



**HAL**  
open science

## Surveillance dynamique de mesures

Sophie Miallaret

► **To cite this version:**

Sophie Miallaret. Surveillance dynamique de mesures. Statistiques [math.ST]. Université Clermont Auvergne [2017-2020], 2019. Français. NNT : 2019CLFAC091 . tel-02866517

**HAL Id: tel-02866517**

**<https://theses.hal.science/tel-02866517>**

Submitted on 12 Jun 2020

**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.



# Doctorat Université Clermont Auvergne

École doctorale : Sciences Fondamentales

## THÈSE

pour obtenir le grade de docteur délivré par

## l'Université Clermont Auvergne

Spécialité doctorale “Mathématiques appliquées -  
statistiques et applications”

*présentée et soutenue publiquement par*

**Sophie Mialaret**

le 11 Décembre 2019

## **Dynamic Monitoring Measures**

Directeurs de thèse : **Arnaud Guillin et Anne-Françoise Yao**

### **Jury**

<b>Sophie Dabo-Niang,</b>	Professeur, Université de Lille	Rapporteur
<b>Jean-Michel Marin,</b>	Professeur, Université de Montpellier	Rapporteur
<b>Hacène Djellout,</b>	Maître de conférence, Université Clermont Auvergne	Examineur
<b>Denys Pommeret,</b>	Professeur, Aix-Marseille Université	Examineur
<b>Laurence Reboul,</b>	Maître de conférence, Aix-Marseille Université	Examineur
<b>Vincent Sapin,</b>	Professeur, Université Clermont Auvergne	Examineur

**Université Clermont Auvergne**

**Laboratoire Mathématiques Blaise Pascal (LMBP)**

UMR CNRS 6620, Clermont Ferrand, France



# Remerciements

Cette thèse met fin à plusieurs années d'études. Le manuscrit que vous avez sous vos yeux ne pourrait pas être là sans l'aide de différentes personnes que je souhaite remercier.

Mes premiers remerciements s'adressent à mes directeurs de thèse, Anne-Françoise et Arnaud pour m'avoir accompagné lors de toutes les étapes de ma thèse. Je tiens à leur exprimer ma profonde gratitude pour leur présence, leur disponibilité et leurs conseils durant ces trois années.

Je remercie Sophie Dabo-Niang et Jean-Michel Marin d'avoir accepté de rapporter ma thèse, ainsi que Hacène Djellout, Denys Pommeret, Laurence Reboul et Vincent Sapin d'avoir accepté d'être examinateur de mes travaux. Merci de votre intérêt pour mon travail.

Je remercie la région Auvergne Rhône Alpe et le Fond Européen de Développement Régional pour le financement accordé qui a permis la réalisation de cette thèse.

Je tiens à remercier Genbio, Deltamu et le CHU d'avoir rendu ce projet possible et de m'avoir fait confiance pour le faire avancer.

Je remercie Jean-Michel Pou pour son enthousiasme et son accompagnement tout au long du projet. Je remercie également Christophe Dubois et Peggy Courtois d'avoir suivi ce travail et pour les conseils qu'ils ont pu me donner.

Je remercie Vincent Sapin de nous avoir donné son avis et ses conseils tout au long du projet. Merci à Régine Quinard et Lise Landrieaux d'avoir pris du temps pour répondre à mes questions et pour le partage d'informations.

Je remercie Julia d'avoir bien voulu se plonger dans nos analyses statistiques et de réussir à traduire nos résultats en termes géologiques.

Cette thèse est la fin de neuf ans d'étude au LMBP et je tiens à remercier tous les enseignants qui m'ont formé pendant cette période et qui m'ont permis d'en arriver là. Je remercie tous les membres du LMBP pour l'environnement chaleureux qu'ils donnent au laboratoire notamment les membres du secrétariat pour leur aide et leur gentillesse.

Je tiens à remercier les doctorants et post-doctorants que j'ai pu croiser durant ces trois

ans de thèse et qui ont su rendre les journées au laboratoire conviviales. Merci à vous pour les pauses partagées mais aussi, les repas, les soirées, les jeux, les verres et les rires. En tout premier lieu à Arnaud, Chaoen et Rénauld avec qui j'ai partagé le bureau 1213 durant trois ans, mais aussi à Franck, Arthur, Valentin, Fernando, Athina, Damien, Minh, Charlie, Maeva, Thomas, Pape, Sébastien, Baptiste et Emeline.

Merci à mes amis avec qui je partage de merveilleux moments, que ce soit ma fin d'étude, les jours de l'an, les karaokés improvisés ou les fameuses « fêtes à Lanobre ». Un merci tout particulier à Flora, Jessica et Sophie, pour avoir traversé l'adolescence et le passage à la vie d'adulte avec moi. Merci à vous d'avoir été là dans les bons et les mauvais moments de ces 15 dernières années.

Un tendre merci à ma famille pour le soutien et l'amour qu'ils m'ont apporté durant toutes ces années. Un grand merci à mes parents pour leurs encouragements et le soutien indéfectible dans les choix que j'ai pu faire. Merci à ma sœur pour sa folie à toute épreuve qui la rend aussi fantastique. Merci à vous d'être présents pour moi quoiqu'il arrive, je n'en serais absolument pas là sans vous.

Je tiens à remercier avec tout mon amour Carole, pour avoir supporté mes sautes d'humeur et mes grincements de dents mais surtout pour l'écoute, le réconfort, la confiance et les rires que tu m'apportes au quotidien. Je n'aurais pas rêvé meilleur soutien que toi durant ces trois années de thèse.

# Résumé

Les mesures sont des actes quotidiens, elles nous donnent beaucoup d'informations et permettent de prendre des décisions. L'analyse des mesures peut nous permettre d'en apprendre plus sur notre environnement, mais l'erreur d'une mesure peut avoir des conséquences importantes dans certains domaines.

Dans une première partie, nous proposons, grâce à l'étude de mesures d'analyses sanguines réalisées au CHU de Clermont-Ferrand, une procédure permettant de détecter les dérives des analyseurs de laboratoires de biologie médicale, se basant sur les mesures d'analyses de patients. Après une analyse descriptive des données, la méthode mise en place, utilisant des méthodes de détection de ruptures de séries temporelles, est testée pour des simulations de ruptures représentant des décalages, des imprécisions ou des dérives d'analyseurs pour différents paramètres biologiques mesurés. La méthode est adaptée pour deux scénarios : lorsque l'on connaît ou non le service hospitalier des patients. L'étude est complétée par une analyse de l'impact de l'incertitude de mesure sur les analyses des patients.

Dans une seconde partie nous étudions des mesures de formes de cendres volcaniques réalisées au Laboratoire Magmas et Volcans de l'Université Clermont Auvergne, dans le but de déterminer un lien entre les lieux de collecte et les formes des particules. Après avoir montré la dépendance entre ces paramètres, nous proposons, grâce une méthode de classification, un regroupement des particules représentant différentes populations dépendantes de la distance entre les lieux de collecte et le cratère du volcan.



# Abstract

The measures are daily actions, they give us a lot of information and allow us to make decisions. The analysis of measures can allow us to learn more about our environment, but the error of a measure can have important consequences in certain areas.

In a first part, we propose, thanks to the study of blood test measurements carried out at the CHU of Clermont-Ferrand, a procedure for detecting deviations from medical biology laboratory analyzers based on patient analysis measurements. After a descriptive analysis of the data, the method put in place, using methods of detection of breaks of time series, is tested for simulations of breaks representing offsets, imprecision or drifts of machine for different measured biological parameters. The method is adapted for two scenarios: when the patient's hospital service is known or not. The study is supplemented by an analysis of the impact of measurement uncertainty on patient analyses.

In a second part we study measurements of volcanic ash forms made at "Laboratoire Magmas et Volcans" of the Clermont Auvergne University, in order to determine a link between the collection locations and the forms of the particles. After showing the dependence between these parameters, we propose, using a classification method, a grouping of particles representing different populations depending on the distance between the collection locations and the volcano crater.





# Table des matières

<b>Remerciements</b>	<b>iii</b>
<b>Résumé</b>	<b>v</b>
<b>Abstract</b>	<b>vii</b>
<b>Table des matières</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Contexte général et enjeux . . . . .	1
1.1.1 Laboratoires de biologie médicale . . . . .	1
1.1.2 Volcanologie . . . . .	2
1.2 Organisation du manuscrit et contributions . . . . .	2
1.2.1 Chapitre 2 : Étude de mesures sanguines réalisées au CHU de Clermont-Ferrand . . . . .	2
1.2.2 Chapitre 3 : Étude de mesures de cendres volcaniques réalisées au Laboratoire Magmas et Volcans de l'Université Clermont Auvergne . . . . .	10
<b>2 Étude de mesures sanguines réalisées au CHU de Clermont Ferrand</b>	<b>17</b>
2.1 Introduction . . . . .	17
2.1.1 Contrôle interne de qualité . . . . .	17
2.1.2 CHU Clermont-Ferrand . . . . .	19
2.2 Données . . . . .	19
2.2.1 Jeux de données . . . . .	19
2.2.2 Paramètres biologiques . . . . .	20
2.2.3 Différences pôle Estaing et Gabriel-Montpied . . . . .	23
2.2.4 Différences par service . . . . .	25
2.3 Étude des paramètres biologiques . . . . .	26
2.3.1 Distributions des paramètres . . . . .	26
2.3.2 Effets des services sur les distributions des paramètres biologiques . . . . .	28
2.3.3 Corrélations des paramètres . . . . .	30
2.3.4 Classification . . . . .	30
2.4 Détection de ruptures dans des séries temporelles . . . . .	36
2.4.1 Méthodes . . . . .	36
2.4.2 Recherche de ruptures . . . . .	39
2.5 Étude des problèmes observés par le CHU . . . . .	45
2.5.1 Problèmes repérés lors des CIQ . . . . .	45
2.5.2 Détections de ruptures et interprétations . . . . .	46

2.6	Simulations et détections de ruptures . . . . .	48
2.6.1	Méthode de détection en ligne . . . . .	49
2.6.2	Simulation des données . . . . .	50
2.6.3	Création de ruptures . . . . .	52
2.6.4	Méthode de détection de ruptures avec l'information service . . . . .	54
2.6.5	Méthodes de détection de ruptures sans l'information service . . . . .	59
2.7	Incertitudes de mesure . . . . .	66
2.7.1	Incertitudes de mesure (IM) au laboratoire . . . . .	66
2.7.2	Distribution des valeurs vraies . . . . .	67
2.7.3	Impact de l'incertitude de mesure sur les valeurs vraies des mesures de patients . . . . .	69
2.8	Conclusion . . . . .	72
2.8.1	Résumé de l'étude . . . . .	72
2.8.2	Perspectives . . . . .	73
<b>3</b>	<b>A characterization of volcanic ash grain particles based on the shape parameters. Example of observations collected around the crater of the Tungurahua volcano in Ecuador.</b>	<b>75</b>
3.1	Introduction . . . . .	75
3.1.1	Geographical setting . . . . .	75
3.1.2	Data collection . . . . .	76
3.2	Statistical data analysis . . . . .	77
3.3	Results . . . . .	80
3.3.1	Distribution of the number of observations . . . . .	80
3.3.2	Basic statistics on the ash grain shape parameters. . . . .	81
3.3.3	The distribution of the shape parameters and Box-Cox transformation . . . . .	82
3.3.4	Relationship between the grain shape parameters and the distance to the crater. . . . .	84
3.3.5	Classification of the ash grain particles using their shape characteristic and distance to the crater . . . . .	87
3.3.6	Statistical modeling . . . . .	94
3.4	Conclusion . . . . .	99
<b>A</b>	<b>Annexes Chapitre 2</b>	<b>101</b>
A.1	Quantiles et moyennes des paramètres biologiques par sexe et catégories d'âge des patients. . . . .	101
A.2	Représentations graphiques des moyennes des paramètres biologiques par âge et sexe des patients calculées sur le jeu de données 2018 . . . . .	103
A.3	Résumés et graphiques des distributions des paramètres biologiques pour les cinq services ayant fait le plus de mesures pour chaque paramètre en 2018 aux pôles Gabriel Montpied et Estaing . . . . .	107
A.4	Résultats pour les paramètres chlore, créatinine, ferritine, protéines, PSA et urée ; méthode avec l'information service. . . . .	111
A.5	Résultats pour les paramètres chlore, créatinine, ferritine, protéines, PSA et urée ; méthode sans l'information service et avec des classifications. . . . .	129
A.6	Étude de l'impact de l'incertitude de mesure sur les valeurs vraies estimées des paramètres biologiques . . . . .	147





# Chapitre 1

## Introduction

### 1.1 Contexte général et enjeux

Des mesures sont effectuées dans tous les domaines : l'industrie, la recherche, la santé ou même dans nos vies quotidiennes. Leurs résultats servent à prendre des décisions : acceptation d'un produit, réglage d'un instrument de mesure, protection de l'environnement, décision médicale... Chaque mesure apporte donc beaucoup d'informations, c'est pourquoi l'exploration des résultats des mesures est importante, en recherche elles peuvent aider à mieux comprendre notre environnement ou son fonctionnement.

L'erreur d'une mesure peut avoir différentes conséquences : économique, risque d'accident,... Certaines conséquences peuvent être très importantes, comme en santé. En effet une erreur de mesure en médecine peut engendrer de mauvaises actions médicales pouvant impacter la santé d'un patient. Le suivi et le contrôle des appareils de mesure sont donc essentiels. La science de la mesure dit métrologie donne de nombreux outils pour cela.

#### 1.1.1 Laboratoires de biologie médicale

Les laboratoires de biologie médicale peuvent réaliser plusieurs milliers de mesures par jour pour différents types d'analyses. Le suivi, le contrôle et la qualité des mesures réalisées sont réglementés par des normes, notamment la norme NF EN ISO 15189 [2] qui est imposée par la loi depuis 2010 [3].

Pour répondre aux exigences de cette norme, les laboratoires ont mis en place différents outils permettant de suivre et d'analyser la qualité des mesures : les contrôles internes de qualité (CIQ), les comparaisons inter laboratoires (CIL), les évaluations externes de qualité (EEQ) ou les CIQ externalisés. Ces dispositifs sont mis en place pour chaque type d'analyses réalisées dans les laboratoires pour prévenir les possibles dérives des analyseurs. Ils sont tous effectués sur des échantillons indépendants des analyses effectuées pour les patients et sont effectués en parallèle de celles-ci à périodicité arbitraire. Généralement, les CIQ sont réalisés

plusieurs fois par jour, pour suivre quotidiennement la qualité des résultats, alors que les autres contrôles sont effectués chaque mois.

L'utilisation d'outils statistiques ou prédictifs sur les mesures des patients permettrait de suivre en temps réel la qualité des mesures et d'ainsi prévenir au plus vite la dérive d'un analyseur et de faire moins de contrôles inutiles.

C'est l'objectif que nous avons en étudiant des jeux de données du Centre Hospitalier Universitaire (CHU) de Clermont-Ferrand

### **1.1.2 Volcanologie**

En volcanologie, l'étude des dépôts de particules de cendres donne beaucoup d'informations sur le transport des cendres, les mécanismes de fragmentation du magma, etc... Plus précisément la morphologie des particules de cendres permet d'étudier les conditions de fragmentation, de transport et de sédimentation des cendres. Les grains de cendres sont mesurés afin d'évaluer la géométrie et la morphologie des particules. Les mesures sont réalisées grâce à l'acquisition d'images des grains.

Le laboratoire Magmas et Volcans de l'université Clermont-Auvergne a pu utiliser l'instrument Morphologie G3 développé par Malvern, permettant de mesurer automatiquement les caractéristiques morphologiques et géométriques en deux dimensions d'une grande quantité de particules de cendres en prenant des clichés des grains pour différents grossissements. [15] [8] Ce sont les résultats de ces mesures que nous étudions dans le but de comprendre la répartition spatiales des particules en fonction de leur forme.

## **1.2 Organisation du manuscrit et contributions**

### **1.2.1 Chapitre 2 : Étude de mesures sanguines réalisées au CHU de Clermont-Ferrand**

Ce premier chapitre traite les mesures de sept paramètres biologiques réalisées au CHU de Clermont-Ferrand. Les méthodes présentées dans ce chapitre vont permettre de suivre et contrôler les mesures des patients en temps réel et ainsi détecter les dérives des analyseurs le plus rapidement possible.

#### **Contrôles internes de qualité**

Aujourd'hui pour suivre et contrôler les analyses quotidiennement, des contrôles internes de qualité (CIQ) sont réalisés plusieurs fois par jour pour chaque analyse effectuée et sur chaque analyseur. Ces contrôles consistent à mesurer plusieurs fois par jour des étalons dont on connaît les valeurs avec la méthode d'analyse contrôlée. Les résultats de ces contrôles sont étudiés et interprétés notamment grâce aux graphiques de Levey Jenngins et aux règles

d’alarme et de rejet de Westgard permettant d’accepter ou de refuser un contrôle effectué. [32]. Si les contrôles sont rejetés, les analyses des patients effectuées avant ce contrôle peuvent être repassées et comparées aux anciennes mesures. Il existe deux seuils pour décider si l’écart entre ces deux mesures est acceptable ou non. L’un est le Delta-Check, calculé à l’aide de la moyenne et de l’écart type des CIQ, l’autre est le TCC, calculé à l’aide du coefficient de variation intra-individuelle du paramètre biologique mesuré et le coefficient de variation calculé avec les CIQ.

## Jeux de données

Pour notre étude, nous avons pu étudier les mesures de sept paramètres biologiques : albumine, chlore, créatinine, ferritine, protéines, PSA et urée. Ces paramètres sont accompagnés de l’âge et du sexe des patients, la date et l’heure de la mesure ainsi que l’analyseur sur lequel elle a été effectuée.

Le principal jeu de données étudié donne les mesures de ces paramètres effectuées en 2018 au CHU de Clermont-Ferrand. il donne aussi le service hospitalier de provenance des patients.

Les mesures étudiées ont été réalisées sur 5 analyseurs différents : VISTA 1, 2, 3, 500 et 1500. VISTA 1, 2 et 3 se trouvent sur le site Gabriel-Montpied et VISTA 500 et 1500 sur le site Estaing.

Les paramètres biologiques sont accompagnés par leurs valeurs de référence. Ces valeurs permettent d’interpréter les résultats d’analyses des patients et dépendent de l’âge et du sexe des personnes. Pour chaque patient, nous pouvons ainsi dire si les taux des paramètres biologiques entrent ou non dans les valeurs de référence, si c’est le cas les valeurs seront dites usuelles.

Nous savons que certains paramètres dépendent fortement du sexe et de l’âge de la personne. L’annexe A.1 décrit les quantiles et les moyennes des paramètres biologiques en fonction de l’âge et du sexe et A.2 représente ces moyennes graphiquement.

Les données proviennent de deux sites différents du CHU : Estaing et Gabriel Montpied. Ces pôles regroupent des services hospitaliers différents et n’ont donc pas le même panel de patients.

325 services sont présents dans le jeu de données. Certains sont rares avec un unique patient et d’autres sont très présents avec au maximum 32995 patients (Accueil urgence de Gabriel Montpied). Chaque service regroupe des panels de patients différents, les proportions d’hommes et de femmes peuvent varier ainsi que la moyenne d’âge. Ces différences de panels engendrent des distributions très distinctes par service pour certains paramètres, comme l’albumine (figure 1.1) .



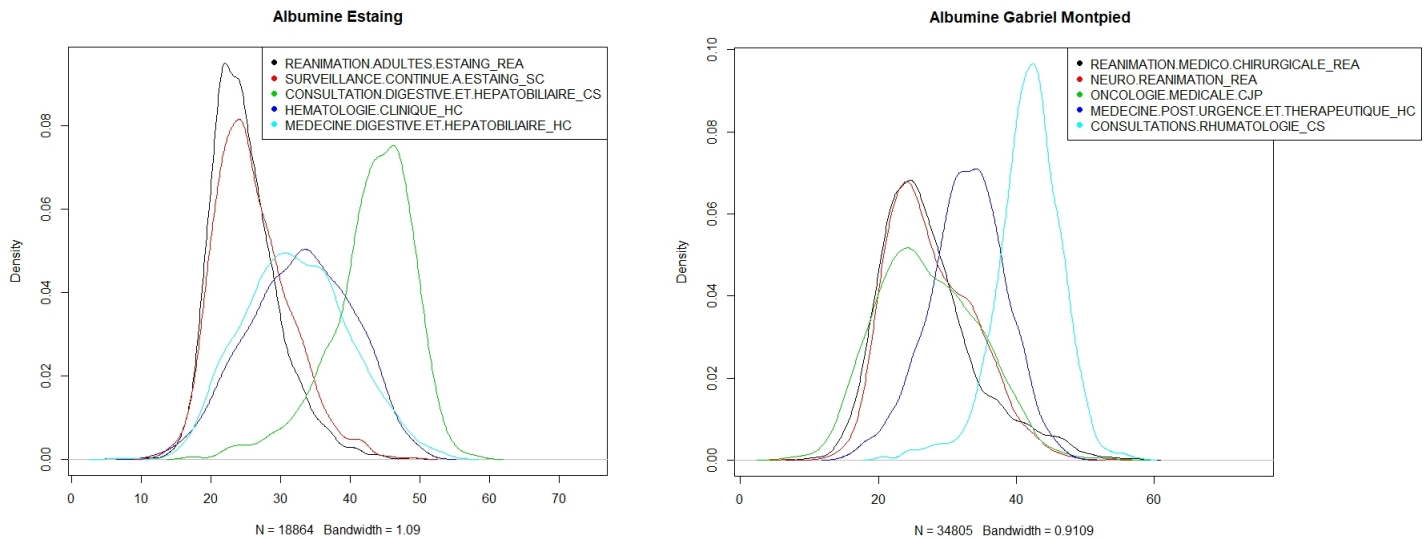


FIGURE 1.1 – Distributions de l’Albumine pour les cinq services ayant fait le plus de mesures d’Albumine en 2018 aux pôles Estaing et Gabriel Montpied

### Problèmes observés sur les analyseurs lors des CIQ

Le CHU a fourni une liste de dates correspondant aux jours où un problème a été observé sur un analyseur lors des CIQ. Ces dates sont données avec le paramètre biologique et l’analyseur correspondants. Pour les données de 2018, vingt dates correspondent à des paramètres biologiques étudiées lors de notre étude.

Nous avons étudié les mesures réalisées lors des jours ayant eu un problème, pour voir si ceux-ci avaient eu un impact sur les données des patients. Pour cela, nous avons calculé la variance, la moyenne et le pourcentage de mesures en dehors des valeurs de référence du paramètre biologique pour chaque jour. Malheureusement, cela n’a pas permis de mettre en évidence les jours problématiques.

Les distributions des paramètres biologiques de ces jours ont été comparées à une distribution de référence estimée sur le reste du jeu de données. Les comparaisons sont réalisées grâce à un test de Wilcoxon. Deux distributions sont différentes de leurs références, celle des protéines du 12/02/2018 de l’analyseur VISTA500 et celle de l’urée du 31/07/2018 de l’analyseur VISTA3.

Ayant la date et l’heure de chaque mesure, les échantillons peuvent être analysés en tant que séries temporelles. Les séries analysées correspondent à plusieurs jours de mesures autour des dates problématiques.

Des méthodes de détection de ruptures dites hors ligne ont été appliquées sur ces séries. Différentes méthodes ont été utilisées, permettant de détecter des changements de moyenne,

de variance ou des pentes. Les mesures arrivant service par service, les détections de ruptures peuvent localiser ces changements de services. Pour éviter cela, chaque individu est centré et réduit avec la moyenne et l'écart type correspondant à son service.

Les détections ainsi réalisées ont localisé des ruptures sur un seul des échantillons étudiés : les mesures de protéines du 10/02/2018 au 14/02/2018, le problème étant repéré par les CIQ le 12/02. Les changements trouvés sont : le 12/02 à 06h31, 09h55 et le 13/02 à 12h14, elles sont représentées sur la figure 1.2

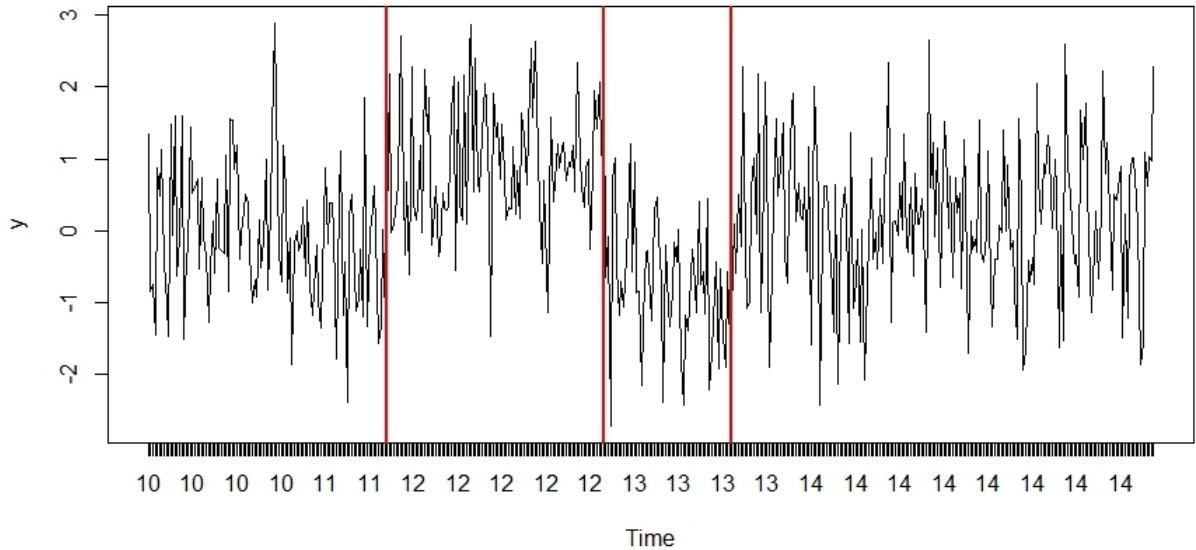


FIGURE 1.2 – Ruptures détections sur la série temporelle : protéines, centrée réduite selon les services, VISTA500, du 10/02/2018 au 14/02/2018

### Simulations et créations de ruptures

Malheureusement nous n'avons pas pu identifier les conséquences que peuvent avoir les problèmes provenant des analyseurs sur les mesures des patients. Nous avons donc émis l'hypothèse que lorsqu'un analyseur dysfonctionnerait cela engendrerait un bruit gaussien ajouté aux mesures des patients avec une certaine moyenne ou variance, ou même une moyenne qui pourrait évoluer en fonction du temps.

Afin de tester les différentes méthodes utilisées nous avons simulé des données ayant les mêmes distributions que les données de patients et ajouté une rupture pour tenter de la détecter le plus tôt possible. Les simulations de données prennent en compte l'arrivée des échantillons service par service au laboratoire en utilisant la matrice de transitions des services calculée sur les données patients.

Trois types de bruit sont testés représentant respectivement un décalage, une imprécision ou une dérive de la machine :  $\mathcal{N}(\mu, 1)$ ,  $\mathcal{N}(0, \sigma)$  ou  $\mathcal{N}(0.5T, 1)$  où  $T$  représente le temps de la

rupture, donc  $T = 1, \dots, 50$ . Les bruits sont ajoutés aux séries temporelles au temps  $i = 70$ . Pour que les ruptures simulées soient représentatives des écarts considérés comme de réels impacts biologiques,  $\mu$  et  $\sigma$  prendront d'abord les valeurs les plus basses de Delta-Check et TCC correspondant au paramètre biologique étudié puis varieront pour déterminer à partir duquel les détections localisent au mieux la rupture.

Dans cette partie des méthodes de détection hors ligne sont utilisées, mais aussi des méthodes en ligne. Ces dernières permettent de détecter en temps réel des changements de moyenne, de variance ou des pentes, en donnant le moment de la rupture estimée ainsi que le temps de détection qu'il a fallu pour la trouver.

### Méthode avec l'information service

Dans un premier temps nous émettons l'hypothèse que l'on connaît le service d'où proviennent les patients. Avant de réaliser les détections de ruptures, les données sont centrées et réduites en fonction du service, donc avec les moyennes et les écarts-types de chacun. Cela permet de perdre l'effet des services et de ne plus détecter leurs changements.

Les détections de ruptures hors ligne sont utilisées sur l'ensemble de l'échantillon. Les méthodes sont : segmentation binaire en moyenne et variance (BinMoy et BinVar), segment neighbourhood en moyenne et variance (SegMoy et SegVar), PELT, Pettitt (Pett) et breakpoint (Br).

Les détections de ruptures en ligne analysent l'échantillon valeur par valeur, les 20 premières mesures sont considérées comme l'échantillon de référence, le premier test est donc réalisé à la 20ème mesure. Les tests statistiques utilisés sont : Student (Stu), Bartlett (Bar), rapport de vraisemblance généralisé (GLR), Mann-Whitney (MW), Mood (Moo), Lepage (LP), Kolmogorov-Smirnov (KS) et Cramer-Von-Mises (CVM).

Pour chaque rupture testée, 100 simulations et détections sont réalisées. Pour les méthodes hors ligne, on regarde la première rupture détectée et la dernière. Pour les méthodes en ligne, on regarde la rupture détectée et le temps qu'il a fallu pour la trouver.

Suite aux 100 simulations, la rupture débutant à  $i = 70$ , nous regardons le pourcentage de simulations ayant trouvées la première rupture entre  $i = 60$  et  $i = 80$ , puis entre 60 et 90 ainsi que 60 et 100. On calcule aussi le pourcentage de simulations ayant trouvées la dernière rupture entre 110 et 130 pour les méthodes hors ligne et le temps de détection moyen pour les méthodes en ligne. Ces informations vont nous permettre de voir quelle méthode détecte le mieux la rupture et quelle est la plus rapide.

La figure 1.3 donne les résultats de détection d'un bruit  $\mathcal{N}(\mu; 1)$  en représentant les pourcentages de simulations d'albumine où les méthodes en ligne ont trouvé la première rupture entre  $i = 60$  et  $i = 80$ , ainsi que le temps de détection moyen, pour  $\mu$  variant de 0 à 15.

Pour chaque  $\mu$  différent, le pourcentage et le temps de détection moyen ont été calculés sur 50 simulations.

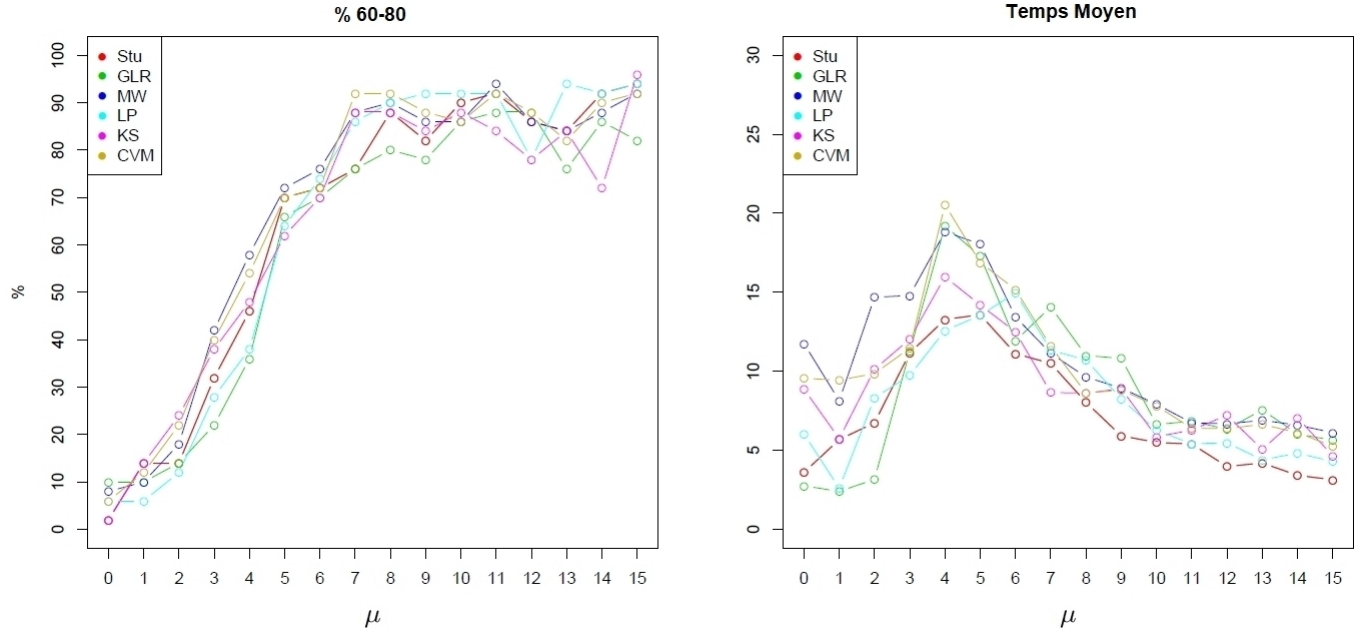


FIGURE 1.3 – Pourcentages de simulations d’Albumine avec ajout d’un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 0 à 15 ; Méthode avec l’information service

Nous nous intéressons maintenant aux fausses alarmes des méthodes en ligne. Si une méthode détecte bien les ruptures, elle doit aussi détecter le moins possible de fausses ruptures et donc de fausses alarmes en pratique dans les laboratoires.

Afin d’évaluer les détections de fausses ruptures, nous avons simulé, sans ajout de rupture, 200 échantillons d’une taille qui correspond aux nombres de mesures réalisées par paramètre biologique sur une journée. Les détections de ruptures en ligne ont été appliquées sur ces échantillons mais sans interruption, à chaque rupture détectée, la méthode se réinitialise à la mesure suivante. Cela permet d’évaluer le nombre de fausses alarmes que peuvent signaler les méthodes sur une journée de mesures au laboratoire.

Le tableau 2.14 donne les pourcentages de simulations ayant détectées 0, 1, 2 ou 3 ruptures pour 200 simulations d’albumine de 60 individus. On peut donc dire que dans environ 90% des cas, il n’y aurait pas de fausses alarmes lors d’une journée de mesures d’albumine en laboratoire et que s’il y en avait, elles seraient de 1, 2 ou 3 par jour.

Méthode	% 0	% 1	% 2	% 3
Stu	88,5	9,5	2	0
Bar	89	4	7	0
GLR	89,5	3	7,5	0
MW	90,5	7,5	2	0
Moo	92	7	1	0
LP	91,5	7,5	1	0
KS	91	8	0,5	0,5
CVM	91	7	2	0

TABLE 1.1 – Pourcentages de simulations où les méthodes de détection en ligne (avec services) ont trouvé 0, 1, 2 ou 3 ruptures sur 200 simulations d’albumine de 60 individus

### Méthode sans l’information service

Dans un second temps nous émettons l’hypothèse que nous n’avons pas l’information service, pour pallier à cela des classifications seront couplées aux détections de ruptures.

Avant de réaliser les détections de ruptures hors ligne, on effectue une classification de l’échantillon de données. Chaque groupe ainsi créé donne une nouvelle série temporelle et ce sont sur ces séries que l’on applique les détections de ruptures.

Pour les détections en ligne, on réalise une première classification sur l’échantillon référent qui correspondra aux 20 premiers individus. Des détections de ruptures en ligne seront effectuées sur chaque groupe représenté sous la forme d’une série temporelle. Si une rupture est trouvée, alors l’algorithme est arrêté et l’on peut donner le moment de rupture et le temps de détection. Sinon, on remet à jour la classification avec les deux individus suivant en plus, on relance les détections de ruptures et ainsi de suite. L’algorithme s’arrête lorsqu’une rupture est détectée sur une série temporelle ou lorsque l’on arrive au bout de l’échantillon.

Deux méthodes de classification sont utilisées. La première est une méthode probabiliste basée sur un modèle de mélange gaussien, elle prend en compte uniquement l’échantillon des mesures du paramètre biologique. La deuxième méthode est une classification ascendante hiérarchique (CAH) réalisée sur les mesures du paramètre biologique, le sexe et l’âge des patients.

Les simulations et les résultats des détections de ruptures sont étudiés de la même façon que pour la méthode utilisant l’information service.

La figure 1.4 donne les résultats de détection d’un bruit  $\mathcal{N}(\mu; 1)$  en représentant les pourcentages de simulations d’albumine où les méthodes en ligne ont trouvé la première rupture entre  $i = 60$  et  $i = 80$ , ainsi que le temps de détection moyen, pour  $\mu$  variant de 0 à 15. Pour chaque  $\mu$  différent, le pourcentage et le temps de détection moyen ont été calculés sur 50 simulations.

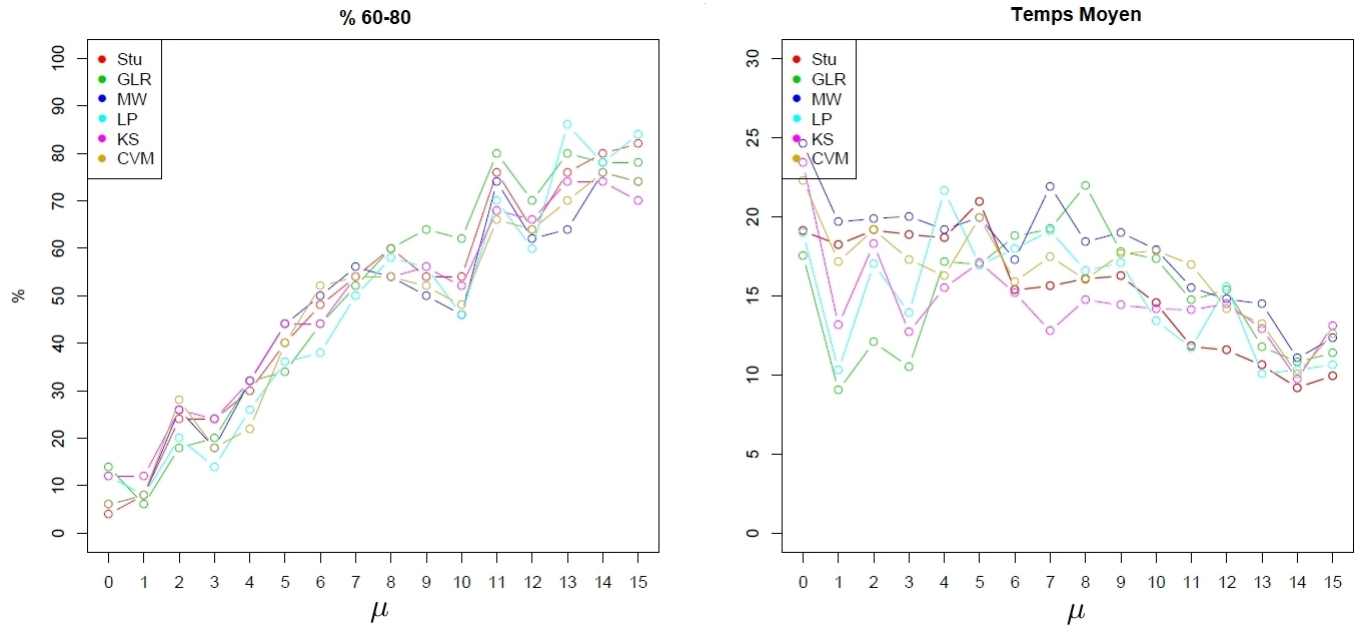


FIGURE 1.4 – Pourcentages de simulations d’Albumine avec ajout d’un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 0 à 15 ; Méthode avec classification

Les taux de fausses détections sont évalués de la même façon qu’auparavant. Le tableau 1.2 donne les pourcentages de simulations ayant détectées 0 ou 1 rupture pour 200 échantillons d’albumine simulés de 60 individus. On peut donc dire que dans 83 à 95% des cas, il n’y aurait pas de fausses alarmes lors d’une journée de mesures d’albumine en laboratoire et que s’il y en avait, elles seraient de 1 par jour.

Méthode	% 0	% 1
Stu	86	14
Bar	95	5
GLR	86,5	13,5
MW	83,5	16,5
Moo	94	6
LP	87,5	12,5
KS	85,5	14,5
CVM	83,5	16,5

TABLE 1.2 – Pourcentages de simulations où les méthodes de détection en ligne (avec classification) ont trouvé 0, 1, 2 ou 3 ruptures sur 200 simulations d’albumine de 60 individus

### Incertitude de mesure

Chaque année le CHU calcule les incertitudes de mesure pour chaque analyseur et chaque méthode d’analyses. Ces incertitudes, données en pourcentage, permettent de quantifier l’er-

reur aléatoire ou défaut de fidélité et l'erreur systématique ou erreur de justesse. Cette partie étudie l'effet des incertitudes de mesure sur les analyses des patients.

On suppose que les données des patients sont des données dites contaminées par les incertitudes de mesure. Dans un premier temps, nous estimons donc les densités des valeurs vraies des paramètres biologiques des patients pour chaque analyseur en utilisant des méthodes de déconvolution traitant les erreurs de mesure. Une fois ces densités de valeurs vraies estimées, nous simulons des échantillons suivant ces densités grâce à une méthode du rejet.

Pour étudier l'impact des incertitudes sur les valeurs vraies des mesures, nous avons comparé, par des tests de Kolmogorov-Smirnov, les distributions des valeurs avec ces mêmes valeurs contaminées par un pourcentage d'incertitude de mesure.

Les valeurs contaminées sont données par une valeur suivant la loi  $\mathcal{N}(m, \%IM \times m)$  où  $m$  est la valeur vraie et  $\%IM$  le pourcentage d'incertitude de mesure étudié.

Puis, nous avons ensuite étudié les échantillons simulés sous forme de séries temporelles pour déterminer à partir de quel pourcentage les incertitudes de mesure créent une rupture et donc un changement de moyenne, de variance ou de distribution sur les séries temporelles. L'incertitude de mesure est ajoutée pour 50 individus à  $i = 70$ .

Les pourcentages de simulations ayant trouvées une rupture entre  $i = 60$  et  $i = 80$  permettront de déterminer à partir de quel pourcentage l'incertitude de mesure crée des changements de moyenne ou de variance sur les échantillons.

### **1.2.2 Chapitre 3 : Étude de mesures de cendres volcaniques réalisées au Laboratoire Magmas et Volcans de l'Université Clermont Auvergne**

Dans ce chapitre nous traitons les mesures de forme de particules de cendres d'un volcan d'équateur, le Tungurahua. Ce travail fait suite aux travaux de Jean Luc Le Pennec, Sébastien Leibrandt et Julia Eychenne au laboratoire Magmas et Volcans de l'Université Clermont Auvergne [15] [8]. L'objectif de l'analyse de données présentée dans ce chapitre est de développer un outil pour comprendre la répartition spatiale des particules en fonction de leurs paramètres de forme.

#### **Jeu de données**

Les données décrivent la forme des particules de cendres. Elles ont été acquises grâce à un morpho-granulomètre G3 qui permet de mesurer automatiquement les caractéristiques morphologiques et géométriques d'une grande quantité de cendres.

Les paramètres étudiés sont : aspect ratio, circularité, convexité et solidité, ce sont des variables comprises entre 0 et 1. Le paramètre aspect ratio, mesure la forme plus ou moins ronde de la particule, plus celle-ci est allongée, plus l'aspect ratio tend vers 0. Le paramètre de

circularité mesure également la forme plus ou moins arrondie de la particule mais prend aussi en compte l'irrégularité du contour. Les paramètres convexité et solidité mesurent la tendance de la particule à être plus ou moins concave donc l'irrégularité du contour, plus la particule est lisse, plus ils tendent vers 1.

Les particules de cendres étudiées proviennent de 22 sites de collecte différents, F1 à F22, situés plus ou moins loin du cratère du volcan et sont séparées en trois tailles de grains : 75 à 90 $\mu\text{m}$ , 250 à 300 $\mu\text{m}$  et 710 à 850 $\mu\text{m}$ . Le lieu de collecte ainsi que la distance au cratère de ce lieu sont associés à chaque particule.

### Transformation des variables

Les distributions des paramètres montrent des asymétries. Ces spécificités ne conviennent pas à l'analyse statistique standard. Pour obtenir des distributions symétriques, on procède à une transformation Box-Cox. La figure 1.5 donne les distributions de la convexité pour les trois tailles de grains avec et sans transformation Box-Cox. Les analyses statistiques suivantes seront basées sur ces transformations.

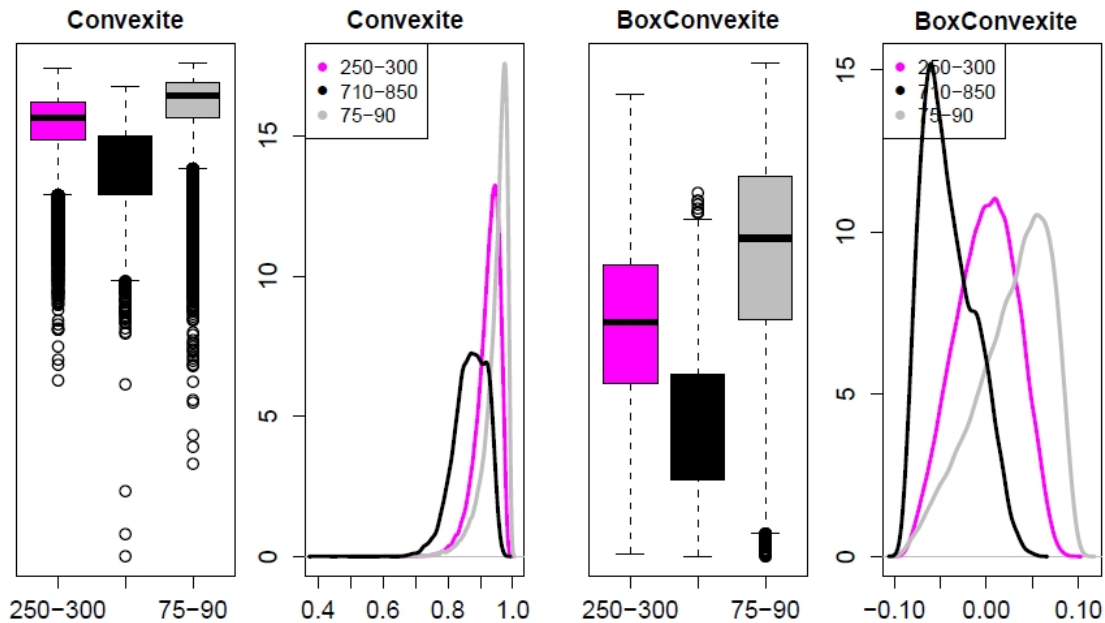


FIGURE 1.5 – Distribution du paramètre convexité pour les trois tailles de grains avec et sans transformation Box-Cox

Une analyse bivariée des corrélations entre les paramètres de forme et la localisation des particules a conduit à la dépendance de chaque paramètre de forme avec la localisation. Cependant la structure de l'ensemble de données ne permet pas d'obtenir un modèle expliquant ces relations.



## Classification

Une classification K-means est réalisée pour chaque taille de grains avec les paramètres de forme et la distance au cratère, 5 groupes sont gardés. La variable distance au cratère permet de prendre en compte, grâce à une variable quantitative, la répartition spatiale des particules.

Les groupes ainsi créés sont référés par étiquette par nombre et couleur : 1 (rose), 2(bleu), 3(vert), 4(violet), 5(brun), cela créé une nouvelle variable qualitative "groupe".

La figure 1.6 donne les distributions des variables de forme et la distance au cratère pour chaque groupe des particules de taille  $75-90\mu\text{m}$  ainsi que la répartition en pourcentage des groupes sur chaque lieu de collecte des particules.

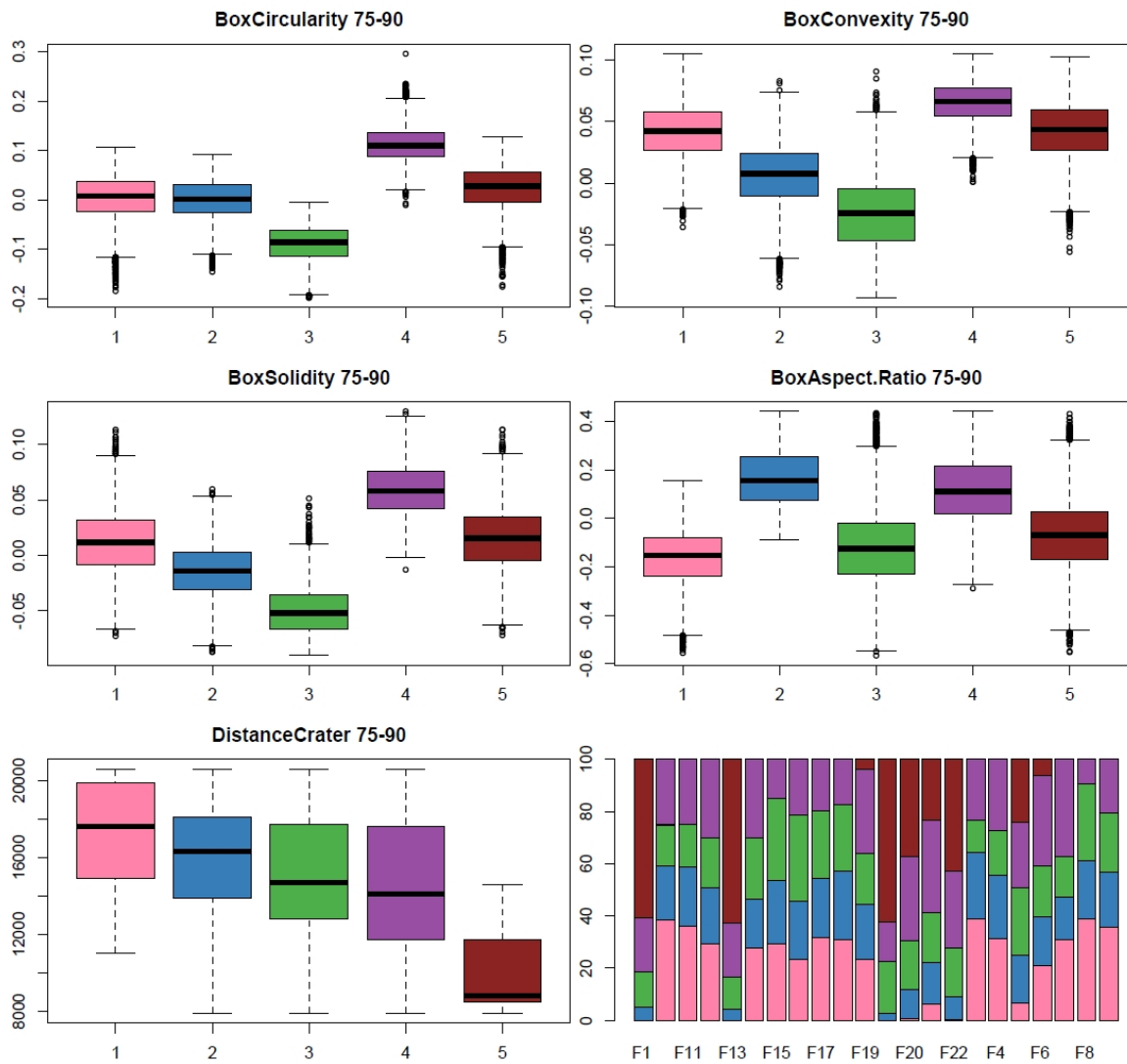


FIGURE 1.6 – Distribution des paramètres par groupe issus d'une classification K-means pour les particules de taille  $75-90\mu\text{m}$  et répartition des groupes par lieu de collecte des particules

Pour faciliter l'interprétation des groupes et mettre en évidence le lien entre les groupes, les paramètres de forme et la distance au cratère, les observations sont représentées sur des graphiques d'analyses en composantes principales (ACP). L'ACP résume l'interaction entre les variables de forme et la distance au cratère. Chaque graphique représente deux axes, chacun exprimant un pourcentage d'information (inertie) de l'ensemble des données. La figure 1.7 donne les représentations des individus sur les axes 1 et 2 puis 2 et 3, les points sont colorés selon le groupe auquel ils appartiennent. Les trois premiers axes de l'ACP représentent 94.9% d'informations.

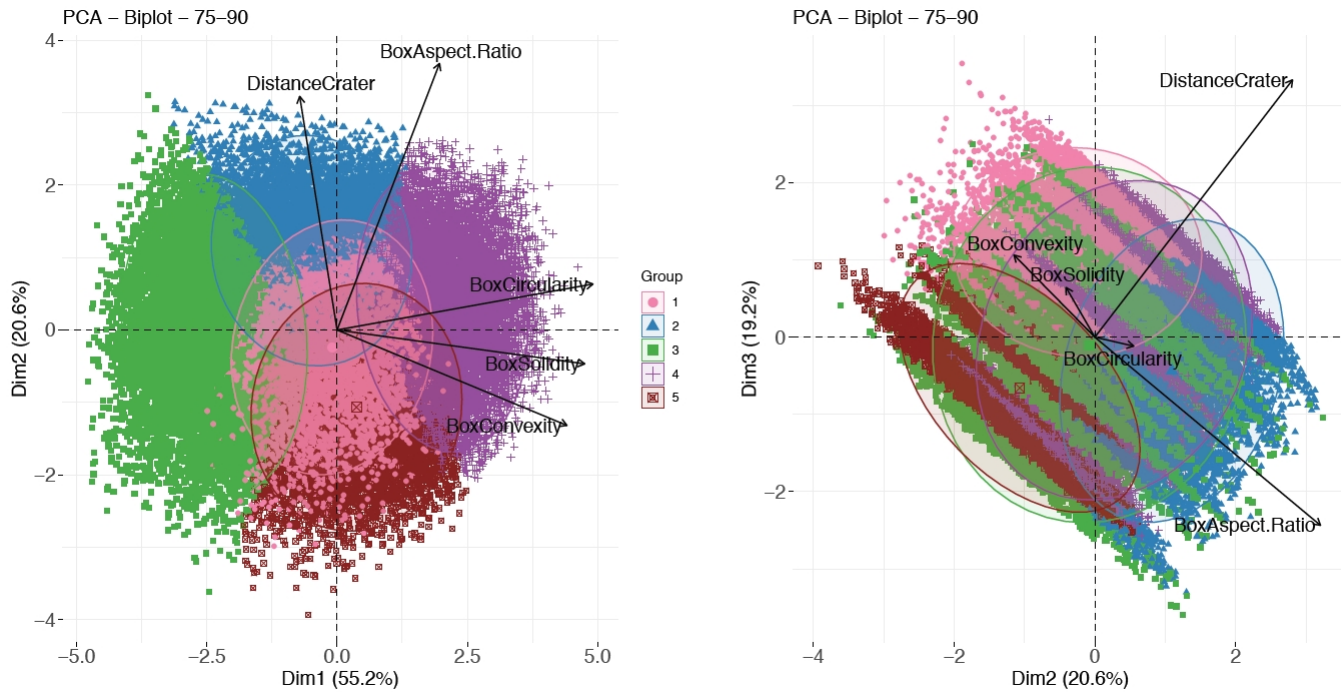


FIGURE 1.7 – Représentation de l'analyse en composantes principales des particules de taille 75-90 μm, avec représentation en couleur des groupes issus d'une classification K-means.)

Le premier axe fait ressortir les informations sur la circularité, la convexité et la solidité. Les scores positifs sur cet axe indiquent des valeurs élevées pour ces trois paramètres de forme et l'inverse se produit pour des scores négatifs. L'axe 2 fait ressortir les informations aspect ratio et distance au cratère.

