategy, with an Improved Analysis	131
4 Lower Bounds	137
5 Technical Details	144

6	Additional Comments for the Literature Review	155
7	Conclusion	160
Chapter 6	Asymptotically Optimal Adaptive Top-Two Algorithms in the Fixed-Confidence Setting	161
1	Introduction	162
2	A Fixed Point Property	165
3	Asymptotically Optimal Adaptive Algorithms	174
4	Side Note: On the Asymptotic Optimality of Track-and-Stop	184
5	Conclusion	186
	Bibliography	187
	Index	191

Résumé

Contributions à une théorie de l'exploration pure en statistique séquentielle

Cette thèse, à la croisée entre les domaines de l'intelligence artificielle, de la statistique séquentielle et de l'optimisation, s'intéresse au problème d'identification du meilleur bras (en espérance) dans les bandits non structurés à K bras. Ce problème possède deux approches dont les niveaux de compréhension sont très différents.

Le cadre à confiance fixée est le mieux compris : des stratégies asymptotiquement optimales sont connues, et l'on s'intéresse à l'obtention de garanties non asymptotiques pour des stratégies (si possible) simples et naturelles. Avec des bandits Gaussiens, nous proposons l'analyse à risque fini d'une nouvelle stratégie (asymptotiquement optimale) grâce aux propriétés de régularité de ce modèle. Cette stratégie modifie subtilement la règle d'attribution des tirages de l'algorithme Track-and-Stop en une règle plus prudente et interprétable. Dans le contexte plus général d'un modèle exponentiel, nous proposons l'ébauche d'une analyse de l'asymptotique optimalité d'algorithmes de type Top-Two adaptatifs, dont les règles de choix de tirages sont particulièrement simples.

Par ailleurs, dans le cadre à budget fixé, où l'existence d'une hypothétique complexité reste à démontrer, nous proposons des généralisations à des modèles non-paramétriques des bornes (supérieures et inférieures) connues jusqu'à présent pour des modèles très spécifiques. Les bornes obtenues font intervenir des quantités de théorie de l'information plus précises que les écarts entre les moyennes qui apparaissaient précédemment. Ces nouvelles quantités pourraient être la clé pour mesurer la complexité de l'identification de meilleur bras à budget fixé.

Mots-clés. Problèmes de bandits · Identification de meilleur bras · Statistiques séquentielles · Apprentissage statistique · Intelligence artificielle

Liste des publications

Le présent manuscrit s'appuie sur les deux publications suivantes :



A. Barrier, A. Garivier, and T. Kocák. A Non-Asymptotic Approach to Best-Arm Identification for Gaussian Bandits. In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, pages 10078–10109. PMLR, 2022



A. Barrier, A. Garivier, and G. Stoltz. On Best-Arm Identification with a Fixed Budget in Non-Parametric Multi-Armed Bandits. In *Proceedings of the 34th International Conference on Algorithmic Learning Theory*, volume 201 of *Proceedings of Machine Learning Research*, pages 136–181. PMLR, 2023

Les chapitres 3 et 4 reprennent le contenu du papier AISTATS 2022, qui porte sur l'étude d'une nouvelle stratégie d'identification de meilleur bras à confiance fixée, *Exploration-Biased-Sampling*, pour des variables gaussiennes. Des garanties à risque fini sont notamment démontrées, en utilisant de nouveaux résultats de régularité du problème d'optimisation définissant les poids optimaux.

Ensuite, le chapitre 5 présente les travaux issus de la publication ALT 2023. L'objectif est de généraliser les bornes connues en identification de meilleur bras à budget fixé à des modèles généraux, possiblement non paramétriques, en faisant intervenir de nouvelles quantités de théorie de l'information plus précises que les gaps.

Le chapitre 6, qui s'intéresse aux propriétés asymptotiques des algorithmes top-two adaptatifs en confiance fixée, est une ébauche contenant des résultats préliminaires non publiés.

Abstract

Contributions to a Theory of Pure Exploration in Sequential Statistics

This thesis lies in the fields of artificial intelligence, sequential statistics and optimization. We focus on the problem of best (in expectation) arm identification in unstructured multi-armed bandits. This problem has two approaches with very different levels of understanding.

The fixed-confidence framework is the best understood: asymptotically optimal strategies are known, and we are interested in obtaining non-asymptotic guarantees for (if possible) simple and natural strategies. Working with Gaussian bandits, we propose a finite risk analysis of a new (asymptotically optimal) strategy using the regularity properties of this model. This strategy slightly modifies the sampling rule of the Track-and-Stop algorithm into a more conservative and interpretable rule. In the more general context of an exponential model, we propose a preliminary analysis of the asymptotic optimality of adaptive Top-Two algorithms, whose sampling rules are particularly simple.

Independently, in the fixed-budget framework, for which the existence of a hypothetical complexity remains to be demonstrated, we propose generalizations to non-parametric models of the existing bounds (upper and lower) that were available so far only for very specific models. The obtained bounds involve more precise information-theoretic quantities than the gaps (differences between the means) which appeared previously. These new quantities could be the key to measuring the complexity of fixed-budget best-arm identification.

Keywords. Multi-Armed Bandits · Best-Arm Identification · Sequential Statistics · Statistical Learning · Machine Learning

Origin of the materials

This manuscript is based on the following two publications:



A. Barrier, A. Garivier, and T. Kocák. A Non-Asymptotic Approach to Best-Arm Identification for Gaussian Bandits. In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, pages 10078–10109. PMLR, 2022



A. Barrier, A. Garivier, and G. Stoltz. On Best-Arm Identification with a Fixed Budget in Non-Parametric Multi-Armed Bandits. In *Proceedings of the 34th International Conference on Algorithmic Learning Theory*, volume 201 of *Proceedings of Machine Learning Research*, pages 136–181. PMLR, 2023

Chapters 3 and 4 follow the contents of the AISTATS 2022 paper, which deals with the study of a new fixed-confidence best-arm identification strategy, *Exploration-Biased-Sampling*, for Gaussian variables. In particular, finite-risk guarantees are demonstrated, using new regularity results for the optimization problem defining the optimal weights.

Then, Chapter 5 presents the work resulting from the ALT 2023 publication. The aim is to generalize the known bounds in best-arm identification with a fixed budget to general, possibly non-parametric, models, by involving new quantities of information theory that are more precise than the gaps.

Finally, Chapter 6, which deals with the asymptotic properties of adaptive top-two algorithms in the fixed-confidence setting, is a draft containing preliminary unpublished results.

CHAPITRE 1

Vue d'ensemble des résultats

Ce chapitre initial résume l'ensemble des contributions de la thèse en les replaçant dans le contexte bibliographique de l'identification de meilleur bras. Après une présentation du problème, nous expliquons comment une profonde compréhension du modèle Gaussien dans le cadre à confiance fixée nous a permis de définir une stratégie avec des garanties non asymptotiques. Puis, dans le cadre à budget fixé, nous énonçons des généralisations de bornes existantes à des modèles quelconques en introduisant de nouvelles mesures de complexité à base de quantités de théorie de l'information.

Sommaire

1	Exploration pure dans les problèmes de bandits	10
2	Identification de meilleur bras à confiance fixée	11
1	Borne inférieure pour un modèle exponentiel	12
2	Résolution du problème d'optimisation définissant $T(\underline{\mu})$	14
3	L'algorithme Track-and-Stop	17
4	Des garanties non asymptotiques	18
5	Les algorithmes Top-Two	19
3	Identification de meilleur bras à budget fixé	21
1	Bornes connues de la littérature	21
2	Généralisations des bornes existantes à des modèles quelconques	22
3	Existence d'une complexité	25

1.1. Exploration pure dans les problèmes de bandits

Dans cette thèse, on s'intéresse au problème statistique suivant. Un joueur (un *learner*) fait face à $K \geq 2$ distributions de probabilité ν_1, \dots, ν_K inconnues. À chaque instant, il peut observer une réalisation de la distribution de son choix, en utilisant les observations passées. Son objectif est d'identifier rapidement la distribution de plus grande espérance. Il s'agit d'un problème d'optimisation séquentielle aux nombreuses applications.

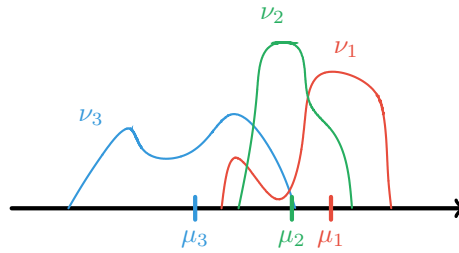


Figure 1.1: μ_a désigne l'espérance de ν_a . Considérant les trois distributions dont nous représentons les fonctions de masse, l'objectif du joueur est d'identifier ν_1 comme étant la meilleure distribution.

Par exemple, en médecine, si l'on souhaite déterminer quel médicament est le plus efficace pour traiter une maladie parmi un ensemble de K médicaments, on peut organiser un essai clinique où l'on administre à chaque malade un des médicaments et l'on observe s'il guérit ou non. On peut alors procéder de manière séquentielle en choisissant le médicament administré à un patient en fonction de toutes les observations précédentes. Si, pour les premiers patients, on voudra observer les médicaments uniformément en l'attente de données fiables, il faudra en cours d'essai clinique se concentrer au fur et à mesure sur les médicaments les plus prometteurs empiriquement. L'efficacité de l'essai clinique dépendra de la stratégie adoptée, et donc de l'arbitrage choisi entre exploration de l'ensemble des distributions et focalisation sur les distributions les plus prometteuses.

Problèmes de bandits. Le cadre et la terminologie utilisés pour modéliser la situation sont ceux des bandits (Lattimore and Szepesvári, 2020) : on se donne un ensemble $\underline{\nu} = (\nu_1, \dots, \nu_K)$ de distributions appartenant à un modèle donné \mathcal{D} , où $K \geq 2$ est fixé. Bien que les motivations initiales concernaient des essais cliniques (Thompson, 1933), le terme de *problème de bandit* est utilisé pour désigner $\underline{\nu}$ et provient du nom de machines à sous, les *bandits manchots*, dans les casinos. Les indices des distributions (et par extension leurs distributions associées) sont quant à eux les *bras* du problème $\underline{\nu}$. On désigne par $\underline{\mu} = (\mu_1, \dots, \mu_K)$ le vecteur de moyennes de $\underline{\nu}$, où $\mu_a = \mathbb{E}(\nu_a)$ est l'espérance de ν_a . Dans l'ensemble de la thèse, nous travaillons avec le modèle élémentaire des bandits *non structurés*, i.e., dont les distributions sont supposées indépendantes. Parmi les modèles de distributions possibles, on peut notamment citer :

- le modèle $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ des distributions Gaussiennes avec variance commune $\sigma^2 > 0$, le modèle $\mathcal{D}_{\mathcal{B}}$ des distributions de Bernoulli, ou plus généralement les modèles exponentiels,
- le modèle non paramétrique $\mathcal{P}[0, 1]$ des lois à valeurs dans $[0, 1]$, ou plus généralement tout modèle \mathcal{D}_{σ^2} constitué de variables σ^2 -sous-Gaussiennes, i.e., de lois ν satisfaisant la borne suivante sur la fonction des moments, où $X \sim \nu$:

$$\forall \lambda \in \mathbb{R}, \quad \mathbb{E} \left[e^{\lambda(X - \mathbb{E}[X])} \right] \leq e^{\frac{\lambda^2 \sigma^2}{2}} .$$

1.2. IDENTIFICATION DE MEILLEUR BRAS À CONFIANCE FIXÉE

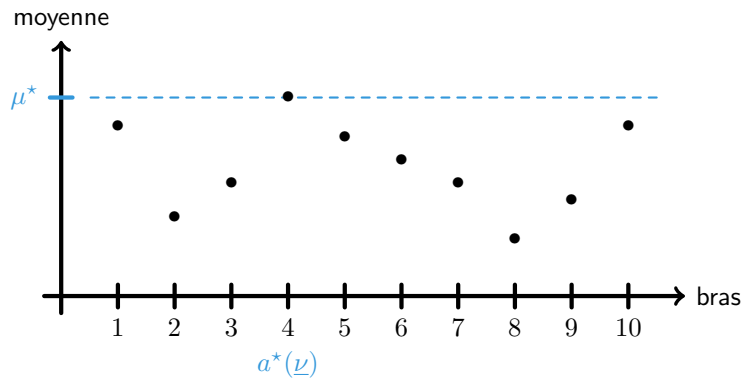


Figure 1.2: Bras optimal et moyenne optimale d'un problème de bandit $\underline{\nu}$.

Identification de meilleur bras. Le problème présenté dans le premier paragraphe est connu sous le nom d'*identification de meilleur bras* ou d'*exploration pure* : l'objectif est d'identifier le bras, supposé unique et noté $a^*(\underline{\nu})$, dont l'espérance est la plus élevée :

$$\{a^*(\underline{\nu})\} \stackrel{\text{def}}{=} \operatorname{argmax}_{1 \leq a \leq K} \mu_a.$$

Dans l'ensemble de ce chapitre, nous considérons des problèmes de bandits ayant un unique bras *optimal*. Pour cela, le joueur donne, après un nombre (potentiellement aléatoire) τ d'observations, une estimation $\hat{a}_\tau \in [K] \stackrel{\text{def}}{=} \{1, \dots, K\}$ du meilleur bras $a^*(\underline{\nu})$. Ce problème peut être traité sous deux approches différentes en fonction des applications.

- Une première possibilité est de travailler à *confiance fixée* (Even-Dar et al., 2006) : le joueur ne peut pas décider de s'arrêter tant qu'il n'est pas certain d'avoir identifié le bras optimal avec un niveau de risque inférieur à un seuil $\delta \in (0, 1)$ fixé. On dit alors que sa stratégie est δ -correcte. Le calcul du niveau de risque d'une stratégie dépend des hypothèses sur les distributions, i.e., sur le modèle \mathcal{D} auxquelles elles appartiennent. Ces hypothèses sont connues du joueur. Parmi toutes les stratégies δ -correctes, on voudrait trouver celles qui minimisent le nombre nécessaire d'observations.
- La deuxième approche est donnée par le cadre du *budget fixé* (Audibert et al., 2010) : le joueur est obligé de donner son estimation après un nombre total d'observations T fixé. On cherche les stratégies qui minimisent la probabilité d'erreur, c'est-à-dire la probabilité d'avoir une mauvaise estimation du bras optimal.

De nombreuses applications existent pour les deux approches, selon le nombre d'observations possibles : le budget fixé est adapté pour des essais cliniques où le nombre de patients à disposition est difficilement extensible, alors que la confiance fixée peut être utilisée, par exemple, par des sites web souhaitant tester diverses configurations et disposant de nombreux internautes (test A/B).

De manière assez surprenante, les deux approches conduisent à des connaissances théoriques et empiriques très différentes. Nous allons voir que les stratégies d'apprentissage diffèrent sensiblement.

1.2. Identification de meilleur bras à confiance fixée

Nous étudions d'abord l'identification de meilleur bras à confiance fixée : étant donné $\delta \in (0, 1)$ fixé, la stratégie du joueur doit s'arrêter après un temps aléatoire τ_δ de sorte de pouvoir garantir que son estimation \hat{a}_{τ_δ} est correcte avec risque δ . L'objectif est de trouver des stratégies δ -correctes qui nécessitent le moins d'observations possibles.

Stratégie. Pour commencer, nous définissons formellement ce qu'est une stratégie. On se donne une suite $(U_t)_{t \geq 0}$ de variables aléatoires i.i.d. de loi $\mathcal{U}([0, 1])$, indépendante de toute autre source d'aléa, et on note Y_t l'observation au temps $t \geq 1$. La stratégie du joueur est définie par la donnée :

- d'une *règle d'échantillonnage*, qui décide du choix du bras $A_t \in [K]$ à observer à l'étape $t \geq 1$, en fonction des observations passées $I_{t-1} \stackrel{\text{def}}{=} (U_0, Y_1, U_1, Y_2, U_2, \dots, Y_{t-1}, U_{t-1})$; A_t est donc \mathcal{F}_{t-1} -mesurable, où $\mathcal{F}_{t-1} = \sigma(I_{t-1})$.
- d'une *règle d'arrêt* τ_δ , qui est un temps d'arrêt par rapport à $(\mathcal{F}_t)_{t \geq 0}$,
- d'une *règle de décision* \hat{a}_{τ_δ} qui est $\mathcal{F}_{\tau_\delta}$ -mesurable et décide quel bras est estimé comme étant optimal (il s'agit le plus souvent du bras de meilleure moyenne empirique).

La stratégie est dite δ -correcte si elle vérifie

$$\forall \underline{\nu} \text{ dans } \mathcal{D}, \quad \mathbb{P}_{\underline{\nu}}(\tau_\delta < +\infty, \hat{a}_{\tau_\delta} \neq a^*(\underline{\nu})) \leq \delta,$$

où $\mathbb{P}_{\underline{\nu}}$ désigne la probabilité sous le problème de bandit $\underline{\nu}$.

Notations. Pour un problème de bandit $\underline{\nu}$ fixé, on note $\Delta_a(\underline{\nu}) = \mu^* - \mu_a$ le *gap* du bras a , c'est-à-dire l'écart entre la moyenne optimale et la moyenne du bras a . Pour une stratégie donnée face au problème $\underline{\nu}$, on note $N_a(t)$ et $\hat{\mu}_a(t)$ le nombre de tirages et la moyenne empirique¹ du bras a à la fin de l'étape t :

$$N_a(t) \stackrel{\text{def}}{=} \sum_{s \in [t]} \mathbb{I}\{A_s = a\} \quad \text{and} \quad \hat{\mu}_a(t) \stackrel{\text{def}}{=} \frac{1}{N_a(t)} \sum_{s \in [t]} Y_s \mathbb{I}\{A_s = a\},$$

où $\mathbb{I}\{E\}$ est l'indicatrice de l'évènement E .

1.2.1. Borne inférieure pour un modèle exponentiel

Un bon critère pour mesurer la performance de ces stratégies est le nombre moyen de tirages $\mathbb{E}_{\underline{\nu}}[\tau_\delta]$. Nous allons voir qu'une notion d'optimalité existe en effet pour cette quantité. Dans un premier temps, [Garivier and Kaufmann \(2016\)](#) ont montré une borne inférieure sur le nombre moyen de tirages d'une stratégie δ -correcte pour un modèle exponentiel. Elle repose sur une inégalité de théorie de l'information faisant intervenir la divergence de Kullback-Leibler.

Remarque. Les premières stratégies d'identification de meilleur bras à confiance fixée, comme la stratégie d'élimination de [Even-Dar et al. \(2006\)](#), garantissent des bornes sur τ_δ avec forte probabilité. Cependant, il n'y a pas de borne inférieure pour ce critère.

La divergence de Kullback-Leibler. La *divergence de Kullback-Leibler* est une pseudo-distance sur les distributions de probabilité qui joue un rôle important en théorie de l'information. Elle est définie, pour toute paire de probabilités \mathbb{P} et \mathbb{Q} définies sur un espace mesurable (Ω, \mathcal{A}) , par

$$\text{KL}(\mathbb{P}, \mathbb{Q}) \stackrel{\text{def}}{=} \begin{cases} \int_{\Omega} \log \frac{d\mathbb{P}}{d\mathbb{Q}} d\mathbb{P} = \mathbb{E}_{X \sim \mathbb{P}} \left[\log \frac{d\mathbb{P}}{d\mathbb{Q}}(X) \right] & \text{si } \mathbb{P} \ll \mathbb{Q}, \\ +\infty & \text{sinon,} \end{cases}$$

où $\frac{d\mathbb{P}}{d\mathbb{Q}}$ est la dérivée de Radon-Nikodym de \mathbb{P} par rapport à \mathbb{Q} lorsque \mathbb{P} est absolument continue par rapport à \mathbb{Q} . On peut vérifier que la KL est toujours positive mais n'est pas symétrique, et que $\text{KL}(\mathbb{P}, \mathbb{Q}) = 0$ si et seulement si $\mathbb{P} = \mathbb{Q}$.

¹Toutes les stratégies observent chaque bras une fois initialement, donc $\hat{\mu}_a(t)$ est bien défini pour $t \geq K$.

L'inégalité fondamentale. Les bornes inférieures de la littérature sur les bandits utilisent régulièrement des changements de mesure : on considère un problème de bandit $\underline{\zeta}$, dit *alternatif*, dont le bras optimal diffère de celui de $\underline{\nu}$ et on souhaite quantifier combien d'observations sont nécessaires pour discriminer $\underline{\nu}$ et $\underline{\zeta}$. Dans notre cadre, il est pratique d'utiliser l'inégalité de théorie de l'information suivante, qui évite l'utilisation de changements de mesure plus explicites. Une nouvelle preuve de ce résultat, utilisant une martingale arrêtée, est donnée en Section 2.2.2.

Lemme 1.1. [Kaufmann et al., 2016, Lemma 1]

Soient $\underline{\nu}$ et $\underline{\zeta}$ deux problèmes de bandits. Soit une stratégie d'identification de meilleur bras telle que le temps d'arrêt τ est $\mathbb{P}_{\underline{\nu}}$ -intégrable. Alors, pour tout événement E dans \mathcal{F}_τ , on a

$$\sum_{a \in [K]} \mathbb{E}_{\underline{\nu}}[N_a(\tau)] \text{KL}(\nu_a, \zeta_a) \geq \text{KL}(\text{Ber}(\mathbb{P}_{\underline{\nu}}(E)), \text{Ber}(\mathbb{P}_{\underline{\zeta}}(E))).$$

Modèles exponentiels. On désigne par modèle exponentiel, et on note \mathcal{D}_{exp} , toute famille exponentielle canonique à un paramètre (voir Lehmann and Casella, 1998, Section 1.5). Il s'agit d'un ensemble de distributions ν_θ indexées par $\theta \in \Theta \subset \mathbb{R}$, toutes absolument continues par rapport à une mesure ρ sur \mathbb{R} , de densités données par

$$x \mapsto \frac{d\nu_\theta}{d\rho}(x) \stackrel{\text{def}}{=} \exp(\theta x - b(\theta)),$$

pour une fonction de normalisation b au moins deux fois dérivable. Citons par exemple les distributions de Bernoulli, binomiales, de Poisson, et Gaussiennes de variance commune. On peut montrer que $\mathbb{E}(\nu_\theta) = b'(\theta)$ pour tout $\theta \in \Theta$, et donc que les distributions du modèle sont caractérisées par leur espérance. La divergence de Kullback-Leibler peut aussi être paramétrée par les espérances : on note, pour tout $\theta, \theta' \in \Theta$,

$$d(\mathbb{E}(\nu_\theta), \mathbb{E}(\nu_{\theta'})) \stackrel{\text{def}}{=} \text{KL}(\nu_\theta, \nu_{\theta'}).$$

Cela définit une divergence d qui est strictement convexe et différentiable sur $\mathcal{M} \times \mathcal{M}$, où $\mathcal{M} = (\mu_-, \mu_+)$ est l'intervalle supposé ouvert des moyennes du modèle. En particulier, d est continue, telle que $d(\mu, \mu') = 0$ si et seulement si $\mu = \mu'$, et, pour tout $\mu \in \mathcal{M}$, $d(\mu, \cdot)$ et $d(\cdot, \mu)$ sont strictement décroissantes sur $(\mu_-, \mu]$, et strictement croissantes sur $[\mu, \mu_+)$.

La divergence de Kullback-Leibler entre variables de Bernoulli jouant un rôle particulier dans les bornes inférieures, on la note kl et on a l'expression suivante :

$$\forall p, q \in (0, 1), \quad \text{kl}(p, q) \stackrel{\text{def}}{=} \text{KL}(\text{Ber}(p), \text{Ber}(q)) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}.$$

Borne inférieure pour un modèle exponentiel. En utilisant le Lemme 1.1 et la propriété de stratégie δ -correcte, Garivier and Kaufmann (2016) ont prouvé la borne inférieure suivante pour un modèle exponentiel, où les problèmes de bandits $\underline{\nu}$ sont caractérisés par leur moyennes $\underline{\mu}$ comme expliqué ci-dessus, et où l'on définit le simplexe

$$\Sigma_K \stackrel{\text{def}}{=} \left\{ \underline{\nu} \in [0, 1]^K : \nu_1 + \dots + \nu_K = 1 \right\},$$

et l'ensemble des bandits *alternatifs* au problème $\underline{\mu}$:

$$\text{Alt}(\underline{\mu}) \stackrel{\text{def}}{=} \left\{ \lambda \text{ in } \mathcal{D}_{\text{exp}} : a^*(\lambda) \neq a^*(\underline{\mu}) \right\}.$$

Théorème 1.1. [Garivier and Kaufmann, 2016, Theorem 1]

Soit \mathcal{D}_{exp} un modèle exponentiel. Soit $\delta \in (0, \frac{1}{2})$. Pour toute stratégie δ -correcte, et pour tout problème de bandit $\underline{\mu}$ dans \mathcal{D}_{exp} avec un unique bras optimal,

$$\mathbb{E}_{\underline{\mu}}[\tau_{\delta}] \geq T(\underline{\mu}) \text{kl}(\delta, 1 - \delta) \geq T(\underline{\mu}) \log \frac{1}{2.4\delta}, \quad (1.1)$$

où $T(\underline{\mu})$ est le temps caractéristique de $\underline{\mu}$, défini par

$$T(\underline{\mu})^{-1} \stackrel{\text{def}}{=} \sup_{\underline{v} \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a). \quad (1.2)$$

Un vecteur de poids optimal. Garivier and Kaufmann (2016) ont montré que le problème d'optimisation (1.2) admet un unique maximiseur $\underline{w}(\underline{\mu})$, appelé le *vecteur de poids optimal* de $\underline{\mu}$:

$$\{\underline{w}(\underline{\mu})\} \stackrel{\text{def}}{=} \underset{\underline{v} \in \Sigma_K}{\text{argmax}} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a).$$

Pour obtenir une stratégie dont la performance est proche de la borne inférieure (1.1), il faut que les proportions moyennes de tirages $(\mathbb{E}_{\underline{\mu}}[N_a(\tau_{\delta})]/\mathbb{E}_{\underline{\mu}}[\tau_{\delta}])_{a \in [K]}$ de cette stratégie soient proches de $\underline{w}(\underline{\mu})$. Autrement dit, si la stratégie ne tire pas les bras selon des fréquences bien précises, elle ne pourra pas être performante.

1.2.2. Résolution du problème d'optimisation définissant $T(\underline{\mu})$

Considérons toujours un modèle exponentiel \mathcal{D}_{exp} . Nous allons voir que la solution du problème d'optimisation (1.2) définissant $T(\underline{\mu})$ est cruciale pour déterminer l'efficacité d'une stratégie.

Coûts de transport. Le problème d'optimisation (1.2) définissant $T(\underline{\mu})$ peut être écrit sous la forme suivante, où $a^* = a^*(\underline{\mu})$:

$$T(\underline{\mu})^{-1} = \sup_{\underline{v} \in \Sigma_K} \min_{a \neq a^*} v_{a^*} d(\mu_{a^*}, \bar{\mu}_{a^*, a, \underline{v}}) + \underbrace{v_a d(\mu_a, \bar{\mu}_{a^*, a, \underline{v}})}_{\stackrel{\text{def}}{=} \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})},$$

$$\text{où } \bar{\mu}_{a^*, a, \underline{v}} \stackrel{\text{def}}{=} \frac{v_{a^*} \mu_{a^*} + v_a \mu_a}{v_{a^*} + v_a}.$$

La quantité $\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})$ s'interprète comme un coût de transport représentant la difficulté de changer les distributions des bras a^* et a de sorte que a devienne optimal, étant données les fréquences de tirage \underline{v} . La preuve de l'unicité de $\underline{w}(\underline{\mu})$ assure que le vecteur de poids optimal égalise les coûts de transport :

$$\forall a \neq a^*, \quad \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w}(\underline{\mu})) = T(\underline{\mu})^{-1}.$$

Approximation du vecteur de poids optimal. Il est important de pouvoir approcher la valeur du vecteur de poids optimal $\underline{w}(\underline{\mu})$ et de comprendre sa dépendance vis-à-vis du paramètre $\underline{\mu}$. Cela peut notamment permettre de définir des stratégies (voir la présentation de Track-and-Stop dans la section suivante). Garivier and Kaufmann (2016) ont montré que le calcul de $\underline{w}(\underline{\mu})$ revenait à déterminer la racine d'une fonction croissante à une variable. On peut donc approcher $\underline{w}(\underline{\mu})$ avec une précision arbitraire en utilisant une méthode de dichotomie. Cette caractérisation de $\underline{w}(\underline{\mu})$ a également permis d'obtenir la continuité de $\underline{\mu} \mapsto \underline{w}(\underline{\mu})$. Cependant, d'autres caractérisations pourraient

permettre d'obtenir des méthodes d'approximation plus efficaces d'un point de vue computationnel et de meilleurs résultats de régularité. Dans cette thèse, nous travaillons sur ce problème d'optimisation dans le cas spécifique d'un modèle Gaussien au Chapitre 3 et, indépendamment, dans un cadre exponentiel général au Chapitre 6.

Caractérisation pour des variables Gaussiennes. Dans le Chapitre 3, on s'intéresse au problème d'optimisation pour des variables aléatoires Gaussiennes standards². Dans ce cas, le problème s'écrit simplement :

$$T(\underline{\mu})^{-1} = \sup_{v \in \Sigma_K} \frac{1}{2} \min_{a \neq a^*} \frac{v_{a^*} v_a}{v_{a^*} + v_a} \Delta_a^2,$$

où l'on note, tant qu'il n'y a pas d'ambiguïté, $\Delta_a = \Delta_a(\underline{\mu})$. Nous montrons que le vecteur de poids optimal $\underline{w}(\underline{\mu})$ est caractérisé par la racine de la fonction $\phi_{\underline{\mu}}$ définie par :

$$\forall r \in \left(\frac{1}{\Delta_{\min}^2}, +\infty \right), \quad \phi_{\underline{\mu}}(r) \stackrel{\text{def}}{=} \sum_{a \neq a^*} \frac{1}{(r \Delta_a^2 - 1)^2} - 1.$$

On a en effet la proposition suivante.

Proposition 1.2. [voir Proposition 3.2]

Soit $\underline{\mu}$ un problème de bandit dont les bras sont Gaussiens standards. Soient $\underline{\Delta} = \underline{\Delta}(\underline{\mu})$, $\underline{w} = \underline{w}(\underline{\mu})$, $T = T(\underline{\mu})$, et soit $r = r(\underline{\mu})$ la solution de $\phi_{\underline{\mu}}(r) = 0$. Alors

$$\begin{aligned} w_{a^*} &= \frac{1}{1 + \sum_{a \neq a^*} \frac{1}{r \Delta_a^2 - 1}}, \\ \forall a \neq a^*, \quad w_a &= \frac{w_{a^*}}{r \Delta_a^2 - 1}, \\ \text{et} \quad T &= 2 \frac{r}{w_{a^*}}. \end{aligned}$$

Cette caractérisation a plusieurs conséquences pour les modèles Gaussiens à variance commune.

1. Tout d'abord, la fonction $\phi_{\underline{\mu}}$ étant convexe et décroissante, l'application d'une méthode de Newton permet d'approcher r , et donc $\underline{w}(\underline{\mu})$ et $T(\underline{\mu})$, avec une vitesse de convergence quadratique des itérées. Cela accélère la méthode proposée initialement par [Garivier and Kaufmann \(2016\)](#). Voir Section 3.3 pour plus de détails.
2. Nous proposons également de nouvelles bornes pour les valeurs possibles de la coordonnée $w_{a^*}(\underline{\mu})$ du vecteur de poids optimal et du temps caractéristique $T(\underline{\mu})$:

$$\begin{aligned} \frac{1}{1 + \sqrt{K-1}} &\leq w_{a^*}(\underline{\mu}) \leq \frac{1}{2}, \\ \text{et} \quad \max \left(\frac{8}{\Delta_{\min}^2}, 4 \frac{1 + \sqrt{K-1}}{\Delta^2} \right) &\leq T(\underline{\mu}) \leq 2 \frac{(1 + \sqrt{K-1})^2}{\Delta_{\min}^2}, \end{aligned}$$

où $\overline{\Delta^2} \stackrel{\text{def}}{=} \frac{1}{K-1} \sum_{a \neq a^*} \Delta_a^2$ est la moyenne des carrés des gaps et $\Delta_{\min} \stackrel{\text{def}}{=} \min_{a \neq a^*} \Delta_a$ est le gap minimal. Toutes ces bornes sont atteintes par des problèmes spécifiques.

3. Nous établissons des résultats de monotonie lorsque l'on bouge l'un (ou plusieurs) des bras de $\underline{\mu}$. Par exemple, en augmentant la valeur de la moyenne d'un bras sous-optimal a (voir

²Ou plus généralement de variance commune $\sigma^2 > 0$.

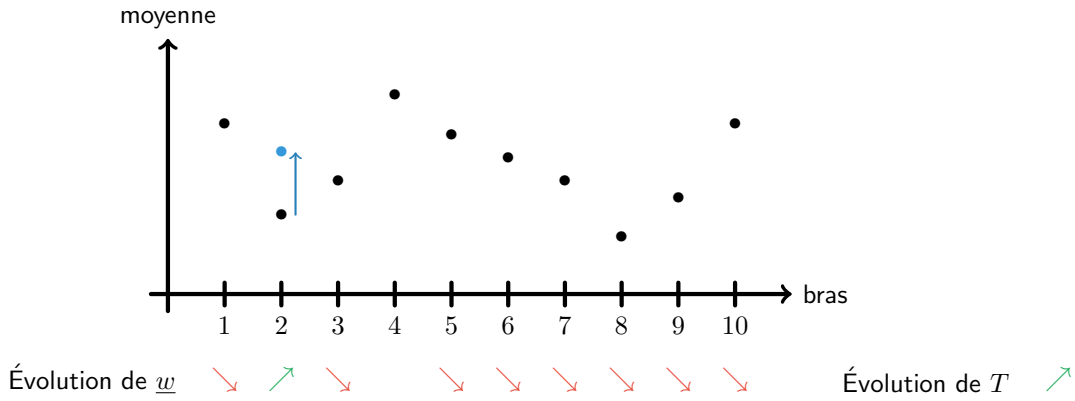


Figure 1.3: En **augmentant la moyenne d'un bras sous-optimal**, son poids optimal augmente tandis que les poids des autres bras sous-optimaux diminuent, et le temps caractéristique augmente.

Figure 1.3), on augmente son poids $w_a(\underline{\mu})$ tout en diminuant le poids des autres poids sous-optimaux. Comme le gap entre le bras a et le bras optimal diminue, le problème d'identification se complexifie et donc la valeur de $T(\underline{\mu})$ augmente. La Section 3.4 donne l'ensemble des résultats de monotonicit  obtenus, qui vont se r v ler importants pour la d finition de strat gies exploratrices comme la strat gie Exploration-Biased-Sampling propos e au Chapitre 4.

4. Enfin, nous d montrons des r sultats quantitatifs de r gularit  des solutions au probl me d'optimisation dans la Section 3.5. On montre notamment que les fonctions $\underline{\mu} \mapsto \underline{w}(\underline{\mu})$ et $\underline{\mu} \mapsto T(\underline{\mu})$ sont localement lipschitziennes :

Th or me 1.2. [voir Theorem 3.9]

Soient $\underline{\mu}$ et $\underline{\mu}'$ deux probl mes de bandit dont les bras sont Gaussiens standards, et de m me bras optimaux. Soient $\underline{\Delta}$, \underline{w} et T (respectivement $\underline{\Delta}'$, \underline{w}' et T') d finis comme dans la Proposition 1.2. Supposons que, pour un $\varepsilon \in [0, 1/7]$, on ait

$$\forall a \neq a^*, \quad (1 - \varepsilon)\Delta_a^2 \leq \Delta_a'^2 \leq (1 + \varepsilon)\Delta_a^2.$$

Alors

$$\forall a \in [K], \quad (1 - 10\varepsilon)w_a \leq w_a' \leq (1 + 10\varepsilon)w_a,$$

et $(1 - 3\varepsilon)T \leq T' \leq (1 + 6\varepsilon)T.$

Cela pr cise le r sultat de continuit  de $\underline{\mu} \mapsto \underline{w}(\underline{\mu})$ obtenu³ par [Garivier and Kaufmann \(2016\)](#). La r gularit  du probl me d'optimisation permet d'envisager de nouvelles techniques de preuves de bornes non asymptotiques, comme celles pr sent es au Chapitre 4.

Caract risation pour un mod le exponentiel quelconque. Les techniques d velopp es dans le Chapitre 3 utilisent la forme particuli rement simple de la divergence de Kullback-Leibler entre deux Gaussiennes, et il semble compliqu  de les g n raliser   tout mod le exponentiel. Dans la Section 6.2, nous donnons cependant une nouvelle caract risation pour un mod le exponentiel \mathcal{D}_{exp} quelconque. Celle-ci repose sur l'introduction d'une quantit  d finie pour $\underline{v} \in \text{int}(\Sigma_K)$, l'int rieur de Σ_K , par

$$T_{\underline{v}} \stackrel{\text{def}}{=} \sum_{b \neq a^*} \frac{1}{d(\mu_b, \bar{\mu}_{a^*}, b, \underline{v})}.$$

³Dans le cadre plus g n ral d'un mod le exponentiel.

La caractérisation de $\underline{w}(\underline{\mu})$ faite par [Garivier and Kaufmann \(2016\)](#) assure que $T_{\underline{w}} = T(\underline{\mu})$, et donc, étant donné un vecteur de poids \underline{v} , $T_{\underline{v}}$ est une approximation du temps caractéristique. Cela permet de définir une transformation $\underline{W} : \text{int}(\Sigma_K) \rightarrow \text{int}(\Sigma_K)$ par

$$\forall \underline{v} \in \text{int}(\Sigma_K), \forall a \in [K], \quad W_a(\underline{v}) \stackrel{\text{def}}{=} \begin{cases} v_a \frac{T_{\underline{v}}^{-1}}{\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})} & \text{si } a \neq a^*, \\ 1 - \sum_{b \neq a^*} v_b \frac{T_{\underline{v}}^{-1}}{\text{TC}_{b \rightarrow a^*}(\underline{\mu}, \underline{v})} & \text{si } a = a^*. \end{cases}$$

On peut alors vérifier que $\underline{w}(\underline{\mu})$ est l'unique point fixe de cette transformation.

Proposition 1.3. [voir Proposition 6.2]

Soit \mathcal{D}_{exp} un modèle exponentiel et $\underline{\mu}$ dans \mathcal{D}_{exp} . Alors \underline{W} a pour unique point fixe $\underline{w}(\underline{\mu})$.

Cette caractérisation a elle aussi des conséquences importantes.

1. Des algorithmes itératifs, utilisant la transformation \underline{W} , permettent de faire converger des séquences de vecteurs vers le vecteur de poids optimal. Si les résultats théoriques de convergence restent à établir, nous vérifions que cette convergence a lieu numériquement, quelque soit le problème de bandit et le modèle, avec une convergence assez rapide, qui pourrait là encore accélérer la méthode initiale de [Garivier and Kaufmann \(2016\)](#).
2. Aussi, une interprétation de la transformation \underline{W} mène naturellement à la définition d'un nouvel algorithme de type Top-Two adaptatif, que nous décrirons plus en détails en Section 1.2.5.

1.2.3. L'algorithme Track-and-Stop

Pour obtenir des garanties proches de la borne inférieure (1.1), nous avons vu qu'une stratégie doit observer les bras avec des fréquences proches du vecteur de poids optimal (inconnu) $\underline{w}(\underline{\mu})$. L'idée de l'algorithme Track-and-Stop de [Garivier and Kaufmann \(2016\)](#) est d'estimer ce vecteur de poids à chaque étape t en utilisant la moyenne empirique disponible $\hat{\underline{\mu}}(t-1)$.

Règle d'échantillonnage. À l'étape t , Track-and-Stop choisit donc d'observer le bras qui est le plus en retard par rapport au vecteur de poids $\hat{\underline{w}}(t-1) = \underline{w}(\hat{\underline{\mu}}(t-1))$, on parle de D-tracking, ou à une version cumulative $\frac{1}{t-1} \sum_{s \in [t-1]} \hat{w}_a(s)$ plus prudente, on parle de C-tracking :

$$A_t \in \begin{cases} \operatorname{argmin}_{a \in [K]} N_a(t-1) - \sum_{s \in [t-1]} \hat{w}_a(s) & \text{(C-tracking),} \\ \operatorname{argmin}_{a \in [K]} N_a(t-1) - (t-1)\hat{w}_a(t-1) & \text{(D-tracking).} \end{cases}$$

Le biais introduit par l'utilisation du vecteur $\hat{\underline{w}}(t-1)$ plutôt que $\underline{w}(\underline{\mu})$ peut conduire à un sous-échantillonnage de bras de moyennes assez élevées, notamment lorsque les premières estimations sont mauvaises. Pour y remédier, un mécanisme d'exploration forcée vient remplacer la règle d'échantillonnage lorsque l'un des bras n'a pas assez été observé : si le bras le moins observé a été échantillonné moins de \sqrt{t} fois (taux arbitraire), il est tiré automatiquement.

Règle d'arrêt et décision. Le choix des règles d'arrêt est aussi crucial : celles-ci doivent assurer que la stratégie est δ -correcte sans pour autant nécessiter trop d'observations, ce qui diminuerait la performance de la stratégie. Il semblerait qu'il y ait un consensus autour de la règle d'arrêt du Global-Likelihood-Ratio introduite également par [Garivier and Kaufmann \(2016\)](#), qui consiste à stopper l'algorithme dès lors que $Z(t)$ dépasse un certain seuil $\beta(t, \delta)$, où

$$Z(t) \stackrel{\text{def}}{=} \max_{a \in [K]} \min_{b \neq a} Z_{a,b}(t),$$

avec $Z_{a,b}(t)$ le ratio du log-likelihood sous les hypothèses $H_1 : \mu_a > \mu_b$ et $H_0 : \mu_a \leq \mu_b$, étant données les observations à la fin de l'étape t (voir Section 2.2.5 pour plus de détails).

Ils ont notamment établi que cette règle d'arrêt assurait que la stratégie est δ -correcte pour le seuil suivant, quelque soit la règle d'échantillonnage.

Théorème 1.3. [Garivier and Kaufmann, 2016, Proposition 12]

Soit un modèle exponentiel \mathcal{D}_{exp} . Soit $\delta \in (0, 1)$ et $\alpha > 1$. Il existe une constante $R = R(\alpha, K)$ telle que, quelque soit la règle d'échantillonnage, la règle d'arrêt du Global-Likelihood-Ratio avec le seuil

$$\beta(t, \delta) \stackrel{\text{def}}{=} \log \frac{Rt^\alpha}{\delta}, \quad (1.3)$$

et avec l'estimation du meilleur bras empirique, assure que la stratégie est δ -correcte.

S'il est encore possible de diminuer légèrement le seuil, il semblerait que l'on ne puisse pas améliorer sensiblement le temps d'arrêt en changeant la condition donnée par le Global-Likelihood-Ratio.

Optimalité asymptotique. L'algorithme Track-and-Stop est la première stratégie *asymptotiquement optimale*, pour laquelle les garanties lorsque le risque δ tend vers 0 correspondent avec la borne inférieure :

$$\lim_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_\delta]}{\log \frac{1}{\delta}} \leq T(\underline{\mu}).$$

Nous donnons une preuve de ce résultat en Section 6.4 (voir Theorem 6.17).

1.2.4. Des garanties non asymptotiques

L'existence d'une stratégie asymptotiquement optimale questionne sur l'obtention de bornes non asymptotiques. Dans le cadre du modèle $\mathcal{D}_{\mathcal{N}_1}^{[0,1]}$ des distributions Gaussiennes standards avec moyennes dans $[0, 1]$, nous définissons au Chapitre 4 une nouvelle stratégie appelée Exploration-Biased-Sampling, pour laquelle nous démontrons des bornes à risque fini.

Règle d'échantillonnage. L'idée de cette nouvelle stratégie est de modifier légèrement le fonctionnement de Track-and-Stop, en instaurant une région de confiance $\mathcal{CR}_{\underline{\mu}}(t-1)$ autour de $\hat{\underline{\mu}}(t-1)$ afin de définir un vecteur de poids qui favorise et maximise l'exploration, c'est-à-dire qui assure un taux d'observation minimal supérieur à celui donné par le vecteur de poids optimal $\underline{w}(\underline{\mu})$. Pour cela on définit

$$\tilde{\underline{\mu}}(t-1) \in \underset{\underline{\rho} \in \mathcal{CR}_{\underline{\mu}}(t-1)}{\operatorname{argmax}} w_{\min}(\underline{\rho}),$$

où $w_{\min}(\underline{\rho}) \stackrel{\text{def}}{=} \min_{a \in [K]} w_a(\underline{\rho})$ est la composante minimale de $\underline{w}(\underline{\rho})$. Ainsi, dès lors que le problème de bandit $\underline{\mu}$ appartient à la région de confiance $\mathcal{CR}_{\underline{\mu}}(t-1)$, utiliser le vecteur $\tilde{\underline{\mu}}(t-1)$ plutôt que $\underline{w}(\underline{\mu})$ augmente la fréquence minimale d'observation des bras.

La règle d'échantillonnage d'Exploration-Biased-Sampling consiste donc à suivre le poids $\tilde{\underline{w}}(t-1)$ de manière directe (D-tracking) ou cumulative (C-tracking) de la même manière que Track-and-Stop. Pour ce faire, il faut justifier que le vecteur $\tilde{\underline{w}}(t-1)$ est calculable. En utilisant les résultats de monotonie de la Section 3.4, nous montrons que cela est possible dès lors que la région de confiance est un produit d'intervalles de confiance. Nous choisissons

$$\mathcal{CR}_{\underline{\mu}}(t) \stackrel{\text{def}}{=} \prod_{a \in [K]} \left[\hat{\mu}_a(t) \pm C_{\frac{\gamma}{K}}(N_a(t)) \right], \quad \text{où} \quad C_\gamma(s) \stackrel{\text{def}}{=} 2 \sqrt{\frac{\log(\frac{4s}{\gamma})}{s}},$$

où $\gamma \in (0, 1)$ est un paramètre de l'algorithme fixé. À mesure que le nombre d'observations augmente, la région de confiance se rétrécit jusqu'à $\{\underline{\mu}\}$ et le biais d'exploration décroît, i.e., $\tilde{\underline{w}}(t) \rightarrow \underline{w}(\underline{\mu})$.

Borne non asymptotique. Les garanties théoriques d'Exploration-Biased-Sampling sont présentées dans la Section 4.3. Une première remarque intéressante est que cette stratégie a un taux naturel d'exploration : tous les bras sont tirés au moins de l'ordre de $\simeq \sqrt{t}$ fois après t étapes. Cette stratégie ne nécessite donc pas besoin d'exploration forcée (à l'inverse de Track-and-Stop). La stratégie est asymptotiquement optimale (comme Track-and-Stop), les poids estimés convergeant vers ceux de μ au fur et à mesure que les régions de confiance rétrécissent. Nous démontrons également des garanties non asymptotiques, i.e., à risque fini, lorsque nous conservons la règle d'arrêt Global-Likelihood-Ratio de Track-and-Stop.

Théorème 1.4. [voir Theorem 4.5]

Soient $\gamma \in (0, 1)$, $\alpha \in [1, 2]$, $\eta \in (0, 1]$, et soit $\underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}^{[0,1]}$. Il existe un évènement \mathcal{E} de probabilité au moins $1 - \gamma$ et $\delta_0 \stackrel{\text{def}}{=} \delta_0(\underline{\mu}, K, \gamma, \eta, \alpha) > 0$ tels que, pour tout $0 < \delta \leq \delta_0$, l'algorithme Exploration-Biased-Sampling avec le seuil (1.3) vérifie

$$\forall t > (1 + \eta)T(\underline{\mu}) \log \frac{1}{\delta}, \quad \mathbb{P}_{\underline{\mu}}(\tau_\delta > t \cap \mathcal{E}) \leq 2Kt \exp\left(-\frac{tw_{\min}(\underline{\mu})}{4T(\underline{\mu})^2} \frac{1}{\log^{\frac{2}{3}} \frac{1}{\delta}}\right),$$

et

$$\mathbb{E}_{\underline{\mu}}[\tau_\delta \mathbb{I}\{\mathcal{E}\}] \leq (1 + \eta)T(\underline{\mu}) \log \frac{1}{\delta} + \frac{2^7 K T(\underline{\mu})^4}{w_{\min}(\underline{\mu})^2} \exp\left(-\frac{w_{\min}(\underline{\mu})}{4T(\underline{\mu})} \log^{\frac{1}{3}} \frac{1}{\delta}\right) \log^2 \frac{1}{\delta}.$$

Pour obtenir ces garanties, nous utilisons notamment les résultats de régularité démontrés pour les poids optimaux en Section 3.5, couplés à la stabilité de la procédure de tracking.

Correction des lacunes de Track-and-Stop. La stratégie de tracking de Track-and-Stop possède quelques limitations. Outre la nécessité de recourir à une exploration forcée, le fort bruitage des observations lors des premiers tours peut rendre les proportions estimées très volatiles, et donc la procédure très aléatoire. De plus, la stratégie ne présente pas le comportement attendu de tirer les bras uniformément tant que trop peu d'informations ont été collectées. Les régions de confiance et la procédure introduite par Exploration-Biased-Sampling permettent de corriger ces défauts : l'algorithme commence par une phase d'uniforme exploration (de longueur variable), puis les poids évoluent de manière assez stable. En contrepartie, cela entraîne un fort biais en faveur de l'exploration qui, en pratique, se manifeste par un temps de prise de décision toujours plus long que Track-and-Stop (mais d'ordre de grandeur comparable même pour des valeurs modérées de risque).

1.2.5. Les algorithmes Top-Two

De nouvelles stratégies tentent de combler les faiblesses de Track-and-Stop, notamment du point de vue du coût computationnel. Track-and-Stop nécessite en effet le calcul d'un vecteur de poids optimal à chaque étape. Les algorithmes de type top-two introduits par Russo (2016) utilisent des règles d'échantillonnages simples : le bras à tirer est choisi entre un meneur L_t (un *leader*) et un concurrent C_t (un *challenger*). De nombreux choix sont possibles, citons notamment (Jourdan et al., 2022) le meilleur bras empirique pour le leader :

$$L_t \in \operatorname{argmax}_{a \in [K]} \hat{\mu}_a(t-1),$$

puis, pour le challenger, le bras dont le coût de transport empirique avec le leader est le plus faible :

$$C_t \in \operatorname{argmin}_{a \neq L_t} \operatorname{TC}_{a \rightarrow L_t} \left(\hat{\mu}(t-1), \frac{N(t-1)}{t-1} \right).$$

Avec un paramètre β non adaptatif. Initialement, ces algorithmes ont été étudiés lorsque le bras à observer est choisi comme étant le leader avec une probabilité $\beta \in (0, 1)$ fixe. Comme un bon choix de leader mènera à choisir la plupart du temps le meilleur bras comme leader, cela signifie que ces algorithmes tirent en moyenne le meilleur bras avec proportion β , ce qui ne garantit pas un choix optimal puisque a priori $\beta \neq w_{a^*}(\underline{\mu})$. Ainsi, les analyses actuelles montrent que certains algorithmes de type Top-Two sont β -asymptotiquement optimaux, car ils vérifient pour des modèles exponentiels :

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \leq T_{\beta}(\underline{\mu}), \quad \text{où} \quad T_{\beta}(\underline{\mu})^{-1} \stackrel{\text{def}}{=} \sup_{\substack{v \in \Sigma_K \\ v_{a^*} = \beta}} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a),$$

ce qui correspond à la borne inférieure asymptotique pour ces algorithmes. Si l'on peut montrer que des choix raisonnables de β (comme $\beta = \frac{1}{2}$) n'engendrent pas une perte de garantie importante, il n'est pas possible d'obtenir l'optimalité asymptotique des algorithmes Top-Two non adaptatifs.

Avec des proportions adaptatives. Une question naturelle consiste donc à comprendre si l'on peut apprendre, au cours des observations, les proportions idéales de tirages entre le leader et le challenger. La transformation \underline{W} présentée en Section 1.2.2 permet en fait de mettre la main sur de telles proportions, et donc de proposer des algorithmes adaptatifs. Nous montrons en effet que

$$\beta_{L,C} = \frac{w_L d(\mu_L, \bar{\mu}_{L,C,\underline{w}})}{\text{TC}_{C \rightarrow L}(\underline{\mu}, \underline{w})} = \frac{w_L d(\mu_L, \bar{\mu}_{L,C,\underline{w}})}{w_L d(\mu_L, \bar{\mu}_{L,C,\underline{w}}) + w_C d(\mu_C, \bar{\mu}_{L,C,\underline{w}})},$$

est la proportion optimale de tirages du leader L face au challenger C . Cette quantité peut être estimée en remplaçant \underline{w} par les proportions empiriques de tirages $\frac{N(t)}{t}$.

Pour un modèle Gaussien avec variances communes, [You et al. \(2023\)](#) ont montré que certains de ces algorithmes Top-Two adaptatifs étaient asymptotiquement optimaux. Cependant, l'analyse reste à établir pour des modèles exponentiels, pour lesquels nous conjecturons l'optimalité asymptotique.

Conjecture 1.4. [voir Conjecture 6.4]

Soit \mathcal{D}_{exp} un modèle exponentiel. L'algorithme de type top-two adaptatif utilisant le meilleur bras empirique comme leader, puis le challenger minimisant le coût de transport avec le leader, vérifie, pour tout problème de bandit $\underline{\mu}$ dans \mathcal{D}_{exp} avec moyennes distinctes :

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \leq T(\underline{\mu}).$$

Nous explorons cette conjecture dans la Section 6.3. Les résultats établis dans un cadre Gaussien par [You et al. \(2023\)](#) donnent une feuille de route, et nous expliquons comment adapter leurs arguments dans le cadre général. Deux difficultés majeures sont identifiées : d'une part la propriété d'exploration suffisante (le fait que les bras soient tirés au moins de l'ordre de \sqrt{t} à l'étape t), et d'autre part une relation limite vérifiée par les fréquences empiriques. Si nous conjecturons, en donnant quelques intuitions, que l'exploration suffisante est également vérifiée dans le cadre général (ou peut être forcée si ce n'est pas le cas), la difficulté majeure semble désormais d'obtenir une relation limite similaire au cas Gaussien.

1.3. Identification de meilleur bras à budget fixé

Nous explorons ensuite le cadre à budget fixé, où la stratégie retourne son estimation \hat{a}_T du bras optimal après un nombre déterministe d'observations T . L'objectif est de trouver des stratégies qui minimisent la probabilité d'erreur :

$$\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})).$$

Une décroissance exponentielle. En utilisant l'inégalité de Hoeffding avec des modèles sous-gaussiens, on peut vérifier que la probabilité d'erreur converge vers 0 à vitesse exponentielle pour la stratégie qui tire tous les bras uniformément et retourne le meilleur bras empirique :

$$\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \leq (K - 1) \exp(C(\underline{\nu})T),$$

pour une constante $C(\underline{\nu}) < 0$. On s'intéresse alors à borner le taux de décroissance :

$$\frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})).$$

D'une part, on souhaite obtenir des bornes inférieures, valides pour des suites de stratégies "raisonnables", des quantités

$$\ell(\underline{\nu}) \stackrel{\text{def}}{=} \liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})).$$

D'autre part, pour une suite de stratégies fixée, on cherche une borne supérieure sur la quantité

$$u(\underline{\nu}) \stackrel{\text{def}}{=} \limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})).$$

Notons que $\ell(\underline{\nu}) \leq u(\underline{\nu})$ et que ces taux sont négatifs.

Nous allons voir que le cadre à budget fixé est pour l'instant moins bien compris que celui à confiance fixée : les bornes que l'on connaît sur $\ell(\underline{\nu})$ et $u(\underline{\nu})$ ne permettent pas d'identifier un taux de décroissance précis, et donc une notion de complexité optimale comme en confiance fixée.

1.3.1. Bornes connues de la littérature

Borne supérieure : l'algorithme Successive-Rejects. L'une des premières stratégies étudiée à budget fixé est un algorithme d'élimination proposé par [Audibert et al. \(2010\)](#) nommé Successive-Rejects. La stratégie consiste à diviser l'exploration en $K - 1$ tours. Une liste de bras candidats est initialisée à l'ensemble des bras, et à chaque tour l'ensemble des bras restants sont tirés uniformément, puis le pire bras empirique est éliminé des candidats. [Audibert et al. \(2010\)](#) ont considérées les longueurs de phases, soigneusement choisies et déterministes (non adaptatives), suivantes :

$$\ell_1 \stackrel{\text{def}}{=} \frac{T}{\overline{\log K}}, \quad \text{et} \quad \forall r \in \{2, \dots, K - 1\}, \quad \ell_r \stackrel{\text{def}}{=} \frac{T}{(K - r + 2) \overline{\log K}}, \quad (1.4)$$

où

$$\overline{\log K} \stackrel{\text{def}}{=} \frac{1}{2} + \sum_{k=2}^K \frac{1}{k},$$

Ils ont démontré une borne supérieure sur le taux exponentiel limite de décroissance impliquant les gaps. Pour tout problème de bandit $\underline{\nu}$ avec un unique bras optimal dans un modèle σ^2 -sous-Gaussien,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{1}{4\sigma^2 \overline{\log K}} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}, \quad (1.5)$$

où $\overline{\log K}$ est de l'ordre de $\log K$, et les gaps sont ordonnés selon $0 = \Delta_{(1)} < \Delta_{(2)} \leq \dots \leq \Delta_{(K)}$.

Bornes inférieures. La borne supérieure (1.5) dépend cruciallement du choix des longueurs de phases, qui a en fait été guidé par le souhait d'obtenir un même terme de complexité que celui obtenu dans une borne inférieure pour un modèle de Bernoulli avec des moyennes dans $[p, 1 - p]$. Une preuve technique leur a permis de montrer que, intuitivement, pour tout problème de bandit $\underline{\nu}$ avec un unique bras optimal :

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq -\frac{5}{p(1-p)} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}. \quad (1.6)$$

Le choix des longueurs de phases de *Successive-Rejects* a ainsi permis d'obtenir la même quantité $\min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}$ dans la borne supérieure, au prix de l'apparition d'un facteur supplémentaire $\overline{\log K}$ dû à la contrainte de budget. Autrement dit, les deux bornes ne correspondent pas.

Si d'autres bornes inférieures ont été proposées dans le cas très spécifique d'un modèle Gaussien (pour lequel la divergence de Kullback-Leibler est symétrique), ou dans le cas de $K = 2$ bras (voir [Kaufmann et al., 2016](#)), les bornes inférieures et supérieures connues aujourd'hui ne correspondent pas en toute généralité. Pour résumer, la compréhension du problème à budget fixé est donc bien moindre qu'à confiance fixée :

- les bornes supérieure et inférieure ne correspondent pas, avec un écart multiplicatif de l'ordre de $\log K$,
- les bornes existantes sont formulées uniquement à partir des gaps entre les distributions (et non pas de quantités plus précises comme la divergence de Kullback-Leibler),
- les bornes (inférieures) sont valables pour des modèles très spécifiques (Gaussiens et Bernoulli).

1.3.2. Généralisations des bornes existantes à des modèles quelconques

Si le premier point de la liste ci-dessus semble délicat à traiter (voir Section 1.3.3), les techniques de preuves connues peuvent être généralisées à de nombreux modèles (incluant par exemple des modèles exponentiels et non paramétriques), comme nous le montrons au Chapitre 5. Les généralisations de ces bornes font intervenir de nouvelles quantités de théorie de l'information, plus informatives que les gaps, et l'obtention des bornes inférieures se fait grâce à des hypothèses naturelles sur le comportement attendu d'une bonne stratégie.

De nouvelles quantités de théorie de l'information. La borne inférieure de théorie d'information donnée par le Lemme 1.1 peut également être utilisée dans le cadre à budget fixé. Cependant, pour obtenir des bornes sur le taux exponentiel de décroissance, les rôles de $\underline{\nu}$ et de l'alternative $\underline{\zeta}$ doivent être inversés, i.e., on utilise l'inégalité

$$\sum_{a \in [K]} \mathbb{E}_{\underline{\zeta}}[N_a(T)] \text{KL}(\zeta_a, \nu_a) \geq \text{kl}(\mathbb{P}_{\underline{\nu}}(E), \mathbb{P}_{\underline{\zeta}}(E)).$$

pour tout évènement $E \in \mathcal{F}_T$. Si en confiance fixée, le temps caractéristique s'exprimait en fonction des $(\text{KL}(\nu_a, \zeta_a))_{a \in [K]}$, comme par exemple dans (1.2), il semblerait en conséquence que les bornes sur les taux de décroissance à budget fixé s'expriment plutôt en fonction des $(\text{KL}(\zeta_a, \nu_a))_{a \in [K]}$, autrement dit les arguments des divergences de Kullback-Leibler sont inversés.

Pour un modèle général \mathcal{D} , on introduit alors, pour $\nu \in \mathcal{D}$ et $x \in \mathbb{R}$, les quantités

$$\begin{aligned} \mathcal{L}_{\inf}^<(x, \nu) &= \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) < x \}, \\ \text{et } \mathcal{L}_{\inf}^>(x, \nu) &= \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) > x \}. \end{aligned}$$

Nous allons montrer que ces quantités interviennent dans des généralisations des bornes existantes en remplacement des gaps. Elles pourraient ainsi être la clé pour mesurer la complexité du cadre à budget fixé.

1.3. IDENTIFICATION DE MEILLEUR BRAS À BUDGET FIXÉ

Remarque 1.5. Une application de l'inégalité de Pinsker permet de montrer, lorsque les distributions du modèle sont à valeurs dans $[0, 1]$, que

$$\forall x \leq \mathbb{E}(\nu), \quad \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \geq 2(x - \mathbb{E}(\nu))^2, \quad \text{et} \quad \forall x \geq \mathbb{E}(\nu), \quad \mathcal{L}_{\text{inf}}^{\geq}(x, \nu) \geq 2(x - \mathbb{E}(\nu))^2.$$

Dans le cas d'un modèle exponentiel \mathcal{D}_{exp} , ces quantités s'écrivent simplement sous la forme de la divergence de Kullback-Leibler d du modèle :

$$\forall x \leq \mathbb{E}(\nu), \quad \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = d(x, \mathbb{E}(\nu)), \quad \text{et} \quad \forall x \geq \mathbb{E}(\nu), \quad \mathcal{L}_{\text{inf}}^{\geq}(x, \nu) = d(x, \mathbb{E}(\nu)).$$

Une analyse plus précise de Successive-Rejects. Nous proposons une nouvelle analyse de Successive-Rejects qui consiste à remplacer l'application de l'inégalité de Hoeffding dans l'analyse de Audibert et al. (2010) par une simple application de la borne de Cramér-Chernoff. Il en résulte le remplacement des gaps dans la borne (1.5) par une quantité dépendant des transformées de Fenchel-Legendre $(\phi_{\nu_a}^*)_{a \in [K]}$ des log-fonctions génératrices des moments des distributions de $\underline{\nu}$.

Remarque. On définit, pour une distribution ν ,

$$\forall \lambda \in \mathbb{R}, \quad \phi_{\nu}(\lambda) = \log \int_{\mathbb{R}} e^{\lambda x} d\nu(x) \quad \text{et} \quad \forall x \in \mathbb{R}, \quad \phi_{\nu}^*(x) = \sup_{\lambda \in \mathbb{R}} \{\lambda x - \phi_{\nu}(\lambda)\},$$

et l'on rappelle que ϕ_{ν}^* intervient dans le contrôle des déviations d'un N -échantillon de loi ν par la borne de Cramér-Chernoff :

$$\begin{aligned} \forall x \leq \mathbb{E}(\nu), \quad & \mathbb{P}(\bar{X}_N \leq x) \leq \exp(-N \phi_{\nu}^*(x)), \\ \text{et} \quad \forall x \geq \mathbb{E}(\nu), \quad & \mathbb{P}(\bar{X}_N \geq x) \leq \exp(-N \phi_{\nu}^*(x)). \end{aligned}$$

On obtient le résultat suivant, où l'on définit, pour tout $\nu, \nu' \in \mathcal{D}$ avec $\mathbb{E}(\nu') < \mathbb{E}(\nu)$,

$$\Phi(\nu', \nu) \stackrel{\text{def}}{=} \inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \{\phi_{\nu'}^*(x) + \phi_{\nu}^*(x)\}.$$

Proposition 1.6. [voir Corollary 5.4]

Soit un modèle \mathcal{D} . La séquence de stratégies Successive-Rejects utilisant les longueurs de phase (1.4) vérifie, pour tout problème de bandit $\underline{\nu}$ dans \mathcal{D} avec un unique bras optimal,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\log K} \min_{2 \leq k \leq K} \frac{\Phi(\nu_{\sigma_k}, \nu^*)}{k}. \quad (1.7)$$

où l'on range les bras a par ordre croissant de $\Phi(\nu_a, \nu^*)$, i.e., on considère la permutation σ telle que

$$0 = \Phi(\nu_{\sigma_1}, \nu^*) < \Phi(\nu_{\sigma_2}, \nu^*) \leq \dots \leq \Phi(\nu_{\sigma_{K-1}}, \nu^*) \leq \Phi(\nu_{\sigma_K}, \nu^*).$$

Sous des conditions de régularité assez légères du modèle \mathcal{D} considéré (voir Section 5.4.2), nous montrons que la quantité $\Phi(\nu', \nu)$ est égale à la quantité

$$\mathcal{L}(\nu', \nu) \stackrel{\text{def}}{=} \inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \left\{ \mathcal{L}_{\text{inf}}^{\geq}(x, \nu') + \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \right\}, \quad (1.8)$$

ce qui permet de réécrire la borne supérieure (1.7) en fonction des quantités $\mathcal{L}_{\text{inf}}^{\leq}$ et $\mathcal{L}_{\text{inf}}^{\geq}$:

$$\forall \underline{\nu} \text{ dans } \mathcal{D}, \quad \limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\log K} \min_{2 \leq k \leq K} \frac{\mathcal{L}(\nu_{\sigma_k}, \nu^*)}{k}. \quad (1.9)$$

Pour un modèle sous-Gaussien, cela améliore la borne supérieure (1.5) puisque l'on peut prouver, via l'inégalité de Pinsker, que $\mathcal{L}(\nu', \nu) \geq \frac{1}{4\sigma^2} (\mathbb{E}(\nu) - \mathbb{E}(\nu'))^2$.

De nouvelles bornes inférieures. Nous montrons également plusieurs bornes inférieures, qui dépendent de différentes hypothèses sur les suites de stratégies considérées. Une hypothèse générale est la *consistence*, c'est-à-dire le fait que l'identification a lieu asymptotiquement presque sûrement quelque soit le problème de bandit, i.e., que $\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \xrightarrow{T \rightarrow +\infty} 0$ pour tout $\underline{\nu}$ dans \mathcal{D} .

Une première borne est obtenue en faisant les hypothèses supplémentaires suivantes :

- *équilibre pour le pire bras* : le bras de plus petite espérance est tiré avec fréquence au plus $\frac{1}{K}$ asymptotiquement,
- *exploitation astucieuse de l'élagage des bras sous-optimaux* : la suite de stratégies, doublement indexée par K et T , est plus efficace pour identifier le meilleur bras si l'on enlève un des bras sous-optimaux du problème.

Sous ces deux hypothèses, nous montrons la borne suivante.

Théorème 1.5. [voir Theorem 5.10]

Soit un modèle \mathcal{D} . Soit une suite de stratégies, doublement indexée, qui est consistante, équilibrée pour le pire bras, et qui exploite astucieusement l'élagage des bras sous-optimaux. Alors, pour tout problème de bandit $\underline{\nu}$ dans \mathcal{D} avec des moyennes distinctes,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \frac{\mathcal{L}_{\text{inf}}^<(\mu_{(k)}, \nu^*)}{k}, \quad (1.10)$$

où les bras sont ordonnés selon $\mu_{(1)} > \mu_{(2)} > \dots > \mu_{(K)}$.

Ce résultat généralise la borne inférieure (1.6) de Audibert et al. (2010) à tout modèle \mathcal{D} . Dans le cas particulier d'un modèle de Bernoulli, l'inégalité de Pinsker permet en effet (cf. Remarque 1.5) de déduire (1.6) de la borne (1.10).

Cette borne inférieure ne s'exprime pas en termes d'infima de combinaisons de $\mathcal{L}_{\text{inf}}^>$ et $\mathcal{L}_{\text{inf}}^<$, i.e., sous la forme \mathcal{L} introduite en (1.8) et que l'on retrouve dans la borne supérieure de Successive-Rejects (1.9). Sans parvenir à obtenir une borne inférieure impliquant la quantité \mathcal{L} , nous proposons tout de même une borne inférieure pour des modèles dits *normaux*⁴ et des suites de stratégies dites *monotones*, qui limitent (asymptotiquement) les fréquences de tirages des bras sous-optimaux :

$$\forall \underline{\nu} \text{ dans } \mathcal{D}, \forall a \in [K], \quad \limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\nu}}[N_{(a)}(T)]}{T} \leq \frac{1}{a}.$$

Théorème 1.6. [Voir Theorem 5.13]

Considérons un modèle normal \mathcal{D} . Soit une séquence de stratégies consistante et monotone sur \mathcal{D} . Alors, pour tout problème de bandit $\underline{\nu}$ dans \mathcal{D} avec des moyennes distinctes,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \min_{2 \leq j \leq k} \inf_{x \in [\mu_{(j)}, \mu_{(j-1)})} \left\{ \frac{\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})}{j-1} + \frac{\mathcal{L}_{\text{inf}}^<(x, \nu^*)}{j} \right\}.$$

Cette borne implique notamment la borne plus lisible suivante

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \inf_{x \in [\mu_{(k)}, \mu_{(k-1)})} \left\{ \frac{\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})}{k-1} + \frac{\mathcal{L}_{\text{inf}}^<(x, \nu^*)}{k} \right\}, \quad (1.11)$$

⁴Qui nécessitent uniquement que l'infimum définissant $\mathcal{L}_{\text{inf}}^>(x, \nu)$ soit atteint par une suite de distributions dont les espérances tendent vers x .

1.3. IDENTIFICATION DE MEILLEUR BRAS À BUDGET FIXÉ

et, de surcroît, la borne inférieure (1.10) en prenant $x = \mu_{(k)}$. La borne (1.11) va dans la direction du terme de complexité qui apparaît dans la borne supérieure (1.9) avec la quantité \mathcal{L} , cependant l'infima est pris sur de plus petits intervalles.

Remarque. Les bornes présentées au Chapitre 5 n'améliorent pas l'écart entre les bornes inférieure et supérieure : il y a toujours (au moins) une constante multiplicative de l'ordre de $\overline{\log K} \simeq \log K$ entre les deux.

1.3.3. Existence d'une complexité

La différence de facteur de l'ordre $\log K$ entre les bornes inférieure (1.6) et supérieure (1.5) amène à se demander s'il existe, comme en confiance fixée, une notion de complexité dans le cadre à budget fixé. Si des résultats minimax comme la borne inférieure de [Carpentier and Locatelli \(2016\)](#) montrent que les bornes inférieures ne sont peut-être pas encore optimales, on se demande également s'il existe une stratégie *uniformément* optimale, qui puisse atteindre une borne inférieure asymptotique simultanément sur l'ensemble des problèmes de bandits.

Des travaux récents tentent de démontrer qu'une telle stratégie n'existe pas pour des modèles assez vastes. [Degenne \(2023\)](#) a notamment démontré que, pour un modèle Gaussien, il n'existe pas de stratégie qui soit, uniformément sur chaque problème de bandit, aussi performante que la meilleure stratégie tirant les bras avec des proportions fixes contre ce problème de bandit. Ce comportement, très différent du cadre à confiance fixée, pourrait s'expliquer par la difficulté que peuvent avoir les stratégies à se fier aux observations dans ce cadre.

Notation

This symbol in the external margin highlights an important notation or convention for the remaining of the thesis.



All tuples (vectors of \mathbb{R}^K , list of distributions, ...) are underlined.

Remark. For the sake of clarity, some technical notations of Chapters 3 to 5 are not aggregated here. For the same reasons, most of the pieces of notation of Chapter 6 are not included, as the chapter presents work-in-progress ideas with abundant notation.

Generalities

This first table contains general notation. Each piece of notation is associated with a reference to the page where it is defined or where it appears first.

Notation	Designation	Page
K	Positive integer (number of arms)	32
$\underline{\nu}$	Vector of \mathbb{R}^K	44
$[K]$	Set $\{1, \dots, K\}$	33
$\text{int}(A)$	Interior of a set A	48
sgn	Sign function	50
Σ_K	Simplex of dimension $K - 1$	44
\mathfrak{S}_K	Set of permutations of $[K]$	62
$\overline{\log K}$	$\frac{1}{2} + \sum_{k=2}^K \frac{1}{k}$	61
Probabilities		
$X \sim \nu$	The random variable X has law ν	33
$\mathbb{E}(\nu)$	Expectation of a distribution ν	33
$\mathbb{E}[X]$	Expectation of a random variable X	33
$\mathbb{I}\{E\}$	Indicator function of an event E	35
\mathbb{P}^X	Push-forward measure of a random variable X under probability \mathbb{P}	40
$\mathbb{P} \ll \mathbb{Q}$	\mathbb{Q} absolutely dominates \mathbb{P}	42
$\text{Leb}_{[0,1]}$	Lebesgue measure restricted to $[0, 1]$	41
$\text{Supp}(\nu)$	Support of distribution ν	145
$m(\nu)$	$\inf \text{Supp}(\nu)$	145
$M(\nu)$	$\sup \text{Supp}(\nu)$	145
δ_x	Dirac mass at x	146

Bandit problems

Letters a, b, c always refer to arms, that is elements of $[K]$. In subindices for sums and infima, or with quantifiers, we sometimes omit to explicitly mention $[K]$ for simplicity: for example, given a fixed arm b , $\sum_{a \neq b}$ denotes the sum over arms $a \in [K] \setminus \{b\}$.

Notation	Designation	Page
K	Number of arms	32
ν_a	Probability distribution of arm a	32
μ_a	Mean of probability distribution of arm a	33
$\underline{\nu}$	Bandit problem, set of K distributions (ν_1, \dots, ν_K)	33
$\underline{\mu}$	Mean vector of bandit problem $\underline{\nu}$	44
$a^*(\underline{\nu})$	Set of best arms of $\underline{\nu}$, or unique best arm of $\underline{\nu}$	33
$a^*(\underline{\mu})$		44
ν^*	Optimal distribution of $\underline{\nu}$	33
μ^*	Optimal mean of $\underline{\nu}$	33
$\mathbb{P}_{\underline{\nu}}$	Probability under bandit problem $\underline{\nu}$	35
$\mathbb{P}_{\underline{\mu}}$		44
$\mathbb{E}_{\underline{\nu}}$	Expectation under bandit problem $\underline{\nu}$	35
$\mathbb{E}_{\underline{\mu}}$		44
$\Delta_a, \Delta_a(\underline{\nu})$	Gap of arm a in $\underline{\nu}$	36
Δ_{\min}	Minimal positive gap of $\underline{\nu}$	39
Δ_{\max}	Maximal gap of $\underline{\nu}$	74
$\overline{\Delta}^2$	Average squared gap of $\underline{\nu}$	48
Alternative		
ζ	Another bandit problem (often an alternative to $\underline{\nu}$)	40
λ	Mean of ζ	44
$\text{Alt}(\underline{\nu})$	Set of alternative bandits problems to $\underline{\nu}$	66
$\text{Alt}(\underline{\mu})$		44
Optimization problem		
$T(\underline{\mu})$	Characteristic time of $\underline{\mu}$	44
$T_\beta(\underline{\mu})$	Characteristic time of $\underline{\mu}$ under the constraint that the best arm is pulled a fraction β of times	57
$w, w(\underline{\mu})$	Optimal weight vector of $\underline{\mu}$	47
$w_{\min}, w_{\min}(\underline{\mu})$	Minimal component of the optimal weight vector of $\underline{\mu}$	74
$\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{\nu})$	Transportation cost of arm a to a^* in $\underline{\nu}$, given draw proportions $\underline{\nu}$	46
$\bar{\mu}_{a^*, a, \underline{\nu}}$	Weighted mean of μ^* and μ_a with weights ν_{a^*} and ν_a	46
(a)	a -th arm of $\underline{\nu}$ in terms of means	37
$\underline{\nu}^\sigma$	Permuted bandit instance associated to $\underline{\nu}$ and permutation σ	57
$H(\underline{\nu})$	Complexity function	69
$H_\Sigma(\underline{\nu})$	$\sum_{a \neq a^*(\underline{\nu})} \frac{1}{\Delta_a^2}$	64

Strategies

We use a special font when referring to algorithms and procedures, like for instance Track-and-Stop or sampling-rule. Algorithm names start with a capital letter, while general procedures do not.

Notation	Designation	Page
δ	Confidence parameter	33
T	Budget parameter	34
A_t	Observed arm at time step t	35
Y_t	Observation at time step t	35
U_t	External randomization at time step t	35
I_t	Available history at time step t	35
\mathcal{F}_t	σ -algebra generated by I_{t+1}	35
$N_a(t)$	Number of observations of arm a at time step t	35
$\hat{\mu}_a(t)$	Empirical mean of arm a at time step t	35
$X_{a,n}$	n -th observation of arm a	37
$\hat{\mu}_{a,n}$	Mean of the n first observations of arm a	37
Stopping and decision rules		
τ	Stopping time of a δ -correct strategy	35
τ_δ		35
$\hat{a}_\tau, \hat{a}_{\tau_\delta}, \hat{a}_T$	Estimated best arm	35
$Z_{a,b}(t)$	Generalized log-likelihood ratio	50
$Z(t)$	Test statistic for the Global-Likelihood-Ratio stopping rule	50
$\beta(t, \delta)$	Threshold function	50
R	Constant ensuring δ -correctness of threshold (2.27)	51
α	Constant in threshold (2.27)	51
Elimination strategies		
r	Round	36
S_{r-1}	Set of candidates arms remaining at the beginning of round r	36
a_r	Arm removed at the end of round r	36
n_r	Number of pulls of the arm eliminated at the end of phase r	38
ℓ_r	Length of round r	60
Track-and-Stop		
$\hat{w}(t)$	Plug-in estimate of $w(\underline{\mu})$ at time step t	49
$U(t)$	Set of under-sampled arms at time step t	49
Exploration-Biased-Sampling		
$\mathcal{CR}_\mu(t)$	Confidence region for $\underline{\mu}$ after t observations	101
$C_\gamma(s)$	Constant defining the width of $\mathcal{CR}_\mu(t)$	101
$\tilde{\mu}(t)$	Biased exploring bandit inside $\mathcal{CR}_\mu(t)$	102
$\tilde{w}(t)$	Optimal weight of $\tilde{\mu}(t)$	102
top-two algorithms		
L_t	Leader at time step t	55
C_t	Challenger at time step t	55
β	Probability parameter of choosing the leader in a top-two non-adaptive algorithm	57

Information-theoretic quantities

Notation	Designation	Page
$\text{KL}(\mathbb{P}, \mathbb{Q})$	Kullback-Leibler divergence between probability measures \mathbb{P} and \mathbb{Q}	40
d	Mean-parametrized Kullback-Leibler divergence of an exponential family	43
kl	Mean-parametrized Kullback-Leibler divergence of the Bernoulli family	44
$\mathcal{L}_{\text{inf}}^<(x, \nu)$	$\inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) < x\}$	67
$\mathcal{L}_{\text{inf}}^{\leq}(x, \nu)$	$\inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) \leq x\}$	126
$\mathcal{L}_{\text{inf}}^>(x, \nu)$	$\inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) > x\}$	67
$\mathcal{L}_{\text{inf}}^{\geq}(x, \nu)$	$\inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) \geq x\}$	126
$\mathcal{L}(\nu, \nu')$	$\inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \left\{ \mathcal{L}_{\text{inf}}^{\geq}(x, \nu') + \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \right\}$	67
$\mathcal{K}_{\text{inf}}^<(\nu, x)$	$\inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) < x\}$	46
$\mathcal{K}_{\text{inf}}^>(\nu, x)$	$\inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) > x\}$	46

Models of distributions

Notation	Designation	Page
\mathcal{D}	General set of distributions all admitting a finite mean	33
$\mathcal{P}[0, 1]$	Distributions taking values in $[0, 1]$	33
$\mathcal{D}_{\mathcal{B}}$	Bernoulli distributions	33
$\mathcal{B}_{[p, 1-p]}$	Bernoulli distributions with means in $[p, 1-p]$	64
$\mathcal{D}_{\mathcal{N}_{\sigma^2}}$	Gaussian distributions with known variance $\sigma^2 > 0$	33
$\mathcal{D}_{\mathcal{N}_{\sigma^2}}^{\mathcal{M}}$	Distributions of $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ with means in \mathcal{M}	96
$\mathcal{D}_{\mathcal{N}_1}^{[0,1]}$	Standard Gaussian distributions with mean in $[0, 1]$	89
\mathcal{D}_{σ^2}	σ^2 -sub-Gaussian distributions	33

Exponential models

\mathcal{D}_{exp}	One-canonical exponential family	33
Θ	Natural space parameter	43
θ	Parameter	43
b	Normalizing function	43
$\mathcal{M} = (\mu_-, \mu_+)$	Open interval of the expectations of distributions	43
d	Mean-parametrized Kullback-Leibler divergence	43
kl	Mean-parametrized Kullback-Leibler divergence of the Bernoulli family	44
$C_{\mathcal{D}}$	Regularity constant associated to a model \mathcal{D}	65

CHAPTER 2

Pure Exploration in Multi-Armed Bandits

In this preliminary chapter, we first introduce and motivate the statistical problems considered in this thesis, which constitute a sub-domain of the multi-armed bandit literature known as best-arm identification. Then we review the state-of-the-art theoretical knowledge around those problems and highlight our contributions (to be all presented in detail in Chapters 3 to 6).

In the fixed-confidence setting, we explain that the problem is well-understood in the asymptotic regime and we present new results for a Gaussian model which allow us to define a new strategy with non-asymptotic guarantees. We also present partial results for the asymptotic analysis of adaptive top-two algorithms.

Then, we highlight the difficulties faced by the literature to obtain a precise comprehension of the fixed-budget setting. In that setting, we present generalizations of existing bounds to possibly non-parametric models. Those bounds are based on new information-theoretic quantities.

Contents

1	A Sequential Learning Problem	32
2	Best-Arm Identification with a Fixed-Confidence	35
1	Naive Approaches	36
2	The Fundamental Inequality	40
3	Lower Bound for Exponential Models	43
4	The Sample Complexity Optimization Problem	45
5	The Track-and-Stop Algorithm: an Asymptotically Optimal Strategy	49
6	Non-Asymptotic Guarantees	53
7	Towards Computationally More Efficient and More Natural Strategies	54
3	Best-Arm Identification with a Fixed-Budget	58
1	An Exponential Decay Rate	58
2	Successive-Rejects-Type Strategies	59
3	Lower Bounds	62
4	Comparing Upper and Lower Bounds: the Challenges of the Fixed-Budget Setting	66
5	On Non-Matching Bounds and Minimax Results	68

2.1. A Sequential Learning Problem

In this thesis, we study the following sequential learning task. A learner faces a set of $K \geq 2$ unknown real probability distributions ν_1, \dots, ν_K . At each time step, the learner observes an independent realization of one of the distributions, that she actively chooses according to previous observations. Her objective is to quickly gather evidence so as to identify the distribution of highest mean. This statistical task is a sequential and active optimization problem, which are sub-domains of machine learning.

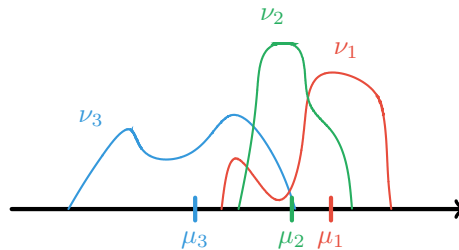


Figure 2.1: μ_a denotes the expectation of ν_a . For the three distributions for which we represent the probability mass functions, the goal of the learner is to identify ν_1 as the best distribution.

Motivations. Before going any further, we give two concrete applications of this learning task.

Application 2.1. [Medical trial]

Assume that you get K candidate treatments to cure a disease. You organize a medical trial to find the treatment that responds best to the disease¹. You can allocate one treatment to each participant of the trial and observe how she responds to this treatment. You may proceed sequentially so as to choose the next treatment to allocate based on past observations.

Application 2.2. [Online advertising, A/B testing]

You own a webpage that you want to monetize by including an advertisement in it. Your income will be proportional to the number of clicks on the ad. Given K possible ads, you want to find the ad which maximizes the commitment of users. To do so, for each new viewer, you can choose an ad to display, based on previously collected data, and check whether or not he clicks on it.

A famous variant of this problem is A/B testing: you want to find which one among two versions of a webpage maximizes customer satisfaction.

First intuitions on the learner's strategy. The learner's strategy has an impact on the collected data: she chooses the next distribution to observe with the knowledge of previous observations. This is why we say that the learner is *active* and we introduced the problem as a *sequential* learning task. It is in contrast with more "classical" statistical problems for which we cannot influence the data collected. In our problem, how the learner relies on past observations might evolve over time, as intuition suggests:

- during the first steps, the learner has access to a limited amount of observations of each distribution, hence she cannot estimate precisely the associated means. Consequently, it sounds natural to mistrust those observations and proceed to a uniform exploration of the distributions (i.e., observe all of them equally often),
- as more and more observations are gathered, the accuracy of the estimations increases, and the learner should allocate more observations to empirically promising distributions in order to better differentiate them.

