



HAL
open science

Markov Substitute Processes : a statistical model for linguistics

Thomas Mainguy

► **To cite this version:**

Thomas Mainguy. Markov Substitute Processes : a statistical model for linguistics. General Mathematics [math.GM]. Université Pierre et Marie Curie - Paris VI, 2014. English. NNT : 2014PA066354 . tel-01127344

HAL Id: tel-01127344

<https://theses.hal.science/tel-01127344>

Submitted on 7 Mar 2015

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.

Université Pierre et Marie Curie
École doctorale de sciences mathématiques de Paris Centre

*Département de Mathématiques et Applications de l'École Normale Supérieure
Équipe d'accueil : Probabilités et Statistiques*

**Processus de substitution
Markoviens**

un modèle statistique pour la linguistique

Markov Substitute Processes
a statistical model for linguistics

par Thomas MAINGUY

Thèse de doctorat de Mathématiques

dirigée par Olivier CATONI

Présentée et soutenue publiquement le 11 décembre 2014
devant un jury composé de

M. Pierre ALQUIER	ENSAE	Examineur
M. Gérard BIAU	UPMC	Examineur
M. Olivier CAPPÉ	Telecom Paris Tech-CNRS	Examineur
M. Olivier CATONI	ENS	Examineur
M. Antoine CHAMBAZ	Université Paris Ouest Nanterre	Examineur
Mme Elisabeth GASSIAT	Université Paris XI Orsay	Rapporteur
Mme Catherine MATIAS	CNRS	Examineur

au vu des rapports de

Mme Elisabeth GASSIAT Université Paris XI Orsay
M. Ramon VAN HANDEL Princeton University

Remerciements

Mes premiers remerciements vont à Olivier CATONI, pour m'avoir proposé ce sujet à la fin de mon stage de M2, rebondissant de manière inattendue sur une remarque en passant. Pour avoir su me guider comme il l'a fait, malgré les incertitudes et les difficultés initiales, m'avoir poussé à chercher toujours plus loin, et m'avoir montré ce qu'est le travail de chercheur.

Mes remerciements vont aussi à tous ceux qui m'ont mené à cette thèse, les excellents professeurs de probabilités et statistiques que j'ai eu lors de ma scolarité à l'École et à l'université d'Orsay, Wendelin WERNER, Elisabeth GASSIAT, Pascal MASSART, Gilles STOLTZ, ainsi que les linguistes du DEC qui m'ont initié à cette science qui m'attirait depuis tant de temps sans que je le sache, en particulier Dominique SPORTICHE et Edward STABLER (UCLA). Sans oublier mes professeurs de classes préparatoires Hélène BENHAMOU et Nicolas TOSEL, qui m'ont montré la voie des mathématiques et sans qui je n'aurais sans doute pas eu la chance d'étudier à l'École Normale.

Merci enfin à tous ces non-mathématiciens sans qui tout ceci n'aurait pas été possible, pour leur soutien et leur intérêt, à mes parents, qui m'ont guidé et soutenu, poussé et motivé, depuis toujours et particulièrement dans les passages difficiles de cette thèse, toujours là pour écouter mes déboires et victoires, même s'ils ne comprenaient pas tout... Merci aussi à tous mes camarades du club d'escrime, Raphaël, Samuel, Romain, et tous les autres, pour me suivre dans cette passion. Et à Émeline, qui m'a donné tant d'énergie et de motivation. Merci, du fond du cœur.

Abstract:

This thesis proposes a new approach to natural language processing. Rather than trying to estimate directly the probability distribution of a random sentence, we will detect syntactic structures in the language, which can be used to modify and create new sentences from an initial sample.

The study of syntactic structures will be done using Markov substitute sets, sets of strings that can be freely substituted in any sentence without affecting the whole distribution. These sets define the notion of Markov substitute processes, modelling conditional independence of certain substrings (given by the sets) with respect to their context. This point of view splits the issue of language analysis into two parts, a model selection stage where Markov substitute sets are selected, and a parameter estimation stage where the actual frequencies for each set are estimated.

We show that these substitute processes form exponential families of distributions, when the language structure (the Markov substitute sets) is fixed. On the other hand, when the language structure is unknown, we propose methods to identify Markov substitute sets from a statistical sample, and to estimate the parameters of the distribution. Markov substitute sets show some connections with context-free grammars, that can be used to help the analysis. We then proceed to build invariant dynamics for Markov substitute processes. They can among other things be used to effectively compute the maximum likelihood estimate. Indeed, Markov substitute models can be seen as the thermodynamical limit of the invariant measure of crossing-over dynamics.

Keywords: Markov processes; Natural language processing; Metropolis algorithm; PAC-Bayes hypothesis testing; Formal grammars; Statistical machine learning

Résumé :

Ce travail de thèse propose une nouvelle approche au traitement des langues naturelles. Plutôt qu'essayer d'estimer directement la probabilité d'une phrase quelconque, nous identifions des structures syntaxiques dans le langage, qui peuvent être utilisées pour modifier et créer de nouvelles phrases à partir d'un échantillon initial.

L'étude des structures syntaxiques est accomplie avec des ensembles de substitution Markoviens, ensembles de chaînes de caractères qui peuvent être échangées sans affecter la distribution. Ces ensembles définissent des processus de substitution Markoviens qui modélisent l'indépendance conditionnelle de certaines chaînes vis-à-vis de leur contexte. Ce point de vue décompose l'analyse du langage en deux parties, une phase de sélection de modèle, où les ensembles de substitution sont sélectionnés, et une phase d'estimation des paramètres, où les fréquences pour chaque ensemble sont estimées.

Nous montrons que ces processus constituent des familles exponentielles quand la structure du langage est fixée. Lorsque la structure du langage est inconnue, nous proposons des méthodes pour identifier des ensembles de substitution à partir d'un échantillon, et pour estimer les paramètres de la distribution. Les ensembles de substitution ont quelques relations avec les grammaires hors-contexte, qui peuvent être utilisées pour aider l'analyse. Nous construisons alors des dynamiques invariantes pour les processus de substitution. Elles peuvent être utilisées pour calculer l'estimateur du maximum de vraisemblance. En effet, les processus de substitution peuvent être vus comme la limite thermodynamique de la mesure invariante d'une dynamique de crossing-over.

Mots-clefs : Processus markoviens ; Analyse des langues naturelles ; Algorithme Metropolis ; Tests d'hypothèse PAC-Bayesiens ; Grammaires formelles ; Apprentissage statistique

Contents

Introduction	1
0.1 Context	1
0.1.1 Natural language processing	1
0.1.2 Remarks on the shortcomings of standard models	1
0.1.3 Syntactic analysis	2
0.1.4 Markov substitute processes and toric grammars in a few words	4
0.2 Overview of the main results	7
0.2.1 Toric grammars	7
0.2.2 Markov substitute sets and processes	9
0.2.3 Model selection: finding a collection of Markov substitute sets	11
0.2.4 Parameter estimation	14
0.2.5 Invariant dynamics	15
1 Toric grammars and communication models	17
1.1 Introduction to a new communication model	17
1.2 First definitions	19
1.2.1 Toric grammars	19
1.2.2 A roadmap towards a communication model	23
1.3 Operations on toric grammars	24
1.3.1 Non stochastic syntax splitting and merging	24
1.3.2 Random split and merge processes	28
1.3.3 Splitting rules and label identification	32
1.4 Parsing and generalization	35
1.5 Expectation of a random toric grammar	37
1.6 Language models	38
1.6.1 Communication model	38
1.6.2 Comparison with other models	39
1.7 A small experiment	43
1.8 The story so far...	48
1.A Proofs	49

1.A.1	Bound on the length of splitting and production processes	49
1.A.2	Parsing Relations	51
1.A.3	Convergence to the expectation of a random toric grammar	54
1.B	Language produced by a toric grammar	60
2	Markov substitute Sets	63
2.1	Presentation of the model	63
2.1.1	Motivation	63
2.1.2	What is a Markov substitute set	63
2.1.3	Weak Markov substitute sets	65
2.1.4	Basic properties of Markov substitute sets	66
2.1.5	Interpretation in terms of random parsing	67
2.2	Invariant dynamics	69
2.2.1	Metropolis invariant dynamics on sentences	69
2.2.2	Reflecting a reversible dynamics on the boundary of a finite domain	70
2.2.3	Compound dynamics	71
2.2.4	A simple example of recursive structure	71
2.2.5	Crossing-over reversible dynamics on texts	73
2.3	Exponential families of Markov substitute processes	74
2.4	Testing Markov substitute sets	83
2.4.1	Alternative construction of a test function	84
2.5	Weakening the Markov substitute assumption	87
2.6	Testing using a parse process	88
2.6.1	Probability of false rejection	89
2.6.2	Probability of false acceptance of the hypothesis	91
2.7	Testing Markov substitute sets without parsing	96
2.7.1	Definition of the test and probability of false rejection	100
2.7.2	Testing for the weak Markov substitute property	102
2.7.3	Computation of the test	103
2.7.4	Some numerical examples	107
2.7.5	Probability of false acceptance of the test	111
2.8	Estimation of the substitute measure	116
3	Markov substitute sets and language	119
3.1	Production rules and Markov substitute sets	119
3.1.1	Markov grammars	119
3.1.2	Parse trees	122
3.1.3	General introduction to parsing	123
3.2	Substitute measures and reversible dynamics	125
3.2.1	Parametrization of substitute measures	125

3.2.2	Estimating substitute measures	127
3.2.3	Simulating substitute measures	129
3.2.4	Reversible dynamics for the language distribution	131
3.2.5	Crossing-over dynamics	134
3.3	Crossing-over dynamics and the maximum likelihood estimator . . .	135
3.4	Building a Markov ruleset	144
3.4.1	Adding new rules	144
3.4.2	Reducing the context space	145
3.4.3	Saturation	146
3.4.4	Identification of labels	150
3.5	Toric grammars	151
3.5.1	Split and merge processes using parsing	151
3.5.2	Reversible split and merge process	152
3.6	Estimating the language distribution	155
3.A	Support of the split and merge process and substitutions	157
3.B	Parsing using a ruleset	159
3.B.1	Parsing general syntagms	161
3.B.2	Parsing inside syntagms of certain type	162
3.C	Building crossing-over tree kernels	163
3.D	Building general Metropolis reversible dynamics	165
4	Conclusion	167
4.1	Possible uses of Markov substitute sets	167
4.2	Further considerations on the scope of the model	169
	Bibliography	174

Introduction

0.1 Context

0.1.1 Natural language processing

The goal of this thesis is to propose some new theoretical tools for the statistical analysis of natural languages. This will lead us to introduce new models with a quite general purpose, that could be applied to approximate any process on a finite alphabet with a complex dependence structure.

A standard approach to language processing is to model a language as a distribution on sentences. The set of grammatical sentences is then given by the support of the distribution. As an alternative, computing the probability of a sentence provides an indication about its correctness with more gradations.

In such a framework, central questions are to represent and estimate the language distribution, to compute the probability of a sentence and to simulate from the estimated language. Since the support of a natural language is made of a huge number of grammatically correct sentences, an adaptive statistical approach, using a collection of parametric models of various dimensions, and going through the two steps of model selection and parameter estimation is required.

0.1.2 Remarks on the shortcomings of standard models

One approach to language processing is to use Markov chain based models. In these models, a chain of words is seen as the result of a production process that outputs words one at a time, depending only on the near past. In classical Markov chains, the n th word depends only on the $n - 1$ th, for example:

$$\mathbb{P}(\omega_n | \omega_1 \dots \omega_{n-1}) = \mathbb{P}(\omega_n | \omega_{n-1}).$$

Many variations on this pattern exists, as hidden Markov chains (where we only see a function of the words of a non observed Markov chain), or n -grams, where the m th word depends on the n previous ones. Context trees may also be used to model the dependence on the past context in a more flexible way, using a

suffix tree. More generally, the conditional probabilities may be defined with the help of any context function f with a finite range, according to the formula

$$\mathbb{P}(\omega_n | \omega_1 \dots \omega_{n-1}) = \mathbb{P}(\omega_n | f(\omega_1 \dots \omega_{n-1})).$$

Markov chains have been thoroughly studied, and as such are a good starting point to study many processes that output a chain of words – including natural languages.

However, Markov chains are missing a few important properties pertaining to the structure of natural languages. The description of a sentence likelihood as a product of forward conditional probabilities does not fit easily the long range backward and forward interactions described by linguistics. As a result, increasing the number of parameters (for instance the size of the context tree in a context tree model) to better fit the data, bumps into the curse of dimensionality, the number of parameters required for a decent fit precluding the possibility of a successful statistical estimation of those parameters.

0.1.3 Syntactic analysis

Linguists have developed several models to grasp the structure of a sentence, such as context-free grammars, minimalist grammars, head-driven phrase structure grammars, and many more. Many of these models use the notion of syntactic tree, analysing the structure of a sentence as a tree, each subtree representing a syntactic sub-unit of the sentence. For example, a (simplified) syntactic tree for the sentence “a cat chased the mouse” could be

[.VP [.DP [.D a] [.NP cat]] [.V' [.V chased] [.DP [.D the] [.NP mouse]]]]

In this example, the subtree “the mouse” forms a syntactic unit, a determiner phrase, as “chased the mouse”, which is a verbal phrase.

The notion of phrases, or syntactic constituents, is fundamental in linguistics, as the structure of sentences seems to, in a first approach at least, follow simple rules if we consider them to be constituted by nested elements.

The first observation, which is actually the main tool used by syntacticians to identify these categories, is that some strings seem to be mostly interchangeable in a given language, without loss of grammaticality. For example, the two strings “a cat” and “the dog” are mostly interchangeable. In any sentence where “a cat” appears, we can substitute it by “the dog” and keep a grammatical sentence: they are both constituents of the same type, DPs. In the same way, “chased the mouse” and “slept” are both constituents of the same type (called V’).

Syntactic relations can apparently reach over huge stretches of the sentence, while actually ignoring only a simple sub-tree. For example, relatives can be added

or subtracted and are effectively invisible to the rest of the sentence: “the dog (that Uncle Jim gave to the son of my neighbour) barked all night”.

Other considerations, such as the fact that some constituents can “move” to other specific parts of the tree, add even more weight to such structures. In our work, however, we will mainly be interested in rigid trees, without movement. Since movement is rarely optional (at least in languages with poor flecional morphology), a rigid structure is not so great a restriction. A classic model using such rigid trees are context-free grammars, in which sentences are generated by rewriting non-terminal symbols. In our cat-and-mice example, DP may for example be rewritten as D and NP, D being in turn rewritten as “the” or “a”, etc.

There is another property of natural languages which is quite dear to linguists, namely the recursivity of language, as typically exemplified by “I tell you that John said that Mary thought that (...)”. Such recursivity is very easily obtained by tree-like structures: in the case of context-free grammars, it suffices that a non-terminal symbol may be rewritten (after a number of interations) as itself.

One last property of natural languages that is closely linked to tree-structure and syntactic constituents is dependency. As we hinted before, whole sub-trees can be virtually invisible from the rest of the sentence. One notion that goes in the same direction is that of syntactic head. From outside of a constituent, only its “head” (usually one word) is visible (to check grammaticality). To go back to our example of the barking dog, the head of the constituent “the dog that Uncle Jim gave to the son of my neighbour” is “the dog” (the actual head is theoretically “the”, but let’s not delve into the linguistics depths here). As such, we could replace the rest of the syntagm by anything we could want (for example, “the dog of my aunt” and, provided that the result is still a correct nominal syntagm, put it in any sentence were the original syntagm was grammatical, and it would still be (“the dog of my aunt barked all night”).

This indicates that legitimate substitutions may be conditioned by relations between constituents and contexts depending only on appropriate surface labels.

Of course, the theoretical properties of natural language structures are still under heavy research, and the models we briefly exposed here are far from being the latest in terms of linguistic advancement. Even some notions presented here are currently questioned, as the rigid structure of head-context dependencies. However, one cannot deny that theories such as \bar{X} , mainly used for the presentation here, are very good first approximations for the analysis of natural languages. They may fail to grasp some intricate subtleties, but still do a good job at providing a rough global picture.

The Markov models presentend in section 0.1.1 on page 1 fail to grasp the recursive nature of natural languages, as well as many syntactic properties, such as far-reaching relations. As such, many efforts have been made to incorporate

syntactic structures in the models (Della Pietra et al. 1994 [DPDPG⁺94]; B. Roark 2001 [Roa01]; M. Tan et al. 2012 [TZZW12], to cite only a few). The question of learning the structure of the syntax has been also considered, for example by A. Clark [Cla14].

0.1.4 Markov substitute processes and toric grammars in a few words

Our proposal follows the trend described in the previous section, but with a change of perspective. While trying some variations on the theme of Markov models on data, we realized that modeling dependencies in natural languages through an improvement of the Markov model was a daunting task. Therefore at some point, we decided to depart from conventional modeling.

In the classical statistical approach, the tabula rasa is in some sense that all words are independent. In statistical words, this corresponds to taking as the default model the independent word model, or Markov chain of order zero. Starting from there, we may build models that describe some possible dependencies adding progressively to the set of parameters of the independent model. This works only with limited success, the structure of natural languages being too flexible to be grasped by rigid parametric models, defined prior to examining data. Using adaptive model selection may help (as in context tree weighting estimators), but does not solve the problem entirely.

So we propose to go in the opposite direction. This means taking as our tabular rasa the hypothesis that all words are dependent. Accordingly, our default model is the multinomial distribution on the language support. Using this default model, we can only output from any i.i.d. data sample the trivial distribution estimate consisting in the empirical sample distribution.

Starting from the multinomial model, we incorporate incremental conditional independence assumptions to decrease progressively the dimension of the model. These conditional independence assumptions are related to the fact that some constituents may be substituted for one another independently of the value of the context. As we remove parameters in this way, we can build distribution estimators whose support is spread further and further beyond the support of the empirical sample distribution.

We first implemented this idea by describing, with the help of a new type of grammar, that we called toric grammars, what we called sample level kernel estimates. These are language distribution estimates that produce a sentence distribution by recombining the sentences of the statistical sample, performing randomly a set of allowed constituent substitutions. This can be described as applying to the sample a Markov kernel with interactions between sentences, hence

the name sample level kernel estimate. Technically, the distribution estimate is the invariant measure of this Markov kernel, restricted to the communicating class of the statistical sample. This can be seen as some extension of classical kernel density estimates, where a smoothing Markov kernel is applied independently at each point of the statistical sample.

At first, we were not able to provide a principled statistical theory to choose this sample level kernel, (that is to choose the substitution rules), and relied on some heuristic criterion that gave satisfactory empirical results. We then realized that we could relate the sample level kernel estimates of the first chapter to a full fledged statistical model, that in fact appears as a very natural extension of Markov random field models. We called this model Markov substitute processes. In the second chapter, we describe those models in terms of conditional independence assumptions and address the two questions of model selection through multiple statistical tests, and of parameter estimation.

We show that once a set of conditional independence assumptions has been chosen, this results in a parametric model that can be described as an exponential family. The construction of the corresponding potential function is given by a theorem that extends in some way the theorem of Hammersley and Clifford [Bes74] to our model but that does not provide an efficient algorithm to compute the potential in practice (except for toy examples). We show on the other side that the conditional independence assumptions can be tested by a collection of simultaneous tests with a controlled complexity. The model complexity in these tests is handled in a data dependent way, using PAC-Bayesian theorems (a theory first proposed by D. McAllester, see [McA98, McA99]).

We then explored the combinatorics of our conditional independence assumptions. Each assumption is described by a substitute set : a list of constituents that can be substituted independently of the context, according to a substitute probability measure. A Markov substitute model is thus described by a finite family of Markov substitute sets, the set of parameters being the parameters of the substitute measures. Interaction between Markov substitute sets means that we need much less parameters than the number of actual members of each set to describe the substitute measures.

Any family of substitute sets defines a non trivial Markov substitute model, that is a non void exponential family (this is one of our results). Moreover, there is a combinatorics on Markov substitute sets : the fact that some sets are Markov substitute sets implies that other sets are also Markov substitute sets. This combinatorics is related to context-free grammars. Starting from a finite family of Markov substitute sets, we can build a context-free grammar, such that the languages that are generated choosing each non terminal symbol in turn as the start symbol forms a broader collection of Markov substitute sets. This broader collec-

tion is not the largest one. In fact, if we include one point sets in the definition of Markov substitute sets, it is possible to show that the family of maximal Markov substitute sets forms a partition of the sentence space. Anyhow we do not have an efficient algorithms to compute this partition, so that parsing algorithms for context-free grammars are at present our best offer to compute a large amount of Markov substitute sets.

This large collection can be used to build more efficient sample level kernel estimates. On the other hand, the support of our language distribution is not necessarily itself a language generated by a context-free grammar, and can be much bigger. Context-free grammars are used as a tool to manipulate a large family of substitute sets, allowing to simulate faster Monte Carlo sampling dynamics and to compute more efficiently the probability of a single sentence. The interplay between grammars and Markov substitute processes is the subject of the third chapter.

In this chapter, we also describe properly balanced sample level kernel estimates, based on random crossing over between sample sentences or based on more sophisticated and faster mixing sentence recombinations using grammatical parsing. By properly balanced, we mean kernels that are reversible with respect to the sample distribution (that a tensor product of the language distribution).

Since we are dealing with models forming exponential families, we know that the maximum likelihood estimate in each such model has good off-model properties. Namely, it will converge to the projection according to the Kullback-Leibler distance of the actual data distribution on the model, at a rate that depends on the rate of estimation of the expectation of each component of the potential function. The problem is that, as long as we are not able to compute explicitly the potential function, we are not able to compute the maximum likelihood estimator either.

To solve this problem, we show that properly balanced sample level kernel estimates have a thermodynamical limit, when applied to n replicas of the statistical sample, when n tends to infinity, that is equal to the maximum likelihood estimate of the parameters of the corresponding Markov substitute model. This means that, although we cannot provide an efficient algorithm to compute the potential function that describes a Markov substitute process model as an exponential family, we still have a feasible Monte-Carlo algorithm to compute the maximum likelihood estimator in this model. This shows that the model can be used not only on data that would exactly match the required assumptions, but that the model can be used on any data distribution to provide an approximation of this distribution, in the same way as high order Markov chains can be used to approximate non Markovian stationary ergodic processes.

0.2 Overview of the main results

We will now present a quick overview of the notions encountered and the main results of this work. In what follows, we will consider a text sample $S_{1:n}$ of n sentences on a dictionary D . It will sometimes be convenient, since we do not attempt to model interactions between sentences in this work, to see this sample (and any text), as an empirical distribution on $D^+ = \bigcup_{j=1}^{\infty} D^j$, so the order of sentences does not count.

0.2.1 Toric grammars

We began with what we called a communication model, in the form of a kernel modifying sentences and texts according to a certain type of interaction. This point of view was inspired by the linguistic notion of syntactic categories, seen as chains of words that can be substituted regardless of their contexts.

The main tool we developed to build these sample level kernels was toric grammars, a notion closely related to context-free grammars. Indeed, the set \mathfrak{G} of toric grammars is defined in definition 1.2.1 on page 22 as the set of positive measures on context-free rewriting rules \mathcal{E} . Using this definition, we can see the set of texts \mathfrak{T} as a subset of the toric grammar set \mathfrak{G} , whose support is included in the rewriting rules for the start symbol. For ease of notation, we write the non-terminal symbols as $]_i$, $i \in \mathbb{N}$, and for the left-hand side of production rules, we replace the usual notation $]_i \rightarrow$ by $[_i$. The start symbol $]_0$ is supposed to never occur in the right-hand side of the rules (as in some normal forms of context-free grammars), so only the symbol $[_0$ appears.

The main difference between toric grammars and context-free grammars could be roughly summarized by the fact that the context-free rewriting rules of toric grammars are more viewed as context-free substitution rules (replace any string that can be generated by a given non-terminal symbol by any other, regardless of the context), than rewriting rules.

These grammars give us a method to modify a text (seen as an empirical measure on sentences, and thus also a toric grammar) by substituting substrings by others. In order to do so, we have to make a distinction between what we chose to call a reference grammar \mathcal{R} , used to describe the model of the language considered, and “regular” toric grammars \mathcal{G} , obtained from parsing a certain text using the rules of a reference grammar in a possibly random manner.

The reference grammar can thus be seen as a collection of context-free rewriting rules, *without a start symbol*. For these, (in chapter 1, at least), the weights are not important. From this reference grammar, we will define in section 1.3.2 on page 28 a random split-and-merge process G_t , that will swap the occurrences of constituents of same type in the text.

This process is implemented in two steps, each step being a time homogeneous stopped Markov chain. The first step, the splitting process S_t , is the decomposition of the initial text into the rules of \mathcal{R} used to produce the text, giving a regular toric grammar. It is stopped at stopping time τ , when no more authorized splits are available. The second step, the production process P_t , is the converse operation of using the same rules back to generate a new text, effectively mixing the constituents of the initial text into a new one. It is stopped at stopping time σ , when there is no possible move left.

As such, the split-and-merge process G_t is a Markov chain on grammars, a grammar on even time being a text produced by the grammar of the previous time using a production process, and a grammar on odd time being a parse of the previous text, obtained with a splitting process.

We show in proposition 1.3.4 on page 30 that both the splitting and the production processes take finite time, with a deterministic bound depending on the size of the initial text. This justifies the definition of the split-and-merge process, whose transitions are then computable in finite time.

Proposition

Let $(S_t)_{t \leq \tau}$, $(P_t)_{t \leq \sigma}$ and (G_t) be a splitting process, a production process and the corresponding split and merge process, starting from $G_0 = \mathcal{T} \in \mathfrak{T}$. For any $\mathcal{G} \in \mathfrak{G}$, any $\mathcal{T}' \in \mathfrak{T}$, such that $\sum_{t \in \mathbb{N}} \mathbb{P}(G_{2t+1} = \mathcal{G}) > 0$ and $\sum_{t \in \mathbb{N}} \mathbb{P}(G_{2t} = \mathcal{T}') > 0$,

$$\begin{aligned} \mathbb{P}\left(\tau \leq 2\left[\mathcal{T}(DS^*) - \mathcal{T}([{}_0 S^*])\right] \mid S_0 = \mathcal{T}'\right) &= 1, \\ \mathbb{P}\left(\sigma \leq 2\left[\mathcal{T}(DS^*) - \mathcal{T}([{}_0 S^*])\right] \mid P_0 = \mathcal{G}\right) &= 1. \end{aligned}$$

In other words, the length of all the splitting and production processes involved in the split and merge process have a uniform bound, given by twice the difference between the number of words and the number of sentences in the original text.

Moreover, we prove in proposition 1.4.2 on page 36 that this process is weakly reversible.

Proposition

Given a parsing process S_t based on a reference grammar $\mathcal{R} \in \mathfrak{G}$ and a production process P_t , the corresponding split and merge process G_t is weakly reversible, in the sense that for any $\mathcal{T} \in \mathfrak{T}$, any $\mathcal{G} \in \bigcup_{t \in \mathbb{N}} \text{supp}(\mathbb{P}_{G_{2t+1}})$,

$$\mathbb{P}(G_1 = \mathcal{G} \mid G_0 = \mathcal{T}) > 0 \iff \mathbb{P}(G_2 = \mathcal{T} \mid G_1 = \mathcal{G}) > 0.$$

Consequently, for any $\mathcal{T}, \mathcal{T}' \in \mathfrak{T}$ and any $\mathcal{G}, \mathcal{G}' \in \bigcup_{t \in \mathbb{N}} \text{supp}(\mathbb{P}_{G_{2t+1}})$,

$$\begin{aligned} \mathbb{P}(G_2 = \mathcal{T}' \mid G_0 = \mathcal{T}) > 0 &\iff \mathbb{P}(G_2 = \mathcal{T}' \mid G_0 = \mathcal{T}') > 0, \\ \mathbb{P}(G_3 = \mathcal{G}' \mid G_1 = \mathcal{G}) > 0 &\iff \mathbb{P}(G_3 = \mathcal{G}' \mid G_1 = \mathcal{G}') > 0. \end{aligned}$$

In other words, the two processes G_{2t} and G_{2t+1} are weakly reversible time homogeneous Markov chains.

As we will also prove that the set of reachable states from any starting point is finite, this proposition shows that the split and merge process defines two recurrent Markov chains, one on texts and one on grammars. These two chains partition their respective state spaces into positive recurrent communicating classes.

This means that the split and merge process can be seen as a communication process, where each speaker, upon hearing a given text, learns a grammar that could have produced it, and utters back a new text to a new speaker, and so forth. What this result shows is that if the reference grammar (which can be seen as an internal language model) is common to all speakers, the asymptotic language, that is, the distribution of $\lim_{\frac{1}{T} \sum_{t \leq T} G_{2t+1}}$, is uniquely determined by the recurrent class of the initial text.

The main difference of this approach from usual context-free grammars is that the rules are weighted by a counting measure, counting how many times each rule is to be applied. This means that, for example, the size, vocabulary, etc. remains constant between texts in the same communicating class. A second difference is that the process G_t , going from texts to texts, is iterated. As such, the texts produced are produced by the context-free rules of the reference grammar, complemented by a set of rewriting rules for the start symbol, that may not be the same for all texts in the same communicating class. We will see in section 1.6.2 on page 39 that in some cases this gives a wider generated language than what can be obtained using a context-free grammar.

0.2.2 Markov substitute sets and processes

The study of the first chapter leaves the question of the statistical estimation of the reference grammar essentially unsolved, although encouraging preliminary experimental results could be obtained using some heuristic rules of thumb to learn the reference grammar.

This question led us to relate our toric grammars to more classical statistical models, described in chapters 2 and 3.

In chapter 2, we begin by defining in definition 2.1.1 on page 64 the notion of Markov substitute sets, as a possible formalization of syntactic constituents. These sets are sets of strings that can be substituted in any sentence, regardless of their contexts. In more formal terms, a Markov substitute set B for a distribution \mathbb{P}_S has a substitute measure q_B such that, for any context (x, y) , and any string $b \in B$,

$$\mathbb{P}_S(xby) = \mathbb{P}_S(xBy)q_B(b).$$

While this notion is rooted in linguistic considerations, it is in no way restricted to linguistics, and could have applications in other domains. It can indeed be seen as

a generalization of the (one dimensional) Markov field conditional independence property, that corresponds to the case when B is of the form $w_1 D w_2$, where w_1 and w_2 are two boundary words.

To any collection of Markov substitute sets \mathcal{B} corresponds an equivalence relation that relates two strings $s \sim_{\mathcal{B}} s'$ if and only if s and s' can be reached from one another by a sequence of substitutions authorized by \mathcal{B} . The components of D^+ for this equivalence relation are all Markov substitute sets.

Conversely, we can define a \mathcal{B} -Markov substitute process (for a given finite collection of finite sets \mathcal{B}) as any process S such that each set of \mathcal{B} is a Markov substitute set (technically, we have to require also that each element of each set of \mathcal{B} has a positive probability to appear in a sentence). We will prove in section 2.3 on page 74 that Markov substitute processes actually exist for any collection of Markov substitute sets, and that the set of \mathcal{B} -Markov substitute processes is an exponential family (technically a union of exponential families, depending on the support of the process).

Proposition

Given any finite set \mathcal{B} of finite subsets of D^+ , there is a finite set of pairs, \mathcal{P} , that we can choose such that each one is included in a member of \mathcal{B} , such that the sets of \mathcal{B} -Markov and \mathcal{P} -Markov substitute processes are the same, and such that \mathcal{P} is minimal for the inclusion relation (removing a pair from \mathcal{P} would break the above property).

For any \mathcal{B} -Markov substitute process S , we can decompose its support as a union of components $\mathcal{C} \subset D^+ / \sim_{\mathcal{B}}$, such that $\text{supp}(\mathbb{P}_S) = \bigcup \mathcal{C}$ and such that

$$\mathcal{B} = \bigcup_{C \in \mathcal{C}} \mathcal{B}_C.$$

For any such subset, there is a free, non empty subset of pairs $\mathcal{F} \subset \mathcal{P}$, a matrix $(e_{j,i}, j \in \mathcal{P} \setminus \mathcal{F}, i \in \mathcal{F})$, an index set $I(\mathcal{C})$ containing \mathcal{F} and an energy function $(U_i(s), i \in I(\mathcal{C}), s \in C)$, such that the set of \mathcal{B} -Markov processes whose support is $\bigcup \mathcal{C}$ is the linear exponential family

$$\mathfrak{M}_{\mathcal{C}}(\mathcal{B}) = \left\{ \left(p(s) = \frac{\exp\left(-\sum_{i \in I(\mathcal{C})} \beta_i U_i(s)\right)}{\sum_{s' \in \bigcup \mathcal{C}} \exp\left(-\sum_{i \in I(\mathcal{C})} \beta_i U_i(s')\right)}, s \in \bigcup \mathcal{C} \right), \beta \in \mathfrak{B} \subset \mathbb{R}^{I(\mathcal{C})} \right\},$$

where

$$\mathfrak{B} = \left\{ \beta \in \mathbb{R}^{I(\mathcal{C})}, \sum_{s \in \bigcup \mathcal{C}} \exp\left(-\sum_{i \in I(\mathcal{C})} \beta_i U_i(s)\right) < \infty \right\},$$

and such that for any member p of this family, the substitute measure of any pair $i = \{y_{i,0}, y'_{i,1}\} \in \mathcal{F}$ (taken in a suitable order compatible with the definition of U), is given by

$$q_i(y_{i,\sigma}) = \frac{\exp(\sigma\beta_i)}{1 + \exp(\sigma\beta_i)}, \quad \sigma \in \{0, 1\},$$

whereas for $j \in \mathcal{P} \setminus \mathcal{F}$,

$$q_j(y_{j,\sigma}) = \frac{\exp(\sigma\beta_j)}{1 + \exp(\sigma\beta_j)}, \quad \sigma \in \{0, 1\},$$

where

$$\beta_j = \sum_{i \in \mathcal{F}(C)} e_{j,i} \beta_i.$$

We can see here that the problem of analysing any language using this framework can be decomposed into two stages:

- first a model selection stage, where a good collection \mathcal{B} of Markov substitute sets is to be selected,
- and then a parameter estimation stage, in which the substitute measures q_B , for each $B \in \mathcal{B}$ are to be estimated. (We could in principle estimate only the parameters of the exponential family, but the above theorem does not provide an efficient algorithm to compute the energy function U explicitly, so that this will not be feasible, except for some toy models.)

Since Markov substitute sets are a possible formalization of syntactic constituents, we will see in section 3.1 on page 119 that it is possible to relate them to context-free rules. Conversely, a set of context-free rules defines a collection of Markov substitute sets, as the collection of the languages generated by each non-terminal symbol. This duality gives us quite efficient ways to analyse any sentence in terms of Markov substitute sets, using parsing algorithms, as described in section 3.B on page 159.

0.2.3 Model selection: finding a collection of Markov substitute sets

The question is to find Markov substitute sets, or, equivalently, context-free rules.

A good starting point is that two strings form a Markov substitute set if their relative frequencies do not depend on their context. This observation can be extended to two sets of strings, for which we know that each already forms a

Markov substitute set. This can be used to build several test functions to test whether the union of two Markov substitute sets $B_1 \cup B_2$ is still a Markov substitute set. We explore two main categories of possible test functions.

The first category involves the simulation of a parse (X, Y) of a sentence S , so that $S = \alpha(X, Y)$ (meaning that S is the string where Y appears in context X), and $Y \in B_1 \cup B_2 = B$. These tests are based on the empirical mean of test functions of the form

$$F_\theta(X, Y, p) = \mathbf{1}(X \in \theta) \left[\mathbf{1}(Y \in B_1) - p \mathbf{1}(Y \in B_1 \cup B_2) \right],$$

where $\theta \in \Theta$ is a context type.

The second category is built by computing explicitly conditional expectations with respect to the raw (non parsed) sentence sample. It is of the type

$$F_{B_1, B_2, \theta}(s, p) = \mathbb{E} \left\{ \left[\mathbf{1}(Y_{B_1 \cup B_2} \in B_1) - p \mathbf{1}(Y_{B_1 \cup B_2} \in B_1 \cup B_2) \right] \mathbf{1}(X_{B_1 \cup B_2} \in \theta) \mid S = s \right\}.$$

The goal is to find test functions $F_{B_1, B_2, \theta}(S, p)$ such that, regardless of the context type θ ,

$$\mathbb{E} \left(F_{B_1, B_2, \theta}(S, p) \right) = 0$$

when $p = q_{B_1 \cup B_2}(B_1)$, in the case when $B_1 \cup B_2$ is indeed a Markov substitute set. In other words, if we define

$$p_+ = \sup \left\{ \mathbb{P}(Y \in B_1 \mid Y \in B, X = x); x \in (D^*)^2, \mathbb{P}(Y \in B, X = x) > 0 \right\},$$

$$p_- = \inf \left\{ \mathbb{P}(Y \in B_1 \mid Y \in B, X = x); x \in (D^*)^2, \mathbb{P}(Y \in B, X = x) > 0 \right\},$$

the test function allows us to test whether $p_+ = p_-$. A third interpretation is that the test functions define $p(B_1, B_2, \theta)$ so that

$$\mathbb{E} \left[F(S, p(B_1, B_2, \theta)) \right] = 0,$$

and we want to test whether $p(B_1, B_2, \theta)$ is independent of θ .

For both of these types of test function, we can define similar tests, for which we can prove accurate bounds for the simultaneous probability of both false rejection and false acceptance.

For the first type of test function, for example, we can prove the following result (propositions 2.6.1 and 2.6.3 on page 89 and on page 92)

Proposition

Let $\mu \in \mathcal{M}_+^1$ be a probability measure depending only on $(X_i, \mathbf{1}(Y_i \in B), i \in \llbracket 1, n \rrbracket)$. Let us define

$$F_{\theta, i}(p) = F_\theta(X_i, Y_i, p).$$

Let Λ be a finite subset of $]0, 1[$. With probability at least $1 - 2\varepsilon$,

$$B_-(p_+) \stackrel{\text{def}}{=} \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \int \sum_{i=1}^n \log(1 + \lambda F_{\theta, i}(p_+)) d\rho(\theta) - \mathcal{K}(\rho, \mu) - \log(|\Lambda|/\varepsilon) \leq 0,$$

$$B_+(p_-) \stackrel{\text{def}}{=} \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \int \sum_{i=1}^n \log(1 - \lambda F_{\theta, i}(p_-)) d\rho(\theta) - \mathcal{K}(\rho, \mu) - \log(|\Lambda|/\varepsilon) \leq 0,$$

where

$$\mathcal{K}(\rho, \mu) = \begin{cases} \int \log\left(\frac{d\rho}{d\mu}\right) d\rho, & \text{when } \rho \ll \mu, \\ +\infty, & \text{otherwise.} \end{cases}$$

Therefore, if we reject the hypothesis that B is a Markov substitute set when

$$\inf_{p \in [0, 1]} \max\{B_-(p), B_+(p)\} > 0,$$

the probability of false rejection is at most 2ε .

Let us assume in addition that the number of context types that can be encountered in the language is bounded by m , and that μ is the uniform probability measure on the finite set of encountered context types $\hat{\Theta}$. Let us put

$$\delta = \frac{\log(mn) + 2\log(\varepsilon^{-1}) + \log(|\Lambda|)}{n},$$

and

$$\chi = \sup_{x \in [n^{-1/2}, n^{1/2}]} \inf_{\lambda \in \Lambda} \cosh\left[\log\left(\frac{\lambda}{(1-\lambda)x}\right)\right].$$

Let us assume that there are θ_+ and $\theta_- \in \Theta$ such that $\bar{p}_+ = p(\theta_+)$, $\bar{p}_- = p(\theta_-)$, $q_+ = \mathbb{E}[\mathbf{1}(Y \in B)\mathbf{1}(X \in \theta_+)]$ and $q_- = \mathbb{E}[\mathbf{1}(Y \in B)\mathbf{1}(X \in \theta_-)]$ satisfy

$$\min\{q_-, q_+\} \geq 8\chi^2\delta,$$

$$\bar{p}_+ - \bar{p}_- \geq 2\chi \sqrt{\frac{\bar{p}_+(1-\bar{p}_+)\delta}{q_+} \left(1 + \frac{4\chi^2\delta}{q_+}\right) + \frac{(2+\sqrt{2})\delta}{q_+}}$$

$$+ 2\chi \sqrt{\frac{\bar{p}_-(1-\bar{p}_-)\delta}{q_-} \left(1 + \frac{4\chi^2\delta}{q_-}\right) + \frac{(2+\sqrt{2})\delta}{q_-}}.$$

Then with probability at least $1 - 2\varepsilon$

$$\inf_{p \in [0, 1]} \max\{B_-(p), B_+(p)\} > 0.$$

Therefore, the test has in this case a probability of false acceptance at most equal to 2ε .

A similar result about the test functions that do not simulate the parse is proved in propositions 2.7.3 and 2.7.7 on page 101 and on page 112. The test proposed in this case also comprises a uniform bound on all tests on pairs of Markov substitute sets chosen in a finite collection \mathcal{B} . We propose a novel PAC-Bayes complexity measure, based on the Kullback divergence between a posterior distribution and another observable reference distribution, that plays the role of a prior, but involves only the contexts appearing in the statistical sample. This observable complexity factor is bounding the more classical divergence of the posterior from a prior, according to eq. (2.7.4) on page 98.

We can remark at this point that these tests give a way to test if the union of two Markov substitute sets is still one, without having to actually compute the whole extent of each set. This is interesting, because, if we take the rewriting rule view, the new Markov substitute set can then be described as a simple rewriting rule $B_1 \cup B_2 \longrightarrow B_1|B_2$. This mechanism gives actually a method to build recursively a complete toric grammar, for which the distribution (provided no test failed) is a Markov substitute process.

0.2.4 Parameter estimation

The second chapter proposes tests to select a particular model described by a collection of Markov substitute sets. The next question is, given this collection, how to estimate the parameters q_B for each Markov substitute set in the collection.

While the prospect of estimating the parameter $q_B(x)$ for each member x of B is quite daunting, we will see that we can use the grammar structure to greatly simplify this estimation. In order to understand this result, we first need to see that collection of Markov substitute sets and reference grammars are closely related. Indeed, we can define for any string e of symbols the language set B_e as the set of strings that can be generated from e . We will see in section 3.1 on page 119 that this set B_e is a Markov substitute set, as soon as each $B_i \stackrel{\text{def}}{=} B_{]_i}$ are Markov substitute sets themselves. We call such a grammar a Markov grammar. Conversely, if we build a collection of Markov substitute sets using the tests of chapter 2, we can define a corresponding Markov grammar, as described in section 3.4 on page 144.

Another useful notion is that of parse trees $t \in \mathcal{T}$, defined in section 3.1.2 on page 122. They are parenthetized strings in $\check{D}^* = (D \cup \{(i,)_i, i \in \mathbb{N}\})^*$, such that the content y of each pair parentheses $(i y)_i$ is a member of B_i , that is, a string in the language generated by the grammar with $]_i$ as start symbol. The structure of a given tree gives a very simple decomposition of the string in terms of Markov substitute sets. The surface structure $\varsigma(t)$ of t , which is simply the string t where all pairs of parentheses (and their contents) $(i y)_i$ are replaced by the corresponding non-terminal symbol $]_i$, gives one language that can generate $\varphi(t)$ (t without

parentheses), that is, $\varphi(t) \in B_{\varsigma(t)}$. Using these notions, we will show the following result, lemma 3.2.1 on page 125, giving a formula for any substitute measure:

Lemma

For any $s \in D^+$, and any $t \in \mathcal{T}$ such that $s = \varphi(t)$,

$$q_{B_{\varsigma(t)}}(s) = A_{\varsigma(t)} \prod_{[j \in \mathcal{A}]} [A_e q_{B_j}(B_e)]^{\chi(t,j,e)},$$

where

$$\chi(t, j, e) = \sum_{x \in (\check{D}^*)^2} \sum_{y \in \mathcal{T}} \mathbb{1}[t = \alpha(x, (jy)_j)] \mathbb{1}(\varsigma(y) = e).$$

The advantage of this result is that we need only to compute the parameters A_e and $q_{B_j}(B_e)$ for expressions e that are actually in the grammar, and not for all members of the Markov substitute sets.

The sections 3.2.2 and 3.2.3 on page 127 and on page 129 use this result to propose methods to estimate there parameters and simulate q_B .

0.2.5 Invariant dynamics

The last part of chapter 3 is about three types of invariant dynamics, that we can define with the previous tools.

The first one is a dynamics on sentences, and requires the knowledge of the substitute measures q_B . This dynamics is defined in section 3.2.4 on page 131, and simply identifies members of a Markov substitute sets B in the sentence, and replaces it by another member of B according to q_B .

The two other dynamics work on whole texts, but do not require actual knowledge of the substitute measures q_B .

The second method, described in section 3.5.1 on page 151, is actually a new take on the split and merge process defined in chapter 1, using a Metropolis algorithm to make it reversible.

The third method, described in section 3.2.5 on page 134, is a sort of crossing-over dynamics, that will simply swap elements of the same set found in the text.

These three methods are all made reversible using a slightly more general (to our knowledge) method for the Metropolis algorithm, that allows for an arbitrary number of intermediate steps between the argument and the result. This method is described in section 3.D on page 165.

This last dynamics also gives an interesting view on Markov substitute processes. We mentioned that these could in fact be described as an exponential family. This means that the maximum likelihood estimate for any sample will converge to the projection according to the Kullback Leibler distance on the model.

The problem is that since the energy function is hard to compute, the projection is also difficult to find.

However, we will be able to prove in section 3.3 on page 135 that the crossing-over dynamics from any sample does indeed converge to the projection (with some additional hypotheses, of course).

If we write K_m a crossing-over dynamics on texts of size nm , we can define, for any text $s_{1:n}$, S^m a random uniform shuffle of m copies of $s_{1:n}$, and the distribution

$$p_m = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{u=1}^t \mathbb{P}_{S^m} K_m^u.$$

We can then define the random measure

$$N_m = \sum_{i=1}^{nm} \delta_{\bar{S}_{m,i}} \in \mathcal{M}_+(\mathcal{D}),$$

where $(\bar{S}_{m,i}, 1 \leq i \leq nm)$ is a random variable distributed according to p_m . Simply put, N_m is the empirical distribution of the crossing-over of m copies of the initial text, randomly shuffled.

We then will prove in proposition 3.3.3 on page 143 that

Proposition

Under some hypotheses described at the beginning of section 3.3 on page 135, the maximum likelihood

$$\sup_{\beta \in \mathbb{R}^I} \sum_{j=1}^n \log [p_\beta(s_j)]$$

is reached at some $\beta_ \in \mathbb{R}^I$. For any $\varepsilon > 0$, there are $\eta > 0$ and $M > 0$ such that for any $m \geq M$,*

$$\mathbb{P} [d(N_m, p_{\beta_*}) \geq \varepsilon] \leq \exp(-nm\eta),$$

where $d(\mu, \nu) = \sqrt{\sum_{s \in \mathcal{D}} (\mu(s) - \nu(s))^2}$. Moreover,

$$\lim_{m \rightarrow \infty} \mathbb{E}(N_m / (nm)) = p_{\beta_*}.$$

This means that we can actually project any sample on the Markov substitute model using a crossing-over dynamics, and thus that Markov substitute processes can be seen as thermodynamical limits of crossing-over dynamics.

Chapter 1

Toric grammars and communication models

1.1 Introduction to a new communication model

In the well known kernel approach to density estimation on a measurable space \mathcal{X} , the probability distribution \mathbb{P} of a random variable $X \in \mathcal{X}$ is estimated from a statistical sample (X_1, \dots, X_n) of n independent copies of X as $\frac{1}{n} \sum_{i=1}^n k(X_i, dx)$, where k is a suitable Markov kernel (also known as a conditional probability kernel). This kernel estimate can be seen as a modification of the empirical measure $\bar{\mathbb{P}} = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$.

Finding sensible kernel estimates or sensible parametric models in the context of natural language processing is a challenge. Therefore, we propose here another route, that we will describe as an alternative way of producing a modification of the empirical measure. The idea is to recombine repeatedly a set of sentences. Let us describe for this a general framework, concerned with an arbitrary countable state space \mathcal{X} .

Let

$$\bar{\mathcal{P}}_n = \left\{ \frac{1}{n} \sum_{i=1}^n \delta_{x_i}, x_i \in \mathcal{X} \right\}$$

be the set of empirical measures of all possible samples of size n . Let us consider a parametric family $\{q_\theta, \theta \in \Theta\}$ of Markov kernels on $\bar{\mathcal{P}}_n$. Let us assume for simplicity that for any $P \in \bar{\mathcal{P}}_n$, the reachable set $\{Q \in \bar{\mathcal{P}}_n, \sum_{t \in \mathbb{N}} q_\theta^t(P, Q) > 0\}$ is finite, where q_θ^t is q_θ composed t times with itself, $t \in \mathbb{N}$, so that for instance $q_\theta^2(P, Q) = \sum_{P' \in \bar{\mathcal{P}}_n} q_\theta(P, P')q_\theta(P', Q)$. In this case we can define the Markov kernel

$$\hat{q}_\theta(P, Q) = \lim_{k \rightarrow \infty} \frac{1}{k} \sum_{t=1}^k q_\theta^t(P, Q).$$

It is such that for any $P \in \overline{\mathcal{P}}_n$, $\hat{q}_\theta(P, \cdot)$ is an invariant measure of q_θ . More generally $q_\theta \hat{q}_\theta = \hat{q}_\theta q_\theta = \hat{q}_\theta$. The distribution $\hat{q}_\theta(P, \cdot) \in \mathcal{M}_+^1(\overline{\mathcal{P}}_n)$ induces a marginal distribution $\hat{Q}_{\theta, P}$ on \mathcal{X} through the formula

$$\hat{Q}_{\theta, P} = \sum_{Q \in \overline{\mathcal{P}}_n} \hat{q}_\theta(P, Q) Q. \quad (1.1.1)$$

We will here be concerned with estimators of the form $\hat{\mathbb{P}} = \hat{Q}_{\theta, \overline{\mathbb{P}}}$, if θ is fixed in advance, or of the form $\hat{Q}_{\hat{\theta}, \overline{\mathbb{P}}}$, if $\hat{\theta}$ is an estimator of the parameter θ depending also on $\overline{\mathbb{P}}$.

Another interpretation of our framework is to consider q_θ as a communication model. One speaker hears a set of sentences described by its empirical distribution $P \in \overline{\mathcal{P}}_n$ (which means that he will not make use of the special order in which he has heard them). He uses those sentences to learn the corresponding language. Then he teaches another speaker what he has learnt by outputting another random set of sentences, distributed according to $q_\theta(P, \cdot)$. The language model (as opposed to the communication model q_θ), is $\hat{Q}_{\theta, P}$, the average sentence distribution along an infinite chain of communicating speakers.

If we start from a recurrent state P , and we assume that θ is known, we obtain a communication model where the target sentence distribution $\hat{Q}_{\theta, P}$ can be learnt without error from the set of sentences output by any involved speaker. Indeed $\hat{Q}_{\theta, P} = \hat{Q}_{\theta, Q}$ for any Q in the communicating class of P , which in this situation is also the reachable set from P .

This error free estimation behaviour is desirable for a communication model. It tells us that the language can be transmitted from speaker to speaker without distortion, a desirable feature in the case of a large number of speakers. The model may also account for weak stimulus learning, the fact that human beings learn language through a limited number of examples compared with the variety of new sentences they are able to formulate. Indeed, whereas the size of the support of $P \in \overline{\mathcal{P}}_n$ (the number of sentences heard by one speaker) is constant and equal to n , the support of the language model $\hat{Q}_{\theta, P}$ may be much larger. We will actually give a toy example where the number of sentences in the language is exponential with n .

Remark on the other hand that if the parameter θ is not known, but rather estimated at each step from the heard sample, the error free property is obviously not true anymore. This view of communication by iterated learning, where an internal grammar θ is learned at each generation and used to produce a learning sample for the next generation, is not new, and has been studied, for example, by Griffiths and Kalish [GK07] to form models of language evolution.

In the language transmission interpretation, we may evaluate the interest of the model by studying whether it can model a large family of sentence distributions.

This richness will depend on the number of recurrent communicating classes of the communication Markov model q_θ , since any invariant distribution $\hat{q}_\theta(P, \cdot)$ is a convex combination of the unique invariant measures supported by each recurrent communicating class. The situation is even simpler in the case when all $P \in \overline{\mathcal{P}}_n$ are recurrent states (a fact we will be able to prove in our particular model). In this case $\hat{q}_\theta(P, \cdot)$ is the unique invariant measure supported by the recurrent communicating class to which P belongs.

The parameter θ of our model will be a new kind of grammar, closely related to context-free grammars, in that it will use rewriting rules, but used to generate sentences in a different way. These toric grammars will be defined in section 1.2 below as a weighted collection of rewriting rules, not unlike stochastic context-free grammars.

We will relate the building of a communication kernel presented in this chapter with the construction of more classical statistical models in the following chapters of this work.

1.2 First definitions

1.2.1 Toric grammars

We consider some finite dictionary of words D , and a random sentence S , that is a random sequence of words of random finite length. We will use the notation

$$D^+ = \bigcup_{j=1}^{\infty} D^j$$

for the set of sequences of words of finite positive length. As explained in the introduction, S , or equivalently its probability distribution $\mathbb{P}_S \in \mathcal{M}_+^1(D^+)$, will be our mathematical representation of a language.

From this language, we will observe a statistical sample of sentences S_1, \dots, S_n , made of n independent copies of S .

In this chapter, we will build estimators by applying a Markov kernel to the empirical sample distribution. To describe this construction, let us introduce the state space of counting measures of weight n as

$$\overline{\mathcal{P}}_n = \left\{ \sum_{i=1}^n \delta_{s_i}, s_i \in D^+ \right\}.$$

(Let us remark that in this definition, where δ_{s_i} is the Dirac mass at s_i , the same sentence can be repeated more than once.)

We will call $\overline{\mathcal{P}}_n$ the set of texts of length n . Let us notice that for us, texts are unordered sets of sentences (with possible repetitions). The question of generating meaningful ordered sequences of sentences is also of interest, but will not be addressed in this study.

In order to define a language estimator, we will first define a Markov kernel q (that is a conditional probability kernel) on the state space $\overline{\mathcal{P}}_n$. We will construct this kernel in such a way that the reachable sets $\{Q \in \overline{\mathcal{P}}_n, \sum_{t=0}^{\infty} q^t(P, Q) > 0\}$ are finite (for any starting point P), so that we can define a limit Markov kernel

$$\hat{q}(P, Q) = \lim_{k \rightarrow \infty} \frac{1}{k} \sum_{t=1}^k q^t(P, Q),$$

and take as our language estimator

$$\hat{\mathbb{P}} = \frac{1}{n} \sum_{Q \in \overline{\mathcal{P}}_n} \hat{q}\left(\sum_{i=1}^n \delta_{S_i}, Q\right) Q.$$

In order to define the communication kernel q , we will describe random transformations on texts, related to the notion of context-free grammars. Let us start with an informal presentation. The communication kernel will perform random recombinations of sentences.

Our point of view is to see a context-free grammar as the result of some fragmentation process applied to a set of sentences. Let us explain this on a simple example. Consider the sentence

This is my friend Peter.

Imagine we would like to represent this sentence as the result of pasting the expression *my friend* in its context, because we think language is built by cutting and pasting expressions drawn from some large sets of memorized expressions. We can do this by introducing the simple context-free grammar

$$\begin{aligned} \boxed{0} &\rightarrow \textit{This is } \boxed{1} \textit{ Peter.} \\ \boxed{1} &\rightarrow \textit{my friend} \end{aligned}$$

where we have used numbered framed boxes for non terminal symbols, the start symbol being $\boxed{0}$. The two rules mean that we can rewrite the start symbol $\boxed{0}$ to obtain the right-hand side of the first rule, and that we can then rewrite the non terminal symbol $\boxed{1}$ as the right-hand side of the second rule.

Since we want to see the rules of the grammar as the result of some splitting operation, we are going to use more symmetric notations. Instead of considering that we have described our original sentence with the help of two rules and two non

terminal symbols $\boxed{0}$ and $\boxed{1}$, we may as well consider that we have split our original sentence into two new sentences using *three* non terminal symbols, namely $\boxed{0} \rightarrow$, $\boxed{1}$ and $\boxed{1} \rightarrow$. To emphasize this interpretation, we can adopt more symmetric notations and write these three non terminal symbols as $[_0,]_1$ and $[_1$. With these new notations, the representation of our original sentence is now

$$\begin{array}{l} [_0 \textit{ This is }]_1 \textit{ Peter} . \\ [_1 \textit{ my friend} \end{array}$$

In this new representation, the rewriting rules can be replaced by merge operations of the type

$$a]_i c + [_i b \mapsto abc$$

We can make this merge operations even more symmetric, if we consider that each expression can be represented by any of its circular permutations. Indeed, each expression contains exactly one non terminal symbol of the form $[_i$, and therefore is uniquely defined by any of its circular permutations (since, due to this feature, we can define the permutation in which the opening bracket $[_i$ comes first as the canonical form, and recover it from any other circular permutation). Using this convention, we can write $a]_i c$ as $ca]_i$ (note that a starts with an opening bracket and contains no other opening bracket, whereas c does not contain any opening bracket) and describe the merge operation as

$$ca]_i + [_i b \mapsto cab,$$

or, renaming ca as a , simply as

$$a]_i + [_i b \mapsto ab.$$

Let us formalize what we have explained so far. Let D be some dictionary of words (which can be for the sake of this mathematical description any finite set, representing the words of the natural language to be modeled). Let us form the symbol set $S = D \cup \{[_i,]_i, i \in \mathbb{N}\}$. Let us define the set of circular permutations of a sequence of symbols as

$$\mathfrak{S}(w_0, \dots, w_{\ell-1}) = \{(w_{(i+j \bmod \ell)}, i = 0, \dots, \ell - 1), j = 0, \dots, \ell - 1\},$$

so that for instance $\mathfrak{S}(w_0, w_1, w_2) = \{w_0 w_1 w_2, w_1 w_2 w_0, w_2 w_0 w_1\}$. Let us define its support (the set of symbols included in the sequence) as

$$\text{supp}(w_0, \dots, w_{\ell-1}) = \{w_0, \dots, w_{\ell-1}\}.$$

Let $A^+ = \bigcup_{n=1}^{+\infty} A^n$, $[_+ = \{[_i, i \in \mathbb{N} \setminus \{0\}\}$, and consider the set of expressions

$$\mathcal{E} = \left\{ e \in \mathfrak{S}([_i a), i \in \mathbb{N}, a \in (D \cup]_+)^+ \setminus]_+ \right\}.$$

In plain words, an expression is a circular permutation of a finite sequence of symbols starting with an opening bracket, containing no other opening bracket and not reduced to an opening bracket followed by a closing bracket.