Pour les particules de taille 75-90 μm, le groupe 4 est caractérisé par des valeurs plus élevées de circularité, convexité et solidité et des valeurs intermédiaires de la distance au cratère et d'aspect ratio. Le groupe 3 est à l'opposé du groupe 4 en ce qui concerne la circularité, la convexité et la solidité. Le groupe 2 a des valeurs intermédiaires de circularité, de convexité et de solidité et des valeurs plus élevées de distance au cratère et d'aspect ratio. Le groupe 1 a des valeurs intermédiaires de circularité, convexité, solidité et d'aspect ratio et des valeurs

plus élevées de distance au cratère. Le groupe 5 a des valeurs intermédiaires de circularité, convexité, solidité et d'aspect ratio ainsi que des valeurs basses de distance au cratère.

Une analyse des correspondances (CA) est ensuite réalisée pour faire ressortir les correspondances entre les groupes issus de la classification et les lieux de collecte des particules. La CA résume les relations en donnant des scores permettant de représenter les valeurs des variables qualitatives sur des graphiques 2D. En utilisant ces scores, nous pouvons classer les lieux. La figure 1.8 donne la représentation graphique des deux premiers axes de la CA pour les particules de taille 75-90 $\mu\text{m}$  et le dendrogramme classant les lieux. Les lieux étant dans une même classe ont des répartitions de groupes de particules issus de la classification k-means similaires.

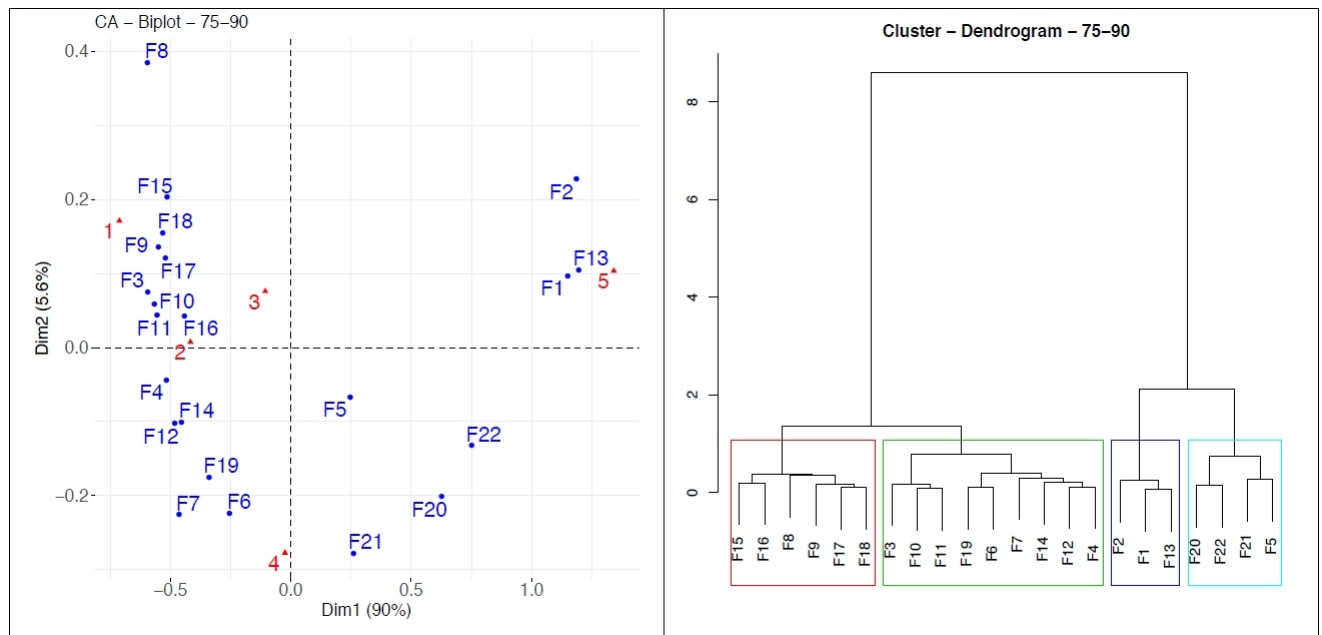


FIGURE 1.8 – Représentation de l'analyse des correspondances des particules de taille 75-90 $\mu\text{m}$ , et du dendrogramme classant les lieux

Les particules du groupes 5 sont principalement situées dans les lieux F1, F13, F2, F20, F21, F22 et F5 et non à des endroits éloignés du cratère. Dans ces endroits il n'y a pas ou peu de particules avec le profil du groupe 1, en particulier sur les lieux F1, F13 et F2.

## Modélisation

La dernière étape de cette étude est de modéliser le lien entre les groupes et les paramètres de formes ainsi que les variables de localisation. Pour cela, nous donnons les arbres de régression expliquant les groupes à l'aide des variables quantitatives de forme et ses groupes de lieux créés avec l'analyse des correspondances. Une stratégie fondée sur le sous-échantillonnage a été réalisée pour l'étape de validation du modèle, et un « taux de bien classés » a été déterminé.

Ce taux décrit la performance du modèle. Nous donnons également la matrice de confusion qui fournit un « taux de bien classés » par groupe.



## Chapitre 2

# Étude de mesures sanguines réalisées au CHU de Clermont Ferrand

### 2.1 Introduction

#### 2.1.1 Contrôle interne de qualité

Les CIQ sont réalisés au sein des laboratoires à l'aide d'échantillons de contrôle dans le but de surveiller la maîtrise des processus analytiques ainsi que la fiabilité des résultats et d'assurer la qualité des examens réalisés.

La mise en œuvre des CIQ vise à : vérifier la conformité analytique des résultats, détecter et corriger des erreurs ainsi que de les prévenir par l'observation d'un certain nombre de phénomènes (dérives, tendances, variabilités, ...) [7]

Pour chaque type d'examen réalisé dans le laboratoire, il y a trois niveaux de contrôle qui consistent à mesurer trois étalons plusieurs fois par jour. Ces trois niveaux représentent différentes valeurs atteignables par les paramètres biologiques mesurés lors des examens.

Les résultats de ces mesures sont étudiés et interprétés quotidiennement par les techniciens et biologistes. Ils sont représentés chronologiquement dans des tableaux de Levey-Jennings sous forme de graphiques. Ces tableaux représentent la moyenne des CIQ ainsi que les limites de décision qui sont égales à 1, 2 et 3 écarts types.

L'interprétation des graphiques de Levey-Jennings est gérée à l'aide des règles de Westgard [32]. Il s'agit de règles de rejet ou d'alarme : les règles de rejet indiquent que les résultats sont considérés comme « non conforme » et les règles d'alarme permettent de dégager une tendance.

Règles de rejet :

- 13s : 1 valeur éloignée de plus de 3 écarts-types de la moyenne
- 22s : 2 valeurs consécutives éloignées de plus de 2 écarts-types du même côté de la

moyenne

- R4s : 2 valeurs consécutives éloignées l'une de l'autre de plus de 4 écarts-types.

Règles d'alarme :

- 12s : 1 valeur éloignée de plus de 2 écarts-types de la moyenne
- 41s : 4 valeurs consécutives éloignées de plus de 1 écart-type du même côté de la moyenne
- 10x : 10 valeurs consécutives situées du même côté de la moyenne

Si une règle se déclenche, différentes actions peuvent être réalisées. Le contrôle peut être repassé, le point peut être accepté (en cas de règle d'alarme) ou refusé. Dans ce dernier cas, des actions pour résoudre le « problème » sont mises en place (désactivation de la méthode, calibration de la machine, ...) [28]

Suite à un problème lors des CIQ, une étude d'impact peut être réalisée sur les données des patients, elle peut nécessiter la reprise d'un échantillonnage des patients dosés depuis les derniers CIQ conformes sur un autre automate ou sur le même une fois fonctionnel. Si l'écart entre la nouvelle mesure et l'ancienne est considéré comme acceptable, alors la mesure du patient ne sera pas modifiée, sinon une troisième mesure sera effectuée et le compte rendu du patient sera modifiée.

L'écart acceptable est défini par la norme 5725-6 [1], il est appelé Delta-Check :

$$Delta - Check = \frac{2.8 \times \text{écart-type des CIQ}}{\text{moyenne des CIQ}} \times 100 \quad (2.1)$$

Un deuxième écart acceptable peut être pris en compte : le taux de changement critique (TCC), il permet de dire si la modification mesurée chez un individu est cliniquement significative. Ce taux est calculé à l'aide du coefficient de variation intra-individuelle RICOS [33], mais cette valeur n'existe pas pour tous les paramètres biologiques mesurés au CHU.

$$TCC = \sqrt{2} \times Z \times \sqrt{CV_a^2 + CV_w^2} \quad (2.2)$$

avec  $Z = 1.96$ ,

$CV_w$  : variation intra-individuelle déterminée sur des patients en bonne santé, issue d'une base de données mise à jour régulièrement [25] [22]

$CV_a$  : imprécision analytique, coefficient de variation calculé sur les CIQ.

Il y a donc un TCC et un Delta-Check par niveau de CIQ, le niveau 1 est utilisé si la mesure du patient est inférieure à la valeur moyenne des CIQ niveau 1, le niveau 2 si la mesure se trouve dans l'intervalle des valeurs moyennes des CIQ niveaux 1 et 2 et le 3 si la mesure est supérieure à la moyenne des CIQ niveau 3.

Si l'écart entre les deux mesures n'est pas considéré comme acceptable (supérieure à Delta-Check ou à TCC) alors la mesure est changée, sauf si le biologiste estime que la différence n'a pas d'impact sur l'interprétation biologique.

### **2.1.2 CHU Clermont-Ferrand**

Le CHU de Clermont-Ferrand se distingue en deux pôles : Gabriel-Montpied et Estaing.

Le pôle Gabriel-Montpied regroupe : urgences, cardiologie médicale et chirurgicale, RMNDO (rhumatologie médecine physique et de réadaptation, ophtalmologie, . . .), chirurgie, RHEUNNIRS (néphrologie, hémodialyse, pneumologie, endocrinologie, réanimation, urologie, . . .) et psychiatrie enfant et adulte. Le pôle Estaing regroupe : FEE (gynécologie obstétrique et reproduction humaine, pédiatrie générale et multidisciplinaire, chirurgie infantile, urgence pédiatrique, . . .), spécialités médicales et chirurgicales (chirurgie digestive et hépatobiliaire, chirurgie maxillo-faciale et chirurgie plastique, dermatologie et oncologie cutanée, Soins palliatifs, . . .) et ressources interventionnelles. Les analyses biomédicales peuvent être réalisées sur les deux sites.

Dans le cadre de cette thèse, les données étudiées proviennent des deux pôles du CHU et sont des mesures sanguines de patients.

## **2.2 Données**

### **2.2.1 Jeux de données**

Pour notre étude, nous avons pu étudier les mesures de sept paramètres biologiques sélectionnés par deux biologistes médicales. Ces paramètres sont accompagnés de l'âge et du sexe des patients, la date et l'heure de la mesure ainsi que l'analyseur sur lequel elle a été effectuée. Nous avons pu étudier deux jeux de données, l'un donnant les mesures pour six paramètres biologiques du 1er juillet 2016 au 31 juin 2017 et le deuxième donnant les mesures pour sept paramètres biologiques entre le 1er janvier 2018 et le 31 décembre 2018. Ce dernier jeu de données est complété par l'information Service qui donne le service où le patient fut hospitalisé.

Les mesures étudiées ont été réalisées sur 5 analyseurs différents : VISTA 1, 2, 3, 500 et 1500. VISTA 1, 2 et 3 se trouvent sur le site Gabriel-Montpied et VISTA 500 et 1500 sur le site Estaing.

Les paramètres biologiques sont accompagnés par leurs valeurs de référence (Tableau 2.1). Ces valeurs permettent d'interpréter les résultats d'analyses des patients et dépendent de l'âge et du sexe des personnes. Pour chaque patient, nous pouvons ainsi dire si les taux des paramètres biologiques entrent ou non dans les valeurs de référence, si c'est le cas les valeurs



seront dites usuelles.

Le CHU a fourni une liste de dates correspondant aux jours où un problème a été observé sur un analyseur lors des CIQ. Ces dates sont données avec le paramètre biologique et l'analyseur sur lequel les contrôles ont été effectués. Pour les données de juillet 2016 à Juin 2017, cinq dates correspondent à des paramètres biologiques étudiées lors de notre étude et 20 dates pour les données de 2018.

Le jeu de données de juillet 2016 à Juin 2017 contient 303545 patients, celui de 2018 290536. Les patients peuvent avoir des mesures pour plusieurs paramètres biologiques ou pour un seul.

### 2.2.2 Paramètres biologiques

Les variables biologiques étudiées sont des paramètres sanguins : l'albumine, le chlore, la créatinine, la ferritine, les protéines totales, l'urée et le PSA (antigène spécifique de la prostate).

— Albumine :

L'albumine est synthétisée par le foie et représente plus de 50% des protéines totales du plasma. Elle constitue la principale protéine de transport dans le sang. La diminution du taux d'albumine s'observe dans les atteintes hépatiques, les états de nutrition, dans les situations d'expansion des liquides biologiques ainsi qu'au court de pertes (syndrome néphrotique, brûlures étendues,...). [16]

— Chlore :

Le chlore régule la distribution des fluides extracellulaires, il est absorbé au niveau du tractus gastro-intestinal et l'excès est éliminé par les reins. L'hypochlorémie est due à une réduction de l'apport du chlore dans l'alimentation, des vomissements, des diarrhées, une diminution de réabsorption rénale ou des diurétiques. Quant à l'hyperchlorémie, elle est observée en cas de déshydratation, d'insuffisance rénale, de certaines formes d'acidose, d'apport chloré important par l'alimentation ou par voie parentérale et d'intoxication aux salicylés. [17]

Les valeurs de référence du chlore varient très peu et ne dépendent pas du sexe (Tableau 2.2.2).

— Créatinine :

La créatinine est un produit de dégradation de la créatine et est essentiellement éliminée par les reins. Sa variabilité intra-individuelle dépend du régime alimentaire et de l'exercice physique. Ainsi le jeûne, le régime végétarien la diminue et à l'inverse les régimes riches en protides et la consommation de tabac l'augmentent.

Le dosage de la créatinine permet le diagnostic ou le suivi d'affections rénales aiguës et chroniques ainsi que la surveillance de dialyses rénales. Une diminution du taux peut être observée en cas de myopathie. Une augmentation peut être le signe d'insuffisances rénales, prématurés, prééclampsies ou peut être observée en cas de leucémie, goutte, hyperthyroïdie, acromégalie et gigantisme, diabète, hypertension et insuffisance cardiaque.[18] Le taux de créatinine étant très variable selon l'âge ou le sexe des personnes, les valeurs de référence évoluent donc en fonction de ces deux paramètres (Tableau 2.2.2).

— Ferritine :

La ferritine est la protéine de mise en réserve du fer dans l'organisme. Elle a deux fonctions : la réserve et la détoxification du fer. Elle est un bon reflet des réserves en fer de l'organisme.

Le taux de ferritine varie en fonction de l'âge et du sexe des personnes comme le montre les valeurs de référence (Tableau 2.1). Il permet d'évaluer les réserves en fer de l'organisme et ainsi détecter des carences ou d'autres diagnostics si la mesure est accompagnée par d'autres bilans.[5]

— Protéines totales :

Les protéines plasmatiques interviennent principalement dans le maintien de la pression oncotique sanguine, dans le transport non spécifique de substances comme le fer, l'hémoglobine, les phospholipides ou des médicaments, dans la coagulation, dans l'immunité humorale ainsi que dans les systèmes tampons sanguins.

La concentration de protéines plasmatiques varie en fonction de l'état d'hydratation de l'organisme. Chez le nouveau-né, le taux est inférieur à celui de l'adulte, il augmente progressivement durant l'enfance et l'adolescence pour se stabiliser à l'âge adulte et diminuer légèrement chez les personnes âgées (Tableau 2.1). [19]

— Urée :

L'urée est un produit azoté issu du catabolisme des protéines. Elle est synthétisée au niveau du foie. Le niveau d'urée dans le plasma dépend de la balance entre la production et l'excrétion rénale.

Une augmentation de l'urée sanguine peut être le signe de : hémorragies gastro-intestinales, atteintes rénales, syndromes urémiques, insuffisance cardiaque, . . . Une diminution peut signifier : une hépatite toxique, une tumeur hépatique ou de l'alcoolisme. [21]

Les valeurs de référence de l'urée varient en fonction de l'âge des patients (Tableau 2.1).

— PSA (antigène spécifique de la prostate) :

Des taux de PSA élevés reflètent généralement une affection de la prostate. Le PSA étant

présent dans les glandes para-urétrales et anales, ainsi que dans les tissus mammaires sains ou cancéreux, on peut trouver, chez la femme, de faibles concentrations sanguines en PSA.

Le dosage de PSA est utile pour le suivi d'affection et de l'efficacité du traitement chez les patients atteints du cancer de la prostate ou les patients sous hormonothérapie. [20]

Le taux de PSA ne varie ni en fonction du sexe, ni en fonction de l'âge (Tableau 2.1).

### Bornes de référence

	Max âge		Unité âge	Sexe	Limite Basse (Inclusive)	Limite Haute (Inclusive)
	Jours (Exclusif)		j=jour m=mois a=an			
<b>Albumine</b>	14	7	j	?	30	43
	365.25	1	a	?	27	48
	2922	8	a	?	37	47
	5478.75	15	a	?	39	49
	6939.75	19	a	Femme	37	51
	6939.75	19	a	Homme	40	53
	6939.75	19	a	?	37	53
	43800	120	a	?	35	52
<b>Chlore</b>	7	7	j	?	97	108
	30	30	j	?	97	108
	182.58	6	m	?	97	108
	365.25	1	a	?	97	106
	6570	18	a	?	97	107
	43800	120	a	?	98	107
<b>Créatinine</b>	14	14	j	?	27.2	77.6
	730,5	2	a	?	8.4	30.7
	1826,25	5	a	?	19.6	36
	4383	12	a	?	25.9	50.9
	5478,75	15	a	?	37.9	68.3
	6939,75	19	a	Homme	52.4	90.8
	6939,75	19	a	Femme	41.2	70.5
	6939,75	19	a	?	41.2	90.8
	43800	120	a	Homme	59.2	104
	43800	120	a	Femme	45.1	84
	43800	120	a	?	45.1	104

### Bornes de référence

	Max âge  Jours (Exclusif)		Unité âge  j=jour m=mois a=an	Sexe	Limite Basse (Inclusive)	Limite Haute (Inclusive)
<b>Ferritine</b>	30,43	1	m	Homme	47	438
	30,43	1	m	Femme	47	554
	30,43	1	m	?	47	554
	182,58	6	m	Homme	47	449
	182,58	6	m	Femme	47	379
	182,58	6	m	?	47	449
	365,25	1	a	Homme	47	120
	365,25	1	a	Femme	47	85
	365,25	1	a	?	47	120
	1826,5	5	a	?	47	101
	6939.75	19	a	Homme	47	358
	6939.75	19	a	Femme	47	110
	6939.75	19	a	?	?	?
	43800	120	a	Homme	26	388
	43800	120	a	Femme	8	252
43800	120	a	?	?	?	
<b>Protéines</b>	14	14	j	?	54	86
	365.25	1	a	?	45	73
	2191.5	6	a	?	63	78
	3287.25	9	a	?	66	80
	6939.75	19	a	?	67	84
	43800	120	a	?	64	82
<b>PSA</b>	6570	18	a	?	?	?
	43800	120	a	?	?	4
<b>Urée</b>	14	14	j	?	1	8.3
	3652.5	1	a	?	3.2	8
	6939.75	19	a	Male	2.6	7.6
	6939.75	19	a	Female	2.6	6.9
	6939.75	19	a	?	2.6	7.6
	43800	120	a	?	2.5	6.4

TABLE 2.1 – Bornes de référence des paramètres biologiques

### 2.2.3 Différences pôle Estaing et Gabriel-Montpied

Les données analysées dans cette partie sont des mesures de paramètres biologiques réalisées au CHU de Clermont-Ferrand sur deux sites : Gabriel-Montpied et Estaing. Ces deux

sites regroupent des pôles différents du CHU, ils n'ont donc pas les mêmes panels de patients.

La première différence se trouve au niveau de l'âge. La figure 2.1 représente les histogrammes de l'âge des patients pour les deux sites. Les services de pédiatrie se trouvant à Estaing, la moyenne d'âge est plus basse (53.42 ans pour Estaing, 64.4 ans pour Gabriel-Montpied) (Tableau 2.2), Le pourcentage de mineurs est aussi plus bas à Estaing (15.97% à Estaing et 2.09% à Gabriel-Montpied).

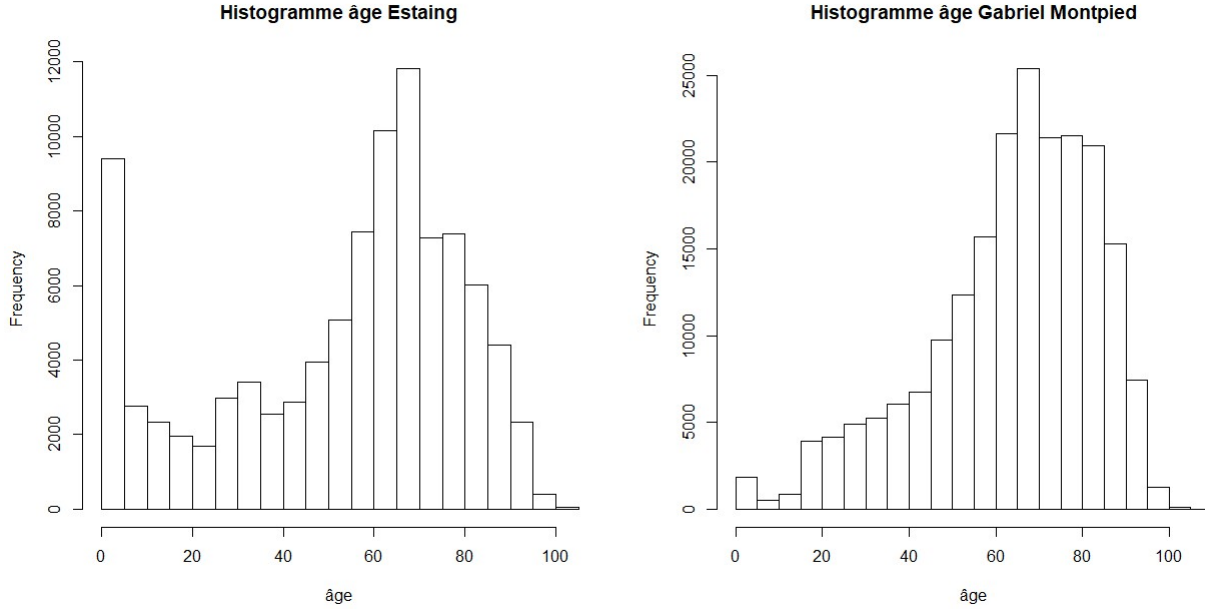


FIGURE 2.1 – Histogrammes de l'âge, Estaing / Garbiel-Montpied

	Min	25%	50%	75%	Max	Mean
Estaing	0	35	61	73	105	53.42
Gabriel Montpied	0	54	67	79	110	64.2

TABLE 2.2 – Quantiles et moyennes de l'âge, Estaing / Garbiel-Montpied

Les mesures sont réalisées en majorité le matin (Figure 2.2). En effet, à Estaing 79% des mesures sont réalisées entre 7h et 12h, ce pourcentage est de 54% pour Gabriel-Montpied. Le site Estaing étant fermé de 20h à 6h, un panel de patients plus jeunes apparaît à Gabriel-Montpied. De jour (6h-20h), la moyenne d'âge est de 65,4 et le pourcentage de patients mineurs est de 1.43%, la nuit ce pourcentage passe à 5.60% et la moyenne diminue à 57.75 ans.

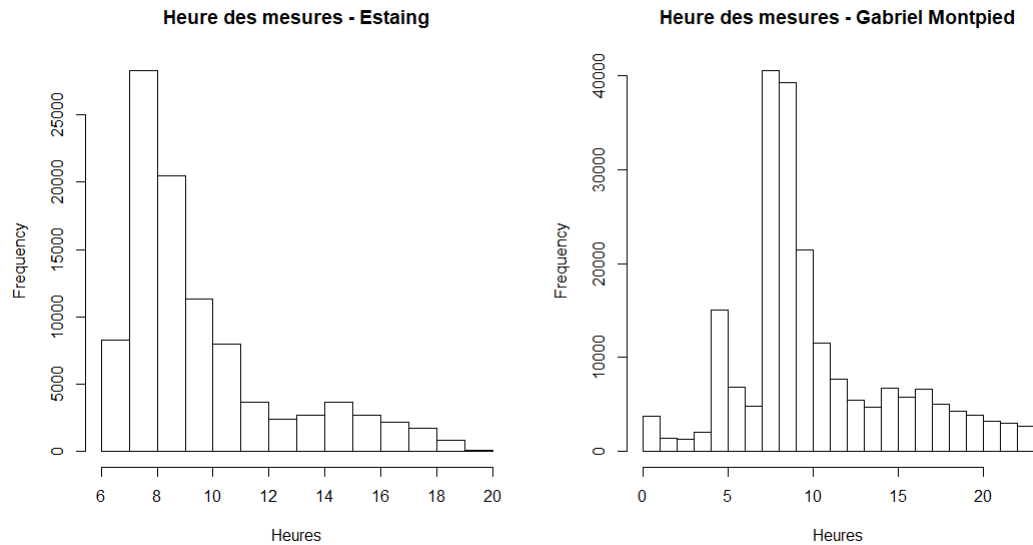


FIGURE 2.2 – Histogrammes des heures des mesures Estaing / Garbiel-Montpied

À Estaing les patients sont à 46.2% des femmes et à 53.8% des hommes, à Gabriel-Montpied ils sont à 44% des femmes et 56% des hommes. Le pourcentage de femmes plus élevé à Estaing est expliqué par la présence des services : gynécologie obstétrique,....

Ces différentes informations ont été calculées sur le jeu de données de juillet 2016 à juin, 2017, mais les résultats sont les mêmes en 2018.

## 2.2.4 Différences par service

Les pôles Gabriel Montpied et Estaing sont segmentés en différents services. Dans le jeu de données de 2018, on peut décompter 325 services distincts. Certains services sont rares et ne comptent qu'un patient correspondant aux sept paramètres étudiés ici comme l'addictologie, neuroradiologie GM ou convalescence B, à l'inverse, certains services regroupent beaucoup de patients.

Les services peuvent regrouper des panels de patients différents, le tableau 2.3 donne le nombre de patients (dans notre jeu de données) la moyenne d'âge, et les pourcentages d'hommes et de femmes, pour les 15 principaux services de Garbiel Montpied puis d'Estaing.

	Services	Nombre de patients	Moyenne âge	% Femme	% Homme
Gabriel Montpied	ACCUEIL.URGENCES _URG	32995	55,84	50,22	49,78
	UNITE.DE.SOINS.INTENSIFS.CARDIOLOGIE _SI	9167	68,01	30,57	69,43
	REANIMATION.MEDICO.CHIRURGICALE _REA	6122	57,22	21,92	78,08
	MEDECINE.POST.URGENCE.ET.THERAPEUTIQUE _HC	5186	75,05	54,20	45,80
	CARDIOLOGIE.B2 _HC	5161	78,20	43,98	56,02
	CARDIOLOGIE.A _HC	4863	70,39	43,41	56,59
	ONCOLOGIE.MEDICALE.CJP	4667	63,67	53,07	46,93
	NEURO.REANIMATION _REA	4546	56,87	32,78	67,22
	REANIMATION.CHIRURGIE.CARDIO.VASCULAIRE _REA	4540	67,44	28,19	71,81
	PNEUMOLOGIE _HC	4413	66,25	33,85	66,15
	CHIRURGIE.CJP	4292	61,50	56,83	43,17
	CHIRURGIE.CARDIO.VASCULAIRE _HC	4260	67,09	27,11	72,89
	NEPHROLOGIE _HC	4240	62,55	35,71	64,29
	REANIMATION..MEDICALE _REA	4036	64,18	29,93	70,07
	RHUMATOLOGIE _HC	3920	72,28	57,76	42,24
Estaing	CHIRURGIE.DIGESTIVE.ET.HEPATOBIILAIRE. _HC	7702	64,00	44,60	55,40
	HEMATOLOGIE.CLINIQUE _HC	7076	59,92	36,07	63,93
	MEDECINE.DIGESTIVE.ET.HEPATOBIILAIRE _HC	6310	61,67	36,40	63,60
	CLINIQUE.DE.DURTOL	6055	66,83	31,42	68,58
	SOINS.INTENSIFS.DHEMATOLOGIE.CLINIQUE _SI	5058	54,03	39,50	60,50
	REANIMATION.ADULTES.ESTAING _REA	5000	61,67	34,62	65,38
	MEDECINE.INTERNE.ESTAING _HC	3625	65,82	49,41	50,59
	HOPITAL.DE.JOUR.HEMATOLOGIE.CLINIQUE _HJ	3197	58,23	34,75	65,25
	CONSULTATIONS.HEMATOLOGIE.CS	2700	60,30	45,70	54,30
	SURVEILLANCE.CONTINUE.A.ESTAING _SC	2498	62,15	40,47	59,53
	CONSULTATION.DIGESTIVE.ET.HEPATOBIILAIRE _CS	2223	53,98	40,17	59,83
	CONSULTATIONS.MEDICO.CHIRURGICALES.PEDIATRIE _CS	2206	8,38	39,21	60,79
	HJ.PEDIATRIE.GENERALE.et.HEMATO.ONCO.PEDIATRIQUE _HJ	1977	9,60	40,21	59,79
	HEMATO.ONCOLOGIE.PEDIATRIQUE _SC	1875	8,97	42,51	57,49
	DERMATOLOGIE.ET.ONCOLOGIE.CUTANEE _HC	1696	71,04	49,35	50,65

TABLE 2.3 – Nombres de patients, moyennes d'âge et pourcentages d'hommes et de femmes pour les 15 principaux services du pôle Gabriel Montpied puis d'Estaing

## 2.3 Étude des paramètres biologiques

### 2.3.1 Distributions des paramètres

Les paramètres, albumine, chlore, créatinine, protéines et urée, sont mesurés sur les 5 analyseurs (VISTA1, VISTA2, VISTA 3, VISTA 500 et VISTA 1500). La ferritine n'est jamais mesurée sur VISTA1 et le PSA est mesuré seulement sur les analyseurs VISTA2 et VISTA3.

Le tableau 2.4 indique le nombre de mesures réalisées pour chaque paramètre biologique et analyseur pour une année ainsi que la moyenne par jour. Tout d'abord, il les indique pour toutes les mesures réalisées puis pour les mesures entrant dans les valeurs de référence, dites valeurs usuelles. On peut remarquer que le CHU effectue beaucoup moins de mesures de Ferritine et de PSA. Au pôle Gabriel-Montpied, l'analyseur le plus utilisé est le VISTA1 et à Estaing, c'est le VISTA1500.

		Toutes valeurs		Valeurs usuelles	
		Sur 1 an	Moyenne par jour	Sur 1 an	Moyenne par jour
Albumine	VISTA1	4016	11	1161	3
	VISTA2	13738	37	5239	14
	VISTA3	17051	46	6610	18
	VISTA500	6310	17	2138	6
	VISTA1500	12554	34	4439	12
Chlore	VISTA1	83633	228	53110	145
	VISTA2	56814	155	33651	92
	VISTA3	51676	141	32905	90
	VISTA500	25613	70	17426	47
	VISTA1500	54607	149	39756	106
Créatinine	VISTA1	82996	226	49296	134
	VISTA2	56851	155	33105	90
	VISTA3	51517	140	29376	80
	VISTA500	25378	69	16455	45
	VISTA1500	82996	147	34021	92
Ferritine	VISTA1	0	0	0	0
	VISTA2	5123	14	3000	8
	VISTA3	5525	14	3172	9
	VISTA500	2649	7	1538	4
	VISTA1500	5371	15	3120	9
Protéines	VISTA1	81776	223	56921	155
	VISTA2	55684	152	38344	104
	VISTA3	50280	137	33790	92
	VISTA500	24515	67	13593	37
	VISTA1500	49852	136	29063	79
PSA	VISTA1	0	0	0	0
	VISTA2	1165	3	840	2
	VISTA3	1246	3	872	2
	VISTA500	0	0	0	0
	VISTA1500	0	0	0	0
Urée	VISTA1	81535	222	39241	107
	VISTA2	55701	152	27797	76
	VISTA3	50436	137	23959	65
	VISTA500	23725	65	13575	37
	VISTA1500	48847	133	27797	76

TABLE 2.4 – Nombres de mesures réalisées pour chaque paramètre biologique et analyseur sur un an (année 2018) ainsi que la moyenne par jour – Toutes valeurs et valeurs usuelles.

Le tableau 2.5 indique les quantiles et moyennes des paramètres biologiques. Les paramètres peuvent prendre des valeurs très élevées, ces valeurs restent probables et sont sûrement dues à certaines pathologies des patients. Nous savons que certains paramètres dépendent fortement du sexe et de l'âge de la personne. L'annexe A.1 décrit les quantiles et les moyennes des paramètres biologiques en fonction de l'âge et du sexe et A.2 représente ces moyennes



graphiquement. Cela permet de voir les différences dues à ces deux paramètres. Les catégories d'âge de l'annexe A.1 sont celles données dans les valeurs de référence, elles ne sont donc pas les mêmes pour chaque paramètre puisqu'ils ne dépendent pas tous de la même façon de l'âge.

	Min	25%	50%	Moyenne	75%	Max
Albumine	3.00	25.40	31.70	32.02	38.20	71.70
Chlore	51	102	105,00	105,00	108,00	165
Créatinine	12.40	56.60	72.50	96.88	97.00	2580.00
Ferritine	0.70	40.08	128.75	518.40	380.02	281077.00
Protéines	11.00	62.80	69.30	68.66	75.00	152.50
Urée	0.400	4.200	6.000	8.202	9.200	115.800
PSA	0.01	0.68	1.73	22.15	5.00	3600.00

TABLE 2.5 – Quantiles et moyennes des paramètres biologiques calculés sur le jeu de données 2018

### 2.3.2 Effets des services sur les distributions des paramètres biologiques

Comme expliqué précédemment, les services peuvent avoir des panels de patients très variés, ce qui engendre des distributions différentes de paramètres biologiques selon le service d'où proviennent les patients. La figure 2.3 et le tableau 2.6 donnent les distributions et résumés de l'albumine dans les cinq services ayant fait le plus de mesures d'albumine aux pôles Estaing et Gabriel Montpied. On peut remarquer que dans des services de réanimation les taux d'albumine sont plus bas que dans les services comme la rhumatologie ou les consultations digestives et hépatobiliaires.

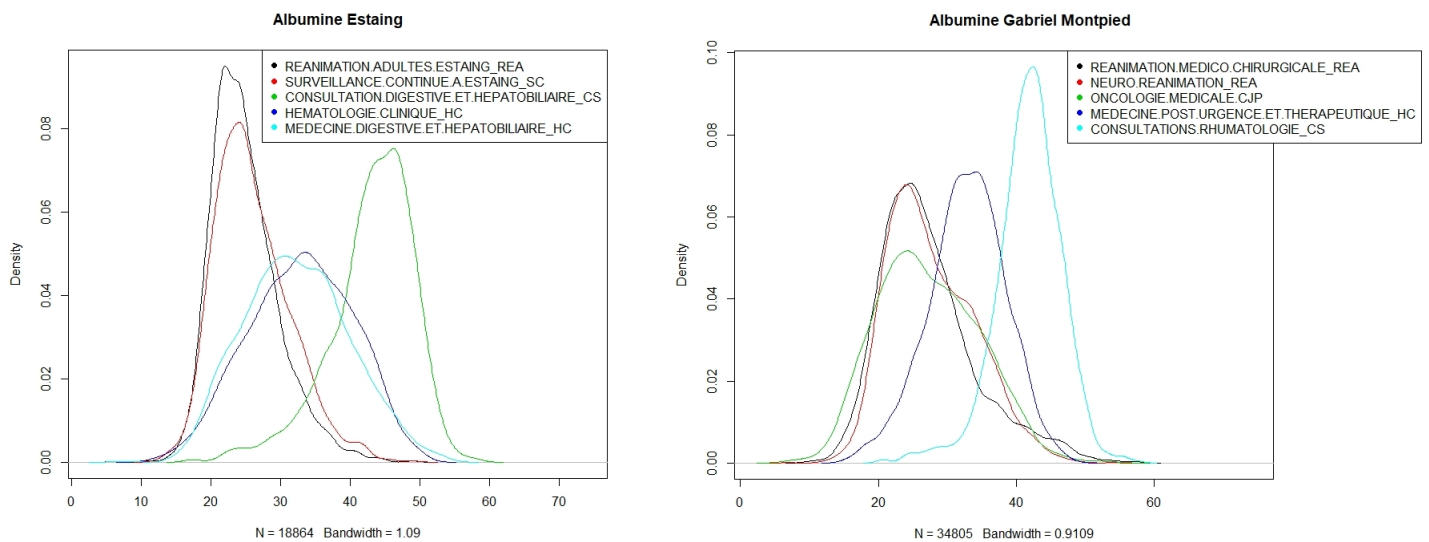


FIGURE 2.3 – Distributions de l'Albumine pour les cinq services ayant fait le plus de mesures d'Albumine en 2018 aux pôles Estaing et Gabriel Montpied

		Min	25%	50%	Moyenne	75%	Max
Gabriel Montpied	REANIMATION.MEDICO.CHIRURGICALE_REA	9,80	22,30	25,90	27,20	30,60	57,70
	NEURO.REANIMATION_REA	6,90	22,90	26,70	27,73	32,30	54,10
	ONCOLOGIE.MEDICALE.CJP	6,80	21,80	26,70	27,32	32,60	54,80
	MEDECINE.POST.URGENCE.ET.THERAPEUTIQUE_HC	15,40	29,50	33,00	32,85	36,70	48,10
	CONSULTATIONS.RHUMATOLOGIE_CS	20,60	39,20	42,10	41,81	44,73	57,60
Estaing	REANIMATION.ADULTES.ESTAING_REA	5,70	21,50	24,20	24,91	27,50	49,40
	SURVEILLANCE.CONTINUE.A.ESTAING_SC	12,00	22,10	25,10	26,00	29,20	49,50
	CONSULTATION.DIGESTIVE.ET.HEPATOBILIAIRE_CS	17,00	40,10	43,90	43,08	47,30	58,70
	HEMATOLOGIE.CLINIQUE_HC	11,20	27,70	33,20	32,90	38,60	50,50
	MEDECINE.DIGESTIVE.ET.HEPATOBILIAIRE_HC	7,50	26,90	31,90	32,33	37,20	53,50

TABLE 2.6 – Résumés de l’Albumine pour les cinq services ayant fait le plus de mesures d’Albumine en 2018 aux pôles Gabriel Montpied et Estaing

L’annexe A.3 donne les mêmes informations pour les autres paramètres biologiques étudiés : chlore, créatinine, ferritine, protéines, PSA et urée. On remarque que le chlore est le paramètre qui varie le moins en fonction du service. Le PSA peut atteindre des valeurs extrêmes dans certains services, il y a beaucoup moins de valeurs de PSA que pour les autres paramètres, et donc peu de mesures par service.

Les différences des distributions des paramètres biologiques selon les services s’expliquent soit par le panel de patients présents, par exemple un service de pédiatrie aura une moyenne de créatinine plus basse, puisque les enfants ont des taux beaucoup plus bas que les adultes, soit par les caractéristiques médicales présentes dans les services, en effet, ils peuvent regrouper des pathologies spécifiques qui pour certaines entraînent des variations des paramètres biologiques.

La figure 2.7 indique les moyennes des paramètres biologiques par analyseurs, calculées via une méthode Bootstrap. Pour la créatinine, la ferritine, les protéines et l’urée, les moyennes sont plus basses sur les analyseurs VISTA1500 et VISTA500 que les autres, elles sont donc plus basses sur le site Estaing qu’à Gabriel-Montpied.

	Albumine	Chlore	Créatinine	Ferritine	Protéines	PSA	Urée
VISTA1	30.46	104.89	103.78	-	69.41	-	8.77
VISTA2	32.14	105.46	105.32	476.33	69.81	18.23	8.45
VISTA3	32.36	105.00	107.53	465.01	69.31	19.56	8.90
VISTA500	31.79	104.62	73.24	417.38	66.11	-	6.78
VISTA1500	32.04	104.63	77.87	388.34	66.78	-	6.91

TABLE 2.7 – Moyennes des paramètres biologiques par analyseur, calculées en Bootstrap

Les différentes populations observées sur les deux sites du CHU s’expliquent par les différents services présents sur chacun. Certains services, très spécifiques, ne se situent que sur un pôle, comme la gynécologie obstétrique et la pédiatrie à Estaing ou l’hémodialyse à Gabriel-Montpied.

On n’observe pas de différence d’albumine selon les lieux car l’albumine diffèrent sur des services tel que la réanimation, hors un tel service se trouve à Gabriel Montpied et à Estaing. Quant au chlore, il varie très peu selon le service, l’âge et le sexe des patients, ce qui explique qu’il n’y est pas de différence pour ce paramètre entre les deux pôles.

### 2.3.3 Corrélations des paramètres

	Albumine	Chlore	Créatinine	Ferritine	Protéines	Urée	PSA
Albumine	1,000	-0,045	0,001	-0,189	0,603	-0,165	-0,165
Chlore	-0,045	1,000	-0,092	-0,091	-0,173	-0,092	0,011
Créatinine	0,001	-0,092	1,000	0,033	-0,006	0,662	0,006
Ferritine	-0,189	-0,091	0,033	1,000	-0,198	0,087	0,588
Protéines	0,603	-0,173	-0,006	-0,198	1,000	-0,045	-0,087
Urée	-0,165	-0,092	0,662	0,087	-0,045	1,000	0,046
PSA	-0,165	0,011	0,006	0,588	-0,087	0,046	1,000

TABLE 2.8 – Corrélations des paramètres biologiques

Le tableau 2.8 indique les corrélations entre les paramètres biologiques et l’âge. Les paramètres biologiques les plus corrélés sont donc : l’urée et la créatinine (0.66), l’albumine et les protéines (0.60) et la ferritine et le PSA (0.588).

### 2.3.4 Classification

Les classifications sont des méthodes d’analyses de données permettant d’organiser les données en groupes homogènes. Elles permettent de repérer des groupes d’individus ayant des caractéristiques communes. L’objectif peut être de créer des regroupements d’individus ou de construire une règle de classement lorsque l’on connaît déjà les groupes.

Dans notre étude, le but des classifications est de repérer des groupes de patients ayant les mêmes caractéristiques pour les paramètres biologiques étudiés. On étudiera ensuite la répartition de l’âge, du sexe et des analyseurs dans les groupes ainsi que celle des groupes dans chaque service.

La classification est réalisée sur les données de 2018 afin de pouvoir prendre en compte les services. Pour la réaliser, nous ne pouvons étudier uniquement les patients ayant des résultats pour tous les paramètres biologiques, cela représente 82 individus si l’on prend les sept paramètres. Ce nombre est dû au peu de mesures de PSA, sans ce paramètre on a 4487 individus. La classification est donc réalisée sur : albumine, chlore, créatinine, ferritine, protéines et urée. Les patients étudiés proviennent de 136 services différents, dont seulement 18 comptent plus de 50 patients sur ce jeu de données.

## K-Means

La classification réalisée est une méthode non supervisée permettant de rechercher une segmentation ou une typologie caractérisant l'ensemble des observations. Nous utilisons la méthode k-means.

Soit  $k$  le nombre de classes  $C_i$  ( $i = 1, \dots, k$ ) recherchées, l'algorithme kmeans va identifier  $k$  centres de gravité  $\hat{\mu}_i$  minimisant la distance entre eux et les points assignés à leur classe associée :

$$\frac{1}{N} \sum_x (x - \mu_{\hat{\omega}(x)})^t \cdot (x - \mu_{\hat{\omega}(x)}) \quad (2.3)$$

où  $\hat{\omega}(x)$  est la classe d'assignation de la donnée  $x$ , soit la classe dont le centre de gravité est le plus proche de  $x$  :

$$\hat{\omega}(x) = \arg \min_i (x - \mu_i)^t \cdot (x - \mu_i) \quad (2.4)$$

L'algorithme de la méthode k-means est le suivant :

---

**Algorithm 1** K-means

---

Données entrées :  $k$  et  $x$

Positionnement initial des centres de gravité des classes :  $\hat{\mu}_i^{(0)}$  avec  $i = 1, \dots, k$

$t = 1$

Tant que le critère d'arrêt n'est pas satisfait :

— Assignation des points : assigner chaque point au groupe dont le centre est le plus proche :

$$\hat{\omega}(x)^{(t)} = \arg \min_i (x - \mu_i^{(t-1)})^t \cdot (x - \mu_i^{(t-1)})$$

— Soit  $S_i^{(t)}$  l'ensemble des points assignés à la classe  $i$  :  $S_i^{(t)} = \{x \text{ tq } \hat{\omega}(x)^{(t)} = i\}$

— Mise à jour des centres de gravité des  $k$  classes :  $\mu_k(t) = \frac{1}{|S_i^{(t)}|} \sum_{x \in S_i^{(t)}} x$

—  $t = t + 1$

---

Il existe différents critères d'arrêt possibles tel qu'un nombre maximum d'itérations ou un seuil  $s$  pour la différence entre les centres de gravité entre deux itérations :

$$s = \sum_i (\mu_i^{(t)} - \mu_i^{(t-1)})^t \cdot (\mu_i^{(t)} - \mu_i^{(t-1)}) \quad (2.5)$$

## Résultats

La figure 2.4 et le tableau 2.9 donnent les distributions et les résumés des paramètres biologiques sur lesquels ont été réalisés la classification ainsi que l'âge. La figure 2.6 donne les proportions de sexe et d'analyseurs dans chaque groupe et la figure 2.7 donne les proportions de groupes dans les services représentant plus de 50 individus dans le jeu de données utilisé pour la classification.

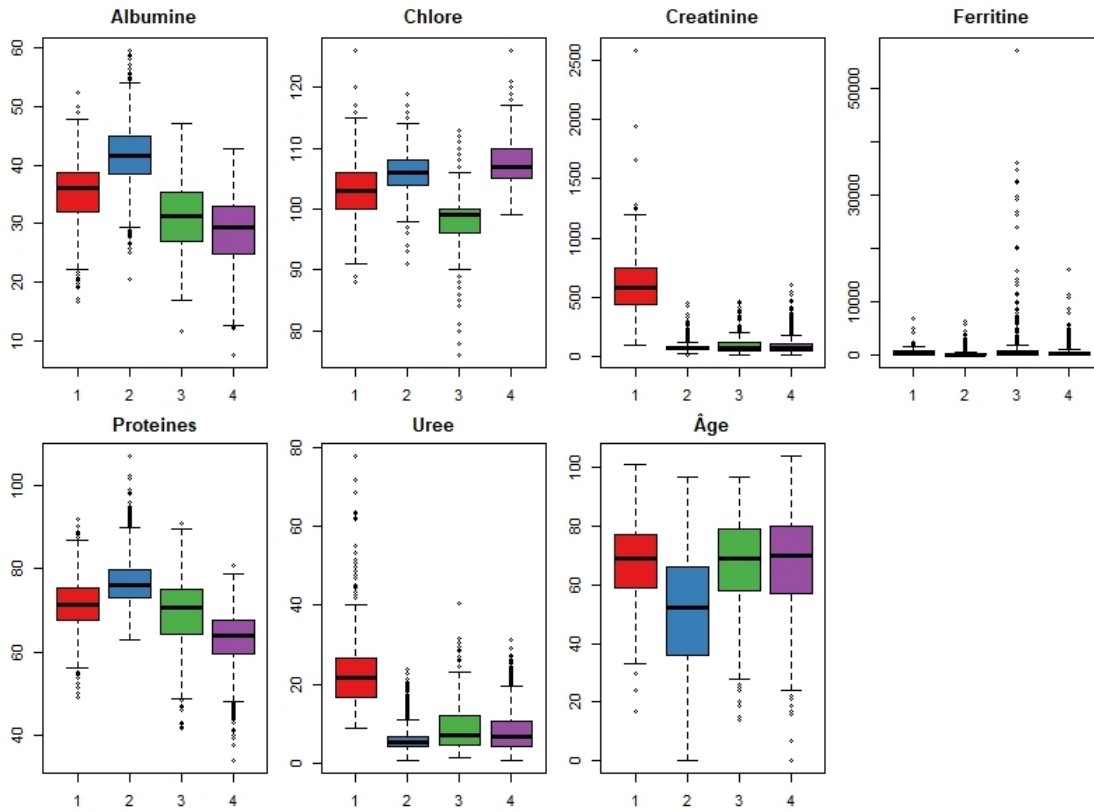


FIGURE 2.4 – Distributions des paramètres biologiques par classe résultant d'une classification k-means

Groupe	Albumine				Chlore			
	1	2	3	4	1	2	3	4
Min	16,70	20,40	11,60	7,50	88,00	91,00	76,00	99,00
25 %	32,00	38,60	26,90	24,90	100,00	104,00	96,00	105,00
50 %	36,00	41,60	31,40	29,45	103,00	106,00	99,00	107,00
Moyenne	35,19	41,78	31,15	28,69	103,03	106,05	97,87	107,34
75 %	38,80	44,90	35,30	33,08	106,00	108,00	100,00	110,00
Max	52,40	59,60	47,10	42,80	126,00	119,00	113,00	126,00

Groupe	Créatinine				Ferritine			
	1	2	3	4	1	2	3	4
Min	95,50	17,90	19,30	12,50	6,30	2,00	6,90	2,40
25 %	438,00	58,88	55,90	55,90	171,48	40,70	125,00	74,90
50 %	582,50	69,30	73,30	73,05	394,55	94,50	340,10	227,80
Moyenne	602,12	77,84	99,63	96,88	554,69	206,25	1838,95	530,58
75 %	744,25	84,80	116,00	106,00	775,70	225,23	869,90	514,30
Max	2580,00	455,00	464,00	607,00	6885,00	6443,60	57208,50	16241,20

Groupe	Protéines				Urée			
	1	2	3	4	1	2	3	4
Min	49,00	62,90	41,60	33,90	8,60	0,50	1,20	0,60
25 %	67,50	72,90	64,40	59,60	16,70	4,10	4,50	4,20
50 %	71,45	76,10	70,60	64,00	21,45	5,20	7,00	6,50
Moyenne	71,21	76,55	69,63	63,01	23,30	5,96	8,85	8,05
75 %	75,33	79,70	74,90	67,60	26,43	6,80	12,10	10,40
Max	91,90	106,90	90,90	80,70	77,90	23,70	40,50	31,20

Groupe	Âge			
	1	2	3	4
Min	17,00	0,00	14,00	0,00
25 %	59,00	36,00	58,00	57,00
50 %	69,00	52,50	69,00	70,00
Moyenne	68,21	51,13	67,71	66,50
75 %	77,00	66,00	79,00	80,00
Max	101	97	97	104

TABLE 2.9 – Résumés des paramètres biologiques par classe résultant d’une classification k-means

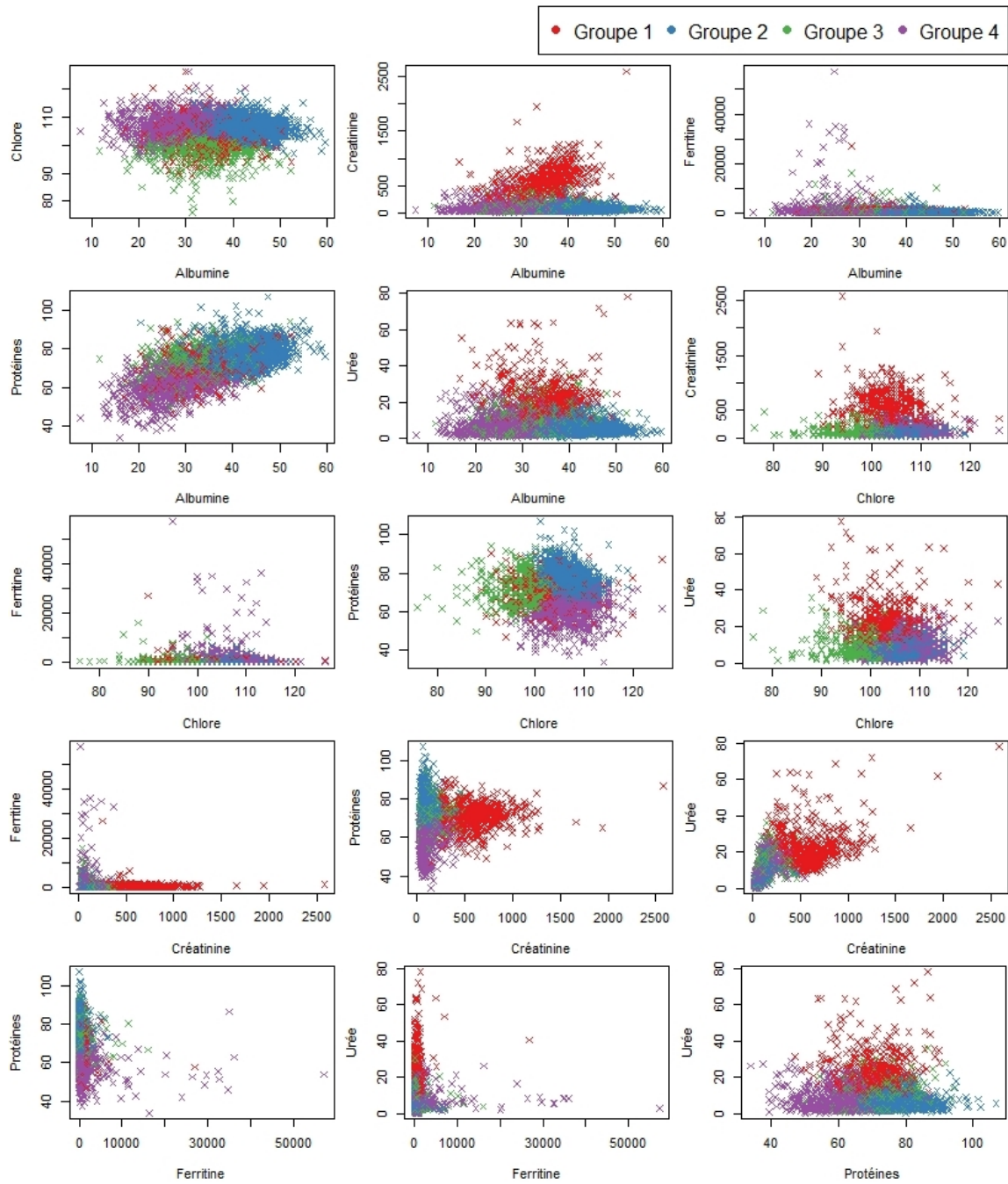


FIGURE 2.5 – Représentations des groupes résultants d’une classification k-means sur les graphiques biplot des paramètres biologiques

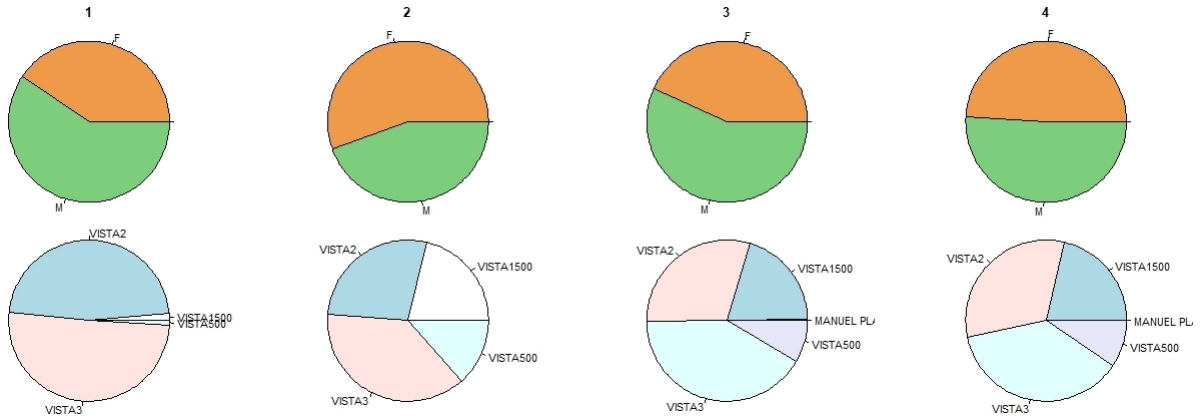


FIGURE 2.6 – Proportions des variables sexe et analyseurs dans les groupes résultant d'une classification k-means

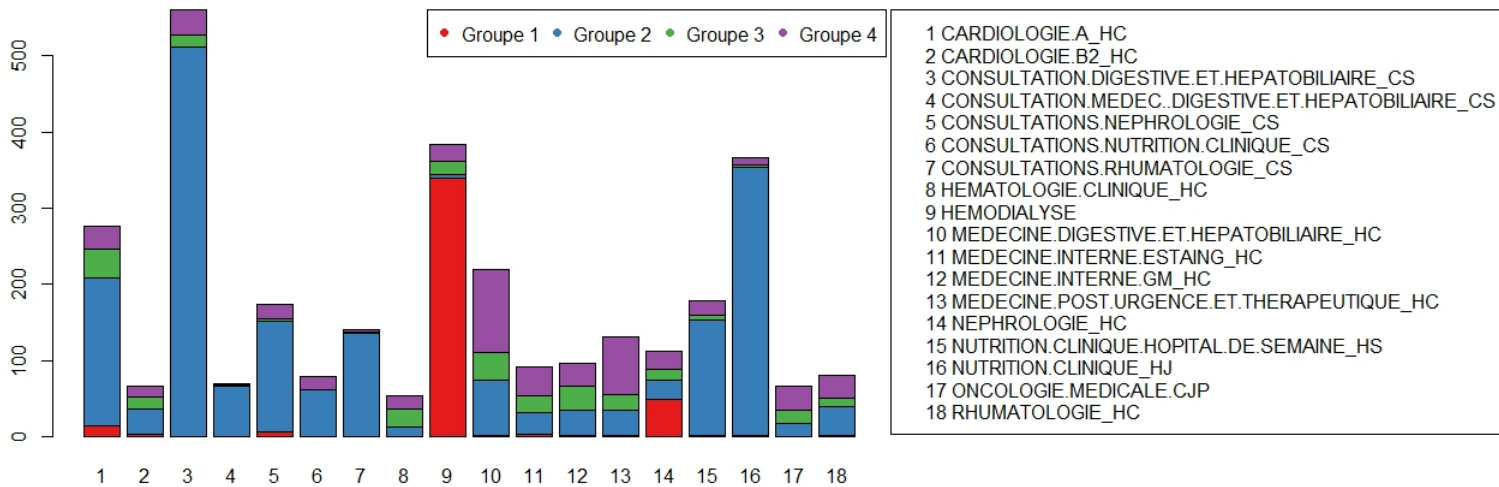


FIGURE 2.7 – Répartition des groupes dans les services ayant plus de 50 individus dans le jeu de données correspondant à la classification K-means

Le groupe 1 regroupe des individus ayant des valeurs d'urée et de créatinine élevées, ces mesures ont été réalisées principalement sur les analyseurs VISTA2 et VISTA3, donc à Gabriel Montpied. Ce groupe est essentiellement présent en hémodialyse mais on le retrouve aussi en cardiologie A HC, en consultation néphrologie CS et en néphrologie HC.