¹In a meaning that needs to be defined: a good treatment should cure the disease while avoiding side effects.

This highlights a major difficulty when dealing with those identification problems: how much can the learner trust previous observations? Does she need to focus on under-sampled distributions (for which she has poor estimates) or on promising distributions (which behave well empirically)?

Multi-armed bandits. The problem introduced in this section is known in the literature as the problem of *best-arm identification* in *multi-armed bandits*. Although the first motivations for those questions concerned medical applications (see [Thompson, 1933](#)), this vocabulary comes from casinos and the slot machines sometimes called “one-armed bandits”. Considering a set of K machines, a player might try to identify the machine with the best average payoff.

The problem instance $\underline{\nu} = (\nu_1, \dots, \nu_K)$ is called a K -armed bandit instance and the indices of its distributions, which belong to $[K] \stackrel{\text{def}}{=} \{1, \dots, K\}$ are its *arms*² (hence at each round we choose an arm to observe). All distributions of $\underline{\nu}$ belong to some model denoted by \mathcal{D} , which is a subset of the set of real distributions which admit a finite mean.

Example. We may cite, for example,

- the model $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ of Gaussian distributions with common variance $\sigma^2 > 0$,
- the model $\mathcal{D}_{\mathcal{B}}$ of Bernoulli distributions with means in $(0, 1)$,
- exponential models, denoted by \mathcal{D}_{exp} (see details in page 43),
- the model \mathcal{D}_{σ^2} of σ^2 -sub-Gaussian distributions, that is the set of distributions ν such that:

$$\forall \lambda \in \mathbb{R}, \quad \mathbb{E} \left[e^{\lambda(X - \mathbb{E}[X])} \right] \leq e^{\frac{\lambda^2 \sigma^2}{2}}, \quad (2.1)$$

where $X \sim \nu$,

- the model $\mathcal{P}[0, 1]$ of bounded distributions taking values in $[0, 1]$.

Multi-armed bandit problems have been vastly studied in the literature. The interest in this model comes from the wide variety of applications (see [Lattimore and Szepesvári, 2020](#), Chapter 1 for additional motivations), which has led to the study of many objectives (including regret minimization and best-arm identification), the consideration of different utility functions (mean, quantile, CVaR, etc.), of various structures when there is an underlying dependency between the means (linear bandits, combinatorial bandits, graph bandits, etc.), and also extensions to an infinite number of arms.

Best-arm identification. In the rest of this chapter, unless otherwise specified, we only consider bandit instances that have a unique *best arm* (or *optimal arm*) denoted by $a^*(\underline{\nu})$:

$$\{a^*(\underline{\nu})\} \stackrel{\text{def}}{=} \operatorname{argmax}_{a \in [K]} \mu_a,$$

where $\mu_a \stackrel{\text{def}}{=} \mathbb{E}(\nu_a)$ is the mean of arm a . The corresponding distribution and its mean will often be denoted by ν^* and μ^* :

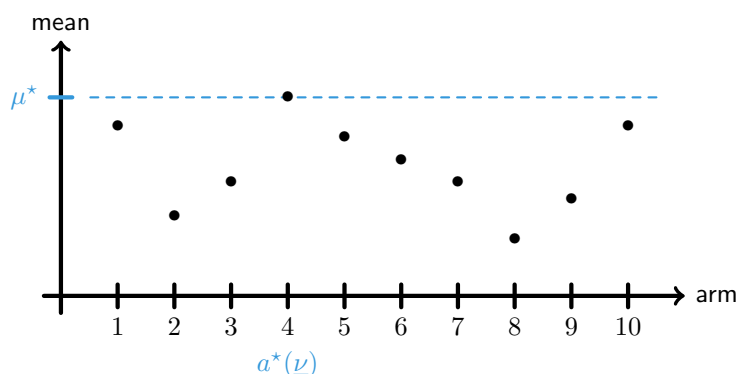
$$\nu^* \stackrel{\text{def}}{=} \nu_{a^*(\underline{\nu})} \quad \text{and} \quad \mu^* \stackrel{\text{def}}{=} \mu_{a^*(\underline{\nu})},$$

and other arms are said to be *sub-optimal*s.

The best-arm identification task consists in identifying $a^*(\underline{\nu})$. As we will detail in the next sections, the problem can be considered under two different approaches:

- in the *fixed-confidence* setting, we introduce some confidence parameter $\delta \in (0, 1)$ and we want to design strategies that will stop after some (random) finite number of observations and give an estimate of the best arm $a^*(\underline{\nu})$ which is correct with probability at least $1 - \delta$. Such

²As the indices characterize the distribution, a slight abuse consists in using arms to directly denote the distributions.

Figure 2.2: Optimal arm and mean of a bandit instance ν .

strategies are called δ -correct. Of course, by accessing a very large number of observations, we will be able to design δ -correct strategies, yet we want to avoid unnecessary observations. This is why, in that setting, we try to find δ -correct strategies that minimize the expected number of queries. This setting was originally considered by [Even-Dar et al. \(2006\)](#).

- in the *fixed-budget* setting, we have to give an estimate of the best arm after a given total of T observations. The objective is to design strategies that minimize the probability of misidentification. This setting is studied in [Audibert et al. \(2010\)](#).

Surprisingly, the two settings are not equivalent: they do not lead to the same strategies, and they are not equally understood in the literature (as we will see in this introductory chapter).

Example. Obtaining a δ -correct estimate of the best arm might require a long number of queries, which is not adapted to every situation:

- during a medical trial ([Application 2.1](#)), the panel of participants is limited, hence you have a maximal number of possible observations and might model the situation as a fixed-budget best-arm identification problem,
- in online advertising ([Application 2.2](#)), if your webpage is visited frequently enough, you have access to a potentially large number of internet users, so that you can wait until obtaining a δ -correct estimate of the best add, and then choose this add to be displayed to all future users.

The problem of best-arm identification is part of the *pure exploration* problems ([Lattimore and Szepesvári, 2020](#), Chapter 33): observing sub-optimal distributions is not directly penalized (contrarily to the problem of regret minimization) and even required to increase the confidence on the prediction. However, exploration might be done carefully:

- in the fixed-confidence setting, observing sub-optimal distributions more than necessary will be costly for the total number of queries that we want to minimize,
- in the fixed-budget setting, too many observations of the worst arms will imply that less budget will be devoted to comparing the best distributions, hence increasing the probability of error.

Before being seen as bandit problems, both settings of the best-arm identification problem have been considered in the “*ranking and selection*” literature (see [Hong et al., 2021](#) for a survey).

2.2. Best-Arm Identification with a Fixed-Confidence

In this section, we describe the fixed-confidence setting introduced by [Even-Dar et al. \(2006\)](#) and the main results obtained in the literature for that setting. The objective is to design strategies that after some random number of steps τ return an estimate $\hat{a}_\tau \in [K]$, which is equal to the best arm $a^*(\underline{\nu})$ with probability at least $1 - \delta$, where $\delta \in (0, 1)$ is a fixed confidence parameter.

We begin with a mathematical definition of what a strategy is in that setting. We denote by Y_t the observation at step $t \geq 1$ and consider an independent sequence $(U_t)_{t \geq 0}$ of i.i.d. random variables of law $\mathcal{U}([0, 1])$ to be used as an external randomization by the strategy.

Formal definition of a strategy. In the fixed-confidence setting, a strategy is defined by

- A *sampling rule*, which consists in choosing the arm $A_t \in [K]$ to observe at each time step $t \geq 1$. This arm A_t depends on the previous observations Y_1, \dots, Y_{t-1} , but also possibly on some external randomization that we capture by the random variable U_{t-1} . A_t is thus \mathcal{F}_{t-1} -measurable, where

$$\mathcal{F}_{t-1} \stackrel{\text{def}}{=} \sigma(I_{t-1}), \quad \text{with } I_{t-1} \stackrel{\text{def}}{=} (U_0, Y_1, U_1, Y_2, U_2, \dots, Y_{t-1}, U_{t-1}).$$

I_{t-1} corresponds to the information available at the end of time step $t - 1$, i.e., to the history.

- A *stopping rule* τ , which is a stopping time with respect to the filtration $(\mathcal{F}_t)_{t \geq 0}$.
- A *decision rule* \hat{a}_τ which is \mathcal{F}_τ -measurable.

The general structure of a strategy is presented in [Algorithm 1](#). The stopping rule τ highly depends on the confidence parameter δ , this is why we will denote it as τ_δ from now on. For the sampling and decision rules, most of the presented strategies in this chapter are independent³ from δ .

For a given strategy facing a bandit problem $\underline{\nu}$, let $N_a(t)$ and $\hat{\mu}_a(t)$ denote the number of pulls and the empirical mean⁴ of arm a at step t :

$$N_a(t) \stackrel{\text{def}}{=} \sum_{s \in [t]} \mathbb{I}\{A_s = a\} \quad \text{and} \quad \hat{\mu}_a(t) \stackrel{\text{def}}{=} \frac{1}{N_a(t)} \sum_{s \in [t]} Y_s \mathbb{I}\{A_s = a\},$$

where $\mathbb{I}\{E\}$ denotes the indicator (or characteristic) function of event E .

A most natural decision rule, used by many strategies, is to return the arm with the best (i.e., the largest) empirical estimate $\hat{\mu}_a(\tau_\delta)$, see [Algorithm 2](#).

For the sake of clarity, probabilities and expectations will be indexed by the bandit problem under which observations are done: $\mathbb{P}_\underline{\nu}$ (respectively, $\mathbb{E}_\underline{\nu}$) denotes the probability (respectively, expectation) under the bandit problem $\underline{\nu}$. We now give a concrete definition of δ -correctness.

Definition 2.3. [δ -correct strategy]

Let $\delta \in (0, 1)$. A strategy is δ -correct on a model \mathcal{D} if for all K -armed bandit problems $\underline{\nu}$ in \mathcal{D} with a unique optimal arm,

$$\mathbb{P}_\underline{\nu}(\tau_\delta < +\infty, \hat{a}_{\tau_\delta} \neq a^*(\underline{\nu})) \leq \delta. \quad (2.2)$$

Remark. The condition (2.2) means that the strategy returns a bad arm with controlled probability. As the strategy may never stop, it is not equivalent to ask that the strategy returns a^* with high probability:

$$\mathbb{P}_\underline{\nu}(\tau_\delta < +\infty, \hat{a}_{\tau_\delta} = a^*(\underline{\nu})) \geq 1 - \delta. \quad (2.3)$$

Note that both conditions (2.2) and (2.3) are used in the literature, which is not confusing as we will see that efficient strategies satisfy $\tau_\delta < +\infty$ $\mathbb{P}_\underline{\nu}$ -a.s..

³All but elimination algorithms presented in the next section.

⁴All considered strategies initially pull each arm once, hence $\hat{\mu}_a(t)$ is well-defined for $t \geq K$.

Algorithm 1: General structure of a fixed-confidence strategy

Input: confidence parameter δ
sampling-rule, stopping-condition, decision-rule

Output: stopping time τ_δ
estimated best arm \hat{a}_{τ_δ}

```
1 Observe each arm once // initialization
2  $t \leftarrow K$ 
3 while stopping-condition( $I_t, \delta$ ) is not satisfied do
4   | Increase  $t$  by 1
5   |  $A_t \leftarrow$  sampling-rule( $I_{t-1}$ )
6   | Observe  $Y_t \sim \nu_{A_t}$ 
7  $\tau_\delta \leftarrow t$ 
8  $\hat{a}_{\tau_\delta} \leftarrow$  decision-rule( $I_{\tau_\delta}$ )
```

Algorithm 2: Empirical-Best decision rule

Input: history of observations I_{τ_δ}

Output: estimated best arm \hat{a}_{τ_δ}

```
1 Choose  $\hat{a}_{\tau_\delta} \in \operatorname{argmax}_{a \in [K]} \hat{\mu}_a(\tau_\delta)$ 
```

We will see that the δ -correct condition is ensured by a careful choice of the stopping rule, while the sampling rule determines the performance of the strategy. We explore possible choices of both procedures in Sections 2.2.1, 2.2.5 and 2.2.7.

2.2.1. Naive Approaches

In this section, we present the first δ -correct strategies introduced in the literature by [Even-Dar et al. \(2006\)](#). They considered *elimination* (or *racing*) algorithms, which are strategies with a particular structure: the exploration is split in rounds (or phases). The strategy maintains a list of candidate arms, starting with all arms, and drops (at least) one arm at the end of each phase. We will only consider elimination strategies for which, in each phase r , the random set S_{r-1} of active candidates is uniformly explored, and the worst empirical arm a_r is removed (hence strategies have $K - 1$ rounds). See Algorithm 3 for details.

Remark. Elimination algorithms are the only strategies considered in this thesis that do not use the Empirical-Best decision rule (Algorithm 2): they recommend the remaining arm of S_{K-1} , which might have a smaller empirical estimate than one of the arms eliminated in rounds 1 to $K - 2$ (even if this is quite unlikely).

We consider here the model $\mathcal{P}[0, 1]$ of bounded distributions taking values in $[0, 1]$, but the results generalize to any σ^2 -sub-Gaussian model \mathcal{D}_{σ^2} . [Even-Dar et al. \(2006\)](#) proposed a δ -correct elimination algorithm for which they obtained guarantees on the number of pulls depending on the *gaps* of bandit problem $\underline{\nu}$. The gap $\Delta_a(\underline{\nu})$ of arm a is the expected difference between the rewards under the optimal distribution ν^* and under ν_a :

$$\Delta_a(\underline{\nu}) \stackrel{\text{def}}{=} \mu^* - \mu_a.$$

Algorithm 3: General structure of an elimination algorithm

Input: confidence parameter δ
 stopping-condition for each phase
Output: stopping time τ_δ
 estimated best arm \hat{a}_{τ_δ}

```

1  $t \leftarrow 0$ 
2  $S_0 \leftarrow [K]$ 
3 for each phase  $r \in [K - 1]$  do
4     while stopping-condition( $r, I_t$ ) does not hold do
5         Observe each arm of  $S_{r-1}$  once
6         Increase  $t$  by the cardinality of  $S_{r-1}$ 
7         Choose  $a_r \in \operatorname{argmin}_{a \in S_{r-1}} \hat{\mu}_a(t)$ 
8          $S_r \leftarrow S_{r-1} \setminus \{a_r\}$ 
9  $\tau_\delta \leftarrow t$ 
10 Define  $\hat{a}_{\tau_\delta}$  as the unique element of  $S_{K-1}$ 
    
```

In the sequel, when there is no confusion, we set $a^* = a^*(\underline{\nu})$ and $\Delta_a = \Delta_a(\underline{\nu})$ for all $a \in [K]$. We will also order arms by their means, using (reverse) notation of order statistics:

$$\mu_{(1)} > \mu_{(2)} \geq \mu_{(3)} \geq \cdots \geq \mu_{(K)}.$$

As a consequence, note that $a^*(\underline{\nu}) = (1)$.

Letters a, b, c always refer to arms, that is, to elements of $[K]$. In subindices for sums and infima, or with quantifiers, we sometimes omit to explicitly mention $[K]$ for simplicity: for example, given a fixed arm b , $\sum_{a \neq b}$ denotes the sum over arms $a \in [K] \setminus \{b\}$.

To be efficient in the identification task, the racing Algorithm 3 should not eliminate the best arm at the end of all phases. If remaining arms have been pulled n times at the end of a phase, such elimination might occur if two n -sample averages of distributions ν^* and ν_a , for some sub-optimal arm a , are in reverse order compared to the expectations of the underlying distributions. We will first focus on the probability of those undesired events, which will then allow us to derive phase lengths that ensure the δ -correctness of the elimination Algorithm 3. Before doing so, we introduce convenient notation offered by optional skipping.

Optional skipping. Let $(X_{a,n})_{a \in [K], n \geq 1}$ be independent random variables such that $X_{a,n} \sim \nu_a$ for all $a \in [K]$ and $n \geq 1$. *Optional skipping* (see Garivier et al., 2022, Section 4.1 for an extensive presentation) is a useful tool in multi-armed bandit problems based on Doob's optional skipping (Doob, 1953, Theorem 5.2), which allows assuming that, given $A_t = a$, the observation at time step t is $Y_t = X_{a, N_a(t)}$. As a consequence, we notably get

$$\hat{\mu}_a(t) = \frac{1}{N_a(t)} \sum_{n=1}^{N_a(t)} X_{a,n} \stackrel{\text{def}}{=} \hat{\mu}_{a, N_a(t)},$$

making a link between $\hat{\mu}_a(t)$ which is a quantity defined by the global time step t to some $\hat{\mu}_{a,n}$ which depends on local steps of arm a . It might be highlighted that the sequence $(X_{a,n})_{a \in [K], n \geq 1}$ is assumed to be drawn regardless of the strategy. In the analysis, it is sometimes useful to rely on the existence of $\hat{\mu}_{a,n}$ even if arm a has been totally pulled less than n times. See, e.g., Equation (2.7).

Discriminate between two arms. Assume that you have to discriminate the means of arms a^* and a with confidence δ . How many observations do you need if you pull both arms equally? One can prove that at most $n = \Delta_a^{-2} \log \frac{1}{\delta}$ pulls of each arm are required using Hoeffding's inequality.

Remark. We recall that Hoeffding's inequality states that, if N is an integer and X_1, \dots, X_N are independent random variables such that

- either each X_s belongs to $[a_s, b_s]$ almost surely, with $\sigma_s = \frac{b_s - a_s}{2}$,
- or, more generally⁵, each X_s is σ_s^2 -sub-Gaussian, see (2.1),

then

$$\forall t > 0, \quad \mathbb{P}\left(\sum_{s=1}^N X_s - \mathbb{E}[X_s] \geq t\right) \leq \exp\left(-\frac{t^2}{2 \sum_{s=1}^N \sigma_s^2}\right). \quad (2.4)$$

For the choice of n given above, the misidentification probability is indeed controlled by δ :

$$\begin{aligned} \mathbb{P}_{\mathcal{L}}(\hat{\mu}_{a,n} \geq \hat{\mu}_{a^*,n}) &= \mathbb{P}_{\mathcal{L}}(n(\hat{\mu}_{a,n} - \hat{\mu}_{a^*,n} + \Delta_a) \geq n\Delta_a) \\ &\leq \exp(-\Delta_a^2 n) \\ &\leq \delta, \end{aligned} \quad (2.5)$$

by applying Hoeffding's inequality (2.4) to

$$n(\hat{\mu}_{a,n} - \hat{\mu}_{a^*,n} + \Delta_a) = \sum_{s=1}^n (X_{a,s} - \mathbb{E}[X_{a,s}]) - \sum_{s=1}^n (X_{a^*,s} - \mathbb{E}[X_{a^*,s}]),$$

which is the centered sum of a $2n$ -sample of $\frac{1}{4}$ -sub-Gaussian variables (as each $X_{a,s}$ and $X_{a^*,s}$ belongs to $[0, 1]$ almost surely). Using this simple observation will allow us to design δ -correct elimination algorithms.

With known gaps. Assume that we know all values of the gaps, but not the mapping $a \mapsto \Delta_a$. Then we can use inequality (2.5) to fix deterministic phase lengths so as to control the probability of misidentification by δ . Indeed, let n_r be the determined and pre-defined total number of pulls at the end of round r of an arm that was still running in that phase (i.e., belonging to S_{r-1}). The algorithm fails to identify the best arm if a^* is eliminated at some phase, hence the probability of misidentification can be decomposed as follows:

$$\begin{aligned} \mathbb{P}_{\mathcal{L}}(\hat{a}_{\tau_\delta} \neq a^*) &= \sum_{r \in [K-1]} \mathbb{P}_{\mathcal{L}}(a^* \text{ is eliminated at round } r) \\ &= \sum_{r \in [K-1]} \mathbb{P}_{\mathcal{L}}(a_r = a^*). \end{aligned} \quad (2.6)$$

Fixing $r \in [K-1]$, we will choose n_r such that $\mathbb{P}_{\mathcal{L}}(a_r = a^*) \leq \frac{\delta}{K}$, which will make the strategy δ -correct. We observe that S_{r-1} contains $K-r+1$ arms, hence at least one arm among the r arms $\{(K), (K-1), \dots, (K-r+1)\}$. We deduce the following bound, using optional skipping and, for

⁵Hoeffding's lemma exactly claims that, if a random variable belongs almost surely to an interval of length 2σ , then it is σ^2 -sub-Gaussian.

the third inequality, bound (2.5):

$$\begin{aligned}
 \mathbb{P}_{\mathcal{L}}(a_r = a^*) &= \mathbb{P}_{\mathcal{L}}\left(a^* \in S_{r-1} \text{ and } \hat{\mu}_{a^*,n_r} \leq \min_{a \in S_{r-1} \setminus \{a^*\}} \hat{\mu}_{a,n_r}\right) \\
 &\leq \mathbb{P}_{\mathcal{L}}\left(\exists a \in \{(K), (K-1), \dots, (K-r+1)\} : \hat{\mu}_{a^*,n_r} \leq \hat{\mu}_{a,n_r}\right) \quad (2.7) \\
 &\leq \sum_{a=K-r+1}^K \mathbb{P}_{\mathcal{L}}\left(\hat{\mu}_{a^*,n_r} \leq \hat{\mu}_{(a),n_r}\right) \\
 &\leq \sum_{a=K-r+1}^K \exp\left(-\Delta_{(a)}^2 n_r\right) \\
 &\leq K \exp\left(-\Delta_{(K-r+1)}^2 n_r\right) \quad (2.8) \\
 &\leq \frac{\delta}{K}
 \end{aligned}$$

as soon as $n_r \geq \left\lceil \frac{2}{\Delta_{(K-r+1)}^2} \log \frac{K}{\delta} \right\rceil$.

Using this lower bound to define the phase lengths we obtain a δ -correct elimination algorithm (relying on the knowledge of the gaps) for which the number of queries is deterministic: as the arm eliminated in phase r is pulled n_r times and the surviving arm is pulled n_{K-1} times, we get

$$\tau_{\delta} = \sum_{r \in [K-1]} n_r + n_{K-1} = \sum_{a \neq a^*} \left\lceil \frac{2 \log \frac{K}{\delta}}{\Delta_a^2} \right\rceil + \left\lceil \frac{2 \log \frac{K}{\delta}}{\Delta_{\min}^2} \right\rceil, \quad (2.9)$$

where $\Delta_{\min} \stackrel{\text{def}}{=} \min_{a \neq a^*} \Delta_a = \Delta_{(2)}$. By observing that

$$2 \log \frac{K}{\delta} \sum_{a \neq a^*} \frac{1}{\Delta_a^2} \leq \tau_{\delta} \leq 8 \log \frac{K}{\delta} \sum_{a \neq a^*} \frac{1}{\Delta_a^2},$$

the sample complexity is of order

$$\mathcal{O}\left(\log \frac{K}{\delta} \sum_{a \neq a^*} \frac{1}{\Delta_a^2}\right).$$

Remark. As we observed, this strategy has a constant stopping time τ_{δ} . Elimination algorithms have been similarly studied in fixed-budget best-arm identification ([Audibert et al., 2010](#); [Karnin et al., 2013](#)), as we will discuss in Section 2.3.

Without the knowledge of gaps. Without the knowledge of gaps, an elimination strategy has no choice but to consider adaptive phase lengths in order to be δ -correct. [Even-Dar et al. \(2006\)](#) developed such an algorithm, using Algorithm 3 with a stopping criterion for each phase depending on δ , K and t and based on anytime confidence intervals (we will discuss more this concept in the paragraph on stopping conditions in page 50). They proved that their strategy satisfies

$$\tau_{\delta} = \mathcal{O}\left(\sum_{a \neq a^*} \frac{\log \frac{K}{\delta \Delta_a}}{\Delta_a^2}\right) = \mathcal{O}\left(\log\left(\frac{K}{\delta \Delta_{\min}}\right) \sum_{a \neq a^*} \frac{1}{\Delta_a^2}\right) \quad (2.10)$$

with probability at least $1 - \delta$.

2.2.2. The Fundamental Inequality

When measuring the performance of δ -correct strategies, one may be interested in deriving bounds on the number of pulls either with high probability, as in bound (2.10), or in expectation. From now on, we focus on bounds obtained in expectation, as Sections 2.2.3 and 2.2.4 will prove that this criterion seems more suitable.

[Garivier and Kaufmann \(2016\)](#) proved that the performance in expectation of δ -correct strategies cannot be arbitrarily good: they obtained an instance-dependent lower bound on the expected number of pulls of any δ -correct strategy that we will present in Section 2.2.3. The result is based on a fundamental inequality whose applications go beyond the fixed-confidence setting. This fundamental inequality relies on change-of-measure arguments and on information-theoretic properties of the Kullback-Leibler divergence that we briefly introduce in the next paragraph.

The Kullback-Leibler divergence. The *Kullback-Leibler divergence* is a pseudo-distance on probability distributions that appears to play a special role in information theory. It is defined, for all pairs of probability measures \mathbb{P} and \mathbb{Q} defined on a measurable space (Ω, \mathcal{A}) , by

$$\text{KL}(\mathbb{P}, \mathbb{Q}) \stackrel{\text{def}}{=} \begin{cases} \int_{\Omega} \log \frac{d\mathbb{P}}{d\mathbb{Q}} d\mathbb{P} = \mathbb{E}_{X \sim \mathbb{P}} \left[\log \frac{d\mathbb{P}}{d\mathbb{Q}}(X) \right] & \text{if } \mathbb{P} \ll \mathbb{Q}, \\ +\infty & \text{otherwise,} \end{cases}$$

where $\frac{d\mathbb{P}}{d\mathbb{Q}}$ denotes the Radon-Nikodym derivative of \mathbb{P} with respect to \mathbb{Q} when \mathbb{P} is absolutely continuous with respect to \mathbb{Q} . One can easily check that $\text{KL}(\mathbb{P}, \mathbb{Q})$ is a non-negative and non-symmetric quantity, and that it is null if and only if $\mathbb{P} = \mathbb{Q}$.

The Kullback-Leibler divergence satisfies a useful inequality to derive lower bounds, namely the *data-processing inequality* (or *contraction of entropy*), which states that the divergence between push-forward measures is smaller than the divergence between the original measures. More precisely, if \mathbb{P}, \mathbb{Q} are two probability distributions on a measurable space (Ω, \mathcal{A}) and $X : (\Omega, \mathcal{A}) \rightarrow (\Omega', \mathcal{A}')$ is a random variable, then

$$\text{KL}(\mathbb{P}^X, \mathbb{Q}^X) \leq \text{KL}(\mathbb{P}, \mathbb{Q}), \quad (2.11)$$

where \mathbb{P}^X (respectively, \mathbb{Q}^X) denotes the push-forward measure of X under \mathbb{P} (respectively, \mathbb{Q}). See [Ali and Silvey \(1966\)](#) for proof of this statement, and, e.g., [Cover and Thomas \(2006\)](#) for additional information on the Kullback-Leibler divergence.

The fundamental inequality. Instance-dependent lower bounds can be obtained by performing changes of measure, an argument extensively used in the bandit literature (see, e.g., [Lai and Robbins, 1985](#); [Mannor and Tsitsiklis, 2004](#)). This will allow us to quantify how many steps are necessary to discriminate $\underline{\nu}$ from a “close” alternative bandit problem $\underline{\zeta}$ which has a different optimal arm. [Kaufmann et al. \(2016\)](#) derived the following inequality for best-arm identification problems, which is convenient to use as it “hides” the change of measure within a high-level informational argument.

Lemma 2.4. [[Kaufmann et al., 2016, Lemma 1](#)]

Let $\underline{\nu}$ and $\underline{\zeta}$ be K -armed bandit problems. Consider any best-arm identification strategy and assume that the stopping-time τ is $\mathbb{P}_{\underline{\nu}}$ -integrable. Then, for all events E in \mathcal{F}_{τ} , we have

$$\sum_{a \in [K]} \mathbb{E}_{\underline{\nu}}[N_a(\tau)] \text{KL}(\nu_a, \zeta_a) \geq \text{KL}(\text{Ber}(\mathbb{P}_{\underline{\nu}}(E)), \text{Ber}(\mathbb{P}_{\underline{\zeta}}(E))). \quad (2.12)$$

Remark. As we will see, this inequality can be applied both in the fixed-confidence and the fixed-budget settings, but it will have stronger implications in the first one. We will explain in Remark 2.18 that this is mainly due to the fact that, in Equation (2.12), expectations on the number of pulls are relative to bandit problem $\underline{\nu}$; which implies that $\underline{\nu}$ and $\underline{\zeta}$ have very different roles in this inequality (in addition to the fact that the Kullback-Leibler divergence is non-symmetric).

The proof of Kaufmann et al. (2016) is based on explicit changes of measure and Wald's lemma, while Garivier et al. (2019) provided a shorter proof using information-theoretic properties of the Kullback-Leibler divergence when τ is a deterministic time. We will unify both proof schemes by applying a martingale stopping theorem.

Stochastic transition kernels. In the following proof, we will use the fact that the law of I_t under $\mathbb{P}_{\underline{\nu}}$ can be defined iteratively by *regular stochastic transition kernels*, from which we now give a definition. Let (Ω, \mathcal{A}) and (Ω', \mathcal{A}') be two measurable spaces. We say that $K : \Omega \times \mathcal{A}' \rightarrow [0, 1]$ is a *regular transition kernel* from (Ω, \mathcal{A}) to (Ω', \mathcal{A}') if $K(\omega, \cdot)$ is a probability distribution on (Ω', \mathcal{A}') and if $K(\cdot, B)$ is \mathcal{A} -measurable for all $B \in \mathcal{A}'$. For a probability measure \mathbb{P} on (Ω, \mathcal{A}) , we define the probability distribution $K\mathbb{P}$ on $(\Omega \times \Omega', \mathcal{A} \otimes \mathcal{A}')$ by

$$\forall (A, B) \in \mathcal{A} \times \mathcal{A}', \quad K\mathbb{P}(A, B) \stackrel{\text{def}}{=} \int_{\Omega} \mathbb{I}\{\omega \in A\} K(\omega, B) d\mathbb{P}(\omega).$$

Fixing $t \geq 1$, we get that

$$\mathbb{P}_{\underline{\nu}}^{I_t} = K_{\underline{\nu}}^{I_{t-1}} \mathbb{P}_{\underline{\nu}}^{I_{t-1}}, \quad \text{where} \quad K_{\underline{\nu}}^t(I_{t-1}, \cdot) = \sum_{a \in [K]} \mathbb{I}\{A_t = a\} \nu_a \otimes \text{Leb}_{[0,1]} = \nu_{A_t} \otimes \text{Leb}_{[0,1]} \quad (2.13)$$

is a regular transition kernel, with $\text{Leb}_{[0,1]}$ the Lebesgue measure on $[0, 1]$.

Proof. Once again, we rely on optional skipping (see page 10) and consider the sequence of rewards $(X_{a,n})_{a \in [K], n \geq 1}$. Without loss of generality, we assume that $\text{KL}(\nu_a, \zeta_a)$ is finite⁷ for each arm $a \in [K]$. As a consequence, note that $\nu_a \ll \zeta_a$ and that $\log \frac{d\nu_a}{d\zeta_a}(X_{a,1})$ is $\mathbb{P}_{\underline{\nu}}$ -integrable for all $a \in [K]$.

We will prove that

$$\sum_{a \in [K]} \mathbb{E}_{\underline{\nu}}[N_a(\tau)] \text{KL}(\nu_a, \zeta_a) = \text{KL}(\mathbb{P}_{\underline{\nu}}^{I_\tau}, \mathbb{P}_{\underline{\zeta}}^{I_\tau}). \quad (2.14)$$

The result will follow by combining this equality with an application of the data-processing inequality (2.11): $\mathbb{I}\{E\}$ being \mathcal{F}_τ -measurable, it is a random variable depending only on I_τ , hence:

$$\text{KL}(\mathbb{P}_{\underline{\nu}}^{I_\tau}, \mathbb{P}_{\underline{\zeta}}^{I_\tau}) \geq \text{KL}(\mathbb{P}_{\underline{\nu}}^{\mathbb{I}\{E\}}, \mathbb{P}_{\underline{\zeta}}^{\mathbb{I}\{E\}}) = \text{KL}(\text{Ber}(\mathbb{P}_{\underline{\nu}}(E)), \text{Ber}(\mathbb{P}_{\underline{\zeta}}(E))).$$

We now prove Equation (2.14). We define $M_0 = 0$ and

$$\forall t \geq 1, \quad M_t = \log \frac{d\mathbb{P}_{\underline{\nu}}^{I_t}}{d\mathbb{P}_{\underline{\zeta}}^{I_t}} - \sum_{a \in [K]} N_a(t) \text{KL}(\nu_a, \zeta_a).$$

The process $(M_t)_{t \geq 0}$ is adapted to $(\mathcal{F}_t)_{t \geq 0}$. We will show that $(M_t)_{t \geq 0}$ is a martingale by using the following chain rule property (see, e.g., Wainwright, 2019, Exercice 3.2 or Stoltz, 2022, Lecture 5), which is a consequence of the Radon-Nikodym theorem.

⁶Note that the existence of a measurable space on which $\mathbb{P}_{\underline{\nu}}$ is defined is ensured by Kolmogorov's extension theorem.

⁷If $\text{KL}(\nu_a, \zeta_a) = +\infty$, either the left-hand-side of Equation (2.12) is infinite, or $\mathbb{E}_{\underline{\nu}}[N_a(\tau)] = 0$ so that arm a is almost surely never pulled and the proof still applies.

Proposition 2.5. Let K, L be two regular transition kernels from (Ω, \mathcal{A}) to (Ω', \mathcal{A}') and \mathbb{P}, \mathbb{Q} be two probability distributions on (Ω, \mathcal{A}) with $\mathbb{P} \ll \mathbb{Q}$, and assume that

- $K(\omega, \cdot) \ll L(\omega, \cdot)$ for \mathbb{Q} -almost all $\omega \in \Omega$,
- there exists a version of $(\omega, \omega') \mapsto \frac{dK(\omega, \cdot)}{dL(\omega, \cdot)}(\omega')$ which is $\mathcal{A} \otimes \mathcal{A}'$ -bi-measurable.

Then $K\mathbb{P} \ll L\mathbb{Q}$ and the following chain rule property for densities applies:

$$\frac{dK\mathbb{P}}{dL\mathbb{Q}}(\omega, \omega') = \frac{d\mathbb{P}}{d\mathbb{Q}}(\omega) \cdot \frac{dK(\omega, \cdot)}{dL(\omega, \cdot)}(\omega') \quad \text{for } L\mathbb{Q}\text{-almost all } (\omega, \omega').$$

Fix $t \geq 1$. The transition kernels $K = K_{\underline{\mu}}^t$ and $L = K_{\underline{\zeta}}^t$, defined in (2.13), satisfy the hypothesis of the above proposition, as $\nu_a \ll \zeta_a$ for all $a \in [K]$, and as for all $h_{t-1} = (u_0, y_1, \dots, u_{t-1})$ and y_t, u_t :

$$\frac{dK_{\underline{\mu}}^t(h_{t-1}, \cdot)}{dK_{\underline{\zeta}}^t(h_{t-1}, \cdot)}(y_t, u_t) = \sum_{a \in [K]} \mathbb{I}\{A_t(h_{t-1}) = a\} \frac{d\nu_a}{d\zeta_a}(y_t),$$

where $A_t(h_{t-1})$ denotes the arm chosen by the strategy knowing history h_{t-1} . We can apply the proposition with $\mathbb{P} = \mathbb{P}_{\underline{\nu}}^{t-1}$ and $\mathbb{Q} = \mathbb{P}_{\underline{\zeta}}^{t-1}$, as one can show that $\mathbb{P} \ll \mathbb{Q}$ by induction. Using (2.13), and recalling that U_t is an auxiliary random variable independent from all other randomization, this entails

$$\begin{aligned} M_t &= \log \frac{d\mathbb{P}_{\underline{\nu}}^{t-1}}{d\mathbb{P}_{\underline{\zeta}}^{t-1}} + \log \frac{dK_{\underline{\mu}}^t(I_{t-1}, \cdot)}{dK_{\underline{\zeta}}^t(I_{t-1}, \cdot)}(Y_t, U_t) - \sum_{a \in [K]} (N_a(t-1) + \mathbb{I}\{A_t = a\}) \text{KL}(\nu_a, \zeta_a) \\ &= M_{t-1} + \sum_{a \in [K]} \mathbb{I}\{A_t = a\} \left(\log \frac{d\nu_a}{d\zeta_a}(Y_t) - \text{KL}(\nu_a, \zeta_a) \right) \\ &= M_{t-1} + \sum_{a \in [K]} \mathbb{I}\{A_t = a\} \left(\log \frac{d\nu_a}{d\zeta_a}(X_{a, N_a(t)}) - \text{KL}(\nu_a, \zeta_a) \right). \end{aligned}$$

We deduce from this expression that the sequence $(M_t)_{t \geq 0}$ is a martingale, as for all $t \geq 1$

$$\begin{aligned} \mathbb{E}_{\underline{\nu}}[M_t | \mathcal{F}_{t-1}] &= M_{t-1} + \sum_{a \in [K]} \mathbb{I}\{A_t = a\} \mathbb{E}_{\underline{\nu}} \left[\left(\log \frac{d\nu_a}{d\zeta_a}(X_{a, N_a(t)}) - \text{KL}(\nu_a, \zeta_a) \right) | \mathcal{F}_{t-1} \right] \\ &= M_{t-1} + \sum_{a \in [K]} \mathbb{I}\{A_t = a\} \underbrace{\left(\mathbb{E}_{\underline{\nu}} \left[\log \frac{d\nu_a}{d\zeta_a}(X_{a,1}) \right] - \text{KL}(\nu_a, \zeta_a) \right)}_{=0} = M_{t-1}, \end{aligned}$$

where we used that A_t is \mathcal{F}_{t-1} -measurable and that $X_{a, N_a(t)} | \mathcal{F}_{t-1} \sim \nu_a$ has same law as $X_{a,1}$.

Note that τ was assumed to be $\mathbb{P}_{\underline{\nu}}$ -integrable and that the conditional expectation of the absolute increments of $(M_t)_{t \geq 0}$ are almost surely bounded: for any $t \geq 1$ we have

$$\begin{aligned} \mathbb{E}_{\underline{\nu}}[|M_t - M_{t-1}| | \mathcal{F}_{t-1}] &\leq \sum_{a \in [K]} \mathbb{I}\{A_t = a\} \mathbb{E}_{\underline{\nu}} \left[\left| \log \frac{d\nu_a}{d\zeta_a}(X_{a, N_a(t)}) - \text{KL}(\nu_a, \zeta_a) \right| | \mathcal{F}_{t-1} \right] \\ &\leq \sup_{a \in [K]} \left\{ \text{KL}(\nu_a, \zeta_a) + \mathbb{E}_{\underline{\nu}} \left[\left| \log \frac{d\nu_a}{d\zeta_a}(X_{a,1}) \right| \right] \right\} < +\infty, \end{aligned}$$

where we used that $\log \frac{d\nu_a}{d\zeta_a}(X_{a,1})$ is $\mathbb{P}_{\underline{\nu}}$ -integrable by assumption.

By Doob's optional stopping theorem (Grimmett and Stirzaker, 2001, Section 12.5), we obtain that

$$\mathbb{E}_{\underline{\nu}}[M_\tau] = \mathbb{E}_{\underline{\nu}}[M_0] = 0.$$

Rearranging the terms directly leads to Equation (2.14), which concludes the proof. \square

2.2.3. Lower Bound for Exponential Models

Garivier and Kaufmann (2016) used the fundamental inequality (2.12) to derive a lower bound valid for common models of distributions like the model $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ of Gaussian variables with a fixed variance $\sigma^2 > 0$, or the model $\mathcal{D}_{\mathcal{B}}$ of Bernoulli distributions with means in $(0, 1)$, and more generally for any canonical one-parameter exponential family. Before stating the lower bound, we give a quick description of these families.

Canonical one-parameter exponential families. We follow largely the exposition by Lehmann and Casella (1998, Section 1.5); more details, including the proofs of the stated properties, may be found in the monograph by Lehmann and Casella (1998). A (regular) canonical one-parameter exponential family \mathcal{D}_{exp} is a set of distributions ν_θ indexed by $\theta \in \Theta$, all absolutely continuous with respect to some measure ρ on \mathbb{R} , with densities given by

$$x \mapsto \frac{d\nu_\theta}{d\rho}(x) \stackrel{\text{def}}{=} \exp(\theta x - b(\theta)), \quad (2.15)$$

for some smooth enough normalizing function b . More precisely, b is assumed to be twice differentiable. We also assume that Θ is the natural parameter space, i.e., that Θ contains all possible parameters for ρ ,

$$\Theta = \left\{ \theta \in \mathbb{R} : \int_{\mathbb{R}} \exp(\theta y) d\rho(y) < +\infty \right\},$$

and that Θ is an open interval (this latter fact is what regularity stands for). A closed-form expression of b is: for all $\theta \in \Theta$,


$$b(\theta) = \log \int_{\mathbb{R}} e^{\theta y} d\rho(y). \quad (2.16)$$

The derivative b' of b is a continuous function, by assumption, and it may be shown that it is increasing, so that b' is a one-to-one mapping with a continuous inverse $(b')^{-1}$. In addition, it can be seen, by a differentiation under the integral sign, that $\mathbb{E}(\nu_\theta) = b'(\theta)$ for all $\theta \in \Theta$. Therefore, the distributions in \mathcal{D}_{exp} may be rather parameterized by their expectations. We denote by $\mathcal{M} = b'(\Theta)$ the open interval of the expectations of distributions in \mathcal{D}_{exp} , and let μ_- and μ_+ be its lower and upper ends:

$$\mathcal{M} = (\mu_-, \mu_+).$$

For each $x \in \mathcal{M}$, there exists a unique distribution in \mathcal{D}_{exp} with expectation x , namely, $\nu_{(b')^{-1}(x)}$.

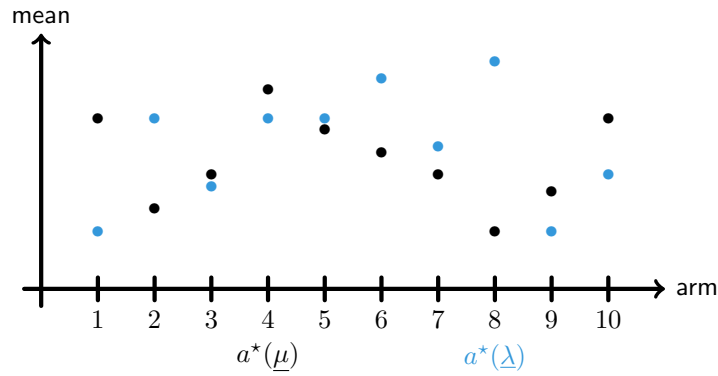
Example. Bernoulli distributions, binomial distributions, Poisson distributions and Gaussian distributions with common variance $\sigma^2 > 0$ all are canonical one-parameter exponential families.

In this thesis, any canonical one-parameter exponential family will be simply referred to as an *exponential family*, as we will not consider other exponential families. 

Kullback-Leibler divergences for \mathcal{D}_{exp} . For an exponential model \mathcal{D}_{exp} , we may also parameterize the Kullback-Leibler divergence function by the expectations: we define, for all $\theta, \theta' \in \Theta$,

$$d(\mathbb{E}(\nu_\theta), \mathbb{E}(\nu_{\theta'})) \stackrel{\text{def}}{=} \text{KL}(\nu_\theta, \nu_{\theta'}) = (\theta - \theta')b'(\theta) - b(\theta) + b(\theta'). \quad (2.17)$$

This defines a divergence d which is strictly convex and differentiable on the open set $\mathcal{M} \times \mathcal{M}$. In particular, d is continuous, is such that $d(\mu, \mu') = 0$ if and only if $\mu = \mu'$, and, for all $\mu \in \mathcal{M}$, both $d(\mu, \cdot)$ and $d(\cdot, \mu)$ are decreasing on $(\mu_-, \mu]$, and increasing on $[\mu, \mu_+)$.


 Figure 2.3: A bandit instance $\underline{\mu}$ and an alternative $\underline{\lambda}$.

As we will see, the mean-parametrized Kullback-Leibler divergence of the model of Bernoulli distributions $\mathcal{D}_{\mathcal{B}}$ plays a particular role. To that end, let us denote it by kl . Its closed-form expression reads:

$$\forall p, q \in (0, 1), \quad \text{kl}(p, q) \stackrel{\text{def}}{=} p \log \frac{p}{q} + (1 - p) \log \frac{1 - p}{1 - q}.$$

Lower bound for an exponential model. We are now able to state and prove the lower bound obtained by [Garivier and Kaufmann \(2016\)](#) for an exponential model \mathcal{D}_{exp} . As explained in the paragraph on exponential families, distributions are characterized by their means, hence in the sequel we identify $\underline{\nu}$ and its mean vector $\underline{\mu} \in \mathcal{M}^K$; for instance, $\mathbb{P}_{\underline{\mu}} = \mathbb{P}_{\underline{\nu}}$, $\mathbb{E}_{\underline{\mu}} = \mathbb{E}_{\underline{\nu}}$, and $a^*(\underline{\mu}) = a^*(\underline{\nu})$.

Before stating the lower bound, we introduce the simplex

$$\Sigma_K \stackrel{\text{def}}{=} \left\{ \underline{v} \in [0, 1]^K : v_1 + \dots + v_K = 1 \right\},$$

and the set of *alternative bandits*⁸ to a problem $\underline{\mu}$, which are bandit problems with a different optimal arm than $\underline{\mu}$ (see Figure 2.3):

$$\text{Alt}(\underline{\mu}) \stackrel{\text{def}}{=} \left\{ \underline{\lambda} \text{ in } \mathcal{D}_{\text{exp}} : a^*(\underline{\lambda}) \neq a^*(\underline{\mu}) \right\}. \quad (2.18)$$

Theorem 2.6. [[Garivier and Kaufmann, 2016, Theorem 1](#)]

Let \mathcal{D}_{exp} be an exponential model. Let $\delta \in (0, \frac{1}{2})$. For all δ -correct strategy and all bandit problems $\underline{\mu}$ in \mathcal{D}_{exp} with a unique optimal arm,

$$\mathbb{E}_{\underline{\mu}}[\tau_{\delta}] \geq T(\underline{\mu}) \text{kl}(\delta, 1 - \delta) \geq T(\underline{\mu}) \log \frac{1}{2.4\delta}, \quad (2.19)$$

where $T(\underline{\mu})$ is the characteristic time of $\underline{\mu}$, defined as

$$T(\underline{\mu})^{-1} \stackrel{\text{def}}{=} \sup_{\underline{v} \in \Sigma_K} \inf_{\underline{\lambda} \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a). \quad (2.20)$$

Example. For the model $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ of Gaussian variables with variance $\sigma^2 > 0$, the Kullback-Leibler divergence enjoys the simple closed-form expression

$$\forall \mu, \mu' \in \mathbb{R}, \quad d(\mu, \mu') = \frac{(\mu' - \mu)^2}{2\sigma^2}, \quad (2.21)$$

⁸Again, alternative bandits $\underline{\zeta}$ will be parameterized by their expectation $\underline{\lambda}$.

from which it can be proved that $T(\underline{\mu})$ satisfies the following inequalities:

$$\sum_{a \neq a^*} \frac{1}{\Delta_a^2} \leq \sum_{a \neq a^*} \frac{1}{\Delta_a^2} + \frac{1}{\Delta_2^2} \leq \frac{T(\underline{\mu})}{2\sigma^2} \leq 2 \left(\sum_{a \neq a^*} \frac{1}{\Delta_a^2} + \frac{1}{\Delta_2^2} \right) \leq 4 \sum_{a \neq a^*} \frac{1}{\Delta_a^2}.$$

This should be compared to the high probability upper bound (2.10) obtained by [Even-Dar et al. \(2006\)](#) with an adaptive elimination algorithm, which applies up to a variance term for the model $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ and gets an additional factor $\log \frac{K}{\Delta_{\min}}$.

Proof. Let $\underline{\mu}$ be a K -armed bandit problem with a unique optimal arm, and fix an alternative bandit instance $\underline{\lambda} \in \text{Alt}(\underline{\mu})$. If $\mathbb{E}_{\underline{\mu}}[\tau_\delta] = +\infty$ the result holds. We assume in the rest of the proof that $\mathbb{E}_{\underline{\mu}}[\tau_\delta] < +\infty$, and, as a consequence, $\tau_\delta < +\infty$ $\mathbb{P}_{\underline{\mu}}$ -a.s..

Considering $E = \{\tau_\delta = +\infty \text{ or } \hat{a}_{\tau_\delta} \neq a^*(\underline{\mu})\} \in \mathcal{F}_{\tau_\delta}$, and recalling that the strategy is δ -correct, we get, on the one hand, as $\tau_\delta < +\infty$ $\mathbb{P}_{\underline{\mu}}$ -a.s.,

$$\mathbb{P}_{\underline{\mu}}(E) = \mathbb{P}_{\underline{\mu}}(\tau_\delta = +\infty \text{ or } \hat{a}_{\tau_\delta}) = \mathbb{P}_{\underline{\mu}}(\tau_\delta < +\infty, \hat{a}_{\tau_\delta} \neq a^*(\underline{\mu})) \leq \delta,$$

and on the other hand, as $a^*(\underline{\mu}) \neq a^*(\underline{\lambda})$,

$$\mathbb{P}_{\underline{\lambda}}(E) = \mathbb{P}_{\underline{\lambda}}(\tau_\delta = +\infty \text{ or } \hat{a}_{\tau_\delta} \neq a^*(\underline{\mu})) \geq \mathbb{P}_{\underline{\lambda}}(\tau_\delta = +\infty \text{ or } \hat{a}_{\tau_\delta} = a^*(\underline{\lambda})) \geq 1 - \delta.$$

Injecting those two inequalities in the fundamental inequality (2.12), and using the monotonicity properties, recalled after Equation (2.17), of the kl divergence with $\delta < \frac{1}{2}$, we obtain

$$\begin{aligned} \mathbb{E}_{\underline{\mu}}[\tau_\delta] \times \sum_{a \in [K]} \frac{\mathbb{E}_{\underline{\mu}}[N_a(\tau_\delta)]}{\mathbb{E}_{\underline{\mu}}[\tau_\delta]} d(\mu_a, \lambda_a) &= \sum_{a \in [K]} \mathbb{E}_{\underline{\mu}}[N_a(\tau_\delta)] d(\mu_a, \lambda_a) \\ &\geq \text{kl} \left(\mathbb{P}_{\underline{\mu}}(\hat{a}_{\tau_\delta} \neq a^*(\underline{\mu})), \mathbb{P}_{\underline{\lambda}}(\hat{a}_{\tau_\delta} \neq a^*(\underline{\mu})) \right) \\ &\geq \text{kl}(\delta, 1 - \delta). \end{aligned}$$

This inequality holds for any alternative bandit $\underline{\lambda}$, hence, taking the infimum over $\text{Alt}(\underline{\mu})$ entails

$$\begin{aligned} \text{kl}(\delta, 1 - \delta) &\leq \mathbb{E}_{\underline{\mu}}[\tau_\delta] \times \inf_{\underline{\lambda} \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} \frac{\mathbb{E}_{\underline{\mu}}[N_a(\tau_\delta)]}{\mathbb{E}_{\underline{\mu}}[\tau_\delta]} d(\mu_a, \lambda_a) \\ &\leq \mathbb{E}_{\underline{\mu}}[\tau_\delta] \times \sup_{v \in \Sigma_K} \inf_{\underline{\lambda} \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a), \end{aligned} \quad (2.22)$$

using in the last inequality that $\tau_\delta = \sum_{a \in [K]} N_a(\tau_\delta)$. This concludes the proof of the first inequality, while the second can be derived by a mere calculus. \square

2.2.4. The Sample Complexity Optimization Problem

As in the previous subsection, we consider here an exponential model \mathcal{D}_{exp} . As we will see, the solution of the optimization problem (2.20) defining $T(\underline{\mu})$ is of high interest to design good strategies.

Transportation costs. Given a fixed proportion $v \in \Sigma_K$, reaching the infimum over $\text{Alt}(\underline{\mu})$ in optimization problem (2.20) is done by considering bandit instances for which only two arms differ from $\underline{\mu}$, namely the best arm and one of the sub-optimal arms (see Figure 2.4):

$$\begin{aligned} \inf_{\underline{\lambda} \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a) &= \min_{a \neq a^*} \inf_{\substack{\lambda_{a^*}, \lambda_a \in \mathcal{M} \\ \lambda_{a^*} < \lambda_a}} v_{a^*} d(\mu_{a^*}, \lambda_{a^*}) + v_a d(\mu_a, \lambda_a) \\ &= \min_{a \neq a^*} \inf_{x \in [\mu_a, \mu_{a^*}]} v_{a^*} d(\mu_{a^*}, x) + v_a d(\mu_a, x), \end{aligned} \quad (2.23)$$

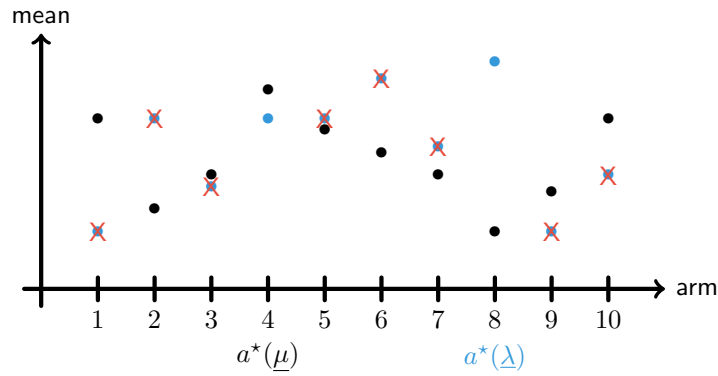


Figure 2.4: To decrease its transportation cost with $\underline{\mu}$, an alternative bandit problem $\underline{\lambda}$ should only move two arms, namely the optimal arm of $\underline{\mu}$ and one sub-optimal arm of $\underline{\mu}$.

by using the continuity and monotonicity properties of d , where $a^* = a^*(\underline{\mu})$ and we recall that $\mu^* = \mu_{a^*}$ is the optimal mean of $\underline{\mu}$. The quantity

$$\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v}) \stackrel{\text{def}}{=} \inf_{x \in [\mu_a, \mu^*]} v_{a^*} d(\mu^*, x) + v_a d(\mu_a, x)$$

can be seen as a *transportation cost* representing the difficulty of changing the distributions of arms a and a^* so as to make arm a optimal, given the pulling frequencies \underline{v} . The monotonicity properties of d ensures that the continuous function $x \mapsto v_{a^*} d(\mu^*, x) + v_a d(\mu_a, x)$ reaches its infimum in $[\mu_a, \mu^*]$, and using the closed-form (2.17) of d , we see that its derivative vanishes if and only if

$$v_{a^*}(x - \mu^*) + v_a(x - \mu_a) = 0,$$

or, equivalently, if and only if $x = \bar{\mu}_{a^*, a, \underline{v}}$, where

$$\bar{\mu}_{a^*, a, \underline{v}} \stackrel{\text{def}}{=} \frac{v_{a^*} \mu^* + v_a \mu_a}{v_{a^*} + v_a},$$

so that

$$\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v}) = v_{a^*} d(\mu^*, \bar{\mu}_{a^*, a, \underline{v}}) + v_a d(\mu_a, \bar{\mu}_{a^*, a, \underline{v}}).$$

Hence we proved that optimization problem (2.20) rewrites

$$T(\underline{\mu})^{-1} = \sup_{\underline{v} \in \Sigma_K} \min_{a \neq a^*} \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v}) = \sup_{\underline{v} \in \Sigma_K} \min_{a \neq a^*} v_{a^*} d(\mu^*, \bar{\mu}_{a^*, a, \underline{v}}) + v_a d(\mu_a, \bar{\mu}_{a^*, a, \underline{v}}). \quad (2.24)$$

Remark 2.7. When considering a general model \mathcal{D} (even non-parametric), the transportation cost becomes

$$\text{TC}_{a \rightarrow a^*}(\underline{\nu}, \underline{v}) \stackrel{\text{def}}{=} \inf_{x \in [\mu_a, \mu^*]} v_{a^*} \mathcal{K}_{\text{inf}}^<(\nu^*, x) + v_a \mathcal{K}_{\text{inf}}^>(\nu_a, x),$$

where we define, for $\nu \in \mathcal{D}$ and $x \in \mathbb{R}$,

$$\begin{aligned} \mathcal{K}_{\text{inf}}^<(\nu, x) &\stackrel{\text{def}}{=} \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) < x \}, \\ \text{and } \mathcal{K}_{\text{inf}}^>(\nu, x) &\stackrel{\text{def}}{=} \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) > x \}. \end{aligned}$$

A generalization of Theorem 2.6 to any model \mathcal{D} involving those costs is straightforward (see [Agrawal et al., 2020](#)).

An optimal weight vector. Garivier and Kaufmann (2016) proved that the sample complexity optimization problem (2.20) admits a unique maximizer $\underline{w}(\underline{\mu})$, called the *optimal weight vector* of $\underline{\mu}$:

$$\{\underline{w}(\underline{\mu})\} \stackrel{\text{def}}{=} \operatorname{argmax}_{\underline{v} \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a) = \operatorname{argmax}_{\underline{v} \in \Sigma_K} \min_{a \neq a^*} \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v}).$$

The main arguments are summarized below⁹.

- Existence: it can be seen that $\underline{v} \mapsto \min_{a \neq a^*} \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})$ is a continuous function since each $\underline{v} \mapsto \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})$ is continuous. Hence, it admits a maximum on the compact set Σ_K .
- Unicity:
 1. First, we prove that, given a fixed proportion $w_{a^*} \in (0, 1)$, there exists a unique vector $\underline{w} \in \Sigma_K$ which equalizes all transportation costs, i.e., such that

$$\exists y \in \mathbb{R}_+, \forall a \neq a^*, \quad \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w}) = y.$$

This result is a consequence of both the constraint $\underline{w} \in \Sigma_K$ and the fact that, given a fixed v_{a^*} , each cost $\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})$ is an increasing and continuous function of v_a .

2. Then, we remark that a maximizer \underline{w} has to equalize the transportation costs. Let us assume that this does not hold, i.e., that the set $\mathcal{A}_{\min} = \operatorname{argmin}_{a \neq a^*} \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w})$ does not contain all sub-optimal arms. Let us transform \underline{w} into a new vector $\underline{w}' \in \Sigma_K$ such that $w'_a > w_a$ for each $a \in \mathcal{A}_{\min}$, $w_{a^*} = w'_{a^*}$, and $w'_a < w_a$ otherwise. If the weights are modified slightly enough, this entails¹⁰ that the minimal costs increase while the other costs only decrease sufficiently low, i.e., that

$$\min_{a \neq a^*} \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w}) < \min_{a \neq a^*} \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w}'),$$

or, in other words, that \underline{w} is not a maximizer, which is a contradiction.

3. From the two first points, potential maximizers belong to a set $\{\underline{w}^\beta : \beta \in (0, 1)\}$ of weight vectors that are parameterized by the proportion $\beta = w_{a^*}^\beta$ associated to the optimal weight. By defining $y(\beta)$ as the (common) value of the transportation costs under proportions \underline{w}^β , it might be seen that $y(\beta)$ is a differentiable function, such that $y(\beta) \rightarrow 0$ when $\beta \rightarrow 0$ and $\beta \rightarrow 1$, and whose derivative vanishes only once, at some $\beta^* \in (0, 1)$. Hence the optimal maximizer is unique and equal to \underline{w}^{β^*} .

Note that the proof of unicity indicates that the optimal weight vector $\underline{w}(\underline{\mu})$ equalizes the transportation costs:

$$\forall a \neq a^*, \quad \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w}(\underline{\mu})) = T(\underline{\mu})^{-1}.$$

If we want a strategy for which inequality (2.22) is an equality, then its average proportions of draws $(\mathbb{E}_{\underline{\mu}}[N_a(\tau_\delta)] / \mathbb{E}_{\underline{\mu}}[\tau_\delta])_{a \in [K]}$ should exactly be $\underline{w}(\underline{\mu})$. To put it differently, if those average proportions are not close to $\underline{w}(\underline{\mu})$, then the performance of the strategy will be far from lower bound (2.19).

Computability of the optimal weight vector. Garivier and Kaufmann (2016) proved that solving the optimization problem (2.20) reduces to determining the root of a one-variable increasing function. By applying a bisection method, it is then possible to compute $\underline{w}(\underline{\mu})$ with arbitrary precision. We will denote by `Optimal-Weights` an algorithm computing this optimal vector.

Using their procedure, they also proved some regularity results concerning the solution of $\underline{w}(\underline{\mu})$:

- the function $\underline{\mu} \mapsto \underline{w}(\underline{\mu})$ is continuous (at problems having a unique optimal arm),
- all arms have to be linearly pulled: $w_{\min}(\underline{\mu}) \stackrel{\text{def}}{=} \min_{a \in [K]} w_a(\underline{\mu}) > 0$,

⁹Note, however, that the structure of the proof does not exactly follow the path given here, due to technicalities (see Garivier and Kaufmann, 2016, Appendix A.2).

¹⁰It is, again, a consequence of the increasing and continuity properties of each $\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})$ relatively to v_a .

- if arms are ordered so that $\mu_{(1)} > \mu_{(2)} \geq \dots \geq \mu_{(K)}$, then $w_{(2)}(\underline{\mu}) \geq \dots \geq w_{(K)}(\underline{\mu})$, but there may be instances (e.g., for Bernoulli instances) for which $w_{(1)}(\underline{\mu}) < w_{(2)}(\underline{\mu})$.

Remark 2.8. It is easy to find instances for which $w_{(1)}(\underline{\mu}) \gg w_{(2)}(\underline{\mu})$. As a consequence, for such instances, elimination algorithms presented in Section 2.2.1 cannot get close to the lower bound (2.19), as the last two remaining arms are pulled equally often.

Contribution. We make progress in the understanding of the solution of optimization problem (2.20) in two directions: first, in Chapter 3, we obtain *regularity results* for the specific case of a Gaussian model, and then, in Section 6.2 we characterize, for all exponential models, the solution $\underline{w}(\underline{\mu})$ as the unique fixed point of some transformation. Those theoretical contributions might be used to improve fixed-confidence best-arm identification algorithms, as we will explain in Contributions 2.13 and 2.14.

The solution for a Gaussian model. Considering the specific model $\mathcal{D}_{\mathcal{N}_1}$ of standard Gaussian distributions (or more generally any model of Gaussian distributions with common variance $\sigma^2 > 0$ and bounded means), we will present in Chapter 3 a new method for computing the optimal weight vector $\underline{w}(\underline{\mu})$. The proposed procedure speeds up the resolution of the computation by using Newton's method on a convex function. In addition, we obtain new quantitative regularity results concerning the solution of optimization problem (2.20):

- First, we prove that $a^* = \operatorname{argmax}_{a \in [K]} w_a(\underline{\mu})$ for this model, or equivalently $w_{(1)}(\underline{\mu}) > w_{(2)}(\underline{\mu})$. Optimal weights are thus ordered as the means for a Gaussian model,
- Then, we derive the following bounds for $w_{a^*}(\underline{\mu})$ and $T(\underline{\mu})$:

$$\frac{1}{1 + \sqrt{K-1}} \leq w_{a^*}(\underline{\mu}) \leq \frac{1}{2},$$

$$\max\left(\frac{8}{\Delta_{\min}^2}, 4 \frac{1 + \sqrt{K-1}}{\Delta^2}\right) \leq T(\underline{\mu}) \leq 2 \frac{(1 + \sqrt{K-1})^2}{\Delta_{\min}^2},$$

where $\overline{\Delta^2} \stackrel{\text{def}}{=} \frac{1}{K-1} \sum_{a \neq a^*} \Delta_a^2$ is the average squared gap and Δ_{\min} has been defined in page 39. For each inequality, we characterize instances for which the equality holds, which shows the tightness of those bounds in all generality.

- Also, we study the variations of $\underline{w}(\underline{\mu})$ and $T(\underline{\mu})$ when moving one or several arms of $\underline{\mu}$. For instance (see Figure 2.5), increasing the mean of a sub-optimal arm will increase its associated optimal weight and increase the characteristic time $T(\underline{\mu})$.
- Finally, we prove that both \underline{w} and T are locally Lipschitz functions, giving a quantification of the continuity result obtained by [Garivier and Kaufmann \(2016\)](#).

Remark. All those results are specific to the Gaussian case. In particular, similar regularity results of \underline{w} and T for other models are still to be determined.

A fixed point property for exponential models. For a general exponential model \mathcal{D}_{exp} , we use sufficient and necessary conditions satisfied by the optimal weight vector $\underline{w}(\underline{\mu})$ to prove in Section 6.2 that $\underline{w}(\underline{\mu})$ is the unique fixed point of a transformation $\underline{W} : \text{int}(\Sigma_K) \rightarrow \text{int}(\Sigma_K)$, which roughly speaking modifies a weight vector so as to decrease the difference between the transportation costs. This could entail the design of new empirical procedures that speed up the computation of $\underline{w}(\underline{\mu})$.

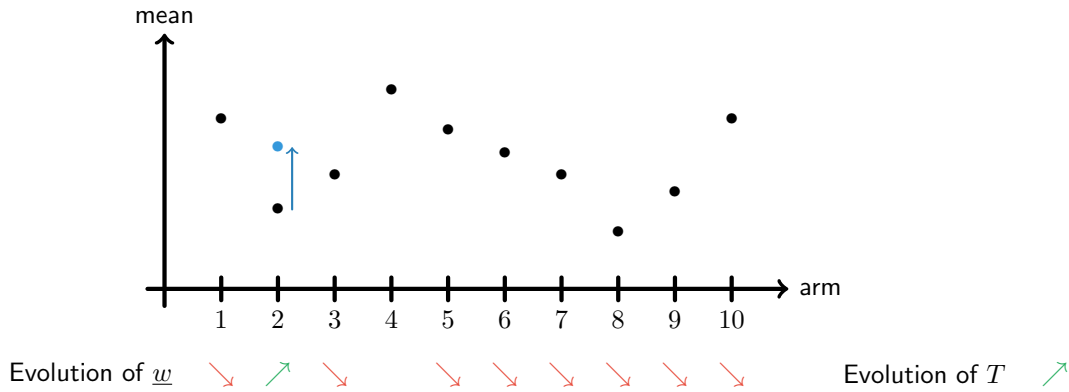


Figure 2.5: By increasing the mean of a sub-optimal arm, its associated optimal weight increases while those of other sub-optimal arms decrease, and the characteristic time increases.

2.2.5. The Track-and-Stop Algorithm: an Asymptotically Optimal Strategy

To obtain strategies that could reach lower bound (2.19), we have seen that we have no other choice than finding sampling rules which ensure that the empirical proportions of draws are close to $\underline{w}(\underline{\mu})$. Of course, the learner does not know $\underline{\mu}$, so she cannot directly compute $\underline{w}(\underline{\mu})$. However, by continuity of \underline{w} , a strategy can estimate it by using the current empirical estimate $\hat{\underline{\mu}}(t)$ available at step t . This is the main idea of the Track-and-Stop algorithm that we will now describe.

Sampling rule. At the beginning of time step t , we have access to $\hat{\underline{\mu}}(t-1)$ the maximum likelihood estimator of $\underline{\mu}$. Hence, using the Optimal-Weights procedure, we can compute the plug-in estimate $\hat{\underline{w}}(t-1) \stackrel{\text{def}}{=} \underline{w}(\hat{\underline{\mu}}(t-1))$ which is a good estimate of $\underline{w}(\underline{\mu})$ when $\hat{\underline{\mu}}(t)$ is close to $\underline{\mu}$, thanks to the continuity of \underline{w} . The sampling rule consists in *tracking* this proportion $\hat{\underline{w}}(t-1)$: it compares $\hat{\underline{w}}(t-1)$ with the empirical proportions $\frac{N_a(t-1)}{(t-1)}$ and chooses the arm which has been most under-sampled. There are two approaches:

- we can work in a cumulative way, by tracking the average $\frac{1}{t-1} \sum_{s \in [t-1]} \hat{\underline{w}}(s)$. This option, referred to as *C-tracking*, is the simplest to analyze theoretically as the variability between the two consecutive tracked weights is low,
- or we can use a direct approach by tracking the current weight $\hat{\underline{w}}(t-1)$. This option called *D-tracking* performs empirically better than C-tracking.

Unfortunately, due to the randomness, it might happen that some good arms with poor initial observations stay under-sampled and hence do not improve their estimates, which makes the strategy inefficient. To counter this issue¹¹, the strategy forces exploration by observing the least pulled arm among the set $U(t-1)$ of arms that are under-sampled with respect to the sub-linear rate \sqrt{t} :

$$U(t-1) \stackrel{\text{def}}{=} \left\{ a \in [K] : N_a(t-1) < \sqrt{t-1} - \frac{K}{2} \right\}.$$

The sub-linear rate \sqrt{t} is in fact arbitrary and ensures guarantees on the minimal convergence rate of the empirical mean to $\underline{\mu}$ (see Section 6.3.3). The sampling rule of Track-and-Stop is presented in Algorithm 4.

Remark. Note that we will present, in Contribution 2.13, a new sampling rule which modifies the tracking of Track-and-Stop in a natural way and which enjoys the same exploration rate \sqrt{t} without any force exploration mechanism.

¹¹We present here the way forced exploration is ensured with D-tracking. For C-tracking, the mechanism is different but we omit its presentation for the sake of simplicity.

Algorithm 4: Track-and-Stop sampling rule at step $t > K$

Input: history of observations I_{t-1}
Output: next arm to observe A_t

```

1  $U(t-1) \leftarrow \{a \in [K] : N_a(t-1) < \sqrt{t-1} - \frac{K}{2}\}$ 
2  $\hat{w}(t-1) \leftarrow \text{Optimal-Weights}(\hat{\mu}(t-1))$ 
3 if  $U(t-1) \neq \emptyset$  then
4   |  $A_t \leftarrow \underset{a \in U(t-1)}{\text{argmin}} N_a(t-1)$ 
5 else
6   | /* C-tracking */
7   | Choose  $A_t \in \underset{a \in [K]}{\text{argmin}} N_a(t-1) - \sum_{s \in [t-1]} \hat{w}_a(s)$ 
8   | /* D-tracking */
9   | Choose  $A_t \in \underset{a \in [K]}{\text{argmin}} N_a(t-1) - (t-1)\hat{w}_a(t-1)$ 

```

Stopping condition. Given observations up to time t , the algorithm should stop if it has statistical evidence that it found the optimal arm with risk δ . The problem can be seen as a statistical test. [Garivier and Kaufmann \(2016\)](#) introduced the generalized log-likelihood ratio (see [Chernoff, 1959](#)) between arms a and b :

$$Z_{a,b}(t) \stackrel{\text{def}}{=} \log \frac{\sup_{\substack{\lambda_a, \lambda_b \in \mathcal{M} \\ \lambda_a \geq \lambda_b}} \prod_{s=1}^{N_a(t)} p_{\lambda_a}(X_{a,s}) \prod_{s=1}^{N_b(t)} p_{\lambda_b}(X_{b,s})}{\sup_{\substack{\lambda_a, \lambda_b \in \mathcal{M} \\ \lambda_a \leq \lambda_b}} \prod_{s=1}^{N_a(t)} p_{\lambda_a}(X_{a,s}) \prod_{s=1}^{N_b(t)} p_{\lambda_b}(X_{b,s})},$$

where $p_\lambda = \frac{d\zeta}{d\rho}$ if ζ is the distribution of \mathcal{D}_{exp} of mean λ , and whose expression is given in (2.15). It may be shown that $Z_{a,b}(t)$ is the empirical transportation cost from arm b to arm a up to a factor t and the sign of $\hat{\mu}_a(t) - \hat{\mu}_b(t)$:

$$Z_{a,b}(t) = t \text{TC}_{b \rightarrow a} \left(\hat{\mu}(t), \frac{N(t)}{t} \right) \text{sgn} \left(\hat{\mu}_a(t) - \hat{\mu}_b(t) \right).$$

If $Z_{a,b}(t)$ is large enough then we can confidently reject $H_0: \mu_a \leq \mu_b$ and conclude to $H_1: \mu_a > \mu_b$. When considering all arms, we define¹²

$$\begin{aligned} Z(t) &\stackrel{\text{def}}{=} \max_{a \in [K]} \min_{b \neq a} Z_{a,b}(t) \\ &= \min_{b \neq a^*(\hat{\mu}(t))} Z_{a^*(\hat{\mu}(t)), b}(t) \\ &= t \min_{b \neq a^*(\hat{\mu}(t))} \text{TC}_{b \rightarrow a^*(\hat{\mu}(t))} \left(\hat{\mu}(t), \frac{N(t)}{t} \right) \end{aligned} \quad (2.25)$$

$$= \inf_{\lambda \in \text{Alt}(\hat{\mu}(t))} \sum_{a \in [K]} N_a(t) d(\mu_a, \lambda_a). \quad (2.26)$$

The Global-Likelihood-Ratio stopping rule (or *Chernoff stopping rule*), see Algorithm 5, associated to a given threshold function $\beta(t, \delta)$ consists in stopping the strategy as soon as $Z(t)$ exceeds

¹²Where $a^*(\hat{\mu}(t))$ is any of the optimal arms of $\hat{\mu}(t)$.

Algorithm 5: Global-Likelihood-Ratio stopping rule at step $t > K$

Input: history of observations I_t
 threshold function $\beta(t, \delta)$

- 1 $Z(t) \leftarrow \max_{a \in [K]} \min_{b \neq a} Z_{a,b}(t)$ // $Z_{a,b}(t)$ is defined in (4.7)
- 2 **if** $Z(t) > \beta(t, \delta)$ **then**
- 3 | Stop
- 4 **else**
- 5 | Continue

$\beta(t, \delta)$:

$$\tau_\delta \stackrel{\text{def}}{=} \inf \left\{ t \geq 1 : Z(t) > \beta(t, \delta) \right\}.$$

Of course, the associated decision rule will be to return the best empirical arm at time τ_δ (see the Empirical-Best Procedure 2). The stopping rule problem now reduces to the research of threshold functions which ensure δ -correctness, a statistical problem that does not depend on the sampling rule. Using Equation (2.26), we get $Z(t) \leq \frac{t}{T(\hat{\mu}(t))}$ and expect $Z(t) \sim \frac{t}{T(\mu)}$ for large values of t . For a time-independent threshold $\beta(t, \delta) = \beta_\delta$, this indicates that the strategy should stop after approximately $\tau_\delta \sim T(\mu)\beta_\delta$ steps, and the δ -correctness condition might force $\beta_\delta \gtrsim \log \frac{1}{\delta}$ at the sight of lower bound (2.19).


Obtaining δ -correctness is based on time-uniform bounds that require time-dependent thresholds. Garivier and Kaufmann (2016) proposed the following threshold using deviation results.

Theorem 2.9. [Garivier and Kaufmann, 2016, Proposition 12]

Consider an exponential model \mathcal{D}_{exp} . Let $\delta \in (0, 1)$ and $\alpha > 1$. There exists a constant $R = R(\alpha, K)$ such that, whatever the sampling rule, using the Global-Likelihood-Ratio stopping rule (Algorithm 5) with threshold

$$\beta(t, \delta) \stackrel{\text{def}}{=} \log \frac{Rt^\alpha}{\delta}, \quad (2.27)$$

and the Empirical-Best decision rule (Algorithm 2) ensures that the strategy is δ -correct.

The performance of a strategy highly depends on both its sampling and stopping rules. While many efficient sampling rules have been studied in the literature (see Section 2.2.7), the use of the Global-Likelihood-Ratio stopping rule appears to be consensual: it does not seem that it might be significantly improved, apart from slightly reducing the threshold. In the remainder of this chapter, except for elimination strategies, all considered fixed-confidence best-arm identification strategies will use the Global-Likelihood-Ratio stopping rule with threshold (2.27) for a given value of α and the Empirical-Best decision rule. The name of the strategy will be defined as the name of the sampling rule. 

Remark. • With threshold (2.27), the stopping time τ_δ is almost surely finite. Indeed, intuitively, for t large enough and under a good sampling rule, based on Equation (2.25) we get

$$Z(t) \simeq \frac{t}{T(\mu)},$$

hence $Z(t)$ grows much faster than $\beta(t, \delta)$.

- Obtaining δ -correctness with threshold (2.27) crucially depends on a condition on R that does not scale pretty well with the number of arms K . However, this is not an issue numerically:

more convenient thresholds, which are closer to $\log \frac{1}{\delta}$ like

$$\beta(t, \delta) = \log \frac{1 + \log t}{\delta},$$

are known to be empirically δ -correct (although there is no theoretical proof of this observation).

More generally, theoretical algorithms are often adapted in experiments to versions that perform better in practice (see, e.g., the use of D-tracking instead of C-tracking for the stopping rule of Track-and-Stop).

Open question. Finding theoretical arguments to obtain better thresholds — as close to $\log \frac{1}{\delta}$ as possible —, is an important statistical question, whose answer will get implications beyond bandit problems. Recent techniques using mixture martingales (see [Garivier and Kaufmann, 2021](#); [Kaufmann and Koolen, 2021](#) and very recently [Chowdhury et al., 2023](#)) proposed thresholds with better asymptotical behaviors which are however model-dependent and quite less explicit.

Asymptotic optimality. The Track-and-Stop sampling rule has been chosen so that the empirical frequencies of pulls converge to $\underline{w}(\underline{\mu})$, while the design of the Global-Likelihood-Ratio stopping-rule with threshold (2.27) ensures the δ -correctness. We may wonder how close the performance of the corresponding strategy is from lower bound (2.19). It turns out that, in the *asymptotic* regime when δ goes to 0, the upper bound of Track-and-Stop matches the lower bound.

Theorem 2.10. Consider an exponential model \mathcal{D}_{exp} . The Track-and-Stop strategy, with threshold (2.27) for a fixed $\alpha > 1$, satisfies, for all bandit problems $\underline{\mu}$ in \mathcal{D}_{exp} with a unique optimal arm,

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \leq T(\underline{\mu}).$$

Remark 2.11. The original result of [Garivier and Kaufmann \(2016, Theorem 14\)](#) only ensures that

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \leq \alpha T(\underline{\mu}).$$

The multiplicative factor α can be avoided using recent proof techniques, as we explain in Section 6.4.

This result indicates that the lower bound (2.19) cannot be improved in the asymptotic regime. The Track-and-Stop strategy is hence *asymptotically optimal* for all exponential models.

Definition 2.12. [asymptotically optimal strategy]

A strategy is said to be *asymptotically optimal* on a model \mathcal{D} if for all bandit problems $\underline{\mu}$ in \mathcal{D} with a unique optimal arm,

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \leq T(\underline{\mu}).$$

Empirical performance. The Track-and-Stop algorithm is not only a theoretical contribution, but it also proved to be numerically efficient: the strategy outperforms its competitors in a wide variety of settings.

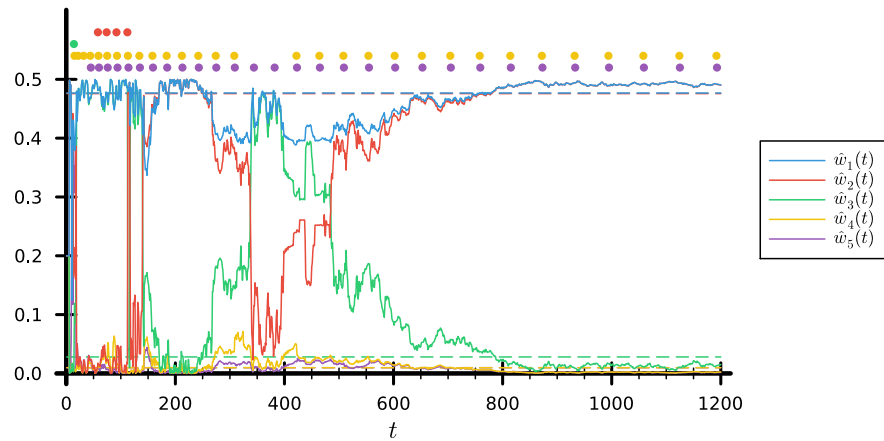


Figure 2.6: Instability of the targeted weights $\hat{w}(t)$ during the first 1200 rounds of a simulation of Track-and-Stop with parameters $\delta = 0.01$, $\mu = (0.9, 0.8, 0.6, 0.4, 0.4)$ and D-tracking. Dots correspond to the use of forced exploration for under-sampled arms, while the values of the optimal weights $w(\mu) = (0.477, 0.476, 0.028, 0.010, 0.010)$ are dashed.

2.2.6. Non-Asymptotic Guarantees

Some questions raised by Track-and-Stop. The Track-and-Stop strategy suffers from certain shortcomings illustrated in Figure 2.6:

- first, the sampling rule appears to be pretty unstable, especially at the beginning: the target frequencies can vary significantly as the estimated means fluctuate before stabilizing around their expectations,
- second, Track-and-Stop does not exhibit the intuitively desirable behavior to sample uniformly in the beginning, until sufficient information has been gathered for significant differences between the arms to emerge. This is in contrast with elimination strategies, which are sub-optimal but intuitively appealing.
- also, the forced exploration appears very arbitrary, with a rate of \sqrt{t} that has no other justification than lying somewhere between constant and linear functions.

Altogether, these issues lead for example to unpredictable and irregular behaviors at the beginning of multiple A/B testing cases with many arms extremely close to optimal.

Non-asymptotic guarantees.

Open question. As we have seen, Track-and-Stop is asymptotically optimal. The next natural direction is to obtain guarantees in the *non-asymptotic* regime, for fixed values of δ . A close look into the proofs of [Garivier and Kaufmann \(2016\)](#) shows that the theoretical guarantees proved so far are asymptotic in nature. Even if the strategy works pretty well empirically, one needs to develop different proof techniques to get non-asymptotic guarantees.

Note that lower bound (2.19) might not be well-suited for the moderate regime (values of δ that are not too small): [Simchowitz et al. \(2017\)](#) proved a lower bound for moderate values of δ that does not depend on the risk δ . This moderate setting has led to strategies that are sub-optimal by a multiplicative constant but are proved to satisfy explicit non-asymptotic bounds ([Karnin et al., 2013](#); [Jamieson et al., 2014](#); [Chen et al., 2017](#)). Yet, can we design asymptotically optimal strategies with finite risk bounds that match the asymptotic complexity?

Contribution 2.13. Due to the instability of the tracking procedure of Track-and-Stop that we described in the previous paragraph, it seems difficult to obtain theoretical non-asymptotic bounds. In order to address this problem, we present a novel algorithm named Exploration-Biased-Sampling in Chapter 4. This algorithm effectively resolves all the limitations discussed earlier and attains certain non-asymptotic assurances, specifically for a Gaussian model with a shared variance $\sigma^2 > 0$ and constrained means.

The exploration is conducted differently, in a statistically natural way that softens the fluctuations of empirical means and avoids arbitrary parameters. The idea is to introduce, at time step t , a confidence region $\mathcal{CR}_\mu(t)$ for $\underline{\mu}$, find the bandit problem $\tilde{\mu}(t)$ “maximizing exploration” (in a way that we define) inside $\mathcal{CR}_\mu(t)$, and track its associated optimal weight vector $\tilde{w}(t) \stackrel{\text{def}}{=} w(\tilde{\mu}(t))$. Based on the lemmas of Section 3.4, which give the evolution of the weights when moving one or several arms of a bandit problem, we give a procedure that allows computing $\tilde{\mu}(t)$. As the confidence regions shrink to $\{\underline{\mu}\}$ with time, the targeted weights will converge to $w(\underline{\mu})$ as for the Track-and-Stop algorithm. A major benefit is that the procedure results in a stabilized sampling strategy, which is much easier to follow and understand (see Figure 2.7).

This stabilization, together with the careful analysis of the quantitative regularity of the solution to the optimization problem (2.20) developed in Section 3, allows us to propose a non-asymptotic analysis of Exploration-Biased-Sampling with finite risk bounds. A simplification of the obtained bound reads that for all standard Gaussian bandit problems $\underline{\nu}$ with means in $[0, 1]$, there exist an event \mathcal{E} of high probability (independent from δ) and $\delta_0 > 0$ such that algorithm Exploration-Biased-Sampling with the threshold of Equation (2.27) satisfies

$$\forall \delta \in (0, \delta_0], \quad \mathbb{E}_\mu[\tau_\delta \mathbb{I}\{\mathcal{E}\}] \lesssim T(\underline{\mu}) \log \frac{1}{\delta} + o_{\delta \rightarrow 0}(1), \quad (2.28)$$

See Theorem 4.5 for a precise statement, with a closed-form expression for $o_{\delta \rightarrow 0}(1)$; note also that δ_0 depends, among others, on the probability of \mathcal{E} . We observe that bound (2.28) matches the asymptotic complexity. Similarly to Track-and-Stop, it can be shown that our strategy Exploration-Biased-Sampling is asymptotically optimal, although it is not a direct consequence of our non-asymptotic bound, as we considered $\mathbb{E}_\mu[\tau_\delta \mathbb{I}\{\mathcal{E}\}]$ instead of $\mathbb{E}_\mu[\tau_\delta]$.

In the same direction, [Degenne et al. \(2019\)](#) obtained a highly general and remarkable non-asymptotic bound for a pure exploration algorithm. Independently, [Wang et al. \(2021\)](#) obtained a sampling rule based on a Frank-Wolfe method for which they proved finite risk analysis and asymptotic optimality. Our bound has a better asymptotic behavior, but a worse behavior in the regime where gaps go to zero¹³ (see [Jourdan and Degenne, 2023](#), Table 1).