This definition mirrors the fact that a given rule of a context-free grammar has exactly one $\boxed{i} \rightarrow$ (the left side), and the right side of the rule cannot be just a non terminal symbol \boxed{j} . Indeed, if we had allowed $\boxed{i} \rightarrow \boxed{j}$, or with our notations $[_i]_j$, we could as well have replaced j by i everywhere.

We can now define the notion of toric grammars, which will be the central tool in all of this chapter, both to describe the language structure, and to define the transformations on texts.

Definition 1.2.1

The set of toric grammars is the set \mathfrak{G} of positive measures \mathcal{G} on \mathcal{E} with finite support such that for any circular permutation $e' \in \mathfrak{S}(e)$ of any expression $e \in \mathcal{E}$, $\mathcal{G}(e') = \mathcal{G}(e)$.

In other words, a toric grammar \mathcal{G} is a positive measure with finite support on the set of expressions \mathcal{E} satisfying

$$\mathcal{G}(e) = |\mathfrak{S}(e)|^{-1} \mathcal{G}(\mathfrak{S}(e)).$$

Let us remark that, in our definition of toric grammars, on top of choosing some special notations for context-free grammars, we also introduced positive weights, so that it is more the support of a toric grammar than the grammar itself that corresponds to the usual notion of context-free grammar, while a toric grammar in its entirety corresponds to a stochastic context-free grammar.

The weights will serve to keep track of word frequencies through the process of splitting a set of sentences to obtain a toric grammar.

Our aim is indeed to build a toric grammar from a text. To be consistent with our definition of grammars, we will also define texts as positive measures (only on whole sentences, so that no non-terminal other than $[_0$ is used). Let us give a formal definition. We will forget the sentence order, a text will be an unordered set of sentences with possible repetitions.

Definition 1.2.2

The set \mathfrak{T} of texts is the set of toric grammars with integer weights supported by $\mathfrak{S}([_0 D^+)$, that is the set of toric grammars with integer weights using only one non terminal symbol, the start symbol $[_0$.

In this definition, it should be understood that

$$[_0 D^+ = \left\{ \left([_0, w_1, \dots, w_k \right), \text{ where } k \in \mathbb{N} \setminus \{0\} \text{ and } w_i \in D, 1 \leq i \leq k \right\},$$

and that

$$\mathfrak{S}([_0 D^+) = \bigcup_{e \in [_0 D^+} \mathfrak{S}(e).$$

1.2.2 A roadmap towards a communication model

We will use toric grammars as intermediate steps to define the transition probabilities of our communication model on texts. To this purpose, we will first introduce some general types of transformations on toric grammars (reminding the reader that in our formalism texts are some special subset of toric grammars).

It will turn out that two types of expressions, global expressions and local expressions, will play different roles. Let us define them respectively as

$$\begin{aligned}\mathcal{E}_g &= \mathcal{E} \cap \mathfrak{G}([\![_0 S^+), \\ \mathcal{E}_\ell &= \mathcal{E} \cap \mathfrak{G}([\![_+ S^+),\end{aligned}$$

where we remind that $[\![_+ = \{[\![_i, i \in \mathbb{N} \setminus \{0\}\}$ and $S^+ = \bigcup_{j=1}^{\infty} S^j$. Roughly speaking, global expressions correspond to full sentences, whereas local expressions correspond to partial constituents. Any toric grammar $\mathcal{G} \in \mathfrak{G}$ can be accordingly decomposed into $\mathcal{G} = \mathcal{G}_g + \mathcal{G}_\ell$, where $\mathcal{G}_g(A) = \mathcal{G}(A \cap \mathcal{E}_g)$ and $\mathcal{G}_\ell(A) = \mathcal{G}(A \cap \mathcal{E}_\ell)$, for any subset $A \subset \mathcal{E}$.

The transitions of the communication chain with kernel $q_\theta(\mathcal{T}, \mathcal{T}')$ will be defined in two steps. The first step consists in learning from the text \mathcal{T} a toric grammar \mathcal{G} . To this purpose we will split the sentences of \mathcal{T} into syntactic constituents. The second step consists in merging the constituents again to produce a random new text \mathcal{T}' . The parameter $\theta = \mathcal{R}$ of the communication kernel q_θ , will also be a toric grammar. The role of this reference grammar \mathcal{R} will be to provide a stock of local expressions to be used when computing \mathcal{G} from \mathcal{T} . We will discuss later the question of the estimation of \mathcal{R} itself, (and bring satisfactory statistical answers to this issue only in the forthcoming chapters). For the time being, we will assume that the reference grammar \mathcal{R} is a parameter of the communication chain, known to all involved speakers.

We could have defined a communication kernel $q_{\widehat{\mathcal{R}}(\mathcal{T})}(\mathcal{T}, \mathcal{T}')$, where the reference grammar $\widehat{\mathcal{R}}(\mathcal{T})$ itself is estimated at each step from the current text \mathcal{T} , but we would have obtained a model with weaker properties, where, in particular, all the states are not necessarily recurrent states. On the other hand, the proof that the reachable set from any starting point is finite still holds for this modified model, so that it does provide an alternative way of defining a language model as described in the introduction.

We will still need an estimator $\widehat{\mathcal{R}}(\mathcal{T})$ of the reference grammar, in order to provide a language estimator $\widehat{Q}_{\widehat{\mathcal{R}}(\mathcal{T}), \mathcal{T}'}$, where we are using the notations of eq. (1.1.1) on page 18. We will propose a first tentative method for the estimation $\widehat{\mathcal{R}}(\mathcal{T})$ of the reference grammar in this chapter, and study this question in more details in the two following ones.

1.3 Operations on toric grammars

1.3.1 Non stochastic syntax splitting and merging

Let us now describe the model, starting with the description of some non random grammar transformations. We already introduced a model for grammars that includes texts as a special case. We have now to describe how to generate a toric grammar from a text, with, or without, the help of a reference grammar to learn the local component of the grammar. The mechanism producing a grammar from a text will be some sort of random parse algorithm (or rather tentative parse algorithm).

All of this will be achieved by two transformations on toric grammars that will respectively *split* and *merge* expressions (syntagms) of a toric grammar into smaller or bigger ones. We will first describe the sets of possible splits and merges from a given grammar. This will serve as a basis to define random transitions from one grammar to another in subsequent sections.

Let us first introduce some elementary operations involving toric grammars.

$$\begin{aligned} e \oplus f &= \sum_{s \in \mathfrak{G}(e)} \delta_s + \sum_{s \in \mathfrak{G}(f)} \delta_s, & e, f \in \mathcal{E}, \\ e \ominus f &= \sum_{s \in \mathfrak{G}(e)} \delta_s - \sum_{s \in \mathfrak{G}(f)} \delta_s, & e, f \in \mathcal{E}, \\ \rho \otimes e &= \rho \sum_{s \in \mathfrak{G}(e)} \delta_s, & \rho \in \mathbb{R}, e \in \mathcal{E}, \end{aligned}$$

The first operation builds a toric grammar containing expressions e and f with weights 1, and the third one builds a toric grammar containing expression e with weight ρ .

We can generalize these notations to be able to take the sum of a toric grammar and an expression, as well as the sum of two toric grammars.

$$\begin{aligned} \mathcal{G} \oplus e &= \mathcal{G} + \sum_{s \in \mathfrak{G}(e)} \delta_s, & \mathcal{G} \in \mathfrak{G}, e \in \mathcal{E} \\ \mathcal{G} \ominus e &= \mathcal{G} - \sum_{s \in \mathfrak{G}(e)} \delta_s, & \mathcal{G} \in \mathfrak{G}, e \in \mathcal{E} \\ \mathcal{G} \oplus \mathcal{G}' &= \mathcal{G} + \mathcal{G}', & \mathcal{G}, \mathcal{G}' \in \mathfrak{G}. \end{aligned}$$

With these notations, a split is described as

$$\mathcal{G}' = \mathcal{G} \ominus ab \oplus a]_i \oplus]_i b, \quad \mathcal{G}, \mathcal{G}' \in \mathfrak{G},$$

the fact that $\mathcal{G}, \mathcal{G}' \in \mathfrak{G}$ implying that

$$i \in \mathbb{N} \setminus \{0\}, ab, a]_i,]_i b \in \mathcal{E} \text{ and } \mathcal{G}(ab) \geq 1.$$

The (partial) order relation $\mathcal{G} \leq \mathcal{G}'$ will also be defined by the rule

$$\mathcal{G} \leq \mathcal{G}' \iff \mathcal{G}' - \mathcal{G} \in \mathfrak{G},$$

or equivalently

$$\mathcal{G} \leq \mathcal{G}' \iff \mathcal{G}' - \mathcal{G} \in \mathcal{M}_+(\mathcal{E}).$$

Let us resume our example. Starting from the one sentence text

$$\mathcal{T} = 1 \otimes [{}_0 \text{ This is my friend Peter } .$$

we get after splitting the grammar

$$\mathcal{G} = [{}_0 \text{ This is }]_1 \text{ Peter } . \oplus [{}_1 \text{ my friend}$$

which can also be written as

$$\mathcal{G} = \text{ Peter } . [{}_0 \text{ This is }]_1 \oplus [{}_1 \text{ my friend}$$

In this example, as well as in the following, punctuation marks are treated as words, so that here the required dictionary has to include the set

$$\{is, friend, my, Peter, This, .\}.$$

Splitting a sentence providing a new label for each split does not create generalization, since it allows only to merge back two expressions that came from the same split. To create a grammar capable of yielding new sentences, we need some label identification scheme. We will perform label identification through the more general process of label remapping, identification being a consequence of the fact that the map may not be one to one. Let

$$\mathfrak{F} = \{f : \mathbb{N} \rightarrow \mathbb{N} \text{ such that } f(0) = 0\}$$

be the set of label maps. For any symbol $]_i$ or $[_i$, let us define $f(]_i) =]_{f(i)}$ and $f([_i) = [_{f(i)}$. Let us also define for any word $w \in D$, $f(w) = w$ and for any expression $e = (w_0, \dots, w_{\ell-1})$, $f(e) = (f(w_0), \dots, f(w_{\ell-1}))$. Since any grammar $\mathcal{G} \in \mathfrak{G}$ is a measure on the set of expressions \mathcal{E} , we can define its image measure by f , considered as a map from \mathcal{E} to \mathcal{E} . We will put $f(\mathcal{G}) = \mathcal{G} \circ f^{-1}$, meaning that $f(\mathcal{G})(A) = \mathcal{G}(f^{-1}(A))$, for any subset $A \subset \mathcal{E}$.

Definition 1.3.1

Two label maps f and $g \in \mathfrak{F}$ are said to be isomorphic if there is a one to one label map $h \in \mathfrak{F}$ such that $g = h \circ f$. In this case $h^{-1} \in \mathfrak{F}$ and $f = h^{-1} \circ g$. Two grammars \mathcal{G} and $\mathcal{G}' \in \mathfrak{G}$ are said to be isomorphic if there is a one to one label map $f \in \mathfrak{F}$ such that $f(\mathcal{G}) = \mathcal{G}'$. In this case, $f^{-1}(\mathcal{G}') = \mathcal{G}$ and we will write $\mathcal{G} \equiv \mathcal{G}'$. If f and g are two isomorphic label maps, then for any toric grammar $\mathcal{G} \in \mathfrak{G}$, $f(\mathcal{G})$ and $g(\mathcal{G})$ are isomorphic grammars.

In the following of this work, to ease notations and simplify exposition, we will freely identify isomorphic label maps and isomorphic grammars and often speak of them as if they were equal.

This being put, we proceed with the introduction of a set of grammar transformations β that consist in a split with possible label remapping. The *split* will be the core component for generating a toric grammar from a text, by splitting the sentences in smaller parts (syntagms).

Definition 1.3.2 (Splitting rule)

For any $\mathcal{G} \in \mathfrak{G}$, let us consider

$$\beta(\mathcal{G}) = \left\{ f(\mathcal{G}'), f \in \mathfrak{F}, \mathcal{G}' \in \mathfrak{G}, \mathcal{G}' = \mathcal{G} \ominus ab \oplus a]_i \oplus [_i b \right\} \subset \mathfrak{G}.$$

Let us remark that in this definition, necessarily, $ab, a]_i, [_i b \in \mathcal{E}$, $i \in \mathbb{N} \setminus \{0\}$, $1 \otimes ab \leq \mathcal{G}$, and $a]_i \oplus [_i b \leq \mathcal{G}'$. Let us put

$$\beta^*(\mathcal{G}) = \bigcup_{n=0}^{+\infty} \underbrace{\beta \circ \dots \circ \beta}_{n \text{ times}}(\mathcal{G}),$$

the set of grammars that can be constructed from repeated invocations of β .

Lemma 1.3.1

Let us recall that $S = D \cup \{ [_i,]_i, i \in \mathbb{N} \}$ and let us put $S^* = \bigcup_{n=0}^{+\infty} S^n$. For any text $\mathcal{T} \in \mathfrak{T}$, and any $\mathcal{G} \in \beta^*(\mathcal{T})$, \mathcal{G} is a toric grammar with integer weights,

$$\begin{aligned} \mathcal{G}([_i S^*) &= \mathcal{G}([_i S^*), & i \in \mathbb{N} \setminus \{0\}, \\ \mathcal{G}(wS^*) &\leq \mathcal{T}(wS^*), & w \in (D \cup \{ [_0 \})^+, \\ \mathcal{G}(wS^*) &= \mathcal{T}(wS^*), & w \in (D \cup \{ [_0 \}), \end{aligned}$$

and in particular

$$\mathcal{G}([_0 S^*) = \mathcal{T}([_0 S^*).$$

This means that in any toric grammar obtained by splitting a text, the weights of expressions containing the two forms $]_i$ and $[_i$ of a label are balanced, the word frequencies are the same in the grammar and in the text, and the number of sentences contained in the text is given by the total weight of expressions containing the start symbol $[_0$ in the grammar.

PROOF. For the first assertion, an induction on the number of applications of β yields the result, since

$$\mathcal{T}([_i S^*) = \mathcal{T}(\lfloor_i S^*) = 0, i \in \mathbb{N} \setminus \{0\},$$

and, for any $\mathcal{G}' = \mathcal{G} \ominus ab \oplus a \rfloor_i \oplus \lfloor_i b$, and any label $j \in \mathbb{N} \setminus \{0, i\}$,

$$\mathcal{G}'([_j S^*) = \mathcal{G}([_j S^*), \quad (1.3.1)$$

$$\mathcal{G}'(\lfloor_j S^*) = \mathcal{G}(\lfloor_j S^*), \quad (1.3.2)$$

whereas

$$\mathcal{G}'([_i S^*) = \mathcal{G}([_i S^*) + 1, \quad (1.3.3)$$

$$\mathcal{G}'(\lfloor_i S^*) = \mathcal{G}(\lfloor_i S^*) + 1. \quad (1.3.4)$$

For the second assertion, it suffices to remark that the weight of expressions beginning with a given sequence of words is not increased by application of β . Indeed, any word sequence $w \in (D \cup \{[_0\}^+)$ appears the same number of times at the beginning of an expression of $1 \otimes ab$ and of $a \rfloor_i \oplus \lfloor_i b$ if w is not split between a and b and appears a smaller number of times if w is split between a and b . When w is a single word, it cannot be split, so we have equality, as stated in the third assertion of the lemma. The subsequent application of a label map f to \mathcal{G}' does not change the counts involved in the lemma. \square

This lemma is important, because we will subsequently want to impose restrictions on the splitting rule based on word frequencies. Our choice to define a new type of grammar as a positive measure on symbol sequences was made to keep track of word frequencies throughout the construction.

Let us now describe the reverse of a splitting transformation, that we will call a merge transformation. This transformation will be central in generating new texts from a toric grammar, by merging the syntagms into bigger ones, ending with a full sentence.

Definition 1.3.3 (Merge rule)

For any toric grammar $\mathcal{G} \in \mathfrak{G}$ we consider the following set of allowed merge transformations

$$\alpha(\mathcal{G}) = \left\{ \mathcal{G}' \in \mathfrak{G}, \mathcal{G}' = \mathcal{G} \ominus a \rfloor_i \ominus \lfloor_i b \oplus ab \right\}.$$

Let us remark that in this definition, necessarily $i \in \mathbb{N} \setminus \{0\}$, $a \rfloor_i, \lfloor_i b, ab \in \mathcal{E}$, and $a \rfloor_i \oplus \lfloor_i b \leq \mathcal{G}$.

The merge transformation is indeed the reverse of the *split*, in the sense that:

Lemma 1.3.2

For any $\mathcal{G}, \mathcal{G}' \in \beta^*(\mathfrak{T})$, $\mathcal{G}' \in \beta(\mathcal{G})$ if, and only if, there is $f \in \mathfrak{F}$ such that $f(\mathcal{G}) \in \alpha(\mathcal{G}')$.

PROOF. Let us suppose that $\mathcal{G}' = f(\mathcal{G} \oplus a]_i \oplus [{}_i b \ominus ab)$ is in $\beta(\mathcal{G})$. In this case, $\mathcal{G}' = f(\mathcal{G}) \oplus f(a]_i) \oplus f([{}_i b) \ominus f(ab)$, so that both $f(a]_i)$ and $f([{}_i b)$ are in $\text{supp}(\mathcal{G}')$, $f(ab) \in \text{supp}(f(\mathcal{G}))$, and consequently $f(a]_i), f([{}_i b)$ and $f(ab) \in \mathcal{E}$. Moreover $f(\mathcal{G}) = \mathcal{G}' \oplus f(a)f(b) \ominus f(a)]_{f(i)} \ominus [{}_{f(i)} f(b)$, so that $f(\mathcal{G}) \in \alpha(\mathcal{G}')$.

On the other hand, if for some $f \in \mathfrak{F}$, $f(\mathcal{G}) \in \alpha(\mathcal{G}')$, $f(\mathcal{G}) = \mathcal{G}' \oplus ab \ominus a]_i \ominus [{}_i b$. Since $ab \in \text{supp}(f(\mathcal{G}))$, there is $e \in \mathcal{E}$ such that $f(e) = ab$. But this implies that there is $c, d \in S^+$ such that $a = f(c)$ and $b = f(d)$. We can then if needed modify f outside $\{j \in \mathbb{N} : [{}_j S^* \in \text{supp}(\mathcal{G})\}$, to make sure that $i \in f(\mathbb{N})$. Let $f(j) = i$. We now get that $f(\mathcal{G}) = \mathcal{G}' \oplus f(c)f(d) \ominus f(c)]_{f(j)} \ominus [{}_{f(j)} f(d)$, so that finally $\mathcal{G}' = f(\mathcal{G} \oplus c]_j \oplus [{}_j d \ominus cd)$, proving that $\mathcal{G}' \in \beta(\mathcal{G})$. \square

Another useful property of the merge rule is given by the following lemma:

Lemma 1.3.3

For any $f \in \mathfrak{F}$ and any $\mathcal{G} \in \mathfrak{G}$, $f(\alpha(\mathcal{G})) \subset \alpha(f(\mathcal{G}))$.

PROOF. Indeed, any $\mathcal{G}' \in f(\alpha(\mathcal{G}))$ is of the form

$$\begin{aligned} \mathcal{G}' &= f(\mathcal{G} \oplus ab \ominus a]_i \ominus [{}_i b) \\ &= f(\mathcal{G}) \oplus f(a)f(b) \ominus f(a)]_{f(i)} \ominus [{}_{f(i)} b \in \alpha(f(\mathcal{G})). \end{aligned} \quad \square$$

Unfortunately, repeating the merge transformation will not provide a text in all circumstances. Indeed, we can end up with some expressions of the type $[{}_i a]_i b$. However, since an expression is allowed to contain only one opening bracket, we are sure that $[{}_0 \notin \text{supp}([{}_i a]_i b)$, and that all sentences (global expressions, beginning with $[{}_0)$ are “complete”, in the sense that they do not contain any other non-terminal symbol.

To continue the discussion, we will switch to a random context, where split and merge transformations are performed according to some probability measure.

1.3.2 Random split and merge processes

The grammars we described so far are obtained using splitting rules. Texts can be reconstructed using merge transformations. The splitting rules as well as the merge rules allow for multiple choices at each step. We will account for this by introducing random processes where these choices are made at random.

We will describe two types of random grammar transformations. Each of these will appear as a finite length Markov chain, where the length of the chain is given by a uniformly bounded stopping time.

- The learning process (or splitting process) will start with a text and build a grammar through iterated splits;
- the production process will start with a grammar and produce a text through iterated merge operations.

These two types of processes may be combined into a split and merge process, going back and forth between texts and toric grammars.

Let us give more formal definitions. Learning and parsing processes will be some special kinds of splitting processes, to be defined hereafter.

Definition 1.3.4 (Splitting process)

Given some restricted splitting rule $\beta_r : \mathfrak{G} \rightarrow 2^{\mathfrak{G}}$ from the set of grammars to the set of subsets of \mathfrak{G} , such that for any $\mathcal{G} \in \mathfrak{G}$, $\beta_r(\mathcal{G}) \subset \beta(\mathcal{G})$, a splitting process is a time homogeneous stopped Markov chain $S_t, 0 \leq t \leq \tau$ defined on \mathfrak{G} such that

$$\tau = \inf\{t \in \mathbb{N} : \beta_r(S_t) = \emptyset\},$$

$$\mathbb{P}(S_t = \mathcal{G}' \mid S_{t-1} = \mathcal{G}) > 0 \iff \mathcal{G}' \in \beta_r(\mathcal{G}).$$

Definition 1.3.5 (Production process)

A production process is a time homogenous stopped Markov chain $P_t, 0 \leq t \leq \sigma$ defined on \mathfrak{G} such that

$$\sigma = \inf\{t \in \mathbb{N}, \alpha(P_t) = \emptyset\},$$

and

$$\mathbb{P}(P_t = \mathcal{G}' \mid P_{t-1} = \mathcal{G}) > 0 \iff \mathcal{G}' \in \alpha(\mathcal{G}).$$

Definition 1.3.6 (Split and Merge process)

Given a splitting process $S_t, t \in \mathbb{N}$ and a production process $P_t, t \in \mathbb{N}$, a split and merge process is a Markov chain $G_t \in \mathfrak{G}, t \in \mathbb{N}$, with transitions

$$\mathbb{P}(G_{2t+1} = \mathcal{G}' \mid G_{2t} = \mathcal{G}) = \mathbb{P}(S_\tau = \mathcal{G}' \mid S_0 = \mathcal{G}), \quad t \in \mathbb{N},$$

$$\mathbb{P}(G_{2t} = \mathcal{G}' \mid G_{2t-1} = \mathcal{G}) = \mathbb{P}(P_\sigma = \mathcal{G}' \mid P_0 = \mathcal{G}, P_\sigma \in \mathfrak{T}), \quad t \in \mathbb{N} \setminus \{0\},$$

whose initial distribution is a probability measure on texts, so that almost surely $G_0 \in \mathfrak{T}$.

Let us remark that we have to impose the condition that $P_\sigma \in \mathfrak{T}$, because the production process does not produce a true text with probability one. On the other hand it can yield back G_{2t-2} with positive probability when started at G_{2t-1} , as will be proved later on. Therefore $\mathbb{P}(P_\sigma \in \mathfrak{T} \mid P_0 = \mathcal{G}) > 0$ for any \mathcal{G} such that $\mathbb{P}(G_{2t-1} = \mathcal{G}) > 0$. One way to simulate $\mathbb{P}_{G_{2t} \mid G_{2t-1}}$ is to use a rejection method, simulating repeatedly from the production process until a true text is produced. In the experiments we made, $\mathbb{P}(P_\sigma \in \mathfrak{T} \mid P_0 = \mathcal{G})$ was close to one and rejection a rare event. We will describe the relations between these processes and parsing rules issued from the statistical model of the following chapters in section 3.5 on page 151.

Proposition 1.3.4

Let S_t , P_t and G_t be a splitting process, a production process and the corresponding split and merge process, starting from $G_0 = \mathcal{T} \in \mathfrak{T}$. For any $\mathcal{G} \in \mathfrak{G}$, any $\mathcal{T}' \in \mathfrak{T}$, such that $\sum_{t \in \mathbb{N}} \mathbb{P}(G_{2t+1} = \mathcal{G}) > 0$ and $\sum_{t \in \mathbb{N}} \mathbb{P}(G_{2t} = \mathcal{T}') > 0$,

$$\mathbb{P}\left(\tau \leq 2[\mathcal{T}(DS^*) - \mathcal{T}([{}_0 S^*)] \mid S_0 = \mathcal{T}'\right) = 1, \quad (1.3.5)$$

$$\mathbb{P}\left(\sigma \leq 2[\mathcal{T}(DS^*) - \mathcal{T}([{}_0 S^*)] \mid P_0 = \mathcal{G}\right) = 1. \quad (1.3.6)$$

In other words, the length of all the splitting and production processes involved in the split and merge process have a uniform bound, given by twice the difference between the number of words and the number of sentences in the original text.

PROOF. This proof is a bit lengthy and is based on some invariants in the split and merge operations. It has been put off to section 1.A.1 on page 49. \square

Proposition 1.3.5

If G_t is a split and merge process starting almost surely from the text $G_0 = \mathcal{T} \in \mathfrak{T}$, there is a finite subset of toric grammars $\mathfrak{G}_{\mathcal{T}}$ such that with probability equal to one there is for each time t a grammar G'_t isomorphic to G_t such that $G'_t \in \mathfrak{G}_{\mathcal{T}}$. Thus, after identification of isomorphic grammars, we can analyze the split and merge process as a finite state Markov chain, since the reachable set from any starting point is finite. We should however keep in mind that the finite state space $\mathfrak{G}_{\mathcal{T}}$ depends on the initial state \mathcal{T} , so the state space is still infinite, although any trajectory will almost surely stay in a finite subset of reachable states.

PROOF. Let us assume that the labels of \mathcal{G} are taken from $\llbracket 0, W_\ell(\mathcal{G}) \rrbracket$, where, as defined in the appendix, $W_\ell(\mathcal{G}) = \sum_{i=1}^{\infty} \mathcal{G}([{}_i S^*)$ is the number of labels (with their multiplicities) used in the grammar \mathcal{G} . Consequently $\mathcal{G}([{}_i S^*) = 0$ for $i > W_\ell(\mathcal{G})$. This can be achieved, up to grammar isomorphisms, by applying to \mathcal{G} a suitable

label map.

Let us define the set of canonical expressions as

$$\mathcal{E}_c = \mathcal{E} \cap \left(\bigcup_{i \in \mathbb{N}} [{}_i S^* \right),$$

and the canonical decomposition of \mathcal{G} as

$$\mathcal{G} = \sum_{e \in \mathcal{E}_c} \mathcal{G}(e) \otimes e.$$

We see that \mathcal{G} can be described by the concatenation of the canonical expressions, each repeated a number of times equal to its weight, to form a sequence of symbols of length $W_s(\mathcal{G})$. From the proof of the previous proposition, we know that

$$W_s(\mathcal{G}) \leq M = 5W_w(\mathcal{T}) - 3W_e(\mathcal{T}) = 5\mathcal{T}(DS^*) - 3\mathcal{T}([{}_0 S^*),$$

(with notations defined in the appendix). We can represent \mathcal{G} by a sequence of exactly M symbols by padding with trailing $[{}_0$ symbols the representation described above. Let us give an example

$$\mathcal{G} = 2 \otimes [{}_0 w_1]_1 w_2 \oplus [{}_1 w_3 \oplus [{}_1 w_4$$

can be coded as

$$[{}_0 w_1]_1 w_2 [{}_0 w_1]_1 w_2 [{}_1 w_3 [{}_1 w_4 [{}_0 [{}_0 [{}_0$$

in the case when $M = 15$. Let us consider the set of symbols

$$S_{\mathcal{T}} = D \cup \left\{ [{}_0, [{}_i,]_i, 0 < i \leq 2 \left[\mathcal{T}(DS^*) - \mathcal{T}([{}_0 S^*) \right] \right\}.$$

Since \mathcal{G} uses only those symbols, we see from the proposed coding of \mathcal{G} that it can take at most $|S_{\mathcal{T}}|^M$ different values. Since

$$|S_{\mathcal{T}}| = |D| + 1 + 4 \left[\mathcal{T}(DS^*) - \mathcal{T}([{}_0 S^*) \right],$$

we have proved that

$$|\mathfrak{G}_{\mathcal{T}}| \leq \left(|D| + 1 + 4 \left[\mathcal{T}(DS^*) - \mathcal{T}([{}_0 S^*) \right] \right)^{5\mathcal{T}(DS^*) - 3\mathcal{T}([{}_0 S^*)}.$$

Let us notice that this bound, while being finite, is very large, and probably quite loose in practice. \square

1.3.3 Splitting rules and label identification

In the previous section, we introduced some class of random processes, and studied some of their general properties. In this section, we are going to describe some more specific schemes and go further in the description of split and merge processes that can learn toric grammars in a satisfactory way.

The choice of splitting rules and label identification rules has a decisive influence on the way syntactic categories and syntactic rules are learnt by the split and merge process. While it is necessary as a starting point to consider rules learnt from the text to be parsed itself, it will also be fruitful to consider the case when a previously learnt reference grammar $\mathcal{R} \in \mathfrak{G}$ can be used to govern the splits.

To make things easier to grasp, let us explain on some example the basics of syntactic generalization by label identification. Let us start with the simple text with two sentences.

$$G_0 = \mathcal{T} = [{}_0 \text{ This is my friend Peter } . \oplus [{}_0 \text{ This is my neighbour John } .$$

If we split “my friend” and “my neighbour” in the two sentences using the same label, we will form after two splits the grammar

$$G_1 = [{}_0 \text{ This is }]_1 \text{ Peter } . \oplus [{}_0 \text{ This is }]_1 \text{ John } . \\ \oplus [{}_1 \text{ my friend } \oplus [{}_1 \text{ my neighbour } .$$

If no more splits are allowed and we therefore reached the stopping time of the splitting process, so that $\tau = 2$, we can proceed to the production process, and reach after two more steps the new text G_2 that can either be $G_2 = G_0$ or

$$G_2 = [{}_0 \text{ This is my neighbour Peter } . \oplus [{}_0 \text{ This is my friend John } .$$

Now is a good time to remind the reader of the distinction made in section 1.2.2 on page 23 about local and global expressions.

Legitimate local expressions will be provided by the reference grammar \mathcal{R} , whereas global expressions will be deduced from the text itself. This approach will be particularly efficient in the case when the set of local expressions is smaller than the set of global expressions.

We will need two different kinds of split processes, one to learn the reference grammar from a text and the other one to perform the first part of the transitions of the communication Markov chain.

These split processes may be viewed as performing some parsing of the text they are applied to. Here, we do not use parsing as it is usually used to discover whether a sentence is correct or not, we use it instead to discover new expressions.

We will start by defining the parsing rules to be used in the communication chain. We will call them *narrow* parsing rules. We will then proceed to the

definition of a *broad* parsing rule suitable for a first go at learning the reference grammar $\widehat{\mathcal{R}}(\mathcal{T})$ from a text.

Definition 1.3.7

Let us define the narrow parsing rule with reference grammar \mathcal{R} as

$$\beta_n(\mathcal{G}, \mathcal{R}) = \left\{ \mathcal{G}' \in \mathfrak{G} : \mathcal{G}' = \mathcal{G} \oplus a]_i \oplus []_i b \ominus ab, \right. \\ \left. ab \in \mathcal{E}_g, \mathcal{R}([]_i b) > 0 \right\}, \quad \mathcal{G} \in \mathfrak{G}.$$

Let us remark that, due to the definition of the set of expressions \mathcal{E} and of toric grammars $\mathfrak{G} \subset \mathcal{M}_+(\mathcal{E})$, the fact that \mathcal{G} and $\mathcal{G}' \in \mathfrak{G}$ implies that $i \in \mathbb{N} \setminus \{0\}$ in this definition, since necessarily $a]_i, []_i b \in \mathcal{E}$. It implies also that $[]_0 \in \text{supp}(a)$, a condition equivalent to $ab \in \mathcal{E}_g$.

The narrow parsing rule depends on \mathcal{R} only through $\text{supp}(\mathcal{R}) \cap \mathcal{E}_i$.

Let us define the broad parsing rule as

$$\beta_b(\mathcal{G}, \mathcal{R}) = \left\{ \mathcal{G}' \in \mathfrak{G} : \mathcal{G}' = \mathcal{G} \oplus a]_i \oplus []_i b \ominus ab, \right. \\ \left. \mathcal{R}(a]_i) + \mathcal{R}([]_i b) > 0, \mathcal{R}(aS^*) \leq \mu_1 \mathcal{R}([]_0 S^*), \right. \\ \left. \text{and } \mathcal{R}(bS^*) \leq \mu_2 \mathcal{R}([]_0 S^*) \right\}, \quad \mathcal{G}, \mathcal{R} \in \mathfrak{G},$$

where $\mu_1, \mu_2 \in \mathbb{R}_+$ are two positive real parameters.

Since the reference grammar is under construction during broad parsing, we will mainly use this rule with $\mathcal{R} = \mathcal{G}$, as will be explained later. The same learning parameters μ_1 and μ_2 are present here and in the innovation rule to be described next. They serve to split expressions into sufficiently infrequent halves, in order to constrain the model.

Let us define now maximal sequences, a notion that will be needed to define learning rules.

Definition 1.3.8

Given some toric grammar \mathcal{G} , we will say that $a \in S^+$ is \mathcal{G} -maximal and write $a \in \text{max}(\mathcal{G})$ when

$$\mathcal{G}(aS^*) > \max\{\mathcal{G}(awS^*), \mathcal{G}(waS^*), w \in S\}.$$

In other words, a is a maximal subsequence among the subsequences with the same weight in \mathcal{G} . Note that if a is \mathcal{G} -maximal, usually $\mathcal{G}(a) = 0$ (meaning that a is not an expression of the grammar, but only a subexpression) and if moreover the grammar \mathcal{G} has integer weights (which will be the case if it has been produced by a split and merge process), then $\mathcal{G}(aS^*) \geq 2$.

Definition 1.3.9 (Innovation rule)

Using the notations $[_+ = \{[_i, i \in \mathbb{N} \setminus \{0\}\}$ and $]_+ = \{]_i, i \in \mathbb{N} \setminus \{0\}\}$, let us define the innovation rule with reference grammar \mathcal{R} as

$$\begin{aligned} \beta_i(\mathcal{G}, \mathcal{R}) &= \left\{ \mathcal{G}' \in \mathfrak{G} : \mathcal{G}' = \mathcal{G} \oplus a]_i \oplus [_i b \ominus ab, \right. \\ &\quad \mathcal{R}([_i S^*) = 0, \{a, b\} \cap \max(\mathcal{R}) \neq \emptyset, \\ &\quad \left. \mathcal{R}(aS^*) \leq \mu_1 \mathcal{R}([_0 S^*), \text{ and } \mathcal{R}(bS^*) \leq \mu_2 \mathcal{R}([_0 S^*) \right\}. \end{aligned}$$

Here again, the rule will be used while learning the reference grammar with $\mathcal{R} = \mathcal{G}$.

We will now introduce a label map that identifies the labels appearing in the same context.

Definition 1.3.10 (Label identification through context)

Given some toric grammar $\mathcal{G} \in \mathfrak{G}$, let us consider the relation $C \in (\mathbb{N} \setminus \{0\})^2$ defined as

$$C = \left\{ (i, j) \in (\mathbb{N} \setminus \{0\})^2 : \sum_{a \in S^*} \mathcal{G}(a]_i) \mathcal{G}(a]_j) + \mathcal{G}([_i a) \mathcal{G}([_j a) > 0 \right\}.$$

The smallest equivalence relation containing C defines a partition of $\mathbb{N} \setminus \{0\}$ into equivalence classes. Let $(A_k)_{k \in \mathbb{N} \setminus \{0\}}$ be an arbitrary indexing of this partition. Each positive integer falls in a unique class of the partition, so that the relation $i \in A_{\underline{\chi}_{\mathcal{G}}(i)}$ defines a label map $\underline{\chi}_{\mathcal{G}} : \mathbb{N} \rightarrow \mathbb{N}$ in a non ambiguous way. The choice of the indexing of the partition $(A_k)_{k \in \mathbb{N} \setminus \{0\}}$ does not matter, since two different choices lead to two isomorphic label maps. When applying $\underline{\chi}_{\mathcal{G}}$ to \mathcal{G} itself, we will use the short notation $\underline{\chi}(\mathcal{G}) \stackrel{\text{def}}{=} \underline{\chi}_{\mathcal{G}}(\mathcal{G})$.

Let us consider the evolution of the number of labels used by \mathcal{G} :

$$L(\mathcal{G}) = \left| \{i \in \mathbb{N} : \mathcal{G}([_i S^*) > 0\} \right|.$$

It is easy to see that $L(\underline{\chi}(\mathcal{G})) \leq L(\mathcal{G})$ and that $\underline{\chi}(\mathcal{G}) \equiv \mathcal{G}$ (where the symbol \equiv means isomorphic) if and only if $L(\underline{\chi}_{\mathcal{G}}(\mathcal{G})) = L(\mathcal{G})$. Accordingly there is $k \in \mathbb{N}$ such that $\underline{\chi}^{k+1}(\mathcal{G}) \equiv \underline{\chi}^k(\mathcal{G})$, and we can take it to be the smallest integer such that $L(\underline{\chi}^{k+1}(\mathcal{G})) = L(\underline{\chi}^k(\mathcal{G}))$. Consequently, k is such that for any $n \geq k$, $\underline{\chi}^n(\mathcal{G}) \equiv \underline{\chi}^k(\mathcal{G})$. We will define $\chi(\mathcal{G}) = \underline{\chi}^k(\mathcal{G})$, up to grammar isomorphisms (so that $\chi(\mathcal{G})$ belongs to \mathfrak{G}/\equiv rather than to \mathfrak{G} itself).

A characterisation in terms of more elementary label maps will be established in proposition 1.A.6 on page 58. This characterization provides an algorithm to compute χ in practice.

We are now ready to define a learning rule.

Definition 1.3.11

Let us define the learning rule

$$\beta_\ell(\mathcal{G}) = \begin{cases} \beta_i(\mathcal{G}, \mathcal{G}), & \text{when } \beta_b(\mathcal{G}, \mathcal{G}) = \emptyset, \\ \{\chi(\mathcal{G}') : \mathcal{G}' \in \beta_b(\mathcal{G}, \mathcal{G})\}, & \text{otherwise.} \end{cases}$$

We will define two kinds of splitting processes, based on two different choices of the restricted splitting rule β_r .

Definition 1.3.12 (Learning process)

A learning process is a splitting process with restricted splitting rule

$$\beta_r(\mathcal{G}) = \beta_\ell(\mathcal{G}).$$

Definition 1.3.13 (Parsing process)

A parsing process with reference grammar $\mathcal{R} \in \mathfrak{G}$ is a splitting process with restricted splitting rule

$$\beta_r(\mathcal{G}) = \beta_n(\mathcal{G}, \mathcal{R}).$$

Before we reach the aim of this chapter and describe our statistical language model, we need to explore some of the properties of the production, learning and parsing processes introduced so far.

1.4 Parsing and generalization

Let us introduce some notations for the output of parsing, learning and production processes.

Definition 1.4.1

Let S_t be a parsing process, with reference grammar $\mathcal{R} \in \mathfrak{G}$. We will use the following notation for the distribution of S_τ .

$$\mathbf{G}_{\mathcal{T}, \mathcal{R}} = \mathbb{P}_{S_\tau | S_0 = \mathcal{T}}, \quad \mathcal{T} \in \mathfrak{T}.$$

We will also use a short notation for the distribution of the output of a production process.

$$\mathbf{T}_{\mathcal{G}} = \mathbb{P}_{P_\sigma | P_0 = \mathcal{G}, P_\sigma \in \mathfrak{T}}, \quad \mathcal{G} \in \mathfrak{G}.$$

Eventually, $\mathbf{G}_{\mathcal{T}}$ will be the probability distribution of the output of a learning process S_t , according to the definition

$$\mathbf{G}_{\mathcal{T}} = \mathbb{P}_{S_\tau | S_0 = \mathcal{T}}, \quad \mathcal{T} \in \mathfrak{T}.$$

At this point we obviously may consider different notions of parsing that we have to connect together. Namely, we would like to make a link between the following statements:

- $\mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0$, the grammar \mathcal{G} can produce the text \mathcal{T} ;
- $\mathbb{G}_{\mathcal{T}, \mathcal{R}}(\mathcal{G}) > 0$, the text \mathcal{T} can generate the grammar \mathcal{G} when parsed with the help of the grammar \mathcal{R} ;
- $\mathbb{G}_{\mathcal{T}}(\mathcal{G}) > 0$, the grammar \mathcal{G} can be learnt from the text \mathcal{T} .

Lemma 1.4.1

The previous parse notions are related in the following way. For any $\mathcal{G}, \mathcal{R} \in \mathfrak{G}$, and any $\mathcal{T} \in \mathfrak{T}$,

$$\begin{aligned} \mathbb{G}_{\mathcal{T}}(\mathcal{G}) > 0 &\implies \mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0, \\ \mathbb{G}_{\mathcal{T}, \mathcal{R}}(\mathcal{G}) > 0 &\implies \mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0, \\ \mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0 &\implies \mathbb{G}_{\mathcal{T}, \mathcal{G}}(\mathcal{G}) > 0. \end{aligned}$$

Consequently, for any $\mathcal{G}, \mathcal{R} \in \mathfrak{G}$ such that $(\text{supp}(\mathcal{G}) \cap \mathcal{E}_l) \subset \text{supp}(\mathcal{R})$, and any $\mathcal{T} \in \mathfrak{T}$,

$$\mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0 \iff \mathbb{G}_{\mathcal{T}, \mathcal{R}}(\mathcal{G}) > 0.$$

PROOF. This is one of the core lemmas of this chapter. The proof is given in section 1.A.2 on page 51, on account of its length. \square

It has the following important implication.

Proposition 1.4.2

Given a parsing process S_t based on a reference grammar $\mathcal{R} \in \mathfrak{G}$ and a production process P_t , the corresponding split and merge process G_t is weakly reversible, in the sense that for any $\mathcal{T} \in \mathfrak{T}$, any $\mathcal{G} \in \bigcup_{t \in \mathbb{N}} \text{supp}(\mathbb{P}_{G_{2t+1}})$,

$$\mathbb{P}(G_1 = \mathcal{G} \mid G_0 = \mathcal{T}) > 0 \iff \mathbb{P}(G_2 = \mathcal{T} \mid G_1 = \mathcal{G}) > 0.$$

Consequently, for any $\mathcal{T}, \mathcal{T}' \in \mathfrak{T}$ and any $\mathcal{G}, \mathcal{G}' \in \bigcup_{t \in \mathbb{N}} \text{supp}(\mathbb{P}_{G_{2t+1}})$,

$$\begin{aligned} \mathbb{P}(G_2 = \mathcal{T}' \mid G_0 = \mathcal{T}) > 0 &\iff \mathbb{P}(G_2 = \mathcal{T} \mid G_0 = \mathcal{T}') > 0, \\ \mathbb{P}(G_3 = \mathcal{G}' \mid G_1 = \mathcal{G}) > 0 &\iff \mathbb{P}(G_3 = \mathcal{G} \mid G_1 = \mathcal{G}') > 0. \end{aligned}$$

In other words, the two processes G_{2t} and G_{2t+1} are weakly reversible time homogeneous Markov chains. As we already proved that the set of reachable states from any starting point is finite, it shows that they are recurrent Markov chains: they partition their respective state spaces into positive recurrent communicating classes.

PROOF. Let us remark first that

$$\begin{aligned}\mathbb{P}(G_1 = \mathcal{G} \mid G_0 = \mathcal{T}) > 0 &\iff \mathbb{G}_{\mathcal{T}, \mathcal{R}}(\mathcal{G}) > 0 \\ \mathbb{P}(G_2 = \mathcal{T} \mid G_1 = \mathcal{G}) > 0 &\iff \mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0.\end{aligned}$$

Moreover, since $\mathcal{G} \in \text{supp}(\mathbb{P}_{G_{2t+1}})$ for some $t \in \mathbb{N}$, there is $\mathcal{T}' \in \mathfrak{T}$ such that $\mathbb{G}_{\mathcal{T}', \mathcal{R}}(\mathcal{G}) > 0$, implying that $\text{supp}(\mathcal{G}) \cap \mathcal{E}_l \subset \text{supp}(\mathcal{R})$. This ends the proof according to the last statement of the previous lemma. \square

We will see in section 3.5 on page 151 a method to build split and merge processes that are actually reversible, and not only weakly so. This method involves a slightly modified Metropolis algorithm, that allows for intermediate steps (here, the grammar G_{2t+1}).

1.5 Expectation of a random toric grammar

In section 1.4 on page 35, given some text $\mathcal{T} \in \mathfrak{T}$, we defined a random distribution on toric grammars $\mathbb{G}_{\mathcal{T}}$ that we would like to use to learn a grammar from a text. The most obvious way to do this is to draw a toric grammar at random according to the distribution $\mathbb{G}_{\mathcal{T}}$, and we already saw an algorithm, described by a Markov chain and a stopping time, to do this.

The distribution $\mathbb{G}_{\mathcal{T}}$ will be spread in general on many grammars. This is a kind of instability that we would like to avoid, if possible. A natural way to get rid of this instability would be to simulate the expectation of $\mathbb{G}_{\mathcal{T}}$. To do this, we are facing a problem: the usual definition of the expectation of $\mathbb{G}_{\mathcal{T}}$, that is

$$\int \mathcal{G} \, d\mathbb{G}_{\mathcal{T}}(\mathcal{G}),$$

although well defined from a mathematical point of view, is a meaningless toric grammar, due to the possible fluctuations of the label mapping. To get a meaningful notion of expectation, we need to define in a meaningful way the sum of two toric grammars. We will achieve this in two steps.

Let us introduce first the *disjoint sum* of two toric grammars. We will do this with the help of two disjoint label maps. Let us define the *even* and *odd* label maps f_e and f_o as

$$f_e(i) = 2i, \quad f_o(i) = \max\{0, 2i - 1\}, \quad i \in \mathbb{N}.$$

Definition 1.5.1

The disjoint sum of two toric grammars $\mathcal{G}, \mathcal{G}' \in \mathfrak{G}$ is defined as

$$\mathcal{G} \boxplus \mathcal{G}' = f_e(\mathcal{G}) + f_o(\mathcal{G}').$$

Definition 1.5.2

Given a probability measure $\mathbb{G} \in \mathcal{M}_+^1(\mathfrak{G})$ with finite support, we define the mean of \mathbb{G} as

$$\oint \mathcal{G} \, d\mathbb{G}(\mathcal{G}) = \chi \left(\boxplus_{\mathcal{G} \in \mathfrak{G}} \mathbb{G}(\mathcal{G}) \mathcal{G} \right).$$

Lemma 1.5.1

If G_i is an i.i.d. sequence of random grammars distributed according to \mathbb{G} , then almost surely

$$\lim_{n \rightarrow +\infty} \frac{1}{n} \chi \left(\boxplus_{i=1}^n G_i \right) = \oint \mathcal{G} \, d\mathbb{G}(\mathcal{G}).$$

PROOF. The proof of this result is quite lengthy, and postponed till section 1.A.3 on page 54. \square

1.6 Language models

1.6.1 Communication model

We are now ready to define the language model announced in the introduction. Given a reference grammar \mathcal{R} , and the corresponding split and merge process $(G_t)_{t \in \mathbb{N}}$ with reference \mathcal{R} , we define the communication kernel $q_{\mathcal{R}}(\mathcal{T}, \mathcal{T}')$ on \mathfrak{T}^2 as

$$q_{\mathcal{R}}(\mathcal{T}, \mathcal{T}') = \mathbb{P}(G_2 = \mathcal{T}' \mid G_0 = \mathcal{T}).$$

According to proposition 1.3.5 on page 30 and proposition 1.4.2 on page 36, $q_{\mathcal{R}}$ has finite reachable sets and is weakly reversible, so that all texts $\mathcal{T} \in \mathfrak{T}$ are positive recurrent states of the communication kernel $q_{\mathcal{R}}$.

Thus to each text $\mathcal{T} \in \mathfrak{T}$ corresponds a unique invariant text distribution $\hat{q}_{\mathcal{R}}(\mathcal{T}, \cdot)$, as explained in the introduction. As all states are positive recurrent, $\hat{q}_{\mathcal{R}}(\mathcal{T}, \cdot)$ is the unique invariant measure of $q_{\mathcal{R}}$ on the communicating class containing \mathcal{T} . Moreover, from the ergodic theorem,

$$\mathbb{P} \left(\hat{q}_{\mathcal{R}}(\mathcal{T}, \cdot) = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{j=1}^t \delta_{G_{2j}} \mid G_0 = \mathcal{T} \right) = 1,$$

showing that $\hat{q}_{\mathcal{R}}(\mathcal{T}, \cdot)$ can be computed by an almost surely convergent Monte-Carlo simulation. Eventually, we deduce from the invariant probability measure on texts $\hat{q}_{\mathcal{R}}(\mathcal{T}, \cdot)$, a probability measure on sentences $\hat{Q}_{\mathcal{R}, \mathcal{T}}$ as explained in the introduction, according to the formula

$$\hat{Q}_{\mathcal{R}, \mathcal{T}} = \mathcal{T}([_0S^*])^{-1} \sum_{\mathcal{T}' \in \mathfrak{T}} \hat{q}_{\mathcal{R}}(\mathcal{T}, \mathcal{T}') \mathcal{T}'.$$

(This is the same formula as eq. (1.1.1) on page 18 in the introduction, taking into account the fact that texts in the support of $\hat{q}_{\mathcal{R}}(\mathcal{T}, \cdot)$ are non normalized empirical measures with the same total mass equal to $\mathcal{T}([_0S^*])$, the number of sentences in the text \mathcal{T} .)

To obtain a true language estimator, there remains to estimate \mathcal{R} by some estimator $\widehat{\mathcal{R}}(\mathcal{T})$. We could do this as described in section 1.5 on page 37, putting

$$\widehat{\mathcal{R}}(\mathcal{T}) = \oint \mathcal{G} \, d\mathbb{G}_{\mathcal{T}}(\mathcal{G}).$$

Let us remark that, according to lemma 1.5.1 on the preceding page, $\widehat{\mathcal{R}}(\mathcal{T})$ can be computed from repeated simulations from the distribution $\mathbb{G}_{\mathcal{T}}$.

1.6.2 Comparison with other models

Comparison with context-free grammars

Given a toric grammar $\mathcal{G} \in \beta^*(\mathfrak{T})$, we may consider the split and merge process G_t with reference grammar \mathcal{G} starting at $G_1 = \mathcal{G}$ (so here we start at time 1 with an initial state that is a grammar, instead of starting at time 0 with an initial state that is a text). Due to the weak reversibility of proposition 1.4.2 on page 36, G_2 almost surely falls in the same recurrent communicating class of $t \mapsto G_{2t}$, and the unique invariant probability measure supported by this recurrent communicating class defines a probability measure $\overline{\mathbb{T}}_{\mathcal{G}}$ on texts, and therefore a stochastic language model. This way of defining the language generated by the grammar \mathcal{G} can be compared to the usual definition of the language generated by a context-free grammar. Indeed, the support of \mathcal{G} is a context-free grammar, so this is meaningful to consider the language generated by this grammar and to compare it with the support of our stochastic language model.

None of these two sets of sentences is contained in the other one. In our stochastic model, the number of times a rule can be used is bounded, so if the recursive use of some rules is possible, the deterministic language will in this sense be larger. On the other hand, the stochastic model uses both production and parsing to build new sentences, whereas the deterministic model uses only

production rules. In this respect, the stochastic model may, at least in some cases, define a much broader language, as we will show on the following example.

Let us take as dictionary the set

$$D = \{+, =\} \cup \llbracket 1, N \rrbracket,$$

where $\llbracket 1, N \rrbracket = \{i \in \mathbb{N}, 1 \leq i \leq N\}$, and consider the toric grammar

$$\mathcal{G} = N^2 \otimes "[_0]_N = N" \oplus \bigoplus_{i=1}^N N \otimes "[_i i]" \oplus \bigoplus_{i=2}^N N(i-1) \otimes "[_i]_{i-1} + 1",$$

(where we put expressions into double quotes " to ease the reading of this definition.) Let us also consider the text

$$\mathcal{T} = N \otimes \bigoplus_{i=1}^N "[_0 \underbrace{i+1+\dots+1}_{N-i \text{ times}} = N".$$

It is easy to check that $\mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0$, (so that $\mathcal{G} \in \beta^*(\mathcal{T})$), that indeed the support of \mathcal{T} is the language generated by $\text{supp}(\mathcal{G})$, seen as a context-free grammar, and that the stochastic language $\overline{\mathbb{T}}_{\mathcal{G}}$ generated by \mathcal{G} is able to produce with positive probability a set of sentences

$$\text{supp}^2(\overline{\mathbb{T}}_{\mathcal{G}}) \stackrel{\text{def}}{=} \bigcup_{\mathcal{T} \in \text{supp}(\overline{\mathbb{T}}_{\mathcal{G}})} \text{supp}(\mathcal{T}),$$

equal to

$$\text{supp}^2(\overline{\mathbb{T}}_{\mathcal{G}}) = \left\{ "[_0 x_1 + \dots + x_i = x_{i+1} + \dots + x_j", \right. \\ \left. 1 \leq i < j \leq 2N, x_k \in \llbracket 1, N \rrbracket, 1 \leq k \leq j, \sum_{k=1}^i x_k = \sum_{k=i+1}^j x_k = N \right\}.$$

Here, the number of sentences produced by the underlying context-free grammar is $|\text{supp}(\mathcal{T})| = N$, whereas the number of sentences produced by our stochastic language model is $|\text{supp}^2(\overline{\mathbb{T}}_{\mathcal{G}})| = 2^{2(N-1)}$. Thus, in this small example based on arithmetic expressions (admittedly closer to a computing language than it is to a natural language), our new definition of the generated language induces a huge increase in the number of generated sentences.

Note that with usual context-free grammar notations, $\text{supp}(\mathcal{G})$ would have been described as

$$\begin{aligned} \boxed{0} &\rightarrow \boxed{N} = N \\ \boxed{i} &\rightarrow i, & i = 1, \dots, N, \\ \boxed{i} &\rightarrow \boxed{i-1} + 1, & i = 2, \dots, N, \end{aligned}$$

where $\boxed{0}$ is the start symbol and \boxed{i} , $i = 1, \dots, N$, are other non terminal symbols.

To count the number of elements in $\text{supp}^2(\overline{\mathbb{T}}_{\mathcal{G}})$, one can remark that the number of ways N can be written as $\sum_{k=1}^i x_k$ with an arbitrary number of terms is also the number of increasing integer sequences $0 < s_1 < \dots < s_{i-1} < N$ of arbitrary length, which is also the number of subsets $\{s_1, \dots, s_{i-1}\}$ of $\{1, \dots, N-1\}$, that is 2^{N-1} .

Intuitively speaking, the underlying context-free grammar $\text{supp}(\mathcal{G})$ is limited to producing a small set of global expressions of the form $i + 1 + \dots + 1 = N$, whereas the stochastic language model incorporates some crude logical reasoning that is capable of deducing from them a large set of new global expressions.

Let us remark also that, when we start as here from a text made of true arithmetic statements, the language generated by our language model is also made of true arithmetic statements. This shows that our approach to language modeling is capable of some sort of logical reasoning.

Comparison with Markov models

The kind of reasoning illustrated in the previous section is related to the fact that we analyse global syntactic structures represented by the global expressions of our toric grammars.

In order to give another point of comparison, we would like in this section to make a qualitative comparison with Markov models, that do not share this feature. To make a parallel between toric grammars and Markov models, we are going to show how a Markov model could be described in terms of toric grammars and label identification rules. More striking relations between our models and Markov models will however be outlined in forthcoming chapters.

To build a Markov model in our framework, we have to use a deterministic splitting (or parsing) rule. This is because in a Markov model, conditional probabilities are specified from left to right in a rigid data independent way. Let us introduce the Markov splitting rule

$$\beta_m(\mathcal{G}) = \left\{ \mathcal{G}' \in \mathfrak{G}, \mathcal{G}' = \mathcal{G} \ominus [{}_0 aw]_i \oplus [{}_0 a]_j \oplus [{}_j w]_i, \right. \\ \left. i, j \in \mathbb{N} \setminus \{0\}, a \in D^+, w \in D, \mathcal{G}([{}_j S^*] = 0) \right\}.$$

We will describe now label identification rules using concepts introduced in section 1.A.3 on page 54. Let us say that the pair of labels $p \in (\mathbb{N} \setminus \{0\})^2$ is \mathcal{G} -Markov if there is $w \in D$ such that $\mathcal{G}(w]_{p_1} S^*) \mathcal{G}(w]_{p_2} S^*) > 0$. Let us say that the sequence of pairs p^1, \dots, p^k is \mathcal{G} -Markov if p^j is $\xi_{p^1, \dots, p^{j-1}}(\mathcal{G})$ -Markov. It can be proved as in the case of congruent sequences that if σ is a permutation and p is \mathcal{G} -Markov, then $p \circ \sigma$ is also \mathcal{G} -Markov. It can also be proved that if p and q are maximal \mathcal{G} -Markov

sequences, then $\xi_p \equiv \xi_q$, and therefore $\xi_p(\mathcal{G}) \equiv \xi_q(\mathcal{G})$. We will call $\xi_p(\mathcal{G}) \in \mathfrak{G} / \equiv$ the Markov closure of \mathcal{G} and use the notation $\xi_p(\mathcal{G}) \stackrel{\text{def}}{=} \mu(\mathcal{G})$, where $\mu(\mathcal{G})$ is the Markov pendent of $\chi(\mathcal{G})$ in the construction of toric grammars.