Le groupe 2 rassemble des individus avec des taux d'albumine et de protéines plus élevés et d'urée plus bas. La moyenne d'âge de ce groupe est aussi plus basse que les autres. Ce groupe est présent dans tous les services mais en grande majorité dans les services de consultations, de nutrition et en cardiologie A HC.

Le groupe 3 inclut les individus ayant des taux de ferritine très élevés et des taux de chlore bas. Peu présent dans les services on peut noter sa proportion très faible voire inexistante dans



les services tel que les consultations ou la nutrition.

Le groupe 4 rassemble des individus couplant des taux d'albumine et de protéines bas et des taux de chlore élevés. Cette population est présente dans tous les services, cependant très peu en consultation médecine digestive et hépatobiliaire et rhumatologie, ce groupe apparaît plus dans les services de médecine.

La figure 2.7 montre que les services de consultation ont le même profil que les services de nutrition, ceux de médecine ont aussi un profil similaire. Le service se démarquant est l'hémodialyse avec une population d'individus ayant des taux de créatinine et d'urée très haut.

Cette méthode peut être utilisée avec d'autres paramètres biologiques et donc d'autres services, elle permet de rechercher une segmentation des patients mais aussi de repérer les services ayant des profils similaires, pouvant ensuite être modéliser à l'aide d'une régression logistique ou d'un arbre de décision.

## 2.4 Détection de ruptures dans des séries temporelles

Les données étudiées sont des mesures réalisées en laboratoire, chaque mesure est associée à la date et l'heure à laquelle elle est réalisée. Ces données peuvent donc être étudiées sous forme de séries temporelles.

Une série temporelle ou chronologique est une série ordonnée de données  $(t_i, X_i)_{1 \leq i \leq n}$  où  $t_i$  est le temps et  $X_i$  un nombre. On représente souvent une série temporelle dans un repère où l'axe des abscisses représente le temps et l'axe des ordonnées les valeurs observées. Il est possible d'avoir des ruptures dans les séries. Il existe des méthodes statistiques de détection de ruptures, ces méthodes ont pour but d'estimer les instants où la série présente des changements dans la distribution (moyenne, variance, ...).

Un changement dans la distribution des données patients pourrait mettre en évidence un problème intervenu sur l'analyseur qui aurait impacté les analyses.

### 2.4.1 Méthodes

Les méthodes de détection de ruptures présentées ici sont des méthodes dites *offline*, c'est-à-dire que l'on dispose de l'ensemble des données. Il existe aussi des détections dites en temps réel ou *online* où les données arrivent progressivement.

Les trois premières méthodes présentées ici (Segmentation binaire, Segment Neighbourhood et PELT) sont issues du package développé pour le logiciel R : « Changepoint » [13] [12]. Ce package contient différentes méthodes de détection de ruptures *offline*.

## Introduction au package Changepoint

La détection de ruptures revient à identifier des points dans un ensemble de données où les propriétés statistiques changent. On suppose une séquence de données ordonnées  $y_{1:n} = (y_1, \dots, y_n)$ . Ce modèle aura  $m$  points de ruptures de positions :  $\tau_{1:m} = (\tau_1, \dots, \tau_m)$ . Chaque position est un entier compris entre 1 et  $n - 1$ . On définit  $\tau_0 = 0$  et  $\tau_m = n$  et on suppose que les points de ruptures sont ordonnés. Les points de ruptures fractionnent donc les données en  $m$  segments. Une approche couramment utilisée pour détecter plusieurs points de ruptures consiste à minimiser :

$$\sum_{i=1}^{m+1} [\mathcal{C}(y_{(\tau_{i-1}+1):\tau_i})] + \beta f(m). \quad (2.6)$$

Ici  $\mathcal{C}$  est une fonction de coût pour un segment et  $\beta f(m)$  est une pénalité. La fonction de coût la plus utilisée est la log-vraisemblance négative. Akaike's Information Criterion ( $\beta = 2p$ ) et Schwarz Information Criterion ( $\beta = p \log(n)$ ), avec  $p$  le nombre de paramètres supplémentaires introduits par l'ajout d'un point de rupture, sont deux exemples de pénalités. Les méthodes du package Changepoint permettent de détecter des points de ruptures en moyenne ou en variance.

## Segmentation Binaire

La méthode segmentation binaire consiste, tout d'abord, à appliquer la méthode du point unique à l'ensemble des données, c'est-à-dire à tester l'existence d'un  $\tau$  satisfaisant :

$$\mathcal{C}(y_{1:\tau}) + \mathcal{C}(y_{(\tau+j):n}) + \beta < \mathcal{C}(y_{1:n}) \quad (2.7)$$

Si aucun point n'est détecté, la méthode s'arrête et aucune rupture n'est observée. Si un point est détecté, les données sont séparées en deux segments : avant et après le point. La méthode de détection est alors appliquée sur les deux nouveaux segments. Le processus est répété jusqu'à ce qu'aucun point ne soit détecté. L'avantage de cette méthode est qu'elle est efficace en calcul ( $\mathcal{O}(n \log(n))$ ) [14]

## Segment Neighbourhood

Cette méthode explore l'ensemble de l'espace à l'aide d'une programmation dynamique. Elle commence par définir une limite supérieure de la taille de l'espace de segmentation (c'est-à-dire le nombre maximum de point de rupture) qui est noté  $Q$ . Ensuite, elle calcule la fonction coût pour tous les segments possibles. Toutes les segmentations possibles avec des points de ruptures compris entre 0 et  $Q$  sont considérées. C'est une méthode exacte et elle a la capacité d'inclure une pénalité arbitraire  $\beta f(m)$  mais elle est coûteuse en calcul ( $\mathcal{O}(Qn^2)$ ). [13] [14]

## PELT

L'algorithme PELT est similaire à celui de Segment Neighbourhood puisqu'il fournit une segmentation exacte mais il se montre plus efficace en temps de calcul. Le temps de calcul est contrôlé par l'hypothèse que les points de changements sont répartis dans l'ensemble des données plutôt qu'en une seule portion. . L'algorithme de programmation dynamique PELT permet de minimiser :

$$\sum_{i=1}^{m+1} [\mathcal{C}(y_{(\tau_{i-1}+1):\tau_i}) + \beta] \quad (2.8)$$

Cela correspond à l'équation 2.6 où  $f(m) = m$ . L'algorithme utilise une méthode de partitionnement optimal mais réalise un élagage des points de changements possibles ce qui permet de réduire le temps de calcul. La méthode de partitionnement optimale et l'algorithme PELT sont détaillés dans l'article [14].

## Test Pettitt

Le test pettitt est un test non paramétrique permettant de tester l'évolution de la tendance d'une série temporelle et ainsi de détecter des ruptures en moyenne. Ce test est réalisé via le package, développé pour le logiciel R, « Trend » [24]. Ce test est un test sur le rang dont la statistique du test est :

$$\hat{U} = \max_{1 \leq t \leq N} (|U_{t,N}|) \quad (2.9)$$

$$\text{avec } U_{t,N} = \sum_{i=1}^t \sum_{j=t+1}^N \text{signe}(x_i - x_j) \text{ et } \text{signe}(x_i - x_j) = \begin{cases} +1 & \text{si } x > 0 \\ 0 & \text{si } x = 0 \\ -1 & \text{si } x < 0 \end{cases} .$$

La probabilité pour que cette statistique soit supérieure à une certaine valeur  $k$  est donnée par :

$$P(U \geq k) = 2e^{\frac{-6k^2}{n^3+n^2}} \quad (2.10)$$

On rejette l'hypothèse nulle d'absence d'une rupture au seuil  $\alpha$  si  $P(U > \hat{U}) < \alpha$ . Dans le cas contraire, le point de rupture est estimé lorsque que  $U_{t,N}$  atteint son maximum.[23]

## BreackPoint

Cette méthode recherche les changements de structure dans les régressions linéaires à l'aide d'un algorithme de programmation dynamique.

On considère un modèle de régression standard :

$$y_i = x_i^t b_j + u_i \quad (2.11)$$

où, au temps  $i$ ,  $y_i$  est l'observation de la variable dépendante,  $x_i$  est un vecteur de régresseurs, avec le premier composant généralement égal à l'unité, et  $\beta_i$  le vecteur des coefficients de régression, qui peuvent varier au fil du temps.

La méthode teste l'hypothèse que les coefficients de régression restent constants :

$$H_0 : \beta_i = \beta_0 \quad (i = 1, \dots, n) \quad (2.12)$$

contre l'alternative qu'il existe au moins une variation des coefficients au cours du temps. S'il y a  $m$  points de changements, il existe  $m+1$  segments dans lesquels les coefficients de régressions sont constants et le modèle (2.11) peut être réécrit comme suit :

$$y_i = x_i' b_j + u_i \quad (2.13)$$

avec ( $i = i_{j-1} + 1, \dots, i_j, j = 1, \dots, m + 1$ ). où  $j$  est l'indice des segments et  $\mathcal{I} = \{i_1, \dots, i_m\}$  désigne l'ensemble des points de ruptures, par convention  $i_0 = 1$  et  $i_{m+1} = n$ .

Les points de ruptures sont estimés en minimisant la somme résiduelle des carrés (RSS) de l'équation ci-dessus pour toutes les segmentations possibles jusqu'au nombre de maximal de ruptures ainsi que le BIC associé :

$$RSS(i_1, \dots, i_m) = \sum_{j=1}^{m+1} r_{ss}(i_{j-1} + 1, i_k) \quad (2.14)$$

où  $r_{ss}(i_{j-1} + 1, i_j)$  est la somme résiduelle minimale habituelle des carrés du  $j$ -ième segment. Le problème est maintenant de trouver les points de ruptures ( $\hat{i}_1, \dots, \hat{i}_m$ ) qui minimisent la fonction objective :

$$(\hat{i}_1, \dots, \hat{i}_m) = \underset{(i_1, \dots, i_m)}{\operatorname{argmin}} RSS(i_1, \dots, i_m) \quad (2.15)$$

L'obtention des minimiseurs globaux en (2.15) par une recherche approfondie dans la grille serait de l'ordre de  $O(nm)$  et serait fastidieuse sur le plan du calcul pour  $m > 2$ . Ceux-ci peuvent être trouvés beaucoup plus facilement par une approche de programmation dynamique qui est de l'ordre  $O(n^2)$  pour un certain nombre de changements  $m$ . [34] [4] [35]

### 2.4.2 Recherche de ruptures

Les premières détections de ruptures ont été réalisées sans a priori, nous n'avions pas les dates des problèmes observés par le CHU et donc aucune information d'où se trouvaient des ruptures s'il y en avait.

Le but étant de repérer les dérives ou les ruptures des machines sur les données patients, les détections ont donc été réalisées par machine et par paramètre biologique séparément.

## Différences des échantillons

Les méthodes de détections sont réalisées sur différents échantillons. Un échantillon peut prendre en compte toutes les valeurs ou seulement les valeurs dites usuelles (valeurs comprises entre les bornes de référence (Tableau 2.1)). Les figures 2.8 et 2.9 représentent, en séries temporelles, des mesures d'Urée réalisées sur l'analyseur VISTA500 sur 5 jours, du 17 octobre 2016 au 21 octobre 2016, prenant en compte toutes les valeurs pour la première et seulement les valeurs usuelles pour la deuxième. Le fait de travailler seulement avec les mesures entrant dans les valeurs de références répond à l'hypothèse que ces données auraient des distributions plus stables puisqu'elles n'incluent pas les valeurs extrêmes qui pourraient être dues à certaines pathologies des patients.

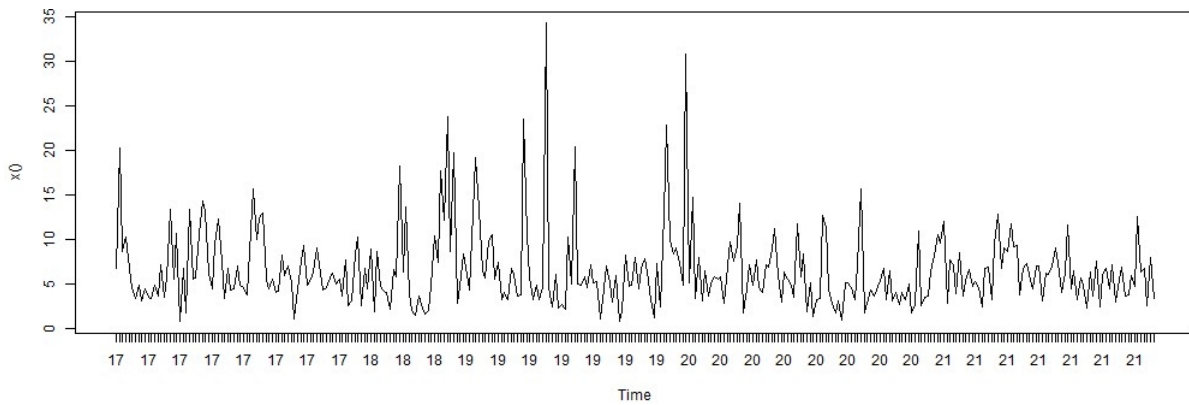


FIGURE 2.8 – Série temporelle : Urée, analyseur VISTA500, toutes valeurs du 17/10/16 au 21/10/16

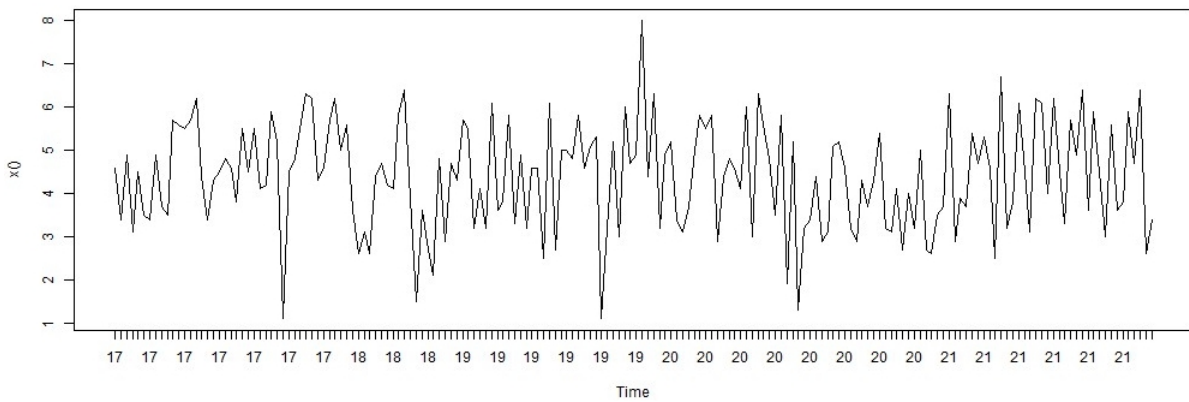


FIGURE 2.9 – Série temporelle : Urée, analyseur VISTA500, valeurs usuelles du 17/10/16 au 21/10/16

Les échantillons peuvent être composés seulement de femmes, d'hommes ou des deux, ou

encore ils peuvent inclure ou non les enfants. Ces différences peuvent être prises en compte si les distributions des paramètres biologiques varient selon l'âge ou le sexe des patients.

Par exemple, les valeurs usuelles de créatinine des enfants sont très différentes de celles des adultes, ce qui est observable sur les séries temporelles et peut entraîner de mauvaises détections de ruptures. Les figures 2.10 et 2.11 représentent, en série temporelle, des mesures de Créatinine entrant dans les valeurs de référence réalisées durant 16 jours, du 1er août 2016 au 16 août 2016, sur l'analyseur VISTA2 respectivement avec et sans enfant.

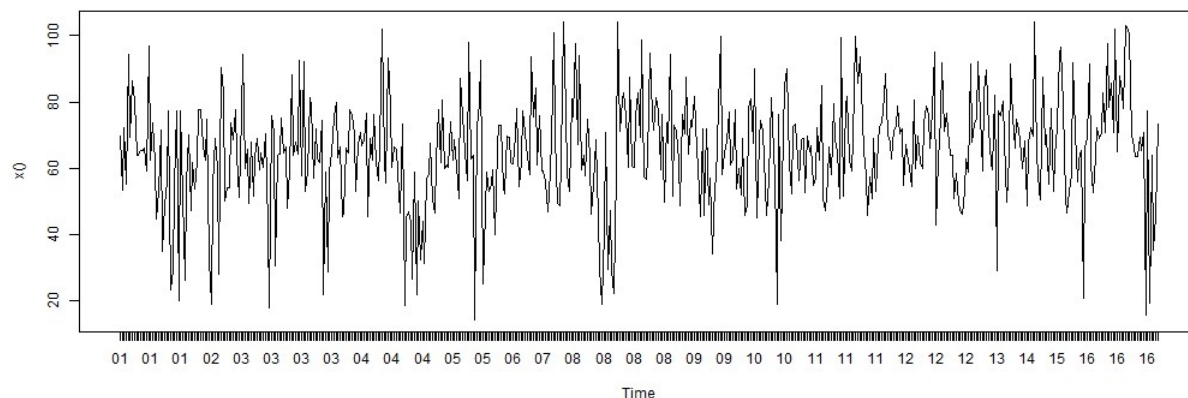


FIGURE 2.10 – Série temporelle : Créatinine, analyseur VISTA2, avec enfants, valeurs usuelles du 01/08/2016 au 16/10/16

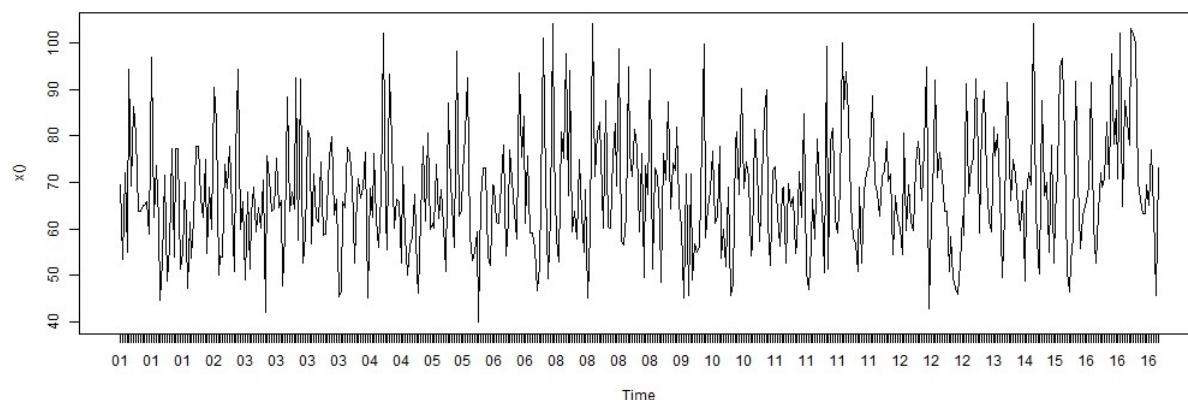


FIGURE 2.11 – Série temporelle : Créatinine, analyseur VISTA2, sans enfant, valeurs usuelles du 01/08/2016 au 16/10/16

Des détections de ruptures sont donc réalisées avec les différentes méthodes présentées au point 2.4.1 sur les différents échantillons possibles. On peut déjà noter qu'aucune rupture en variance n'a été trouvée sur les échantillons prenant uniquement les valeurs dites usuelles

(entrant dans les valeurs de référence).

### Exemple de détections de ruptures

Voici un exemple d'application des différentes méthodes présentées. L'échantillon étudié est constitué de cinq jours de mesures d'Albumine d'hommes et de femmes réalisées sur l'analyseur VISTA1500. Ce qui fait 158 mesures réalisées du 16/11/2018 au 20/11/2018 (figure 2.12).

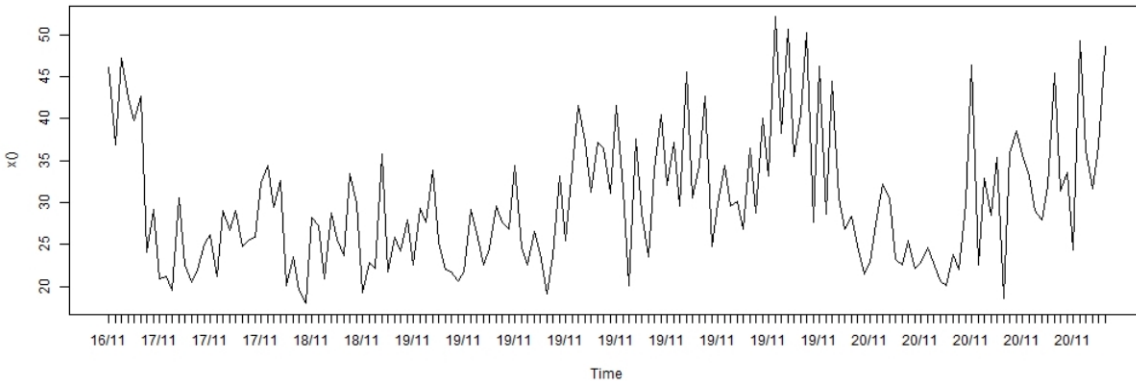


FIGURE 2.12 – Série temporelle : Albumine, analyseur VISTA1500 du 16/11/2018 au 20/11/2018

La première méthode utilisée est la segmentation binaire. Cette méthode permet de détecter des ruptures en variance ou en moyenne, ici aucune rupture en variance n'est détectée mais on trouve 5 ruptures en moyenne : les 16/11 à 18h29, 19/11 à 08h23, 11h12 et 18h50 et 20/11 07h30. Ces ruptures sont représentées sur la figure 2.13.

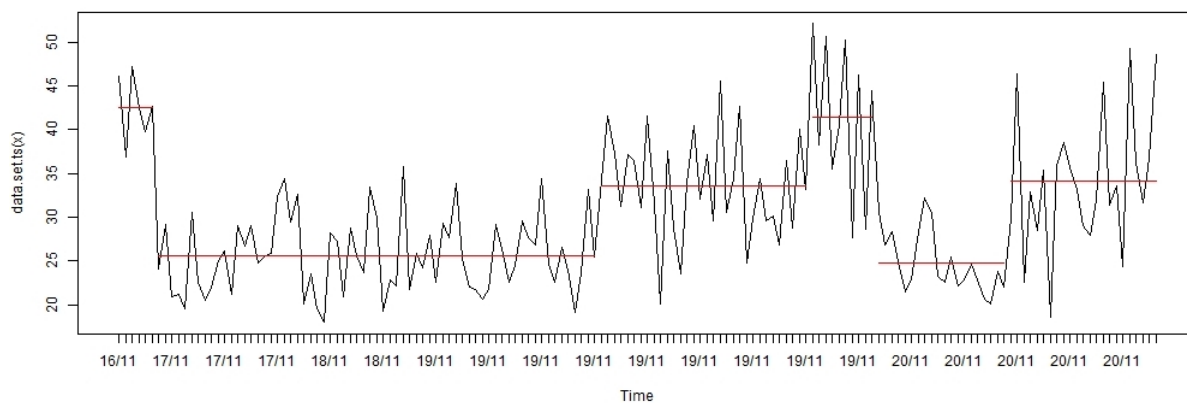


FIGURE 2.13 – Détection de ruptures par Segmentation Binaire sur série temporelle : Albumine, analyseur VISTA1500 du 16/11/2018 au 20/11/2018

La méthode suivante est la Segment Neighbourhood, cette méthode peut aussi détecter des ruptures en moyenne ou en variance mais n'a pas trouvé de rupture en variance. On retrouve 4 ruptures en moyenne : les 16/11 à 18h29, 19/11 08h58 et 18h58 et 20/11 07h30 (figure 2.14).

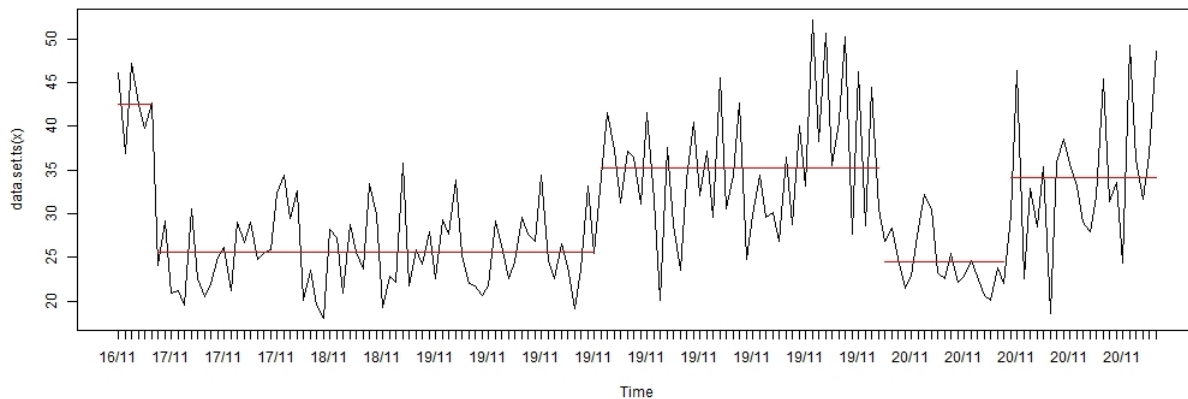


FIGURE 2.14 – Détection de ruptures par Segment Neighbourhood sur série temporelle : Albumine, analyseur VISTA1500 du 16/11/2018 au 20/11/2018

La méthode suivante est PELT. On retrouve 5 ruptures : les 16/11 à 18h29, 19/11 à 08h23, 11h12 et 18h50 et 20/11 07h30. (figure 2.15).

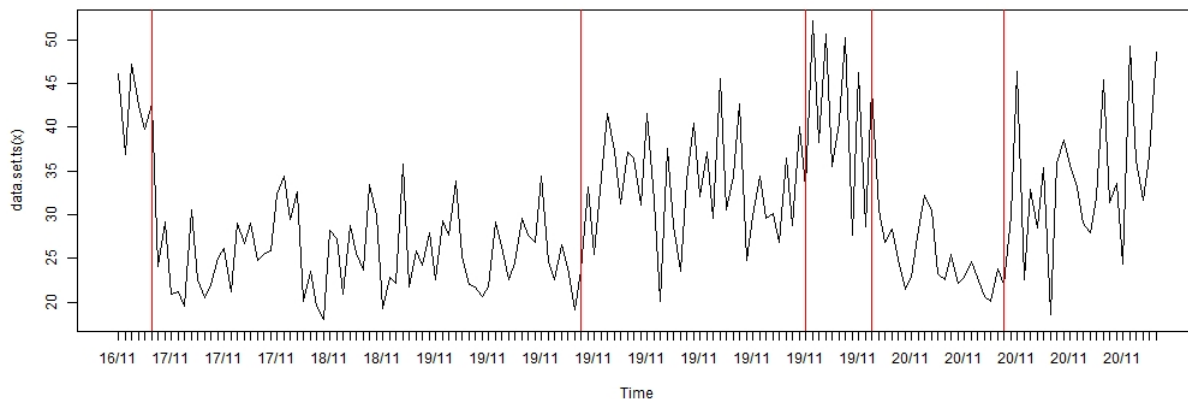


FIGURE 2.15 – Détection de ruptures par PELT sur série temporelle : Albumine, analyseur VISTA1500 du 16/11/2018 au 20/11/2018

La méthode suivante est le test Pettitt. Ce test ne peut trouver qu'un unique point de rupture par échantillon, il est donc utilisé sur une fenêtre glissante. Une taille de fenêtre est choisie autour d'un point, cette fenêtre se déplace autour de chaque point les uns après les autres et à chaque fenêtre on réalise un test de Pettitt. Trois ruptures sont ainsi détectées : les 16/11 à 18h29, 17/11 à 07h09 et 19/11 à 08h23 et 18h50 (figure 2.16).



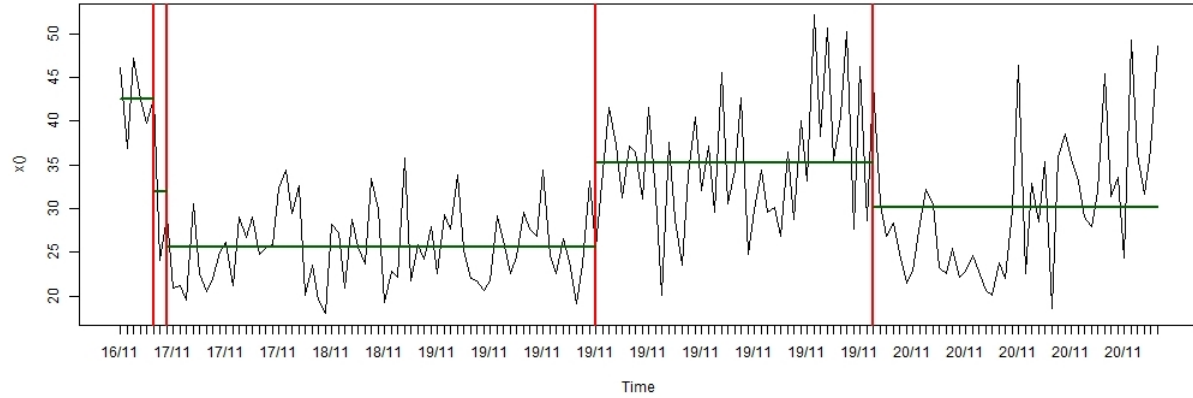


FIGURE 2.16 – Détection de ruptures par tests Pettitt sur série temporelle : Albumine, analyseur VISTA1500 du 16/11/2018 au 20/11/2018

La méthode suivante est Breackpoint. On retrouve 5 ruptures : les 16/11 à 18h29, 19/11 à 08h23, 11h12 et 18h50 et 20/11 07h30. Ces ruptures sont représentées sur la figure 2.17, les ruptures sont en rouge, les lignes vertes représentent les intervalles de confiance autour des ruptures.

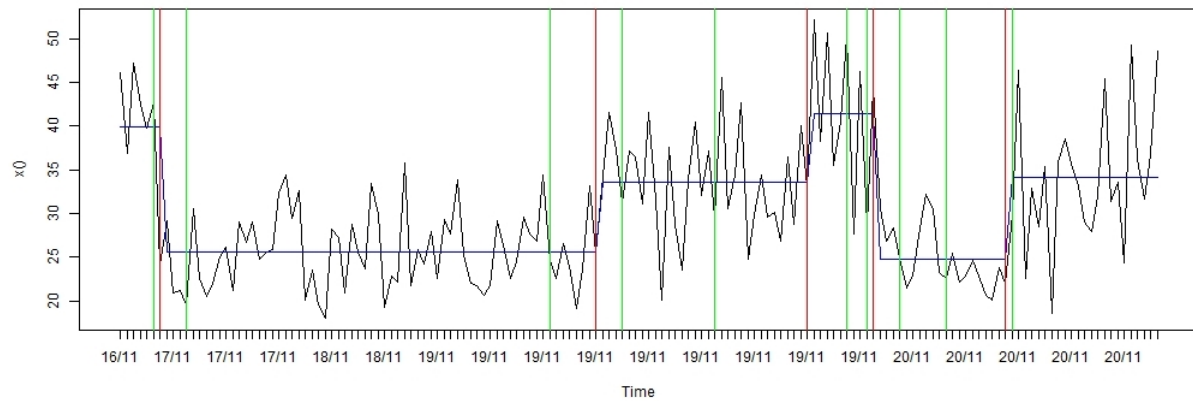


FIGURE 2.17 – Détection de ruptures par Breackpoint sur série temporelle : Albumine, analyseur VISTA1500 du 16/11/2018 au 20/11/2018

On peut ainsi voir que les différentes méthodes peuvent trouver différents points de ruptures.

### Interprétation des détections de ruptures

Les détections de ruptures réalisées sur l'exemple ci-dessus nous donnent 5 principales ruptures : 16/11 à 18h29, 19/11 à 08h23, 11h12 et 18h50 et 20/11 07h30 qui créent une

segmentation de l'échantillon. Si l'on regarde les services d'où proviennent les mesures, on peut remarquer que les ruptures mettent en évidence des changements de services où la distribution de l'albumine est très différente. En effet le premier segment est composé de 7 mesures issues de 7 services différents, mais le deuxième segment inclut, sur 68 mesures, 33 de réanimation adulte et 20 de surveillance continue, ce sont deux services où les taux d'albumine sont bas. Le troisième segment est composé en majorité de mesures de soin intensif d'hématologie clinique et d'hématologie clinique, le quatrième de mesures issues de 8 services différents, le cinquième comporte sur 21 mesures, 13 de réanimation et 6 de surveillance continue et le dernier segment contient des mesures de 14 services.

Les échantillons sanguins arrivent au laboratoire service par service ce qui explique les groupes que l'on vient de remarquer. L'albumine ayant des distributions différentes selon les services, les détections de ruptures sur les échantillons de ce paramètre mettent donc en évidence des changements de services, plus que de potentielles problèmes dus aux machines.

Certains paramètres comme le chlore ne varient pas ou très peu selon les services. Des détections de ruptures sur ces échantillons n'ont trouvé que très peu de ruptures contrairement aux échantillons de paramètres tel que l'albumine ou la créatinine.

## 2.5 Étude des problèmes observés par le CHU

### 2.5.1 Problèmes repérés lors des CIQ

Les problèmes observés par le CHU font suite aux CIQ réalisés plusieurs fois par jour pour chaque analyseur et pour tous les paramètres biologiques mesurés. Chaque problème est donné avec la date, le paramètre biologique et l'analyseur concernés.

En 2018, le CHU a pu observer des problèmes qui concernaient les paramètres biologiques que nous étudions, grâce aux contrôles internes de qualité pour 20 jours. 17 d'entre eux concernent l'albumine, 2 l'urée et 1 les protéines.

Pour rappel, si un problème est détecté lors d'un CIQ, les analyses des patients sont étudiées par une étude d'impact. Suite à cela, certaines analyses de patients peuvent être mesurées une nouvelle fois, dans le cas où le biologiste conclut qu'il y a eu impact sur interprétation, la mesure du patient est modifiée.

Pour trois des vingt problèmes observés, des analyses de patients ont dû être reprises mais aucun impact sur interprétation n'a été détecté et donc les données n'ont pas été modifiées.

Nous avons cherché à savoir si les problèmes observés lors des CIQ ont eu des conséquences sur les données des patients qui n'auraient pas été repérées lors des études d'impact ou qui n'auraient pas eu d'impact sur interprétation biologique.

Pour commencer à comparer les jours où les CIQ ont détecté un problème, nous avons calculé pour chaque jour par paramètre et par analyseur les variances, les moyennes et le pourcentage de valeurs n'entrant pas dans les valeurs de référence des paramètres biologiques. Mais aucune de ces nouvelles variables ne met en évidence les problèmes détectés par les CIQ.

Les distributions des paramètres biologiques de ces jours ont été comparées à une distribution de référence estimée sur le reste du jeu de données. Les comparaisons sont réalisées grâce à un test de Wilcoxon. Deux distributions sont différentes de leurs références, celle des protéines du 12/02/2018 de l'analyseur VISTA500 et celle de l'urée du 31/07/2018 de l'analyseur VISTA3.

## 2.5.2 Détections de ruptures et interprétations

Les données correspondant aux dates, analyseurs et paramètres biologiques sont maintenant étudiées sous forme de série temporelles. Les séries correspondent à plusieurs jours de mesures autour de la date de problème afin d'observer l'évolution du paramètres. Les différentes méthodes de détection de ruptures présentées au paragraphe 2.4.1 ont été utilisées. Pour les deux séries d'urée concernées par des problèmes, aucune rupture n'a été détectée avec les différentes méthodes utilisées. Pour les séries d'albumine, des ruptures ont été détectées à chaque fois, mais en regardant les services d'où provenaient les mesures, on peut émettre l'hypothèse que ces ruptures mettent en évidence des changements de services. La série temporelle de protéines montre aussi des ruptures.

Pour éviter la détection des changements de services, les données ont été centrées et réduites en fonction des services. Pour chacun d'eux nous calculons la moyenne et l'écart type des paramètres, cela a permis de centrer et de réduire les mesures en fonction de leur provenance. Les figures 2.18 et 2.19 donnent un exemple pour une série temporelle d'Albumine réalisée sur l'analyseur VISTA 1500 du 03/02/2018 au 09/02/2018.

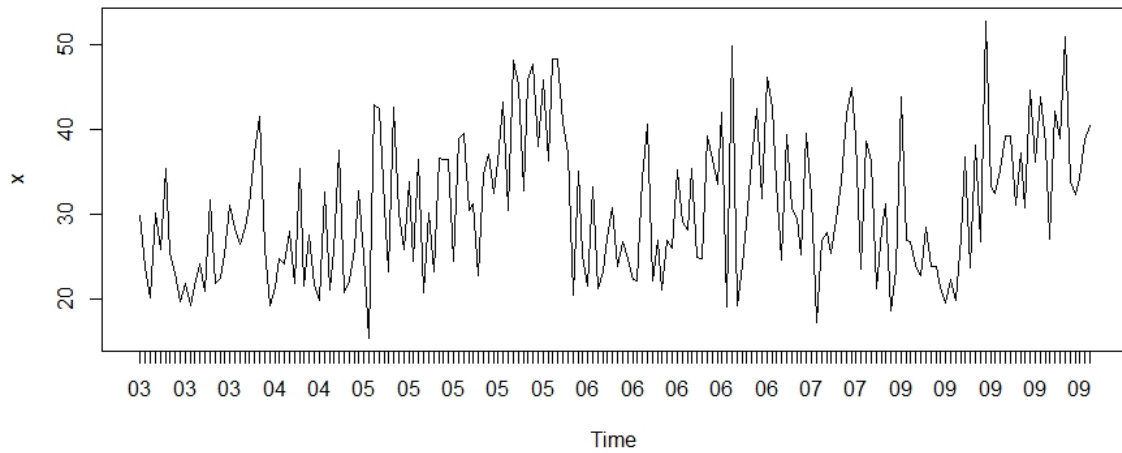


FIGURE 2.18 – Série temporelle : Albumine, VISTA1500, du 03/02/2018 au 09/02/2018

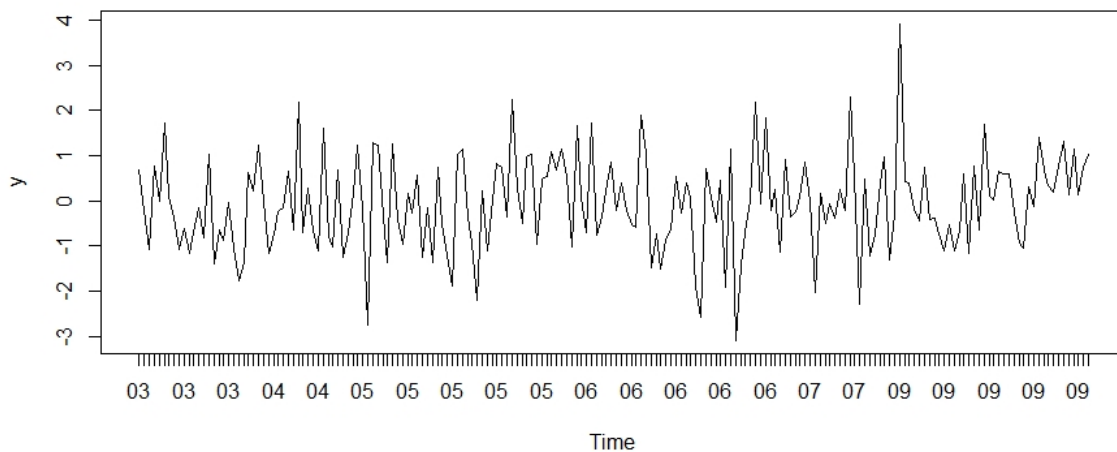


FIGURE 2.19 – Série temporelle : Albumine, centrée et réduite selon les services, VISTA1500, du 03/02/2018 au 09/02/2018

Suite à cela, de nouvelles détections de ruptures ont été réalisées. Aucune rupture n'a été détectée sur les nouvelles séries temporelles d'albumine et d'urée. Les CIQ ont mis en évidence un problème le 12/02/2018 sur l'analyseur V500 pour les mesures de protéines. Les détections de ruptures sur la série temporelle représentant les mesures de protéines du 10/02/2018 au 14/02/2018 ont localisé trois ruptures : le 12/02 à 06h31, à 09h55 et le 13/02 à 12h14, elles sont représentées sur la figure 2.20.

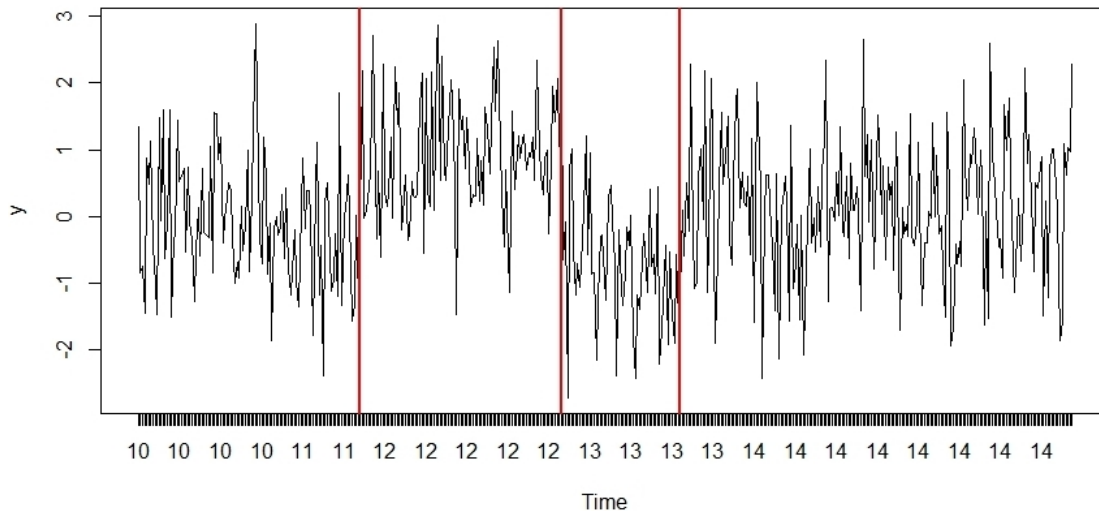


FIGURE 2.20 – Ruptures détectées sur la série temporelle : protéines, centrée et réduite selon les services, VISTA500, du 10/02/2018 au 14/02/2018

Les comparaisons de distribution et les détections de ruptures ont mis en évidence des comportements de données de patients différents lors de deux jours où le CHU a observé un problème grâce aux CIQ. Les mesures de Protéines du 12/02/2018 ont donc une distribution différente des autres jours ce qui se manifeste par des ruptures en séries temporelles, même lorsque les mesures sont centrées et réduites selon les services. Les mesures d'urée du 31/07/2018 ont une distribution différente des autres jours mais cela n'est pas visible lors d'une détection de ruptures sur les mesures centrées et réduites selon les services.

Pour les autres jours, d'après nos analyses, les problèmes observés lors des CIQ n'ont pas eu de répercussion sur les données des patients.

## 2.6 Simulations et détections de ruptures

Malheureusement nous n'avons pas pu identifier les conséquences que peuvent avoir les problèmes provenant des analyseurs sur les mesures des patients. Nous avons donc émis l'hypothèse que lorsqu'un analyseur dysfonctionnerait cela engendrerait un décalage, une imprécision ou une dérive traduits par un bruit gaussien ajouté aux mesures des patients avec une certaine moyenne ou variance, ou même une moyenne qui pourrait évoluer en fonction du temps.

Afin de tester les différentes méthodes utilisées, nous avons simulé des données ayant les mêmes distributions que les données patients et ajouté une rupture pour tenter de la détecter le plus tôt possible.

Dans cette partie, des méthodes de détection hors ligne sont utilisées mais aussi des

méthodes en ligne permettant de détecter les ruptures en temps réelle. Elles permettent de détecter les ruptures le plus tôt possible pour que le problème soit pris en charge rapidement et qu'il impacte le moins de données possible. .

### 2.6.1 Méthode de détection en ligne

Les méthodes de détection de ruptures dites séquentielles ou en ligne traitent les données au fur et à mesure de leurs acquisitions, par opposition aux méthodes offline présentées à la section 2.4.1. Dans ce contexte en temps réel, l'enjeu principal de la détection est de décider si un changement a eu lieu plutôt que de localiser la rupture. La décision est prise d'après la comparaison des dernières observations réalisées avec le début du signal.

La méthode présentée est issue du package "CPM" développé pour le logiciel R. [26] Elle est décrite en deux phases, la première permet de détecter une rupture sur un échantillon donné et la deuxième utilise la phase 1 séquentiellement.

#### Phase 1 :

Dans cette phase, on a une séquence d'observations de taille  $n : x_1, \dots, x_n$  qui contient ou non un point de changement. S'il n'y en a pas, les observations sont indépendamment et identiquement distribuées selon une certaine distribution  $F_0$ . S'il existe un point de changement à un moment  $k$ , alors avant cet instant les observations suivent une distribution  $F_0$  et après une distribution  $F_1$ , où  $F_0 \neq F_1$ .

La décision de présence d'un changement à un moment  $k$  conduit donc à choisir entre deux hypothèses :

$$\begin{aligned} H_0 : X_i &\sim F_0(x; \theta_0), \quad i = 1, \dots, n \\ H_1 : X_i &\sim \begin{cases} F_0(x; \theta_0) & i = 1, 2, \dots, k \\ F_1(x; \theta_1) & i = k + 1, \dots, n \end{cases} \end{aligned} \quad (2.16)$$

où  $\theta_i$  représente le potentiel paramètre inconnu de chaque distribution. Ce problème peut être résolu à l'aide d'un test d'hypothèse à deux échantillons. Le choix de la statistique de test se fait en fonction de ce que l'on suppose sur les échantillons et du type de changement qu'ils peuvent subir. Si l'on suppose des échantillons gaussiens on peut utiliser un test de Student, Bartlett ou le rapport de vraisemblance généralisée (GLR). Des tests non paramétriques peuvent aussi être utilisés tels que : Mann-Whitney, Mood, Kolmogorov-Smirnoff, Lepage, ou Cramer-Von-Mises. Les tests de Bartlett et Mood seront plus sensibles à des changements en variance qu'en moyenne.

Après avoir choisi le test, la valeur du test  $\tilde{D}_{n,k}$  est calculée pour tous les moments  $k$  de la série temporelle donnée. La statistique finale est alors :

$$D_n = \max_{k=2, \dots, n-1} D_{k,n} = \max_{k=2, \dots, n-1} \left| \frac{\tilde{D}_{k,n} - \mu_{\tilde{D}_{k,n}}}{\sigma_{\tilde{D}_{k,n}}} \right| \quad (2.17)$$

L'hypothèse nulle est alors rejetée si  $D_n > h_n$ .  $h_n$  est un seuil choisi, lié au taux de faux positifs  $\alpha$  (probabilité de détecter un changement alors qu'il n'y en a pas). Il convient de choisir  $h_n$  comme le quantile supérieur  $\alpha$  de la distribution de  $D_n$  sous l'hypothèse nulle.

Finalement la meilleure estimation du point de changement sera la valeur  $k$  qui maximise  $D_{k,n}$  :

$$\hat{\mathcal{T}} = \arg \max_k D_{k,n} \quad (2.18)$$

### Phase II :

L'approche décrite à la phase I peut être utilisée de manière séquentielle lorsque les observations sont reçues au fur et à mesure. Soit  $x_t$  la  $t$ -ième observation,  $t = 1, 2, \dots$ . Pour chaque observation reçue  $x_t$ , l'approche CPM traite  $x_1, \dots, x_t$  comme étant une séquence de longueur fixe et calcul  $D_t$  associé en utilisant la méthode expliquée ci-dessus. Un changement est détecté si  $D_t > h_t$ , dans ce cas, l'estimation du point de changement  $\tilde{\mathcal{T}}$  est donnée ainsi que le temps de détection :  $t - \tilde{\mathcal{T}}$ , sinon l'observation suivante  $x_{t+1}$  est reçue et  $D_{t+1}$  est calculé et ainsi de suite.

Dans ce cadre séquentiel,  $h_t$  est choisi de sorte que la probabilité d'encourir une erreur de type 1 est constante au fil du temps.

$$\begin{aligned} P(D_1 > h_1) &= \alpha \\ P(D_t > h_t | D_{t-1} \leq h_{t-1}, \dots, D_1 \leq h_1) &= \alpha, \quad t > 1 \end{aligned} \quad (2.19)$$

Dans ce cas, en supposant qu'aucun changement n'est lieu, le nombre moyen d'observations reçues avant qu'une fausse détection positive se produise est égal à  $\frac{1}{\alpha}$ . Cette quantité est appelée longueur moyenne d'exécution, *ARL0*. En général, la distribution conditionnelle dans l'équation 2.19 est intraitable sur le plan analytique et la simulation de Monte Carlo est utilisée pour calculer les séquences requises de valeurs de  $h_t$  correspondant à un choix donné de  $\alpha$ . Il s'agit d'une procédure de calcul coûteuse, mais elle n'a besoin d'être effectuée qu'une seule fois. Le package CPM contient des séquences précompilées de seuils qui correspondent à une variété de choix de  $\alpha$ . [27]

### **2.6.2 Simulation des données**

Les données des patients permettent d'estimer les fonctions de répartition des lois des paramètres biologiques. Nous utiliserons la méthode d'inversion des fonctions de répartition pour simuler les données :

Soit  $X$  une variable aléatoire réelle de fonction de répartition  $F$ . Posons,  $0 \leq t \leq 1$ ,

$$F^{-1}(t) = \inf\{x, F(x) \geq t\} \quad (2.20)$$

Alors, si  $U \sim \mathcal{U}([0; 1])$  alors  $F^{-1}(U)$  a la même loi que  $X$ .

Une fonction de répartition du paramètre biologique étudié est estimé pour chaque service.

Soit  $(X_n, S_n)_{n \in \mathbb{N}}$  une chaîne de Markov donnant un échantillon de données simulées de valeurs du paramètre biologique étudié et des services.

Les valeurs du paramètre biologique dépendent du service du patient, nous simulons donc les lois conditionnelles de  $X_n$  sachant  $S_n$

$$P(X_{i+1} = x', S_{i+1} = s' | X_i = x, S_i = s) = P(X_{i+1} = x' | S_{i+1} = s')P(S_{i+1} = s' | S_i = s) \quad (2.21)$$

La matrice de transition des services donne les probabilités des valeurs de  $S_{k+1}$  lorsque l'on connaît celle de  $S_k$  :

$$P_{x,y} = P(S_{i+1} = y | S_i = x) = \frac{\text{Nombre de transition } x \rightarrow y}{\text{Nombre de transition total}} \quad (2.22)$$

Ce qui signifie que le service  $S_{i+1}$  est simulé selon les probabilités de transition de  $s_i$ , données par la matrice de transition des services. Une fois que le service est connu,  $x_{i+1}$  est simulé grâce à la méthode d'inversion des fonctions de répartition.  $s_1$  est une valeur suivant une loi discrète ayant comme probabilités les fréquences des services dans les données réelles. La matrice de transition permet de traduire l'arrivée des échantillons des patients service par service.