2.2.7. Towards Computationally More Efficient and More Natural Strategies

The sampling rule of the Track-and-Stop strategy requires computing, at each step, the solution of optimization problem (2.20). While a call to `Optimal-Weights` is not too costly experimentally (approximately proportional in the number of arms), obtaining easier sampling rules that do not need to solve the optimization problem at each step could allow us to obtain much more efficient strategies. In fact, one wants to find the most natural and simple sampling rules that ensure asymptotic optimality. In this section, we present sampling rules that go in that direction. We will still consider a general exponential model \mathcal{D}_{exp} and a bandit problem $\underline{\mu}$ in \mathcal{D}_{exp} with a unique optimal arm.

¹³However, we do not pretend that those bounds are comparable as we only considered $\mathbb{E}_\mu[\tau_\delta \mathbb{I}\{\mathcal{E}\}]$.

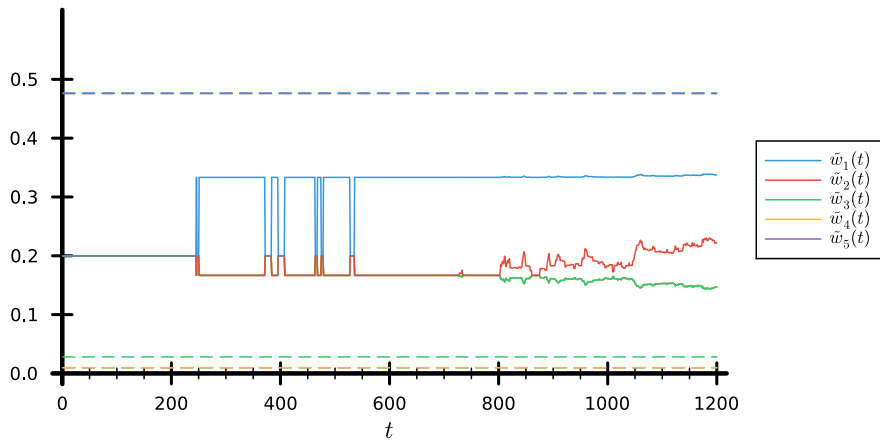


Figure 2.7: Stability of the targeted weights $\tilde{w}(t)$ during the first 1200 rounds of a simulation of Exploration-Biased-Sampling with parameters $\delta = 0.01$, $\mu = (0.9, 0.8, 0.6, 0.4, 0.4)$ and D-tracking. The values of the optimal weights $\underline{w}(\mu) = (0.477, 0.476, 0.028, 0.010, 0.010)$ are dashed. See Figure 2.6 for a comparison with Track-and-Stop.

Gradient-based sampling rules. The solution of optimization problem (2.20) can be sequentially approximated using (sub-)gradient optimization methods like the Frank-Wolfe algorithm. Two such sampling strategies, Lazy-Mirror-Ascent and Frank-Wolfe-based-Sampling, were respectively proposed by Ménard (2019) and Wang et al. (2021) with improved efficiency compared to Track-and-Stop, and proved to be asymptotically optimal.

top-two algorithms. A promising set of strategies are top-two algorithms which come with very simple sampling rules: at time step t , the algorithm chooses the next arm to sample A_t between two arms, namely a *leader* L_t and a *challenger* C_t . Originally, Russo (2016, 2020) considered an algorithm for which the choice of leader and challenger is done using Thompson Sampling:

- given some prior Π_0 on the value of $a^*(\mu)$, the leader is chosen according to the posterior Π_{t-1} computed given I_{t-1} ,
- the challenger is obtained similarly by drawing according to Π_{t-1} until obtaining a different arm than the leader.

The leader and challenger can also be chosen by non-Bayesian procedures (see, e.g., Jourdan et al., 2022). For the leader, we may cite the Empirical-Best (EB) and Upper-Confidence-Bound (UCB) procedures, which respectively select the arm with best empirical mean or best upper confidence bound on the mean. It might be natural to choose the challenger among arms $a \neq L_t$ minimizing the empirical transportation cost to the leader, which we call the Transportation-Cost (TC) challenger, or a version favoring exploration by adding a penalization to the cost of over-sampled arms, leading to the Transportation-Cost-Penalized (TCP) challenger. The previously mentioned leaders and challengers are gathered in Algorithms 7 and 8. The list is far from being exhaustive (see also Qin et al., 2017; Shang et al., 2020 for additional Bayesian procedures).

Remark. In the special case of best-arm identification, the LUCBtop-two algorithm, introduced a few years earlier by Kalyanakrishnan et al. (2012), can be seen as a Top-Two algorithm, in which the leader is pulled at even steps and the challenger is pulled at odd steps.

Algorithm 6: top-two sampling rule at step $t > K$

Input: history of observations I_{t-1}
 leader, challenger, sample-arm procedures

Output: next arm to observe A_t

- 1 $L_t \leftarrow \text{leader}(I_{t-1})$
 - 2 $C_t \leftarrow \text{challenger}(I_{t-1}, L_t)$
 - 3 $A_t \leftarrow \text{sample-arm}(I_{t-1}, L_t, C_t)$
-

Algorithm 7: leader procedures for top-two algorithms

Input: history of observations I_{t-1}
 optional: bonus function g , prior probability Π_0

Output: leader L_t

/ Empirical-Best (EB) leader */*

- 1 Choose $L_t \in \operatorname{argmax}_{a \in [K]} \hat{\mu}_a(t-1)$

/ Upper-Confidence-Bound (UCB) leader */*

- 2 Choose $L_t \in \operatorname{argmax}_{a \in [K]} \hat{\mu}_a(t-1) + \sqrt{\frac{g(t-1)}{N_a(t-1)}}$

/ Thompson-Sampling (TS) leader */*

- 3 Choose $L_t \sim \Pi_{t-1}$
-

Algorithm 8: challenger procedures for top-two algorithms

Input: history of observations I_{t-1}

leader L_t

Output: challenger C_t

/ Transportation-Cost (TC) challenger */*

- 1 Choose $C_t \in \operatorname{argmin}_{a \neq L_t} \text{TC}_{a \rightarrow L_t} \left(\hat{\mu}(t-1), \frac{N(t-1)}{t-1} \right)$

/ Transportation-Cost-Penalized (TCP) challenger */*

- 2 Choose

$$C_t \in \operatorname{argmin}_{a \neq L_t} \text{TC}_{a \rightarrow L_t} \left(\hat{\mu}(t-1), \frac{N(t-1)}{t-1} \right) + \log N_a(t-1)$$

/ Thompson-Sampling (TS) challenger */*


- 3 **repeat**

- 4 | Draw $C_t \sim \Pi_{t-1}$

- 5 **until** $C_t \neq L_t$
-

To define top-two strategies, it remains to choose which arm to sample between the leader and the challenger (see Algorithm 6). A simple way of doing so is to fix some probability parameter $\beta \in (0, 1)$ and choose the leader with that probability¹⁴:

$$A_t = \begin{cases} L_t & \text{with probability } \beta, \\ C_t & \text{otherwise.} \end{cases} \quad (2.29)$$

We will refer to top-two instances using this sampling rule as TT-leader-challenger- β . For instance, TT-EB-TCP- β denotes the top-two algorithm used with the Empirical-Best leader, the Transportation-Cost-Penalized challenger, and sampling rule (2.29). 

As a good leader might satisfy $L_t = a^*$ except for a sub-linear number of time steps, using sampling rule (2.29) comes with some limitation: it implies that the best arm will be pulled a fraction β of the time, which can be far from the optimal frequency $w_{a^*}(\underline{\mu})$. As a consequence, an easy adaption of lower bound (2.19) implies the following lower bound (see Russo, 2016) for such strategies:

$$\mathbb{E}_{\underline{\mu}}[\tau_\delta] \geq T_\beta(\underline{\mu}) \text{kl}(\delta, 1 - \delta), \quad \text{where } T_\beta(\underline{\mu})^{-1} \stackrel{\text{def}}{=} \sup_{\substack{v \in \Sigma_K \\ v_{a^*} = \beta}} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a), \quad (2.30)$$

and β -asymptotically optimal strategies are algorithms that pull the best arm with proportion β and having an upper bound asymptotically matching this lower bound. Of course, by uniqueness of $w(\underline{\mu})$, we get $T_\beta(\underline{\mu}) = T(\underline{\mu})$ if and only if $\beta = w_{a^*}(\underline{\mu})$. In other words, those algorithms can achieve asymptotic optimality only if $\beta = w_{a^*}(\underline{\mu})$. An arbitrary choice of parameter β is however not too dramatic: Russo (2016) proved that for $\beta = \frac{1}{2}$, one gets:

$$T_{\frac{1}{2}}(\underline{\mu}) \leq 2T(\underline{\mu}),$$

hence $\frac{1}{2}$ -asymptotically optimal strategies only loose a multiplicative factor 2 in theoretical performance compared to asymptotically optimal strategies.

Noticeably, to ensure β -asymptotical optimality of TT-leader-challenger- β algorithms, the conditions on the leader and challenger procedures are quite mild (see Jourdan et al., 2022); the combination of those proposed in Algorithms 7 and 8 are suitable.

Remark. The TT-EB-TC- β algorithm does not present any exploratory mechanism in the leader or challenger definitions. Even if it is β -asymptotically optimal, empirical performance at fixed values of δ suffers from some outliers at which the strategy fails to stop quickly. Strategies with an explorative procedure like TT-EB-TCP- β or TT-UCB-TC- β do not present this misbehavior and thus are preferred in practice.

Remark. Jourdan and Degenne (2023) also proved a non-asymptotic bound for Gaussian variables and the TT-UCB-TC- β algorithm, which asymptotically matches the β -optimality up to a factor $\frac{1}{\beta}$.

Adaptivity of top-two algorithms. Until recently, it has been possible to analyze top-two algorithms only with the use of sampling rule (2.29). It is quite frustrating that those top-two algorithms can only achieve β -optimality. To tackle this issue, we must rely on an adaptive sampling rule, like, for instance, adaptive pulling frequencies $(\beta(I_t, L_t, C_t))_{t \geq 1}$ of the leader. How to design such frequencies might be driven by the targeted proportions that we want to asymptotically reach, namely $w(\underline{\mu})$. See details in Section 6.2.

Recently, You et al. (2023) proposed an adaptive top-two algorithm based on the original Bayesian leader and challenger of Russo (2016) called Top-Two-Thompson-Sampling. Their analysis proved the asymptotical optimality of the strategy for standard Gaussian variables. The proof can be adapted to obtain the same guarantees for adaptive versions of TT-EB-TC and TT-EB-TCP .

¹⁴One can also use a tracking procedure of proportions β instead of a sampling.

Contribution 2.14. Based on the transformation \underline{W} introduced in Section 6.2.1, we propose in Section 6.2.3 a new challenger procedure together with adaptive pulling frequencies.

We also present partial results that generalize the analysis of You et al. (2023) in order to prove that those algorithms are asymptotically optimal for all exponential models (not only Gaussian distributions with known and common variance). See Section 6.3 for details. Yet, some work remains to complete the analysis.

2.3. Best-Arm Identification with a Fixed-Budget

We now move to the problem of best-arm identification with a fixed-budget. While the identification task seems similar to the fixed-confidence setting, it turns out that this setting is much less understood. Strategies are defined similarly to the fixed-confidence setting, except that the stopping time τ is not random and is equal to a given deterministic, known-in-advance, budget $T \geq 1$ (see Algorithm 9). The strategy returns an estimate \hat{a}_T of $a^*(\nu)$, and its quality is measured by the *probability of misidentification* or *probability of error*

$$\mathbb{P}_\nu(\hat{a}_T \neq a^*(\nu)).$$

Algorithm 9: Structure of a fixed-budget strategy

Input: budget parameter T

sampling-rule and decision-rule

Output: estimated best arm \hat{a}_T

```

1 for  $t \in [T]$  do
2    $A_t \leftarrow \text{sampling-rule}(I_{t-1})$ 
3   Observe  $Y_t \sim \nu_{A_t}$ 
4  $\hat{a}_T \leftarrow \text{decision-rule}(I_T)$ 
    
```

2.3.1. An Exponential Decay Rate

The simplest sampling, **Uniform-Sampling**, consists of pulling all arms equally often and returning the best empirical estimate. We will see, by using Hoeffding's inequality, that the error probability vanishes exponentially fast as T goes to $+\infty$ for a σ^2 -sub-Gaussian model \mathcal{D}_{σ^2} .

Analysis of the Uniform-Sampling strategy. We can upper bound the probability of error when comparing arms a^* and $a \neq a^*$ by applying Hoeffding's inequality (2.4), similarly to what we did in Equation (2.5) with the Successive-Elimination strategy. This gives, as all arms are pulled $\lfloor \frac{T}{K} \rfloor$ times

$$\begin{aligned}
 \mathbb{P}_\nu(\hat{a}_T = a) &\leq \mathbb{P}_\nu\left(\hat{\mu}_a(T) \geq \max_{b \neq a} \hat{\mu}_b(T)\right) \\
 &\leq \mathbb{P}_\nu\left(\hat{\mu}_a(T) \geq \hat{\mu}_{a^*}(T)\right) \\
 &\leq \mathbb{P}_\nu\left(\left|\frac{T}{K}\right| \left(\hat{\mu}_a(T) - \hat{\mu}_{a^*}(T) + \Delta_a\right) \geq \left|\frac{T}{K}\right| \Delta_a\right) \\
 &\leq \exp\left(-\frac{\Delta_a^2}{4\sigma^2} \left|\frac{T}{K}\right|\right). \tag{2.31}
 \end{aligned}$$

Applying a union bound, we obtain the following bound for the probability of misidentification:

$$\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \leq (K - 1) \max_{a \neq a^*(\underline{\nu})} \mathbb{P}_{\underline{\nu}}(\hat{a}_T = a) = (K - 1) \exp\left(-\frac{\Delta_{\min}^2}{4\sigma^2} \left\lfloor \frac{T}{K} \right\rfloor\right), \quad (2.32)$$

where we recall that $\Delta_{\min} = \min_{a \neq a^*(\underline{\nu})} \Delta_a$.

Asymptotic rate of the exponential decay. We proved that the probability of misidentification of the `Uniform-Sampling` strategy goes to 0 exponentially fast with respect to the budget T . From now on, we will thus be interested in obtaining the best rate for this exponential decay, that is, in minimizing the quantity

$$\frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})).$$

In the sequel, we will mainly work in the *asymptotic* regime where the budget T goes to $+\infty$ and hence consider *sequences* of strategies¹⁵ indexed by their budgets T . On the one hand, we focus on obtaining instance-dependent lower bounds, valid for “reasonable” sequences of strategies, on the quantities

$$\ell(\underline{\nu}) \stackrel{\text{def}}{=} \liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})),$$

On the other hand, for a given sequence of strategies, we look for upper bounds of

$$u(\underline{\nu}) \stackrel{\text{def}}{=} \limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})).$$

Note that $\ell(\underline{\nu}) \leq u(\underline{\nu})$ and that those rates are negative.

Example. Given Equation (2.32), we proved the upper bound on the asymptotic rate of the sequence of `Uniform-Sampling` strategies:

$$\forall \underline{\nu} \text{ in } \mathcal{D}_{\sigma^2}, \quad \limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{\Delta_{\min}^2}{4\sigma^2 K}. \quad (2.33)$$

2.3.2. Successive-Rejects-Type Strategies

Can we do better than `Uniform-Sampling`? Upper bound (2.33) only depends on the minimal gap Δ_{\min} and the number of arms, that is, one can change the values of the means of the $K - 2$ worst arms without modifying the bound. A simple look into Equation (2.31) shows that the exponential decay of the probability of recommending sub-optimal arm a depends crucially on its gap Δ_a . In order to equalize this exponential rate among sub-optimal arms, one can be tempted to pull arm a a fraction proportional to $\frac{1}{\Delta_a^2}$ of the budget. Unfortunately, this requires good estimates of the gaps that the strategy cannot access.

This idea of having a (possibly pre-defined) different number of pulls for each arm was already studied in the fixed-confidence setting with `Successive-Elimination` strategies (see Section 2.2.1). Those strategies have also been considered in the fixed-budget setting by [Audibert et al. \(2010\)](#) under the name `Successive-Rejects`, as we now discuss. We continue to work with a σ^2 -sub-Gaussian model \mathcal{D}_{σ^2} .

¹⁵In fact, we did the same in the fixed-confidence regime when considering asymptotic behaviors when δ goes to 0.

Algorithm 10: Successive-Rejects algorithm

Input: budget parameter T

 phase lengths $(\ell_r)_{r \in [K-1]}$ such that $\sum_{r \in [K-1]} \ell_r = T$
Output: estimated best arm \hat{a}_T

```

1  $t \leftarrow 0$ 
2  $S_0 \leftarrow [K]$ 
3 for each round  $r \in [K - 1]$  do
4     Observe each arm  $\lfloor \frac{\ell_r}{K-r+1} \rfloor$  times
5     Increase  $t$  by  $\ell_r$ 
6     Choose  $a_r \in \operatorname{argmin}_{a \in S_{r-1}} \hat{\mu}_a(t)$ 
7      $S_r \leftarrow S_{r-1} \setminus \{a_r\}$ 
8 Define  $\hat{a}_T$  as the unique element of  $S_{K-1}$ 
    
```

Successive-Rejects strategies. We keep the notation of the Successive-Elimination algorithm (see Algorithm 3). Namely, the strategy works in $K - 1$ phases: at round r , candidate arms in a set S_{r-1} are pulled uniformly and the worst empirical arm of S_{r-1} (since the beginning) is dropped to obtain the set of next candidates S_r . We will consider versions of Successive-Rejects in which the lengths of the phases are set beforehand; they are denoted by $\ell_1, \dots, \ell_{K-1} \geq 1$ and satisfy $\ell_1 + \dots + \ell_{K-1} = T$. More precisely, during phase $r \in [K - 1]$, the strategy draws $\lfloor \frac{\ell_r}{K-r+1} \rfloor$ times each of the $K - r + 1$ arms in S_{r-1} (and does not use the few remaining time steps, if there are some). At the end of phase r , an arm of S_{r-1} has been pulled

$$n_r = \left\lfloor \frac{\ell_1}{K} \right\rfloor + \dots + \left\lfloor \frac{\ell_r}{K - r + 1} \right\rfloor$$

times since the beginning of the first phase. The description of the strategy is summarized in Algorithm 10.

Combining adapted versions of decomposition (2.6) and inequality (2.8) in the analysis of Successive-Elimination to σ^2 -sub-Gaussian variables as in (2.31), we prove the following bound on the probability of error of Successive-Rejects

$$\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*) \leq K \sum_{r \in [K-1]} \exp\left(-\frac{\Delta_{(K-r+1)}^2}{4\sigma^2} n_r\right) \leq K^2 \max_{r \in [K-1]} \exp\left(-\frac{\Delta_{(K-r+1)}^2}{4\sigma^2} n_r\right),$$

where we recall that we use the (reverse) notation of order statistics:

$$\mu_{(1)} > \mu_{(2)} \geq \mu_{(3)} \geq \dots \geq \mu_{(K)}.$$

The bound can be rewritten, in terms of rates (recall that n_r depends on T),

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*) \leq -\frac{1}{4\sigma^2} \min_{r \in [K-1]} \left\{ \Delta_{(K-r+1)}^2 \liminf_{T \rightarrow +\infty} \frac{n_r}{T} \right\}. \quad (2.34)$$

In order to exploit this general upper bound, one needs to carefully choose the phase lengths $\ell_1, \dots, \ell_{K-1}$ (or equivalently the number of pulls n_1, \dots, n_{K-1}).

With the knowledge of the gaps. As in the fixed-confidence setting (see page 38), we can select the phase lengths giving the best theoretical bound with the knowledge of the set of gaps. This boils down to equalizing the quantities appearing in the minimum of inequality (2.34), i.e., to take $n_r \propto \Delta_{(K-r+1)}^{-2}$. In order to not exceed the budget T , we should take

$$\forall r \in [K-1], \quad n_r \stackrel{\text{def}}{=} \left[\left(\sum_{a \neq a^*} \frac{1}{\Delta_{(a)}^2} + \frac{1}{\Delta_{(2)}^2} \right)^{-1} \frac{T}{\Delta_{(K-r+1)}^{-2}} \right].$$

The bound (2.34) reads

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*) \leq -\frac{1}{4\sigma^2} \left(\sum_{a \neq a^*} \frac{1}{\Delta_{(a)}^2} + \frac{1}{\Delta_{(2)}^2} \right).$$

Remark. The complexity appearing in the right-hand side was already involved in the expression (2.9) giving the value of τ_δ for the Successive-Elimination algorithm.

With the phase lengths of Audibert et al. (2010). How much do we lose without the knowledge of the gaps? The Successive-Rejects strategy was studied by Audibert et al. (2010) with the pre-defined phase lengths

$$\ell_1 \stackrel{\text{def}}{=} \frac{T}{\log K}, \quad \text{and} \quad \forall r \in \{2, \dots, K-1\}, \quad \ell_r \stackrel{\text{def}}{=} \frac{T}{(K-r+2) \log K}, \quad (2.35)$$

where we define

$$\overline{\log K} \stackrel{\text{def}}{=} \frac{1}{2} + \sum_{k=2}^K \frac{1}{k},$$

which ensures that the phase lengths sum to T . In a nutshell, the first phase is the longest and is devoted to a uniform exploration of all arms, and then the $K-2$ next rounds share the rest of the budget with a slight increase in length from one phase to the next (see Figure 2.8).

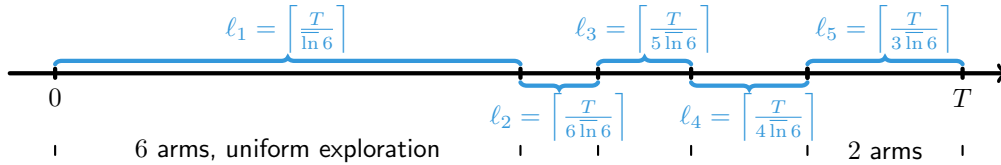


Figure 2.8: Successive-Rejects phase lengths of Audibert et al. (2010) for $K = 6$ arms.

Substituting the numbers of pulls $(n_r)_{r \in [K-1]}$ induced by the choice of phase lengths (2.35) into upper bound (2.34) yields the following result.

Theorem 2.15. [Audibert et al., 2010, Theorem 2]

The sequence of Successive-Rejects strategies with phase lengths given by Equation (2.35) satisfies, for all $\underline{\nu}$ in \mathcal{D}_{σ^2} with a unique optimal arm:

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{1}{4\sigma^2 \log K} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}. \quad (2.36)$$

Remark. • Audibert et al. (2010) also obtained a lower bound (that we describe in page 64) involving the quantity $\min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}$. The phase lengths (2.35) are in fact chosen so as to obtain the same gap-based quantity in the upper bound, which results in the apparition of the normalizing constant $\overline{\log K}$.

- In all generality, bounds (2.33) and (2.36) are not comparable: taking a bandit problem with all gaps equal will lead to a better bound for `Uniform-Sampling` while taking one tiny gap compared to the others gives an advantage to `Successive-Rejects` (as intuition suggests). Still, bound (2.36) captures the bandit structure by considering the value of all gaps and is smaller than (2.33) in a lot of regimes. Experimentally, `Uniform-Sampling` does not outperform `Successive-Rejects` even in the first situation.
- [Karnin et al. \(2013\)](#) proposed `Sequential-Halving`, a variant of `Successive-Rejects` in which half of the remaining arms are dropped at each round. Its theoretical guarantees are similar to the ones for `Successive-Rejects` (up to a factor $\log 2$ in the exponential rate), while its experimental behavior shows slight improvement with a large number of arms.

Despite being the first non-trivial proposed strategy for the fixed-budget setting, the instance-dependent guarantee of `Successive-Rejects` has not been significantly improved. It seems, however, that this gap-based upper bound might be non-optimal, as we will discuss in Section 2.3.4.

Contribution. The upper bound (2.36) for the `Successive-Rejects` strategy only applies to sub-Gaussian models. We will show in Chapter 5.3 how to generalize it to other models (see Contribution 2.22 for more details). This generalization will reveal a new information-theoretic quantity replacing gaps that we present in page 67.

2.3.3. Lower Bounds

We now discuss existing lower bounds in the fixed-budget setting.

Avoiding bad strategies. Considering the class of all strategies precludes obtaining lower bounds that hold simultaneously for all bandit instances. Indeed, a strategy that always recommends a given arm (independently of all observations) will make no error for problems admitting that arm as the best, but of course, such strategies are poor, as they completely fail on any other bandit instance. To prevent this, it is possible to modify the objective and lower bound, for a given instance $\underline{\nu}$, the maximal rate of error probability among permuted instances $(\underline{\nu}^\sigma)_{\sigma \in \mathfrak{S}_K}$, where \mathfrak{S}_K denotes the set of permutations of $[K]$ and $\underline{\nu}^\sigma \stackrel{\text{def}}{=} (\nu_{\sigma(1)}, \dots, \nu_{\sigma(K)})$:

$$\max_{\sigma \in \mathfrak{S}_K} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}^\sigma}(\hat{a}_T \neq a^*(\underline{\nu}^\sigma)).$$

This is not really satisfying in terms of writing, but sometimes handy (see [Audibert et al., 2010](#)). Yet, a more natural way of proceeding is to restrict the considered class of strategies¹⁶.

(Exponentially) consistent sequences of strategies. Recall that, in the asymptotic point of view, we consider sequences of strategies indexed by $T \geq 1$ given a value of $K \geq 2$. We will assume that these sequences are “reasonable” in the sense below. The probability $\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu}))$ of misidentifying the unique optimal arm may vanish asymptotically (and even vanish exponentially fast) for all bandit problems—in not too large a model \mathcal{D} —, as illustrated in Section 2.3.1. We will therefore only be interested in such sequences of strategies, called (exponentially) consistent.

¹⁶As in the regret literature, where lower bounds are derived for uniformly fast convergent strategies (see, e.g., [Burnetas and Katehakis, 1996](#)).

Definition 2.16. [(exponentially) consistent sequence of strategies]

Fix $K \geq 2$. A sequence of strategies is *consistent*, respectively, *exponentially consistent*, on a model \mathcal{D} if for all problems $\underline{\nu}$ in \mathcal{D} with a unique optimal arm,

$$\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \xrightarrow{T \rightarrow +\infty} 0, \quad \text{respectively,} \quad \limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) < 0.$$

The fundamental inequality. The fundamental inequality (2.12) that we presented in the fixed-confidence setting can also be used in fixed-budget best-arm identification to obtain lower bounds on the rate of the exponential decay. It entails the following asymptotic lemma.

Lemma 2.17. Fix $K \geq 2$ and a model \mathcal{D} . Consider a consistent sequence of strategies on \mathcal{D} , and two bandit problems $\underline{\nu}$ and $\underline{\zeta}$ in \mathcal{D} with unique and distinct optimal arms such that $a^*(\underline{\zeta}) \neq a^*(\underline{\nu})$. Then

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \limsup_{T \rightarrow +\infty} \sum_{a \in [K]} \frac{\mathbb{E}_{\underline{\zeta}}[N_a(T)]}{T} \text{KL}(\zeta_a, \nu_a). \quad (2.37)$$

Proof. We will first prove the following non-asymptotic lower bound, which does not make any assumption on the strategy:

$$\log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \frac{\sum_{a \in [K]} \mathbb{E}_{\underline{\zeta}}[N_a(T)] \text{KL}(\zeta_a, \nu_a) + \log 2}{\mathbb{P}_{\underline{\zeta}}(\hat{a}_T = a^*(\underline{\zeta}))}. \quad (2.38)$$

We apply the fundamental inequality (2.12) with $E = \{\hat{a}_T \neq a^*(\underline{\nu})\}$, and by exchanging the roles of $\underline{\nu}$ and $\underline{\zeta}$:

$$\sum_{a \in [K]} \mathbb{E}_{\underline{\zeta}}[N_a(T)] \text{KL}(\zeta_a, \nu_a) \geq \text{kl}\left(\mathbb{P}_{\underline{\zeta}}(\hat{a}_T \neq a^*(\underline{\nu})), \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu}))\right).$$

Combining this inequality with the following bound satisfied by the Kullback-Leibler divergence of the Bernoulli model¹⁷:

$$\forall p, q \in (0, 1), \quad \text{kl}(p, q) \geq p \log \frac{1}{q} - \log 2, \quad \text{i.e.,} \quad \log q \geq - \frac{\text{kl}(p, q) + \log 2}{p}, \quad (2.39)$$

and the fact that $\mathbb{P}_{\underline{\zeta}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq \mathbb{P}_{\underline{\zeta}}(\hat{a}_T = a^*(\underline{\zeta}))$ as $a^*(\underline{\nu}) \neq a^*(\underline{\zeta})$, we obtain

$$\begin{aligned} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) &\geq - \frac{\sum_{a \in [K]} \mathbb{E}_{\underline{\zeta}}[N_a(T)] \text{KL}(\zeta_a, \nu_a) + \log 2}{\mathbb{P}_{\underline{\zeta}}(\hat{a}_T \neq a^*(\underline{\nu}))} \\ &\geq - \frac{\sum_{a \in [K]} \mathbb{E}_{\underline{\zeta}}[N_a(T)] \text{KL}(\zeta_a, \nu_a) + \log 2}{\mathbb{P}_{\underline{\zeta}}(\hat{a}_T = a^*(\underline{\zeta}))}, \end{aligned}$$

which concludes the proof of (2.38).

Inequality (2.37) follows simply by considering asymptotics of (2.38)¹⁸, and the fact that the considered sequence of strategies is consistent on \mathcal{D} :

$$\mathbb{P}_{\underline{\zeta}}(\hat{a}_T = a^*(\underline{\zeta})) \xrightarrow{T \rightarrow +\infty} 1, \quad \text{and} \quad \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \xrightarrow{T \rightarrow +\infty} 0. \quad \square$$

¹⁷As $\text{kl}(p, q) = p \log \frac{1}{q} + (1-p) \log \frac{1}{1-q} + \underbrace{p \log p + (1-p) \log(1-p)}_{\geq -\log 2}$.

¹⁸One can directly use the asymptotic of the kl and obtain the result without inequality (2.39).

Remark 2.18. By using inequality (2.39), we are forced to exchange the roles of $\underline{\nu}$ and $\underline{\zeta}$ when applying the fundamental inequality (2.12). This is why expectations of pulls are now under the alternative bandit problem $\underline{\zeta}$. The task of transforming those inequalities into informative lower bounds gets harder: not only do we need to find alternative bandit problems close to $\underline{\nu}$ in terms of Kullback-Leibler divergences, but we also need to simultaneously control the pull frequencies under $\underline{\zeta}$. Note also that, as a consequence, we cannot use similar techniques as in the fixed-confidence proof of lower bound (2.19), where in inequality (2.22) we replaced the unknown frequency vector $\mathbb{E}_{\underline{\nu}}[N(\tau_\delta)]/\mathbb{E}_{\underline{\nu}}[\tau_\delta]$ by the “best” vector $\underline{\nu} \in \Sigma_K$, which was only possible because this frequency vector did not depend on $\underline{\zeta}$; hence, at least at first sight, we cannot conclude the existence of an optimal weight vector associated to a bandit problem $\underline{\nu}$ as in the fixed-confidence case. As we will discuss later, it does not seem possible to avoid those issues unless for models for which the Kullback-Leibler divergence is symmetric.

Lower bound of Audibert et al. (2010). Audibert et al. (2010) proposed the following lower bound for the model

$$\mathcal{B}_{[p, 1-p]} \stackrel{\text{def}}{=} \{\text{Ber}(x) : x \in [p, 1-p]\}$$

of Bernoulli distributions $\text{Ber}(x)$ with parameters x in $[p, 1-p]$ for some $p \in (0, \frac{1}{2})$,

Theorem 2.19. [Audibert et al., 2010, Theorem 4]

Consider the model $\mathcal{B}_{[p, 1-p]}$ for some $p \in (0, \frac{1}{2})$. For all strategies, and for all $\underline{\nu}$ in $\mathcal{B}_{[p, 1-p]}$ with a unique optimal arm,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \left(\max_{\sigma \in \mathfrak{S}_K} \mathbb{P}_{\underline{\nu}^\sigma}(\hat{a}_T \neq a^*(\underline{\nu}^\sigma)) \right) \geq -\frac{5}{p(1-p)} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}. \quad (2.40)$$

The bound involves the same quantity $\min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}$ as in the upper bound (2.36) of Successive-Rejects. Yet, these two bounds do not match as there is a factor $\overline{\log K} \simeq \log K$ between both. We will compare upper and lower bounds deeply in Section 2.3.4.

The Bretagnolle-Huber technique by Kaufmann et al. (2016, Section 5.2). A specific lower bound was obtained by Kaufmann et al. (2016) for the model $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ of Gaussian distributions with common variance $\sigma^2 > 0$, based on the Bretagnolle-Huber inequality (Bretagnolle and Huber, 1979) instead of inequality (2.39):

$$\forall p, q \in (0, 1), \quad p + 1 - q \geq \frac{1}{2} \exp(-\text{kl}(p, q)).$$

Remark. We used inequality (2.39) in the regime $p \rightarrow 1$ and $q \rightarrow 0$. The Bretagnolle-Huber inequality should be applied in the reverse regime where $p \rightarrow 0$ and $q \rightarrow 1$, giving a similar lower bound to (2.39) except that the logarithm depends on both p and q :

$$\log(2(p + 1 - q)) \geq -\text{kl}(p, q).$$

Kaufmann et al. (2016) proved that the rate was driven by the sum of the inverse squared gaps

$$H_\Sigma(\underline{\nu}) \stackrel{\text{def}}{=} \sum_{a \neq a^*(\underline{\nu})} \frac{1}{\Delta_a^2}. \quad (2.41)$$

More precisely, an asymptotic version of their result is the following.

Proposition 2.20. [Asymptotic version of Kaufmann et al., 2016, Theorem 16]

For all sequences of strategies and for all bandit problems $\underline{\nu}$ in $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ with a unique optimal arm, there exists a set of alternative bandit instances $(\underline{\nu}^{(k)})_{k \neq a^*(\underline{\nu})}$ in $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$, where each $\underline{\nu}^{(k)}$ admits k as a best arm and satisfies $H_{\Sigma}(\underline{\nu}^{(k)}) \leq H_{\Sigma}(\underline{\nu})$, and for which

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \max \left\{ \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})), \max_{k \neq a^*(\underline{\nu})} \mathbb{P}_{\underline{\nu}^{(k)}}(\hat{a}_T \neq k) \right\} \geq -\frac{2}{\sigma^2} H_{\Sigma}(\underline{\nu})^{-1}. \quad (2.42)$$

We will give a proof and a precise interpretation of this result in Section 5.6.2, but it conveys the intuition that

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \gtrsim -\frac{2}{\sigma^2} H_{\Sigma}(\underline{\nu})^{-1}.$$

The proposed lower bounds (2.40) and (2.42) both involve a variance term but rely on gap-based complexity terms that differ, namely, $H_{\Sigma}(\underline{\nu})$ and

$$H_{\max}(\underline{\nu}) = \left(\min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k} \right)^{-1} = \max_{2 \leq k \leq K} \frac{k}{\Delta_{(k)}^2}.$$

The second quantity also appears in the upper bound (2.36), somewhat in an arbitrary way as it crucially depends on the phase lengths, and both quantities can be linked as follow

$$H_{\max}(\underline{\nu}) \leq H_{\Sigma}(\underline{\nu}) \leq \log(2K) H_{\max}(\underline{\nu}),$$

with the two inequalities being tight in all generality (there exist instances for which equalities hold).

Contribution. In Chapter 5, we will introduce new information-theoretic complexity measures that replace H_{\max} and H_{Σ} . Those quantities will generalize the existing lower and upper bounds to more models including, e.g., non-parametric models. See Contribution 2.22 for more details.

Generalizations to other models. In order to understand if a complexity measure is a good candidate for being involved in the complexity of the fixed-budget setting, one may wonder if the proof techniques might be applied to more general models. Indeed, the lower bounds presented so far in this section are specific to Bernoulli and Gaussian.

Concerning the bound (2.42) of Kaufmann et al. (2016), a close look at the proof reveals that it heavily relies on a property even stronger than the symmetry of the Kullback-Leibler divergence for this model. In particular, generalizations beyond the Gaussian case appear to be infeasible.

Contribution. A detailed discussion on this statement might be found in Section 5.6.2, together with a proof of Proposition 2.20.

For the lower bound (2.40) of Audibert et al. (2010), however, a key inequality in their proof follows from the Kullback-Leibler $-\chi^2$ -divergence bound:

$$\forall x, y \in [p, 1-p], \quad \text{kl}(x, y) \leq \frac{(x-y)^2}{2p(1-p)}. \quad (2.43)$$

The construction may actually be generalized to models \mathcal{D} with $C_{\mathcal{D}} > 0$ such that

$$\forall \nu, \nu' \text{ in } \mathcal{D}, \quad \text{KL}(\nu, \nu') \leq C_{\mathcal{D}} (\mathbb{E}(\nu) - \mathbb{E}(\nu'))^2. \quad (2.44)$$

This is a property that clearly holds for some exponential families: on top of the restricted Bernoulli model discussed above, where $C_{\mathcal{B}_{[p, 1-p]}} = \frac{1}{2p(1-p)}$ by (2.43), we may cite the model \mathcal{D}_{σ^2} of Gaussian distributions with variance σ^2 , for which $C_{\mathcal{D}_{\sigma^2}} = \frac{1}{2\sigma^2}$.

A generalization of Theorem 2.19 to such models reads as follows.

Proposition 2.21. Consider a model \mathcal{D} for which property (2.44) holds with $C_{\mathcal{D}} > 0$. Then

$$\forall \underline{\nu} \text{ in } \mathcal{D}, \quad \liminf_{T \rightarrow +\infty} \frac{1}{T} \log \left(\max_{\sigma \in \mathfrak{S}_K} \mathbb{P}_{\underline{\nu}^\sigma} (\hat{a}_T \neq a^*(\underline{\nu}^\sigma)) \right) \geq -5 C_{\mathcal{D}} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}.$$

2.3.4. Comparing Upper and Lower Bounds: the Challenges of the Fixed-Budget Setting

We now compare the lower bounds proposed in the previous section with the upper bound of the Successive-Rejects strategy. This comparison indicates important directions of research for the fixed-budget setting.

Questions. This setting is much less understood than the fixed-confidence setting.

1. We actually do not know if there exists a complexity for the fixed-budget setting: there is a $\sim \log K$ multiplicative gap between the lower and upper bounds (2.40) and (2.36). Can we obtain matching lower and upper bounds even for simple models like Gaussian or Bernoulli models? Is there an optimal strategy that asymptotically reaches the lower bounds simultaneously in all instances? We will discuss those questions in the next section.
2. The bounds so far only depend on the gaps between arms, whereas they involve the Kullback-Leibler divergence in the fixed-confidence setting, which provides a more precise quantification of the difficulty in terms of the geometry of information of the problem. Can we obtain bounds that rely on more informative quantities than gaps?
3. Most of the proposed bounds are specific to a few models (even the generalization proposed in Proposition 2.21 is limited to models with a particular structure). Can we generalize the proof techniques to more general models, like an exponential model or even a non-parametric model?

From now on, let us forget about the first point that we will consider in the next section, and see how to tackle the two other questions. It turns out that the existing literature for the fixed-budget setting offered so far a non-parametric lower bound, in the case of $K = 2$ arms.

The non-parametric bound of Kaufmann et al. (2016) for $K = 2$ arms. Namely, in a general, possibly non-parametric, model \mathcal{D} , Kaufmann et al. (2016, Theorem 12) stated a lower bound for all 2-armed bandit problems $\underline{\nu} = (\nu_1, \nu_2)$, and for all consistent sequence of strategies:

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}} (\hat{a}_T \neq a^*(\underline{\nu})) \geq - \inf_{\underline{\zeta} \in \text{Alt}(\underline{\nu})} \max \{ \text{KL}(\zeta_1, \nu_1), \text{KL}(\zeta_2, \nu_2) \}, \quad (2.45)$$

where we recall that $\text{Alt}(\underline{\nu})$ is the set of alternative bandit problems¹⁹ to $\underline{\nu}$:

$$\text{Alt}(\underline{\nu}) \stackrel{\text{def}}{=} \{ \underline{\zeta} \text{ in } \mathcal{D} : a^*(\underline{\zeta}) \neq a^*(\underline{\nu}) \}.$$

The bound can be directly derived by taking the infimum over alternative bandit problems $\underline{\zeta}$ in lower bound (2.37):

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}} (\hat{a}_T \neq a^*(\underline{\nu})) \geq - \inf_{\underline{\zeta} \in \text{Alt}(\underline{\nu})} \limsup_{T \rightarrow +\infty} \underbrace{\sum_{a=1}^2 \frac{\mathbb{E}_{\underline{\zeta}} [N_a(T)]}{T} \text{KL}(\zeta_a, \nu_a)}_{\leq \max \{ \text{KL}(\zeta_1, \nu_1), \text{KL}(\zeta_2, \nu_2) \}}.$$

¹⁹Already defined in the fixed-confidence setting with mean-parametrized instances, see (2.18).

New information-theoretic quantities. We have seen in lower bound (2.37) that the arguments of the involved Kullback-Leibler divergences were in reversed order compared to the fixed-confidence setting.

Remark. Recall that we have seen in Remark 2.7 that the key quantities for the non-parametric study of best-arm identification with fixed confidence are defined based on Kullback-Leibler divergences $\mathcal{K}_{\text{inf}}^<$ and $\mathcal{K}_{\text{inf}}^>$ with arguments in reverse order. Note also that optimal bound regret-minimization literature only depends on the $\mathcal{K}_{\text{inf}}^<$ (see, e.g., [Honda and Takemura, 2015](#); [Garivier et al., 2022](#)).

Therefore, let us introduce, for a distribution $\nu \in \mathcal{D}$ and a real number $x \in \mathbb{R}$,

$$\mathcal{L}_{\text{inf}}^<(x, \nu) = \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) < x\},$$

and symmetrically, by considering rather distributions ζ with expectations larger than x ,

$$\mathcal{L}_{\text{inf}}^>(x, \nu) = \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) > x\}.$$

The $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^>$ quantities might be key in measuring the complexity of best-arm identification under a fixed budget and be involved in generalizations of gap-bounds to general models. As a first illustration, we note that we may actually rewrite lower bound (2.45) in terms of $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^>$ quantities²⁰:

$$\inf_{\underline{\zeta} \in \text{Alt}(\underline{\nu})} \max\{\text{KL}(\zeta_1, \nu_1), \text{KL}(\zeta_2, \nu_2)\} = \inf_{x \in [\mu_{(2)}, \mu^*]} \max\{\mathcal{L}_{\text{inf}}^>(x, \nu_{(2)}), \mathcal{L}_{\text{inf}}^<(x, \nu^*)\}.$$

Contribution 2.22. In Chapter 5, we will prove upper and lower bounds based on the information-theoretic quantities $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^>$. Those bounds will be applicable to a wide variety of models, including non-parametric models, and will imply all existing ones presented above. More precisely:

- for the upper bound, we prove that the quantity

$$\mathcal{L}(\nu', \nu) \stackrel{\text{def}}{=} \inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \left\{ \mathcal{L}_{\text{inf}}^>(x, \nu') + \mathcal{L}_{\text{inf}}^<(x, \nu) \right\}, \quad (2.46)$$

defined for ν, ν' in \mathcal{D} such that $\mathbb{E}(\nu') < \mathbb{E}(\nu)$, can replace the role of gaps in the Successive-Rejects analysis of [Audibert et al. \(2010\)](#). Given a bandit problem $\underline{\nu}$ with a unique optimal arm, we may rank the arms a in non-decreasing order of $\mathcal{L}(\nu_a, \nu^*)$, i.e., consider the permutation σ such that

$$0 = \mathcal{L}(\nu_{\sigma_1}, \nu^*) < \mathcal{L}(\nu_{\sigma_2}, \nu^*) \leq \dots \leq \mathcal{L}(\nu_{\sigma_{K-1}}, \nu^*) \leq \mathcal{L}(\nu_{\sigma_K}, \nu^*),$$

and we prove that for all so-called *regular* models \mathcal{D} , Successive-Rejects guarantees that

$$\forall \underline{\nu} \in \mathcal{D}, \quad \limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\log K} \min_{2 \leq k \leq K} \frac{\mathcal{L}(\nu_{\sigma_k}, \nu^*)}{k}. \quad (2.47)$$

For a sub-Gaussian model, this result improves the upper bound (2.36) as it can be proved using Pinsker's inequality that $\mathcal{L}(\nu', \nu) \geq \frac{1}{4\sigma^2} (\mathbb{E}(\nu) - \mathbb{E}(\nu'))^2$.

- we also derive several lower bounds, which depend on various natural assumptions on the considered strategies. For instance, as long as the sequence of strategies asymptotically pulls the arm associated with the smallest expectation less than a fraction $\frac{1}{K}$ of the time, and

²⁰See the proof of Theorem 5.14 for details about this equality, which is obtained similarly to the transportation costs calculation (2.23) of the fixed-confidence setting.

is more efficient in the identification task when one sub-optimal arm is removed from the bandit problem, then the lower bound (2.40) can be generalized to any model \mathcal{D} as:

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \frac{\mathcal{L}_{\text{inf}}^<(\mu^{(k)}, \nu^*)}{k}.$$

This lower bound is not stated in terms of infima of combinations of $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^<$, i.e., in terms of the function \mathcal{L} introduced in (2.46). While it seems challenging to obtain lower bounds involving \mathcal{L} , we rather obtained the following bound for so-called *monotonous* sequences of strategies, which pull more often arms with higher means:

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \inf_{x \in [\mu^{(k)}, \mu^{(k-1)})} \left\{ \frac{\mathcal{L}_{\text{inf}}^>(x, \nu^{(k)})}{k-1} + \frac{\mathcal{L}_{\text{inf}}^<(x, \nu^*)}{k} \right\}.$$

This bound goes in the direction of the complexity appearing in upper bound (2.47), involving \mathcal{L} quantities, however, the infima are taken over smaller intervals.

The bounds presented in Chapter 5 do not close the gap between lower and upper bounds: there is still (at least) a multiplicative factor $\overline{\log K} \simeq \log K$ between them. The question of reducing this gap is still open and will be discussed in the next section.

2.3.5. On Non-Matching Bounds and Minimax Results

We recall that the lower bound (2.40) and the upper bound (2.36) differ in particular by a factor proportional to $\overline{\log K}$. As a consequence, we need to understand how to fill the gap between those two bounds. Is it even possible?

The minimax lower bound of Carpentier and Locatelli (2016). Carpentier and Locatelli (2016) discuss the gap between lower bound (2.40) and upper bound (2.36) in the case of the Bernoulli model $\mathcal{B}_{[1/4, 3/4]}$. They improve the lower bound (2.40) by a factor of $\log K$, but not simultaneously for all bandit problems $\underline{\nu}$: they obtain the improvement just for one bandit problem $\underline{\nu}$. Their lower bound result (formally stated and discussed in Section 5.6.1) can be asymptotically stated as the existence of $\underline{\nu}$ in $\mathcal{B}_{[1/4, 3/4]}$ and of an increasing sequence of budgets $(T_n)_{n \geq 1}$ such that

$$\exists \underline{\nu} \text{ in } \mathcal{B}_{[1/4, 3/4]}, \quad \liminf_{n \rightarrow +\infty} \frac{1}{T_n} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_{T_n} \neq a^*(\underline{\nu})) \geq - \frac{400}{\log K} H_{\Sigma}(\underline{\nu})^{-1}. \quad (2.48)$$

This result is different in nature from the uniform instance-dependent lower bounds considered in the previous section and in Chapter 5, which hold simultaneously for all bandit problems of a given model. Yet, it gives directions to answer the questions raised above:

- either we need to obtain new lower bounds with the $\log K$ factor so as to match upper bound (2.36),
- or this might not be possible, which would mean that lower and upper bounds cannot perfectly match for large enough models, i.e., that there does not exist an optimal strategy that reaches the best rate (among all strategies) uniformly on all instances.

Recent works have been focussing their efforts on proving that the second scenario occurs. Before discussing those results, let us explain why the fixed-budget setting seems harder than the fixed-confidence setting.

Adaptive strategies and the difficulty of the fixed-budget setting. The upper bound (2.36) was obtained for a Successive-Rejects strategy with pre-defined lengths phase, so, without any attempt to adapt to collected data. It could be useful to rely on those observations to choose the round lengths or design other types of strategies. If adaptivity might complicate the theoretical arguments to obtain upper bounds, one first interesting step would be to design adaptive sampling rules that perform well empirically. This task does not seem to be easy, as we expose now.

- One needs to have a high precision on the estimates of the means to confidently rely on observed data in order to adapt sampling rules. However, the misidentification event $\{\hat{a}_T \neq a^*(\underline{\nu})\}$ that we want to control is the event that (at least) one sub-optimal empirical estimate is above the empirical estimate of the best arm:

$$\left\{ \exists a \neq a^*(\underline{\nu}), \hat{\mu}_a(T) \geq \hat{\mu}_{a^*(\underline{\nu})}(T) \right\}.$$

This holds when there is a failure in the precision of (at least) one of the estimates: either a sub-optimal means is over-estimated, or the optimal mean is under-estimated.

- Using forced exploration (which is equivalent to beginning the process by a round of uniform exploration) does not seem possible in this setting: even a careful choice of the forced exploration strategy will lead to imprecisions in the estimates, and the strategy might not exploit those estimates correctly (see the first item). This is why Audibert et al. (2010) rather proposed another strategy called Upper-Confidence-Bound-Exploration, a UCB-based algorithm with an exploration bonus. To ensure good theoretical guarantees of those strategies, one requires the knowledge of a complexity term, namely $H_\Sigma(\underline{\nu})$, in order to set the scaling of the bonus (see also Gabillon et al., 2012).

Existence of optimal sequences of strategies. Komyama et al. (2022) conjectured that decreasing the error probability under an instance $\underline{\nu}$ should increase the error probability under another instance $\underline{\zeta}$. In other words, that no strategy might perform uniformly well under all bandit instances. Given a fixed complexity function $H(\underline{\nu}) > 0$, the probability of error can be written as

$$\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) = \exp\left(-\frac{T}{R(\underline{\nu}, T)H(\underline{\nu})}\right),$$

for some corrective term $R(\underline{\nu}, T) \in [0, +\infty]$. If a strategy is such that $R(\underline{\nu}, T) \leq 1$, then $H(\underline{\nu})$ is a good upper bound of the complexity of $\underline{\nu}$ as the probability of error is less than

$$\exp\left(-\frac{T}{H(\underline{\nu})}\right).$$

We can extend this statement to sequences of strategies such that

$$R(\underline{\nu}) \stackrel{\text{def}}{=} \limsup_{T \rightarrow +\infty} R(\underline{\nu}, T) \leq 1.$$

In addition, if $R_{\text{inf}}(\underline{\nu}) \stackrel{\text{def}}{=} \inf_{\text{strategies}} R(\underline{\nu}) > 0$, then $H(\underline{\nu})$ is a good lower bound of the complexity as the probability of error is asymptotically higher than

$$\exp\left(-\frac{T}{R_{\text{inf}}(\underline{\nu})H(\underline{\nu})}\right).$$

Can we find a “good” complexity measure H , for which

$$R_{\text{inf}} \stackrel{\text{def}}{=} \inf_{\underline{\nu} \text{ in } \mathcal{D}} R_{\text{inf}}(\underline{\nu}) > 0,$$

and such that there exists a uniformly optimal sequence of strategies, for which

$$R \stackrel{\text{def}}{=} \sup_{\underline{\nu} \text{ in } \mathcal{D}} R_{\text{sup}}(\underline{\nu}) \leq 1?$$

[Komiyama et al. \(2022\)](#) conjectured that such complexity does not exist for large enough models (including an exponential model). Given a complexity H , they proposed an intractable sequence of strategies that can reach the best uniform upper rate R_{inf} among all strategies. For the complexity H_{Σ} defined in (2.41), they also introduced an algorithm approaching this strategy by using a pre-trained neural network. It results in a strategy that performs empirically better than `Successive-Rejects` on some instances (with an improved behavior at the beginning, the uniform phase of exploration being reduced), yet without strong theoretical guarantees.

The recent results of [Degenne \(2023\)](#) partly confirmed the conjecture of [Komiyama et al. \(2022\)](#) on a theoretical aspect: for the model of standard Gaussian distributions, they consider the oracle complexity for which $R_{\text{inf}} = 1$, which is the exact rate of the best strategy pulling arms with fixed proportions, which is equal to

$$H(\underline{\nu}) \stackrel{\text{def}}{=} \left(\sup_{\underline{\nu} \in \Sigma_K} \inf_{\zeta \in \text{Alt}(\underline{\nu})} \sum_{a \in [K]} v_a \text{KL}(\zeta_a, \nu_a) \right)^{-1}. \quad (2.49)$$

They proved that the best rate $\inf_{\text{strategies}} R$ goes to $+\infty$ as $K \rightarrow +\infty$. In other words, for large enough K , we cannot design an algorithm that will be able to track the optimal proportions of $\underline{\nu}$ for all bandit problems $\underline{\nu}$, unlike the fixed-confidence setting. This contradicts in particular the conjecture in the conclusion of [Garivier and Kaufmann \(2016\)](#) according to which the complexity of the fixed-budget best-arm identification might be (2.49).

Remark. [Degenne \(2023\)](#) obtained this result by considering similar arguments than the minimax lower bound (2.48) of [Carpentier and Locatelli \(2016\)](#).

CHAPTER 3

About the Fixed-Confidence Sample Complexity Optimization Problem for Gaussian Variables

In this chapter, we study, for a Gaussian model, the solution of the sample complexity optimization problem (2.20) in best-arm identification with fixed-confidence. We present a new characterization of the solution which allows the design of a new procedure to compute the optimal weight vector and to obtain precise regularity properties concerning its dependency on the bandit instance. Materials are extracted from Section 3 (and associated appendix) of the conference paper:



A. Barrier, A. Garivier, and T. Kocák. A Non-Asymptotic Approach to Best-Arm Identification for Gaussian Bandits. In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, pages 10078–10109. PMLR, 2022

Note that the lemmas of Section 3.4 contain unpublished additional statements.

Contents

1	Introduction	72
2	Solving the Optimization Problem	74
3	Bounds and Computation of the Problem Characteristics	77
4	Monotonicity of the max–min Problem	80
1	Increasing the Mean of a Sub-Optimal Arm	80
2	Increasing the Mean of the Best Arm	83
3	Increasing the Mean of the Worst Arms	85
5	Regularity Properties	86
1	Regularity of \underline{w} and T	87
2	Regularity of g	89
6	Conclusion	91

3.1. Introduction

In this chapter, we consider the “multi-armed bandit” framework, a collection of $K \geq 2$ independent probability distributions $\underline{\nu} = (\nu_1, \dots, \nu_K)$ called *arms*, of unknown means $\underline{\mu} = (\mu_a)_{a \in [K]}$, and belonging to some model \mathcal{D} . These arms are sampled sequentially and independently: at every discrete time step $t \in \mathbb{N}^*$, an agent chooses an arm $A_t \in [K]$ based on past information and observes an independent draw Y_t from distribution ν_{A_t} . Multi-armed bandits are used as models for many situations in which one needs to find the best among a set of options, using noisy observations.

The *best-arm identification* problem consists in identifying the arm with highest mean of a bandit problem $\underline{\nu}$ (we only consider, in the rest of this section, bandit problems with a unique optimal arm):

$$\{a^*(\underline{\nu})\} = \operatorname{argmax}_{a \in [K]} \mu_a.$$

In the *fixed-confidence* setting (see [Even-Dar et al., 2006](#)), a confidence parameter $\delta \in (0, 1)$ is given, and the objective is to design strategies that, after some random number of steps τ_δ return an estimate $\hat{a}_{\tau_\delta} \in [K]$, which is equal to the best arm $a^*(\underline{\nu})$ with probability at least $1 - \delta$. The aim is to find a strategy that minimizes the expected number of samplings $\mathbb{E}_{\underline{\nu}}[\tau_\delta]$ among all δ -correct strategies, which are strategies satisfying

$$\forall \underline{\nu} \text{ in } \mathcal{D}, \quad \mathbb{P}_{\underline{\nu}}(\tau_\delta < +\infty, \hat{a}_{\tau_\delta} \neq a^*(\underline{\nu})) \leq \delta.$$

The sample complexity optimization problem. The sample complexity of δ -correct strategies cannot be arbitrarily good. For an exponential model \mathcal{D}_{exp} , for which instances $\underline{\nu}$ are characterized by their means $\underline{\mu}$, it has been proved by [Garivier and Kaufmann \(2016\)](#) that all strategies satisfy

$$\forall \underline{\mu} \in \mathcal{D}_{\text{exp}}, \quad \mathbb{E}_{\underline{\mu}}[\tau_\delta] \geq T(\underline{\mu}) \log \frac{1}{2.4\delta}, \quad (3.1)$$

where $T(\underline{\mu})$ is the *characteristic time* of $\underline{\mu}$, defined by the following optimization problem as

$$T(\underline{\mu})^{-1} \stackrel{\text{def}}{=} \sup_{\underline{\nu} \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a). \quad (3.2)$$

where d denotes the mean-parameterized Kullback-Leibler divergence of the model \mathcal{D}_{exp} ,

$$\Sigma_K = \left\{ \underline{\nu} \in [0, 1]^K : v_1 + \dots + v_K = 1 \right\} \quad \text{and} \quad \text{Alt}(\underline{\mu}) = \left\{ \lambda \text{ in } \mathcal{D}_{\text{exp}} : a^*(\lambda) \neq a^*(\underline{\mu}) \right\}.$$

An optimal weight vector. [Garivier and Kaufmann \(2016\)](#) proved that optimization problem (3.2) admits a unique maximizer $\underline{w}(\underline{\mu})$, called the *optimal weight vector*:

$$\left\{ \underline{w}(\underline{\mu}) \right\} \stackrel{\text{def}}{=} \operatorname{argmax}_{\underline{\nu} \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a).$$

This maximizer plays a very important role in the design of efficient fixed-confidence strategies: whatever the value of the risk δ , the more a strategy pulls arms in proportions close to $\underline{w}(\underline{\mu})$, the closest from the lower bound (3.1) its performance will be. By giving a procedure to compute the optimal weight vector when the means are known, [Garivier and Kaufmann \(2016\)](#) designed Track-and-Stop, which in a nutshell boils down to tracking the current optimal weight vector $\underline{w}(\hat{\underline{\mu}}(t))$, where $\hat{\underline{\mu}}(t)$ is the empirical estimate of $\underline{\mu}$ at time step t . Track-and-Stop is the first *asymptotically optimal strategy*, for which the asymptotic upper bound matches lower bound (3.1):

$$\forall \underline{\mu} \text{ in } \mathcal{D}_{\text{exp}}, \quad \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_\delta]}{\log \frac{1}{\delta}} \leq T(\underline{\mu}).$$

The existence of such strategies ensures the asymptotic tightness of the lower bound (3.1).

Computing the solution $\underline{w}(\underline{\mu})$. Garivier and Kaufmann (2016) proved that, with the knowledge of $\underline{\mu}$, solving the optimization problem (3.2) reduces to determining the root of a one-variable increasing function. By applying a bisection method, one may then compute $\underline{w}(\underline{\mu})$ with arbitrary precision. The procedure, referred to as `Optimal-Weights`, has consequences on the execution time of strategies that require a computation at each time step. It might be interesting to obtain more efficient methods to compute this optimal weight vector in order to speed up strategies such as `Track-and-Stop`.

Regularity and monotonicity of \underline{w} and T . For a general exponential model, Garivier and Kaufmann (2016) obtained a few properties concerning the mappings $\underline{\mu} \mapsto \underline{w}(\underline{\mu})$ and $\underline{\mu} \mapsto T(\underline{\mu})$. They proved that $\underline{\mu} \mapsto \underline{w}(\underline{\mu})$, and hence $\underline{\mu} \mapsto T(\underline{\mu})$, are continuous (at problems having a unique optimal arm), that $\underline{w}(\underline{\mu})$ charges all arms:

$$\min_{a \in [K]} w_a(\underline{\mu}) > 0,$$

and that, if arms are ordered so that $\mu_{(1)} > \mu_{(2)} \geq \dots \geq \mu_{(K)}$, then

$$w_{(2)}(\underline{\mu}) \geq \dots \geq w_{(K)}(\underline{\mu}),$$

but there might be instances (e.g., for a Bernoulli model) for which $w_{(1)}(\underline{\mu}) < w_{(2)}(\underline{\mu})$. Apart from those basic properties, only a little is known about quantitative regularity results, or monotonicity properties of $\underline{w}(\underline{\mu})$ when, for instance, moving the mean of one of the arms. Obtaining such results would have implications in the design of more efficient and natural (in several meanings) strategies, see, e.g., Chapter 4.

The special case of Gaussian variables. Among all exponential models, the model $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ of Gaussian variables with common variance $\sigma^2 > 0$ provides a suitable setting for analysis. Its Kullback-Leibler divergence enjoys the simple closed-form expression

$$\forall \mu, \mu' \in \mathbb{R}, \quad d(\mu, \mu') = \frac{(\mu' - \mu)^2}{2\sigma^2}, \quad (3.3)$$

which allows to rewrite the optimization problem (3.2) as

$$T(\underline{\mu})^{-1} = \sup_{\underline{v} \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a \frac{(\mu_a - \lambda_a)^2}{2\sigma^2}. \quad (3.4)$$

As this max–min optimization problem has a handy objective function, it might be possible to obtain new procedures for the computation of its solution $\underline{w}(\underline{\mu})$, together with monotonicity and regularity results of \underline{w} and T .

Outline and contributions. This chapter is devoted to the study of the Gaussian optimization problem (3.4). For the clarity of the presentation, we only consider the model $\mathcal{D}_{\mathcal{N}_1}$ of standard Gaussian variables, but all results naturally extend to $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$.

In Section 3.2, we describe a new procedure for solving optimization problem (3.4). This results in an accelerated algorithm for its numerical resolution that we present in Section 3.3, allowing a significant speed-up for the `Track-and-Stop` algorithm in the Gaussian case. In addition, we use our procedure to deduce monotonicity properties in Section 3.4 and develop a careful analysis in Section 3.5 of the quantitative regularity of the solution to the optimization problem (3.4). Those results will be essential for defining the `Exploration-Biased-Sampling` strategy of Chapter 4 and proving its non-asymptotic guarantees.



Notation. Let, in the rest of this chapter, $\underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}$ be a fixed bandit parameter with a unique optimal arm. For the simplicity of the presentation, we set

$$a^* = a^*(\underline{\mu}), \quad \underline{w} = \underline{w}(\underline{\mu}), \quad \text{and} \quad T = T(\underline{\mu}).$$

We recall that, for $\underline{v} \in \Sigma_K$, the quantity

$$g(\underline{\mu}, \underline{v}) \stackrel{\text{def}}{=} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a \frac{(\mu_a - \lambda_a)^2}{2}$$

can be seen as a minimum of transportation costs of the sub-optimal arms, see Equation (2.24) and more generally Section 2.2.4, which in the Gaussian case reads

$$g(\underline{\mu}, \underline{v}) = \min_{a \neq a^*} \underbrace{\frac{v_{a^*} v_a}{v_{a^*} + v_a}}_{=\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})} \frac{\Delta_a^2}{2}, \quad (3.5)$$

by use of (3.3), where $\Delta_a = \Delta_a(\underline{\mu}) \stackrel{\text{def}}{=} \mu_{a^*} - \mu_a$ is the *gap* of arm a . Thus, the optimization problem (3.4) can be equivalently written as

$$T(\underline{\mu})^{-1} = \sup_{\underline{v} \in \Sigma_K} g(\underline{\mu}, \underline{v}) = \frac{1}{2} \sup_{\underline{v} \in \Sigma_K} \min_{a \neq a^*} \frac{v_{a^*} v_a}{v_{a^*} + v_a} \Delta_a^2, \quad (3.6)$$

and the optimal weight vector \underline{w} satisfies:

$$\{\underline{w}\} = \operatorname{argmax}_{\underline{v} \in \Sigma_K} g(\underline{\mu}, \underline{v}) = \operatorname{argmax}_{\underline{v} \in \Sigma_K} \frac{1}{2} \min_{a \neq a^*} \frac{v_{a^*} v_a}{v_{a^*} + v_a} \Delta_a^2. \quad (3.7)$$

[Garivier and Kaufmann \(2016\)](#) proved that, at the optimum \underline{w} , all transportation costs are equal (see also the proof of Proposition 3.2), so that, for all arms $a \neq a^*$,

$$T(\underline{\mu})^{-1} = g(\underline{\mu}, \underline{w}) = \frac{1}{2} \frac{w_{a^*} w_a}{w_{a^*} + w_a} \Delta_a^2. \quad (3.8)$$

Remark. We might sometimes use that $g(\underline{\mu}, \underline{v})$ is defined for bandit problems $\underline{\mu}$ that admit several optimal arms (with a^* begin any of the optimal arms). In that case, $g(\underline{\mu}, \underline{v}) = 0$, all vectors \underline{v} are solutions of optimization problem (3.6) and $T = +\infty$. However, unless explicitly stated, we only consider problems with a unique optimal arm.

Finally, we define the following quantities:

$$\Delta_{\min} \stackrel{\text{def}}{=} \min_{a \neq a^*} \Delta_a, \quad \Delta_{\max} \stackrel{\text{def}}{=} \max_{a \in [K]} \Delta_a, \quad \text{and} \quad w_{\min} \stackrel{\text{def}}{=} \min_{a \in [K]} w_a.$$

3.2. Solving the Optimization Problem

In this section, we give a new characterization of the solution \underline{w} of optimization problem (3.6). We show that \underline{w} can be derived using the root of the function $\phi_{\underline{\mu}}$ defined by:

$$\forall r \in \left(\frac{1}{\Delta_{\min}^2}, +\infty \right), \quad \phi_{\underline{\mu}}(r) \stackrel{\text{def}}{=} \sum_{a \neq a^*} \frac{1}{(r \Delta_a^2 - 1)^2} - 1.$$

The following properties of $\phi_{\underline{\mu}}$ are straightforward.

Lemma 3.1. $\phi_{\underline{\mu}}$ is strictly convex, decreasing on $(\frac{1}{\Delta_{\min}^2}, +\infty)$, and thus has a unique root.

To compute the solution \underline{w} of optimization problem (3.6), it is sufficient to solve $\phi_{\underline{\mu}}(r) = 0$. The next result explains how to link r and \underline{w} .

Proposition 3.2. Let $r = r(\underline{\mu})$ be the solution of $\phi_{\underline{\mu}}(r) = 0$. Then

$$w_{a^*} = \frac{1}{1 + \sum_{a \neq a^*} \frac{1}{r\Delta_a^2 - 1}}, \quad (3.9)$$

$$\forall a \neq a^*, \quad w_a = \frac{w_{a^*}}{r\Delta_a^2 - 1}, \quad (3.10)$$

$$\text{and} \quad T = 2 \frac{r}{w_{a^*}}. \quad (3.11)$$

Besides,

$$w_{a^*}^2 = \sum_{a \neq a^*} w_a^2. \quad (3.12)$$

Before proving this result, we make two observations. Other important consequences of Proposition 3.2 will be derived in the next sections.

Monotonicity of the weights. In the case of $K = 2$ arms, we recall¹ that $\underline{w}(\underline{\mu}) = (0.5, 0.5)$. For $K \geq 3$, we obtain that the optimal weights are monotonous with respect to the means.

Corollary 3.3. Assume that $K \geq 3$. Then

$$\forall a, b \in [K], \quad \mu_a > \mu_b \implies w_a > w_b.$$

Proof. When a and b are sub-optimal, the result is a direct consequence of Equation (3.10). It remains to see that $w_{a^*} > \max_{a \neq a^*} w_a$ using Equation (3.12) with the fact that all weights are positive and that $K \geq 3$. \square

Optimal ratio of sub-optimal arms. Equation (3.10) also implies that

$$\forall a, b \neq a^*, \quad \frac{w_a}{w_b} = \frac{\Delta_b^2 - \frac{1}{r}}{\Delta_a^2 - \frac{1}{r}}.$$

Intuitively, it requires about $\frac{1}{\Delta_a^2}$ samplings of arms a^* and a before being able to distinguish them, so that one could expect $\frac{w_a}{w_b}$ to be $\frac{\Delta_b^2}{\Delta_a^2}$. This would be the case if the comparisons between arms were independent. In our problem, sampling the best arm benefits the comparison with all arms, so it is worth sampling the optimal arm a little more than any single comparison would require, and hence each sub-optimal arm a little less. As a result, the ratio $\frac{w_a}{w_b}$ is closer to 1, and the factor can be seen as a “discount” on each squared gap for sharing the comparisons.

We close this section by proving Proposition 3.2.

¹This can be obtained by Proposition 3.2, or by directly solving optimization problem (3.6).

Proof of Proposition 3.2. Let us define, for some $\beta \in [0, 1]$:

$$G(\beta) \stackrel{\text{def}}{=} \max_{\substack{\underline{v} \in \Sigma_K \\ v_{a^*} = \beta}} \min_{a \neq a^*} \frac{\beta v_a}{\beta + v_a} \Delta_a^2, \quad (3.13)$$

so that, by Equation (3.8),

$$T^{-1} = \max_{\underline{v} \in \Sigma_K} g(\underline{\mu}, \underline{v}) = \frac{1}{2} \max_{\beta \in [0,1]} G(\beta). \quad (3.14)$$

Remark. $G(\beta) = 2(T_\beta)^{-1}$ where T_β is the characteristic time of strategies allowing a proportion β of their budget to the best arm, like non-adaptive top-two algorithms, see Equation (2.30).

By unicity of the optimal weight vector \underline{w} , see (3.7), we know that:

$$\{w_{a^*}\} = \operatorname{argmax}_{\beta \in (0,1)} G(\beta). \quad (3.15)$$

Fix $\beta \in [0, 1]$. The maximum in Equation (3.13) is reached for a vector \underline{v} such that $v_{a^*} = \beta$ and all the costs $(\frac{\beta v_a}{\beta + v_a} \Delta_a^2)_{a \neq a^*}$ are equal: $G(\beta)$ is such that

$$\forall a \neq a^*, \quad G(\beta) = \frac{\beta v_a}{\beta + v_a} \Delta_a^2,$$

and hence

$$\forall a \neq a^*, \quad v_a = \frac{\beta G(\beta)}{\beta \Delta_a^2 - G(\beta)}. \quad (3.16)$$

The fact that $\underline{v} \in \Sigma_K$ yields:

$$\Phi(\beta, G(\beta)) \stackrel{\text{def}}{=} \beta + \sum_{a \neq a^*} \frac{\beta G(\beta)}{\beta \Delta_a^2 - G(\beta)} - 1 = 0. \quad (3.17)$$

By the implicit function theorem, there exists a mapping $\beta \in [0, 1] \mapsto G(\beta)$ such that

$$\begin{aligned} \Phi(\beta, G(\beta)) &= 0, \\ \text{and} \quad G'(\beta) &= - \frac{\frac{\partial \Phi}{\partial \beta}(\beta, G(\beta))}{\frac{\partial \Phi}{\partial G}(\beta, G(\beta))} \\ &= - \frac{1 + \sum_{a \neq a^*} \frac{G(\beta)(\beta \Delta_a^2 - G(\beta)) - \beta G(\beta) \Delta_a^2}{(\beta \Delta_a^2 - G(\beta))^2}}{\beta^2 \sum_{a \neq a^*} \frac{\Delta_a^2}{(\beta \Delta_a^2 - G(\beta))^2}} \\ &= - \frac{1 - \sum_{a \neq a^*} \frac{1}{(\frac{\beta}{G(\beta)} \Delta_a^2 - 1)^2}}{\beta^2 \sum_{a \neq a^*} \frac{\Delta_a^2}{(\beta \Delta_a^2 - G(\beta))^2}}. \end{aligned}$$

Hence $\beta \mapsto G(\beta)$ is a smooth non-negative function with a continuous derivative. By Equation (3.13), it vanishes when $\beta \rightarrow 0$ and $\beta \rightarrow 1$, and hence its maximum is reached at a point β^* where $G'(\beta^*) = 0$. Defining

$$r \stackrel{\text{def}}{=} \frac{\beta^*}{G(\beta^*)},$$

we get that r is the unique solution of $\phi_{\underline{\mu}}(r) = 0$ using the relation

$$G'(\beta^*) = 0 \iff 1 - \sum_{a \neq a^*} \frac{1}{\left(\frac{\beta^*}{G(\beta^*)} \Delta_a^2 - 1\right)^2} = 0.$$

The claimed results follow by remarking that, by Equation (3.15), β^* is (unique and) equal to w_{a^*} . Equations (3.9), (3.10) and (3.11) can be respectively derived from (3.17), (3.16) and (3.14). Lastly, we obtain Equation (3.12) by combining Equation (3.10) and the characterization $\phi_{\underline{\mu}}(r) = 0$:

$$\sum_{a \neq a^*} w_a^2 = w_{a^*}^2 \sum_{a \neq a^*} \frac{1}{(r \Delta_a^2 - 1)^2} = w_{a^*}^2 (\phi_{\underline{\mu}}(r) + 1) = w_{a^*}^2.$$

□

Remark. We largely used the explicit (and nice) expression 3.3 of the Kullback-Leibler divergence between standard Gaussian variables. Obtaining a similar proof in the general case of an exponential model might not allow us to obtain simple relations as in Proposition 3.2.

3.3. Bounds and Computation of the Problem Characteristics

By Proposition 3.2, it suffices to compute r to obtain the values of both T and \underline{w} . As $\phi_{\underline{\mu}}$ is a strictly convex and strictly decreasing function, Newton's iterates initialized with a value $r_0 < r$ converge to r from below at quadratic speed (the number of correct digits roughly doubles at every step). This implies that a few iterations are sufficient to guarantee machine precision. The cost of the algorithm can hence be considered proportional to that of evaluating $\phi_{\underline{\mu}}(r)$, which is linear in the number of arms. The procedure called **Gaussian-Optimal-Weights** is summarized in Algorithm 11, where we use the close form expression of the derivative of $\phi_{\underline{\mu}}$:

$$\phi'_{\underline{\mu}} : r \in \left(\frac{1}{\Delta_{\min}^2}, +\infty \right) \mapsto -2 \sum_{a \neq a^*} \frac{\Delta_a^2}{(r \Delta_a^2 - 1)^3}.$$

It remains to show that it is possible to find $r_0 < r$, and possibly close to r . The next proposition offers such a lower bound as simple functions of the gaps. This also yields tight bounds on the optimal weight vector \underline{w} and the characteristic time T .

Proposition 3.4. *Let r be the solution of $\phi_{\underline{\mu}}(r) = 0$. Then the following holds:*

$$\max\left(\frac{2}{\Delta_{\min}^2}, \frac{1 + \sqrt{K-1}}{\bar{\Delta}^2}\right) \leq r \leq \frac{1 + \sqrt{K-1}}{\Delta_{\min}^2}, \quad (3.18)$$

$$\frac{1}{1 + \sqrt{K-1}} \leq w_{a^*} \leq \frac{1}{2}, \quad (3.19)$$

$$\text{and} \quad \max\left(\frac{8}{\Delta_{\min}^2}, 4 \frac{1 + \sqrt{K-1}}{\bar{\Delta}^2}\right) \leq T \leq 2 \frac{(1 + \sqrt{K-1})^2}{\Delta_{\min}^2}, \quad (3.20)$$

where $\bar{\Delta}^2 \stackrel{\text{def}}{=} \frac{1}{K-1} \sum_{a \neq a^*} \Delta_a^2$ is the average squared gap.

Note that all of these inequalities can be reached for certain parameters $\underline{\mu}$, as discussed after the proof of the proposition.

Algorithm 11: Gaussian-Optimal-Weights

Input: bandit μ in $\mathcal{D}_{\mathcal{N}_1}$
 initialization r_0
 tolerance parameter tol (typically 10^{-10})
Output: optimal weight vector \underline{w}
 characteristic time T

```

1  $r \leftarrow r_0$ 
2 while  $|\phi_{\underline{\mu}}(r)| \geq \text{tol}$  do
3    $r \leftarrow r - \frac{\phi_{\underline{\mu}}(r)}{\phi'_{\underline{\mu}}(r)}$ 
4  $w_{a^*} \leftarrow \left(1 + \sum_{a \neq a^*} \frac{1}{r\Delta_a^2 - 1}\right)^{-1}$ 
5 for  $a \neq a^*$  do
6    $w_a \leftarrow \frac{w_{a^*}}{r\Delta_a^2 - 1}$ 
7  $T \leftarrow 2 \frac{r}{w_{a^*}}$ 
    
```

Proof. For all sub-optimal arms a , we set

$$q_a \stackrel{\text{def}}{=} \frac{1}{r\Delta_a^2 - 1}.$$

By definition of r , we get

$$0 = \phi_{\underline{\mu}}(r) = \sum_{a \neq a^*} q_a^2 - 1. \quad (3.21)$$

Hence the $(q_a^2)_{a \neq a^*}$ are positive and sum to 1. This implies, in particular, that $q_a \leq 1$ for all $a \neq a^*$, with strict inequality when $K \geq 3$.

We first prove inequalities (3.19). We recall that Equation (3.9) reads, in terms of the $(q_a)_{a \neq a^*}$:

$$w_{a^*} = \left(1 + \sum_{a \neq a^*} q_a\right)^{-1}.$$

On the one hand, the upper bound of (3.19) follows from the fact that the $(q_a)_{a \neq a^*}$ are less than or equal to 1 and relation (3.21):

$$w_{a^*} \leq \left(1 + \sum_{a \neq a^*} q_a^2\right)^{-1} = (1 + 1)^{-1} = \frac{1}{2},$$

and on the other hand, we derive the lower bound of (3.19) by the Cauchy-Schwarz inequality:

$$w_{a^*} \geq \left(1 + \sqrt{(K-1) \sum_{a \neq a^*} q_a^2}\right)^{-1} = \frac{1}{1 + \sqrt{K-1}}.$$

We now prove inequalities (3.18). For the lower bound, note first that, since $q_a \leq 1$ or equivalently $r\Delta_a^2 \geq 2$ for every $a \neq a^*$, we get

$$r \geq \frac{2}{\Delta_{\min}^2}. \quad (3.22)$$

Then, as $\overline{\Delta^2} = \frac{1}{K-1} \sum_{a \neq a^*} \Delta_a^2$, the convexity of $x \mapsto \frac{1}{(rx-1)^2}$ ensures that

$$\frac{1}{K-1} \sum_{a \neq a^*} \frac{1}{\left(\frac{1+\sqrt{K-1}}{\Delta^2} \Delta_a^2 - 1\right)^2} \geq \frac{1}{\left(\frac{1+\sqrt{K-1}}{\Delta^2} \overline{\Delta^2} - 1\right)^2} = \frac{1}{K-1},$$

and hence

$$\phi_{\underline{\mu}}\left(\frac{1+\sqrt{K-1}}{\Delta^2}\right) = \sum_{a \neq a^*} \frac{1}{\left(\frac{1+\sqrt{K-1}}{\Delta^2} \Delta_a^2 - 1\right)^2} + 1 \geq 0.$$

By decreasing of $\phi_{\underline{\mu}}$ (see Lemma 3.1), this entails that

$$r \geq \frac{1 + \sqrt{K-1}}{\Delta^2},$$

which, together with (3.22), concludes the proof of the lower bound of (3.18). The decreasing of $\phi_{\underline{\mu}}$ is also helpful for the upper bound of (3.18): as

$$\phi_{\underline{\mu}}\left(\frac{1+\sqrt{K-1}}{\Delta_{\min}^2}\right) = \sum_{a \neq 1} \frac{1}{\left(\frac{1+\sqrt{K-1}}{\Delta_{\min}^2} \Delta_a^2 - 1\right)^2} - 1 \leq 0,$$

we get

$$r \leq \frac{1 + \sqrt{K-1}}{\Delta_{\min}^2}.$$

Finally, Equation (3.20) is derived by combining inequalities (3.18) and (3.19) with Equation (3.11). \square

Tightness of the bounds. To conclude this section, we discuss the tightness of the inequalities obtained in Proposition 3.4.

- When $K = 2$, we note that lower and upper bounds match in inequalities (3.18), (3.19) and (3.20). However, the bounds do not provide any additional information as $\underline{w} = (0.5, 0.5)$ whatever the value of $\underline{\mu}$.
- In fact, equalities $r = \frac{2}{\Delta_{\min}^2}$, $w_{a^*} = \frac{1}{2}$ and $T = \frac{8}{\Delta_{\min}^2}$ occur if and only if $K = 2$. This is because the $(q_a)_{a \neq a^*}$ are positive and sum to 1, hence they cannot equal 1 unless $K = 2$: as we proved the associated inequalities by injecting that $q_a \leq 1$ for all $a \neq a^*$, there is equality only when $K = 2$. The presence of additional arms increases r and T while decreases w_{a^*} .
- If there is at least $K \geq 3$ arms, then the remaining equalities

$$w_{a^*} = \frac{1}{1 + \sqrt{K-1}}, \quad r = \frac{1 + \sqrt{K-1}}{\Delta^2} = \frac{1 + \sqrt{K-1}}{\Delta_{\min}^2}, \quad \text{and} \quad T = 2 \frac{(1 + \sqrt{K-1})^2}{\Delta_{\min}^2},$$

are reached if and only if $\Delta_{\min} = \Delta_{\max}$, that is, if $\Delta_2 = \dots = \Delta_K$. This might be obtained by studying the equality cases in the proof above, using the equality case of the Cauchy-Schwarz inequality for w_{a^*} , the strict convexity of $x \mapsto \frac{1}{(rx-1)^2}$ and the decreasing of $\phi_{\underline{\mu}}$ for r , and finally Equation (3.20) for T . If the condition holds, note that T grows linearly with K .

3.4. Monotonicity of the max–min Problem

We now show monotonicity results for the characteristic time T and the optimal weight vector \underline{w} when moving some arm(s) of $\underline{\mu}$. When $K = 2$, we have seen that $\underline{w} = (0.5, 0.5)$ whatever the bandit problem $\underline{\mu}$, and the variations of T are also known since $T = \frac{8}{\Delta^2}$ where Δ is the unique gap of $\underline{\mu}$.



In the rest of the section, we assume that $K \geq 3$. Let $\underline{\mu}'$ denote another bandit problem in $\mathcal{D}_{\mathcal{N}_1}$ sharing the same unique optimal arm a^* as $\underline{\mu}$, and define $\underline{\Delta}'$ its gap vector, \underline{w}' its optimal weight vector, T' its characteristic time and r' the root of $\phi_{\underline{\mu}'}$. Those results will be useful to define the Exploration–Biased–Sampling strategy in Chapter 4.

3.4.1. Increasing the Mean of a Sub-Optimal Arm

The following lemma gives monotonicity properties of \underline{w} and T when increasing the mean of a sub-optimal arm, which corresponds to decreasing its gap.

Lemma 3.5. *Assume that $\Delta'_b < \Delta_b$ for a fixed $b \neq a^*$ while $\Delta'_a = \Delta_a$ for all $a \notin \{a^*, b\}$. Then*

1. $w'_b > w_b$,
2. $w'_a < w_a$ for all $a \notin \{a^*, b\}$,
3. $T' > T$,
4. in addition:
 - $w'_{a^*} > w_{a^*}$ if arm b is (one of) the second best arm(s) of μ and μ' ,
 - $w'_{a^*} < w_{a^*}$ if arm b is (one of) the worst arm(s) of μ and μ' .

The monotonicity results of the lemma are summarized in Figure 3.1. As an example, we give in Figure 3.2 the evolution of the optimal weights when modifying the value of one sub-optimal arm. This illustrates, in particular, the point 4 of the lemma: $w_{a^*}(\underline{\mu}^x)$ decreases (respectively increases) when the moved arm is the worst arm (respectively the second best arm), which corresponds to $x \leq 2$ (respectively $x > 3$) in the figure. Otherwise, if $x \in [2, 3]$, the variation of $w_{a^*}(\underline{\mu}^x)$ changes at some point that is marked by a dashed line on the figure, and we cannot conclude a monotonicity property in that case.

Proof.

1. Since $\Delta'_a \leq \Delta_a$ for all $a \neq a^*$, we get

$$\phi_{\underline{\mu}'}(r) = \sum_{a \neq a^*} \frac{1}{(r\Delta'_a - 1)^2} + 1 > \sum_{a \neq a^*} \frac{1}{(r\Delta_a - 1)^2} + 1 = \phi_{\underline{\mu}}(r) = 0,$$

hence it holds that $r' > r$ by decreasing of $\phi_{\underline{\mu}'}$ (see Lemma 3.1). This implies that

$$\forall a \notin \{a^*, b\}, \quad \frac{1}{r'\Delta'_a - 1} = \frac{1}{r\Delta_a - 1} < \frac{1}{r\Delta'_a - 1}. \quad (3.23)$$

As $K \geq 3$, such an arm a exists and hence, as $\phi_{\underline{\mu}}(r) = 0 = \phi_{\underline{\mu}'}(r')$:

$$\frac{1}{r'\Delta'_b - 1} > \frac{1}{r\Delta_b - 1}, \quad \text{i.e.,} \quad r'\Delta'_b - 1 < r\Delta_b - 1. \quad (3.24)$$

Combining (3.23) and (3.24) with Equation (3.10) entails that:

$$\begin{aligned} \forall a \notin \{a^*, b\}, \quad \frac{w'_a}{w'_b} &= \frac{r'\Delta'_b - 1}{r'\Delta'_a - 1} < \frac{r\Delta_b - 1}{r\Delta_a - 1} = \frac{w_a}{w_b}, \\ \text{and} \quad \frac{w'_{a^*}}{w'_b} &= r'\Delta'_b - 1 < r\Delta_b - 1 = \frac{w_{a^*}}{w_b}. \end{aligned}$$

3.4. MONOTONICITY OF THE max-min PROBLEM

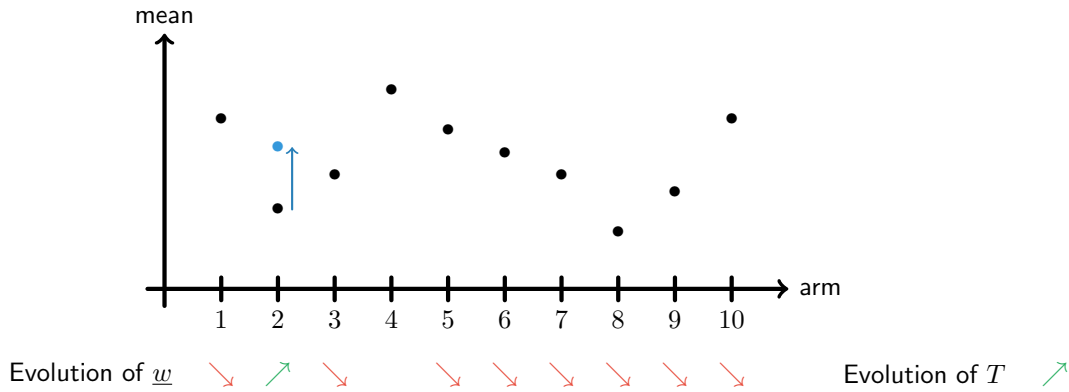


Figure 3.1: By **increasing the mean of a sub-optimal arm**, its associated optimal weight increases while those of other sub-optimal arms decrease, and the characteristic time increases.