Let S_t , $0 \leq t \leq \tau$ be a splitting process based on the restricted splitting rule

$$\beta_r(\mathcal{G}) = \left\{ \mu(\mathcal{G}'), \mathcal{G}' \in \beta_m(\mathcal{G}) \right\}.$$

It is not very difficult to check that the support of S_τ is contained in a single isomorphic class of grammars, so that, up to label remapping the result of this splitting process is deterministic. More specifically, starting from a text

$$\mathcal{T} = \bigoplus_{j=1}^n [{}_0 w_1^j \dots w_{\ell(j)}^j],$$

where $w_i^j \in D \setminus \{.\}$, $1 \leq i < \ell(j)$, $1 \leq j \leq n$, and $w_{\ell(j)}^j = .$, $1 \leq j \leq n$ so that all sentences end with a period, we obtain a grammar isomorphic to

$$\mathcal{G} = \bigoplus_{j=1}^n \left([{}_0 w_1^j]_{w_1^j} \bigoplus_{i=2}^{\ell(j)-1} [{}_{w_{i-1}^j} w_i^j]_{w_i^j} \oplus [{}_{w_{\ell(j)-1}^j} w_{\ell(j)}^j] \right),$$

where we have used words as labels instead of integers, since in this model, due to the label identification rule, labels are functions of words (namely $[_w]$ is the non terminal symbol following the word $w \in D$).

We can now define a Markov production mechanism, to replace the production process. It is described as a Markov chain X_i , $i \in \mathbb{N}$, where $X_i \in D \cup \{\Delta\}$, where $\Delta \notin D$ is a padding symbol used to embed finite sentences into infinite sequences of symbols, all equal to Δ for indices larger than the sentence length. The distribution of the Markov chain X_i is as follows. Its initial distribution is

$$\mathbb{P}(X_0 = w) = \frac{\mathcal{G}([{}_0 w]_w)}{\mathcal{G}([{}_0 S^*])},$$

and its transition probabilities are

$$\begin{aligned} \mathbb{P}(X_i = \Delta \mid X_{i-1} = .) &= 1, \\ \mathbb{P}(X_i = . \mid X_{i-1} = w) &= \frac{\mathcal{G}([{}_w .])}{\mathcal{G}([{}_w S^*])}, & w \in D \setminus \{.\} \\ \mathbb{P}(X_i = w' \mid X_{i-1} = w) &= \frac{\mathcal{G}([{}_w w']_{w'})}{\mathcal{G}([{}_w S^*])}, & w, w' \in D \setminus \{.\}. \end{aligned}$$

Roughly speaking, the difference with the production process P_t defined previously is that in the production process the production rules are drawn at random without replacement whereas here, the production rules are drawn with replacement.

It is easy to see that the initial distribution and transition probabilities of the Markov chain X_i are the empirical initial distribution and empirical transition probabilities of the training text \mathcal{T} .

In conclusion, to build a Markov model using the same framework as for toric grammars, we had to modify two steps in a dramatic way:

- we had to change the splitting process, and replace the random splitting process of toric grammars with a non random splitting process which chains forward transitions in a linear way;
- we had to change in a dramatic way the label identification rule to replace the *forward and backward global condition* of toric grammars with a *backward only local condition*.

(The modification of the production process is less crucial and boils down to drawing production rules with or without replacement.)

We hope that this discussion of Markov models will help the reader realize that our model proposal is indeed really different from the Markov model at sentence level. We could have extended easily the discussion to Markov models of higher order, or to more general context tree models. We let the reader figure out the details. All these more sophisticated models show the same differences from toric grammars: a more rigid splitting process and local backward label identification rules.

1.7 A small experiment

Let us end this chapter with a small example. Here we use a small text that is meant to mimic what could be found in a tutorial to learn English as a foreign language. We have added a more elaborate sentence at the end of the text to show its impact. More systematic experiments are yet to be carried out, although the conception of this model was guided by experimental trial and errors with models starting with variable length Markov chains, before we tried global rules leading to grammars.

This is the training text \mathcal{T} (each line shows an expression, starting with its weight) :

```
1 [0 He is a clever guy .
1 [0 He is doing some shopping .
1 [0 He is laughing .
```

1 [0 He is not interested in sports .
 1 [0 He is walking .
 1 [0 He likes to walk in the streets .
 1 [0 I am driving a car .
 1 [0 I am riding a horse too .
 1 [0 I am running .
 1 [0 Paul is crossing the street .
 1 [0 Paul is driving a car .
 1 [0 Paul is riding a horse .
 1 [0 Paul is walking .
 1 [0 Peter is walking .
 1 [0 While I was walking , I saw Paul crossing the street .

And now, the new sentences produced by the model (that is by $\hat{Q}_{\mathcal{R}, \mathcal{T}}$, approximated on 50 iterations of the communication chain with kernel $q_{\hat{\mathcal{R}}}$).

1 [0 Paul is driving a car too .
 1 [0 Paul is doing some shopping .
 1 [0 Paul is laughing .
 1 [0 Paul is riding a horse too .
 1 [0 Paul is running too .
 1 [0 Paul is running .
 1 [0 Paul is not interested in sports too .
 1 [0 Paul is not interested in sports .
 1 [0 Paul is a clever guy too .
 1 [0 Paul is a clever guy .
 1 [0 Paul is walking too .
 1 [0 Peter is driving a car too .
 1 [0 Peter is driving a car .
 1 [0 Peter is doing some shopping .
 1 [0 Peter is laughing .
 1 [0 Peter is riding a horse too .
 1 [0 Peter is riding a horse .
 1 [0 Peter is running too .
 1 [0 Peter is running .
 1 [0 Peter is not interested in sports .
 1 [0 Peter is a clever guy .
 1 [0 Peter is crossing the street .
 1 [0 He is driving a car too .
 1 [0 He is driving a car .
 1 [0 He is riding a horse too .
 1 [0 He is riding a horse .
 1 [0 He is running too .
 1 [0 He is running .

1 [0 He is not interested in sports too .
1 [0 He is crossing the street too .
1 [0 He is crossing the street .
1 [0 He is walking too .
1 [0 I am driving a car too .
1 [0 I am doing some shopping .
1 [0 I am laughing too .
1 [0 I am laughing .
1 [0 I am riding a horse .
1 [0 I am not interested in sports .
1 [0 I am a clever guy .
1 [0 I am crossing the street too .
1 [0 I am crossing the street .
1 [0 I am walking too .
1 [0 I am walking .
1 [0 While I was driving a car , I saw Paul doing some shopping too .
1 [0 While I was driving a car , I saw Paul doing some shopping .
1 [0 While I was driving a car , I saw Paul riding a horse .
1 [0 While I was driving a car , I saw Paul crossing the street .
1 [0 While I was driving a car , I saw Paul walking .
1 [0 While I was driving a car , I saw Peter riding a horse .
1 [0 While I was doing some shopping , I saw Paul riding a horse .
1 [0 While I was doing some shopping , I saw Paul walking .
1 [0 While I was laughing too , I saw Peter crossing the street .
1 [0 While I was laughing , I saw Peter riding a horse .
1 [0 While I was riding a horse , I saw Paul driving a car too .
1 [0 While I was riding a horse , I saw Paul driving a car .
1 [0 While I was riding a horse , I saw Paul laughing .
1 [0 While I was riding a horse , I saw Paul running .
1 [0 While I was riding a horse , I saw Paul walking .
1 [0 While I was riding a horse , I saw Peter not interested in sports .
1 [0 While I was running , I saw Paul laughing .
1 [0 While I was running , I saw Paul not interested in sports .
1 [0 While I was running , I saw Paul a clever guy .
1 [0 While I was running , I saw Paul walking .
1 [0 While I was not interested in sports , I saw Paul driving a car .
1 [0 While I was not interested in sports , I saw Paul riding a horse .
1 [0 While I was a clever guy , I saw Paul running .
1 [0 While I was a clever guy , I saw Paul crossing the street .
1 [0 While I was a clever guy , I saw Paul walking .
1 [0 While I was crossing the street , I saw Paul riding a horse .
1 [0 While I was crossing the street , I saw Paul running .
1 [0 While I was crossing the street , I saw Paul crossing the street .

1 [0 While I was crossing the street , I saw Paul walking .
 1 [0 While I was crossing the street , I saw Peter walking .
 1 [0 While I was walking , I saw Paul driving a car .
 1 [0 While I was walking , I saw Paul laughing .
 1 [0 While I was walking , I saw Paul riding a horse .
 1 [0 While I was walking , I saw Paul running .
 1 [0 While I was walking , I saw Paul not interested in sports .
 1 [0 While I was walking , I saw Paul crossing the street too .
 1 [0 While I was walking , I saw Paul walking .
 1 [0 While I was walking , I saw Peter not interested in sports .
 1 [0 While I was walking , I saw Peter walking .

The reference grammar was learnt first, and was computed from 10 samples of $G_{\mathcal{G}}$. (We did not normalize the weights, since we were interested in the support of the local expressions only.)

10 [0 He likes to walk]6]3 streets .
 2 [0]1]8 clever guy .
 2 [0]1 doing some shopping .
 2 [0]1 laughing .
 2 [0]1 not interested]6 sports .
 2 [0]1 riding]8 horse .
 2 [0]1 riding]8 horse]2 .
 2 [0]1 running .
 24 [0]7 am]5 .
 28 [0 Paul is]5 .
 40 [0 He is]5 .
 4 [0]1 crossing]3 street .
 4 [0]1 driving]8 car .
 5 [0]4 is]5 .
 6 [0]1 walking .
 7 [0 Peter is]5 .
 8 [0 While]7 was]5 ,]7 saw]4]5 .
 10 [1 He is
 2 [1 Peter is
 2 [1 While]7 was]5 ,]7 saw]4
 6 [1]7 am
 8 [1 Paul is
 2 [2 too
 30 [3 the
 14 [4 Paul
 1 [4 Peter
 16 [5 crossing]3 street
 16 [5 driving]8 car

```

16 [5 riding ]8 horse
34 [5 walking
8 [5 ]5 too
8 [5 ]8 clever guy
8 [5 doing some shopping
8 [5 laughing
8 [5 not interested ]6 sports
8 [5 running
20 [6 in
50 [7 I
50 [8 a

```

Although we did not yet make the software development effort required to test large text corpora, we learnt a few interesting things from what we already tried:

- As it is, the model requires the inclusion of a sufficient number of simple and redundant sentences to start generalizing. At this stage, we do not know whether this could be avoided by changing the learning rules. We made quite a few attempts in this direction. All of them resulted in the production of grammatical nonsense. Breaking the global constraints that are enforced by the model seems to have a dramatic effect on grammatical coherence. This could be a clue that these global conservation rules reflect some fundamental feature of the syntactic structure of natural languages. Including a bunch of “simple” sentences made of frequent words may be seen as introducing a pinch of supervision in the learning process.
- The constraints on subexpressions frequencies in the learning rule definition 1.3.7 on page 33 and 1.3.9 were added to avoid some unwanted generalizations. For instance here we took $\mu_1 \mathcal{R}([_0 S^*) = \mu_2 \mathcal{R}([_0 S^*) = 5$. If we had chosen 10 instead of 5, sentences of the kind

```
[0 While I was walking , I saw He crossing the street .
```

would have emerged, where the pronoun “He” is substituted to a noun in the wrong place. We deliberately wrote the training text in such a way that “He” is more frequent than any noun, since we expect that to be true for any reasonable large corpus. Doing so, we were able to rule out the wrong construction by lowering the frequency constraint to avoid the unwanted substitution.

- Despite all the limitations of this small example, it shows that the model is able to find out non trivial new constructs, like

```
[0 While I was laughing too, I saw Peter crossing the street.
```

where it has discovered that “too” could be added to the subordinate clause opening the sentence. We are quite pleased to see that such things could be learnt along very general label identification rules, while all the generalized sentences remain, if not all grammatically correct, at least all grammatically plausible. Of course this judgement is purely subjective. But since we have no mathematical or otherwise quantitative definition of what natural languages are, we have to be content with a subjective evaluation of models.

Studying how this learning model scales with large corpora is still a work to be done (it will require from us that we optimize our code so that it can run efficiently on large data sets).

1.8 The story so far...

We have built in this chapter a new statistical framework for the syntactic analysis of natural languages.

The main idea pervading our approach is that trying to estimate the distribution of an isolated random sentence is hopeless. Instead we propose to build a Markov chain on sets of sentences (called texts in this work), with non trivial recurrent communicating classes and to define our language model as the invariant measures of this Markov chain on each of these recurrent communicating classes. At each step, the Markov chain recombines the set of sentences constituting its current state, using cut and paste operations described by grammar rules. In this way we define the probability distribution of an isolated random sentence only in an indirect way. We replace the hard question of generating a random sentence by the hopefully simpler one of recombining a set of sentences in a way that keep the desired distribution invariant.

The main result of this chapter is the construction of the split and merge process with reference grammar \mathcal{R} . It has non trivial mathematical properties proving that it can be simulated using a bounded number of operations at each step, and that the state space is divided into recurrent communicating classes each including a finite number of states.

Nonetheless, the model needs some further refinements. For example, we did not (yet) propose a simple way to define the actual probabilities on the split and merge process. But more importantly, our method to estimate the reference grammar through the grammar expectation

$$\widehat{\mathcal{R}}(\mathcal{T}) = \int \mathcal{G} \, dG_{\mathcal{T}}(\mathcal{G})$$

is clearly lacking. The next chapter will be dedicated to introducing the notion of Markov substitute sets, which will be used to estimate the reference grammar

within a more traditional statistical framework, enjoying proved mathematical properties.

1.A Proofs

1.A.1 Bound on the length of splitting and production processes

PROOF OF PROPOSITION 1.3.4 ON PAGE 30. Let us define the length of an expression $e \in S^k \cap \mathcal{E}$ as $\ell(e) = k$. Let us introduce some remarkable weights associated with a grammar $\mathcal{G} \in \beta^*(\mathfrak{T})$.

$$\begin{aligned} W_s(\mathcal{G}) &= \sum_{e \in \mathcal{E}} \mathcal{G}(e), \\ W_e(\mathcal{G}) &= \sum_{e \in \mathcal{E}} \mathcal{G}(e) \ell(e)^{-1}, \\ W_l(\mathcal{G}) &= \sum_{i=1}^{+\infty} \mathcal{G}([_i S^*), \\ W_w(\mathcal{G}) &= \sum_{w \in D} \mathcal{G}(w S^*). \end{aligned}$$

Let us define the set of canonical expressions as

$$\mathcal{E}_c = \mathcal{E} \cap \left(\bigcup_{i \in \mathbb{N}} [_i S^* \right).$$

Using previously introduced notations, we can write the grammar as

$$\mathcal{G} = \sum_{e \in \mathcal{E}_c} \mathcal{G}(e) \otimes e.$$

We will call this the canonical decomposition of \mathcal{G} . The two weights $W_s(\mathcal{G})$ and $W_e(\mathcal{G})$ are better understood in terms of this canonical decomposition. They can be expressed as

$$\begin{aligned} W_s(\mathcal{G}) &= \sum_{e \in \mathcal{E}_c} \mathcal{G}(e) \ell(e), \\ W_e(\mathcal{G}) &= \sum_{e \in \mathcal{E}_c} \mathcal{G}(e). \end{aligned}$$

This shows that $W_s(\mathcal{G})$ counts the “number of symbols” in the canonical decomposition of \mathcal{G} , whereas $W_e(\mathcal{G})$ counts the number of expressions (that is $\mathcal{G}(\mathcal{E}_c)$), the

weight put by the grammar on canonical expressions). We can also see from the definitions that $W_l(\mathcal{G})$ counts the number of canonical expressions starting with a positive (that is non terminal) label, that we will call for short the number of labels, and that $W_w(\mathcal{G})$ counts the number of words.

Since a split increases the number of canonical expressions by one, the number of symbols in canonical expressions by two, the number of labels by one, and keeps the number of words constant, whereas a merge decreases these quantities in the same proportions, the following quantities are invariant in all the toric grammars involved: for any $\mathcal{G} \in \mathfrak{G}$ such that $\sum_{t \in \mathbb{N}} \mathbb{P}(G_t = \mathcal{G}) > 0$,

$$\begin{aligned} W_s(\mathcal{G}) - 2W_e(\mathcal{G}) &= W_s(\mathcal{T}) - 2W_e(\mathcal{T}), \\ W_e(\mathcal{G}) - W_l(\mathcal{G}) &= W_e(\mathcal{T}) - W_l(\mathcal{T}) = W_e(\mathcal{T}), \\ W_w(\mathcal{G}) &= W_w(\mathcal{T}). \end{aligned}$$

Moreover, for the same reasons, for any $\mathcal{T}' \in \mathfrak{T}$ and $\mathcal{G} \in \mathfrak{G}$ such that $\sum_{t \in \mathbb{N}} \mathbb{P}(G_{2t} = \mathcal{T}') > 0$ and $\sum_{t \in \mathbb{N}} \mathbb{P}(G_{2t+1} = \mathcal{G}) > 0$,

$$\begin{aligned} \mathbb{P}(\tau = W_l(S_\tau) \mid S_0 = \mathcal{T}') &= 1, \\ \mathbb{P}(\sigma = W_l(\mathcal{G}) \mid P_0 = \mathcal{G}, P_\sigma \in \mathfrak{T}) &= 1. \end{aligned}$$

Thus, we will prove the lemma if we can bound $W_l(\mathcal{G})$ (or equivalently $W_l(S_\tau)$ when $S_0 = \mathcal{T}'$, since S_τ almost surely satisfies the conditions imposed on \mathcal{G}). We can then remark that

$$\begin{aligned} \sum_{e \in \mathcal{E}_c} \mathcal{G}(e) \mathbf{1}[\ell(e) \geq 3] &\leq \sum_{e \in \mathcal{E}_c} \mathcal{G}(e) [\ell(e) - 2] = W_s(\mathcal{G}) - 2W_e(\mathcal{G}), \\ \sum_{e \in \mathcal{E}_c} \mathcal{G}(e) \mathbf{1}[\ell(e) = 2] &= \sum_{e \in \mathcal{E}^e} \mathcal{G}(e) \mathbf{1}[\ell(e) = 2] \sum_{w \in D} \mathbf{1}(e \in wS^*) \\ &\leq \sum_{e \in \mathcal{E}^e} \mathcal{G}(e) \sum_{w \in D} \mathbf{1}(e \in wS^*) = W_w(\mathcal{G}), \end{aligned}$$

because any canonical expression of length 2 is of the form $e = [i]w$, with $i \in \mathbb{N}$ and $w \in D$, so that for any $e \in \mathcal{E}_c$ of length 2,

$$\sum_{e' \in \mathfrak{G}(e)} \sum_{w \in D} \mathbf{1}(e' \in wS^*) = 1.$$

Thus

$$W_e(\mathcal{G}) \leq W_w(\mathcal{G}) + W_s(\mathcal{G}) - 2W_e(\mathcal{G}),$$

and consequently we can bound $W_l(\mathcal{G})$ by the split and merge invariant bound

$$W_l(\mathcal{G}) \leq W_l(\mathcal{G}) - W_e(\mathcal{G}) + W_w(\mathcal{G}) + W_s(\mathcal{G}) - 2W_e(\mathcal{G}).$$

This, added to the fact that $W_l(\mathcal{T}) = 0$ and $W_s(\mathcal{T}) = W_w(\mathcal{T}) + W_e(\mathcal{T})$, proves that

$$W_l(\mathcal{G}) \leq 2[W_w(\mathcal{T}) - W_e(\mathcal{T})].$$

This ends the proof, since $W_w(\mathcal{T}) = \mathcal{T}(DS^*)$ and $W_e(\mathcal{T}) = \mathcal{T}([_0S^*])$. \square

1.A.2 Parsing Relations

PROOF OF LEMMA 1.4.1 ON PAGE 36. The implication

$$\mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0 \implies \mathbb{G}_{\mathcal{T},\mathcal{G}}(\mathcal{G}) > 0$$

is less trivial than it may seem. Indeed we can reverse the path of the splitting process S_t , be it a parsing or a learning process, to obtain a path followed with positive probability by the production process, but reversing the production process does not give a parsing process. Let us illustrate this difficulty on a simple example. Consider

$$\mathcal{T} = 1 \otimes [abcd] \quad \text{and} \quad \mathcal{G} = [{}_0a]_1 \oplus [{}_1b]_2 \oplus [{}_2c]_3 \oplus [{}_3d].$$

The production path

$$\mathcal{G}, \quad [{}_0ab]_2 \oplus [{}_2c]_3 \oplus [{}_3d], \quad [{}_0ab]_2 \oplus [{}_2cd], \quad \mathcal{T}$$

has positive probability. The reverse path may have a positive probability for the learning process but not for the parsing process with reference \mathcal{G} , since none of the expressions $[{}_0ab]_2$ or $[{}_2cd]$ belongs to the support of \mathcal{G} . To parse \mathcal{T} according to \mathcal{G} , one can instead follow with positive probability such a path as

$$\mathcal{T}, \quad [{}_0abc]_3 \oplus [{}_3d], \quad [{}_0ab]_2 \oplus [{}_2c]_3 \oplus [{}_3d], \quad \mathcal{G}.$$

To prove the lemma, we will have to show that it is always possible to find such an alternative parsing path. This property is fundamental to our approach, since it proves that the toric grammars we build can be used to parse the texts they can produce.

Let us start with the easiest part of the proof. Assume that $\mathbb{G}_{\mathcal{T},\mathcal{R}}(\mathcal{G}) > 0$. This means that there is a path $\mathcal{G}_0, \dots, \mathcal{G}_k$ such that $\mathcal{G}_0 = \mathcal{T}$, $\mathcal{G}_k = \mathcal{G}$, and at each step $\mathcal{G}_t \in \beta_n(\mathcal{G}_{t-1}, \mathcal{R})$. Anyhow it is easy to check that

$$\mathcal{G}_t \in \beta_n(\mathcal{G}_{t-1}, \mathcal{R}) \implies \mathcal{G}_{t-1} \in \alpha(\mathcal{G}_t),$$

so that the reverse path is followed with a positive probability by the production process. This means that $\mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0$.

In the case of the learning process, if $\mathbb{G}_{\mathcal{G}}(\mathcal{G}) > 0$, there is a path $\mathcal{G}_t, 0 \leq t \leq k$, such that $\mathcal{G}_t \in \beta_{\ell}(\mathcal{G}_{t-1})$, $\mathcal{G}_0 = \mathcal{T}$ and $\mathcal{G}_k = \mathcal{G}$, consequently there is a label map $f_t \in \mathfrak{F}$ such that $f_t(\mathcal{G}_{t-1}) \in \alpha(\mathcal{G}_t)$. We can then remark that

$$f_k \circ \cdots \circ f_t(\mathcal{G}_{t-1}) \in \alpha(f_k \circ \cdots \circ f_{t+1}(\mathcal{G}_t)),$$

because as already proved before in lemma 1.3.3 on page 28, $f(\alpha(\mathcal{G})) \subset \alpha(f(\mathcal{G}))$. Let us consider the path $\tilde{\mathcal{G}}_t = f_k \circ \cdots \circ f_{k-t+1}(\mathcal{G}_{k-t})$. It begins at $\tilde{\mathcal{G}}_0 = \mathcal{G}_k = \mathcal{G}$ and ends at $\tilde{\mathcal{G}}_k = f_k \circ \cdots \circ f_1(\mathcal{G}_0) = \mathcal{T}$. According to the previous remark, this path is followed by the production process with positive probability, proving that $\mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0$.

Let us now come to the proof of the third implication of the lemma. For this let us assume now that $\mathbb{T}_{\mathcal{G}}(\mathcal{T}) > 0$. Consider a path $\mathcal{G}_0, \dots, \mathcal{G}_k$ such that $\mathcal{G}_0 = \mathcal{G}, \dots, \mathcal{G}_k = \mathcal{T}$ and $\mathcal{G}_t \in \alpha(\mathcal{G}_{t-1})$. We are going to define some *decorated* path $\tilde{\mathcal{G}}_0, \dots, \tilde{\mathcal{G}}_k$ with some added parentheses. Introduce a new set of symbols $B = \{(i,)_i, i \in \mathbb{N} \setminus \{0\}\}$ and assume that it is disjoint from the other symbols used so far, so that $B \cap S = \emptyset$. Consider the set of toric grammars $\tilde{\mathfrak{G}}$ based on the enlarged dictionary $D \cup B$, and the projection $\pi : \tilde{\mathfrak{G}} \rightarrow \mathfrak{G}$ defined with the help of the canonical decomposition of toric grammars as

$$\pi \left(\sum_{e \in \tilde{\mathcal{E}}_c} \mathcal{G}(e) \otimes e \right) = \sum_{e \in \tilde{\mathcal{E}}_c} \mathcal{G}(e) \otimes \pi(e),$$

where $\tilde{\mathcal{E}}_c$ is the set of canonical expressions based on the enlarged dictionary $D \cup B$, and where $\pi(e)$ is obtained by removing from the sequence of symbols e the symbols belonging to the decoration set B (that is the parentheses).

Let us put $\tilde{\mathcal{G}}_0 = \mathcal{G}$ and define $\tilde{\mathcal{G}}_t$ for $t = 1, \dots, k$ by induction. We will check on the go that $\pi(\tilde{\mathcal{G}}_t) = \mathcal{G}_t$. It is obviously true for $\tilde{\mathcal{G}}_0$, because $\tilde{\mathcal{G}}_0 \in \mathfrak{G}$, so that $\pi(\tilde{\mathcal{G}}_0) = \tilde{\mathcal{G}}_0 = \mathcal{G}_0$. That said, let us describe the construction of $\tilde{\mathcal{G}}_t$, assuming that $\tilde{\mathcal{G}}_{t-1}$ is already defined, and satisfies $\pi(\tilde{\mathcal{G}}_{t-1}) = \mathcal{G}_{t-1}$. Consider the sequence of symbols a and $b \in S^*$ and the index $i \in \mathbb{N} \setminus \{0\}$ such that

$$\mathcal{G}_t = \mathcal{G}_{t-1} \oplus ab \ominus a]_i \ominus []_i b.$$

Since $\pi(\tilde{\mathcal{G}}_{t-1}) = \mathcal{G}_{t-1}$, and since $a]_i \oplus []_i b \leq \mathcal{G}_{t-1}$, there are $\tilde{a} \in \tilde{S}^*$ and $\tilde{b} \in \tilde{S}^*$ such that $\pi(\tilde{a}) = a$, $\pi(\tilde{b}) = b$, and $\tilde{a}]_i \oplus []_i \tilde{b} \leq \tilde{\mathcal{G}}_{t-1}$. (The choice of \tilde{a} and \tilde{b} may not be unique, in which case we can make any arbitrary choice). Let us define

$$\tilde{\mathcal{G}}_t = \tilde{\mathcal{G}}_{t-1} \oplus \tilde{a}(\tilde{b})_i \ominus \tilde{a}]_i \ominus []_i \tilde{b}.$$

Since $\pi(\tilde{a}(\tilde{b})_i) = \pi(\tilde{a}\tilde{b}) = ab$,

$$\pi(\tilde{\mathcal{G}}_t) = \pi(\tilde{\mathcal{G}}_{t-1}) \oplus \pi(\tilde{a}(\tilde{b})_i) \ominus \pi(\tilde{a}]_i) \ominus \pi([]_i \tilde{b}) = \tilde{\mathcal{G}}_{t-1} \oplus ab \ominus a]_i \ominus []_i b = \mathcal{G}_t,$$

where we have used the obvious fact that π is linear.

We are now going to define another mapping between grammars that allows to recover \mathcal{G} from any $\tilde{\mathcal{G}}_t$ (obviously the decorations were added to keep track of \mathcal{G}). Let us define $\psi : \tilde{\mathfrak{S}}' \rightarrow \mathfrak{S}$ on the set of decorated grammars $\tilde{\mathfrak{S}}'$ which are supported by expressions where the parentheses $(\cdot)_i$ are matched (at the same level) by the formula

$$\psi \left(\sum_{e \in \tilde{\mathcal{E}}_c} \tilde{\mathcal{G}}(e) \otimes e \right) = \sum_{e \in \tilde{\mathcal{E}}_c} \tilde{\mathcal{G}}(e) \psi(e),$$

where $\psi(e)$ is defined by the rules

$$\psi(e) = \begin{cases} \psi([{}_i a]_j c) + \psi([{}_j b]), & \text{if } e = [{}_i a({}_j b)_j] c, \text{ with } a, b, c \in \tilde{S}^* \\ \psi(e) = 1 \otimes e, & \text{otherwise.} \end{cases}$$

It is easy to check that this definition is not ambiguous and that

$$\psi(e) = \psi'(e) \oplus \bigoplus_{(ia)_i \in \text{supp}(e)} [{}_i \psi'(a),$$

where $\psi'(e)$ is the expression obtained from e by replacing all the sequences between outer parentheses pairs $({}_j a)_j$ by $]_j$. This is may be easier to grasp on some example:

$$\psi([{}_0 a({}_1 b({}_2 c)_2 d)_1] e({}_3 f({}_4 g)_4)_3 h) = [{}_0 a]_1 e]_3 h \oplus [{}_1 b]_2 d \oplus [{}_2 c \oplus [{}_3 f]_4 \oplus [{}_4 g.$$

It is easy to check by induction that $\tilde{\mathcal{G}}_t \in \tilde{\mathfrak{S}}'$. Moreover, we have that $\psi(\tilde{\mathcal{G}}_t) = \mathcal{G}$. Indeed $\psi(\tilde{\mathcal{G}}_0) = \psi(\mathcal{G}) = \mathcal{G}$ and

$$\psi(\tilde{G}_t) = \psi(\tilde{G}_{t-1}) \oplus \psi(\tilde{a}({}_i \tilde{b})_i) \ominus \psi(\tilde{a}]_i) \ominus \psi([{}_i \tilde{b}) = \psi(\tilde{G}_{t-1}),$$

since ψ is linear and $\psi(\tilde{a}({}_i \tilde{b})_i) = \psi(\tilde{a}]_i) \oplus \psi([{}_i \tilde{b})$.

We are now going to define a continuation for the path $(\tilde{\mathcal{G}}_t, 0 \leq t \leq k)$ that will bring us back to \mathcal{G} .

We will maintain during our inductive construction two properties:

$$\begin{aligned} \psi(\tilde{\mathcal{G}}_t) &= \mathcal{G}, \\ \text{and } \text{supp}(\tilde{\mathcal{G}}_t) \cap \tilde{\mathcal{E}}_l &\subset \mathcal{E}_l, \end{aligned}$$

where $\tilde{\mathcal{E}}_l$ is the set of local decorated expressions, so that

$$\tilde{\mathcal{E}}_l = \{e \in \tilde{\mathcal{E}} : [{}_0 \notin \text{supp}(e)\}.$$

We already proved that the first property is satisfied by $\tilde{\mathcal{G}}_k$. As $\pi(\tilde{\mathcal{G}}_k) = \mathcal{G}_k = \mathcal{T}$, $\text{supp}(\tilde{\mathcal{G}}_k) \cap \tilde{\mathcal{E}}_l = \emptyset$, so that the second condition is also satisfied. Let us assume that, for some $t > k$, $\tilde{\mathcal{G}}_{t-1}$ has been defined and satisfies the two conditions above, and let us proceed to the construction of $\tilde{\mathcal{G}}_t$.

As long as $\tilde{\mathcal{G}}_{t-1} \notin \mathfrak{G}$, (and this will be the case for $t < 2k$), find some canonical expression $e \in \tilde{\mathcal{E}}_c \setminus \mathcal{E}_c$, such that $\tilde{\mathcal{G}}_{t-1}(e) \geq 1$. From our induction hypotheses, we see that necessarily $[_0 \in \text{supp}(e)$. Our continuation will be such that each such expression has matching parentheses with matching labels, and we will check this on the go while building it by induction. Among those matching pairs of parentheses, there is necessarily at least one inner pair. We can for instance choose the one starting with the last opening parenthesis $(_j$ of the sequence e . This choice makes it obvious that the subsequence of e enclosed between $(_j$ and $)_j$ contains no further parentheses.

Since ψ is linear and preserves positive measures,

$$\mathcal{G} \ominus \psi(e) = \psi(\tilde{\mathcal{G}}_{t-1}) \ominus \psi(e) = \psi(\tilde{\mathcal{G}}_{t-1} \ominus e) \geq 0.$$

On the other hand, e has the form $e = [_0 a (_j b)_j c$, where $\psi(b) = b$ (since $(_j)_j$ is an inner pair of parentheses in e). As $\psi(e) = \psi([_0 a]_j c) + \psi([_j b)$ and $\psi([_j b) = [_j b$, this shows that $[_j b \leq \mathcal{G}$, and therefore that $\mathcal{G}([_j b) > 0$. Let us now define

$$\tilde{\mathcal{G}}_t = \tilde{\mathcal{G}}_{t-1} \ominus e \oplus [_j b \oplus [_0 a]_j c.$$

Applying ψ to $\tilde{\mathcal{G}}_t$, we see as previously that $\psi(\tilde{\mathcal{G}}_t) = \psi(\tilde{\mathcal{G}}_{t-1}) = \mathcal{G}$. As $\tilde{\mathcal{G}}_k$ contains k pairs of parentheses, and we consume one pair at each step $t > k$, we see that $\tilde{\mathcal{G}}_{2k}$ contains no more parentheses, so that $\tilde{\mathcal{G}}_{2k} \in \mathfrak{G}$ and $\tilde{\mathcal{G}}_{2k} = \psi(\tilde{\mathcal{G}}_{2k}) = \mathcal{G}$. Let us put now $\mathcal{G}_t = \pi(\tilde{\mathcal{G}}_t)$, for $t = k + 1, \dots, 2k$. We see that

$$\mathcal{G}_t = \mathcal{G}_{t-1} \ominus [_0 abc \oplus [_0 a]_j c \oplus [_j b,$$

where $[_0 abc, [_0 a]_j c \in \mathcal{E}$ and $\mathcal{G}([_j b) > 0$, so that $\mathcal{G}_t \in \beta_n(\mathcal{G}_{t-1}, \mathcal{G})$, therefore $\mathcal{G}_k = \mathcal{T}, \dots, \mathcal{G}_{2k} = \mathcal{G}$ is a path of positive probability under the parsing process with reference \mathcal{G} , leading from \mathcal{T} to \mathcal{G} , in other words, $\mathbb{G}_{\mathcal{T}, \mathcal{G}}(\mathcal{G}) > 0$ as required. \square

1.A.3 Convergence to the expectation of a random toric grammar

We will here prove lemma 1.5.1 on page 38. This result is based on the fact that the operation

$$(\mathcal{G}, \mathcal{G}') \mapsto \chi(\mathcal{G} \boxplus \mathcal{G}')$$

is associative.

Let us begin the proof by several definitions and lemmas.

For any grammar $\mathcal{G} \in \mathfrak{G}$ and any pair of indices $p = (p^1, p^2) \in (\mathbb{N} \setminus \{0\})^2$, we will say that p is \mathcal{G} -congruent when there is $a \in S^*$ such that $\mathcal{G}(a]_{p^1})\mathcal{G}(a]_{p^2}) > 0$ or $\mathcal{G}([_{p^1} a)\mathcal{G}([_{p^2} a) > 0$.

Let us define the label map ξ_p as

$$\xi_p(i) = \begin{cases} i, & \text{when } i \notin \{p^1, p^2\}, \\ \min\{p^1, p^2\}, & \text{when } i \in \{p^1, p^2\}. \end{cases}$$

For any sequence $p_1, \dots, p_k \in (\mathbb{N} \setminus \{0\})^{2k}$ of pairs of indices, let us define the label map ξ_{p_1, \dots, p_k} as

$$\xi_{p_1, \dots, p_k} = \xi_{\xi_{p_1, \dots, p_{k-1}}(p_k)} \circ \xi_{p_1, \dots, p_{k-1}},$$

where $f((i, j)) = (f(i), f(j))$, for any $(i, j) \in (\mathbb{N} \setminus \{0\})^2$.

Let us say that $(p_1, \dots, p_k) \in (\mathbb{N} \setminus \{0\})^{2k}$ is \mathcal{G} -congruent, if $\xi_{p_1, \dots, p_{j-1}}(p_j)$ is $\xi_{p_1, \dots, p_{j-1}}(\mathcal{G})$ -congruent for any $j \leq k$, and that it is maximal \mathcal{G} -congruent if it is \mathcal{G} -congruent and any \mathcal{G} -congruent sequence of the form $(p_1, \dots, p_k, p_{k+1})$ is such that

$$\xi_{p_1, \dots, p_k}(p_{k+1}^1) = \xi_{p_1, \dots, p_k}(p_{k+1}^2),$$

or equivalently such that $\xi_{p_1, \dots, p_{k+1}} = \xi_{p_1, \dots, p_k}$.

Lemma 1.A.1

For any sequence $(p_1, \dots, p_\ell) \in (\mathbb{N} \setminus \{0\})^{2\ell}$, for any $k < \ell$,

$$\xi_{p_1, \dots, p_\ell} = \xi_{\xi_{p_1, \dots, p_k}(p_{k+1}, \dots, p_\ell)} \circ \xi_{p_1, \dots, p_k}.$$

PROOF. By induction on ℓ for k fixed. This is true from the definition for $\ell = k+1$. Assuming we have established the lemma for $\ell - 1$, we can write

$$\begin{aligned} \xi_{p_1, \dots, p_\ell} &= \xi_{\xi_{p_1, \dots, p_{\ell-1}}(p_\ell)} \circ \xi_{p_1, \dots, p_{\ell-1}} \\ &= \xi_{\xi_{\xi_{p_1, \dots, p_k}(p_{k+1}, \dots, p_{\ell-1})} \circ \xi_{p_1, \dots, p_k}(p_\ell)} \circ \xi_{\xi_{p_1, \dots, p_k}(p_{k+1}, \dots, p_{\ell-1})} \circ \xi_{p_1, \dots, p_k} \\ &= \xi_{\xi_{p_1, \dots, p_k}(p_{k+1}, \dots, p_\ell)} \circ \xi_{p_1, \dots, p_k}, \quad \square \end{aligned}$$

Lemma 1.A.2

For any permutation σ of $\{1, \dots, k\}$,

$$\xi_{p_1, \dots, p_k} \equiv \xi_{p_{\sigma(1)}, \dots, p_{\sigma(k)}}$$

PROOF. Let us consider the smallest equivalence relation containing the set $\{p_1, \dots, p_k\}$. Let $\bar{\pi}_k : \mathbb{N} \setminus \{0\} \rightarrow \mathcal{C}$ be the corresponding projection of each label to its component. Let us define the label map π_k by $\pi_k(0) = 0$ and

$$\pi_k(i) = \min \bar{\pi}_k(i), \quad i > 0.$$

We are going to prove by induction on k that $\xi_{p_1, \dots, p_k} \equiv \pi_k$. Since π_k is invariant by permutation of the sequence (p_1, \dots, p_k) , this will prove the lemma.

Let us remark now that $\xi_{p_1, \dots, p_k} \equiv \pi_k$ if and only if

$$\xi_{p_1, \dots, p_k}(i) = \xi_{p_1, \dots, p_k}(j) \iff \pi_k(i) = \pi_k(j), \quad i, j > 0.$$

So we are going to prove this equivalence. It is easy to see from the previous lemma that for any integer $m = 1, \dots, k$,

$$\xi_{p_1, \dots, p_k}(p_m^1) = \xi_{p_1, \dots, p_k}(p_m^2). \quad (1.A.1)$$

Indeed,

$$\xi_{p_1, \dots, p_k} = \xi_{\xi_{p_1, \dots, p_m}(p_{m+1}, \dots, p_k)} \circ \xi_{\xi_{p_1, \dots, p_{m-1}}(p_m)} \circ \xi_{p_1, \dots, p_{m-1}},$$

so that, changing p_m for $\xi_{p_1, \dots, p_{m-1}}(p_m)$, we are back to proving the result when $m = k = 1$, where it is obvious from the definitions.

Now, eq. (1.A.1) on the current page and the minimality of π_k implies that

$$\pi_k(i) = \pi_k(j) \implies \xi_{p_1, \dots, p_k}(i) = \xi_{p_1, \dots, p_k}(j), \quad i, j > 0.$$

Let us assume conversely that $\xi_{p_1, \dots, p_k}(i) = \xi_{p_1, \dots, p_k}(j)$ and let

$$m = \min \{ \ell : \xi_{p_1, \dots, p_\ell}(i) = \xi_{p_1, \dots, p_\ell}(j) \}.$$

Since $\xi_{p_1, \dots, p_m} = \xi_{\xi_{p_1, \dots, p_{m-1}}(p_m)} \circ \xi_{p_1, \dots, p_{m-1}}$, we see that necessarily

$$\xi_{p_1, \dots, p_{m-1}}(\{i, j\}) = \xi_{p_1, \dots, p_{m-1}}(\{p_m^1, p_m^2\}),$$

and that this set contains two distinct elements. Exchanging the role of i and j if necessary, we can assume without loss of generality that

$$\xi_{p_1, \dots, p_{m-1}}((i, j)) = \xi_{p_1, \dots, p_{m-1}}(p_m).$$

From the induction hypothesis, this implies that $\pi_{m-1}((i, j)) = \pi_{m-1}(p_m)$. Since the equivalence relation defined by π_{m-1} is a subset of the equivalence relation defined by π_k , this implies that $\pi_k((i, j)) = \pi_k(p_m)$. Since moreover we have that $\pi_k(p_m^1) = \pi_k(p_m^2)$, this implies that $\pi_k(i) = \pi_k(j)$. \square

Lemma 1.A.3

For any $f \in \mathfrak{F}$, any sequence of pairs of positive labels p_1, \dots, p_k , there is a label map $g \in \mathfrak{F}$ such that

$$\xi_{f(p_1, \dots, p_k)} \circ f = g \circ \xi_{p_1, \dots, p_k}.$$

PROOF. We have to prove that

$$\xi_{p_1, \dots, p_k}(i) = \xi_{p_1, \dots, p_k}(j) \implies \xi_{f(p_1, \dots, p_k)} \circ f(i) = \xi_{f(p_1, \dots, p_k)} \circ f(j), \quad i, j > 0.$$

From the proof of the previous lemma, it is enough to check that the right-hand side holds when $(i, j) = p_m$, $m = 1, \dots, k$, which is then obvious. \square

Lemma 1.A.4

If $f \in \mathfrak{F}$ and (p_1, \dots, p_k) is \mathcal{G} -congruent, then the sequence $(f(p_1), \dots, f(p_k))$ is also $f(\mathcal{G})$ -congruent.

PROOF. Assume that for some $a \in S^*$

$$\xi_{p_1, \dots, p_{m-1}}(\mathcal{G})(a]_{\xi_{p_1, \dots, p_{m-1}}(p_m^1)}) > 0.$$

Then, $\xi_{f(p_1, \dots, p_{m-1})} \circ f = g \circ \xi_{p_1, \dots, p_{m-1}}$, and

$$\begin{aligned} \xi_{f(p_1, \dots, p_{m-1})} \circ f(\mathcal{G})(g(a)]_{\xi_{f(p_1, \dots, p_{m-1})} \circ f(p_m^1)}) \\ &= g \circ \xi_{p_1, \dots, p_{m-1}}(\mathcal{G})(g(a)]_{g \circ \xi_{p_1, \dots, p_{m-1}}(p_m^1)}) \\ &= \xi_{p_1, \dots, p_m}(\mathcal{G})(g^{-1} \circ g(a)]_{\xi_{p_1, \dots, p_{m-1}}(p_m^1)}) \\ &\geq \xi_{p_1, \dots, p_{m-1}}(\mathcal{G})(a)]_{\xi_{p_1, \dots, p_{m-1}}(p_m^1)}) > 0. \end{aligned}$$

The same is true when p_m^1 is replaced with p_m^2 and when $a]_{\xi_{p_1, \dots, p_{m-1}}(p_m^1)}$ is replaced with $a]_{\xi_{p_1, \dots, p_{m-1}}(p_m^2)}$.

The lemma is a straightforward consequence of these remarks and the definition of a congruent sequence. \square

Lemma 1.A.5

If (p_1, \dots, p_k) and (q_1, \dots, q_ℓ) are both \mathcal{G} -congruent, then

$$(p_1, \dots, p_k, q_1, \dots, q_\ell)$$

is \mathcal{G} -congruent.

PROOF. According to the previous lemma, we know that $\xi_{p_1, \dots, p_k}(q_1, \dots, q_\ell)$ is $\xi_{p_1, \dots, p_k}(\mathcal{G})$ -congruent. Coming back to the definition this proves that

$$\xi_{\xi_{p_1, \dots, p_k}(q_1, \dots, q_{\ell-1})} \circ \xi_{p_1, \dots, p_k}(q_\ell)$$

is

$$\xi_{\xi_{p_1, \dots, p_k}(q_1, \dots, q_{\ell-1})} \circ \xi_{p_1, \dots, p_k}(\mathcal{G})\text{-congruent.}$$

In lemma 1.A.1 on page 55 we have moreover proved that

$$\xi_{p_1, \dots, p_k, q_1, \dots, q_{\ell-1}} = \xi_{\xi_{p_1, \dots, p_k}(q_1, \dots, q_{\ell-1})} \circ \xi_{p_1, \dots, p_k}.$$

This identity applied to the above statement shows that $(p_1, \dots, p_k, q_1, \dots, q_\ell)$ satisfies the definition of a \mathcal{G} -congruent sequence. \square

Proposition 1.A.6

If (p_1, \dots, p_k) and (q_1, \dots, q_ℓ) are both maximal \mathcal{G} -congruent, then

$$\xi_{p_1, \dots, p_k}(\mathcal{G}) \equiv \xi_{q_1, \dots, q_\ell}(\mathcal{G}) \equiv \chi(\mathcal{G}).$$

PROOF. From the previous lemma, $(p_1, \dots, p_k, q_1, \dots, q_\ell)$ is \mathcal{G} -congruent. Since p is maximal, $\xi_{p_1, \dots, p_k, q_1, \dots, q_\ell} = \xi_{p_1, \dots, p_k}$. In the same way $\xi_{q_1, \dots, q_\ell, p_1, \dots, p_k} = \xi_{q_1, \dots, q_\ell}$. We have seen moreover in a previous lemma that

$$\xi_{p_1, \dots, p_k, q_1, \dots, q_\ell} \equiv \xi_{q_1, \dots, q_\ell, p_1, \dots, p_k}.$$

This proves that $\xi_{p_1, \dots, p_k} \equiv \xi_{q_1, \dots, q_\ell}$.

We see from the definition of χ (see definition 1.3.10 on page 34) that there is some maximal \mathcal{G} -congruent sequence r_1, \dots, r_m such that $\chi(\mathcal{G}) = \xi_{r_1, \dots, r_m}(\mathcal{G})$. Therefore

$$\chi(\mathcal{G}) \equiv \xi_{p_1, \dots, p_k}(\mathcal{G}) \equiv \xi_{q_1, \dots, q_\ell}(\mathcal{G}). \quad \square$$

Proposition 1.A.7

For any $\mathcal{G}, \mathcal{G}' \in \mathfrak{G}$,

$$\chi(\chi(\mathcal{G}) \boxplus \mathcal{G}') = \chi(\mathcal{G} \boxplus \mathcal{G}').$$

Consequently, for any $\mathcal{G}, \mathcal{G}', \mathcal{G}'' \in \mathfrak{G}$,

$$\chi(\chi(\mathcal{G} \boxplus \mathcal{G}') \boxplus \mathcal{G}'') = \chi(\mathcal{G} \boxplus \mathcal{G}' \boxplus \mathcal{G}'').$$

PROOF. Let us assume that \mathcal{G} , \mathcal{G}' and $\chi(\mathcal{G})$ use disjoint label sets, so that

$$\begin{aligned}\chi(\chi(\mathcal{G}) \boxplus \mathcal{G}') &\equiv \chi(\chi(\mathcal{G}) + \mathcal{G}'), \\ \chi(\mathcal{G} \boxplus \mathcal{G}') &\equiv \chi(\mathcal{G} + \mathcal{G}').\end{aligned}$$

Let p_1, \dots, p_k be some maximal \mathcal{G} -congruent sequence. It is then obviously also $(\mathcal{G} + \mathcal{G}')$ -congruent, and since label sets are disjoint,

$$\xi_{p_1, \dots, p_k}(\mathcal{G}) + \mathcal{G}' = \xi_{p_1, \dots, p_k}(\mathcal{G} + \mathcal{G}').$$

Let us continue the sequence p_1, \dots, p_k to form a maximal $(\mathcal{G} + \mathcal{G}')$ -congruent sequence p_1, \dots, p_ℓ . Let (q_{k+1}, \dots, q_ℓ) be defined as

$$q_m = \xi_{p_1, \dots, p_k}(p_{k+m}).$$

It now follows from the definitions that the sequence (q_{k+1}, \dots, q_ℓ) is a maximal $\xi_{p_1, \dots, p_k}(\mathcal{G} + \mathcal{G}')$ -congruent sequence, and a maximal $(\xi_{p_1, \dots, p_k}(\mathcal{G}) + \mathcal{G}')$ -congruent sequence by consequence. Consequently

$$\begin{aligned}\chi(\chi(\mathcal{G}) + \mathcal{G}') &\equiv \xi_{q_{k+1}, \dots, q_\ell}(\xi_{p_1, \dots, p_k}(\mathcal{G}) + \mathcal{G}') \\ &= \xi_{q_{k+1}, \dots, q_\ell} \circ \xi_{p_1, \dots, p_k}(\mathcal{G} + \mathcal{G}') = \xi_{\xi_{p_1, \dots, p_k}(p_{k+1}, \dots, p_\ell)} \circ \xi_{p_1, \dots, p_k}(\mathcal{G} + \mathcal{G}') \\ &= \xi_{p_1, \dots, p_\ell}(\mathcal{G} + \mathcal{G}') \equiv \chi(\mathcal{G} + \mathcal{G}'),\end{aligned}$$

proving the proposition. \square

PROOF OF LEMMA 1.5.1 ON PAGE 38. Let π be the projection of \mathfrak{G} on \mathfrak{G}/\equiv .

From the law of large numbers, we have that, for all $\mathcal{G} \in \mathfrak{G}$,

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}(G_i \equiv \mathcal{G}) \xrightarrow[n \rightarrow \infty]{} \mathbb{G}(\pi(\mathcal{G})).$$

Let us now remark that $\bigsqcup_{i=1}^n n^{-1}G_i = \bigsqcup_{\overline{\mathcal{G}} \in \mathfrak{G}/\equiv} \bigsqcup_{G_i \in \overline{\mathcal{G}}} n^{-1}G_i$. Thus

$$\begin{aligned}\frac{1}{n} \chi \left(\bigsqcup_{i=1}^n G_i \right) &= \chi \left(\bigsqcup_{\overline{\mathcal{G}} \in \mathfrak{G}/\equiv} \chi \left(\bigsqcup_{G_i \in \overline{\mathcal{G}}} n^{-1}G_i \right) \right) \\ &= \chi \left(\bigsqcup_{\overline{\mathcal{G}} \in \mathfrak{G}/\equiv} \left(\sum_{i=1}^n n^{-1} \mathbf{1}(G_i \in \overline{\mathcal{G}}) \right) \chi(\overline{\mathcal{G}}) \right) \\ &= \chi \left(\bigsqcup_{\overline{\mathcal{G}} \in \mathfrak{G}/\equiv} \left(\sum_{i=1}^n n^{-1} \mathbf{1}(G_i \in \overline{\mathcal{G}}) \right) \overline{\mathcal{G}} \right).\end{aligned}$$

We used here proposition 1.A.7 on page 58, together with the fact that for any numbers $a, b \in \mathbb{R}_+$,

$$\chi[(a\mathcal{G}) \boxplus (b\mathcal{G})] = (a+b)\chi(\mathcal{G}),$$

which comes from the following reasoning: Suppose that

$$\{1, \dots, d\} = \{i; \mathcal{G}([_i S^*] > 0)\},$$

and let $p_i = (2i, 2i-1)$. Since each p_i is $(a\mathcal{G}) \boxplus (b\mathcal{G})$ -congruent, (p_1, \dots, p_d) is also $(a\mathcal{G}) \boxplus (b\mathcal{G})$ -congruent, from lemma 1.A.5 on page 57. It is quite straightforward to see that

$$\xi_{p_1, \dots, p_d}[(a\mathcal{G}) \boxplus (b\mathcal{G})] \equiv (a+b)\mathcal{G}.$$

This implies that

$$\chi[(a\mathcal{G}) \boxplus (b\mathcal{G})] = \chi \circ \xi_{p_1, \dots, p_d}[(a\mathcal{G}) \boxplus (b\mathcal{G})] = \chi[(a+b)\mathcal{G}] = (a+b)\chi(\mathcal{G}).$$

To take the limit inside χ , we need to prove that χ is continuous in a suitable sense. Actually, $\mathcal{G} \mapsto \chi(\mathcal{G})$ is continuous on sets of fixed support, and this is what is required to conclude.

Indeed, for any sequence (\mathcal{G}_i) with fixed support for n large enough, there is a fixed label map f (depending on the support) such that for n large enough $\chi(\mathcal{G}_i) = f(\mathcal{G}_i)$, and the result follows from the fact that $\mathcal{G} \mapsto f(\mathcal{G})$ is continuous; since $f(\mathcal{G})(A) = \mathcal{G}(f^{-1}(A))$.

Consequently

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \chi \left(\bigoplus_{i=1}^n G_i \right) &= \chi \left(\bigoplus_{\mathcal{G} \in \mathfrak{G}/\equiv} \lim_{n \rightarrow \infty} \left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}(G_i \in \overline{\mathcal{G}}) \right) \overline{\mathcal{G}} \right) \\ &= \chi \left(\bigoplus_{\mathcal{G} \in \mathfrak{G}/\equiv} G(\overline{\mathcal{G}}) \overline{\mathcal{G}} \right) = \chi \left(\bigoplus_{\mathcal{G} \in \mathfrak{G}/\equiv} G(\overline{\mathcal{G}}) \chi(\overline{\mathcal{G}}) \right) \\ &= \chi \left(\bigoplus_{\mathcal{G} \in \mathfrak{G}/\equiv} \bigoplus_{\mathcal{G} \in \overline{\mathcal{G}}} G(\mathcal{G}) \chi(\mathcal{G}) \right) = \chi \left(\bigoplus_{\mathcal{G} \in \mathfrak{G}} G(\mathcal{G}) \chi(\mathcal{G}) \right) \\ &= \chi \left(\bigoplus_{\mathcal{G} \in \mathfrak{G}} G(\mathcal{G}) \mathcal{G} \right) = \int \mathcal{G} dG(\mathcal{G}). \quad \square \end{aligned}$$

1.B Language produced by a toric grammar

In this appendix, we make a deterministic study of the language produced by a toric grammar $\mathcal{G} \in \beta^*(\mathfrak{T})$. More precisely, we are interested in the support of the distribution $\mathbb{T}_{\mathcal{G}}$ of the final state of the production process.

Lemma 1.B.1

Let $\mathcal{T} \in \mathfrak{T}$ be some text and $\mathcal{G} \in \beta^*(\mathcal{T})$ be some grammar obtained by splitting this text a finite number of times. The number of splits performed can be read in \mathcal{G} and is equal to

$$n = \sum_{i=1}^{+\infty} \mathcal{G} \left(]_i S^* \right).$$

Let us put $\bar{\alpha}(\mathcal{G}) = \alpha^n(\mathcal{G})$. Then, $\mathcal{T} \in \bar{\alpha}(\mathcal{G}) \subset \mathfrak{T}$, moreover $\bar{\alpha}(\mathcal{G}) = \text{supp}(\mathbb{T}_{\mathcal{G}})$.

PROOF. The grammar \mathcal{G} is obtained by making a succession of splits. Each of those splits add one $[_i$ and one $]_i$ to the grammar, whereas in the original text there are no $[_i$ nor $]_i$, except for the $[_0$ at the beginning of each sentence. Since application of an element of \mathfrak{F} does not change the number of such symbols, they may be used to count the number of splits performed.

Let us take then a sequence of toric grammars $\mathcal{T} = \mathcal{G}_0, \dots, \mathcal{G}_n = \mathcal{G}$, such that $\mathcal{G}_k \in \beta(\mathcal{G}_{k-1})$. From lemma 1.3.2 on page 28, there is a sequence $f_1, \dots, f_n \in \mathfrak{F}$ such that $f_k(\mathcal{G}_{k-1}) \in \alpha(\mathcal{G}_k)$. Let us prove by induction that for any $k = 0, \dots, n$,

$$f_k \circ \dots \circ f_1(\mathcal{T}) \in \alpha^k(\mathcal{G}_k).$$

Indeed, this is true for $k = 0$, since $\mathcal{G}_0 = \mathcal{T}$. Moreover, assuming that the assertion holds for $k - 1$, we deduce that

$$f_k \circ \dots \circ f_1(\mathcal{T}) \in f_k \left(\alpha^{k-1}(\mathcal{G}_{k-1}) \right) \subset \alpha^{k-1} \left(f_k(\mathcal{G}_{k-1}) \right) \subset \alpha^k(\mathcal{G}_k).$$

showing that if the assertion holds for k , it also holds for $k + 1$. For $k = n$, we obtain that

$$f_n \circ \dots \circ f_1(\mathcal{T}) \in \alpha^n(\mathcal{G}_n).$$

As $f_n \circ \dots \circ f_1(\mathcal{T}) = \mathcal{T}$, since \mathcal{T} is a text, and $\mathcal{G}_n = \mathcal{G}$, we get that $\mathcal{T} \in \alpha^n(\mathcal{G})$.

Let us consider now $\mathcal{G}' \in \bar{\alpha}(\mathcal{G})$. Let $(\mathcal{G} = \mathcal{G}_0, \dots, \mathcal{G}_n = \mathcal{G}')$ the chain of grammars leading to \mathcal{G}' . Then for any $k = 0, \dots, n$,

$$\sum_{i=1}^{+\infty} \mathcal{G}_k \left(]_i S^* \right) = n - k,$$

since $\mathcal{G}_k \in \alpha(\mathcal{G}_{k-1})$ and each merge takes away one $[_i$ and one $]_i$. This implies that $\sum_{i=1}^{+\infty} \mathcal{G}' \left(]_i S^* \right) = 0$, and thus $\mathcal{G}' \in \mathfrak{T}$.

Note that, as remarked above, repeated merges may create elements of the type $[_i a]_i b$. However, this will not happen if n successful merges can be performed. Indeed in the case when expressions of the form $[_i a]_i b$ remain unmatched during the merge process, we will get $\alpha(\mathcal{G}_k) = \emptyset$ for some $k < n$. \square

Chapter 2

Markov substitute Sets

2.1 Presentation of the model

2.1.1 Motivation

In the previous chapter, we constructed a split and merge process that rearranged texts according to a known reference grammar. We saw that when this reference grammar is known, the distribution of texts can be estimated without error via a Monte-Carlo simulation.