Les simulations réalisées sont basées sur les 25 services ayant demandés le plus de mesures du paramètre biologique étudié. La figure 2.21 donne un exemple de simulation d'un échantillon de 200 individus d'Albumine sur l'analyseur VISTA 500 et la figure 2.22 donne un échantillon réel de 200 individus d'Albumine sur l'analyseur VISTA 500

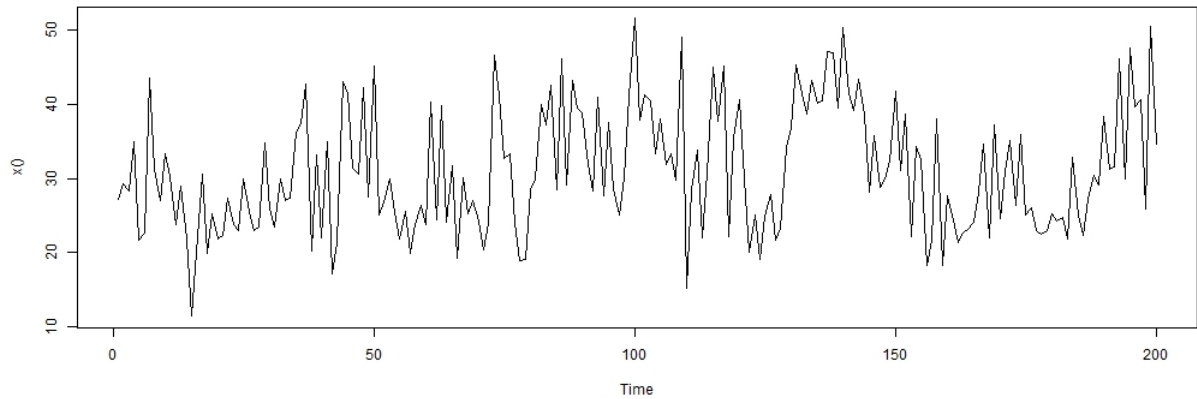


FIGURE 2.21 – Exemple d'une simulation d'un échantillon de 200 valeurs d'albumine ; analyseur VISTA500



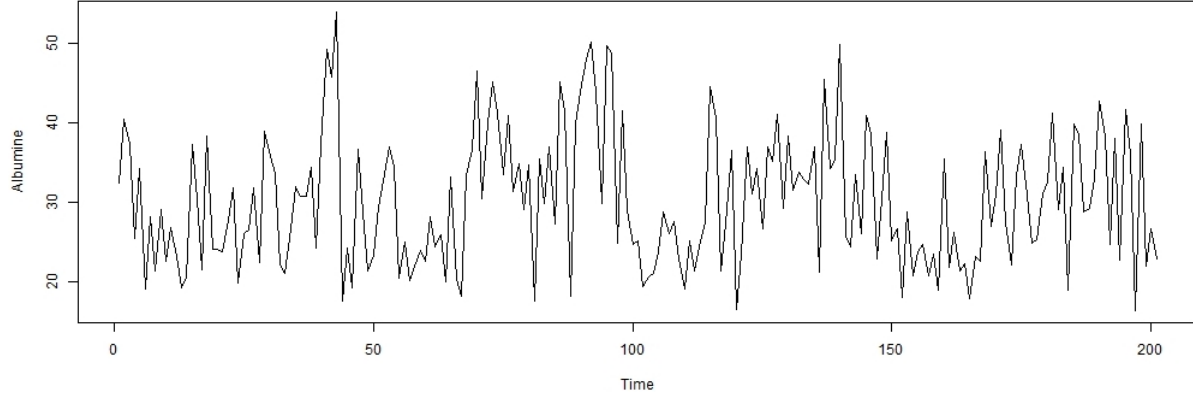


FIGURE 2.22 – Exemple d’un réel échantillon d’Albumine de 200 individus; analyseur VISTA500

### 2.6.3 Création de ruptures

Pour créer une rupture sur les données, nous avons ajouté un bruit gaussien de moyenne  $\mu$  et d’écart-type  $\sigma$  aux simulations  $X_n$ , qui seront visualisées et étudiées comme des séries temporelles  $(t_i, X_i)_{1 \leq i \leq n}$ . Ces séries sont de 200 individus, le bruit est ajouté au temps  $i = 70$  et sera de longueur 50. Trois types de bruit sont testés :  $\mathcal{N}(\mu, 1)$ ,  $\mathcal{N}(0, \sigma)$  ou  $\mathcal{N}(0.5T, 1)$  où  $T$  représente le temps de la rupture, donc  $T = 1, \dots, 50$ , représentant respectivement un décalage, une imprécision ou une dérive de l’analyseur. Pour que les ruptures simulées soient représentatives des écarts considérés comme de réels impacts biologiques,  $\mu$  et  $\sigma$  prendront les valeurs les plus basses de Delta-Check et TCC (tableau 2.10) pour chaque paramètre biologique.

	Delta-Check			TCC		
	Niveau 1	Niveau2	Niveau 3	Niveau 1	Niveau2	Niveau 3
Albumine	8.8	8.5	8.8	12.47	12.25	12.47
Chlore	2.4	2.5	2.8	4.09	4.16	4.37
Créatinine	10.5	6.3	4.3	19.50	17.65	17.10
Ferritine	7.6	8.0	9.0	40.07	40.15	40.35
Protéines	4.1	3.7	3.5	8.63	8.46	8.37
PSA	8.5	6.8	7.3	50.87	50.62	50.69
Urée	9.4	7.2	6.9	34.78	34.31	34.24

TABLE 2.10 – Delta-Check et TCC des trois niveaux de CIQ pour les sept paramètres biologiques étudiés

La figure 2.23 donne l’exemple d’une simulation d’Albumine sur l’analyseur VISTA2. Les figures 2.24, 2.25 et 2.26 donnent la même simulation mais avec l’ajout d’un bruit gaussien, respectivement,  $\mathcal{N}(8.6; 1)$ ,  $\mathcal{N}(0; 8.6)$  et  $\mathcal{N}(0.5T; 1)$ .

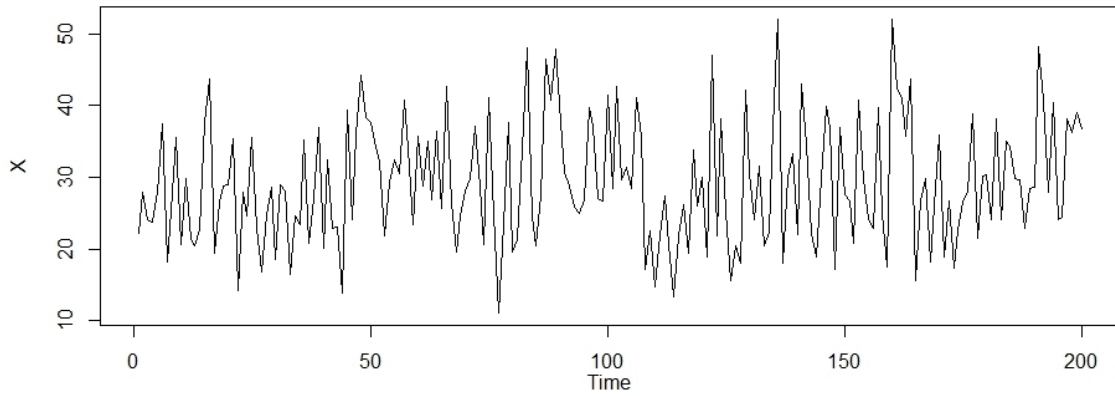


FIGURE 2.23 – Simulation d'un échantillon d'Albumine sur analyseur VISTA 2

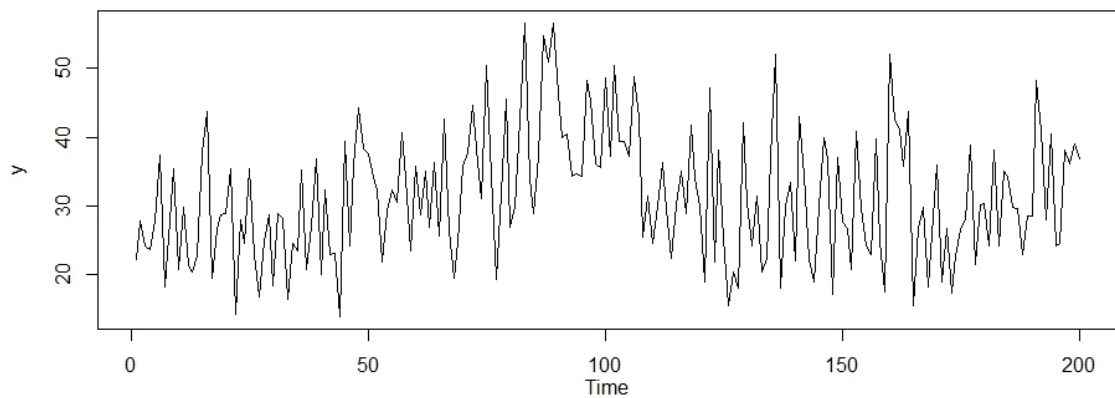


FIGURE 2.24 – Simulation d'un échantillon d'Albumine sur analyseur VISTA 2 avec ajout d'un bruit gaussien  $\mathcal{N}(8.6; 1)$  de taille 50 à  $i=70$

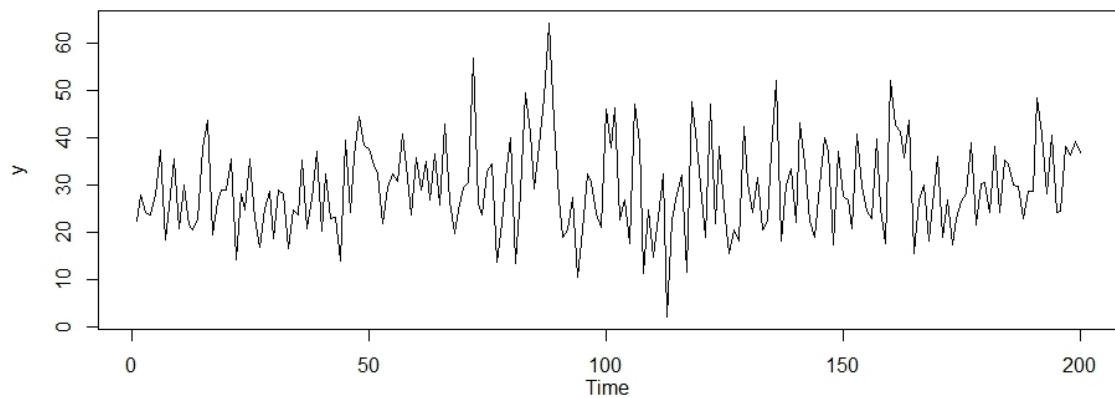


FIGURE 2.25 – Simulation d'un échantillon d'Albumine sur analyseur VISTA 2 avec ajout d'un bruit gaussien  $\mathcal{N}(0; 8.6)$  de taille 50 à  $i=70$

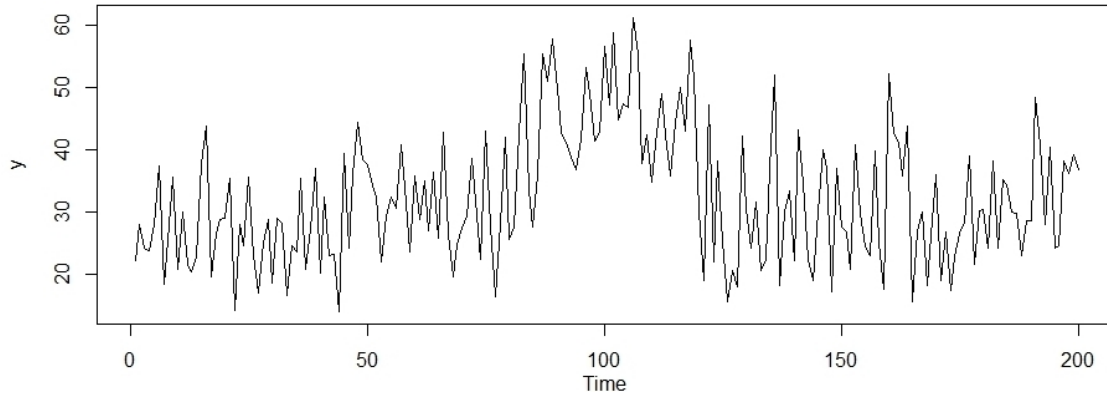


FIGURE 2.26 – Simulation d’un échantillon d’Albumine sur analyseur VISTA 2 avec ajout d’un bruit gaussien  $\mathcal{N}(0.5T; 1)$  de taille 50 à  $i=70$

#### 2.6.4 Méthode de détection de ruptures avec l’information service

Dans cette méthode nous émettons l’hypothèse que l’on connaît les services de provenance des patients. Avant de réaliser les détections de ruptures, les données sont centrées et réduites en fonction du service, donc avec les moyennes et les écarts-types de chacun. Cela permet de perdre l’effet des services et de ne plus détecter leurs changements comme nous avons déjà observé.

Les détections de ruptures hors ligne sont utilisées sur l’ensemble de l’échantillon. Les méthodes sont : segmentation binaire en moyenne et variance (BinMoy et BinVar), segment neighbourhood en moyenne et variance (SegMoy et SegVar), PELT, pettitt (Pett) et break-point (Br).

Les détections de ruptures en ligne analysent l’échantillon valeur par valeur, les 20 premières mesures sont considérées comme l’échantillon de référence, le premier test est donc réalisé à la 20ème mesure. Les tests statistiques utilisés sont : Student (Stu), Bartlett (Bar), rapport de vraisemblance généralisé (GLR), Mann-Whitney (MW), Mood (Moo), Lepage (LP), Kolmogorov-Smirnov (KS) et Cramer-Von-Mises (CVM).

Pour chaque rupture testée, 100 simulations et détections sont réalisées. Pour les méthodes hors ligne, on regarde la première rupture détectée et la dernière, les méthodes segmentation binaire et segment neighbourhood sont forcées de trouver deux ruptures. Pour les méthodes en ligne, on regarde la rupture détectée et le temps qu’il a fallu pour la trouver.

Suite aux 100 simulations, la rupture débutant à  $i = 70$ , nous regardons le pourcentage de simulations ayant trouvées la première rupture entre  $i = 60$  et  $i = 80$ , puis entre 60 et 90 ainsi que 60 et 100. On calcule aussi le pourcentage de simulations ayant trouvées la dernière rupture entre 110 et 130 pour les méthodes hors ligne et le temps de détection moyen pour

les méthodes en ligne. Ces informations vont nous permettre de voir quelle méthode détecte le mieux la rupture et quelle est la plus rapide.

Pour de meilleurs résultats et moins de fausses alarmes lors des détections de ruptures, les valeurs trop extrêmes et rares de trois paramètres biologiques ont été retirées des échantillons. Ainsi pour la créatinine les valeurs au-dessus de 250 sont retirées, pour la ferritine ce sont les valeurs à partir de 800 et celles au-dessus de 40 pour l'urée.

## Résultats

Les tableaux 2.11, 2.12 et 2.13 donnent les résultats de toutes les méthodes pour 100 simulations d'Albumine sur VISTA 2, respectivement, pour les bruits  $\mathcal{N}(\mu, 1)$ ,  $\mathcal{N}(0, \sigma)$  et  $\mathcal{N}(0.5T, 1)$  avec  $\mu$  et  $\sigma$  prenant comme valeurs 8.6 et 12.25.

Les méthodes segment Neighbourhood en moyenne et breackpoint détectent très bien les ruptures en moyenne et la meilleure méthode en ligne pour ces ruptures est le test Lepage. On remarque tout de même que dans environ 10% des simulations, les méthodes en ligne ne détectent pas la bonne rupture.

Les ruptures en variance sont plus compliquées à détecter, les méthodes donnent de meilleurs résultats quand  $\sigma$  vaut 12.25. La meilleure méthode hors ligne est segment neighbourhood et Mood pour les méthodes en ligne.

Au vu des résultats des détections des ruptures en pente, il est compliqué de détecter le début de la rupture mais entre 60 et 100, le bruit est devenu assez important pour le détecter à 100% en hors ligne et à 90% en ligne. On peut remarquer qu'ici les temps moyens de détection sont plus grand qu'auparavant.

		N( 8,6 ; 1)				N (12,25 ; 1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	87	87	87	94	100	100	100	97
Ligne	BinVar	6	6	6	1	39	41	41	37
	SegMoy	99	99	99	100	100	100	100	100
	SegVar	16	18	19	10	61	63	65	57
	Pelt	92	92	92	94	94	94	94	95
	Pett	78	78	78	91	95	95	95	97
	Br	97	97	97	98	99	99	99	100
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	86	86	86	7.01	90	90	90	4.2
	Bar	14	16	21	14.95	22	26	27	22.41
	GLR	83	84	84	9.92	84	84	84	7.15
	MW	84	84	84	9.27	90	90	90	6.95
	Moo	25	26	28	4.26	68	72	76	9.74
	LP	91	91	91	8.87	92	92	92	5.07
	KS	90	90	90	7.92	86	86	86	6.99
	CVM	87	87	87	8.71	90	90	90	6.69

TABLE 2.11 – Résultats des détections de ruptures sur 100 simulations d'Albumine, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec l'information service

		N(0 ; 8,6)				N (0 ; 12,25)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	0	1	1	1	2	3	4	1
Ligne	BinVar	55	59	60	59	90	91	91	87
	SegMoy	12	27	46	16	22	44	68	21
	SegVar	67	72	74	68	95	97	97	90
	Pelt	54	66	73	55	89	94	96	85
	Pett	11	13	20	7	10	17	24	13
	Br	4	4	5	1	8	14	16	6
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	46	54	70	4.2	64	77	81	3.24
	Bar	62	68	71	14.81	69	71	71	9.42
	GLR	51	56	60	16.03	75	77	78	12.47
	MW	18	24	30	5.75	17	24	29	5.77
	Moo	59	74	78	10.45	81	84	88	8.43
	LP	46	65	70	7.33	77	86	90	8.83
	KS	21	30	35	7.62	31	39	46	9
	CVM	19	26	34	5.87	26	38	41	6.61

TABLE 2.12 – Résultats des détections de ruptures sur 100 simulations d'Albumine, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$ ; méthode avec l'information service

Méthode		$\mathcal{N}(0,5T; 1)$			
		%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	46	85	97	96
Ligne	BinVar	16	57	95	99
	SegMoy	17	76	100	100
	SegVar	5	42	94	99
	Pelt	63	88	89	96
	Pett	86	90	91	93
	Br	61	95	96	98

Méthode		%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	68	90	90	8.59
	Bar	26	32	35	38.89
	GLR	59	82	82	12.11
	MW	71	83	83	12.93
	Moo	17	52	79	8.44
	LP	61	89	89	8.55
	KS	60	81	81	10.38
	CVM	68	81	81	12.58

TABLE 2.13 – Résultats des détections de ruptures sur 100 simulations d’Albumine, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$ ; méthode avec l’information service

Les résultats des paramètres biologiques chlore, créatinine ferritine, protéines, PSA et urée sont dans l’annexe A.4. Malgré un TCC grand, les ruptures en moyenne, variance ou pente de ferritine ne sont pas bien détectées. La ferritine peut atteindre de grandes valeurs et sa variation est importante même à l’intérieur des services, ce qui peut expliquer ces mauvaises détections. À l’inverse, on a de très bons résultats pour les trois types de ruptures, que se soit pour le Delta-Check ou le TCC, pour l’Urée et les PSA.

En général, les ruptures en moyenne avec le TCC sont bien détectées mais ce n’est pas toujours le cas avec le Delta-Check et les ruptures en variance sont rarement détectées. La détection des pentes demande plus de temps et varient selon les échelles des paramètres biologiques. En effet, pour les paramètres avec une plus grande variance, les bruits en pente à  $0.5T$  auront moins d’impact et seront par conséquence détectables beaucoup plus tard que pour les paramètres biologiques avec une variance plus petite.

Afin d’évaluer à partir de quel  $\mu$  ou  $\sigma$  les méthodes détectent bien les différents types de bruit, des détections de ruptures ont été réalisées sur des simulations d’Albumine avec un bruit  $\mathcal{N}(\mu; 1)$  et  $\mathcal{N}(0; \sigma)$ , pour  $\mu$  variant de 0 à 15 et  $\sigma$  variant de 5 à 20. Les figures 2.27 et 2.28 donnent les pourcentages de simulations où les méthodes en ligne ont trouvé la première rupture entre  $i = 60$  et  $i = 80$ , ainsi que le temps de détection moyen. Pour chaque  $\mu$  ou  $\sigma$  différent, le pourcentage et le temps de détection moyen ont été calculés sur 50 simulations.

Ces mêmes résultats sont données pour les paramètres biologiques : chlore, créatinine, ferritine, protéines, PSA et urée dans l’annexe A.4.

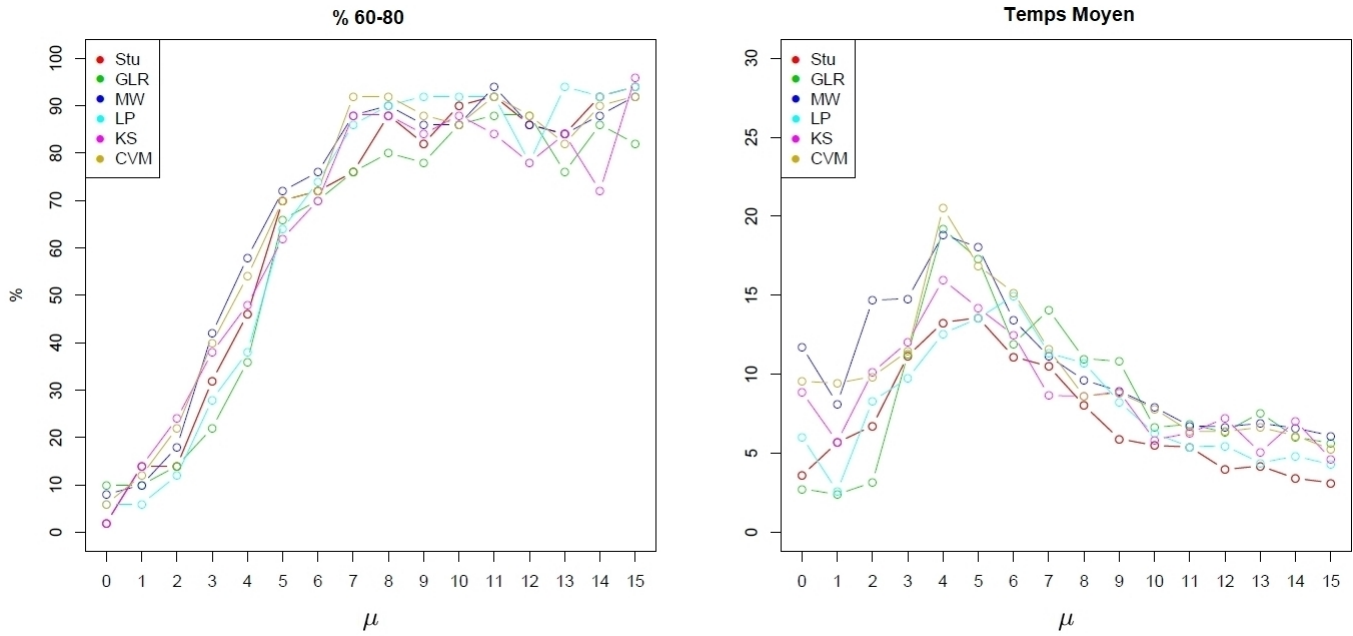


FIGURE 2.27 – Pourcentages de simulations d’Albumine avec ajout d’un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 0 à 15 ; Méthode avec l’information service

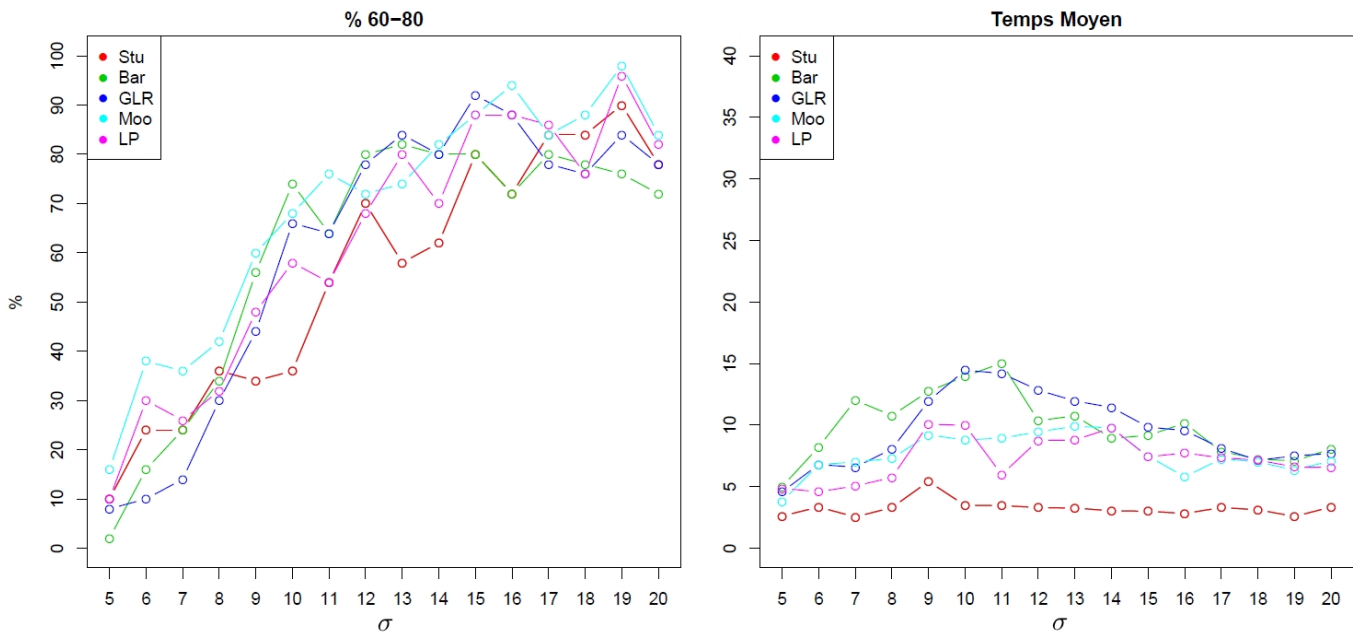


FIGURE 2.28 – Pourcentages de simulations d’Albumine avec ajout d’un bruit  $\mathcal{N}(0; \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 5 à 50 Méthode avec l’information service

## Fausse alarmes

Nous nous intéressons maintenant aux fausses alarmes des méthodes en ligne. Si une méthode détecte bien les ruptures, elle doit aussi détecter le moins possible de fausses ruptures et donc de fausses alarmes en pratique dans les laboratoires.

Afin d'évaluer les détections de fausses ruptures, nous avons simulé, sans ajout de rupture, 200 échantillons d'une taille qui correspond aux nombres des mesures réalisées par paramètre biologique sur une journée. Les détections de ruptures en ligne ont été appliquées sur ces échantillons mais sans interruption, à chaque rupture détectée, la méthode se réinitialise à la mesure suivante. Cela permet d'évaluer le nombre de fausses alarmes que peuvent signaler les méthodes sur une journée de mesure au laboratoire.

Le tableau 2.14 donne les pourcentages de simulations ayant détectées 0, 1, 2 ou 3 ruptures pour 200 simulations d'albumine de 60 individus. On peut donc dire que dans environ 90% des cas, il n'y aurait pas de fausses alarmes lors d'une journée de mesures d'albumine en laboratoire et que s'il y en a, elles seraient de 1, 2 ou 3 par jour.

Méthode	0	1	2	3
Stu	88,5	9,5	2	0
Bar	89	4	7	0
GLR	89,5	3	7,5	0
MW	90,5	7,5	2	0
Moo	92	7	1	0
LP	91,5	7,5	1	0
KS	91	8	0,5	0,5
CVM	91	7	2	0

TABLE 2.14 – Pourcentages de simulations où les méthodes de détection en ligne en continue (avec services) ont trouvé 0, 1, 2 ou 3 ruptures sur 200 simulations d'albumine de 60 individus

Les résultats pour les paramètres chlore, créatinine, ferritine, protéines, PSA et urée sont donnés dans l'annexe A.4.

### 2.6.5 Méthodes de détection de ruptures sans l'information service

Dans cette partie nous émettons l'hypothèse que nous n'avons pas l'information service, pour pallier à cela des classifications sont couplées aux détections de ruptures.

Avant de réaliser les détections de ruptures hors ligne, on effectue une classification de l'échantillon de données. Chaque groupe créé donne ainsi une nouvelle série temporelle et ce sont sur ces séries que l'on applique les détections de ruptures.

Pour les détections en ligne, on réalise une première classification sur l'échantillon référent qui correspondra aux 20 premiers individus. Des détections de ruptures en ligne sont effectuées



sur chaque groupe représenté sous forme d'une série temporelle. Si une rupture est trouvée (la plus tôt en cas de multiples ruptures), alors l'algorithme est arrêté et on peut donner le moment de rupture et le temps de détection. Sinon, on remet à jour la classification avec les deux individus suivant en plus, on relance les détections de ruptures et ainsi de suite. L'algorithme s'arrête lorsqu'une rupture est détectée sur une série temporelle ou lorsque l'on arrive à la fin de l'échantillon.

Les mêmes méthodes de détection de ruptures hors ligne et en ligne que données au paragraphe 2.6.4 sont utilisées.

Pour chaque rupture testée, 100 simulations et détections sont réalisées. Pour les détections hors ligne, on regarde la première rupture détectée et la dernière, les méthodes segmentation binaire et segment neighbourhood sont forcées de trouver deux ruptures. Pour les méthodes en ligne, on regarde la rupture détectée et le temps qu'il a fallu pour la détecter.

Suite aux 100 simulations, la rupture débutant à  $i = 70$ , nous regardons le pourcentage de simulations ayant trouvées la première rupture entre  $i = 60$  et  $i = 80$ , puis entre 60 et 90 ainsi que 60 et 100. On calcule aussi le pourcentage de simulations ayant trouvées la dernière rupture entre 110 et 130 pour les méthodes hors ligne et le temps de détection moyen pour les méthodes en ligne. Ces informations vont nous permettre de voir quelle méthode détecte le mieux la rupture et quelle est la plus rapide.

Deux méthodes de classification sont utilisées. La première est une méthode probabiliste basée sur un modèle de mélange gaussien, elle prend en compte uniquement l'échantillon des mesures du paramètre biologique. Elle est utilisée grâce au package "mclust" développé pour le logiciel R. [11]

Le modèle de mélange consiste à supposer que les données proviennent d'une source contenant plusieurs sous-populations. La forme générale d'un modèle à  $g$  composants est :

$$f(x_i, \Psi) = \sum_{k=1}^G \pi_k f_k(x_i, \theta_k) \quad (2.23)$$

où  $\Psi = (\pi_1, \dots, \pi_n, \theta_1, \dots, \theta_G)$  sont les paramètres du modèle de mélange et  $f_k(\cdot)$  les densités des composants.  $\pi_1, \dots, \pi_G$  sont les poids du mélange ou probabilités, on a donc :  $\pi_k < 1$  et  $\sum_{k=1}^G \pi_k = 1$

Soit  $G$  le nombre de groupe fixé, les paramètres du modèle de mélange  $\Psi$  sont généralement estimés par la maximisation de la log vraisemblance :

$$l(\Psi, x_1, \dots, x_n) = \sum_{i=1}^n \log(f(x_i, \Psi)) \quad (2.24)$$

Comme cette maximisation est compliquée, l'estimateur du maximum de vraisemblance (EMV) est obtenu grâce à un algorithme EM (Expectation-Maximization).

Une densité du mélange fini est associée à un groupe. Le modèle de mélange gaussien suppose une distribution gaussienne pour chaque composant :  $f_k(x_i, \theta_k) \sim \mathcal{N}(\mu_k, \Sigma_k)$  Les paramétrisations des matrices de covariances peuvent être obtenues par décomposition propre de la forme  $k = \lambda_k D_k A_k D_k$ , où  $\lambda_k$  est un scalaire contrôlant le volume de l'ellipsoïde,  $A_k$  est une matrice diagonale spécifiant la forme des contours de densité avec  $\det(A_k) = 1$ , et  $D_k$  est une matrice orthogonale qui détermine l'orientation de l'ellipsoïde correspondante.

Les modèles de mélange sont testés pour  $G = 1, \dots, 9$ . Le choix du meilleur modèle (donc du nombre de composants) mais également des paramétrages de covariances sont traités par le critère d'information BIC :

$$BIC_{\mathcal{M},G} = 2l_{\mathcal{M},G}(x|\hat{\Psi}) - \nu \log(n) \quad (2.25)$$

où  $l_{\mathcal{M},G}(x|\hat{\Psi})$  est la log-vraisemblance à l'EMV  $\hat{\Psi}$  pour le modèle avec les paramètres  $\mathcal{M}$  et  $G$ ,  $n$  est la taille de l'échantillon et  $\nu$  est le nombre de paramètres estimés. La paire  $(\mathcal{M}, G)$  qui maximise  $BIC_{\mathcal{M},G}$  est choisi.[10] [29]

La deuxième méthode est une classification ascendante hiérarchique (CAH) réalisée sur l'échantillon des mesures du paramètre biologique, le sexe et l'âge des patients. Le principe de la CAH est de calculer un tableau de distances entre les individus à classer. L'algorithme cherche, à chaque étape, à constituer des classes par agrégation des deux éléments les plus proches. Il existe différents algorithmes, ici c'est la méthode Ward qui est utilisée. Elle permet d'avoir une perte d'inertie inter-classe la plus faible possible à chaque étape, la perte d'inertie étant :

$$P = \frac{m_A m_B}{m_A + m_B} d^2(A, B) \quad (2.26)$$

$m_i$  est le poids de la classe  $i$  (nombre d'individus appartenant à la classe) et  $d^2$  la distance euclidienne au carré. L'algorithme s'arrête avec l'obtention d'une seule classe. Les regroupements successifs sont représentés sous la forme d'un arbre binaire ou dendrogramme. Le choix du nombre de classe est fait a posteriori en utilisant le dendrogramme. Avec la méthode de Ward, la hauteur des branches est proportionnelle à la perte d'inertie inter-classe, on coupe avant une forte perte d'inertie, ici le nombre de classes sera de deux.

Les détections de ruptures ont été réalisées avec les deux classifications mais on ne donnera que les meilleurs résultats pour chaque paramètre biologique. Ainsi, les résultats pour l'albumine et les protéines sont issus de la méthode avec la classification probabiliste et les résultats pour le chlore, la créatinine, la ferritine, le PSA et l'urée sont issus de la méthode avec la CAH.

Comme pour la méthode prenant en compte l'information service, les valeurs trop extrêmes et rares de trois paramètres biologiques ont été retirées des échantillons afin d'avoir de meilleurs résultats et moins de fausses alarmes. Ainsi pour la créatinine les valeurs au dessus de 250

sont retirées, pour la ferritine ce sont les valeurs à partir de 800 et celles au-dessus de 40 pour l'urée.

## Résultats

Les tableaux 2.15, 2.16 et 2.17 donnent les résultats de toutes les méthodes de détection de ruptures pour 100 simulations d'Albumine sur VISTA 2, respectivement, pour les bruits  $\mathcal{N}(\mu, 1)$ ,  $\mathcal{N}(0, \sigma)$  et  $\mathcal{N}(0.5T, 1)$  avec  $\mu$  et  $\sigma$  prenant comme valeurs 8.6 et 12.25.

On remarque que les résultats sont moins bons que lorsque l'on connaît le service de provenance les patients.

Pour la détection du bruit en moyenne ( $\mathcal{N}(\mu, 1)$ ), segment Neighbourhood en moyenne et breackpoint restent les meilleures parmi les méthodes hors ligne et Lepage pour les méthodes en ligne.

Les ruptures en variance ( $\mathcal{N}(0, \sigma)$ ) sont toujours plus compliquées à détecter, la meilleure méthode hors ligne est segment neighbourhood en variance et Student pour les méthodes en ligne.

La détection de ruptures en pente donne de moins bons résultats qu'avec l'information service, on peut atteindre 82% entre 60 et 100 avec la méthode Breackpoint et 82% avec la méthode GLR.

		N(8,6 ; 1)				N (12,25 ; 1)			
Méthode		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110_130
Hors Ligne	BinMoy	62	63	64	65	71	72	72	68
	BinVar	0	0	1	0	5	7	8	5
	SegMoy	71	74	76	72	78	80	81	77
	SegVar	2	4	5	0	6	10	12	7
	Pelt	52	57	58	57	69	73	74	62
	Pett	47	49	52	51	58	59	60	68
	Br	57	63	64	64	78	82	84	74
Méthode		%60-80	%60-90	%60-100	TpsMoyen	%60-80	%60-90	%60-100	TpsMoyen
En Ligne	Stu	65	71	71	15.46	79	80	80	9.3
	Bar	7	9	12	10.56	13	17	21	17.27
	GLR	68	70	70	18.14	77	79	79	11.14
	MW	65	66	66	17.99	72	73	73	13.04
	Moo	24	35	40	8.31	38	52	59	10.64
	LP	58	72	74	13.92	79	85	85	11.18
	KS	59	64	66	15.3	66	69	69	11.18
CVM	61	65	66	17.52	69	72	72	12.09	

TABLE 2.15 – Résultats des détections de ruptures sur 100 simulations d'Albumine, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec classification

		N(0 ; 8,6)				N (0 ; 12,25 )			
Methode		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors Ligne	BinMoy	17	24	29	10	12	22	30	14
	BinVar	10	11	11	9	52	56	56	48
	SegMoy	23	41	48	21	16	36	53	25
	SegVar	18	24	24	18	69	76	77	64
	Pelt	47	61	62	40	66	71	74	59
	Pett	18	27	33	10	4	12	17	13
	Br	10	13	14	5	3	7	10	6
Methode		%60-80	%60-90	%60-100	TpsMoyen	%60-80	%60-90	%60-100	TpsMoyen
En Ligne	Stu	29	46	49	11.5	43	55	61	8.02
	Bar	32	39	43	14.43	72	77	78	16.44
	GLR	24	32	34	13.29	56	63	64	18.62
	MW	17	27	31	14.14	15	25	33	13.96
	Moo	30	44	49	13.5	68	79	85	8.68
	LP	25	38	44	11.1	48	58	66	10.59
	KS	15	29	33	13.12	12	19	27	16.18
	CVM	16	27	30	13.1	15	25	34	13.59

TABLE 2.16 – Résultats des détections de ruptures sur 100 simulations d’Albumine, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$ ; méthode avec classification

		N(0,5T ; 1)			
Méthode		%60-80	%60-90	%60-100	%110-130
Hors Ligne	BinMoy	21	43	57	62
	BinVar	2	18	40	46
	SegMoy	13	40	60	70
	SegVar	2	20	53	71
	Pelt	28	56	66	74
	Pett	39	50	58	59
	Br	34	65	82	85
Méthode		%60-80	%60-90	%60-100	TpsMoyen
En Ligne	Stu	45	73	78	12.53
	Bar	26	31	41	37.6
	GLR	41	72	82	13.93
	MW	47	70	71	18.96
	Moo	13	36	72	8.73
	LP	32	67	80	11.43
	KS	41	68	72	14.46
	CVM	44	69	73	16.91

TABLE 2.17 – Résultats des détections de ruptures sur 100 simulation d’Albumine, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$ ; méthode avec classification

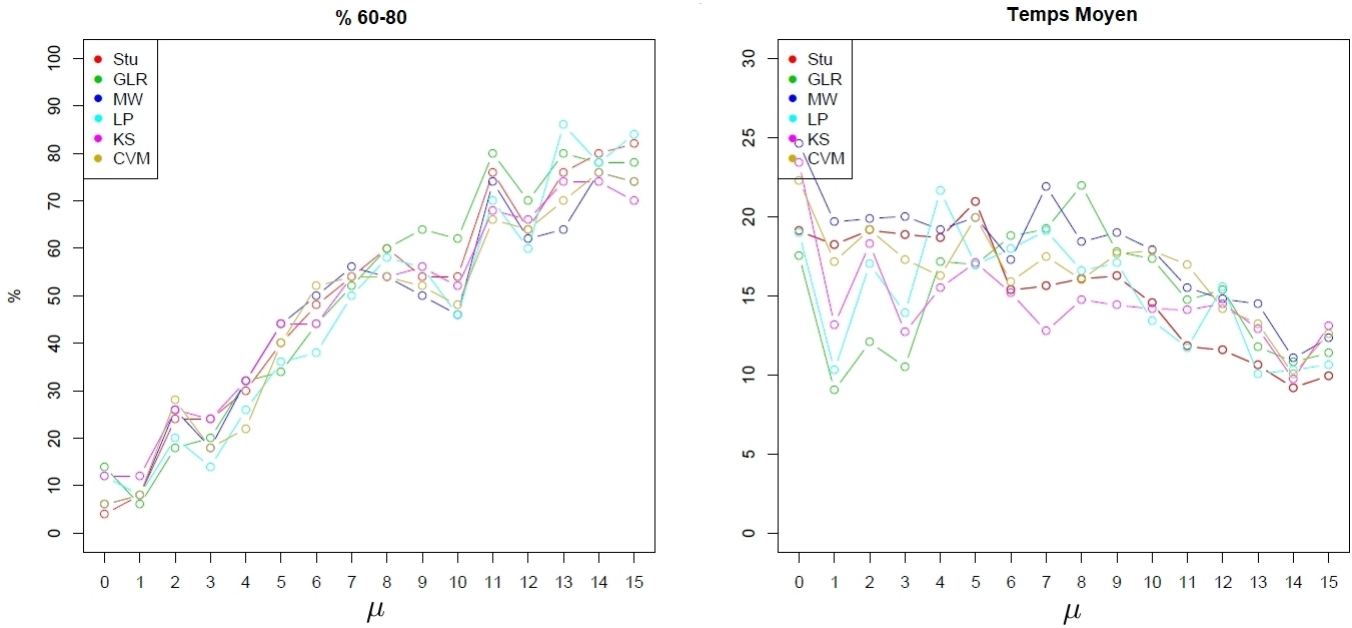


FIGURE 2.29 – Pourcentages de simulations d’Albumine avec ajout d’un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 0 à 15 ; Méthode avec classification

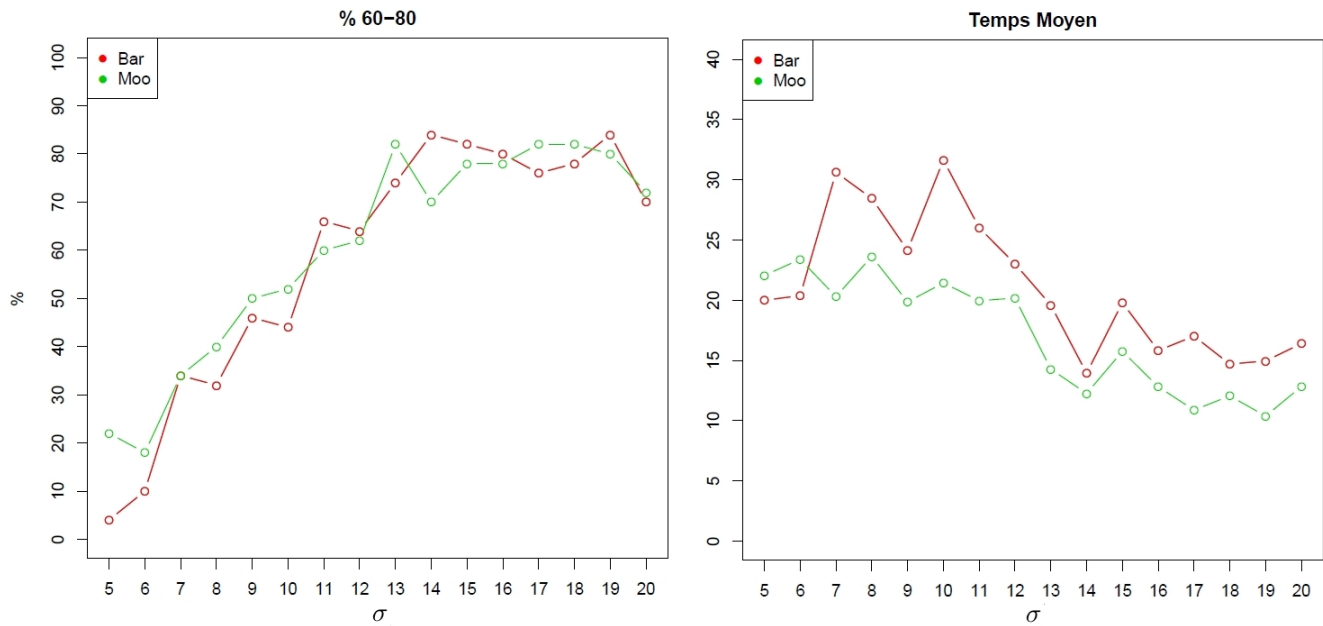


FIGURE 2.30 – Pourcentages de simulations d’Albumine avec ajout d’un bruit  $\mathcal{N}(0; \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 0 à 15 ; Méthode avec classification

Les résultats des paramètres biologiques chlore, créatinine ferritine, protéines, PSA et urée

sont dans l'annexe A.5. Malgré un TCC grand, les ruptures en moyenne, variance ou pente de ferritine sont encore mal détectées mais on obtient de meilleurs résultats qu'avec la méthode où l'on a l'information service. À l'inverse, on a de très bons résultats pour les trois types de ruptures que se soit pour le Delta-Check ou le TCC pour l'Urée et le PSA mais ils sont globalement moins bons qu'avec l'information service.

En général, les ruptures en moyenne avec le TCC sont bien détectées mais ce n'est pas toujours le cas avec le Delta-Check, les ruptures en variance sont rarement détectées.

### Fausses alarmes

Afin d'évaluer les détections de fausses ruptures, nous avons simulé, sans ajout de rupture, 200 échantillons d'une taille correspondant aux nombres de mesures réalisées par paramètre biologique sur une journée. Les détections de ruptures en ligne ont été appliquées sur ces échantillons de façon continue, à chaque rupture détectée, la méthode se réinitialise à la mesure suivante. Cela permet d'évaluer le nombre de fausses détections révélées par les méthodes en une journée au laboratoire.

Le tableau 2.18 donne les pourcentages de simulations ayant détectées 0 ou 1 rupture pour 200 échantillons d'albumine simulés de 60 individus. On peut donc dire que dans 83 à 95% des cas, il n'y aurait pas de fausses alarmes lors d'une journée de mesures d'albumine en laboratoire et que s'il y en avait, elles seraient de 1 par jour.

Méthode	0	1
Stu	86	14
Bar	95	5
GLR	86,5	13,5
MW	83,5	16,5
Moo	94	6
LP	87,5	12,5
KS	85,5	14,5
CVM	83,5	16,5

TABLE 2.18 – Pourcentages de simulations où les méthodes de détection en ligne en continue (avec classification) ont trouvé 0 ou 1 rupture sur 200 simulations d'albumine de 60 individus

Les résultats pour les paramètres chlore, créatinine, ferritine, protéines, PSA et urée sont donnés dans l'annexe A.5.

## 2.7 Incertitudes de mesure

### 2.7.1 Incertitudes de mesure (IM) au laboratoire

L'estimation des incertitudes de mesure sur les résultats dans les laboratoires de biologie médicale est imposée par la norme NF EN ISO 15189, elles doivent être calculées chaque année pour tous les niveaux de chaque type d'examen. L'incertitude de mesure permet de caractériser la dispersion des valeurs attribuées à un mesurande, c'est un indicateur de qualité d'un résultat et de fiabilité.

Plusieurs méthodes pour calculer l'incertitude de mesure existent, le CHU de Clermont-Ferrand utilise la méthode "CIQ, EEQ". Elle repose sur l'utilisation des contrôles internes de qualité et des données externes telles que des CIQ externalisés ou des évaluations externes de qualité (EEQ).

Les CIQ externalisés permettent de confronter les résultats des CIQ d'autres laboratoires, ayant utilisés les mêmes lots d'échantillons de contrôles, par leur moyenne, généralement mensuellement. Les EEQ sont des procédures d'évaluation des performances par le biais d'une comparaison inter-laboratoire. Ce sont les comparaisons d'une référence externe entre plusieurs laboratoires d'analyses. [6]

L'incertitude de mesure  $u(C)$  est calculée en prenant en compte les résultats des CIQ et des EEQ ou des CIQ externalisés (CIQE) :

$$u(C) = \sqrt{u^2(CIQ) + u^2(EEQ \text{ ou } CIQE)} \quad (2.27)$$

$u^2(CIQ)$  quantifie l'erreur aléatoire ou défaut de fidélité, la fidélité donne la reproductibilité intra-laboratoire et est obtenue par l'analyse d'un même échantillon, ainsi on a :  $u^2(CIQ) =$  la variance de l'ensemble des résultats du CIQ

$u^2(EEQ \text{ ou } CIQE)$  quantifie l'erreur systématique ou erreur de justesse, la justesse représente le biais entre les valeurs mesurées répétées et une valeurs de références. Il est calculé à l'aide des EEQ ou des CIQE.

Soit  $E_i = X_{lab,i} - X_{ref,i}$  où  $X_{lab}$  est le résultat du laboratoire et  $X_{ref}$  est la valeur assignée de la comparaison. Soit  $n$  le nombre de comparaisons réalisées, on a :

$$\bar{E} = \frac{\sum_i E_i}{n} \quad \text{et} \quad \widehat{\sigma}_E = \sqrt{\frac{\sum_i (E_i - \bar{E})^2}{n - 1}} \quad (2.28)$$

On obtient :

$$u^2(EEQ \text{ ou } CIQE) = \sqrt{\left(\frac{|\bar{E}|}{\sqrt{3}}\right)^2 + \widehat{\sigma}_E^2} \quad (2.29)$$

Le tableau 2.19 donne les incertitudes de mesure des trois niveaux de contrôle pour chaque analyseur pour les paramètres biologiques étudiés.

		Niveau					Niveau		
		1	2	3			1	2	3
Albumine	VISTA1	5,80	5,72	6,52	Ferritine	VISTA2	8,75	6,90	6,82
	VISTA2	7,53	7,10	7,10		VISTA3	7,91	7,18	7,64
	VISTA3	6,04	6,49	7,01		VISTA500	6,03	5,39	5,69
	VISTA500	6,92	7,24	6,94		VISTA1500	6,55	6,85	7,36
	VISTA1500	6,87	7,14	8,15					
Chlore	VISTA1	2,68	2,40	2,52	Protéines	VISTA1	4,21	4,11	3,95
	VISTA2	2,69	2,40	2,62		VISTA2	3,01	2,80	2,84
	VISTA3	2,08	1,81	1,97		VISTA3	4,64	4,23	4,10
	VISTA500	2,47	2,11	2,62		VISTA500	2,72	3,27	3,04
	VISTA1500	1,80	1,76	2,17		VISTA1500	3,45	2,91	3,38
Créatinine	VISTA1	8,11	5,30	4,22	PSA	VISTA2	6,81	5,54	6,27
	VISTA2	9,20	5,98	4,96		VISTA3	7,85	5,78	6,33
	VISTA3	10,53	6,48	4,81	Urée	VISTA1	9,34	7,39	7,22
	VISTA500	8,95	6,09	4,57		VISTA2	8,53	6,93	6,54
	VISTA1500	10,98	6,31	5,21		VISTA3	7,51	6,48	6,22
					VISTA500	8,00	7,53	6,93	
					VISTA1500	7,66	5,07	5,95	

TABLE 2.19 – Incertitudes de mesure calculées pour les paramètres biologiques étudiés

Les incertitudes sont données en pourcentage, pour chaque observation, l'incertitude correspond donc au pourcentage IM de la mesure réalisée.

### 2.7.2 Distribution des valeurs vraies

Nous recherchons la distribution des valeurs vraies des mesures de patients. Pour cela, nous allons estimer la fonction de densité à partir des données dites contaminées par l'incertitude de mesure. Nous utiliserons le package "decon", développé pour le logiciel R, utilisant des méthodes de déconvolution à noyau pour traiter les problèmes d'erreurs de mesure. [31].

On suppose que les données observées sont des données contaminées générées par un modèle d'erreurs additives. Les méthodes utilisées permettent de prendre en compte des erreurs de mesure homoscédastiques ou hétéroscédastiques. Les fonctions de densité sont estimées grâce à une adaptation de l'algorithme de transformation de Fourier rapide pour l'estimation de densité à l'estimation du noyau de déconvolution. L'estimation du paramètre de lissage peut être fait grâce à une méthode de bootstrap ou par une méthode empirique basée sur les travaux de Fan (1991). [9] [30] Dans notre cas les distributions d'erreurs de mesure sont considérées comme normales et hétéroscédastiques, elles varient donc selon chaque observation.

Les incertitudes de mesure sont données par niveau de CIQ, nous allons d'abord séparer les données par niveau. Les valeurs inférieures ou égales à la moyenne des CIQ niveau 1 auront



l'incertitude de ce niveau, celles comprises entre la moyenne des CIQ des niveaux 1 et 2 auront l'incertitude du niveau 2 et enfin les valeurs supérieures à la moyenne des CIQ niveau 2 auront l'incertitude du niveau 3. La figure 2.31 donne les distributions des trois niveaux de l'albumine pour l'analyseur VISTA2.

L'estimation des fonctions de densité des valeurs vraies est donc réalisée par niveau avec des erreurs de mesure hétéroscédastiques. L'écart type des distributions des erreurs pour chaque mesure  $m$  est :  $\sigma = \%IM \times m$ . La figure 2.31 donne les estimations des fonctions de densité des valeurs vraies de l'albumine sur VISTA2 des trois niveaux.

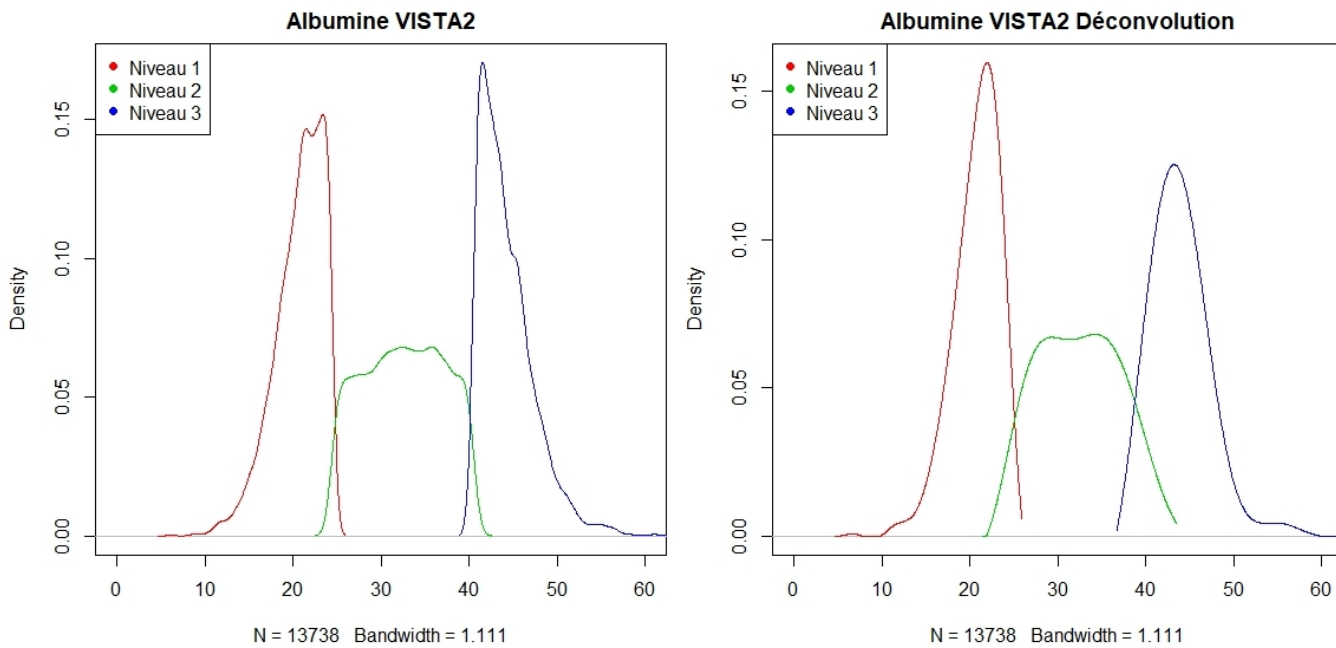


FIGURE 2.31 – Estimation des fonctions de densité des trois niveaux de CIQ pour l'albumine sur VISTA2 sur données contaminées puis par déconvolution de l'incertitude de mesure.