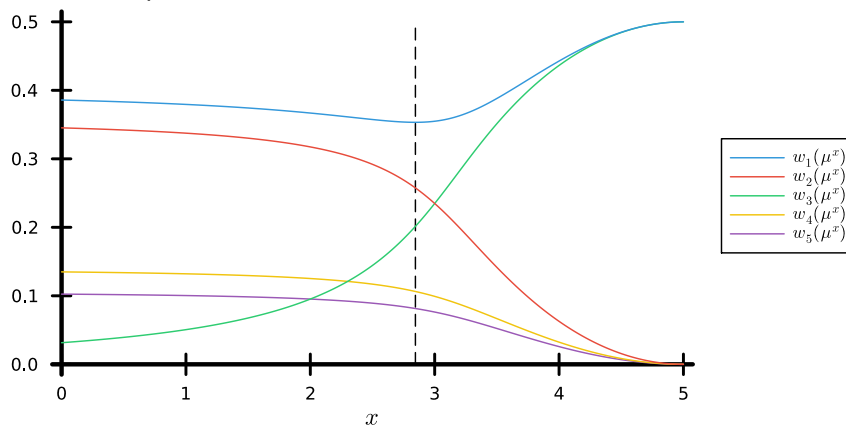


Figure 3.2: Evolution of the optimal weight $w(\underline{\mu}^x)$ for $x \in [0, 5]$, where $\underline{\mu}^x = (5, 4, x, 2, 3, 2)$. The dashed line corresponds to the minimizer of $w_1(\underline{\mu}^x)$.

This leads to

$$\frac{1 - w'_b}{w'_b} = \sum_{a \neq b} \frac{w'_a}{w'_b} < \sum_{a \neq b} \frac{w_a}{w_b} = \frac{1 - w_b}{w_b},$$

and thus, finally, to $w'_b > w_b$.

- Given a bandit problem $\underline{\rho}$ with unique optimal arm a^* , its optimal weights and the root of $\phi_{\underline{\rho}}$ can be seen as functions of its squared gap vector $\underline{d} \stackrel{\text{def}}{=} \underline{\Delta}(\underline{\rho})^2$: we set

$$\forall a \notin \{a^*, b\}, \quad F_a(\underline{d}) \stackrel{\text{def}}{=} \frac{1}{w_a(\underline{\rho})} = \frac{r(\underline{d}) d_a - 1}{w_{a^*}(\underline{\rho})} = (r(\underline{d}) d_a - 1) + \sum_{c \neq a^*} \frac{r(\underline{d}) d_a - 1}{r(\underline{d}) d_c - 1}.$$

where $r(\underline{d})$ is the unique solution of $\phi_{\underline{\rho}}(r) = 0$, and where the right-equalities are derived from Equations (3.9) and (3.10).

We are interested in the variations of F_a with respect to d_b (at points \underline{d} corresponding to bandit problems with optimal arm a^* , i.e., such that $d_{a^*} = 0$ and $d_a > 0$ for $a \neq a^*$). Note that, by an application of the implicit function theorem to the relation $\phi_{\underline{\rho}}(r(\underline{d})) = 0$, we know that r is differentiable with respect to d_b , with

$$\frac{\partial r}{\partial d_b}(\underline{d}) = - \frac{\frac{r(\underline{d})}{(r(\underline{d}) d_b - 1)^3}}{\sum_{c \neq a^*} \frac{d_c}{(r(\underline{d}) d_c - 1)^3}} < 0. \quad (3.25)$$

Remark. Thanks to the first part of the proof, we already knew that r is decreasing with respect to d_b (or equivalently to the gap Δ_b).

As a consequence, the functions $(F_a)_{a \notin \{a^*, b\}}$ are differentiable with respect to d_b , and the proof of the second point of the lemma will follow by proving that

$$\forall a \notin \{a^*, b\}, \quad \frac{\partial F_a}{\partial d_b}(\underline{d}) < 0. \quad (3.26)$$

Indeed, this will lead, as $\Delta'_b < \Delta_b$ and $\Delta'_a = \Delta_a$ for all $a \notin \{a^*, b\}$, to

$$\forall a \notin \{a^*, b\}, \quad w'_a = \frac{1}{F_a(\underline{\Delta}'^2)} < \frac{1}{F_a(\underline{\Delta}^2)} = w_a.$$

We now prove (3.26). In the following calculations, we omit the dependency of the quantities on \underline{d} to facilitate the reading, even if the dependency with respect to the squared gaps \underline{d} is crucial for, e.g., $r = r(\underline{d})$. We get, for all $a \notin \{a^*, b\}$:

$$\begin{aligned} \frac{\partial F_a}{\partial d_b} &= \frac{\partial r}{\partial d_b} d_a + \sum_{c \neq a^*} \left[\frac{\partial r}{\partial d_b} d_c - \frac{rd_c - 1}{(rd_c - 1)^2} \left(\frac{\partial r}{\partial d_b} d_c \right) \right] - \frac{rd_a - 1}{(rd_b - 1)^2} r \\ &= \frac{\partial r}{\partial d_b} d_a \left(1 + \sum_{c \neq a^*} \frac{1}{rd_c - 1} - \frac{rd_c}{(rd_c - 1)^2} \right) + \frac{\partial r}{\partial d_b} \sum_{c \neq a^*} \frac{d_c}{(rd_c - 1)^2} - \frac{rd_a - 1}{(rd_b - 1)^2} r \\ &= \frac{\partial r}{\partial d_b} d_a \sum_{c \neq a^*} \underbrace{\frac{1 + (rd_c - 1) - rd_c}{(rd_c - 1)^2}}_{=0} + \frac{\partial r}{\partial d_b} \sum_{c \neq a^*} \frac{d_c}{(rd_c - 1)^2} - \frac{rd_a - 1}{(rd_b - 1)^2} r \\ &= \frac{\partial r}{\partial d_b} \sum_{c \neq a^*} \frac{d_c}{(rd_c - 1)^2} - \frac{rd_a - 1}{(rd_b - 1)^2} r, \end{aligned}$$

where, in the third equality, we used that, by definition of r ,

$$1 = \sum_{c \neq a^*} \frac{1}{(rd_c - 1)^2}.$$

As $\frac{\partial r}{\partial d_b}$ is negative by (3.25), this implies that $\frac{\partial F_a}{\partial d_b} < 0$, i.e., (3.26) holds.

3. Using Equations (3.8) and (3.5), we get:

$$T'^{-1} = \frac{1}{2} \min_{a \neq a^*} \frac{w'_{a^*} w'_a}{w'_{a^*} + w'_a} \Delta_a'^2 \leq \frac{1}{2} \min_{a \neq a^*} \frac{w'_{a^*} w'_a}{w'_{a^*} + w'_a} \Delta_a^2 < \frac{1}{2} \min_{a \neq a^*} \frac{w_{a^*} w_a}{w_{a^*} + w_a} \Delta_a^2 = T^{-1},$$

where the first inequality uses that $\Delta'_a \leq \Delta_b$ for all $a \neq a^*$, and the second inequality is a consequence of the unicity of the optimal weight vector \underline{w} and the fact that (as obtained in the previous points) $\underline{w} \neq \underline{w}'$.

4. As in the second point of the proof, we can look at the variations of the function

$$F_{a^*}(\underline{d}) \stackrel{\text{def}}{=} \frac{1}{w_{a^*}(\underline{\rho})} = 1 + \sum_{c \neq a^*} \frac{1}{r(\underline{d}) d_c - 1},$$

where the second equality is derived from Equation (3.9).

3.4. MONOTONICITY OF THE max-min PROBLEM

- We want to prove that $w'_{a^*} > w_{a^*}$ when b is the second best arm of $\underline{\mu}$ (hence of $\underline{\mu}'$), which boils down to proving that the partial derivative $\frac{\partial F_{a^*}}{\partial d_b}$ is positive when $d_b < \min_{a \notin \{a^*, b\}} d_a$. Indeed, this will imply, as $\Delta'_b < \Delta_b$ and $\Delta'_a = \Delta_a$ for all $a \notin \{a^*, b\}$, that

$$w'_{a^*} = \frac{1}{F_{a^*}(\underline{\Delta}'^2)} > \frac{1}{F_{a^*}(\underline{\Delta}^2)} = w_{a^*}.$$

Computing the derivative of F_{a^*} with respect to d_b leads, by using (3.25), to

$$\begin{aligned} \frac{\partial F_{a^*}}{\partial d_b} &= - \sum_{c \neq a^*} \frac{\frac{\partial r}{\partial d_b} d_c}{(rd_c - 1)^2} - \frac{r}{(rd_c - 1)^2} \\ &= \underbrace{\frac{r}{(rd_b - 1)^2}}_{>0} \left[\frac{1}{rd_b - 1} \cdot \frac{\sum_{c \neq a^*} \frac{d_c}{(rd_c - 1)^2}}{\sum_{c \neq a^*} \frac{d_c}{(rd_c - 1)^3}} - 1 \right], \end{aligned}$$

which is positive when $d_b < \min_{a \notin \{a^*, b\}} d_a$, as

$$(rd_b - 1) \cdot \sum_{c \neq a^*} \frac{d_c}{(rd_c - 1)^3} = \sum_{c \neq a^*} \frac{(rd_b - 1)d_c}{(rd_c - 1)^3} < \sum_{c \neq a^*} \frac{d_c}{(rd_c - 1)^2}.$$

- To prove that $w'_{a^*} < w_{a^*}$ when b is the worst best arm of $\underline{\mu}'$ (hence of $\underline{\mu}$), we prove that the derivative $\frac{\partial F_{a^*}}{\partial d_b}$ is negative as soon as $d_b > \max_{a \notin \{a^*, b\}} d_a$. By the previous calculations, this indeed holds as

$$(rd_b - 1) \cdot \sum_{c \neq a^*} \frac{d_c}{(rd_c - 1)^3} = \sum_{c \neq a^*} \frac{(rd_b - 1)d_c}{(rd_c - 1)^3} > \sum_{c \neq a^*} \frac{d_c}{(rd_c - 1)^2}.$$

□

3.4.2. Increasing the Mean of the Best Arm

This second lemma gives the monotonicity of w and T when increasing the mean of the best arm, which corresponds to increasing all gaps by a constant value.

Lemma 3.6. *Assume that $\Delta'_a = \Delta_a + d$ for every $a \neq a^*$ and some $d > 0$. Then*

1. *if $b \in \operatorname{argmin}_{a \neq a^*} \Delta_a$ is (one of) the second best arm(s) of μ , then $w'_b \leq w_b$,*
2. *$w'_{\min} \geq w_{\min}$,*
3. *$T' < T$,*
4. *we get*

$$\lim_{d \rightarrow +\infty} w'_{a^*} = \frac{1}{1 + \sqrt{K-1}}, \quad \text{and} \quad \forall a \neq a^*, \quad \lim_{d \rightarrow +\infty} w'_a = \frac{1}{K-1 + \sqrt{K-1}}.$$

In addition, all inequalities are strict whenever gaps $(\Delta_a)_{a \neq a^}$ are not all equal.*

The monotonicity results of the lemma are illustrated in Figure 3.3. In Figure 3.4, we observe the evolution of the optimal weights when modifying the value of the optimal arm. It shows that we cannot give a general variation for sub-optimal arms that are not either the second or the worst arm: the weight of the third arm increases at the beginning and decreases. We observe that the weight of the optimal arm decreases and conjecture that this holds in all generality.

Before proving Lemma 3.6, we recall that the problem is scaling-invariant.

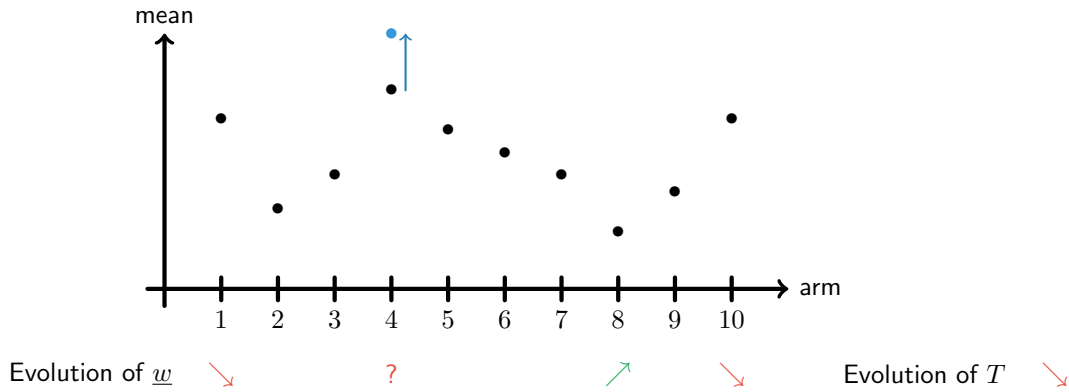


Figure 3.3: By **increasing the mean of the best arm**, the weight of the worst arm(s) increases, while those of the second best arm(s) decrease, and the characteristic time decreases. We also conjecture that the weight of the optimal arm decreases.

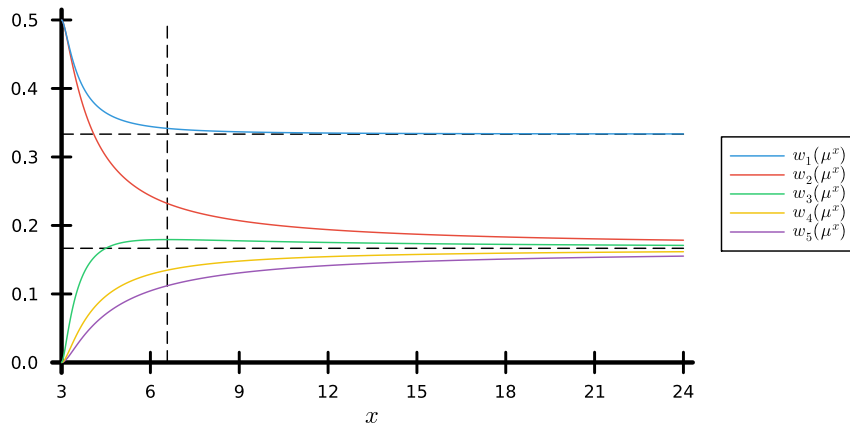


Figure 3.4: Evolution of the optimal weight $\underline{w}(\underline{\mu}^x)$ for $x \in [3, 24]$, where $\underline{\mu}^x = (x, 3, 2.7, 2.3, 2)$. The vertical dashed line corresponds to the maximizer of $w_3(\underline{\mu}^x)$, while the horizontal dashed lines are the limit values $\frac{1}{1+\sqrt{K-1}}$ and $\frac{1}{K-1+\sqrt{K-1}}$.

Lemma 3.7. *If there exists $\kappa > 0$ such that $\Delta'_a = \kappa \Delta_a$ for all $a \neq a^*$, then $\underline{w}' = \underline{w}$.*

Proof. By definition, \underline{w} and \underline{w}' are solutions of the same optimization problem up to the multiplicative constant κ , as (3.7) gives:

$$\{\underline{w}'\} = \operatorname{argmax}_{\underline{v} \in \Sigma_K} \frac{1}{2} \min_{a \neq a^*} \frac{v_{a^*} v_a}{v_{a^*} + v_a} \Delta_a^2 = \operatorname{argmax}_{\underline{v} \in \Sigma_K} \frac{1}{2} \min_{a \neq a^*} \frac{v_{a^*} v_a}{v_{a^*} + v_a} (\kappa \Delta_a)^2 = \{\underline{w}\}. \quad \square$$

Proof of Lemma 3.6. We first treat the easy case where all gaps $(\Delta_a)_{a \neq a^*}$ are equal to a common gap Δ . By the equality cases of Proposition 3.4, we know that

$$w'_{a^*} = \frac{1}{1 + \sqrt{K-1}} = w_{a^*},$$

and

$$T' = \frac{2(1 + \sqrt{K-1})^2}{(\Delta + d)^2} = \frac{\Delta^2}{(\Delta + d)^2} \frac{2(1 + \sqrt{K-1})^2}{\Delta^2} = \frac{\Delta^2}{(\Delta + d)^2} T < T.$$

As the optimal weights of the sub-optimal arms are equal by Equation (3.10), this ensures that

$$\forall a \neq a^*, \quad w'_a = \frac{1}{K-1 + \sqrt{K-1}} = w_a.$$

3.4. MONOTONICITY OF THE max-min PROBLEM

We now treat the case where at least two gaps are distinct, i.e., there exist sub-optimal arms a and b such that $\Delta_a \neq \Delta_b$.

1. Let $b \in \operatorname{argmin}_{a \neq a^*} \Delta_a$. Let us rescale bandit problem $\underline{\mu}'$ to obtain the same gap for arm b than in $\underline{\mu}$, by multiplying the gaps of $\underline{\mu}'$ by constant $\kappa = \frac{\Delta_b}{\Delta_b + d}$. Let $\underline{\mu}''$ denote the obtained bandit (taking arbitrarily $\underline{\mu}''_{a^*} = \underline{\mu}_{a^*}$), and set $\underline{\Delta}'' = \underline{\Delta}(\underline{\mu}'') = \kappa \underline{\Delta}'$ and $\underline{w}'' = \underline{w}(\underline{\mu}'')$. By Lemma 3.7, we have $\underline{w}'' = \underline{w}'$. Then, on the one hand

$$\Delta''_b = \kappa \Delta'_b = \frac{\Delta_b}{\Delta_b + d} (\Delta_b + d) = \Delta_b,$$

and for all sub-optimal arms $a \neq b$, by non-decreasing of $x \mapsto \frac{x}{x+d}$:

$$\Delta''_a = \kappa \Delta'_a = \frac{\Delta_b}{\Delta_b + d} (\Delta_a + d) \geq \frac{\Delta_a}{\Delta_a + d} (\Delta_a + d) = \Delta_a,$$

with strict inequality if $\Delta_a > \Delta_b$, that is, for at least one arm a as gaps are not all equal.

Applying point 2 of Lemma 3.5 to every sub-optimal arm $a \neq b$, we go from $\underline{\mu}$ to $\underline{\mu}''$ and the weight of arm b is non-decreasing, and increases during steps for which arm a is such that $\Delta''_a > \Delta_a$. Hence $w'_b = w(\underline{\mu}'') > w_b$.

2. The result can be obtained similarly to point 1 by rescaling bandit $\underline{\mu}'$ by constant $\kappa = \frac{\Delta_{\min}}{\Delta_{\min} + d}$.
3. Using Equations (3.8) and (3.5), we get:

$$T'^{-1} = \frac{1}{2} \min_{a \neq a^*} \frac{w'_{a^*} w'_a}{w'_{a^*} + w'_a} \Delta'^2_a > \frac{1}{2} \min_{a \neq a^*} \frac{w_{a^*} w_a}{w_{a^*} + w_a} \Delta'^2_a > \frac{1}{2} \min_{a \neq a^*} \frac{w_{a^*} w_a}{w_{a^*} + w_a} \Delta_a^2 = T^{-1},$$

where the first inequality is a consequence of the unicity of the optimal weight vector \underline{w} and the fact that (as obtained in the previous points) $\underline{w} \neq \underline{w}'$, and the second uses that $\Delta'_a > \Delta_a$ for all $a \neq a^*$.

4. Using the rescaling argument of the second point of this proof, we get that the limit, when d goes to ∞ , of bandit $\underline{\mu}''$ is the bandit problem $\underline{\mu}_{\lim}$ with best arm a^* and constant gap $\Delta = \Delta_b$. By continuity of the optimal weight vector, this implies that

$$\forall a \in [K], \quad \lim_{d \rightarrow +\infty} w'_a = w'_a(\underline{\mu}_{\lim}).$$

Injecting the values of the optimal weight vector for a bandit problem with constant gaps, recalled at the beginning of this proof, gives the result. \square

3.4.3. Increasing the Mean of the Worst Arms

Lastly, we look at the monotonicity of \underline{w} and T when increasing the worst means.

Lemma 3.8. *Let $B = \operatorname{argmin}_{a \in [K]} \mu_a$ (respectively $B' = \operatorname{argmin}_{a \in [K]} \mu'_a$) be the set of the worst arms of $\underline{\mu}$ (respectively $\underline{\mu}'$) and assume that $B \subset B'$ and $\Delta'_{\max} < \Delta_{\max}$, while $\Delta'_a = \Delta_a$ for all $a \notin B'$.*

1. $w'_{a^*} \leq w_{a^*}$,
2. if c is a sub-optimal arm such that $\Delta'_c = \Delta_c$, then $w'_{a^*} \leq w_{a^*}$,
3. $w'_{\min} \geq w_{\min}$,
4. $T' > T$,
5. we get

$$\lim_{\Delta'_{\max} \rightarrow \Delta_{\min}} w'_{a^*} = \frac{1}{1 + \sqrt{K-1}}, \quad \text{and} \quad \forall a \neq a^*, \quad \lim_{\Delta'_{\max} \rightarrow \Delta_{\min}} w'_a = \frac{1}{K-1 + \sqrt{K-1}}.$$

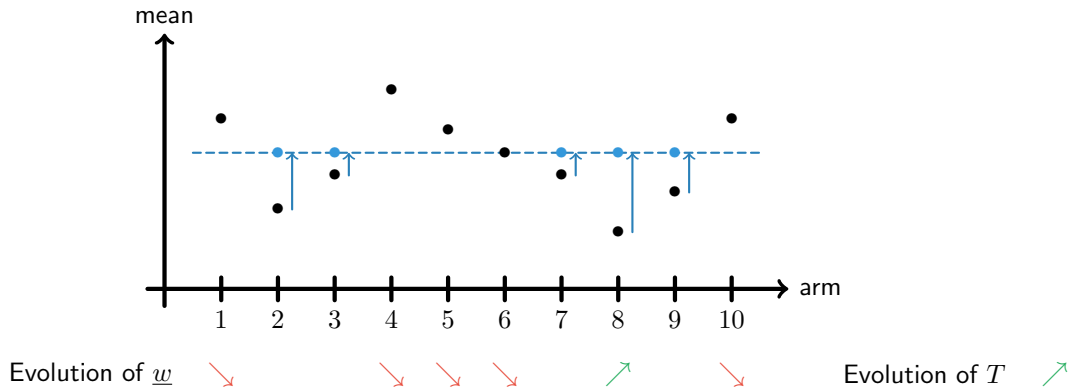


Figure 3.5: By increasing the means of the worst arms to a common value, the minimal optimal weight and the characteristic time increase, while the optimal weights of unmoved arms decrease.

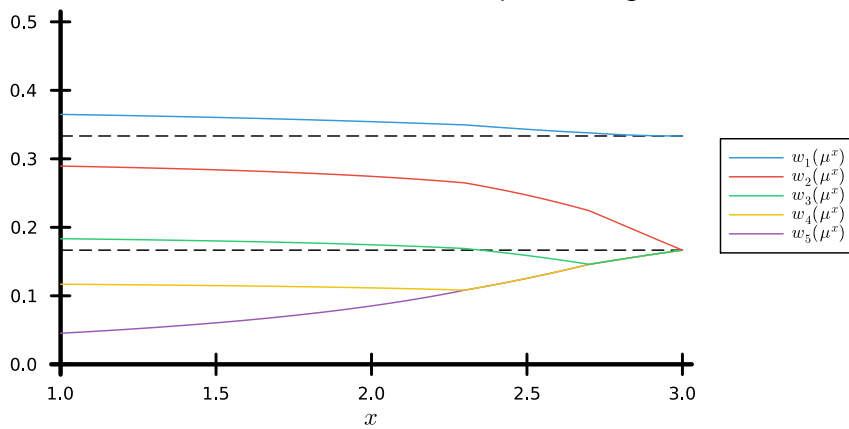


Figure 3.6: Evolution of the optimal weight vector $\underline{w}(\underline{\mu}^x)$ for $x \in [1, 3]$, where $\underline{\mu}^x = (5, 3, \max(2.7, x), \max(2.3, x), x)$. The dashed lines are the limit values $\frac{1}{1+\sqrt{K-1}}$ and $\frac{1}{K-1+\sqrt{K-1}}$.

The monotonicity results of the lemma are illustrated in Figure 3.5. In Figure 3.6, we observe the evolution of the optimal weights when modifying the minimal mean of the arms. It shows in particular the limit behaviors and the fact that while the mean of an arm is not modified, its weight decreases.

Proof. One can mimic the proof of Lemma 3.5 to prove that, if we equally increase the means of all the worst arms of $\underline{\mu}$ to a value which is smaller or equal to the second worst mean of $\underline{\mu}$, then:

1. w_{\min} increases,
2. w_a decreases for all $a \notin B \cup \{a^*\}$,
3. w_{a^*} decreases,
4. T increases.

By applying this observation a finite number of steps (one needs to split the reasoning each time a new arm is added to the list of worst arms), we deduce all but the last point.

Limit behaviors of point 5 might be obtained similarly to Lemma 3.6. □

3.5. Regularity Properties

In this section, we provide quantitative regularity results for the functions \underline{w} , T and g . The results presented in this section will prove to be essential to the non-asymptotic analysis of the Exploration-Biased-Sampling strategy in Chapter 4.

3.5.1. Regularity of \underline{w} and T

In this section, we show explicit bounds on the regularity of $\underline{\rho} \mapsto \underline{w}(\underline{\rho})$ and $\underline{\rho} \mapsto T(\underline{\rho})$. We keep the notation of the last section.

Theorem 3.9. *Assume that, for some $\varepsilon \in [0, 1/7]$,*

$$\forall a \neq a^*, \quad (1 - \varepsilon)\Delta_a^2 \leq \Delta_a'^2 \leq (1 + \varepsilon)\Delta_a^2. \quad (3.27)$$

Then

$$\forall a \in [K], \quad (1 - 10\varepsilon)w_a \leq w_a' \leq (1 + 10\varepsilon)w_a, \quad (3.28)$$

$$\text{and} \quad (1 - 3\varepsilon)T \leq T' \leq (1 + 6\varepsilon)T. \quad (3.29)$$

Proof. All along the proof, we will use the following inequalities, that might be easily checked,

$$\forall x \in \mathbb{R}, \quad \frac{1}{1+x} \leq 1 - x, \quad (3.30)$$

$$\forall u \in [0, 1/2], \quad \frac{1}{1-u} \leq 1 + 2u. \quad (3.31)$$

We will first prove the two following supporting results:

$$\frac{r}{1+\varepsilon} \leq r' \leq \frac{r}{1-\varepsilon}, \quad (3.32)$$

$$\text{and} \quad \frac{1-3\varepsilon}{w_{a^*}} \leq \frac{1-3\varepsilon}{1-\varepsilon} \frac{1}{w_{a^*}} \leq \frac{1}{w_{a^*}'} \leq \frac{1+5\varepsilon}{1+\varepsilon} \frac{1}{w_{a^*}} \leq \frac{1+5\varepsilon}{w_{a^*}}. \quad (3.33)$$

On the one hand, we get, using the right-inequality in (3.27),

$$\phi_{\underline{\mu}'}\left(\frac{r}{1+\varepsilon}\right) = \sum_{a \neq a^*} \frac{1}{\left(\frac{r}{1+\varepsilon}\Delta_a'^2 - 1\right)^2} - 1 \geq \sum_{a \neq a^*} \frac{1}{\left(\frac{r}{1+\varepsilon}\Delta_a^2(1+\varepsilon) - 1\right)^2} - 1 = \phi_{\underline{\mu}}(r) = 0,$$

and on the other hand, using the left inequality,

$$\phi_{\underline{\mu}'}\left(\frac{r}{1-\varepsilon}\right) = \sum_{a \neq a^*} \frac{1}{\left(\frac{r}{1-\varepsilon}\Delta_a'^2 - 1\right)^2} - 1 \leq \sum_{a \neq a^*} \frac{1}{\left(\frac{r}{1-\varepsilon}\Delta_a^2(1-\varepsilon) - 1\right)^2} - 1 = \phi_{\underline{\mu}}(r) = 0.$$

By decreasing of $\phi_{\underline{\mu}'}$ (Lemma 3.1) and definition of r' , we deduce that (3.32) holds.

We move to the proof of (3.33). For all $a \neq a^*$, we get

$$r'\Delta_a'^2 \leq \frac{1+\varepsilon}{1-\varepsilon}r\Delta_a^2 = (1+\eta)r\Delta_a^2, \quad \text{where} \quad \eta \stackrel{\text{def}}{=} \frac{1+\varepsilon}{1-\varepsilon} - 1 = \frac{2\varepsilon}{1-\varepsilon},$$

hence, by inequality (3.30),

$$\begin{aligned} \frac{1}{r'\Delta_a'^2 - 1} &\geq \frac{1}{(1+\eta)r\Delta_a^2 - 1} = \frac{1}{(r\Delta_a^2 - 1)\left(1 + \frac{\eta r\Delta_a^2}{r\Delta_a^2 - 1}\right)} \\ &\geq \frac{1}{r\Delta_a^2 - 1} \left(1 - \frac{\eta r\Delta_a^2}{r\Delta_a^2 - 1}\right) = \frac{1}{r\Delta_a^2 - 1} - \eta \frac{1}{r\Delta_a^2 - 1} - \eta \frac{1}{(r\Delta_a^2 - 1)^2}. \end{aligned} \quad (3.34)$$

By Equation (3.9), we get

$$\begin{aligned}
 \frac{1}{w'_{a^*}} &= 1 + \sum_{a \neq a^*} \frac{1}{r' \Delta_a'^2 - 1} \\
 &\geq 1 + (1 - \eta) \sum_{a \neq a^*} \frac{1}{r \Delta_a^2 - 1} - \eta \underbrace{\sum_{a \neq a^*} \frac{1}{(r \Delta_a^2 - 1)^2}}_{=\phi_{\underline{\mu}}(r)+1=1} \\
 &= (1 - \eta) \frac{1}{w_{a^*}} = \frac{1 - 3\varepsilon}{1 - \varepsilon} \frac{1}{w_{a^*}},
 \end{aligned}$$

which gives the second inequality of (3.33).

We move to the proof of the third inequality. Let $a \neq a^*$. We get

$$r' \Delta_a'^2 \geq \frac{1 - \varepsilon}{1 + \varepsilon} r \Delta_a^2 = (1 - \eta) r \Delta_a^2, \quad \text{where} \quad \eta \stackrel{\text{def}}{=} 1 - \frac{1 - \varepsilon}{1 + \varepsilon} = \frac{2\varepsilon}{1 + \varepsilon}.$$

Note that $\eta \leq 1/4$ as $\varepsilon \leq 1/7$. Using that $r \geq \frac{2}{\Delta_{\min}}$ by inequality (3.18), we get $r \Delta_a^2 \geq 2$, hence, by decreasing of $x \mapsto \frac{x}{x-1}$ on $(2, +\infty)$,

$$\eta \frac{r \Delta_a^2}{r \Delta_a^2 - 1} \leq \frac{1}{2}.$$

By using inequality (3.31), we have

$$\begin{aligned}
 \frac{1}{r' \Delta_a'^2 - 1} &\leq \frac{1}{(1 - \eta) r \Delta_a^2 - 1} = \frac{1}{(r \Delta_a^2 - 1) \left(1 - \frac{\eta r \Delta_a^2}{r \Delta_a^2 - 1}\right)} \\
 &\leq \frac{1}{r \Delta_a^2 - 1} \left(1 + 2 \frac{\eta r \Delta_a^2}{r \Delta_a^2 - 1}\right) = \frac{1}{r \Delta_a^2 - 1} + 2\eta \frac{1}{r \Delta_a^2 - 1} + 2\eta \frac{1}{(r \Delta_a^2 - 1)^2}.
 \end{aligned} \tag{3.35}$$

Consequently,

$$\begin{aligned}
 \frac{1}{w'_{a^*}} &= 1 + \sum_{a \neq a^*} \frac{1}{r' \Delta_a'^2 - 1} \\
 &\leq 1 + (1 + 2\eta) \sum_{a \neq a^*} \frac{1}{r \Delta_a^2 - 1} + 2\eta \underbrace{\sum_{a \neq a^*} \frac{1}{(r \Delta_a^2 - 1)^2}}_{\phi_{\underline{\mu}}(r)+1=1} \\
 &= (1 + 2\eta) \frac{1}{w_{a^*}} = \frac{1 + 5\varepsilon}{1 + \varepsilon} \frac{1}{w_{a^*}},
 \end{aligned}$$

which concludes the proof of (3.33).

We now deduce the results from inequalities (3.32) and (3.33). In all the following equations, we use at least one of the inequalities of (3.32) and (3.33), together with various bounds satisfied by $\varepsilon \leq \frac{1}{7}$. First, taking the inverse of (3.33), and using (3.30) for the left-inequality and (3.31) for the right-inequality, we get

$$(1 - 5\varepsilon)w_{a^*} \leq \frac{w_{a^*}}{1 + 5\varepsilon} \leq w'_{a^*} \leq \frac{w_{a^*}}{1 - 3\varepsilon} \leq (1 + 6\varepsilon)w_{a^*}.$$

Equation (3.10), with respectively (3.34) and (3.35), yield for all $a \neq a^*$:

$$w'_a = \frac{w'_{a^*}}{r'\Delta'_a{}^2 - 1} \geq \frac{\frac{w_{a^*}}{1+5\varepsilon}}{(r\Delta_a^2 - 1) \underbrace{\left(1 + \frac{2\varepsilon}{1-\varepsilon} \frac{r\Delta_a^2}{\Delta_a^2 - 1}\right)}_{\leq 2}} = \frac{1-\varepsilon}{(1+5\varepsilon)(1+3\varepsilon)} w_a \geq (1-10\varepsilon)w_a,$$

$$\text{and } w'_a = \frac{w'_{a^*}}{r'\Delta'_a{}^2 - 1} \leq \frac{\frac{w_{a^*}}{1-3\varepsilon}}{(r\Delta_a^2 - 1) \underbrace{\left(1 - \frac{2\varepsilon}{1+\varepsilon} \frac{r\Delta_a^2}{\Delta_a^2 - 1}\right)}_{\leq 2}} = \frac{1+\varepsilon}{(1-3\varepsilon)^2} w_a \leq (1+10\varepsilon)w_a.$$

This concludes the proof of (3.28). Finally, we obtain (3.29) by using Equation (3.11), as, on the one hand,

$$T' = \frac{2r'}{w'_{a^*}} \geq 2 \cdot \frac{r}{1+\varepsilon} \cdot \frac{1-3\varepsilon}{1-\varepsilon} \frac{1}{w_{a^*}} = \frac{1-3\varepsilon}{1-\varepsilon^2} \cdot \frac{2r}{w_{a^*}} \geq (1-3\varepsilon)T,$$

and, on the other hand,

$$T' = \frac{2r'}{w'_{a^*}} \leq 2 \cdot \frac{r}{1-\varepsilon} \cdot \frac{1+5\varepsilon}{1+\varepsilon} \frac{1}{w_{a^*}} = \frac{1+5\varepsilon}{1-\varepsilon^2} \cdot \frac{2r}{w_{a^*}} \leq (1+6\varepsilon)T. \quad \square$$

3.5.2. Regularity of g

The regularity of g might also be studied. Controlling the variations of g requires getting an upper bound on the gap values. For simplicity, we work here with the model $\mathcal{D}_{\mathcal{N}_1}^{[0,1]}$ of standard Gaussian variables with means in $[0, 1]$, but one might generalize it to all sub-models of $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ of standard Gaussian bandit problems with bounded gaps. We get the following property:

Proposition 3.10. *Let $\underline{\mu}$ and $\underline{\mu}'$ be two bandit problems in $\mathcal{D}_{\mathcal{N}_1}^{[0,1]}$, with potentially several best arms, and let $\underline{v} \in \Sigma_K$. Setting*

$$\varepsilon = \max_{a \in [K]} |\mu_a - \mu'_a|, \quad \text{and} \quad \eta = \max_{a \in [K]} \frac{|w_a(\underline{\mu}) - v_a|}{w_a(\underline{\mu})},$$

we get

$$g(\underline{\mu}', \underline{v}) \geq \frac{(1-\eta)^2}{1+\eta} \left(g(\underline{\mu}, \underline{w}(\underline{\mu})) - \frac{\varepsilon}{2} \right).$$

We will prove Proposition 3.10 by combining two Lemmas. The first lemma states that $g(\cdot, \underline{u})$ is $\frac{1}{2}$ -Lipschitz for the infinity norm.

Lemma 3.11. *Let $\underline{\mu}$ and $\underline{\mu}'$ be two bandit problems in $\mathcal{D}_{\mathcal{N}_1}^{[0,1]}$, with potentially several best arms. Then, for all optimal weight vectors $\underline{u} \in \Sigma_K$,*

$$|g(\underline{\mu}', \underline{u}) - g(\underline{\mu}, \underline{u})| \leq \frac{1}{2} \max_{a \in [K]} |\mu_a - \mu'_a|.$$

Proof. As $\underline{\mu}$ and $\underline{\mu}'$ play similar roles, we only prove that

$$g(\underline{\mu}', \underline{u}) \geq g(\underline{\mu}, \underline{u}) - \frac{\varepsilon}{2}, \quad \text{where} \quad \varepsilon = \max_{a \in [K]} |\mu_a - \mu'_a|.$$

- Assume first that $\underline{\mu}$ and $\underline{\mu}'$ have a common best arm denoted by a^* . Then by Equation (3.5),

$$\begin{aligned}
 g(\underline{\mu}', \underline{u}) - g(\underline{\mu}, \underline{u}) &= \frac{1}{2} \min_{a \neq a^*} \frac{u_{a^*} u_a}{u_{a^*} + u_a} \Delta_a'^2 - \frac{1}{2} \min_{b \neq a^*} \frac{u_{a^*} u_b}{u_{a^*} + u_b} \Delta_b^2 \\
 &= \frac{1}{2} \min_{a \neq a^*} \max_{b \neq a^*} \frac{u_{a^*} u_a}{u_{a^*} + u_a} \Delta_a'^2 - \frac{u_{a^*} u_b}{u_{a^*} + u_b} \Delta_b^2 \\
 &\geq \frac{1}{2} \min_{a \neq a^*} \frac{u_{a^*} u_a}{u_{a^*} + u_a} \left(\Delta_a'^2 - \Delta_a^2 \right). \tag{3.36}
 \end{aligned}$$

by taking $b = a$. Fix $a \neq a^*$. One has:

$$|\Delta_a - \Delta_a'| = |(\mu_{a^*} - \mu_{a^*}') - (\mu_a - \mu_a')| \leq |\mu_{a^*} - \mu_{a^*}'| + |\mu_a - \mu_a'| \leq 2\varepsilon,$$

from which we obtain, using that the gaps are bounded by 1

$$|\Delta_a^2 - \Delta_a'^2| = |\Delta_a - \Delta_a'|(\Delta_a + \Delta_a') \leq 4\varepsilon.$$

As \underline{u} is an optimal weight vector, it satisfies $u_a \leq u_{a^*} \leq \frac{1}{2}$ by Corollary 3.3 and bound (3.19), so that:

$$\frac{u_{a^*} u_a}{u_{a^*} + u_a} \leq \frac{1}{2} \frac{u_a}{u_{a^*} + u_a} \leq \frac{1}{2} \frac{u_a}{2u_a} = \frac{1}{4}.$$

Finally,

$$\frac{u_{a^*} u_a}{u_{a^*} + u_a} \left(\Delta_a'^2 - \Delta_a^2 \right) \geq -\varepsilon.$$

This concludes the proof by injecting this inequality into (3.36).

- In case $\underline{\mu}$ and $\underline{\mu}'$ do not have a common best arm, define the family of bandits $(\underline{\mu}^{(t)})_{t \in [0,1]}$ by

$$\forall t \in [0, 1], \forall a \in [K], \quad \mu_a^{(t)} = (1-t)\mu_a + t\mu_a'.$$

One can check that $\underline{\mu} = \underline{\mu}^{(0)}$, $\underline{\mu}' = \underline{\mu}^{(1)}$ and

$$\forall t_1, t_2 \in [0, 1], \quad \max_{a \in [K]} |\mu_a^{(t_1)} - \mu_a^{(t_2)}| \leq |t_1 - t_2| \cdot \varepsilon \leq \varepsilon. \tag{3.37}$$

Select the subdivision $0 = t_0 < t_1 < \dots < t_N = 1$ of times at which the optimal arms of $\underline{\mu}^{(t)}$ are modified. Note that $N \geq 2$ as $\underline{\mu}$ and $\underline{\mu}'$ do not have a common best arm. Note that by continuity:

- for any $1 \leq n \leq N-1$, $\underline{\mu}^{(t_n)}$ has at least two best arms so that $g(\underline{\mu}^{(t_n)}, \underline{u}) = 0$,
- $\underline{\mu}^{(1)}$ and $\underline{\mu}$ have a common best arm,
- $\underline{\mu}^{(N-1)}$ and $\underline{\mu}'$ have a common best arm.

Thus, applying the first part of the proof,

$$\begin{aligned}
 g(\underline{\mu}', \underline{u}) - g(\underline{\mu}, \underline{u}) &= \underbrace{g(\underline{\mu}', \underline{u}) - g(\underline{\mu}^{(t_{N-1})}, \underline{u})}_{=0} + \underbrace{g(\underline{\mu}^{(t_1)}, \underline{u}) - g(\underline{\mu}, \underline{u})}_{=0} \\
 &\geq g(\underline{\mu}^{(t_1)}, \underline{u}) - g(\underline{\mu}, \underline{u}) \\
 &\geq -\frac{1}{2} \max_{a \in [K]} |\mu_a^{(t_1)} - \mu_a| \\
 &\geq -\frac{\varepsilon}{2},
 \end{aligned}$$

where we used (3.37) in the last inequality. This concludes the proof. \square

3.6. CONCLUSION

The second lemma controls the variations of $g(\underline{\mu}', \cdot)$. Note that this lemma does not require the use of bounded gaps.

Lemma 3.12. *Let $\underline{\mu}'$ be a bandit problem in $\mathcal{D}_{\mathcal{N}_1}$ with potentially several arms. Let $\underline{u}, \underline{v} \in \Sigma_K$ be such that, for a fixed $\eta \in [0, 1]$,*

$$\max_{a \in [K]} \frac{|u_a - v_a|}{u_a} \leq \eta. \quad (3.38)$$

Then:

$$g(\underline{\mu}', \underline{v}) \geq \frac{(1 - \eta)^2}{1 + \eta} g(\underline{\mu}', \underline{u}).$$

Proof. Let a^* is (one of) the best arm(s) of $\underline{\mu}'$. Note that condition (3.38) can be written as

$$\forall a \in [K], \quad (1 - \eta)u_a \leq v_a \leq (1 + \eta)u_a.$$

Hence, for all $a \neq a^*$,

$$\frac{v_{a^*}v_a}{v_{a^*} + v_a} \geq \frac{(1 - \eta)^2}{(1 + \eta)} \cdot \frac{u_1u_a}{u_1 + u_a},$$

which entails

$$g(\underline{\mu}', \underline{v}) = \min_{a \neq a^*} \frac{v_{a^*}v_a}{v_{a^*} + v_a} \Delta'_a{}^2 \geq \frac{(1 - \eta)^2}{1 + \eta} \min_{a \neq a^*} \frac{u_1u_a}{u_1 + u_a} \Delta'_a{}^2 = \frac{(1 - \eta)^2}{1 + \eta} g(\underline{\mu}', \underline{u}). \quad \square$$

Proof of Proposition 3.10. The result follows directly by Lemmas 3.12 and 3.11 with $\underline{u} = \underline{w}(\underline{\mu})$:

$$g(\underline{\mu}', \underline{v}) \geq \frac{(1 - \eta)^2}{1 + \eta} g(\underline{\mu}', \underline{w}(\underline{\mu})) \geq \frac{(1 - \eta)^2}{1 + \eta} \left(g(\underline{\mu}, \underline{w}(\underline{\mu})) - \frac{\varepsilon}{2} \right). \quad \square$$

3.6. Conclusion

In this chapter, we studied the sample complexity optimization 3.2 for a Gaussian model. The fact that this optimization problem only depends on gaps allowed us to obtain a new characterization of the solution $\underline{w}(\underline{\mu})$. This characterization implies new general bounds for $\underline{w}(\underline{\mu})$ and $T(\underline{\mu})$, together with monotonicity properties that will be useful for the definition of the strategy Exploration-Biased-Sampling in Chapter 4.

We also investigated the regularity of the solution, proving quantitative results that complement the continuity result of [Garivier and Kaufmann \(2016\)](#). Obtaining such results in the context of a general exponential model could be an interesting perspective for future work, but requires dealing with a less explicit and easy-to-handle Kullback-Leibler divergence.

CHAPTER 4

A Fixed-Confidence Strategy with Non-Asymptotic Guarantees

In this chapter, we propose a new strategy for the problem of best-arm identification with fixed-confidence of Gaussian variables. This strategy, called `Exploration-Biased-Sampling`, is not only asymptotically optimal: it is to the best of our knowledge the first strategy with non-asymptotic bounds that asymptotically matches the sample complexity. Its sampling rule is built in order to naturally encourage exploration and uses the results of Chapter 3. The content of the present chapter is extracted from Sections 1, 2 and 4 (and associated appendices) of the conference paper



A. Barrier, A. Garivier, and T. Kocák. A Non-Asymptotic Approach to Best-Arm Identification for Gaussian Bandits. In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, pages 10078–10109. PMLR, 2022

Contents

1	Introduction	94
2	The <code>Exploration-Biased-Sampling</code> Strategy	97
1	Conservative Tracking	97
2	The Strategy	100
3	Theoretical Results	103
4	Numerical Experiments	104
5	Proof of the Non-Asymptotic Bounds of Theorem 4.5	110
1	Step 1: Controlling the Difference between Vectors \underline{w} and $\underline{N}(t)/t$	111
2	Step 2: a Useful Inclusion of Events	113
3	Step 3: Bounding $\mathbb{P}_{\underline{\mu}}(\tau_{\delta} > t \cap \mathcal{E})$ and $\mathbb{E}_{\underline{\mu}}[\tau_{\delta}\mathbb{I}\{\mathcal{E}\}]$	114
6	Technical Results	116
1	Proof of Lemma 4.4	116
2	Technical Details for the Proof of Theorem 4.5	117
3	Almost Sure Asymptotic Bound	120
7	Conclusion	121

4.1. Introduction

Many modern systems of automatic decisions (from recommender systems to clinical trials, through auto-ML and parameter tuning) require finding the best among a set of options, using noisy observations obtained by successive calls to a random mechanism (see, e.g., [Lattimore and Szepesvári, 2020](#)). The simplest formal model for such situations is the *standard Gaussian multi-armed bandit*, a collection of $K \geq 2$ independent Gaussian distributions called *arms* of unknown means $\underline{\mu} = (\mu_a)_{a \in [K]} \in \mathbb{R}^K$ and variances all equal to 1. They are sampled sequentially and independently: at every discrete time step $t \in \mathbb{N}^*$, a learner chooses an arm $A_t \in [K]$ based on past information, and observes an independent draw Y_t from distribution $\mathcal{N}(\mu_{A_t}, 1)$.

Best-arm identification. The *best-arm identification* problem consists in identifying the arm with highest mean of $\underline{\mu}$. Unless otherwise specified, we only consider bandit problems with a unique optimal arm:

$$\{a^*(\underline{\mu})\} \stackrel{\text{def}}{=} \operatorname{argmax}_{a \in [K]} \mu_a.$$

The corresponding distribution mean will be denoted by μ^* . In the *fixed-confidence* setting (see [Even-Dar et al., 2006](#); [Kalyanakrishnan et al., 2012](#)), a confidence parameter $\delta \in (0, 1)$ is given, and the objective is to design strategies that will stop after some (random) finite number of observations and give an estimate of the best arm $a^*(\underline{\mu})$ which is correct with probability at least $1 - \delta$. Such strategies are called δ -correct, and their performance is measured by how quick they are to take a decision.

Formally, a strategy is defined by

- a *sampling rule*, which consists in choosing the arm $A_t \in [K]$ to observe at each time step $t \geq 1$. This arm A_t depends on the previous observations Y_1, \dots, Y_{t-1} , but also possibly on some external randomization that we capture by the random variable U_{t-1} . A_t is thus \mathcal{F}_{t-1} -measurable, where $\mathcal{F}_{t-1} \stackrel{\text{def}}{=} \sigma(I_{t-1})$ with $I_{t-1} \stackrel{\text{def}}{=} (U_0, Y_1, U_1, Y_2, U_2, \dots, Y_{t-1}, U_{t-1})$. I_{t-1} corresponds to the information available at the end of the time step $t - 1$.
- a *stopping rule* τ_δ , which is a stopping time with respect to the filtration $(\mathcal{F}_t)_{t \geq 0}$,
- a *decision rule* \hat{a}_{τ_δ} which is $\mathcal{F}_{\tau_\delta}$ -measurable.

The general structure of a strategy is presented in Algorithm 12. A strategy is δ -correct if for all bandit problems $\underline{\mu}$ in the model $\mathcal{D}_{\mathcal{N}_1}$ of standard Gaussian variables,

$$\mathbb{P}_{\underline{\mu}}(\tau_\delta < +\infty, \hat{a}_{\tau_\delta} \neq a^*(\underline{\mu})) \leq \delta,$$

where $\mathbb{P}_{\underline{\mu}}$ is the probability distribution under bandit problem $\underline{\mu}$.

Lower bound. The aim of the problem is to find strategies that minimize the *sample complexity*, that is, the expected number of samplings $\mathbb{E}_{\underline{\mu}}[\tau_\delta]$. The sample complexity of δ -correct strategies cannot be arbitrarily good: it has been proved by [Garivier and Kaufmann \(2016\)](#) that they obey the lower bound

$$\forall \underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}, \quad \mathbb{E}_{\underline{\mu}}[\tau_\delta] \geq T(\underline{\mu}) \log \frac{1}{2.4\delta}, \quad (4.1)$$

where the *characteristic time* $T(\underline{\mu})$ is the solution of the following optimization problem

$$T(\underline{\mu})^{-1} = \sup_{\underline{v} \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a \frac{(\mu_a - \lambda_a)^2}{2}, \quad (4.2)$$

where

$$\Sigma_K = \left\{ \underline{v} \in [0, 1]^K : v_1 + \dots + v_K = 1 \right\} \quad \text{and} \quad \text{Alt}(\underline{\mu}) = \left\{ \underline{\lambda} \text{ in } \mathcal{D}_{\mathcal{N}_1}^{[0,1]} : a^*(\underline{\lambda}) \neq a^*(\underline{\mu}) \right\}.$$

Algorithm 12: General structure of a fixed-confidence strategy

Input: confidence parameter δ
sampling-rule, stopping-condition, decision-rule

Output: stopping time τ_δ
estimated best arm \hat{a}_{τ_δ}

```
1 Observe each arm once // initialization
2  $t \leftarrow K$ 
3 while stopping-condition( $I_t, \delta$ ) is not satisfied do
4   | Increase  $t$  by 1
5   |  $A_t \leftarrow$  sampling-rule( $I_{t-1}$ )
6   | Observe  $Y_t \sim \mathcal{N}(\mu_{A_t}, 1)$ 
7  $\tau_\delta \leftarrow t$ 
8  $\hat{a}_{\tau_\delta} \leftarrow$  decision-rule( $I_{\tau_\delta}$ )
```

An optimal weight vector. The information-theoretic analysis of [Garivier and Kaufmann \(2016\)](#) also highlights the nature of the optimal sampling strategy: whatever the value of the risk δ , one should sample the arms with frequencies proportional to $\underline{v} = \underline{w}(\underline{\mu})$, the (unique and well-defined) maximizer in the right-hand side of Equation (4.2). This observation allowed the authors to introduce Track-and-Stop, the first asymptotically optimal strategy, which satisfies,

$$\forall \underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}, \quad \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_\delta]}{\log \frac{1}{\delta}} = T(\underline{\mu}),$$

proving, by passing, that the lower bound (4.1) is tight. The algorithm works as follows: at every time step t , the optimal frequencies relative to an estimate $\hat{\underline{\mu}}(t)$ of the mean parameter $\underline{\mu}$ are computed and used to determine which action is to be selected next: we pick the action that lays the most behind its estimated optimal frequency unless one action was severely under-sampled (in which case its exploration is forced). Some improvements were proposed: for example, [Ménard \(2019\)](#) proved that it is not necessary to solve the optimization problem in every time step.

The shortcomings of Track-and-Stop. The Track-and-Stop algorithm is not only a theoretical contribution, but it also proved to be numerically efficient, far exceeding its competitors in a wide variety of settings. It was improved in different directions ([Degenne and Koolen, 2019](#); [Degenne et al., 2019](#); [Shang et al., 2020](#)), and also provides a simple template for extensions, for bandit problems with structure ([Kocák and Garivier, 2020](#)), as long as the optimization problem (4.2) can be solved. Yet, Track-and-Stop suffers from certain shortcomings. First, a close look into the proofs shows that the theoretical guarantees proved so far are really asymptotic in nature. Second, the forced exploration appears very arbitrary, with a rate of \sqrt{t} that has no other justification than lying somewhere between constant and linear functions. Third, the sampling strategy appears to be pretty unstable, especially at the beginning: the target frequencies can vary significantly as the estimated means fluctuate before stabilizing around their expectations. Fourth, Track-and-Stop does not present the intuitively desirable behavior to sample uniformly in the beginning, until sufficient information has been gathered for significant differences between the arms to emerge. This is in contrast with strategies like Racing ([Kaufmann and Kalyanakrishnan, 2013](#)), which are sub-optimal but intuitively appealing. Altogether, these issues lead for example to unpredictable and irregular conduct at the beginning of multiple A/B testing cases with many arms very close to optimal.

Towards non-asymptotic bounds. While the proven optimality of Track-and-Stop is purely asymptotic, a different approach is followed in (Karnin et al., 2013; Jamieson et al., 2014; Chen et al., 2017) for moderate values of δ . The proposed strategies are sub-optimal by a multiplicative constant but are proven to satisfy explicit non-asymptotic bounds. More recently, Degenne et al. (2019) obtained a general non-asymptotic bound, a remarkable but hardly comparable result in particular settings.

Outline and contributions. In this chapter, we try to make a link between both approaches by introducing a new strategy, Exploration-Biased-Sampling, with non-asymptotic guarantees which, in the regime $\delta \rightarrow 0$, correspond to the lower bound on the sample complexity. Additionally, Exploration-Biased-Sampling solves the issues of Track-and-Stop mentioned above. We work with the model $\mathcal{D}_{\mathcal{N}_1}^{[0,1]}$ of standard Gaussian variables with means in $[0, 1]$, but all results can be generalized to¹ the model $\mathcal{D}_{\mathcal{N}_{\sigma^2}}^{\mathcal{M}}$ of Gaussian variables with means in a bounded set \mathcal{M} and common variance σ^2 .

We present in Section 4.2 our new Exploration-Biased-Sampling strategy. The exploration is conducted very differently, in a statistically natural way that softens the fluctuations of empirical means and avoids arbitrary parameters. It results in a stabilized sampling strategy, that is much easier to follow and understand. Its theoretical properties and guarantees are stated in Section 4.3, including a non-asymptotic analysis with finite risk bounds for which the proof is presented in Section 4.5. These results have required developing a careful analysis of the quantitative regularity of the solution to the optimization problem (4.2) that was presented in Chapter 3. Lastly, we illustrate the performance and behavior of our strategy by numerical experiments in Section 4.4.

Note that, independently, Wang et al. (2021) obtained a sampling rule based on a Frank-Wolfe method for which they proved finite risk analysis and asymptotic optimality. Our finite risk bound has a better asymptotic behavior, but a worse behavior in the regime where gaps go to zero.



Notation. For the simplicity of the presentation, when there is no confusion we set

$$a^* = a^*(\underline{\mu}), \quad \Delta_a = \Delta_a(\underline{\mu}), \quad \underline{w} = \underline{w}(\underline{\mu}), \quad \text{and} \quad T = T(\underline{\mu}),$$

where $\Delta_a(\underline{\mu}) = \mu^* - \mu_a$ is the gap of arm a . We recall (see Section 3.1) that the optimization problem rewrites:

$$T(\underline{\mu})^{-1} = \sup_{\underline{v} \in \Sigma_K} g(\underline{\mu}, \underline{v}), \quad \text{where} \quad g(\underline{\mu}, \underline{v}) = \underbrace{\frac{v_{a^*} v_a}{v_{a^*} + v_a} \frac{\Delta_a^2}{2}}_{=TC_{a \rightarrow a^*}(\underline{\mu}, \underline{v})},$$

and hence the optimal weight vector \underline{w} satisfies:

$$\{\underline{w}\} = \operatorname{argmax}_{\underline{v} \in \Sigma_K} g(\underline{\mu}, \underline{v}) = \operatorname{argmax}_{\underline{v} \in \Sigma_K} \frac{1}{2} \min_{a \neq a^*} \frac{v_{a^*} v_a}{v_{a^*} + v_a} \Delta_a^2.$$

Garivier and Kaufmann (2016) proved that, at the optimum \underline{w} , all transportation costs are equal, so that, for all arms $a \neq a^*$,

$$T(\underline{\mu})^{-1} = g(\underline{\mu}, \underline{w}) = \frac{1}{2} \frac{w_{a^*} w_a}{w_{a^*} + w_a} \Delta_a^2. \quad (4.3)$$

Remark. Unless explicitly stated, we only consider instances with a unique optimal arm. However, we might sometimes use that $g(\underline{\mu}, \underline{v})$ is also defined for bandit problems $\underline{\mu}$ that admit several optimal

¹The algorithms discussed here can be used with a sub-Gaussian model with a known upper bound on the variances. However, for such models, the sample complexity bounds proved in this chapter apply but are not necessarily optimal.

arms (with a^* being any of the optimal arms). In that case, $g(\underline{\mu}, \underline{v}) = 0$, $T(\underline{\mu}) = +\infty$, and we define the optimal weight vector $\underline{w}(\underline{\mu})$, somewhat arbitrarily, by

$$\underline{w}_a(\underline{\mu}) = \frac{\mathbb{I}\{a \in \mathcal{A}^*(\underline{\mu})\}}{\text{card } \mathcal{A}^*(\underline{\mu})},$$

where $\mathcal{A}^*(\underline{\mu}) = \text{argmax}_{a \in [K]} \mu_a$ is the set of optimal arms of $\underline{\mu}$.

For a given strategy facing a bandit problem $\underline{\mu}$, let $N_a(t)$ and $\hat{\mu}_a(t)$ denote the number of pulls and the empirical mean² of arm a at step t :

$$N_a(t) \stackrel{\text{def}}{=} \sum_{s \in [t]} \mathbb{I}\{A_s = a\} \quad \text{and} \quad \hat{\mu}_a(t) \stackrel{\text{def}}{=} \frac{1}{N_a(t)} \sum_{s \in [t]} Y_s \mathbb{I}\{A_s = a\}.$$

Without loss of generality (see the paragraph on optional skipping page 37), we assume that the observation at time step t is $Y_t = X_{A_t, N_{A_t}(t)}$, where $(X_{a,n})_{a \in [K], n \geq 1}$ are independent random variables such that $X_{a,n} \sim \mathcal{N}(\mu_{A_t}, 1)$ for all $a \in [K]$ and $n \geq 1$. As a consequence, we notably get

$$\hat{\mu}_a(t) = \frac{1}{N_a(t)} \sum_{n=1}^{N_a(t)} X_{a,n} \stackrel{\text{def}}{=} \hat{\mu}_{a, N_a(t)}. \quad (4.4)$$

4.2. The Exploration-Biased-Sampling Strategy

In this section, we introduce our new strategy called Exploration-Biased-Sampling. Instead of Track-and-Stop's greedy choice of actions based on a plug-in estimate of $\underline{\mu}$, it relies on a specific estimator that is biased toward uniform exploration. We fix $\underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}^{[0,1]}$ with a unique optimal arm and define the following quantities, for which the dependency with respect to $\underline{\mu}$ is omitted if there is no confusion:

$$\Delta_{\min} \stackrel{\text{def}}{=} \min_{a \neq a^*} \Delta_a > 0, \quad \Delta_{\max} \stackrel{\text{def}}{=} \max_{a \in [K]} \Delta_a, \quad \text{and} \quad w_{\min} \stackrel{\text{def}}{=} \min_{a \in [K]} w_a.$$

We recall that w_{\min} is the optimal weight value of the worst arm(s) of $\underline{\mu}$.

4.2.1. Conservative Tracking

The main idea of the algorithm is to design a sampling policy of arms that naturally encourages exploration without forcing it like Track-and-Stop does. To do so, the objective is to “wrap” the optimal weight vector $\underline{w}(\underline{\mu})$ “from above”, by ensuring that we never under-estimate its minimal value. Indeed, even an arm with a low mean needs to be sampled sufficiently often until one is very confident that it is sub-optimal. The idea is to construct a confidence region $\mathcal{CR}_{\underline{\mu}} \subset [0, 1]^K$ for $\underline{\mu}$ on which one can efficiently find a bandit $\tilde{\underline{\mu}} \in \mathcal{CR}_{\underline{\mu}}$ maximizing the minimal weight w_{\min} :

$$\tilde{\underline{\mu}} \in \underset{\underline{\rho} \in \mathcal{CR}_{\underline{\mu}}}{\text{argmax}} w_{\min}(\underline{\rho}). \quad (4.5)$$

As long as $\underline{\mu}$ belongs to the confidence region $\mathcal{CR}_{\underline{\mu}}$, choosing the target weights $\underline{w}(\tilde{\underline{\mu}})$ guarantees that every arm is explored sufficiently, as $w_{\min}(\tilde{\underline{\mu}}) \geq w_{\min}(\underline{\mu})$. The exploration bias decreases with the number of observations, as $\mathcal{CR}_{\underline{\mu}}$ shrinks to $\{\underline{\mu}\}$, and in the end arms are sampled with frequencies close to the optimal weight vector $\underline{w}(\underline{\mu})$.

²As strategies initially observe each arm once, $\hat{\mu}_a(t)$ is well-defined for $t \geq K$.

This approach to exploration requires two ingredients:

- the exploration-biased bandit $\tilde{\mu}$ needs to be efficiently computable. It turns out to be the case if the confidence region is a product of confidence intervals on each arm (a mild requirement since the arms are independent). We propose Algorithm 13, an efficient procedure for computing $\tilde{\mu}$. Intuitively, maximizing w_{\min} over $\mathcal{CR}(\underline{\mu})$ requires increasing and equalizing all the positive gaps as much as possible. The associated bandit will indeed be the one for which it is harder to identify the second-best arm and thus it will require to sample the worst arms more frequently. This gives a candidate bandit for each potential best arm, and our algorithm compares those candidates. Figure 4.1 illustrates on an example the principle of Algorithm 13, whose correctness is proved in Proposition 4.1. The algorithm requires Gaussian-Optimal-Weights (Algorithm 11 of Chapter 3), an efficient procedure for solving optimization problem (4.2).
- the regularity of the mapping $\underline{\rho} \mapsto \underline{w}(\underline{\rho})$ needs to be explicitly known. Indeed, the confidence region will decrease with the number of observations, and $\tilde{\mu}$ will come close to $\underline{\mu}$. The continuity proved by Garivier and Kaufmann (2016) for the asymptotic optimality of Track-and-Stop is not sufficient: the first quantitative bounds were given in Section 3.5.

One can remark that as long as the confidence intervals have a non-empty intersection, which means the observations do not permit to exclude that any of them is optimal, the exploration-biased weights returned by Algorithm 13 are uniform and the arms are sampled in a round-robin way (as in a Racing or Successive Elimination algorithm like in Even-Dar et al., 2006).

Proposition 4.1. *Let $\mathcal{CR} = \prod_{a \in [K]} [\mu_a^-, \mu_a^+] \subset [0, 1]^K$ be a confidence region and*

$$(\tilde{\mu}, \tilde{w}) \leftarrow \text{Exploration-Biased-Weights}(\mathcal{CR}).$$

Then $\tilde{w} = \underline{w}(\tilde{\mu})$ and $\tilde{\mu}$ satisfies Equation (4.5).

The proof relies on the results of Section 3.4.

Proof. We assume that $K \geq 3$, otherwise $\underline{w}(\underline{\rho}) = (\frac{1}{2}, \frac{1}{2})$ for all $\underline{\rho}$ and the result is clear.

With the notation of Algorithm 13, we first observe that $\tilde{w} = \underline{w}(\tilde{\mu})$. When $\text{minUB} \geq \text{maxLB}$ the algorithm returns a constant bandit and $\tilde{w} = (\frac{1}{K}, \dots, \frac{1}{K})$ is its optimal weight vector by convention. As all optimal weight vectors $\underline{w}(\underline{\rho})$ are such that $w_{\min}(\underline{\rho}) \leq \frac{1}{K}$, we obtain that $\tilde{\mu}$ satisfies (4.5).

Now assume that $\text{minUB} < \text{maxLB}$, i.e., at least two confidence intervals are disjoint, and fix $\underline{\rho} \in \mathcal{CR}$. If $\underline{\rho}$ has several optimal arms, then $w_{\min}(\underline{\rho}) = 0$ so that trivially $w_{\min}(\underline{\rho}) \leq w_{\min}(\tilde{\mu})$. Assume now that $\underline{\rho}$ has a unique optimal arm denoted by a . Note that $a \in \text{PotentialBest}$, so that we will show that

$$w_{\min}(\underline{\rho}) \leq w_{\min}(\tilde{\mu}^{\text{test}(a)}) \leq \tilde{w}_{\min}.$$

The latest inequality stems from the choice of $\tilde{\mu}$ by the procedure into the set of potential most exploring bandits $\{\tilde{\mu}^{\text{test}(a)} : a \in \text{PotentialBest}\}$. To obtain the first inequality, we will transform the means of $\underline{\rho}$ to those of $\tilde{\mu}^{\text{test}(a)}$ by modifications that only increase w_{\min} , thanks to the results of Section 3.4. We recall that the value of w_{\min} is the optimal weight value of any of the worst arms. The procedure, illustrated in Figure 4.2, is the following:

1. Transform $\underline{\rho}$ into $\underline{\rho}^{(1)}$ by increasing arm a so that $\underline{\rho}_a^{(1)} = \mu_a^+$. Using Lemma 3.6, one has

$$w_{\min}(\underline{\rho}^{(1)}) \geq w_{\min}(\underline{\rho}).$$

2. Transform $\underline{\rho}^{(1)}$ into $\underline{\rho}^{(2)}$ by decreasing, for each arm $b \neq a$, μ_b to $\max(\mu_b^-, \rho_{\min})$, where $\rho_{\min} = \min_{c \in [K]} \rho_c$. By several applications³ of Lemma 3.5, one has

$$w_{\min}(\underline{\rho}^{(2)}) \geq w_{\min}(\underline{\rho}^{(1)}).$$

³Note that, as the modified means do not go below ρ_{\min} , the worst arms of $\underline{\rho}$ stay in the set of worst arms.

Algorithm 13: Exploration-Biased-Weights

Input: confidence region $\mathcal{CR} = \prod_{a \in [K]} [\mu_a^-, \mu_a^+]$
Output: exploration-biased bandit $\tilde{\underline{\mu}} \in \mathcal{CR}$

 exploration-biased optimal weight vector $\tilde{\underline{w}} = \underline{w}(\tilde{\underline{\mu}})$

```

1 maxLB  $\leftarrow \max_{a \in [K]} \mu_a^-$ 
2 minUB  $\leftarrow \min_{a \in [K]} \mu_a^+$ 
3 if minUB  $\geq$  maxLB then
4      $\tilde{\underline{\mu}} \leftarrow (\text{minUB}, \dots, \text{minUB})$ 
5      $\tilde{\underline{w}} \leftarrow (\frac{1}{K}, \dots, \frac{1}{K})$ 
6 else
7     PotentialBest  $\leftarrow \{a \in [K] : \mu_a^+ > \text{maxLB}\}$ 
8      $\tilde{\underline{w}} \leftarrow (0, \dots, 0)$ 
9     for  $a \in \text{PotentialBest}$  do
10         $\tilde{\mu}_a^{\text{test}(a)} \leftarrow \mu_a^+$ 
11        for  $b \in [K] \setminus \{a\}$  do
12             $\tilde{\mu}_b^{\text{test}(a)} \leftarrow \max(\mu_b^-, \text{minUB})$ 
13         $\underline{w}^{\text{test}(a)} \leftarrow \text{Gaussian-Optimal-Weights}(\tilde{\underline{\mu}}^{\text{test}(a)})$ 
14        if  $\min_{b \in [K]} w_b^{\text{test}(a)} > \min_{b \in [K]} \tilde{w}_b$  then
15             $\tilde{\underline{w}} \leftarrow \underline{w}^{\text{test}(a)}$ 
16             $\tilde{\underline{\mu}} \leftarrow \tilde{\underline{\mu}}^{\text{test}(a)}$ 
    
```

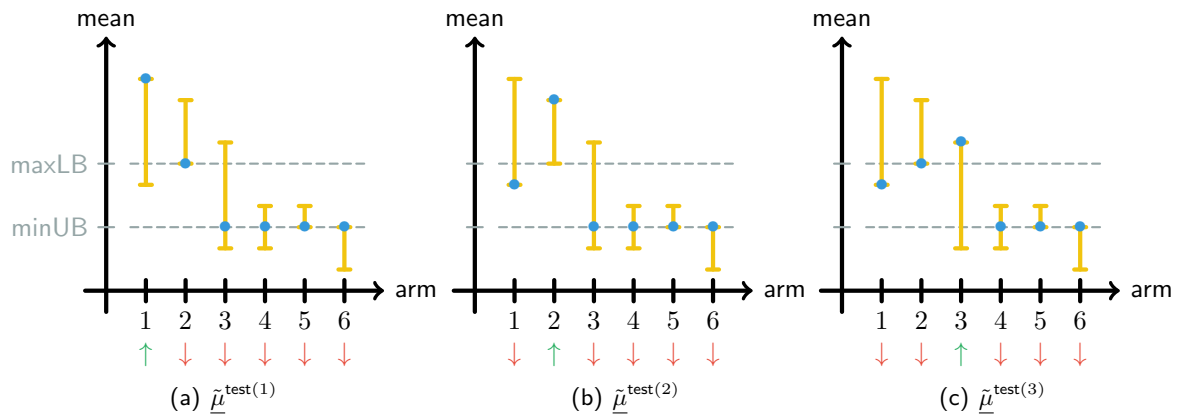


Figure 4.1: List of bandits $(\tilde{\underline{\mu}}^{\text{test}(a)})_{a \in \text{PotentialBest}}$ tried by Algorithm 13 for the example confidence region in yellow. As only the three first arms are above maxLB, they are the potential candidates for being the best arm: $\text{PotentialBest} = \{1, 2, 3\}$. For each potential best arm, we associate a unique bandit instance that “maximizes exploration” by putting the mean of that arm as high as possible (↑) and all other arms as lower as possible (↓) while staying above minUB.

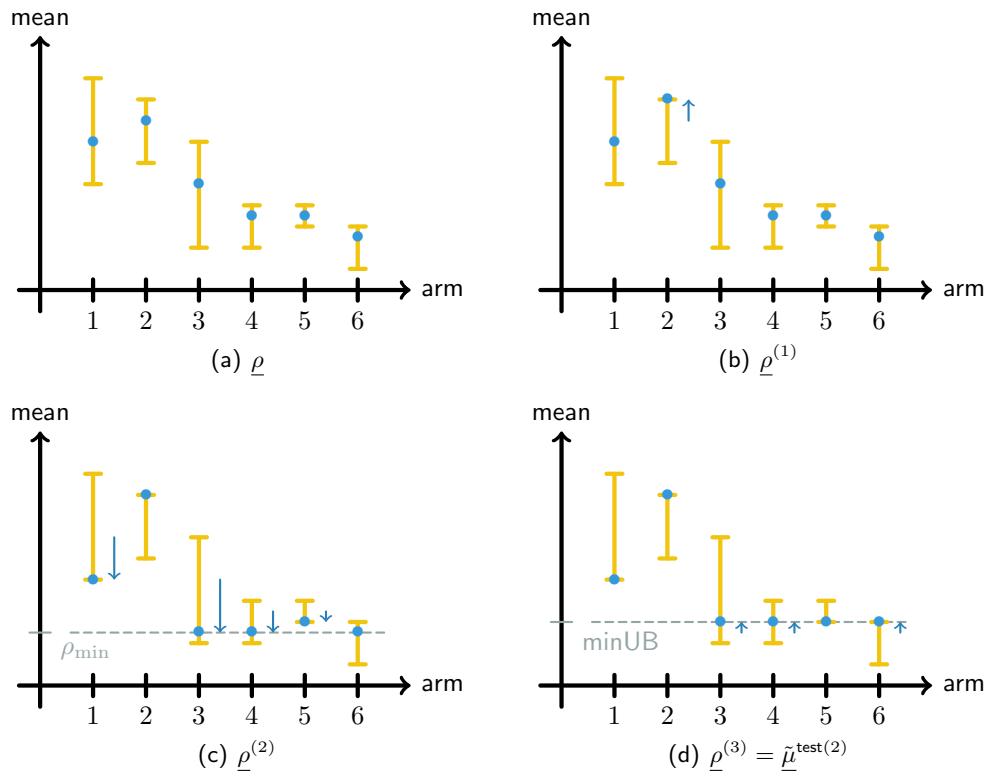


Figure 4.2: Transformations in the proof of Proposition 4.1, for some instance bandit $\underline{\rho}$.

3. Transform $\underline{\rho}^{(2)}$ into $\underline{\rho}^{(3)}$ by increasing all the worst arms to minUB. By Lemma 3.8, one has

$$w_{\min}(\underline{\rho}^{(3)}) \geq w_{\min}(\underline{\rho}^{(2)}).$$

Observing that we now get $\underline{\rho}^{(3)} = \tilde{\underline{\mu}}^{\text{test}(a)}$, we proved that

$$w_{\min}(\underline{\rho}) \leq w_{\min}(\tilde{\underline{\mu}}^{\text{test}(a)}).$$

In conclusion,

$$\max_{\underline{\rho} \in \mathcal{CR}} w_{\min}(\underline{\rho}) = \max_{a \in \text{PotentialBest}} w_{\min}(\tilde{\underline{\mu}}^{\text{test}(a)}) = w_{\min}(\tilde{\underline{\mu}}),$$

where the last inequality comes from the procedure defining $\tilde{\underline{\mu}}$. □

4.2.2. The Strategy

We are now able to introduce our strategy called Exploration-Biased-Sampling. Given a risk $\delta \in (0, 1)$ and a threshold function $\beta(t, \delta)$, we compute at each time confidence intervals for each μ_a that will ensure $\underline{\mu}$ to belong to each associated confidence region with probability at least $1 - \gamma$, where $\gamma \in (0, 1)$ is a fixed parameter. We can then ensure enough exploration by biasing the optimal weight vector $\underline{w}(\underline{\mu})$ using Algorithm 13.

Confidence regions. Confidence regions are designed to satisfy two requirements. First, we need products of confidence intervals in order to use Algorithm 13, and then, we will require a time-uniform

confidence guarantee as a key ingredient for the non-asymptotic analysis of Exploration-Biased-Sampling. For $\gamma \in (0, 1)$ and $t \geq K$, we define

$$\mathcal{CR}_{\underline{\mu}}(t) \stackrel{\text{def}}{=} \prod_{a \in [K]} \left[\hat{\mu}_a(t) \pm C_{\frac{\gamma}{K}}(N_a(t)) \right], \quad \text{where} \quad C_{\gamma}(s) \stackrel{\text{def}}{=} 2\sqrt{\frac{\log\left(\frac{4s}{\gamma}\right)}{s}}. \quad (4.6)$$

The following lemma states a time-uniform γ -confidence guarantee for $\underline{\mu}$.

Lemma 4.2. *Fix $\gamma \in (0, 1)$. For all bandit problems $\underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}^{[0,1]}$ (with potentially several optimal arms), it holds that*

$$\mathbb{P}_{\underline{\mu}}\left(\exists t \geq K : \underline{\mu} \notin \mathcal{CR}_{\underline{\mu}}(t)\right) \leq \gamma.$$

Proof. By union bound we only have to show that, for all $a \in [K]$,

$$\mathbb{P}_{\underline{\mu}}\left(\exists t \geq K : |\hat{\mu}_a(t) - \mu_a| \geq C_{\frac{\gamma}{K}}(N_a(t))\right) \leq \frac{\gamma}{K}.$$

Let $a \in [K]$. We recall that $N_a(t) \geq 1$ for $t \geq K$, as all arms are observed once at the beginning. Thus, using Equation (4.4), we have

$$\begin{aligned} \mathbb{P}_{\underline{\mu}}\left(\exists t \geq K : |\hat{\mu}_a(t) - \mu_a| \geq C_{\frac{\gamma}{K}}(N_a(t))\right) &= \mathbb{P}_{\underline{\mu}}\left(\exists t \geq K : |\hat{\mu}_{a, N_a(t)} - \mu_a| \geq C_{\frac{\gamma}{K}}(N_a(t))\right) \\ &= \mathbb{P}_{\underline{\mu}}\left(\exists n \in \mathbb{N}^* : |\hat{\mu}_{a, n} - \mu_a| \geq C_{\frac{\gamma}{K}}(n)\right). \end{aligned}$$

By a peeling argument (see, e.g., [Boucheron et al., 2013](#)), this leads to

$$\begin{aligned} \mathbb{P}_{\underline{\mu}}\left(\exists n \in \mathbb{N}^* : |\hat{\mu}_{a, n} - \mu_a| \geq C_{\frac{\gamma}{K}}(n)\right) &\leq \sum_{k \geq 0} \mathbb{P}_{\underline{\mu}}\left(\exists n \in [2^k, 2^{k+1}] : \left| \frac{1}{n} \sum_{s \in [n]} (X_{a, s} - \mu_a) \right| \geq C_{\frac{\gamma}{K}}(n)\right) \\ &= \sum_{k \geq 0} \mathbb{P}_{\underline{\mu}}\left(\exists n \in [2^k, 2^{k+1}] : \left| \sum_{s \in [n]} X_{a, s} - \mu_a \right| \geq n C_{\frac{\gamma}{K}}(n)\right) \\ &\leq \sum_{k \geq 0} \mathbb{P}_{\underline{\mu}}\left(\exists n \in [0, 2^{k+1}] : \left| \sum_{s \in [n]} X_{a, s} - \mu_a \right| \geq 2^k C_{\frac{\gamma}{K}}(2^k)\right) \\ &\leq 2 \sum_{k \geq 0} \exp\left(-\frac{(2^k C_{\frac{\gamma}{K}}(2^k))^2}{2 \times 2^{k+1}}\right) \\ &= 2 \sum_{k \geq 0} \exp\left(-\log \frac{2^{k+2} K}{\gamma}\right) \\ &= 2 \frac{\gamma}{K} \sum_{k \geq 0} \frac{1}{2^{k+2}} \\ &= \frac{\gamma}{K}, \end{aligned}$$

where the second inequality is obtained using the fact that $n \mapsto n C_{\frac{\gamma}{K}}(n)$ is non-decreasing and the second inequality is a well-known concentration bound for the sum of sub-Gaussian variables (see [Lattimore and Szepesvári, 2020](#), Theorem 9.2). \square

Sampling rule. The sampling rule of Exploration-Biased-Sampling is summarized in Algorithm 14. As explained in [Garivier and Kaufmann \(2016\)](#), one can either follow the exploration-biased weights directly (D-tracking) or their cumulative sums (C-tracking). The theoretical results of the strategy will be derived with the use of C-tracking, while we will run the experiments with both options, as D-tracking appears to perform slightly better.

Algorithm 14: Exploration-Biased-Sampling sampling rule at step $t > K$

Input: history of observations I_{t-1}
 confidence parameter γ
Output: next arm to observe A_t

- 1 $\mathcal{CR}_\mu(t-1) \stackrel{\text{def}}{=} \prod_{a \in [K]} \left[\hat{\mu}_a(t-1) \pm C_{\frac{\gamma}{K}}(N_a(t-1)) \right]$
- 2 $(\tilde{\mu}(t-1), \tilde{w}(t-1)) \leftarrow \text{Exploration-Biased-Weights}(\mathcal{CR}_\mu(t))$
 /* C-tracking */
- 3 Choose $A_t \in \operatorname{argmin}_{a \in [K]} N_a(t-1) - \sum_{s \in [t-1]} \tilde{w}_a(s)$
 /* D-tracking */
- 4 Choose $A_t \in \operatorname{argmin}_{a \in [K]} N_a(t-1) - (t-1)\tilde{w}_a(t-1)$

Stopping and decision rule. Following [Garivier and Kaufmann \(2016\)](#), our stopping rule relies on the statistic

$$Z(t) \stackrel{\text{def}}{=} \max_{a \in [K]} \min_{b \neq a} Z_{a,b}(t),$$

where $Z_{a,b}(t)$ is the Generalized Likelihood Ratio statistic (see [Chernoff, 1959](#)), which is defined, for the standard Gaussian model, as

$$Z_{a,b}(t) \stackrel{\text{def}}{=} \frac{1}{2} \frac{N_a(t)N_b(t)}{N_a(t) + N_b(t)} (\hat{\mu}_a(t) - \hat{\mu}_b(t))^2 \operatorname{sgn}(\hat{\mu}_a(t) - \hat{\mu}_b(t)). \quad (4.7)$$

The stopping rule consists of stopping the procedure if $Z(t)$ exceeds some threshold $\beta(t, \delta)$, as described in Algorithm 15. See page 50 for more details about the Global-Likelihood-Ratio stopping rule. The best empirical arm is then recommended by the strategy: we choose

$$\hat{a}_{\tau_\delta} \in \operatorname{argmax}_{a \in [K]} \hat{\mu}_a(\tau_\delta).$$

Algorithm 15: Global-Likelihood-Ratio stopping rule at step $t > K$

Input: history of observations I_t
 threshold function $\beta(t, \delta)$

- 1 $Z(t) \leftarrow \max_{a \in [K]} \min_{b \neq a} Z_{a,b}(t)$ // $Z_{a,b}(t)$ is defined in (4.7)
- 2
- 3 **if** $Z(t) > \beta(t, \delta)$ **then**
- 4 | Stop
- 5 **else**
- 6 | Continue

4.3. Theoretical Results

We present the theoretical guarantees of the Exploration-Biased-Sampling algorithm.

δ -correctness. Garivier and Kaufmann (2016) proved that, whatever the sampling rule, the use of the Global-Likelihood-Ratio stopping rule ensures the δ -correct property for some suitable threshold.

Proposition 4.3. [Garivier and Kaufmann, 2016, Proposition 12]

Let $\delta \in (0, 1)$ and $\alpha > 1$. There exists a constant $R = R(\alpha, K)$ such that, whatever the sampling rule, using the Global-Likelihood-Ratio stopping rule (Algorithm 15) with threshold

$$\beta(t, \delta) \stackrel{\text{def}}{=} \log \frac{Rt^\alpha}{\delta}, \quad (4.8)$$

and recommending the best empirical arm ensure that the strategy is δ -correct:

$$\forall \underline{\mu} \text{ in } \mathcal{D}_{\mathcal{N}_1}^{[0,1]}, \quad \mathbb{P}_{\underline{\mu}}(\tau_\delta < +\infty, \hat{a}_{\tau_\delta} \neq a^*(\underline{\mu})) \leq \delta.$$

The Exploration-Biased-Sampling strategy is hence δ -correct if used with threshold (4.8), whatever the choice of parameter $\gamma \in (0, 1)$.

Sufficient exploration. Interestingly, it happens that the choice of the confidence regions (4.6) naturally leads to a minimal exploration rate of the arms of order \sqrt{t} .

Lemma 4.4. Let $\gamma \in (0, 1)$. The Exploration-Biased-Sampling strategy satisfies, for all bandit problems $\underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}^{[0,1]}$,

$$\forall t \geq 0, \forall a \in [K], \quad N_a(t) \geq \frac{2}{K}\sqrt{t} - K.$$

What is surprising is that this is exactly the arbitrary rate used by Track-and-Stop for forced exploration. The proof of this lemma can be found in Section 4.6.1. Other practical advantages of Exploration-Biased-Sampling over Track-and-Stop are discussed in Section 4.4.