The question of building such a reference grammar \mathcal{R} along justified statistical principles was left open. The model of toric grammars, as we saw, is based on the assumption that any elements ${}_i a$ and ${}_i b$ in the grammar can be substituted in any sentence, while keeping grammaticality. The goal of this chapter is to propose some ways to define properly this kind of property, and to test it.

2.1.2 What is a Markov substitute set

We will work for a while with a probability distribution of sentences using the dictionary of words D , leaving the notion of text aside. We will use the notation $D^* = \bigcup_{j=0}^{\infty} D^j = \{\epsilon\} \cup D^+$, where ϵ stands for the empty string. Let $S \in D^+$ be a random sentence drawn from the language distribution and S_i , $1 \leq i \leq n$, a statistical sample made of n independent copies of S , our observed sample.

Given a right and left context $x = (x_1, x_2) \in D^* \times D^*$, and an expression $y \in D^+$ (that is any non void finite string of words), it will be useful to introduce the insertion operator α as

$$\alpha(x, y) = x_1 y x_2 \in D^+,$$

the concatenation of the strings of words x_1 , y , and x_2 .

Let us remark that for any $x \in (D^*)^2$, the map $y \mapsto \alpha(x, y)$ is one to one, whereas the map $x \mapsto \alpha(x, y)$ is not, since for instance $\alpha((y, y'), y) = \alpha((\emptyset, yy'), y)$.

Definition 2.1.1

We will say that the set $B \subset D^+$ is a Markov substitute set of S when

$$\sum_{x \in (D^*)^2} \mathbb{P}(S = \alpha(x, y)) > 0, \quad y \in B, \quad (2.1.1)$$

and when there exists a probability measure $q_B \in \mathcal{M}_+^1(B)$ on B , called the substitute measure of B , such that for any $x \in (D^*)^2$ and any $y \in B$,

$$\mathbb{P}_S[\alpha(x, y)] = \mathbb{P}_S[\alpha(x, B)]q_B(y), \quad (2.1.2)$$

where $\alpha(x, B) = \{\alpha(x, y), y \in B\}$.

For any context $x \in (D^*)^2$, such that $\mathbb{P}_S[\alpha(x, B)] > 0$, the substitute measure is equal to

$$q_B(y) = \frac{\mathbb{P}_S[\alpha(x, y)]}{\mathbb{P}_S[\alpha(x, B)]}.$$

Therefore $\text{supp}(q_B) = B$. An equivalent characterization of the Markov substitute property is that B is a Markov substitute set if and only if it satisfies eq. (2.1.1) on this page and for any $y, y' \in B$, any $x, x' \in (D^*)^2$,

$$\mathbb{P}_S[\alpha(x, y)] \mathbb{P}_S[\alpha(x', y')] = \mathbb{P}_S[\alpha(x, y')] \mathbb{P}_S[\alpha(x', y)]. \quad (2.1.3)$$

To justify calling this a Markov substitute set, we can make a link with Markov chains. Indeed, in the case when $S = (Z_1, \dots, Z_L)$ is a finite length time homogeneous Markov chain, where $L > 2$ is a non random integer, and where $\text{supp}(\mathbb{P}_{Z_1}) = D$, for any $w_1, w_3 \in D$, the set

$$B_{w_1, w_3} = \{(w_1, w_2, w_3) \in D^3, \mathbb{P}(Z_2 = w_2, Z_3 = w_3 | Z_1 = w_1) > 0\},$$

is a Markov substitute set. In the same way, let us consider a time homogeneous Markov chain $(Z_n, n \in \mathbb{N} \setminus \{0\})$, a subset $C \subset D$ and let us take for its stopping time $\tau = \inf\{k > 1, Z_k \in C\}$. Let us assume that $\mathbb{P}(\tau < \infty | Z_1 = z) = 1$ for any $z \in D$. In this situation, the chain $S = Z_{1:\tau} \stackrel{\text{def}}{=} (Z_1, \dots, Z_\tau)$, stopped at time τ , is a process in D^+ . For any $w, w' \in D$, such that

$$\sum_{1 \leq j < k} \mathbb{P}(Z_j = w, Z_k = w', j \leq \tau) > 0$$

the set

$$B_{w, w'} = \left\{ w_{1:k} \in D^k, k > 1, w_1 = w, w_k = w', w_j \in D \setminus C, 1 < j < k, \right. \\ \left. \mathbb{P}_{Z_{2:k} | Z_1 = w_1}(w_{2:k}) > 0 \right\}$$

is a Markov substitute set.

Indeed, in this case, there are $j < k \leq t$ such that

$$\mathbb{P}(Z_j = w, Z_k = w', \tau = t) > 0,$$

so that there is a chain $y_{1:t} \in D^t$ such that $y_{2:t-1} \in (D \setminus C)^{t-2}$, $y_j = w$, $y_k = w'$, $y_t \in C$ and

$$\mathbb{P}(Z_{1:\tau} = y_{1:t}) > 0.$$

Since

$$\begin{aligned} \mathbb{P}(Z_{1:\tau} = y_{1:t}) &= \mathbb{P}(Z_{1:t} = y_{1:t}) \\ &= \mathbb{P}(Z_{1:j} = y_{1:j}) \mathbb{P}(Z_{2:k-j+1} = y_{j+1:k} \mid Z_1 = w) \mathbb{P}(Z_{2:t-k+1} = y_{k+1:t} \mid Z_1 = w'), \end{aligned}$$

this implies that $y_{j:k} \in B_{w,w'}$, and that for any $w_{1:\ell} \in B_{w,w'}$,

$$\begin{aligned} \mathbb{P}[S = \alpha((y_{1:j-1}, y_{k+1:t}), w_{1,\ell})] \\ = \mathbb{P}(Z_{1:j} = y_{1:j}) \mathbb{P}(Z_{2:\ell} = w_{2:\ell} \mid Z_1 = w) \mathbb{P}(Z_{2:t-k+1} = y_{k+1:t} \mid Z_1 = w') > 0, \end{aligned}$$

so that the set $B_{w,w'}$ satisfies eq. (2.1.1) on the preceding page.

Moreover for any $(x_{1:j}, x_{j+1:m}) \in (D^*)^2$,

$$\begin{aligned} \mathbb{P}_S(\alpha(x, w)) &= \mathbb{P}(Z_{1:j+1} = \gamma(x_{1:j}, w), \tau > j) \mathbb{P}(Z_{2:k} = w_{2:\ell} \mid Z_1 = w, \tau \geq \ell) \\ &\quad \times \mathbb{P}(Z_{2,m-j+1} = x_{j+1,m} \mid Z_1 = w', \tau \geq m - j + 1). \end{aligned}$$

(where γ is the concatenation operator) so that $B_{w,w'}$ satisfies also eq. (2.1.2) on the facing page, where

$$q_{B_{w,w'}}(w_{1:\ell}) = \frac{\mathbb{P}(Z_{2,k} = w_{2,\ell} \mid Z_1 = w, \tau \geq \ell)}{\sum_{y \in (D \setminus C)^*} \mathbb{P}(Z_{2:\ell(y)+1} = \gamma(w, y) \mid Z_1 = w)}.$$

We gave these examples to show the connections with the usual Markov property and to provide an example of Markov substitute sets containing strings of variable lengths. We will see other examples of Markov substitute models a little later.

2.1.3 Weak Markov substitute sets

In natural language analysis, we may be mainly interested in the support of the distribution of S . If this is the case, we can weaken our definition of Markov substitute set to deal only with support shapes.

Definition 2.1.2

We will say that the set $B \subset D^+$ is a weak Markov substitute set of S when it satisfies eq. (2.1.1) on page 64 and when for any context $x \in (D^*)^2$ such that $\mathbb{P}_S[\alpha(x, B)] > 0$, $\mathbb{P}_S[\alpha(x, y)] > 0$ for any $y \in B$.

2.1.4 Basic properties of Markov substitute sets**Proposition 2.1.1**

Any one point set $\{y\}$, where $y \in D^+$ satisfies eq. (2.1.1) on page 64, is a Markov substitute set, whose substitute measure is the Dirac mass at y .

A subset of a Markov substitute set is itself a Markov substitute set.

If B and C are Markov substitute sets such that $B \cap C \neq \emptyset$, then $B \cup C$ is also a Markov substitute set.

As such, a set B is a Markov substitute set if and only if, for any $y, y' \in B$, the pair $\{y, y'\}$ is a Markov substitute pair.

Less restrictively, a set B is a Markov substitute set if and only if there is a connected undirected spanning graph $\mathcal{G} \subset B^2$ such that for any $(y, y') \in \mathcal{G}$, $\{y, y'\}$ is a Markov substitute pair.

This is obvious from the characterization by eq. (2.1.2) on page 64 for the first point and eq. (2.1.3) on page 64 for the others.

These properties lead us to define the relation

$$y \sim_S y' \iff \{y, y'\} \text{ is a Markov substitute pair.}$$

This is an equivalence relation and D^+/\sim_S forms a partition of D^+ into maximal Markov substitute sets.

Proposition 2.1.2

If B is a Markov substitute set, then for any $x \in (D^*)^2$, $\alpha(x, B)$ is also a Markov substitute set, as soon as it satisfies

$$\sum_{z \in (D^*)^2} \mathbb{P}[S \in \alpha(z, \alpha(x, B))] > 0.$$

PROOF. From the definition of the Markov substitute property, we see that for any $z \in (D^*)^2$, any $y \in B$,

$$\mathbb{P}_S[\alpha(z, \alpha(x, y))] = \mathbb{P}_S[\alpha(z_1 x_1, x_2 z_2, y)] = \mathbb{P}_S[\alpha(z, \alpha(x, B))] q_B(y),$$

so that $\alpha(x, B)$ is a Markov substitute set and $q_{\alpha(x, B)}(\alpha(x, y)) = q_B(y)$. \square

Proposition 2.1.3

If B_j , $1 \leq j \leq k$ are Markov substitute sets (including possibly some one point sets), then

$$B = \gamma(B_1, \dots, B_k) \stackrel{\text{def}}{=} \{y = y_1 \dots y_k, y_j \in B_j, 1 \leq j \leq \ell\},$$

is also a Markov substitute set, as soon as it satisfies

$$\sum_{x \in (D^*)^2} \mathbb{P}(\alpha(x, B)) > 0.$$

PROOF. For any $y_{1:k}, y'_{1:k} \in \gamma(B_1, \dots, B_k)$,

$$\{y'_{1:j-1}y_{j:k}, y'_{1:j}y_{j+1:k}\} = \alpha(y'_{1:j-1}, y_{j+1:k}, \{y_j, y'_j\})$$

and is therefore a Markov substitute pair, so that $y_{1:k}$ and $y'_{1:k}$ are connected by the equivalence relation \sim_S defined above and themselves form a Markov substitute pair. \square

The same properties are true for weak Markov substitute sets. They show the key role played by Markov substitute pairs.

2.1.5 Interpretation in terms of random parsing

Let us consider some finite subset $B \in D^+$ satisfying eq. (2.1.1) on page 64. We define the set of B splits of s , for any sentence $s \in D^+$ as

$$\mathcal{S}(s, B) = \{(x, y), x \in (D^*)^2, y \in B, \alpha(x, y) = s\}.$$

Let us consider some conditional probability kernel

$$(\pi(s, x, y), s \in D^+, x \in (D^*)^2, y \in (B \cup \{\epsilon\}))$$

such that $\pi(s, \cdot) \in \mathcal{M}_+^1((D^*)^2 \times D^*)$, and such that

$$\mathcal{S}(s, B) \subset \text{supp}(\pi(s, \cdot)) \subset \mathcal{S}(s, B) \cup \{(s, \epsilon), \epsilon\}. \quad (2.1.4)$$

Introducing the split counts

$$c(s) = |\mathcal{S}(s, B)|,$$

we can for instance take

$$\pi(s, x, y) = \begin{cases} c(s)^{-1}, & (x, y) \in \mathcal{S}(s, B), \\ 1 - \sum_{(x', y') \in \mathcal{S}(s, B)} \pi(s, x', y'), & x = (s, \epsilon), y = \epsilon. \end{cases}$$

Let us define the random B -parse X, Y of the random sentence S on the same probability space by its conditional distribution

$$\mathbb{P}_{X, Y|S=s}(x, y) = \begin{cases} \min_{y' \in B} \pi(\alpha(x, y'), x, y'), & (x, y) \in \mathcal{S}(s, B), \\ 1 - \sum_{(x', y') \in \mathcal{S}(s, B)} \mathbb{P}_{X, Y|S=s}(x', y'), & x = (s, \epsilon), y = \epsilon. \end{cases}$$

Let us remark that the random B -parse can be simulated using a rejection method. Indeed, let us put

$$w(x, y) = \min_{y' \in B} \frac{\pi(\alpha(x, y'), x, y')}{\pi(\alpha(x, y), x, y)}.$$

We can first draw the pair (X', Y') according to the conditional probability distribution $\mathbb{P}_{X', Y'|S} = \pi(S, \cdot)$, and then draw (X, Y) according to the distribution

$$\mathbb{P}_{(X, Y)|S, X', Y'} = w(X', Y')\delta_{(X', Y')} + (1 - w(X', Y'))\delta_{((S, \epsilon), \epsilon)}.$$

Proposition 2.1.4

The set B is a Markov substitute set for S if and only if the random B -parse (X, Y) satisfies one of the three following conditions :

$$\mathbb{P}_{X, Y|Y \in B} = \mathbb{P}_{X|Y \in B} \otimes \mathbb{P}_{Y|Y \in B}, \quad (2.1.5)$$

$$\mathbb{P}_{X|Y=y} = \mathbb{P}_{X|Y=y'}, \quad y, y' \in B, \quad (2.1.6)$$

$$\mathbb{P}_{Y|X=x, Y \in B} = \mathbb{P}_{Y|Y \in B}, \quad x \in (D^*)^2, \mathbb{P}_{X|Y \in B}(x) > 0. \quad (2.1.7)$$

Moreover, when B is a Markov substitute set, the substitute measure q_B is equal to $\mathbb{P}_{Y|Y \in B}$.

PROOF. The three properties are quite obviously equivalent, and the fact that, for any $x \in (D^*)^2$ and any $y \in B$,

$$\begin{aligned} \mathbb{P}_{X, Y}(x, y) &= \min_{y' \in B} \pi(\alpha(x, y'), x, y') \mathbb{P}_S[\alpha(x, y)] \\ &= \min_{y' \in B} \pi(\alpha(x, y'), x, y') \mathbb{P}_S[\alpha(x, B)] q_B(y), \end{aligned}$$

ends the proof. □

2.2 Invariant dynamics

Markov substitute sets can be used to define various reversible (and consequently invariant) dynamics.

Moreover, when \mathbb{P}_S has a sufficient number of Markov substitute sets, we will be able to define an irreducible (in other words ergodic) dynamics that will have a single invariant measure.

The idea is to replace an element from a Markov substitute set by another one. Such a transformation stays necessarily in the support of the language, and, if the substitution weights are properly related to the substitute measure, this dynamics can be invariant, and even reversible. We will see here simple invariant dynamics, both on single sentences and on texts, that will only modify one element of a Markov substitute set at a time. More evolved invariant dynamics will be defined and studied in chapter 3.

2.2.1 Metropolis invariant dynamics on sentences

We will here have a first go at describing invariant dynamics on sentences, related to the Markov substitute property.

Let us consider a family $B_j, 1 \leq j \leq t$ of Markov substitute sets, and let us assume that we know for each B_i the substitute measure $q_{\bar{B}_i}$ of a Markov substitute set \bar{B}_i containing B_i . We do not need necessarily to know q_{B_i} itself, since we can deduce it from the relation $q_{B_i}(y) = q_{\bar{B}_i}(y)/q_{\bar{B}_i}(B_i)\mathbf{1}(y \in B_i)$.

We will locate a member of a Markov substitute set in the sentence and replace it by another element of the same Markov substitute set, according to the relevant substitute measure. Let us remark that we accordingly will only use the substitute measures $q_{\bar{B}_i}$ such that

$$\sum_{x \in (D^*)^2} \mathbb{E}(\pi(S, x, \bar{B}_i, i)) > 0.$$

Definition 2.2.1

A splitting kernel is a kernel

$$\pi : D^+ \rightarrow \mathcal{M}_+^1((D^*)^2 \times D^* \times \{1, \dots, t\})$$

such that for any $s \in D^+$,

$$\pi(s, x, y, j) > 0 \implies y \in B_j \cup \{\epsilon\} \text{ (the empty string), and } \alpha(x, y) = s.$$

For now, we will suppose that we have access to such kernels. Actual constructions will be presented in chapter 3.

Definition 2.2.2

Given a splitting kernel π , we can define the substitute dynamics

$$\begin{aligned} k_\pi(s, s') &= \sum_{x \in (D^*)^2, (y, y') \in (D^*)^2, j} \pi(s, x, y, j) q_{\overline{B}_j}(y') \left(\frac{\pi(s', x, y', j)}{\pi(s, x, y, j)} \wedge 1 \right), \quad s \neq s', \\ k_\pi(s, s) &= 1 - \sum_{s' \in D^+ \setminus \{s\}} k_\pi(s, s'). \end{aligned} \quad (2.2.1)$$

The simulation of S' such that $\mathbb{P}_{S'|S=s}(s') = k_\pi(s, s')$ can be implemented in the following way. First draw (X, Y, j) such that $\mathbb{P}_{X, Y, j|S} = \pi$, then draw Y' according to $q_{\overline{B}_j}$ and set

$$S' = \begin{cases} \alpha(X, Y') & \text{with probability } \left(\frac{\pi(\alpha(X, Y'), X, Y', j)}{\pi(S, X, Y, j)} \wedge 1 \right), \\ S & \text{otherwise.} \end{cases}$$

Proposition 2.2.1

Any substitute dynamics defined as in definition 2.2.2 on this page is reversible (and therefore invariant) with respect to \mathbb{P}_S .

PROOF. We can see that k_π is reversible by writing

$$\begin{aligned} \mathbb{P}_S(s)k(s, s') &= \sum_{x \in (D^*)^2, (y, y') \in (D^*)^2, j} \mathbb{P}_S(\alpha(x, \overline{B}_j)) \\ &\quad \times q_{\overline{B}_j}(y)q_{\overline{B}_j}(y') [\pi(s, x, y, j) \wedge \pi(s', x, y', j)], \end{aligned}$$

which is obviously a symmetric expression in s and s' . \square

Remark that, while all these operations are defined on single sentences, it is quite straightforward to extend them to texts, that is, an empirical distribution on sentences, by applying them to each sentence of the text.

2.2.2 Reflecting a reversible dynamics on the boundary of a finite domain

We can reflect a \mathbb{P}_S reversible dynamics k on the boundary of any finite domain $\mathcal{D} \subset D^+$, by defining the reflected dynamics as

$$k_{\mathcal{D}}(s, s') = \begin{cases} k(s, s'), & \text{when } s, s' \in \mathcal{D}, s \neq s', \\ 0, & \text{when } s \in \mathcal{D}, s' \notin \mathcal{D}, \\ k(s, s) + \sum_{s'' \in D^+ \setminus \mathcal{D}} k(s, s''), & \text{when } s = s' \in \mathcal{D}, \\ \delta_s(s'), & s \notin \mathcal{D}. \end{cases}$$

It is immediate to see that $k_{\mathcal{D}}$ is still reversible with respect to \mathbb{P}_S .

2.2.3 Compound dynamics

Definition 2.2.3

Let us consider a finite family of \mathbb{P}_S reversible dynamics, k_1, \dots, k_t and a probability distribution $\xi \in \mathcal{M}_+^1(\{1, \dots, t\})$ on the t first integers. We define the compound substitute dynamics k_ξ as the reversible kernel

$$k_\xi(s, s') = \sum_{j=1}^t \xi(j) k_j(s, s').$$

This compound kernel can be easily implemented by drawing at random an index j according to the distribution ξ and then drawing s' according to $k_{B_j}(s, s')$. It is obviously reversible with respect to \mathbb{P}_S .

Definition 2.2.4

We can then define the accelerated compound dynamics

$$\begin{aligned} \bar{k}_\xi(s, s') &= \left(\sum_{s'' \in D^+ \setminus \{s\}} k_\xi(s, s'') \vee \sum_{s'' \in D^+ \setminus \{s'\}} k_\xi(s', s'') \right)^{-1} k_\xi(s, s'), \quad s \neq s' \in D^+, \\ \bar{k}_\xi(s, s) &= 1 - \sum_{s' \in D^+ \setminus \{s\}} \bar{k}_\xi(s, s'), \quad s \in D^+. \end{aligned}$$

Let us remark that for states $s \in D^+$ such that

$$\sum_{s'' \in D^+ \setminus \{s\}} k_\xi(s, s'') = \max_{s' \in D^+, k(s, s') > 0} \sum_{s'' \in D^+ \setminus \{s'\}} k_\xi(s', s''),$$

the accelerated dynamics jumps with probability one, meaning that $\bar{k}_\xi(s, s) = 0$.

2.2.4 A simple example of recursive structure

In this subsection, we will present a small example, showing that the Markov substitute property is suitable to model recursive structures.

Let us consider on $D = \{a, b\}$ the language distribution

$$\mathbb{P}_S(ab^n) = \left(\frac{1}{2}\right)^{n+1}.$$

It is easy to check that $\{a, ab\}$ is a Markov substitute set for this language, with

$$q_B = \frac{2}{3}\delta_a + \frac{1}{3}\delta_{ab}.$$

A random B -parse as defined in section 2.1.5 on page 67 can be defined in this case as

$$\mathbb{P}(X = (\epsilon, b^n), Y = a | S = ab^n) = \mathbb{P}(X = (\epsilon, b^{n-1}), Y = ab | S = ab^n) = \frac{1}{2}.$$

The invariant dynamics in this setup is to either remove one b with probability $\frac{1}{3}$ or add one with probability $\frac{1}{6}$, as shown in fig. 2.1 on this page.

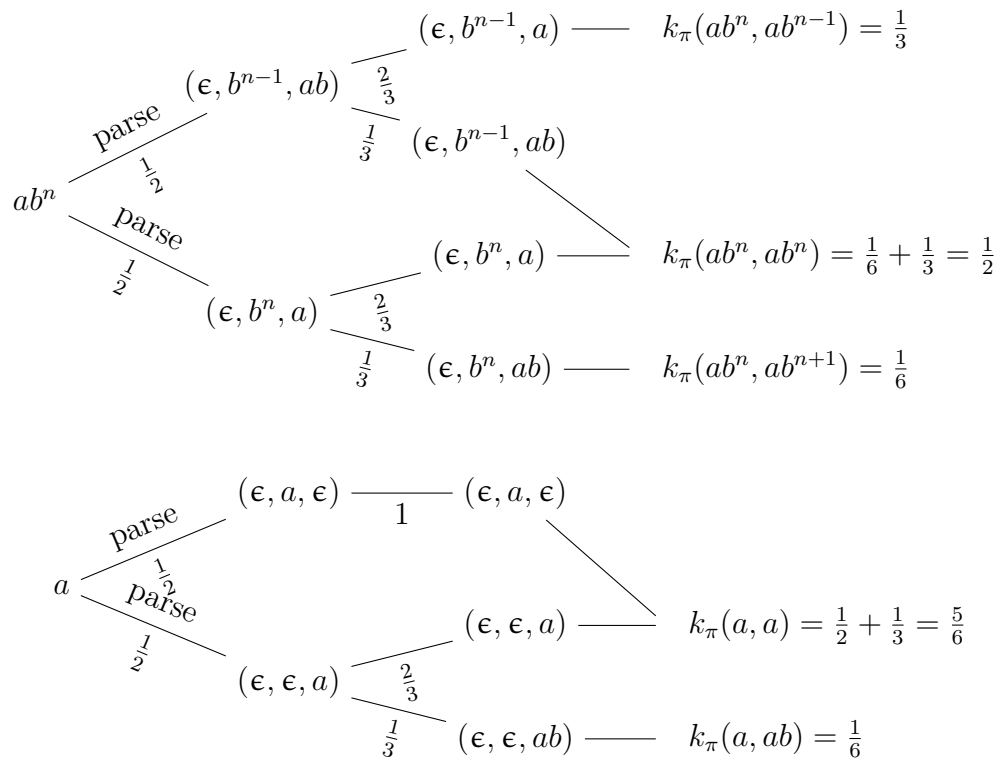


Figure 2.1: Invariant dynamics for ab^n .

It is quite straightforward to see that this dynamics is indeed invariant.

It is important to realize on this basic example that substitute dynamics have the power to generate languages with a recursive structure.

Here, the full language $ab^n, n \in \mathbb{N}$ is generated from any starting point ab^k , because the invariant dynamics is ergodic (since it has a single communicating class). A statistical sample larger than a single sentence would still be necessary to discover that $\{a, ab\}$ is a Markov substitute pair.

On the other hand, we can define invariant dynamics implementing any set of rewriting rules, generating languages with arbitrarily complex recursive structures.

2.2.5 Crossing-over reversible dynamics on texts

The previous method required that we know the substitute measures $q_{\overline{B}_i}$, which must be estimated. We will now present another way to build an invariant dynamics, on texts this time, that does not require to know the substitute measures. We will here see texts as a concatenation of symbols, supposing that each sentence ends with a dedicated symbol (a period, if you prefer), so that the sentences of a text are uniquely defined. This dynamics will simply swap two elements of the same Markov substitute set.

Accordingly, we introduce $T = \gamma(S_1, \dots, S_n) \in D^+$, obtained by concatenating the sentences of our statistical sample.

To describe crossing-over dynamics, it will be useful to generalize the merge operation to multiple contexts. A two-fold context x being an element of $(D^*)^3$, we define

$$\alpha((x_1, x_2, x_3), y) = (x_1 y x_2, x_3).$$

Quite obviously, all formulas on α may be generalized on two-fold contexts, in particular the definition of Markov substitute sets

$$\mathbb{P}_S\left(\alpha(\alpha(x, y), z)\right) = \mathbb{P}_S\left(\alpha(\alpha(x, B), z)\right)q_B(y).$$

Definition 2.2.5

A double splitting kernel on a text is an kernel

$$\pi : D^+ \rightarrow \mathcal{M}_+^1\left(\left((D^*)^3\right) \times D^* \times D^* \times \{1, \dots, t\}\right)$$

such that for any $t \in D^+$,

$$\pi(t, x, y_1, y_2, j) > 0 \implies y_1, y_2 \in B_j \cup \{\epsilon\}, x \in (D^*)^3, \alpha(\alpha(x, y_1), y_2) = t.$$

Definition 2.2.6

Given a double splitting kernel π , we can define the crossing-over dynamics on texts

$$\begin{aligned} \sigma_\pi(t, t') &= \sum_{x \in (D^*)^3, (y_1, y_2) \in (D^*)^2, j} \pi(t, x, y_1, y_2, j) & (2.2.2) \\ &\quad \times \left(\frac{\pi(t', x, y_2, y_1, j)}{\pi(t, x, y_1, y_2, j)} \wedge 1 \right), & t \neq t' \end{aligned}$$

$$\sigma_\pi(t, t) = 1 - \sum_{t' \in D^+ \setminus \{t\}} \sigma_\pi(t, t'). \quad (2.2.3)$$

The simulation of T' such that $\mathbb{P}_{T'|T=t}(t') = \sigma_\pi(t, t')$ can be implemented in the following way. First draw $P = (X, Y_1, Y_2, j)$ such that $\mathbb{P}_{P|T} = \pi$, then set

$$T' = \begin{cases} \alpha(\alpha(X, Y_2), Y_1), & \text{w. p. } \left(\frac{\pi(T', X, Y_2, Y_1, j, n_1, n_2)}{\pi(T, X, Y_1, Y_2, j, n_1, n_2)} \wedge 1 \right), \\ T & \text{otherwise.} \end{cases}$$

Proposition 2.2.2

Any crossing-over dynamics defined as in definition 2.2.6 on the previous page is reversible with respect to $\mathbb{P}_T = (\mathbb{P}_S)^{\otimes n}$.

PROOF. It is obvious that $\sigma_\pi(t, t') = \sigma_\pi(t', t)$ for any texts t, t' , so we only have to prove that $\mathbb{P}_T(t) = \mathbb{P}_T(t')$ as soon as $\sigma_\pi(t, t') > 0$.

In this case, there is a configuration $(x, y_1, y_2, j) \in (D^*)^3 \times B_j^2 \times \{j\}$ such that $t = \alpha(\alpha(x, y_1), y_2)$ and $t' = \alpha(\alpha(x, y_2), y_1)$. We can then remark that

$$\begin{aligned} \mathbb{P}(T = t) &= q_{B_j}(y_2) \sum_{y \in B_j} \mathbb{P}[T = \alpha(\alpha(x, y_1), y)] \\ &= q_{B_j}(y_1) q_{B_j}(y_2) \sum_{y, y' \in B_j} \mathbb{P}[\alpha(\alpha(x, y), y')] \\ &= q_{B_j}(y_1) \sum_{y \in B_j} \mathbb{P}[T = \alpha(\alpha(x, y_2), y)] \\ &= \mathbb{P}[T = \alpha(\alpha(x, y_2), y_1)] \\ &= \mathbb{P}(T = t'). \end{aligned} \quad \square$$

2.3 Exponential families of Markov substitute processes

As there may be compatibility issues, knowing whether there exists a language distribution with a given family of Markov substitute sets is not obvious. We are going in this section to give a theoretical answer to this question, by which we mean a non constructive description of the Markov substitute processes having a prescribed set of Markov substitute sets. We will show more precisely that the language probability distributions having a prescribed finite collection of Markov substitute sets form a union of exponential families (depending on their support), and that these exponential families are non empty, meaning that the compatibility issue between Markov substitute sets never results in a void model.

Definition 2.3.1

For any given family \mathcal{B} of subsets of D^+ , we will say that the random process $S \in D^+$ is a \mathcal{B} -Markov substitute process if all the members $B \in \mathcal{B}$ are Markov substitute sets for \mathbb{P}_S .

The purpose of this section is to describe the set of \mathcal{B} -Markov substitute processes.

We will first introduce the special subfamily of independent Markov substitute processes, that will be useful to show that the set of \mathcal{B} -Markov substitute processes is not empty.

Given a strict sub-probability measure $\xi \in \mathcal{M}_+(D)$, let us put

$$r = 1 - \sum_{w \in D} \xi(w) > 0,$$

and let us define the independent process \tilde{S}_ξ by its distribution

$$\mathbb{P}(\tilde{S}_\xi = w_{1:k}) = \frac{r}{1-r} \prod_{j=1}^k \xi(w_j).$$

Let us remark that it is such that $\mathbb{P}[\ell(\tilde{S}) = L] = r(1-r)^{L-1}$. It is easy to see from the definition that $\text{supp}(\xi)^+$ is a Markov substitute set for the independent model \tilde{S} and that its substitute measure is equal to the distribution of \tilde{S}_ξ itself.

We will start with a description of the possible supports of \mathcal{B} -Markov substitute processes.

To do so, let us introduce the equivalence relation $\sim_{\mathcal{B}}$ defined as

$$\begin{aligned} y \sim_{\mathcal{B}} y' &\iff \exists(x_j, y_j, y'_j), x_j \in (D^*)^2, y_j, y'_j \in B_j \in \mathcal{B}, 0 \leq j \leq J, \\ y &= \alpha(x_0, y_0), y' = \alpha(x_J, y'_J), \alpha(x_{j-1}, y'_{j-1}) = \alpha(x_j, y_j), 0 < j \leq J. \end{aligned} \quad (2.3.1)$$

Let us remark that any component of $D^+/\sim_{\mathcal{B}}$ is a Markov substitute set for any \mathcal{B} -Markov substitute process, due to the basic properties of Markov substitute sets.

For any $C \in D^+/\sim_{\mathcal{B}}$, let us define the subset \mathcal{B}_C of \mathcal{B} present in C as

$$\mathcal{B}_C = \left\{ B \in \mathcal{B}, \sum_{x \in (D^*)^2} \mathbf{1}(\alpha(x, B) \cap C \neq \emptyset) > 0 \right\}.$$

Proposition 2.3.1

For any family \mathcal{B} of subsets of D^+ , for any \mathcal{B} -Markov substitute process S , there is a subset $\mathcal{C}_S \subset D^+/\sim_{\mathcal{B}}$ such that the support $\text{supp}(\mathbb{P}_S)$ is of the form

$$\text{supp}(\mathbb{P}_S) = \bigcup_{C \in \mathcal{C}_S} C \quad (2.3.2)$$

and such that

$$\mathcal{B} = \bigcup_{C \in \mathcal{C}_S} \mathcal{B}_C. \quad (2.3.3)$$

Conversely, if $\mathcal{C} \subset D^+ / \sim_{\mathcal{B}}$ is such that

$$\mathcal{B} = \bigcup_{C \in \mathcal{C}} \mathcal{B}_C, \quad (2.3.4)$$

given any strict sub-probability measure $\xi \in \mathcal{M}_+(D)$, such that $\text{supp}(\xi) = D$, and any probability measure $\mu \in \mathcal{M}_+^1(\mathcal{C})$, such that $\text{supp}(\mu) = \mathcal{C}$, the process S defined as

$$\mathbb{P}_S = \sum_{C \in \mathcal{C}} \mu(C) \mathbb{P}_{\tilde{S}_\xi | \tilde{S}_\xi \in C} \quad (2.3.5)$$

is a \mathcal{B} -Markov substitute process such that

$$\text{supp}(\mathbb{P}_S) = \bigcup_{C \in \mathcal{C}} C. \quad (2.3.6)$$

PROOF. Any component $C \in D^+ / \sim_{\mathcal{B}}$ such that $\text{supp}(\mathbb{P}_S) \cap C \neq \emptyset$ is a Markov substitute set. Consequently, $C \subset \text{supp}(\mathbb{P}_S)$, since by considering an empty context x in eq. (2.1.2) on page 64,

$$\mathbb{P}(S = s) = \mathbb{P}(C) q_C(s) > 0, \quad s \in C.$$

Moreover, condition eq. (2.3.3) on the current page is required to ensure that all members of \mathcal{B} satisfies eq. (2.1.1) on page 64.

The support of \tilde{S}_ξ being D^+ , it is clear that the support of S is as described in eq. (2.3.6) on the current page. This implies that any member $B \in \mathcal{B}$ satisfies eq. (2.1.1) on page 64.

Let us consider now any $B \in \mathcal{B}$, any $x \in (D^*)^2$, and any $y \in B$. There is $C \in D^+ / \sim_{\mathcal{B}}$ such that $\alpha(x, y) \in C$. Consequently, $\alpha(x, B) \in C$, since all members of this set are connected as described in eq. (2.3.1) on the previous page. We can write

$$\begin{aligned} \mathbb{P}_S(\alpha(x, y)) &= \frac{\mu(C) \mathbb{P}(\tilde{S}_\xi = \alpha(x, y))}{\mathbb{P}(\tilde{S}_\xi \in C)} \\ &= \frac{\mu(C) \mathbb{P}(\tilde{S}_\xi = \alpha(x, B)) \mathbb{P}(\tilde{S}_\xi = y | \tilde{S}_\xi \in B)}{\mathbb{P}(\tilde{S}_\xi \in C)} \\ &= \mathbb{P}_S(\alpha(x, B)) \mathbb{P}(\tilde{S}_\xi = y | \tilde{S}_\xi \in B), \end{aligned}$$

proving that B satisfies eq. (2.1.2) on page 64, with $q_B = \mathbb{P}_{\tilde{S}_\xi | \tilde{S}_\xi \in B}$. Consequently, S is a \mathcal{B} -Markov process as claimed in the proposition. In this proof, we have used the fact that the Markov substitute property is stable by conditioning. \square

We are now going to show that the \mathcal{B} -Markov processes form an exponential family, although we will not provide an efficient algorithm to compute the corresponding energy function (or in other terms sufficient statistics). In this respect, the method can be considered as non constructive.

Let \mathcal{B} be any finite family of finite subsets of D^+ .

The \mathcal{B} -Markov substitute processes remain the same if we replace \mathcal{B} by any other family of sets such that the equivalence relation $\sim_{\mathcal{B}}$ remains unchanged. Indeed, the fact that S is a \mathcal{B} -Markov substitute process is equivalent to the fact that any subset of $D^+/\sim_{\mathcal{B}}$ is a Markov substitute set, since it implies it and the reverse implication is trivial (since each $B \in \mathcal{B}$ is included in one of the components of $D^+/\sim_{\mathcal{B}}$).

This property allows to propose some reorganization of \mathcal{B} . First of all, we can split any $B \in \mathcal{B}$ into pairs. Namely, if $B = \{y_j, 1 \leq j \leq k\}$, we can replace B with the $k - 1$ sets $\{y_1, y_j\}, 1 < j \leq k$.

After this transformation, \mathcal{B} , being a finite set of pairs, can be considered as an undirected graph on D^+ (with a finite number of edges). We can then remove successively any pair included in a cycle, and prune in this way \mathcal{B} until we obtain an acyclic graph (this will leave the connected components of \mathcal{B} , and therefore also $\sim_{\mathcal{B}}$, unchanged). We can go further and remove more pairs as long as this removal does not change $\sim_{\mathcal{B}}$, until we obtain a minimal set of pairs. However, this last removal operation is not constructive since it is not clear that an efficient algorithm could check the required property.

Let us now decompose the distribution of any \mathcal{B} -Markov substitute process into

$$\mathbb{P}_S(s) = \sum_{C \in D^+/\sim_{\mathcal{B}}} \mathbb{P}(S \in C) \mathbb{P}(S = s | S \in C).$$

Let us choose in each $C \in D^+/\sim_{\mathcal{B}}$ a reference point $s_C \in C$, and let us write explicitly

$$\mathcal{B} = \left\{ \{y_{i,0}, y_{i,1}\}, 1 \leq i \leq I \right\},$$

where we have chosen an explicit order or indexation for each pair of \mathcal{B} .

Let us choose for any $s \in C \in D^+/\sim_{\mathcal{B}}$, such that $\mathbb{P}(C) > 0$ a path

$$\left((x_j, i_j, \sigma_j), 0 \leq j \leq J, x_j \in (D^*)^2, 1 \leq i_j \leq I, \sigma_j \in \{0, 1\} \right) \quad (2.3.7)$$

such that

$$\begin{aligned} \alpha(x_1, y_{i_0, \sigma_0}) &= s_C, \\ \alpha(x_J, y_{i_J, 1 - \sigma_J}) &= s, \\ \alpha(x_{j-1}, y_{i_{j-1}, 1 - \sigma_{j-1}}) &= \alpha(x_j, y_{i_j, \sigma_j}), \end{aligned} \quad (2.3.8)$$

and let us put

$$\begin{aligned}\beta_i &= \log \left(\frac{q_{\{y_{i,0}, y_{i,1}\}}(y_{i,1})}{q_{\{y_{i,0}, y_{i,1}\}}(y_{i,0})} \right), \quad 1 \leq i \leq I, \\ a_{s,i} &= \sum_{j=0}^J \mathbb{1}(i_j = i) \left(\mathbb{1}(\sigma_j = 0) - \mathbb{1}(\sigma_j = 1) \right), \quad s \in C, 1 \leq i \leq I, \\ b_s &= \log \left(\frac{\mathbb{P}(S = s)}{\mathbb{P}(S = s_C)} \right).\end{aligned}$$

The vector b_s defines $\mathbb{P}_{S|S \in C}$, according to the relation

$$\mathbb{P}(S = s | S \in C) = \frac{\exp(b_s)}{\sum_{s' \in C} \exp(b_{s'})}.$$

Moreover, due to the Markov substitute property,

$$b_s = \sum_{i=1}^I a_{s,i} \beta_i, \quad (2.3.9)$$

since

$$\frac{\mathbb{P}_S(s)}{\mathbb{P}_S(s_C)} = \prod_{j=0}^J \frac{q_{\{y_{j,0}, y_{j,1}\}}(y_{j,1-\sigma_j})}{q_{\{y_{j,0}, y_{j,1}\}}(y_{j,\sigma_j})}. \quad (2.3.10)$$

Let us remark now that S is a \mathcal{B} -Markov substitute process, if and only if

$$\mathcal{B} = \bigcup_{C \in D^+ / \sim_{\emptyset}, \mathbb{P}(S \in C) > 0} \mathcal{B}_C, \quad (2.3.11)$$

and eq. (2.3.10) on the current page holds not only for the chosen path, but for any path connecting s_C and s as prescribed above, for some set of substitute measures q_B , $B \in \mathcal{B}$, with full support $\text{supp}(q_B) = B$. Indeed the first condition eq. (2.3.11) on this page is required by eq. (2.1.1) on page 64 and the second condition is obviously necessary too. The two conditions are also sufficient. Indeed, in this case, for any $y \in B = \{y_{i,0}, y_{i,1}\} \in \mathcal{B}$, there is C such that $\mathbb{P}(C) > 0$ and there is $x \in (D^*)^+$ such that $\alpha(x, B) \cap C \neq \emptyset$. Therefore, $\alpha(x, B) \subset C$, because obviously $\alpha(x, y_{i,0}) \sim_{\mathcal{B}} \alpha(x, y_{i,1})$. According to eq. (2.3.10) on this page, $\mathbb{P}(S = s) > 0$ for any $s \in C$, so that $\mathbb{P}(S = \alpha(x, y_{i,0})) > 0$ and $\mathbb{P}(S = \alpha(x, y_{i,1})) > 0$, implying that $\mathbb{P}(S = \alpha(x, y)) > 0$, and therefore that eq. (2.1.1) on page 64 is satisfied. Moreover, given $x \in (D^*)^2$ and $i, 1 \leq i \leq I$, such that $\mathbb{P}(S = \alpha(x, \{y_{i,0}, y_{i,1}\})) > 0$, considering the component $C \in D^+ / \sim_{\mathcal{B}}$ such that $\alpha(x, y_{i,0}) \in C$, this implies that

$\alpha(x, y_{i,1}) \in C$ also, and that we can connect s_C with $\alpha(x, y_{i,1})$ adding to the path chosen to connect s_C with $\alpha(x, y_{i,0})$ the step $(x, i, 0)$ at the end of the path. Applying eq. (2.3.10) on the facing page to these two paths proves that

$$\frac{\mathbb{P}_S(\alpha(x, y_{i,1}))}{\mathbb{P}_S(\alpha(x, y_{i,0}))} = \frac{q_{\{y_{i,0}, y_{i,1}\}}(y_{i,1})}{q_{\{y_{i,0}, y_{i,1}\}}(y_{i,0})},$$

and therefore that S is a \mathcal{B} -Markov process.

Now we can remark that eq. (2.3.10) on the preceding page will hold for any path, and not only the chosen one, if and only if, for any loop $\ell \in \mathcal{L}$ connecting some s_C , such that $\mathbb{P}(S \in C) > 0$ to itself, that is satisfying eq. (2.3.8) on page 77, with $s = s_C$, the corresponding vector of weights

$$c_{\ell, i} = \sum_{j=0}^J \mathbf{1}(i_j = i) (\mathbf{1}(\sigma_j = 0) - \mathbf{1}(\sigma_j = 1)), \quad 1 \leq i \leq I,$$

is such that

$$\sum_{i=1}^I c_{\ell, i} \beta_i = 0. \quad (2.3.12)$$

Let us remark that \mathcal{L} is countable but may be infinite. By reindexing \mathcal{B}_C if necessary, we may assume that the vectors $(c_{\ell, i}, d < i \leq I)$ forms a basis of the vector space generated by the vectors $c_{\ell, i}, 1 \leq i \leq I$, where d is some integer. We can therefore find a real matrix $(e_{j, i}, d < j \leq I, 1 \leq i \leq I)$ such that

$$c_{\ell, i} = \sum_{j=d+1}^I c_{\ell, j} e_{j, i}, \quad \ell \in \mathcal{L}, \quad 1 \leq i \leq I.$$

Let us remark that the submatrix $(e_{j, i}, d < j \leq I, d < i \leq I)$ is equal to the identity matrix. Substituting in eq. (2.3.12) on the current page, we get

$$\sum_{j=d+1}^I c_{\ell, j} \sum_{i=1}^I e_{j, i} \beta_i = 0.$$

Since by construction $c_{\ell, i}, d < j \leq I$ are independent, eq. (2.3.12) on this page is equivalent to

$$\sum_{i=1}^I e_{j, i} \beta_i = 0, \quad d < j \leq I,$$

which can be rewritten according to the fact that $(e_{j, i}, d < j \leq I, d < i \leq I)$ is the identity, as

$$\beta_j = - \sum_{i=1}^d e_{j, i} \beta_i, \quad d < j \leq I.$$

Plugging this condition into eq. (2.3.9) on page 78, we get that S is a \mathcal{B} -Markov process if and only if

$$b_s = - \sum_{i=1}^d \beta_i U_i(s), \quad s \in \text{supp}(\mathbb{P}_S),$$

where

$$U_i(s) = - \left(a_{s,i} - \sum_{j=d+1}^I a_{s,j} e_{j,i} \right), \quad s \in \text{supp}(\mathbb{P}_S), \quad 1 \leq i \leq d.$$

Let us insist on the fact that the energy $U_i(s)$ depends only on $\text{supp}(\mathbb{P}_S)$, and not on the coefficients β_i that defines the substitute measures of the members of \mathcal{B} under \mathbb{P}_S .

We conclude that the process S is a \mathcal{B} -Markov substitute process if and only if its support satisfies eq. (2.3.11) on page 78 and there is a vector of real coefficients $(\beta_i \in \mathbb{R}, 1 \leq i \leq d)$ such that

$$\mathbb{P}(S = s) = \sum_{C \in D^+ / \sim_{\mathcal{B}}} \mathbf{1}(s \in C) \mathbb{P}(S \in C) \frac{\exp\left(- \sum_{i=1}^d \beta_i U_i(s)\right)}{\sum_{s' \in C} \exp\left(- \sum_{i=1}^d \beta_i U_i(s')\right)},$$

where the energy function $U_i(s)$, depends only on $\text{supp}(\mathbb{P}_S)$. This can be written as an exponential family on $\text{supp}(\mathbb{P}_S)$. More precisely, for any $\mathcal{C} \in D^+ / \sim_{\mathcal{B}}$ such that

$$\mathcal{B} = \bigcup_{C \in \mathcal{C}} \mathcal{B}_C,$$

the set of \mathcal{B} -Markov processes whose support is $\bigcup \mathcal{C}$ is the exponential family

$$\left\{ \left(\frac{\exp\left(- \sum_{C \in \mathcal{C} \setminus \{C_0\}} \beta_C \mathbf{1}(s \in C) - \sum_{i=1}^d \beta_i U_i(s)\right)}{\sum_{C \in \mathcal{C}} \sum_{s' \in C} \exp\left(- \sum_{C \in \mathcal{C} \setminus \{C_0\}} \beta_C \mathbf{1}(s' \in C) - \sum_{i=1}^d \beta_i U_i(s')\right)} \right), s \in \bigcup \mathcal{C} \right\}, \beta_C, \beta_i \in \mathbb{R},$$

where β_C and β_i should also be such that the denominator is not infinite.

Moreover, this exponential family contains the independent Markov substitute processes described previously, so that necessarily $d \geq 1$, as soon as \mathcal{B} is not trivial, that is contains at least a member that is not a one point set.

Let us summarize what we have proved in the following proposition.

Proposition 2.3.2

Given any finite set \mathcal{B} of finite subsets of D^+ , there is a finite set of pairs, \mathcal{P} , that we can choose such that each one is included in a member of \mathcal{B} , such that the sets of \mathcal{B} -Markov and \mathcal{P} -Markov substitute processes are the same, and such that \mathcal{P} is minimal for the inclusion relation (removing a pair from \mathcal{P} would break the above property).

For any \mathcal{B} -Markov substitute process S , we can decompose its support as a union of components $\mathcal{C} \subset D^+/\sim_{\mathcal{B}}$, such that $\text{supp}(\mathbb{P}_S) = \bigcup \mathcal{C}$ and such that

$$\mathcal{B} = \bigcup_{C \in \mathcal{C}} \mathcal{B}_C.$$

For any such subset, there is a free, non empty subset of pairs $\mathcal{F} \subset \mathcal{P}$, a matrix $(e_{j,i}, j \in \mathcal{P} \setminus \mathcal{F}, i \in \mathcal{F})$, an index set $I(\mathcal{C})$ containing \mathcal{F} and an energy function $(U_i(s), i \in I(\mathcal{C}), s \in C)$, such that the set of \mathcal{B} -Markov processes whose support is $\bigcup \mathcal{C}$ is the linear exponential family

$$\mathfrak{M}_{\mathcal{C}}(\mathcal{B}) = \left\{ \left(p(s) = \frac{\exp\left(-\sum_{i \in I(\mathcal{C})} \beta_i U_i(s)\right)}{\sum_{s' \in \bigcup \mathcal{C}} \exp\left(-\sum_{i \in I(\mathcal{C})} \beta_i U_i(s')\right)}, s \in \bigcup \mathcal{C} \right), \beta \in \mathfrak{B} \subset \mathbb{R}^{I(\mathcal{C})} \right\},$$

where

$$\mathfrak{B} = \left\{ \beta \in \mathbb{R}^{I(\mathcal{C})}, \sum_{s \in \bigcup \mathcal{C}} \exp\left(-\sum_{i \in I(\mathcal{C})} \beta_i U_i(s)\right) < \infty \right\},$$

and such that for any member p of this family, the substitute measure of any pair $i = \{y_{i,0}, y'_{i,1}\} \in \mathcal{F}$ (taken in a suitable order compatible with the definition of U), is given by

$$q_i(y_{i,\sigma}) = \frac{\exp(\sigma \beta_i)}{1 + \exp(\beta_i)}, \quad \sigma \in \{0, 1\},$$

whereas for $j \in \mathcal{P} \setminus \mathcal{F}$,

$$q_j(y_{j,\sigma}) = \frac{\exp(\sigma \beta_j)}{1 + \exp(\beta_j)}, \quad \sigma \in \{0, 1\},$$

where

$$\beta_j = \sum_{i \in \mathcal{F}(C)} e_{j,i} \beta_i.$$

This result implies that the maximum likelihood estimator in $\mathfrak{M}_{\mathcal{C}}(\mathcal{B})$ provides an asymptotically efficient estimator of the parameters β_i , $i \in I(\mathcal{C})$. We will describe in section 3.3 on page 135 exponential families of Markov substitute processes restricted to finite subdomains \mathcal{D} of D^+ and provide a way to compute with the help of a Monte-Carlo simulation the maximum likelihood estimator in this finite setting (where $\mathfrak{B} = \mathbb{R}^{I(\mathcal{C})}$).

Let us give an example showing that it is possible to have $\mathcal{P} \neq \mathcal{F}$ and that \mathcal{F} and therefore the form of the energy may depend on \mathcal{C} .

Let us consider the family $\mathcal{B} = \{B_1, B_2\}$, where $B_1 = \{ab, a\}$, $B_2 = \{bc, c\}$, and where $a, b, c \in D$.

Let us consider

$$\begin{aligned} C_1 &= \{ab^n c, n \in \mathbb{N}\}, \\ C_2 &= \{b^m cab^n, m, n \in \mathbb{N}\}, \\ C_3 &= \{b^k cab^m cab^n, k, m, n \in \mathbb{N}\}. \end{aligned}$$

It is easy to check that $C_1, C_2, C_3 \in D^+/\sim_{\mathcal{B}}$, and that $\{C_1\}$, $\{C_2\}$ and $\{C_3\}$ are all valid choices for \mathcal{C} , since both B_1 and B_2 are present in those three components. On C_1 , we have the loop

$$(a)c \rightarrow (ab)c \rightarrow a(bc) \rightarrow a(c),$$

and the induced constraint that $q_{B_1}(ab) = q_{B_2}(bc)$, so that

$$\mathfrak{M}_{\{C_1\}} = \left\{ \left(p(ab^n c) = r(1-r)^n \right), r \in]0, 1[\right\}.$$

Moreover we can take $\mathcal{F}(\{C_1\}) = \{B_1\}$, and although $\mathcal{B}_{C_1} = \mathcal{B}$, we could remove either B_1 or B_2 from \mathcal{B} while maintaining the fact that $C_1 \in D^+/\sim_{\mathcal{B}}$ and without changing $\mathfrak{M}_{\{C_1\}}$.

On C_2 , we have no non trivial loop, so that $\mathcal{F}(\{C_2\}) = \mathcal{B}$ and

$$\mathfrak{M}_{\{C_2\}} = \left\{ \left(p(b^m cab^n) = r_1 r_2 (1-r_1)^m (1-r_2)^m \right), r_1, r_2 \in]0, 1[\right\}.$$

On C_3 , we have the same kind of non trivial loop as on C_1 , imposing the constraint $q_{B_1} = q_{B_2}$, so that

$$\mathfrak{M}_{\{C_3\}} = \left\{ \left(p(b^k cab^m cab^n) = r(1-r)^{k+m+n}, k, m, n \in \mathbb{N} \right), r \in]0, 1[\right\},$$

but this time, removing B_1 or B_2 from \mathcal{B} would disconnect the component C_3 that would no longer belong to $D^+/\sim_{\mathcal{B}}$.

On $\{C_1, C_2\}$, the situation is again different

$$\mathfrak{M}_{\{C_1, C_2\}} = \left\{ \left(p(ab^n c) = r_1 r_2 (1 - r_2)^n, \right. \right. \\ \left. \left. p(b^m c a b^n) = (1 - r_1) r_2 (1 - r_2)^{m+n}, m, n \in \mathbb{N} \right), r_1, r_2 \in]0, 1[\right\}.$$

It is interesting to remark that the trace of $\mathfrak{M}_{\{C_1, C_2\}}$ on C_2 (the set of conditional probabilities on C_2 of distributions in $\mathfrak{M}_{\{C_1, C_2\}}$), is not equal to $\mathfrak{M}_{\{C_2\}}$: the presence of C_1 in the support imposes further constraints on the dependence between substitute measures.

2.4 Testing Markov substitute sets

We will discuss in this section the question of building statistical tests to decide whether sets of expressions are Markov substitute sets.

Let us consider a family of sets of contexts $\Theta \subset 2^{(D^*)^2}$, and assume that it contains all single contexts $\{x\}$, $x \in (D^*)^2$.

Let us assume that we already know that B_1 and B_2 are two Markov substitute sets. If they have a non void intersection, their union is also a Markov substitute set. However, this may be hard to check, so we assume that we do not know whether they are disjoint or not, and want to test whether their union $B = B_1 \cup B_2$ is a Markov substitute set.

To perform this test, we assume that we can check efficiently whether a given $y \in D^+$ belongs to B_1 or to B_2 , using a parsing algorithm. We consider the random variables

$$F_{B_1, B_2, \theta}(X_B, Y_B, p) = \left[\mathbf{1}(Y_B \in B_1) - p \mathbf{1}(Y_B \in B) \right] \mathbf{1}(X_B \in \theta), \\ \theta \in \Theta, p \in [0, 1], \quad (2.4.1)$$

where (X_B, Y_B) is the B -parse defined in section 2.1.5 on page 67. It is easy to deduce from the properties stated in this section that B is a Markov substitute set if and only if, there is $p \in [0, 1]$ such that for any $\theta \in \Theta$,

$$\mathbb{E} \left[F_{B_1, B_2, \theta}(X_B, Y_B, p) \right] = 0.$$

If we do not want to actually simulate the parse process (X_B, Y_B) , we can work with $\mathbb{E} \left[F_{B_1, B_2, \theta}(X_B, Y_B, p) \middle| S \right]$ instead of working with $F_{B_1, B_2, \theta}(X_B, Y_B, p)$ itself. In

other words, we may consider

$$\begin{aligned}
F_{B_1, B_2, \theta}(S, p) &= \mathbb{E} \left[\left(\mathbb{1}(Y_B \in B_1) - p \mathbb{1}(Y_B \in B) \right) \mathbb{1}(X_B \in \theta) \middle| S \right] \\
&= \sum_{x \in (D^*)^2} \sum_{y \in B_1 \cup B_2} \mathbb{1}(x \in \theta) \min_{y' \in B_1 \cup B_2} \pi(\alpha(x, y'), x, y') \\
&\quad \times \left[\mathbb{1}(y \in B_1)(1-p) - \mathbb{1}(y \in B_2)p \right] \mathbb{1}(S = \alpha(x, y)). \quad (2.4.2)
\end{aligned}$$

Using this last formula, we can see that

$$\begin{aligned}
\mathbb{E} \left[F_{B_1, B_2, \theta}(S, p) \right] &= \sum_{x \in \theta} \min_{y' \in B_1 \cup B_2} \pi(\alpha(x, y'), x, y') \\
&\quad \times \left[\mathbb{P}_S(\alpha(x, B_1))(1-p) - p \mathbb{P}_S(\alpha(x, B_2)) \right].
\end{aligned}$$

2.4.1 Alternative construction of a test function

For any $\theta \in \Theta$, let us consider some probability measure $\nu_\theta \in \mathcal{M}_+^1((D^*)^2)$, such that $\theta \subset \text{supp}(\nu_\theta)$ (in the case when θ is a finite set, we can for instance consider the uniform measure on θ).

Let X_θ be some random variable distributed according to ν_θ and independent of S . Since $y \mapsto \alpha(x, y)$ is one to one, we can define a random variable $Y_\theta \in D^*$ by the property $\alpha(X_\theta, Y_\theta) = S$ when $S \in \alpha(X_\theta, D^+)$, and $Y_\theta = \epsilon$ otherwise.

Let us define the test functions

$$\begin{aligned}
F_{B_1, B_2, \theta}(S, p) &= \sum_{x \in (D^*)^2} \sum_{y \in D^+} \nu_\theta(x) \mathbb{1}(x \in \theta) \mathbb{1}(S = \alpha(x, y)) \\
&\quad \times \left[\mathbb{1}(y \in B_1) - p \mathbb{1}(y \in B_1 \cup B_2) \right] \\
&= \mathbb{E} \left[\left(\mathbb{1}(Y_\theta \in B_1) - p \mathbb{1}(Y_\theta \in B_1 \cup B_2) \right) \mathbb{1}(X_\theta \in \theta) \middle| S \right]. \quad (2.4.3)
\end{aligned}$$

We see on this expression that this definition of the test function may not be very efficient when θ is large, because of the factor $\nu_\theta(x)$, but that it will not be affected by the size of $B_1 \cup B_2$, that is not even required to be a finite set, whereas the definition given in eq. (2.4.2) on the current page may be inefficient when $B_1 \cup B_2$ is large, and does not allow $B_1 \cup B_2$ to contain an infinite number of elements.