Les fonctions de densité estimées vont nous permettre de simuler des échantillons de valeurs vraies afin d'étudier l'impact des incertitudes de mesure sur ces valeurs.

Le niveau de chaque individu simulé suit une loi discrète avec, comme probabilités, les fréquences de chaque niveau dans les données réelles. Une fois le niveau choisi, l'individu est simulé selon la fonction de densité de son niveau à l'aide d'une méthode du rejet :

Soit  $f$  une densité de probabilité sur  $\mathbb{R}$ . Si le vecteur  $(X, Y) \in \mathbb{R}^2$  suit une loi uniforme sur  $\{(x, y) \in \mathbb{R}^2, 0 < y < f(x)\}$ , alors  $X$  a pour densité  $f$ .

La figure 2.32 donne la distribution des valeurs d'albumine sur VISTA2 dites contaminées ainsi que celle des simulations réalisées sur les fonctions de densité estimées après déconvolution de l'incertitude de mesure.

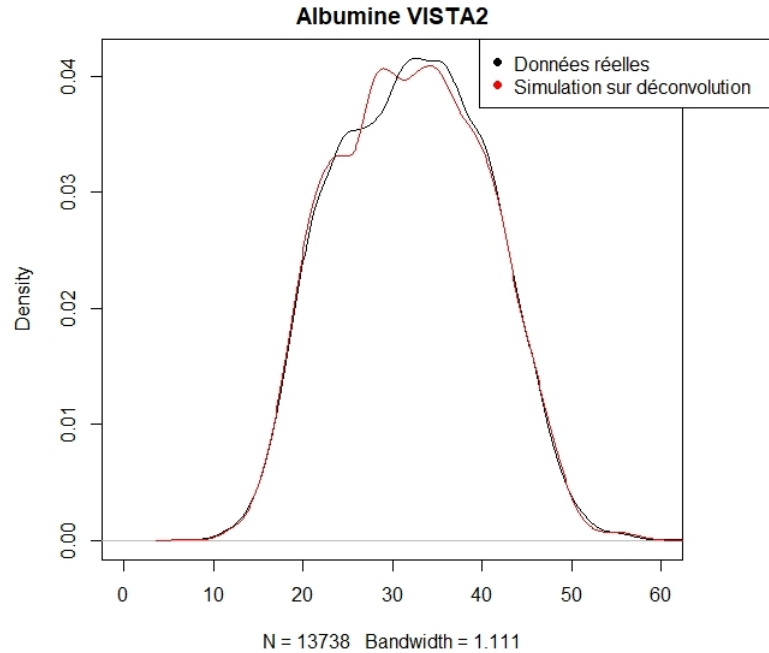


FIGURE 2.32 – Distributions d’albumine sur VISTA2 sur données contaminées et des simulations sur les fonctions de densité estimées après déconvolution de l’incertitude de mesure.

### 2.7.3 Impact de l’incertitude de mesure sur les valeurs vraies des mesures de patients

Pour étudier l’impact des incertitudes sur les valeurs vraies des mesures, nous avons comparé, par des tests de Kolmogorov-Smirnov, les distributions des valeurs avec ces mêmes valeurs contaminées par l’incertitude de mesure.

Un échantillon de valeurs contaminées  $m_1^*, \dots, m_n^*$  est simulé à partir d’un échantillon de valeurs vraies  $m_1, \dots, m_n$ . Toute valeur contaminée  $m_i^*$  est une valeur suivant la loi :  $\mathcal{N}(m_i, \sigma)$ , avec  $m_i$  la valeur vraie,  $\sigma = \%IM \times m_i$  et  $\%IM$  le pourcentage d’incertitude de mesure étudié.

Le tableau 2.20 donne les pourcentages d’incertitudes de mesure à partir desquels les lois sont différentes, au-dessus de ce seuil les lois sont différentes alors qu’en dessous les valeurs contaminées ont la même distribution d’après des tests de Kolmogorov-Smirnov. Les tests ont été effectués pour différentes tailles d’échantillons : 50, 100, 500, 2000 et 10000.

Paramètre	Taille échantillon	V1	V2	V3	V500	V1500
Albumine	50	93	95	91	94	89
	100	55	54	53	54	54
	500	26	25	26	26	26
	2000	17	16	16	19	17
Chlore	50	14	14	15	12	13
	100	8	8	9	7	8
	500	4	4	4	3	4
	2000	2,1	1,9	2,1	1,4	1,9
	10000	1,1	1	1,1	0,8	1,1
Créatinine	50	132	143	143	154	147
	100	72	75	77	84	77
	500	32	33	34	39	35
	2000	20	21	21	24	22
	10000	12	13	13	15	14
Ferritine	50		281	280	276	279
	100		133	134	134	131
	500		67	68	68	67
	2000		48	47	48	48
Protéines	50	45	44	47	55	51
	100	27	27	28	33	31
	500	13	12	13	16	15
	2000	8	8	8	10	9
	10000	5	5	5	6	6
PSA	50		438	750		
	100		198	221		
	500		87	92		
Urée	50	172	179	180	187	175
	100	91	92	94	98	96
	500	42	42	42	44	43
	2000	26	26	27	27	27
	10000	15	15	15	14	15

TABLE 2.20 – Pourcentages d’incertitudes de mesure à partir desquels les lois des valeurs vraies des paramètres biologiques et ces mêmes valeurs contaminées par l’incertitude sont différentes

Nous avons ensuite étudié l’impact de l’incertitude de mesure sur les valeurs vraies sous forme de séries temporelles pour déterminer à partir de quel pourcentage les incertitudes de mesure créent une rupture et donc un changement de moyenne, de variance ou de distribution sur les séries temporelles.

Pour cela, nous avons simulé des échantillons de valeurs vraies et ajouté l’incertitude de mesure pour 50 individus à  $i = 70$ , 75 simulations ont été réalisées pour chaque pourcentage d’incertitude étudié. Des détections de ruptures hors ligne ont ensuite été réalisées sur chaque échantillon.

La figure 2.33 donne les pourcentages de simulations d'albumine sur VISTA2 ayant trouvé une rupture entre  $i = 60$  et  $i = 80$  pour les pourcentages d'incertitudes allant de 0 à 30. On peut voir que les méthodes de détection hors ligne commencent à trouver une rupture quand l'incertitude de mesure dépasse 10.

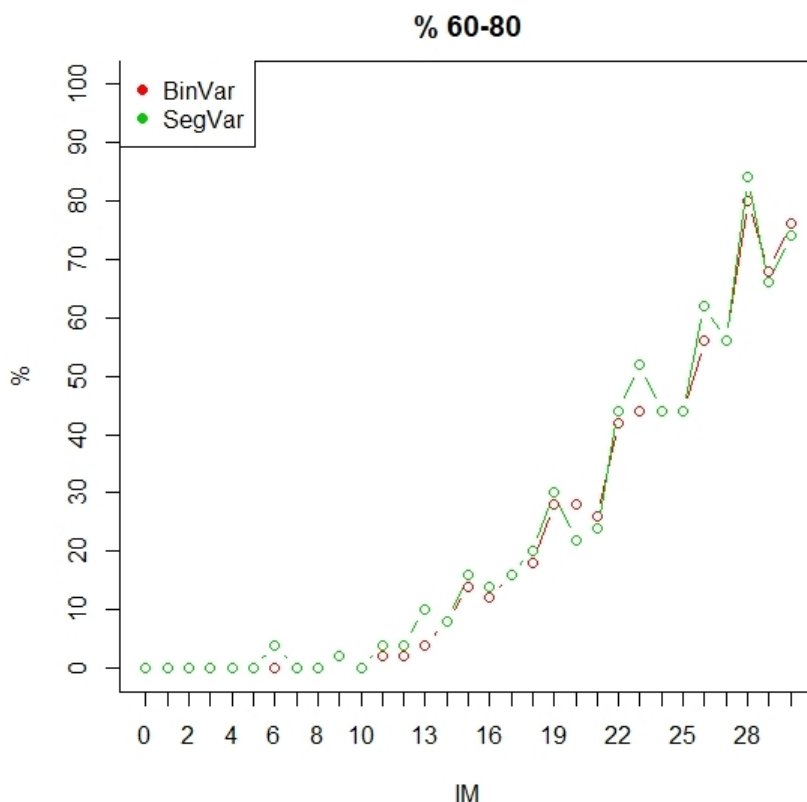


FIGURE 2.33 – Pourcentages de simulations d'albumine sur les fonctions de densité des valeurs vraies avec ajout d'un bruit  $\mathcal{N}(m, \%IM \times m)$  de taille 50 à  $i = 70$ , où les méthodes hors ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$ .

Les résultats montrent que les pourcentages d'incertitudes de mesure calculés par le CHU sont plus petits que ceux déterminés comme ayant un impact sur les distributions ou sur les séries temporelles des valeurs vraies.

Il n'y a que pour le chlore que ces pourcentages sont proches. En effet le pourcentage d'incertitude calculé par le CHU le plus grand est de 2.69 et l'on peut voir que les pourcentages d'incertitudes commencent à avoir un impact sur les séries temporelles à partir de 2.5 ou 3% et sur les distributions à partir de 1% pour des échantillons de taille 10000 et 1.9% pour les échantillons de taille 2000.

## 2.8 Conclusion

### 2.8.1 Résumé de l'étude

Dans ce travail de recherche, nous proposons une procédure permettant de détecter les dérives des analyseurs de laboratoires de biologie médicale en se basant sur les mesures réalisées sur des prélèvements sanguins de patients. L'étude est réalisée pour différents paramètres biologiques mesurées sur plusieurs analyseurs au CHU de Clermont-Ferrand.

Après avoir présenté les contrôles internes de qualité mis en place quotidiennement dans les laboratoires et étudié les différents panels de patients par service, un premier travail statistique d'analyse de données des patients est détaillé, donnant notamment une classification des patients et l'étude des répartitions dans les services des différentes populations trouvées. Un modèle permettant de simuler des données fictives suivant le même schéma que les données du CHU est donné.

### Stratégies de détection des dérives des analyseurs

Nous avons émis l'hypothèse que les dérives des analyseurs se traduisent par l'ajout d'un bruit gaussien sur les résultats des patients. Trois scénarios sont imaginés : une rupture en moyenne ( $\mathcal{N}(\mu; 1)$ ), une rupture en variance ( $\mathcal{N}(0; \sigma)$ ) et une rupture en moyenne évoluant en fonction du temps ( $\mathcal{N}(0, 5T; 1)$ ).

Deux algorithmes sont présentés pour détecter les différents types de dérives simulées, basés sur des méthodes de détection de ruptures de séries temporelles hors ligne et en ligne. Le premier utilise l'information des services en centrant et réduisant chaque mesure par la moyenne et l'écart type du service de provenance du patient, avant d'appliquer des méthodes de détection de ruptures. Le deuxième n'utilise pas cette information mais réalise des classifications, les méthodes de détection de ruptures sont ensuite appliquées sur chaque groupe. Les résultats sont donnés pour chaque type de dérives simulées.

Les probabilités de fausses alarmes sont estimées pour chaque méthode mise en place. Elles sont estimées sur des échantillons de taille correspondant à une journée de mesures en laboratoire. Cela permet donc de contrôler le nombre de fausses alarmes déclenchées par les méthodes lors d'une journée de mesures au laboratoire.

### Incertitudes de mesure

L'impact des incertitudes de mesure des analyseurs sur les données des patients est étudié. Pour cela, nous estimons d'abord les fonctions de densité des valeurs vraies des mesures, soit les valeurs non contaminées par les incertitudes de mesure, grâce à des méthodes de déconvolution. L'impact des incertitudes est ensuite étudié sur des simulations réalisées sur les fonctions de densité des valeurs vraies estimées. L'incertitude à partir de laquelle les données des patients

sont impactées est estimée pour tous les paramètres biologiques étudiés. Cette étude montre que les incertitudes de mesure calculées par le CHU n'ont pas d'impact sur les distributions et les séries temporelles des mesures des patients.

### 2.8.2 Perspectives

Le deuxième algorithme proposé, n'utilisant pas les services hospitaliers de provenance des patients, est une solution pouvant être mise en place dans des laboratoires où les patients ne sont pas aussi bien segmentés que dans un hôpital. En effet la segmentation des patients par service hospitalier permet de comprendre les valeurs parfois hautes ou basses de certains groupes de patients. Sans cela, les classifications permettent que les groupes de patients « atypiques » ne déclenchent pas la détection d'une fausse rupture tout en gardant les détections de réelles dérives impactant tous les mesures.

Mises en place dans les laboratoires, les stratégies proposées dans cette étude permettraient de donner des alarmes dans le but de lancer un contrôle de qualité. Aujourd'hui les stratégies de suivi et de contrôle des analyses biomédicales dans les laboratoires entraînent, en général, pour chaque paramètre biologique mesuré et pour chaque analyseur utilisé, trois tests de trois niveaux par jour, soit 90 CIQ par mois pour les contrôles routiniers.

Sur une année, très peu de problèmes ont été observés avec les CIQ et encore moins ont eu un impact sur les données des patients. De plus l'estimation des probabilités de fausses alarmes, réalisée pour chaque méthode donnée, montre que les stratégies ainsi mises en place n'engendreraient que très peu d'alarmes dans les cas où aucun problème n'apparaîtrait. Pour les méthodes utilisant les services de provenance des patients, nous avons estimé que des fausses alarmes s'enclencheraient environ 10 jours sur 100 soit seulement pour 3 jours par mois et que le nombre de fausses alarmes par jour ne dépasserait pas trois. Avec un tel scénario, le nombre de contrôle à effectuer serait au maximum de 9 par mois, ce qui diviserait par 10 le nombre de contrôles.

Aujourd'hui dans les laboratoires, les contrôles internes sont effectués à des heures imposées. Théoriquement, la dérive d'un analyseur peut n'être observée par les CIQ que plusieurs heures après avoir débutée. Les stratégies que nous proposons ont l'avantage d'étudier en temps réel les analyses effectuées et donc avertir au plus vite d'un problème potentiel. Le nombre de mesures de patients affectées serait donc plus petit avec nos méthodes.

Les types de dérives étudiées dans cette étude sont des hypothèses émises sur les réels impacts d'un problème sur les mesures effectuées pour les patients. Nous n'avons malheureusement pas pu observer de réels problèmes. Ces observations devraient être effectuées dans le but de déterminer le scénario et la méthode proposés qui correspondraient le mieux et de pouvoir les mettre en place.

Les algorithmes présentés dans ce travail étudient les paramètres biologiques indépendamment les uns des autres. Or, un même analyseur mesure plusieurs paramètres, on peut faire l'hypothèse qu'un dysfonctionnement devrait donc impacter plusieurs paramètres en même temps. Une stratégie permettant d'étudier les paramètres simultanément devrait permettre de détecter les problèmes plus rapidement grâce à la quantité d'informations plus importantes. On peut imaginer des détections de ruptures multiples ou l'analyse des paramètres centrés et réduits sur une unique série temporelle.

L'étude simultanée de plusieurs paramètres pour contrôler les analyseurs pourrait diminuer le risque de fausses alarmes. En effet, une alarme donnée pour un paramètre n'est peut-être pas la conséquence d'un réel problème si les mesures d'un autre paramètre biologique effectuées au même moment sur l'analyseur signalé ne rencontrent aucune anomalie.

De plus, ce même type d'étude pourrait aider à comprendre l'origine du problème sur l'analyseur. En effet, les analyseurs mesurent les échantillons à l'aide de différents modules. Un module est utilisé pour plusieurs paramètres. La détection de ruptures simultanée sur ces paramètres peut donc indiquer le module à l'origine du problème.

Les stratégies proposées dans cette recherche et celles proposées pour la suite devraient être testées pour différents paramètres biologiques résultant de plusieurs techniques de mesures et variant plus ou moins selon les panels de patients.

## Chapitre 3

# A characterization of volcanic ash grain particles based on the shape parameters. Example of observations collected around the crater of the Tungurahua volcano in Ecuador.

Sophie Miallaret<sup>1</sup>   Julia Eychenne<sup>2</sup>   Jean-Luc Le Pennec<sup>2</sup>   Anne-Fraçoise Yao<sup>1</sup>

<sup>1</sup> Université Clermont-Auvergne, Laboratoire de Mathématiques, CNRS UMR 6620, Campus des cézeaux 63171 Aubière Cedex, France

<sup>2</sup> Université Clermont-Auvergne, CNRS, IRD, OPGC, Laboratoire Magmas et Volcans, 63000 Clermont-Ferrand, France

### 3.1 Introduction

#### 3.1.1 Geographical setting

Tungurahua is an andesitic volcano, which means that it has a structure consisting flows of tephra or pyroclastic formed during the various eruptive stages. It is located in the central chain of the Andes 120 km south of Quito in Ecuador. It rises to an altitude of 5023 m and is surrounded at its base by three rivers, called Chambo, Patate and Pastaza. A plateau named the plateau of Quéro, at about 3 km of altitude, extends to the west of the volcano, at the foot north is located a tourist town Baños of 20000 inhabitants, the lower parts of the northwest west and southwest flank are also cultivated and inhabited, as is the Quéro plateau.

Tungurahua has frequent explosive activity. Over the past 1,300 years, the Tungurahua has been in operation on average every 80 to 100 years, with major periods of activity occurring



in 1640, 1773, 1883, 1918. Since 1999, it has entered an active phase, in 2006 its activity increased, in 2010 it erupts with large explosions and experienced continuous volcanic activity between 2012 and 2013. The local population had to be evacuated several times [8].

### 3.1.2 Data collection

This study follows the work of Julia Eychenne [8], Jean-Luc Le Penneec and Sébastien Leibrandt [15] and is based on data from ash measurements of the Tungurahua. The data describe the shape of the ash particles. They were acquired thanks to a Morphologi G3SE [15] which allows to measure automatically the morphological and geometric characteristics of a large quantity of grains. Thanks to a mobile plate, it can fully scan the glass plate and thus analyze thousands or even hundreds of thousands of particles by taking pictures of the plate for different magnifications. The time required to scan the plate ranges from 10 to 20 minutes. Software associated with the tool allows you to recombine the different shots. The images thus acquired allow the software to calculate different parameters of the ash grains such as the aspect ratio, the circularity, the convexity and the solidity, these parameters vary between 0 and 1. The aspect ratio parameter measures the more or less elongated shape of the particle, the longer it is elongated the more the parameter tends to 0. The circularity measures the more or less rounded shape but also takes into account the irregularity of the contour. The parameters of convexity and solidity measure the tendency of the particle to be more or less concave, that is, the irregularity of the contour. The more convex the particle, the more these parameters tend to 1.[8].

The dataset concerns information about ash samples described in the work of Lepenneec and Eychenne [[8]] who were interested with three volcanic ash grain size classes : 75-90  $\mu m$ , 250-300  $\mu m$  and 750-800  $\mu m$ . Then, the ash samples were collected from 22 sites around Tungurahua (Figure 3.1), sieved to be grouped by size classes. They were then analyzed by the G3 SE Morpho-granulometer.

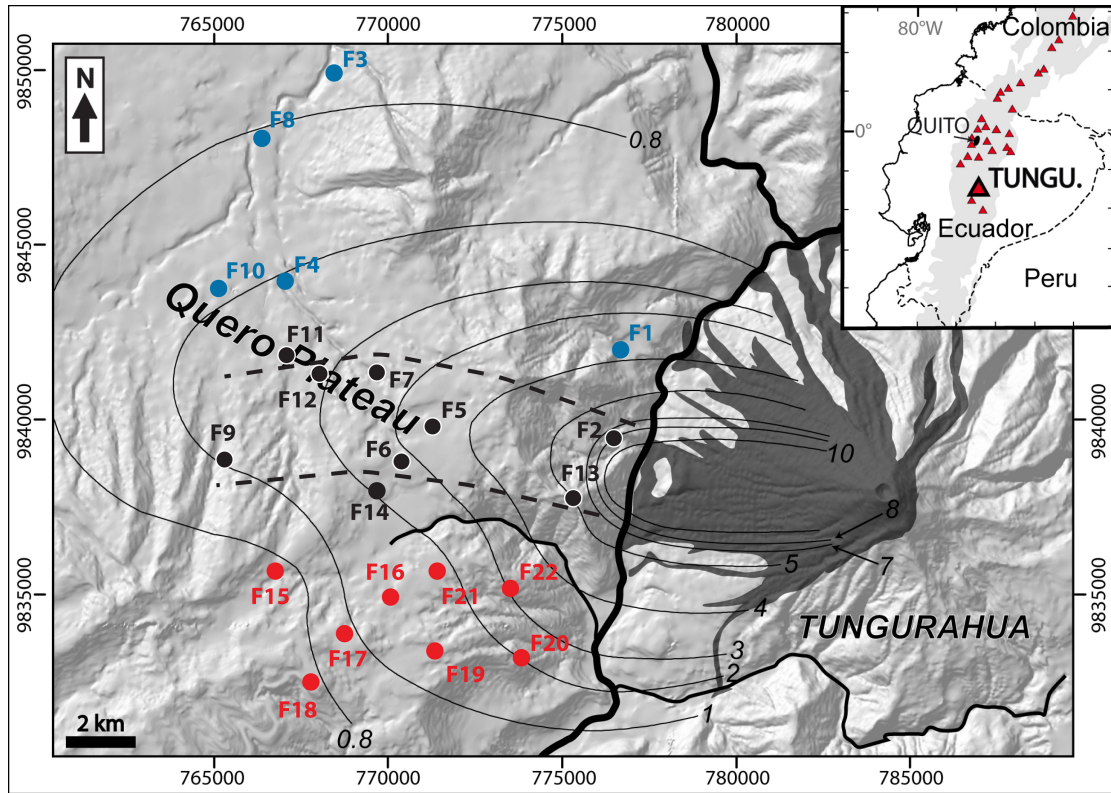


FIGURE 3.1 – Shaded relief map of Tungurahua edifice, dots are sampling sites where the for ash grain size were collected (The labels F1 to F22 represent the locations)

### 3.2 Statistical data analysis

We have carried out a statistical analysis in order to characterize the ash particles by their shape parameters. Namely, we aim at building statistical tools that help the understanding of the spatial repartition and shapes of the particles for each grain size class. Then, for each grain size class, our strategy is to 1) visualize shape parameters and how they are correlated 2) achieve the first point by taking into account the spatial location around the volcano 3) build models to characterize of the relationships which will be brought out by the two first points. This study was conducted using R software packages. Several variables are described by Lepennec and Eychenne [[8]]. In this work, we are interested with the 5 shape variables given in (Table 3.1) which characterize the morphologies of the ash particles.

Variable	Description
Circularity	Ratio of the perimeter of the equivalent circle and the perimeter of the particle
Convexity	Ratio of the perimeter of the convex envelope (equivalent to an elastic stretched around the particle) and the perimeter of the particle
Solidity	Particle area report on convex envelope area
Aspect Ratio	width to length ratio
Elongation	1-Aspect Ratio

TABLE 3.1 – Shape parameters measures on the ash particles

Because of the natural relationship between Aspect Ratio and Elongation, we finally deal on the one hand with five quantitative variables (or parameters) : the 4 shape parameters and the distance to the crater and on the other hand with a qualitative variable, the location variable with 22 levels : label from F1 to F22.

### Steps of the statistical analysis.

The first step of the statistical study concerned the distribution of each variable. The distribution of the variable location is based on a representation of the number observations per level (location) when those of the grain shape parameters are represented both by the histogram and the density of probability of the observations. The latter brought out an asymmetric in each distribution. Such a specificity does not suit to standard statistical analysis. An alternative consists then to proceed to a Box-Cox transformation to get symmetric distributions. The following steps are mainly based on the transformed variables.

In a second step, we proceeded to a bivariate (one-to-one) correlation study by first representing the correlation matrix in order to measure the strength of the linear correlation between the shape parameters (as well as with their transformed version) and also between them and distance to the crater. A Kruskal-Wallis dependence test was carried out to measure the dependence between each shape variable and the (qualitative) location variable. In each case, the test led to conclude to a dependence between the shape parameter and the location. But, the structure of the dataset did not allow us to get a useful model to understand the relationships we deal with. Then, we have proceeded to an intermediate step of clustering.

In the third step, we have clustered the particles using their observations of their quantitative parameters. The clustering approach considered here is the k-means. We have fixed the number of clusters to five. This step has led to a new qualitative variable which we denote by Group. In the following, we study the link between this variable and the quantitative variables (in the fourth step) as well as with the location variable in the other hand (in Step five).

The step four is devoted to Principal Component Analysis (PCA) in order to represent simultaneous interdependency between the shape variables and the distance to the crater. This

method gives scores that enables several 2D-plots orthogonal graphical representations. Each graphic, 2D-plot, representation is naturally based on two axes (dimensions), each expresses a percentage of information (inertia) of the dataset. The axis are ranked according to their contribution to global inertia (=100%) of the dataset. For example, in our study, the contributions are respectively 55.2% for dimension 1, 20.6% for dimension 2 and 19.2% dimension 3 for the grain size class 75-90. The PCA also provides coordinates for each particle (based on its shape characteristics and distance to the crater) which is then represented by a dotted point. Here, the points are colored according to the group (stated in the third step) they belong to. We notice that only the significant projections are interpretable. The quality of projection is measured by the square of the cosine between the initial vector and its projection on the PCA dimension or plane.

In the fifth step, A Correspondence Analysis (CA) was conducted in order to bring out the correspondence between the levels of the (qualitative) variables Group and Location. Similarly to the PCA, the CA summarizes the relationships between both variables by bringing scores. It allows the graphical representation of these scores in several 2D-plots. These scores allow us to proceed to aggregate of the locations in clusters which leads to a new variable, denoted by G.location. This new variable has enable to build some models.

The step six was conducted to model the link between the groups and the shape parameters, the distance to the crater and the variable G.location (referred as g.loc). A strategy based on subsampling (training and test approach) was achieved for the model validation step, “good classification rate” was determined. This rate describes the performance of the model. We also give the confusion matrix which provide “good classification rate” per group. The model of interest is the regression tree. It provide more information on how dependency between the variable Group and the other variables is structured.

A regression tree is made of nodes and leaves (terminal nodes) and edges. See for example Figure 3.11. The numbers on white square label the nodes. Each node corresponds to a question with a binary response about a variable : **yes**, labels the edge on the left of the node when **no** labels the other edge.

If the variable is quantitative or qualitative with more than two levels, an algorithm based on within-leave variance minimizing combing with a pruning step leads to a threshold or levels merging for such a variable. Only the significant information is kept in the final tree. The results give here are based on training and test sampling (both samples have the same size) approach.

Example of Figure 3.11. The tree is structured as follows :

1. Level 1 : Label 1 is for the first node (also called root node) it asks the question : “is Circularity<0.89”. The answer “yes” leads to node 2 (on the left) when “no” leads to 3 (on the right).

2. Level 2 : Nodes 2 and 3.

- Node 2 asks the question : “is Circularity $\geq 0.79$ ”. The answer “yes” leads to the node 4 and “no” leads to the node 5 .
- Node 3 asks the question : “is g.loc = F15/16/8/9/17,F20/22/21/5,F3/10/11/19/6/7/14/12/4”. The answers yes or no respectively leads to nodes 6 and 7 .

3. Level 3 :

- The node 6 is a is terminate node or leaf, due to the prune step, all following nodes (for example nodes 12 and 13 ) connected to this node do not appear in the final tree.
- The node 7 asks the question : “is Circular  $\geq 0.9$ ”. The answers yes or no respectively leads to nodes 14 and 15.

### 3.3 Results

#### 3.3.1 Distribution of the number of observations

The distribution of the number of particles per grain size class (75-90, 250-300, 710-850) and location is given in Table 3.2 and Figure 3.2. The percentage per class (in red on the barplot in Figure 3.2) shows that only between 2.8 and 6.6% of the ash grain particles (collected on all the locations) belongs to the class 710-850 when the percentage is almost the same for both other grain size classes : 75-90 and 250-300.

	75-90	250-300	710-850		75-90	250-300	710-850
F1	4458	3842	237	F12	2351	2462	252
F2	3009	3276	314	F13	3014	2691	281
F3	3000	2874	344	F14	2035	3348	381
F4	2735	2945	308	F15	2182	2454	256
F5	2653	2601	282	F16	2266	3266	372
F6	2975	3170	357	F17	2114	3345	439
F7	2467	3434	383	F18	1787	3128	260
F8	2513	3209	326	F19	2486	3016	380
F9	2027	2698	381	F20	3544	3094	345
F10	2620	2701	305	F21	2762	2542	318
F11	3013	3089	358	F22	2724	2792	276
Total	31470	33839	3595	Total	58735	65977	7155

TABLE 3.2 – The distributions of the sample grain size per class-size and location

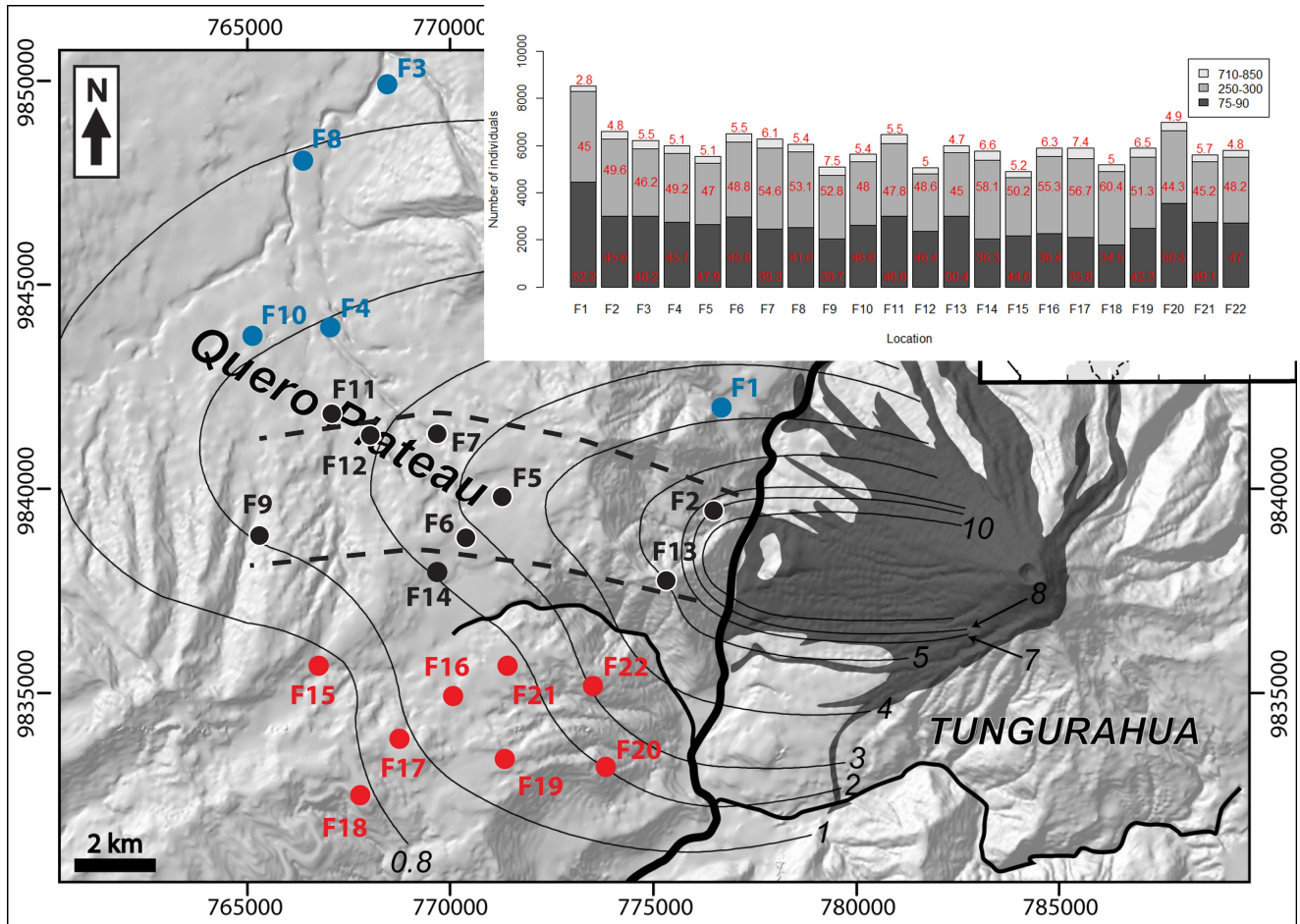


FIGURE 3.2 – The spatial repartition of the ash particles

### 3.3.2 Basic statistics on the ash grain shape parameters.

Basic characteristics of the distribution of shape parameters (presented in Table 3.1) are summarized in the Table 3.3.

	Circularity	Convexity	Solidity	Aspect Ratio	Elongation
Minimum	0,336	0,395	0,378	0,118	0
1st Quartile	0,801	0,91	0,903	0,645	0,17
Median	0,847	0,94	0,932	0,741	0,259
Mean	0,837	0,931	0,922	0,732	0,268
3rd Quartile	0,883	0,963	0,952	0,83	0,355
Maximum	1	0,998	0,998	1	0,882

TABLE 3.3 – Basic statistics about the grain shape parameters

### 3.3.3 The distribution of the shape parameters and Box-Cox transformation

On Figure 3.3, we give the distributions of the sample ash grain shape parameters per class size. This graphics show that all the distributions are asymmetric distributions. In this condition, the several classical statistical methods as those we used (for example mean computing, PCA or linear regressions) in this study fail. Actually, this is mainly due to the fact that in this condition, the extreme values have a great influence on the results which then, hide the main information concerning the majority of the observations. A strategy to fix this problem and get symmetric distributions consists in applying a Box-Cox transformation to the variable of interest. We have proceeded to such a transformation for each grain shape parameter. The corresponding transformed distributions are given on Figure 3.4. A first characteristic is : most of the particles of the sample of grain size class 710-800 have lower values of circularity and convexity (either transformed version or not) than in the two other grain size classes.

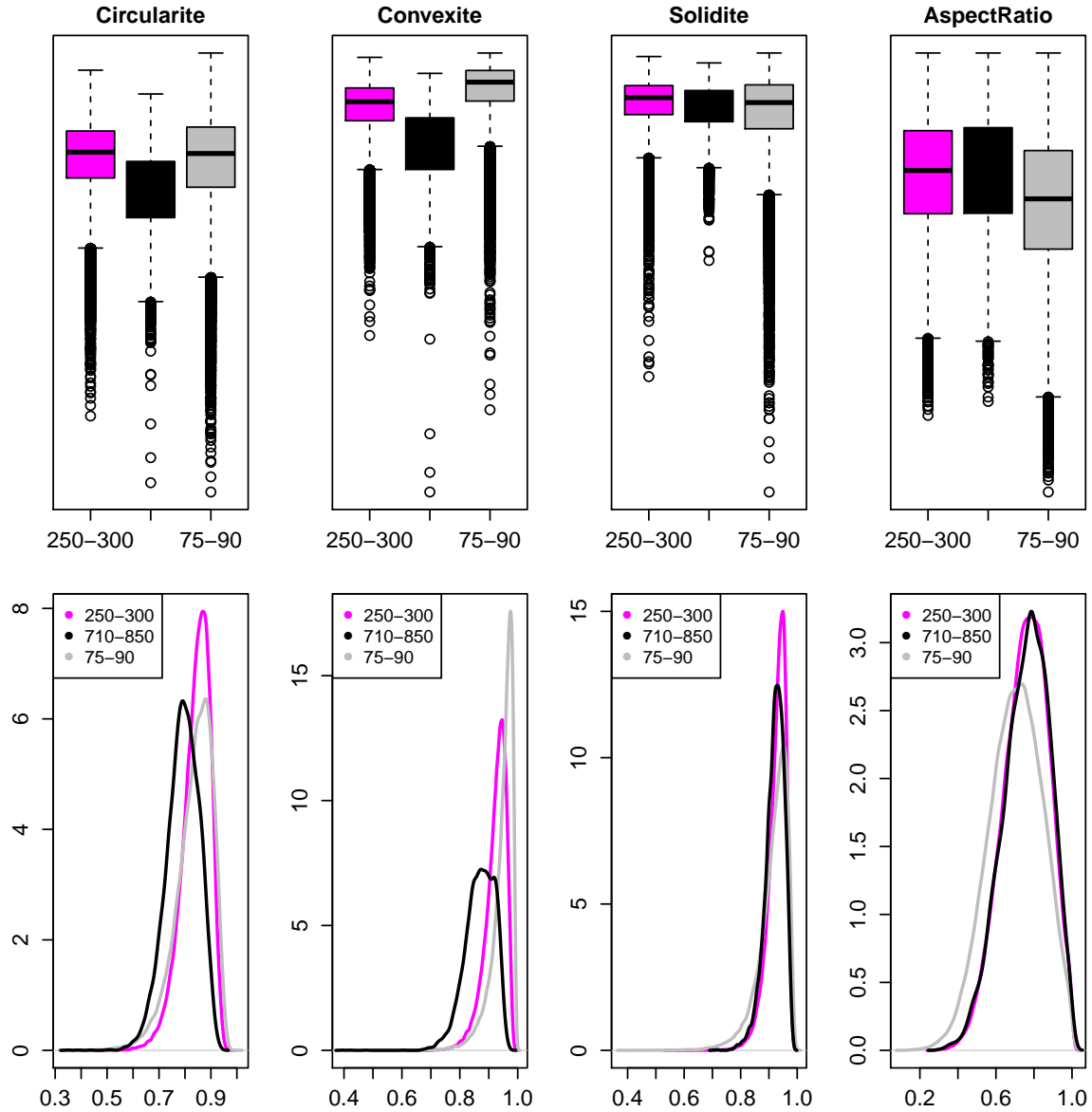


FIGURE 3.3 – Distributions of the initial the grain shape parameters versus grain size class



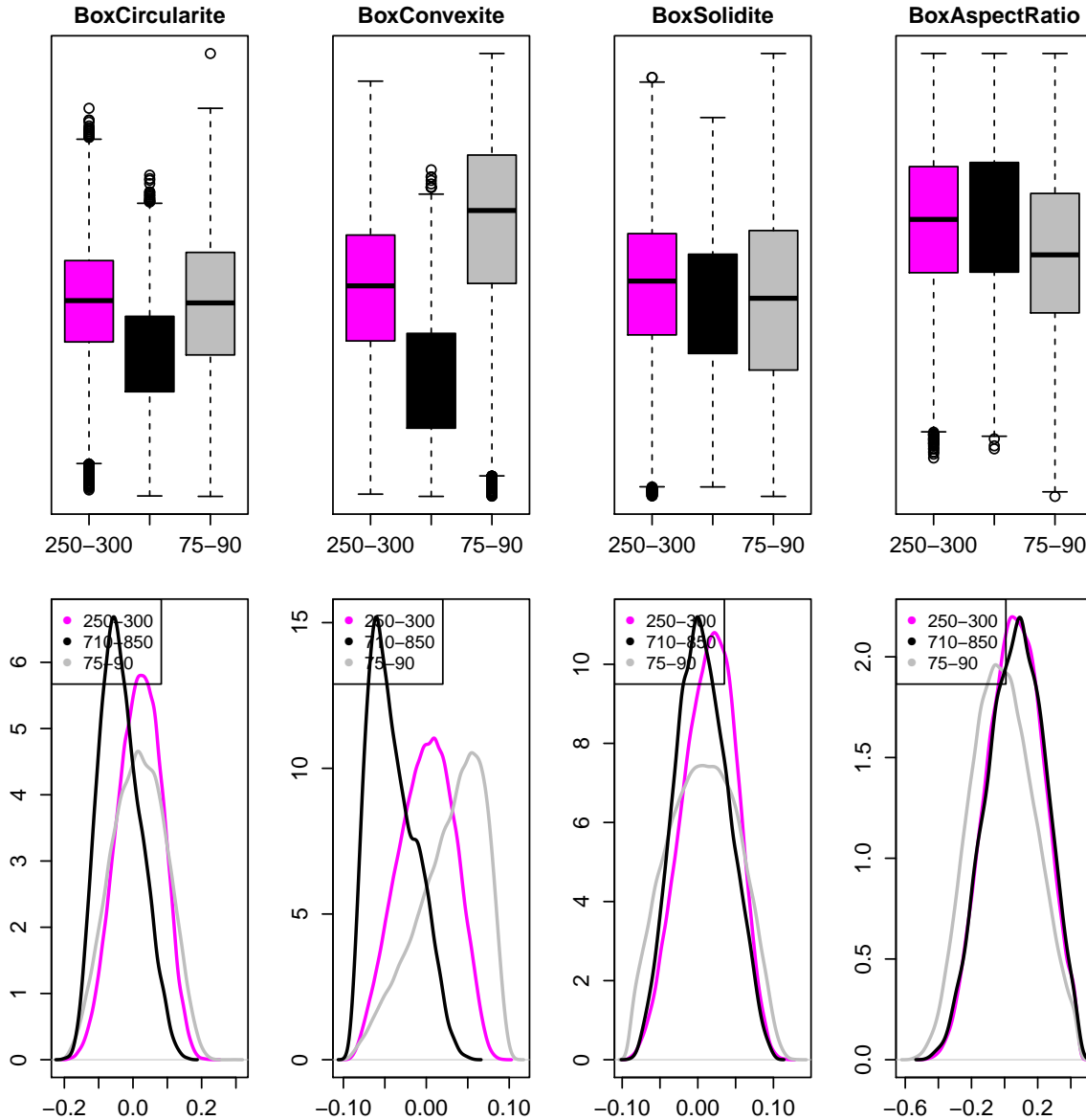


FIGURE 3.4 – Distribution of the transformed shape parameters by Box-Cox transformation

### 3.3.4 Relationship between the grain shape parameters and the distance to the crater.

We carried out a representation of the correlation matrix as a first step toward the study of the correlation between shape parameters as well as those between each of them and the variable Distance to the crater. The corresponding linear coefficients of correlation are given in Figure 3.5. All the significant linear coefficients of correlation are positive either non-transformed or the transformed versions and have the same scale, within each group (non-transformed or transformed). Obviously, if there is a stronger linear dependence between the parameters : Circularity, Convexity and Solidity. The correlations of these three parameters with the Aspect

Ratio are weaker (linear coefficients of correlation smaller 0.4) when those with the distance to the crater are almost null.

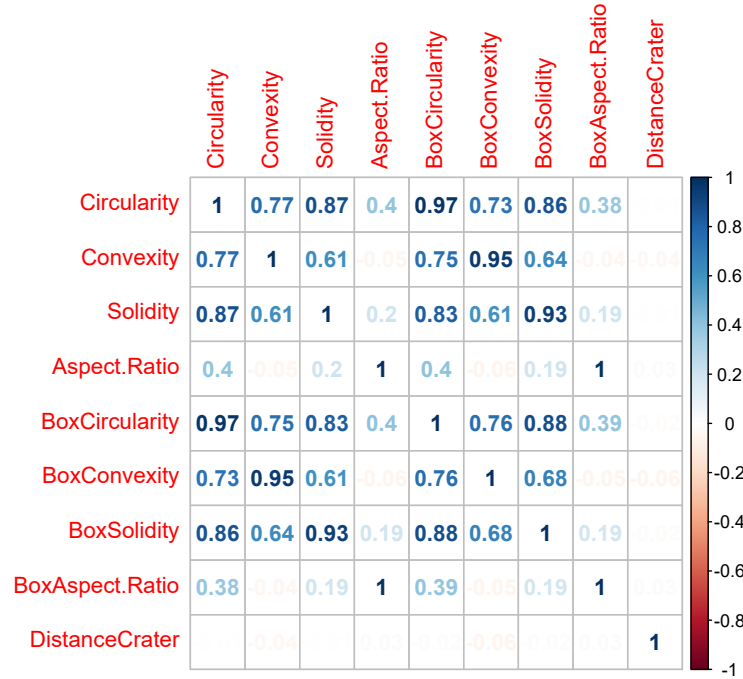


FIGURE 3.5 – Representation of linear coefficients of correlation between the shape parameters

Then, we were interested in the behavior these dependencies within grain size class. We have computed the matrix of correlation within each grain size class. The results are given in Figure 3.6 where obviously for example, the linear coefficient correlation (increases) between the Circularity and Convexity is respectively equal to 0.78, 0.86, 0.91 respectively for class 70-90, 250-300 and 710-850. This order still the same when replacing the variables by their the transformed version where it is respectively equal to : 0.78, 0.85 and 0.89. When, the linear correlation between Solidity and Convexity seems to invariant from a grain size class to another. Concerning the correlations with the distance to the crater, it seems less close to zero when splitting the sample into the three grain size classes, particularly for grain size class 710-850.

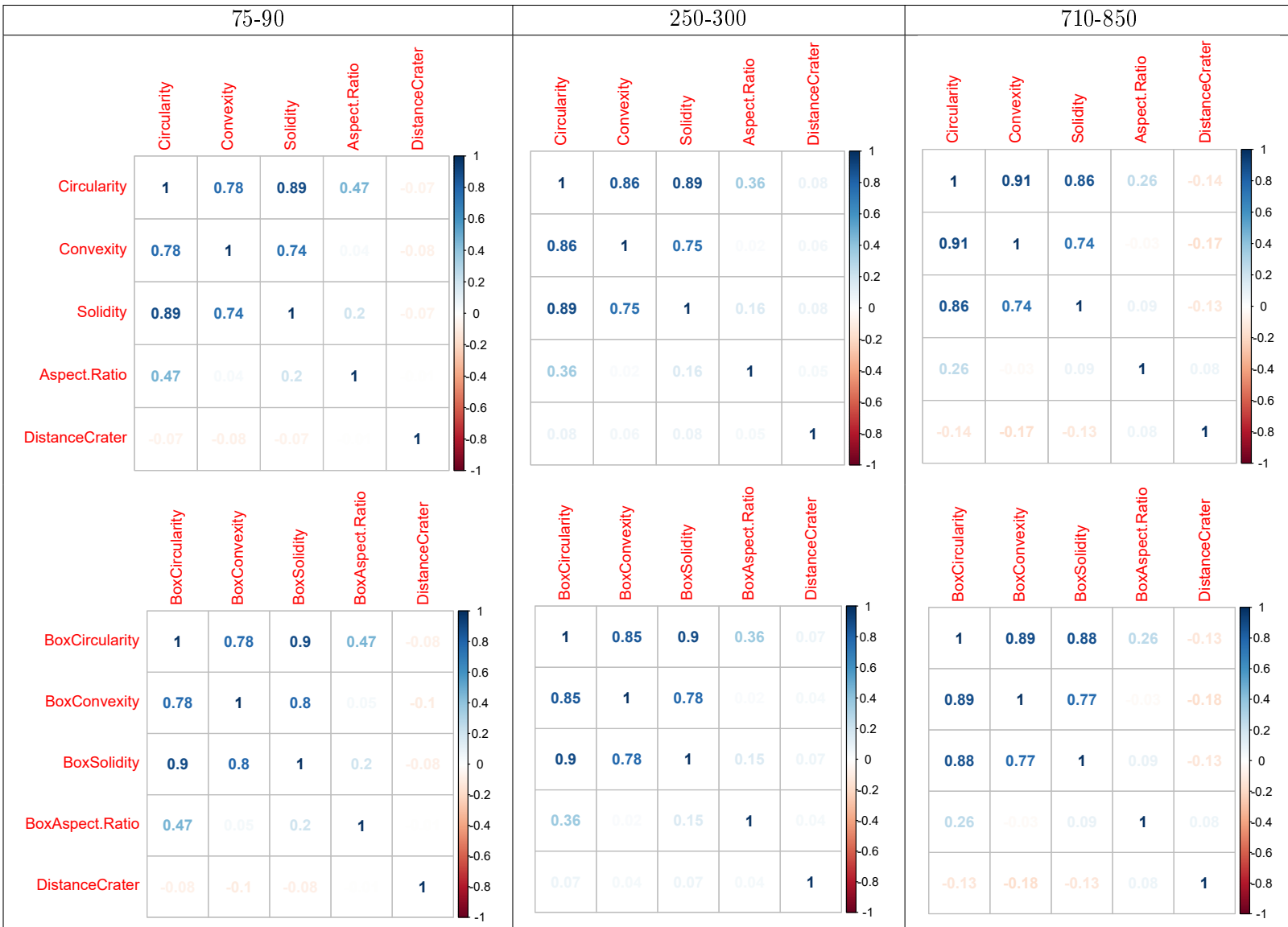


FIGURE 3.6 – Representation of linear coefficients of correlation between the shape parameters

But, a Kruskal-Wallis test between the variable Location (the  $F_i$ 's) and the grain shape parameters has led to a  $p$ -values  $< 10^{-4}$ . Then, we concluded that there is a dependency between the shape parameters and the location variable. Unfortunately, a direct statistical analysis between the shape parameters and the variables Distance to the crater and Location was not significant. But, since both the location parameter and the distance to the crater are naturally dependent, for convenience, we have chosen to keep the distance to the crater in the following study. We have proceeded to a clustering of the sample of observations of each grain size class in order to detect more (detailed) relationship between these variables. These groups should bring out more information to understand the relationships we deal with.

### **3.3.5 Classification of the ash grain particles using their shape characteristic and distance to the crater**

We have clustered the sample of each grain size class of particles into five groups by a K-means approach. Actually, we have proceeded to several numbers of group classification, the five groups clustering appeared to be a trade off. Of course, one can use an automatic approach for each class (that could be the subject of some other work). The characteristics of the five groups obtained by the k-means algorithm conducted with the Box-Cox transformed variables are given on Figure 3.7. The groups are referred both labels by numbers and colors : 1 (pink), 2(blue), 3(green), 4(violet), 5(brown). This procedure leads to a new qualitative variable which we denote by Group. The levels of this qualitative variable are the numbers of the groups. The distributions of this variable for each grain size class and location are given in Figure 3.7.

### **Representation of the clusters in the Principal Component Analysis graphics**

Several characteristics of the groups are brought out by the representations of Figure 3.7. But, to ease the interpretation of these results and highlight the link between the groups and the grain shape parameters (Links between the groups and their corresponding shape parameters and distance to the crater), we represent the observations in the Principal Component Analysis (PCA) graphics in Figure 3.8.

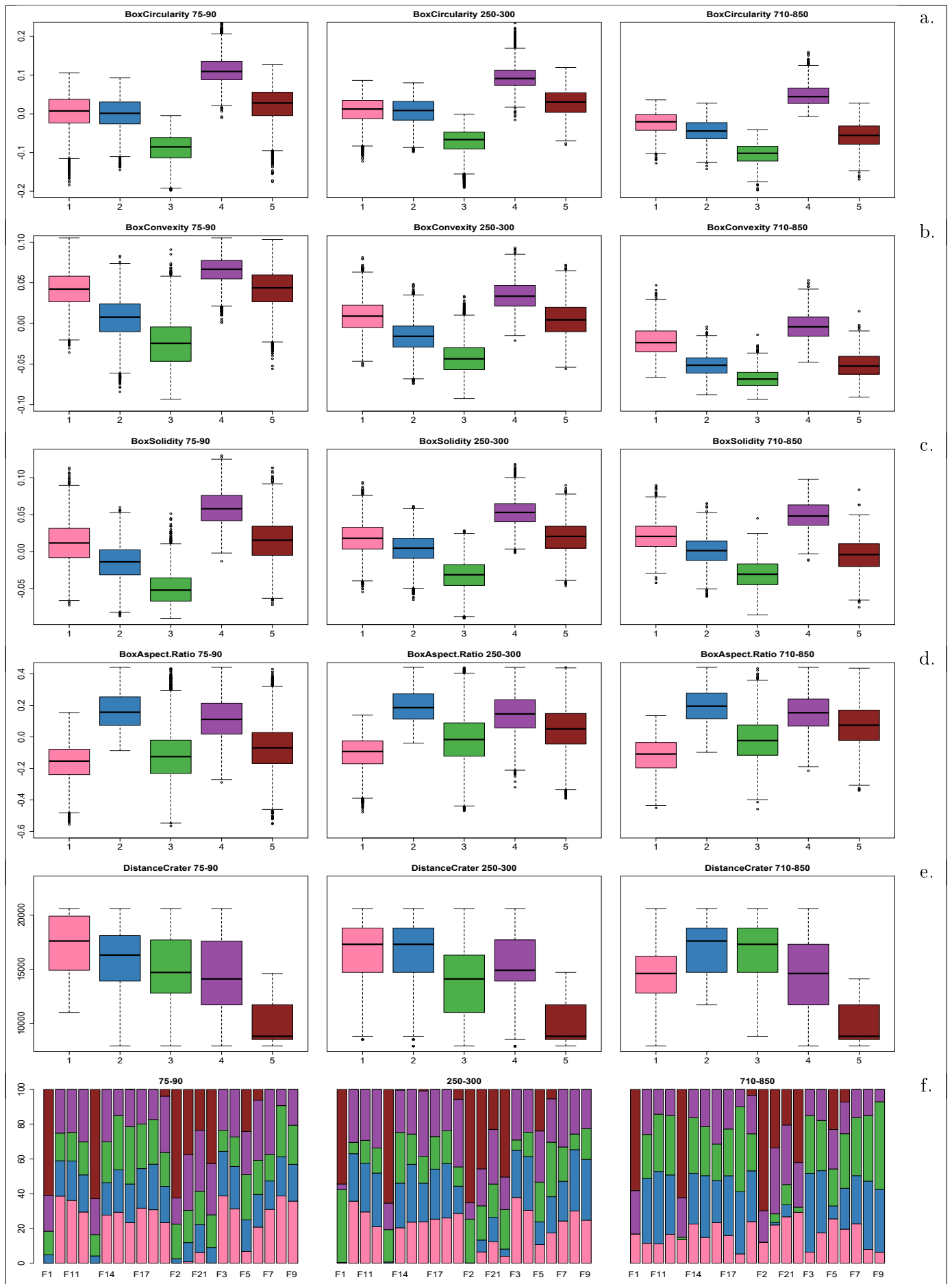


FIGURE 3.7 – Characteristic of the Groups

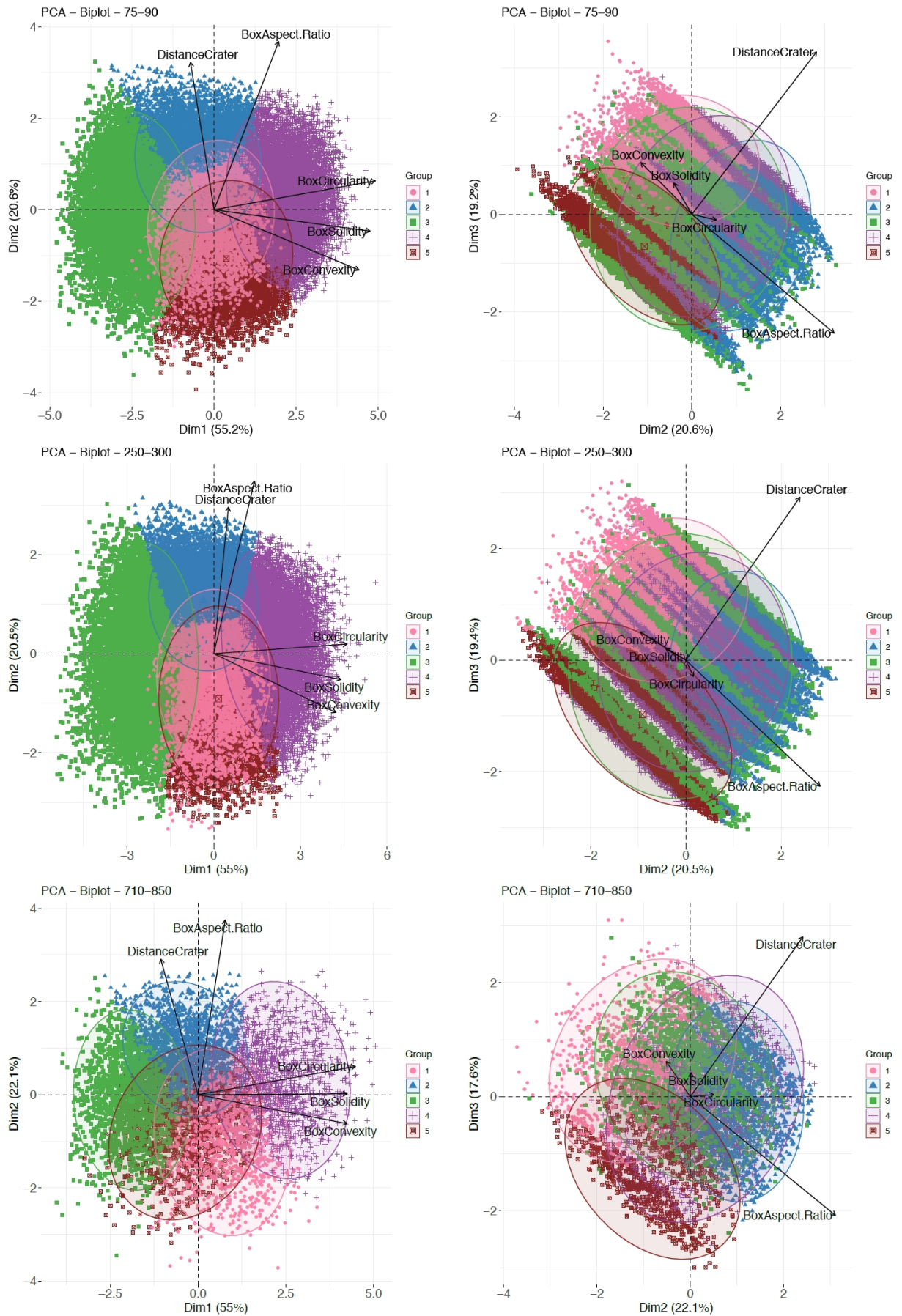


FIGURE 3.8 – Projection of observations of the grain shape parameters on the two first planes of PCA. Each dot is colored according to the group it belongs to as in Figure 3.7. 89

## Interpretation of the results of the PCA

**Dim 1** Whatever the grain size class, the first dimension (Dim1) brings out information about the parameters Circularity, Convexity and Solidity. Positive scores of the projection on Dim1 (on the right side of the first plane : Dim1-Dim2) indicate simultaneously high values of all these shape parameters and the reverse occurs for negative scores (on the left side of the first plane : Dim1-Dim2) of the projection on Dim1. This dimension concerns 55.2% of the information about the grain size class 75-90 when this percentage is of 55% for the two other the grain size classes.