Non-asymptotic bound. Our main result is to obtain high probability bounds for the sample complexity of Exploration-Biased-Sampling in finite horizon.

Theorem 4.5. Fix $\gamma \in (0, 1)$, $\alpha \in [1, 2]$, $\eta \in (0, 1)$, and let $\underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}^{[0,1]}$. There exists an event \mathcal{E} of probability at least $1 - \gamma$ and $\delta_0 \stackrel{\text{def}}{=} \delta_0(\underline{\mu}, K, \gamma, \eta, \alpha) > 0$ such that, for all $0 < \delta \leq \delta_0$, algorithm Exploration-Biased-Sampling with threshold (4.8) satisfies

$$\forall t > (1 + \eta)T(\underline{\mu}) \log \frac{1}{\delta}, \quad \mathbb{P}_{\underline{\mu}}(\tau_\delta > t \cap \mathcal{E}) \leq 2Kt \exp\left(-\frac{tw_{\min}(\underline{\mu})}{4T(\underline{\mu})^2} \frac{1}{\log^{\frac{2}{3}} \frac{1}{\delta}}\right), \quad (4.9)$$

and

$$\mathbb{E}_{\underline{\mu}}[\tau_\delta \mathbb{I}\{\mathcal{E}\}] \leq (1 + \eta)T(\underline{\mu}) \log \frac{1}{\delta} + \frac{2^7 K T(\underline{\mu})^4}{w_{\min}(\underline{\mu})^2} \exp\left(-\frac{w_{\min}(\underline{\mu})}{4T(\underline{\mu})} \log^{\frac{1}{3}} \frac{1}{\delta}\right) \log^2 \frac{1}{\delta}. \quad (4.10)$$

Section 4.5 will be devoted to the proof of this theorem. The proof will highlight why, contrary to Track-and-Stop, the exploration strategy of Exploration-Biased-Sampling is adequate for obtaining non-asymptotic bounds. Note that:

- the proof of Theorem 4.5 provides an explicit expression for δ_0 ,
- the second term of bound (4.10) vanishes when δ decreases to 0, and hence negligible with respect to the first term: the sample complexity is therefore arbitrarily close to lower bound (4.1),
- The dependency on $w_{\min}(\underline{\mu})$ of the bounds might be replaced by the minimal gap Δ_{\min} , as Lemma 4.12 in Section 4.6.1 ensures that

$$\forall \underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}^{[0,1]}, \quad w_{\min}(\underline{\mu}) \geq \frac{\Delta_{\min}(\underline{\mu})}{2K}.$$

Asymptotic optimality. We additionally prove that, from an asymptotic point of view, our strategy presents the same guarantees as Track-and-Stop (see also Theorem 4.17 in Section 4.6.3).

Theorem 4.6. *Let $\gamma \in (0, 1)$, $\alpha \in (1, e/2]$. Algorithm Exploration-Biased-Sampling with threshold (4.8) satisfies, for all bandit problems $\underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}^{[0,1]}$,*

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log 1/\delta} \leq \alpha T(\underline{\mu}). \quad (4.11)$$

The proof of this result can be found in Section 4.6.3.

Remark.

- The multiplicative factor α in bound (4.11) can be avoided using recent proof techniques, as we explain in Section 6.4 (see also Remark 2.11).
- It is worth mentioning that the guarantees of Exploration-Biased-Sampling presented in this section hold true not only for Gaussian arms but more generally for 1-sub-Gaussian arms with means in $[0, 1]$ (in which case, of course, a better lower bound might hold). Indeed, these proofs only rely on sub-Gaussian deviation bounds.

4.4. Numerical Experiments

In this section, we discuss the behavior and performance of Exploration-Biased-Sampling for practical values of risk δ . We propose a comparison with Track-and-Stop, Chernoff-Racing and LUCB++, and begin by recalling a quick description of those strategies:

- the Track-and-Stop strategy (see 2.2.5) tracks the optimal weight vector $\underline{w}(\underline{\mu})$ by using the plugin estimate $\underline{w}(\hat{\underline{\mu}}(t))$. Some exploration rate is forced to ensure that bad initial observations do not lead to an under-sampling of some arms. The stopping rule is the same as the one presented for Exploration-Biased-Sampling,
- the Chernoff-Racing algorithm is an elimination algorithm (see Algorithm 3 in Section 2.2.1): the strategy maintains a list of candidate arms, starting with all arms, and divides the exploration into rounds, during which each arm of the list is observed once. At the end of each round, a decision is made to keep or eliminate the current worst arm from the active set. Several decision rules are possible, we will use the Chernoff rule presented in (Garivier and Kaufmann, 2016), which eliminates the worst empirical arm a_r at the end of round r if

$$Z_{\ell_r, a_r}(t_r) = \frac{r}{4} (\hat{\mu}_{\ell_r}(t_r) - \hat{\mu}_{a_r}(t_r))^2 > \beta(t_r, \delta),$$

where ℓ_r (respectively t_r) is the best empirical arm (respectively the number of times steps) at the end of round r .

4.4. NUMERICAL EXPERIMENTS

- the LUCB++ strategy, introduced by [Simchowitz et al. \(2017\)](#) (see also [Kalyanakrishnan et al., 2012](#); [Howard et al., 2021](#)) samples two arms at each round: the one with the current best estimate and the one in the remaining arms with the highest optimistic index $U_a(t)$ which is an upper confidence bound⁴:

$$U_a(t) \stackrel{\text{def}}{=} \hat{\mu}_a(t) + \sqrt{\frac{3}{N_a(t)} \log\left(\frac{2K \log(N_a(t))}{\delta}\right)}.$$

For the fairness of the comparison, we will take the same stopping condition as the strategies Track-and-Stop and Exploration-Biased-Sampling.

For all strategies, we ran our experiments with the same threshold, given by

$$\beta(t, \delta) = \log\left(\frac{1 + \log(t)}{\delta}\right),$$

and we use, for the Exploration-Biased-Sampling strategy, the confidence lengths

$$C_\gamma(s) = \sqrt{\frac{\log(\frac{s}{\gamma})}{s}}.$$

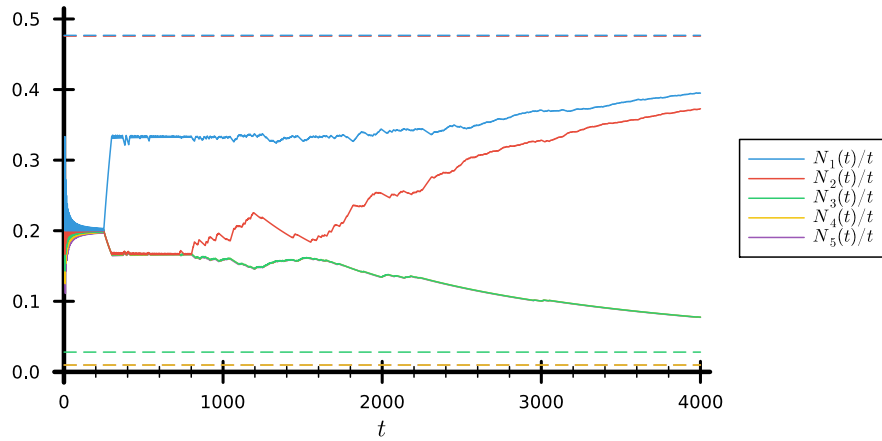
These choices are more aggressive than what the theoretical analysis suggests: yet, empirically, they appear to guarantee the desired failure rate. Using the larger intervals of Section 4.2 would have increased the number of rounds with uniform exploration, and using larger thresholds unnecessarily delays the stopping for all strategies.

Improving the stability of Track-and-Stop. In Section 4.1, we highlighted the weaknesses of Track-and-Stop (see page 95), especially the forced exploration parameter and the non-interpretable and unstable sampling strategy during the first rounds. In Figures 4.3 and 4.4 we see the improvements of Exploration-Biased-Sampling concerning those behaviors. During the first rounds, as for a racing algorithm, a uniform sampling phase is observed as the learner has not collected enough information (the confidence intervals on all arms are not separated), which is the expected behavior. Then the best arms are sampled more and more often, but still in a more cautious way than Track-and-Stop. We observe in Figure 4.4 the stability of the sampling strategies compared to Track-and-Stop during the first rounds: the targeted weights of Exploration-Biased-Sampling are stable and separate from each other cautiously (note that the three last arms still have the same weight at time 1200) whereas, for Track-and-Stop, we observe an important variation of the targeted weights with time. As a matter of fact, there is a clear discontinuity each time the estimated best arm changes, as we can see with the red and green arms. We also remark that Track-and-Stop uses forced exploration at regular rounds (for the yellow and purple dots), which is unnecessary for Exploration-Biased-Sampling as a natural exploration is always performed (Lemma 4.4).

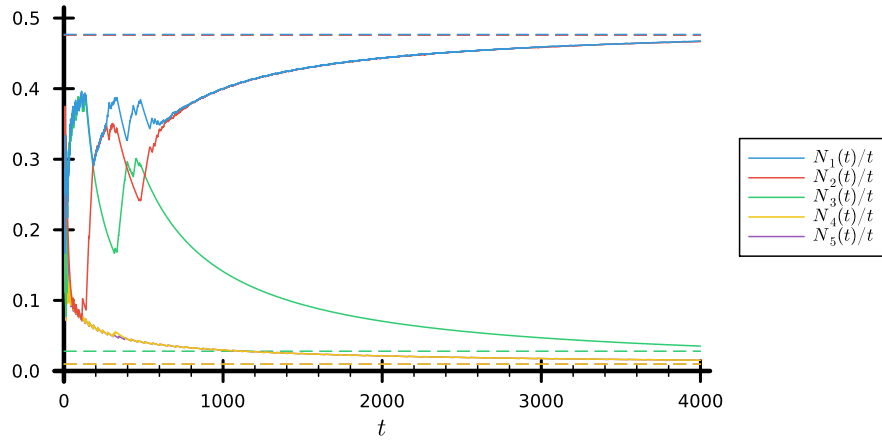
Comparisons of the strategies. The cost of the cautiousness of the algorithm (the exploration-biased weights) is that in terms of pure performance, it takes a little longer for the proportions of draws of Exploration-Biased-Sampling to converge to the optimal weight vector. This results in a slightly larger stopping time than Track-and-Stop that occurs for all bandit parameters⁵. In other words, Exploration-Biased-Sampling does not improve the numerical efficiency of Track-and-Stop. This can be observed in Table 4.1, where we present the performances of the strategies with two scenarios and a set of parameters. Exploration-Biased-Sampling globally performs

⁴Constant $\sqrt{3}$ appeared to be empirically optimal.

⁵Note that the cautiousness of our strategy is required to obtain the non-asymptotic bounds of Theorem 4.5.



(a) Exploration-Biased-Sampling



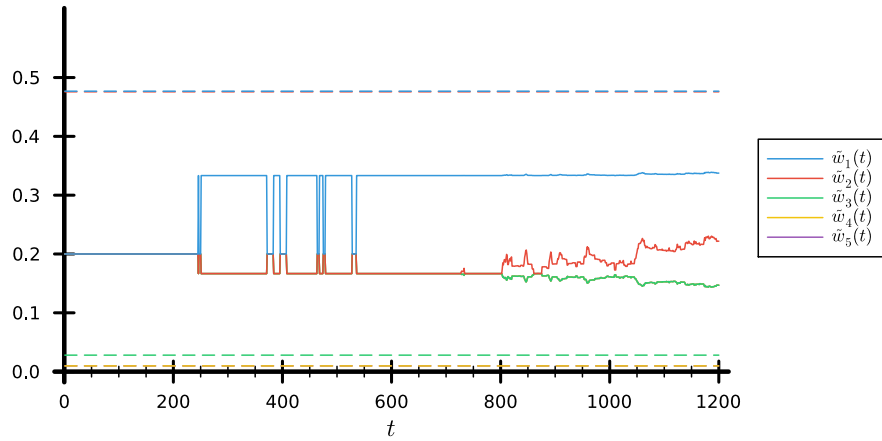
(b) Track-and-Stop

Figure 4.3: Evolution of the sampling frequencies $\frac{N(t)}{t}$ on a single simulation of the strategies Exploration-Biased-Sampling and Track-and-Stop with D-tracking. The parameters are $\delta = 0.01$, $\gamma = 0.2$ and $\underline{\mu} = (0.9, 0.8, 0.6, 0.4, 0.4)$. The values of $\underline{w}(\underline{\mu}) = (0.477, 0.476, 0.028, 0.010, 0.010)$ are dashed.

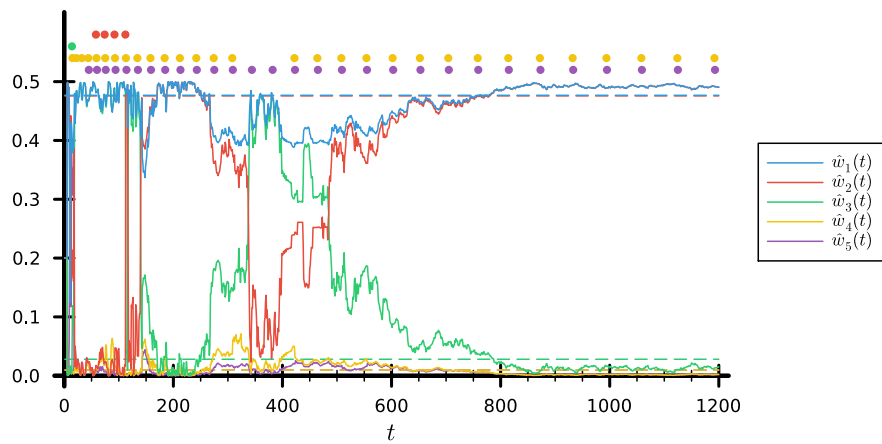
Table 4.1: Empirical expected number of draws $\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]$, averaged over 1000 experiments, with $\underline{\mu}^{(1)} = (0.9, 0.8, 0.6, 0.4, 0.4)$, $\underline{w}(\underline{\mu}^{(1)}) = (0.477, 0.476, 0.028, 0.010, 0.010)$ and $\underline{\mu}^{(2)} = (0.9, 0.5, 0.45, 0.4)$, $\underline{w}(\underline{\mu}^{(2)}) = (0.375, 0.286, 0.195, 0.144)$.

Instance	δ	γ	$T \text{kl}(\delta, 1 - \delta)$	EBS-C	TaS-C	EBS-D	TaS-D	CR	LUCB++
$\underline{\mu}^{(1)}$	0.1	0.05	1476	4727	3597	4191	3477	3124	3353
$\underline{\mu}^{(1)}$	0.01	0.05	3782	7363	5664	6330	5584	5419	5549
$\underline{\mu}^{(1)}$	0.01	0.2	3782	7090	5664	6136	5584	5419	5372
$\underline{\mu}^{(1)}$	10^{-5}	0.2	9669	13801	12181	12376	11439	11557	11644
$\underline{\mu}^{(2)}$	0.1	0.05	135	476	367	470	322	405	365
$\underline{\mu}^{(2)}$	0.01	0.05	347	708	588	699	485	542	565

4.4. NUMERICAL EXPERIMENTS



(a) Exploration-Biased-Sampling



(b) Track-and-Stop

Figure 4.4: Evolution of the Targeted Weights $\tilde{w}(t)$ (respectively $\hat{w}(t)$) During the First 1200 Rounds on a Simulation of Exploration-Biased-Sampling (respectively Track-and-Stop). ($\delta = 0.01$, $\gamma = 0.2$, $\underline{\mu} = (0.9, 0.8, 0.6, 0.4, 0.4)$)

correctly but we see that the other strategies are always a little more efficient. Note that when increasing γ , the confidence intervals reduce so that the targeted weights are closer to \underline{w} , improving the performance of the algorithm. For similar reasons, the initial cautiousness of the strategy disappears in the long term, thus when δ is very small the relative performance of Track-and-Stop and Exploration-Biased-Sampling gets closer. Of course, Exploration-Biased-Sampling overperforms Chernoff-Racing in the long run when the optimal weight vector is far from the sampling proportions of Chernoff-Racing (e.g., when $w_1 \gg w_2$).

Chernoff-Racing shows great performance with both $\underline{\mu}^{(1)}$ and $\underline{\mu}^{(2)}$. This strategy samples the two last arms of the race equally often, thus can be optimal only when $\underline{w}(\underline{\mu})$ has its two highest components of similar value, e.g. when the two best arms are well separated from the others: this is the case of bandit $\underline{\mu}^{(1)}$. For $\underline{\mu}^{(2)}$, any strategy performs well as the problem is easy. However, Chernoff-Racing (whose theoretical analysis remains to be written) leads to a few more misidentifications in our experiments that might be linked to the stopping rule we chose here; for fairness reasons, it was taken as identical to that of the other algorithms. LUCB++ presents similar performance with Chernoff-Racing, which can be explained by the similar behavior of the strategies: LUCB++ samples half of the time the best arm asymptotically, and the worst arms are eliminated one by one once their indexes fall under the two best estimates.

Finally, note that D-tracking shows better performance than C-tracking, either for Exploration-Biased-Sampling and Track-and-Stop: D-tracking indeed benefits directly from the current estimate of $\underline{\mu}$ (thus the empirical proportions of draws converge faster to $\underline{w}(\underline{\mu})$), while the impact is diluted in time with C-tracking. However, we did not prove theoretical guarantees for D-tracking.

Dependency on the confidence parameter δ . We present numerical experiments to compare the dependency on parameter δ of Exploration-Biased-Sampling, Track-and-Stop and Uniform-Sampling (that samples arms uniformly).

In Figure 4.5, we plot for each strategy and several bandit parameters the estimate of $\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]$ for different values of δ (using the same threshold β as in the experiments of Section 4 and $\gamma = 0.1$ for Exploration-Biased-Sampling). We also plot in black the lower bound of [Garivier and Kaufmann \(2016\)](#), which is of order $T(\underline{\mu}) \log \frac{1}{\delta}$ when δ goes to 0.

In terms of performance, we observe that Exploration-Biased-Sampling is always between Uniform-Sampling and Track-and-Stop. More precisely, there are different behaviors:

- For some problems (like for instance $\underline{\mu}^{(1)}$), Exploration-Biased-Sampling behaves almost like (but always a little worse than) Track-and-Stop. For this instance, the uniform sampling phase of Exploration-Biased-Sampling is relatively small compared to the required number of samples so that Exploration-Biased-Sampling has time to shrink its confidence regions close to parameter $\underline{\mu}$ and thus behaves like Track-and-Stop.
- When the problem is easier (with large gaps, see $\underline{\mu}^{(2)}$), the sample complexity is very low and the confidence regions of Exploration-Biased-Sampling do not have enough time to shrink. It results in a performance close to Uniform-Sampling.
- When the problem is difficult (with small gaps, see $\underline{\mu}^{(3)}$), it takes a large number of samples for Exploration-Biased-Sampling before leaving the uniform exploration phase, and this results in a behavior close to Uniform-Sampling for moderate values of δ . When δ decreases, there is a separation between Exploration-Biased-Sampling and Uniform-Sampling as more and more simulations reach the non-uniform sampling phase of our strategy. For even smaller values of δ , one can expect that Exploration-Biased-Sampling will come closer to Track-and-Stop than Uniform-Sampling, for the same reasons as before: the confidence regions of Exploration-Biased-Sampling have more time to shrink. This is what we observe with bandit instance $\underline{\mu}^{(4)}$, for which Exploration-Biased-Sampling behaves like Uniform-Sampling for moderate values of δ and like Track-and-Stop for small values of δ .

4.4. NUMERICAL EXPERIMENTS

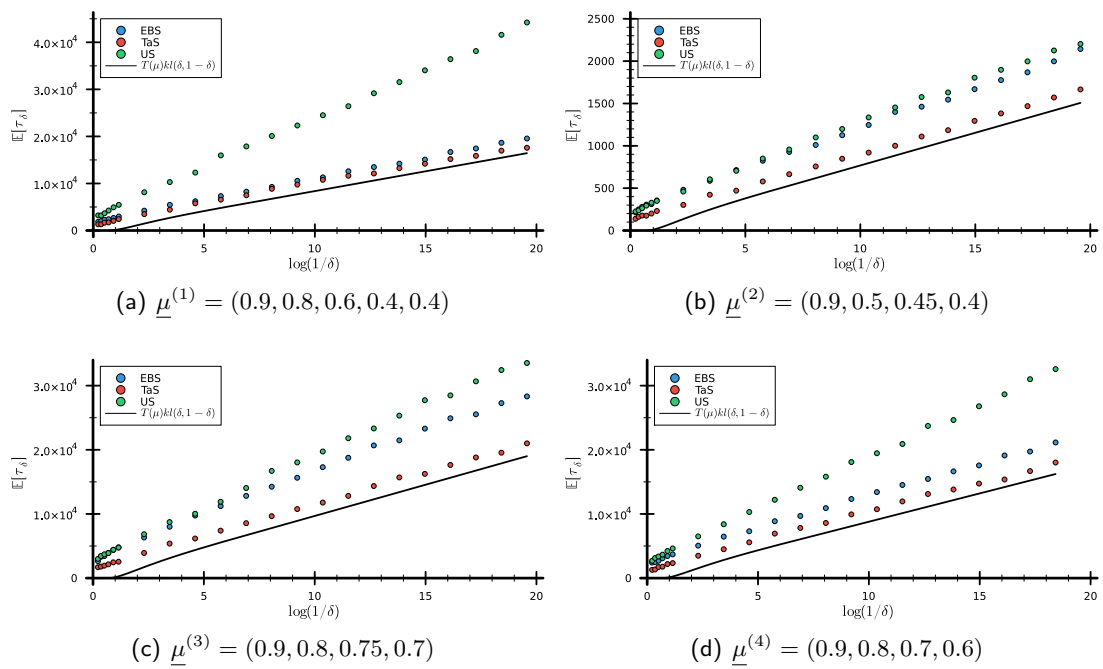


Figure 4.5: Empirical expected number of draws $\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]$ as a function of δ , averaged over 500 experiments.

4.5. Proof of the Non-Asymptotic Bounds of Theorem 4.5

The aim of this section is to prove Theorem 4.5. Let $\gamma \in (0, 1)$ and $\underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}^{[0,1]}$. We assume, without loss of generality, that $a^*(\underline{\mu}) = 1$. We also write, for simplicity, $\underline{\Delta} = \underline{\Delta}(\underline{\mu})$, $\underline{w} = \underline{w}(\underline{\mu})$ and $T = T(\underline{\mu})$.

We recall that the confidence regions are defined in (4.6) by

$$\forall t \geq K, \quad \mathcal{CR}_{\underline{\mu}}(t) = \prod_{a \in [K]} [\hat{\mu}_a(t) \pm \ell_a(t)],$$

$$\text{where } \forall a \in [K], \quad \ell_a(t) \stackrel{\text{def}}{=} C_{\frac{\gamma}{K}}(N_a(t)) = 2\sqrt{\frac{\log\left(\frac{4KN_a(t)}{\gamma}\right)}{N_a(t)}}.$$

Let \mathcal{E} denotes the event that $\underline{\mu}$ belongs to all confidence regions:

$$\mathcal{E} = \bigcap_{t=K}^{\tau_\delta} \left\{ \underline{\mu} \in \mathcal{CR}_{\underline{\mu}}(t) \right\}.$$

By the choice of the confidence regions and Lemma 4.2, we have

$$\mathbb{P}_{\underline{\mu}}(\mathcal{E}) \geq 1 - \gamma.$$

Furthermore, the design of Exploration-Biased-Sampling ensures that, under event \mathcal{E} , all arms are observed with some minimal linear rate, specified by Lemma 4.7 and proved in Section 4.6.2.

Lemma 4.7. *On event \mathcal{E} , one has:*

$$\forall t \in \mathbb{N}^*, \quad \min_{a \in [K]} N_a(t) \geq tw_{\min} - K. \quad (4.12)$$

Inequality (4.12) implies the more handy bound:

$$\forall t \geq \frac{2K}{w_{\min}}, \quad \min_{a \in [K]} N_a(t) \geq \frac{tw_{\min}}{2}. \quad (4.13)$$

Proof Outline. The proof is organized in 3 steps:

1. We first show that, on event \mathcal{E} , the optimal vector \underline{w} and the sampling frequency vector $\underline{N}(t)/t$ are very close for any $t \geq T_1$, where T_1 is a (problem-dependent) constant. To do so, we will make use of the regularity results of Section 3.5 and the fact that the confidence regions shrink with time.
2. Then, we control the event $(\tau_\delta > t) \cap \mathcal{E}$ for $t > T \log \frac{1}{\delta}$ by another event for which we can easily bound the probability using Hoeffding's inequality. This inclusion relies once again on the regularity results of Section 3.5 and on conditions on δ , in particular, we will require to have $T \log \frac{1}{\delta} \geq T_1$ with T_1 obtained at step 1.
3. Finally, we derive the two bounds of the theorem from Hoeffding's inequality and elementary calculations.

The proof uses some technical lemmas introduced and shown in Section 4.6.2.

4.5.1. Step 1: Controlling the Difference between Vectors \underline{w} and $\underline{N}(t)/t$

In this step we assume that event \mathcal{E} occurs.

Let $t \geq \frac{2K}{w_{\min}}$. Equation (4.13) implies that

$$\forall a \in [K], \quad \ell_a(t) = 2\sqrt{\frac{\log\left(\frac{4N_a(t)K}{\gamma}\right)}{N_a(t)}} \leq L(t) \stackrel{\text{def}}{=} \sqrt{8\frac{\log\left(\frac{4tK}{\gamma}\right)}{tw_{\min}}}.$$

$L(t)$ is an arm-independent bound on the half-length of the confidence interval of each μ_a . In other words, $\|\tilde{\underline{\mu}}(t) - \underline{\mu}\|_{\infty} \leq L(t)$ as we are on event \mathcal{E} . Note that $L(t)$ is deterministic and goes to 0 as t goes to $+\infty$. This control of $\|\tilde{\underline{\mu}}(t) - \underline{\mu}\|_{\infty} \leq L(t)$ together with Theorem 3.9 allows to control the difference between \underline{w} and $\tilde{\underline{w}}(t)$ for t large enough, as the following Lemma claims.

Lemma 4.8. *Let*

$$T_0 = \max\left(\frac{224^2}{\Delta_{\min}^2 w_{\min}} \log\left(\frac{2 \times 224^2 eK}{\Delta_{\min}^2 w_{\min} \gamma}\right), \frac{2K}{w_{\min}}\right). \quad (4.14)$$

Then for every $t \geq T_0$, one has, introducing $\varepsilon_t = \frac{80L(t)}{\Delta_{\min}}$,

$$\forall a \in [K], \quad w_a(1 - \varepsilon_t) \leq \tilde{w}_a(t) \leq w_a(1 + \varepsilon_t). \quad (4.15)$$

Proof. Let $t \geq \frac{2K}{w_{\min}}$ and assume that t is such that $4L(t) < \Delta_{\min}$. On event \mathcal{E} , one has

$$\underline{\mu} \in \mathcal{CR}_{\underline{\mu}}(t) = \prod_{a \in [K]} [\mu_a^-(t), \bar{\mu}_a(t)],$$

hence for any $a \neq 1$:

$$\mu_1^-(t) - \bar{\mu}_a(t) \geq \mu_1 - 2L(t) - (\mu_a + 2L(t)) \geq \Delta_a - 4L(t) > 0,$$

so that the confidence interval for μ_1 is strictly above all other confidence intervals. Hence $\tilde{\underline{\mu}}(t)$ has a unique optimal arm which is arm 1.

For each arm $a \neq 1$, define $\tilde{\Delta}_a(t) \stackrel{\text{def}}{=} \Delta_a(\tilde{\underline{\mu}}(t)) = \tilde{\mu}_1(t) - \tilde{\mu}_a(t)$. Then

$$\begin{aligned} \tilde{\Delta}_a(t)^2 &\leq (\Delta_a + 2L(t))^2 = \Delta_a^2 \left(1 + \frac{4L(t)}{\Delta_a} + \frac{4L(t)^2}{\Delta_a^2}\right) \leq \Delta_a^2 \left(1 + \frac{8L(t)}{\Delta_{\min}}\right) \\ \text{and } \tilde{\Delta}_a(t)^2 &\geq (\Delta_a - 2L(t))^2 = \Delta_a^2 \left(1 - \frac{4L(t)}{\Delta_a} + \frac{4L(t)^2}{\Delta_a^2}\right) \geq \Delta_a^2 \left(1 - \frac{8L(t)}{\Delta_{\min}}\right). \end{aligned}$$

If t is such that $\frac{8L(t)}{\Delta_{\min}} \leq \frac{1}{7}$ (this condition is stronger than $4L(t) < \Delta_{\min}$), we can apply Theorem 3.9 which gives

$$\forall a \in [K], \quad w_a(1 - \varepsilon_t) \leq \tilde{w}_a(t) \leq w_a(1 + \varepsilon_t).$$

It remains to understand when the condition $\frac{8L(t)}{\Delta_{\min}} \leq \frac{1}{7}$ holds. We have:

$$\frac{8L(t)}{\Delta_{\min}} \leq \frac{1}{7} \iff \frac{\log\left(\frac{4tK}{\gamma}\right)}{t} \leq \frac{\Delta_{\min}^2 w_{\min}}{(7 \times 8)^2 \times 8} = \frac{\Delta_{\min}^2 w_{\min}}{2 \times 112^2}$$

and this inequality is satisfied, by Lemma 4.15, for

$$t \geq \frac{224^2}{\Delta_{\min}^2 w_{\min}} \log \left(\frac{2 \times 224^2 eK}{\Delta_{\min}^2 w_{\min} \gamma} \right).$$

Combining with the initial condition $t \geq \frac{2K}{w_{\min}}$ leads to the definition of T_0 . \square

As the C-tracking procedure ensures that each $\frac{N_a(t)}{t}$ is roughly the Cesàro sum of the $(\tilde{w}_a(s))_{0 \leq s \leq t-1}$ (see Lemma 4.14), and as $\varepsilon_t \xrightarrow{t \rightarrow +\infty} 0$, we are able to control the difference between \underline{w} and $\frac{N(t)}{t}$ after a deterministic time T_1 .

Lemma 4.9. Fix $\eta \in (0, 1)$ and let

$$T_1 \stackrel{\text{def}}{=} \frac{\max(640^2, 8K)}{\eta^2 \Delta_{\min}^2 w_{\min}^2} \log \left(\frac{2 \times 640^2 eK}{\eta^2 \Delta_{\min}^2 w_{\min} \gamma} \right). \quad (4.16)$$

Then for any $t \geq T_1$ one has:

$$\forall a \in [K], \quad w_a(1 - \eta) \leq \frac{N_a(t)}{t} \leq w_a(1 + \eta).$$

Proof. Let T_0 be defined by Equation (4.14). Let $t > T_0$ and $a \in [K]$. Equation (4.15) of Lemma 4.8 gives:

$$\left| \sum_{s=0}^{t-1} \tilde{w}_a(s) - tw_a \right| \leq \sum_{s=0}^{T_0-1} |\tilde{w}_a(s) - w_a| + \sum_{s=T_0}^{t-1} |\tilde{w}_a(s) - w_a| \leq T_0 + w_a \sum_{s=T_0}^{t-1} \varepsilon_s.$$

By definition of ε_t one has:

$$\sum_{s=T_0}^{t-1} \varepsilon_s = \frac{80\sqrt{8}}{\Delta_{\min} \sqrt{w_{\min}}} \sum_{s=T_0}^{t-1} \sqrt{\frac{\log\left(\frac{4sK}{\gamma}\right)}{s}} \leq \frac{80\sqrt{8} \sqrt{\log\left(\frac{4tK}{\gamma}\right)}}{\Delta_{\min} \sqrt{w_{\min}}} \sum_{s=T_0}^{t-1} \frac{1}{\sqrt{s}} \leq \frac{80\sqrt{8} \sqrt{t \log\left(\frac{4tK}{\gamma}\right)}}{\Delta_{\min} \sqrt{w_{\min}}},$$

so that we have, using Lemma 4.14:

$$\begin{aligned} \left| \frac{N_a(t)}{t} - w_a \right| &\leq \frac{1}{t} \left[\left| N_a(t) - \sum_{s=0}^{t-1} \tilde{w}_a(s) \right| + \left| \sum_{s=0}^{t-1} \tilde{w}_a(s) - tw_a \right| \right] \\ &\leq \frac{K + T_0}{t} + w_a \frac{80\sqrt{8} \sqrt{\log\left(\frac{4tK}{\gamma}\right)}}{\Delta_{\min} \sqrt{w_{\min} t}} \\ &\leq w_a \left(\frac{K + T_0}{tw_{\min}} + \frac{80\sqrt{8} \sqrt{\log\left(\frac{4tK}{\gamma}\right)}}{\Delta_{\min} \sqrt{w_{\min} t}} \right). \end{aligned}$$

Thus the conclusion of the lemma holds when:

$$\max \left(\frac{K + T_0}{tw_{\min}}, \frac{80\sqrt{8} \sqrt{\log(4tK/\gamma)}}{\Delta_{\min} \sqrt{w_{\min} t}} \right) \leq \frac{\eta}{2},$$

and this inequality is satisfied, using Lemma 4.15, when:

$$t \geq \max \left(\frac{2}{\eta} \frac{K + T_0}{w_{\min}}, \frac{640^2}{\eta^2 \Delta_{\min}^2 w_{\min}} \log \left(\frac{2 \times 640^2 eK}{\eta^2 \Delta_{\min}^2 w_{\min} \gamma} \right) \right).$$

The definition of T_0 implies

$$K + T_0 \leq \frac{4 \max(112^2, K)}{\Delta_{\min}^2 w_{\min}} \log \left(\frac{2 \times 224^2 e K}{\Delta_{\min}^2 w_{\min} \gamma} \right),$$

hence the inequality still holds for

$$t \geq \max \left(\frac{8 \max(112^2, K)}{\eta \Delta_{\min}^2 w_{\min}^2} \log \left(\frac{2 \times 224^2 e K}{\Delta_{\min}^2 w_{\min} \gamma} \right), \frac{640^2}{\eta^2 \Delta_{\min}^2 w_{\min}} \log \left(\frac{2 \times 640^2 e K}{\eta^2 \Delta_{\min}^2 w_{\min} \gamma} \right) \right),$$

and T_1 is greater than this lower bound. \square

4.5.2. Step 2: a Useful Inclusion of Events

We want to control the event $\{\tau_\delta > t\} \cap \mathcal{E}$ for $t > T \log \frac{1}{\delta}$. For δ small enough, we have the following inclusion of events.

Lemma 4.10. Fix $\eta \in (0, 0.15]$ and let δ be such that

$$T \log \frac{1}{\delta} \geq T_1, \quad (\text{C1})$$

where T_1 is defined by Equation (4.16), and

$$\log \frac{1}{\delta} > \frac{4}{\eta} \log \left(\frac{8eTR^{1/2}}{\eta} \right). \quad (\text{C2})$$

Then for all $C \in (0, 1]$:

$$\forall t \geq (1 + C) \frac{(1 + \eta)^2}{(1 - \eta)^2} T \log \frac{1}{\delta}, \quad \{\tau_\delta > t\} \cap \mathcal{E} \subseteq \left\{ \|\tilde{\underline{\mu}}(t) - \underline{\mu}\|_\infty \geq \frac{C}{T} \right\} \cap \mathcal{E}.$$

Remark. Latter, we will use this Lemma with $C = \frac{1}{\log^{\frac{1}{3}} \frac{1}{\delta}}$.

Proof. Assume in the following that $T \log \frac{1}{\delta} \geq T_1$ and let $t \geq T \log \frac{1}{\delta}$. By definition of T_1 and Lemma 4.9, one has

$$\max_{a \in [K]} \left| \frac{w_a - \frac{N_a(t)}{t}}{w_a} \right| \leq \eta.$$

Then using Proposition 3.10 and Equation (4.3):

$$\begin{aligned} \{\tau_\delta > t\} \cap \mathcal{E} &\subseteq \left\{ Z(t) = t g \left(\hat{\underline{\mu}}(t), \frac{N(t)}{t} \right) \leq \beta(t, \delta) \right\} \cap \mathcal{E} \\ &\subseteq \left\{ t \frac{(1 - \eta)^2}{1 + \eta} \left(g(\underline{\mu}, \underline{w}) - \frac{\|\tilde{\underline{\mu}}(t) - \underline{\mu}\|_\infty}{2} \right) \leq \beta(t, \delta) \right\} \cap \mathcal{E} \\ &\subseteq \left\{ \frac{\|\tilde{\underline{\mu}}(t) - \underline{\mu}\|_\infty}{2} \geq \frac{1}{T} - \frac{1 + \eta}{(1 - \eta)^2} \frac{\beta(t, \delta)}{t} \right\} \cap \mathcal{E}. \end{aligned}$$

Consider now

$$f(t) = \frac{1 + \eta}{(1 - \eta)^2} \frac{\beta(t, \delta)}{t} = \frac{1 + \eta}{(1 - \eta)^2} \frac{\log \left(\frac{Rt^\alpha}{\delta} \right)}{t}.$$

As $\alpha \leq 2$, one can check that f is decreasing on $(4, +\infty)$. Let us show that

$$\forall C \in (0, 1], \quad f\left((1+C)\frac{(1+\eta)^2}{(1-\eta)^2}T \log \frac{1}{\delta}\right) \leq \frac{1}{(1+C)T}. \quad (4.17)$$

Fix $C \in (0, 1]$. As $\alpha \leq 2$ and as $\eta \leq 0.15$ is such that $\frac{(1+\eta)^2}{(1-\eta)^2} \leq 2$, we have:

$$\begin{aligned} f\left((1+C)\frac{(1+\eta)^2}{(1-\eta)^2}T \log \frac{1}{\delta}\right) &\leq \frac{1+\eta}{(1-\eta)^2} \frac{\log\left(\frac{R(4T \log \frac{1}{\delta})^2}{\delta}\right)}{(1+C)\frac{(1+\eta)^2}{(1-\eta)^2}T \log \frac{1}{\delta}} \\ &\leq \frac{1}{(1+C)T} \frac{1}{1+\eta} \left(1 + 2 \frac{\log(4R^{1/2}T \log \frac{1}{\delta})}{\log \frac{1}{\delta}}\right), \end{aligned}$$

hence inequality (4.17) is satisfied if

$$\log(4R^{1/2}T \log \frac{1}{\delta}) \leq \frac{\eta}{2} \log \frac{1}{\delta},$$

which is the case, by Lemma 4.15, when:

$$\log \frac{1}{\delta} > \frac{4}{\eta} \log\left(\frac{8eTR^{1/2}}{\eta}\right).$$

Finally when inequality (4.17) holds we have for $t \geq (1+C)\frac{(1+\eta)^2}{(1-\eta)^2}T \log \frac{1}{\delta}$:

$$\{\tau_\delta > t\} \cap \mathcal{E} \subseteq \left\{ \|\tilde{\mu}(t) - \underline{\mu}\|_\infty \geq \frac{2}{T} - \frac{2}{(1+C)T} \right\} \cap \mathcal{E} \subseteq \left\{ \|\tilde{\mu}(t) - \underline{\mu}\|_\infty \geq \frac{C}{T} \right\} \cap \mathcal{E}$$

where we use $C \leq 1$ in the last inclusion. \square

4.5.3. Step 3: Bounding $\mathbb{P}_\mu(\tau_\delta > t \cap \mathcal{E})$ and $\mathbb{E}_\mu[\tau_\delta \mathbb{I}\{\mathcal{E}\}]$.

Fix $\eta \in (0, 1]$ and assume in the following that conditions (C1) and (C2) of Lemma 4.10 are satisfied with $\eta' = \frac{\eta}{7} \leq 0.15$. We set $\zeta = \frac{(1+\eta')^2}{(1-\eta')^2}$. Let $C \in (0, 1]$, $t > (1+C)\zeta T \log \frac{1}{\delta}$ and define

$$\mathcal{E}_t = \left\{ \|\tilde{\mu}(t) - \underline{\mu}\|_\infty \geq \frac{C}{T} \right\} \cap \mathcal{E}.$$

Lemmas 4.10 and 4.16 – a consequence of Hoeffding's inequality – (note that Condition (C1) ensures that $t \geq \frac{2K}{w_{\min}}$) give the bound:

$$\mathbb{P}_\mu(\tau_\delta > t \cap \mathcal{E}) \leq \mathbb{P}_\mu(\mathcal{E}_t) \leq 2Kt \exp\left(-\frac{tw_{\min}}{4T^2}C^2\right). \quad (4.18)$$

By taking $C = \frac{1}{\log^{\frac{1}{3}} \frac{1}{\delta}}$, we obtained so far that

$$\forall t > \left(1 + \frac{1}{\log^{\frac{1}{3}} \frac{1}{\delta}}\right) \zeta T \log \frac{1}{\delta}, \quad \mathbb{P}_\mu(\tau_\delta > t \cap \mathcal{E}) \leq 2Kt \exp\left(-\frac{tw_{\min}}{4T^2} \frac{1}{\log^{\frac{2}{3}} \frac{1}{\delta}}\right),$$

giving Bound (4.9) as long as

$$\left(1 + \frac{1}{\log^{\frac{1}{3}} \frac{1}{\delta}}\right) \zeta \leq 1 + \eta.$$

Note that $\zeta \leq 1 + 6\eta'$ as $\eta' \leq 0.15$ so that when

$$\frac{1}{\log^{\frac{1}{3}} \frac{1}{\delta}} \leq \frac{\eta'}{2} \iff \log \frac{1}{\delta} \geq \frac{8 \times 7^3}{\eta^3}, \quad (\text{C3})$$

the condition holds as

$$\left(1 + \frac{1}{\log^{\frac{1}{3}} \frac{1}{\delta}}\right) \zeta \leq \left(1 + \frac{\eta'}{2}\right) (1 + 6\eta') \leq 1 + 6.6\eta' \leq 1 + \eta.$$

It remains to focus on the bound of $\mathbb{E}_{\underline{\mu}}[\tau_{\delta} \mathbb{I}\{\mathcal{E}\}]$. Using Equation (4.18) we have:

$$\begin{aligned} \mathbb{E}_{\underline{\mu}}[\tau_{\delta} \mathbb{I}\{\mathcal{E}\}] &= \sum_{t=0}^{\lfloor (1+C)\zeta T \log \frac{1}{\delta} \rfloor} \mathbb{P}_{\underline{\mu}}(\tau_{\delta} > t \cap \mathcal{E}) + \sum_{t > (1+C)\zeta T \log \frac{1}{\delta}} \mathbb{P}_{\underline{\mu}}(\tau_{\delta} > t \cap \mathcal{E}) \\ &\leq (1+C)\zeta T \log \frac{1}{\delta} + 1 + 2K \sum_{t > (1+C)\zeta T \log \frac{1}{\delta}} t \exp\left(-\frac{tw_{\min}}{4T^2} C^2\right). \end{aligned}$$

By defining

$$S(C) = \sum_{t > C\zeta T \log \frac{1}{\delta}} t \exp\left(-\frac{tw_{\min}}{4T^2} C^2\right),$$

we obtain, with some technical calculations (see Section 4.6.2), the following bound.

Lemma 4.11. *One has*

$$S(C) \leq \frac{32T^4}{w_{\min}^2} \exp\left(-\frac{w_{\min}}{4T} C^2 \log \frac{1}{\delta}\right) \left(\frac{\log \frac{1}{\delta}}{C^2} + \frac{1}{C^4}\right).$$

Once again, taking $C = \frac{1}{\log^{\frac{1}{3}} \frac{1}{\delta}}$ leads to

$$S(C) \leq \frac{32T^4}{w_{\min}^2} \exp\left(-\frac{w_{\min}}{4T} \log^{\frac{1}{3}} \frac{1}{\delta}\right) \left(\log^{\frac{5}{3}} \frac{1}{\delta} + \log^{\frac{4}{3}} \frac{1}{\delta}\right) \leq \frac{64T^4}{w_{\min}^2} \exp\left(-\frac{w_{\min}}{4T} \log^{\frac{1}{3}} \frac{1}{\delta}\right) \log^2 \frac{1}{\delta},$$

thus

$$\mathbb{E}_{\underline{\mu}}[\tau_{\delta} \mathbb{I}\{\mathcal{E}\}] \leq \zeta \left(1 + \frac{1}{\log^{\frac{1}{3}} \frac{1}{\delta}}\right) T \log \frac{1}{\delta} + 1 + \frac{2^7 K T^4}{w_{\min}^2} \exp\left(-\frac{w_{\min}}{4T} \log^{\frac{1}{3}} \frac{1}{\delta}\right) \log^2 \frac{1}{\delta}.$$

Under Condition (C3) we get

$$\zeta \left(1 + \frac{1}{\log^{\frac{1}{3}} \frac{1}{\delta}}\right) T \log \frac{1}{\delta} + 1 \leq (1 + 6.6\eta') T \log \frac{1}{\delta} + 1 \leq (1 + \eta) T \log \frac{1}{\delta},$$

and obtain the Bound (4.10) claimed in the theorem.

Combining conditions (C1), (C2) and (C3) together, one can define δ_0 satisfying:

$$\log \frac{1}{\delta_0} \geq \frac{7^3 \times \max(2 \times 160^2, K)}{\eta^3 \Delta_{\min} w_{\min}^2} \log\left(\frac{7^2 \times 2 \times 640^2 e K R^{1/2}}{\eta^2 \Delta_{\min}^2 w_{\min} \gamma}\right),$$

with some simplifications allowed by Equation (3.20) of Proposition 3.4.

4.6. Technical Results

4.6.1. Proof of Lemma 4.4

We will prove the lemma using two supporting results. The following lemma gives a lower bound of $w_{\min}(\underline{\mu})$ in terms of the minimal gap $\Delta_{\min}(\underline{\mu})$, based on the study of optimization problem (4.2) in Chapter 3.

Lemma 4.12. *For all $\underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}^{[0,1]}$ with a unique optimal arm, one has*

$$w_{\min}(\underline{\mu}) \geq \frac{\Delta_{\min}(\underline{\mu})}{2K}.$$

Proof. Let $\underline{w} = \underline{w}(\underline{\mu})$, $w_{\min} = w_{\min}(\underline{\mu})$ and $\underline{\Delta} = \underline{\Delta}(\underline{\mu})$. We have, by Equation (3.10), recalling that the weight w_{\min} is the optimal weight of the worst arm(s) of $\underline{\mu}$, i.e., the arm(s) with highest gap(s),

$$w_{\min} = \frac{w_{\max}}{r\Delta_{\max} - 1}.$$

Applying inequalities (3.18) and (3.19), and using the fact that $\Delta_{\max} \leq 1$ as means belong to $[0, 1]$,

$$w_{\min} \geq \frac{1}{\sqrt{K-1} + 1} \cdot \frac{1}{\frac{\sqrt{K-1} + 1}{\Delta_{\min}} \Delta_{\max} - 1} \geq \frac{\Delta_{\min}}{(\sqrt{K-1} + 1)^2} \geq \frac{\Delta_{\min}}{2K}. \quad \square$$

The next lemma provides a lower bound on the minimal gap of the optimistic bandit $\tilde{\underline{\mu}}$ computed by Algorithm 13.

Lemma 4.13. *Let $\mathcal{CR} = \prod_{a \in [K]} [\mu_a^-, \mu_a^+]$ be a confidence region with $\mu_a^- < \mu_a^+$ for all $a \in [K]$, and assume that*

$$\max_{a \in [K]} \mu_a^- = \max LB > \min UB = \min_{a \in [K]} \mu_a^+.$$

Then, if $(\tilde{\underline{\mu}}, \underline{v}) \leftarrow \text{Exploration-Biased-Weights}(\mathcal{CR})$, we get

$$\Delta_{\min}(\tilde{\underline{\mu}}) \geq \min_{a \in [K]} \mu_a^+ - \mu_a^-.$$

Proof. We proceed by contradiction: let us assume that $\tilde{\underline{\mu}}$ is such that

$$\Delta_{\min}(\tilde{\underline{\mu}}) < \min_{a \in [K]} \mu_a^+ - \mu_a^-.$$

By the two hypotheses and the algorithm's procedure, it is clear that $\tilde{\underline{\mu}}$ has a unique best arm. Without loss of generality let us arrange the arms so that $\tilde{\mu}_1 > \tilde{\mu}_2 \geq \tilde{\mu}_3 \geq \dots \geq \tilde{\mu}_K$. Note that $\Delta_{\min}(\tilde{\underline{\mu}}) = \tilde{\mu}_1 - \tilde{\mu}_2$.

As 1 is the best arm, once again the algorithm's procedure ensures that $\tilde{\mu}_1 = \mu_1^+$. In addition, our assumption implies $\Delta_{\min}(\tilde{\underline{\mu}}) < \mu_1^+ - \mu_1^-$, giving $\tilde{\mu}_2 > \mu_1^-$. Recall that $\tilde{\mu}_2 = \max(\mu_2^-, \min UB)$, so we split our analysis into the two possible cases:

4.6. TECHNICAL RESULTS

- if $\tilde{\mu}_2 = \mu_2^-$, then we cannot have $\mu_2^+ \leq \mu_1^+ = \tilde{\mu}_1$ otherwise $\Delta_{\min}(\tilde{\mu}) > \mu_2^+ - \mu_2^-$, which is impossible.
Then $\mu_2^+ > \mu_1^+$. By defining $\underline{\rho} = (\tilde{\mu}_2, \mu_2^+, \tilde{\mu}_3, \dots, \tilde{\mu}_K)$, one has $\underline{\rho} \in \mathcal{CR}$ and $w_{\min}(\underline{\rho}) > w_{\min}(\tilde{\mu})$ by Lemma 3.6. Thus $\tilde{\mu}$ cannot maximize w_{\min} over \mathcal{CR} which is in contradiction with Proposition 4.1.
- if $\tilde{\mu}_2 = \text{minUB}$, then $\tilde{\mu}_2 = \tilde{\mu}_3 = \dots = \tilde{\mu}_K$ and thus all confidence intervals share a common point equal to $\tilde{\mu}_2$ (recall that $\tilde{\mu}_2 \in [\mu_1^-, \mu_1^+]$), which is a contradiction with $\text{maxLB} > \text{minUB}$. \square

We can now prove Lemma 4.4.

Proof of Lemma 4.4. Let $t \geq 0$. We want to lower bound $\tilde{w}_{\min}(t)$.

- If at time t one has $\tilde{w}(t) = (\frac{1}{K}, \dots, \frac{1}{K})$, that is, all confidence intervals share a common value, then $\tilde{w}_{\min}(t) = \frac{1}{K}$.
- Otherwise, by the construction of Algorithms 13 and 14 we know that $t \geq K$ and the confidence region $\mathcal{CR}(t)$ is such that at least two confidence intervals are separated. In that case, the optimistic bandit $\tilde{\mu}(t)$ has a unique optimal arm, hence by Lemma 4.12

$$\tilde{w}_{\min}(t) \geq \frac{\tilde{\Delta}_{\min}(t)}{2K}.$$

Applying Lemma 4.13 leads to (recalling that $N_a(t) \geq 1$ as all arms are pulled once at the beginning)

$$\tilde{\Delta}_{\min}^{(t)} \geq \min_{a \in [K]} 2\ell_a(t) \geq 4 \min_{a \in [K]} \sqrt{\frac{\log\left(\frac{4N_a(t)K}{\gamma}\right)}{N_a(t)}} \geq 4\sqrt{\frac{\log\left(\frac{4K}{\gamma}\right)}{t}} \geq 4\sqrt{\frac{\log 8}{t}} \geq \frac{4}{\sqrt{t}}.$$

Combining the two last equations, we obtain

$$\tilde{w}_{\min}(t) \geq \frac{2}{K} \frac{1}{\sqrt{t}}.$$

Hence, in both cases, we get:

$$\forall t \geq 0, \quad \tilde{w}_{\min}(t) \geq \min\left(\frac{2}{K} \frac{1}{\sqrt{t}}, \frac{1}{K}\right) \geq \frac{1}{K} \frac{1}{\sqrt{t}}.$$

Using Lemma 4.14, this implies that for all $t \geq 0$ and $a \in [K]$

$$N_a(t) \geq \sum_{s=0}^{t-1} \tilde{w}_a(s) - (K-1) \geq \sum_{s=2}^{t-1} \tilde{w}_{\min}(s) - K \geq \frac{1}{K} \sum_{s=2}^{t-1} \frac{1}{\sqrt{s}} - K \geq \frac{1}{K} \int_1^t \frac{ds}{\sqrt{s}} - K \geq \frac{2}{K} \sqrt{t} - K.$$

\square

4.6.2. Technical Details for the Proof of Theorem 4.5

Proof of Lemma 4.7. We will use the following deterministic Lemma.

Lemma 4.14. *One has:*

$$\forall t > 0, \forall a \in [K], \quad \left| N_a(t) - \sum_{s=0}^{t-1} \tilde{w}_a(s) \right| \leq K - 1.$$

Proof. Apply [Garivier and Kaufmann \(2016, Lemma 15\)](#) with $p(s) = \tilde{w}(s)$. □

Proof of Lemma 4.7. The claim is true for $t \leq K$. Otherwise, fix $t > K$ and $a \in [K]$. For all $0 \leq s \leq K-1$, one has $\tilde{w}_a(s) = \frac{1}{K}$ by convention⁶, and thus $\tilde{w}_a(s) \geq w_{\min}$. As $\underline{\mu} \in \mathcal{CR}_{\underline{\mu}}(s)$ for all $K \leq s \leq t$ on event \mathcal{E} , Proposition 4.1 ensures that

$$\forall K \leq s \leq t, \quad \tilde{w}_a(s) \geq \tilde{w}_{\min}(s) = \max_{\rho \in \mathcal{CR}_{\underline{\mu}}(s)} w_{\min}(\rho) \geq w_{\min}.$$

Hence, by Lemma 4.14

$$N_a(t) \geq \sum_{s=0}^{t-1} \tilde{w}_a(s) - (K-1) \geq tw_{\min} - (K-1) \geq tw_{\min} - K. \quad \square$$

A Technical Lemma. The following result is a direct consequence of [Garivier and Kaufmann \(2016, Lemma 18\)](#).

Lemma 4.15. For any $c_1, c_2 > 0$,

$$x = \frac{2}{c_1} \log\left(\frac{c_2 e}{c_1}\right)$$

is such that $c_1 x \geq \log(c_2 x)$.

Deviation Bound. We prove the following simple consequence of Hoeffding's inequality.

Lemma 4.16. For all $t \geq \frac{2K}{w_{\min}}$ and $x > 0$, one has

$$\mathbb{P}_{\underline{\mu}}\left(\max_{a \in [K]} |\hat{\mu}_a(t) - \mu_a| > x \cap \mathcal{E}\right) \leq 2Kt \exp\left(-\frac{tw_{\min}}{4} x^2\right).$$

Proof. Fix $t \geq \frac{2K}{w_{\min}}$ and $x > 0$. With $T = \frac{tw_{\min}}{2}$, we get, using Equations (4.13) and (4.4),

$$\begin{aligned} \forall a \in [K], \quad \mathbb{P}_{\underline{\mu}}\left(|\hat{\mu}_a(t) - \mu_a| > x \cap \mathcal{E}\right) &= \sum_{s=T}^t \mathbb{P}_{\underline{\mu}}\left(|\hat{\mu}_a(t) - \mu_a| > x \cap \mathcal{E} \cap N_a(t) = s\right) \\ &\leq \sum_{s=T}^t \mathbb{P}_{\underline{\mu}}\left(|\hat{\mu}_{a,s} - \mu_a| > x\right) \\ &\leq \sum_{s=T}^t 2 \exp\left(-\frac{s}{2} x^2\right) \\ &\leq 2t \exp\left(-\frac{T}{2} x^2\right), \end{aligned}$$

where the second inequality uses Hoeffding's inequality (2.4). The result follows by union bound. □

Proof of Lemma 4.11. We have

$$S(C) = \sum_{t > (1+C)\zeta T \log \frac{1}{\delta}} t \exp\left(-\frac{tw_{\min}}{4T^2} C^2\right) = \sum_{t > B} f(t)$$

⁶As all arms are drawn once during the K first rounds, the only request is $\sum_{s=0}^{K-1} \tilde{w}_a(s) = 1$.

4.6. TECHNICAL RESULTS

where $f : t \mapsto t \exp(-At)$, $A = \frac{w_{\min}}{4T^2} C^2$ and $B = (1 + C)\zeta T \log \frac{1}{\delta}$. f is increasing until $1/A$ and then decreasing. Let $n_0 = \lfloor \frac{1}{A} \rfloor$. We will show that $S(C) \leq 2 \int_B^{+\infty} f(t) dt$.

- If $B > n_0$ then f is decreasing on $[B, +\infty[$ and one has $S(C) \leq \int_B^{+\infty} f(t) dt$.
- Otherwise, one has:

$$\begin{aligned} S(C) &= \sum_{t=\lceil B \rceil}^{n_0-1} f(t) + f(n_0) + f(n_0 + 1) + \sum_{t>n_0+1} f(t) \\ &\leq \sum_{t=\lceil B \rceil}^{n_0-1} \int_t^{t+1} f(t) dt + f(n_0) + f(n_0 + 1) + \sum_{t>n_0+1} \int_{t-1}^t f(t) dt \\ &\leq \int_{\lceil B \rceil}^{+\infty} f(t) dt + f(n_0) + f(n_0 + 1) \end{aligned}$$

where in the second inequality, we use the increasing (respectively the decreasing) of f on $[B, n_0]$ (respectively on $[n_0 + 1, +\infty[$). The result will be true if

$$f(n_0) + f(n_0 + 1) \leq \int_B^{+\infty} f(t) dt.$$

We have:

$$\begin{aligned} f(n_0) + f(n_0 + 1) &= \left\lfloor \frac{1}{A} \right\rfloor e^{-A \lfloor \frac{1}{A} \rfloor} + \left\lceil \frac{1}{A} \right\rceil e^{-A \lceil \frac{1}{A} \rceil} \\ &\leq \left(\left\lfloor \frac{1}{A} \right\rfloor + \left\lceil \frac{1}{A} \right\rceil \right) e^{-A \lfloor \frac{1}{A} \rfloor} \\ &\leq \left(\left\lfloor \frac{1}{A} \right\rfloor \frac{1}{A} + \frac{1}{A^2} \right) e^{-A \lfloor \frac{1}{A} \rfloor} \quad \text{as } A < \frac{1}{2} \\ &= \int_{\lfloor \frac{1}{A} \rfloor}^{+\infty} f(t) dt \leq \int_B^{+\infty} f(t) dt \quad \text{as } B \leq \left\lfloor \frac{1}{A} \right\rfloor = n_0. \end{aligned}$$

where in the last inequality, we used the simple calculation

$$\int_Y^{+\infty} t \exp(-tX) dt = \exp(-YX) \left(\frac{Y}{X} + \frac{1}{X^2} \right)$$

for $X, Y > 0$.

In both cases, we have:

$$S(C) \leq 2 \int_{(1+C)\zeta T \log \frac{1}{\delta}}^{\infty} t \exp\left(-\frac{tw_{\min}}{4T^2} C^2\right) dt$$

and using the same calculation as before

$$S(C) \leq 2 \exp\left(-\frac{\zeta w_{\min}}{4T} (1 + C) C^2 \log \frac{1}{\delta}\right) \left(\frac{4(1 + C)\zeta T^3 \log \frac{1}{\delta}}{w_{\min} C^2} + \frac{16T^4}{w_{\min}^2 C^4} \right).$$

Bounding $C \in (0, 1]$ and $\zeta \in [1, 2]$ (remind that $\zeta \leq 1 + 6\eta'$):

$$\begin{aligned} S(C) &\leq 2 \exp\left(-\frac{w_{\min}}{4T} C^2 \log \frac{1}{\delta}\right) \left(\frac{16T^3 \log \frac{1}{\delta}}{w_{\min} C^2} + \frac{16T^4}{w_{\min}^2 C^4} \right) \\ &\leq \frac{32T^4}{w_{\min}^2} \exp\left(-\frac{w_{\min}}{4T} C^2 \log \frac{1}{\delta}\right) \left(\frac{\log \frac{1}{\delta}}{C^2} + \frac{1}{C^4} \right). \end{aligned}$$

□

4.6.3. Almost Sure Asymptotic Bound

Almost Sure Asymptotic Bound. Exploration-Biased-Sampling satisfies the following almost sure asymptotic bound.

Theorem 4.17. Let $\gamma \in (0, 1)$ and $\alpha \in [1, e/2]$. Algorithm Exploration-Biased-Sampling with threshold (4.8) satisfies

$$\forall \underline{\mu} \text{ in } \mathcal{D}_{\mathcal{N}_1}^{[0,1]}, \quad \limsup_{\delta \rightarrow 0} \frac{\tau_\delta}{\log \frac{1}{\delta}} \leq \alpha T(\underline{\mu}) \quad \mathbb{P}_{\underline{\mu}}\text{-a.s.}$$

This result was obtained for the Track-and-Stop algorithm by [Garivier and Kaufmann \(2016, Proposition 13\)](#). The adaptation to Exploration-Biased-Sampling is straightforward, as soon as we prove the almost sure convergence of the empirical means and frequencies of pulls.

Proposition 4.18. For all $\gamma \in (0, 1)$, the sampling rule of Exploration-Biased-Sampling satisfies, for all $\underline{\mu} \in \mathcal{D}_{\mathcal{N}_1}^{[0,1]}$,

$$\lim_{t \rightarrow +\infty} \hat{\underline{\mu}}(t) = \underline{\mu} \quad \mathbb{P}_{\underline{\mu}}\text{-a.s.} \quad \text{and} \quad \lim_{t \rightarrow +\infty} \frac{N(t)}{t} = w(\underline{\mu}) \quad \mathbb{P}_{\underline{\mu}}\text{-a.s.}$$

Proof. The first convergence is a simple application of the law of large numbers, given that the number of pulls of all arms diverges thanks to Lemma 4.4. Note also that, as

$$\forall a \in [K], \quad |\tilde{\mu}_a(t) - \hat{\mu}_a(t)| \leq C_{\frac{\gamma}{K}}(N_a(t)) = 2\sqrt{\frac{\log\left(\frac{4N_a(t)K}{\gamma}\right)}{N_a(t)}} \xrightarrow{t \rightarrow +\infty} 0,$$

we have that $\tilde{\underline{\mu}}(t)$ and $\hat{\underline{\mu}}(t)$ get the same limit:

$$\lim_{t \rightarrow +\infty} \tilde{\underline{\mu}}(t) = \underline{\mu} \quad \mathbb{P}_{\underline{\mu}}\text{-a.s.}$$

Thus, by continuity of the optimal weights \underline{w} at $\underline{\mu}$:

$$\lim_{t \rightarrow +\infty} \tilde{\underline{w}}(t) = \underline{w}(\underline{\mu}) \quad \mathbb{P}_{\underline{\mu}}\text{-a.s.}$$

Finally, for all $t \in \mathbb{N}^*$ and $a \in [K]$, we have:

$$\begin{aligned} \left| \frac{N_a(t)}{t} - w_a(\underline{\mu}) \right| &\leq \frac{1}{t} \left| N_a(t) - \sum_{s=0}^{t-1} \tilde{w}_a(s) \right| + \left| \frac{1}{t} \sum_{s=0}^{t-1} (\tilde{w}_a(s) - w_a(\underline{\mu})) \right| \\ &\leq \frac{K-1}{t} + \left| \frac{1}{t} \sum_{s=0}^{t-1} (\tilde{w}_a(s) - w_a(\underline{\mu})) \right| \\ &\xrightarrow{t \rightarrow +\infty} 0 \quad \mathbb{P}_{\underline{\mu}}\text{-a.s.}, \end{aligned}$$

using Lemma 4.14 for the second inequality, and Cesàro Lemma for the convergence. \square

Asymptotic optimality. We now prove Theorem 4.6 on the asymptotic optimality of our strategy. Once again this is a direct adaptation of [Garivier and Kaufmann \(2016, Theorem 14\)](#). One can apply the original proof as long as the two lemmas shown in this section are satisfied.

4.7. CONCLUSION

Proof. Let $T \geq 0$. We get, for all $t \geq \sqrt{T} = h(T)^2$ and $a \in [K]$,

$$\begin{aligned}
 \left| \frac{N_a(t)}{t} - w_a(\underline{\mu}) \right| &\leq \frac{1}{t} \left| N_a(t) - \sum_{s=0}^{t-1} \tilde{w}_a(s) \right| + \left| \frac{1}{t} \sum_{s=0}^{t-1} (\tilde{w}_a(s) - w_a(\underline{\mu})) \right| \\
 &\leq \frac{K-1}{t} + \frac{h(T)}{t} + \left| \frac{1}{t} \sum_{s=h(T)}^{t-1} (\tilde{w}_a(s) - w_a(\underline{\mu})) \right| && \text{by Lemma 4.14} \\
 &\leq \frac{K-1}{T^{1/2}} + \frac{1}{T^{1/4}} + \varepsilon && \text{by definition of } \mathcal{E}_T \\
 &\leq \frac{K}{T^{1/4}} + \varepsilon, \leq 3\varepsilon,
 \end{aligned}$$

where we used Lemma 4.14 in the second inequality, and the definition of \mathcal{E}_T for the third inequality. This gives the result whenever $T \geq T_\varepsilon \stackrel{\text{def}}{=} \left(\frac{K}{2\varepsilon}\right)^4$. \square

4.7. Conclusion

We introduced **Exploration-Biased-Sampling**, a new strategy for the problem of best arm identification with fixed confidence. In addition to asymptotic optimal results, we proved non-asymptotic bounds for this strategy in the case of (sub-)Gaussian bandits. Those finite risk bounds were made possible by a new analysis of the sample complexity optimization problem presented in Chapter 3, and by the design of our strategy which tackles some shortcomings of **Track-and-Stop**: the procedure ensures exploration in an unforced way and stabilizes the sampling strategy, observing uniformly before having a high certainty that one arm is better than another.

Improving the guarantees of Exploration-Biased-Sampling. Although our new strategy enjoys interesting stability properties, it is always a bit worse than **Track-and-Stop** in terms of pure performance (sample complexity). This is a consequence of the biased tracking of the strategy which implies a slower convergence of the pulling frequencies than **Track-and-Stop**, but was necessary to obtain our non-asymptotic bound. A future direction of research might consist in carefully modifying the exploration mechanism of **Exploration-Biased-Sampling** in order to obtain a more efficient strategy that still benefits from a finite risk analysis. The non-asymptotic guarantees of Theorem 4.5 also come with a few limitations, which leaves room for improvement:

- The finite risk bound is given on an event \mathcal{E} of high probability (depending on the external parameter γ for the confidence regions). What happens when \mathcal{E} does not occur?
- The analysis is valid for values of confidence δ that are in practice extremely small. Can we analyze it for more moderate values of δ ?

Non-Gaussian models. It would be interesting but it remains out of reach to generalize this approach to non-Gaussian models: this requires extending our results on the sample-complexity optimization problem, a technically challenging task for which the simple and clean arguments developed here are likely to be replaced by much more involved derivations if this is possible. In addition, it will be necessary to modify the confidence intervals on the arm means in a way that ensures exploration. Another direction of improvement will be to investigate if similar analysis and strategies are possible for the problem of ε -best arm identification.

CHAPTER 5

A Non-Parametric Theory of Fixed-Budget Best-Arm Identification

In this chapter, we study non-parametric generalizations of existing bounds in fixed-budget best-arm identification. We consider general models \mathcal{D} for distributions over the arms; an overarching example is the model $\mathcal{D} = \mathcal{P}[0, 1]$ of all probability distributions over $[0, 1]$. We propose upper bounds on the average log-probability of misidentifying the optimal arm based on information-theoretic quantities that we name $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$ and $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu)$ and that correspond to infima over Kullback-Leibler divergences between some distributions in \mathcal{D} and a given distribution ν . This is made possible by a refined analysis of the Successive-Rejects strategy of [Audibert et al. \(2010\)](#). We finally provide lower bounds on the same average log-probability, also in terms of the same new information-theoretic quantities; these lower bounds are larger when the (natural) assumptions on the considered strategies are stronger. All these new upper and lower bounds generalize existing bounds based, e.g., on gaps between distributions. The content of this chapter is extracted from the conference paper



A. Barrier, A. Garivier, and G. Stoltz. On Best-Arm Identification with a Fixed Budget in Non-Parametric Multi-Armed Bandits. In *Proceedings of the 34th International Conference on Algorithmic Learning Theory*, volume 201 of *Proceedings of Machine Learning Research*, pages 136–181. PMLR, 2023

Contents

1	Introduction	124
2	Overview of the Results and more Extended Literature Review	126
1	The Key new Quantities: $\mathcal{L}_{\text{inf}}^{\leq}$ and $\mathcal{L}_{\text{inf}}^{\geq}$, as well as $\mathcal{L}_{\text{inf}}^{\leq}$ and $\mathcal{L}_{\text{inf}}^{\geq}$	126
2	Overview of the Results	127
3	Re-Derivation of Existing Bounds	128
4	Discussion of the (Lack of) Optimality of the new Bounds Exhibited	130
3	Upper Bound for the Successive-Rejects Strategy, with an Improved Analysis	131
1	General Analysis	131
2	On Links between Φ and the Quantities $\mathcal{L}_{\text{inf}}^{\leq}$, $\mathcal{L}_{\text{inf}}^{\leq}$, $\mathcal{L}_{\text{inf}}^{\geq}$ and $\mathcal{L}_{\text{inf}}^{\geq}$	137
4	Lower Bounds	137
1	Common Restriction: Consistence	138
2	A Lower Bound Revisiting and Extending the one by Audibert et al. (2010)	139
3	A Larger Lower Bound, for a more Restrictive Class of Strategies	141
4	A General Lower Bound, Valid for any Strategy	143

5	Technical Details	144
1	Properties of the $\mathcal{L}_{\text{inf}}^<$, $\mathcal{L}_{\text{inf}}^{\leq}$, $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^{\geq}$	144
2	Reminder: the Cramér-Chernoff Bound	147
3	Proofs and Details for Section 5.4.2: Rewriting of Φ as \mathcal{L}	148
4	Proof of the Normality of the Models $\mathcal{P}[0, 1]$ and \mathcal{D}_{exp}	154
6	Additional Comments for the Literature Review	155
1	The Minimax Lower Bound of Carpentier and Locatelli (2016)	155
2	The Bretagnolle-Huber Technique by Kaufmann et al. (2016)	157
7	Conclusion	160

5.1. Introduction

We consider a class \mathcal{D} of distributions over \mathbb{R} with finite first moments, which we refer to as the model \mathcal{D} . A K -armed bandit problem in \mathcal{D} is a K -tuple $\underline{\nu} = (\nu_1, \dots, \nu_K)$ of distributions in \mathcal{D} . We denote by (μ_1, \dots, μ_K) the K -tuple of their expectations. An agent sequentially interacts with $\underline{\nu}$: at each step $t \geq 1$, she selects an arm A_t and receives a reward Y_t drawn from the distribution ν_{A_t} . This is the only feedback that she obtains.

While regret minimization has been vastly studied (see [Lattimore and Szepesvári, 2020](#)), another relevant objective is *best-arm identification*, that is, identifying the distribution with highest expectation. In the fixed-confidence setting, this identification is performed under the constraint that a given confidence level $1 - \delta$ is respected, while minimizing the expected number of pulls of the arms (the expected sample complexity). This setting is fairly well understood (see [Lattimore and Szepesvári, 2020](#), Chapter 33 for a review). A turning point in this literature was achieved by [Garivier and Kaufmann \(2016\)](#), who provided matching upper and lower bounds on the expected number of pulls of the arms in the case of canonical one-parameter exponential families. Since then, improvements have been made in several directions, including for example non-asymptotic bounds ([Degenne et al., 2019](#)) and the problem of ε -best-arm identification ([Garivier and Kaufmann, 2021](#)). The first generalization to non-parametric models in this fixed-confidence setting was achieved by [Jourdan and Degenne \(2023\)](#), who worked in a concurrent and independent manner from us. Their upper and lower bounds differ by a multiplicative factor of 2 (only).

Fixed-budget best-arm identification. The *fixed-budget setting* seems to be much less understood. Therein, the total number T of pulls of the arms is fixed. After these T pulls, a strategy must issue a recommendation \hat{a}_T . Assuming that $\underline{\nu}$ contains a unique optimal distribution ν^* of index $a^*(\underline{\nu})$, one aims at minimizing $\mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu}))$. We are interested in (upper and lower) bounds that hold for all problems $\underline{\nu}$ in \mathcal{D} , possibly under the restriction that they only contain a unique optimal arm. It may be straightforwardly seen that the probability of error can decay exponentially fast—for instance, by uniformly exploring the arms (pulling each of them about T/K times) and recommending the one with the largest empirical average. This is why the literature (see, for instance, [Audibert et al., 2010](#) and [Lattimore and Szepesvári, 2020](#), Chapter 33) focuses on lower and upper bound functions $\ell \leq U < 0$ of the typical form: *for all bandit problems $\underline{\nu}$ in \mathcal{D} , with a unique optimal arm,*

$$\ell(\underline{\nu}) \leq \liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq \limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq U(\underline{\nu}) < 0,$$

$$\text{or, put differently,} \quad \exp(\ell(\underline{\nu}) T(1 + o(1))) \leq \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq \exp(U(\underline{\nu}) T(1 + o(1))).$$

This problem is generally considered more difficult than the fixed-confidence setting (see, e.g., [Lattimore and Szepesvári, 2020](#), Chapter 33 and [Jourdan and Degenne, 2023](#), Section 6), and even for parametric models like canonical one-parameter exponential models, no strategy with matching upper and lower bounds (i.e., no optimal strategy) is known so far.

Earlier approaches. So far, four main approaches were considered for the problem of best-arm identification with a fixed budget. *First*, the early approach by [Audibert et al. \(2010\)](#) relies on gaps: we define the gap Δ_a of arm a as the difference $\mu^* - \mu_a$ between the largest expectation μ^* in $\underline{\nu}$ and the expectation of the distribution ν_a . They introduce a *Successive-Rejects* strategy and provide gap-based upper bounds for sub-Gaussian models, based on Hoeffding’s inequality. They however propose a lower bound only in the case of a Bernoulli model, not for larger, non-parametric, models. This lower bound was further discussed by [Carpentier and Locatelli \(2016\)](#), in a minimax sense. A *second series of approaches* (see, e.g., [Kaufmann et al., 2016](#)) focused on Gaussian bandits with fixed variances, but their results do not seem to be easily generalized to other models as they rely on specific properties (even stronger than the symmetry of the Kullback-Leibler divergence, namely, that in this model, the Kullback-Leibler divergence only depends on the gap between the expectations of the distributions). A *third approach*, led by [Russo \(2016, 2020\)](#), considered canonical one-parameter exponential families, but for a different target probability. Namely, a Bayesian setting is considered and the quality of a strategy is measured as the posterior probability of identifying the best arm. An optimal non-gap-based complexity is exhibited, together with optimal strategies matching this complexity. However, [Komiyama \(2022\)](#) argues that such an approach is specific to the Bayesian case and is not suited to the frequentist case that we consider. A *fourth approach* is to focus on the case of $K = 2$ arms, see, e.g., [Kaufmann et al. \(2016\)](#). The non-parametric bounds obtained therein do not enjoy any obvious generalization to the case of $K \geq 3$ arms beyond the one stated in Theorem 5.14 and criticized in Section 5.2.3 for only involving pairwise comparisons with the best arm. By considering very specific models, [Kato et al. \(2022\)](#) constructed a strategy that is optimal (only) in the regime where the gap between the 2 arms is small —yet, this gap-based approach does not, by nature, go in the direction of non-parametric bounds.

We will provide more details concerning some of these approaches while presenting and discussing our main results, in Section 5.2.2; see also Section 5.6.

Outline and contributions. In this chapter, we focus our attention on instance-dependent upper and lower bounds, holding for all problems of general models \mathcal{D} , including non-parametric models, and valid for any number K of arms. Put differently, we target a high degree of generality. While admittedly not exhibiting matching upper and lower bounds, we show that the same (new) information-theoretic quantities $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^>$ are at stake in these upper and lower bounds. These information-theoretic quantities are defined, in Section 5.2.1, as infima of Kullback-Leibler divergences and provide a quantification of the difficulty of the identification in terms of the geometry of information of the problem. We also present in Section 5.2.2 an overview of our results, which we carefully compare to existing bounds (restated therein, occasionally with some improvements). We state upper bounds in Section 5.3 and to do so, we provide an improved analysis of the classical *Successive-Rejects* strategy, not relying on gaps through Hoeffding’s lemma. Section 5.4 exhibits several possible lower bounds, which are inversely larger to the strength of the assumptions made on the strategies. These lower bounds generalize known lower bounds in the literature, like the lower bound for Bernoulli models by [Audibert et al. \(2010\)](#), but hold for arbitrary models. They share some similar flavor with the lower bounds by [Lai and Robbins \(1985\)](#) and [Burnetas and Katehakis \(1996\)](#) for the cumulative regret.

Notation. For a given strategy facing a bandit problem $\underline{\mu}$, let $N_a(t)$ and $\hat{\mu}_a(t)$ denote the number of pulls and the empirical mean¹ of arm a at step t :

$$N_a(t) \stackrel{\text{def}}{=} \sum_{s \in [t]} \mathbb{I}\{A_s = a\} \quad \text{and} \quad \hat{\mu}_a(t) \stackrel{\text{def}}{=} \frac{1}{N_a(t)} \sum_{s \in [t]} Y_s \mathbb{I}\{A_s = a\} .$$

¹As strategies initially observe each arm once, $\hat{\mu}_a(t)$ is well-defined for $t \geq K$.



Without loss of generality (see the paragraph on optional skipping page 37), we assume that the observation at time step t is $Y_t = X_{A_t, N_{A_t}(t)}$, where $(X_{a,n})_{a \in [K], n \geq 1}$ are independent random variables such that $X_{a,n} \sim \mathcal{N}(\mu_{A_t}, 1)$ for all $a \in [K]$ and $n \geq 1$. As a consequence, we notably get

$$\hat{\mu}_a(t) = \frac{1}{N_a(t)} \sum_{n=1}^{N_a(t)} X_{a,n} \stackrel{\text{def}}{=} \hat{\mu}_{a, N_a(t)}.$$

5.2. Overview of the Results and more Extended Literature Review

5.2.1. The Key new Quantities: $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^{\leq}$, as well as $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^{\geq}$

In this chapter, we only consider models \mathcal{D} whose distributions all admit an expectation. We denote by $\mathbb{E}(\zeta)$ the expectation of a distribution $\zeta \in \mathcal{D}$. For a distribution $\nu \in \mathcal{D}$ and a real number $x \in \mathbb{R}$, we then introduce

$$\begin{aligned} \mathcal{L}_{\text{inf}}^<(x, \nu) &= \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) < x\} \\ \text{and } \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) &= \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) \leq x\}, \end{aligned}$$

where KL denotes the Kullback-Leibler divergence and with the usual convention that the infimum of an empty set equals $+\infty$. Symmetrically, by considering rather distributions ζ with expectations larger than x , we define

$$\begin{aligned} \mathcal{L}_{\text{inf}}^>(x, \nu) &= \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) > x\} \\ \text{and } \mathcal{L}_{\text{inf}}^{\geq}(x, \nu) &= \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) \geq x\}. \end{aligned}$$

We state some general properties on these quantities in Section 5.5.1 —among others, that $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^{\leq}$, as well as $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^{\geq}$, are almost identical for the model $\mathcal{P}[0, 1]$. The same holds for canonical one-parameter exponential models, as discussed in Section 5.5.3 (see page 151). In our results, lower bounds will be typically expressed with $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^>$ quantities, while upper bounds will rely on $\mathcal{L}_{\text{inf}}^{\leq}$ and $\mathcal{L}_{\text{inf}}^{\geq}$ quantities.

Remark. The key quantities for the non-parametric study of best-arm identification with fixed-confidence by [Jourdan and Degenne \(2023\)](#) are defined based on Kullback-Leibler divergences with arguments in reverse order, namely,

$$\begin{aligned} \mathcal{K}_{\text{inf}}^-(\nu, x) &= \inf\{\text{KL}(\nu, \zeta) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) < x\} = \mathcal{K}_{\text{inf}}(\nu, x) \\ \text{and } \mathcal{K}_{\text{inf}}^+(\nu, x) &= \inf\{\text{KL}(\nu, \zeta) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) > x\}, \end{aligned}$$

where the first quantity was referred to as simply $\mathcal{K}_{\text{inf}}(\nu, x)$ by [Honda and Takemura \(2015\)](#) in the regret-minimization literature (see also Section 5.5.3 and [Garviev et al., 2022](#)). Optimal bounds for regret minimization only depend on $\mathcal{K}_{\text{inf}}(\nu, x)$.

The introduction of $\mathcal{L}_{\text{inf}}^<$, $\mathcal{L}_{\text{inf}}^{\leq}$, $\mathcal{L}_{\text{inf}}^>$, and $\mathcal{L}_{\text{inf}}^{\geq}$ are motivated by the fact that, in fixed-budget best-arm identification, the arguments in the KL are in reverse order compared to the fixed-confidence setting, see Section 5.2.2 (and also Section 2.3). Except for very specific models (e.g., the model \mathcal{D}_{σ^2} of Gaussian distributions with a fixed variance $\sigma^2 > 0$), the Kullback-Leibler divergence is not symmetric, i.e., $\text{KL}(\zeta, \nu)$ and $\text{KL}(\nu, \zeta)$ differ in general. Specific best-arm-identification results were obtained by [Kaufmann et al. \(2016\)](#) for the model \mathcal{D}_{σ^2} , based on the Bretagnolle-Huber inequality ([Bretagnolle and Huber, 1979](#)); they indicate that the sum of the inverse squared gaps would be driving both the lower bound and upper bound functions ℓ and U . However, a close look at the proof reveals that they heavily rely on a property even stronger than the symmetry of KL for this model: details and discussions on this matter are provided in Section 5.6.2. In particular, generalizations beyond the Gaussian case appear to be infeasible.

5.2.2. Overview of the Results

The chapter provides new and more general (possibly non-parametric) bounds on the misidentification errors based on the information-theoretic quantities introduced above. In particular, we consider a version of Chernoff information defined, for ν, ν' in \mathcal{D} with $E(\nu') < E(\nu)$, as

$$\mathcal{L}(\nu', \nu) = \inf_{x \in [E(\nu'), E(\nu)]} \left\{ \mathcal{L}_{\text{inf}}^{\geq}(x, \nu') + \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \right\}. \quad (5.1)$$

Given a bandit problem $\underline{\nu}$ with a unique optimal distribution denoted by ν^* , we may rank the arms a in non-decreasing order of $\mathcal{L}(\nu_a, \nu^*)$, i.e., consider the permutation σ such that

$$0 = \mathcal{L}(\nu_{\sigma_1}, \nu^*) < \mathcal{L}(\nu_{\sigma_2}, \nu^*) \leq \dots \leq \mathcal{L}(\nu_{\sigma_{K-1}}, \nu^*) \leq \mathcal{L}(\nu_{\sigma_K}, \nu^*). \quad (5.2)$$

Upper bound. *Our first main result* (Corollary 5.4 together with Lemma 5.5) considers models \mathcal{D} like $\mathcal{D} = \mathcal{P}[0, 1]$, the set of all probability distributions over $[0, 1]$, or $\mathcal{D} = \mathcal{D}_{\text{exp}}$, any canonical one-parameter exponential family. We study the **Successive-Rejects** strategy, introduced by [Audibert et al. \(2010\)](#), for which arms are rejected one by one at the end of phases of uniform exploration, and state that this strategy is such that for all bandit problems $\underline{\nu}$ in \mathcal{D} with a unique optimal arm,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\overline{\log K}} \min_{2 \leq k \leq K} \frac{\mathcal{L}(\nu_{\sigma_k}, \nu^*)}{k}, \quad (5.3)$$

where $\overline{\log K}$ is defined in (5.18) and is of order $\log K$. The key for this result (Lemma 5.2, of independent interest) is a grid-based application of the Cramér-Chernoff bound to control $\mathbb{P}(\overline{X}_N \leq \overline{Y}_N)$, where \overline{X}_N and \overline{Y}_N are averages of two independent N -samples. This approach can be used to analyze similar algorithms, like **Sequential-Halving** ([Karnin et al., 2013](#)).

Lower bounds. The corresponding lower bounds are stated rather in terms of $\mathcal{L}_{\text{inf}}^{\leq}$ and $\mathcal{L}_{\text{inf}}^{\geq}$ quantities, but Sections 5.5.1 and 5.5.3 (page 151) explain that those quantities are almost the same than $\mathcal{L}_{\text{inf}}^{\leq}$ and $\mathcal{L}_{\text{inf}}^{\geq}$ for regular models like $\mathcal{P}[0, 1]$ and exponential models (for those models, $\mathcal{L}(\nu', \nu)$ could be alternatively defined with $\mathcal{L}_{\text{inf}}^{\leq}$ and $\mathcal{L}_{\text{inf}}^{\geq}$ instead of $\mathcal{L}_{\text{inf}}^{\leq}$ and $\mathcal{L}_{\text{inf}}^{\geq}$, except in a single pathological case of the $\mathcal{P}[0, 1]$ model). We actually state several lower bounds in Section 5.4, that are larger as the assumptions on the strategies considered are more restrictive; as usual, there is a trade-off between the strength of a lower bound and its generality. However, all assumptions considered remain rather mild and are empirically satisfied by **Successive-Rejects**-type strategies: for instance, Definition 5.8 restricts the attention to strategies such that for all bandit problems, the arm associated with the smallest expectation is pulled less than a fraction $\frac{1}{K}$ of the time. Out of all lower bounds exhibited, *our second main result* (Theorem 5.13) holds, as indicated, under mild assumptions on the model and sequences of strategies considered, and reads: for all bandit problems $\underline{\nu}$ with no two same expectations,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq -\min_{2 \leq k \leq K} \inf_{x \in [\mu_{(k)}, \mu_{(k-1)})} \left\{ \frac{\mathcal{L}_{\text{inf}}^{\geq}(x, \nu_{(k)})}{k-1} + \frac{\mathcal{L}_{\text{inf}}^{\leq}(x, \nu^*)}{k} \right\}, \quad (5.4)$$

where $\mu_{(1)} > \mu_{(2)} > \mu_{(3)} > \dots > \mu_{(K)}$ and where $\nu_{(a)}$ denotes the distribution with expectation $\mu_{(a)}$. Here, we considered the notation (k) for order statistics in reverse order.