Given that B_1 and B_2 are two Markov substitute sets, it is easy to see that $B = B_1 \cup B_2$ is a Markov substitute set if and only if there is $p \in [0, 1]$ such that for any $\theta \in \Theta$,

$$\mathbb{E} \left[F_{B_1, B_2, \theta}(S, p) \right] = 0.$$

Indeed, we can write this expectation as

$$\sum_{x \in \theta} \nu_\theta(x) \left[\mathbb{P}_S(\alpha(x, B_1))(1-p) - p \mathbb{P}_S(\alpha(x, B_2)) \right].$$

If we are ready to loose the interpretation of $F_{B_1, B_2, \theta}(S, p)$ in terms of the random variables X_θ and Y_θ , while retaining the property that $|F_{B_1, B_2, \theta}(S, p)| \leq 1$, we can take for some real parameter $\alpha > 1$,

$$\nu_\theta(x) = \nu(x) = \nu(x_1, x_2) = (\alpha - 1)^2 (\ell(x_1) + 2)^{-\alpha} (\ell(x_2) + 2)^{-\alpha},$$

that is such that for any $s \in D^+$,

$$\sum_{x \in (D^*)^2} \sum_{y \in D^+} \nu(x) \mathbb{1}(s = \alpha(x, y)) \leq 1.$$

We can also build tests that are more closely related to the invariant dynamics introduced in section 2.2.1 on page 69. Let us assume that B_1 and B_2 are disjoint Markov substitute sets. Let π be a conditional probability kernel, satisfying eq. (2.1.4) on page 67, for $B = B_1 \cup B_2$.

Consider the family of test functions

$$\begin{aligned} F_{B_1, B_2, \theta}(S, p) &= \sum_{x \in (D^*)^2} \sum_{y_1 \in B_1, y_2 \in B_2} \mathbb{1}(x \in \theta) \left[\pi(\alpha(x, y_1), x, y_1) \wedge \pi(\alpha(x, y_2), x, y_2) \right] \\ &\times \left[\mathbb{1}(S = \alpha(x, y_1)) q_{B_2}(y_2) (1 - p) - \mathbb{1}(S = \alpha(x, y_2)) q_{B_1}(y_1) p \right]. \end{aligned} \quad (2.4.4)$$

The advantage, when compared to eq. (2.4.2) on the facing page or eq. (2.4.3) on the preceding page is that we take the infimum on two values of the kernel π only, and that the coefficient $\nu_\theta(x)$ is not present. The drawback is that we have to know, or at least to estimate the substitute measures q_{B_1} and q_{B_2} .

Proposition 2.4.1

Let us assume that B_1 and B_2 are two Markov substitute sets. Let us also assume that for any $(x, y) \in (D^*)^2 \times (B_1 \cup B_2)$, $\pi(\alpha(x, y), x, y) > 0$. Then $B_1 \cup B_2$ is a Markov substitute set if and only if there is $p \in [0, 1]$ such that for any $\theta \in \Theta$,

$$\mathbb{E}(F_{B_1, B_2, \theta}(S, p)) = 0,$$

so that we can use the random variables $F_{B_1, B_2, \theta}(S, p)$ to define a test.

To see that the proposition is true, we can remark that

$$\begin{aligned} \mathbb{E}(F_{B_1, B_2, \theta}(S, p)) &= \sum_{x \in (D^*)^2} \left(\sum_{(y_1, y_2) \in B_1 \times B_2} \mathbb{1}(x \in \theta) q_{B_1}(y_1) q_{B_2}(y_2) \right. \\ &\times \left. \left[\pi(\alpha(x, y_1), x, y_1) \wedge \pi(\alpha(x, y_2), x, y_2) \right] \right) \\ &\times \left[\mathbb{P}_S(\alpha(x, B_1)) (1 - p) - \mathbb{P}_S(\alpha(x, B_2)) p \right]. \end{aligned}$$

If $B_1 \cup B_2$ is a Markov substitute set, then

$$\begin{aligned}\mathbb{P}_S(\alpha(x, B_1)) &= \mathbb{P}_S(\alpha(x, B_1 \cup B_2)q_{B_1 \cup B_2}(B_1)), \\ \mathbb{P}_S(\alpha(x, B_2)) &= \mathbb{P}_S(\alpha(x, B_1 \cup B_2)q_{B_1 \cup B_2}(B_2)),\end{aligned}$$

so that $\mathbb{E}(F_{B_1, B_2, \theta}(S, q_{B_1 \cup B_2}(B_1))) = 0$.

On the other hand, if for any $\theta \in \Theta$, $\mathbb{E}(F_{B_1, B_2, \theta}(S, p)) = 0$, then for any context $x \in (D^*)^2$,

$$\mathbb{P}_S(\alpha(x, B_1)) - p \mathbb{P}_S(\alpha(x, B_1 \cup B_2)) = 0,$$

implying that $B_1 \cup B_2$ is a Markov substitute set with $q_{B_1 \cup B_2}(B_1) = p$.

To see how $F_{B_1, B_2, \theta}(S, p)$ can be computed by simulations, we can remark that

$$\begin{aligned}F_{B_1, B_2, \theta}(S, p) &= \sum_{x \in (D^*)^2} \sum_{y, y' \in (D^*)^2} \pi(S, x, y) \mathbb{1}(x \in \theta) \\ &\quad \times \left[\mathbb{1}(y \in B_1)(1-p)q_{B_2}(y') - \mathbb{1}(y \in B_2)pq_{B_1}(y') \right] \\ &\quad \times \left(\frac{\pi(\alpha(x, y'), x, y')}{\pi(\alpha(x, y), x, y)} \wedge 1 \right).\end{aligned}$$

Thus, if we draw

$$\mathbb{P}_{(X, Y)|S} = \pi$$

and

$$\mathbb{P}_{Y'|X, Y} = \mathbb{1}(Y \in B_1)q_{B_2} + \mathbb{1}(Y \in B_2)q_{B_1}$$

we obtain that

$$\begin{aligned}F_{B_1, B_2, \theta}(S, p) &= \mathbb{E} \left\{ \mathbb{1}(X \in \theta) \left[\mathbb{1}(Y \in B_1)(1-p) - \mathbb{1}(Y \in B_2)p \right] \right. \\ &\quad \left. \times \left(\frac{\pi(\alpha(X, Y'), X, Y')}{\pi(\alpha(X, Y), X, Y)} \wedge 1 \right) \middle| S \right\}.\end{aligned}$$

Let us introduce the weights

$$w(x, y) = \mathbb{E} \left(\left(\frac{\pi(\alpha(X, Y'), X, Y')}{\pi(\alpha(X, Y), X, Y)} \wedge 1 \right) \middle| X = x, Y = y \right),$$

and the random variables (X'', Y'') defined as

$$\mathbb{P}_{X'', Y''|X, Y} = w(X, Y)\delta_{X, Y} + (1 - w(X, Y))\delta_{(\alpha(X, Y), \epsilon)}, \epsilon.$$

We can rewrite $F_{B_1, B_2, \theta}(S, p)$ as

$$F_{B_1, B_2, \theta}(S, p) = \mathbb{E}\left\{\mathbf{1}(X'' \in \theta)\left[\mathbf{1}(Y'' \in B_1) - p\mathbf{1}(Y'' \in B_1 \cup B_2)\right]\middle| S\right\}.$$

This shows that we could also simulate (X'', Y'') and base the tests on

$$F_{B_1, B_2, \theta}(X'', Y'', p) = \mathbf{1}(X'' \in \theta)\left[\mathbf{1}(Y'' \in B_1) - p\mathbf{1}(Y'' \in B_1 \cup B_2)\right]. \quad (2.4.5)$$

In this way we fall back on the same setting as the previous one, defined by eq. (2.4.1) on page 83.

2.5 Weakening the Markov substitute assumption

Let us remark that if we do not want to assume that $\pi(\alpha(x, y), x, y) > 0$ for all $(x, y) \in (D^*)^2 \times D^*$, we will test a weaker property than the Markov substitute property, as stated in the following proposition.

Proposition 2.5.1

Let us assume that for any $x \in (D^*)^2$, there is $\theta \in \Theta$ such that $\theta = \{x\}$. Let us assume also that the dynamics

$$\begin{aligned} k(s, s') &= \sum_{x \in (D^*)^2, (y, y') \in (D^*)^2} \pi(s, x, y) \\ &\quad \times \left[\mathbf{1}(y \in B_1)q_{B_1}(y') + \mathbf{1}(y \in B_2)q_{B_2}(y') \right] \left(\frac{\pi(s', x, y')}{\pi(s, x, y)} \wedge 1 \right), \quad s \neq s', \\ k(s, s) &= 1 - \sum_{s' \neq s} k(s, s') \end{aligned}$$

is reversible with respect to \mathbb{P}_S . Under these assumptions, the dynamics

$$\begin{aligned} k'(s, s') &= \sum_{x \in (D^*)^2, (y, y') \in (D^*)^2} \pi(s, x, y) \\ &\quad \times \left[pq_{B_1}(y') + (1-p)q_{B_2}(y') \right] \left(\frac{\pi(s', x, y')}{\pi(s, x, y)} \wedge 1 \right), \quad s \neq s', \\ k'(s, s) &= 1 - \sum_{s' \neq s} k(s, s'). \end{aligned}$$

is reversible with respect to the same invariant probability measure if and only if, for all $\theta \in \Theta$, $\mathbb{E}(F_{B_1, B_2, \theta}(S, p)) = 0$, where the test function is defined as in eq. (2.4.4) on page 85.

This comes from the fact that k is reversible if and only if for any any $i \in \{1, 2\}$, any $x \in (D^*)^2$, $(y, y') \in B_i^2$ such that $\pi(\alpha(x, y), x, y) \wedge \pi(\alpha(x, y'), x, y') > 0$,

$$\mathbb{P}(S = \alpha(x, y))q_{B_i}(y') = \mathbb{P}(S = \alpha(x, y'))q_{B_i}(y).$$

Actually, weakening the hypothesis opens perspectives for a more refined modeling. We can for instance implement along this line the masking of some syntactic categories. As an example, in the idiom “ a Trojan horse ” you would like to prevent the possibility of replacing horse with another noun, whereas you would like to be able to do so in other contexts.

2.6 Testing the Markov substitute property using a parse process

In this section, we will propose a test for the Markov substitute property based on test functions of the type

$$F_\theta(X, Y, p) = \mathbb{1}(X \in \theta) [\mathbb{1}(Y \in B_1) - p \mathbb{1}(Y \in B_1 \cup B_2)],$$

where B_1 and B_2 are Markov substitute sets, and the test has to decide whether $B = B_1 \cup B_2$ is also a Markov substitute set. This covers two possible constructions of the parse process, defined either by eq. (2.4.1) on page 83, or by eq. (2.4.5) on the previous page. We make this study as an introduction to tests based on the other type of test functions, where a conditional expectation knowing S is taken. In particular, we will not address here the question of making the test uniform with respect to the range of values that the pair B_1, B_2 may take.

As already explained, in both cases, we want to test whether there is a value of $p \in [0, 1]$ such that for any $\theta \in \Theta$, $\mathbb{E}[F_\theta(X, Y, p)] = 0$. Let us introduce the constants

$$p_+ = \sup\left\{\mathbb{P}(Y \in B_1 | Y \in B, X = x); x \in (D^*)^2, \mathbb{P}(Y \in B, X = x) > 0\right\},$$

$$p_- = \inf\left\{\mathbb{P}(Y \in B_1 | Y \in B, X = x); x \in (D^*)^2, \mathbb{P}(Y \in B, X = x) > 0\right\}.$$

Remark that with these notations, $\mathbb{E}[F_\theta(X, Y, p)] = 0$ when $p_+ = p_- = p$, which corresponds exactly to the Markov substitute property to be tested.

Since $p \mapsto F_\theta(X, Y, p)$ is non-increasing,

$$\mathbb{E}[F_\theta(X, Y, p_+) | X, Y \in B] \leq 0,$$

$$\mathbb{E}[F_\theta(X, Y, p_-) | X, Y \in B] \geq 0,$$

and, observing that F_θ is null when $\mathbb{1}(Y \in B) = 0$, we obtain that

$$\begin{aligned} \mathbb{E}\left[F_\theta(X, Y, p_+) \mid X, \mathbb{1}(Y \in B)\right] &\leq 0, \\ \mathbb{E}\left[F_\theta(X, Y, p_-) \mid X, \mathbb{1}(Y \in B)\right] &\geq 0. \end{aligned} \quad (2.6.1)$$

These inequalities will be used to get bounds for p_+ and p_- , based on concentration inequalities concerning the empirical mean of $F_\theta(X, Y, p)$.

The test will be based on a simulation (X_i, Y_i) , $1 \leq i \leq n$ of the random parse process, where each pair (X_i, Y_i) is independently drawn from the conditional probability distribution $\mathbb{P}_{X, Y \mid S}$. Doing so, we obtain an i.i.d. sample (X_i, Y_i) , $1 \leq i \leq n$.

2.6.1 Probability of false rejection

Proposition 2.6.1

Let $\mu \in \mathcal{M}_+^1$ be a probability measure depending only on $(X_i, \mathbb{1}(Y_i \in B))$, $i \in \llbracket 1, n \rrbracket$. Let us define

$$F_{\theta, i}(p) = F_\theta(X_i, Y_i, p).$$

Let Λ be a finite subset of $]0, 1[$. With probability at least $1 - 2\varepsilon$,

$$\begin{aligned} B_-(p_+) &\stackrel{\text{def}}{=} \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \int \sum_{i=1}^n \log(1 + \lambda F_{\theta, i}(p_+)) d\rho(\theta) - \mathcal{H}(\rho, \mu) - \log(|\Lambda|/\varepsilon) \leq 0, \\ B_+(p_-) &\stackrel{\text{def}}{=} \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \int \sum_{i=1}^n \log(1 - \lambda F_{\theta, i}(p_-)) d\rho(\theta) - \mathcal{H}(\rho, \mu) - \log(|\Lambda|/\varepsilon) \leq 0, \end{aligned}$$

where

$$\mathcal{H}(\rho, \mu) = \begin{cases} \int \log\left(\frac{d\rho}{d\mu}\right) d\rho, & \text{when } \rho \ll \mu, \\ +\infty, & \text{otherwise.} \end{cases}$$

Therefore, if we reject the hypothesis that B is a Markov substitute set when

$$\inf_{p \in [0, 1]} \max\{B_-(p), B_+(p)\} > 0, \quad (2.6.2)$$

the probability of false rejection is at most 2ε .

Remark 2.6.1

We may also strengthen the Markov substitute property, requiring that the substitute measure is not too unbalanced. For any real parameter $\eta \in]0, 1[$, we will say that (B_1, B_2) is an η -Markov substitute pair of sets when $B = B_1 \cup B_2$ is a Markov substitute set such that $q_B(B_1) \in [\eta, 1 - \eta]$. According to the proposition, we can

reject the hypothesis that the pair (B_1, B_2) is an η -Markov substitute pair of sets when

$$\inf_{p \in [\eta, 1-\eta]} \max\{B_-(p), B_+(p)\} > 0, \quad (2.6.3)$$

with a probability of false rejection not greater than 2ε .

PROOF. Let us prove the last statement of the proposition first, assuming that the rest is true. In the case when B is a Markov substitute set, $p_- = p_+$, so that with probability at least $1 - 2\varepsilon$,

$$\inf_{p \in [0,1]} \max\{B_-(p), B_+(p)\} \leq \max\{B_-(p_+), B_+(p_+)\} \leq 0.$$

Consequently

$$\mathbb{P}\left(\inf_{p \in [0,1]} \max\{B_-(p), B_+(p)\} > 0\right) \leq 2\varepsilon.$$

If moreover (B_1, B_2) is an η -Markov substitute pair of sets, then $p_+ \in [\eta, 1 - \eta]$, so that with probability at least $1 - 2\varepsilon$,

$$\inf_{p \in [\eta, 1-\eta]} \max\{B_-(p), B_+(p)\} \leq \max\{B_-(p_+), B_+(p_+)\} \leq 0,$$

proving that

$$\mathbb{P}\left(\inf_{p \in [\eta, 1-\eta]} \max\{B_-(p), B_+(p)\} > 0\right) \leq 2\varepsilon.$$

Let us now proceed to the proof of the first statement of the proposition. Let us put

$$W_\theta(\lambda) = \sum_{i=1}^n \log(1 + \lambda F_{\theta,i}(p_+)).$$

This is legitimate, since $|F_{\theta,i}(p)| \leq 1$, and $\lambda < 1$. From eq. (2.6.1) on the preceding page, and the independence of $(X_i, \mathbf{1}(Y_i \in B), F_{\theta,i}(p_+))$, $1 \leq i \leq n$,

$$\begin{aligned} \mathbb{E}\left[\exp(W_\theta(\lambda)) \mid X_i, \mathbf{1}(Y_i \in B), 1 \leq i \leq n\right] \\ = \prod_{i=1}^n \mathbb{E}\left(1 + \lambda F_{\theta,i}(p_+) \mid X_i, \mathbf{1}(Y_i \in B)\right) \leq 1. \end{aligned}$$

Since μ depends only on X_i , and $\mathbf{1}(Y_i \in B)$, $1 \leq i \leq n$, we get that

$$\begin{aligned} \mathbb{E}\left[\int \exp(W_\theta(\lambda)) d\mu(\theta)\right] \\ = \mathbb{E}\left[\int \mathbb{E}\left[\exp(W_\theta(\lambda)) \mid X_i, \mathbf{1}(Y_i \in B), 1 \leq i \leq n\right] d\mu(\theta)\right] \leq 1. \end{aligned}$$

Together with the remark that for any distribution $\rho \in \mathcal{M}_+^1(\Theta)$,

$$\log\left(\int \exp(h) d\mu\right) \geq \int h d\rho - \mathcal{K}(\rho, \mu), \quad (2.6.4)$$

we get finally that

$$\begin{aligned} \mathbb{E}\left[\exp(B_-(p_+))\right] &= \mathbb{E}\left[\exp\left(\sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \int W_\theta(\lambda) d\rho(\theta) - \mathcal{K}(\rho, \mu) - \log(|\Lambda|/\varepsilon)\right)\right] \\ &\leq \mathbb{E}\left[\sum_{\lambda \in \Lambda} \int \exp\left(W_\theta(\lambda) - \log(|\Lambda|/\varepsilon)\right) d\mu(\theta)\right] \leq \varepsilon, \end{aligned}$$

proving through the Markov inequality that $\mathbb{P}(B_-(p_+) \geq 0) \leq \varepsilon$. The proof concerning $B_+(p_-)$ is the same, mutatis mutandi (changing λ for $-\lambda$ and $+$ for $-$). \square

In practice, we will assume that

$$\sup\left\{\left|\left\{\theta \in \Theta : x \in \theta\right\}\right| : x \in (D^*)^2, \alpha(x, B) \cap \text{supp}(\mathbb{P}_S) \neq \emptyset\right\} \leq m, \quad (2.6.5)$$

and choose for μ the uniform probability measure on

$$\hat{\Theta} = \left\{\theta \in \Theta, \{X_i, 1 \leq i \leq n\} \cap \theta \neq \emptyset\right\}, \quad (2.6.6)$$

a set of size at most mn . In this case $\mathcal{K}(\rho, \mu) \leq \log(|\hat{\Theta}|) \leq \log(mn)$, the former inequality being an equality when ρ is a Dirac mass.

2.6.2 Probability of false acceptance of the hypothesis

We have built the test of eq. (2.6.2) on page 89 by controlling its probability of false rejection. We would like now to study its probability of false acceptance. To this purpose, let us establish lower bounds for $B_-(p_+)$ and $B_+(p_-)$. (We need indeed to prove that these quantities cannot be both negative when the hypothesis is false.)

Lemma 2.6.2

Let us define

$$p(\theta) = \mathbb{P}(Y \in B_1 | X \in \theta, Y \in B).$$

For any $p \in [0, 1]$, any $\lambda \in]-1, 1[$,

$$\begin{aligned} \mathbb{E}\left[\exp\left(-\log\left[1 + \lambda F_\theta(p)\right]\right)\right] &\leq \exp\left[\mathbb{E}\left[\mathbf{1}(X \in \theta, Y \in B)\right]\left[\lambda(p - p(\theta))\right.\right. \\ &\quad \left.\left.+ \frac{\lambda^2}{1 - |\lambda|}\left(p(\theta)(1 - p(\theta)) + (p - p(\theta))^2\right)\right]\right]. \end{aligned}$$

PROOF. Let us remark first that for $-1 < \lambda < 1$, $-1 \leq x \leq 1$,

$$-\log(1 + \lambda x) = \log\left(1 - \lambda x + \frac{\lambda^2 x^2}{1 + \lambda x}\right) \leq \log\left(1 - \lambda x + \frac{\lambda^2 x^2}{1 - |\lambda|}\right).$$

As a consequence,

$$\mathbb{E}\left[\exp\left(-\log[1 + \lambda F_\theta(p)]\right)\right] \leq 1 - \lambda \mathbb{E}(F_\theta(p)) + \frac{\lambda^2}{1 - |\lambda|} \mathbb{E}(F_\theta(p)^2).$$

The rest follows easily from

$$\begin{aligned} \mathbb{E}(F_\theta(p)) &= \mathbb{E}\left[\left(\mathbf{1}(Y \in B_1) - p\right)\mathbf{1}(X \in \theta, Y \in B)\right] \\ &= (p(\theta) - p)\mathbb{E}\left[\mathbf{1}(X \in \theta, Y \in B)\right], \\ \mathbb{E}(F_\theta(p)^2) &= \mathbb{E}\left[\left(\mathbf{1}(Y \in B_1) - p\right)^2 \mathbf{1}(X \in \theta, Y \in B)\right] \\ &= \mathbb{E}\left[\left(\mathbf{1}(Y \in B_1)(1 - 2p) + p^2\right)\mathbf{1}(X \in \theta, Y \in B)\right] \\ &= \left(p(\theta)(1 - 2p) + p^2\right)\mathbb{P}(X \in \theta, Y \in B) \\ &= \left[p(\theta)(1 - p(\theta)) + (p - p(\theta))^2\right]\mathbb{P}(X \in \theta, Y \in B). \quad \square \end{aligned}$$

Proposition 2.6.3

Let us assume that eq. (2.6.5) on the preceding page holds and that μ is the uniform probability measure on the finite set of parameters $\hat{\Theta}$ defined by eq. (2.6.6) on the previous page. Let us put

$$\delta = \frac{\log(mn) + 2\log(\varepsilon^{-1}) + \log(|\Lambda|)}{n},$$

and

$$\chi = \sup_{x \in [n^{-1/2}, n^{1/2}]} \inf_{\lambda \in \Lambda} \cosh\left[\log\left(\frac{\lambda}{(1 - \lambda)x}\right)\right].$$

Let us assume that there are θ_+ and $\theta_- \in \Theta$ such that $\bar{p}_+ = p(\theta_+)$, $\bar{p}_- = p(\theta_-)$, $q_+ = \mathbb{E}[\mathbf{1}(Y \in B)f_{\theta_+}(X)]$ and $q_- = \mathbb{E}[\mathbf{1}(Y \in B)f_{\theta_-}(X)]$ satisfy

$$\begin{aligned} \min\{q_-, q_+\} &\geq 8\chi^2\delta, \\ \bar{p}_+ - \bar{p}_- &\geq 2\chi\sqrt{\frac{\bar{p}_+(1 - \bar{p}_+)\delta}{q_+}}\left(1 + \frac{4\chi^2\delta}{q_+}\right) + \frac{(2 + \sqrt{2})\delta}{q_+} \\ &\quad + 2\chi\sqrt{\frac{\bar{p}_-(1 - \bar{p}_-)\delta}{q_-}}\left(1 + \frac{4\chi^2\delta}{q_-}\right) + \frac{(2 + \sqrt{2})\delta}{q_-}. \end{aligned} \quad (2.6.7)$$

Then with probability at least $1 - 2\varepsilon$

$$\inf_{p \in [0,1]} \max\{B_-(p), B_+(p)\} > 0.$$

Therefore, the test defined by eq. (2.6.2) on page 89 has in this case a probability of false acceptance at most equal to 2ε .

If we assume instead of eq. (2.6.7) on the preceding page that

$$\bar{p}_+ - 2\chi \sqrt{\frac{\bar{p}_+(1 - \bar{p}_+)\delta}{q_+}} \left(1 + \frac{4\chi^2\delta}{q_+}\right) - \frac{(2 + \sqrt{2})\delta}{q_+} > 1 - \eta, \quad (2.6.8)$$

$$\text{or that } \bar{p}_- + 2\chi \sqrt{\frac{\bar{p}_-(1 - \bar{p}_-)\delta}{q_-}} \left(1 + \frac{4\chi^2\delta}{q_-}\right) + \frac{(2 + \sqrt{2})\delta}{q_-} > \eta, \quad (2.6.9)$$

then with probability at least $1 - \varepsilon$,

$$\inf_{p \in [\eta, 1-\eta]} \max\{B_-(p), B_+(p)\} > 0.$$

Therefore the test defined by eq. (2.6.3) on page 90 has a probability of false acceptance not greater than ε in this case.

These two statements are consequences of the following more precise one. Each of the two following inequalities holds with probability at least $1 - \varepsilon$ (so that both hold together with probability at least $1 - 2\varepsilon$):

$$B_- \left(\bar{p}_+ - 2\chi \sqrt{\frac{\bar{p}_+(1 - \bar{p}_+)\delta}{q_+}} \left(1 + \frac{4\chi^2\delta}{q_+}\right) - \frac{(2 + \sqrt{2})\delta}{q_+} \right) > 0 \quad (2.6.10)$$

$$B_+ \left(\bar{p}_- + 2\chi \sqrt{\frac{\bar{p}_-(1 - \bar{p}_-)\delta}{q_-}} \left(1 + \frac{4\chi^2\delta}{q_-}\right) + \frac{(2 + \sqrt{2})\delta}{q_-} \right) > 0. \quad (2.6.11)$$

PROOF. Let us assume for a while the more precise statement at the end of the proposition. Under the hypothesis specified by eq. (2.6.7) on the preceding page, there is $p_* \in [0, 1]$, such that

$$\begin{aligned} \bar{p}_- + 2\chi \sqrt{\frac{\bar{p}_-(1 - \bar{p}_-)\delta}{q_-}} \left(1 + \frac{4\chi^2\delta}{q_-}\right) + \frac{(2 + \sqrt{2})\delta}{q_-} \\ \leq p_* \\ \leq \bar{p}_+ - 2\chi \sqrt{\frac{\bar{p}_+(1 - \bar{p}_+)\delta}{q_+}} \left(1 + \frac{4\chi^2\delta}{q_+}\right) - \frac{(2 + \sqrt{2})\delta}{q_+}. \end{aligned}$$

As a consequence of the fact that $p \mapsto B_-(p)$ is non-increasing and $p \mapsto B_+(p)$ is non-decreasing, this implies from the precise statement at the end of the proposition that with probability at least $1 - 2\varepsilon$, $\min\{B_-(p_*), B_+(p_*)\} > 0$, and therefore that

$$\begin{aligned} & \inf_{p \in [0,1]} \max\{B_-(p), B_+(p)\} \\ & \geq \min\left\{\inf_{p \in [0,p_*]} B_-(p), \inf_{p \in [p_*,1]} B_+(p)\right\} = \min\{B_-(p_*), B_+(p_*)\} > 0. \end{aligned}$$

In the same way, if we assume eq. (2.6.8) on the preceding page, we get with probability at least $1 - \varepsilon$ that $B_-(1 - \eta) > 0$ and if we assume that eq. (2.6.9) on the previous page, we get with probability at least $1 - \varepsilon$ that $B_+(\eta) > 0$.

Let us now prove the precise statement at the end of the proposition. Let us put

$$p = \bar{p}_+ - 2\chi \sqrt{\frac{\bar{p}_+(1 - \bar{p}_+)\delta}{q_+}} \left(1 + \frac{4\chi^2\delta}{q_+}\right) - \frac{(2 + \sqrt{2})\delta}{q_+}.$$

Let us choose $\lambda \in \Lambda$ such that

$$\lambda \in \operatorname{argmin}_{\zeta \in \Lambda} \cosh \left[\log \left(\frac{\zeta}{1 - \zeta} \sqrt{\frac{[\bar{p}_+(1 - \bar{p}_+) + (\bar{p}_+ - p)^2]q_+}{\delta}} \right) \right].$$

Let us remark that $\bar{p}_+ - p \geq \frac{\delta}{q_+}$, so that

$$\sqrt{\frac{1}{n}} \leq \sqrt{\frac{\delta}{q_+}} \leq \sqrt{\frac{[\bar{p}_+(1 - \bar{p}_+) + (\bar{p}_+ - p)^2]q_+}{\delta}} \leq \sqrt{\frac{q_+}{\delta}} \leq \sqrt{n}.$$

This implies that

$$\cosh \left[\log \left(\frac{\lambda}{1 - \lambda} \sqrt{\frac{[\bar{p}_+(1 - \bar{p}_+) + (\bar{p}_+ - p)^2]q_+}{\delta}} \right) \right] \leq \chi.$$

Let us remark also that

$$B_-(p) \geq \sum_{i=1}^n \log[1 + \lambda F_{\theta_+,i}(p)] - \log(mn) - \log(|\Lambda|/\varepsilon).$$

According to lemma 2.6.2 on page 91, with probability at least $1 - \varepsilon$,

$$\begin{aligned} & \sum_i \log(1 + \lambda F_{\theta_+,i}(p)) \\ & \geq -\log(\varepsilon^{-1}) - \lambda(p - \bar{p}_+) - \frac{\lambda^2}{1 - \lambda} \left(\bar{p}_+(1 - \bar{p}_+) + (p - \bar{p}_+)^2 \right). \end{aligned}$$

This implies that with probability at least $1 - \varepsilon$,

$$B_-(p) > nq_+\lambda \left((\bar{p}_+ - p) - \frac{\lambda}{1-\lambda} [\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2] - \frac{\delta}{\lambda q_+} \right).$$

Decomposing $\frac{\delta}{\lambda q_+} = \frac{\delta(1-\lambda)}{\lambda q_+} + \frac{\delta}{q_+}$ and using the identity

$$ax + b/x = 2\sqrt{ab} \cosh \left[\log \left(x \sqrt{\frac{a}{b}} \right) \right], \quad x, a, b \in \mathbb{R}_+^*,$$

we obtain that with probability at least $1 - \varepsilon$,

$$\begin{aligned} B_-(p) &> nq_+\lambda \left\{ (\bar{p}_+ - p) - \frac{\delta}{q_+} - 2\sqrt{\frac{[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]\delta}{q_+}} \right. \\ &\quad \left. \times \cosh \left[\log \left(\frac{\lambda}{1-\lambda} \sqrt{\frac{[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]q_+}{\delta}} \right) \right] \right\} \\ &\geq nq_+\lambda \left(\bar{p}_+ - p - \frac{\delta}{q_+} - 2\chi \sqrt{\frac{[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]\delta}{q_+}} \right). \end{aligned}$$

Let us put $\rho = \frac{4\chi^2\delta}{q_+}$, $v = \bar{p}_+(1-\bar{p}_+)$, $y = \frac{\delta}{q_+}$ and $x = \bar{p}_+ - p$. The sign of the right-hand side of the previous inequality is the same as the sign of

$$\begin{aligned} \left(\bar{p}_+ - p - \frac{\delta}{q_+} \right)^2 - \frac{4\chi^2\delta[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]}{q_+} \\ = x^2 - 2yx + y^2 - \rho v - \rho x^2 \\ = (1-\rho)x^2 - 2yx - (\rho v - y^2). \end{aligned}$$

This quantity is non negative when x is not less than

$$\begin{aligned} \frac{y}{1-\rho} + \sqrt{\frac{y^2}{(1-\rho)^2} + \frac{\rho v - y^2}{1-\rho}} &= \frac{y}{1-\rho} + \sqrt{\frac{\rho y^2}{(1-\rho)^2} + \frac{\rho v}{1-\rho}} \\ &\leq \sqrt{\frac{\rho v}{1-\rho}} + \frac{(1+\sqrt{\rho})y}{1-\rho} \leq \sqrt{\rho v} \left(1 + \frac{\rho}{2(1-\rho)} \right) + \frac{(1+\sqrt{\rho})y}{1-\rho} \\ &\leq \sqrt{\rho v}(1+\rho) + (2+\sqrt{2})y = 2\chi \sqrt{\frac{\bar{p}_+(1-\bar{p}_+)\delta}{q_+}} \left(1 + \frac{4\chi^2\delta}{q_+} \right) + \frac{(2+\sqrt{2})\delta}{q_+}, \end{aligned}$$

where we have used the inequalities $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$, $\sqrt{\frac{1}{1-\rho}} \leq 1 + \frac{\rho}{2(1-\rho)}$ and $\rho \leq 1/2$.

This proves that with probability at least $1 - \varepsilon$,

$$B_- \left(\bar{p}_+ - 2\chi \sqrt{\frac{\bar{p}_+(1-\bar{p}_+)\delta}{q_+}} \left(1 + \frac{4\chi^2\delta}{q_+} \right) - \frac{(2+\sqrt{2})\delta}{q_+} \right) > 0.$$

The proof of eq. (2.6.11) on page 93 is obtained exactly in the same way. \square

2.7 Testing Markov substitute sets without simulating the parse process

In this section, we will embrace a broader scope than testing for a single Markov substitute pair of sets.

We will assume that some partition \mathcal{B} of D^+ into Markov substitute set is known, and we will simultaneously test whether any pair $(B_1, B_2) \in \mathcal{B}^2$ is such that $B = B_1 \cup B_2$ is a Markov substitute set.

Instead of simulating random B -parse processes, which add some random fluctuations to the observations, we will consider the following conditional expectations:

$$F_{B_1, B_2, \theta}(s, p) = \mathbb{E} \left\{ \left[\mathbf{1}(Y_{B_1 \cup B_2} \in B_1) - p \mathbf{1}(Y_{B_1 \cup B_2} \in B_1 \cup B_2) \right] \mathbf{1}(X_{B_1 \cup B_2} \in \theta) \mid S = s \right\}. \quad (2.7.1)$$

Let us notice that similar results could be stated about

$$F_{B_1, B_2, \theta}(s, p) = \mathbb{E} \left\{ \left[\mathbf{1}(X_{B_1 \cup B_2} \in \theta) - p \right] \mathbf{1}(Y_{B_1 \cup B_2} \in B_1) \mid S = s \right\}, \quad (2.7.2)$$

We will not follow this second path, since the first is more closely related to the invariant dynamics we will use to simulate from the model. The results of this section also apply to the alternative choices of test functions described in section 2.4.1 on page 84.

In the case when

$$p = \mathbb{P}(Y_{B_1 \cup B_2} \in B_1 \mid X_{B_1 \cup B_1} \in \theta, Y_{B_1 \cup B_2} \in B_1 \cup B_2),$$

it is easy to see that $\mathbb{E}(F_{B_1, B_2, \theta}(S, p)) = 0$, and consequently that

$$\mathbb{E} \left[\exp \left(\log \left(1 + \lambda F_{B_1, B_2, \theta}(S, p) \right) \right) \right] = 1, \quad \lambda \in]-1, 1[.$$

Let us introduce $\bar{\Theta} = \mathcal{B}^2 \times \Theta$, the parameter space indexing the random variables $F_{B_1, B_2, \theta}(S, p)$ (apart from p that can be viewed as another parameter).

In order to master the complexity of the problem in a data driven way, we will use a shadow sample, as in the proof of Vapnik-Chervonenkis complexity bounds. However, the actual observation of this shadow sample will not be required, because our generalization bounds will eventually be obtained by taking a conditional expectation with respect to the observed sample $S_{1:n}$.

Lemma 2.7.1

Consider an i.i.d sample S_i , $1 \leq i \leq kn$ of size kn and an exchangeable prior distribution $\mu_{S_{1:kn}} \in \mathcal{M}_+^1(\bar{\Theta})$, where we have put for short $S_{1:kn} = (S_i, 1 \leq i \leq kn)$.

For any $\lambda \in]-1, 1[$, any $\theta \in \bar{\Theta}$, any $p \in [0, 1]$,

$$\mathbb{E} \left(\int \exp \left(\sum_{i=1}^n \log(1 + \lambda F_{\theta}(S_i, p)) - \frac{\lambda}{k} \sum_{i=1}^{kn} F_{\theta}(S_i, p) \right) d\mu(\theta) \right) \leq 1.$$

PROOF. Let us put $W_i = \lambda F_{\theta}(S_i, p)$. For any exchangeable function $G(W_{1:kn})$,

$$\begin{aligned} \mathbb{E} \left(G(W_{1:kn}) \prod_{i=1}^n (1 + W_i) \right) &= \mathbb{E} \left(G(W_{1:kn}) \sum_{j_1, \dots, j_n=0}^{k-1} \prod_{i=1}^n \frac{(1 + W_{i+jn})}{k} \right) \\ &= \mathbb{E} \left(G(W_{1:kn}) \prod_{i=1}^n \left(1 + \frac{1}{k} \sum_{j=0}^{k-1} W_{i+jn} \right) \right) \\ &= \mathbb{E} \left(G(W_{1:kn}) \exp \left(\sum_{i=1}^n \log \left(1 + \frac{1}{k} \sum_{j=0}^{k-1} W_{i+jn} \right) \right) \right) \\ &\leq \mathbb{E} \left(G(W_{1:kn}) \exp \left(n \log \left(1 + \frac{1}{nk} \sum_{i=1}^{kn} W_i \right) \right) \right) \\ &\leq \mathbb{E} \left(G(W_{1:kn}) \exp \left(\frac{1}{k} \sum_{i=1}^{kn} W_i \right) \right). \end{aligned}$$

If we take for any value of θ

$$G_{\theta}(W_{1:kn}) = \mu(\theta) \exp \left(-\frac{\lambda}{k} \sum_{i=1}^{kn} F_{\theta}(S_i, p) \right),$$

we get that

$$\begin{aligned} &\mathbb{E} \left(\prod_{i=1}^n (1 + \lambda F_{\theta}(S_i, p)) \exp \left(-\frac{\lambda}{k} \sum_{i=1}^{kn} F_{\theta}(S_i, p) \right) \mu(\theta) \right) \\ &\leq \mathbb{E} \left(\exp \left(\frac{\lambda}{k} \sum_{i=1}^{kn} F_{\theta}(S_i, p) \right) \exp \left(-\frac{\lambda}{k} \sum_{i=1}^{kn} F_{\theta}(S_i, p) \right) \mu(\theta) \right) = \mathbb{E}(\mu(\theta)). \end{aligned}$$

The inequality follows by summing on θ . □

Let us now make a more specific choice for μ .

Let us consider some indicator function $h : \mathcal{B} \times D^+ \rightarrow \{0, 1\}$, and the exchangeable prior on \mathcal{B} defined by

$$\nu(B) = \frac{1}{kn} \sum_{i=1}^{kn} \left(\sum_{B' \in \mathcal{B}} h(B', S_i) \right)^{-1} h(B, S_i).$$

Let us also consider another indicator function $g : \mathcal{B} \times \Theta \times D^+ \rightarrow \{0, 1\}$, and define

$$\xi(\theta | B_1, B_2) = \frac{1}{kn} \sum_{i=1}^{kn} \left(\sum_{\theta' \in \Theta} g(B_1 \cup B_2, \theta', S_i) \right)^{-1} g(B_1 \cup B_2, \theta, S_i).$$

Let us introduce

$$\begin{aligned} \nu_1(B) &= \frac{1}{n} \sum_{i=1}^n \left(\sum_{B' \in \mathcal{B}} h(B', S_i) \right)^{-1} h(B, S_i), \text{ and} \\ \xi_1(\theta | B_1, B_2) &= \frac{1}{n} \sum_{i=1}^n \left(\sum_{\theta' \in \Theta} g(B_1 \cup B_2, \theta', S_i) \right)^{-1} g(B_1 \cup B_2, \theta, S_i). \end{aligned}$$

We can then choose the prior distribution

$$\mu(B_1, B_2, \theta) = \nu(B_1)\nu(B_2)\xi(\theta | B_1, B_2), \quad (2.7.3)$$

and decompose it as $\mu = k^{-3}\mu_1 + (1 - k^{-3})\mu_2$, where

$$\mu_1(B_1, B_2, \theta) = \nu_1(B_1)\nu_1(B_2)\xi_1(\theta | B_1, B_2)$$

depends only on the observable sample $S_{1:n}$. We can remark moreover that

$$\mathcal{K}(\rho, \mu) \leq \mathcal{K}(\rho, \mu_1) + 3 \log(k). \quad (2.7.4)$$

Let us remark that by the law of large numbers, when n is large, μ_1 will get close to μ .

We have more specifically in mind the following weight functions

$$\begin{aligned} h(B, s) &= \mathbb{1}(\mathcal{S}(s, B) \neq \emptyset), \\ g(B, \theta, s) &= \mathbb{1}(\exists (x, y) \in \mathcal{S}(s, B); x \in \theta). \end{aligned}$$

With this choice, the support of ν_1 is made of all Markov substitute sets $B \in \mathcal{B}$ such that one expression of B is observed in the sample $S_{1:n}$. In the same way, the support of $\mu_1(\cdot | B_1, B_2)$ is made of all the context parameters $\theta \in \Theta$ indexing contexts that appear in the sample $S_{1:n}$.

Let us put

$$\begin{aligned}\bar{h}(s) &= \sum_{B \in \mathcal{B}} h(B, s), \\ \bar{g}(s) &= \sup_{B_1, B_2 \in \mathcal{B}} \sum_{\theta \in \Theta} g(B_1 \cup B_2, \theta, s).\end{aligned}$$

We can easily see that as soon as $\mu_1(B_1, B_2, \theta) > 0$, then

$$\mu_1(B_1, B_2, \theta) \geq n^{-3} \left(\max_{i=1, \dots, n} \bar{h}(S_i) \right)^{-2} \left(\max_{i=1, \dots, n} \bar{g}(S_i) \right)^{-1}.$$

We can improve this inequality if we choose a conditional probability kernel π satisfying eq. (2.1.4) on page 67 only when $\bar{h}(s) \leq C_1$ and $\bar{g}(s) \leq C_2$, and such that $\pi(s, (s, \epsilon), \epsilon) = 1$ otherwise. With this choice of kernel, we test the Markov substitute property only for the sentences s satisfying the above constraint. More specifically, we test that

$$\mathbb{P}_S(\alpha(x, y)) = \mathbb{P}_S(\alpha(x, B)) q_B(y), \quad \inf_{y' \in B} \pi(\alpha(x, y'), x, y') > 0.$$

On the other hand, we are now sure that

$$\mu_1(B_1, B_2, \theta) \geq n^{-3} C_1^{-2} C_2^{-1},$$

as soon as $\mu_1(B_1, B_2, \theta) > 0$,

From lemma 2.7.1 on page 97, we get in this context

Proposition 2.7.2

Consider some finite set $\Lambda \subset]0, 1[$, and $-\Lambda = \{-\lambda, \lambda \in \Lambda\}$. With probability at least $1 - 2\varepsilon$, for any $\lambda \in \Lambda \cup (-\Lambda)$, any $p \in \mathcal{P} \subset]0, 1[$, any $\rho \in \mathcal{M}_+^1(\bar{\Theta})$,

$$\begin{aligned}& \sum_{i=1}^n \frac{(k-1)\lambda}{k} \int F_\theta(S_i, p) d\rho(\theta) - \frac{\lambda^2}{2(1-|\lambda|)^2} \int F_\theta(S_i, p)^2 d\rho(\theta) \\ & \leq \int \sum_{i=1}^n \left[\log(1 + \lambda F_\theta(S_i, p)) - \frac{\lambda}{k} F_\theta(S_i, p) \right] d\rho(\theta) \\ & \leq \frac{(k-1)n\lambda}{k} \int \mathbb{E}[F_\theta(S, p)] d\rho(\theta) + \mathcal{K}(\rho, \mu_1) + 3 \log(k) + \log(|\Lambda| |\mathcal{P}| / \varepsilon).\end{aligned}$$

Remark 2.7.1

From the shape of this equation, it makes sense to take k of order $k = 10$ or so.

PROOF. The first inequality is a consequence of the fact that

$$\log(1 + \lambda x) = \lambda x - \int_0^{\lambda x} \frac{(\lambda x - y)}{(1 + y)^2} dy \geq \lambda x - \frac{\lambda^2 x^2}{2(1 - |\lambda|)^2}, \quad x \in [-1, 1], \lambda \in]-1, 1[.$$

In order to prove the second inequality, let us put for short

$$\begin{aligned} W_1(\theta) &= \sum_{i=1}^n \log(1 + \lambda F_\theta(S_i, p)) - \frac{\lambda}{k} \sum_{i=1}^n F_\theta(S_i, p), \\ W_2(\theta) &= -\frac{\lambda}{k} \sum_{i=n+1}^{kn} F_\theta(S_i, p). \end{aligned}$$

According to the previous lemma,

$$\mathbb{E} \left(\exp \left(\sup_{\rho \in \mathcal{M}(\bar{\Theta})} \int [W_1(\theta) + W_2(\theta)] d\rho(\theta) - \mathcal{K}(\rho, \mu_1) - 3 \log(k) \right) \right) \leq 1.$$

Using

$$\begin{aligned} \mathbb{E}(\exp(\cdot)) &\geq \mathbb{E}[\exp(\mathbb{E}(\cdot|S_{1:n}))], \\ \mathbb{E}(\sup_{\rho}(\cdot)|S_{1:n}) &\geq \sup_{\rho} \mathbb{E}(\cdot|S_{1:n}), \\ \mathbb{E}(W_1(\theta)|S_{1:n}) &= W_1(\theta), \\ \mathbb{E}(W_2(\theta)|S_{1:n}) &= \mathbb{E}(W_2(\theta)), \\ \text{and } \mathbb{E}(\mathcal{K}(\rho, \mu_1)|S_{1:n}) &= \mathcal{K}(\rho, \mu_1), \end{aligned}$$

we obtain that

$$\mathbb{E} \left(\exp \left(\sup_{\rho \in \mathcal{M}_+^1(\bar{\Theta})} \int W_1(\theta) + \mathbb{E}(W_2(\theta)) d\rho(\theta) - \mathcal{K}(\rho, \mu_1) - 3 \log(k) \right) \right) \leq 1,$$

from which the proposition is an easy consequence. \square

2.7.1 Definition of the test and probability of false rejection

This leads us to define an alternative to the test proposed in section 2.6 on page 88. Let $p(B_1, B_2, \theta)$ be the value of p such that

$$\mathbb{E}(F_{B_1, B_2, \theta}(S, p)) = 0,$$

so that

$$p(B_1, B_2, \theta) = \mathbb{P}(Y_{B_1 \cup B_2} \in B_1 \mid X_{B_1 \cup B_2} \in \theta, Y_{B_1 \cup B_2} \in B_1 \cup B_2),$$

where $F_{B_1, B_2, \theta}(S, p)$ is defined by eq. (2.7.1) on page 96.

Let us define

$$p_+(B_1, B_2) = \sup \left\{ p(B_1, B_2, \theta) : \theta \in \Theta, \mathbb{P}(X_{B_1 \cup B_2} \in \theta, Y_{B_1 \cup B_2} \in B_1 \cup B_2) > 0 \right\},$$

and let $p_-(B_1, B_2)$ be the infimum of the same set.

Let us introduce the function

$$\psi(z) = \log(1 + z) - k^{-1}z.$$

Proposition 2.7.3

Let Λ be a finite subset of $]0, 1[$.

With probability at least $1 - 2\varepsilon$, for any pair $(B_1, B_2) \in \mathcal{B}^2$,

$$\begin{aligned} B_-(p_+(B_1, B_2)) &\stackrel{\text{def}}{=} \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \int \sum_{i=1}^n \psi \left(\lambda F_{B_1, B_2, \theta}(S_i, p_+(B_1, B_2)) \right) d\rho(\theta) \\ &\quad - \mathcal{K}(\rho, \mu_1) - 3 \log(k) - \log \left(\frac{|\Lambda|}{\varepsilon \nu_1(B_1) \nu_1(B_2)} \right) \leq 0 \\ B_+(p_-(B_1, B_2)) &\stackrel{\text{def}}{=} \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \int \sum_{i=1}^n \psi \left(-\lambda F_{B_1, B_2, \theta}(S_i, p_-(B_1, B_2)) \right) d\rho(\theta) \\ &\quad - \mathcal{K}(\rho, \mu_1) - 3 \log(k) - \log \left(\frac{|\Lambda|}{\varepsilon \nu_1(B_1) \nu_1(B_2)} \right) \leq 0 \end{aligned}$$

Therefore, if we reject the hypothesis that $B_1 \cup B_2$ is a Markov substitute set when

$$\inf_{p \in [0, 1]} \max \{ B_-(p), B_+(p) \} > 0, \quad (2.7.5)$$

the probability of making a false rejection (after testing all pairs in \mathcal{B}^2) is at most 2ε .

In the same way we can reject the hypothesis that $(B_1, B_2) \in \mathcal{B}^2$ is an η -Markov substitute pair of sets when

$$\inf_{p \in [\eta, 1-\eta]} \max \{ B_-(p), B_+(p) \} > 0,$$

with a probability of rejecting one of the true η -Markov pairs (after testing all pairs in \mathcal{B}^2), not greater than 2ε .

Remark 2.7.2

To our knowledge the observable entropy term $\mathcal{K}(\rho, \mu_1)$ is a novel addition to PAC-Bayes theory that was not proposed before in the literature.

This proposition is a direct consequence of the previous one.

We can remark at this point that this result stays true for any test functions $F_\theta(S, p)$ such that $|F_\theta(S, p)| \leq 1$ and $\mathbb{E}(F_\theta(S, p)) = 0$ for $p = p(B_1, B_2, \theta)$, so this kind of test can be used for a wide range of test functions. However, this proposition gives no indication on the probability on false acceptance (which is admittedly the more important part). Section 2.7.5 on page 111 will deal with this issue, but will require a more restrictive form of the test function.

2.7.2 Testing for the weak Markov substitute property

We may be interested in situations where the target language does not have a sufficient number of strong Markov substitute sets, but nevertheless has a bigger number of weak Markov substitute sets. In this case, we will be able to generate the support of the language, running an ill specified Markov substitute set model whose Markov substitute sets are only weak Markov substitute sets for the true language distribution.

This is how we propose to modify the test. To get the right support, we want to select pairs which satisfy $0 < p_-(B_1, B_2) \leq p_+(B_1, B_2) < 1$. Since it will not be possible to test that $p_-(B_1, B_2) = 0$ or $p_+(B_1, B_2) = 1$, we will take some margin and focus on pairs for which

$$\eta \leq p_-(B_1, B_2) \leq p_+(B_1, B_2) \leq 1 - \eta,$$

where $\eta \in]0, 1/2[$ is a fixed parameter. We will also relax the property that $p_-(B_1, B_2) = p_+(B_1, B_2)$ to $p_+(B_1, B_2) - p_-(B_1, B_2) \leq \gamma$, where $\gamma \in [0, 1 - 2\eta]$ is a second parameter. Pairs satisfying those two conditions will be called γ -weak η -Markov substitute pairs of sets. Let us remark that when $\gamma = 1 - 2\eta$, the second condition is void (because it is implied by the first one). Let us also remark that the first condition says that the ratio between the probabilities of observing $y \in B_1$ and $y \in B_2$ in the same context belongs to the interval $[\eta/(1 - \eta), (1 - \eta)/\eta]$.

Corollary 2.7.4

With probability at least $1 - 2\varepsilon$, for any γ -weak η -Markov substitute pair,

$$\inf_{p \in [\eta, 1 - \eta - \gamma]} \max\{B_-(p + \gamma), B_+(p)\} \leq 0,$$

so that if we accept as γ -weak η -Markov pairs all pairs satisfying this condition, the probability of false rejection, that is, the probability that some truly γ -weak η -Markov pair in \mathcal{B}^2 may have been unduly rejected is not greater than 2ε .

2.7.3 Computation of the test

Let us remark first that a more explicit formula is available for the test function. Indeed, the solution of the optimization with respect to ρ is explicit, according to the general formula

$$\sup_{\rho \in \mathcal{M}_+^1(\Theta)} \int h(\theta) d\rho(\theta) - \mathcal{K}(\rho, \mu_1) = \log \left[\int \exp(h(\theta)) d\mu_1(\theta) \right],$$

the supremum being achieved when

$$\frac{d\rho}{d\mu_1} = \frac{\exp(h)}{\int \exp(h) d\mu_1}.$$

According to this general formula

$$B_-(p) = \log \left[\int \exp \left[\sum_{i=1}^n \psi \left(\lambda F_{B_1, B_2, \theta}(S_i, p) \right) \right] d\mu_1(\theta) \right] - \log \left(\frac{k^3 |\Lambda|}{\varepsilon \nu_1(B_1) \nu_1(B_2)} \right), \quad (2.7.6)$$

$$B_+(p) = \log \left[\int \exp \left[\sum_{i=1}^n \psi \left(-\lambda F_{B_1, B_2, \theta}(S_i, p) \right) \right] d\mu_1(\theta) \right] - \log \left(\frac{k^3 |\Lambda|}{\varepsilon \nu_1(B_1) \nu_1(B_2)} \right). \quad (2.7.7)$$

These formulas are indeed very compact and pleasing from this point of view, but they do not show the structure of the test function. For example they do not show whether or how the test function involves some variance and complexity estimates.

To show this structure, we are going in the following lines to compute a more explicit approximation, in the form of a close empirical lower bound for B_- and B_+ . Let us put

$$\begin{aligned} a(\theta, i) &= \mathbb{P}(X_{B_1 \cup B_2} \in \theta, Y_{B_1 \cup B_2} \in B_1 \mid S = S_i), \\ b(\theta, i) &= \mathbb{P}(X_{B_1 \cup B_2} \in \theta, Y_{B_1 \cup B_2} \in B_1 \cup B_2 \mid S = S_i), \\ \delta(\rho) &= \mathcal{K}(\rho, \mu_1) + \log \left(\frac{k^3 |\Lambda|}{\varepsilon \nu_1(B_1) \nu_2(B_2)} \right). \end{aligned}$$

so that

$$F_{B_1, B_2, \theta}(S_i, p) = a(\theta, i) - p b(\theta, i).$$

According to proposition 2.7.2 on page 99,

$$B_-(p) \geq \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \int \left(\sum_{i=1}^n \frac{(k-1)\lambda}{k} [a(\theta, i) - p b(\theta, i)] - \frac{\lambda^2}{2(1-\lambda)^2} [a(\theta, i) - p b(\theta, i)]^2 \right) d\rho(\theta) - \delta(\rho).$$

Let us put, for $k = 1, 2$,

$$\begin{aligned} a_k(\rho) &= \int \sum_{i=1}^n a(\theta, i)^k d\rho(\theta), \\ b_k(\rho) &= \int \sum_{i=1}^n b(\theta, i)^k d\rho(\theta), \\ \hat{p}_1(\rho) &= \frac{a_1(\rho)}{b_1(\rho)}, \\ \hat{p}_2(\rho) &= b_2(\rho)^{-1} \int \sum_{i=1}^n a(\theta, i) b(\theta, i) d\rho(\theta), \\ \hat{v}(\rho) &= b_2(\rho)^{-1} \int \sum_{i=1}^n (a(\theta, i) - \hat{p}_2 b(\theta, i))^2 d\rho(\theta), \\ \alpha &= \frac{k-1}{k}. \end{aligned}$$

We get that

$$B_-(p) \geq \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \alpha b_1(\rho) \lambda (\hat{p}_1(\rho) - p) - \frac{\lambda^2}{2(1-\lambda)^2} b_2(\rho) [\hat{v}(\rho) + (\hat{p}_2(\rho) - p)^2] - \delta(\rho).$$

We deduce that with probability at least $1 - \varepsilon$, $p_+(y_1, y_2)$ is solution in p of the inequalities (indexed by $\lambda \in \Lambda$ and $\rho \in \mathcal{M}_+^1(\Theta)$)

$$\hat{p}_1(\rho) - p - \frac{\lambda b_2(\rho)}{2(1-\lambda)^2 \alpha b_1(\rho)} [\hat{v}(\rho) + (\hat{p}_2(\rho) - p)^2] - \frac{\delta(\rho)}{\lambda \alpha b_1(\rho)} \leq 0.$$

Using the inequality $(1-\lambda)^{-2} \leq (1-2\lambda)^{-1}$, and optimizing the value of λ , we get

$$\hat{p}_1(\rho) - p - \chi \sqrt{\frac{2b_2(\rho) [\hat{v}(\rho) + (\hat{p}_2(\rho) - p)^2] \delta(\rho)}{\alpha^2 b_1(\rho)^2}} - \frac{2\delta(\rho)}{\alpha b_1(\rho)} \leq 0, \quad (2.7.8)$$

where

$$\chi = \inf_{\lambda \in \Lambda} \cosh \left[\log \left(\frac{\lambda}{1-2\lambda} \sqrt{\frac{b_2(\rho) [\hat{v}(\rho) + (\hat{p}_2(\rho) - p)^2]}{2\delta(\rho)}} \right) \right].$$

We are now going to solve this inequality in p . Although this is straightforward, it will get a little technical, so this may be the best place to give an interpretation of the bound. We have an empirical estimate $\hat{p}_1(\rho)$ of some mixture of conditional probabilities (the mixture being related to ρ). In this empirical estimate, $b_1(\rho)$ plays the role of an empirical “effective” sample dimension. This empirical estimate is corrected by a deviation term, to provide a lower bound for $p_+(B_1, B_2)$ with confidence level $1 - \varepsilon$. We also get an upper bound for $p_-(B_1, B_2)$ with confidence level $1 - \varepsilon$,

$$p - \hat{p}_1(\rho) - \chi \sqrt{\frac{2b_2(\rho) [\hat{v}(\rho) + (\hat{p}_2(\rho) - p)^2] \delta(\rho)}{\alpha^2 b_1(\rho)^2}} - \frac{2\delta(\rho)}{\alpha b_1(\rho)} \leq 0,$$

resulting in a final joint statement holding at confidence level $1 - 2\varepsilon$. The deviation term involves a variance estimate given by $\hat{v}(\rho)$ and a complexity (or “effective dimension”) estimate given by $\delta(\rho)$. The confidence level is uniform in ρ , which includes the case where ρ is a Dirac mass at parameter value θ , but a better compromise may be found for a distribution ρ with more spread (as shown by eq. (2.7.6) on page 103, the optimal ρ is indeed a Gibbs distribution and not a Dirac mass, although a Dirac mass may be a close winner in this discrete setting). Optimizing the bound in ρ will give a tighter lower bound for $p_+(B_1, B_2)$, which is itself defined as an upper bound in θ , and for $p_-(B_1, B_2)$ which is defined as a lower bound in θ .