Projections close to the center of the plan correspond to intermediate values of the circularity, convexity and solidity.

**Dim 2** The second dimension (Dim2) corresponds to the information on the Distance to the crater and Aspect ratio. The information concerning this dimension can be visualized both on the first plane (Dim1-Dim2) and the second (Dim2-Dim3). The latter gives a best view of this feature. Positive scores of the projection on Dim 2 (on the top the plane Dim1-Dim2 or right side of the second plane : Dim2-Dim3) indicate simultaneously higher values of the variables Distance to the crater and Aspect ratio and the reverse for negative scores (at the bottom of the plane : Dim1-Dim2 or left side of the second plane : Dim2-Dim3). This dimension brings 20.6% of the information about the class 75-90 when this percentage is of 20.5% for class 250-300 and 22.1% for class 710-850.

Projection close to the center of the plane represent intermediate values of the distance to the crater and aspect ratio.

**Dim 3** The third dimension (Dim3) also corresponds to the information on the distance to the crater and aspect ratio. But, with opposite sign on the scores : positive scores of the projection on Dim 3 (on the top the plane Dim2-Dim3) indicate higher values of the distance to the crater and lower values of the aspect ratio and the reverse for negative scores (at the bottom of the plane : Dim2-Dim3). This dimension brings 19.2% of the information about the grain size class 75-90 when this percentage is of 19.4% grain size class 250-300 and 17.6% grain size class 710-850.

A projection close to the center of the plane represent intermediate values of the distance to the crater or aspect ratio.

## Characterization of the groups

Based on the interpretations, we can characterize the five groups stated in the clustering step.

**Group 4** (violet) : whatever the grain size class, this group is characterized by higher values of Circularity, Convexity and Solidity (referring to the interpretation of Dim1) and intermediate values of the distance to the crater and the aspect ratio (referring to the interpretation Dim2 and Dim3)

**Group 3** (in green) : Whatever the grain size class, this group is at the opposite of Group 4 on Dim1 so that it concerns particles with lower values of convexity, solidity, aspect ratio and circularity. The projection of the observations of this group on plane Dim2-Dim3 reveals higher values of the distance to the crater for the class 710-850 and intermediate values of the distance to the crater for classes 75-90 and 250-300. For each size class, it concerns particles with intermediate values of the aspect ratio.

**Group 2** (in blue) consists of particles with 1) intermediate values of Circularity, Convexity and Solidity, 2) higher values of the distance to the crater for classes 250-300 and 710-850 and intermediate values of the variables for the class 75-90, 3) higher values of the aspect ratio for each class.

**Group 1** (in pink) consists of particles with 1) intermediate values of Circularity, Convexity and Solidity, 2) higher values of the distance to the crater for classes 75-90 and 250-300 and intermediate values of the aspect ratio for the 710-850.

**Group 5** (in brown) consists of particles with 1) intermediate values of Circularity, Convexity and Solidity, 2) Lower values of distance to the crater and intermediate values of the Aspect Ratio.

### **Multivariate analysis and classification : link between the location and the groups.**

In order to visualize the link between the variables Location and Group, we have achieved a Correspondence Analysis (CA) between the locations and the groups. The results are displayed on Figure 3.9 which shows that of the information concerning this relationship are concentrated by Dim 1 (90% for class 75–90, 90.3% for class 250–300 and 86.3% for class 710–850) when Dim 2 represents respectively 5.6%, 5.9% and so 10.1% of the information. We have clustered the location using their CA scores. The new location variable is denoted to g.loc. Example of interpretation, particles of Group 5 are mainly located in F1, F2, F13 unlike the other groups (see also Figure 3.7).



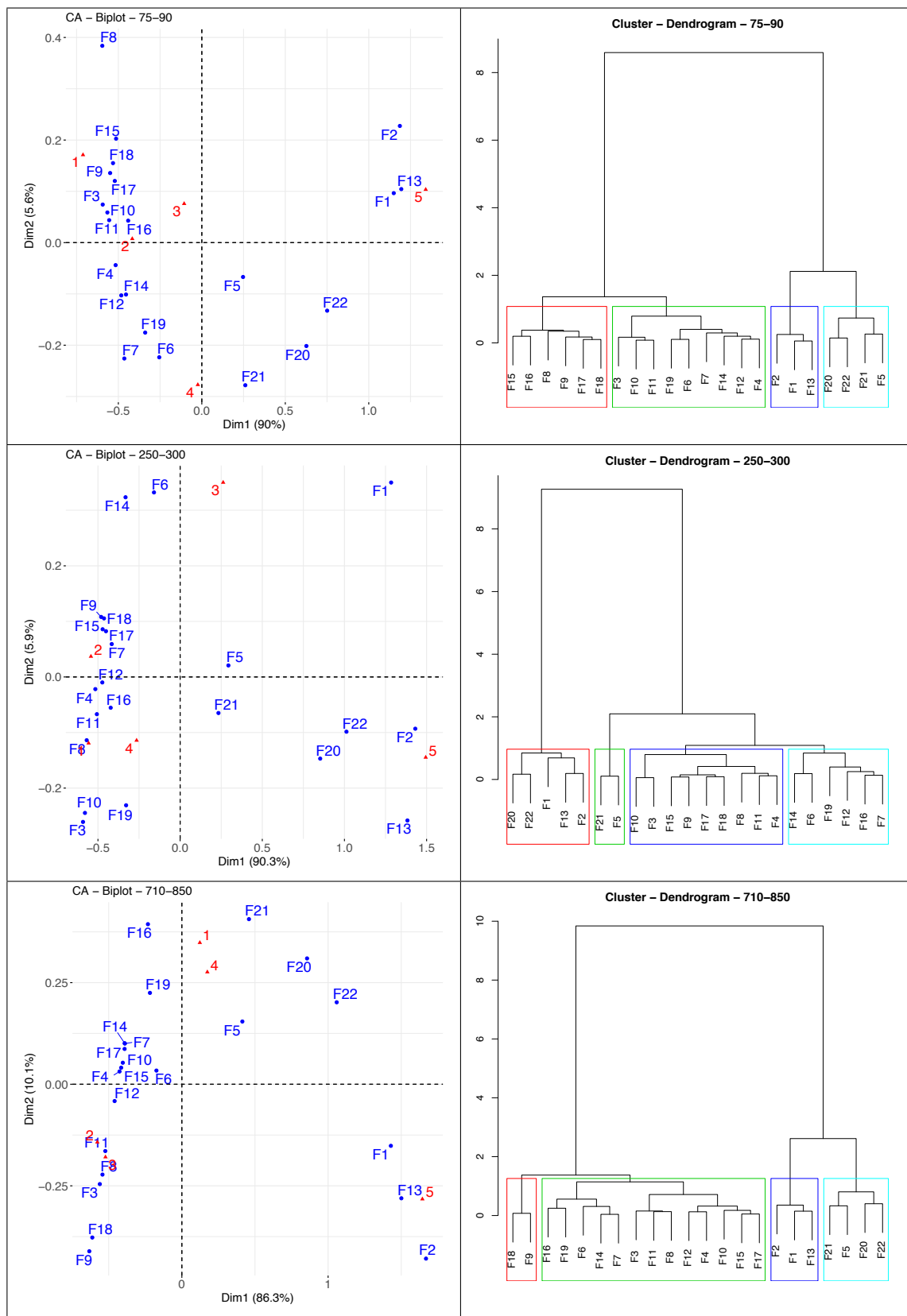


FIGURE 3.9 – Projection of the table of contingency Group versus Location in the first plane of CA. Correspondence between the locations and the groups

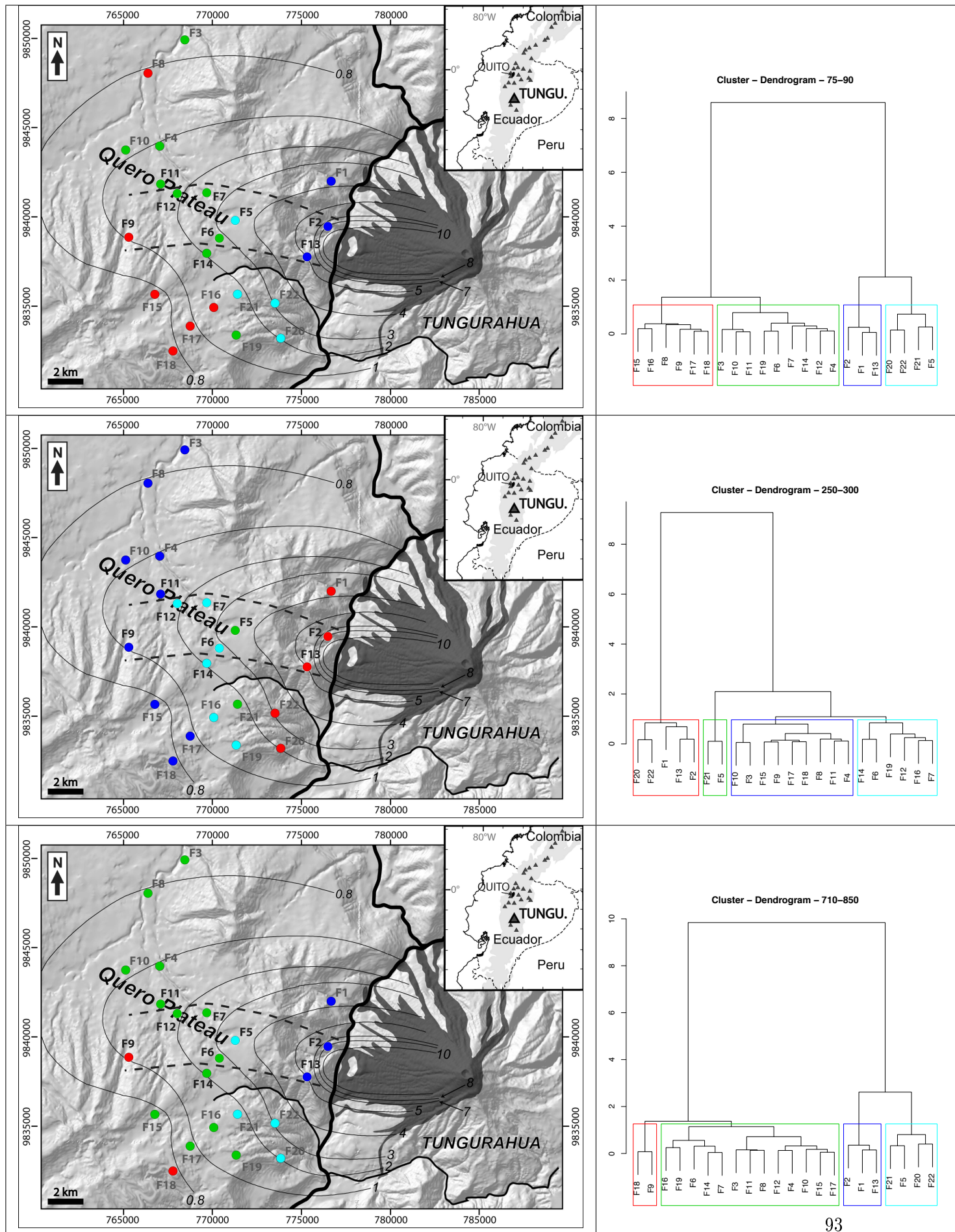


FIGURE 3.10 – Classification of the locations according to their CA scores.

### 3.3.6 Statistical modeling

We have applied regression tree approach in order to build a model which will precise the all relationships previously described between the variable Group and the other variables.

#### Regression tree

We first proceeded to a sampling without replacement of 300 observations per grain size class of the initial dataset. We have applied the regression tree model to the sample of observations of size is then equal to  $3 \times 300$ . Different regression tree models based on different sub-samples of size  $3 \times 300$  were computed. A part from some minor differences (the nearest 0.01 for the percentage), the trees most often were similar to the following one. The results for one of the trees are given in Figure 3.11, Figure 3.12 and Figure 3.13.

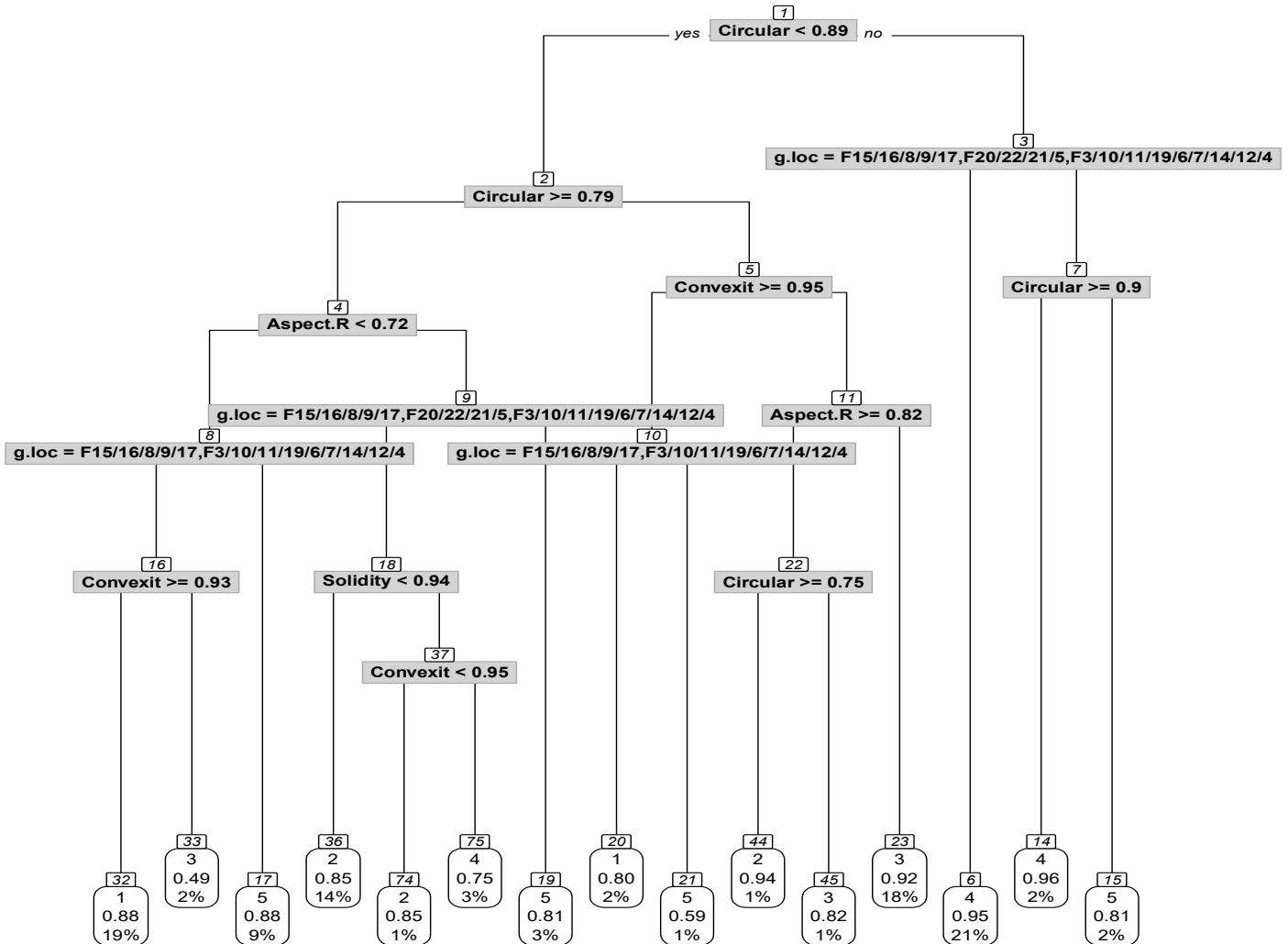


FIGURE 3.11 – Regression tree - Class 75-90

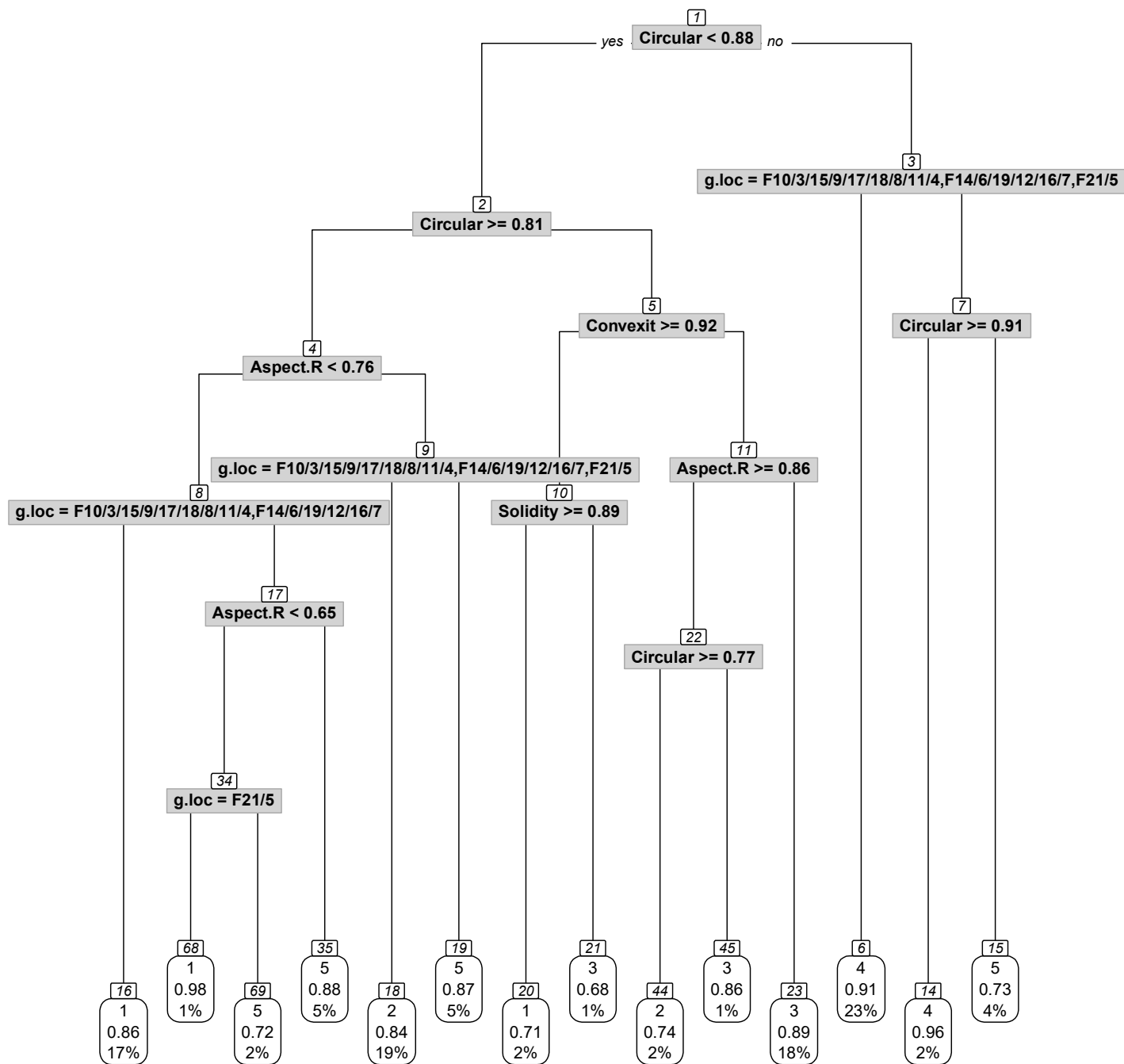


FIGURE 3.12 – Regression tree - Class 250-300

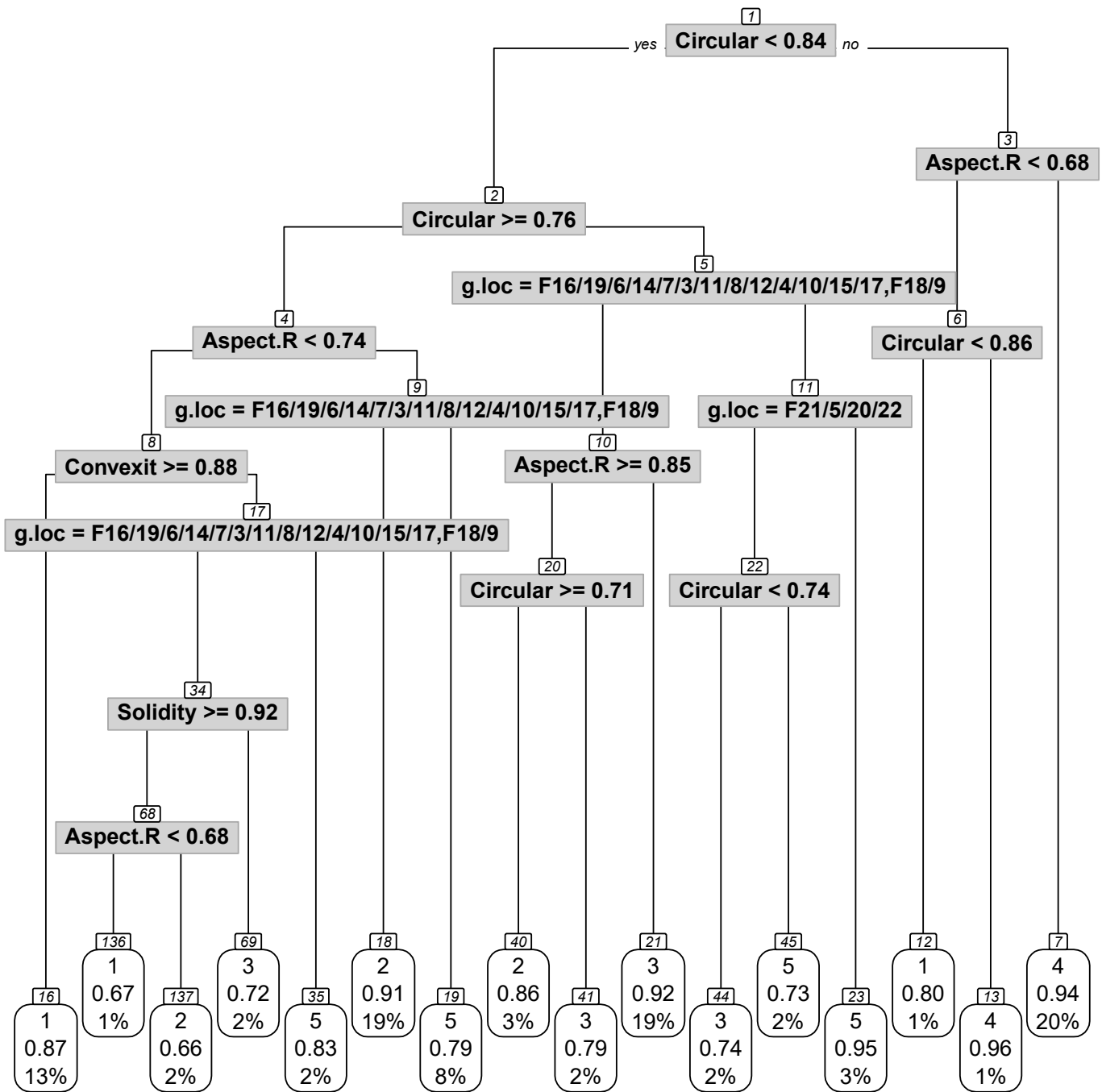


FIGURE 3.13 – Regression tree - Class 710-800

## Interpretation

### Example of interpretation for grain size class 75-90

We characterize the groups by an interpretation of combination of the information about connected nodes starting from the root node.

Example the combination of the information of the connected nodes [1], [3] and [6] are interpreted as follows : a particle of the leaf [6] has  $\text{Circularity} \geq 0.89$  and is located at  $\text{g.loc} = \text{F15/16/8/9/17,F20/22/21/5,F3/10/11/19/6/7/14/12/4}$ . With 95% of good prediction rate, we can say that the particles with this feature belongs to the group 4 and represent 21% of the whole 3x300 sample. This feature appears as a characterization of Group 4.

We characterize the main information in each group as follows.

**Group 1** mainly characterized by the information of the leaf [32]. The information given by the successive nodes [2], [4], [8], [16] represent particles with  $\text{Circularity} < 0.89$  and  $\text{Circularity} \geq 0.79$ ,  $\text{Aspect Ratio} < 0.72$  located at  $\text{g.loc} = \text{F15/16/8/9/17,F3/10/11/19/6/7/14/12/4}$  and  $\text{Convexity} \geq 0.93$ .

They represent 18% of all the observations and the good classification rate is equal to 88% in the group 1.

**Group 2** summarized at the leaf [36] represent particles with  $\text{Circularity} > 0.89$  and  $\text{Circularity} \geq 0.79$ ,  $\text{Aspect Ratio} \geq 0.72$  located at  $\text{g.loc} = \text{F15/16/8/9/17,F20/22/21/5,F3/10/11/19/6/7/14/12/4}$  and  $\text{Solidity} < 0.94$ .

They represent 14% of all the observations and the good classification rate is equal to 0.85 in the group 2.

**Group 3** corresponds to the information of the leaf [23] represent particles with  $\text{Circularity} > 0.89$  and  $\text{Circularity} \geq 0.79$ ,  $\text{Convexity} < 0.95$  and  $\text{Aspect Ratio} < 0.82$ .

They represent 18% of all the observations and the good classification rate is equal to 0.92 in the group 3.

**Group 4** corresponds to the information of the leaf [6] : previously comment.

**Group 5** corresponds to the information of the leaf [17] represent particles with  $\text{Circularity} < 0.89$  and  $\text{Circularity} \geq 0.79$ ,  $\text{Aspect Ratio} < 0.72$  are not present at  $\text{g.loc} = \text{F15/16/8/9/17,F3/10/11/19/6/7/14/12/4}$  and  $\text{Convexity} \geq 0.93$ .

They represent 10% of all the observations and the good classification rate is equal to 0.88 in the group 5.

The information about only a few percentage of particle of each group can be obtained by interpreting for example the leaf [20] for Group 1, leaves [74] and [44] for Group 2, leaves [33] and [45] for Group 3 and leaves [75] and [14] for Group 4, leaves [19], [21] and [15] for Group 5.

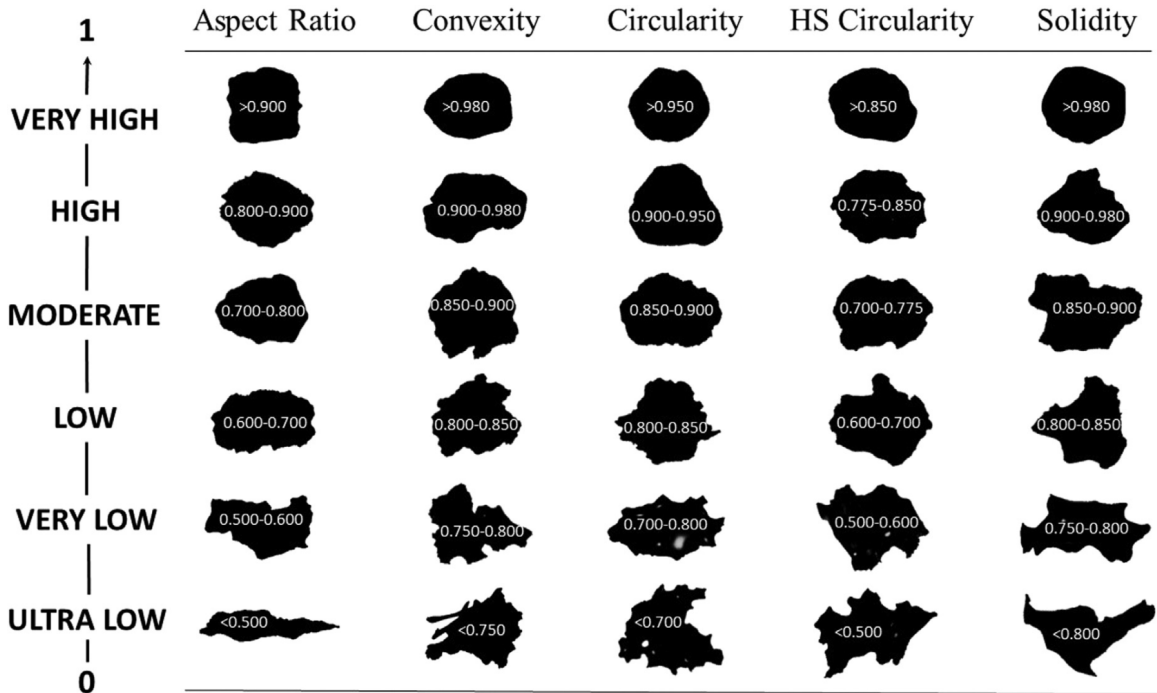


FIGURE 3.14 – Correspondence between the shape parameters and the morphology

### Model validation

In order to evaluate the predictive performance of the model, we have proceed to a training and test approach sampling (both samples have the same size). This has led to the following confusion matrix.

75-90 $\mu\text{m}$					250-300 $\mu\text{m}$					710-800 $\mu\text{m}$							
1	2	3	4	5	1	2	3	4	5	1	2	3	4	5			
1	1185	49	62	50	18	1	1069	100	52	19	31	1	603	13	19	13	51
2	49	1018	23	18	74	2	31	1128	58	79	72	2	55	926	45	22	13
3	71	88	1198	0	15	3	76	93	1108	0	15	3	40	66	975	0	27
4	78	41	0	1476	68	4	56	61	0	1509	40	4	21	19	0	825	11
5	42	43	61	35	838	5	28	27	33	81	834	5	38	35	12	25	546

75-90 $\mu\text{m}$					250-300 $\mu\text{m}$					710-800 $\mu\text{m}$							
	1	2	3	4	5		1	2	3	4	5		1	2	3	4	5
1	<b>17.95</b>	0.74	0.94	0.76	0.27	<b>16.20</b>	1.52	0.79	0.29	0.47	<b>13.70</b>	0.30	0.43	0.30	1.16		
2	0.74	<b>15.42</b>	0.35	0.27	1.12	0.47	<b>17.09</b>	0.88	1.20	1.09	1.25	<b>21.05</b>	1.02	0.50	0.30		
3	1.08	1.33	<b>18.15</b>	0.00	0.23	1.15	1.41	<b>16.79</b>	0.00	0.23	0.91	1.50	<b>22.16</b>	0.00	0.61		
4	1.18	0.62	0.00	<b>22.36</b>	1.03	0.85	0.92	0.00	<b>22.86</b>	0.61	0.48	0.43	0.00	<b>18.75</b>	0.25		
5	0.64	0.65	0.92	0.53	<b>12.70</b>	0.42	0.41	0.50	1.23	<b>12.64</b>	0.86	0.80	0.27	0.57	<b>12.41</b>		

TABLE 3.4 – Confusion matrix based of the effective (the table on the top) and the rate (table on the at the bottom)

### 3.4 Conclusion

This study was carried out on the results of measurements of volcanic ash form from Tungurahua with the aim of determining a link between the ash collection sites and their forms.

Tests initially showed that the shape variables and the locations of particle collection are well dependent, but direct analyses between these variables were not significant.

A classification is therefore made to find groups with characteristics of parameters of different shape depending on the distance to the crater. Five groups are thus established. The study of the repair of groups at each location, using a factor analysis of the matches and a classification of the places via this analysis, showed that the distributions of evolved groups well according to the location and distance to the crater.

Regression trees are proposed to model particle groups with form parameters and site groups. One can imagine later to model the groups of places created following the factorial analysis of the correspondences according to the parameters of form using regression tree or logistic regressions.

To continue this study, the variable representing the locations in the analyses, here distance to the crater, could be replaced by variables representing better the geographical positioning of the locations regarding the crater as it is the case but also according to the other places. For this, we can imagine the abscissa and ordered coordinates of the places in a landmark where the crater of the volcano would be the origin.





# Annexe A

## Annexes Chapitre 2

### A.1 Quantiles et moyennes des paramètres biologiques par sexe et catégories d'âge des patients.

	Age	Sexe	Min	25%	50%	Moyenne	75%	Max
Albumine	0-14j	F	18,5	32,05	37,95	35,75	41,13	50,1
		M	21,6	24,15	33,05	31,96	36,58	45,6
	14j-1a	F	9,8	29	34	32,86	37,7	49,7
		M	13,8	27	30,65	30,73	35,18	46,5
	1-8a	F	7,6	27,33	35,55	34,96	43,48	55,1
		M	3	29,95	37,2	35,71	43,5	55,3
	8-15a	F	8,7	35,5	42,8	40,72	46,6	59,9
		M	10	31,53	40,45	38,22	45,18	54,8
	15-19a	F	11,8	34,2	41,1	39,89	45,9	62,3
		M	13	29,48	36,9	37,36	45,53	70,3
	19a ou +	F	5,7	26,18	32,6	32,42	38,5	70,8
		M	6,9	24,7	30,7	31,37	37,5	71,7

	Age	Sexe	Min	25%	50%	Moyenne	75%	Max
Chlore	0-6m	F	86	105	108	107,91	111	147
		M	69	105	107	107,27	110	165
	6m-1a	F	93	105	107	107,99	109	143
		M	84	105	107	107,57	109	147
	1-18a	F	84	104	106	105,84	108	145
		M	81	103	105	105,41	107	165
	18a ou +	F	52	102	105	104,76	108	151
		M	51	102	105	104,71	108	155

	Age	Sexe	Min	25%	50%	Moyenne	75%	Max
Créatinine	0-14j	F	13	39,7	52,8	55,96	68,2	189
		M	12,39	39,65	54,7	57,78	70,8	306
	14j-2a	F	12,39	20,4	24	26,63	29,1	187
		M	12,39	20,6	24,4	30,42	30,4	357
	2a-5a	F	12,39	23,1	27,4	34,08	32,5	533
		M	12,39	23,6	28	29,8	33	179
	5a-12a	F	12,39	32,9	39,3	46,64	47	572
		M	12,39	31,7	38	42,59	46,63	522
	12a-15a	F	12,39	43,8	51,7	54,09	60,7	150
		M	12,5	43,53	53,5	59,07	64,28	440
	15a-19a	F	12,5	50,4	58,4	62,44	66,7	735
		M	12,39	60,5	71,05	77,19	82,1	1210
	19a ou +	F	12,6	53,3	65,3	85,75	85,5	2110
		M	12,39	66,2	82,2	111,36	111	2580

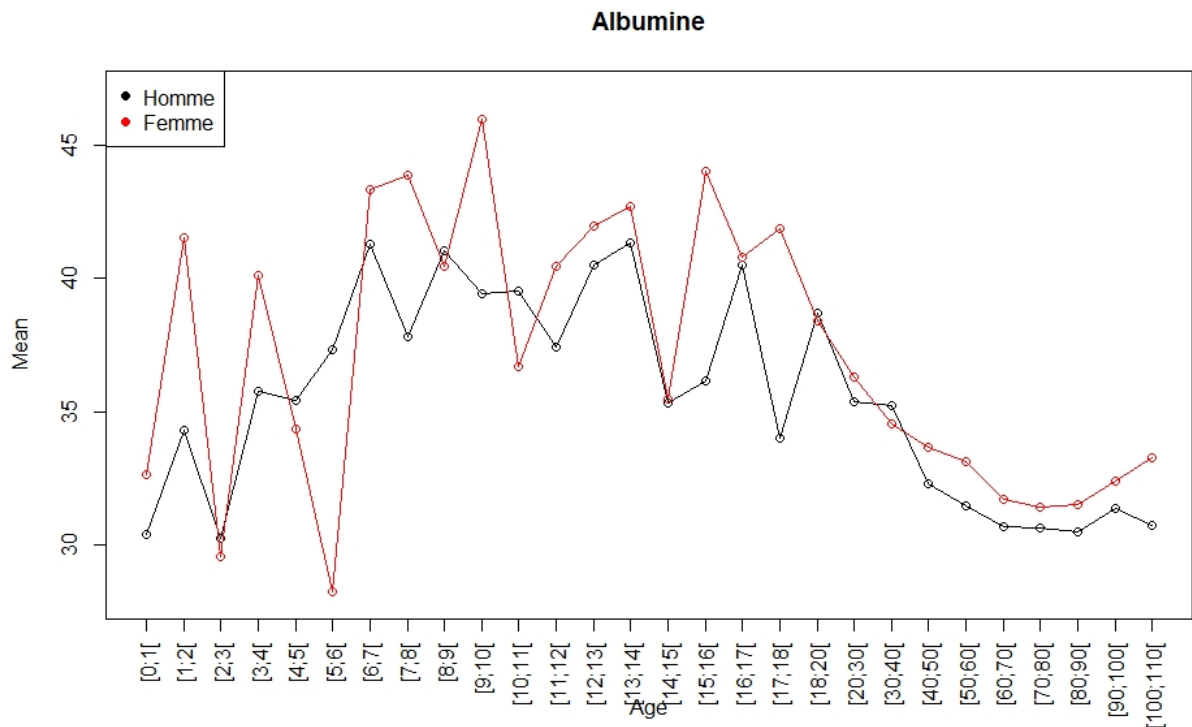
	Age	Sexe	Min	25%	50%	Moyenne	75%	Max
Ferritine	0-1m	F	129	158,45	227,6	270,43	317,43	667,9
		M	56,7	148,3	290,9	332,39	376,95	1116,1
	1m-6m	F	12,6	50,5	84,4	191,09	169,6	2131,3
		M	4,3	33,3	91,95	550,15	206,73	9973,7
	6m-1a	F	7,9	21,93	31,25	45,66	57,95	179,8
		M	4	18,9	26,3	53,04	48,5	719,6
	1a-5a	F	1,1	17	27,2	46,34	52,7	702,8
		M	1,7	15,53	26,5	136,13	53,3	8516,8
	5a-19a	F	1	17,6	33,9	731,95	71,05	281077
		M	2	25,85	42,5	147,04	84,33	7979,1
	19a ou +	F	0,9	31,8	94,9	359,16	288,6	120207
		M	0,7	94,63	236,2	643,51	540,18	168664

	Age	Sexe	Min	25%	50%	Moyenne	75%	Max
Protéines	0-14j	F	21	49	55,3	54,57	60,9	85,6
		M	23,8	47,7	53,8	53,27	59,4	88,2
	14j-1a	F	19,6	49,2	58,1	56,63	64,8	97,8
		M	24,5	49,7	58,3	57,16	64,6	88,6
	1a-6a	F	36,4	65,5	70,2	69,25	74,3	101,3
		M	33,5	63,3	68,6	67,56	72,8	108,2
	6a-9a	F	25,1	66,5	72,1	70,87	76,1	96
		M	38,6	65,48	70,1	69,67	74,6	114
	9a-19a	F	23,1	67,9	73,4	72,5	78	105,5
		M	20,3	68,3	73,7	72,65	78,2	109
	19a ou +	F	11	62,7	68,9	68,36	74,5	141,7
		M	13,6	62,4	68,8	68,23	74,5	152,5

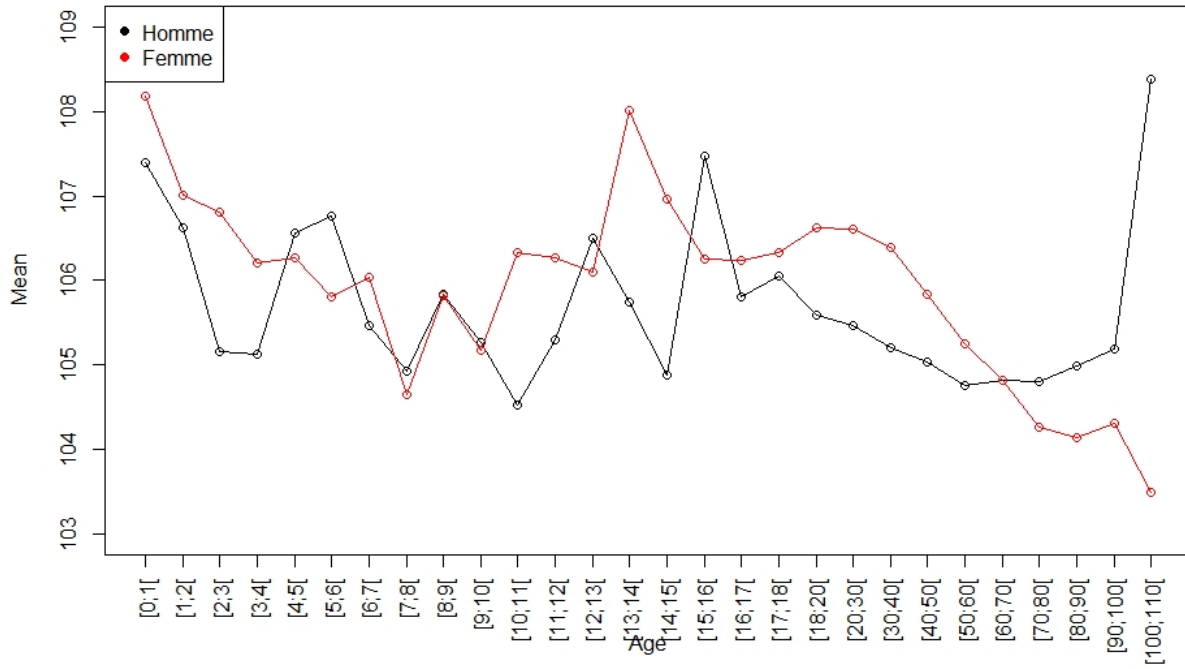
	Age	Sexe	Min	0,25	0,5	Moyenne	0,75	Max
PSA	0 ou +	F	0,01	0,01	0,01	0,11	0,02	1,33
		M	0,01	0,65	1,66	20,21	4,69	3600

	Age	Sexe	Min	25%	50%	Moyenne	75%	Max
Urée	0-14j	F	0,4	2,1	3,4	3,87	5,3	16,5
		M	0,4	2	3,4	3,88	5,1	16,5
	14j-1a	F	0,4	1,7	2,7	3,06	3,8	23,9
		M	0,4	1,7	2,6	2,99	3,8	19,2
	1a-19a	F	0,4	2,9	3,9	4,3	5	48,1
		M	0,4	3,1	4,2	4,65	5,4	38,5
19a ou +	F	0,4	4,1	5,8	7,78	8,9	121,6	
	M	0,4	4,7	6,6	9,06	10,2	115,8	

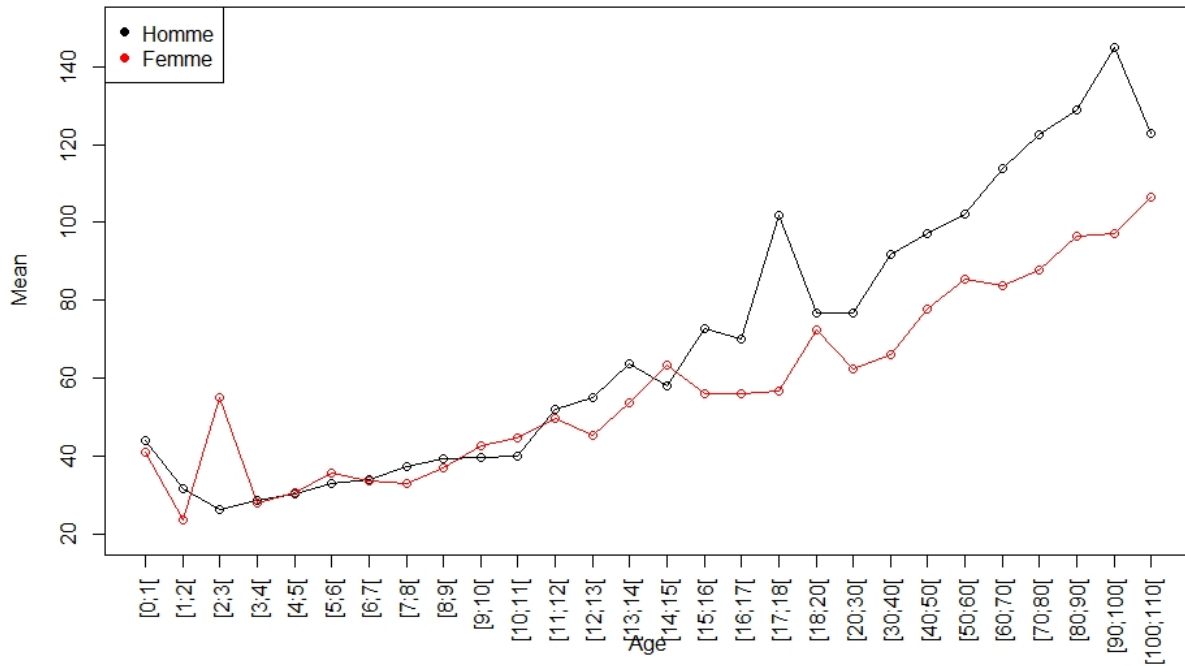
## A.2 Représentations graphiques des moyennes des paramètres biologiques par âge et sexe des patients calculées sur le jeu de données 2018