This lower bound does not match the exhibited upper bound, as is further discussed in Section 5.2.4. Still, we argue that quantities defined as infima over x of $\mathcal{L}_{\text{inf}}^{\geq}(x, \nu_{(k)}) + \mathcal{L}_{\text{inf}}^{\leq}(x, \nu^*)$ should measure how difficult a best-arm-identification problem is under a fixed budget. *This is the main insight of this chapter.*

5.2.3. Re-Derivation of Existing Bounds

We now survey the most important existing bounds and re-derive them from our general bounds. These existing bounds all hold only for sub-Gaussian models and for exponential models when $K \geq 3$, while a non-parametric bound was only available in the case of $K = 2$ arms.

To do so, we will sometimes consider the following weaker version of the lower bound (5.4), obtained by picking $x = \mu_{(k)}$:

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \frac{\mathcal{L}_{\text{inf}}^{\leq}(\mu_{(k)}, \nu^*)}{k}. \quad (5.5)$$

Comparison to the gap-based approaches. Audibert et al. (2010) propose an analysis of the Successive-Rejects strategy based on Hoeffding's inequality, stating that for all bandit problems in $\mathcal{P}[0, 1]$ with a unique optimal arm,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \leq - \frac{1}{\log K} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}, \quad (5.6)$$

where we recall the definition of the gaps $\Delta_{(k)} = \mu^* - \mu_{(k)}$. This bound is a consequence of (Corollary 5.4, a slightly more general form of) the bound (5.3), given Pinsker's inequality (5.26):

$$\mathcal{L}(\nu_{(k)}, \nu^*) \geq \inf_{x \in [\mu_{(k)}, \mu^*]} \left\{ 2(x - \mu_{(k)})^2 + 2(x - \mu^*)^2 \right\} = (\mu^* - \mu_{(k)})^2 = \Delta_{(k)}^2. \quad (5.7)$$

We remark that the bound (5.6) and the lower bound on $\mathcal{L}(\nu_{(k)}, \nu^*)$ may actually be extended to the model of σ^2 -sub-Gaussian distributions, up to considering factors $\frac{1}{4\sigma^2}$. We do not discuss the UCB-E algorithm of Audibert et al. (2010), as its performance and analysis crucially depend on a tuning parameter set with some knowledge of the gaps.

Audibert et al. (2010) also propose a carefully constructed lower bound for the model $\mathcal{B}_{[p, 1-p]} = \{\text{Ber}(x) : x \in [p, 1-p]\}$ of Bernoulli distributions $\text{Ber}(x)$ with parameters x in $[p, 1-p]$ for some $p \in (0, \frac{1}{2})$. A key inequality in their proof follows from the Kullback-Leibler $-\chi^2$ -divergence bound:

$$\forall x, y \in [p, 1-p], \quad \text{KL}(\text{Ber}(x), \text{Ber}(y)) \leq \frac{(x-y)^2}{2p(1-p)}.$$

Their construction may actually be generalized to models \mathcal{D} with $C_{\mathcal{D}} > 0$ such that for all ν, ν' in \mathcal{D} , one has $\text{KL}(\nu, \nu') \leq C_{\mathcal{D}} (\mathbb{E}(\nu) - \mathbb{E}(\nu'))^2$. This is a property that clearly holds for some exponential families: on top of the restricted Bernoulli model discussed above, for which

$$C_{\mathcal{B}_{[p, 1-p]}} = \frac{1}{(2p(1-p))},$$

we may cite the model \mathcal{D}_{σ^2} of Gaussian distributions with variance σ^2 , for which $C_{\mathcal{D}_{\sigma^2}} = 1/(2\sigma^2)$. For models enjoying the existence of such a constant $C_{\mathcal{D}}$, (a straightforward modification of) the analysis by Audibert et al. (2010) entails that for any $\underline{\nu}$ in \mathcal{D} ,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq -5 C_{\mathcal{D}} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}. \quad (5.8)$$

As by the very assumption on the model, $\mathcal{L}_{\text{inf}}^{\leq}(\mu_{(k)}, \nu^*) \leq C_{\mathcal{D}} \Delta_{(k)}^2$, the lower bound (5.5) implies the stated lower bound (5.8), with an improved constant factor.

The lower bound (5.8) and the upper bound (5.6) differ in particular by a factor proportional to $\overline{\log} K$. [Carpentier and Locatelli \(2016\)](#) discuss this gap in the case of the Bernoulli model $\mathcal{B}_{[1/4, 3/4]}$ and improve the lower bound (5.8) by a factor of $\log K$, but not simultaneously for all bandit problems $\underline{\nu}$ (as we aim for); they obtain the improvement just for one bandit problem $\underline{\nu}$. Their lower bound result (formally stated and discussed in Section 5.6.1) is therefore of a totally different nature. More results on how and when given lower bounds with a given complexity measure may, or may not, be improved were stated by [Komiyama et al. \(2022\)](#).

Discussion on the non-parametric bound for $K = 2$ arms of [Kaufmann et al. \(2016\)](#). It turns out that the existing literature for the fixed-budget setting offered so far a non-parametric bound, in the case of $K = 2$ arms. Namely, in a general, possibly non-parametric model \mathcal{D} , [Kaufmann et al. \(2016, Theorem 12\)](#) stated a lower bound for all 2-armed bandit problems $\underline{\nu} = (\nu_1, \nu_2)$:

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \inf_{\substack{\lambda \text{ in } \mathcal{D}: \\ \mathbb{E}(\lambda_{a^*(\underline{\nu})}) < \mathbb{E}(\lambda_{w_*(\underline{\nu})})}} \max\{\text{KL}(\lambda_{w_*(\underline{\nu})}, \nu_{w_*(\underline{\nu})}), \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})})\}, \quad (5.9)$$

where $w_*(\underline{\nu})$ denotes the sub-optimal arm in $\underline{\nu}$ and where the infimum is over all alternative bandit problems (λ_1, λ_2) in \mathcal{D} with a different best arm than $\underline{\nu}$. We note (see the proof of Theorem 5.14) that we may actually rewrite this lower bound in a more readable way, in terms of $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^>$ quantities, illustrating once again that these quantities are key in measuring the complexity of best-arm identification under a fixed budget:

$$\begin{aligned} \inf_{\substack{\lambda \text{ in } \mathcal{D}: \\ \mathbb{E}(\lambda_{a^*(\underline{\nu})}) < \mathbb{E}(\lambda_{w_*(\underline{\nu})})}} \max\{\text{KL}(\lambda_{w_*(\underline{\nu})}, \nu_{w_*(\underline{\nu})}), \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})})\} \\ = \inf_{x \in [\mu_{w_*(\underline{\nu})}, \mu^*]} \left\{ \max\{\mathcal{L}_{\text{inf}}^>(x, \nu_{w_*(\underline{\nu})}), \mathcal{L}_{\text{inf}}^<(x, \nu^*)\} \right\}. \end{aligned} \quad (5.10)$$

The proof technique of [Kaufmann et al. \(2016\)](#) may be applied in a pairwise fashion to generalize the lower bound (5.10) for 2 arms into a lower bound for $K \geq 2$ arms, stated in Theorem 5.14: for all $\underline{\nu}$ in \mathcal{D} with a unique optimal arm,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \min_{k \neq a^*(\underline{\nu})} \inf_{x \in [\mu_k, \mu^*]} \left\{ \max\{\mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu^*)\} \right\}. \quad (5.11)$$

We do not claim that (5.11) is a deep and interesting bound, as it only involves pairwise comparisons with the best arm. In particular, we lack divisions by the ranks of the arms, as in (5.4). This is why we had not stated the result (5.11) of Theorem 5.14 in Section 5.2.2 and mention it only here.

That being said, given that the infima in (5.4) are over more restricted ranges than in (5.11), we can see no obvious ranking between the two bounds, which rather look incomparable.

Bounds for $K = 2$ arms and exponential families, cf. comments after Theorem 12 of [Kaufmann et al. \(2016\)](#). We denote by \mathcal{D}_{exp} the model corresponding to a canonical one-parameter exponential family with expectations defined on an open interval \mathcal{M} (see page 151 in Section 5.5.3 for a reminder on this matter). For such a model, we denote by d the mean-parameterized Kullback-Leibler divergence. By continuity of d , we have that for all ν in \mathcal{D}_{exp} and for all $x \in \mathcal{M}$,

$$\forall x \leq \mathbb{E}(\nu), \quad \mathcal{L}_{\text{inf}}^<(x, \nu) = \mathcal{L}_{\text{inf}}^<(x, \nu) = d(x, \mathbb{E}(\nu)), \quad (5.12)$$

$$\text{and } \forall x \geq \mathbb{E}(\nu), \quad \mathcal{L}_{\text{inf}}^>(x, \nu) = \mathcal{L}_{\text{inf}}^>(x, \nu) = d(x, \mathbb{E}(\nu)). \quad (5.13)$$

Note that all bounds stated in Section 5.2.2 then admit simple reformulations in terms of d . The Chernoff-information-type quantity \mathcal{L} introduced in (5.1) may also be mean-parameterized as follows: for $\mu' < \mu$,

$$L(\mu', \mu) = \min_{x \in [\mu', \mu]} \{d(x, \mu') + d(x, \mu)\}. \quad (5.14)$$

We now explain why we called L (and therefore \mathcal{L}) a version of Chernoff information. The original definition of the Chernoff information $D(\mu', \mu)$ is the value $d(y, \mu)$ for $y \in [\mu', \mu]$ such that $d(y, \mu') = d(y, \mu)$. As mentioned in the comments after Theorem 12 of [Kaufmann et al. \(2016\)](#), D is the quantity at stake in (5.10) for a canonical one-parameter exponential family: given that $d(\cdot, \mu')$ and $d(\cdot, \mu)$ are respectively increasing and decreasing on $[\mu', \mu]$,

$$\min_{x \in [\mu', \mu]} \max\{d(x, \mu'), d(x, \mu)\} = D(\mu', \mu).$$

Therefore, $D(\mu', \mu) \leq L(\mu', \mu) \leq 2D(\mu', \mu)$, which shows that L is related to D , as claimed.

Example 5.1. We state the lower bound (5.5) and the upper bound (5.3) for the model $\mathcal{B}_{[p, 1-p]}$ of Bernoulli distributions with parameters in $[p, 1-p]$, where $p \in (0, \frac{1}{2})$. We denote by

$$\text{kl}(x, y) = x \log \frac{x}{y} + (1-x) \log \frac{1-x}{1-y}, \quad \text{where } x, y \in [p, 1-p]$$

the mean-parameterized Kullback-Leibler divergence of this model. We consider a generic bandit problem $\nu = (\text{Ber}(p_1), \dots, \text{Ber}(p_K))$. We rank the parameters as in (5.4), i.e., introduce the notation $p^* = p_{(1)} > p_{(2)} > \dots > p_{(K)}$. Then, after noticing (see Lemma 5.19 in Section 5.5.3) that this ranking is the same as the one considered in (5.2), the upper bound (5.3) rewrites as

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\nu)) \leq -\frac{1}{\log K} \min_{2 \leq k \leq K} \frac{\min_{x \in [p_{(k)}, p^*]} \{\text{kl}(x, p_{(k)}) + \text{kl}(x, p^*)\}}{k},$$

while the lower bound (5.5) rewrites as

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_\nu(\hat{a}_T \neq a^*(\nu)) \geq -\min_{2 \leq k \leq K} \frac{\text{kl}(p_{(k)}, p^*)}{k}.$$

They should be compared to the upper (5.6) and lower (5.8) bounds of [Audibert et al. \(2010\)](#), respectively.

5.2.4. Discussion of the (Lack of) Optimality of the new Bounds Exhibited

The lower bound (5.4) does not match the upper bound (5.3) because of two aspects. First, the infima in (5.4) are only taken on restricted ranges $[\mu_{(k)}, \mu_{(k-1)})$ and not on the entire intervals $[\mu_{(k)}, \mu^*]$ as in (5.3). Second, the upper bound (5.3) involves a $1/\log K$ factor, while the lower bound (5.4) does not. A similar $1/\log K$ factor was missing between the upper (5.6) and lower (5.8) bounds of [Audibert et al. \(2010\)](#) for Bernoulli models, together with a numerical factor of $5C_{\mathcal{B}_{[p, 1-p]}}$. The non-parametric bounds exhibited in this chapter mainly generalize and extend the known parametric bounds but do not refine the latter in the sense that gaps between upper and lower bounds would be closed.

That being said, we would like to illustrate below one specific example to which extent the gap-based bounds can be looser.

Example of an extreme improvement: distributions with separated supports. For general non-parametric models, gaps are not enough at all to measure complexity as we may well have a finite gap between two distributions ν_1 and ν_2 with $\mu_1 > \mu_2$, but $\mathcal{L}(\nu_2, \nu_1) = +\infty$. This holds, for instance, as soon as ν_1 and ν_2 have closed supports separated by a threshold x_0 (see Figure 5.1), i.e., the closed supports of ν_1 and ν_2 are included in $(-\infty, x_0)$ and $(x_0, +\infty)$, respectively. Indeed,

by mimicking the beginning of the proof of Lemma 5.16, it may be seen that $\mathcal{L}_{\inf}^{\leq}(x, \nu_1) = +\infty$ for $x \leq x_0$ and $\mathcal{L}_{\inf}^{\geq}(x, \nu_2) = +\infty$ if $x \geq x_0$, so that in all cases, the sum $\mathcal{L}_{\inf}^{\geq}(x, \nu_2) + \mathcal{L}_{\inf}^{\leq}(x, \nu_1)$ equals $+\infty$, and thus, $\mathcal{L}(\nu_2, \nu_1) = +\infty$. In our bounds, e.g., the upper bound (5.3), the pair of distributions ν_1, ν_2 will therefore not contribute —as intuition commands: these two distributions are easy to distinguish—, while it does contribute to the earlier gap-based bounds.

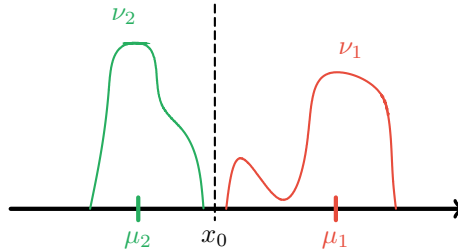


Figure 5.1: When the two distributions have separated supports, there is no confusion to identify which of both has the highest mean.

5.3. Upper Bound for the Successive-Rejects Strategy, with an Improved Analysis

We consider the Successive-Rejects strategy introduced by [Audibert et al. \(2010\)](#), for K arms and a budget T . The strategy works in phases, and the lengths of the phases are set beforehand; they are denoted by $\ell_1, \dots, \ell_{K-1} \geq 1$ and satisfy $\ell_1 + \dots + \ell_{K-1} = T$. The strategy maintains a list of candidate arms, starting with all arms, i.e., $S_0 = [K]$. At the end of each phase $r \in \{1, \dots, K-1\}$, it drops an arm to get S_r , while during phase r , it operates with the $K-r+1$ arms in S_{r-1} .

More precisely, during phase $r \in \{1, \dots, K-1\}$, the strategy draws $\lfloor \frac{\ell_r}{K-r+1} \rfloor$ times each arm in S_{r-1} (and does not use the few remaining time steps, if there are some). At the end of each phase r , the strategy computes the empirical averages \bar{X}_a^r of the payoffs obtained by each arm $a \in S_{r-1}$ since the beginning; i.e., \bar{X}_a^r is an average over

$$N_r = \left\lfloor \frac{\ell_1}{K} \right\rfloor + \dots + \left\lfloor \frac{\ell_r}{K-r+1} \right\rfloor$$

i.i.d. realizations of ν_a . It then drops the arm a_r with smallest empirical average (ties broken arbitrarily). This description is summarized in Algorithm 16.

5.3.1. General Analysis

The key quantities for the general analysis will be the logarithmic moment-generating function ϕ_ν of a distribution $\nu \in \mathcal{D}$, and its Fenchel-Legendre transform ϕ_ν^* :

$$\forall \lambda \in \mathbb{R}, \quad \phi_\nu(\lambda) = \log \int_{\mathbb{R}} e^{\lambda x} d\nu(x) \quad \text{and} \quad \forall x \in \mathbb{R}, \quad \phi_\nu^*(x) = \sup_{\lambda \in \mathbb{R}} \{\lambda x - \phi_\nu(\lambda)\}. \quad (5.15)$$

Based on them, we can now define, for all $\nu, \nu' \in \mathcal{D}$ with $E(\nu') < E(\nu)$,

$$\Phi(\nu', \nu) \stackrel{\text{def}}{=} \inf_{x \in [E(\nu'), E(\nu)]} \{\phi_{\nu'}^*(x) + \phi_\nu^*(x)\}.$$

The following lemma shows that Φ plays a significant role in bounding the probability that two sample averages are in reverse order compared to the expectations of the underlying distributions. It supersedes the use of Hoeffding's inequality in [Audibert et al. \(2010\)](#).

Algorithm 16: Successive-Rejects algorithm

Input: budget parameter T

 phase lengths $(\ell_r)_{r \in [K-1]}$ such that $\sum_{r \in [K-1]} \ell_r = T$
Output: estimated best arm \hat{a}_T

- 1 $t \leftarrow 0$
 - 2 $S_0 \leftarrow [K]$
 - 3 **for** each round $r \in [K-1]$ **do**
 - 4 Observe each arm $\lfloor \frac{\ell_r}{K-r+1} \rfloor$ times
 - 5 Increase t by ℓ_r
 - 6 Choose $a_r \in \underset{a \in S_{r-1}}{\operatorname{argmin}} \hat{\mu}_a(t)$
 - 7 $S_r \leftarrow S_{r-1} \setminus \{a_r\}$
 - 8 Define \hat{a}_T as the unique element of S_{K-1}
-

Lemma 5.2. Fix ν and ν' in \mathcal{D} , with respective expectations $\mu = \mathbb{E}(\nu) > \mu' = \mathbb{E}(\nu')$. For all $N \geq 1$, let \bar{X}_N and \bar{Y}_N be the averages of N -samples with respective distributions ν and ν' . Then,

$$\limsup_{N \rightarrow +\infty} \frac{1}{N} \log \mathbb{P}(\bar{X}_N \leq \bar{Y}_N) \leq - \inf_{x \in [\mu', \mu]} \{ \phi_{\nu'}^*(x) + \phi_{\nu}^*(x) \} \stackrel{\text{def}}{=} -\Phi(\nu', \nu).$$

Proof. The proof consists of two parts. We first show that for any finite grid $\mathcal{G} = \{g_2, \dots, g_{G-1}\}$ in (μ', μ) , to which we add the points $g_1 = \mu'$ and $g_G = \mu$, we have

$$\limsup_{N \rightarrow +\infty} \frac{1}{N} \log \mathbb{P}(\bar{X}_N \leq \bar{Y}_N) \leq - \min \left\{ \phi_{\nu}^*(\mu'), \min_{2 \leq j \leq G-1} \{ \phi_{\nu'}^*(g_{j-1}) + \phi_{\nu}^*(g_j) \}, \phi_{\nu'}^*(\mu) \right\}. \quad (5.16)$$

Indeed, by identifying (when \bar{X}_N and \bar{Y}_N belong to $[\mu', \mu]$) in which interval $[g_{j-1}, g_j]$ lies \bar{X}_N , we note that

$$\{\bar{X}_N \leq \bar{Y}_N\} \subseteq \{\bar{X}_N \leq \mu'\} \cup \{\bar{Y}_N \geq \mu\} \cup \bigcup_{j=2}^{G-1} \{\bar{Y}_N \geq g_{j-1} \text{ and } \bar{X}_N \leq g_j\}.$$

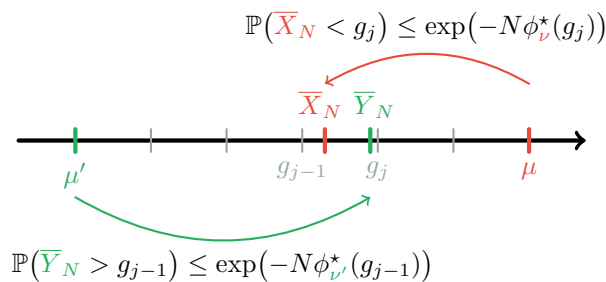


Figure 5.2: Control of the probability of $\{\bar{Y}_N \geq g_{j-1} \text{ and } \bar{X}_N \leq g_j\}$ using a Cramér-Chernoff bound.

First, by independence and by the Cramér-Chernoff inequalities (5.29) and (5.30),

$$\mathbb{P}(\bar{Y}_N \geq g_{j-1} \text{ and } \bar{X}_N \leq g_j) = \mathbb{P}(\bar{Y}_N \geq g_{j-1}) \mathbb{P}(\bar{X}_N \leq g_j) \leq \exp\left(-N(\phi_{\nu'}^*(g_{j-1}) + \phi_{\nu}^*(g_j))\right).$$

Second, again by the Cramér-Chernoff inequalities,

$$\mathbb{P}(\bar{X}_N \leq \mu') \leq \exp(-N \phi_\nu^*(\mu')) \quad \text{and} \quad \mathbb{P}(\bar{Y}_N \geq \mu) \leq \exp(-N \phi_{\nu'}^*(\mu)).$$

By a union bound,

$$\mathbb{P}(\bar{X}_N \leq \bar{Y}_N) \leq \exp(-N \phi_\nu^*(\mu')) + \exp(-N \phi_{\nu'}^*(\mu)) + \sum_{j=2}^{G-1} \exp(-N(\phi_{\nu'}^*(g_{j-1}) + \phi_\nu^*(g_j))).$$

The stated bound (5.16) follows by identifying the (finitely many) terms with the smallest rate in the exponent.

In the second part of the proof, we note that the bound (5.16) holds for any finite grid in (μ', μ) , and we consider a sequence

$$\mathcal{G}^{(n)} = \{g_2^{(n)}, \dots, g_{G_n-1}^{(n)}\}$$

of such finite grids. In particular,

$$\begin{aligned} \limsup_{N \rightarrow +\infty} \frac{1}{N} \log \mathbb{P}(\bar{X}_N \leq \bar{Y}_N) &\leq -\min\left\{\phi_\nu^*(\mu'), \max_{n \geq 1} S_n, \phi_{\nu'}^*(\mu)\right\}, \\ \text{where } S_n &\stackrel{\text{def}}{=} \min_{2 \leq j \leq G_n-1} \left\{\phi_{\nu'}^*(g_{j-1}^{(n)}) + \phi_\nu^*(g_j^{(n)})\right\}. \end{aligned}$$

To obtain the claimed bound, given that (see the end of Section 5.5.2)

$$\phi_\nu^*(\mu) = 0 = \phi_{\nu'}^*(\mu'),$$

it suffices to show that

$$\max_{n \geq 1} S_n \geq \inf_{x \in [\mu', \mu]} \{\phi_{\nu'}^*(x) + \phi_\nu^*(x)\}.$$

To that end, we assume that the steps ε_n of the grids $\mathcal{G}^{(n)}$, which are defined as

$$\varepsilon_n \stackrel{\text{def}}{=} \max_{2 \leq j \leq G_n} |g_j^{(n)} - g_{j-1}^{(n)}|,$$

vanish asymptotically, i.e., $\varepsilon_n \rightarrow 0$. For each grid $\mathcal{G}^{(n)}$, we denote by $x_n^* \in (\mu', \mu)$ the argument of the minimum in the definition of S_n . As a consequence, for each $n \geq 1$,

$$S_n = \phi_{\nu'}^*(x_n^* - \varepsilon_n^*) + \phi_\nu^*(x_n^*),$$

for some $0 < \varepsilon_n^* \leq \varepsilon_n$. The quantity $x_n^* - \varepsilon_n^*$ denotes the point in the grid that is right before x_n^* , and it belongs to $[\mu', \mu)$. We note that we also have $\varepsilon_n^* \rightarrow 0$. In the compact interval $[\mu', \mu]$, the Bolzano-Weierstrass theorem (see, e.g., [Bartle and Sherbert, 2000](#), Section 3.4) ensures the existence of a converging sub-sequence: there exists $x_\infty^* \in [\mu', \mu]$ and a sequence $(n_k)_{k \geq 1}$ of integers such that

$$x_{n_k}^* \xrightarrow[k \rightarrow +\infty]{} x_\infty^*, \quad \text{which also entails} \quad x_{n_k}^* - \varepsilon_{n_k}^* \xrightarrow[k \rightarrow +\infty]{} x_\infty^*.$$

Now, the functions ϕ_ν^* , respectively, $\phi_{\nu'}^*$, are lower semi-continuous, as the suprema over $\lambda \in \mathbb{R}$ of the continuous functions $x \mapsto \lambda x - \varphi_\nu(\lambda)$, respectively, $x \mapsto \lambda x - \varphi_{\nu'}(\lambda)$. Therefore, by these lower semi-continuities,

$$\begin{aligned} \max_{n \geq 1} S_n &\geq \liminf_{k \rightarrow +\infty} \phi_{\nu'}^*(x_{n_k}^* - \varepsilon_{n_k}^*) + \phi_\nu^*(x_{n_k}^*) \geq \phi_{\nu'}^*(x_\infty^*) + \phi_\nu^*(x_\infty^*) \\ &\geq \inf_{x \in [\mu', \mu]} \{\phi_{\nu'}^*(x) + \phi_\nu^*(x)\}. \end{aligned}$$

This concludes the proof. □

General upper bound. The main performance upper bound is stated below in terms of Φ , that is, in terms of Fenchel-Legendre transforms of logarithmic moment-generating functions. Section 5.4.2 will later explain why and when the latter may be replaced by $\mathcal{L}_{\inf}^{\leq}$ and $\mathcal{L}_{\inf}^{\geq}$ quantities, leading to a rewriting $\Phi = \mathcal{L}$ and to the bound claimed in (5.3).

Theorem 5.3. Fix $K \geq 2$ and a model \mathcal{D} . Consider a sequence of Successive-Rejects strategies, indexed by T , such that $N_r/T \rightarrow \gamma_r > 0$ as $T \rightarrow +\infty$ for all $r \in \{1, \dots, K-1\}$. Let $\underline{\nu}$ be a bandit problem in \mathcal{D} with a unique optimal arm and, for each $r \in \{1, \dots, K-1\}$, let \mathcal{A}_r be a subset of arms of cardinality r that does not contain $a^*(\underline{\nu})$. Then

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq - \min_{1 \leq r \leq K-1} \left\{ \gamma_r \min_{k \in \mathcal{A}_r} \Phi(\nu_k, \nu^*) \right\}.$$

The proof mimics the analysis by [Audibert et al. \(2010\)](#), the main modification being the substitution of Hoeffding's inequality by the bound of Lemma 5.2.

Proof. We recall that for $r \in \{1, \dots, K-1\}$, we denoted by $N_r = \lfloor \ell_1/K \rfloor + \dots + \lfloor \ell_r/(K-r+1) \rfloor$ the total number of times an arm still considered in phase r , i.e., belonging to S_{r-1} , was pulled in phases 1 to r . For each arm a , we denote by \bar{Y}_a^r the average of a N_r -sample distributed according to ν_a . By optional skipping (see [Doob, 1953](#), Chapter III, Theorem 5.2, p. 145, or [Chow and Teicher, 1988](#), Section 5.3 for a more recent reference), we may assume, with no loss of generality, that for each $r \in \{1, \dots, K-1\}$,

$$\text{on the event } \{a \in S_{r-1}\}, \quad \bar{X}_a^r = \bar{Y}_a^r. \quad (5.17)$$

We fix a bandit problem $\underline{\nu}$ with a unique optimal arm $a^*(\underline{\nu})$. The Successive-Rejects strategy fails if (and only) if it rejects $a^*(\underline{\nu})$ in one of the phases. This corresponds to the event

$$\{\hat{a}_T \neq a^*(\underline{\nu})\} = \bigcup_{r=1}^{K-1} \{a_r = a^*(\underline{\nu})\} \subseteq \bigcup_{r=1}^{K-1} \left\{ a^*(\underline{\nu}) \in S_{r-1} \text{ and } \forall k \in S_{r-1}, \bar{X}_{a^*(\underline{\nu})}^r \leq \bar{X}_k^r \right\},$$

where we recall that a_r is the arm removed at the end of phase r (we have an inclusion because ties are broken arbitrarily). By optional skipping (5.17),

$$\begin{aligned} & \bigcup_{r=1}^{K-1} \left\{ a^*(\underline{\nu}) \in S_{r-1} \text{ and } \forall k \in S_{r-1}, \bar{X}_{a^*(\underline{\nu})}^r \leq \bar{X}_k^r \right\} \\ &= \bigcup_{r=1}^{K-1} \left\{ a^*(\underline{\nu}) \in S_{r-1} \text{ and } \forall k \in S_{r-1}, \bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r \right\}. \end{aligned}$$

Recall that the set S_{r-1} is a random set; dealing with it therefore requires some care. On the event of interest, S_{r-1} contains $K-r+1$ elements, among which $a^*(\underline{\nu})$. The set \mathcal{A}_r is of cardinality r and does not contain $a^*(\underline{\nu})$. By the pigeonhole principle, S_{r-1} thus necessarily contains one arm in \mathcal{A}_r . As a consequence, for each phase $r \in \{1, \dots, K-1\}$,

$$\left\{ a^*(\underline{\nu}) \in S_{r-1} \text{ and } \forall k \in S_{r-1}, \bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r \right\} \subseteq \bigcup_{k \in \mathcal{A}_r} \left\{ \bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r \right\}.$$

Summarizing the inclusions above, taking unions bounds, and upper bounding the obtained sum in a crude way, we proved so far

$$\mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq \sum_{r=1}^{K-1} \sum_{k \in \mathcal{A}_r} \mathbb{P}(\bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r) \leq K^2 \max_{1 \leq r \leq K-1} \max_{k \in \mathcal{A}_r} \mathbb{P}(\bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r),$$

or equivalently,

$$\begin{aligned} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) &\leq \frac{2}{T} \log K + \max_{1 \leq r \leq K-1} \max_{k \in \mathcal{A}_r} \frac{1}{T} \log \mathbb{P}(\bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r) \\ &= \frac{2}{T} \log K + \max_{1 \leq r \leq K-1} \max_{k \in \mathcal{A}_r} \frac{N_r}{T} \frac{1}{N_r} \log \mathbb{P}(\bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r). \end{aligned}$$

As $N_r/T \rightarrow \gamma_r > 0$ as $T \rightarrow +\infty$, we may apply Lemma 5.2, together with an exchange between the lim sup and the maximum over a finite number of quantities. We obtain

$$\begin{aligned} \limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) &\leq \max_{1 \leq r \leq K-1} \max_{k \in \mathcal{A}_r} \left\{ \gamma_r \left(-\Phi(\nu_k, \nu^*) \right) \right\} \\ &= - \min_{1 \leq r \leq K-1} \left\{ \gamma_r \min_{k \in \mathcal{A}_r} \Phi(\nu_k, \nu^*) \right\}. \end{aligned}$$

This concludes the proof. \square

With the phase lengths of Audibert et al. (2010). We conclude this subsection by stating the bound of Theorem 5.3 for the phase lengths suggested by Audibert et al. (2010), namely,

$$\ell_1 \stackrel{\text{def}}{=} \frac{T}{\overline{\log K}}, \quad \text{and} \quad \forall r \in \{2, \dots, K-1\}, \quad \ell_r \stackrel{\text{def}}{=} \frac{T}{(K-r+2) \overline{\log K}}, \quad (5.18)$$

where we define

$$\overline{\log K} \stackrel{\text{def}}{=} \frac{1}{2} + \sum_{k=2}^K \frac{1}{k},$$

The phase lengths in the example case of $K = 6$ arms are shown in Figure 5.3. Rather than Φ , it

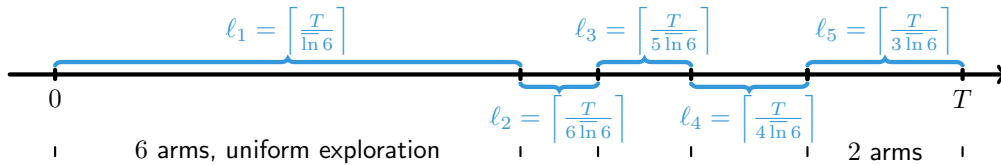


Figure 5.3: Successive-Rejects phase lengths of Audibert et al. (2010) for $K = 6$ arms.

is sometimes handy to rely on more readable quantities, this is why we will consider lower bounds $f(\nu_k, \nu^*)$ on the $\Phi(\nu_k, \nu^*)$ quantities. We may of course use $f = \Phi$ but f can also be, for instance, the squared gaps: in the case of the $\mathcal{P}[0, 1]$ model, Hoeffding's inequality entails that

$$\phi_\nu^*(x) \geq 2(x - \mathbb{E}(\nu))^2, \quad \text{so that} \quad \Phi(\nu_k, \nu^*) \geq \Delta_k^2 \stackrel{\text{def}}{=} f(\nu_k, \nu^*), \quad (5.19)$$

as explained in the proof below.

Remark. Such bounds hold more generally in models consisting of sub-Gaussian distributions. For ease of exposition, the path followed in Section 5.2.2 to show that $\Phi(\nu_k, \nu^*) \geq \Delta_k^2$ was to first note that $\Phi = \mathcal{L}$ when $\mathcal{D} = \mathcal{P}[0, 1]$ (see Lemma 5.5) and then use Pinsker's inequality (5.7). We provide below a slightly more direct but equivalent approach, based on Hoeffding's inequality.

Proof of (5.19). When $\nu \in \mathcal{P}[0, 1]$, Hoeffding's inequality exactly states that

$$\forall \lambda \in \mathbb{R}, \quad \phi_\nu(\lambda) \leq \lambda \mathbb{E}(\nu) + \frac{\lambda^2}{8},$$

$$\text{so that} \quad \forall x \in \mathbb{R}, \quad \phi_\nu^*(x) \geq \sup_{\lambda \in \mathbb{R}} \left\{ \lambda(x - \mathbb{E}(\nu)) - \frac{\lambda^2}{8} \right\} = 2(x - \mathbb{E}(\nu))^2.$$

This corresponds to the first part of (5.19).

For its second part, we consider a pair ν, ν' of distributions in $\mathcal{P}[0, 1]$, we set any $x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]$, and we apply twice the bound of the first part to get

$$\phi_{\nu'}^*(x) + \phi_{\nu}^*(x) \geq 2(x - \mathbb{E}(\nu'))^2 + 2(x - \mathbb{E}(\nu))^2.$$

From the definition of Φ , it follows that

$$\Phi(\nu', \nu) \geq \inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \left\{ 2(x - \mathbb{E}(\nu'))^2 + 2(x - \mathbb{E}(\nu))^2 \right\} = (\mathbb{E}(\nu') - \mathbb{E}(\nu))^2.$$

This corresponds to the second part of (5.19). \square

We now order the arms into $\sigma_1, \dots, \sigma_K$ based on f , namely, we let $\sigma_1 = a^*(\underline{\nu})$ and

$$0 = f(\nu_{\sigma_1}, \nu^*) < f(\nu_{\sigma_2}, \nu^*) \leq \dots \leq f(\nu_{\sigma_{K-1}}, \nu^*) \leq f(\nu_{\sigma_K}, \nu^*), \quad (5.20)$$

and we take $\mathcal{A}_r = \{\sigma_{K-r+1}, \dots, \sigma_K\}$. We obtain the following corollary.

Corollary 5.4. *Fix $K \geq 2$, a model \mathcal{D} , and consider a lower bound f on Φ . The sequence of Successive-Rejects strategies based on the phase lengths (5.18) ensures, that for all bandit problems $\underline{\nu}$ in \mathcal{D} with a unique optimal arm,*

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\log K} \min_{2 \leq k \leq K} \frac{f(\nu_{\sigma_k}, \nu^*)}{k},$$

where arms were reordered as in (5.20).

Proof. To apply Theorem 5.3, we need only to show that the phase lengths of (5.18) are such that N_r/T converges to a positive value, and to identify this limit value γ_r . As $N_1 = \lfloor \ell_1/K \rfloor$, where $\ell_1 = T/\log K$, we immediately have $N_1/T \rightarrow \gamma_1 = 1/(K \log K) > 0$. For $r \in \{2, \dots, K-1\}$,

$$\begin{aligned} \frac{N_r}{T} &= \sum_{p=1}^r \frac{1}{T} \left\lfloor \frac{\ell_p}{K} \right\rfloor = \frac{1}{T} \left(\left\lfloor \frac{T}{K \log K} \right\rfloor + \sum_{p=2}^r \left\lfloor \frac{T}{(K-p+1)(K-p+2) \log K} \right\rfloor \right) \\ &\xrightarrow{T \rightarrow +\infty} \gamma_r \stackrel{\text{def}}{=} \frac{1}{\log K} \left(\frac{1}{K} + \sum_{p=2}^r \frac{1}{K-p+1} - \frac{1}{K-p+2} \right) = \frac{1}{(K-r+1) \log K}. \end{aligned}$$

The bound of Theorem 5.3 reads:

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\log K} \min_{1 \leq r \leq K-1} \left\{ \frac{1}{K-r+1} \min_{k \in \mathcal{A}_r} \Phi(\nu_k, \nu^*) \right\}.$$

It implies, in terms of lower bounds $f(\nu_k, \nu^*) \leq \Phi(\nu_k, \nu^*)$,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\log K} \min_{1 \leq r \leq K-1} \left\{ \frac{1}{K-r+1} \min_{k \in \mathcal{A}_r} f(\nu_k, \nu^*) \right\}. \quad (5.21)$$

The permutation σ in (5.20) and the sets $\mathcal{A}_r = \{\sigma_{K-r+1}, \dots, \sigma_K\}$ were exactly picked, for each $r \in \{1, \dots, K-1\}$, to minimize

$$\min_{k \in \mathcal{B}_r} f(\nu_k, \nu^*)$$

over sets \mathcal{B}_r abiding by the indicated constraints: being of cardinal r and not containing the optimal arm $a^*(\underline{\nu}) = \sigma_1$. We get

$$\min_{k \in \mathcal{A}_r} f(\nu_k, \nu^*) = \min_{K-r+1 \leq k \leq K} f(\nu_{\sigma_k}, \nu^*) = f(\nu_{\sigma_{K-r+1}}, \nu^*),$$

which, together with (5.21), yields the stated bound, up to replacing $K-r+1$ with $r \in \{1, \dots, K-1\}$ by $k \in \{2, \dots, K\}$:

$$-\frac{1}{\log K} \min_{1 \leq r \leq K-1} \left\{ \frac{1}{K-r+1} f(\nu_{\sigma_{K-r+1}}, \nu^*) \right\} = -\frac{1}{\log K} \min_{2 \leq k \leq K} \left\{ \frac{1}{k} f(\nu_{\sigma_k}, \nu^*) \right\}. \quad \square$$

5.3.2. On Links between Φ and the Quantities $\mathcal{L}_{\text{inf}}^<$, $\mathcal{L}_{\text{inf}}^{\leq}$, $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^{\geq}$

The Fenchel-Legendre transform ϕ_ν^* of the logarithmic moment-generating function of ν admits a classical (see, e.g., [Boucheron et al., 2013](#), Exercice 4.13) dual formulation in terms of infima of Kullback-Leibler divergences. The following lemma, proved in Section 5.5.3 (see page 149), reveals that these infima correspond to $\mathcal{L}_{\text{inf}}^{\leq}$ and $\mathcal{L}_{\text{inf}}^{\geq}$ for the model $\mathcal{P}[0, 1]$ of distributions supported on $[0, 1]$.

Lemma 5.5. *Consider the model $\mathcal{D} = \mathcal{P}[0, 1]$. For all $\nu \in \mathcal{P}[0, 1]$,*

$$\forall x \leq \mathbb{E}(\nu), \quad \phi_\nu^*(x) = \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \quad \text{and} \quad \forall x \geq \mathbb{E}(\nu), \quad \phi_\nu^*(x) = \mathcal{L}_{\text{inf}}^{\geq}(x, \nu).$$

Based on this lemma, we have the following rewriting, which is useful to reinterpret the quantities appearing in Theorem 5.3 and Corollary 5.4: $\Phi(\nu', \nu) = \mathcal{L}(\nu', \nu)$ for the model $\mathcal{P}[0, 1]$, i.e.,

$$\inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \{ \phi_{\nu'}^*(x) + \phi_\nu^*(x) \} = \inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \{ \mathcal{L}_{\text{inf}}^{\geq}(x, \nu') + \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \}. \quad (5.22)$$

For canonical one-parameter exponential models \mathcal{D}_{exp} , a slightly weaker version of Lemma 5.5, only holding for x corresponding to expectations in \mathcal{D}_{exp} and provided in page 151 of 5.5.3.151, similarly shows (5.22), i.e., $\Phi = \mathcal{L}$. Conditions on general models for $\Phi = \mathcal{L}$ to hold are discussed in 5.5.3.152.

5.4. Lower Bounds

In most of this section, we restrict our attention to *generic*² K -armed bandit problems $\underline{\nu}$, that are such that $\mu_j \neq \mu_k$ for $j \neq k$. In particular, the best arm $a^*(\underline{\nu})$ is unique.

Definition of a strategy, and of a (doubly-indexed) sequence of strategies. In the fixed-budget setting, a strategy is defined by

- a *sampling rule*, which consists in choosing the arm $A_t \in [K]$ to observe at each time step $t \in [T]$. This arm A_t depends on the previous observations Y_1, \dots, Y_{t-1} , but also possibly on some external randomization that we capture by the random variable U_{t-1} . A_t is thus \mathcal{F}_{t-1} -measurable, where

$$\mathcal{F}_{t-1} \stackrel{\text{def}}{=} \sigma(I_{t-1}), \quad \text{with } I_{t-1} \stackrel{\text{def}}{=} (U_0, Y_1, U_1, Y_2, U_2, \dots, Y_{t-1}, U_{t-1}).$$

I_{t-1} corresponds to the information available at the end of time step $t-1$, i.e., to the history.

- a *decision rule* \hat{a}_T which is \mathcal{F}_T -measurable.

A strategy depends on the budget T and the number K of arms. As we are interested in asymptotic rates when the budget T goes to ∞ , we consider sequences of strategies indexed by T . Our results will depend on assumptions made on the strategies, which will sometimes be stated for doubly-indexed sequences, that is, sequences of strategies indexed by T and K .

²This is probably a new terminology³ for referring to bandit problems with no two same expectations for the distributions over the arms. It comes from measure theory: if expectations were drawn at random according to some diffuse distribution, e.g., a uniform distribution over an interval, or a Gaussian distribution, then, almost surely, no two expectations would be equal.

Outline of this section. As always in lower-bound results, there is a trade-off between how restrictive are the assumptions on the (doubly-indexed) sequences of strategies, and sometimes on the models, and how large the lower bounds are: the more restrictive the assumptions, the larger the lower bounds. We are interested in assumptions on strategies that are natural in the sense that they should be satisfied by *Successive-Rejects*-type strategies. For instance, Theorem 5.14 comes with the least assumptions but provides a bound where there are no divisions by the ranks k of the arms, which Theorems 5.10 and 5.13 do. We may see Theorem 5.10 as a warm-up result: its main aim is to generalize the lower bound by Audibert et al. (2010) to non-parametric models with a (non-constructive) proof that is only a few-line long. Our preferred result is Theorem 5.13, which provides the largest lower bound while putting the heaviest (though natural) constraints on the sequences of strategies.

5.4.1. Common Restriction: Consistence

For our lower bounds, we will consider sequences of strategies, either only indexed by $T \geq 1$ given a value of $K \geq 2$, or doubly indexed by T and K . These sequences will also be assumed to be “reasonable” in the sense below.

Consistent (or exponentially consistent) sequences of strategies. The misidentification probability $\mathbb{P}(\hat{a}_T \neq a^*(\underline{\nu}))$ may vanish asymptotically (and even vanish exponentially fast) for all bandit problems—in not too large a model \mathcal{D} —, as illustrated in Section 5.3. We will therefore only be interested in such sequences of strategies, called (exponentially) consistent. In the sequel and for extra clarity, we index the probabilities by the ambient bandit problem $\underline{\nu}$ considered.

Definition 5.6. [(exponentially) consistant sequence of strategies]

Fix $K \geq 2$. A sequence of strategies indexed by $T \geq 1$ is consistent, respectively, exponentially consistent, on a model \mathcal{D} if for all generic problems $\underline{\nu}$ in \mathcal{D} ,

$$\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \xrightarrow{T \rightarrow +\infty} 0, \quad \text{respectively,} \quad \limsup_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) < 0.$$

By extension, a doubly-indexed sequence of strategies is (exponentially) consistent if for all $K \geq 2$, the associated sequences of strategies are so.

The fundamental inequality. The fundamental inequality by Garivier et al. (2019), together with the very definition of consistency, yields in a straightforward manner our building block for lower bounds. Details of the derivation are provided in Section 2.3.3.

Lemma 5.7. Fix $K \geq 2$ and a model \mathcal{D} . Consider a consistent sequence of strategies on \mathcal{D} , and two generic bandit problems $\underline{\nu}$ and $\underline{\lambda}$ in \mathcal{D} such that $a^*(\underline{\lambda}) \neq a^*(\underline{\nu})$. Then

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \limsup_{T \rightarrow +\infty} \sum_{a=1}^K \frac{\mathbb{E}_{\underline{\lambda}}[N_a(T)]}{T} \text{KL}(\lambda_a, \nu_a).$$

Note that the proof reveals that the inequality actually holds for limits taken along sub-sequences $(T_n)_{n \geq 1}$. Also, we may only relax the assumptions on the bandit models; e.g., they do not need to be generic and it suffices that they have different unique optimal arms (the notion of a generic bandit problem is defined in the first lines of Section 5.4).

5.4.2. A Lower Bound Revisiting and Extending the one by Audibert et al. (2010)

The focus of this subsection is to establish the lower bound (5.5), from which we derived the gap-based lower bound (5.6) by Audibert et al. (2010). The lower bound (5.5) is smaller than the lower bound to be exhibited in the next subsection, but it comes with less restrictive assumptions on the behaviors of the sequences of strategies considered.

Firstly, we only consider sequences of strategies —actually, sequences of sampling schemes— that do not pull too often the worst arm, and which we will refer to as being balanced against the worst arm. Successive-Rejects-type strategies sample the worst arm less than other arms in expectations, and hence, are indeed balanced against the worst arm. To define this constraint formally, we denote by $w_*(\underline{\nu})$ the index of the unique worst arm of a generic bandit problem $\underline{\nu}$.

Definition 5.8. [balanced against the worst arm sequence]

A doubly-indexed sequence of strategies is *balanced against the worst arm* on a model \mathcal{D} if for all $K \geq 2$, for all generic K -armed bandit problems $\underline{\nu}$ in \mathcal{D} ,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \mathbb{E}_{\underline{\nu}}[N_{w_*(\underline{\nu})}(T)] \leq \frac{1}{K}.$$

A second constraint is related to bandit sub-problems. We say that $\underline{\nu}'$ is a sub-problem of a K -armed bandit problem $\underline{\nu}$ if $\underline{\nu}' = (\nu_a)_{a \in \mathcal{A}}$ for a subset $\mathcal{A} \subseteq [K]$ of cardinality greater than or equal to 2; we denote by $\underline{\nu}' \subseteq \underline{\nu}$ this fact. We say in addition that $\underline{\nu}'$ and $\underline{\nu}$ feature the same optimal arm if $\nu'_{a^*(\underline{\nu}')} = \nu_{a^*(\underline{\nu})}$. It should be easier to identify the best arm in $\underline{\nu}'$ than in $\underline{\nu}$, in the sense below, and this defines the fact that a strategy cleverly exploits pruning of sub-optimal arms. Again, Successive-Rejects-type strategies naturally satisfy this constraint.

Definition 5.9. [clever exploitation of the pruning of sub-optimal arms]

A doubly-indexed sequence of strategies *cleverly exploits pruning of sub-optimal arms* on a model \mathcal{D} if for all generic bandit problems $\underline{\nu}$ in \mathcal{D} with $K \geq 2$ arms, for all sub-problems $\underline{\nu}' \subseteq \underline{\nu}$ featuring the same optimal arm,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq \liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}'}(\hat{a}_T \neq a^*(\underline{\nu}')).$$

We use again the order statistics $\mu_{w_*(\underline{\nu})} = \mu_{(K)} < \mu_{(K-1)} < \dots < \mu_{(1)} = \mu_{a^*(\underline{\nu})}$.

Theorem 5.10. Fix a model \mathcal{D} . Consider a doubly-indexed sequence of strategies that is consistent, balanced against the worst arm on \mathcal{D} , and that cleverly exploits the pruning of sub-optimal arms on \mathcal{D} . For all generic bandit problems $\underline{\nu}$ in \mathcal{D} with $K \geq 2$ arms,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \frac{\mathcal{L}_{\inf}^<(\mu_{(k)}, \nu^*)}{k}.$$

Proof. The proof consists of two steps. The first step is to prove that for a generic bandit problem $\underline{\nu}$ in \mathcal{D} with $K \geq 2$ arms, we have,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \frac{\mathcal{L}_{\inf}^<(\mu_{(K)}, \nu^*)}{K}. \quad (5.23)$$

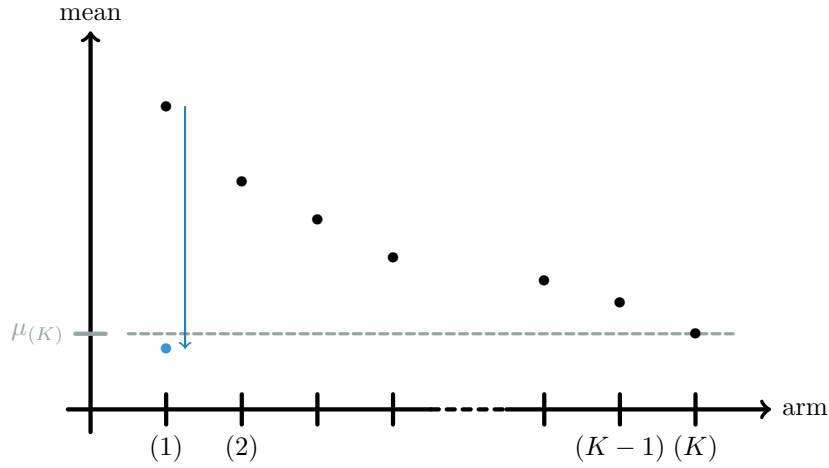


Figure 5.4: The [alternative \$\underline{\lambda}\$](#) of bandit problem $\underline{\nu}$ is obtained by modifying only the best arm of $\underline{\nu}$.

In the second step, we use this lower bound and the very definition of the clever exploitation of the pruning of sub-optimal arms to get the claimed bound.

Step 1: lower bound (5.23). We follow a well-established methodology and consider an alternative bandit problem only differing from $\underline{\nu}$ at one arm, namely, at the best arm. To do so, we set some distribution $\zeta \in \mathcal{D}$ with $\mathbb{E}(\zeta) < \mu_{(K)}$, if some exists, and define the bandit problem $\underline{\lambda} = (\lambda_1, \dots, \lambda_K)$ as

$$\lambda_a = \begin{cases} \zeta & \text{if } a = a^*(\underline{\nu}), \\ \nu_a & \text{if } a \neq a^*(\underline{\nu}). \end{cases}$$

Observe (see Figure 5.4 that $\underline{\lambda}$ is also a generic bandit problem in \mathcal{D} , that $a^*(\underline{\nu})$ is the worst arm in $\underline{\lambda}$ (and also that the second best arm of $\underline{\nu}$ is the optimal arm in $\underline{\lambda}$, but we will not use this specific fact). Therefore, Lemma 5.7 yields, as $\underline{\lambda}$ and $\underline{\nu}$ only differ at arm $a^*(\underline{\nu})$,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\lambda}}[N_{a^*(\underline{\nu})}(T)]}{T} \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu^*),$$

where we recall that $\nu^* = \nu_{a^*(\underline{\nu})}$. Given that $a^*(\underline{\nu})$ is the worst arm of $\underline{\lambda}$, and since by assumption, the sequence of strategies is balanced against the worst arm,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \mathbb{E}_{\underline{\lambda}}[N_{a^*(\underline{\nu})}(T)] \leq \frac{1}{K},$$

proving that

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \frac{\text{KL}(\zeta, \nu^*)}{K}.$$

The claimed inequality (5.23) follows from taking the supremum on the right-hand side over distributions $\zeta \in \mathcal{D}$ with $\mathbb{E}(\zeta) < \mu_{(K)}$.

Step 2: clever exploitation of pruning. For each $k \in \{2, \dots, K-1\}$, define $\underline{\nu}'_{1:k}$ as the subproblem of $\underline{\nu}$ obtained by keeping the k best arms and dropping the $K-k$ worst arms. Use the definition of clever exploitation of pruning of sub-optimal arms and apply (5.23) to $\underline{\nu}'_{1:k}$ to get

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq \liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}'_{1:k}}(\hat{a}_T \neq a^*(\underline{\nu}'_{1:k})) \geq - \frac{\mathcal{L}_{\text{inf}}^<(\mu_{(k)}, \nu^*)}{k}.$$

Taking the maximum of all lower bounds exhibited as k varies between 2 and K , we proved the claimed result. \square

5.4.3. A Larger Lower Bound, for a more Restrictive Class of Strategies

In this section, we derive a slightly stronger version of the lower bound (5.4). This lower bound is larger than the bound exhibited in the previous subsection but relies on stronger assumptions on the strategies considered. Namely, we introduce an assumption of monotonicity, which extends Definition 5.8 to provide frequency constraints on each arm $a \in [K]$.

Definition 5.11. [monotonous sequence of strategies]

Fix $K \geq 2$. A sequence of strategies is *monotonous* on a model \mathcal{D} if for all generic problems $\underline{\nu}$ in \mathcal{D} , for all arms $a \in \{1, \dots, K\}$,

$$\limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\nu}}[N_{(a)}(T)]}{T} \leq \frac{1}{a},$$

where arms are ordered such that $\mu_{(1)} > \mu_{(2)} > \dots > \mu_{(K)}$.

This condition is satisfied as soon as a given arm is not pulled more often, asymptotically and on average, than better-performing arms (note that Definition 5.11 is slightly weaker than this). Successive-Rejects-type strategies naturally satisfy this requirement.

We also rely on the following assumption on the model \mathcal{D} , which essentially indicates that there is “no gap” in \mathcal{D} . Once again, the model $\mathcal{P}[0, 1]$ and canonical one-parameter exponential models \mathcal{D}_{exp} all satisfy this mild requirement (see Section 5.5.4 for the immediate details).

Definition 5.12. A model \mathcal{D} is *normal* if for all $\nu \in \mathcal{D}$, for all $x \geq \mathbb{E}(\nu)$,

$$\begin{aligned} \forall \varepsilon > 0, \quad \mathcal{L}_{\text{inf}}^>(x, \nu) &\stackrel{\text{def}}{=} \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) > x\} \\ &= \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } x + \varepsilon > \mathbb{E}(\zeta) > x\}. \end{aligned}$$

Theorem 5.13. Fix $K \geq 2$ and a normal model \mathcal{D} . Consider a sequence of strategies that is consistent and *monotonous* on \mathcal{D} . For all generic bandit problems $\underline{\nu}$ in \mathcal{D} ,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \min_{2 \leq j \leq k} \inf_{x \in [\mu_{(j)}, \mu_{(j-1)})} \left\{ \frac{\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})}{j-1} + \frac{\mathcal{L}_{\text{inf}}^<(x, \nu^*)}{j} \right\}.$$

Proof. We fix a generic bandit $\underline{\nu}$ in \mathcal{D} and consider the following sets of alternative bandit problems, indexed by triplets (k, j, x) satisfying $2 \leq k \leq K$ and $2 \leq j \leq k$, as well as $x \in [\mu_{(j)}, \mu_{(j-1)})$:

$$\text{Alt}_{k,j,x}(\underline{\nu}) = \left\{ \underline{\lambda} \text{ in } \mathcal{D} : \mathbb{E}(\lambda_{(1)}) < x < \mathbb{E}(\lambda_{(k)}) < \mu_{(j-1)} \text{ and } \lambda_a = \nu_a \text{ for } a \notin \{(1), (k)\} \right\};$$

in particular, an alternative problem $\underline{\lambda}$ in $\text{Alt}_{k,j,x}(\underline{\nu})$ only differ from the original bandit problem $\underline{\nu}$ at the best arm (1) and at the k -th best arm (k). Given $x \in [\mu_{(j)}, \mu_{(j-1)})$ and $\mathbb{E}(\lambda_{(1)}) < x$, arm (1) is at best the j -th best arm of $\underline{\lambda}$, but it can be possibly worse. Similarly, the same condition on x and the fact that $x < \mathbb{E}(\lambda_{(k)})$ implies that arm (k) is exactly the $j-1$ -th best arm of $\underline{\lambda}$. Both facts are illustrated in Figure 5.5.

Thus, by monotonicity of the strategy,

$$\limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\lambda}}[N_{(k)}(T)]}{T} \leq \frac{1}{j-1} \quad \text{and} \quad \limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\lambda}}[N_{(1)}(T)]}{T} \leq \frac{1}{j}.$$

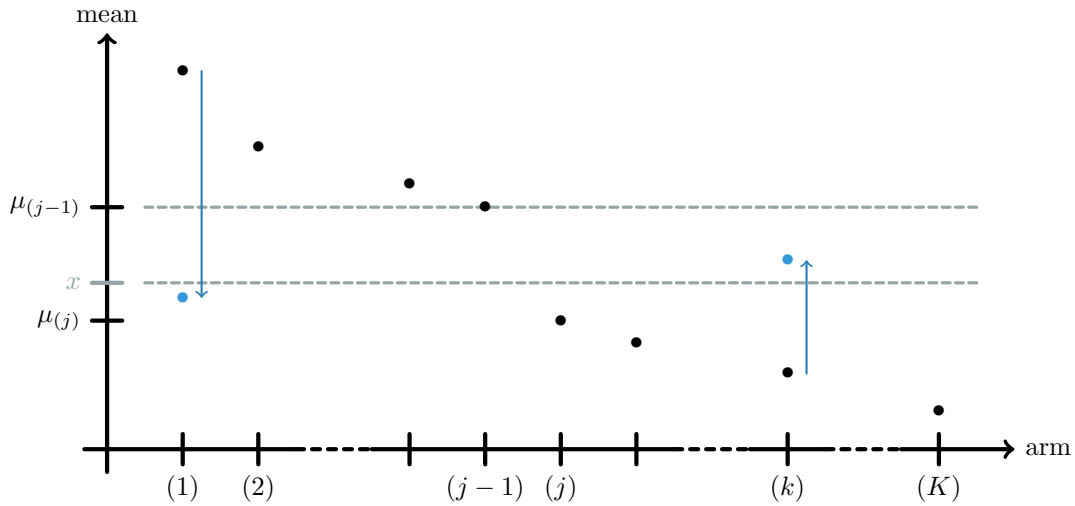


Figure 5.5: Original bandit problem $\underline{\nu}$ (in dark) and modifications made to arms (1) and (k) to obtain an alternative bandit problem $\underline{\lambda} \in \text{Alt}_{k,j,x}(\underline{\nu})$ (in blue): in $\underline{\lambda}$, arm (k) is the $j-1$ -th best arm, while arm (1) = $a^*(\underline{\nu})$ is at best the j -th best arm.

Given that the optimal arm in $\underline{\lambda}$ is different from the optimal arm (1) of $\underline{\nu}$, Lemma 5.7 may be applied; together with the two upper bounds above, it yields

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \left(\frac{\text{KL}(\lambda_{(k)}, \nu_{(k)})}{j-1} + \frac{\text{KL}(\lambda_{(1)}, \nu^*)}{j} \right).$$

We can now take the infimum over all bandit problems $\underline{\lambda} \in \text{Alt}_{k,j,x}(\underline{\nu})$ and obtain the following lower bound, where we define a quantity $\mathcal{I}_{k,j,x}(\underline{\nu})$:

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \inf_{\underline{\lambda} \in \text{Alt}_{k,j,x}(\underline{\nu})} \left\{ \frac{\text{KL}(\lambda_{(k)}, \nu_{(k)})}{j-1} + \frac{\text{KL}(\lambda_{(1)}, \nu^*)}{j} \right\} \stackrel{\text{def}}{=} -\mathcal{I}_{k,j,x}(\underline{\nu}).$$

We prove below that

$$\mathcal{I}_{k,j,x}(\underline{\nu}) = \frac{\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})}{j-1} + \frac{\mathcal{L}_{\text{inf}}^<(x, \nu^*)}{j}, \quad (5.24)$$

from which the lower bound claimed in Theorem 5.13 will follow, by taking the supremum of $-\mathcal{I}_{k,j,x}(\underline{\nu})$ first over $x \in [\mu_{(j)}, \mu_{(j-1)})$, then the maximum over $2 \leq j \leq k$, and finally, the maximum over $2 \leq k \leq K$.

We now prove (5.24). The infimum over $\underline{\lambda} \in \text{Alt}_{k,j,x}(\underline{\nu})$ may be split into two separate infima, respectively over $\lambda_{(k)}$ and $\lambda_{(1)}$; given that each term of the sum of KL only depends either on $\lambda_{(k)}$, or on $\lambda_{(1)}$, but not on both, we may write

$$\begin{aligned} \mathcal{I}_{k,j,x}(\underline{\nu}) &= \inf_{\substack{\lambda_{(1)}, \lambda_{(k)} \in \mathcal{D}: \\ \mathbb{E}(\lambda_{(1)}) < x \\ x < \mathbb{E}(\lambda_{(k)}) < \mu_{(j-1)}}} \left\{ \frac{\text{KL}(\lambda_{(k)}, \nu_{(k)})}{j-1} + \frac{\text{KL}(\lambda_{(1)}, \nu^*)}{j} \right\} \\ &= \frac{1}{j-1} \underbrace{\inf_{\substack{\lambda_{(k)} \in \mathcal{D}: \\ x < \mathbb{E}(\lambda_{(k)}) < \mu_{(j-1)}}} \text{KL}(\lambda_{(k)}, \nu_{(k)})}_{=\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})} + \frac{1}{j} \underbrace{\inf_{\substack{\lambda_{(1)} \in \mathcal{D}: \\ \mathbb{E}(\lambda_{(1)}) < x}} \text{KL}(\lambda_{(1)}, \nu^*)}_{=\mathcal{L}_{\text{inf}}^<(x, \nu^*)}, \end{aligned}$$

where we obtained $\mathcal{L}_{\text{inf}}^<(x, \nu^*)$ by definition while we relied on the normality of the model (Definition 5.12) to obtain $\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})$. We did so with $\varepsilon = \mu_{(j-1)} - x$, which is indeed positive as we considered $x < \mu_{(j-1)}$. \square

5.4.4. A General Lower Bound, Valid for any Strategy

The previous subsections illustrated what may be achieved under restrictions —though natural restrictions— on the classes of strategies considered. For the sake of completeness, we also provide a lower bound relying on no other restriction than consistency; it extends the lower bound (5.9) exhibited by [Kaufmann et al. \(2016\)](#) for $K = 2$ arms, and is formulated in terms of $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^>$.

Theorem 5.14. *Fix $K \geq 2$ and a model \mathcal{D} . Consider a consistent sequence of strategies on \mathcal{D} . For all generic bandit problems $\underline{\nu}$ in \mathcal{D} ,*

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \min_{k \neq a^*(\underline{\nu})} \inf_{x \in [\mu_k, \mu^*]} \max\{\mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu^*)\}.$$

Proof. Let $\underline{\nu}$ be a generic bandit problem. We fix $k \neq a^*(\underline{\nu})$ and $x \in [\mu_k, \mu^*]$, and prove that

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \max\{\mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu^*)\},$$

from which the stated lower bound follows, by taking suprema. To do so, we consider the set of alternative bandit problems

$$\text{Alt}_{k,x}(\underline{\nu}) = \left\{ \underline{\lambda} \text{ in } \mathcal{D} : \mathbb{E}(\lambda_{a^*(\underline{\nu})}) < x < \mathbb{E}(\lambda_k) \text{ and } \lambda_a = \nu_a \text{ for } a \notin \{a^*(\underline{\nu}), k\} \right\};$$

it is composed of bandit problems, only differing from $\underline{\nu}$ at arms $a^*(\underline{\nu})$ and k , and for which arm k is better than arm $a^*(\underline{\nu})$, with associated expectations separated by x . In particular, the optimal arm in $\underline{\lambda}$ is different from the optimal arm $a^*(\underline{\nu})$ of $\underline{\nu}$. Lemma 5.7 may therefore be applied; it states that

$$\begin{aligned} & \liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \\ & \geq - \limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\lambda}}[N_k(T)]}{T} \text{KL}(\lambda_k, \nu_k) + \frac{\mathbb{E}_{\underline{\lambda}}[N_{a^*(\underline{\nu})}(T)]}{T} \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \\ & \geq - \max \left\{ \text{KL}(\lambda_k, \nu_k), \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \right\}, \end{aligned}$$

where we used, for the second inequality, the crude upper bound $N_k(T) + N_{a^*(\underline{\nu})}(T) \leq T$. Taking the supremum of the obtained lower bound over all $\underline{\lambda} \in \text{Alt}_{k,x}(\underline{\nu})$ leads to the following inequality, where we define the short-hand notation $\mathcal{I}_{k,x}(\underline{\nu})$:

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq - \inf_{\underline{\lambda} \in \text{Alt}_{k,x}(\underline{\nu})} \max \left\{ \text{KL}(\lambda_k, \nu_k), \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \right\} \stackrel{\text{def}}{=} -\mathcal{I}_{k,x}(\underline{\nu}).$$

The proof is concluded below by showing that $\mathcal{I}_{k,x}(\underline{\nu}) = \max\{\mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu^*)\}$.

As in the proof of Theorem 5.13, we use a separation of the infima, in the abstract form, for two functions f and g ,

$$\inf_{u,v} \max\{f(u), g(v)\} = \max\left\{ \inf_u f(u), \inf_v g(v) \right\}.$$

Here, by definition of $\text{Alt}_{k,x}(\underline{\nu})$,

$$\begin{aligned} \mathcal{I}_{k,x}(\underline{\nu}) &= \inf_{\substack{\lambda_{a^*(\underline{\nu})}, \lambda_k \in \mathcal{D} \\ \mathbb{E}(\lambda_{a^*(\underline{\nu})}) < x \\ \mathbb{E}(\lambda_k) > x}} \max \left\{ \text{KL}(\lambda_k, \nu_k), \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \right\} \\ &= \max \left\{ \inf_{\substack{\lambda_k \in \mathcal{D} \\ \mathbb{E}(\lambda_k) > x}} \text{KL}(\lambda_k, \nu_k), \inf_{\substack{\lambda_{a^*(\underline{\nu})} \in \mathcal{D} \\ \mathbb{E}(\lambda_{a^*(\underline{\nu})}) < x}} \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \right\} \\ &= \max \left\{ \mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu_{a^*}) \right\}, \end{aligned}$$

which concludes the proof. \square

5.5. Technical Details

5.5.1. Properties of the $\mathcal{L}_{\text{inf}}^<$, $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^{\geq}$

General Properties

We present here general properties, that hold for all models \mathcal{D} , of the $\mathcal{L}_{\text{inf}}^<$, $\mathcal{L}_{\text{inf}}^>$, $\mathcal{L}_{\text{inf}}^>$, and $\mathcal{L}_{\text{inf}}^{\geq}$ quantities. We state some properties for $\mathcal{L}_{\text{inf}}^<$, that all also hold for $\mathcal{L}_{\text{inf}}^>$; the corresponding properties for $\mathcal{L}_{\text{inf}}^{\geq}$ and $\mathcal{L}_{\text{inf}}^{\leq}$ are deduced by symmetry.

The function $\mathcal{L}_{\text{inf}}^<(\cdot, \nu)$ is non-increasing and satisfies $\mathcal{L}_{\text{inf}}^<(x, \nu) = 0$ for all $x > \mathbb{E}(\nu)$, as can be seen by taking $\zeta = \nu$. Also, whenever \mathcal{D} is convex, the function $\mathcal{L}_{\text{inf}}^<$ is jointly convex over $\mathbb{R} \times \mathcal{D}$, as indicated in the lemma below. In particular, $x \mapsto \mathcal{L}_{\text{inf}}^<(x, \nu)$ is continuous on the interior of its domain (the set where it takes finite values).

Lemma 5.15. *When \mathcal{D} is a convex model, all four functions $\mathcal{L}_{\text{inf}}^<$, $\mathcal{L}_{\text{inf}}^>$, $\mathcal{L}_{\text{inf}}^>$, and $\mathcal{L}_{\text{inf}}^{\geq}$ are jointly convex over $\mathbb{R} \times \mathcal{D}$.*

Proof. We provide the proof for $\mathcal{L}_{\text{inf}}^<$, and it may be adapted in a straightforward manner for the other functions.

We set two distributions ν and ν' of \mathcal{D} , two expectation levels μ and μ' in \mathbb{R} , and a weight $\lambda \in (0, 1)$. We want to prove that

$$\mathcal{L}_{\text{inf}}^<(\lambda\mu + (1-\lambda)\mu', \lambda\nu + (1-\lambda)\nu') \leq \lambda\mathcal{L}_{\text{inf}}^<(\mu, \nu) + (1-\lambda)\mathcal{L}_{\text{inf}}^<(\mu', \nu'). \quad (5.25)$$

The desired inequality holds whenever $\mathcal{L}_{\text{inf}}^<(\mu, \nu) = +\infty$ or $\mathcal{L}_{\text{inf}}^<(\mu', \nu') = +\infty$. Otherwise, assuming that both $\mathcal{L}_{\text{inf}}^<(\mu, \nu)$ and $\mathcal{L}_{\text{inf}}^<(\mu', \nu')$ are finite, we set $\delta > 0$ (which we will ultimately let converge to 0) and pick ζ and ζ' in \mathcal{D} such that $\mathbb{E}(\zeta) < \mu$ and $\mathbb{E}(\zeta') < \mu'$, as well as

$$\text{KL}(\zeta, \nu) \leq \mathcal{L}_{\text{inf}}^<(\mu, \nu) + \delta \quad \text{and} \quad \text{KL}(\zeta', \nu') \leq \mathcal{L}_{\text{inf}}^<(\mu', \nu') + \delta.$$

Then, by joint convexity of the Kullback-Leibler divergence:

$$\begin{aligned} \lambda\mathcal{L}_{\text{inf}}^<(\mu, \nu) + (1-\lambda)\mathcal{L}_{\text{inf}}^<(\mu', \nu') + \delta &\geq \lambda\text{KL}(\zeta, \nu) + (1-\lambda)\text{KL}(\zeta', \nu') \\ &\geq \text{KL}(\lambda\zeta + (1-\lambda)\zeta', \lambda\nu + (1-\lambda)\nu') \\ &\geq \mathcal{L}_{\text{inf}}^<(\lambda\mu + (1-\lambda)\mu', \lambda\nu + (1-\lambda)\nu'), \end{aligned}$$

where for the last inequality, we used the definition of $\mathcal{L}_{\text{inf}}^<$ as an infimum and the fact that by convexity, the distribution $\lambda\zeta + (1-\lambda)\zeta'$ belongs to \mathcal{D} , with expectation larger than $\lambda\mu + (1-\lambda)\mu'$. The desired convexity inequality (5.25) follows by letting $\delta \rightarrow 0$. \square

Specific Properties for the Model $\mathcal{P}[0, 1]$

We study some properties of the $\mathcal{L}_{\text{inf}}^<$, $\mathcal{L}_{\text{inf}}^{\leq}$, $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^{\geq}$ quantities for the model $\mathcal{P}[0, 1]$ of all distributions over $[0, 1]$. Similar analysis might be done for, e.g., exponential models, as explained in Section 5.5.3.

Since we are considering distributions over the interval $[0, 1]$, the data-processing inequality (2.11) for Kullback-Leibler divergences ensures that for all $\zeta \in \mathcal{P}[0, 1]$,

$$\text{KL}(\zeta, \nu) \geq \text{KL}(\text{Ber}(\mathbb{E}(\zeta)), \text{Ber}(\mathbb{E}(\nu))) \geq 2(\mathbb{E}(\zeta) - \mathbb{E}(\nu))^2,$$

where $\text{Ber}(p)$ denotes the Bernoulli distribution with parameter p and where we applied Pinsker's inequality for Bernoulli distributions. Therefore, taking the infimum over distributions $\zeta \in \mathcal{P}[0, 1]$ with $\mathbb{E}(\zeta) < x$,

$$\forall x \leq \mathbb{E}(\nu), \quad \mathcal{L}_{\text{inf}}^<(x, \nu) \geq 2(\mathbb{E}(\nu) - x)^2. \quad (5.26)$$

We denote by $m(\nu) = \min(\text{Supp}(\nu)) \geq 0$ the minimum of the closed support $\text{Supp}(\nu)$ of ν ; that is, $m(\nu)$ is the largest value such that $\text{Supp}(\nu) \subseteq [m(\nu), 1]$. We will refer to $m(\nu)$ as the lower end of the support of ν . Though we will not need it immediately, we also define the upper end of the support of ν as $M(\nu) = \max(\text{Supp}(\nu)) \leq 1$; by symmetry, it will be considered when studying $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^{\geq}$ instead of $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^{\leq}$.

The lemma below states that the functions $\mathcal{L}_{\text{inf}}^<(\cdot, \nu)$ and $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$ coincide, except maybe at $m(\nu)$. One may wonder what happens at $x = m(\nu)$. We denote by $\nu\{m(\nu)\}$ the probability mass assigned by ν to the point $m(\nu)$. It follows from the second part of the lemma below that $\mathcal{L}_{\text{inf}}^<(m(\nu), \nu) = \mathcal{L}_{\text{inf}}^{\leq}(m(\nu), \nu)$ if and only if $\{m(\nu)\}$ is not an atom of ν .

Lemma 5.16. *We consider the model $\mathcal{D} = \mathcal{P}[0, 1]$. The function $\mathcal{L}_{\text{inf}}^<(\cdot, \nu)$ is continuous on the interval $(m(\nu), +\infty)$. We also have, on the one hand,*

$$\forall \mu \neq m(\nu), \quad \mathcal{L}_{\text{inf}}^<(\mu, \nu) = \mathcal{L}_{\text{inf}}^{\leq}(\mu, \nu), \quad (5.27)$$

and on the other hand, at $\mu = m(\nu)$,

$$\log \frac{1}{\nu\{m(\nu)\}} = \mathcal{L}_{\text{inf}}^{\leq}(m(\nu), \nu) \leq \mathcal{L}_{\text{inf}}^<(m(\nu), \nu) = +\infty. \quad (5.28)$$

Analogous results hold for $\mathcal{L}_{\text{inf}}^>(\cdot, \nu)$, $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu)$, and $M(\nu)$.

Proof. To prove (5.27), we first identify the interior of the domain of $\mathcal{L}_{\text{inf}}^<$.

Distributions ζ such that $\mathbb{E}(\zeta) < m(\nu)$ cannot be absolutely continuous with respect to ν ; otherwise, they would also give a null probability to values strictly smaller than $m(\nu)$, which contradicts the assumption $\mathbb{E}(\zeta) < m(\nu)$. Hence $\text{KL}(\zeta, \nu) = +\infty$ for these distributions. It follows that $\mathcal{L}_{\text{inf}}^<(\mu, \nu) = \mathcal{L}_{\text{inf}}^{\leq}(\mu, \nu) = +\infty$ for $\mu < m(\nu)$; we note in passing that we also have $\mathcal{L}_{\text{inf}}^<(m(\nu), \nu) = +\infty$.

For $\mu > m(\nu)$, we take $\varepsilon > 0$ with $m(\nu) + \varepsilon < \mu$ and have, by definition of the support of a measure, that $[m(\nu), m(\nu) + \varepsilon]$ has a positive ν -measure denoted by κ . The distribution ζ given by ν conditioned to the interval $[m(\nu), m(\nu) + \varepsilon]$ is absolutely continuous with respect to ν , with density $d\zeta/d\nu = 1/\kappa$ on $[m(\nu), m(\nu) + \varepsilon]$, and 0 elsewhere; therefore, $\text{KL}(\zeta, \nu) = \log(1/\kappa) < +\infty$ and $\mathcal{L}_{\text{inf}}^<(\mu, \nu) < +\infty$.

The interior of the domain of $\mu \mapsto \mathcal{L}_{\text{inf}}^<(\mu, \nu)$ is therefore $(m(\nu), +\infty)$, and we recall that $\mathcal{L}_{\text{inf}}^<(\cdot, \nu)$ is continuous on this interval. We fix some $\mu > m(\nu)$. For all $\varepsilon > 0$, by the very

definitions of all quantities as infima of nested sets, we have

$$\mathcal{L}_{\inf}^{\leq}(\mu - \varepsilon, \nu) \leq \mathcal{L}_{\inf}^{\leq}(\mu, \nu) \leq \mathcal{L}_{\inf}^{\leq}(\mu, \nu).$$

Letting $\varepsilon \rightarrow 0$, we get, by a sandwich argument, that $\mathcal{L}_{\inf}^{\leq}(\mu, \nu) = \mathcal{L}_{\inf}^{\leq}(\mu, \nu)$. This concludes the proof of (5.27).

We turn our attention to (5.28). We already showed above that $\mathcal{L}_{\inf}^{\leq}(m(\nu), \nu) = +\infty$. Now, to compute $\mathcal{L}_{\inf}^{\leq}(\mu, \nu)$, we wonder which are the distributions ζ that are absolutely continuous with respect to ν , and thus, give a null probability to values strictly smaller than $m(\nu)$, and are also such that $\mathbb{E}(\zeta) \leq m(\nu)$: at most one such distribution exists, the Dirac mass at $m(\nu)$, denoted by $\delta_{m(\nu)}$. We then distinguish the cases $\nu\{m(\nu)\} > 0$ and $\nu\{m(\nu)\} = 0$ to establish, respectively, the equalities

$$\mathcal{L}_{\inf}^{\leq}(m(\nu), \nu) = \text{KL}(\delta_{m(\nu)}, \nu) = \log \frac{1}{\nu\{m(\nu)\}} \quad \text{and} \quad \mathcal{L}_{\inf}^{\leq}(m(\nu), \nu) = +\infty = \log \frac{1}{\nu\{m(\nu)\}}.$$

In both cases, the first equality in (5.28) is proved, which concludes the proof. \square

We also have the following result, which is the most important and useful one, as it discussed the quantity that appears in the upper bounds on the average log-probability of misidentification of the optimal arm; see Corollary 5.4 together with Lemma 5.5.

Lemma 5.17. *Let $\nu, \nu' \in \mathcal{P}[0, 1]$ with $\mu = \mathbb{E}(\nu) > \mathbb{E}(\nu') = \mu'$. Then*

$$\inf_{x \in [\mu', \mu]} \mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu') = \inf_{x \in [\mu', \mu]} \mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu')$$

if and only if either $m(\nu) \neq M(\nu')$ or $\nu\{m(\nu)\} \times \nu'\{M(\nu')\} = 0$.

Remark. In other words, the only case for which the two infima differ is when $m(\nu) = M(\nu')$, i.e., the upper end of the support of ν' equals the lower end of the support of ν , and both ν and ν' admit this common value as an atom.

Proof. The first lines of the proof of Lemma 5.16 show that $\mathcal{L}_{\inf}^{\leq}(x, \nu) = \mathcal{L}_{\inf}^{\leq}(x, \nu) = +\infty$ for $x < m(\nu)$. We can symmetrically show that $\mathcal{L}_{\inf}^{\geq}(x, \nu') = \mathcal{L}_{\inf}^{\geq}(x, \nu') = +\infty$ for $x > M(\nu')$. Therefore, $\mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu')$ and $\mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu')$ are infinite whenever x lies outside of $[m(\nu), M(\nu')]$. This implies that

$$\begin{aligned} \inf_{x \in [\mu', \mu]} \mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu') &= \inf_{x \in [\mu', \mu] \cap [m(\nu), M(\nu')]} \mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu') \\ \text{and} \quad \inf_{x \in [\mu', \mu]} \mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu') &= \inf_{x \in [\mu', \mu] \cap [m(\nu), M(\nu')]} \mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu'). \end{aligned}$$

We now split the analysis according to how large the interval \mathcal{I} is, where

$$\mathcal{I} = [\mu', \mu] \cap [m(\nu), M(\nu')] = [\max\{\mu', m(\nu)\}, \min\{\mu, M(\nu')\}].$$

Case 1: \mathcal{I} is empty. In that case, the two infima are over an empty set and both equal $+\infty$.

Case 2: \mathcal{I} has a non-empty interior. When $a \neq b$, the infimum of a convex function over a closed interval $[a, b]$ equals the infimum over (a, b) , whether the function takes finite or infinite values at a and b . Now, the interior of $\mathcal{I} = [a, b]$ equals

$$(a, b) = (\max\{\mu', m(\nu)\}, \min\{\mu, M(\nu')\}) = (\mu', \mu) \cap (m(\nu), M(\nu'))$$

and does not contain neither $m(\nu)$ nor $M(\nu')$. By Lemma 5.16, the functions $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$ and $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$ coincide on $\mathbb{R} \setminus \{m(\nu)\}$. It may be similarly shown that $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu')$ and $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu')$ coincide on $\mathbb{R} \setminus \{M(\nu')\}$. In particular, the functions $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu) + \mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu')$ and $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu) + \mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu')$ coincide on the interior of \mathcal{I} . Their infima over the interior of \mathcal{I} , which, by convexity, are equal to the infima over \mathcal{I} , are therefore equal.

Case 3: \mathcal{I} is a singleton. This case arises if and only if $m(\nu) = M(\nu')$, as by definition, $m(\nu) \leq \mu$ and $M(\nu') \geq \mu'$. We then have $\mathcal{I} = \{m(\nu)\} = \{M(\nu')\}$, and both infima are equal to the values of the sums at $m(\nu) = M(\nu')$. By Lemma 5.16 and by symmetric results for $\mathcal{L}_{\text{inf}}^{\geq}$ and $\mathcal{L}_{\text{inf}}^{\leq}$, on the one hand,

$$\mathcal{L}_{\text{inf}}^{\leq}(m(\nu), \nu) = \mathcal{L}_{\text{inf}}^{\geq}(M(\nu'), \nu') = +\infty,$$

and on the other hand,

$$\mathcal{L}_{\text{inf}}^{\leq}(m(\nu), \nu) + \mathcal{L}_{\text{inf}}^{\geq}(M(\nu'), \nu') = \log \frac{1}{\nu\{m(\nu)\}} + \log \frac{1}{\nu'\{M(\nu')\}}.$$

We get the desired equality if and only if either $\nu\{m(\nu)\} = 0$ or $\nu'\{M(\nu')\} = 0$. \square

5.5.2. Reminder: the Cramér-Chernoff Bound

In this section, we recall the statement of the highly classical Cramér-Chernoff bound: with the notation introduced in Section 5.3, for an N -sample X_1, \dots, X_N , distributed according to ν and of average denoted by \bar{X}_N ,

$$\forall x \leq \mathbb{E}(\nu), \quad \mathbb{P}(\bar{X}_N \leq x) \leq \exp(-N \phi_\nu^*(x)), \quad (5.29)$$

$$\text{and} \quad \forall x \geq \mathbb{E}(\nu), \quad \mathbb{P}(\bar{X}_N \geq x) \leq \exp(-N \phi_\nu^*(x)). \quad (5.30)$$

Such a classical result would in principle not require to be proved here. However, it turns out that we will re-use parts of this proof in later proofs, like the application 5.31 of Jensen's inequality or the variations of ϕ_ν^* discussed at the end of this section. This is why, despite all, we now prove (5.29)–(5.30).

Proof. For all $\lambda < 0$, by Markov's inequality first and then by independence,

$$\begin{aligned} \mathbb{P}(\bar{X}_N \leq x) &= \mathbb{P}(e^{\lambda \bar{X}_N} \geq e^{\lambda x}) \leq e^{-\lambda x} \mathbb{E}[e^{\lambda \bar{X}_N}] = e^{-\lambda x} \left(\mathbb{E}[e^{\lambda X_1/N}] \right)^N \\ &= \exp(-\lambda x + N \phi_\nu(\lambda/N)) = \exp(-N(\lambda' x - \phi_\nu(\lambda'))), \end{aligned}$$

where $\lambda' = \lambda/N$. The bound also holds for $\lambda = \lambda' = 0$ given that $\phi_\nu(0) = 0$. Optimizing over $\lambda \leq 0$ (or, equivalently, over $\lambda' \leq 0$), we proved so far

$$\mathbb{P}(\bar{X}_N \leq x) \leq \exp\left(-N \sup_{\lambda \leq 0} \{\lambda x - \phi_\nu(\lambda)\}\right).$$

Now, by Jensen's inequality,

$$\forall \lambda \in \mathbb{R}, \quad \phi_\nu(\lambda) = \log \mathbb{E}[e^{\lambda X}] \geq \lambda \mathbb{E}[X] = \lambda \mathbb{E}(\nu); \quad (5.31)$$

therefore, for $x \leq \mathbb{E}(\nu)$,

$$\forall \lambda \geq 0, \quad \lambda x - \phi_\nu(\lambda) \leq \lambda(x - \mathbb{E}(\nu)) \leq 0.$$

In particular,

$$0 = -\phi_\nu(0) \leq \sup_{\lambda \leq 0} \{\lambda x - \phi_\nu(\lambda)\} = \sup_{\lambda \in \mathbb{R}} \{\lambda x - \phi_\nu(\lambda)\} \stackrel{\text{def}}{=} \phi_\nu^*(x). \quad (5.32)$$

This concludes the proof of (5.29). The bound (5.30) follows by symmetry. \square

We also note, in passing, that Jensen's inequality entails, for $x = \mathbb{E}(\nu)$, that

$$\forall \lambda \in \mathbb{R}, \quad \lambda \mathbb{E}(\nu) - \phi_\nu(\lambda) \leq \lambda(\mathbb{E}(\nu) - \mathbb{E}(\nu)) = 0,$$

thus showing that $\phi_\nu^*(\mathbb{E}(\nu)) = 0$. The property (5.32) and its counterpart for $x \geq \mathbb{E}(\nu)$ and $\lambda \geq 0$ actually show that ϕ_ν^* is non-increasing on $(-\infty, \mathbb{E}(\nu)]$ and non-decreasing on $[\mathbb{E}(\nu), +\infty)$.