Let us now solve explicitly the quadratic inequality eq. (2.7.8) on the facing page. Let us put

$$\begin{aligned} x &= \hat{p}_2(\rho) - p_+(B_1, B_2), \\ y &= \hat{p}_2(\rho) - \hat{p}_1(\rho) + \frac{2\delta(\rho)}{\alpha b_1(\rho)}, \\ z &= \frac{2\chi^2 b_2(\rho) \delta(\rho)}{\alpha^2 b_1(\rho)^2}, \\ v &= \hat{v}(\rho). \end{aligned}$$

We deduce from the previous discussion that with probability at least $1 - \varepsilon$,

$$(x - y)^2 - z(v + x^2) \leq 0,$$

which can also be written as

$$(1 - z)x^2 - 2yx - (zv - y^2) \leq 0,$$

implying that

$$\begin{aligned} x &\leq \frac{y}{1-z} + \sqrt{\frac{y^2}{(1-z)^2} + \frac{zv-y^2}{1-z}} \\ &= \frac{y}{1-z} + \sqrt{\frac{zy^2}{(1-z)^2} + \frac{zv}{1-z}} \\ &\leq \frac{1+\sqrt{z}}{1-z}y + \sqrt{zv}\left(1 + \frac{z}{2(1-z)}\right). \end{aligned}$$

Thus, if we put

$$\hat{p}_+(\rho) = \hat{p}_2(\rho) - \frac{1+\sqrt{z}}{1-z}y - \sqrt{zv}\left(1 + \frac{z}{2(1-z)}\right),$$

and in the same way (reasoning with $B_+(p)$)

$$\hat{p}_-(\rho) = \hat{p}_2(\rho) + \frac{1+\sqrt{z}}{1-z}[y + 2(\hat{p}_1(\rho) - \hat{p}_2(\rho))] + \sqrt{zv}\left(1 + \frac{z}{2(1-z)}\right),$$

we obtain

Proposition 2.7.5

With the previous notations, with probability at least $1 - 2\varepsilon$, for any $\rho \in \mathcal{M}_+^1(\Theta)$,

$$\begin{aligned} p_-(B_1, B_2) &< \inf_{\rho \in \mathcal{M}_+^1(\Theta)} \hat{p}_-(\rho), \\ \sup_{\rho \in \mathcal{M}_+^1(\Theta)} \hat{p}_+(\rho) &< p_+(B_1, B_2). \end{aligned}$$

Therefore, we can reject the hypothesis that $B_1 \cup B_2$ is a Markov substitute set with a probability of false rejection not greater than 2ε whenever

$$\inf_{\rho \in \mathcal{M}_+^1(\Theta)} \hat{p}_-(\rho) \leq \sup_{\rho \in \mathcal{M}_+^1(\Theta)} \hat{p}_+(\rho).$$

In the same way, we can reject the hypothesis that the pair (B_1, B_2) is a γ -weak η -Markov substitute pair of sets with a probability of false rejection not greater than 2ε whenever

$$\begin{aligned} \sup_{\rho \in \mathcal{M}_+^1(\Theta)} \hat{p}_+(\rho) &\geq 1 - \eta, \quad \inf_{\rho \in \mathcal{M}_+^1(\Theta)} \hat{p}_-(\rho) \leq \eta \\ \text{or} \quad \sup_{\rho \in \mathcal{M}_+^1(\Theta)} \hat{p}_+(\rho) - \inf_{\rho \in \mathcal{M}_+^1(\Theta)} \hat{p}_-(\rho) &\geq \gamma. \end{aligned}$$

We hope that these more explicit approximations of the previous test will help the reader getting a better understanding of the structure of the test. The main thing to understand is that hidden in the tighter test of proposition 2.7.3 on page 101 are lower bound estimators for $p_+(B_1, B_2)$ and upper bound estimators for $p_-(B_1, B_2)$, and that these confidence bounds hold uniformly for any $(B_1, B_2) \in \mathcal{B}^2$, at the given confidence level.

We would also like to point out that the estimates $\hat{p}_+(\rho)$ and $\hat{p}_-(\rho)$ contain an explicit empirical variance estimate, that was somehow hidden in the properties of the influence function ψ in the previous section.

2.7.4 Some numerical examples and considerations on the accuracy of the tests

We will here give some numerical examples of the computations presented in the previous subsections. Let us assume for simplicity that $B_1 = \{y_1\}$, $B_2 = \{y_2\}$, and $\pi(\alpha(x, y), x, y) = 1$ for any

$$(x, y) \in \bigcup_{i=1}^n \mathcal{S}(S_i, \{y_1, y_2\}),$$

meaning that there is at most one way to $\{y_1, y_2\}$ -parse any given sentence of the observed sample. In this simple case $(X_{\{y_1, y_2\}, i}, Y_{\{y_1, y_2\}, i})$ is a deterministic function of S_i , so that $a(\theta, i), b(\theta, i) \in \{0, 1\}$. As a consequence, $a_2(\rho) = a_1(\rho) \stackrel{\text{def}}{=} a(\rho)$, $b_2(\rho) = b_1(\rho) \stackrel{\text{def}}{=} b(\rho)$, $\hat{p}_2(\rho) = \hat{p}_1(\rho) \stackrel{\text{def}}{=} \hat{p}(\rho)$ and $\hat{v}(\rho) = \hat{p}(\rho)(1 - \hat{p}(\rho))$. Making as few approximations as possible, we can take

$$\begin{aligned} \hat{p}_+(\rho) &= \hat{p}(\rho) - \frac{y}{1-z} - \sqrt{\frac{zy^2}{(1-z)^2} + \frac{z\hat{v}(\rho)}{1-z}}, \\ \hat{p}_-(\rho) &= \hat{p}(\rho) + \frac{y}{1-z} + \sqrt{\frac{zy^2}{(1-z)^2} + \frac{z\hat{v}(\rho)}{1-z}}. \end{aligned}$$

If we choose μ as mentioned in eq. (2.7.3) on page 98, we get that

$$\delta(\rho) \leq \log \left(k^3 \left(\max_{1 \leq i \leq n} \bar{h}(S_i) \right)^2 \left(\max_{1 \leq i \leq n} \bar{g}(S_i) \right) n^3 |\Lambda| \varepsilon^{-1} \right).$$

If we want to be very conservative, taking for instance values of $n = 10^4$, $k = 10$, $\max_{1 \leq i \leq n} \bar{h}(S_i) = 20$, $\max_{1 \leq i \leq n} \bar{g}(S_i) = 50$, $|\Lambda| = 10 \log(n)$, we get

$$\begin{aligned} \delta(\rho) &\leq \log(10^3 \times 50^3 \times 10^{15} \times \log(10^3)) < 56, \\ \chi &\leq \cosh(0.1) \leq 1.006, \\ y &\leq \frac{2 \times 10 \times 56}{9 \times b(\rho)} \leq \frac{125}{b(\rho)}, \\ z &\leq \frac{2 \times (1.006)^2 \times 56 \times 10^2}{9^2 \times b(\rho)} \leq \frac{140}{b(\rho)}. \end{aligned}$$

Let us compute the minimal gap we have to find between $\hat{p}_+(\rho_+) = \sup \hat{p}_+(\rho)$ and $\hat{p}_-(\rho_-) = \inf \hat{p}_-(\rho)$ to reject the Markov substitute hypothesis when, let us say, $b(\rho_+) = b(\rho_-) = 1000$. We may bound $\hat{v}(\rho)$ by $1/4$. We get that we need

$$\hat{p}(\rho_-) \leq \hat{p}(\rho_+) - 0.71$$

to reject the hypothesis that the pair is a Markov substitute pair, of in other term, we need $a(\rho_-) \leq a(\rho_+) + 710$. That is not a very tight bound.

However, if we want to be more optimistic, we may think that the union bound precaution leading to the mathematically correct value of $\delta(\rho)$ computed above is too shy, and that we should get rid of it and keep only $\delta(\rho) = \log(\varepsilon^{-1})$ to get more realistic numerical values. If we are ready to take this kind of liberty with the theory, and choose as before $\varepsilon = 10^{-2}$, we get $\delta(\rho) \leq 5$, and consequently

$$\begin{aligned} y &\leq \frac{12}{b(\rho)}, \\ z &\leq \frac{13}{b(\rho)}. \end{aligned}$$

With this more optimistic guess, using again for comparison $b(\rho_+) = b(\rho_-) = 1000$, $\hat{v} \leq 1/4$, we obtain a required gap of

$$\hat{p}(\rho_-) \leq \hat{p}(\rho_+) - 0.14.$$

We may also ask in this context what number of observations we need to reject the hypothesis when we can find two contexts (say with the same number of observations as above) such that $\hat{p}(\rho_+) = 1$ and $\hat{p}(\rho_-) = 0$ (so that in both cases $\hat{v} = 0$). We get that we need to observe at least 50 times y_1 and y_2 in disjoint contexts before deciding at a confidence level of $98/100$ that $\{y_1, y_2\}$ is not a Markov substitute pair (still ignoring the union bound factors).

This is still conservative, since we can make a direct computation in this simple situation, showing that if $\{y_1, y_2\}$ is a Markov substitute pair, then the probability

of observing 7 times y_1 in one context and 7 times y_2 in another context is not greater than $2 \times 2^{-7} \leq 0.016$ (the worst case being when the substitute measure is a Bernoulli with parameter $1/2$).

In this context, it is interesting to see whether the non approximated test of proposition 2.7.3 on page 101 is more satisfactory for small sample sizes. In the simplified framework described in this section the test of proposition 2.7.3 on page 101 can be written as

$$B_-(p) = \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} a(\rho)\psi(\lambda(1-p)) + (b(\rho) - a(\rho))\psi(-\lambda p) + \delta(\rho), \quad (2.7.9)$$

$$B_+(p) = \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} a(\rho)\psi(-\lambda(1-p)) + (b(\rho) - a(\rho))\psi(\lambda p) + \delta(\rho). \quad (2.7.10)$$

The test will in particular reject the hypothesis when

$$\min\{B_-(1/2), B_+(1/2)\} > 0. \quad (2.7.11)$$

Let us see what we get if we forget the entropy and union bound factors. Take as suggested before $\delta(\rho) = \log(\varepsilon^{-1}) = \log(100) \simeq 4.6$, $b(\rho_+) = b(\rho_-) = a(\rho_+) = b$, $a(\rho_-) = 0$, and look for the minimum value of b for which eq. 2.7.11 holds. Since $p = 1/2$, $F_{\{y_1\}, \{y_2\}, \theta}(S, 1/2) \in [-1/2, 1/2]$, so in this case we can use up to $\lambda = 2$. With this choice we obtain that

$$\min\{B_+(1/2), B_-(1/2)\} \geq b(\log(2) - k^{-1}) - \delta(\rho) \geq b(\log(2) - 1/10) - \log(100) > 0$$

when $b \geq 8$, and even $b \geq 7$, if we forget the k^{-1} also (that was contributing to make the bound uniform). So the conclusion is that in this low sample size simple situation, the test of proposition 2.7.3 on page 101 is in some sense exact. If we do not compromise the mathematical properties of the test, keep k^{-1} and take $\delta(\rho) = 63$, we get $b \geq 107$, whereas the approximated test gives $b \geq 584$.

So a first conclusion is that for small sample sizes, the non approximated test of proposition 2.7.3 on page 101, although not so intuitive, is noticeably sharper than its more understandable approximation given in the next section.

The above discussion was interested in the difference between $\hat{p}(\rho_-)$ and $\hat{p}(\rho_+)$. We may also want to study the ratio $\hat{p}(\rho)/(1 - \hat{p}(\rho))$.

Let us consider again the non approximated test. Equations (2.7.9) and (2.7.10) on the current page give us that the test will reject the hypothesis as soon as, for

some p , there are two contexts ρ_+ and ρ_- such that

$$\frac{a(\rho_+)}{b(\rho_+) - a(\rho_+)} \geq \frac{-\log(1 - \lambda p) - \frac{\lambda}{k}p + \frac{\delta}{b(\rho_+) - a(\rho_+)}}{\log(1 + \lambda(1 - p)) - \frac{\lambda}{k}(1 - p)},$$

$$\frac{b(\rho_-) - a(\rho_-)}{a(\rho_-)} \geq \frac{-\log(1 - \lambda(1 - p)) - \frac{\lambda}{k}(1 - p) + \frac{\delta}{b(\rho_-) - a(\rho_-)}}{\log(1 + \lambda p) - \frac{\lambda}{k}p}.$$

Once again, we may consider only the case $p = \frac{1}{2}$ for the rejection, which has the advantage of symmetrizing both equations: if we suppose also that $b(\rho_+) = b(\rho_-)$, the hypothesis will be rejected as soon as both ratios are greater than

$$\frac{\hat{p}}{1 - \hat{p}} \geq \frac{-\log\left(1 - \frac{\lambda}{2}\right) - \frac{\lambda}{2k} + \frac{\delta}{b(\rho)}}{\log\left(1 + \frac{\lambda}{2}\right) - \frac{\lambda}{2k}}.$$

We are now left with the optimization of

$$f(x) = \frac{-\log(1 - x) - x/k}{\log(1 + x) - x/k} + \frac{\Delta}{\log(1 + x) - x/k},$$

for $x \in]0, \frac{1}{2}[$, where $\Delta = \delta/b(\rho)$. We may first simplify it by observing that, if $k = 10$, the x/k term may be neglected (at least for the derivative), so $\operatorname{argmin} f$ should not be far from

$$x_m = \operatorname{argmin} \frac{-\log(1 - x)}{\log(1 + x)} + \frac{\Delta}{\log(1 + x)}.$$

At this point, we can remark that the minimum is attained for $x > \frac{1}{2}$ if $\Delta > \frac{1}{2}$ (approximately), so we may take the value on $\lambda = 1$ in this case, which gives the bound

$$\frac{\hat{p}}{1 - \hat{p}} \geq 1.81 + 2.82\Delta.$$

This bound is always true, but gives suboptimal results when $\Delta < \frac{1}{2}$, which is the case when we have more than $2\delta(\rho)$ observations, which is quite probable (remember that a plausible bound for $\delta(\rho)$ is 56).

Using again $b(\rho_+) = b(\rho_-) = 1000$, $\Delta = .065$, and computing the actual minimum of the bound function, we get a bound around 1.65 ($\lambda \approx 0.45$), better than the 2 of the suboptimal bound ($\lambda = 1$).

Another practical question is to choose a reasonable entropy term $\delta(\rho)$. One could argue heuristically that $\delta(\rho)$ should be taken of the form $\log(\varepsilon^{-1}N)$, where N is the number of Markov substitute pairs and contexts actually tested. For instance, if we take $N = 10^6$, we get that a pair that has been seen $b = 32$ times in disjoint contexts can be rejected as a Markov substitute pair with a confidence level of 98/100.

Another way to think about tuning the value of the entropy term $\delta(\rho)$, is that if we lower it, we are sure to both increase the false rejection rate and decrease the false acceptance rate. Thus, if we want to put the stress on minimizing the false acceptance rate rather than the false rejection rate, lowering $\delta(\rho)$ will probably strike a more favorable balance than the theoretical value tuned to obtain a false rejection rate less than 2ε .

On the other hand, the approximations made to the test in proposition 2.7.5 on page 106 are sharp in the large sample size limit, since we kept the leading term unchanged at each step and made approximations only on second order terms.

2.7.5 Probability of false acceptance of the test

In the following, we are going to study the probability of false acceptance of the test of proposition 2.7.3 on page 101. Let us introduce

$$q(B_1, B_2, \theta) = \mathbb{E}(X_{B_1 \cup B_2} \in \theta, Y_{B_1 \cup B_2} \in B_1 \cup B_2).$$

Lemma 2.7.6

For any $p \in [0, 1]$, any $\lambda \in]-1, 1[$, any $B_1, B_2 \in \mathcal{B}$, any $\theta \in \Theta$,

$$\begin{aligned} \mathbb{E} \left[\exp \left(-\psi \left(\lambda F_{B_1, B_2, \theta}(S, p) \right) \right) \right] &\leq \exp \left(q(\theta) \left[\lambda \frac{k-1}{k} (p - p(\theta)) \right. \right. \\ &\quad \left. \left. + \frac{\lambda^2}{1 - |\lambda|} \left(\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2} \right) \left(p(\theta)(1 - p(\theta)) + (p - p(\theta))^2 \right) \right] \right), \end{aligned}$$

where $\phi(z) = 2z^{-2}(\exp(z) - 1 - z)$ is an increasing function such that $\phi(0) = 1$, and where we have put for short $p(B_1, B_2, \theta) = p(\theta)$ and $q(B_1, B_2, \theta) = q(\theta)$. This obviously implies that, with probability at least $1 - \varepsilon$,

$$\begin{aligned} \sum_{i=1}^n \psi \left(\lambda F_{B_1, B_2, \theta}(S, p) \right) &\geq -\log(\varepsilon^{-1}) - nq(\theta) \left[\lambda \frac{k-1}{k} (p - p(\theta)) \right. \\ &\quad \left. + \frac{\lambda^2}{1 - |\lambda|} \left(\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2} \right) \left(p(\theta)(1 - p(\theta)) + (p - p(\theta))^2 \right) \right]. \end{aligned}$$

PROOF. Let us remark first that

$$\begin{aligned} \exp(-\psi(\lambda x)) &= \exp\left(\frac{\lambda}{k}x - \log(1 + \lambda x)\right) \leq \frac{1 + \frac{\lambda}{k}x + \phi\left(\frac{\lambda}{k}\right)\frac{\lambda^2}{2k^2}x^2}{1 + \lambda x} \\ &= 1 - \frac{k-1}{k}\lambda x + \left(\frac{k-1}{k} + \frac{1}{2k^2}\phi\left(\frac{\lambda}{k}\right)\right)\frac{\lambda^2 x^2}{1 + \lambda x} \\ &\leq 1 - \frac{k-1}{k}\lambda x + \left(\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2}\right)\frac{\lambda^2}{1 - \lambda}x^2, \quad \lambda \in [0, 1[, x \in [-1, 1]. \end{aligned}$$

As a consequence, for any $\lambda \in]-1, 1[$,

$$\begin{aligned} \mathbb{E}\left[\exp\left(-\psi\left(\lambda F_{B_1, B_2, \theta}(S, p)\right)\right)\right] &\leq 1 - \frac{k-1}{k}\lambda \mathbb{E}\left(F_{B_1, B_2, \theta}(S, p)\right) \\ &\quad + \left(\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2}\right)\frac{\lambda^2}{1 - |\lambda|}\mathbb{E}\left(F_{B_1, B_2, \theta}(S, p)^2\right) \\ &\leq \exp\left(-\frac{k-1}{k}\lambda \mathbb{E}\left(F_{B_1, B_2, \theta}(S, p)\right)\right) \\ &\quad + \left(\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2}\right)\frac{\lambda^2}{1 - |\lambda|}\mathbb{E}\left(F_{B_1, B_2, \theta}(S, p)^2\right). \end{aligned}$$

The rest follows easily from

$$\begin{aligned} \mathbb{E}\left(F_{B_1, B_2, \theta}(p)\right) &= (p(\theta) - p)q(\theta) \\ \mathbb{E}\left(F_{B_1, B_2, \theta}(p)^2\right) &= \left(p(\theta)(1 - p(\theta)) + (p - p(\theta))^2\right)q(\theta). \quad \square \end{aligned}$$

We can now state the following proposition:

Proposition 2.7.7

Let us put

$$\begin{aligned} \delta &= \frac{1}{n} \log \left[k^3 n^3 \left(\max_{1 \leq i \leq n} \bar{h}(S_i) \right) \left(\max_{1 \leq i \leq n} \bar{g}(S_i) \right) |\Lambda| \varepsilon^{-2} \right], \\ \chi &= \sup_{x \in [(2n)^{-1/2}, (2n)^{1/2}]} \inf_{\lambda \in \Lambda} \cosh \left[\log \left(\frac{\lambda x}{(1 - \lambda)} \right) \right], \\ a &= \frac{4\chi^2 k}{(k-1)} \left(1 + \frac{\varphi(k^{-1})}{2k(k-1)} \right) \leq 4.47\chi^2 \text{ when } k = 10, \\ b &= \frac{(2 + \sqrt{2})k}{k-1} \leq 3.8 \text{ when } k = 10. \end{aligned}$$

Let us assume that there are two sets $B_1, B_2 \in \mathcal{B}$ and two contexts $\theta_+, \theta_- \in \Theta$ such that $\bar{p}_\epsilon = p(B_1, B_2, \theta_\epsilon)$ and $q_\epsilon = q(B_1, B_2, \theta_\epsilon)$, $\epsilon = \pm 1$, satisfy

$$\begin{aligned} \min\{q_-, q_+\} &\geq \frac{16k\chi^2\delta}{(k-1)}, & \text{and} \\ \bar{p}_+ - \bar{p}_- &\geq \sqrt{\frac{a\bar{p}_+(1-\bar{p}_+)\delta}{q_+} \left(1 + \frac{a\delta}{q_+}\right) + \frac{b\delta}{q_+}} \\ &\quad + \sqrt{\frac{a\bar{p}_-(1-\bar{p}_-)\delta}{q_-} \left(1 + \frac{a\delta}{q_-}\right) + \frac{b\delta}{q_-}}. \end{aligned} \quad (2.7.12)$$

Then with probability at least $1 - 2\varepsilon$

$$\inf_{p \in [0,1]} \max\{B_-(p), B_+(p)\} > 0.$$

Therefore, the test defined by eq. (2.7.5) on page 101 has in this case a probability of false acceptance of the pair $(B_1, B_2) \in \mathcal{B}^2$ as a Markov substitute pair of sets at most equal to 2ε .

More precisely, with probability at least $1 - 2\varepsilon$,

$$B_- \left(\bar{p}_+ - \sqrt{\frac{a\bar{p}_+(1-\bar{p}_+)\delta}{q_+} \left(1 + \frac{a\delta}{q_+}\right) - \frac{b\delta}{q_+}} \right) > 0, \quad (2.7.13)$$

$$B_+ \left(\bar{p}_- + \sqrt{\frac{a\bar{p}_-(1-\bar{p}_-)\delta}{q_-} \left(1 + \frac{a\delta}{q_-}\right) + \frac{b\delta}{q_-}} \right) > 0. \quad (2.7.14)$$

If we assume now in place of eq. (2.7.12) on the current page that

$$\begin{aligned} \bar{p}_+ - 1 + \eta &\geq \sqrt{\frac{a\bar{p}_+(1-\bar{p}_+)\delta}{q_+} \left(1 + \frac{a\delta}{q_+}\right) + \frac{b\delta}{q_+}}, \\ \text{or that } \eta - \bar{p}_- &\geq \sqrt{\frac{a\bar{p}_-(1-\bar{p}_-)\delta}{q_-} \left(1 + \frac{a\delta}{q_-}\right) + \frac{b\delta}{q_-}}, \\ \text{or that } \bar{p}_+ - \bar{p}_- &\geq \gamma + \sqrt{\frac{a\bar{p}_+(1-\bar{p}_+)\delta}{q_+} \left(1 + \frac{a\delta}{q_+}\right) + \frac{b\delta}{q_+}} \\ &\quad + \sqrt{\frac{a\bar{p}_-(1-\bar{p}_-)\delta}{q_-} \left(1 + \frac{a\delta}{q_-}\right) + \frac{b\delta}{q_-}}. \end{aligned}$$

the false acceptance probability of the test (defined in corollary 2.7.4 on page 102) that $(B_1, B_2) \in \mathcal{B}^2$ is an γ -weak η -Markov substitute pair of sets will be not greater than 2ε .

PROOF. The proof is very similar to the one of proposition 2.6.3 on page 92.

Let us put

$$p = \bar{p}_+ - \sqrt{\frac{a\bar{p}_+(1-\bar{p}_+)\delta}{q_+}} \left(1 + \frac{a\delta}{q_+}\right) - \frac{b\delta}{q_+}.$$

We will here choose $\lambda \in \Lambda$ such that

$$\lambda \in \arg \min_{\zeta \in \Lambda} \cosh \left[\log \left(\frac{\zeta}{1-\zeta} \sqrt{\left[\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2} \right] \frac{[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]q_+}{\delta}} \right) \right].$$

Once again, since $\bar{p}_+ - p \geq \frac{\delta}{q_+}$,

$$\sqrt{\frac{1}{n}} \leq \sqrt{\frac{\delta}{q_+}} \leq \sqrt{\frac{[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]q_+}{\delta}} \leq \sqrt{\frac{q_+}{\delta}} \leq \sqrt{n},$$

which implies, together with the fact that $\frac{1}{2} \leq \left[\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2} \right] \leq 2$, that

$$\cosh \left[\log \left(\frac{\lambda}{1-\lambda} \sqrt{\left[\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2} \right] \frac{[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]q_+}{\delta}} \right) \right] \leq \chi.$$

Now, we have that

$$B_-(p) \geq \sum_{i=1}^n \psi \left(\lambda F_{B_1, B_2, \theta_+}(S_i, p) \right) - \log \left[k^3 |\Lambda| \left(\max_{1 \leq i \leq n} \bar{h}(S_i) \right)^2 \left(\max_{1 \leq i \leq n} \bar{g}(S_i) \right) n^3 \varepsilon^{-1} \right]$$

Together with lemma 2.7.6 on page 111, we get that, with probability at least $1 - \varepsilon$

$$B_-(p) \geq n\lambda q_+ \left(\frac{k-1}{k} (\bar{p}_+ - p) - \frac{\lambda}{1-\lambda} \left[\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2} \right] \left(\bar{p}_+(1-\bar{p}_+) + (p - \bar{p}_+)^2 \right) - \frac{\delta}{\lambda q_+} \right).$$

We obtain then that with probability at least $1 - \varepsilon$,

$$\begin{aligned}
B_-(p) &> nq_+\lambda \left\{ \frac{k-1}{k}(\bar{p}_+ - p) - \frac{\delta}{q_+} \right. \\
&\quad \left. - 2\sqrt{\left[\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2} \right] \frac{[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]\delta}{q_+}} \right. \\
&\quad \left. \times \cosh \left[\log \left(\frac{\lambda}{1-\lambda} \sqrt{\left[\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2} \right] \frac{[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]q_+}{\delta}} \right) \right] \right\} \\
&\geq nq_+\lambda \left(\frac{k-1}{k}(\bar{p}_+ - p) - \frac{\delta}{q_+} \right. \\
&\quad \left. - 2\chi \sqrt{\left[\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2} \right] \frac{[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]\delta}{q_+}} \right).
\end{aligned}$$

Let us put $\alpha = \frac{k-1}{k}$ and $\rho = \frac{4\chi^2\delta}{q_+} \left(\alpha + \frac{\phi(k^{-1})}{2k^2} \right)$, $v = \bar{p}_+(1-\bar{p}_+)$, $y = \frac{\delta}{q_+}$ and $x = \bar{p}_+ - p$. The sign of the right-hand side of the previous inequality is the same as the sign of

$$\begin{aligned}
&\left(\frac{k-1}{k}(\bar{p}_+ - p) - \frac{\delta}{q_+} \right)^2 - \left[\frac{k-1}{k} + \frac{\phi(k^{-1})}{2k^2} \right] \frac{4\chi^2\delta[\bar{p}_+(1-\bar{p}_+) + (\bar{p}_+ - p)^2]}{q_+} \\
&\quad = \alpha^2 x^2 - 2\alpha y x + y^2 - \rho v - \rho x^2 \\
&\quad = (\alpha^2 - \rho)x^2 - 2\alpha y x - (\rho v - y^2).
\end{aligned}$$

This quantity is non negative when x is not less than

$$\begin{aligned}
\frac{\alpha y}{\alpha^2 - \rho} + \sqrt{\frac{\alpha^2 y^2}{(\alpha^2 - \rho)^2} + \frac{\rho v - y^2}{\alpha^2 - \rho}} &= \frac{\alpha y}{\alpha^2 - \rho} + \sqrt{\frac{\rho y^2}{(\alpha^2 - \rho)^2} + \frac{\rho v}{\alpha^2 - \rho}} \\
&\leq \sqrt{\frac{\rho v}{\alpha^2 - \rho}} + \frac{(\alpha + \sqrt{\rho})y}{\alpha^2 - \rho} \leq \sqrt{\frac{\rho v}{\alpha^2}} \left(1 + \frac{\rho}{2\alpha^2(1 - \rho\alpha^{-2})} \right) + \frac{(\alpha + \sqrt{\rho})y}{\alpha^2 - \rho} \\
&\leq \sqrt{\frac{\rho v}{\alpha^2}} \left(1 + \frac{\rho}{\alpha^2} \right) + \frac{(2 + \sqrt{2})y}{\alpha},
\end{aligned}$$

where we have used the inequalities $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ and $\rho \leq \alpha^2/2$.

Thus we obtain with probability at least $1 - \varepsilon$ that $B_-(p) > 0$ when

$$\bar{p}_+ - p \geq \sqrt{\frac{\rho v}{\alpha^2}} \left(1 + \frac{\rho}{\alpha^2} \right) + \frac{(2 + \sqrt{2})\delta}{\alpha q_+}.$$

Since $\frac{\rho}{\alpha^2} = \frac{a\delta}{q_+}$ and $\frac{(2 + \sqrt{2})}{\alpha} = b$, this proves in particular that with probability at least $1 - \varepsilon$,

$$B_- \left(\bar{p}_+ - \sqrt{\frac{a\bar{p}_+(1 - \bar{p}_+)\delta}{q_+}} \left(1 + \frac{a\delta}{q_+} \right) - \frac{b\delta}{q_+} \right) > 0,$$

The proof of eq. (2.7.14) on page 113 is obtained exactly in the same way, and the rest of the proposition is a straightforward consequence of eqs. (2.7.13) and (2.7.14) on page 113. \square

2.8 Estimation of the substitute measure using the results of the test

The family of tests we defined here may be used to get an estimation of the substitute measure q_B of a set that passed the test. More precisely, we can get an estimation \hat{p} of p such that $\mathbb{E}(F(X, Y, p)) = 0$. These p are usually linked to the substitute measures, $p = q_{\{y_1, y_2\}}(y_1)$ in a pair test, or $p = q_{B_1 \cup B_2}(B_1)$ in the union test, for example. Note that in the second example, the symmetric test giving $q_{B_1 \cup B_2}(B_2)$ is required to have the substitute measure (and gives us an indication of the size of the intersection $q_{B_1 \cup B_2}(B_1 \cap B_2)$ as a bonus).

Let us remark that all the tests we defined were based on two functions $B_+(p)$ and $B_-(p)$, one non-decreasing and one non-increasing, that could be written as

$$B_{\pm}(p) = \tilde{B}_{\pm}(p) + \log(\varepsilon),$$

where $\tilde{B}_{\pm}(p)$ did not depend on ε . Note that while the functions B_{\pm} were expressed as a supremum on $\rho(\theta)$ and λ , we can simply take in this case the supremum on λ , and $\rho = \delta_{(D^*)^2}$ to maximize the size of the sample. The analysis told us that if the set tested was a Markov substitute set with parameter p , with probability greater than $1 - 2\varepsilon$,

$$B_{\pm}(p) < 0,$$

or equivalently,

$$\tilde{B}_{\pm}(p) < \log(\varepsilon^{-1}) = \ell.$$

Thus, if we fix ℓ , we have, as before, a test $\max(B_+(p), B_-(p)) > 0$. However, we could see this in another light. Indeed, any q such that $B_-(q) > 0$ is smaller than p , or bigger if $B_+(q) > 0$. Studying the graph of B_{\pm} gives us then a confidence interval for p , if we fix ℓ . The middle of the confidence interval is a possible estimator. However, there is another one we can consider.

If we do not fix ℓ , we can consider the value \hat{p} such that $\tilde{B}_-(\hat{p}) = \tilde{B}_+(\hat{p})$. If we put $\ell_0 = \tilde{B}_+(\hat{p})$, we would reject the hypothesis if we wanted a confidence of more than $1 - 2e^{\ell_0}$. As such, the parameter \hat{p} is the only parameter that is in all confidence intervals, for any confidence level that accepts the pair as a Markov substitute set.

We could also use ℓ_0 to sort different candidates, the lower ℓ_0 , the most likely they are to be Markov substitute sets.

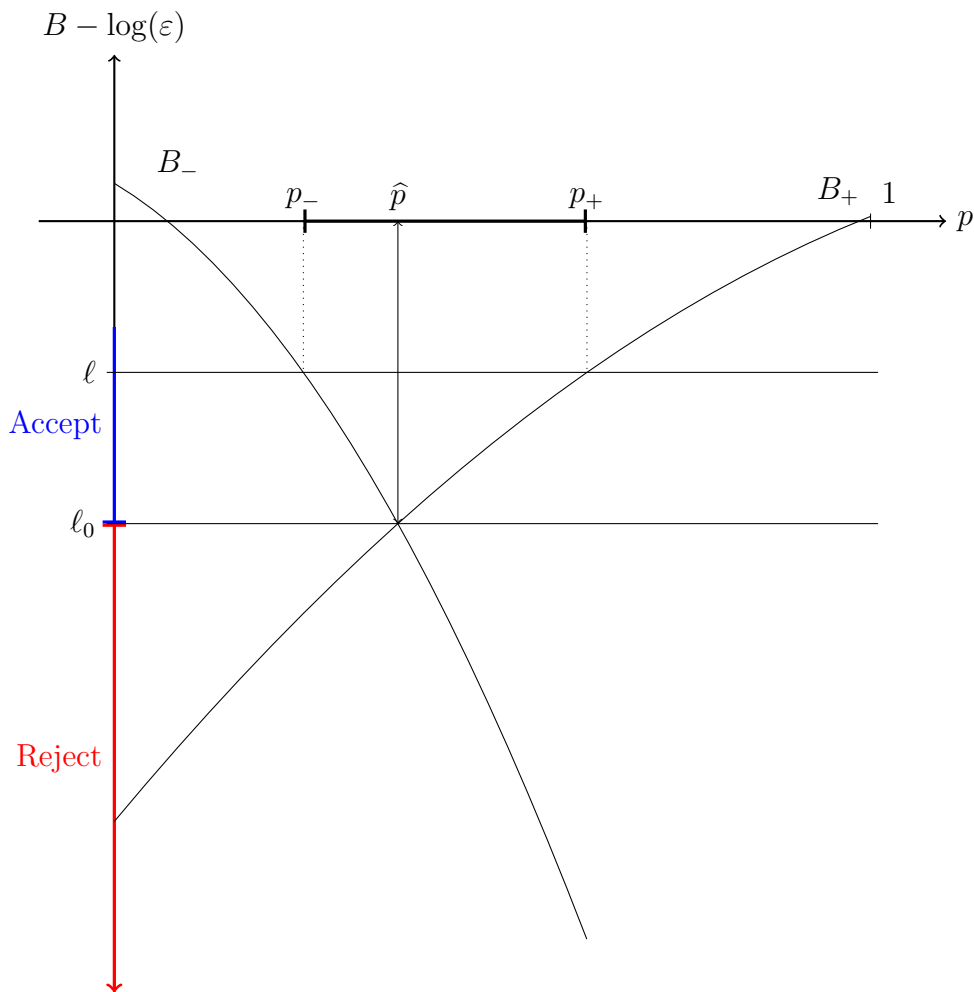


Figure 2.2: Test of the Markov substitute property.

Chapter 3

Markov substitute sets and language

3.1 Production rules and Markov substitute sets

3.1.1 Markov grammars

Consider a set of context-free production rules \mathcal{R} whose terminal symbols belong to an alphabet D and whose non terminal symbols belong to $N = \{]_i, i \in \mathbb{N} \setminus \{0\} \}$. Let us assume that each rule is of the form $]_i \rightarrow e$, where $i \in \mathbb{N} \setminus \{0\}$ and the expression $e \in (D \cup N)^+$, so that no rule uses the empty string.

Production rules alone do not suffice to define a complete context-free grammar, we also need a start symbol. We will here consider multiple start symbols, defining each a different context-free grammar. This approach will enable us to study in a broader sense the properties of our model, and move outside of regular context-free grammars.

Definition 3.1.1

For any set of context-free production rules (ruleset for short) \mathcal{R} , and any $i \in \mathbb{N}$, let \mathcal{R}_i be the context-free grammar whose set of rules is \mathcal{R} and whose start symbol is $]_i$.

For any non empty expression $e \in (D \cup N)^+$, let \mathcal{R}_e be the context-free grammar whose start symbol is $]_0$ and whose set of rules is $\mathcal{R} \cup \{]_0 \rightarrow e \}$.

For any $i \in \mathbb{N} \setminus \{0\}$, let B_i be the language generated by \mathcal{R}_i , and for any expression $e \in (D \cup N)^+$, let B_e be the language generated by \mathcal{R}_e . Let us remark that B_i and B_e may be empty, and that this will be the case except for a finite set of values of i and for the expressions e using only the corresponding non terminal symbols. Elements of the set B_i will be called syntagms of type i .

In the following, we will use for the rule $]_i \rightarrow e$ the more compact notation $]_i e$, that we already used in chapter 1.

Let us remark that for any $i \in \mathbb{N} \setminus \{0\}$,

$$B_i = B_{]_i} = \bigcup_{e,]_i e \in \mathcal{R}} B_e,$$

but that this union is not necessarily a disjoint union.

The context-free property ensures that any elements of a B_i may appear in the same contexts. We will here strengthen this property.

Definition 3.1.2

A ruleset will be called a Markov ruleset if, for any $i \in \mathbb{N} \setminus \{0\}$, B_i is a Markov substitute set.

Lemma 3.1.1

For any Markov ruleset, for any $e \in S^+$, such that $\mathbb{P}(\mathcal{S}(S, B_e) \neq \emptyset) > 0$, B_e is a Markov substitute set.

PROOF. If we write $e = (w_k, 1 \leq k \leq \ell)$, we can write B_e as

$$B_e = \gamma(B_{w_1}, \dots, B_{w_\ell}),$$

and each B_{w_k} is either a singleton (if $w_k \in D$) or $B_{]_i}$, so in any case a Markov substitute set by assumption, so that B_e is also a Markov substitute set, according to proposition 2.1.3 on page 67. \square

The substitute measure of a syntagm set B_e is given by the following lemma.

Lemma 3.1.2

Consider any non terminal expression $e = w_{1:k} \in (D \cup N)^k$ of length k , such that

$$\mathbb{P}(\mathcal{S}(S, B_e) \neq \emptyset) > 0.$$

Consider any string $\gamma(y_{1:k}) \in B_e$, where γ is the concatenation operator, such that $y_j \in B_{w_j}$, $1 \leq j \leq k$. The substitute measure of B_e can be expressed as

$$q_{B_e}(\gamma(y_{1:k})) = A_e \prod_{j=1}^k q_{B_{w_j}}(y_j),$$

where the constant A_e is such that for any functions f_j , $1 \leq j \leq k$,

$$\sum_{y \in B_e} q_{B_e}(y) \sum_{y_j \in B_{w_j}, 1 \leq j \leq k} \mathbf{1}(y = \gamma(y_{1:k})) \prod_{j=1}^k f_j(y_j) = A_e \prod_{j=1}^k \left(\sum_{y_j \in B_{w_j}} q_{B_{w_j}}(y_j) f_j(y_j) \right).$$

Let us remark that the constant A_e may indeed be different from one, as in the case when $e =]_1]_2$, where $B_1 = \{a, ab\}$ and $B_2 = \{bc, c\}$. In this case $B_e = \gamma(B_1, B_2) = \{ac, abc, abbc\}$ has only three elements and not four, so that $\sum_{z \in B_e} p(z) = \sum_{z_1 \in B_1, z_2 \in B_2} p(z_1 z_2) - p(abc)$. However, $A_{j_i} = 1$ for all i .

PROOF. Let us choose $x \in (D^*)^2$ such that $\mathbb{P}_S(\alpha(x, B_e)) > 0$.

For any $y_{1:k} \in B_{w_1} \times \cdots \times B_{w_k}$,

$$\begin{aligned} \mathbb{P}\left(S = \alpha(x, \gamma(y_{1:k}))\right) &= \left(\sum_{y'_k \in B_{w_k}} \mathbb{P}_S\left(\alpha(x, \gamma(y_{1:k-1}, y'_k))\right) \right) q_{B_{w_k}}(y_k) \\ &\quad \vdots \\ &= C_x \bar{q}(y_{1:k}). \end{aligned}$$

where

$$\begin{aligned} C_x &= \left(\sum_{y'_{1:k} \in \prod_{j=1}^k B_{w_j}} \mathbb{P}_S\left(\alpha(x, \gamma(y'_{1:k}))\right) \right) \geq \mathbb{P}_S(\alpha(x, B_e)) > 0, \\ \text{and } \bar{q}(y_{1:k}) &= \left(\prod_{j=1}^k q_{B_{w_j}}(y_j) \right). \end{aligned}$$

This shows that $\bar{q}(y_{1:k})$ depends in fact only on $\gamma(y_{1:k})$ and that $q_{B_e}(\gamma(y_{1:k}))$ is proportional to $\bar{q}(y_{1:k})$, so that for some constant A_e , $q_{B_e} = A_e \bar{q}$.

We can then write for any test functions f_j , $1 \leq j \leq k$,

$$\begin{aligned} &\sum_{y \in B_e} q_{B_e}(y) \sum_{y_j \in B_{w_j}, 1 \leq j \leq k} \mathbf{1}(y = \gamma(y_{1:k})) \prod_{j=1}^k f_j(y_j) \\ &= \sum_{y_{1:k} \in \prod_{j=1}^k B_{w_j}} q_{B_e}(\gamma(y_{1:k})) \prod_{j=1}^k f_j(y_j) = \sum_{y_{1:k} \in \prod_{j=1}^k B_{w_j}} A_e \bar{q}(y_{1:k}) \prod_{j=1}^k f_j(y_j) \\ &= A_e \left(\prod_{j=1}^k \sum_{y_j \in B_{w_j}} q_{B_{w_j}}(y_j) f_j(y_j) \right). \quad \square \end{aligned}$$

This lemma will be used in section 3.2.2 on page 127 to build an estimator of the substitute measures, but we need to introduce first some other notions.

3.1.2 Parse trees

This definition of grammars and language gives production and parsing methods, but the structure of a string is not readily apparent. To manipulate this structure, we will introduce the notion of parse trees.

Definition 3.1.3

For any ruleset \mathcal{R} , let $\check{\mathcal{R}}$ be the ruleset

$$\check{\mathcal{R}} = \{[_i(i e)_i, [_i e \in \mathcal{R}]\}.$$

Let $\check{D} = D \cup \{(i,)_i, i \in \mathbb{N}\}$ be the related extended set of terminal symbols.

Similarly to definition 3.1.1 on page 119, let $\check{\mathcal{R}}_i$ be the context-free grammar with start symbol $[_i$, terminal symbol set \check{D} and set of rules $\check{\mathcal{R}}$. In the same way, let $\check{\mathcal{R}}_e$ be the grammar with start symbol $[_0$, terminal symbol set \check{D} and set of rules $\check{\mathcal{R}} \cup \{[_0 e\}$.

Let \check{B}_i be the language generated by $\check{\mathcal{R}}_i$.

We can in the same way define \check{B}_e , the language generated by $\check{\mathcal{R}}_e$, for any expression $e \in (D \cup N)^+$. Let us remark that, again, $\check{B}_j = \check{B}_{]_j}$.

Let $\mathcal{T} = \bigcup_{e \in (D \cup N)^+} \check{B}_e$ the set of parse trees. When $e \neq e'$, \check{B}_e and $\check{B}_{e'}$ are disjoint. We can therefore, for any parse tree $t \in \mathcal{T}$ define its surface structure $\varsigma(t) \in (D \cup N)^+$ by the condition $t \in \check{B}_{\varsigma(t)}$. It is obtained by replacing each outer pair $(\cdot \cdot \cdot)_j$ of matched parentheses in t and its content by $]_j$.

Members of a particular \check{B}_i will be called parse trees of type i .

The realization of a parse tree is given by the operator $\varphi : \mathcal{T} \rightarrow D^+$ which removes parentheses.

We will say that a tree $t \in \mathcal{T}$ is a parse of another t' (noted $t \succcurlyeq t'$), if they have same realizations $\varphi(t) = \varphi(t')$, and all parentheses in t' are also in t .

The realization operator is such that $\varphi(\check{B}_i) = B_i$. It is however not one to one, since a given string may have different parses.

These new grammars will be called for short tree grammars, and will be used to describe the parse trees. Parse trees of a given type describe the parses of the syntagms of same type.

Let us remark that for any $t \in \check{B}_i$, $\varsigma(t) =]_i$, and $t = (i y)_i$, where $y \in \mathcal{T}$ is such that $[_i \varsigma(y) \in \mathcal{R}$.

In the same way, for any $t \in \mathcal{T}$, and any $x \in (\check{D}^*)^2$, $y \in \mathcal{T}$ and $j \in \mathbb{N} \setminus \{0\}$, such that $t = \alpha(x, (j y)_j)$, $[_j \varsigma(y) \in \mathcal{R}$. This is because the condition $y \in \mathcal{T}$ imposes that in the expression $(j y)_j$ the outer parentheses $(j)_j$ are matched.

3.1.3 General introduction to parsing

The next step is to construct parsing trees from a string.

Definition 3.1.4

A parsing kernel is a probability kernel $\Pi : D^+ \mapsto \mathcal{T}$ such that

$$\text{supp}(\Pi(s, \cdot)) \subset \varphi^{-1}(s).$$

A parsing kernel for a certain syntagm set B_e is a probability kernel $\Pi_e : B_e \mapsto \check{B}_e$ such that

$$\text{supp}(\Pi_e(s, \cdot)) = \varphi^{-1}(s) \cup \check{B}_e.$$

Parsing kernels can be constructed using classical parsing algorithms for context-free grammars such as CYK, as discussed in section 3.B on page 159. They provide a way to build the splitting kernels π used in sections 2.2.1 and 2.2.5 on page 69 and on page 73, by simply drawing at random one or two matching pairs of parentheses in the parsed sentence or text.

A parsing kernel builds a tree from a sentence that is compatible with a ruleset. This operation does not lose any information, and adds a syntactic structure to the sentence. However, the manipulation of natural language corpora may be quite difficult due to its sheer size. For this reason, it may be interesting to find some sort of process that can reduce the size of sentences, while losing as little information as possible. One way to do this is to actually only keep the syntactic structure of a given sentence, by considering the surface structure of a parse of the sentence.

Definition 3.1.5

A parsing reduction kernel is an kernel r from D^+ to S^+ such that

$$r(s, e) > 0 \implies s \in B_e.$$

Given a parsing kernel Π , we can define a parsing reduction kernel r as

$$r(s, e) = \sum_{t:\varsigma(t)=e} \Pi(s, t).$$

In order to implement the reversible dynamics we have defined, we need to have Markov kernels for which both simulating transitions and computing transition probabilities are fast enough. A detailed description of these two operations is provided in section 3.D on page 165. In the following sections, we will use parsing kernels to define the splitting kernel of our Metropolis reversible dynamics. While simulation and computation are both easy when parsing in a given syntagm set, the general parse is more tricky. We will have to access the history of the parse,

and manipulate strings of syntactic trees. This section introduces the required definitions to achieve this, while the actual parsing schemes, based on the CYK algorithm, are described in section 3.B on page 159.

Let us note b and e the begin and end operators on strings, $b(s_{1:n}) = s_1$, $e(s_{1:n}) = s_n$.

We will build the parse using extension kernels, that add one pair of parentheses at a time.

Definition 3.1.6

A parsing extension is a kernel

$$\{\varpi(t, t'), t, t' \in \mathcal{T}, t' \succ t\}.$$

A parsing extension relative to a reference tree T is a kernel

$$\{\rho_T(t, t'), t, t' \in \mathcal{T}, T \succ t' \succ t\}.$$

We will usually consider bottom-up parsing extensions, that is, that add one pair of outer parentheses at a time,

$$\left\{ \varpi(t, t') : t, t' \in \mathcal{T}, \exists i \in \mathbb{N}, x \in \mathcal{T}, a \in S^+, [{}_i a \in \mathcal{R}, y \in \check{B}_a, \right. \\ \left. t = \alpha(x, y), t' = \alpha(x, ({}_i y)_i) \right\}.$$

Parsing extensions define chains of trees, the end of which is our final parse.

Definition 3.1.7

Any parsing extension ϖ defines a stopped Markov chain T_n on parse trees, its stopping time being

$$\tau = \inf\{n, \text{supp}(\varpi(T_n, \cdot)) = \emptyset\}.$$

A general parsing kernel will then be a kernel $\bar{\Pi}$ from syntactic trees (in practice, strings), to sequences of parsing trees.

$$\bar{\Pi}(t, \tilde{t}) = \mathbf{1}(\tilde{t}_0 = t) \mathbf{1}(\text{supp}(\varpi(e(\tilde{t}), \cdot)) = \emptyset) \prod_{i=0}^{\ell(\tilde{t})-1} \varpi(\tilde{t}_i, \tilde{t}_{i+1}), \quad \tilde{t} \in \mathcal{T}^+.$$

The “final” parse being, of course, $e(\tilde{t})$.

In the case of parsing extensions relative to a tree, we define the kernels

$$\zeta(t, \tilde{t}) = \mathbf{1}(e(\tilde{t}) = t) \mathbf{1}(b(\tilde{t}) = \varphi(t)) \prod_{i=0}^{\ell(\tilde{t})-1} \rho_t(\tilde{t}_i, \tilde{t}_{i+1}), \quad \tilde{t} \in \mathcal{T}^+.$$

These parsings will be used to define reversible dynamics as described in section 3.D on page 165. This method supposes that we have, for each kernel π used in the simulation, a “symmetric” kernel π' , so that the quotient $\frac{\pi(t,t')}{\pi'(t',t)}$ is as close to 1 as possible. In this case, we start from a string s , and build a sequence of increasingly bigger parses \tilde{t}_i using a parsing extension. Conversely, we want, from a parsing tree t , to build a sequence of smaller parsing trees \tilde{t}_i , ending at $\varphi(t)$.

Building “parsing reduction kernels”

$$\{\rho(t, t'), t, t' \in \mathcal{T}, t \succ t'\}$$

that remove parentheses from an initial tree would be an intuitive way to achieve this, but we can remark that since we know the end part of the tree sequence ($\varphi(t)$), we can actually begin from the end result, and use regular parsing extensions, the only additional requirement being that we finish at the fixed t . This is the motivation behind the definition of parsing extensions relative to a tree.

Actual construction of these parsing extensions, and discussion of the quotient $\frac{\varpi}{\rho}$ will be found in section 3.B on page 159.

3.2 Substitute measures and reversible dynamics

3.2.1 Parametrization of substitute measures

Of course, Markov substitute sets are mostly interesting in conjunction with their substitute measures q_B . The proof of lemma 3.1.1 on page 120 hints that this requires to compute q_{B_i} . We will now see how to perform this computation from the estimation of a few parameters — one for each production rule — using a parse tree structure.

Lemma 3.2.1

For any $s \in D^+$, and any $t \in \mathcal{T}$ such that $s = \varphi(t)$,

$$q_{B_{\varsigma(t)}}(s) = A_{\varsigma(t)} \prod_{[j \in \mathcal{R}]} [A_e q_{B_j}(B_e)]^{\chi(t, j, e)},$$

where

$$\chi(t, j, e) = \sum_{x \in (\check{D}^*)^2} \sum_{y \in \mathcal{T}} \mathbb{1}[t = \alpha(x, (jy)_j)] \mathbb{1}(\varsigma(y) = e). \quad (3.2.1)$$

PROOF. The result is equivalent to the fact that for any $t \in \mathcal{T}$

$$q_{B_{\varsigma(t)}}(\varphi(t)) = A_{\varsigma(t)} \prod_{[j \in \mathcal{R}]} [A_e q_{B_j}(B_e)]^{\chi(t, j, e)},$$

We will prove it by induction on the length of t . If $\ell(t) = 1$, then $t \in D$, $\varsigma(t) = \varphi(t) = t$, and $q_{B_{\varsigma(t)}}(\varphi(t)) = q_{\{t\}}(t) = 1$, so that the result holds in this case.

Let $\varsigma(t) = \gamma(w_{1:\ell}) \in (D \cup N)^\ell$. We can write the tree t as $t = \gamma(t_{1:\ell})$, where $t_k \in \tilde{B}_{w_k}$. With these notations, $w_k = \varsigma(t_k)$, $\varphi(t) = \gamma(\varphi(t_{1:\ell}))$, and $\varphi(t_k) \in B_{\varsigma(t_k)}$. According to lemma 3.1.2 on page 120 this implies that

$$q_{B_{\varsigma(t)}}(\varphi(t)) = \prod_{k=1}^{\ell} A_{\varsigma(t)} q_{B_{\varsigma(t_k)}}(\varphi(t_k)). \quad (3.2.2)$$

If $w_k = \varsigma(t_k) \in D$, then $t_k = w_k$, and as already seen, $q_{B_{\varsigma(t_k)}}(\varphi(t_k)) = 1$.

In the case when $w_k \in N$, $\varsigma(t_k) =]_i$ and $t_k = ({}_i y_k)_i$, where $[_i \varsigma(y_k) \in \mathcal{R}$. Consequently $\varphi(t_k) = \varphi(y_k) \in B_{\varsigma(y_k)} \subset B_{\varsigma(t_k)}$. Therefore

$$q_{B_{\varsigma(t_k)}}(\varphi(t_k)) = q_{B_{\varsigma(t_k)}}(B_{\varsigma(y_k)}) q_{B_{\varsigma(y_k)}}(\varphi(y_k)).$$

Moreover, since $\varsigma(t_k)$ is of length one, $A_{\varsigma(t_k)} = 1$. Since $\ell(y_k) \leq \ell(t) - 2 < \ell(t)$, we can now apply the induction hypothesis to y_k , to get

$$\begin{aligned} q_{B_{\varsigma(t_k)}}(\varphi(t_k)) &= A_{\varsigma(y_k)} q_{B_{\varsigma(t_k)}}(B_{\varsigma(y_k)}) \prod_{[_j e \in \mathcal{R}} A_e q_{B_j}(B_e) \chi(y_k, j, e) + \mathbb{1}(j = i, e = \varsigma(y_k)) \\ &= \prod_{[_j e \in \mathcal{R}} A_e q_{B_j}(B_e) \chi(t_k, j, e). \end{aligned}$$

This proves the result for t_k , $1 \leq k \leq \ell$, (since $A_{\varsigma(t_k)} = 1$) and implies that

$$q_{B_{\varsigma(t)}}(\varphi(t)) = A_{\varsigma(t)} \prod_{[_j e \in \mathcal{R}} [A_e q_{B_j}(B_e)]^{\sum_{k=1}^{\ell} \chi(t_k, j, e)},$$

proving the result for t , since

$$\chi(t, j, e) = \sum_{k=1}^{\ell} \chi(t_k, j, e). \quad \square$$

This lemma means that we only need to know the parameters $q_{B_j}(B_e)$, and A_e , for $[_j e \in \mathcal{R}$, to know all substitute measures. These parameters can effectively be stored with the ruleset by adding weights to each element $[_j e \in \mathcal{R}$.

Let us remark that the knowledge of these measures reduces the estimation of the distribution \mathbb{P}_S to the estimation of the probability of syntagm sets according to the relation

$$\mathbb{P}_S(s) = \mathbb{P}(S \in B_{\varsigma(t)}) q_{B_{\varsigma(t)}}(s), \quad s \in \varphi(t).$$

3.2.2 Estimating substitute measures

We saw that estimating the substitute measures q_{B_i} can be performed through the estimation of the family of coefficients $(A_e, q_{B_i}(B_e), i \in \mathbb{N} \setminus \{0\}, [{}_i e \in \mathcal{R}_i)$.

Let us start with the estimation of the constants A_e , where $e = w_{1:k} \in (D \cup N)^k$. We will estimate separately

$$N_e = \sum_{y \in B_e} q_{B_e}(y) \sum_{y_j \in B_{w_j}, 1 \leq j \leq k} \mathbb{1}(y = \gamma(y_{1:k})) \prod_{j=1}^k f(y_j),$$

where $f(y) = r(1-r)^{\ell(y)-1}$, $y \in D^+$, for some parameter $r \in (0, 1)$, and each

$$D_{w_j} = \sum_{y_j \in B_{w_j}} q_{B_{w_j}}(y_j) f(y_j),$$

to estimate in the end A_e using the formula

$$A_e = \frac{N_e}{\prod_{j=1}^k D_{w_j}}.$$

Let us remark first that for any $y \in B_e$,

$$\begin{aligned} \sum_{y_j \in B_{w_j}, 1 \leq j \leq k} \mathbb{1}(y = \gamma(y_{1:k})) \prod_{j=1}^k f(y_j) &= \sum_{y_j \in B_{w_j}, 1 \leq j \leq k} \mathbb{1}(y = \gamma(y_{1:k})) r^k (1-r)^{\ell(y)-k} \\ &\leq \binom{\ell(y)-1}{k-1} r^k (1-r)^{\ell(y)-k} \leq r. \end{aligned}$$

This inequality motivated our choice of test function.

To estimate N_e , we may consider the test functions

$$\begin{aligned} F_{\theta,e}(s,p) &= \sum_{x \in \theta} \sum_{y \in B_e} \mathbb{1}(s = \alpha(x,y)) \nu_\theta(x) \\ &\quad \times \left[\sum_{y_j \in B_{w_j}, 1 \leq j \leq k} \mathbb{1}(y = \gamma(y_{1:k})) r^k (1-r)^{\ell(y)-k} - p \right]. \end{aligned}$$

They are such that $|F_{\theta,e}(s,p)| \leq 1$ and

$$\mathbb{E}[F_{\theta,e}(S,p)] = \sum_{x \in \theta} \nu_\theta(x) \mathbb{P}_S(\alpha(x, B_e)) (N_e - p).$$

Therefore, for any θ such that

$$\sum_{x \in \theta} \mathbb{P}(S \in \alpha(B_e)) > 0,$$

N_e is the only value of $p \in [0, 1]$ such that $\mathbb{E}[F_{\theta,e}(S, p)] = 0$.

To estimate each D_{w_j} , we may consider the test functions

$$G_{\theta,w_j}(s, p) = \sum_{x \in \theta} \sum_{y \in B_{w_j}} \mathbf{1}(s = \alpha(x, y)) \nu_{\theta}(x) [r(1-r)^{\ell(y)-1} - p].$$

They are such that $|G_{\theta,w_j}(s, p)| \leq 1$ and

$$\mathbb{E}[G_{\theta,w_j}(S, p)] = \sum_{x \in \theta} \mathbb{P}_S(\alpha(x, B_{w_j})) (D_{w_j} - p).$$

Therefore, for any θ such that

$$\sum_{x \in \theta} \mathbb{P}(S \in \alpha(x, B_{w_j})) > 0,$$

D_{w_j} is the only value of $p \in [0, 1]$ such that $\mathbb{E}[G_{\theta,w_j}(S, p)] = 0$.