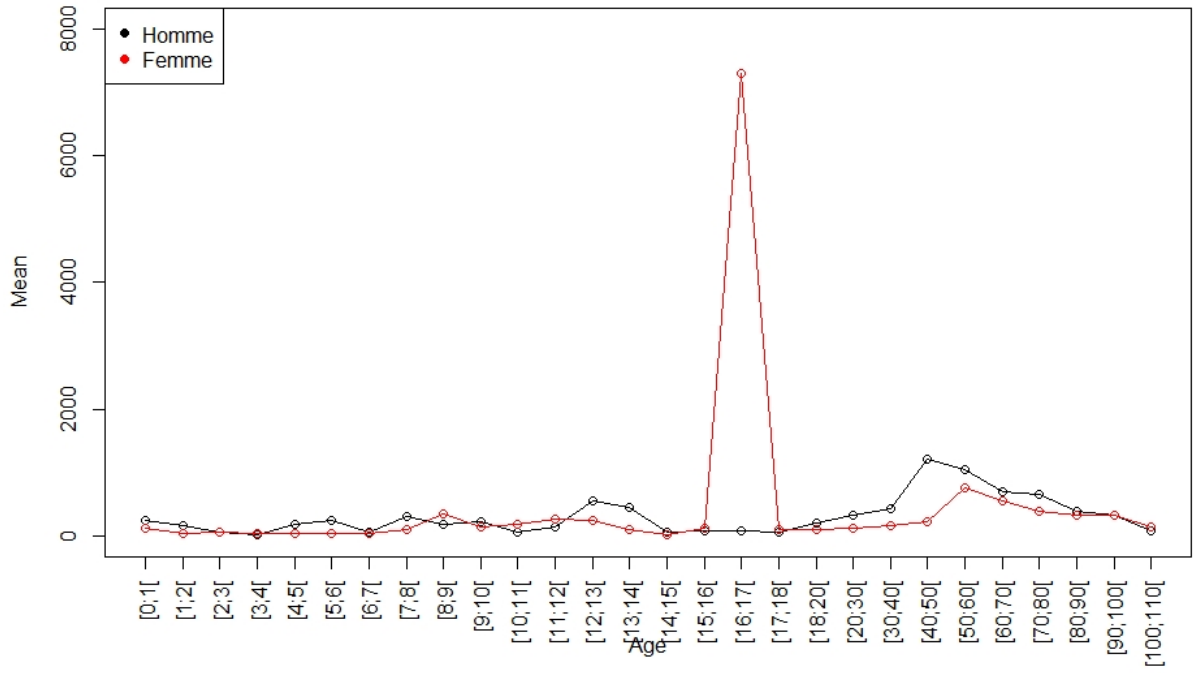
### Chlore



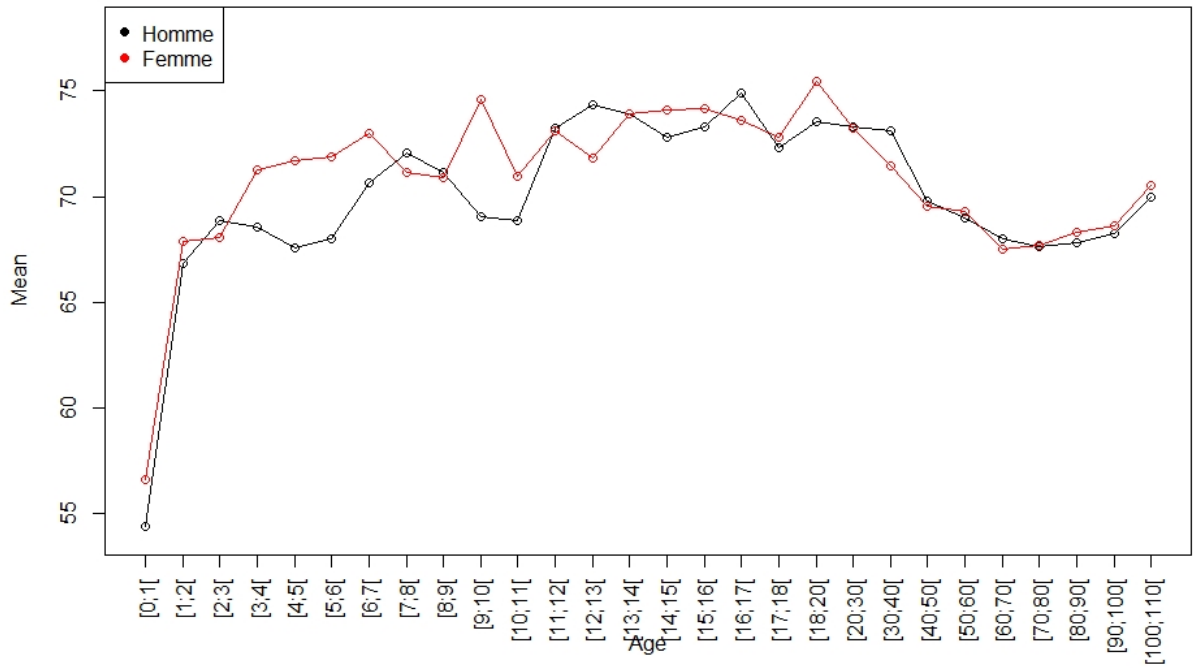
### Creatinine



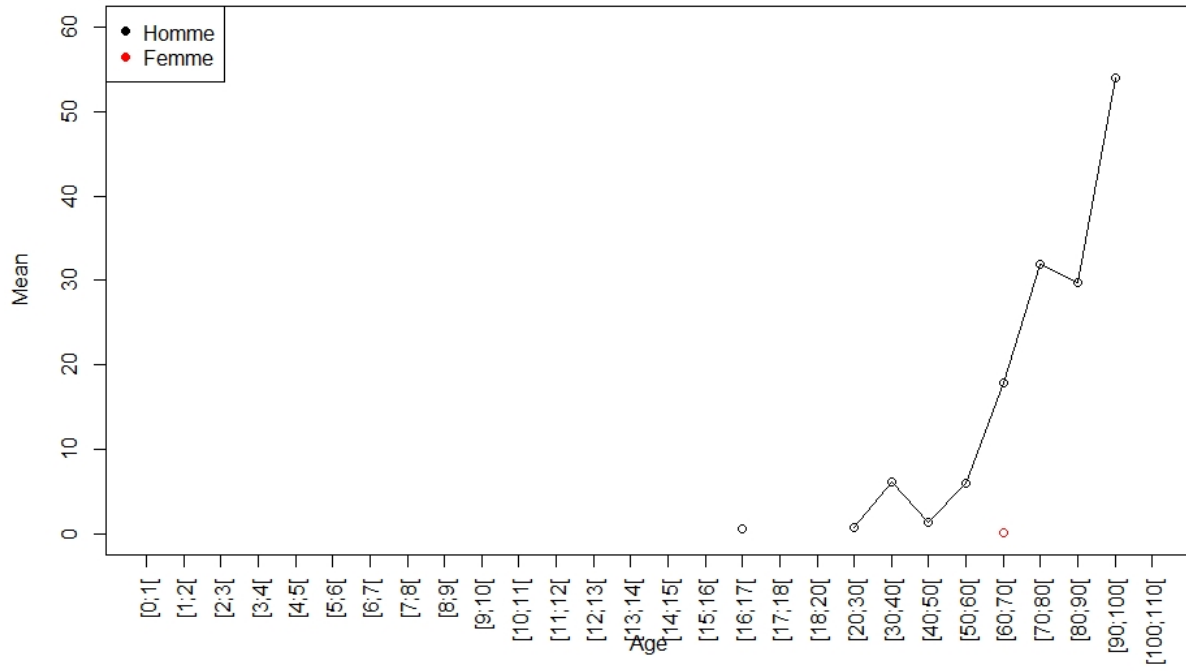
### Ferritine



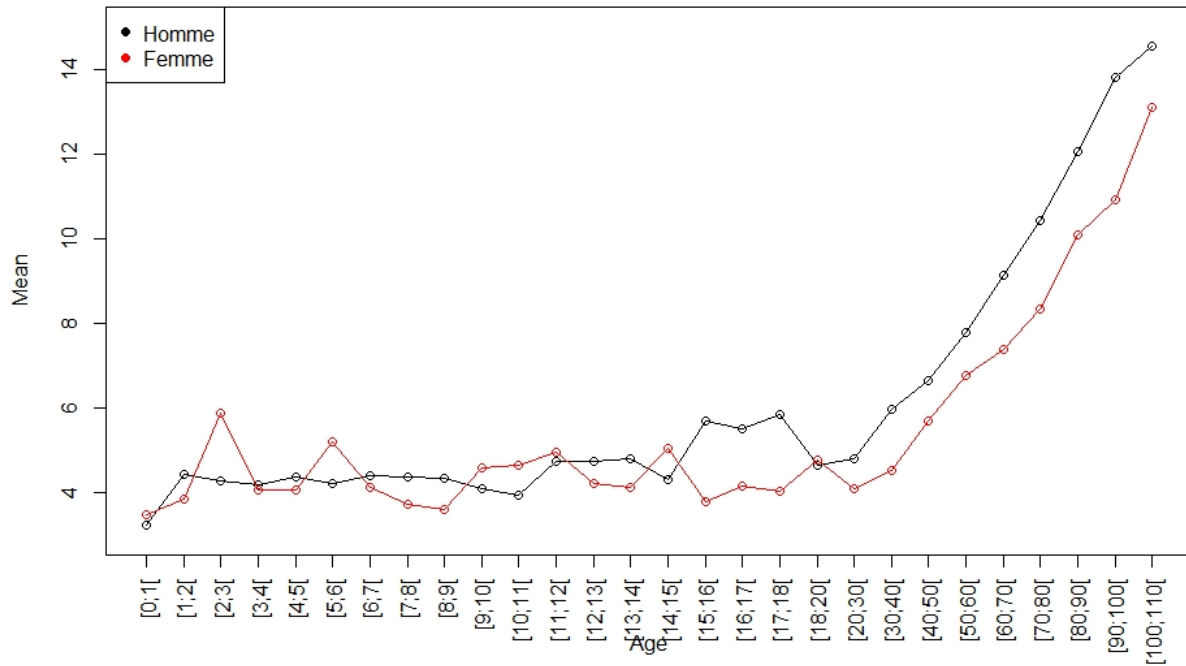
### Protéines



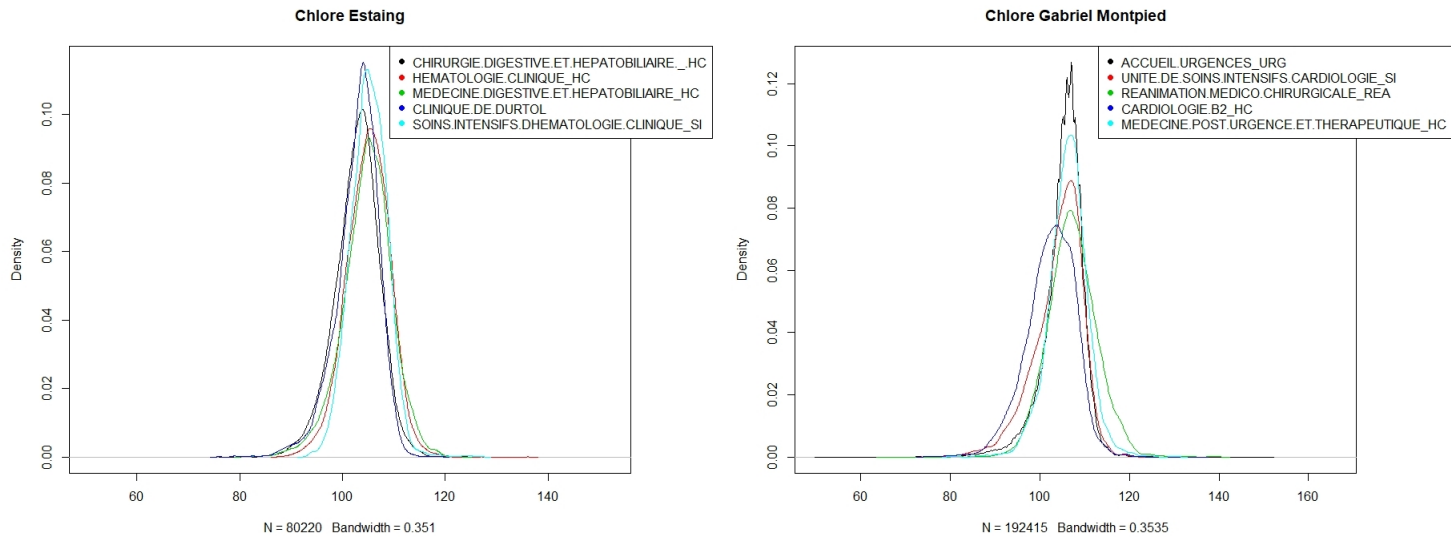
### PSA



### Uree



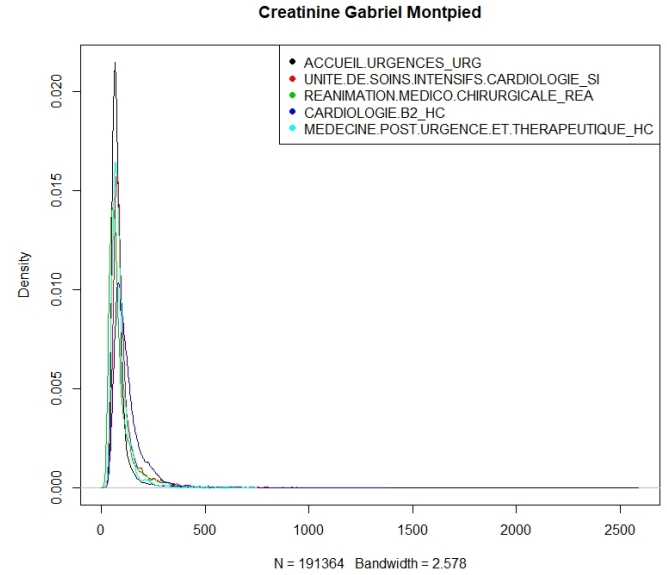
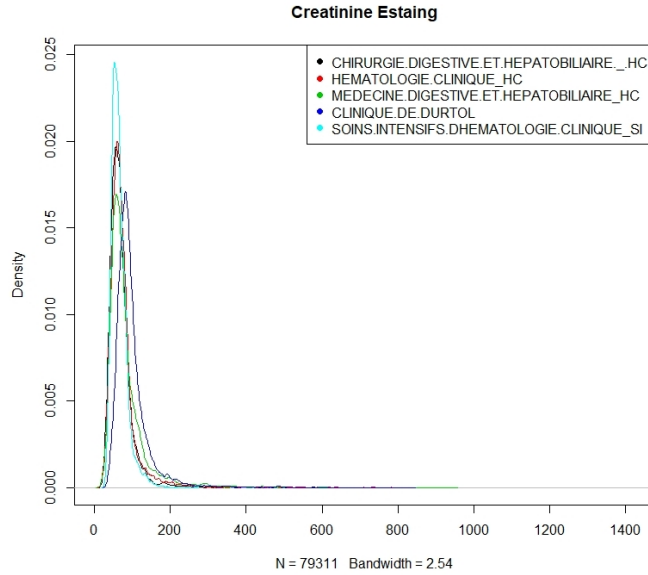
### A.3 Résumés et graphiques des distributions des paramètres biologiques pour les cinq services ayant fait le plus de mesures pour chaque paramètre en 2018 aux pôles Gabriel Montpied et Estaing



		Min	25%	50%	Moyenne	75%	Max
Gabriel Montpied	ACCUEIL.URGENCES._URG	51,00	103,00	106,00	105,25	108,00	151,00
	UNITE.DE.SOINS.INTENSIFS.CARDIOLOGIE._SI	73,00	101,00	105,00	104,08	108,00	125,00
	REANIMATION.MEDICO.CHIRURGICALE._REA	66,00	104,00	107,00	107,25	111,00	140,00
	CARDIOLOGIE.B2._HC	75,00	99,00	103,00	102,42	106,00	135,00
	MEDECINE.POST.URGENCE.ET.THERAPEUTIQUE._HC	80,00	104,00	106,00	106,24	109,00	131,00
Estaing	CHIRURGIE.DIGESTIVE.ET.HEPATOBILIAIRE._HC	78,00	100,00	103,00	102,82	106,00	124,00
	HEMATOLOGIE.CLINIQUE._HC	88,00	102,00	105,00	104,98	108,00	136,00
	MEDECINE.DIGESTIVE.ET.HEPATOBILIAIRE._HC	81,00	102,00	105,00	104,78	108,00	123,00
	CLINIQUE.DE.DURTOL	76,00	101,00	103,00	102,89	106,00	119,00
	SOINS.INTENSIFS.DHEMATOLOGIE.CLINIQUE._SI	93,00	103,00	105,00	105,17	108,00	127,00

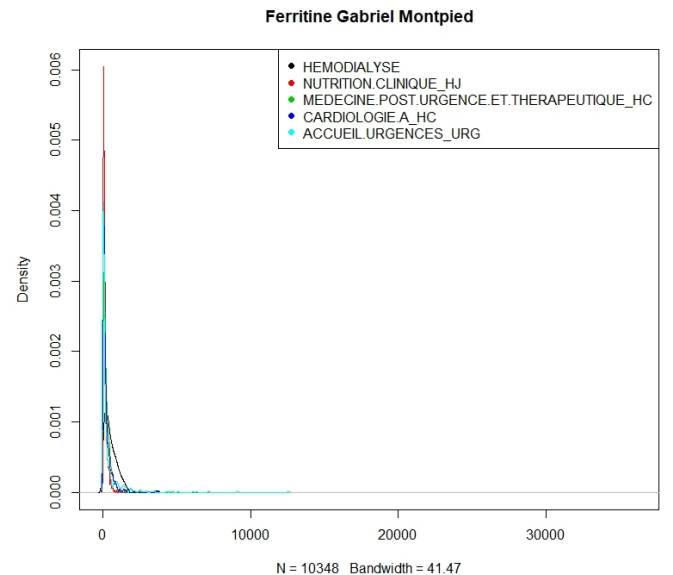
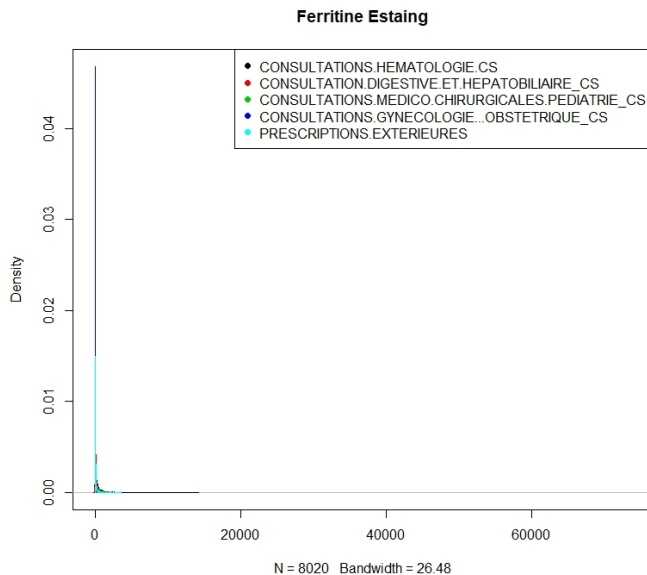
TABLE A.1 – Chlore





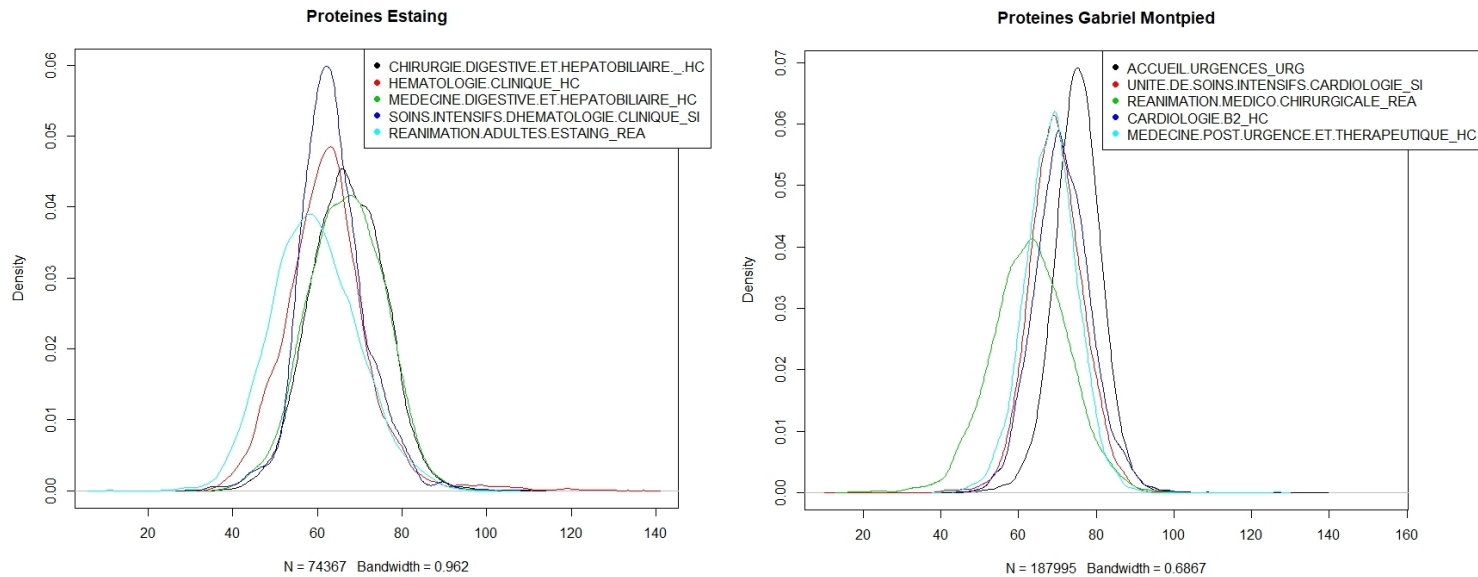
		Min	25%	50%	Moyenne	75%	Max
Gabriel Montpied	ACCUEIL.URGENCES_URG	17,20	60,70	72,70	85,37	88,70	2580,00
	UNITE.DE.SOINS.INTENSIFS.CARDIOLOGIE_SI	25,80	70,30	86,00	108,41	114,00	1150,00
	REANIMATION.MEDICO.CHIRURGICALE_REA	15,30	50,60	68,20	90,21	99,30	881,00
	CARDIOLOGIE.B2_HC	30,80	79,80	105,00	125,33	146,00	1360,00
	MEDECINE.POST.URGENCE.ET.THERAPEUTIQUE_HC	25,50	62,30	77,60	90,07	99,90	737,00
Estaing	CHIRURGIE.DIGESTIVE.ET.HEPATOBIILAIRE_.HC	18,50	50,40	63,30	69,65	78,60	876,00
	HEMATOLOGIE.CLINIQUE_HC	17,40	52,50	65,30	74,73	81,10	782,00
	MEDECINE.DIGESTIVE.ET.HEPATOBIILAIRE_HC	17,20	53,60	68,40	83,56	92,43	943,00
	CLINIQUE.DE.DURTOL	32,90	72,10	87,00	98,36	108,00	834,00
	SOINS.INTENSIFS.DHEMATOLOGIE.CLINIQUE_SI	27,60	50,90	61,20	65,26	74,60	266,00

TABLE A.2 – Créatinine



		Min	25%	50%	Moyenne	75%	Max
Gabriel Montpied	HEMODIALYSE	11,20	193,00	415,80	543,61	789,00	3694,00
	NUTRITION.CLINIQUE_HJ	3,90	39,83	85,10	130,39	162,93	1466,10
	MEDECINE.POST.URGENCE.ET.THERAPEUTIQUE_HC	3,50	65,80	146,60	296,73	326,55	7156,90
	CARDIOLOGIE.A_HC	3,00	34,38	96,35	216,36	264,28	3768,60
Estaing	ACCUEIL.URGENCES_URG	0,70	12,05	55,30	356,14	266,40	12540,20
	CONSULTATIONS.HEMATOLOGIE.CS	1,60	59,40	168,90	401,17	402,20	14030,30
	CONSULTATION.DIGESTIVE.ET.HEPATOBLIAIRE_CS	2,30	37,15	92,70	171,84	202,35	1819,60
	CONSULTATIONS.MEDICO.CHIRURGICALES.PEDIATRIE_CS	2,00	17,90	29,50	81,74	57,10	1813,30
	CONSULTATIONS.GYNECOLOGIE...OBSTETRIQUE_CS	1,50	7,30	12,70	26,86	26,70	716,90
	PRESCRIPTIONS.EXTERIEURES	1,70	16,90	31,10	75,07	72,70	3481,00

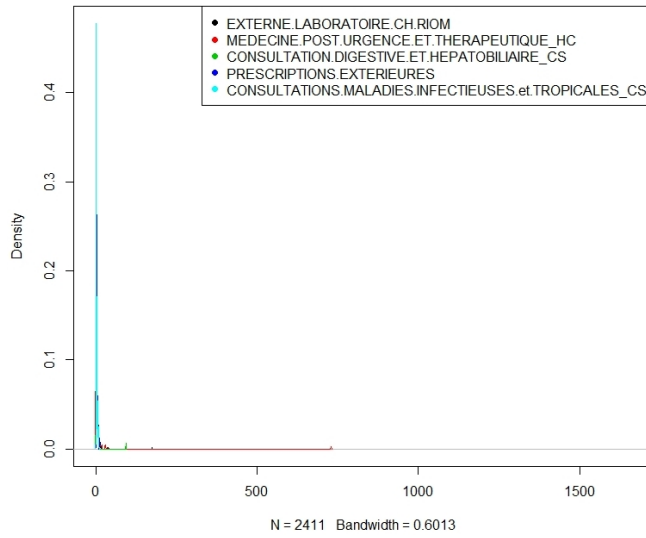
TABLE A.3 – Ferritine



		Min	25%	50%	Moyenne	75%	Max
Gabriel Montpied	ACCUEIL.URGENCES_URG	37,00	71,40	75,30	75,23	79,20	137,70
	UNITE.DE.SOINS.INTENSIFS.CARDIOLOGIE_SI	12,90	65,30	69,60	69,74	74,10	99,40
	REANIMATION.MEDICO.CHIRURGICALE_REA	17,60	56,70	63,10	63,02	69,65	100,50
	CARDIOLOGIE.B2_HC	46,80	66,60	71,10	71,40	76,00	101,50
	MEDECINE.POST.URGENCE.ET.THERAPEUTIQUE_HC	41,00	64,03	68,60	68,40	72,80	126,70
Estaing	CHIRURGIE.DIGESTIVE.ET.HEPATOBLIAIRE_.HC	31,30	61,20	67,00	67,16	73,10	110,00
	HEMATOLOGIE.CLINIQUE_HC	35,00	56,10	62,00	62,33	67,50	137,10
	MEDECINE.DIGESTIVE.ET.HEPATOBLIAIRE_HC	31,80	60,60	66,90	66,73	73,13	98,40
	SOINS.INTENSIFS.DHEMATOLOGIE.CLINIQUE_SI	30,00	58,70	63,00	63,65	68,00	106,80
	REANIMATION.ADULTES.ESTAING_REA	11,00	52,30	59,00	59,40	66,20	102,70

TABLE A.4 – Protéines

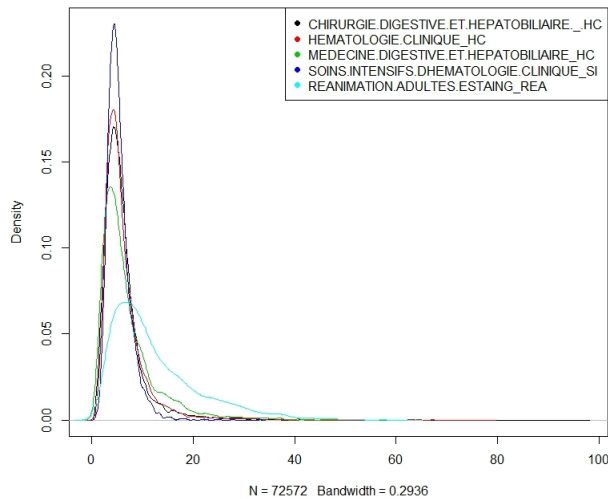
**PSA Gabriel Montpied**



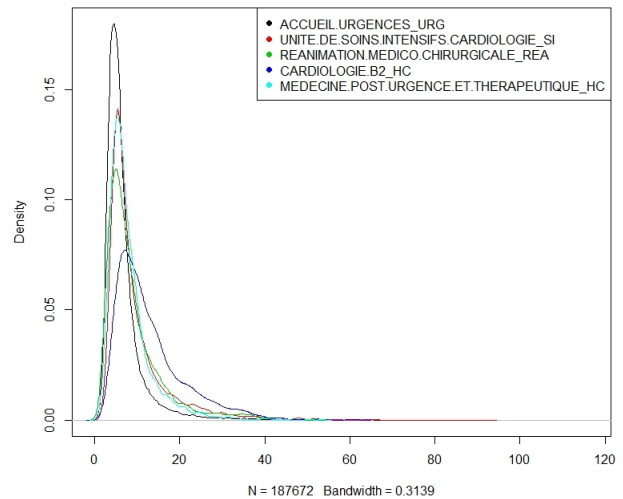
		Min	25%	50%	Moyenne	75%	Max
Gabriel Montpied	EXTERNE.LABORATOIRE.CH.RIOM	0,01	0,78	1,78	3,34	3,80	174,00
	MEDECINE.POST.URGENCE.ET.THERAPEUTIQUE_HC	0,01	0,71	1,44	8,50	3,27	730,00
	CONSULTATION.DIGESTIVE.ET.HEPATOBIILAIRE_CS	0,10	0,50	1,20	2,55	2,36	93,40
	PRESCRIPTIONS.EXTERIEURES	0,05	0,76	1,31	1,93	2,06	9,91
	CONSULTATIONS.MALADIES.INFECTIEUSES.et.TROPICALES_CS	0,01	0,53	0,94	1,87	1,88	9,09

TABLE A.5 – PSA

**Uree Estaing**



**Uree Gabriel Montpied**



		Min	25%	50%	Moyenne	75%	Max
Gabriel Montpied	ACCUEIL.URGENCES_URG	0,50	4,10	5,50	6,78	7,60	81,40
	UNITE.DE.SOINS.INTENSIFS.CARDIOLOGIE_SI	1,10	5,20	7,00	9,62	11,00	92,60
	REANIMATION.MEDICO.CHIRURGICALE_REA	0,60	4,70	7,00	9,02	10,80	56,70
	CARDIOLOGIE.B2_HC	1,50	7,20	10,60	12,97	16,20	61,60
	MEDECINE.POST.URGENCE.ET.THERAPEUTIQUE_HC	0,70	4,90	6,70	8,26	9,70	53,90
Estaing	CHIRURGIE.DIGESTIVE.ET.HEPATOBLIAIRE._HC	0,40	3,70	5,20	6,11	7,10	97,00
	HEMATOLOGIE.CLINIQUE_HC	0,50	3,70	5,10	6,36	7,30	78,40
	MEDECINE.DIGESTIVE.ET.HEPATOBLIAIRE_HC	0,50	3,40	5,30	7,20	8,70	59,30
	SOINS.INTENSIFS.DHEMATOLOGIE.CLINIQUE_SI	0,70	3,80	4,90	5,35	6,40	33,30
	REANIMATION.ADULTES.ESTAING_REA	0,50	6,10	9,90	12,45	16,40	61,40

TABLE A.6 – Urée

#### A.4 Résultats pour les paramètres chlore, créatinine, ferritine, protéines, PSA et urée ; méthode avec l'information service.

##### Chlore

	Méthode	N(2.4; 1)				N(4; 1)			
		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors Ligne	BinMoy	4	4	4	4	24	25	26	27
	BinVar	1	2	3	0	6	8	8	0
	SegMoy	10	12	13	21	49	56	60	63
	SegVar	3	5	7	8	7	10	12	3
	Pelt	35	41	45	39	70	76	79	67
	Pett	24	26	27	37	38	39	39	53
	Br	9	12	15	17	58	62	65	69
En Ligne	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
	Stu	26	35	37	9.96	63	72	75	15.53
	Bar	12	15	16	6.42	8	10	12	6.59
	GLR	21	24	28	8.01	25	32	37	10.61
	MW	35	43	47	13.39	71	77	78	14.44
	Moo	12	18	21	3.28	23	25	30	5.5
	LP	35	42	45	9.22	70	78	81	9.69
	KS	41	49	56	11.99	69	77	80	10.66
	CVM	36	42	46	12.78	73	78	79	13.51

TABLE A.7 – Résultats des détections de ruptures sur 100 simulations de Chlore, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec l'information service

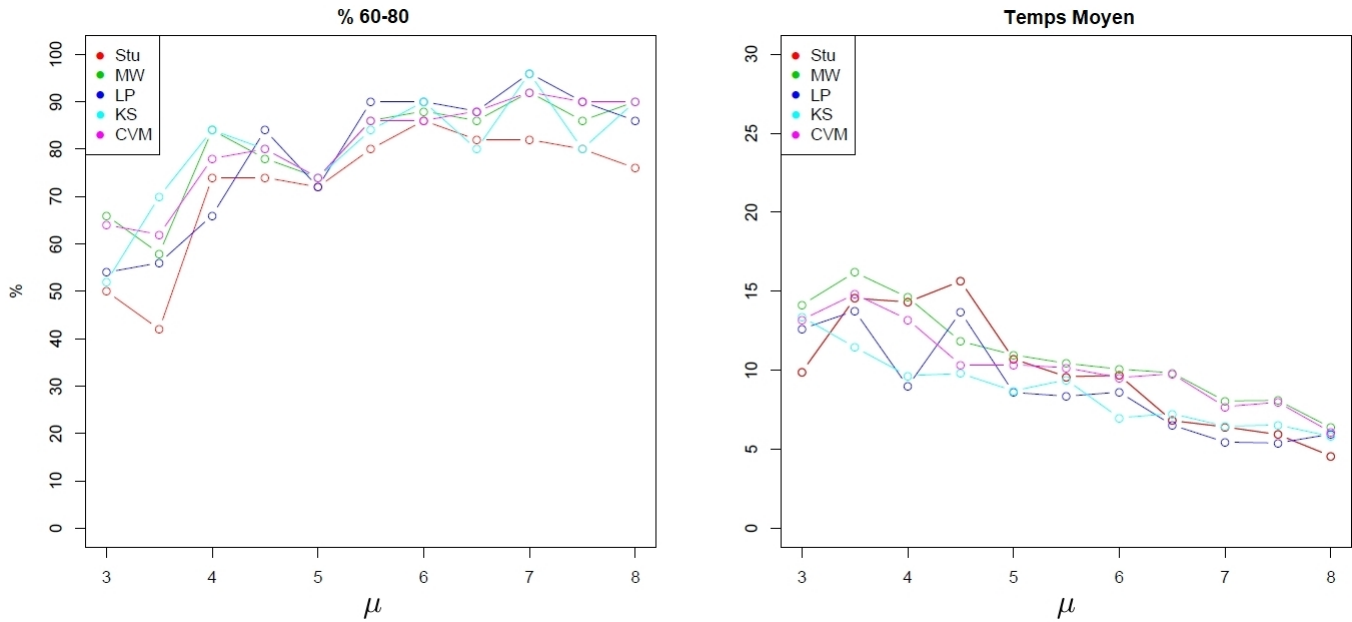


FIGURE A.1 – Pourcentages de simulations de Chlore avec ajout d'un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 3 à 6 par pas de 0.5 ; Méthode avec l'information service

		$N(0; 2.4)$				$N(0; 4)$			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	0	0	0	0	0	0	0	0
Ligne	BinVar	2	3	4	2	6	9	9	3
	SegMoy	8	8	10	0	8	15	18	6
	SegVar	2	5	6	2	15	18	20	11
	Pelt	11	13	16	8	13	18	23	13
	Pett	4	10	12	13	11	16	20	12
	Br	0	0	0	0	0	0	0	0
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	15	17	24	5.67	14	20	25	3
	Bar	8	9	10	5.46	7	10	12	7.03
	GLR	9	12	13	4.55	10	10	12	6.85
	MW	4	5	7	2.55	5	7	11	3.22
	Moo	7	12	16	2.72	24	29	33	7.11
	LP	3	4	8	2.78	23	25	30	4.04
	KS	4	6	9	4.38	14	14	16	4.27
CVM	4	5	8	4.03	5	7	12	3.84	

TABLE A.8 – Résultats des détections de ruptures sur 100 simulations de Chlore, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$ ; méthode avec l'information service

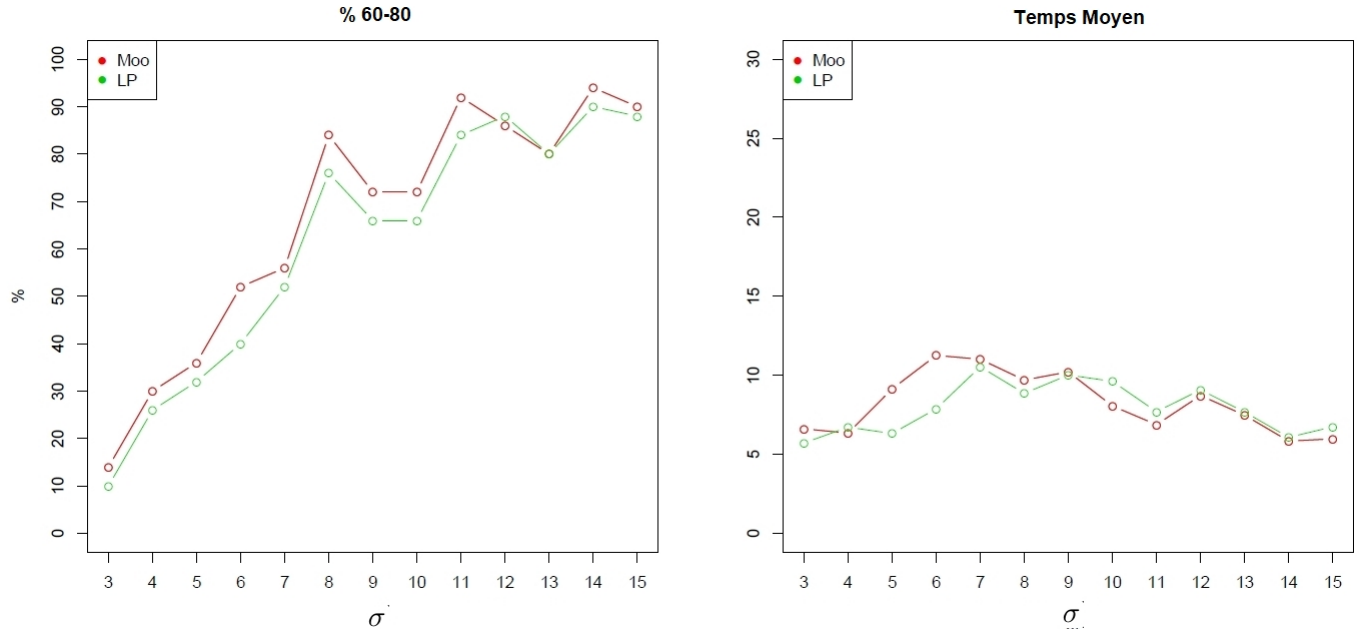


FIGURE A.2 – Pourcentages de simulations de Chlore avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 3 à 15 ; Méthode avec l'information service

Méthode		$\mathcal{N}(0, 5T; 1)$			
		%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	34	82	99	100
Ligne	BinVar	12	71	99	91
	SegMoy	5	73	99	100
	SegVar	3	50	98	95
	Pelt	74	96	96	98
	Pett	88	92	93	98
	Br	72	99	99	100
	Méthode		%60-80	%60-90	%60-100
En Ligne	Stu	71	78	78	8.82
	Bar	22	23	23	21.57
	GLR	45	50	50	7.35
	MW	84	91	91	9.73
	Moo	58	91	95	7.53
	LP	75	89	89	7.16
	KS	81	89	89	7.98
	CVM	82	91	91	9

TABLE A.9 – Résultats des détections de ruptures sur 100 simulations de Chlore, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec l'information service

Méthode	0	1	2	3	4	5	6	7	8
Stu	33	22	23,5	13,5	7	1	0	0	0
Bar	3	2,5	15,5	8,5	17,5	16	16	9	4
GLR	5	2	20,5	6,5	22,5	14	16	4	4
MW	68	11	15,5	3,5	1,5	0,5	0	0	0
Moo	69	16	13,5	1,5	0	0	0	0	0
LP	60	21,5	14,5	2,5	1	0,5	0	0	0
KS	62,5	21,5	11,5	2,5	1	0,5	0,5	0	0
CVM	67	10,5	16,5	3	2	0	1	0	0

TABLE A.10 – Fausses alarmes : Pourcentages de simulations où les méthodes de détection en ligne (avec services) ont trouvé 0,1,2,... ruptures sur 200 simulations de chlore de 190 individus

## Créatinine

	Méthode	$N(4.6; 1)$				$N(17.1; 1)$			
		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	0	0	0	1	5	5	5	3
Ligne	BinVar	4	6	9	1	4	5	5	0
	SegMoy	3	4	7	5	26	31	35	24
	SegVar	4	9	17	4	7	10	14	4
	Pelt	10	12	14	6	66	68	68	58
	Pett	15	19	22	11	39	40	44	44
	Br	1	1	1	1	27	29	33	31
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	19	25	28	3.08	46	52	57	11.25
	Bar	8	12	15	10.39	12	17	21	9.67
	GLR	9	14	17	12.52	25	34	38	13.15
	MW	5	7	10	6.13	70	72	74	21.48
	Moo	5	5	6	5.13	7	8	9	10.04
	LP	4	6	8	3.37	53	59	63	22.6
	KS	6	10	14	7.08	70	82	82	14.6
	CVM	6	8	10	5.46	68	71	73	20.65

TABLE A.11 – Résultats des détections de ruptures sur 100 simulations de Créatinine, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec l'information service

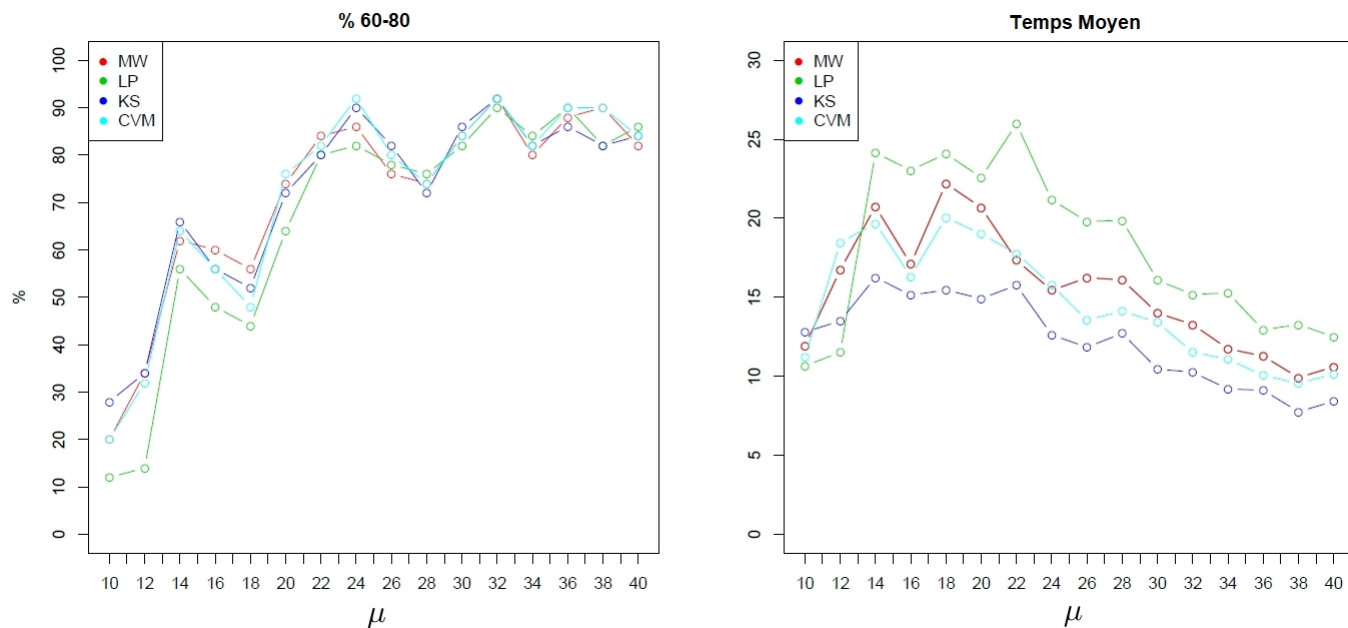


FIGURE A.3 – Pourcentages de simulations de Créatinine avec ajout d'un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 10 à 40 ; Méthode avec l'information service

		N(0 ; 4.6)				N(0 ; 17.1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	1	2	2	0	0	0	0	0
Ligne	BinVar	4	5	6	0	4	4	6	2
	SegMoy	4	6	7	1	5	8	11	4
	SegVar	7	7	8	1	9	11	13	4
	Pelt	8	11	15	4	17	21	27	20
	Pett	15	22	27	13	10	13	20	13
	Br	0	1	2	0	0	1	2	0
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	11	14	18	7.55	7	12	16	3.28
	Bar	13	15	22	13.11	13	13	16	12.53
	GLR	13	16	25	12.16	12	15	17	10.65
	MW	5	10	12	5.89	5	10	13	4.34
	Moo	7	9	10	3.38	16	22	26	7.47
	LP	8	12	12	3.17	10	19	22	3.52
	KS	4	8	9	6.25	5	9	12	4.78
	CVM	7	12	13	5.2	3	8	12	3.81

TABLE A.12 – Résultats des détections de ruptures sur 100 simulations de Créatinine, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$  ; méthode avec l'information service



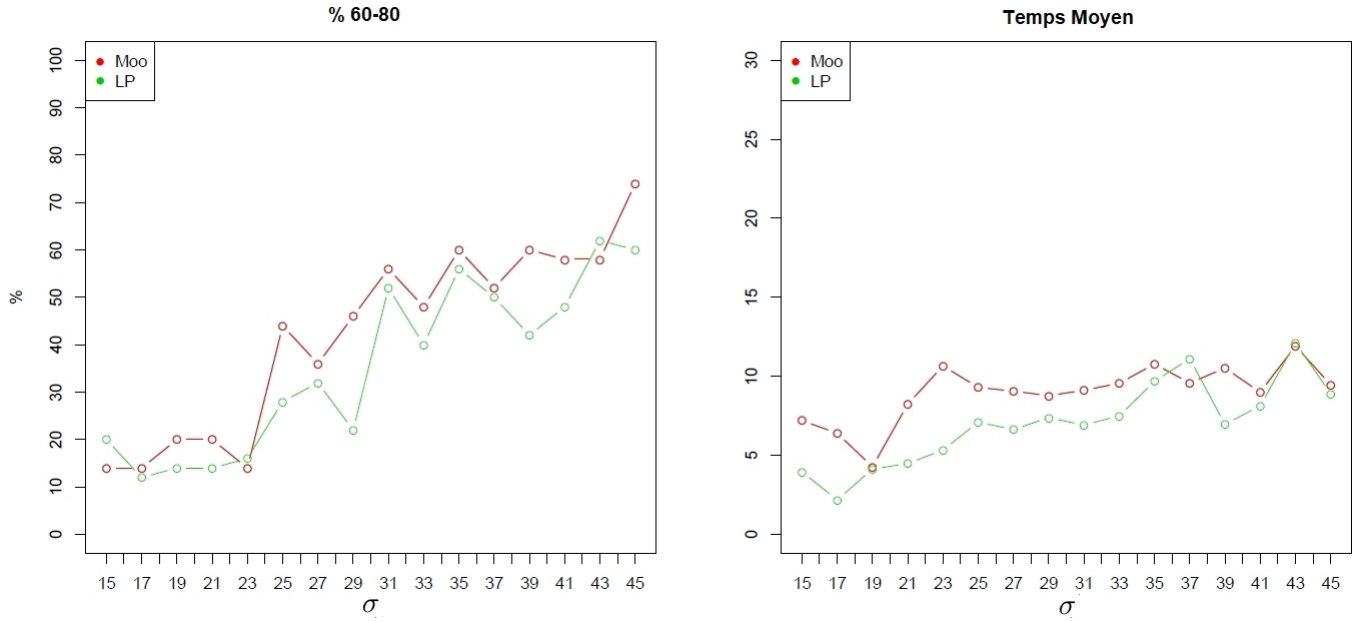


FIGURE A.4 – Pourcentages de simulations de Créatinine avec ajout d’un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 15 à 45 par pas de 2 ; Méthode avec l’information service

Méthode		$\mathcal{N}(0.5T; 1)$			
		%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	0	2	2	3
Ligne	BinVar	3	8	8	2
	SegMoy	3	12	19	20
	SegVar	5	11	13	7
	Pelt	24	44	54	48
	Pett	25	40	47	32
	Br	3	9	14	17
	Méthode		%60-80	%60-90	%60-100
En Ligne	Stu	15	25	34	7.47
	Bar	14	17	19	9.79
	GLR	11	20	26	10.53
	MW	23	46	56	16.74
	Moo	6	7	8	5.32
	LP	17	34	40	16.49
	KS	21	41	55	15.65
	CVM	26	50	59	17.6

TABLE A.13 – Résultats des détections de ruptures sur 100 simulations de Créatinine, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec l’information service

Méthode	0	1	2	3	4	5	6	7	8
Stu	22	20	27	15	10	5	1	0	0
Bar	3,5	4,5	12,5	17	21	19,5	8,5	11	2
GLR	10	7	17	22	16,5	16,5	7	3,5	0,5
MW	70	14,5	12	3	0	0,5	0	0	0
Moo	73,5	14	9,5	1,5	1,5	0	0	0	0
LP	71,5	11	14	2,5	1	0	0	0	0
KS	66	22,5	8,5	2,5	0,5	0	0	0	0
CVM	70,5	14	13,5	1,5	0	0,5	0	0	0

TABLE A.14 – Fausses alarmes : Pourcentages de simulations où les méthodes de détection en ligne (avec services) ont trouvé 0,1,2,... ruptures sur 200 simulations de créatinine de 190 individus

## Ferritine

		N(7,6 ; 1)				N(40 ; 1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	0	0	1	0	0	0	0	0
Ligne	BinVar	0	0	1	0	1	2	2	0
	SegMoy	3	4	6	0	5	7	9	1
	Seg Var	1	2	3	0	2	3	4	0
	Pelt	4	6	12	6	18	22	24	14
	Pett	12	14	20	9	14	19	23	22
	Br	0	0	1	0	3	3	4	2

		%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	6	7	8	5.22	16	19	24	3.19
	Bar	11	14	19	7.26	10	14	18	7.81
	GLR	11	15	21	7.26	12	14	19	8.36
	MW	4	4	5	4.06	10	14	17	10.23
	Moo	7	8	11	7.37	6	7	9	8.51
	LP	9	9	12	5.16	11	12	15	9.41
	KS	4	5	9	4.35	22	27	30	12.15
	CVM	4	4	6	3.64	14	16	20	11.18

TABLE A.15 – Résultats des détections de ruptures sur 100 simulations de Ferritine, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec l'information service

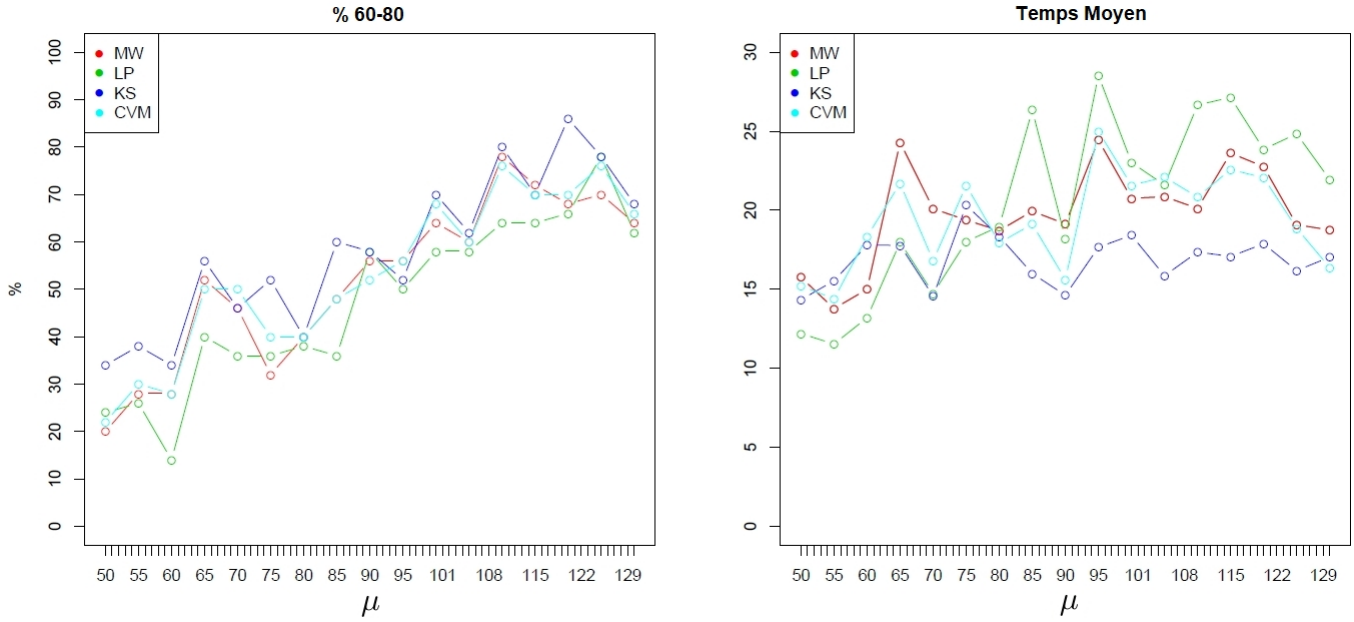


FIGURE A.5 – Pourcentages de simulations de Ferritine avec ajout d’un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 50 à 130 par pas de 5 ; Méthode avec l’information service

Méthode		N(0 ; 7.6)				N(0 ; 40)			
		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	0	0	0	0	2	2	2	0
Ligne	BinVar	2	2	3	0	0	1	2	0
	SegMoy	3	3	4	1	2	4	4	2
	SegVar	3	4	5	1	0	1	3	2
	Pelt	7	11	15	5	5	7	8	7
	Pett	10	17	19	11	17	23	26	9
	Br	0	0	1	0	2	2	2	2
	Méthode		%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100
En Ligne	Stu	9	11	14	2.62	10	13	13	3.51
	Bar	12	16	22	10	8	11	16	9.24
	GLR	11	14	18	8.53	9	13	18	8.58
	MW	1	3	5	4.46	3	6	6	5.27
	Moo	7	11	16	8.12	7	10	10	7.4
	LP	3	4	8	4.35	3	6	6	7.22
	KS	3	5	8	5.2	4	7	9	8.67
	CVM	1	3	4	3.1	4	9	9	4.81

TABLE A.16 – Résultats des détections de ruptures sur 100 simulations de Ferritine, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$  ; méthode avec l’information service

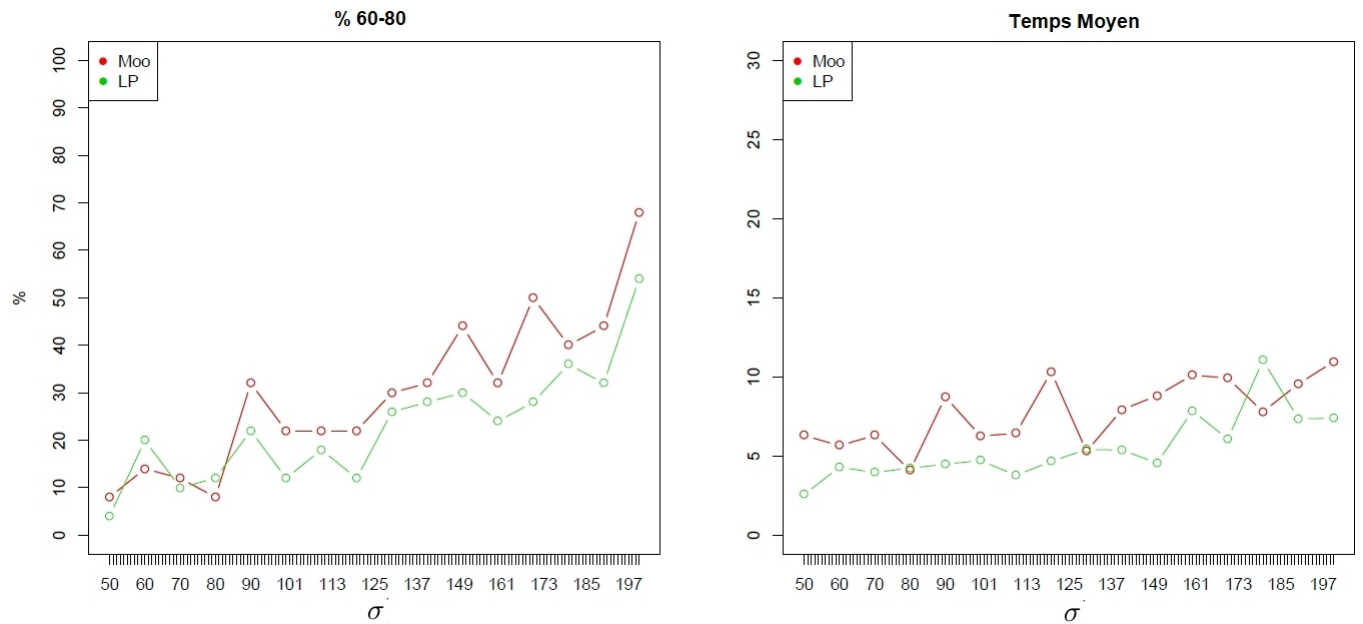


FIGURE A.6 – Pourcentages de simulations de Ferritine avec ajout d’un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 50 à 200 par pas de 10 ; Méthode avec l’information service

Méthode		$\mathcal{N}(0.5T; 1)$			
		%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	0	0	0	0
Ligne	BinVar	1	1	2	0
	SegMoy	4	6	8	2
	SegVar	2	3	7	1
	Pelt	5	6	12	9
	Pett	10	18	28	6
	Br	0	0	0	0
	Méthode		%60-80	%60-90	%60-100
En Ligne	Stu	9	11	16	4.65
	Bar	14	15	21	9.63
	GLR	13	14	20	9.1
	MW	4	5	6	4.12
	Moo	6	6	11	5.64
	LP	6	7	13	5.66
	KS	7	8	9	4.72
	CVM	6	8	10	4.42

TABLE A.17 – Résultats des détections de ruptures sur 100 simulations de Ferritine, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec l’information service

Méthode	0	1
Stu	98,5	1,5
Bar	98	2
GLR	96,5	3,5
MW	100	0
Moo	100	0
LP	100	0
KS	100	0
CVM	100	0

TABLE A.18 – Fausses alarmes : Pourcentages de simulations où les méthodes de détection en ligne (avec services) ont trouvé 0 ou 1 ruptures sur 200 simulations de ferritine de 20 individus

## Protéines

		N(3.5; 1)				N(8.37; 1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	6	7	7	6	66	66	66	66
Ligne	BinVar	2	3	3	0	1	1	1	1
	SegMoy	12	13	13	13	91	94	94	91
	SegVar	4	5	5	0	3	4	4	2
	Pelt	23	24	25	33	83	86	86	80
	Pett	23	25	28	30	57	57	58	70
	Br	13	15	15	17	92	95	95	92
		N(3.5; 1)				N(8.37; 1)			
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	30	33	35	13.56	81	85	85	9.9
	Bar	8	10	10	5.7	9	10	11	7.74
	GLR	16	18	19	11.84	73	77	77	13.86
	MW	30	36	36	14.94	80	82	82	12.51
	Moo	6	7	9	4.28	16	21	29	6.22
	LP	16	19	21	8.54	83	88	89	11.12
	KS	35	43	51	10.37	76	80	82	10.07
	CVM	31	37	37	13.32	80	84	85	11.55

TABLE A.19 – Résultats des détections de ruptures sur 100 simulations de Protéines, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec l'information service

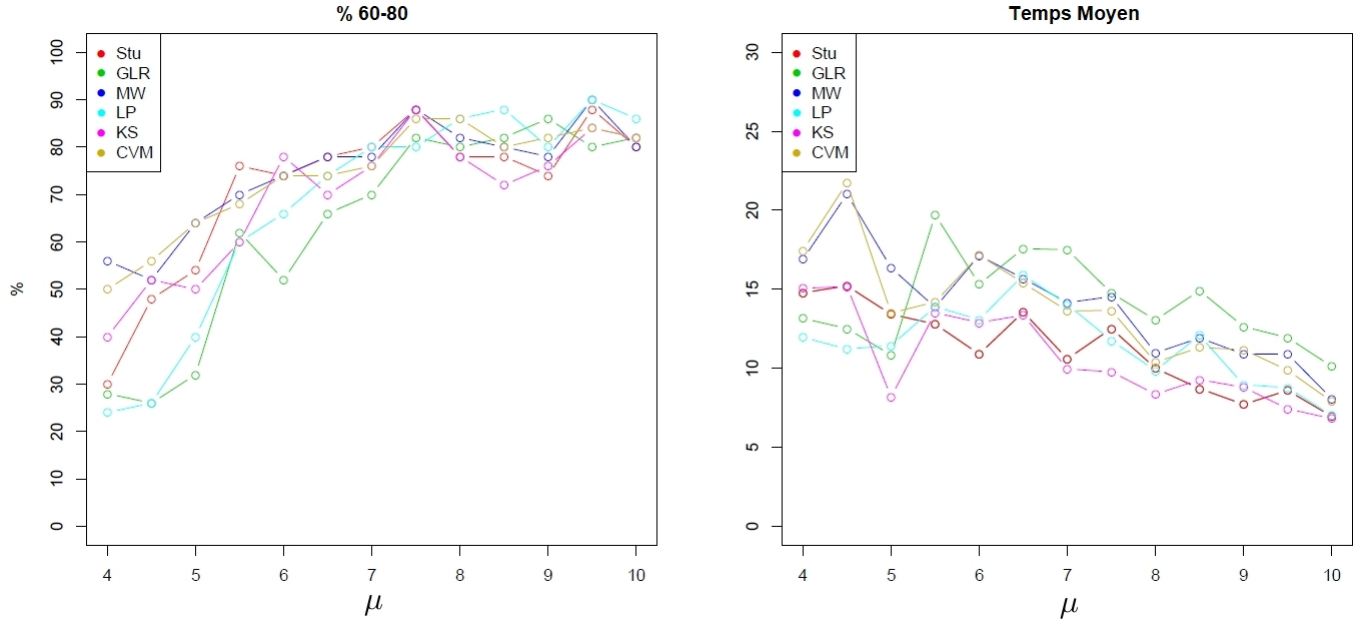


FIGURE A.7 – Pourcentages de simulations de Protéines avec ajout d’un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 4 à 10 par pas de 0.5 ; Méthode avec l’information service