5.5.3. Proofs and Details for Section 5.4.2: Rewriting of Φ as \mathcal{L}

We use the notation of Sections 5.2.1 and 5.3 and discuss conditions on models guaranteeing that $\Phi = \mathcal{L}$, i.e., that (5.22) holds. We do so for $\mathcal{D} = \mathcal{P}[0, 1]$ in and for canonical one-parameter exponential families. Based on these two examples, we provide a set of conditions for general models at the end of the section. A building block of these results is that for all these models \mathcal{D} , the functions $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$ and $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu)$ dominate the Fenchel-Legendre transform ϕ_ν^* defined in (5.15); we prove first.

All proofs of this section are immediate adaptations of a rather standard result, stated, among others, but in a slightly different form (and for the model \mathcal{D} of all real-valued distributions with a first moment), by [Boucheron et al. \(2013, Exercise 4.13\)](#).

Remark. This rewriting of $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$ or $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu)$ as ϕ_ν^* claimed, e.g., by Lemma 5.5, can be seen as a counterpart to a similar rewriting of the \mathcal{K}_{inf} as the supremum of a function of $\lambda \in [0, 1]$. More precisely, we recall (see Remark 1) that the \mathcal{K}_{inf} function is defined, for $\nu \in \mathcal{P}[0, 1]$ and $x \in [0, 1]$, as

$$\mathcal{K}_{\text{inf}}(\nu, x) = \inf\{\text{KL}(\nu, \zeta) : \zeta \in \mathcal{P}[0, 1] \text{ s.t. } \mathbb{E}(\zeta) > x\},$$

and [Honda and Takemura \(2015, Theorem 2\)](#)—see also [Garivier et al., 2022](#), Lemma 18—show that

$$\mathcal{K}_{\text{inf}}(\nu, x) = \sup_{0 \leq \lambda \leq 1} \mathbb{E} \left[\log \left(1 - \lambda \frac{X - x}{1 - x} \right) \right],$$

where X is a random variable distributed according to ν . In both cases, for $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$ or $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu)$, and for \mathcal{K}_{inf} , being able to rewrite the infimum of Kullback-Leibler divergences as a supremum is not unexpected: a given Kullback-Leibler divergence can be formulated as a supremum, see (5.33), and equalities between $\inf \sup$ and $\sup \inf$ holds under suitable assumptions (provided, e.g., by Sion's lemma).

$\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$ and $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu)$ dominate ϕ_ν^*

This domination is a consequence of a variational formula (5.33) for the Kullback-Leibler divergences.

Lemma 5.18. *For all models \mathcal{D} containing distributions with finite first moments, for all distributions $\nu \in \mathcal{D}$,*

$$\forall x \leq \mathbb{E}(\nu), \quad \phi_\nu^*(x) \leq \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \quad \text{and} \quad \forall x \geq \mathbb{E}(\nu), \quad \phi_\nu^*(x) \leq \mathcal{L}_{\text{inf}}^{\geq}(x, \nu).$$

Proof. We rely on a key variational formula for the Kullback-Leibler divergence, see [Boucheron et al. \(2013, Corollary 4.15\)](#): for all distributions ν, ν' over \mathbb{R} ,

$$\begin{aligned} \text{KL}(\nu', \nu) &= \sup \left\{ \mathbb{E}_{\nu'}[Y] - \log \mathbb{E}_\nu[e^Y] : \text{r.v. } Y \in \mathbb{L}^1(\nu') \text{ s.t. } \mathbb{E}_\nu[e^Y] < +\infty \right\}, \\ &= \sup \left\{ \mathbb{E}_{\nu'}[Y] - \log \mathbb{E}_\nu[e^Y] : \text{r.v. } Y \in \mathbb{L}^1(\nu') \right\}, \end{aligned} \quad (5.33)$$

where the supremum is over random variables $Y : \mathbb{R} \rightarrow \mathbb{R}$ with a finite first moment with respect to ν' , and where \mathbb{E}_ν and $\mathbb{E}_{\nu'}$ indicate that expectations are relative to ν and ν' , respectively. In particular, when ν and ν' lie in \mathcal{D} , they admit finite first moments, hence all random variables of the form $Y = \lambda \text{id}_{\mathbb{R}}$ are ν' -integrable, where $\text{id}_{\mathbb{R}}$ denotes the identity function over \mathbb{R} and where $\lambda \in \mathbb{R}$. We have $\mathbb{E}_{\nu'}[Y] = \lambda \mathbb{E}(\nu')$. A consequence of (5.33) and of the definition (5.15) of ϕ_ν^* is therefore that

$$\text{KL}(\nu', \nu) \geq \sup_{\lambda \in \mathbb{R}} \left\{ \lambda \mathbb{E}(\nu') - \log \mathbb{E}_\nu [e^{\lambda \text{id}_{\mathbb{R}}}] \right\} = \phi_\nu^*(\mathbb{E}(\nu')). \quad (5.34)$$

Using the variations of ϕ_ν^* indicated at the end of Section 5.5.2, we see that

$$\phi_\nu^*(\mathbb{E}(\nu')) \geq \phi_\nu^*(x) \quad \text{when } \mathbb{E}(\nu') \leq x \leq \mathbb{E}(\nu) \quad \text{or} \quad \mathbb{E}(\nu') \geq x \geq \mathbb{E}(\nu).$$

Therefore, taking an infimum in (5.34) yields, when $x \leq \mathbb{E}(\nu)$,

$$\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = \inf \{ \text{KL}(\nu', \nu) : \mathbb{E}(\nu') \leq x \} \geq \phi_\nu^*(x),$$

and similarly for the other claimed inequality. \square

The case of $\mathcal{P}[0, 1]$

In this section, we focus on the model $\mathcal{P}[0, 1]$ and prove that the inequalities of Lemma 5.18 are in fact equalities, as claimed by Lemma 5.5. This yields, in particular, the target equality (5.22), as discussed after the statement of Lemma 5.5.

Before proving Lemma 5.5, note that it holds for all $x \in \mathbb{R}$, that is, even outside of the $[0, 1]$ interval, though the proof reveals that when x is smaller than the lower end $m(\nu)$ of the support of ν , we actually have $\phi_\nu^*(x) = \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = +\infty$. The counterpart statement $\phi_\nu^*(x) = \mathcal{L}_{\text{inf}}^{\geq}(x, \nu) = +\infty$ holds for x larger than the upper end $M(\nu)$ of the support of ν . The pieces of notation $m(\nu)$ and $M(\nu)$ were formally defined in Section 5.5.1.

Proof of Lemma 5.5. Note first that by Lemma 5.18, it suffices to prove that

$$\forall x \leq \mathbb{E}(\nu), \quad \phi_\nu^*(x) \geq \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \quad \text{and} \quad \forall x \geq \mathbb{E}(\nu), \quad \phi_\nu^*(x) \geq \mathcal{L}_{\text{inf}}^{\geq}(x, \nu).$$

We only deal with the first inequality, namely $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \leq \phi_\nu^*(x)$ for $x \leq \mathbb{E}(\nu)$, as the other one may be obtained by symmetric arguments.

In the case $x = \mathbb{E}(\nu)$, we have $\phi_\nu^*(\mathbb{E}(\nu)) = 0$, as stated at the end of Section 5.5.2, and $\mathcal{L}_{\text{inf}}^{\leq}(\mathbb{E}(\nu), \nu) = 0$, as can be seen by taking $\zeta = \nu$ in the infimum defining $\mathcal{L}_{\text{inf}}^{\leq}$. We therefore only consider $x < \mathbb{E}(\nu)$ in the sequel. We will rely on the standard fact that, by Hölder's inequality, the logarithmic moment-generating function

$$\phi_\nu : \lambda \in \mathbb{R} \mapsto \log \mathbb{E}_\nu [e^{\lambda \text{id}_{[0,1]}}],$$

is convex, where $\text{id}_{[0,1]}$ denotes the identity function on $[0, 1]$. Also, by two applications of a standard theorem of differentiation under the integral, given that ν is supported by $[0, 1]$, we have that ϕ_ν is continuously differentiable over \mathbb{R} , with derivative

$$\phi'_\nu : \lambda \in \mathbb{R} \mapsto \frac{\mathbb{E}_\nu [\text{id}_{[0,1]} e^{\lambda \text{id}_{[0,1]}}]}{\mathbb{E}_\nu [e^{\lambda \text{id}_{[0,1]}}]}.$$

By convexity of ϕ_ν , this derivative is non-decreasing. Therefore, the limit of ϕ'_ν at $-\infty$ exists; we denote it by ℓ and have that a priori $\ell \in \{-\infty\} \cup \mathbb{R}$. We now prove that actually,

$$\ell \stackrel{\text{def}}{=} \lim_{\lambda \rightarrow -\infty} \phi'_\nu(\lambda) = m(\nu). \quad (5.35)$$

On the one hand, by definition of $m(\nu)$, we have $\text{id}_{[0,1]} \geq m(\nu)$ ν -a.s., which entails $\phi'_\nu(\lambda) \geq m(\nu)$ for all $\lambda \in \mathbb{R}$, and hence, $\ell \geq m(\nu)$. On the other hand, as ϕ'_ν is non-decreasing, it is always larger than its limit ℓ at $-\infty$:

$$\forall \lambda \in \mathbb{R}, \quad \phi'_\nu(\lambda) \geq \ell, \quad \text{thus,} \quad \mathbb{E}_\nu \left[(\text{id}_{[0,1]} - \ell) e^{\lambda \text{id}_{[0,1]}} \right] \geq 0, \quad (5.36)$$

$$\text{or} \quad \mathbb{E}_\nu \left[(\text{id}_{[0,1]} - \ell) e^{\lambda(\text{id}_{[0,1]} - \ell)} \right] \geq 0. \quad (5.37)$$

The last inequality and limit arguments as $\lambda \rightarrow -\infty$ impose that $\text{id}_{[0,1]} - \ell \geq 0$ ν -a.s., which in turn entails that $\ell \leq m(\nu)$. This concludes the proof of (5.35).

The various properties exhibited above for ϕ_ν , including the fact that the derivative ϕ'_ν takes values in $[m(\nu), +\infty)$, entail that the function

$$\Lambda : \lambda \in \mathbb{R} \mapsto \lambda x - \phi_\nu(\lambda)$$

is concave, continuously differentiable, with a non-increasing derivative Λ' taking values in the interval $(-\infty, x - m(\nu)]$ and with limit $x - m(\nu)$ at $-\infty$.

We split the analysis of the case $x < \mathbb{E}(\nu)$ into three sub-cases, depending on the respective positions of x and $m(\nu)$, and recall that we want to show that $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \leq \phi_\nu^*(x)$.

Case 1: $x > m(\nu)$. By Jensen's inequality (5.31) and given that we consider $x < \mathbb{E}(\nu)$, the limit of Λ at $+\infty$ equals $-\infty$. The limit of Λ at $-\infty$ also equals $-\infty$, as the derivative Λ' has limit $x - m(\nu) > 0$ at $-\infty$. By concavity of Λ and the fact that Λ' is continuous, this implies the existence of some $\lambda^* \in \mathbb{R}$ such that

$$\Lambda'(\lambda^*) = x - \phi'_\nu(\lambda^*) = 0 \quad \text{and} \quad \phi_\nu^*(x) = \sup_{\lambda \in \mathbb{R}} \{ \Lambda(\lambda) \} = \Lambda(\lambda^*).$$

Denoting by ζ_{λ^*} the distribution absolutely continuous with respect to ν with density

$$\frac{d\zeta_{\lambda^*}}{d\nu} = \frac{e^{\lambda^* \text{id}_{[0,1]}}}{\mathbb{E}_\nu [e^{\lambda^* \text{id}_{[0,1]}}]} = e^{\lambda^* \text{id}_{[0,1]} - \phi_\nu(\lambda^*)},$$

we have $\mathbb{E}_{\zeta_{\lambda^*}}[\text{id}_{[0,1]}] = \mathbb{E}(\zeta_{\lambda^*}) = \phi'_\nu(\lambda^*) = x$. Therefore, by definition of $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu)$ and of the Kullback-Leibler divergence,

$$\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \leq \text{KL}(\zeta_{\lambda^*}, \nu) = \mathbb{E}_{\zeta_{\lambda^*}} \left[\log \frac{d\zeta_{\lambda^*}}{d\nu} \right] = \lambda^* \mathbb{E}_{\zeta_{\lambda^*}}[\text{id}_{[0,1]}] - \phi_\nu(\lambda^*) = \Lambda(\lambda^*) = \phi_\nu^*(x).$$

Case 2: $x = m(\nu)$. In that case, $\Lambda' \rightarrow 0$ at $-\infty$ and Λ' is non-increasing, thus $\Lambda' \leq 0$ on \mathbb{R} and Λ is non-increasing on \mathbb{R} . Thus,

$$\phi_\nu^*(m(\nu)) = \sup_{\lambda \in \mathbb{R}} \{ \Lambda(\lambda) \} = \lim_{\lambda \rightarrow -\infty} \Lambda(\lambda) = \lim_{\lambda \rightarrow -\infty} -\log \mathbb{E}_\nu \left[e^{\lambda(\text{id}_{[0,1]} - m(\nu))} \right].$$

By monotone convergence based on $\text{id}_{[0,1]} - m(\nu) \geq 0$ ν -a.s.,

$$\lim_{\lambda \rightarrow -\infty} -\log \mathbb{E}_\nu \left[e^{\lambda(\text{id}_{[0,1]} - m(\nu))} \right] = -\log \nu \{ m(\nu) \},$$

whether $\nu \{ m(\nu) \}$ is positive or null. Moreover, Lemma 5.16 states that

$$\mathcal{L}_{\text{inf}}^{\leq}(m(\nu), \nu) = -\log \nu \{ m(\nu) \}.$$

We therefore have $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = \phi_\nu^*(x)$ in this case.

Case 3: $x < m(\nu)$. In that case, as $\Lambda' \rightarrow x - m(\nu) < 0$ at $-\infty$, we get that $\Lambda \rightarrow +\infty$ at $-\infty$, thus $\phi_\nu^*(x) = \sup \Lambda = +\infty$. Now, no distribution $\zeta \in \mathcal{P}[0, 1]$ with $\mathbb{E}(\zeta) \leq x$, if some exists, can be absolutely continuous with respect to ν ; indeed, $x < m(\nu)$ imposes that ζ puts some probability mass to the left of the support of ν . Therefore, $\text{KL}(\zeta, \nu) = +\infty$. All in all, $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu)$ appears as the infimum of either an empty set or of $+\infty$ values, so that $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = +\infty$. In this case as well, $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = \phi_\nu^*(x)$, both being equal to $+\infty$. \square

The case of canonical one-parameter exponential models \mathcal{D}_{exp}

In this section, we show that the target equality (5.22) is satisfied by so-called canonical one-parameter exponential families \mathcal{D}_{exp} . We recall that those models are defined at the beginning of Section 2.2.3 and follow the notation of that Section: ρ is the reference measure, b the normalizing function, Θ is the natural space parameter, and $\mathcal{M} = (\mu_-, \mu_+)$ is the open interval of the expectations of distributions in \mathcal{D}_{exp} .

The mean-parametrized Kullback-Leibler divergence of the model is denoted by d and defined on $\mathcal{M} \times \mathcal{M}$, where. In the following, we extend d to $\mathbb{R} \times \mathbb{R}$ by $+\infty$ values outside of $\mathcal{M} \times \mathcal{M}$.

A direct application of the continuity and monotonicity properties of d is that all functions $\mathcal{L}_{\text{inf}}^<$, $\mathcal{L}_{\text{inf}}^{\leq}$, $\mathcal{L}_{\text{inf}}^>$, $\mathcal{L}_{\text{inf}}^{\geq}$ coincide with d in the sense of the stated equalities (5.12) and (5.13). Indeed and for instance, we have, for $\nu \in \mathcal{D}_{\text{exp}}$ and $x \leq \mathbb{E}(\nu)$ with $x \in \mathcal{M}$:

$$\mathcal{L}_{\text{inf}}^<(x, \nu) = \inf_{\mu < x} \{d(\mu, \nu)\} = \lim_{\substack{\mu \rightarrow x \\ \mu < x}} d(\mu, \nu) = d(x, \nu).$$

When $x \notin \mathcal{M}$, by the convention on the infimum of an empty set, $\mathcal{L}_{\text{inf}}^<(x, \nu) = +\infty$, while by our definition of d outside $\mathcal{M} \times \mathcal{M}$, we also have $d(x, \nu) = +\infty$. But as Lemma 5.20 below illustrates, we will only be interested in the behaviors on $\mathcal{M} \times \mathcal{M}$.

We now state a monotonicity property of the Chernoff-information-type quantity L defined for exponential models in (5.14). This property was referred to in Example 5.1, when indicating that arms can be equivalently ranked in descending expectations or ascending values of $L(\cdot, \mu^*)$.

Lemma 5.19. *Consider a canonical one-parameter exponential family \mathcal{D}_{exp} and fix any $\mu \in \mathcal{M}$. Then $L(\cdot, \mu)$ is non-increasing on $(\mu_-, \mu]$.*

Proof. Fix $\mu_- < \mu_2 \leq \mu_1 \leq \mu$. To get the desired inequality $L(\mu_2, \mu) \geq L(\mu_1, \mu)$, it suffices to show, by (5.14), that

$$\forall y \in [\mu_2, \mu], \quad d(y, \mu_2) + d(y, \mu) \geq \min_{x \in [\mu_1, \mu]} d(x, \mu_1) + d(x, \mu) \stackrel{\text{def}}{=} L(\mu_1, \mu). \quad (5.38)$$

We distinguish two cases. If $\mu_2 \leq \mu_1 \leq y \leq \mu$, then, since $d(y, \cdot)$ is increasing on $(\mu_-, y]$, we have $d(y, \mu_2) \geq d(y, \mu_1)$, from which the inequality (5.38) follows by considering $x = y$. If $\mu_2 \leq y \leq \mu_1 \leq \mu$, then similarly $d(y, \mu) \geq d(\mu_1, \mu)$, which yields

$$\underbrace{d(y, \mu_2) + d(y, \mu)}_{\geq 0} \geq d(\mu_1, \mu) = \underbrace{d(\mu_1, \mu_1)}_{=0} + d(\mu_1, \mu),$$

from which the inequality (5.38) follows by considering $x = \mu_1$. □

A slightly weaker version of Lemma 5.5, sufficient for our purposes. We may now come back to the proof of the target equality (5.22) for canonical one-parameter exponential families. The following slightly weaker version of Lemma 5.5 is enough to yield (5.22), given the rewritings (5.12) and (5.13).

Lemma 5.20. *Consider a canonical one-parameter exponential family $\mathcal{D} = \mathcal{D}_{\text{exp}}$. For all $\nu \in \mathcal{D}_{\text{exp}}$,*

$$\forall x \in \mathcal{M}, \quad \phi_\nu^*(x) = d(x, \mathbb{E}(\nu)).$$

The result of the lemma holds, by conventions, for $x < \mu_-$ or $x > \mu_+$, but does not hold in general for $x \in \{\mu_-, \mu_+\}$.

Proof. By Lemma 5.18, we only need to show that $\phi_\nu^*(x) \geq d(x, \mathbb{E}(\nu))$. Given the definition (5.15) of ϕ_ν^* as a supremum, it suffices to exhibit a $\lambda^* \in \mathbb{R}$ such that

$$d(x, \mathbb{E}(\nu)) = \lambda^* x - \phi_\nu(\lambda^*). \quad (5.39)$$

Let $\theta_1 \in \Theta$ be such that $\nu = \nu_{\theta_1}$ and $\theta_2 = (b')^{-1}(x) \in \Theta$ be such that $\mathbb{E}(\nu_{\theta_2}) = x$. We will prove (5.39) with $\lambda^* = \theta_2 - \theta_1$. Given the closed-form expression of the densities (2.15), the distribution ν_{θ_2} is absolutely continuous with respect to ν_{θ_1} , with density given by $(\theta_2 - \theta_1)\text{id}_{\mathbb{R}} - (b(\theta_2) - b(\theta_1))$. Therefore, by definition of the Kullback-Leibler divergence,

$$\begin{aligned} d(x, \mathbb{E}(\nu)) &= \text{KL}(\nu_{\theta_2}, \nu_{\theta_1}) = \mathbb{E}_{\nu_{\theta_2}} \left[\log \frac{d\nu_{\theta_2}}{d\nu_{\theta_1}} \right] = \mathbb{E}_{\nu_{\theta_2}} \left[(\theta_2 - \theta_1) \text{id}_{\mathbb{R}} - (b(\theta_2) - b(\theta_1)) \right] \\ &= (\theta_2 - \theta_1) \mathbb{E}(\nu_{\theta_2}) - (b(\theta_2) - b(\theta_1)) = \lambda^* x - (b(\theta_2) - b(\theta_1)). \end{aligned} \quad (5.40)$$

To obtain (5.39), it only remains to show that $b(\theta_2) - b(\theta_1) = \phi_\nu(\lambda^*)$. Using the closed-form expressions (2.16) of b at θ_2 and (2.15) of the density at θ_1 , we obtain

$$\begin{aligned} b(\theta_2) &= \log \int_{\mathbb{R}} e^{\theta_2 y} d\rho(y) = b(\theta_1) + \log \int_{\mathbb{R}} e^{(\theta_2 - \theta_1)y} \overbrace{e^{\theta_1 y - b(\theta_1)}}{=d\nu_{\theta_1}(y)=d\nu(y)} d\rho(y) \\ &= b(\theta_1) + \log \int_{\mathbb{R}} e^{\lambda^* y} d\nu(y) = b(\theta_1) + \phi_\nu(\lambda^*), \end{aligned} \quad (5.41)$$

which concludes the proof. \square

Remark. A more direct approach bypassing Lemma 5.18 can be followed with \mathcal{D}_{exp} models, along the following lines. The result (5.41) can be generalized into

$$\forall \theta \in \Theta, \quad \phi_\nu(\theta - \theta_1) = b(\theta) - b(\theta_1). \quad (5.42)$$

As b is differentiable on Θ , the function ϕ_ν is also differentiable; at $\lambda^* = \theta_2 - \theta_1$, we have

$$\phi'_\nu(\lambda^*) = \phi'_\nu(\theta_2 - \theta_1) = b'(\theta_2) = x.$$

Thus, the derivative of the strictly concave function $\Lambda : \lambda \in \mathbb{R} \mapsto \lambda x - \phi_\nu(\lambda)$ vanishes at λ^* , which is therefore the argument of its maximum: $\phi_\nu^*(x) = \Lambda(\lambda^*)$. The closed-form calculation (5.40) and the rewriting (5.42) then lead to Lemma 5.20.

Conditions for general models

In this section, we extend Lemma 5.5, and thus the target equality (5.22), to more general models. We did so by mimicking the proof of Lemma 5.5: the result below can certainly be improved. We extend as follows the definitions of the lower and upper ends $m(\nu)$ and $M(\nu)$ of the closed support $\text{Supp}(\nu)$ of a distribution ν over \mathbb{R} :

$$m(\nu) = \inf(\text{Supp}(\nu)) \in \mathbb{R} \cup \{-\infty\} \quad \text{and} \quad M(\nu) = \sup(\text{Supp}(\nu)) \in \mathbb{R} \cup \{+\infty\}.$$

Lemma 5.21. Consider a model \mathcal{D} containing distributions ν over \mathbb{R} with finite first moments and with exponential moments: $e^{\lambda \text{id}_{\mathbb{R}}} \in \mathbb{L}^1(\nu)$ for all $\lambda \in \mathbb{R}$. Assume that the model \mathcal{D} is stable by exponential reweighting of densities: for all $\nu \in \mathcal{D}$, for all $\lambda \in \mathbb{R}$, the distribution ν_λ with density

$$\frac{d\nu_\lambda}{d\nu} = \frac{e^{\lambda \text{id}_{\mathbb{R}}}}{\mathbb{E}_\nu[e^{\lambda \text{id}_{\mathbb{R}}}]}$$
 with respect to ν (5.43)

also belongs to \mathcal{D} . Assume also that δ_x , the Dirac mass at x , belongs to \mathcal{D} whenever there exists $\nu \in \mathcal{D}$ with $x \in \{m(\nu), M(\nu)\} \cap \mathbb{R}$ and $\nu\{x\} > 0$; put differently, if a distribution $\nu \in \mathcal{D}$ puts some probability mass on an end x of its closed support, then the Dirac mass at x belongs to \mathcal{D} .

Then, for all $\nu \in \mathcal{D}$,

$$\forall x \leq \mathbb{E}(\nu), \quad \phi_\nu^*(x) = \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \quad \text{and} \quad \forall x \geq \mathbb{E}(\nu), \quad \phi_\nu^*(x) = \mathcal{L}_{\text{inf}}^{\geq}(x, \nu).$$

Proof. By symmetry and by Lemma 5.18, we only need to prove that

$$\forall x \leq \mathbb{E}(\nu), \quad \phi_\nu^*(x) \geq \mathcal{L}_{\text{inf}}^{\leq}(x, \nu).$$

For $x = \mathbb{E}(\nu)$, we have $\phi_\nu^*(\mathbb{E}(\nu)) = 0 = \mathcal{L}_{\text{inf}}^{\leq}(\mathbb{E}(\nu), \nu)$, as stated at the end of Section 5.5.2 and by taking $\zeta = \nu$ in the infimum defining $\mathcal{L}_{\text{inf}}^{\leq}$, respectively. Before moving to the case $x < \mathbb{E}(\nu)$, we establish a few properties of ϕ_ν based on the assumptions of Lemma 5.21. All random variables $e^{\lambda \text{id}_{\mathbb{R}}}$ are ν -integrable, for $\lambda \in \mathbb{R}$, which entails, by application of a standard theorem of differentiation under the integral sign together with local domination arguments of the form

$$\forall \lambda \in (\lambda_-, \lambda_+), \quad |\text{id}_{\mathbb{R}} e^{\lambda \text{id}_{\mathbb{R}}}| \leq |\text{id}_{\mathbb{R}}| (e^{\lambda_- \text{id}_{\mathbb{R}}} + e^{\lambda_+ \text{id}_{\mathbb{R}}}) \leq (e^{\text{id}_{\mathbb{R}}} + e^{-\text{id}_{\mathbb{R}}})(e^{\lambda_- \text{id}_{\mathbb{R}}} + e^{\lambda_+ \text{id}_{\mathbb{R}}}),$$

that ϕ_ν is differentiable over \mathbb{R} , with derivative given by

$$\phi'_\nu : \lambda \in \mathbb{R} \mapsto \frac{\mathbb{E}_\nu[\text{id}_{\mathbb{R}} e^{\lambda \text{id}_{\mathbb{R}}}]}{\mathbb{E}_\nu[e^{\lambda \text{id}_{\mathbb{R}}}]}. \tag{5.44}$$

Hölder's inequality still entails that ϕ_ν is convex, thus its derivative ϕ'_ν is non-decreasing; therefore, ϕ'_ν admits a limit $\ell \in \{-\infty\} \cup \mathbb{R}$ at $-\infty$. Actually, we have $\ell = m(\nu)$, as can be seen by combining the following facts. First, by definition, $\text{id}_{\mathbb{R}} \geq m(\nu)$ ν -a.s., thus $\phi'_\nu \geq m(\nu)$, hence $\ell \geq m(\nu)$. As a consequence, if $\ell = -\infty$, then we also have $m(\nu) = -\infty$. Otherwise, if $\ell \in \mathbb{R}$, the same arguments as in (5.36)–(5.37) show that $\text{id}_{\mathbb{R}} - \ell \geq 0$ ν -a.s., i.e., $\ell \leq m(\nu)$.

We may now come back to establishing $\phi_\nu^*(x) \geq \mathcal{L}_{\text{inf}}^{\leq}(x, \nu)$ in the case $x < \mathbb{E}(\nu)$. We consider three sub-cases, depending on the respective positions of x and $m(\nu)$.

Case 1: $x > m(\nu)$. The properties of ϕ_ν ensure, exactly as in Case 1 of the proof of Lemma 5.5, the existence of λ^* such that $\phi'_\nu(\lambda^*) = x$ and $\phi_\nu^*(x) = \lambda^*x - \phi_\nu(\lambda^*)$. Given the assumption (5.43), we may consider the distribution $\nu_{\lambda^*} \in \mathcal{D}$. We note, again exactly as in Case 1 of the proof of Lemma 5.5 and given the closed-form expression (5.44) for ϕ'_ν , that $\mathbb{E}(\nu_{\lambda^*}) = \phi'_\nu(\lambda^*)$, thus $\mathbb{E}(\nu_{\lambda^*}) = x$. Finally, an explicit computation yields

$$\text{KL}(\nu_{\lambda^*}, \nu) = \lambda^* \mathbb{E}(\nu_{\lambda^*}) - \log \mathbb{E}_\nu[e^{\lambda^* \text{id}_{\mathbb{R}}}] = \lambda^* x - \phi_\nu(\lambda^*) = \phi_\nu^*(x).$$

By the defining infimum of $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu)$, we have indeed $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \leq \text{KL}(\nu_{\lambda^*}, \nu) = \phi_\nu^*(x)$.

Case 2: $x = m(\nu)$. In particular, $m(\nu) \in \mathbb{R}$, which allows us to follow the monotone-convergence arguments of Case 2 of the proof of Lemma 5.5 and get the equality $\phi_\nu^*(m(\nu)) = -\log \nu\{m(\nu)\}$. Now, for the second part of this sub-case, we also adapt an argument of the second part of the proof

of Lemma 5.16 (in Section 5.5.1), namely, the fact that either there exists at most one distribution $\zeta \in \mathcal{D}$ absolutely continuous with respect to ν and satisfying $\mathbb{E}(\zeta) \leq m(\nu)$, namely, $\zeta = \delta_{m(\nu)}$, the Dirac mass at $m(\nu)$. The latter is indeed absolutely continuous with respect to ν if and only if $\nu\{m(\nu)\} > 0$. When $\nu\{m(\nu)\} > 0$, we have $\delta_{m(\nu)} \in \mathcal{D}$ by the Dirac assumption of the lemma, so that

$$\mathcal{L}_{\inf}^{\leq}(m(\nu), \nu) = \text{KL}(\delta_{m(\nu)}, \nu) = -\log \nu\{m(\nu)\}.$$

Otherwise, when $\nu\{m(\nu)\} = 0$, the infimum defining $\mathcal{L}_{\inf}^{\leq}(m(\nu), \nu)$ is either over an empty set or of $+\infty$ values, and thus equals $+\infty = -\log \nu\{m(\nu)\}$. In both situations, we obtained $\mathcal{L}_{\inf}^{\leq}(m(\nu), \nu) = \phi_{\nu}^*(m(\nu))$.

Case 3: $x < m(\nu)$. In particular, $m(\nu) \in \mathbb{R}$ in this sub-case as well, which allows us to repeat the exact same arguments as in Case 3 of the proof of Lemma 5.5: we may show that both $\mathcal{L}_{\inf}^{\leq}(x, \nu)$ and $\phi_{\nu}^*(x)$ are equal to $+\infty$. \square

5.5.4. Proof of the Normality of the Models $\mathcal{P}[0, 1]$ and \mathcal{D}_{exp}

In this section, we show that $\mathcal{P}[0, 1]$ and canonical one-parameter exponential models are normal.

Proposition 5.22. *$\mathcal{P}[0, 1]$ is a normal model.*

Proof. We fix $\nu \in \mathcal{P}[0, 1]$, a real $x \geq \mathbb{E}(\nu)$, and $\varepsilon > 0$. Recall the piece of notation $M(\nu)$ for the upper end of the support of ν , as introduced in Section 5.5.1. As in Case 3 of the proof of Lemma 5.5, we note that when $x \geq M(\nu)$, there exists no distribution $\zeta \in \mathcal{P}[0, 1]$ absolutely continuous with respect to ν and such that $\mathbb{E}(\zeta) > x$; hence, both infima in Definition 5.12 equal $+\infty$. We now tackle the case where $\mathbb{E}(\nu) \leq x < M(\nu)$. For all $\delta > 0$, we introduce

$$x'_{\delta} = \min\left\{x + \delta, \frac{x + M(\nu)}{2}\right\} < M(\nu).$$

Case 1 of the proof of Lemma 5.5 and Lemma 5.18 reveal (by symmetry) that for each $\delta > 0$, there exists a distribution $\zeta_{\delta} \in \mathcal{P}[0, 1]$ with expectation x'_{δ} and such that $\mathcal{L}_{\inf}^{\geq}(x'_{\delta}, \nu) = \phi_{\nu}^*(x'_{\delta}) = \text{KL}(\zeta_{\delta}, \nu)$. By Lemma 5.16, $\mathcal{L}_{\inf}^{\geq}(x'_{\delta}, \nu) = \mathcal{L}_{\inf}^{\geq}(x'_{\delta}, \nu)$ and $\mathcal{L}_{\inf}^{\geq}(\cdot, \nu)$ is continuous on $(-\infty, M(\nu))$. Putting all these elements together, we obtain

$$\begin{aligned} \mathcal{L}_{\inf}^{\geq}(x, \nu) &= \lim_{\delta \rightarrow 0} \mathcal{L}_{\inf}^{\geq}(x'_{\delta}, \nu) = \liminf_{\delta \rightarrow 0} \text{KL}(\zeta_{\delta}, \nu) \geq \inf\{\text{KL}(\zeta_{\delta}, \nu) : \delta \in (0, \varepsilon)\} \\ &\geq \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } x + \varepsilon > \mathbb{E}(\zeta) > x\}, \end{aligned}$$

where the first inequality is by the very definition of a \liminf . \square

Proposition 5.23. *All canonical one-parameter exponential models \mathcal{D}_{exp} are normal.*

Proof. The proof consists of rewriting $\mathcal{L}_{\inf}^{\geq}$ as d , as indicated by (5.13), and using the regularity properties for d exhibited in Section 5.5.3 (see page 151). We fix $\nu \in \mathcal{D}_{\text{exp}}$, a real $x \geq \mathbb{E}(\nu)$, and $\varepsilon > 0$. When $x \geq M(\nu)$, the same argument as in the previous proposition shows that both infima equal $+\infty$. For $x < M(\nu)$, we introduce $\delta \in (0, \mu_+ - x)$ and write

$$\begin{aligned} \mathcal{L}_{\inf}^{\geq}(x, \nu) &= d(x, \mathbb{E}(\nu)) = \lim_{\delta \rightarrow 0} d(x + \delta, \mathbb{E}(\nu)) = \inf\{d(x + \delta, \mathbb{E}(\nu)) : \delta \in (0, \varepsilon)\} \\ &= \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } x + \varepsilon > \mathbb{E}(\zeta) > x\}, \end{aligned}$$

where the second and third equalities follow, respectively, by continuity of $d(\cdot, \mathbb{E}(\nu))$ on \mathcal{M} and by the fact that this function is non-decreasing on $(x, \mu_+) \subset [\mathbb{E}(\nu), \mu_+)$, and the final equality is by the rewriting (2.17). \square

5.6. Additional Comments for the Literature Review

This section is devoted to additional discussions concerning the fixed-budget literature. More precisely, we discuss in detail two gap-based lower bounds that we believe are somewhat detached from the spirit of the chapter, namely, the minimax lower bound of [Carpentier and Locatelli \(2016\)](#) in Section 5.6.1 and the Bretagnolle-Huber technique in Section 5.6.2.

5.6.1. The Minimax Lower Bound of [Carpentier and Locatelli \(2016\)](#)

[Carpentier and Locatelli \(2016, Theorem 1\)](#) proved (slightly stronger versions of) the following (non-asymptotic) minimax lower bound. Consider the model $\mathcal{B}_{[1/4, 3/4]}$ of Bernoulli distributions $\text{Ber}(p)$ with parameters $p \in [1/4, 3/4]$. For all sequences of strategies that are consistent on $\mathcal{B}_{[1/4, 3/4]}$, for all $T \geq 0.14 K^4 \log(6KT)$,

$$\exists \underline{\nu} \text{ in } \mathcal{B}_{[1/4, 3/4]}, \quad \frac{1}{T} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq -\frac{400}{\log K} \left(\sum_{a \neq a^*(\underline{\nu})} \frac{1}{\Delta_a^2} \right)^{-1} - \frac{\log 6}{T}, \quad (5.45)$$

where, of course, we may rather use the weaker lower bound based on

$$-\left(\sum_{a \neq a^*(\underline{\nu})} \frac{1}{\Delta_a^2} \right)^{-1} \geq -\min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}.$$

However, the bound (5.45) is different in nature from the lower bounds considered in this chapter, as first and foremost, it only guarantees a $1/\log K$ improvement of the lower bound (5.8) of [Audibert et al. \(2010\)](#) for a single bandit problem $\underline{\nu}$ (actually belonging to a known collection of K bandit problems). This is in strong contrast with the uniform instance-dependent lower bounds presented in this chapter: bounds holding simultaneously for all bandit problems of a given model. Second, the proof of the result (see the simpler proof provided below for Proposition 5.24 stated next) is truly gap-based and does not seem to extend in any obvious way to non-parametric models.

As mentioned above, the proof of (5.45) in [Carpentier and Locatelli \(2016\)](#) uses only K different bandit problems in $\mathcal{B}_{[1/4, 3/4]}$. We may therefore resort to the pigeonhole principle to exchange, in some sense, the “for all $T \geq 0.14K^4 \log(6KT)$ ” and “there exists $\underline{\nu}$ in $\mathcal{B}_{[1/4, 3/4]}$ ” parts. More precisely, we obtain, from (5.45) the following proposition. For the sake of completeness, we provide a self-contained proof of this proposition closely following the original arguments by [Carpentier and Locatelli \(2016\)](#), except for the change-of-measure argument, for which we rather resort to Lemma 5.7. Doing so, we are able to improve the numerical factor 400 that would follow from (5.45) into a smaller factor of 30.

Proposition 5.24. *Fix $K \geq 3$ and consider the model $\mathcal{B}_{[1/4, 3/4]}$ of Bernoulli distributions $\text{Ber}(p)$ with parameters $p \in [1/4, 3/4]$. For all consistent sequences of strategies on $\mathcal{B}_{[1/4, 3/4]}$, there exists an increasing sequence of budgets $(T_n)_{n \geq 1}$ such that*

$$\exists \underline{\nu} \text{ in } \mathcal{B}_{[1/4, 3/4]}, \quad \liminf_{n \rightarrow +\infty} \frac{1}{T_n} \log \mathbb{P}_{\underline{\nu}}(\hat{a}_{T_n} \neq a^*(\underline{\nu})) \geq -\frac{30}{\log K} \left(\sum_{a \neq a^*(\underline{\nu})} \frac{1}{\Delta_a^2} \right)^{-1}. \quad (5.46)$$

Proof. We consider some base Bernoulli bandit problem $\underline{\nu}^{\text{base}} = (\nu_1^{\text{base}}, \dots, \nu_K^{\text{base}})$, where

$$\nu_1^{\text{base}} = \text{Ber}(1/2) \quad \text{and} \quad \forall j \in \{2, \dots, K\}, \quad \nu_j^{\text{base}} = \text{Ber}(p_j),$$

for parameters $p_j \in [1/4, 1/2)$ to be specified later. For each $k \in \{2, \dots, K\}$, we then define the alternative bandit problem $\underline{\nu}^{(k)} = (\nu_1^{(k)}, \dots, \nu_K^{(k)})$ as follows:

$$\nu_j^{(k)} = \begin{cases} \text{Ber}(1 - p_k) & \text{if } j = k, \\ \nu_j^{\text{base}} & \text{if } j \neq k. \end{cases}$$

Given the constraints on the p_j , the unique optimal arm of $\underline{\nu}^{\text{base}}$ is $a^*(\underline{\nu}^{\text{base}}) = 1$, while the unique optimal arm of $\underline{\nu}^{(k)}$ is $a^*(\underline{\nu}^{(k)}) = k$. We introduce, for a given bandit problem $\underline{\nu}$

$$H_\Sigma(\underline{\nu}) \stackrel{\text{def}}{=} \sum_{a \neq a^*(\underline{\nu})} \frac{1}{\Delta_a^2};$$

the right-hand side of (5.46) may be rewritten as $(34/\log K) H_\Sigma(\underline{\nu})^{-1}$. The sub-optimality gaps of the arms of $\underline{\nu}^{\text{base}}$ equal $\Delta_j^{\text{base}} = 1/2 - p_j$ for $j \neq 1$, while the ones of $\underline{\nu}^{(k)}$ equal

$$\begin{aligned} \forall j \neq k, \quad \Delta_j^{(k)} &= 1 - p_k - p_j = (1/2 - p_k) + (1/2 - p_j) = \Delta_k^{\text{base}} + \Delta_j^{\text{base}}, \\ \text{thus} \quad H_\Sigma(\underline{\nu}^{(k)}) &= \sum_{j \neq k} \frac{1}{(\Delta_k^{\text{base}} + \Delta_j^{\text{base}})^2}. \end{aligned} \quad (5.47)$$

The proof is decomposed into two steps. First, we show that for all values of the p_j abiding by the constraints and for all weights u_2, \dots, u_K such that $u_j \geq 0$ for all j and $u_1 + \dots + u_K = 1$, there exists $k^* \in \{2, \dots, K\}$ such that there exists an increasing sequence of budgets $(T_n)_{n \geq 1}$ with

$$\liminf_{n \rightarrow +\infty} \frac{1}{T_n} \log \mathbb{P}_{\underline{\nu}^{(k^*)}}(\hat{a}_{T_n} \neq k^*) \geq -9 u_{k^*} (\Delta_{k^*}^{\text{base}})^2. \quad (5.48)$$

Then, we set specific values of the u_j and p_j to get

$$\forall k \in \{2, \dots, K\}, \quad u_k (\Delta_k^{\text{base}})^2 \leq \frac{10}{3 \log K} H_\Sigma(\underline{\nu}^{(k)})^{-1}. \quad (5.49)$$

Proposition 5.24 follows by combining (5.48) and (5.49).

Part 1: Proof of (5.48). For all $T \geq 1$,

$$\sum_{k=2}^K \frac{\mathbb{E}_{\underline{\nu}^{\text{base}}}[N_k(T)]}{T} \leq 1 = \sum_{k=2}^K u_k;$$

therefore, for all $T \geq 1$, there exists $k_T \in \{2, \dots, K\}$ such that $\mathbb{E}_{\underline{\nu}^{\text{base}}}[N_{k_T}(T)]/T \leq u_{k_T}$. By the pigeonhole principle, there exists $k^* \in \{2, \dots, K\}$ and an (infinite) increasing sequence $(T_n)_{n \geq 1}$ of integers such that $k_{T_n} = k^*$ for all $n \geq 1$. In particular,

$$\limsup_{n \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\nu}^{\text{base}}}[N_{k^*}(T_n)]}{T_n} \leq u_{k^*}.$$

Since $\underline{\nu}^{\text{base}}$ and $\underline{\nu}^{(k^*)}$ only differ at arm k^* , an application of Lemma 5.7 along subsequences (see the initial comments below the lemma) guarantees that

$$\begin{aligned} \liminf_{n \rightarrow +\infty} \frac{1}{T_n} \log \mathbb{P}_{\underline{\nu}^{(k^*)}}(\hat{a}_{T_n} \neq k^*) &\geq - \left(\limsup_{n \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\nu}^{\text{base}}}[N_{k^*}(T_n)]}{T_n} \right) \text{KL}(\text{Ber}(1 - p_{k^*}), \text{Ber}(p_{k^*})) \\ &\geq -u_{k^*} \times 9 (1/2 - p_{k^*})^2 = -9 u_{k^*} (\Delta_{k^*}^{\text{base}})^2, \end{aligned}$$

where, in the last inequality, we used that for all $x \in [1/4, 1/2)$,

$$\text{KL}(\text{Ber}(1-x), \text{Ber}(x)) = (1-x) \log \frac{1-x}{x} + x \log \frac{x}{1-x} \leq 9 \left(\frac{1}{2} - x \right)^2.$$

Part 2: Proof of (5.49). We set, for $j \in \{2, \dots, K\}$,

$$u_j = \frac{U}{(\Delta_j^{\text{base}})^2 H_\Sigma(\nu^{(j)})}, \quad \text{where} \quad U = \left(\sum_{k=2}^K \frac{1}{(\Delta_k^{\text{base}})^2 H_\Sigma(\nu^{(k)})} \right)^{-1}.$$

Then, $u_k (\Delta_k^{\text{base}})^2 = H_\Sigma(\nu^{(k)})^{-1} U$ for all $k \in \{2, \dots, K\}$. To get the desired result, it suffices to guarantee that $U \leq 10/(3 \log K)$. To do so, we consider the same values as in [Carpentier and Locatelli \(2016\)](#) for the p_j , i.e., we set, for $j \in \{2, \dots, K\}$,

$$p_j = \frac{1}{2} - \frac{j}{4K} \quad \text{or, equivalently,} \quad \Delta_j^{\text{base}} = \frac{j}{4K}.$$

We show first that $(\Delta_k^{\text{base}})^2 H_\Sigma(\nu^{(k)}) \leq 2k$, for all $k \in \{2, \dots, K\}$. Indeed, by (5.47) and by lower bounding $\Delta_k^{\text{base}} + \Delta_j^{\text{base}}$ either by Δ_k^{base} or Δ_j^{base} , we get

$$\begin{aligned} (\Delta_k^{\text{base}})^2 H_\Sigma(\nu^{(k)}) &= \sum_{j < k} \frac{(\Delta_k^{\text{base}})^2}{(\Delta_k^{\text{base}} + \Delta_j^{\text{base}})^2} + \sum_{j > k} \frac{(\Delta_k^{\text{base}})^2}{(\Delta_k^{\text{base}} + \Delta_j^{\text{base}})^2} \\ &\leq k-1 + \sum_{j > k} \frac{(\Delta_k^{\text{base}})^2}{(\Delta_j^{\text{base}})^2} = k-1 + \sum_{j > k} \frac{k^2}{j^2} \leq k-1 + k^2 \int_k^K \frac{1}{v^2} dv \leq 2k. \end{aligned}$$

Finally,

$$U \leq \left(\sum_{k=2}^K \frac{1}{2k} \right)^{-1} \leq \left(\int_2^{K+1} \frac{1}{2v} dv \right)^{-1} = 2 (\log(K+1) - \log 2)^{-1} \leq \frac{10}{3 \log K},$$

where the final inequality holds since $K \geq 3$. □

5.6.2. The Bretagnolle-Huber Technique by [Kaufmann et al. \(2016\)](#)

[Kaufmann et al. \(2016, Section 5.2\)](#) provide an interesting series of results relying on the so-called Bretagnolle-Huber inequality recalled below in (5.51); we state one of their lower bounds in Corollary 5.26. But as we argue in this section, the methodology followed seems extremely specific to the case of parametric models where Kullback-Leibler divergences could be controlled (lower bounded and upper bounded) in terms of gaps, like the model \mathcal{D}_{σ^2} of Gaussian distributions with a fixed variance $\sigma^2 > 0$. In particular, we state in Proposition 5.25 what would be the straightforward extension to non-parametric models of the Gaussian results of ([Kaufmann et al., 2016, Section 5.2](#)), and we immediately discuss after this statement why this extension lacks interpretability and interest. Proposition 5.25 considers any sequence of strategies (not necessarily consistent) and provides an asymptotic bound; however, it does not directly control the target probability of error $\mathbb{P}_\nu(\hat{a}_T \neq a^*(\nu))$, but a larger quantity. A proof of Proposition 5.25 is provided at the end of this section.

Proposition 5.25. Fix $K \geq 2$, a model \mathcal{D} , and any sequence of strategies. Let $\underline{\nu}$ be a bandit problem in \mathcal{D} with a unique optimal arm. Consider, for each $k \neq a^*(\underline{\nu})$, a distribution $\zeta_k \in \mathcal{D}$ such that $\mathbb{E}(\zeta_k) > \mu^*$. For $k \neq a^*(\underline{\nu})$, denote by $\underline{\nu}^{(k)}$ the bandit problem obtained from $\underline{\nu}$ by changing the distribution of arm k into ζ_k . For all $T \geq 1$,

$$\frac{1}{T} \log \max \left\{ \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})), \max_{k \neq a^*(\underline{\nu})} \mathbb{P}_{\underline{\nu}^{(k)}}(\hat{a}_T \neq k) \right\} \geq - \left(\sum_{a \neq a^*(\underline{\nu})} \frac{1}{\text{KL}(\nu_a, \zeta_a)} \right)^{-1} - \frac{\log 4}{T}.$$

Lack of interpretability of the bound for general models. To derive an interesting and interpretable bound from this result, one needs to choose carefully the distributions ζ_k . There is a tradeoff between obtaining a large lower bound by choosing ζ_k as close as possible to ν_k in terms of Kullback-Leibler divergences, and controlling the maximum of the misidentification probabilities: when ζ_k gets closer to ν_k while abiding by the constraint $\mathbb{E}(\zeta_k) > \mu^*$, the probability $\mathbb{P}_{\underline{\nu}^{(k)}}(\hat{a}_T \neq k)$ becomes larger, and should even intuitively converge to $1/2$. In any case, the target error $\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu}))$ should get dominated by $\mathbb{P}_{\underline{\nu}^{(k)}}(\hat{a}_T \neq k)$ and the obtained bound is likely to be uninformative on the target error, due to the maximum on the left-hand side. This tradeoff seems to be unsolvable in general unless there exist some specific properties for the Kullback-Leibler divergence of the model, as we illustrate below for a Gaussian model, which was the setting considered by [Kaufmann et al. \(2016, Section 5.2\)](#).

Another intuitive issue with the bound of Proposition 5.25 is that it involves Kullback-Leibler divergences with arguments in reverse order compared to the lower bounds presented in Section 5.4. Indeed, taking the supremum of the lower bound over distributions ζ_k such that $\mathbb{E}(\zeta_k) > \mu^*$ would lead to a complexity in terms of the $\mathcal{K}_{\text{inf}}^>(\nu_k, \mu^*)$, where

$$\mathcal{K}_{\text{inf}}^>(\nu, x) \stackrel{\text{def}}{=} \inf \{ \text{KL}(\nu, \zeta) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) > x \},$$

rather than in terms of the $\mathcal{L}_{\text{inf}}^>(\mu^*, \nu_k)$. Our intuition, given all bounds presented in this chapter, is that the $\mathcal{K}_{\text{inf}}^>(\nu_k, \mu^*)$ would not form the correct notion of complexity for the fixed-budget best-arm identification.

How [Kaufmann et al. \(2016, Section 5.2\)](#) could exploit Proposition 5.25 in the Gaussian case. Yet, in the case of the model \mathcal{D}_{σ^2} of Gaussian distributions with a fixed variance $\sigma^2 > 0$, for which KL is symmetric, Proposition 5.25 admits an interesting corollary, corresponding⁴ to Theorem 16 of [Kaufmann et al. \(2016, Section 5.2\)](#). The corollary actually relies on a strong property of KL in this model: not only is it symmetric, but it only depends on the expectation gaps between its arguments. Namely, for all pairs $\mathcal{N}(\mu, \sigma^2)$ and $\mathcal{N}(\mu', \sigma^2)$ of distributions in \mathcal{D}_{σ^2} , for all $\Delta \in \mathbb{R}$,

$$\text{KL}(\mathcal{N}(\mu, \sigma^2), \mathcal{N}(\mu', \sigma^2)) = \frac{(\mu - \mu')^2}{2\sigma^2} = \text{KL}(\mathcal{N}(\mu + \Delta, \sigma^2), \mathcal{N}(\mu' + \Delta, \sigma^2)). \quad (5.50)$$

We introduce the following short-hand notation:

$$H_{\Sigma}(\underline{\nu}) \stackrel{\text{def}}{=} \sum_{a \neq a^*(\underline{\nu})} \frac{2\sigma^2}{\Delta_a^2}.$$

⁴The maximum of the left-hand side of Corollary 5.26 is present, but somewhat discrete, in the Theorem 16 of [Kaufmann et al. \(2016, Section 5.2\)](#): it corresponds to the ‘‘There exists an alternative bandit problem’’ part of the statement of the latter.

Corollary 5.26. *For all sequences of strategies and for all bandit problems $\underline{\nu}$ in \mathcal{D}_{σ^2} with a unique optimal arm, there exists a set of alternative bandit instances $(\underline{\nu}^{(k)})_{k \neq a^*(\underline{\nu})}$ in \mathcal{D}_{σ^2} , where each $\underline{\nu}^{(k)}$ admits k as a best arm and satisfies $H_{\Sigma}(\underline{\nu}^{(k)}) \leq H_{\Sigma}(\underline{\nu})$, and for which*

$$\frac{1}{T} \log \max \left\{ \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})), \max_{k \neq a^*(\underline{\nu})} \mathbb{P}_{\underline{\nu}^{(k)}}(\hat{a}_T \neq k) \right\} \geq -4 H_{\Sigma}(\underline{\nu})^{-1} - \frac{\log 4}{T}.$$

The proof provided below is highly specific to the Gaussian model and exploits the gap-based rewriting (5.50) of the Kullback-Leibler divergence. The calculations led would only extend to models for which such gap-based rewritings of (upper and lower bounds on) the Kullback-Leibler divergence would be available.

To compare the result of Corollary 5.26 with the bound (5.8) stemming from [Audibert et al. \(2010\)](#), note that

$$H_{\Sigma}(\underline{\nu})^{-1} \geq \frac{2}{\sigma^2} \min_{2 \leq k \leq K} \frac{\Delta_k^2}{k}.$$

Proof. We apply Proposition 5.25 with the distributions $\zeta_k = \mathcal{N}(\mu^* + \Delta_k, \sigma^2)$, for $k \neq a^*(\underline{\nu})$. On the one hand, the bound of Proposition 5.25 involves

$$\sum_{a \neq a^*(\underline{\nu})} \frac{1}{\text{KL}(\nu_a, \zeta_a)} = \sum_{a \neq a^*(\underline{\nu})} \frac{2\sigma^2}{\underbrace{(\mathbb{E}(\nu_a) - \mathbb{E}(\zeta_a))^2}_{\mu^* - \Delta_a \quad \mu^* + \Delta_a}} = \frac{H_{\Sigma}(\underline{\nu})}{4}.$$

On the other hand, for $k \neq a^*(\underline{\nu})$, as the best arm of $\underline{\nu}^{(k)}$ is k , with associated expectation $\mu^* + \Delta_k$,

$$\begin{aligned} H_{\Sigma}(\underline{\nu}^{(k)}) &= \sum_{a \neq k} \frac{2\sigma^2}{(\mu^* + \Delta_k - \mu_a)^2} = \frac{2\sigma^2}{\Delta_k^2} + \sum_{a \notin \{k, a^*(\underline{\nu})\}} \frac{2\sigma^2}{(\mu^* + \Delta_k - \mu_a)^2} \\ &\leq \frac{2\sigma^2}{\Delta_k^2} + \sum_{a \notin \{k, a^*(\underline{\nu})\}} \frac{2\sigma^2}{(\mu^* - \mu_a)^2} = \sum_{a \neq a^*(\underline{\nu})} \frac{2\sigma^2}{\Delta_a^2} = H_{\Sigma}(\underline{\nu}). \end{aligned}$$

These two observations conclude the proof of Corollary 5.26. \square

Proof of Proposition 5.25. We conclude this section with a proof of Proposition 5.25. It relies on the Bretagnolle-Huber inequality ([Bretagnolle and Huber, 1979](#)), which states that, for all $p, q \in [0, 1]$,

$$p + 1 - q \geq \frac{1}{2} \exp\left(-\text{KL}(\text{Ber}(p), \text{Ber}(q))\right). \quad (5.51)$$

Proof. We fix distributions ζ_k abiding by the conditions of the proposition and also fix $T \geq 1$. We will prove below that, for all convex weights $(u_b)_{b \neq a^*(\underline{\nu})}$, i.e., non-negative weights summing up to 1,

$$\frac{1}{T} \log \max \left\{ \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})), \max_{k \neq a^*(\underline{\nu})} \mathbb{P}_{\underline{\nu}^{(k)}}(\hat{a}_T \neq k) \right\} \geq - \max_{b \neq a^*(\underline{\nu})} \{u_b \text{KL}(\nu_b, \zeta_b)\} - \frac{\log 4}{T}, \quad (5.52)$$

from which Proposition 5.25 follows, by optimizing the obtained lower bound, i.e., by taking

$$u_b = \left(\sum_{a \neq a^*(\underline{\nu})} \frac{1}{\text{KL}(\nu_a, \zeta_a)} \right)^{-1} \times \frac{1}{\text{KL}(\nu_b, \zeta_b)}.$$

We now fix convex weights $(u_b)_{b \neq a^*(\underline{\nu})}$ and prove (5.52). As $b \neq a^*(\underline{\nu})$ and b is the unique optimal arm of $\underline{\nu}^{(b)}$, for the first inequality, and by the Bretagnolle-Huber inequality (5.51), for the second inequality,

$$\begin{aligned} \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) + \mathbb{P}_{\underline{\nu}^{(b)}}(\hat{a}_T \neq b) &\geq \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) + \mathbb{P}_{\underline{\nu}^{(b)}}(\hat{a}_T = a^*(\underline{\nu})) \\ &\geq \frac{1}{2} \exp\left(-\text{KL}(\text{Ber}(p_T), \text{Ber}(q_T))\right), \end{aligned}$$

where $p_T \stackrel{\text{def}}{=} \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu}))$ and $q_T \stackrel{\text{def}}{=} \mathbb{P}_{\underline{\nu}^{(b)}}(\hat{a}_T \neq a^*(\underline{\nu}))$. Inequality (2.12) reads, in the present case, as $\underline{\nu}$ and $\underline{\nu}^{(b)}$ only differ at arm b ,

$$\text{KL}(\text{Ber}(p_T), \text{Ber}(q_T)) \leq \mathbb{E}_{\underline{\nu}}[N_b(T)] \text{KL}(\nu_b, \zeta_b).$$

Using $\max\{u, v\} \geq (u + v)/2$ after collecting all bounds obtained so far yields

$$\max\left\{\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})), \mathbb{P}_{\underline{\nu}^{(b)}}(\hat{a}_T \neq b)\right\} \geq \frac{1}{4} \exp\left(-\mathbb{E}_{\underline{\nu}}[N_b(T)] \text{KL}(\nu_b, \zeta_b)\right).$$

We take the maxima over $b \neq a^*(\underline{\nu})$ in both sides, apply logarithms, and conclude the proof of (5.52) by showing that

$$\min_{b \neq a^*(\underline{\nu})} \left\{ \mathbb{E}_{\underline{\nu}}[N_b(T)] \text{KL}(\nu_b, \zeta_b) \right\} \leq \max_{b \neq a^*(\underline{\nu})} \{u_b \text{KL}(\nu_b, \zeta_b)\}. \quad (5.53)$$

Indeed,

$$\sum_{b \neq a^*(\underline{\nu})} \frac{\mathbb{E}_{\underline{\nu}}[N_b(T)]}{T} \leq 1 = \sum_{b \neq a^*(\underline{\nu})} u_b,$$

so that there exists $b^* \neq a^*(\underline{\nu})$ such that $\mathbb{E}_{\underline{\nu}}[N_{b^*}(T)]/T \leq u_{b^*}$. We then have

$$\min_{b \neq a^*(\underline{\nu})} \left\{ \mathbb{E}_{\underline{\nu}}[N_b(T)] \text{KL}(\nu_b, \zeta_b) \right\} \leq u_{b^*} \text{KL}(\nu_{b^*}, \zeta_{b^*}) \leq \max_{b \neq a^*(\underline{\nu})} \{u_b \text{KL}(\nu_b, \zeta_b)\},$$

as desired in (5.53). \square

5.7. Conclusion

In this chapter, we explored the challenging fixed-budget setting of best-arm identification. We introduced new tools, including the information-theoretic quantities $\mathcal{L}_{\text{inf}}^<$, $\mathcal{L}_{\text{inf}}^{\leq}$, $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^{\geq}$, in order to generalize the gap-based bounds of the literature to more general models. Our new general analysis of the Successive-Rejects strategy indicates that the complexity of the fixed-budget setting might be measured by the \mathcal{L} quantity defined as

$$\mathcal{L}(\nu', \nu) = \inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \left\{ \mathcal{L}_{\text{inf}}^{\geq}(x, \nu') + \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \right\}.$$

We then stated several lower bounds depending on various assumptions. Those bounds do not express in terms of the \mathcal{L} quantity —although the bound of Theorem 5.13 invokes quantities that are getting closer to \mathcal{L} —, but we hope that future works could bring new ideas to improve those bounds.

Existence of a complexity As previously explained, there is still a gap of at least a factor $\overline{\log} K$ between the lower and upper bounds presented in this chapter. In fact, we do not know if we can fill this gap and hence prove the existence of a complexity like in the fixed-confidence setting. Recent works (see Komiyama et al., 2022; Degenne, 2023) have been focusing on proving that such optimal complexity does not exist for large enough models (see Section 2.3.5 for discussions).

CHAPTER 6

Asymptotically Optimal Adaptive Top-Two Algorithms in the Fixed-Confidence Setting

This chapter studies adaptive top-two algorithms for a general exponential model in the fixed-confidence setting. It gathers preliminary works that have not been submitted to date. Mainly two conjectures are stated, supported by numerical experiments and elements of proof. Firstly, we obtain new procedures to compute optimal weight vectors that appear to be numerically efficient. Secondly, we give a proof structure to generalize the asymptotic optimality of adaptive top-two algorithms obtained for a Gaussian model by [You et al. \(2023\)](#).

Contents

1	Introduction	162
2	A Fixed Point Property	165
1	The Transformation	166
2	New Empirical Optimal-Weights Procedures	167
3	A New Sampling Rule for top-two Algorithms	169
3	Asymptotically Optimal Adaptive Algorithms	174
1	Step 1: A Sufficient Condition	175
2	Step 2: Tracking versus Sampling	176
3	Step 3: Sufficient Exploration	177
4	Step 4: A Useful Relationship	181
5	Step 5: Quick Convergence of $p_a(t)/p_{a^*}(t)$	182
6	Conclusion	184
4	Side Note: On the Asymptotic Optimality of Track-and-Stop	184
5	Conclusion	186

6.1. Introduction

We consider here the problem of best-arm identification with a fixed-confidence (see Section 2.2 for motivations about this setting and for more details concerning the following introduction). Given some confidence parameter $\delta \in (0, 1)$, our aim is to find a strategy that minimizes the expected number of samplings $\mathbb{E}_{\underline{\nu}}[\tau_\delta]$ among δ -correct strategies, i.e., strategies such that for all bandit problems $\underline{\nu}$,

$$\mathbb{P}_{\underline{\nu}}(\tau_\delta < +\infty, \hat{a}_{\tau_\delta} \neq a^*(\underline{\nu})) \leq \delta.$$

When considering an exponential model \mathcal{D}_{exp} (see the reminder on page 43), we recall that bandit problems $\underline{\nu}$ are characterized by their mean vector $\underline{\mu}$, and a notion of asymptotic optimality (when δ goes to 0) was introduced for the identification problem by [Garivier and Kaufmann \(2016\)](#). They first proved that all δ -correct strategies satisfy, for all bandit problems $\underline{\mu}$ in \mathcal{D}_{exp} ,

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_\delta]}{\log \frac{1}{\delta}} \geq T(\underline{\mu}), \quad (6.1)$$

where $T(\underline{\mu})$ denotes the *characteristic time* of $\underline{\mu}$, defined as

$$T(\underline{\mu})^{-1} \stackrel{\text{def}}{=} \sup_{\underline{v} \in \Sigma_K} \inf_{\underline{\lambda} \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a), \quad (6.2)$$

where d denotes the mean-parameterized Kullback-Leibler divergence of the model \mathcal{D}_{exp} ,

$$\Sigma_K = \left\{ \underline{v} \in [0, 1]^K : v_1 + \dots + v_K = 1 \right\} \quad \text{and} \quad \text{Alt}(\underline{\mu}) = \left\{ \underline{\lambda} \text{ in } \mathcal{D}_{\text{exp}} : a^*(\underline{\lambda}) \neq a^*(\underline{\mu}) \right\}.$$

Additionally, they designed *Track-and-Stop*, the first *asymptotically optimal strategy*, for which the asymptotic upper bound matches the above lower bound:

$$\forall \underline{\mu} \text{ in } \mathcal{D}_{\text{exp}}, \quad \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_\delta]}{\log \frac{1}{\delta}} \leq T(\underline{\mu}).$$

The existence of such strategies ensures the asymptotic tightness of lower bound (6.1).

Computing the solution $\underline{w}(\underline{\mu})$. The information-theoretic analysis of [Garivier and Kaufmann \(2016\)](#) also highlights the nature of optimal sampling strategies: whatever the value of the risk δ , they should sample the arms with frequencies proportional to $\underline{v} = \underline{w}(\underline{\mu})$, the (unique and well-defined) maximizer of optimization problem (6.2), called the *optimal weight vector*. To achieve this, the *Track-and-Stop* sampling rule estimates these proportions by computing, at each time step t , the optimal weight vector $\underline{w}(\hat{\underline{\mu}}(t))$ of the current empirical mean $\hat{\underline{\mu}}(t)$. Indeed, they proved that, with the knowledge of $\underline{\mu}$, solving the optimization problem (6.2) reduces to determining the root of a one-variable increasing function. By applying a bisection method, one may then compute $\underline{w}(\underline{\mu})$ with arbitrary precision. Yet, it can be interesting to obtain more efficient methods to compute this optimal weight vector, as we did for a Gaussian model in Chapter 3.

The shortcomings of *Track-and-Stop*. Computing the solution of optimization problem (6.2) at each time step is not very efficient in terms of computational cost. This is one of the shortcomings of *Track-and-Stop* (see Section 2.2.5), we may also cite, among others, the requirement of forced exploration to ensure that all arms are pulled sufficiently. Some improvements were proposed: for example, [Ménard \(2019\)](#) and [Wang et al. \(2021\)](#) proved that it is not necessary to solve the optimization problem at every time step. Instead, they perform a single gradient step in every round, which enables asymptotic optimality of computationally more efficient algorithms. Another direction of simplification is the study of top-two algorithms (see [Russo, 2016](#)) that we present in depth below.

The importance of sampling rules. The performance of a strategy highly depends on both its sampling and stopping rules. Yet, the analysis of [Garivier and Kaufmann \(2016\)](#) proved (see Theorem 2.9) that the δ -correctness of a strategy can be obtained, whatever the sampling rule is, using the Global-Likelihood-Ratio stopping rule with a carefully designed threshold:

$$\beta(t, \delta) \stackrel{\text{def}}{=} \log \frac{Rt^\alpha}{\delta}, \quad (6.3)$$

where $\alpha > 1$ and R is a constant depending on α and K . As it seems that the stopping rule cannot be significantly improved, the performance of a strategy should mostly be determined by its sampling rule.

top-two algorithms. top-two strategies are promising sets of strategies that come with simple sampling rules: at time step t , the algorithm chooses the next arm to sample A_t between two arms, namely a *leader* L_t and a *challenger* C_t . We recall the general structure of the sampling rule in Algorithm 17. An example of a natural leader is to choose the arm with the current best empirical mean (Algorithm 18), while the challenger might be the arm minimizing its *transportation cost*, defined in Equation (6.6), with the leader, or a penalized version of that cost which encourages exploration (Algorithm 19). See also Section 2.2.7 (and references therein) for additional leader and challenger procedures.

Algorithm 17: top-two sampling rule at time step $t > K$

Input: history of observations I_{t-1}
 leader, challenger, sample-arm procedures

Output: next arm to observe A_t

- 1 $L_t \leftarrow \text{leader}(I_{t-1})$
 - 2 $C_t \leftarrow \text{challenger}(I_{t-1}, L_t)$
 - 3 $A_t \leftarrow \text{sample-arm}(I_{t-1}, L_t, C_t)$
-

Algorithm 18: Empirical-Best leader

Input: history of observations I_{t-1}

Output: leader L_t

- 1 Choose $L_t \in \underset{a \in [K]}{\text{argmax}} \hat{\mu}_a(t-1)$
-

Algorithm 19: challenger procedures for top-two algorithms

Input: history of observations I_{t-1}

leader L_t

Output: challenger C_t

/ Transportation-Cost (TC) challenger */*

- 1 Choose $C_t \in \underset{a \neq L_t}{\text{argmin}} \text{TC}_{a \rightarrow L_t} \left(\hat{\mu}(t-1), \frac{N(t-1)}{t-1} \right)$
// TC_{a→b}(μ, v) is defined in (6.6)

/ Transportation-Cost-Penalized (TCP) challenger */*

- 2 Choose $C_t \in \underset{a \neq L_t}{\text{argmin}} \text{TC}_{a \rightarrow L_t} \left(\hat{\mu}(t-1), \frac{N(t-1)}{t-1} \right) + \log N_a(t-1)$
-

The asymptotic optimality of top-two algorithms has been discussed, but only (until very recently) for non-adaptive strategies, for which the sampling rule chooses the leader with a fixed probability parameter $\beta \in (0, 1)$. As a good leader might satisfy $L_t = a^*$ except for a sub-linear number of time steps, such non-adaptive strategies will pull the best arm a fraction β of the time, which can be far from the optimal frequency $w_{a^*}(\underline{\mu})$. As a consequence, those strategies satisfy the following lower bound (see [Russo, 2016](#)):

$$\forall \underline{\mu} \text{ in } \mathcal{D}_{\text{exp}}, \quad \liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \geq T_{\beta}(\underline{\mu}) \stackrel{\text{def}}{=} \left(\sup_{\substack{v \in \Sigma_K \\ v_{a^*} = \beta}} \inf_{\lambda \in \text{Alt}(\underline{\mu})} \sum_{a \in [K]} v_a d(\mu_a, \lambda_a) \right)^{-1}.$$

It has been proved (see, e.g., [Jourdan et al., 2022](#)) that good choices of leaders and challengers (like those proposed in Algorithms 18 and 19) lead to β -asymptotically optimal strategies, that is, strategies such that, for all generic (hence with distinct means) bandit problems $\underline{\mu}$ in \mathcal{D}_{exp} ,

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \leq T_{\beta}(\underline{\mu}).$$

Adaptive top-two algorithms. Non-adaptive algorithms are hence asymptotically optimal only if $\beta = w_{a^*}(\underline{\mu})$, which is quite frustrating even if it can be proved (see [Russo, 2016](#)) that $T_{\frac{1}{2}}(\underline{\mu}) \leq 2T(\underline{\mu})$, i.e., that the choice of $\beta = \frac{1}{2}$ leads to a loss in guarantee of a multiplicative factor 2 only. To tackle this problem, we might consider adaptive algorithms, for which the sampled arm between the leader and the challenger is not chosen according to an external parameter, but thanks to an adaptive parameter. Based on this idea, [You et al. \(2023\)](#) proved that Top-Two-Thompson-Sampling with an adaptive sampling rule (a top-two strategy with a Bayesian choice of the leader and the challenger, see Section 2.2.7) is asymptotically optimal, but only for a Gaussian model with common variance $\sigma^2 > 0$.

Outline and contributions. The question of extending the asymptotic guarantees of adaptive top-two algorithms to more general models is still open. The objective of this chapter is to present preliminary work in that direction. We focus on exponential models, which form the class of models for which the fixed-confidence best-arm identification problem is best understood (see [Garivier and Kaufmann, 2016](#)).

To begin with, in Section 6.2, we prove that for a general exponential model, the solution $\underline{w}(\underline{\mu})$ of optimization problem (6.2) can be seen as the unique fixed point of some transformation. We use this property to present a new procedure for computing the optimal weight vector $\underline{w}(\underline{\mu})$ that turns out to be empirically efficient. Additionally, an interesting interpretation of the transformation naturally leads to the introduction of a new challenger rule together with an adaptive sampling rule for top-two algorithms. Then, in Section 6.3, we explore how to generalize the Gaussian analysis of [You et al. \(2023\)](#) to a general exponential model \mathcal{D}_{exp} . As a by-product of the analysis of Section 6.3, we give a new proof of the asymptotic optimality of Track-and-Stop in Section 6.4.



Notation. For a given strategy facing a bandit problem $\underline{\mu}$, let $N_a(t)$ and $\hat{\mu}_a(t)$ denote the number of pulls and the empirical mean¹ of arm a at step t :

$$N_a(t) \stackrel{\text{def}}{=} \sum_{s \in [t]} \mathbb{I}\{A_s = a\} \quad \text{and} \quad \hat{\mu}_a(t) \stackrel{\text{def}}{=} \frac{1}{N_a(t)} \sum_{s \in [t]} Y_s \mathbb{I}\{A_s = a\}.$$

¹As strategies initially observe each arm once, $\hat{\mu}_a(t)$ is well-defined for $t \geq K$.

In the rest of this chapter, we fix some exponential model \mathcal{D}_{exp} and rely on the associated notation defined on page 43. In a nutshell, distributions ν_θ of the model are parameterized by elements θ of the natural space parameter Θ , or equivalently by their means $\mu = \mathbb{E}(\nu_\theta)$, which belong to a convex open interval $\mathcal{M} = (\mu_-, \mu_+)$. The log-partition function, denoted by b , is twice differentiable, such that $\mathbb{E}(\nu_\theta) = b'(\theta)$ and b' is invertible, and the mean-parameterized Kullback-Leibler divergence function of the model is defined, for all $\theta, \theta' \in \Theta$, by

$$d(\mathbb{E}(\nu_\theta), \mathbb{E}(\nu_{\theta'})) \stackrel{\text{def}}{=} \text{KL}(\nu_\theta, \nu_{\theta'}) = (\theta - \theta')b'(\theta) - b(\theta) + b(\theta'). \quad (6.4)$$

6.2. A Fixed Point Property

Let, until the end of Section 6.2.2, $\underline{\mu}$ be a fixed bandit instance in \mathcal{D}_{exp} with a unique best arm a^* , and $\underline{w} = \underline{w}(\underline{\mu})$ be its optimal weight vector. We focus on optimization problem (6.2) and recall its formulation in terms of transportation costs (see Section 2.2.4):

$$T(\underline{\mu})^{-1} = \sup_{\underline{v} \in \Sigma_K} \min_{a \neq a^*} \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v}), \quad (6.5)$$

where, for $\underline{v} \in \Sigma_K$ and $a \neq a^*$,

$$\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v}) = v_{a^*} d(\mu^*, \bar{\mu}_{a^*, a, \underline{v}}) + v_a d(\mu_a, \bar{\mu}_{a^*, a, \underline{v}}), \quad \text{where } \bar{\mu}_{a^*, a, \underline{v}} \stackrel{\text{def}}{=} \frac{v_{a^*} \mu^* + v_a \mu_a}{v_{a^*} + v_a}. \quad (6.6)$$

Garivier and Kaufmann (2016) proved that the optimal weight vector \underline{w} solving (6.5) is characterized by the following sufficient and necessary conditions, where $\text{int}(A)$ is the interior of a set A .

Proposition 6.1. *Let $\underline{\mu}$ be a fixed bandit instance in \mathcal{D}_{exp} with a unique best arm a^* . Then for all $\underline{v} \in \text{int}(\Sigma_K)$, the following conditions are equivalent:*

- (i) $\underline{v} = \underline{w}(\underline{\mu})$,
- (ii) for all sub-optimal arms $a \neq a^*$,

$$\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w}) = \left(\sum_{b \neq a^*} \frac{1}{d(\mu_b, \bar{\mu}_{a^*, b, \underline{w}})} \right)^{-1} = T(\underline{\mu})^{-1}, \quad (6.7)$$

- (iii) all transportation costs $(\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w}))_{a \neq a^*}$ are equal and

$$\sum_{b \neq a^*} \frac{d(\mu^*, \bar{\mu}_{a^*, b, \underline{w}})}{d(\mu_b, \bar{\mu}_{a^*, b, \underline{w}})} = 1. \quad (6.8)$$

In Section 2.2.4, we explained the intuition why, for a given proportion $\beta \in (0, 1)$, there exists a unique vector $\underline{w}^\beta = \underline{w}^\beta(\underline{\mu})$ such that $w_{a^*}^\beta = \beta$ and which equalizes the transportation costs $(\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w}^\beta))_{a \neq a^*}$. We also explained that the optimal weight vector \underline{w} belongs to the set of β -optimal weight vectors $\{\underline{w}^\beta : \beta \in (0, 1)\}$, and that the corresponding value of β might be computed by a first-order condition. Conditions (6.7) and (6.8) are equivalent to this condition.

We use Proposition 6.1 to interpret \underline{w} as the unique fixed point of some transformation in Section 6.2.1. We then analyze this transformation to design efficient procedures for the calculation of \underline{w} in Section 6.2.2, and to introduce a new adaptive top-two sampling scheme in Section 6.2.3.

6.2.1. The Transformation

We recall that $\underline{\mu}$ is a fixed and known bandit instance in this section. At the optimum \underline{w} , thank to condition (6.7), all transportation costs are equal and we get

$$T(\underline{\mu}) = \sum_{b \neq a^*} \frac{1}{d(\mu_b, \bar{\mu}_{a^*, b, \underline{w}})}.$$

Defining, for a weight vector $\underline{v} \in \text{int}(\Sigma_K)$,

$$T_{\underline{v}} \stackrel{\text{def}}{=} \sum_{b \neq a^*} \frac{1}{d(\mu_b, \bar{\mu}_{a^*, b, \underline{v}})}, \quad (6.9)$$

the previous equation reads $T_{\underline{w}} = T(\underline{\mu})$. As $\underline{v} \mapsto T_{\underline{v}}$ is continuous, $T_{\underline{v}}$ can be interpreted as an approximation of $T(\underline{\mu})$, at least when \underline{v} is around \underline{w} .

Consider some frequency vector $\underline{v} \in \Sigma_K$. If it is not optimal, one could try to correct the difference between costs $\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w})$ and $T_{\underline{v}}^{-1}$ for all sub-optimal arms a . As the cost $\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w})$ is homogeneous in (v_{a^*}, v_a) , multiplying both the values of v_{a^*} and v_a by $T_{\underline{v}}^{-1} / \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})$ would transform the cost to the exact value $T_{\underline{v}}^{-1}$. Yet, note that if this multiplicative factor is not 1, changing the value of (v_{a^*}, v_a) for a fixed a will result in a vector out of Σ_K . Similarly, it is not possible to change simultaneously the value of the pair for all sub-optimal arm a , as v_{a^*} would be multiplied by different factors. However, we can see what will give the transformation which consists in updating each v_a by the corresponding multiplicative factor and then choosing v_{a^*} so as to get $\underline{v} \in \Sigma_K$. Formally, this defines a transformation $\underline{W} : \text{int}(\Sigma_K) \rightarrow \text{int}(\Sigma_K)$ by

$$\forall \underline{v} \in \text{int}(\Sigma_K), \forall a \in [K], \quad W_a(\underline{v}) \stackrel{\text{def}}{=} \begin{cases} v_a \frac{T_{\underline{v}}^{-1}}{\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})} & \text{if } a \neq a^*, \\ 1 - \sum_{b \neq a^*} v_b \frac{T_{\underline{v}}^{-1}}{\text{TC}_{b \rightarrow a^*}(\underline{\mu}, \underline{v})} & \text{if } a = a^*. \end{cases} \quad (6.10)$$

It is easy to check that this transformation is well-defined. Indeed, for all $\underline{v} \in \text{int}(\Sigma_K)$, we get that $W_{a^*}(\underline{v}) > 0$ by using the definition of the costs (6.6):

$$\sum_{b \neq a^*} \frac{v_b}{\text{TC}_{b \rightarrow a^*}(\underline{\mu}, \underline{v})} < \sum_{b \neq a^*} \frac{v_b}{v_b d(\mu_b, \bar{\mu}_{b, a^*, \underline{v}})} = \sum_{b \neq a^*} \frac{1}{d(\mu_b, \bar{\mu}_{b, a^*, \underline{v}})} = T_{\underline{v}}.$$

By the necessary and sufficient conditions (6.7), the following fixed-point property holds.

Proposition 6.2. *Let $\underline{\mu}$ be a bandit instance in \mathcal{D}_{exp} . Then \underline{w} is the unique fixed point of \underline{W} .*

Proof. Using (6.7), one gets that, for all $a \neq a^*$,

$$W_a(\underline{w}) = w_a \frac{T_{\underline{w}}^{-1}}{\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{w})} = w_a,$$

which ensures that $\underline{W}(\underline{w}) = \underline{w}$.

Reciprocally, if some vector $\underline{v} \in \text{int}(\Sigma_K)$ is such that $\underline{W}(\underline{v}) = \underline{v}$, then

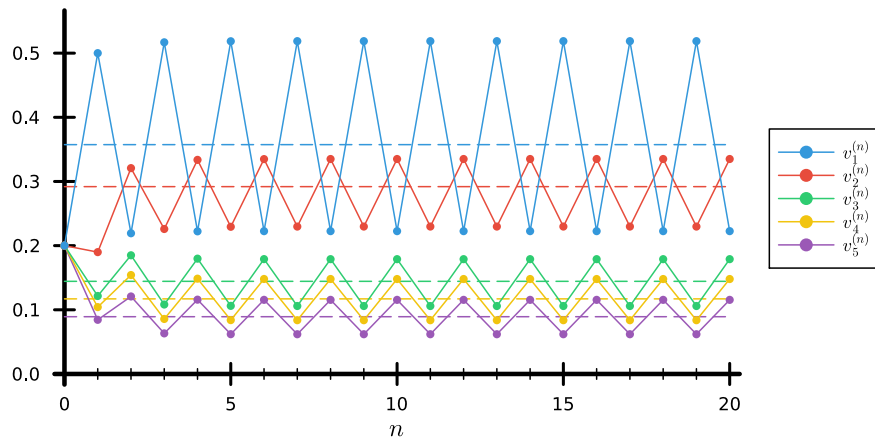
$$\forall a \neq a^*, \quad \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v}) = T_{\underline{v}}^{-1} = \left(\sum_{b \neq a^*} \frac{1}{d(\mu_b, \bar{\mu}_{a^*, b, \underline{v}})} \right)^{-1},$$

hence \underline{v} satisfies condition (6.7), which gives $\underline{v} = \underline{w}$. \square

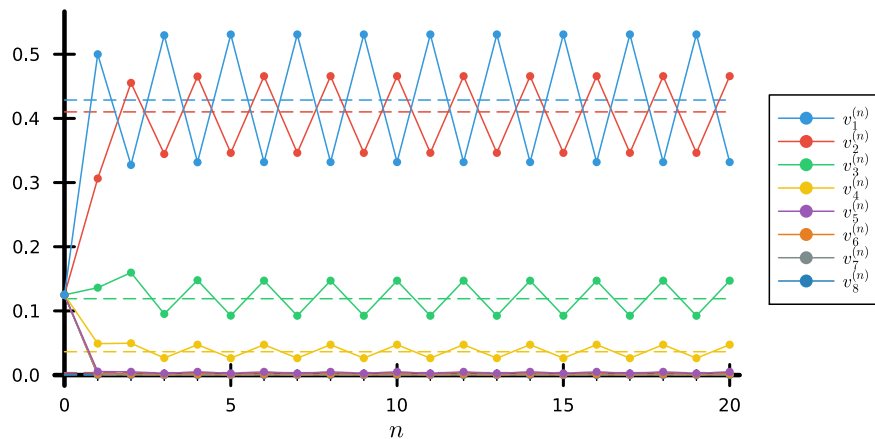
6.2.2. New Empirical Optimal-Weights Procedures

The computation of the optimal weight vector at each time step is the most costly part of strategies like Track-and-Stop. Reducing the complexity of a procedure Optimal-Weights which computes (an approximation of) those weights is then of high interest. We already proposed, independently, such an improvement for Gaussian variables in Chapter 3, and will now see how to use transformation \underline{W} to define a procedure for an exponential model \mathcal{D}_{exp} .

In view of the definition \underline{W} and of its associated fixed point property for \underline{w} , it is natural to wonder whether the iterates $(\underline{v}^{(n)} \stackrel{\text{def}}{=} \underline{W}^n(\underline{v}^{(0)}))_{n \geq 0}$ of some initial vector $\underline{v}^{(0)}$ (e.g., the uniform vector) will converge to \underline{w} , by showing for instance that \underline{W} is a contraction. If so, then we obtain a new Optimal-Weights procedure to approximate \underline{w} , which seems to be computationally efficient.



(a) $\underline{\mu} = (0.9, 0.7, 0.65, 0.63, 0.6)$



(b) $\underline{\mu} = (0.95, 0.93, 0.92, 0.9, 0.8, 0.7, 0.5, 0.4)$

Figure 6.1: For the standard Gaussian model, coordinates of the iterated vectors $(\underline{v}^{(n)})_{n \geq 0}$ with $\underline{v}^{(0)}$ the uniform vector; the values of the corresponding coordinates of \underline{w} are dashed.