Now that we have explained how to estimate A_e , let us come to the estimation of $(q_{B_i}(B_e), i \in \mathbb{N} \setminus \{0\}, \lfloor_i e \in \mathcal{R})$.

Let us choose $\theta \in \Theta$ such that

$$\sum_{x \in \theta} \mathbb{P}(S \in \alpha(x, B_i)) > 0,$$

and consider for any $\lfloor_i e \in \mathcal{G}_i$ the test function defined as

$$F_{i,e,\theta}(s, p) = \sum_{x \in \theta} \sum_{y \in D^+} \mathbf{1}(s = \alpha(x, y)) \nu_{\theta}(x) [\mathbf{1}(y \in B_e) - p \mathbf{1}(y \in B_i)].$$

Then

$$\begin{aligned} \mathbb{E}[F_{i,e,\theta}(S, p)] &= \sum_{x \in \theta} \nu_{\theta}(x) \left(\mathbb{P}_S(\alpha(x, B_e)) - p \mathbb{P}_S(\alpha(x, B_i)) \right) \\ &= \sum_{x \in \theta} \nu_{\theta}(x) \mathbb{P}_S(\alpha(x, B_i)) (q_{B_i}(B_e) - p), \end{aligned}$$

so that $q_{B_i}(B_e)$ is the only value of p such that $\mathbb{E}[F_{i,e,\theta}(S, p)] = 0$. As such, the estimator \hat{p} such that $F_{i,e,\theta}(s, \hat{p}) = 0$ is without bias.

Let us also remark that if we assume only that B_e are Markov substitute sets for all expressions e such that $\lfloor_i e \in \mathcal{R}$, then B_i is itself a Markov substitute set if and only if for any $\lfloor_i e \in \mathcal{G}$, there is $p_{i,e} \in [0, 1]$ such that for any $\theta \in \Theta$, $\mathbb{E}[F_{i,e,\theta}(S, p_{i,e})] = 0$. This can serve to build statistical tests to validate the addition of new rules to the set of rules \mathcal{R} . But we will study these in more details in section 3.4.1 on page 144.

3.2.3 Simulating substitute measures

The next question is the simulation of the substitute measure q_{B_e} , for some expression $e \in S^+$.

Let us consider a parsing kernel $(\Pi_e(s, t), s \in B_e, t \in \check{B}_e)$ for the syntagm set B_e , as defined in definition 3.1.4 on page 123.

Let us define for any $[_i y \in \mathcal{R}$

$$p_i(y) = \frac{A_y q_{B_i}(B_y)}{\sum_{[_i z \in \mathcal{R}} A_z q_{B_i}(B_z)} \leq q_{B_i}(B_y).$$

Consider the subprobability measure on \check{B}_e defined as

$$p_e(t) = \mathbf{1}(\varsigma(t) = e) \prod_{[_j y \in \mathcal{R}} p_j(y)^{\chi(t, j, y)},$$

where χ is defined as in eq. (3.2.1) on page 125. To see that this is indeed a subprobability measure, see lemma 3.2.2 on the following page.

Starting from the obvious equality

$$\sum_{t \in \check{B}_e} p_e(t) \frac{q_{B_e}(\varphi(t)) \Pi_e(s, t)}{p_e(t)} = q_{B_e}(s), \quad s \in B_e,$$

we see that if we draw $T \in \check{B}_e \cup \{\epsilon\}$ according to the subprobability p_e completed with $p_e(\epsilon) = 1 - p_e(\check{B}_e)$, we obtain that for any bounded function f ,

$$\sum_{s \in B_e} q_{B_e}(s) f(s) = \mathbb{E}(w(T) f(T)),$$

where

$$\begin{aligned} w(T) &= \mathbf{1}(T \neq \epsilon) \frac{q_{B_e}(\varphi(T)) \Pi_e(\varphi(T), T)}{p_e(T)} \\ &= \mathbf{1}(T \neq \epsilon) \Pi_e(\varphi(T), T) A_e \prod_{j \in \mathbb{N} \setminus \{0\}} \left(\sum_{[_j y \in \mathcal{R}} A_y q_{B_j}(B_y) \right)^{\sum_{y, [_j y \in \mathcal{R}} \chi(T, j, y)}. \end{aligned}$$

The simulation of T according to p_e may be implemented in the following way. We start with $e_0 = e$. Given $e_k \in S^+$,

- if $e_k \in \mathcal{T}$, put $T = e_k$,

- otherwise $e_k = \alpha(x,]_i)$, $i \in \mathbb{N}$, $x \in \check{D}^* \times S^+$. Draw y according to p_i , completed with $p_i(\epsilon) = 1 - p_i(\check{B}_e)$, and put $e_{k+1} = \alpha(x, (iy)_i)$, when $y \neq \epsilon$ and set $T = \epsilon$ otherwise.

Lemma 3.2.2

For any $k \in \mathbb{N}$, any well-parenthesized expression $t \in (\check{D} \cup N)^*$, if we call $|t|$ the number of opening parentheses in t ,

$$\mathbb{P}(e_k = t) = \mathbf{1}(\zeta(t) = e) \mathbf{1}(|t| = k) \prod_{[j, y \in \mathcal{R}} p_j(y)^{\chi(t, j, y)}.$$

As such, if we define $\tau = \inf k, e_k \in \mathcal{T}$,

$$\mathbb{P}(e_\tau = t) = p_e(t).$$

PROOF. By induction on k :

- Obviously $\chi(e, j, y) = 0$ for all y, j , and

$$\mathbb{P}(e_0 = t) = \mathbf{1}(t = e) = \mathbf{1}(\zeta(t) = e) \mathbf{1}(|t| = 0).$$

- For any $t \in (\check{D} \cup N)^*$, suppose we can write $t = \alpha(x, (iy)_i)$, with $x_2 \in S^+$. Otherwise t has no parentheses and $\mathbb{P}(e_{k+1} = t) = 0$.

$$\begin{aligned} \mathbb{P}(e_{k+1} = t) &= \mathbb{P}(e_k = \alpha(x,]_i)) \cdot p_i(y) \\ &= \mathbf{1}(\zeta(\alpha(x,]_i)) = e) \mathbf{1}(|\alpha(x,]_i)| = k) \\ &\quad \times \prod_{[j, z \in \mathcal{R}} p_j(z) \chi(\alpha(x,]_i), j, z) + \mathbf{1}(j = i, z = y) \\ &= \mathbf{1}(\zeta(t) = e) \mathbf{1}(|t| = k + 1) \prod_{[j, z \in \mathcal{R}} p_j(z)^{\chi(t, j, z)}. \quad \square \end{aligned}$$

This lemma proves that p_e is in fact a probability distribution if the process defined previously finishes in finite time, in which case we can draw T using this method. The process does not finish in finite times if there are looping rewriting rules with probability above some critical value. This will not happen if the parameters are accurately estimated from a true Markov substitute process.

An alternative to the use of importance sampling is the use a reversible dynamics. Let us introduce the conditional probability kernel $(k_e(s, s'), s, s' \in B_e)$, defined for any $s \neq s'$ as

$$k_e(s, s') = \sum_{t, t' \in \check{B}_e} \Pi_e(s, t) p_e(t') w(t, t') \mathbf{1}(s' = \varphi(t')),$$

where

$$\begin{aligned} w(t, t') &= \left(1 \wedge \frac{\Pi_e(\varphi(t'), t')}{\Pi_e(\varphi(t), t)} \cdot \frac{q_{B_e}(\varphi(t'))}{p_e(t')} \cdot \frac{p_e(t)}{q_{B_e}(\varphi(t))} \right) \\ &= \left(1 \wedge \frac{\Pi_e(\varphi(t'), t')}{\Pi_e(\varphi(t), t)} \prod_{j \in \mathbb{N} \setminus \{0\}} \left(\sum_{\substack{y \in \mathcal{R} \\ j, y \in \mathcal{R}}} A_y q_{B_j}(B_y) \right)^{\sum_{j, y \in \mathcal{R}} \chi(t', j, y) - \chi(t, j, y)} \right). \end{aligned}$$

To draw s' according to $k_e(s, \cdot)$, we first draw t according to $\Pi_e(s, t)$, then we draw (independently !) t' according to $p_e(t')$ (again, we suppose that the process finishes in finite time), and we set $s' = \varphi(t')$ with probability $w(t, t')$, and $s' = s$ otherwise.

Proposition 3.2.3

The kernel $k_e(s, s')$ is reversible with respect to q_{B_e} and irreducible on B_e .

PROOF. It is irreducible because obviously $\text{supp}(k_e(s, \cdot)) = B_e$ already. To see that it is reversible we can write

$$q_{B_e}(s)k_e(s, s') = \sum_{t, t' \in \check{B}_e} \left(q_{B_e}(s)\Pi_e(s, t)p_e(t') \wedge q_{B_e}(s')\Pi_e(s', t)p_e(t) \right).$$

This expression is symmetric in s and s' , since it is of the form

$$q_{B_e}(s)k_e(s, s') = \sum_{t, t'} f(s, t, t', s'),$$

where $f(s, t, t', s') = f(s', t', t, s)$. □

3.2.4 Reversible dynamics for the language distribution

The same method can be used to build reversible dynamics for the language distribution on the whole, simply by parsing a sentence in general, allowing any surface structure. Let us consider a general parsing kernel

$$\left(\bar{\Pi}(s, t), s \in D^+, t \in \mathcal{T} \right)$$

such that $\text{supp}(\bar{\Pi}(s, \cdot)) \subset \varphi^{-1}(s)$, and a converse kernel ζ from syntactic trees to sequences of parsing trees. Consider now the conditional probability kernel $(k(s, s'), s, s' \in D^+)$ defined for $s \neq s'$ as

$$k(s, s') = \sum_{\substack{\tilde{t}, \tilde{t}' \in \mathcal{T}^+ \\ t, t' \in \mathcal{T}}} \bar{\Pi}(s, \tilde{t}) \mathbf{1}(t = \tilde{t}_{\ell(\tilde{t})}) p_{\zeta(t)}(t') \zeta(t', \tilde{t}') w(\tilde{t}, t, t', \tilde{t}') \mathbf{1}(s' = \varphi(t')),$$

where

$$\begin{aligned}
w(\tilde{t}, t, t', \tilde{t}') &= \left(1 \wedge \frac{\bar{\Pi}(\varphi(t'), \tilde{t}')}{\zeta(\tilde{t}', \varphi(t'))} \cdot \frac{\zeta(\tilde{t}, \varphi(t))}{\bar{\Pi}(\varphi(t), \tilde{t})} \cdot \frac{q_{B_{\varsigma(t')}}(\varphi(t'))}{p_{\varsigma(t')}(t')} \cdot \frac{p_{\varsigma(t)}(t)}{q_{B_{\varsigma(t)}}(\varphi(t))} \right) \\
&= \left(1 \wedge \frac{\bar{\Pi}(\varphi(t'), \tilde{t}')}{\zeta(\tilde{t}', \varphi(t'))} \cdot \frac{\zeta(\tilde{t}, \varphi(t))}{\bar{\Pi}(\varphi(t), \tilde{t})} \right. \\
&\quad \left. \times \prod_{j \in \mathbb{N} \setminus \{0\}} \left(\sum_{[j, y \in \mathcal{R}} A_y q_{B_j}(B_y) \right)^{\sum_{[j, y \in \mathcal{R}} \chi(t', j, y) - \chi(t, j, y)} \right).
\end{aligned}$$

Proposition 3.2.4

The kernel k is reversible with respect to the language distribution \mathbb{P}_S .

PROOF. We can use the identity $\mathbb{P}_S(\varphi(t)) = \mathbb{P}(S \in B_{\varsigma(t)})q_{B_{\varsigma(t)}}(\varphi(t))$ to remark that

$$\begin{aligned}
\mathbb{P}_S(s)k(s, s') &= \sum_{\substack{\tilde{t}, \tilde{t}' \in \mathcal{T}^+ \\ t, t' \in \mathcal{T}}} \mathbb{P}(S \in B_{\varsigma(t)})\mathbb{1}(\varsigma(t) = \varsigma(t')) \\
&\quad \mathbb{1}(s = \varphi(t))\mathbb{1}(s' = \varphi(t'))\mathbb{1}(t = \tilde{t}_{\ell(\tilde{t})})\mathbb{1}(t' = \tilde{t}'_{\ell(\tilde{t}')}) \\
&\quad \times \left(\bar{\Pi}(\varphi(t), t)\zeta(\tilde{t}', \varphi(t'))p_{\varsigma(t')}(t')q_{B_{\varsigma(t)}}(\varphi(t)) \right. \\
&\quad \left. \wedge \bar{\Pi}(\varphi(t'), t')\zeta(\tilde{t}, \varphi(t))p_{\varsigma(t)}(t)q_{B_{\varsigma(t')}}(\varphi(t')) \right),
\end{aligned}$$

where the symmetry in (s, s') can be seen in the same way as for k_e above. \square

This construction improves on the reversible dynamics of section 2.2.1 on page 69, inasmuch as it performs more substitutions at each step, allowing to substitute a sentence in B_e with a randomly chosen other one in one shot. It can also be seen as an implementation of the first dynamics using the extended Markov substitute set B_e . However, this method still requires the knowledge of q_{B_i} , and requires that the production process will terminate in finite time, something that will happen if the parameters are accurately estimated from a true Markov substitute process, but may fail if parameter estimation is not precise enough or if the data are off model. We will see later in section 3.3 on page 135 how to deal with those termination problems by restricting to a finite state space and using a different kind of estimator (crossing-over dynamics on a replicated sample).

We can now remark that the invariant dynamics k has the following context free property.

Proposition 3.2.5

Let us put $\bar{B}_e = \text{supp}(q_{B_e}k^\infty)$. Then k is reversible with respect to $q_{\bar{B}_e}$.

PROOF. To prove this, we will first prove the following lemma.

Lemma 3.2.6

For any $s, s' \in D^+$ such that $k(s, s') > 0$, $\{s, s'\}$ is a Markov substitute pair and

$$q_{\{s, s'\}}(s)k(s, s') = q_{\{s, s'\}}(s')k(s', s).$$

Indeed, when $k(s, s') > 0$, there is an expression e such that $s, s' \in B_e$, given by the surface structure of the parse tree t of s drawn according to $\Pi(s, \cdot)$. Thus $\{s, s'\}$ forms a Markov substitute pair and

$$q_{\{s, s'\}}(s)k(s, s') = q_{\{s, s'\}}(s')k(s', s). \quad (3.2.3)$$

The lemma implies that for any $s \in B_e$, $\text{supp}(\delta_s k^\infty)$ is a Markov substitute set, and, since $s, s' \in B_e$ implies $k(s, s') > 0$, that

$$\bigcup_{s \in B_e} \text{supp}(\delta_s k^\infty) = \text{supp}(q_{B_e}k^\infty) = \bar{B}_e$$

is also a Markov substitute set. It is stable : $k(s, \bar{B}_e) = 1$ for any $s \in \bar{B}_e$. Moreover, for any $s, s' \in \bar{B}_e$,

$$q_{\bar{B}_e}(s'') = q_{\bar{B}_e}(\{s, s'\})q_{\{s, s'\}}(s''), \quad s'' \in \{s, s'\},$$

proving (still from the above lemma) that k is reversible with respect to $q_{\bar{B}_e}$. \square

This leads us to define a equivalence relation on strings.

Definition 3.2.1

For any pair of strings s and s' , we will say that

$$s \sim s' \iff \sum_n k^n(s, s') > 0.$$

This is an equivalence relation, and the classes will be called syntactic categories.

Remark that the equivalence classes thus defined depend only on the syntagm sets of \mathcal{R} , and not on the actual substitute measures, and form Markov substitute sets.

3.2.5 Crossing-over dynamics

In the same vein as the crossing-over dynamics defined in section 2.2.5 on page 73, we will now present a second approach that requires neither knowledge of substitute measures, nor a structure with no infinite loops, but works on whole texts. This new dynamics will again be a way to accelerate the first crossing-over dynamics by allowing multiple swap at the same time.

Let us extend on texts the previously defined parsing kernel $\bar{\Pi}$ by defining

$$\bar{\Pi}(s_{1:n}, t_{1:n}) = \prod_{k=1}^n \bar{\Pi}(s_k, t_k), \quad s_{1:n} \in (D^+)^n, t_{1:n} \in \mathcal{T}^n.$$

Let us consider some crossing-over tree kernel $(\mathfrak{s}(t, t'), t, t' \in \bar{\mathcal{T}})$, such that $\mathfrak{s}(t, t') > 0$ implies that

$$\begin{aligned} \varsigma(t'_k) &= \varsigma(t_k), & 1 \leq k \leq n, \\ \sum_{k=1}^n \chi(t_k, j, e) &= \sum_{k=1}^n \chi(t'_k, j, e), & [j, e \in \mathcal{B}. \end{aligned} \quad (3.2.4)$$

Let us remark that as a consequence of these properties,

$$\mathbb{P}_S^{\otimes n}(\varphi(t_k), 1 \leq k \leq n) = \mathbb{P}_S^{\otimes n}(\varphi(t'_k), 1 \leq k \leq n). \quad (3.2.5)$$

Such crossing-over tree kernel will be constructed in section 3.C on page 163.

Let us consider the conditional probability kernel K defined for $s_{1:n} \neq s'_{1:n} \in \mathcal{T}$ as

$$\begin{aligned} K(s_{1:n}, s'_{1:n}) &= \sum_{t, t' \in \bar{\mathcal{T}}^+} \bar{\Pi}(s_{1:n}, t) \mathfrak{s}(e(t), e(t')) \zeta(t', s'_{1:n}) \\ &\quad \times \left(\frac{\zeta(t, s_{1:n}) \mathfrak{s}(t, t') \bar{\Pi}(s'_{1:n}, t')}{\bar{\Pi}(s_{1:n}, t) \mathfrak{s}(t', t) \zeta(t', s'_{1:n})} \wedge 1 \right). \end{aligned}$$

Proposition 3.2.7

The crossing-over sample kernel K is reversible with respect to $\mathbb{P}_S^{\otimes n}$.

PROOF. As a consequence of eq. (3.2.5) on the current page, we have to prove that the kernel K is symmetric. This property can easily be deduced from the following identity

$$K(s_{1:n}, s'_{1:n}) = \sum_{t_{1:n}, t'_{1:n} \in \mathcal{T}^n} \left(\bar{\Pi}(s_{1:n}, t_{1:n}) \zeta(t_{1:n}, t'_{1:n}) \wedge \bar{\Pi}(s'_{1:n}, t'_{1:n}) \zeta(t'_{1:n}, t_{1:n}) \right).$$

□

3.3 Crossing-over dynamics and the maximum likelihood estimator

Let us consider some given (deterministic) sample $s_{1:n} \in \mathcal{D}^n$, where $\mathcal{D} \in D^+$ is a finite domain (for example the set of sentences of length no larger than some constant). Let us assume that $s_{1:n}$ is not constant (in other words let us assume that $|\{s_i, 1 \leq i \leq n\}| \geq 2$).

Let \mathcal{B} be a finite set of finite subsets of D^+ , such that

$$\sum_{x \in (D^*)^2} \sum_{j=1}^n \mathbb{1}(s_j = \alpha(x, y)) > 0, \quad y \in B. \quad (3.3.1)$$

This means that each member of each substitute set in \mathcal{B} is present in the sub-strings of the sample $s_{1:n}$.

We will say that $p \in \mathcal{M}_+^1(\mathcal{D})$ is a \mathcal{B} -Markov substitute process on \mathcal{D} if it satisfies eq. (2.1.1) on page 64 and there are substitute measures $q_B \in \mathcal{M}_+^1(B)$ for all $B \in \mathcal{B}$, such that for any $y, y' \in B$ and any $x \in (D^*)^2$ such that $\{\alpha(x, y), \alpha(x, y')\} \subset \mathcal{D}$,

$$p[\alpha(x, y)]q_B(y') = p[\alpha(x, y')]q_B(y).$$

This definition goes with a modification of the equivalence relation defined by eq. (2.3.1) on page 75, where we impose that the path should belong to the restricted domain \mathcal{D} . Namely

$$\begin{aligned} s \sim_{\mathcal{B}, \mathcal{D}} s' &\iff \exists (x_j, y_j, y'_j), x_j \in (D^*)^2, y_j, y'_j \in B_j \in \mathcal{B}, \\ &\alpha(x_j, y_j), \alpha(x_j, y'_j) \in \mathcal{D}, \alpha(x_{j-1}, y'_{j-1}) = \alpha(x_j, y_j), 0 \leq j \leq J, \\ &s = \alpha(x_0, y_0), s' = \alpha(x_J, y'_J). \end{aligned}$$

It is easy to extend proposition 2.3.2 on page 81 to show that the set of \mathcal{B} -Markov substitute processes on \mathcal{D} whose support is the minimal possible support including $\{s_i, 1 \leq i \leq n\}$, is an exponential family of the form

$$\mathfrak{M}_{\mathcal{C}}(\mathcal{D}, \mathcal{B}) = \left\{ \left(p_{\beta}(s) = Z_{\beta}^{-1} \exp \left(- \sum_{i=1}^I \beta_i U_i(s) \right), s \in \bigcup \mathcal{C} \right), \beta \in \mathbb{R}^I \right\},$$

where

$$\mathcal{C} = \left\{ C \in \mathcal{D} / \sim_{\mathcal{B}, \mathcal{D}}; \sum_{i=1}^n \mathbb{1}(s_i \in C) > 0 \right\}.$$

We would like in this section to solve the question of finding the maximum likelihood estimator of $s_{1:n}$ in the set of distributions $\mathfrak{M}_{\mathcal{G}}(\mathcal{D}, \mathcal{B})$. In other words we would like to solve

$$\inf_{\beta \in \mathbb{R}^I} \sum_{i=1}^n -\log[p_{\beta}(s_i)].$$

This is a convex minimization problem, and by taking derivatives, we can see that when the minimum is reached, it is characterized by the fact that

$$\int_{\mathcal{D}} U_i dp_{\beta_*} = \frac{1}{n} \sum_{j=1}^n U_i(s_j), \quad 1 \leq i \leq I.$$

We will prove in this section that the minimum is indeed reached when eq. (3.3.1) on the previous page is satisfied. In practice however, obtaining an explicit parametrization of p_{β} (through an explicit computation of the potential functions U_i), is only possible in very simple situations, so that this result is only of theoretical interest.

We will show on the other hand that a crossing-over dynamics can be used to compute a Monte-Carlo approximation of p_{β_*} .

Let K_m be a symmetric kernel on the product \mathcal{D}^{nm} , such that

$$K_m(s_{1:nm}, s'_{1:nm}) > 0 \implies \left(\sum_{j=1}^{nm} U_i(s_j) - U_i(s'_j) = 0, \quad 1 \leq i \leq I \right).$$

Note that the dynamics K defined in section 3.2.5 on page 134 verifies this property. Let us assume moreover that K_m is permutation invariant in the sense that for any permutation σ ,

$$K_m(s_{1:nm} \circ \sigma, s'_{1:nm} \circ \sigma) = K_m(s_{1:nm}, s'_{1:nm}), \quad s_{1:nm} \in \mathcal{D}^{nm}, s'_{1:nm} \in \mathcal{D}^{nm},$$

where $s_{1:nm} \circ \sigma = (s_{\sigma(i)}, 1 \leq i \leq nm)$. Let us assume also that

$$K_m \left[(\alpha(x, y), \alpha(x', y'), s_{3:nm}), (\alpha(x, y'), \alpha(x', y), s_{3:nm}) \right] > 0, \\ x, x' \in (D^*)^2, y, y' \in B \in \mathcal{B}, \alpha(x, y), \alpha(x, y'), \alpha(x', y), \alpha(x', y') \in \mathcal{D}.$$

This means that K_m performs individual swaps with a positive probability.

It is quite easy to modify the dynamics of section 3.2.5 on page 134 to ensure these properties, by applying a random permutation to the sentences for the first one, and setting the probability of each crossing over to be strictly less than one for the second.

We will show that K_m can be used to sample (approximately) from p_{β_*} . Let $S^m \in (D^+)^{mn}$ be defined as

$$S_{\sigma(kn+i)}^m = s_i, \quad 0 \leq k < m, 1 \leq i \leq n,$$

where σ is a uniform random permutation of $\{1, \dots, nm\}$ independent of everything else, so that S^m is a random uniform shuffle of m copies of s .

Let us consider for each m the distribution on \mathcal{D}^{nm} defined as

$$p_m = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{u=1}^t \mathbb{P}_{S^m} K_m^u.$$

Since the initial distribution \mathbb{P}_{S^m} is exchangeable, and since K_m is permutation invariant, p_m is an exchangeable probability measure. Moreover, since K_m is a symmetric kernel, and since the starting distribution is a uniform shuffle of a single configuration, p_m is uniform on its support, its support being permutation invariant. Let us consider a random variable $(\bar{S}_{m,i}, 1 \leq i \leq nm)$ distributed according to p_m , and let us put

$$N_m = \sum_{i=1}^{nm} \delta_{\bar{S}_{m,i}} \in \mathcal{M}_+(\mathcal{D}).$$

This is a random measure with integer weights. According to what we said above about p_m ,

$$\mathbb{P}(N_m = \nu) = \frac{(nm)!}{Z_m \prod_{s \in D^+} \nu(s)!}, \quad \nu \in \text{supp}(\mathbb{P}_{N_m}),$$

where the partition function Z is defined as

$$Z_m = \sum_{\nu \in \text{supp}(\mathbb{P}_{N_m})} \frac{(nm)!}{\prod_{s \in D^+} \nu(s)!},$$

and where we set by convention that $0! = 1$. This comes from the fact that the map $\Phi(s_{1:nm}) = \sum_{i=1}^{nm} \delta_{s_i}$ is such that $\Phi^{-1}(\nu)$ is made of $\frac{(nm)!}{\prod_{s \in D^+} \nu(s)!}$ distinct permutations

of the same sequence of sentences, that are equiprobable under p_m , as soon as $\nu \in \text{supp}(\mathbb{P}_{N_m})$.

Let us recall Stirling's approximation formula (see for example [Fel68], p.54)

$$\exp\left(\frac{1}{12n+1}\right) < \frac{n!}{\sqrt{2\pi n} n^{n+1/2} \exp(-n)} < \exp\left(\frac{1}{12n}\right), \quad n \geq 1.$$

Since $|\text{supp}(\mathbb{P}_{N_m})| \leq (nm+1)^{|\mathcal{D}|}$, the following lemma is a consequence of Stirling's approximation formula.

Lemma 3.3.1

Let $\mu_m \in \text{supp}(\mathbb{P}_{N_m/(nm)})$ be such that

$$H(\mu_m) = \max \left\{ H(\nu), \nu \in \text{supp}(\mathbb{P}_{N_m/(nm)}) \right\},$$

where

$$H(\mu) = - \sum_{s \in \mathcal{D}} \mu(s) \log[\mu(s)], \quad \mu \in \mathcal{M}_+^1(\mathcal{D}),$$

is the Shannon entropy.

There is a real positive constant c (depending only on $|\mathcal{D}|$) such that

$$(nm + 1)^{-c} \leq \frac{\mathbb{P}(N_m = \nu)}{\exp \left\{ nm [H(\nu/(nm)) - H(\mu_m)] \right\}} \leq (nm + 1)^c. \quad (3.3.2)$$

PROOF. Let us consider some $\nu \in \text{supp}(\mathbb{P}_{N_m})$, and let us put $\mathcal{D}_\nu = \text{supp}(\nu)$, $d_\nu = |\mathcal{D}_\nu|$ and $d = |\mathcal{D}|$. From Stirling formula, we get

$$\begin{aligned} & \frac{(nm)!}{\prod_{s \in \mathcal{D}_\nu} \nu(s)!} \\ & \leq (2\pi)^{-(d_\nu-1)/2} \frac{(nm)^{nm}}{\prod_{s \in \mathcal{D}_\nu} \nu(s)^{\nu(s)}} \left(\frac{nm}{\prod_{s \in \mathcal{D}_\nu} \nu(s)} \right)^{1/2} \exp \left(\frac{1}{12nm} - \sum_{s \in \mathcal{D}_\nu} \frac{1}{12(\nu(s) + 1)} \right) \\ & \leq (nm)^{1/2} \exp [nmH(\nu/(nm))]. \end{aligned}$$

On the other hand

$$\begin{aligned} & \frac{(nm)!}{\prod_{s \in \mathcal{D}_\nu} \nu(s)!} \\ & \geq (2\pi)^{-(d_\nu-1)/2} \frac{(nm)^{nm}}{\prod_{s \in \mathcal{D}_\nu} \nu(s)^{\nu(s)}} \left(\frac{nm}{\prod_{s \in \mathcal{D}_\nu} \nu(s)} \right)^{1/2} \exp \left(\frac{1}{12(nm + 1)} - \sum_{s \in \mathcal{D}_\nu} \frac{1}{12\nu(s)} \right). \end{aligned}$$

Let us then remark that, introducing the uniform probability measure $U_{\mathcal{D}_\nu}$ on \mathcal{D}_ν ,

$$\begin{aligned}
 \prod_{s \in \mathcal{D}_\nu} \nu(s) &= \exp\left(\sum_{s \in \mathcal{D}_\nu} \log(\nu(s))\right) \\
 &= \exp\left[-d_\nu \frac{1}{d_\nu} \sum_{s \in \mathcal{D}_\nu} \log\left(\frac{nm}{d_\nu \nu(s)}\right) + d_\nu \log\left(\frac{nm}{d_\nu}\right)\right] \\
 &= \exp\left[d_\nu \left(\log\left(\frac{nm}{d_\nu}\right) - \mathcal{K}(U_{\mathcal{D}_\nu}, \nu/(nm))\right)\right] \\
 &\leq \exp\left[d_\nu \log\left(\frac{nm}{d_\nu}\right)\right] \\
 &\leq \exp\left[d \max\left\{1, \log\left(\frac{nm}{d}\right)\right\}\right] \\
 &\leq \exp\left[d \log\left(\frac{enm}{d}\right)\right] \\
 &\leq \left(\frac{nm}{d}\right)^d \exp(d).
 \end{aligned}$$

Combining the last two equations, we get

$$\begin{aligned}
 \frac{(nm)!}{\prod_{s \in \mathcal{D}_\nu} \nu(s)!} &\geq d^{1/2} \left(\frac{2\pi nm}{d}\right)^{-(d-1)/2} \exp(-7d/12) \exp[nmH(\nu/(nm))] \\
 &\geq \exp\left(\frac{d}{2} \left(\log(d) - \log(2\pi) - 7/12\right) + \frac{1}{2} \log(2\pi)\right) (nm)^{-(d-1)/2} \exp[nmH(\nu/(nm))] \\
 &\geq (nm)^{-(\max\{d, 10\}-1)/2} \exp[nmH(\nu/(nm))].
 \end{aligned}$$

Accordingly, there is a constant $a > 0$, depending only on $|\mathcal{D}|$, such that

$$(nm)^{-a} \leq \frac{(nm)!}{\prod_{s \in \mathcal{D}_\nu} \nu(s)!} \exp[-nmH(\nu/(nm))] \leq (nm)^{1/2}. \quad (3.3.3)$$

Let us remark now that

$$\left| \text{supp}(\mathbb{P}_{N_m}) \right| \leq (nm + 1)^d.$$

This implies that

$$Z_m \leq (nm + 1)^d (nm)^{1/2} \exp(nmH(\mu_m)).$$

On the other hand

$$Z_m \geq (nm + 1)^{-a} \exp(nmH(\mu_m)),$$

so that

$$(nm + 1)^{-(c+d+1/2)} \leq \frac{\mathbb{P}(N_m = \nu)}{\exp\left(nm\left[H\left(\nu/(nm)\right) - H(\mu_m)\right]\right)} \leq (nm + 1)^{a+1/2},$$

ending the proof of the above lemma. \square

To use this lemma, it remains to prove that $\text{supp}(N_m)$ is essentially defined by the constraints

$$\text{supp}\left(\mathbb{P}_{N_m/(mn)}\right) \simeq \left\{ \nu \in \mathcal{M}_+^1\left(\bigcup \mathcal{C}\right), \int U_i(s) d\nu(s) = \frac{1}{n} \sum_{j=1}^n U_i(s_j), 1 \leq i \leq I \right\},$$

This fact is not so easy to prove directly. To come to a conclusion, we are going to describe a non random path in $\text{supp}\left(\mathbb{P}_{N_m/(mn)}\right)$, starting at the initial probability measure

$$\bar{p} = \frac{1}{n} \sum_{i=1}^n \delta_{s_i},$$

and converging to some $p_\beta \in \mathfrak{M}_\mathcal{C}(\mathcal{D}, \mathcal{B})$.

Let us consider the set of vectors

$$\mathcal{Q} = \left\{ \delta_{\alpha(x,y')} + \delta_{\alpha(x',y)} - \delta_{\alpha(x,y)} - \delta_{\alpha(x',y')}, x, x' \in (D^*)^2, (y, y') \in B^2 \in \mathcal{B}, \right. \\ \left. \alpha(x, y), \alpha(x, y'), \alpha(x', y), \alpha(x', y') \in \mathcal{D} \right\}.$$

Since \mathcal{D} is a finite subset of D^+ , \mathcal{Q} is also a finite set. Let us consider the convex set

$$A = \left\{ \bar{p} + \sum_{v \in \mathcal{Q}} \gamma(v)v, \gamma \in \mathbb{R}^{\mathcal{Q}} \right\} \cap \mathcal{M}_+^1(\mathcal{D}).$$

Let us prove first that we can find a positive constant $\zeta > 0$, a finite sequence $v_j \in \mathcal{Q}$, $1 \leq j \leq J$, such that for any large enough value of m ,

$$\mu_k = \bar{p} + \sum_{j=1}^k \frac{[nm\zeta]}{nm} v_j \in A, \quad 0 \leq k \leq J,$$

and $\{s \in \mathcal{D}, \mu_{J_1}(s) \geq \zeta\} = \bigcup \mathcal{C}$.

Indeed, for any $s \in C \in \mathcal{C}$, there is i , $1 \leq i \leq n$ and a sequence

$$(x_k, y_k, y'_k), 0 \leq k \leq J, x_k \in (D^*)^2, (y_k, y'_k) \in B^2 \in \mathcal{B}^2, \\ \alpha(x_k, y_k) \in \mathcal{D}, \alpha(x_{k-1}, y'_{k-1}) = \alpha(x_k, y_k), 1 \leq k \leq J, \\ \alpha(x_0, y_0) = s_i, \alpha(x_j, y'_j) = s.$$

For any k , $0 \leq k \leq J$, there is i_k , $1 \leq i_k \leq n$ and $x'_k \in (D^*)^2$ such that $s_{i_k} = \alpha(x'_k, y'_k)$. Considering

$$v_k = \delta_{\alpha(x_k, y'_k)} + \delta_{\alpha(x'_k, y_k)} - \delta_{\alpha(x_k, y_k)} - \delta_{s_{i_k}},$$

it is easy to check that for ζ small enough and m large enough,

$$\mu_k = \bar{p} + \frac{[nm\zeta]}{nm} \sum_{j=0}^k v_j \in \mathcal{M}_+^1(\mathcal{D}),$$

(and belongs to A .) and that

$$\mu_J(s_i) \geq \zeta, \text{ and } \mu_J(s) \geq \zeta.$$

Indeed, to create μ_J , we remove some mass from $\{s_i, 1 \leq i \leq n\}$ to carry a fraction not less than ζ of it to s while spilling the rest around. As the amount of mass that is moved around goes to zero when ζ decreases and m increases, it is possible to keep some positive mass at each s_i , $1 \leq i \leq n$.

We can repeat this process, starting from the end point μ_J , and increasing m if necessary, to extend progressively the support to the whole $\cup \mathcal{C}$, along a sequence μ_k , $1 \leq k \leq J$. We build in this way a probability measure

$$\nu_m = \mu_J = \bar{p} + \sum_{j=1}^J \frac{[nm\zeta]}{nm} v_j \in \text{supp}(\mathbb{P}_{(nm)^{-1}N_m}),$$

such that $\nu_m(s) \geq \zeta$, $s \in \cup \mathcal{C}$, and such that

$$\lim_{m \rightarrow \infty} \nu_m = \nu_\infty = \bar{p} + \sum_{j=1}^J \zeta v_j \in A,$$

with $\nu_\infty(s) \geq \zeta$, $s \in \cup \mathcal{C}$.

Now let us consider $\mu_* \in A$ such that

$$H(\mu_*) = \sup\{H(\mu), \mu \in A\}.$$

Since the Shannon entropy H is strictly convex, it reaches its maximum on the compact convex set A at a single point, so that μ_* is uniquely defined by this property.

Let us prove that $\text{supp}(\mu_*) = \cup \mathcal{C}$. Since $H(\mu_*)$ is maximum,

$$\begin{aligned} \lim_{\substack{\lambda \rightarrow 0 \\ \lambda > 0}} \sum_{s \in \cup \mathcal{C}} (\mu_*(s) - \nu_\infty(s)) \log[\lambda \nu_\infty(s) + (1 - \lambda) \mu_*(s)] \\ = \lim_{\substack{\lambda \rightarrow 0 \\ \lambda > 0}} \frac{\partial}{\partial \lambda} H(\lambda \nu_\infty + (1 - \lambda) \mu_*) \leq 0. \end{aligned}$$

Since the left-hand side of this equation is equal to $+\infty$ as soon as $\mu_*(s) = 0$ for some $s \in \bigcup \mathcal{C}$, this proves our claim (here we need the fact that $\nu_\infty(s) > 0$, that we proved before). Let us remark that we could have used any ν_m with m large enough instead of ν_∞ in the previous equation.

Lemma 3.3.2

For any $x, x' \in (D^*)^2$, and $(y, y') \in B^2 \in \mathcal{B}^2$, such that

$$\alpha(x, y), \alpha(x, y'), \alpha(x', y), \alpha(x', y') \in \bigcup \mathcal{C},$$

$$\frac{\mu_*[\alpha(x, y')]}{\mu_*[\alpha(x, y)]} = \frac{\mu_*[\alpha(x', y')]}{\mu_*[\alpha(x', y)]}.$$

PROOF. Consider $v = \delta_{\alpha(x, y)} + \delta_{\alpha(x', y')} - \delta_{\alpha(x, y')} - \delta_{\alpha(x', y)}$. For any small enough positive or negative value of γ , $\mu_* + \gamma v \in A$, therefore

$$0 = \left. \frac{\partial}{\partial \gamma} \right|_{\gamma=0} H(\mu_* + \gamma v)$$

$$= \log[\mu(\alpha(x, y'))] + \log[\mu(\alpha(x', y))] - \log[\mu(\alpha(x, y))] - \log[\mu(\alpha(x', y'))]. \quad \square$$

Since $\{s_i, 1 \leq i \leq n\} \subset \text{supp}(\mu_*)$, μ_* satisfies eq. (2.1.1) on page 64, therefore the above lemma implies that $\mu_* \in \mathfrak{M}_{\mathcal{C}}(\mathcal{D}, \mathcal{B})$. Moreover, as it is the case for any probability measure in A ,

$$\sum_{s \in \mathcal{D}} \mu_*(s) U_i(s) = \sum_{s \in \mathcal{D}} \bar{p}(s) U_i(s),$$

proving that β_* exists and that $\mu_* = p_{\beta_*}$.

Since $\nu, \mu_* \in A$ and since \mathcal{D} is symmetric, we can write μ_* as

$$\mu_* = \nu_\infty + \sum_{j=1}^J \gamma_j v_j, v_j \in \mathcal{D}, \gamma_j > 0, 1 \leq j \leq J, \quad (3.3.4)$$

where v_j are independent vectors. Consider the lattice

$$L_m = \left(\nu_m + \frac{1}{nm} \sum_{j=1}^J \mathbb{Z} v_j \right) \cap A.$$

Since $\text{supp}(\nu_\infty) = \text{supp}(\mu_*) = \bigcup \mathcal{C}$, it is easy to check that for m large enough, there is a connected path in this lattice joining ν_m to the point of the lattice nearest to μ_* . This a consequence of the fact that the segment joining ν_∞ to μ_* is in the

interior of $A \cap \left(\nu_\infty + \sum_{j=1}^J \mathbb{R}v_j \right)$ viewed as a subset of $\nu_\infty + \sum_{j=1}^J \mathbb{R}v_j$. Since the nearest point to μ_* in L_m converges to μ_* according to eq. (3.3.4) on the preceding page, and since $L_m \in \text{supp}(\mathbb{P}_{N_m/(nm)})$, we have proved that

$$\lim_{m \rightarrow \infty} d\left(\text{supp}(\mathbb{P}_{N_m/(nm)}), \mu_*\right) = 0,$$

where d is (for instance) the Euclidean distance.

Proposition 3.3.3

Under the hypotheses described at the beginning of this section, the maximum likelihood

$$\sup_{\beta \in \mathbb{R}^I} \sum_{j=1}^n \log[p_\beta(s_j)]$$

is reached at some $\beta_ \in \mathbb{R}^I$. For any $\varepsilon > 0$, there are $\eta > 0$ and $M > 0$ such that for any $m \geq M$,*

$$\mathbb{P}\left[d(N_m, p_{\beta_*}) \geq \varepsilon\right] \leq \exp(-nm\eta),$$

where $d(\mu, \nu) = \sqrt{\sum_{s \in \mathcal{D}} (\mu(s) - \nu(s))^2}$. Moreover,

$$\lim_{m \rightarrow \infty} \mathbb{E}(N_m/(nm)) = p_{\beta_*}.$$

PROOF. For any $\varepsilon > 0$, there is $\eta > 0$ such that

$$\{\mu \in A, d(\mu, p_{\beta_*}) \geq \varepsilon\} \subset \{\mu \in A, H(\mu) \leq H(p_{\beta_*}) - 3\eta\},$$

because H is strictly concave.

For m large enough, there is a measure $\mu \in \text{supp}(\mathbb{P}_{N_m/(nm)})$ whose entropy is such that $H(\mu) \geq H(p_{\beta_*}) - \eta$. Thus

$$\begin{aligned} \mathbb{P}\left[d(N_m/(nm), p_{\beta_*}) \geq \varepsilon\right] &\leq \mathbb{P}\left[H(N_m/(nm)) \leq H(p_{\beta_*}) - 3\eta\right] \\ &\leq \mathbb{P}\left[H(N_m/(nm)) \leq H(\mu) - 2\eta\right] \\ &\leq \exp(-nm\eta). \end{aligned}$$

The last inequality is a consequence of eq. (3.3.2) on page 138 and of the fact that

$$\left|\text{supp}(\mathbb{P}_{N_m/(nm)})\right| \leq (nm + 1)^{|\mathcal{D}|}$$

is polynomial and not exponential in mn . To study the convergence of the expectation of $N_m/(nm)$, we can remark that the Euclidean distance between two probability distributions cannot be greater than 2, so that, for m large enough,

$$d[\mathbb{E}(N_{mn}/(nm)), p_{\beta_*}] \leq \mathbb{E}[d(N_m/(nm), p_{\beta_*})] \leq 2^{1/2}\varepsilon + \exp(-nm\eta),$$

proving that

$$\lim_{m \rightarrow \infty} \mathbb{E}(N_m/(nm)) = p_{\beta_*}. \quad \square$$

3.4 Building a Markov ruleset

Let us now discuss the question of building a Markov ruleset. If we index the dictionary $D = (w_i)_{i=1}^M$, the trivial set of rules

$$\mathcal{R}_0 = \bigoplus_i 1 \otimes [{}_i w_i$$

is of course Markov, for any distribution. The question now is, given a Markov grammar \mathcal{R} , how to make it grow.

3.4.1 Adding new rules

If we are given a Markov grammar \mathcal{R} , with the corresponding substitute measures, we already know that B_e is a Markov substitute set for any expression e . We could now want to add new rules $[{}_i e$ to the grammar, so that the new grammar

$$\mathcal{R}' = \mathcal{R} \bigoplus [{}_i e$$

is still Markov. Let us note

$$I(\mathcal{R}) = \{i \in \mathbb{N}, \exists e \in S^+, [{}_i e \in \mathcal{R}\}.$$

As a first step, we will only add rules $[{}_i e$ where $i \notin I(\mathcal{R})$. In this case, $B'_i = B_i$, $i \in I(\mathcal{R})$ will still be Markov substitute sets, and we only have to check the property for the new syntagm sets. In practice, this means that we want to find expressions e_1, e_2 so that $B'_i = B_{e_1} \cup B_{e_2}$ is still a Markov substitute set, in which case we can add $[{}_i e_1 \oplus [{}_i e_2$ to the grammar.

Note that if we have at some point added a rule $[{}_i]_j$, we should simply identify the labels i and j instead.

Now, suppose we have a splitting kernel

$$\pi : D^+ \mapsto \mathcal{M}_+^1((D^*)^2 \times D^*)$$

such that

$$\pi(s, x, y) > 0 \text{ and } y \neq \epsilon \implies \exists \check{x} \in (\check{D}^*)^2, \check{y} \in \mathcal{T} : \varsigma(\check{y}) \in \{e_1, e_2\},$$

and such that

$$\mathcal{S}(s, B'_i) \subset \text{supp}(\pi(s, \cdot)). \quad (3.4.1)$$

If we introduce for B'_i the test functions

$$\begin{aligned} F_\theta(S, p) = & \sum_{x \in (D^*)^2} \sum_{y_1, y_2 \in D^+} \mathbf{1}(x \in \theta) \left[\pi(\alpha(x, y_1), x, y_1) \wedge \pi(\alpha(x, y_2), x, y_2) \right] \\ & \times \left[(1-p) \mathbf{1}(S = \alpha(x, y_1)) \mathbf{1}(y_1 \in B_{e_1}) q_{B'_i \setminus B_{e_1}}(y_2) \right. \\ & \left. - p \mathbf{1}(S = \alpha(x, y_2)) \mathbf{1}(y_2 \in B'_i \setminus B_{e_1}) q_{B_{e_1}}(y_1) \right], \end{aligned}$$

we have that B'_i is a Markov substitute set if and only if, for some value of $p \in]0, 1[$,

$$\forall \theta \in \Theta, \mathbb{E}(F_\theta(S, p)) = 0.$$

When this is the case, we have

$$q_{B'_i}(B_{e_1}) = p. \quad (3.4.2)$$

This mean that we can use the random variables $F_\theta(S, p)$ to define a test.

Remark at this point that eq. (3.4.1) on the current page is only necessary to test if the new set is a “true” Markov substitute set. Whether the condition is verified or not, it will not change the fact that if $\mathbb{E}(F_\theta(S, p)) = 0$, the invariant dynamics defined with this splitting kernel will still be invariant. Roughly speaking, the set would be a Markov substitute set only in the contexts weighted by π .

The test function presented in section 3.2.2 on page 127 may also be used to test these new rules, however, the introduction of the measure ν_θ can make the test less efficient when θ is large.

3.4.2 Reducing the context space

The test we described here relies on the fact that $B = C_1 \cup C_2$ is a Markov substitute set if and only if for some $p \in]0, 1[$,

$$\forall \theta \in \Theta, \mathbb{E}(F_\theta(S, p)) = 0,$$

assuming that the context set Θ contains all singletons $\{x\}, x \in D^{*2}$.

However, we can use here the reduction kernel r to reduce the context space, by observing that for any context (x_1, x_2) , any pair of expressions (e_1, e_2) such that $x_i \in B_{e_i}$, any string y ,

$$\mathbb{P}_S(\gamma(x_1, y, x_2)) = \left(\sum_{\substack{x'_1 \in B_{e_1}, \\ x'_2 \in B_{e_2}}} \mathbb{P}_S(\gamma(x'_1, y, x'_2)) \right) q_{B_{e_1}}(x_1) q_{B_{e_2}}(x_2). \quad (3.4.3)$$

Let us assume that we have a collection of expressions $(e_i)_i, i \in I$, such that $D^* = \cup_{i \in I} B_{e_i}$ and $(\{B_{e_i}, i \in I\})^2 \subset \Theta$, and do not assume any more that Θ contains all one point sets.

Proposition 3.4.1

In this case, the union $B = C_1 \cup C_2$ is a Markov substitute set, if and only if, for some $p \in]0, 1[$,

$$\forall \theta \in \Theta, \mathbb{E}(F_\theta(S, p)) = 0,$$

PROOF. We see from the definition of $F_\theta(S, p)$ that $\mathbb{E}(F_\theta(S, p)) = 0$ if and only if

$$\sum_{x \in \theta} \sum_{y \in C_1} \mathbb{P}_S(\alpha(x, y)) = p \sum_{x \in \theta} \sum_{y \in C_1 \cup C_2} \mathbb{P}_S(\alpha(x, y)), \quad (3.4.4)$$

whereas $C_1 \cup C_2$ is a Markov substitute set if and only if for any $x \in (D^*)^2$,

$$\sum_{y \in C_1} \mathbb{P}_S(\alpha(x, y)) = p \sum_{y \in C_1 \cup C_2} \mathbb{P}_S(\alpha(x, y)). \quad (3.4.5)$$

Now, for any $x = (x_1, x_2) \in (D^*)^2$, there are e_1, e_2 such that $x_1 \in B_{e_1}, x_2 \in B_{e_2}$, and we see that the combination of eq. (3.4.4) on the current page applied to $\theta = B_{e_1} \times B_{e_2}$ and of eq. (3.4.3) on this page implies eq. (3.4.5) on the current page. \square

To implement this, we can put the whole of $\{B_e, \exists s \in D^+, r(s, e) > 0\}^2$ in Θ , since we can make Θ as big as we like, knowing that we only have to test on the contexts that we actually see, and compute $\mathbb{1}(x \in B_{e_1} \times B_{e_2})$ as

$$\mathbb{1}(r(x_i, e_i) > 0, i = 1, 2).$$

3.4.3 Saturation

We actually can go further than simply adding rules to the ruleset. Consider for a given ruleset \mathcal{R} , the collection of subsets $\mathcal{B} = \{B_i, i \in \mathbb{N}\}$, and, as constructed in section 2.3 on page 74, the graph in D^+ so that $(\alpha(x, y), \alpha(x, y'))$ is an edge

if, and only if, y and y' are in the same B_i . This graph defines an equivalence relation $\sim_{\mathcal{R}}$, which is obviously the same as the one defined in definition 3.2.1 on page 133.

Now we can remark that we can change the actual graph (that is, the ruleset) without changing the actual equivalence relation. An “sparse” graph, with few edges and small B_i sets, was considered in section 2.3 on page 74, where the goal was to simplify as much as possible the graph to ease computation. A “fuller” graph with many edges and big sets, on the other hand, accelerates the invariant dynamics, since they only jump along the edges of the graph. We will here attempt to build, from a given ruleset, another one, with the same equivalence relation, but with richer syntagm sets.

Let us first consider a small example to see how this works. Say we have in the rules

$$\begin{aligned} & [{}_1ab \\ & \oplus [{}_2b \oplus [{}_2c. \end{aligned}$$

Such a set of rules will not be able to recognise B_{ac} as a syntagm set, while it is a Markov substitute set, easily obtained from these rules. Indeed, the set of rules

$$\begin{aligned} & [{}_1a]_2 \\ & \oplus [{}_2b \oplus [{}_2c, \end{aligned}$$

is still Markov (if the first one was, obviously), and does recognise ac . The missing part is $[{}_1a]_2$, which may be obtained by parsing $[{}_1ab$ by $[{}_2b$.

The question whether we should keep the original $[{}_1ab$ is open. It clutters a bit the reference grammar, and any ab may still be parsed by $[{}_1a]_2$ and $[{}_2b$, but it may speed up the parses, and allows more leeway in the order in which the different elements are parsed. Moreover, it may be smart to keep these elements to maximize the number of parses we can find.

For example, consider the (supposed Markov) set of rules

$$\begin{aligned} & [{}_1abc \oplus [{}_1d \\ & \oplus [{}_2ab \oplus [{}_2e \\ & \oplus [{}_3bc \oplus [{}_3f. \end{aligned}$$

The following set would still be Markov:

$$\begin{aligned} & [{}_1]_2c \oplus [{}_1a]_3 \oplus [{}_1d(\oplus [{}_1abc) \\ & \oplus [{}_2ab \oplus [{}_2e \\ & \oplus [{}_3bc \oplus [{}_3f. \end{aligned}$$

It should be clear here that both $[_1]_2c$ and $[_1a]_3$ come from the same $[_1abc]$. If we did not keep the original $[_1abc]$, we would only get, say, $[_1]_2c$, and af would not be recognised as a syntagm. Running several classic parses (and deleting parsed elements), and merging the results, will give some of the missed elements, but not all (for example, those needing both $[_1]_2c$ and $[_1a]_3$ to parse them, say $[_4]_2ca[_3]$).

Remark however that the redundant elements can always be removed afterwards, by removing all elements that may be parsed by the grammar.

As mentioned, we do not want to modify the rules so that they generalize more, as in the previous section. Rather, we want to keep the same syntactic classes. This leads us to the following definitions.

Definition 3.4.1

A ruleset \mathcal{R} is weakly smaller than another ruleset \mathcal{R}' when, for any syntagm type i , the set B_i of syntagms for \mathcal{R} is contained in a set of syntagms B'_j of \mathcal{R}' .

$$\mathcal{R} \lesssim \mathcal{R}' \iff \forall i \exists j; B_i \subset B'_j.$$

Weak equality is of course the case when there is a bijection $f : \mathbb{N} \mapsto \mathbb{N}$ such that

$$\forall i; B_i = B'_{f(i)}.$$

The weak order on ruleset is obviously a reflexive partial order.

Definition 3.4.2

A ruleset \mathcal{R} is compatible with another \mathcal{R}' , when, the set of syntagms of any type is included in a equivalence class of $\sim_{\mathcal{R}'}$:

$$\mathcal{R} \Leftarrow \mathcal{R}' \iff \forall i \forall a, b \in B_i; a \sim_{\mathcal{R}'} b.$$

It is quite obvious that a ruleset weakly smaller than a compatible one is still compatible, and that the relation \Leftarrow is transitive.

It follows from this definition and from lemma 3.2.6 on page 133 that any ruleset compatible with a Markov ruleset, is also Markov.

Proposition 3.4.2

Let us take two rulesets $\mathcal{R}, \mathcal{R}'$. Then we have the following

$$\begin{aligned} \mathcal{R} \lesssim \mathcal{R}' &\implies \left(\forall s, s' \in D^+, s \sim_{\mathcal{R}} s' \implies s \sim_{\mathcal{R}'} s' \right) \\ \mathcal{R} \Leftarrow \mathcal{R}' &\implies \left(\forall s, s' \in D^+, s \sim_{\mathcal{R}} s' \implies s \sim_{\mathcal{R}'} s' \right). \end{aligned}$$

This means that a weakly bigger, compatible ruleset will define the same equivalence relation.

PROOF. Suppose $\mathcal{R} \lesssim \mathcal{R}'$, and $s \sim_{\mathcal{R}} s'$. By definition, we have a chain (x_j, y_j, y'_j) , $y_j, y'_j \in B_{k_j}$, so that $s = \alpha(x_0, y_0), \alpha(x_J, y'_J) = s'$ and $\alpha(x_{j-1}, y'_{j-1}) = \alpha(x_j, y_j)$. By hypothesis, there are k'_j so that $B_{k_j} \subset B'_{k'_j}$, and $y_j, y'_j \in B_{k_j}$. This means that $s \sim_{\mathcal{R}'} s'$.

Suppose now that $\mathcal{R} \Leftarrow \mathcal{R}'$, and $s \sim_{\mathcal{R}} s'$. Again, we have a chain (x_j, y_j, y'_j) , $y_j, y'_j \in B_{k_j}$, so that $s = \alpha(x_0, y_0), \alpha(x_J, y'_J) = s'$ and $\alpha(x_{j-1}, y'_{j-1}) = \alpha(x_j, y_j)$. This time, we know that $y_j \sim_{\mathcal{R}'} y'_j$, so that again $s \sim_{\mathcal{R}'} s'$. \square

The goal now will be to build, from a given ruleset, a bigger, compatible one. This can be done in the following way.

Proposition 3.4.3

Let \mathcal{R} be a Markov ruleset. Consider the sequence of grammars $(\mathcal{R}_i)_i$:

- $\mathcal{R}_0 = \mathcal{R}$,
- \mathcal{R}_{i+1} is a splitting of $\cup_{j \leq i} \mathcal{R}_j$ with reference grammar $\cup_{j \leq i} \mathcal{R}_j$.

Then $\cup_i \mathcal{R}_i$ is weakly bigger and compatible, and will be called the saturated ruleset, noted $\overline{\mathcal{R}}$.

PROOF. Let us begin with the fact that $\overline{\mathcal{R}} \gtrsim \mathcal{R}$. This comes from the fact that $\mathcal{R}_{i+1} \gtrsim \mathcal{R}_i$, which in turns comes from the fact that, if we put B_i the syntagm sets of \mathcal{R}_i , and B'_i those of \mathcal{R}_{i+1} , $B_i \subset B'_i$.

Now let us consider the compatibility. We will prove the weaker result:

Lemma 3.4.4

For any ruleset \mathcal{R} , for any $[_i \alpha(x, y), [_j y \in \mathcal{R}$, the ruleset $\mathcal{R}' = \mathcal{R} \oplus [_i \alpha(x,]_j)$ is compatible with \mathcal{R} .

PROOF. Let us note B_i the syntagms sets for \mathcal{R} , and B'_i the syntagms sets for \mathcal{R}' . Let us begin by taking a syntagm $s \in B_y$, and a parse \check{s} , $\Pi_y(s, \check{s}) > 0$.

Take any syntagm $a \in B'_i$, and a parse \check{a} , $\Pi_i(a, \check{a}) > 0$.

Consider now all c, t such that $\check{a} = \alpha(c, ({}_j t)_j)$ and $[_j \varsigma(t) = [_i \alpha(a,]_j)$. We will prove by induction on the number of such splits that a is equivalent to a syntagm of type i for \mathcal{G} .

If there is none, then for all c, t such that $\check{a} = \alpha(c, ({}_j t)_j)$, $[_j \varsigma(t) \in \mathcal{R}$, and $a \in B_i$. This initializes the induction.

If not, let us consider one of the deepest such x, y . This means that for all x', y' , if $t = \alpha(x', ({}_k y')_k)$, $[_k \varsigma(y') \in \mathcal{R}$, so $\varphi(t) \in B_k$. Moreover, since $\varsigma(t) \notin \mathcal{R}$, $\varsigma(t) = \alpha(x,]_j)$, and we can take (x', y') such that $t = \alpha(x', y')$, $\varsigma(x') = \varsigma(x)$, and $\varphi(y') \in B_y$. Obviously $\varphi(y') \sim s$, so that, if $\check{a}' = \alpha(c, ({}_j \alpha(x', \check{s}))_j)$, \check{a}' is still a syntagm for \mathcal{R}' , and $a = \varphi(\alpha(c, ({}_j \alpha(x', y'))_j)) \sim \varphi(\check{a}')$. Since \check{a}' has one

less (c, t) , its realization is equivalent by hypothesis to a syntagm for \mathcal{R} , and so is a .

Thus any \mathcal{R}' -syntagm of type i is equivalent to a \mathcal{R} -syntagm of type i . This ends the proof. \square

Since the ruleset $\overline{\mathcal{R}}$ is build by repetitive application of splits as in the lemma, the result is compatible with the first ruleset. \square

Proposition 3.4.5

The process defined in proposition 3.4.3 on the preceding page is stationary in finite time.

PROOF. The saturated ruleset is obviously finite, since the lengths of its elements are bounded by the maximum length of the elements of the original ruleset \mathcal{R}_0 . The dictionary is finite since no new indices for brackets are created.

Since each step adds at least a new element to the saturated ruleset (or the process ends), the process finishes a.s. in finite times. We could even compute an explicit bound using this argument, but chances are it is much larger than the best possible one. \square

3.4.4 Identification of labels

Nowhere in this study did we identify labels, except in the obvious case $[_i]_j$. However, if, at any point, we discover an element $s \in B_i \cap B_j$, while parsing, or simulating, for example, we can at once indentify the labels i and j .

Proposition 3.4.6

For any ruleset \mathcal{R} , if $B_i \cap B_j \neq \emptyset$ for some $i \neq j$, then if we define the ruleset \mathcal{R}' as the ruleset where all instances of j in \mathcal{R} was replaced by i , $\mathcal{R}' \Leftarrow \mathcal{R}$.

PROOF. Let us take $s \in B_i \cap B_j$. For any pair of parse trees $({}_k a)_k, ({}_k b)_k \in \check{B}_k$, we will show, by induction on their maximal depth, that $\varphi(a) \sim \varphi(b)$.

- If their depth is 1, that is, $a, b \in D^+$, either $k \neq i$ and $\varphi(a), \varphi(b) \in B_k$, hence the result, or $k = i$. In this case, by hypothesis, $\varphi(a) \sim s \sim \varphi(b)$.
- Suppose their depth is less than n . Say $a = \alpha(x, ({}_j y)_j)$. By hypothesis, there is $y' \in \check{B}_j$ such that $\varphi(a) \sim \varphi(\alpha(x, ({}_j y')_j))$. This means that there is $t \in \check{B}'_k$, $\varphi(a) \sim \varphi(t)$, and for any $x, y, t = \alpha(x, ({}_j y)_j) \implies y \in \check{B}_k$. The same holds true for b , with a parse t' .

Now, if $k \neq i$, $t, t' \in \check{B}_k$, and $\varphi(a) \sim \varphi(t) \sim \varphi(t') \sim \varphi(b)$. If $k = i$, then $t, t' \in \check{B}_i \cup \check{B}_j$, and $\varphi(a) \sim \varphi(t) \sim s \sim \varphi(t') \sim \varphi(b)$. \square

3.5 Toric grammars

The formalism we used to describe rulesets corresponds exactly to the definition of toric grammars introduced in chapter 1, more precisely, reference grammars (since we only consider local expressions), minus permutations. We will, in this section, use the term reference grammar in place of ruleset.

The construction described in section 3.4 on page 144 may be used to obtain the reference grammar needed for the communication model. Let us now study the properties of this model.

3.5.1 Split and merge processes using parsing

The parsing framework we defined in section 3.1.3 on page 123 can be used to define probabilities for the splitting part of a split-and-merge process, provided the parsing extension kernels are well chosen.

Definition 3.5.1

For any bottom-up parsing extension kernel ϖ that depends only on the surface structures, that is, for any $x, x' \in \check{B}_e \times \check{B}_{e'}, (iy)_i, (iy')_i \in B_i$,

$$\varpi\left(\alpha(x, y), \alpha(x, (iy)_i)\right) = \varpi\left(\alpha(x', y'), \alpha(x', (iy')_i)\right), \quad (3.5.1)$$

we can define a splitting rule

$$\beta_{\varpi}(\mathcal{G}) = \left\{ \mathcal{G}' \in \mathfrak{G} : \mathcal{G}' = \mathcal{G} \oplus \alpha((e, e'),]_i) \oplus [{}_i a \ominus \alpha((e, e'), a), \right. \\ \left. [{}_i a \in \mathcal{R}, \exists x \in \check{B}_e \times \check{B}_{e'}, y \in B_a, \varpi\left(\alpha(x, y), \alpha(x, (iy)_i)\right) > 0 \right\} \subset \beta_n(\mathcal{G}, \mathcal{R}).$$

We can even define a distribution on $\beta_{\varpi}(\mathcal{G})$ as the unique value of ϖ ,

$$\mathbb{P}\left(\mathcal{G} \oplus \alpha(e,]_i) \oplus [{}_i a \ominus \alpha(e, a)\right) = \varpi\left(\alpha(x, y), \alpha(x, (iy)_i)\right).$$

The split process thus defined will be noted \mathcal{S} .

We can complete the definition of split-and-merge process by using the uniform distribution on the merges $\alpha(\mathcal{G})$.

Using this definition, there is a one-to-one correspondence between split processes and a general parsing kernel.

Lemma 3.5.1

Given any string s , and a split $\mathcal{G} \in \beta_n^*(s, \mathcal{R})$, there is a bijection

$$T : \beta_n^*(s, \mathcal{R}) \longrightarrow \{t \in \mathcal{T}, s \in B_t\}.$$

PROOF. Let us take a sequence of splits $s = \mathcal{G}_0, \dots, \mathcal{G}_n$, $\mathcal{G}_i \in \beta_n(\mathcal{G}_{i-1}, \mathcal{R})$. We can reverse this process by building a sequence $(\check{\mathcal{G}}_i)$, such that $\varphi(\check{\mathcal{G}}_i) = \mathcal{G}_i$. We will first let $\check{\mathcal{G}}_n = \mathcal{G}_n$, and, given $\check{\mathcal{G}}_i$, if

$$\mathcal{G}_i = \mathcal{G}_{i-1} \oplus \alpha(e,]_i) \oplus [{}_i a \ominus \alpha(e, a),$$

put

$$\check{\mathcal{G}}_{i-1} = \check{\mathcal{G}}_i \ominus \alpha(e,]_i) \ominus [{}_i \check{a} \oplus \alpha(e, ({}_i \check{a})_i),$$

where \check{a} is the expression corresponding to a in $\check{\mathcal{G}}_i$.

This construction verifies the condition $\varphi(\check{\mathcal{G}}_i) = \mathcal{G}_i$, so in particular we have that $t = \varphi(\check{\mathcal{G}}_0) \in \mathcal{T}$ and $\varphi(t) = s$.

We will now prove that this t does not depend on the actual sequence of splits. This comes from the fact that the surface structure $\varsigma(y)$ of any substring y of t must be in \mathcal{G}_n . This means that we can define $T(\mathcal{G}) = t$. All that is left to prove is that this is a bijection.

In order to do so, consider the converse construction. Take a parse tree t . There is a sequence of parse trees $s = t_0, \dots, t_n = t$, such that there is, for any i , a context $x \in \mathcal{T}^2$, an expression a , $[{}_i a \in \mathcal{R}$ and a tree $y \in B_a$, with $t_i = \alpha(x, y)$, $t_{i+1} = \alpha(x, ({}_i y)_i)$. This sequence defines a sequence of splits, with the rule

$$\mathcal{G}_{i+1} = \mathcal{G}_i \oplus \alpha(x,]_i) \oplus [{}_i y \ominus \alpha(x, y),$$

as mentioned before, and the end result \mathcal{G} does not depend on the actual sequence, since

$$\mathcal{G}(t) = \varsigma(t) \oplus \bigoplus_{[{}_j e \in \mathcal{R}} \chi(t, j, e) \otimes [{}_j e \quad (3.5.2)$$

This obviously defines the inverse of T . □

Remark that this lemma states also that, as soon as we have a parsing kernel such that

$$\Pi(s, t) > 0 \Leftrightarrow \beta_n(\mathcal{G}(t), \mathcal{R}) = \emptyset, t \succ s, \quad (3.5.3)$$

we can simply use this parsing kernel to simulate a split process.

3.5.2 Reversible split and merge process

The results above showed that the (generalized) split and merge process could be used to explore the communication classes of k . However, such a process will probably distort the distribution, and can only be used to study the support of the language. We will here propose a method to build a reversible crossing-over dynamics, using the general Metropolis method described in section 3.D on page 165, using split and merge processes.

We will here work on texts, that is, on sums of Dirac masses at possibly repeated sentences. As such, we will note \mathfrak{T} the set of texts of strings, $\bar{\mathfrak{T}}$ that of texts of parse trees, \mathfrak{G} for texts of expressions, $\bar{\mathfrak{G}}$ for texts of well-parenthesized expressions. The length of texts, that is, their mass, is fixed at n .

We begin with a text $S_{1:n} \in \mathfrak{T}$, and a reference grammar $\mathcal{R} \in \mathfrak{G}$. The basic idea is to split the text according to a parse of its sentences, and then to merge it back, effectively crossing over constituents.

Let us then consider a general parsing kernel $\bar{\Pi}$. We also have the splitting kernel from parsing trees to toric grammars $t \mapsto \mathcal{G}(t)$ defined in eq. (3.5.2) on the facing page. We now want to build a production process that will merge back the elements of \mathcal{G} into a text.

Remark first that any text that can be obtained by regular merges (where we merge any $[_i\alpha(x,]_j)$ with any $[_jy)$ can also be obtained by merging only $[_0\alpha(x,]_j)$ with $[_jy)$, in left-to-right order. Indeed, given a merge result, we can simply follow the order of merges given by taking outer parentheses first, in left-to-right order. If we moreover index the global brackets $[_0$ of the grammar, the order of parentheses are fixed.

Consider now the merge transformations on grammars

$$\left\{ \mathcal{M}_k(\mathcal{G}, \mathcal{G}'); \mathcal{G}' = \mathcal{G} \ominus [_0^k\alpha(x,]_i) \ominus [_iy \oplus [_0^k\alpha(x, (iy)_i) \right\}.$$

We can take, for example,

$$\mathcal{M}_k\left(\mathcal{G}, \mathcal{G} \ominus [_0^k\alpha(x,]_i) \ominus [_iy \oplus [_0^k\alpha(x, (iy)_i)\right) = \frac{\mathbf{1}(x_1 \in \check{D}^*)}{\mathcal{G}([_iS^*)},$$

which corresponds to taking the closing brackets from left-to-right order in the k th sentence of \mathcal{G} (with weights), and a matching $[_iy$.

If we put, for short, $\mathcal{M}(\mathcal{G}) = \mathcal{M}_k(\mathcal{G})$, where the merge is on the first incomplete sentence, that is, $k = \min\{m, \mathcal{G}([_0^m\check{D}^*) = 0\}$, the Markov chain defined by the kernel \mathcal{M} is stationary in finite time, at most $\sum_{i \in \mathbb{N} \setminus \{0\}} \mathcal{G}([_iS^*)$, since each steps removes one $[_i$. We can then define the production process as

$$\mathcal{P} = \mathcal{M}^\infty.$$

The trick of indexing the sentences allows us to have a deterministic order of the merges, at least for the closing brackets, which means that for any two grammars $\mathcal{G}, \mathcal{G}'$, if $\mathcal{P}(\mathcal{G}, \mathcal{G}') > 0$,

$$\mathcal{P}(\mathcal{G}, \mathcal{G}') = \prod_i \frac{\mathcal{G}'([_iS^*)!}{\mathcal{G}([_iS^*)!}.$$

Remark that the result of \mathcal{P} is not necessarily a text (that is, a grammar with no $[_i$, $i > 0$ left). We will not explore how we could recover from this failure, and

merge back the remaining elements, since we think this would be unnecessarily complicated. We prefer to simply reject the result if $\mathcal{G}'(\lfloor_+ D^+) \neq 0$ and repeat the merge until success. Failures will only happen when the grammar is recursive, and be actually rare unless the grammar is strongly recursive. As we mentioned in section 1.3.2 on page 28, this was indeed a rare occurrence in our experiments.

Since the Metropolis algorithm used to build the reversible dynamics already uses some sort of rejection process, we can even integrate it to the chain using the projection kernel

$$\begin{aligned} \tau : \bar{\mathfrak{G}} &\longrightarrow \mathcal{M}_1(\bar{\mathfrak{X}}) \\ \mathcal{G} &\longmapsto \begin{cases} \delta_{\mathcal{G}} & \text{if } \mathcal{G} \in \bar{\mathfrak{X}} \\ \delta_{\epsilon^n} & \text{otherwise,} \end{cases} \end{aligned}$$

and conversely $\tau^{-1}(t) = \delta_t$, $t \in \bar{\mathfrak{X}}$.

We finally obtain the following chain, with which we can define a reversible dynamics on texts, using section 3.D on page 165:

$$\begin{array}{ccccccc} \bar{\mathfrak{X}} & \xrightarrow{\bar{\Pi}} & \bar{\mathfrak{X}}^+ & \longrightarrow & \bar{\mathfrak{X}} & \xrightarrow{\tau^{-1}} & \bar{\mathfrak{G}} \\ & & & & & & \searrow \mathcal{S} \\ & & & & & & \mathfrak{G} \\ & & & & & & \swarrow \mathcal{P} \\ \bar{\mathfrak{X}} & \longleftarrow & \bar{\mathfrak{X}}^+ & \longleftarrow & \bar{\mathfrak{X}} & \xleftarrow{\tau} & \bar{\mathfrak{G}} \end{array}$$

Non-labelled arrows of course represent projections on the last or first element of the string.

Proposition 3.5.2

The dynamics on texts defined by

$$\begin{aligned} \mathcal{SM}(t, t') = & \sum_{\substack{\tilde{t}, \tilde{t}' \in \bar{\mathfrak{X}}^+ \\ \check{t}, \check{t}' \in \bar{\mathfrak{X}} \\ \check{\mathcal{G}}, \check{\mathcal{G}}' \in \bar{\mathfrak{G}} \\ \mathcal{G} \in \mathfrak{G}}} \bar{\Pi}(t, \tilde{t}) \mathbf{1}(e(\tilde{t}) = \check{t} = \check{\mathcal{G}}) \mathcal{S}(\check{\mathcal{G}}, \mathcal{G}) \mathcal{P}(\mathcal{G}, \check{\mathcal{G}}') \tau(\check{\mathcal{G}}', \check{t}') \zeta(\check{t}', \tilde{t}') \mathbf{1}(t' = b(\tilde{t}')) \\ & \times \left(\frac{\zeta(\check{t}, \tilde{t})}{\bar{\Pi}(t, \tilde{t})} \frac{\tau(\check{\mathcal{G}}, \check{t})}{\mathbf{1}(\check{t} = \check{\mathcal{G}})} \frac{\mathcal{P}(\mathcal{G}, \check{\mathcal{G}})}{\mathcal{S}(\check{\mathcal{G}}, \mathcal{G})} \frac{\mathcal{S}(\check{\mathcal{G}}, \mathcal{G}')}{\mathcal{P}(\mathcal{G}, \check{\mathcal{G}}')} \frac{\mathbf{1}(\check{t}' = \check{\mathcal{G}}')}{\tau(\check{\mathcal{G}}', \check{t}')} \frac{\bar{\Pi}(t', \tilde{t}')}{\zeta(\check{t}', \tilde{t}')} \wedge \mathbf{1} \right) \end{aligned}$$

is reversible for $\mathbb{P}_S^{\otimes n}$.

PROOF. As always, the expression of $\mathcal{SM}(t, t')$ is symmetric in t, t' , so we only have to prove that $\mathbb{P}_S^{\otimes n}(t) = \mathbb{P}_S^{\otimes n}(t')$.

Given any sequence $\tilde{t}, \tilde{t}' \in \bar{\mathfrak{T}}^+, \check{t}, \check{t}' \in \bar{\mathfrak{T}}, \check{\mathcal{G}}, \check{\mathcal{G}}' \in \bar{\mathfrak{G}}, \mathcal{G} \in \mathfrak{G}$ with a positive contribution to the sum, according to lemma 3.2.1 on page 125,

$$\begin{aligned}
\mathbb{P}_S^{\otimes n}(t) &= \prod_{k=1}^n \mathbb{P}_S(B_{\zeta(\check{t}_k)}) q_{B_{\zeta(\check{t}_k)}}(t_k) \\
&= \prod_{k=1}^n \left(\mathbb{P}_S(B_{\zeta(\check{t}_k)}) A_{\zeta(\check{t}_k)} \prod_{[j,e \in \mathcal{R}]} [A_e q_{B_j}(B_e)]^{\chi(\check{t}_k, j, e)} \right) \\
&= \prod_{k=1}^n \left(\sum_{[a \in \mathcal{G}]_0^k} \mathbb{P}_S(B_a) A_a \prod_{[j,e \in \mathcal{R}]} [A_e q_{B_j}(B_e)]^{\mathcal{G}([j,e])} \right) \\
&= \prod_{k=1}^n \left(\mathbb{P}_S(B_{\zeta(\check{t}'_k)}) A_{\zeta(\check{t}'_k)} \prod_{[j,e \in \mathcal{R}]} [A_e q_{B_j}(B_e)]^{\chi(\check{t}'_k, j, e)} \right) \\
&= \prod_{k=1}^n \mathbb{P}_S(B_{\zeta(\check{t}'_k)}) q_{B_{\zeta(\check{t}'_k)}}(t'_k) \\
&= \mathbb{P}_S^{\otimes n}(t'). \quad \square
\end{aligned}$$

3.6 Estimating the language distribution

Let us remark that if B is a Markov substitute set such that $B \cap \text{supp}(\mathbb{P}_S) \neq \emptyset$, then $B \subset \text{supp}(\mathbb{P}_S)$ and $\mathbb{P}_{S|S \in B} = q_B$. As such, we can aim to identify a collection of Markov substitute sets that will cover the entire distribution, and from their relative frequencies, deduce the entire distribution.

If we put for any $s \in D^+$, $C(s) = \text{supp}(\delta_s \sum_{j=0}^{\infty} k^j)$ the syntactic category of s , then $C(s)$ is a Markov substitute set and

$$\mathbb{P}(S) = \mathbb{E}(q_{C(S)}),$$

so that

$$\hat{P} = \frac{1}{n} \sum_{i=1}^n q_{C(S_i)}$$

is an unbiased estimate of \mathbb{P}_S .

This formula can be used to estimate the probability of any sentence, provided we have access to the substitute measures $q_{C(s)}$.

Let us remark here that while

$$\text{supp}(\delta_s k) = \bigcup_{e \in S^+, s \in B_e} B_e$$

is a context-free language (with multiple rewriting rules for the start symbols), since $E_s = \{e \in S^+, s \in B_e\}$ is finite, this is not the case anymore for $C(s)$, because nothing prevents the number of surface structures of strings in $C(s)$ to be infinite. Differently put, the syntactic category of a string s is not necessarily a context-free language, though it is still a countable union of context-free languages

$$C(s) = \bigcup_{e \in E(s)} B_e,$$

with

$$E(s) = \left\{ e \in S^+, \sum_{j=0}^{\infty} k^j(s, B_e) > 0 \right\}.$$

For any fixed, finite set E , $C(E) = \bigcup_{e \in E} B_e$ is still a context-free language, and as such its substitute measure $q_{C(E)}$ may be estimated using the tools presented in section 3.2.2 on page 127. This means that we can estimate $E(s)$ by

$$\widehat{E}(s) = \left\{ e \in S^+, \sum_{j=0}^N k^j(s, B_e) > 0 \right\},$$

which is finite and whose substitute measure may be estimated. The error

$$\mathbb{P}_S \left(\bigcup_{s \in \text{supp}(\mathbb{P}_s)} C(s) \setminus \bigcup_{i=1}^n C(\widehat{E}(S_i)) \right)$$

may be estimated using the missing mass estimator (see [Goo53, MO03]).

Another view is to remark that, for any string s , and any expression $e \in S^+$, such that $s \in B_e$,

$$\mathbb{P}_S(s) = \mathbb{P}_S(B_e) q_{B_e}(s).$$

The substitute measure q_{B_e} , as we have already seen, is easy to estimate. To estimate $\mathbb{P}_S(B_e)$, we want to compute

$$\mathbb{1}(s \in B_e).$$

This function is a very simple parsing test, and is also very easy to compute. Thus we can estimate the probability of a given sentence s by drawing one or better yet, multiple parsing trees t_1, \dots, t_m , $t_k \succcurlyeq s$, which gives several expressions $e_k = \varsigma(t_k)$, and set

$$\widehat{\mathbb{P}}(s) = \frac{1}{m} \sum_{k=1}^m \frac{1}{n} \sum_{i=1}^n \mathbb{1}(S_i \in B_{e_k}) q_{B_{e_k}}(s).$$

We thus obtain an estimator without bias of $\mathbb{P}(s)$, on the condition that we use a split sample scheme to learn the substitute measures on a different sample.

This approach can also be used as a grammaticality test. In this test we can choose not to use the substitute measures $q_{B_{e_k}}$, since we know they must have by construction a full support. We can indeed define the set of global structures of the language as

$$\{e \in S^+, \mathbb{P}(B_e) > 0\}.$$

We can then consider the set of maximal surface structures \mathcal{E}_g for the inclusion relation $B_{e'} \subset B_e$. The support of \mathbb{P}_S can then be characterized by the property

$$\mathbb{P}_S(s) > 0 \iff \exists e \in \mathcal{E}_g, s \in B_e.$$

For any string $s \in D^+$, the set $\mathcal{E}_g(s) = \mathcal{E}_g \cap \{e, s \in B_e\}$ is easy to compute from parsing, and \mathcal{E}_g can then easily be estimated by

$$\widehat{\mathcal{E}}_g = \bigcup_i \mathcal{E}_g(s),$$

and the error

$$\mathbb{P}_S \left(\bigcup_{e \in \mathcal{E}_g} B_e \setminus \bigcup_{e \in \widehat{\mathcal{E}}_g} B_e \right)$$

may be once again estimated with the missing mass estimator.

3.A Support of the split and merge process and substitutions

We will study in this section the action of the split and merge process on the syntactic classes.

Let us consider a grammar \mathcal{R}_0 , a corresponding reversible dynamics k_0 , and its associated equivalence relation \sim_0 . These can obviously be extended to texts. Take a compatible reference grammar \mathcal{R} , and let SMP be a split-and-merge process with reference \mathcal{R} , as defined in section 3.5.1 on page 151. We take a compatible grammar, and not simply the same grammar, for the first result.

Proposition 3.A.1

For any pair of texts $(t, t') \in \mathfrak{T}$,

$$SMP(t, t') > 0 \implies t \sim_0 t'.$$

In particular, the split and merge process stays in the support of the language.

PROOF. We will focus in this proof on a single sentence, since the substitute process is independent on the sentences. Let us take a sentence s , and its result, s' . We may consider s' as a syntactic tree \check{s}' , given by the merge operation sequence, and we can do the same with s , its structure \check{s} being given by the parse algorithm used to split. The only common point between \check{s} and \check{s}' is their surface structure.

We want to prove that substitutions can link any two \mathcal{R} -correct syntagms \check{s} and \check{s}' , as soon as they share the same surface structure.

Since \mathcal{R} is compatible, any two trees in the same \check{B}_i are equivalent. Since \check{s} and \check{s}' share the same global structure, any tree inside a pair of outer parentheses of \check{s} has a counterpart in \check{s}' , and each may be substituted to its counterparts in \check{s}' , which gives the result. \square

The converse is obviously not true. The simple compatibility is not sufficient anymore, since the null grammar is compatible with any grammar, but has a trivial split and merge process. So let us now suppose that \mathcal{R} defines the same equivalence relation. Under no further hypothesis, this is still not sufficient, since we need at least all elements of B_i to be in the text to hope to obtain them after a pass of the split and merge process.

However, if we define a slightly different split and merge process, we can get a converse result.

Definition 3.A.1

We will define a generalized split and merge process GSMP as the result of a split process, followed by a merge from the obtained global expressions, using the elements of \mathcal{R} , without any restriction on the number of uses. Such a merge will be noted GMP.

The actual distribution of the merges may be taken as uniform on all possible merges, but since we are here only interested in supports, the actual distribution is not crucial.

Obviously, such a process is independent on each sentence of the text, and proposition 3.A.1 on the preceding page is still true with this new process.

Such a definition, however, is still not sufficient to have the reverse implication, as shown in the following example. Consider the reference grammar (compatible with itself)

$$\begin{aligned} & [{}_1\alpha \oplus [{}_1ab \oplus \\ & [{}_2\gamma \oplus [{}_2bc. \end{aligned}$$

Obviously, $k^2(\alpha c, a\gamma) > 0$. However,

$$GSMP(\alpha c, \cdot) = \frac{1}{2}\delta(\alpha c) + \frac{1}{2}\delta(abc),$$

so that $GSMP(\alpha c, a\gamma) = 0$.

Nonetheless, we have the following, if k is the invariant dynamics defined from \mathcal{R} .

Proposition 3.A.2

For any pair of sentences (s, s') ,

$$k(s, s') > 0 \implies GSMP(s, s') > 0,$$

and as such,

$$\sum_k GSMP^k(s, s') > 0 \iff s \sim s'.$$

This means that, more than staying in the support of the language, the generalized split and merge process will visit all sentences in each of the communication classes of k , provided it begins in it.

PROOF. The proof is actually quite straightforward: if $k(s, s') > 0$, there are two $t, t' \in B_e$ for a certain expression e , and two chains \tilde{t} and \tilde{t}' from s to t and s' to t' respectively so that $\Pi(s, \tilde{t}) > 0$ and $\zeta(\tilde{t}', s') > 0$. Obviously then $\mathcal{S}(s, e) > 0$, and $GMP(e, s') > 0$, so $GSMP(s, s') > 0$. \square

3.B Parsing using a ruleset

We will now describe how to build the different parses we used in chapter 3. In order to do so, we will use a well-known parsing algorithm for context-free grammar, the CYK algorithm.

The version of the CYK algorithm we will (briefly) outline here is that proposed in [LL09], which works on a grammar in the Binary Normal Form, that is, a grammar for which all right-hand sides of rules are of length at most 2.

Let us then take a set of context-free rules $(D,]_R, \mathcal{R})$, where the rule index set $R \subset \mathbb{N}$ is finite. Its size $|\mathcal{R}|$ is by definition the length of the concatenation of all the rules in \mathcal{R} . Its transformation to its binary normal form \mathcal{R}' may be obtained in time $O(|\mathcal{R}|)$, and the resulting grammar is of size $O(|\mathcal{R}|)$, and the size of its non-terminal set is also $O(|\mathcal{R}|)$. The transformation is obtained by replacing any rule of length more than 2, $]_i \omega_1 \dots \omega_n$ by the rules

$$]_i \omega_1]_{i_1}; \dots;]_{i_{j-1}} \omega_j]_{i_j} \dots;]_{i_{n-1}} \omega_{n-1} \omega_n.$$

Let us suppose that we also built the unit graph

$$U = \{ (]_a, e), e \in D \cup]_+, (]_a e \in \mathcal{R} \},$$

which is built in time and space $O(|\mathcal{R}|)$, simply by adding the correct edges for each unit rule. This makes the computation of $U^*(\downarrow_a)$ (the set of reachable states from \downarrow_a in the graph U) possible in time $O(|\mathcal{R}|)$.

Given a string $s_{1:n}$ to parse, the algorithm will actually compute the sets

$$T_v = \{x \in D \cup \downarrow_+, v \in B_x\}$$

for all substring v of s . These sets will be stored in a table T of size n^2 . Each entry in this table is itself a boolean array indexed by pairs (\downarrow_i, k) , where $i \in R$ and $k \in \llbracket 0, n \rrbracket$ (so the array is of size $|R| \times (n+1)$), and whose diagonal is filled with the terminal symbols s_1, \dots, s_n .

The table is then filled so that $T(i, j)(\downarrow_a, 0) = 1$ when there is one symbol b such that $s_{i:j} \in B_a \cup B_b$, $\downarrow_a b \in \mathcal{R}$, and $T(i, j)(\downarrow_a, k) = 1$ when there are two symbols $b, c \in D \cup \downarrow_+$, such that $s_{i:k} \in B_b$, and $s_{k+1:j} \in B_c$, where $\downarrow_a bc \in \mathcal{R}$. In particular, $T(i, j)(\downarrow_a, k) = 1$ for some k means that $s_{i:j} \in B_a$.

The table is then filled so that $s_{i:j} \in B_a$, if, and only if, $s_{i:k} \in B_{e_1}$, and $s_{k+1:j} \in B_{e_2}$, where $\downarrow_a e_1 e_2 \in \mathcal{R}$ if, and only if, $T(i, j)(\downarrow_a, k) = 1$.

This is done by the following algorithm, where all cells in T are initially set to 0:

```

for i=1, ..., n do
  T(i, i) += s_i
  T(i, i) += (U*(T(i, i)), 0)
done
for j=2, ..., n do
  for i=1, ..., j-1 do
    for k=i, ..., j-1 do
      for all rule [a XY do
        if T(i, k)(X, *) = 1 and T(k+1, j)(Y, *) = 1 then
          T(i, j) += (a, k)
        done
      done
    done
  done
  T(i, j) += (U*(T(i, j)), 0)
done
done

```

with the notation $T(i, j) += A$ meaning that we change to 1 the value of each cell of $T(i, j)$ whose index is in A . Again, $U*(T(i, j))$ is the set of symbols produced by the unitary rules of U whose left-hand side is equal to $T(i, j)$ (seen as a list of symbols). On the other hand, $T(i, k)(X, *) = 1$ is the test that at least for one r , $T(i, k)(X, r) = 1$.

This algorithm is in $O(n^4|\mathcal{R}|)$ time, and we could make it $O(n^3|\mathcal{R}|)$ (classical complexity) if we do not need to keep the index k . However, the k is necessary to be able to draw a certain parse from the table, which we'll need in section 3.B.2 on the next page. The space requirement is $O(n^3|\mathcal{R}|)$, $O(n^2|\mathcal{R}|)$ if we do not need the k .

With such a table, we can now present our parsing algorithms. As mentioned before, we are working on rulesets, and not complete context-free grammars. We are missing the start symbol. This means that we have two different approaches to the parsing. In one case, we seek a parsing kernel Π_e that stays inside the context-free language \check{B}_e , which is the classic framework. In the other, we are not in a context-free language, but in the whole \mathcal{T} . We do not have, in this case, a start symbol (or rather, we have an infinity of them).

3.B.1 Parsing general syntagms

Parsing in \mathcal{T} does not require the surface structure of the result to be of a particular form. This means that we can build the parse bottom-up.

Let us begin with the observation that any expression e such that $t \in B_e$ can be described by a succession of cells $(T(i_k, j_k, r_k) > 0)_{1 \leq k \leq K}$, so that $i_{k+1} = i_k + j_k$, $i_0 = 1$, $i_K + j_K = \ell(s)$. Indeed, in this case, $e =]_{i(r_1)} \dots]_{i(r_K)}$. The lengths j_k are additional informations, corresponding, for $s \preceq t \in \check{B}_e$, to the lengths of the contents of each outer parentheses of t .

Definition 3.B.1

For any parse tree $t \in \mathcal{T}$, let us put

$$E(t) = \{t' = \alpha(x, (iy)_i), (x, y) \in \mathcal{T}^3, t = \alpha(x, y), \varsigma(y) \in \mathcal{R}\}.$$

The general parsing extension is defined as

$$\varpi(t, t') = \frac{\mathbf{1}(t' \in E(t))}{|E(t)|}.$$

Proposition 3.B.1

The parsing kernel Π defined from the general parsing extension kernel is such that

$$\text{supp}(\Pi(t, \cdot)) = \{t' \in \mathcal{T}, \varphi(t') = \varphi(t), E(t') = \emptyset\},$$

that is, Π charges all the maximal parses of t .

PROOF. It is quite obvious that if $\Pi(t, t') > 0$, $t' \succcurlyeq t$ and $E(t') = \emptyset$.

Conversely, let us prove that all maximal parses are charged by Π . Let us take such a parse t' . Consider the sequence of trees $t = t_0 \succcurlyeq \dots \succcurlyeq t_n = t'$, where we add the parentheses of t' one pair at a time, following a left-to-right, bottom-up path. It is quite obvious that $\varpi(t, t') > 0$ and $\text{supp}(\varpi(t', \cdot)) = \emptyset$, so $\Pi(t, t') > 0$. \square

As before, we will be more interested in the general parsing kernel $\bar{\Pi}$ defined from ϖ , which keeps track of the whole sequence of trees built using ϖ .

The set $E(t)$ is quite easy to compute from the parse table of $\varphi(t)$, so the simulation and computation of $\bar{\Pi}(s, \cdot)$ is easy.

For the definition of reversible dynamics, we also need a converse kernel ζ that removes parentheses. Given a tree $T \in \mathcal{T}$, writing $s = \varphi(T)$, we will define the parsing extensions relative to T .

Definition 3.B.2

For any parse tree $t \in \mathcal{T}$, let us put

$$B_T(t) = \left\{ t' = \alpha(x, (iy)_i), (x, y) \in \mathcal{T}^3, t = \alpha(x, y), \varsigma(y) \in \mathcal{R}, T \succcurlyeq t' \right\}.$$

The general parsing extension is defined as

$$\rho_T(t, t') = \frac{\mathbb{1}(t' \in B_T(t))}{|B_T(t)|}.$$

Once again, the set $B_T(t)$ is quite easy to compute from the parse table of $\varphi(t)$, so the simulation and computation of $\zeta(T, \cdot) = \rho_T^\infty(T, \cdot)$ is easy.

The Metropolis algorithm will consider, for a sequence of trees \tilde{t}_i , $0 \leq i \leq m$ the ratios

$$\frac{\varpi(\tilde{t}_i, \tilde{t}_{i+1})}{\rho_{\tilde{t}_m}(\tilde{t}_i, \tilde{t}_{i+1})} = \frac{|E(\tilde{t}_i)|}{|B_{\tilde{t}_m}(\tilde{t}_i)|} \mathbb{1}(\tilde{t}_{i+1} \in E(\tilde{t}_i)).$$

While these ratios are not necessarily close to 1, they do not depend on \tilde{t}_{i+1} (apart the obvious requirement $\tilde{t}_{i+1} \in E(\tilde{t}_i)$). The actual values of these ratios depend on the ambiguity of the grammar, and a sentence for which only one maximal parse exists will have all the ratios equal to 1.

3.B.2 Parsing inside syntagms of certain type

Section 3.2.3 on page 129 required that we had a parsing Π_e inside of a particular language B_e . This is actually the main goal of standard parsing algorithms such as CYK, but we need to go a little further, for we need an easy way to draw such a parse, and to compute $\Pi_e(s, t)$ for any string s and tree t . This can easily be done by keeping in the table cell $T(i, j)$ pointers to the table cells $T(i, i+k)$ and $T(i+k, j)$ used to fill it, which can be done by adding k to the indices of the array in $T(i, j)$, as described before.

We have to consider, for any expression e , the set of incomplete parse trees from e , that is

$$T_e = \left\{ t \in \check{D}^*, \exists t' \in \check{B}_e, t' \succcurlyeq t \right\}.$$

Definition 3.B.3

For any incomplete parse tree $t \in T_e$, let us put

$$D(t) = \left\{ t' = \alpha(x, ({}_k(iy)_i(jy')_j)_k), x \in \epsilon \times \check{D}^*, y, y' \in D^+, \right. \\ \left. t = \alpha(x, ({}_kyy')_k), \exists (i\check{y})_i \in B_i, (j\check{y}')_j \in B_j, \right. \\ \left. [{}_k]_i]_j \in \mathcal{R}, \varphi(\check{y}) = y, \varphi(\check{y}') = y' \right\}.$$

The parsing extension from e is defined as

$$\varpi(t, t') = \frac{\mathbb{1}(t' \in D(t))}{|D(t)|}.$$

It is quite obvious that if $t \in T_e$, $D_t \subset T_e$. The set $D(t)$ is, again, easy to compute from the table, on the condition that pointers are kept as described above.

The parsing kernel may then be defined as

$$\Pi_e(s, t) = \varpi^\infty(({}_0s)_0, t) \\ = \exp\left(- \sum_{\substack{x \in (\check{D}^*)^2, \\ y \in \mathcal{T}, k \in \mathbb{N}}} \mathbb{1}(t = \alpha(x, ({}_ky)_k)) \log(|D(\alpha(x, ({}_k\varphi(y))_k))|)\right).$$

3.C Building crossing-over tree kernels

The construction of reversible dynamics defined in section 3.2.5 on page 134 requires a crossing-over tree kernel $(\mathfrak{s}(t, t'), t, t' \in \check{\mathfrak{T}})$. We will present here one way of building it.

In order to do so, it will be useful to index the parentheses of tree texts. We will then use a new parenthesized dictionary $\check{D} = D \cup \{({}_i^k)_i^k, i \in \mathbb{N}, k \in \mathbb{N}\}$. The i will still be called labels, and the new k , indices. We build as before the set of trees $\check{\mathcal{T}}$, with the additional requirement that matching parentheses have the same index, and that no two opening or closing parentheses have the same index. We of course define all the other sets as before, in particular $\check{\mathfrak{T}} = \check{\mathcal{T}}^n$.

For any text $t \in \check{\mathfrak{T}}$, let $t^\mathbb{1}$ be the same tree with the parentheses indexed in the order of the opening ones (left-to-right, depth first order). For any indexed tree text $\dot{t} \in \check{\mathfrak{T}}$, any permutation $\sigma \in \mathfrak{S}_m$, where m is the number of parentheses in \dot{t} , let us write as \dot{t}^σ the text where the pair of parentheses in \dot{t} with index k is instead indexed by $\sigma(k)$, that is, if $\dot{t}_i = \alpha(x, ({}_ky)_k)$, $\dot{t}_i^\sigma = \alpha(x, ({}^{\sigma(k)}y)^{\sigma(k)})$. It is quite obvious that $(\dot{t}^\sigma)^{\sigma'} = \dot{t}^{\sigma \circ \sigma'}$, and we will expand this operator on non-indexed tree texts with the notation $t^\sigma = (t^\mathbb{1})^\sigma$.

For any indexed tree text \dot{t} with m pairs of parentheses, let $\phi(t)$ be the tree itself, without parenthesis indices, and $\sigma_t \in \mathfrak{S}_m$ be the permutation such that $\dot{t} = \phi(t)^{\sigma_t}$.

Let us begin by defining an indexing kernel

$$\mathfrak{J}(t, \dot{t}) = \frac{1}{m!}, \quad t \in \bar{\mathfrak{T}}, \dot{t} \in \dot{\mathfrak{T}}, \dot{t} = t^\sigma, \sigma \in \mathfrak{S}_m.$$

Consider, for any indexed tree text $\dot{t} \in \dot{\mathfrak{T}}$, the sets

$$\begin{aligned} A_K(\dot{t}) &= \left\{ \dot{t}' \in \dot{\mathfrak{T}}, \right. \\ &\quad \left. \begin{aligned} \dot{t}'_u &= \alpha(x, \binom{k}{j} y_j^k), \dot{t}'_v = \alpha(x', \binom{k'}{j'} y_{j'}^{k'}), \dot{t}'_u = \alpha(x, \binom{k'}{j'} y_{j'}^{k'}), \dot{t}'_v = \alpha(x', \binom{k}{j} y_j^k), \\ u, v &\in \mathbb{N}, j \in \mathbb{N} \setminus \{0\}, x, x' \in (\dot{D}^*)^2, y, y' \in \dot{\mathcal{J}}, \sigma_t^{-1}(k') \geq \sigma_t^{-1}(k) = K \end{aligned} \right\}, \\ B_K(\dot{t}) &= \left\{ \dot{t}' \in \dot{\mathfrak{T}}, \right. \\ &\quad \left. \begin{aligned} \dot{t}'_i &= \alpha\left(\alpha(x, \binom{k}{j} y_j^k), \binom{k'}{j'} y_{j'}^{k'}\right), \dot{t}'_i = \alpha\left(\alpha(x, \binom{k'}{j'} y_{j'}^{k'}), \binom{k}{j} y_j^k\right), \\ i &\in \mathbb{N}, j \in \mathbb{N} \setminus \{0\}, x, x' \in (\dot{D}^*)^2, y, y' \in \dot{\mathcal{J}}, \sigma_t^{-1}(k') \geq \sigma_t^{-1}(k) = K \end{aligned} \right\}, \\ C_K(\dot{t}) &= A_K(\dot{t}) \cup B_K(\dot{t}). \end{aligned}$$

We can now define the individual crossing-over kernels

$$\mathfrak{s}_K(\dot{t}, \dot{t}') = \frac{\mathbb{1}(\dot{t}' \in C_K(\dot{t}))}{|C_K(\dot{t})|}.$$

These kernels simply swap the contents of the K th pair of parentheses (that has index $\sigma(K)$) with the contents of a pair of parentheses (with same label) further in the text (but not inside the parentheses). We could add one to the numerator if we wanted to allow no crossing-over to take place (and complete the kernel by giving a positive probability to $\dot{t} = \dot{t}'$).

We can now define the crossing-over kernel $\mathfrak{s} = \mathfrak{s}_1 \circ \dots \circ \mathfrak{s}_m$. Simulation of \mathfrak{s} , as seen, is easy. Moreover, it is apparent from the definition that for any pair of indexed tree texts \dot{t}, \dot{t}' , $\mathfrak{s}(\dot{t}, \dot{t}') = \mathfrak{s}(\dot{t}^\sigma, \dot{t}'^\sigma)$.

Hence, to calculate $\mathfrak{s}(\dot{t}, \dot{t}')$, we can go back to the case where $\sigma_t = \mathbb{1}$. In this case, if $\mathfrak{s}(\dot{t}, \dot{t}') > 0$, the pair of parentheses that were swapped with the k th pair in \dot{t} was that of index $\sigma_t^{-1}(k)$, and computing $\mathfrak{s}(\dot{t}, \dot{t}')$ is easy.

It is quite easy to check that $\mathfrak{s} \circ \mathfrak{J}$ verifies the eq. (3.2.4) on page 134.

The invariant dynamics is then

$$K(s_{1:n}, s'_{1:n}) = \sum_{t, t' \in \bar{\mathfrak{X}}^+, i, i' \in \dot{\bar{\mathfrak{X}}}} \bar{\Pi}(s_{1:n}, t) \mathfrak{J}(e(t), t) \mathfrak{s}(t, t') \mathbf{1}(\mathfrak{J}(e(t)', t') > 0) \zeta(t', s'_{1:n}) \times \left(\frac{\zeta(t, s_{1:n})}{\bar{\Pi}(s_{1:n}, t)} \frac{\mathfrak{J}(e(t)', t')}{\mathfrak{J}(e(t), t)} \frac{\mathfrak{s}(t, t')}{\mathfrak{s}(t', t)} \frac{\bar{\Pi}(s'_{1:n}, t')}{\zeta(t', s'_{1:n})} \wedge 1 \right).$$

(remember that $e(t)$ is the end tree of the sequence of trees t)

Remark that $\mathfrak{s}(t, t') > 0$ implies that t and t' have the same number of parentheses, so that $\frac{\mathfrak{J}(e(t)', t')}{\mathfrak{J}(e(t), t)} = 1$.

3.D Building general Metropolis reversible dynamics

The two reversible dynamics defined in section 2.2 on page 69 use a standard Metropolis algorithm, with a single intermediate step, while those in sections 3.2.4 and 3.2.5 on page 131 and on page 134 use multiple intermediate steps. We will here present a general method to define reversible dynamics, using as many intermediate steps as needed.

Let A_1, \dots, A_n be a sequence of sets, and consider accordingly a sequence of Markov kernels $\pi_i \in \mathcal{M}_1(A_i)^{A_{i-1}}$, $\pi'_i \in \mathcal{M}_1(A_{i-1})^{A_i}$ and $\zeta \in \mathcal{M}_1(A_n)^{A_n}$.

$$A_0 \longrightarrow A_1 \longrightarrow \dots \longrightarrow A_n \longrightarrow A_n \longrightarrow \dots \longrightarrow A_1 \longrightarrow A_0$$

$$s_0 \xrightarrow{\pi_1} s_1 \xrightarrow{\pi_2} \dots \xrightarrow{\pi_n} s_n \xrightarrow{\zeta} s'_n \xrightarrow{\pi'_n} \dots \xrightarrow{\pi'_2} s'_1 \xrightarrow{\pi'_1} s'_0$$

For any $S = (s_0, s_1, \dots, s_n) \in \prod_{i=0}^n A_i$, let $\check{S} = (s_n, \dots, s_1, s_0)$.

Consider a distribution \mathbb{P} on A_0 , and a function $q : \prod_{i=0}^n A_i \times \prod_{i=n}^1 A_i \rightarrow \mathbb{R}$ such that, for any $S, S' \in \prod_{i=0}^n A_i$

$$\prod_{i=0}^n \pi_i(S_{i-1}, S_i) \zeta(S_n, S'_n) \pi'_i(S'_i, S'_{i-1}) > 0 \implies \mathbb{P}(S_0) q(S, \check{S}') = \mathbb{P}(S'_0) q(S', \check{S}). \tag{3.D.1}$$

Proposition 3.D.1

Under eq. (3.D.1) on the preceding page, the Markov kernel $(K(s, s')_s, s' \in A_0)$

defined by

$$K(s, s') = \sum_{\substack{s_i, s'_i \in A_i \\ 0 \leq i \leq n}} \mathbb{1}(s = s_0, s' = s'_0) q(s_0, s_1, \dots, s_n, s'_n, \dots, s'_1, s'_0) \\ \prod_{i=0}^n \pi_i(s_{i-1}, s_i) \pi'_i(s'_i, s'_{i-1}) \left(\prod_{i=0}^n \frac{\pi'_i(s_i, s_{i-1})}{\pi_i(s_{i-1}, s_i)} \cdot \frac{\pi_i(s'_{i-1}, s'_i)}{\pi'_i(s'_i, s'_{i-1})} \wedge \mathbb{1} \right),$$

is reversible with respect to \mathbb{P} .

PROOF. Simply remark that

$$\mathbb{P}_S(s) K(s, s') = \sum_{\substack{s_i, s'_i \in A_i \\ 0 \leq i \leq n}} \mathbb{P}(s_0) q(s_{0:n}, s'_{n:0}) \mathbb{1}(s = s_0, s' = s'_0) \\ \times \left(\prod_{i=0}^n \pi'_i(s_i, s_{i-1}) \cdot \pi_i(s'_{i-1}, s'_i) \wedge \prod_{i=0}^n \pi_i(s_{i-1}, s_i) \cdot \pi'_i(s'_i, s'_{i-1}) \right)$$

is obviously symmetric in s, s' . □

Chapter 4

Conclusion

We would like here to comment on a few possible uses of the objects presented in this work. We will focus on language analysis based on a corpus of independent texts $S_{1:n}^i$, followed by some considerations on the scope of our model, and on possible extensions and open questions.

4.1 Outline of possible uses of Markov substitute sets in natural language processing

The central notion of this work is that of Markov substitute sets, which corresponds to the notion of syntactic constituents in linguistics. Chapter 2 proposes some ways of testing whether a certain set is a Markov substitute set.

The first step of our language analysis from a statistical sample is to build the reference grammar, or a Markov ruleset, which is roughly equivalent to identifying a collection of Markov substitute sets.

We begin by initializing a collection of Markov substitute sets

$$\mathcal{B}_0 = (\{w\}, w \in S),$$

and their substitute measures $q_{\{w\}} = \delta_w$, which would correspond to the trivial reference grammar

$$\mathcal{R}_0 = \bigoplus_i 1 \otimes [{}_i w_i.$$

We then use the sample to make the reference grammar grow, in the way described in section 3.4 on page 144. This construction is recursive, using at each step a new sample (to avoid dealing with dependency issues). Each iteration works as follows:

- Test as many rules as possible that we could add to the reference grammar, using for example the simultaneous test proposed in section 2.7 on page 96, and add those that passed the test,
- Saturate the resulting reference grammar,
- Estimate the substitute measures for each category of rules, as proposed in section 3.2.2 on page 127.

Each step will use the parsing kernels corresponding to the current reference grammar, whose constructions are described in section 3.B on page 159.

When we arrive at a point where no more Markov substitute sets can be identified, or when we do not have anymore samples, or when a missing mass estimator tells us that the grammar we built so far produces surface structures covering almost all of the support of the language distribution, we stop.

The resulting reference grammar can then give us some information on the syntactic structure of the language, and help us define various reversible dynamics on it.

We saw three main types of reversible dynamics. The first one, the most powerful, but with an unknown off-model behaviour, is the substitute dynamics k , using knowledge of the substitute measures to replace instances of them in the initial sample. The second one is the crossing-over dynamics K , that avoids the necessity of knowing the substitute measures at the price of keeping the same number of each word as in the initial sample. When applied to a number of replicas of the statistical sample growing to infinity, the crossing-over dynamics computes the maximum likelihood estimator in a linear exponential family, and therefore has a provable off-model robust behavior. The last dynamics, the split and merge process, is quite similar to the crossing-over dynamics, and uses toric grammars to deconstruct a text in its smaller syntactic components and rebuild it.

Knowledge of a Markov reference grammar gives us parsing tools, that permit us to analyse any string, and, ultimately, to estimate its probability in the context-free language of sentences of the same type B_e . Estimating the frequencies of each type of sentences gives us finally a estimator of the whole language.

The support of the estimated language,

$$\widehat{\mathcal{L}} = \bigcup_{e \in \widehat{\mathcal{E}}_g} B_e,$$

is potentially larger than a context-free language, as a union of a possibly infinite number of such languages. However, with a finite sample, we will never be

able to obtain more than a finite number of surface structures. Nonetheless, the dynamics k defines a slightly bigger support for the estimator of the language,

$$\widehat{\mathcal{L}} = \bigcup_{i=1}^n C(S_i),$$

which may be larger than a context-free language, though still recursively enumerable. The actual position of $C(s)$ in the hierarchy of formal languages is still an open question.

The split and merge process was also used in chapter 1 to define a communication model, which defines another notion of language, as the invariant measure on each class of the communication chain.

If the reference grammar is fixed, we saw that each text was recurrent, and thus that the language could be obtained back, without error, by a Monte-Carlo simulation, from any text on its communicating class.

However, we could ask ourselves what happens when the reference grammar is learnt anew at each generation. We saw that if the reference grammar is indeed Markov for the language, the split and merge process stays in its support (and is even at least weakly reversible). The communicating classes could nonetheless be fragmented if the reference grammar is too poor.

In our study, we supposed that our tests for the Markov substitute property were never wrong (at least, had a false acceptance probability close to zero). This is obviously optimistic, and it would be interesting to study more precisely the distortions introduced to the language when the reference grammar is not Markov.

Some other uses of the model can also be considered. An interesting possibility is to define sub-models, focused on a certain type of texts, or topics, by restricting the Markov substitute sets learnt in a large corpus to the expressions appearing in a more specialized one.

Estimating the substitute measures on a smaller corpus could also be an interesting tool, for example to obtain a statistical description of the style of a given author.

4.2 Further considerations on the scope of the model

The structure of the model presented in this work appears as an extension of the conditional independence assumptions defining Markov random fields.

We saw that a \mathcal{B} -Markov process can be parametrized as an exponential family, with some sort of Markov field structure on a graph on D^+ . Two vertices are neighbours if they differ only by two members of a Markov substitute set, the

ratio of the weights depending only on the non-common part of the two nodes:

$$\frac{\mathbb{P}_S(\alpha(x, y))}{\mathbb{P}_S(\alpha(x, y'))} = \frac{q(y)}{q(y')}.$$

In contrast to the description of the potential function of a Markov random field, depending on a decomposition on cliques, we proposed here a description in terms of minimal “free” edge types and anchor point on each connected component. The number of free edges reflects the number of free parameters in the choice of the substitute measures.

The edge structure considered here was defined using the Markov substitute property, however, we could also obtain linear exponential families using a different definition of the edges, that is of the probability ratios between pairs of states that are constrained to be identical.

We could in particular add restrictive conditions on the context. We could also constrain the probability ratio of states that are linked by a relation that is not a substitution, but that is based on a different type of transformation, for instance a movement of the form $\alpha(\alpha(x, \epsilon), y) \sim \alpha(\alpha(x, y), \epsilon)$.

We proved that the Markov substitute model is the thermodynamic limit of a crossing-over dynamics. This crossing-over dynamics can be interpreted as a communication process, which can be used to transmit a language distribution from speaker to speaker in an error free way, using a limited amount of memory, compared to the size of the support of the language.

This crossing-over dynamics is reversible with respect to the uniform measure on memorized sentences. The fact that the language can be maintained from one speaker to another by simple knowledge of the grammar, that is, the collection of Markov substitute sets (without their substitute measures), hints that language can be “stored” quite efficiently. This could give some insight for the “poverty of stimulus” problem put forward by linguists. In other words, you do not need to learn the parameters (the substitute measures), you need only to learn the model (the set of substitute sets), and then, from a relatively small amount of memorized sentences, you can generate a huge number of new sentences using random crossing-overs, while outputting sentences with a probability distribution close to the model distribution using the optimal model parameters (in terms of likelihood based on what you have heard from others).

To come back to the two steps of the estimation of the language:

- learning the model (i.e. a generating set of Markov substitute sets),
- learning the parameters (the substitute measures),

the second step is simply solved by the crossing-over dynamics.

Various scenarios can be imagined for learning the first step, ranging from theoretical answers based on simultaneous tests, as described in this thesis, or on penalized likelihood criteria available in a more generic way for collections of exponential families, to more supervised learning scenarios, where a teacher, supervising the output of a crossing-over dynamics, raises a red flag each time a wrong sentence is produced (allowing the pupil to cancel the Markov substitute pair involved in the crossing-over from its list of Markov substitute sets). However, children seem to already know the language when they first go to school. Studies in psycholinguistics, especially in children, hints that language learning is mostly based on positive information, and negative feedback is, in the first stages at least, virtually non-existent. Nonetheless, some sort of supervised school-like training could be a way to refine (and homogenize) the model selection step of language learning.

Bibliography

- [Bak79] James K Baker. Trainable grammars for speech recognition. *The Journal of the Acoustical Society of America*, 65(S1):S132–S132, 1979.
- [Bes74] Julian Besag. Spatial interaction and the statistical analysis of lattice systems. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 192–236, 1974.
- [CG98] Zhiyi Chi and Stuart Geman. Estimation of probabilistic context-free grammars. *Computational linguistics*, 24(2):299–305, 1998.
- [Chi99] Zhiyi Chi. Statistical properties of probabilistic context-free grammars. *Computational Linguistics*, 25(1):131–160, 1999.
- [Cho56] Noam Chomsky. Three models for the description of language. *Information Theory, IRE Transactions on*, 2(3):113–124, 1956.
- [Cho65] Noam Chomsky. *Aspects of the Theory of Syntax*. Number 11. MIT press, 1965.
- [Cho95] Noam Chomsky. *The minimalist program*, volume 28. Cambridge Univ Press, 1995.
- [Cla14] Alexander Clark. Learning trees from strings: A strong learning algorithm for some context-free grammars. *Journal of Machine Learning Research*, 14:3537–3559, 2014.
- [CS12] Shay B. Cohen and Noah A. Smith. Empirical risk minimization for probabilistic grammars: Sample complexity and hardness of learning. *Computational Linguistics*, 38(3):479–526, 2012.
- [DPDPG⁺94] Stephen Della Pietra, Vincent Della Pietra, John Gillett, John Lafferty, Harry Printz, and Luboš Ureš. *Inference and estimation of a long-range trigram model*. Springer, 1994.

- [Fel68] William Feller. *An Introduction to Probability Theory and Its Applications: Volume One*. John Wiley & Sons, third edition, 1968.
- [GK07] Thomas L. Griffiths and Michael L. Kalish. Language evolution by iterated learning with bayesian agents. *Cognitive Science*, 31(3):441–480, 2007.
- [Goo53] Irving J. Good. The population frequencies of species and the estimation of population parameters. *Biometrika*, 40(3-4):237–264, 1953.
- [LL09] Martin Lange and Hans Leiß. To CNF or not to CNF? an efficient yet presentable version of the CYK algorithm. *Informatica Didactica*, 8:2008–2010, 2009.
- [LY90] Karim Lari and Steve J. Young. The estimation of stochastic context-free grammars using the inside-outside algorithm. *Computer speech & language*, 4(1):35–56, 1990.
- [McA98] David A. McAllester. Some pac-bayesian theorems. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 230–234. ACM, 1998.
- [McA99] David A. McAllester. Pac-bayesian model averaging. In *Proceedings of the twelfth annual conference on Computational learning theory*, pages 164–170. ACM, 1999.
- [McA03] David A. McAllester. Pac-bayesian stochastic model selection. *Machine Learning*, 51(1):5–21, 2003.
- [MO03] David McAllester and Luis Ortiz. Concentration inequalities for the missing mass and for histogram rule error. *The Journal of Machine Learning Research*, 4:895–911, 2003.
- [Nor98] James R. Norris. *Markov chains*. Cambridge university press, 1998.
- [Roa01] Brian Roark. Probabilistic top-down parsing and language modeling. *Computational linguistics*, 27(2):249–276, 2001.
- [Sta09] Edward P. Stabler. Mathematics of language learning. *Histoire, Épistémologie, Langage*, 31(1):127–145, 2009.
- [Tom85] Masaru Tomita. *Efficient parsing for natural language: a fast algorithm for practical systems*. Kluwer Academic Publishers, 1985.

- [Tom87] Masaru Tomita. An efficient augmented-context-free parsing algorithm. *Computational linguistics*, 13(1-2):31–46, 1987.
- [TZZW12] Ming Tan, Wenli Zhou, Lei Zheng, and Shaojun Wang. A scalable distributed syntactic, semantic, and lexical language model. *Computational Linguistics*, 38(3):631–671, 2012.
- [You67] Daniel H. Younger. Recognition and parsing of context-free languages in time n^3 . *Information and control*, 10(2):189–208, 1967.