\underline{W} is not a contraction. Unfortunately, the convergence does not hold in all generality. Even in the simple model of standard Gaussian variables, the iterates may end up oscillating between two vectors and not converge to \underline{w} , see Figure 6.1. Let us highlights a few remarks concerning the numerical behavior that we observed whatever the Gaussian bandit instance. First, we observe that, when the

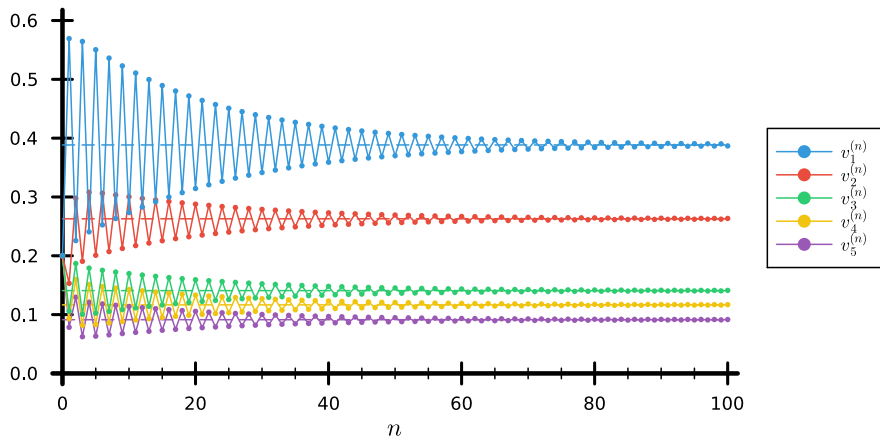


Figure 6.2: For the Bernoulli model and $\underline{\mu} = (0.9, 0.7, 0.65, 0.63, 0.6)$, coordinates of the iterated vectors $(\underline{v}^{(n)})_{n \geq 0}$ with $\underline{v}^{(0)}$ the uniform vector.

weight of the optimal arm is over-estimated, those of sub-optimal arms are under-estimated, and vice versa. Then, the process quickly stabilizes into an oscillation between two vectors, which, strangely, do not satisfy any remarkable property, in particular:

- neither the optimal weight vector is the average of the two oscillations vectors, and, as a consequence, even the Cesàro sums of the iterates would not converge to \underline{w} ,
- nor those two vectors equalize the transportation costs of sub-optimal arms. Nevertheless, they are numerically quite close to vectors satisfying this property, i.e., belonging to the set $\{\underline{w}^\beta(\underline{\mu}) : \beta \in (0, 1)\}$ of β -optimal weight vectors. This explains why the weights of sub-optimal arms are approximately multiplied by the same value at each iteration (see Figure 6.1).

Remark. For Bernoulli instances, we observe in our experiments that the magnitude of the oscillations decreases slowly, and the iterates converge to the optimal weight vector, as can be seen in Figure 6.2. However, the convergence is not very fast and smooth.

Soft update rule. In order to obtain convergence, one can consider less aggressive update rules. For instance, let us focus on the sequence defined by

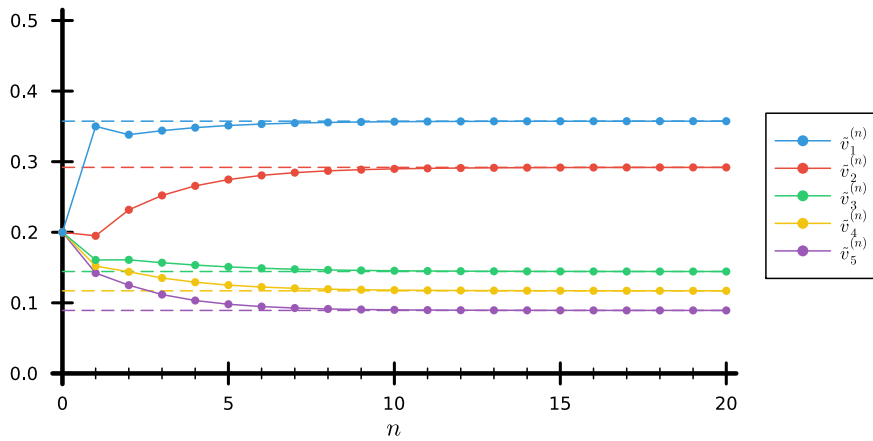
$$\forall n \geq 0, \quad \tilde{\underline{v}}^{(n+1)} \stackrel{\text{def}}{=} \frac{\tilde{\underline{v}}^{(n)} + W(\tilde{\underline{v}}^{(n)})}{2}. \quad (6.11)$$

For this sequence, we observed on all our simulations that the convergence to \underline{w} is quick and unambiguous whatever the model (see Figure 6.3 for a Gaussian model, the figures are similar for, e.g., Bernoulli or Poisson variables). Yet, how to guarantee this convergence theoretically is still to be understood. It might be possible to show that the associated transformation is a contraction, as it numerically seems to be the case for the ℓ^2 -norm.

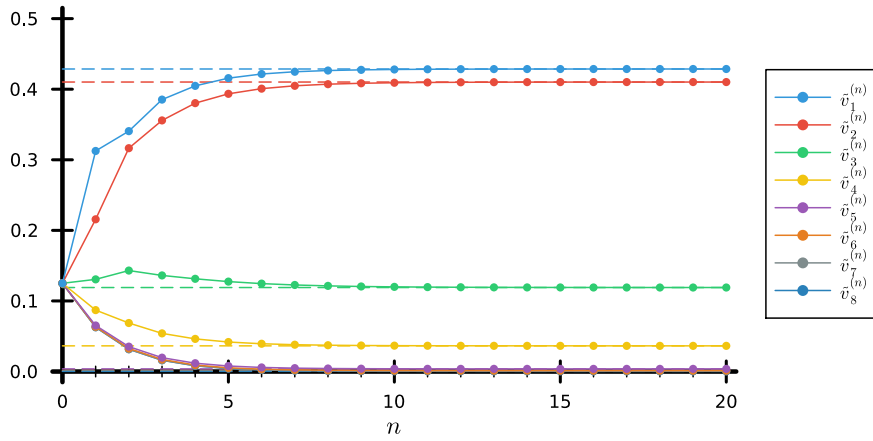
Conjecture 6.3. For all exponential models \mathcal{D}_{exp} , and all bandit problems $\underline{\mu}$ in \mathcal{D}_{exp} with a unique optimal arm, the sequence of iterates $(\tilde{\underline{v}}^{(n)})_{n \geq 0}$ defined by (6.11) converges to \underline{w} .

Remark. One can state other update rules, like modifying only the weights of one sub-optimal arm together with the weight of the optimal arm, choosing this sub-optimal arm circularly or randomly. The empirical performance of those procedures is also pretty good, as far as we observed in some numerical experiments.

6.2. A FIXED POINT PROPERTY



(a) $\underline{\mu} = (0.9, 0.7, 0.65, 0.63, 0.6)$



(b) $\underline{\mu} = (0.95, 0.93, 0.92, 0.9, 0.8, 0.7, 0.5, 0.4)$

Figure 6.3: For the standard Gaussian model, coordinates of the vectors $(\tilde{v}^{(n)})_{n \geq 0}$ with $\tilde{v}^{(0)}$ the uniform vector; the values of the corresponding coordinates of \underline{w} are dashed.

6.2.3. A New Sampling Rule for top-two Algorithms

We now give another point of view on the transformation \underline{W} , that will allow us to define a top-two sampling rule. We write, for all sub-optimal arms a :

$$W_a(\underline{v}) = v_a \frac{T_{\underline{v}}^{-1}}{\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})} \cdot \frac{d(\mu_a, \bar{\mu}_{a^*, a, \underline{v}})}{d(\mu_a, \bar{\mu}_{a^*, a, \underline{v}})} = \frac{T_{\underline{v}}^{-1}}{d(\mu_a, \bar{\mu}_{a^*, a, \underline{v}})} \cdot \frac{v_a d(\mu_a, \bar{\mu}_{a^*, a, \underline{v}})}{\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})} = \chi_a(\underline{v}) \cdot (1 - \beta_a(\underline{v})),$$

where we define

$$\chi_a(\underline{v}) \stackrel{\text{def}}{=} \frac{T_{\underline{v}}^{-1}}{d(\mu_a, \bar{\mu}_{a^*, a, \underline{v}})} \quad \text{and} \quad \beta_a(\underline{v}) \stackrel{\text{def}}{=} \frac{v_{a^*} d(\mu_{a^*}, \bar{\mu}_{a^*, a, \underline{v}})}{\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})}. \quad (6.12)$$

$\chi_a(\underline{v})$ is the relative contribution of arm a in $T_{\underline{v}}^{-1}$, while $\beta_a(\underline{v})$ is the proportion of transportation cost $\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v})$ due to the transport of arm a^* . We note, in passing, that vector $\underline{\chi}$ is a probability vector, which will be useful in Section 6.2.3.

Example. For the model $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ of Gaussian variables with common variance $\sigma^2 > 0$, the closed-form expression (2.21) of the Kullback-Leibler divergence gives

$$\begin{aligned}\beta_a(\underline{v}) &= \frac{v_{a^*} d(\mu^*, \bar{\mu}_{a^*, a, \underline{v}})}{v_{a^*} d(\mu^*, \bar{\mu}_{a^*, a, \underline{v}}) + v_a d(\mu_a, \bar{\mu}_{a^*, a, \underline{v}})} \\ &= \frac{v_{a^*} \left(\frac{v_a}{v_a + v_{a^*}} (\mu^* - \mu_a) \right)^2}{v_{a^*} \left(\frac{v_a}{v_a + v_{a^*}} (\mu^* - \mu_a) \right)^2 + v_a \left(\frac{v_{a^*}}{v_a + v_{a^*}} (\mu_a - \mu^*) \right)^2} \\ &= \frac{v_a}{v_a + v_{a^*}}.\end{aligned}$$

We note that this quantity does not depend on $\underline{\mu}$, which is a specificity of the model.

The definition of transformation \underline{W} now reads

$$\forall \underline{v} \in \text{int}(\Sigma_K), \forall a \in [K], \quad W_a(\underline{v}) = \begin{cases} \chi_a(\underline{v})(1 - \beta_a(\underline{v})) & \text{if } a \neq a^*, \\ \sum_{b \neq a^*} \chi_b(\underline{v})\beta_b(\underline{v}) & \text{if } a = a^*. \end{cases}$$

We note that the knowledge of $\underline{\chi}(\underline{w})$ and $\beta_a(\underline{w})$ is sufficient to generate the optimal weights $\underline{w}(\underline{\mu})$. , we select a couple of best-arm and challenger (a^*, a) where the challenger a is sampled according to the probability $\underline{\chi}(\underline{w})$. Then we choose between the optimal arm and its challenger by keeping the leader with probability $\beta_a(\underline{w})$. This results in a sample from the optimal weight vector \underline{w} thanks to the fixed point property 6.2. Interestingly, we recognize here a top-two procedure that first samples a couple of leader and challenger and then selects an arm from this couple.

Sampling rule. To give a precise description of the top-two sampling rule, we need to write explicitly the dependency on $\underline{\mu}$ of the previously defined quantities $T_{\underline{v}}$, $\chi_a(\underline{v})$, $\beta_a(\underline{v})$ and $W_a(\underline{v})$, as the bandit instance is unknown to the strategy. To that end, we define $T_{\underline{v}}(\underline{\mu})$, $\chi_a(\underline{\mu}, \underline{v})$ and $W_a(\underline{\mu}, \underline{v})$ respectively by Equations (6.9), (6.12) and (6.10), and we set:

$$\forall a \neq b \in [K], \quad \beta_{b,a}(\underline{\mu}, \underline{v}) \stackrel{\text{def}}{=} \begin{cases} \frac{v_b d(\mu_b, \bar{\mu}_{b, a, \underline{v}})}{v_b d(\mu_b, \bar{\mu}_{b, a, \underline{v}}) + v_a d(\mu_a, \bar{\mu}_{b, a, \underline{v}})} & \text{if } \mu_b \neq \mu_a, \\ \frac{1}{2} & \text{if } \mu_b = \mu_a, \end{cases}$$

which corresponds to the relative contribution of arm b to the transportation cost between arms a and b . Note that the dependency of $\beta_{b,a}(\underline{\mu}, \underline{v})$ on vector \underline{v} is only through the ratio $x_{b,a} = \frac{v_b}{v_b + v_a}$, as

$$\beta_{b,a}(\underline{\mu}, \underline{v}) = \frac{1}{1 + \frac{v_a d(\mu_a, \bar{\mu}_{b, a, \underline{v}})}{v_b d(\mu_b, \bar{\mu}_{b, a, \underline{v}})}} = \frac{1}{1 + \frac{1 - x_{b,a} d(\mu_a, (1 - x_{b,a})\mu_a + x_{b,a}\mu_b)}{x_{b,a} d(\mu_b, (1 - x_{b,a})\mu_a + x_{b,a}\mu_b)}}.$$

Hence we will sometimes see $\beta_{b,a}(\underline{\mu}, \cdot)$ as a one-variable function:

$$\forall x \in (0, 1), \quad \beta_{b,a}(\underline{\mu}, x) \stackrel{\text{def}}{=} \frac{1}{1 + \frac{1 - x d(\mu_a, (1 - x)\mu_a + x\mu_b)}{x d(\mu_b, (1 - x)\mu_a + x\mu_b)}}. \quad (6.13)$$

Once again (see, e.g., the Track-and-Stop strategy), at each time step t , we consider the plug-in estimates of those quantities when replacing $\underline{\mu}$ by its empirical estimate $\hat{\underline{\mu}}(t)$. We choose

Algorithm 20: $\underline{\chi}$ challenger

Input: history of observations I_{t-1}

 leader L_t
Output: challenger C_t

- 1 $\underline{p} \leftarrow \underline{p}(t-1)$
 - 2 $\hat{\underline{\mu}} \leftarrow \hat{\underline{\mu}}(t-1)$
 - 3 **for** $a \neq L_t$ **do**
 - 4 $\chi_a \leftarrow \frac{T_{\underline{p}}^{-1}(\hat{\underline{\mu}})}{d(\hat{\underline{\mu}}_a, \hat{\underline{\mu}}_{L_t, a, \underline{p}})}$
 - 5 Choose C_t according to distribution $\underline{\chi}$
-

Algorithm 21: Adaptive sample-arm (sampling version)

Input: history of observations I_{t-1}

 leader L_t

 challenger C_t
Output: sampled arm A_t

- 1 $A_t \leftarrow \begin{cases} L_t & \text{with probability } \beta_{L_t, C_t}(\hat{\underline{\mu}}(t-1), \underline{p}(t-1)), \\ C_t & \text{otherwise.} \end{cases}$
-

L_t as the Empirical-Best leader (see Algorithm 18), and as previously explained, the challenger procedure consists in² choosing C_t according to vector

$$\hat{\underline{\chi}}(t-1) \stackrel{\text{def}}{=} \underline{\chi}(\hat{\underline{\mu}}(t-1), \underline{p}(t-1)), \quad \text{where } \underline{p}(t-1) \stackrel{\text{def}}{=} \frac{N(t-1)}{t-1}$$

denotes the empirical proportions of draws after time step $t-1$. The leader L_t , respectively the challenger C_t , is then pulled with probability

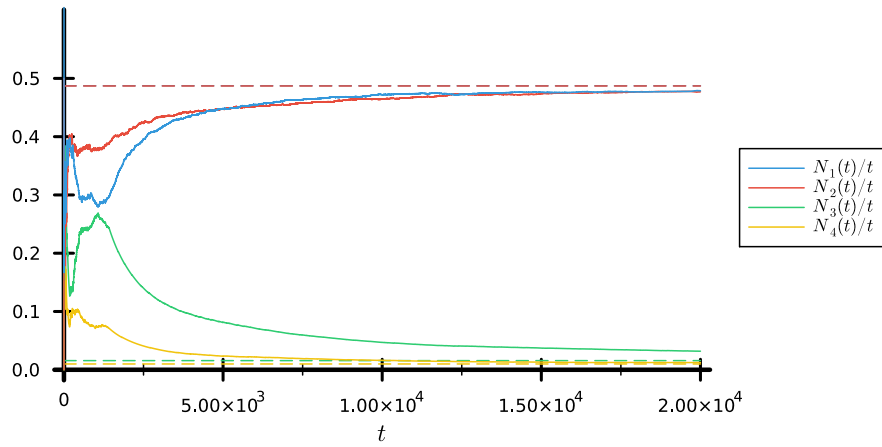
$$\hat{\beta}_{L_t, C_t}(t) \stackrel{\text{def}}{=} \beta_{L_t, C_t}(\hat{\underline{\mu}}(t-1), \underline{p}(t-1)), \quad \text{respectively } \hat{\beta}_{C_t, L_t}(t) \stackrel{\text{def}}{=} \beta_{C_t, L_t}(\hat{\underline{\mu}}(t-1), \underline{p}(t-1)).$$

See Algorithms 20 and 21 for details. This adaptive strategy will be denoted by TT-EB- $\underline{\chi}$.

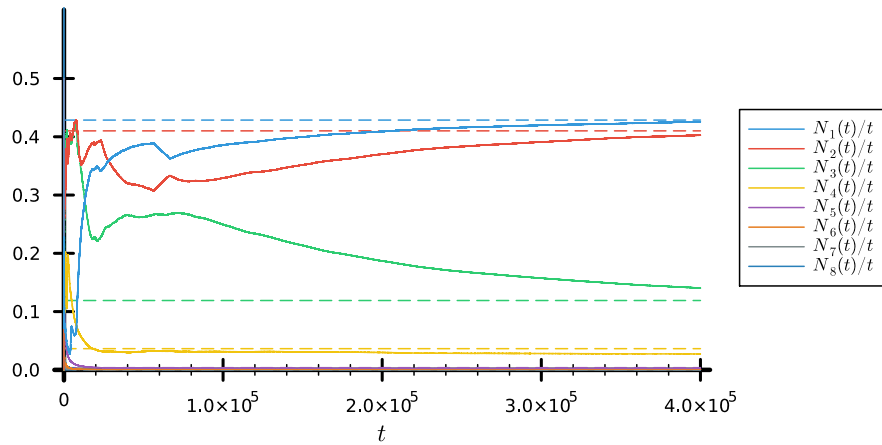
Remark. The adaptive proportions $\beta_{b,a}(\underline{\mu}, \underline{v})$ have also been obtained by [You et al. \(2023\)](#) for a general exponential model.

Experiments. The convergence of the sampling frequencies of the strategy TT-EB- $\underline{\chi}$ to the optimal weight vector $\underline{w}(\underline{\mu})$ is illustrated on Figures 6.4 and 6.5. Whatever the model and parameter, numerical experiments show convergence to the optimal frequencies. We observe that the evolution of the sampling frequencies is quite smooth compared to, e.g., Track-and-Stop or the top-two algorithm TT-EB-TC with adaptive proportions, for which the challenger is chosen as the arm minimizing the transportation cost with the leader.

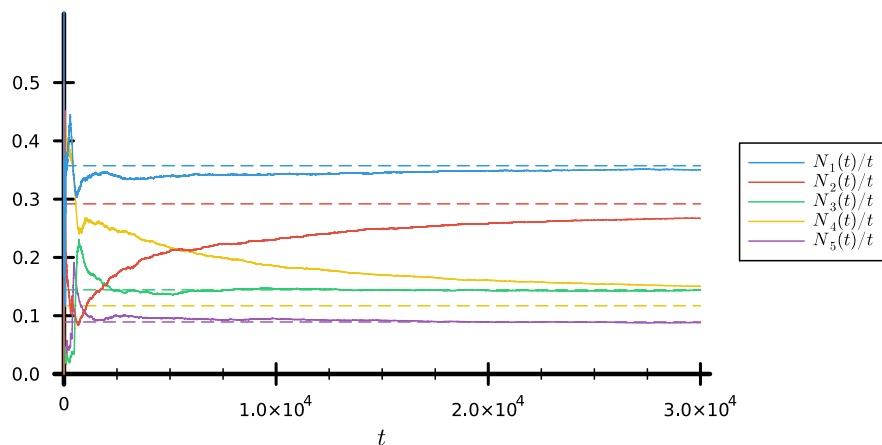
²To be well-defined, $\underline{\chi}$ requires $\hat{\underline{\mu}}(t-1)$ to have a unique optimal arm. If not, we can simply choose A_t among all empirically optimal arms (but note that this event occurs with null probability for a continuous model).



(a) $\underline{\mu} = (0.9, 0.7, 0.4, 0.3)$



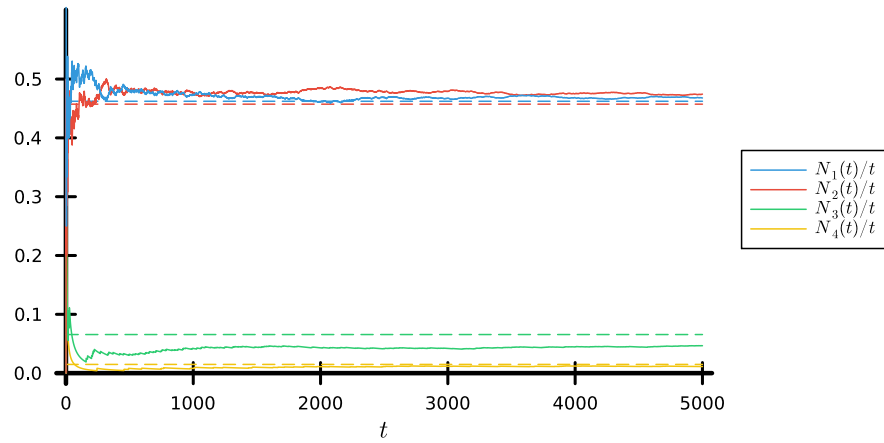
(b) $\underline{\mu} = (0.95, 0.93, 0.92, 0.9, 0.8, 0.7, 0.5, 0.4)$



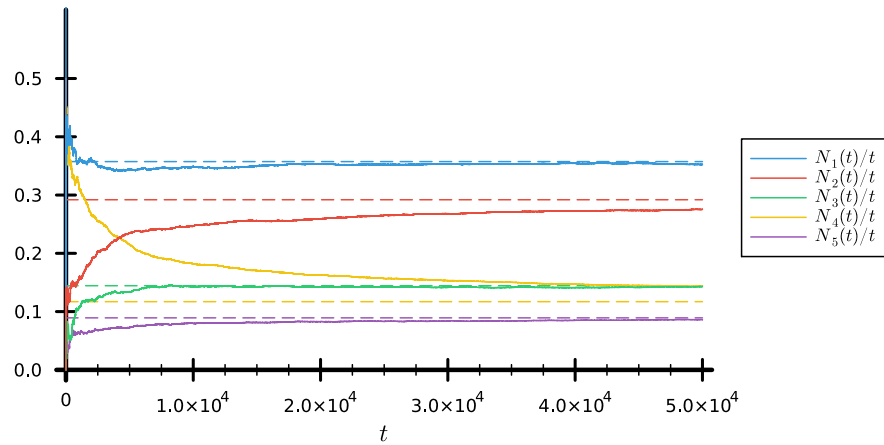
(c) $\underline{\mu} = (0.9, 0.7, 0.65, 0.63, 0.6)$

Figure 6.4: Evolution of the sampling frequencies on a simulation of TT-EB- $\underline{\chi}$ with standard Gaussian variables. The values of the corresponding coordinates of $\underline{w}(\underline{\mu})$ are dashed.

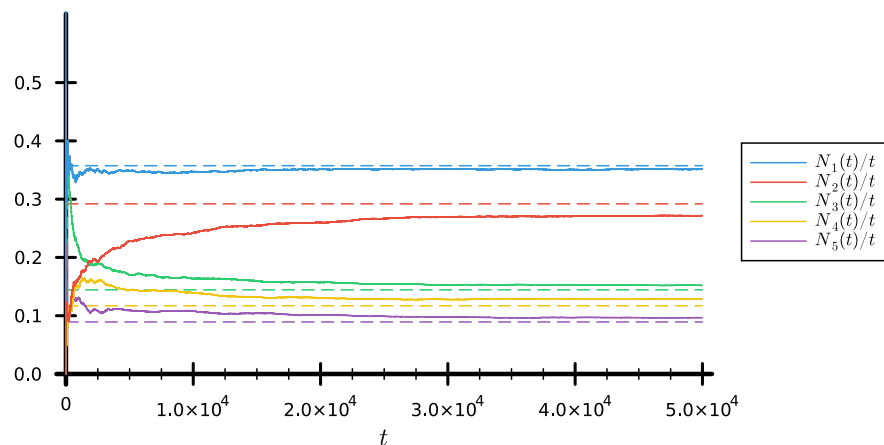
6.2. A FIXED POINT PROPERTY



(a) Bernoulli model, $\underline{\mu} = (0.9, 0.7, 0.5, 0.1)$



(b) Bernoulli model, $\underline{\mu} = (0.9, 0.7, 0.65, 0.63, 0.6)$



(c) Poisson model, $\underline{\mu} = (0.9, 0.7, 0.65, 0.63, 0.6)$

Figure 6.5: Evolution of the sampling frequencies on a simulation of TT-EB- $\underline{\chi}$ with the Bernoulli and Poisson models. The values of the corresponding coordinates of $w(\underline{\mu})$ are dashed.

6.3. Asymptotically Optimal Adaptive Algorithms

The objective of this section is to discuss some adaptive top-two algorithms that could be asymptotically optimal for a general exponential model \mathcal{D}_{exp} . The result was recently obtained for a Gaussian model by You et al. (2023) for the Top-Two-Thompson-Sampling algorithm, but their analysis also holds for, e.g., the pair of Empirical-Best leader and Transportation-Cost or Transportation-Cost-Penalized challenger.

Considered leaders and challengers. In order to get the simplest possible analysis, we focus in this section on the Empirical-Best leader given in Algorithm 18 and the Transportation-Cost and Transportation-Cost-Penalized challengers given in Algorithm 19. The adaptive versions of those top-two algorithms are denoted by TT-EB-TC and TT-EB-TCP (similarly to the acronyms of Jourdan et al., 2022, see also Section 2.2.7). The conjecture is given below.

Conjecture 6.4. *Let \mathcal{D}_{exp} be an exponential model. The TT-EB-TC and TT-EB-TCP algorithms with the adaptive sampling rule of Algorithm 21 satisfy, for all bandit problems $\underline{\mu}$ in \mathcal{D}_{exp} with distinct means:*

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \leq T(\underline{\mu}).$$

Remark. When studying top-two algorithms, it is a common hypothesis to work with bandits instance for which all means differ (see, e.g., Jourdan et al., 2022). We will refer to such instances as *generic* bandits. For other instances, there might be a lack of exploration, especially when the first estimates of the best arm are poor.

In the rest of this section, we provide some first arguments to prove Conjecture 6.4. We largely follow the proof structure of You et al. (2023, Appendix A) for Gaussian variables (except the additional step 2), and try to generalize their statements. As we will see, this might imply slight modifications of the considered strategies TT-EB-TC and TT-EB-TCP, which do not deeply affect the definition of the sampling rules.

Overview of the argumentation. In the sequel, we fix some generic bandit problem $\underline{\mu}$ in \mathcal{D}_{exp} . Let a^* and \underline{w} respectively denote its unique optimal arm and its optimal weight vector.

We split the analysis into several steps:

1. In the first step, we introduce notation and recall that asymptotic optimality can be proved by showing the convergence of the ratios $p_a(t)/p_{a^*}(t)$ to w_a/w_{a^*} for all sub-optimal arms a ,
2. Then, in step 2, we propose to slightly modify the sampling rule in order to facilitate the control of the empirical numbers of pulls.
3. To ensure the convergence of empirical estimates, we have to prove that the strategy sufficiently explores all arms. While we still do not know how to prove the sufficient exploration for TT-EB-TC and TT-EB-TCP with a general exponential model, we give a possible path of reasoning.
4. In order to prove the convergence of the ratios $p_a(t)/p_{a^*}(t)$, You et al. (2023) showed another useful convergence in the Gaussian case. We explain in step 4 that the generalization of this convergence seems to be the difficult part of our analysis.
5. Finally, we quickly explain in step 5 how one could prove the convergence of the ratios $p_a(t)/p_{a^*}(t)$ if the requirements of steps 3 (sufficient exploration) and 4 (supporting convergence) were satisfied.

6.3.1. Step 1: A Sufficient Condition

Let us introduce a notion of *quick convergence* for random variables.

Definition 6.5. [quick convergence]

Let \mathbb{L}^1 denotes the set of real valued random variables X such that $\mathbb{E}[|X|] < +\infty$.

We say that a sequence of real random variables $(X(t))_{t \in \mathbb{N}}$ *quickly converges*:

- to $x \in \mathbb{R}$, a fact which we denote $X(t) \xrightarrow[t \rightarrow +\infty]{q.} x$, if

$$\forall \varepsilon > 0, \exists T \in \mathbb{L}^1, \quad \forall t \geq T, \quad |X(t) - x| \leq \varepsilon.$$

- to $+\infty$, a fact which we denote $X(t) \xrightarrow[t \rightarrow +\infty]{q.} +\infty$, if

$$\forall c > 0, \exists T \in \mathbb{L}^1, \quad \forall t \geq T, \quad X(t) \geq c.$$

- to $-\infty$, a fact which we denote $X(t) \xrightarrow[t \rightarrow +\infty]{q.} -\infty$, if $-X(t) \xrightarrow[t \rightarrow +\infty]{q.} +\infty$.

We say that a sequence of random vectors $(\underline{X}(t))_{t \in \mathbb{N}}$ in \mathbb{R}^K quickly converges to $\underline{x} \in \mathbb{R}^K$, a fact which we denote $\underline{X}(t) \xrightarrow[t \rightarrow +\infty]{q.} \underline{x}$, if $X_a(t) \xrightarrow[t \rightarrow +\infty]{q.} x_a$ for all $a \in [K]$.

Proving the asymptotic optimality of a strategy using the Global-Likelihood-Ratio stopping rule can be reduced to showing that the empirical frequencies of pulls quickly converge to \underline{w} , as the following result states.

Proposition 6.6. *Let \mathcal{D}_{exp} be an exponential model. Consider a strategy for which the stopping rule is the Global-Likelihood-Ratio stopping rule (Algorithm 5) with threshold (6.3) for some $\alpha > 1$. If $\underline{\mu}$ is a bandit problem in $\mathcal{D}_{\mathcal{N}_{\sigma^2}}$ such that the sampling rule (applied infinitely without stopping rule) satisfies*

$$\underline{p}(t) \xrightarrow[t \rightarrow +\infty]{q.} \underline{w},$$

then

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \leq T(\underline{\mu}).$$

This result was originally stated for a Gaussian model by [Qin et al. \(2017\)](#), but the generalization to exponential models is straightforward³.

In fact, the quick convergence of $\underline{p}(t)$ to \underline{w} is equivalent to the quick convergence of the ratios $(p_a(t)/p_{a^*}(t))_{a \neq a^*}$:

$$\underline{p}(t) \xrightarrow[t \rightarrow +\infty]{q.} \underline{w} \quad \iff \quad \forall a \neq a^*, \quad \frac{p_a(t)}{p_{a^*}(t)} \xrightarrow[t \rightarrow +\infty]{q.} \frac{w_a}{w_{a^*}}.$$

This equivalence is easy to prove, using that $\underline{p}(t)$ and \underline{w} belong to Σ_K for the indirect sense. Hence

$$\forall a \neq a^*, \quad \frac{p_a(t)}{p_{a^*}(t)} \xrightarrow[t \rightarrow +\infty]{q.} \frac{w_a}{w_{a^*}} \quad \implies \quad \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log 1/\delta} \leq T(\underline{\mu}).$$

Objective. To prove the conjecture, we now focus on demonstrating that

$$\forall a \neq a^*, \quad \frac{p_a(t)}{p_{a^*}(t)} \xrightarrow[t \rightarrow +\infty]{q.} \frac{w_a}{w_{a^*}}. \quad (6.14)$$

³And even to more general models, up to technicalities, see [Jourdan et al. \(2022, Theorem 2 in Appendix C.5\)](#).

Remark. You et al. (2023) actually worked with a stronger condition than quick convergence: instead of considering random variables in \mathbb{L}^1 , they require those variables to be in $\cap_{p \geq 1} \mathbb{L}^p$, following tools developed by Qin and Russo (2022). Yet, quick convergence is sufficient for what we want to prove.

6.3.2. Step 2: Tracking versus Sampling

Due to the randomization of the strategy (and especially the random choice between the leader and the challenger), the vector $\underline{p}(t)$ might largely deviate from its expected value. When working with (sub-)Gaussian variables, this deviation might be controlled so that it does not really complicate the arguments (see You et al., 2023). With a general exponential model, we might ease the analysis by replacing the sampling with a tracking: no concentration result will be required. The technique was considered by Jourdan and Degenne (2023) to obtain non-asymptotic guarantees for a Top-Two algorithm.

Tracking. When choosing between the leader and the challenger, one can replace the sampling choice with a C-tracking of the targeted proportions. We now explain the modification in detail.

We will need to define tracking procedures associated with fixed pairs of leader and challenger. To that end, for a pair of arms $a \neq b$ and a time step t , we consider the event

$$\mathcal{C}^{a,b}(t) \stackrel{\text{def}}{=} \{L_t = a \text{ and } C_t = b\},$$

which indicates a time step at which the leader and challenger are respectively arms a and b . We count the cumulative number of pulls of arms a and b under those events and denote them by:

$$N_a^{a,b}(t) \stackrel{\text{def}}{=} \sum_{s \in [t]} \mathbb{I}\{A_s = a\} \mathbb{I}\{\mathcal{C}^{a,b}(s)\} \quad \text{and} \quad N_b^{a,b}(t) \stackrel{\text{def}}{=} \sum_{s \in [t]} \mathbb{I}\{A_s = b\} \mathbb{I}\{\mathcal{C}^{a,b}(s)\},$$

and we also define the average target proportions of arms a and b under those events:

$$\begin{aligned} \bar{\beta}_a^{a,b}(t) &\stackrel{\text{def}}{=} \frac{1}{N_a^{a,b}(t)} \sum_{s \in [t]} \hat{\beta}_{a,b}(s) \mathbb{I}\{\mathcal{C}^{a,b}(s)\}, \\ \text{and} \quad \bar{\beta}_b^{a,b}(t) &\stackrel{\text{def}}{=} \frac{1}{N_a^{a,b}(t)} \sum_{s \in [t]} \hat{\beta}_{b,a}(s) \mathbb{I}\{\mathcal{C}^{a,b}(s)\} = 1 - \bar{\beta}_a^{a,b}(t), \\ \text{where} \quad N^{a,b}(t) &\stackrel{\text{def}}{=} N_a^{a,b}(t) + N_b^{a,b}(t) = \sum_{s \in [t]} \mathbb{I}\{\mathcal{C}^{a,b}(s)\}. \end{aligned}$$

We modify the choice of arm A_t between the leader L_t and the challenger C_t so as to track the associated average target proportions, see Algorithm 22.

Algorithm 22: Adaptive sample-arm (tracking version)

Input: history of observations I_{t-1}

leader L_t

challenger C_t

Output: sampled arm A_t

$$1 \quad A_t \leftarrow \begin{cases} L_t & \text{if } \frac{N_{L_t}^{L_t, C_t}(t-1)}{N_{L_t}^{L_t, C_t}(t-1) + N_{C_t}^{C_t, L_t}(t-1)} < \bar{\beta}_{L_t}^{L_t, C_t}(t), \\ C_t & \text{otherwise.} \end{cases}$$

With that modification, we get precise control of the quantities $N_a^{a,b}(t)$ and $N_b^{a,b}(t)$:

Lemma 6.7. For all pairs (a, b) of arms with $a \neq b$, and $t \geq 0$, we get

$$\left| N_a^{a,b}(t) - N^{a,b}(t) \bar{\beta}_a^{a,b}(t) \right| = \left| N_b^{a,b}(t) - N^{a,b}(t) \bar{\beta}_b^{a,b}(t) \right| \leq 1.$$

Before proving this result, we explain why it might be useful. Given a good choice of leader, we hope that $L_t = a^*$ except for a sub-linear number of time steps. If this property holds, then

$$\forall a \neq a^*, \quad p_a(t) \simeq \frac{N_a^{a^*,a}(t)}{t} \simeq \frac{N^{a^*,a}(t)}{t} \bar{\beta}_a^{a^*,a}(t),$$

and proving the condition (6.14) (or maybe $\underline{p}(t) \xrightarrow[t \rightarrow +\infty]{q.} \underline{w}$) might be done by studying the asymptotic behaviour of $N^{a^*,a}(t)/t$ and $\bar{\beta}_a^{a^*,a}(t)$.

Proof. The equality comes from the relationship:

$$N_b^{a,b}(t) - N^{a,b}(t) \bar{\beta}_{b,a}(t) = N^{a,b}(t) - N_a^{a,b}(t) - N^{a,b}(t)(1 - \bar{\beta}_{a,b}(t)) = N^{a,b}(t) \bar{\beta}_{a,b}(t) - N_a^{a,b}(t).$$

To obtain the inequality, we proceed by induction on t . Let us assume that the result holds at time step $t - 1$. If at time step t , we have $L_t \neq a$ or $C_t \neq b$, then all considered quantities are not updated and the result holds by induction. Now, if $L_t = a$ and $C_t = b$, let us assume that arm $A_t = a$ (one can adapt the argument when $A_t = b$) and prove that

$$-1 \leq N_a^{a,b}(t) - N^{a,b}(t) \bar{\beta}_a^{a,b}(t) \leq 1. \quad (6.15)$$

By definition of A_t (see Algorithm 22), we know, as $A_t = a$, that

$$N_a^{a,b}(t-1) \leq N^{a,b}(t-1) \bar{\beta}_a^{a,b}(t),$$

hence, using that $N_a^{a,b}(t) = N_a^{a,b}(t-1) + 1$ and $N^{a,b}(t) = N^{a,b}(t-1) + 1$:

$$N_a^{a,b}(t) - N^{a,b}(t) \bar{\beta}_a^{a,b}(t) = N_a^{a,b}(t-1) + 1 - N^{a,b}(t-1) \bar{\beta}_a^{a,b}(t) - \bar{\beta}_a^{a,b}(t) \leq 1,$$

which gives the second inequality of (6.15). For the first inequality, remark that

$$N^{a,b}(t) \bar{\beta}_a^{a,b}(t) = N^{a,b}(t-1) \bar{\beta}_a^{a,b}(t-1) + \beta_{a,b}(t),$$

which, together with the induction property, leads to

$$N_a^{a,b}(t) - N^{a,b}(t) \bar{\beta}_a^{a,b}(t) = \underbrace{N_a^{a,b}(t-1) + 1 - N^{a,b}(t-1) \bar{\beta}_a^{a,b}(t-1)}_{\geq 0 \text{ by induction}} - \beta_{a,b}(t) \geq -1.$$

This concludes the proof. □

6.3.3. Step 3: Sufficient Exploration

A typical methodology in fixed-confidence best-arm identification (see, e.g., the use of forced exploration for Track-and-Stop in page 49) is to ensure that empirical quantities converge to their expected values by requiring that the number of pulls of arms diverges almost surely, with a minimal sub-linear rate that we take (arbitrarily) of order \sqrt{t} .

Definition 6.8. [sufficient exploration]

We say that a strategy provides sufficient exploration with rate \sqrt{t} if

$$\exists C > 0, \exists T \in \mathbb{L}^1, \quad \forall t \geq T, \quad \min_{a \in [K]} N_a(t) \geq C\sqrt{t}. \quad (6.16)$$

Jourdan et al. (2022) proved that both TT-EB-TC- β and TT-EB-TCP- β , which are non-adaptive versions of TT-EB-TC and TT-EB-TCP, provide sufficient exploration for an exponential model. Their result does not generalize to adaptive versions for which the probability of choosing the leader instead of the challenger might be arbitrarily close to 0 or 1. You et al. (2023) explained how to override this limitation in the Gaussian case. In the last paragraph of this section, we will explain why we do think that the result still holds for an exponential model, without providing explicit proof.

Forced exploration. Meanwhile, we can use a tool that ensures that sufficient exploration is provided: we might simply force exploration (like, e.g., the Track-and-Stop strategy) and hence assume that this step is satisfied for free. We recall the process quickly (see also page 49). We define, at a given time step t , the set $U(t-1)$ of arms that are under-sampled with respect to the sub-linear rate \sqrt{t} :

$$U(t-1) \stackrel{\text{def}}{=} \left\{ a \in [K] : N_a(t-1) < \sqrt{\frac{t}{K}} + 1 \right\}.$$

If this set is not empty, we overrule the choice of A_t given by the top-two sampling rule and replace it by picking the least pulled arm. This implies that the following lemma holds.

Lemma 6.9. *The TT-EB-TC and TT-EB-TCP algorithms, with the forced exploration process defined above, provide sufficient exploration with rate \sqrt{t} :*

$$\forall t \geq 0, \quad \min_{a \in [K]} N_a(t) \geq \sqrt{\frac{t}{K}}.$$

Remark.

- With forced exploration, we do not have to restrict our attention to a generic bandit problem $\underline{\mu}$. The analysis stands for all bandit problems with a unique optimal arm.
- With forced exploration, the penalization of the Transportation-Cost-Penalized challenger is useless. Yet, as we aim to prove that forced exploration is not necessary at least for TT-EB-TCP, we still consider this strategy in this section.

Convergence of the empirical means. The sufficient exploration property implies an important property which is the quick convergence of the empirical means to $\underline{\mu}$.

Proposition 6.10. *Consider a strategy for which sufficient exploration holds with rate \sqrt{t} . Then*

$$\underline{\hat{\mu}}(t) \xrightarrow[t \rightarrow +\infty]{q.} \underline{\mu}.$$

Remark. This result can be used, together with Proposition 6.6 to prove the asymptotic optimality of the Track-and-Stop and Exploration-Biased-Sampling algorithms, as we detail in Section 4.

The proof is an adaptation of Garivier and Kaufmann (2016, Lemma 19).

Proof. Let $C > 0$ and $T \in \mathbb{L}^1$ satisfying (6.16). Fix $\varepsilon > 0$ and define, for $t \geq 1$:

$$\mathcal{E}_t = \bigcap_{s \geq t} \bigcap_{a \in [K]} \left\{ |\hat{\mu}_a(s) - \mu_a| \leq \varepsilon \right\}.$$

We want to prove that

$$\mathbb{E}_{\underline{\mu}}[T_\varepsilon] < +\infty, \quad \text{where } T_\varepsilon \stackrel{\text{def}}{=} \inf \left\{ t \geq 1 : \mathcal{E}_t \right\}. \quad (6.17)$$

We get

$$\begin{aligned} \mathbb{E}_{\underline{\mu}}[T_\varepsilon] &= \sum_{t \geq 0} \mathbb{P}_{\underline{\mu}}(T_\varepsilon > t) \\ &= \sum_{t \geq 0} \mathbb{P}_{\underline{\mu}}(T_\varepsilon > t, T > t) + \mathbb{P}_{\underline{\mu}}(T_\varepsilon > t, t \geq T) \\ &= \sum_{t \geq 0} \mathbb{P}_{\underline{\mu}}(T > t) + \mathbb{P}_{\underline{\mu}}(T_\varepsilon > t, t \geq T) \\ &= \mathbb{E}_{\underline{\mu}}[T] + \sum_{t \geq 0} \mathbb{P}_{\underline{\mu}}(\mathcal{E}_t^c, t \geq T) \end{aligned} \quad (6.18)$$

Fix $t \geq 0$. We have, using optional skipping and the definition of t

$$\begin{aligned} \mathbb{P}_{\underline{\mu}}(\mathcal{E}_t^c, t \geq T) &= \mathbb{P}_{\underline{\mu}}(\exists s \geq t, \exists a \in [K], |\hat{\mu}_a(s) - \mu_a| > \varepsilon, t \geq T) \\ &\leq \mathbb{P}_{\underline{\mu}}(\exists n \geq C\sqrt{t}, \exists a \in [K], |\hat{\mu}_{a,n} - \mu_a| > \varepsilon, t \geq T) \\ &\leq \sum_{a \in [K]} \sum_{n \geq C\sqrt{t}} \mathbb{P}_{\underline{\mu}}(|\hat{\mu}_{a,n} - \mu_a| > \varepsilon, t \geq T) \\ &\leq \sum_{a \in [K]} \sum_{n \geq C\sqrt{t}} \mathbb{P}_{\underline{\mu}}(|\hat{\mu}_{a,n} - \mu_a| > \varepsilon). \end{aligned}$$

For $a \in [K]$ and $n \geq 0$, a Cramér-Chernoff bound gives

$$\begin{aligned} \mathbb{P}_{\underline{\mu}}(|\hat{\mu}_{a,n} - \mu_a| > \varepsilon) &= \mathbb{P}_{\underline{\mu}}(\hat{\mu}_{a,n} < \mu_a - \varepsilon) + \mathbb{P}_{\underline{\mu}}(\hat{\mu}_{a,n} > \mu_a + \varepsilon) \\ &\leq \exp(-nd(\mu_a - \varepsilon, \mu_a)) + \exp(-nd(\mu_a + \varepsilon, \mu_a)). \end{aligned}$$

Hence we get that

$$\begin{aligned} \mathbb{P}_{\underline{\mu}}(\mathcal{E}_t^c, t \geq T) &\leq \sum_{a \in [K]} \frac{\exp(-C\sqrt{t}d(\mu_a - \varepsilon, \mu_a))}{1 - \exp(-d(\mu_a - \varepsilon, \mu_a))} + \frac{\exp(-C\sqrt{t}d(\mu_a + \varepsilon, \mu_a))}{1 - \exp(-d(\mu_a + \varepsilon, \mu_a))} \\ &\leq 2KC_1 \exp(-C_2\sqrt{t}), \end{aligned} \quad (6.19)$$

where we set

$$C_1 \stackrel{\text{def}}{=} \max_{a \in [K]} \left(\max \left(\frac{1}{1 - \exp(-d(\mu_a - \varepsilon, \mu_a))}, \frac{1}{1 - \exp(-d(\mu_a + \varepsilon, \mu_a))} \right) \right) < +\infty,$$

and $C_2 \stackrel{\text{def}}{=} C \cdot \min_{a \in [K]} \left(\min(d(\mu_a - \varepsilon, \mu_a), d(\mu_a + \varepsilon, \mu_a)) \right) > 0.$

Injecting (6.19) into (6.18) finally leads to $\mathbb{E}_{\underline{\mu}}[T_\varepsilon] < +\infty$ and concludes the proof of (6.17). \square

Avoiding forced exploration. We hope to avoid the use of forced exploration with additional arguments. Especially since this sufficient exploration occurs for non-adaptive versions of TT-EB-TC and TT-EB-TCP, and it seems that the adaptive versions are kind of encouraging exploration: given a pair (ℓ, c) of leader and challenger, pulling arm ℓ (respectively c) will increase the probability of pulling c (respectively ℓ) at next time step the pair of leader and challenger will be the same. This is a consequence of the monotonicity of the function $x \in (0, 1) \mapsto \beta_{\ell, c}(\underline{\mu}, x)$ defined in (6.13), where we recall that x is interpreted as the proportion $\frac{v_\ell}{v_\ell + v_c}$ associated to a weight vector \underline{v} .

Lemma 6.11. *Let $a, b \in [K]$ be such that $\mu_a \neq \mu_b$. Then $\beta_{b, a}(\underline{\mu}, \cdot)$ is a decreasing (continuous) function such that*

$$\lim_{x \rightarrow 1} \beta_{b, a}(\underline{\mu}, x) = 0, \quad \text{and} \quad \lim_{x \rightarrow 0} \beta_{b, a}(\underline{\mu}, x) = 1.$$

Proof. The continuity of $\beta_{b, a}(\underline{\mu}, \cdot)$ is a direct consequence of the continuity properties of d (see page 43), since

$$\forall x \in (0, 1), \quad \beta_{b, a}(\underline{\mu}, x) \stackrel{\text{def}}{=} \frac{1}{1 + \frac{1-x}{x} \frac{d(\mu_a, (1-x)\mu_a + x\mu_b)}{d(\mu_b, (1-x)\mu_a + x\mu_b)}}$$

To prove the decreasing of $\beta_{b, a}(\underline{\mu}, \cdot)$, we use the strict convexity of d on its domain. Let $0 < x < x' < 1$. By setting $\bar{\mu}_x = (1-x)\mu_a + x\mu_b$, we get $\bar{\mu}_x \in (\mu_a, \bar{\mu}_{x'})$ or $\bar{\mu}_x \in (\bar{\mu}_{x'}, \mu_a)$, depending on the relative positions of μ_a and μ_b , hence there exists $\alpha \in (0, 1)$ such that

$$\bar{\mu}_x = \alpha \bar{\mu}_{x'} + (1-\alpha)\mu_a.$$

In fact, it is easy to see that $\alpha = \frac{x}{x'}$. By strict convexity of d , we get that

$$d(\mu_a, \bar{\mu}_x) \leq \alpha d(\mu_a, \bar{\mu}_{x'}) + (1-\alpha) \underbrace{d(\mu_a, \mu_a)}_{=0} = \frac{x}{x'} d(\mu_a, \bar{\mu}_{x'}),$$

hence we proved that

$$\frac{d(\mu_a, \bar{\mu}_x)}{x} < \frac{d(\mu_a, \bar{\mu}_{x'})}{x'}.$$

Similarly, one can show that

$$\frac{d(\mu_b, \bar{\mu}_x)}{1-x} > \frac{d(\mu_b, \bar{\mu}_{x'})}{1-x'},$$

which entails that $\beta_{b, a}(\underline{\mu}, x) > \beta_{b, a}(\underline{\mu}, x')$. This concludes the proof of the decreasing of $\beta_{b, a}(\underline{\mu}, \cdot)$.

We now prove the limit behaviors. We recall the expression of the mean-parameterized Kullback-Leibler divergence of the model defined in Equation (6.4):

$$\forall \mu, \mu' \in \mathcal{M}, \quad d(\mu, \mu') \stackrel{\text{def}}{=} \text{KL}(\nu_\theta, \nu_{\theta'}) = (\theta - \theta')b'(\theta) - b(\theta) + b(\theta'),$$

where $\theta = (b')^{-1}(\mu)$ and $\theta' = (b')^{-1}(\mu')$. As b is twice differentiable, a first-order Taylor expansion of $d(\mu, \cdot)$ ensures that for all $\mu \in \mathcal{M}$,

$$d(\mu, \mu') = o_{\mu' \rightarrow \mu}(\mu' - \mu).$$

Injecting this asymptotic behaviour in the expression (6.13) of $\beta_{b,a}(\underline{\mu}, x)$ leads to

$$\begin{aligned}\beta_{b,a}(\underline{\mu}, x) &= \frac{1}{1 + \frac{1-x}{x} \frac{d(\mu_a, (1-x)\mu_a + x\mu_b)}{d(\mu_b, (1-x)\mu_a + x\mu_b)}} \\ &= \frac{1}{1 + \frac{1-x}{x} \frac{o_{x \rightarrow 0}(x(\mu_b - \mu_a))}{d(\mu_b, (1-x)\mu_a + x\mu_b)}} \\ &= \frac{1}{1 + \frac{(\mu_b - \mu_a)}{d(\mu_b, \mu_a)} o_{x \rightarrow 0}(1)} \xrightarrow{x \rightarrow 0} 1.\end{aligned}$$

The second result can be derived similarly. \square

6.3.4. Step 4: A Useful Relationship

To obtain the required convergence (6.14), an important argument in the final step will be to use that, if there exists $a \neq a^*$ such that $p_a(t)/p_{a^*}(t)$ is a little bit above its target w_a/w_{a^*} , then if t is large enough, there exists another sub-optimal arm b such that $p_b(t)/p_{a^*}(t)$ is slightly under w_b/w_{a^*} . This property occurs, for instance (see the next step for precise calculations), as soon as

$$\sum_{a \neq a^*} \frac{p_a(t)}{p_{a^*}(t)} \xrightarrow[t \rightarrow +\infty]{\text{q.}} \sum_{a \neq a^*} \frac{w_a}{w_{a^*}}. \quad (6.20)$$

In the Gaussian case, [You et al. \(2023\)](#) were able to prove that

$$\sum_{a \neq a^*} \frac{p_a(t)^2}{p_{a^*}(t)^2} \xrightarrow[t \rightarrow +\infty]{\text{q.}} 1 = \sum_{a \neq a^*} \frac{w_a^2}{w_{a^*}^2},$$

and uses this relation to get the desired property. Yet, this clearly uses the fact that the optimal weight vector \underline{w} satisfies a property that is highly specific to Gaussian variables:

$$w_{a^*}^2 = \sum_{a \neq a^*} w_a^2,$$

that can be obtained directly from Equation (6.8) or by the analysis of Chapter 3, see Equation (3.12).

It seems that this Gaussian argument is the most challenging to extend the analysis of [You et al. \(2023\)](#) to all exponential models. A general relation that might be obtained to replace this Gaussian property is that the frequencies of pulls satisfy the following quick convergence:

$$\sum_{b \neq a^*} \frac{d(\mu^*, \bar{\mu}_{a^*, b, \underline{p}(t)})}{d(\mu_b, \bar{\mu}_{a^*, b, \underline{p}(t)})} \xrightarrow[t \rightarrow +\infty]{\text{q.}} 1 = \sum_{b \neq a^*} \frac{d(\mu^*, \bar{\mu}_{a^*, b, \underline{w}})}{d(\mu_b, \bar{\mu}_{a^*, b, \underline{w}})},$$

where the equality stands from Equation (6.8).

$p_{a^*}(t)$ **bounded away from 0**. An important implication of such convergence is that the ratio $p_a(t)/p_{a^*}(t)$ does not blow up, i.e., that $p_{a^*}(t)$ is bounded away from 0 after a sufficiently large enough time:

$$\exists p_{\min} > 0, \exists T \in \mathbb{L}^1, \quad \forall t \geq T, \quad p_{a^*}(t) \geq p_{\min}.$$

If another proof structure is required to replace this step, it might be possible to obtain this property⁴ by using Lemma 6.11, which, as previously discussed, indicates that the strategy somewhat encourages exploration.

⁴Even without forced exploration.

6.3.5. Step 5: Quick Convergence of $p_a(t)/p_{a^*}(t)$

In this last step we assume that the convergence (6.20) is satisfied:

$$\sum_{a \neq a^*} \frac{p_a(t)}{p_{a^*}(t)} \xrightarrow[t \rightarrow +\infty]{q.} \sum_{a \neq a^*} \frac{w_a}{w_{a^*}}. \quad (6.21)$$

In this step, we prove (6.14) with similar arguments to [You et al. \(2023\)](#). Some simplifications are made by the use of the deterministic challenger rules considered in this work (Transportation-Cost and Transportation-Cost-Penalized) instead of the Bayesian challenger of [You et al. \(2023\)](#).

Lemma 6.12. *For all $\varepsilon > 0$, there exists $T \in \mathbb{L}^1$ such that for all $t \geq T$ and $a \neq a^*$,*

$$\frac{p_a(t)}{p_{a^*}(t)} \geq \frac{w_a + \varepsilon}{w_{a^*}} \implies A_{t+1} \neq a.$$

This result states that, after a sufficiently long time, if arm a is over-pulled at time step t , it has no chance to be pulled at $t + 1$.

For the sake of completeness, we will give a proof of that result that was already known (see, e.g., [Jourdan et al., 2022](#), Page 7). We need to introduce notation and preliminary results. To simplify the reading, we set:

$$\forall t \geq 0, \forall a \neq b, \quad \widehat{\text{TC}}_{a \rightarrow b}(t) \stackrel{\text{def}}{=} \text{TC}_{a \rightarrow b}(\hat{\underline{\mu}}(t), \hat{\underline{p}}(t))$$

For $a \neq a^*$, we also define

$$f_a(x) \stackrel{\text{def}}{=} d\left(\mu^*, \frac{x}{1+x}\mu_a + \frac{1}{1+x}\mu^*\right) + x d\left(\mu_a, \frac{x}{1+x}\mu_a + \frac{1}{1+x}\mu^*\right),$$

and note that

$$\text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{v}) = v_{a^*} f_a\left(\frac{v_a}{v_{a^*}}\right),$$

As a consequence of Proposition 6.1, we have that

$$\forall a \neq a^*, \quad f_a\left(\frac{w_a}{w_{a^*}}\right) = \frac{T(\underline{\mu})^{-1}}{w_{a^*}}. \quad (6.22)$$

We will use the fact that f_a is an increasing function.

Lemma 6.13. [[Garivier and Kaufmann, 2016, Appendix 5.2](#)]

For all $a \neq a^$, $x \in (0, +\infty) \mapsto f_a(x)$ is an increasing function.*

We also use the sufficient exploration property (Lemma 6.9) to ensure that the empirical costs are quite close to their theoretical value after a sufficiently long time:

Lemma 6.14. *For all $\eta > 0$, there exists $T \in \mathbb{L}^1$ such that*

$$\forall t \geq T, \forall a \neq a^*, \quad (1 - \eta)\widehat{\text{TC}}_{a \rightarrow a^*}(t) \leq \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{p}(t)) \leq (1 + \eta)\widehat{\text{TC}}_{a \rightarrow a^*}(t).$$

Proof. This is a consequence of Proposition 6.10 and the continuity of the transportation costs. \square

Proof of Lemma 6.12. For the sake of simplicity, we only prove the result for the Transportation-Cost challenger. The same proof applies for the Transportation-Cost-Penalized challenger, up to slightly modifying Lemma 6.14 with the penalized empirical costs.

Let $\varepsilon > 0$. Thanks to the previous steps, we get the existence of some integrable times satisfying the following properties:

- By assumption (6.21), there exists $T_1 \in \mathbb{L}^1$ such that

$$\forall t \geq T_1, \quad \left| \sum_{a \neq a^*} \frac{p_a(t)}{p_{a^*}(t)} - \sum_{a \neq a^*} \frac{w_a}{w_{a^*}} \right| \leq \frac{\varepsilon}{w_{a^*}}. \quad (6.23)$$

- In addition, note that by Proposition 6.10, we get the existence of $T_2 \in \mathbb{L}^1$ such that for all $t \geq T_2$

$$\forall t \geq T_2, \forall a \in [K], \quad |\hat{\mu}_a(t) - \mu_a| < \frac{\Delta_{\min}}{2},$$

where $\Delta_{\min} = \min_{a \neq a^*} \Delta_a$. As a consequence $L_t = a^*$ for all $t \geq T_2$.

- Also, by Lemma 6.14, for a fixed $\eta > 0$ there exists $T_3 \in \mathbb{L}^1$ such that for all $t \geq T_3$,

$$\forall a \neq a^*, \quad (1 - \eta) \text{TC}_{a \rightarrow a^*}(t) \leq \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{p}(t)) \leq (1 + \eta) \text{TC}_{a \rightarrow a^*}(t). \quad (6.24)$$

We now set $T \stackrel{\text{def}}{=} \max(T_1, T_2, T_3) \in \mathbb{L}^1$ where T_3 is associated with a value of η to be chosen later. Let $t \geq T$ and assume that there exists $a \neq a^*$ such that

$$\frac{p_a(t)}{p_{a^*}(t)} \geq \frac{w_a + \varepsilon}{w_{a^*}}. \quad (6.25)$$

We will prove that $A_{t+1} \neq a$. As $L_{t+1} = a^*$, we only need to prove that $C_{t+1} \neq a$. By definition of C_{t+1} (which is an arm minimizing its empirical transportation cost with a^*), it suffices to prove that

$$\widehat{\text{TC}}_{b \rightarrow a^*}(t) < \widehat{\text{TC}}_{a \rightarrow a^*}(t). \quad (6.26)$$

By Equation (6.23), the condition (6.25) implies that

$$\exists b \neq a^*, \quad \frac{p_b(t)}{p_{a^*}(t)} \leq \frac{w_b}{w_{a^*}}. \quad (6.27)$$

Using the monotonicity property of f_a and f_b (Lemma 6.13), this leads to

$$\begin{aligned} (1 + \eta) \widehat{\text{TC}}_{a \rightarrow a^*}(t) &\geq \text{TC}_{a \rightarrow a^*}(\underline{\mu}, \underline{p}(t)) && \text{by (6.24)} \\ &= p_{a^*}(t) \cdot f_a\left(\frac{p_a(t)}{p_{a^*}(t)}\right) \\ &\geq p_{a^*}(t) \cdot f_a\left(\frac{w_a + \varepsilon}{w_{a^*}}\right) && \text{by (6.25)} \\ &= p_{a^*}(t) \cdot (1 + \varepsilon_a) \cdot f_a\left(\frac{w_a}{w_{a^*}}\right) \\ &= p_{a^*}(t) \cdot (1 + \varepsilon_a) \cdot f_b\left(\frac{w_b}{w_{a^*}}\right) && \text{by (6.22)} \\ &\geq (1 + \varepsilon_a) \cdot p_{a^*}(t) \cdot f_b\left(\frac{p_b(t)}{p_{a^*}(t)}\right) && \text{by (6.27)} \\ &= (1 + \varepsilon_a) \cdot \text{TC}_{b \rightarrow a^*}(\underline{\mu}, \underline{p}(t)) && \text{by (6.24)} \\ &\geq (1 + \varepsilon_a)(1 - \eta) \widehat{\text{TC}}_{b \rightarrow a^*}(t), \end{aligned}$$

where we defined

$$\varepsilon_a \stackrel{\text{def}}{=} \frac{f_a\left(\frac{w_a + \varepsilon}{w_{a^*}}\right)}{f_a\left(\frac{w_a}{w_{a^*}}\right)} - 1,$$

which is positive by Lemma 6.13. Hence we proved that

$$\widehat{\text{TC}}_{a \rightarrow a^*}(t) \geq \frac{(1 - \eta)}{(1 + \eta)}(1 + \varepsilon_a) \widehat{\text{TC}}_{b \rightarrow a^*}(t).$$

Taking η small enough, such that

$$\frac{(1 - \eta)}{(1 + \eta)} \left(1 + \min_{a \neq a^*} \varepsilon_a\right) > 1,$$

leads to inequality (6.26), which enforces $C_t \neq a$ and completes the proof. \square

Using Lemma 6.12, it is not hard to prove the convergence of $p_a(t)/p_{a^*}(t)$. The analysis of [You et al. \(2023\)](#) might be followed, by proving first that, for all $\varepsilon > 0$, there exists $T \in \mathbb{L}^1$ such that

$$\forall t \geq T, \forall a \neq a^*, \quad \frac{p_a(t)}{p_{a^*}(t)} \leq \frac{w_a + \varepsilon}{w_{a^*}},$$

and then deduce the required objective (6.14) by using assumption (6.21).

6.3.6. Conclusion

In this section, we proposed some ideas to generalize the Gaussian analysis of [You et al. \(2023\)](#) and prove Conjecture 6.4. In steps 2 and 3, we slightly modified the definition of the strategies TT-EB-TC and TT-EB-TCP in order to simplify their analysis for an exponential model. In step 5, we proved that the Gaussian arguments still apply for those models, provided that some hypothetical condition holds. It seems that the last biggest challenge is to prove that this condition is satisfied (see discussions in step 4).

About the χ challenger. The analysis made in this section might be extended to the χ challenger presented in Section 6.2.3, but the randomization of the challenger complicates the analysis as well as technical issues with non-continuous models. For this challenger, a first step might be to obtain the convergence of the empirical frequencies of pulls with the knowledge of $\underline{\mu}$, stated in the following conjecture.

Conjecture 6.15. *Consider a version of the TT-EB- χ algorithm which knows the values of $\underline{\mu}$ and use it instead of $\hat{\mu}(t-1)$ to compute $\hat{\chi}(t-1)$ and $\hat{\beta}_{L_t, C_t}(t)$. Then*

$$\underline{p}(t) \xrightarrow[t \rightarrow +\infty]{} \underline{w} \quad \mathbb{P}_{\underline{\mu}}\text{-a.s.}$$

6.4. Side Note: On the Asymptotic Optimality of Track-and-Stop

In this section, we give a proof of the asymptotic optimality of Track-and-Stop. When introducing their strategy, [Garivier and Kaufmann \(2016\)](#) proved that Track-and-Stop satisfies the following asymptotic bound, which does not exactly correspond to asymptotic optimality in the sense of Definition 2.12.

Theorem 6.16. [Garivier and Kaufmann, 2016, Theorem 14]

Consider an exponential model \mathcal{D}_{exp} . The Track-and-Stop strategy, with threshold (6.3) for a fixed $\alpha > 1$, satisfies

$$\forall \underline{\mu} \text{ in } \mathcal{D}_{\text{exp}}, \quad \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \leq \alpha T(\underline{\mu}).$$

As a by-product of some of the results stated in Section 6.4, we show that Track-and-Stop is asymptotically optimal whatever the choice of $\alpha > 1$.

Theorem 6.17. Consider an exponential model \mathcal{D}_{exp} . The Track-and-Stop strategy, with C-tracking and the Global-Likelihood-Ratio stopping rule with threshold (6.3) for a fixed $\alpha > 1$, is asymptotically optimal:

$$\forall \underline{\mu} \text{ in } \mathcal{D}_{\text{exp}}, \quad \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\underline{\mu}}[\tau_{\delta}]}{\log \frac{1}{\delta}} \leq T(\underline{\mu}).$$

Remark. The proposition also applies for the Exploration-Biased-Sampling strategy studied in Chapter 4. One can follow the same proof structure, recalling that this strategy naturally has sufficient exploration with rate \sqrt{t} .

The result is deduced from the general Propositions 6.6 and 6.10.

Proof. We will prove that the hypothesis of Proposition 6.6 is satisfied, i.e., that

$$\underline{p}(t) \xrightarrow[t \rightarrow +\infty]{\text{q.}} \underline{w}.$$

We rely on the following lemma, which controls the values of the number of pulls.

Lemma 6.18. [Garivier and Kaufmann, 2016, Lemma 7]

The Track-and-Stop algorithm with C-tracking satisfies, for all $t \geq 0$ and $a \in [K]$,

$$N_a(t) \geq \sqrt{t + K^2} - 2K, \tag{6.28}$$

$$\text{and } \left| p_a(t) - \frac{1}{t} \sum_{s \in [t]} \hat{w}(s) \right| \leq K \frac{1 + \sqrt{t}}{t}. \tag{6.29}$$

We deduce from (6.28) that the strategy provides sufficient exploration with rate \sqrt{t} , for some constant $C \in (0, 1)$. By Proposition 6.10, this implies that

$$\hat{\underline{\mu}}(t) \xrightarrow[t \rightarrow +\infty]{\text{q.}} \underline{\mu},$$

and by continuity of $\hat{\underline{\mu}} \mapsto \underline{w}(\hat{\underline{\mu}})$ at $\underline{\mu}$, that

$$\hat{\underline{w}}(t) = \underline{w}(\hat{\underline{\mu}}(t)) \xrightarrow[t \rightarrow +\infty]{\text{q.}} \underline{w}.$$

It can be easily seen that this implies the quick convergence of the Cesàro sum:

$$\frac{1}{t} \sum_{s \in [t]} \hat{w}(s) \xrightarrow[t \rightarrow +\infty]{\text{q.}} \underline{w}.$$

By Equation (6.29), this implies that

$$\underline{p}(t) \xrightarrow[t \rightarrow +\infty]{\mathfrak{q}} \underline{w},$$

and concludes the proof. \square

6.5. Conclusion

In this chapter, we first considered a new transformation of weights which is of high interest for the sample complexity optimization problem, as its unique fixed point is the optimal weight vector. We derived procedures for the computation of the optimal weight vector that revealed empirically correct and efficient, but we are still looking for theoretical arguments that to guarantee the convergence.

Then, we came to grips with adaptive top-two algorithms in order to prove their asymptotic optimality for a general exponential model. Following the Gaussian proof structure of [You et al. \(2023\)](#), we showed that some arguments might easily be generalized or bypassed (in a preliminary study, by using, e.g., forced exploration). However, a crucial convergence argument presented in step 4 is missing in the general case, and constitutes the main remaining challenge, unless different proof structures avoiding this argument are studied. Indeed, in future work, we might also investigate another path based on the three following steps:

- prove Conjecture 6.3, that is, prove the convergence of the empirical procedures proposed in Section 6.2.2 for the computation of $\underline{w}(\underline{\mu})$ with the knowledge of $\underline{\mu}$,
- independently, prove that the TT-EB-TC and/or TT-EB-TCP strategies have sufficient exploration, using the partial ideas introduced in Section 6.3.3,
- finally, prove the asymptotic optimality by combining the results of the two first phases.

Bibliography

- S. Agrawal, S. Juneja, and P. Glynn. Optimal δ -Correct Best-Arm Selection for Heavy-Tailed Distributions. In *Proceedings of the 31st International Conference on Algorithmic Learning Theory*, volume 117 of *Proceedings of Machine Learning Research*, pages 61–110. PMLR, 2020.
- S.M. Ali and S.D Silvey. A General Class of Coefficients of Divergence of One Distribution from Another. *Journal of the Royal Statistical Society*, 28:131–142, 1966.
- J.-Y. Audibert, S. Bubeck, and R. Munos. Best Arm Identification in Multi-Armed Bandits. In *Proceedings of the 23th Conference on Learning Theory*, 2010.
- A. Barrier, A. Garivier, and T. Kocák. A Non-Asymptotic Approach to Best-Arm Identification for Gaussian Bandits. In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, pages 10078–10109. PMLR, 2022.
- A. Barrier, A. Garivier, and G. Stoltz. On Best-Arm Identification with a Fixed Budget in Non-Parametric Multi-Armed Bandits. In *Proceedings of the 34th International Conference on Algorithmic Learning Theory*, volume 201 of *Proceedings of Machine Learning Research*, pages 136–181. PMLR, 2023.
- R.G. Bartle and D.R. Sherbert. *Introduction to Real Analysis*. John Wiley & Sons, 3rd edition, 2000.
- S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.
- J. Bretagnolle and C. Huber. Estimation des densités : risque minimax. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 47, 1979.
- A.N. Burnetas and M.N. Katehakis. Optimal Adaptive Policies for Sequential Allocation Problems. *Advances in Applied Mathematics*, 17(2):122–142, 1996.
- A. Carpentier and A. Locatelli. Tight (Lower) Bounds for the Fixed Budget Best Arm Identification Bandit Problem. In *Proceedings of the 29th Conference on Learning Theory*, volume 49, pages 590–604. PMLR, 2016.
- L. Chen, J. Li, and M. Qiao. Towards Instance Optimal Bounds for Best Arm Identification. In *Proceedings of the 2017 Conference on Learning Theory*, volume 65, pages 535–592. PMLR, 2017.
- H. Chernoff. Sequential Design of Experiments. *The Annals of Mathematical Statistics*, 30(3): 755–770, 1959.
- Y. Chow and H. Teicher. *Probability Theory*. Springer, 1988.
- S.-R. Chowdhury, P. Saux, O.-A. Maillard, and A. Gopalan. Bregman Deviations of Generic Exponential Families. In *Proceedings of the 36th Conference on Learning Theory*, volume 195, pages 394–449. PMLR, 2023.

BIBLIOGRAPHY

- T.M. Cover and J.A. Thomas. *Elements of Information Theory*. Wiley, 2nd edition, 2006.
- R. Degenne. On the Existence of a Complexity in Fixed Budget Bandit Identification. In *Proceedings of the 36th Conference on Learning Theory*, volume 195, pages 1131–1154. PMLR, 2023.
- R. Degenne and W.M. Koolen. Pure Exploration with Multiple Correct Answers. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- R. Degenne, W.M. Koolen, and P. Ménard. Non-Asymptotic Pure Exploration by Solving Games. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- J.L. Doob. *Stochastic Processes*. Wiley Publications in Statistics. John Wiley & Sons, 1953.
- E. Even-Dar, S. Mannor, and Y. Mansour. Action Elimination and Stopping Conditions for the Multi-Armed Bandit and Reinforcement Learning Problems. *Journal of Machine Learning Research*, 7(39):1079–1105, 2006.
- V. Gabillon, M. Ghavamzadeh, and A. Lazaric. Best Arm Identification: A Unified Approach to Fixed Budget and Fixed Confidence. In *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.
- A. Garivier and É. Kaufmann. Optimal Best Arm Identification with Fixed Confidence. In *Proceedings of the 29th Conference on Learning Theory*, volume 49, pages 998–1027. PMLR, 2016.
- A. Garivier and É. Kaufmann. Non-Asymptotic Sequential Tests for Overlapping Hypotheses Applied to Near-Optimal Arm Identification in Bandit Models. *Sequential Analysis*, 40(1):61–96, 2021.
- A. Garivier, P. Ménard, and G. Stoltz. Explore First, Exploit Next: The True Shape of Regret in Bandit Problems. *Mathematics of Operations Research*, 44(2):377–399, 2019.
- A. Garivier, H. Hadji, P. Ménard, and G. Stoltz. KL-UCB-switch: Optimal Regret Bounds for Stochastic Bandits from Both a Distribution-Dependent and a Distribution-Free Viewpoints. *Journal of Machine Learning Research*, 23(179):1–66, 2022.
- G.R. Grimmett and D.R. Stirzaker. *Probability and Random Processes*. Oxford University Press, 3rd edition, 2001.
- J. Honda and A. Takemura. Non-Asymptotic Analysis of a New Bandit Algorithm for Semi-Bounded Rewards. *Journal of Machine Learning Research*, 16(113):3721–3756, 2015.
- L.J. Hong, W. Fan, and J. Luo. Review on Ranking and Selection: A New Perspective. *Frontiers of Engineering Management*, 8(3):321–343, 2021.
- S.R. Howard, A. Ramdas, J. McAuliffe, and J. Sekhon. Time-Uniform, Nonparametric, Nonasymptotic Confidence Sequences. *The Annals of Statistics*, 49(2):1055–1080, 2021.
- K. Jamieson, M. Malloy, R. Nowak, and S. Bubeck. *lil' UCB* : An Optimal Exploration Algorithm for Multi-Armed Bandits. In *Proceedings of the 27th Conference on Learning Theory*, volume 35, pages 423–439. PMLR, 2014.
- M. Jourdan and R. Degenne. Non-Asymptotic Analysis of a UCB-Based Top Two Algorithm, 2023. Preprint, arXiv:2210.05431.
- M. Jourdan, R. Degenne, D. Baudry, R. de Heide, and É. Kaufmann. Top Two Algorithms Revisited. In *Advances in Neural Information Processing Systems*, volume 35. Curran Associates, Inc., 2022.

BIBLIOGRAPHY

- S. Kalyan Krishnan, A. Tewari, P. Auer, and P. Stone. PAC Subset Selection in Stochastic Multi-armed Bandits. In *Proceedings of the 29th International Conference on Machine Learning*, 2012.
- Z. Karnin, T. Koren, and O. Somekh. Almost Optimal Exploration in Multi-Armed Bandits. In *Proceedings of the 30th International Conference on Machine Learning*, volume 28, pages 1238–1246. PMLR, 2013.
- M. Kato, K. Ariu, M. Imaizumi, M. Nomura, and C. Qin. Optimal Best Arm Identification in Two-Armed Bandits with a Fixed Budget under a Small Gap, 2022. Preprint, arXiv:2201.04469.
- É. Kaufmann and S. Kalyan Krishnan. Information Complexity in Bandit Subset Selection. In *Proceedings of the 26th Annual Conference on Learning Theory*, volume 30, pages 228–251. PMLR, 2013.
- É. Kaufmann and W.M. Koolen. Mixture Martingales Revisited with Applications to Sequential Tests and Confidence Intervals. *Journal of Machine Learning Research*, 22(246):1–44, 2021.
- É. Kaufmann, O. Cappé, and A. Garivier. On the Complexity of Best-Arm Identification in Multi-Armed Bandit Models. *Journal of Machine Learning Research*, 17(1):1–42, 2016.
- T. Kocák and A. Garivier. Best Arm Identification in Spectral Bandits. In *Proceedings of the 29th International Joint Conference on Artificial Intelligence*, volume 3, pages 2220–2226. International Joint Conferences on Artificial Intelligence Organization, 2020.
- J. Komiyama. Suboptimal Performance of the Bayes Optimal Algorithm in Frequentist Best Arm Identification, 2022. Preprint, arXiv:2202.05193.
- J. Komiyama, T. Tsuchiya, and J. Honda. Minimax Optimal Algorithms for Fixed-Budget Best Arm Identification. In *Advances in Neural Information Processing Systems*, volume 35. Curran Associates, Inc., 2022.
- T.L. Lai and H. Robbins. Asymptotically Efficient Adaptive Allocation Rules. *Advances in Applied Mathematics*, 6(1):4–22, 1985.
- T. Lattimore and C. Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- E.L. Lehmann and G. Casella. *Theory of Point Estimation*. Springer Texts in Statistics. Springer, 2nd edition, 1998.
- S. Mannor and J.N. Tsitsiklis. The Sample Complexity of Exploration in the Multi-Armed Bandit Problem. *Journal of Machine Learning Research*, 5:623–648, 2004.
- P. Ménard. Gradient Ascent for Active Exploration in Bandit Problems, 2019. Preprint, arXiv:1905.08165.
- C. Qin and D. Russo. Electronic Companion to Adaptivity and Confounding in Multi-Armed Bandit Experiments, 2022. SSRN:4115833.
- C. Qin, D. Klabjan, and D. Russo. Improving the Expected Improvement Algorithm. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- D. Russo. Simple Bayesian Algorithms for Best Arm Identification. In *Proceedings of the 29th Conference on Learning Theory*, volume 49, pages 1417–1418. PMLR, 2016.
- D. Russo. Simple Bayesian Algorithms for Best Arm Identification. *Operations Research*, 68(6): 1625–1647, 2020.

BIBLIOGRAPHY

- X. Shang, R. de Heide, É. Kaufmann, P. Ménard, and M. Valko. Fixed-Confidence Guarantees for Bayesian Best-Arm Identification. In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 1823–1832. PMLR, 2020.
- M. Simchowitz, K. Jamieson, and B. Recht. The Simulator: Understanding Adaptive Sampling in the Moderate-Confidence Regime. In *Proceedings of the 2017 Conference on Learning Theory*, volume 65, pages 1794–1834. PMLR, 2017.
- G. Stoltz. Lecture Notes on “Learning & Sequential Optimization”, 2022.
- W.R. Thompson. On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of Two Samples. *Biometrika*, 25(3/4):285–294, 1933.
- M.J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019.
- P.-A. Wang, R.-C. Tzeng, and A. Proutiere. Fast Pure Exploration via Frank-Wolfe. In *Advances in Neural Information Processing Systems*, volume 34, pages 5810–5821. Curran Associates, Inc., 2021.
- W. You, C. Qin, Z. Wang, and S. Yang. Information-Directed Selection for Top-Two Algorithms. In *Proceedings of the 36th Conference on Learning Theory*, volume 195, pages 2850–2851. PMLR, 2023.

Index

A/B testing	32, 53	elimination algorithm	36, 131
algorithm		candidate arms	36, 60, 131
Chernoff-Racing	104	phase	36, 60, 131
Exploration-Biased-Sampling	54, 96,	round	36, 60
97, 104		empirical mean	35
Frank-Wolfe-based-Sampling	55	exponential family	43, 72
Lazy-Mirror-Ascent	55	divergence	43
LUCB++	55, 104	natural parameter space	43
Successive-Elimination	36	normalizing function	43
Sequential-Halving	62, 127	external randomization	35
Successive-Rejects	60, 125, 131	finite risk bounds	53
top-two	55, 163	fixed-point property	166
TT-EB-TC	57, 174	forced exploration	49, 53, 69
TT-EB-TCP	57, 174	fundamental inequality	40, 45, 63
TT-EB-TCP- β	57	gap	36, 74, 125
TT-EB- χ	171	generalized log-likelihood ratio	50
TT-UCB-TC- β	57	history	35
Track-and-Stop	49, 104, 162	Hoeffding's inequality	38, 58, 125
Uniform-Sampling	58	information-theoretic quantities	125
alternative bandits	44, 66	Kullback-Leibler – χ^2 -divergence bound	65, 128
arm	33, 94	Kullback-Leibler divergence	40, 73, 126
best	33, 72, 94	learner	32, 49, 94, 124
optimal	33	medical trial	32, 34
sub-optimal	33	misidentification probability	38, 58
bandit		model	33, 124
instance	33	σ^2 -sub-Gaussian	33
multi-armed	33, 94, 124	Bernoulli	33
one-armed	33	distributions on $[0, 1]$	33, 145
best-arm identification	33	exponential	33, 43, 129, 162
fixed-budget	34, 58, 124	Gaussian	33, 73, 94
fixed-confidence	33, 35, 72, 94, 162	multi-armed bandit	72
Bretagnolle-Huber inequality	64, 126, 159	Newton's method	77
budget parameter	34, 58	number of pulls	35
chain rule	42	online advertising	32, 34
change of measure	40	optimal weight vector	47, 54, 72, 95, 162
characteristic time	44, 72, 94, 162	optimization problem	44, 45, 72
complexity function	69	optional skipping	37, 97, 126, 134
confidence interval	39	optional stopping theorem	42
confidence parameter	33, 35, 162		
confidence region	54, 97, 100		
Cramér-Chernoff bound	132, 147		
data-processing inequality	40, 145		

INDEX

peeling argument	101	Empirical-Best	35, 51, 102
Pinsker's inequality	67, 128, 145	sampling rule	35, 49, 94, 102, 137, 163
procedure		stopping rule	35, 50, 94
Gaussian-Optimal-Weights	77, 98	Global-Likelihood-Ratio	50, 102
Optimal-Weights	47, 49, 73, 167	sufficient exploration	103, 177
push-forward measure	40	threshold	50
ranking and selection	34	Top-Two algorithm	
reverse order statistics	37, 60, 127	adaptive	57, 164
sample complexity	39, 72, 94	challenger	55, 163
sequence of strategies	59, 137	χ	171
consistent	63, 138	Transportation-Cost	55, 163
doubly-indexed	137	Transportation-Cost-Penalized	55,
exponentially consistent	63, 138	163	
monotonous	68, 141	Thompson-Sampling	55
sequential learning	32	leader	55, 163
simplex	44	Empirical-Best	55, 163
strategy	35, 35 , 137	Thompson-Sampling	55
δ -correct	34, 35 , 72, 94, 103, 162	Upper-Confidence-Bound	55
Exploration-Biased-Sampling	185	tracking	49, 176
Track-and-Stop	184	C-tracking	49, 102
adaptive	69	conservative	97
asymptotically optimal	52	D-tracking	49, 102
decision rule	35, 94, 137	transition kernel	41
		transportation cost	46, 163

Contributions à une théorie de l'exploration pure en statistique séquentielle

Résumé. Cette thèse, à la croisée entre les domaines de l'intelligence artificielle, de la statistique séquentielle et de l'optimisation, s'intéresse au problème d'identification du meilleur bras (en espérance) dans les bandits non structurés à K bras. Ce problème possède deux approches dont les niveaux de compréhension sont très différents.

Le cadre à confiance fixée est le mieux compris : des stratégies asymptotiquement optimales sont connues, et l'on s'intéresse à l'obtention de garanties non asymptotiques pour des stratégies (si possible) simples et naturelles. Avec des bandits Gaussiens, nous proposons l'analyse à risque fini d'une nouvelle stratégie (asymptotiquement optimale) grâce aux propriétés de régularité de ce modèle. Cette stratégie modifie subtilement la règle d'attribution des tirages de l'algorithme Track-and-Stop en une règle plus prudente et interprétable. Dans le contexte plus général d'un modèle exponentiel, nous proposons l'ébauche d'une analyse de l'asymptotique optimalité d'algorithmes de type Top-Two adaptatifs, dont les règles de choix de tirages sont particulièrement simples.

Par ailleurs, dans le cadre à budget fixé, où l'existence d'une hypothétique complexité reste à démontrer, nous proposons des généralisations à des modèles non-paramétriques des bornes (supérieures et inférieures) connues jusqu'à présent pour des modèles très spécifiques. Les bornes obtenues font intervenir des quantités de théorie de l'information plus précises que les écarts entre les moyennes qui apparaissaient précédemment. Ces nouvelles quantités pourraient être la clé pour mesurer la complexité de l'identification de meilleur bras à budget fixé.

Mots-clés. Problèmes de bandits · Identification de meilleur bras · Statistiques séquentielles · Apprentissage statistique · Intelligence artificielle

Contributions to a Theory of Pure Exploration in Sequential Statistics

Abstract. This thesis lies in the fields of artificial intelligence, sequential statistics and optimization. We focus on the problem of best (in expectation) arm identification in unstructured multi-armed bandits. This problem has two approaches with very different levels of understanding.

The fixed-confidence framework is the best understood: asymptotically optimal strategies are known, and we are interested in obtaining non-asymptotic guarantees for (if possible) simple and natural strategies. Working with Gaussian bandits, we propose a finite risk analysis of a new (asymptotically optimal) strategy using the regularity properties of this model. This strategy slightly modifies the sampling rule of the Track-and-Stop algorithm into a more conservative and interpretable rule. In the more general context of an exponential model, we propose a preliminary analysis of the asymptotic optimality of adaptive Top-Two algorithms, whose sampling rules are particularly simple.

Independently, in the fixed-budget framework, for which the existence of a hypothetical complexity remains to be demonstrated, we propose generalizations to non-parametric models of the existing bounds (upper and lower) that were available so far only for very specific models. The obtained bounds involve more precise information-theoretic quantities than the gaps (differences between the means) which appeared previously. These new quantities could be the key to measuring the complexity of fixed-budget best-arm identification.

Keywords. Multi-Armed Bandits · Best-Arm Identification · Sequential Statistics · Statistical Learning · Machine Learning