Méthode		N(0 ; 3.5)				N(0 ; 8.37)			
		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors Ligne	BinMoy	0	0	0	0	0	1	1	0
	BinVar	0	0	2	0	28	31	31	21
	SegMoy	1	2	3	5	16	31	38	7
	SegVar	0	1	4	3	39	46	47	34
	Pelt	6	10	16	6	43	55	62	34
	Pett	7	15	23	8	16	20	24	10
	Br	0	0	0	0	1	4	5	2
Méthode		%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	6	9	17	2.67	32	47	55	2.87
	Bar	2	4	9	6.46	40	48	50	14.11
	GLR	2	2	2	6.17	31	39	41	9.43
	MW	2	3	7	4.51	8	14	20	3.85
	Moo	6	11	14	4.74	42	54	57	10.21
	LP	6	8	15	3.14	37	47	55	6.23
	KS	6	8	10	3.45	19	28	35	4.34
	CVM	2	3	7	3.42	11	15	21	3.75

TABLE A.20 – Résultats des détections de ruptures sur 100 simulations de Protéines, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$  ; méthode avec l’information service

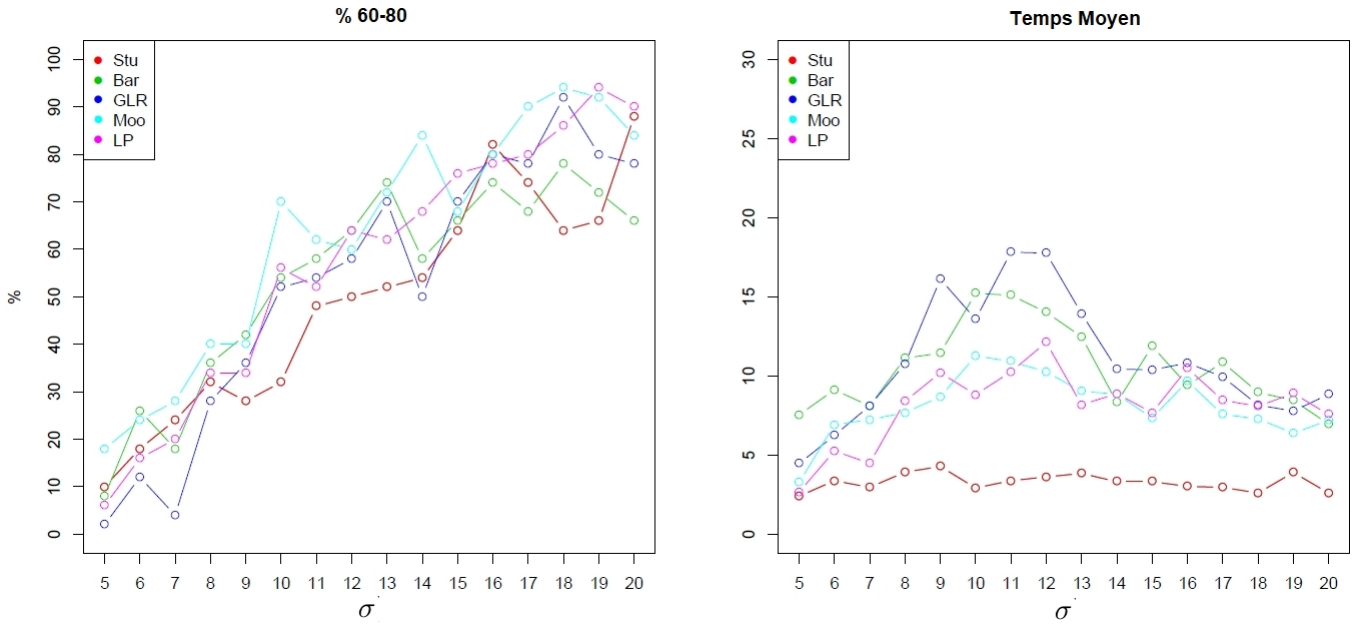


FIGURE A.8 – Pourcentages de simulations de Protéines avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 5 à 20 ; Méthode avec l'information service

		$\mathcal{N}(0.5T; 1)$			
Méthode		%60-80	%60-90	%60-100	%110-130
Hors Ligne	BinMoy	51	83	96	96
	BinVar	12	48	75	82
	SegMoy	16	74	96	100
	SegVar	6	38	81	95
	Pelt	47	84	91	97
	Pett	69	84	85	96
	Br	48	89	97	98
Méthode		%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	56	83	86	9.91
	Bar	27	32	36	34.16
	GLR	53	81	82	11.27
	MW	60	83	84	11.33
	Moo	13	46	78	7.65
	LP	49	84	89	9.05
	KS	45	82	83	8.95
CVM	59	83	84	11.44	

TABLE A.21 – Résultats des détections de ruptures sur 100 simulations de Protéines, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec l'information service

Méthode	0	1	2	3	4	5	6
Stu	60	16	16,5	4	3	0,5	0
Bar	45,5	10	32	4,5	6,5	0,5	1
GLR	51,5	7	31,5	4,5	4	1	0,5
MW	71	10,5	15	2	1,5	0	0
Moo	73,5	14	11,5	0,5	0,5	0	0
LP	75,5	10	12	1	1,5	0	0
KS	65	20	12	1,5	1,5	0	0
CVM	71	9,5	16,5	1	2	0	0

TABLE A.22 – Fausses alarmes : Pourcentages de simulations où les méthodes de détection en ligne (avec services) ont trouvé 0,1,2,... ruptures sur 200 simulations de protéines de 190 individus

## PSA

		N(6.8 ; 1)				N(50.62 ; 1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	97	97	97	95	98	100	100	100
Ligne	BinVar	74	78	79	77	100	100	100	100
	SegMoy	98	99	100	99	97	100	100	96
	SegVar	83	89	90	85	100	100	100	100
	Pelt	89	89	89	89	99	99	99	96
	Pett	84	85	85	88	93	93	93	98
	Br	100	100	100	99	100	100	100	100
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	60	61	61	4.55	57	57	57	4.96
	Bar	5	5	5	8.46	3	3	3	8.71
	GLR	8	8	8	9.05	5	5	5	9.19
	MW	84	84	84	8.68	90	90	90	5.94
	Moo	36	47	49	7.11	86	87	87	5.45
	LP	88	90	90	6.47	88	88	88	4.89
	KS	86	86	86	7.1	85	85	85	5.28
	CVM	84	84	84	8.42	89	89	89	5.85

TABLE A.23 – Résultats des détections de ruptures sur 100 simulations de PSA, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec l'information service



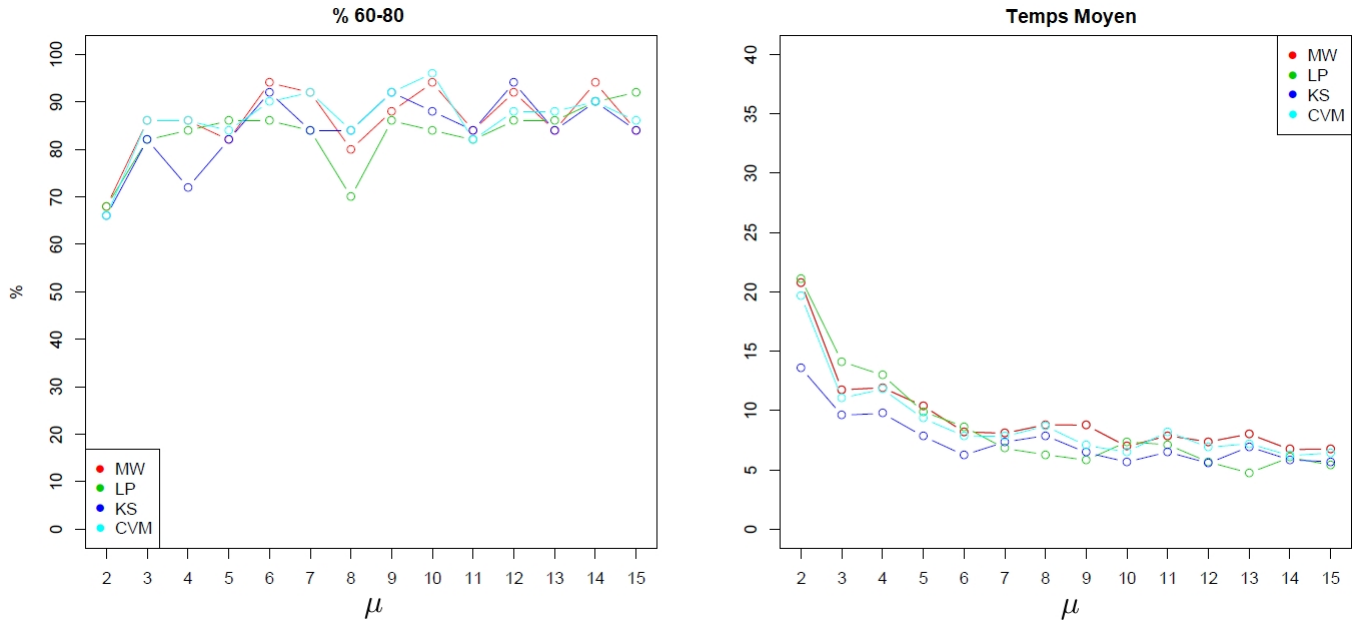


FIGURE A.9 – Pourcentages de simulations de PSA avec ajout d'un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 2 à 15 ; Méthode avec l'information service

		N(0 ; 6.8)				N(0 ; 50.62)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	3	7	8	4	28	44	65	31
Ligne	BinVar	72	75	75	65	100	100	100	100
	SegMoy	28	52	68	14	22	42	60	21
	SegVar	78	81	82	68	100	100	100	100
	Pelt	81	91	92	75	100	100	100	98
	Pett	17	19	22	14	24	30	30	16
	Br	7	12	15	8	23	32	44	14
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	46	56	60	3.02	61	61	61	3.65
	Bar	4	4	4	9.17	2	2	2	8.73
	GLR	9	9	9	9.13	4	4	4	9.04
	MW	21	31	39	4.9	33	50	60	7.02
	Moo	74	77	80	8.72	88	88	89	5
	LP	64	75	77	10.14	88	88	88	5.48
	KS	28	42	49	7.08	76	83	84	11.01
	CVM	22	33	41	6.34	64	82	86	12.48

TABLE A.24 – Résultats des détections de ruptures sur 100 simulations de PSA, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$  ; méthode avec l'information service

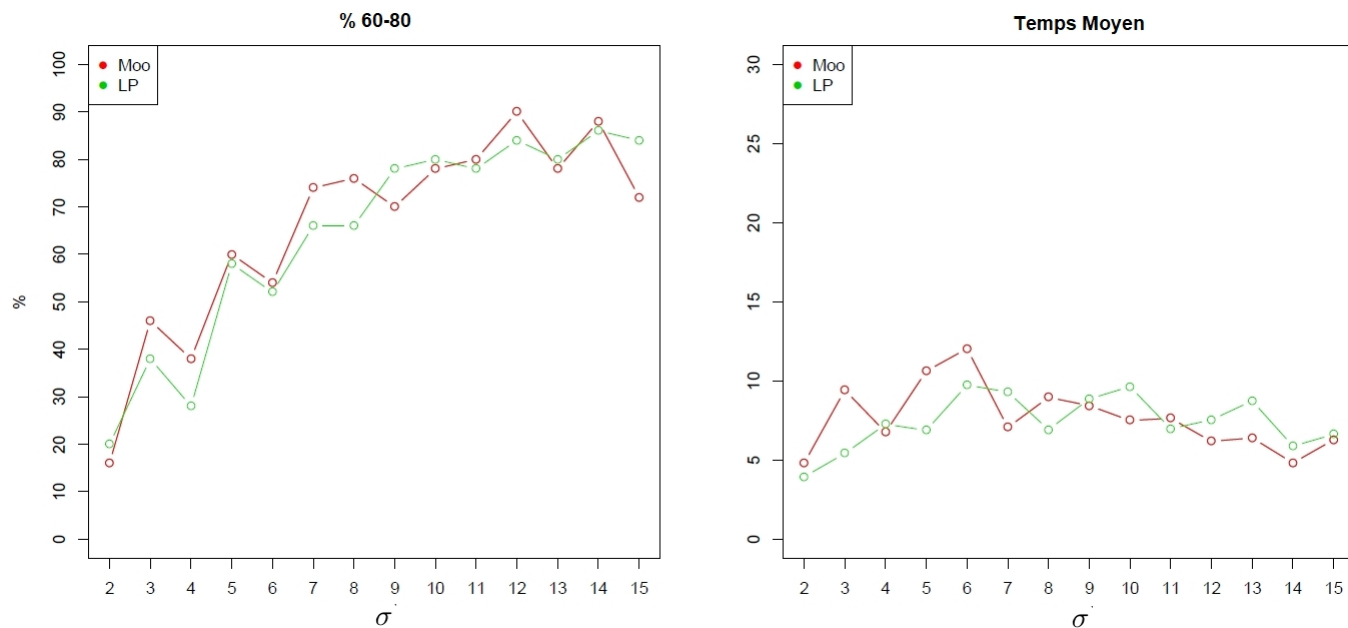


FIGURE A.10 – Pourcentages de simulations de PSA avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 2 à 15 ; Méthode avec l'information service

Méthode		$N(0.5T; 1)$			
		%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	44	72	94	100
Ligne	BinVar	18	80	100	100
	SegMoy	13	53	83	100
	SegVar	12	67	99	100
	Pelt	78	87	87	96
	Pett	84	86	86	93
	Br	59	93	99	100
	Méthode		%60-80	%60-90	%60-100
En Ligne	Stu	49	61	62	5.55
	Bar	5	5	5	9.32
	GLR	9	9	9	10.11
	MW	75	79	79	12.15
	Moo	24	57	78	7.15
	LP	68	80	81	9.75
	KS	74	78	78	9.42
	CVM	76	81	81	11.67

TABLE A.25 – Résultats des détections de ruptures sur 100 simulations de PSA, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec l'information service

Méthode	0	1
Stu	98,5	1,5
Bar	79	21
GLR	80,5	19,5
MW	100	0
Moo	99,5	0,5
LP	100	0
KS	100	0
CVM	100	0

TABLE A.26 – Fausses alarmes : Pourcentages de simulations où les méthodes de détection en ligne (avec services) ont trouvé 0,1,2,... ruptures sur 200 simulations de PSA de 20 individus

## Urée

		N(6,9 ; 1)				N(34,24 ; 1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	85	86	86	80	100	100	100	100
Ligne	BinVar	17	26	29	15	100	100	100	100
	SegMoy	93	94	96	93	100	100	100	100
	SegVar	25	35	40	22	100	100	100	100
	Pelt	75	76	76	77	97	97	97	95
	Pett	82	82	82	90	98	98	98	99
	Br	94	95	95	94	100	100	100	100
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	64	68	68	6.8	64	64	64	3.42
	Bar	6	8	9	13.18	19	19	19	9.87
	GLR	24	24	24	10.8	25	25	25	10.91
	MW	90	90	90	8.65	81	81	81	6.57
	Moo	10	13	15	7.83	81	81	81	4.91
	LP	80	81	81	9.67	85	85	85	4.07
	KS	86	87	87	6.87	80	80	80	5.48
	CVM	89	89	89	8.39	83	83	83	5.69

TABLE A.27 – Résultats des détections de ruptures sur 100 simulations de Urée, VISTA 2, avec ajout de bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec l'information service

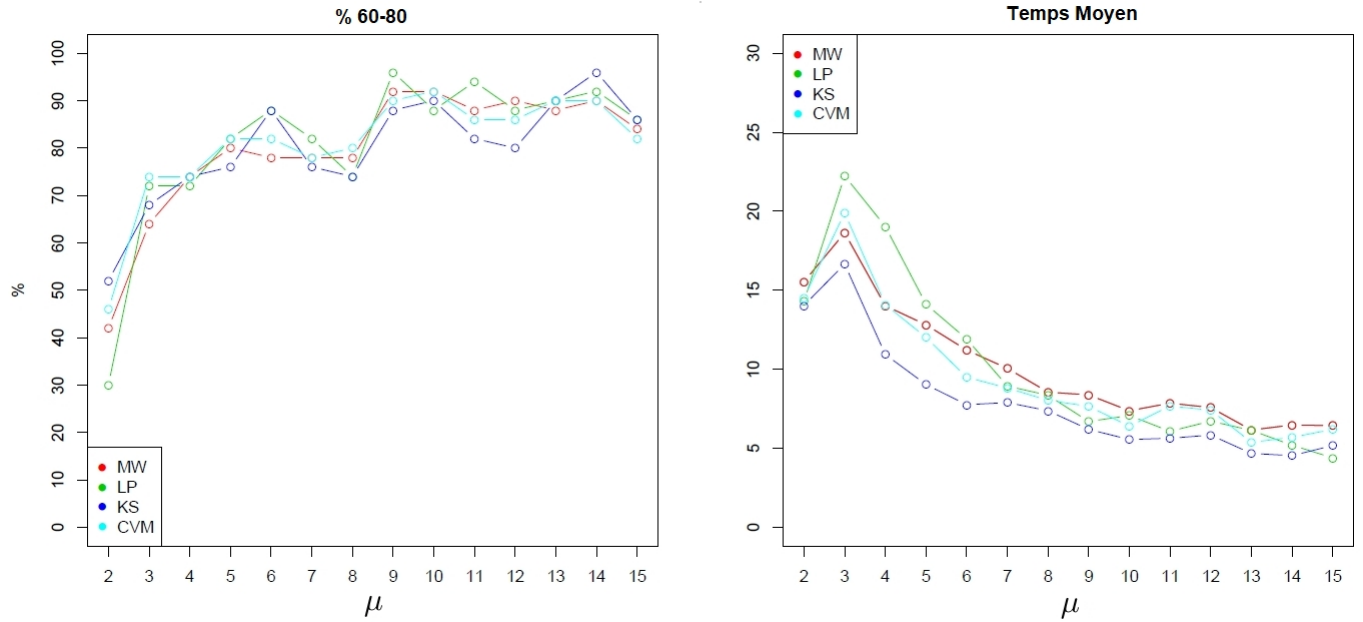


FIGURE A.11 – Pourcentages de simulations d’Urée avec un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne données ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 2 à 15 ; Méthode avec service

		$N(0; 6,9)$				$N(0; 34,24)$			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	0	1	1	0	28	45	68	34
Ligne	BinVar	46	49	49	37	100	100	100	98
	SegMoy	16	32	39	12	22	42	63	30
	Seg Var	44	51	53	40	100	100	100	99
	Pelt	79	91	92	75	100	100	100	100
	Pett	13	18	23	11	13	16	26	23
	Br	3	8	8	2	18	29	35	23
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	31	44	50	5.34	64	65	65	3.8
	Bar	12	17	17	13.02	19	19	19	9.26
	GLR	19	24	24	12.46	27	27	27	9.14
	MW	17	23	30	5.38	38	54	59	5.01
	Moo	68	83	83	10.51	90	90	90	4.84
	LP	63	78	83	9.13	91	91	91	5.27
	KS	37	49	58	11.23	85	87	87	9.8
	CVM	22	33	41	7.23	80	89	90	12.03

TABLE A.28 – Résultats des détections de ruptures sur 100 simulations de Urée, VISTA 2, avec ajout de bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$  ; méthode avec l’information service

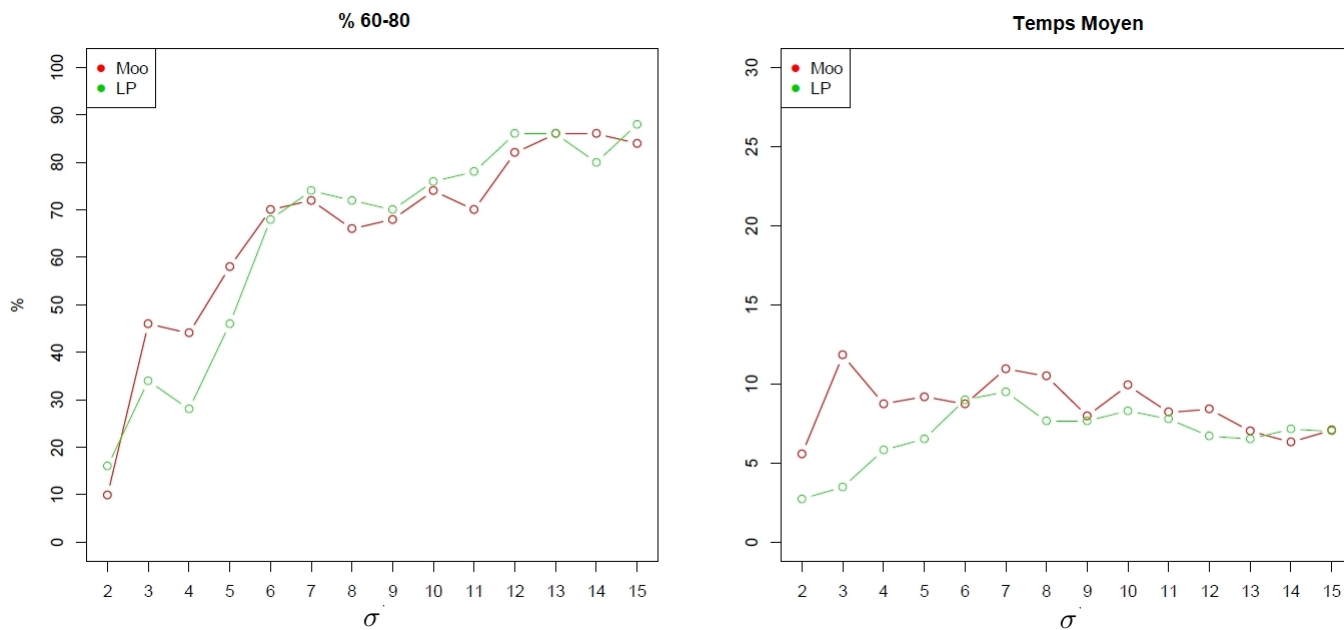


FIGURE A.12 – Pourcentages de simulations d'Urée avec un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne données ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 2 à 15 ; Méthode avec service

		$\mathcal{N}(0.5T; 1)$			
Méthode		%60-80	%60-90	%60-100	%110-130
Hors Ligne	BinMoy	43	77	97	97
	BinVar	13	55	92	97
	SegMoy	12	64	95	100
	SegVar	7	43	92	100
	Pelt	81	89	89	90
	Pett	94	94	94	91
	Br	66	98	99	96
Méthode		%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	56	66	66	6.62
	Bar	5	6	7	12.81
	GLR	17	19	19	10.46
	MW	90	90	90	12.02
	Moo	12	52	80	7.89
	LP	77	85	85	11.11
	KS	84	85	85	9.57
	CVM	87	87	87	11.5

TABLE A.29 – Résultats des détections de ruptures sur 100 simulations de Urée, VISTA 2, avec ajout de bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec l'information service

Méthode	0	1	2	3	4	5	6	7	8
Stu	18	20,5	23,5	17	13,5	5	1,5	1	0
Bar	0,5	2	7,5	9	14	16,5	15	18	9
GLR	1	2,5	7,5	13	15,5	17	17	17	8,5
MW	70,5	11,5	13,5	2,5	1	1	0	0	0
Moo	58,5	18,5	18	2,5	2	0	0,5	0	0
LP	65	18	11,5	3,5	1	1	0	0	0
KS	61,5	23,5	10,5	4	0,5	0	0	0	0
CVM	69	12,5	12,5	3,5	2	0,5	0	0	0

TABLE A.30 – Fausses alarmes : Pourcentages de simulations où les méthodes de détection en ligne (avec services) ont trouvé 0,1,2,... ruptures sur 200 simulations de Urée de 190 individus

### A.5 Résultats pour les paramètres chlore, créatinine, ferritine, protéines, PSA et urée ; méthode sans l'information service et avec des classifications.

#### Chlore

	Méthode	N(2,4; 1)				N(4; 1)			
		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	21	31	34	25	43	46	48	49
Ligne	BinVar	4	6	8	3	3	4	5	4
	SegMoy	21	34	39	28	48	55	57	50
	SegVar	4	7	9	3	3	3	5	4
	Pelt	32	41	46	31	54	67	68	54
	Pett	29	36	44	25	28	30	32	48
	Br	18	26	30	25	55	64	66	58
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	34	42	48	26.51	59	68	72	22.11
	Bar	29	32	32	28.97	25	27	33	24.6
	GLR	35	37	40	24.53	49	54	58	21.9
	MW	43	52	61	27.88	67	76	80	19.13
	Moo	16	24	26	13.99	33	38	42	19.52
	LP	32	41	46	14.06	64	72	78	18.16
	KS	40	52	64	25.27	65	78	85	15.29
	CVM	38	46	56	27.54	66	75	82	19.99

TABLE A.31 – Résultats des détections de ruptures sur 100 simulations de Chlore, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec classifications

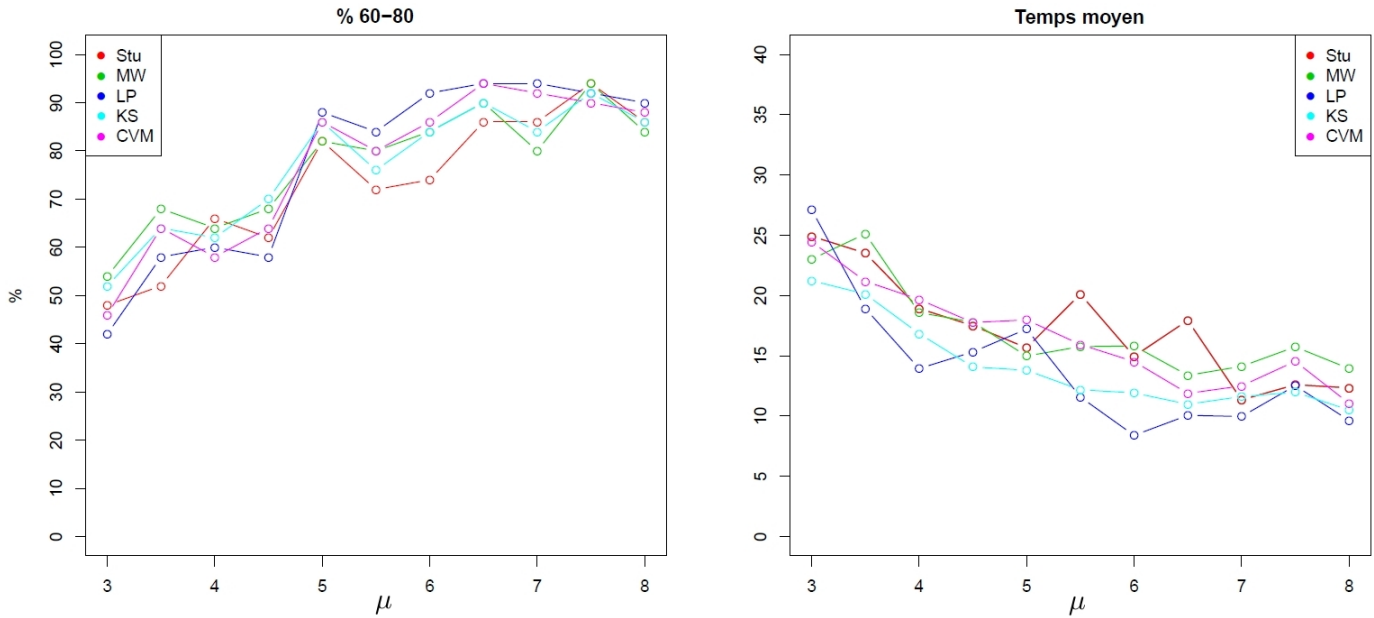


FIGURE A.13 – Pourcentages de simulations de Chlore avec ajout d'un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 3 à 6 par pas de 0.5 ; Méthode avec classifications

	Méthode	N(0 ; 2,4)				N(0 ; 4)			
		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	7	12	15	18	8	16	21	15
Ligne	BinVar	2	5	5	3	3	3	7	0
	SegMoy	11	19	25	12	18	26	34	13
	SegVar	3	8	8	5	4	5	9	2
	Pelt	10	13	15	12	11	17	25	12
	Pett	5	12	20	15	9	15	22	18
	Br	1	1	3	2	0	2	3	1
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	21	25	33	14.36	23	33	41	16.26
	Bar	27	32	33	18.01	24	28	29	23.97
	GLR	28	29	29	16.99	26	32	34	25.39
	MW	9	13	21	16.14	12	19	26	16.21
	Moo	18	28	33	14.06	29	37	44	19.73
	LP	14	19	29	17.22	23	34	41	14.66
	KS	16	27	37	26.47	16	21	29	24.63
	CVM	10	17	24	21.69	15	20	30	22.28

TABLE A.32 – Résultats des détections de ruptures sur 100 simulations de Chlore, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$  ; méthode avec classifications

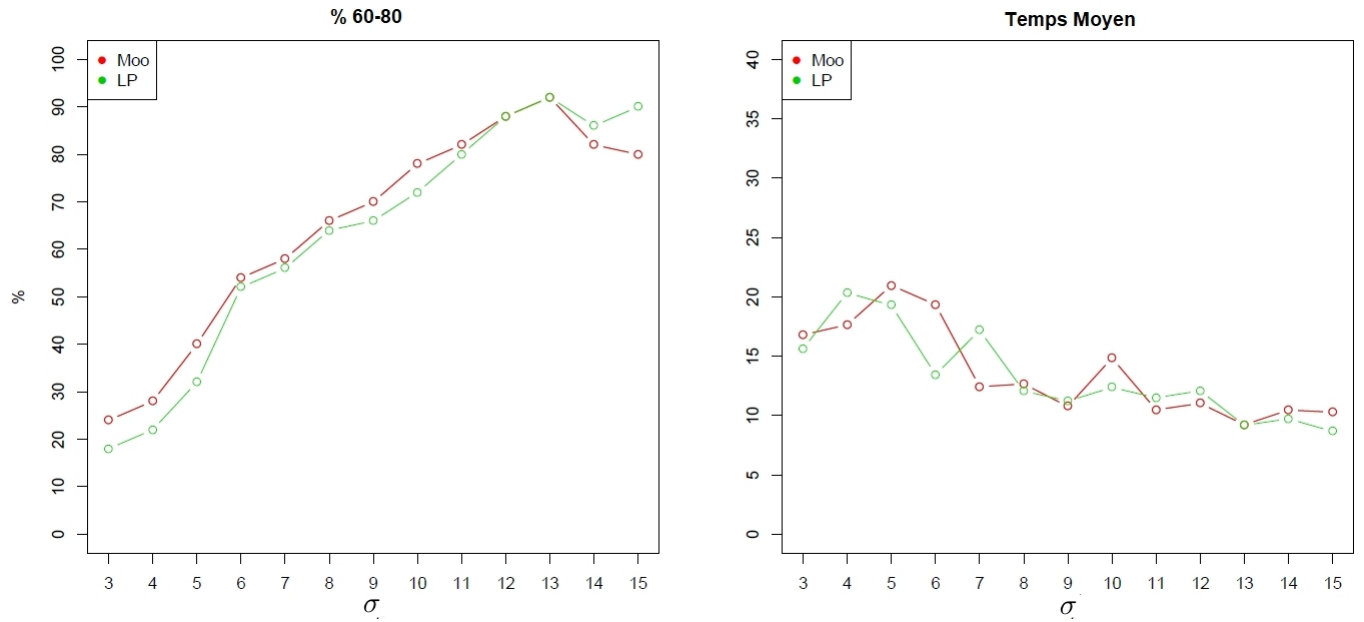


FIGURE A.14 – Pourcentages de simulations de Chlore avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 3 à 15 ; Méthode avec classifications

Méthode		$\mathcal{N}(0, 5T; 1)$			
		%60-80	%60-90	%60-100	%110-130
Hors Ligne	BinMoy	31	42	52	62
	BinVar	18	18	18	10
	SegMoy	43	61	72	79
	SegVar	45	47	47	4
	Pelt	40	66	82	91
	Pett	13	40	56	72
	Br	44	49	52	57
Méthode		%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	75	97	98	9.85
	Bar	34	37	38	36.01
	GLR	63	74	74	18.77
	MW	84	95	95	13.46
	Moo	40	81	90	10.41
	LP	74	94	94	9.52
	KS	78	91	91	11.59
	CVM	82	94	94	12.23

TABLE A.33 – Résultats des détections de ruptures sur 100 simulations de Chlore, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec classifications



Méthode	0	1	2	3
StuR	29	61	10	0
BarR	0	13	80	7
GLRR	0	17,5	77	5,5
MWR	54	40,5	5,5	0
MooR	29	63,5	7,5	0
LPR	26,5	69	4,5	0
KSR	40,5	51	8,5	0
CVMR	42	53,5	4,5	0

TABLE A.34 – Pourcentages de simulations où les méthodes de détection en ligne (avec classifications) ont trouvé 0,1,2,... ruptures sur 200 simulations de chlore de 190 individus

## Créatinine

		$N(4,6; 1)$				$N(17,1; 1)$			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110_130
Hors	BinMoy	12	14	17	8	16	25	32	19
Ligne	BinVar	1	1	1	0	1	1	1	0
	SegMoy	11	14	21	7	17	23	32	15
	SegVar	1	1	1	1	2	2	4	1
	Pelt	15	18	20	11	37	42	45	28
	Pett	16	23	26	12	24	28	32	37
	Br	2	3	4	2	33	34	34	28

		%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	44	49	60	9.03	64	73	76	5.77
	Bar	31	42	50	13.59	41	51	58	14.63
	GLR	27	39	49	13.45	51	57	62	12.82
	MW	13	16	18	22.79	64	69	72	20.51
	Moo	20	25	31	18.3	13	18	21	17.71
	LP	17	25	35	15.85	44	54	62	15.77
	KS	22	26	32	14.48	59	69	73	15.01
	CVM	16	20	21	21.1	60	65	71	20.45

TABLE A.35 – Résultats des détections de ruptures sur 100 simulations de Créatinine, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec classifications

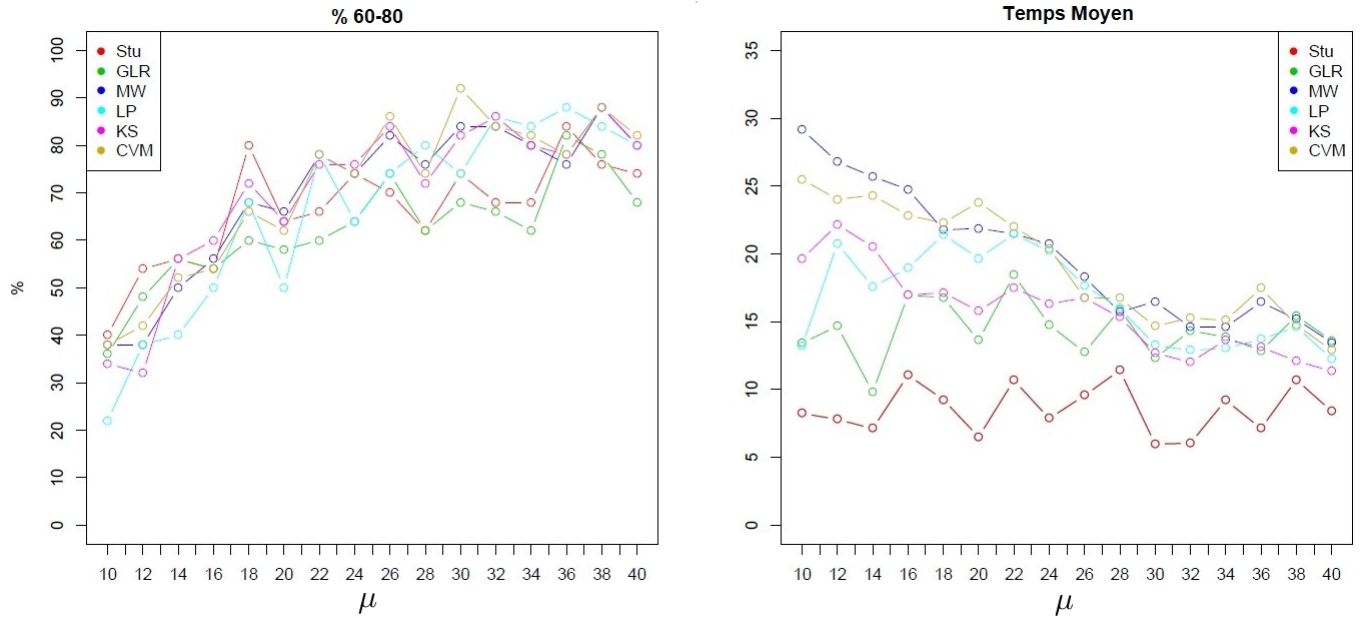


FIGURE A.15 – Pourcentages de simulations de Créatinine avec ajout d'un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 10 à 40 ; Méthode avec classifications

		N(0 ; 4,6)				N (0 ; 17,1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110_130
Hors	BinMoy	9	15	17	11	12	17	19	9
Ligne	BinVar	2	3	3	0	2	2	4	1
	SegMoy	12	18	22	7	13	16	17	13
	Seg Var	2	4	4	3	2	3	6	2
	Pelt	9	11	15	9	10	15	21	10
	Pett	18	21	23	9	4	7	16	18
	Br	4	4	4	4	4	4	9	5
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	34	47	56	7.18	38	54	62	9.27
	Bar	33	48	54	14.22	44	54	61	14.86
	GLR	25	37	49	13.79	40	45	54	16.76
	MW	7	13	16	22.04	15	22	28	20.86
	Moo	9	17	22	18.27	25	43	50	19.68
	LP	12	20	27	11.1	21	34	41	13.11
	KS	11	16	24	17.69	17	23	29	19.22
	CVM	5	12	16	19.28	15	18	26	20.78

TABLE A.36 – Résultats des détections de ruptures sur 100 simulations de Créatinine, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$  ; méthode avec classifications

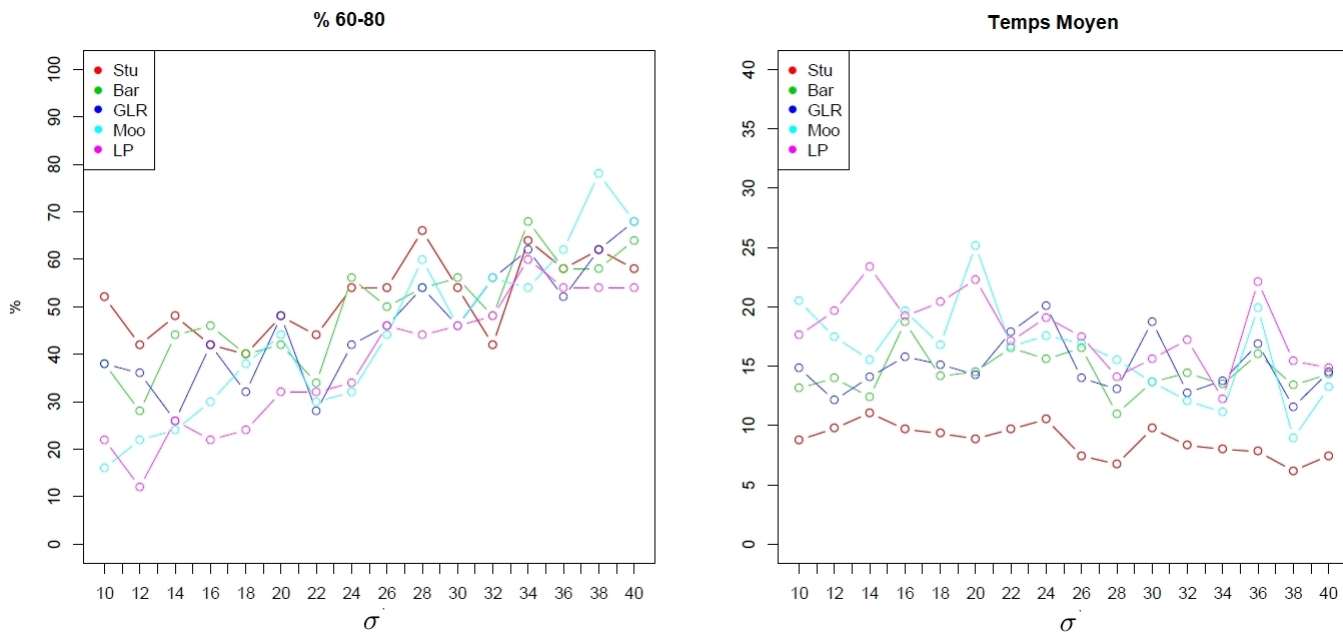


FIGURE A.16 – Pourcentages de simulations de Créatinine avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 15 à 45 par pas de 2 ; Méthode avec classifications

		$\mathcal{N}(0, 5T; 1)$			
Méthode		%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	11	18	24	12
Ligne	BinVar	0	0	1	0
	SegMoy	9	16	24	13
	SegVar	1	1	3	0
	Pelt	16	27	36	32
	Pett	15	25	31	30
	Br	12	19	23	23
Méthode		%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	41	54	63	9.01
	Bar	32	41	49	15.65
	GLR	31	41	54	14.15
	MW	27	47	60	25.3
	Moo	19	24	27	19.96
	LP	29	42	52	20.97
	KS	32	49	65	17
	CVM	28	48	65	21.37

TABLE A.37 – Résultats des détections de ruptures sur 100 simulations de Créatinine, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec classifications

Méthode	0	1	2	3
StuR	1	26,5	67	5,5
BarR	0	41,5	55,5	3
GLRR	1	53,5	43,5	2
MWR	35	56	9	0
MooR	34,5	57	8,5	0
LPR	27,5	61,5	11	0
KSR	32,5	57,5	9,5	0,5
CVMR	38,5	54	7,5	0

TABLE A.38 – Pourcentages de simulations où les méthodes de détection en ligne (avec classifications) ont trouvé 0,1,2,... ruptures sur 200 simulations de créatinine de 190 individus

## Ferritine

		N(7,6 ; 1)				N(40 ; 1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110_130
Hors	BinMoy	11	12	13	8	15	18	18	12
Ligne	BinVar	0	0	0	0	0	0	0	0
	SegMoy	9	18	22	11	21	26	29	15
	Seg Var	0	0	0	0	0	0	1	0
	Pelt	5	7	10	8	22	26	27	20
	Pett	19	22	27	11	30	33	36	13
	Br	1	1	3	1	15	17	17	10
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	50	56	63	9.01	59	71	74	5.65
	Bar	36	43	53	10.94	39	53	59	9.73
	GLR	37	49	53	14.06	41	53	59	11.59
	MW	10	15	18	21.66	29	38	42	26.18
	Moo	15	21	24	15.9	26	30	32	24.7
	LP	23	31	36	13.32	35	44	50	20.52
	KS	21	31	36	16.37	33	49	51	20.44
	CVM	14	20	24	22.95	28	37	39	25.25

TABLE A.39 – Résultats des détections de ruptures sur 100 simulations de Ferritine, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec classifications

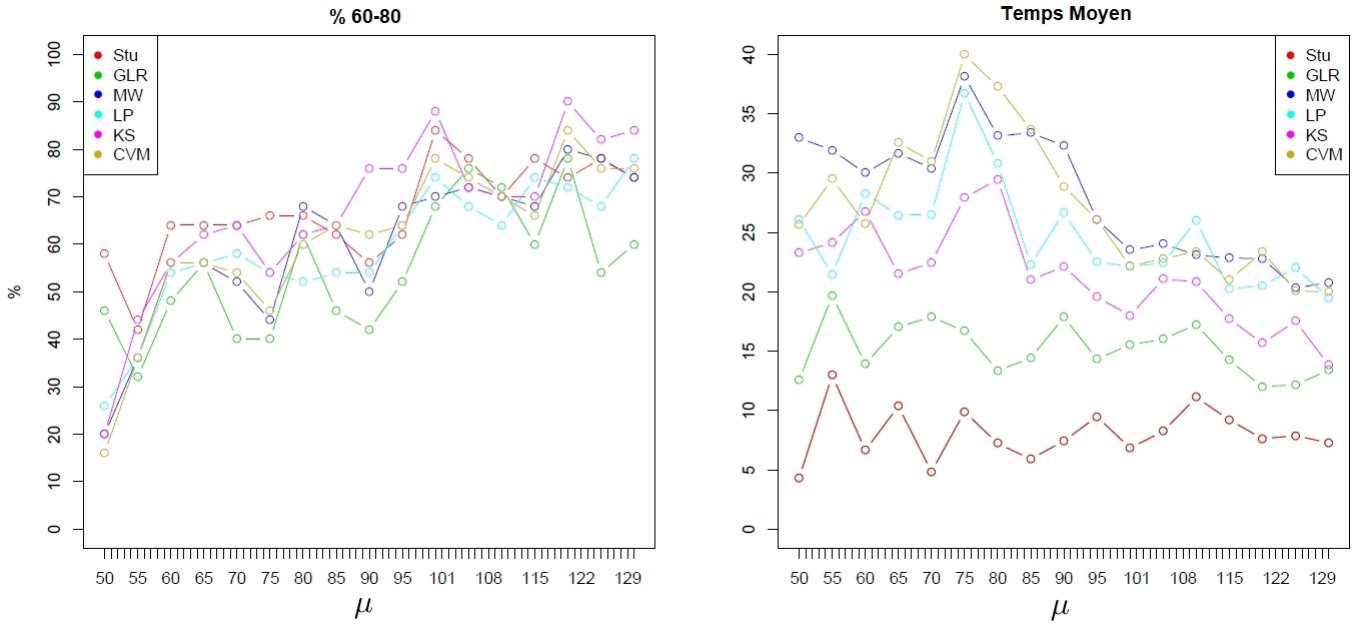


FIGURE A.17 – Pourcentages de simulations de Ferritine avec ajout d’un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 50 à 130 par pas de 5 ; Méthode avec classifications

		N(0 ; 7,6)				N(0 ; 40)			
Méthode		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110_130
Hors	BinMoy	10	11	12	3	9	9	12	9
Ligne	BinVar	0	0	0	0	0	0	0	0
	SegMoy	12	19	24	7	11	14	19	7
	SegVar	0	0	0	0	1	1	1	0
	Pelt	8	12	15	5	22	26	31	17
	Pett	14	21	24	10	16	19	24	15
	Br	0	1	1	1	3	3	3	1

Méthode		%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	53	68	76	7.78	57	58	64	10.32
	Bar	38	53	59	11.99	39	43	52	13.76
	GLR	43	56	68	12.15	44	51	63	13.75
	MW	10	16	19	24.28	17	20	26	19.68
	Moo	14	23	28	16.99	31	37	44	13.48
	LP	10	20	29	17.67	23	25	33	13.67
	KS	9	19	27	18.24	22	25	29	19.65
	CVM	8	13	18	21.46	17	19	25	21.22

TABLE A.40 – Résultats des détections de ruptures sur 100 simulations de Ferritine, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$  ; méthode avec classifications

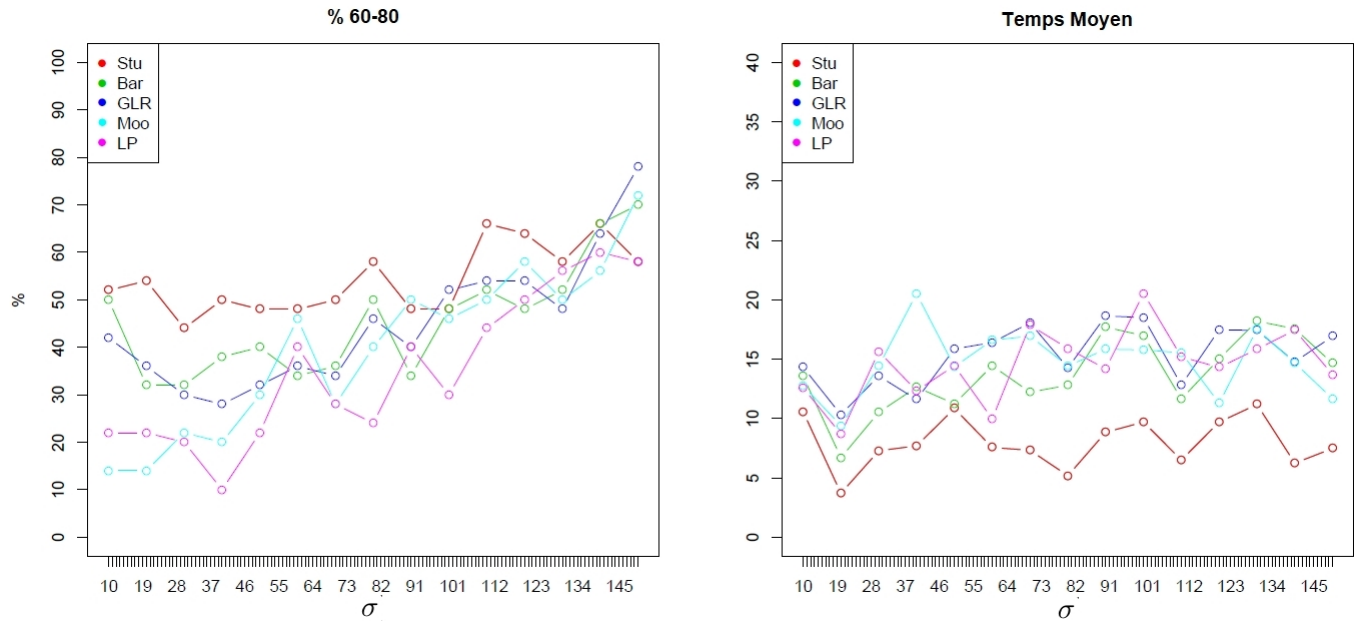


FIGURE A.18 – Pourcentages de simulations de Ferritine avec ajout d’un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 50 à 200 par pas de 10 ; Méthode avec classifications

Méthode		$\mathcal{N}(0,5T; 1)$			
		%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	10	13	17	5
Ligne	BinVar	0	0	0	0
	SegMoy	13	16	21	8
	SegVar	0	0	0	0
	Pelt	16	21	23	4
	Pett	22	28	34	13
	Br	2	5	5	0

Méthode		%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	45	48	58	7.96
	Bar	35	46	57	10.37
	GLR	37	45	55	12.29
	MW	7	12	15	19.8
	Moo	13	19	24	16.5
	LP	17	21	30	15.25
	KS	9	14	24	16.77
	CVM	8	10	13	17.81

TABLE A.41 – Résultats des détections de ruptures sur 100 simulations de Ferritine, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec classifications

Méthode	0
StuR	100
BarR	100
GLRR	100
MWR	100
MooR	100
LPR	100
KSR	100
CVMR	100

TABLE A.42 – Pourcentages de simulations où les méthodes de détection en ligne (avec classifications) ont trouvé 0 ou 1 ruptures sur 200 simulations de ferritine de 20 individus

## Protéines

		N(3,5; 1)				N(8,37; 1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110_130
Hors	BinMoy	22	24	28	27	80	80	81	72
Ligne	BinVar	2	3	3	0	1	1	2	0
	SegMoy	29	29	34	23	84	87	91	84
	SegVar	4	6	8	0	1	1	2	0
	Pelt	34	36	41	25	68	72	73	68
	Pett	28	32	38	25	65	65	65	62
	Br	19	21	24	11	76	79	82	76
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	36	39	44	14.75	72	77	77	12.42
	Bar	9	15	18	7.14	8	14	18	12.27
	GLR	24	33	38	11.96	74	80	83	15.84
	MW	37	39	46	17.29	69	75	75	12.69
	Moo	10	14	19	6.49	36	43	46	11.69
	LP	31	34	43	12.49	75	85	87	10.56
	KS	41	49	56	11.1	73	79	81	10.91
	CVM	34	39	46	14.47	76	80	80	12.91

TABLE A.43 – Résultats des détections de ruptures sur 100 simulations de Protéines, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec classifications

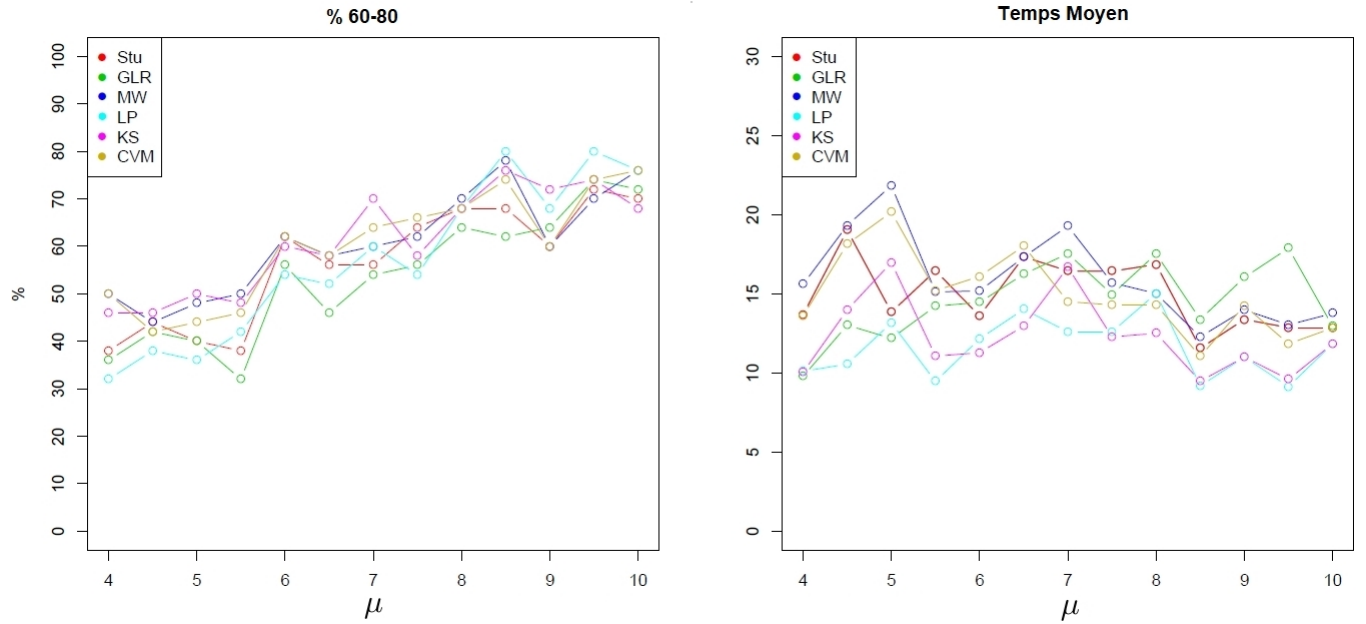


FIGURE A.19 – Pourcentages de simulations de Protéines avec ajout d'un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 4 à 10 par pas de 0.5 ; Méthode avec classifications

	Méthode	N(0 ; 3,5)				N ( 0 ; 8,37 )			
		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110_130
Hors	BinMoy	5	7	9	5	10	15	21	14
Ligne	BinVar	0	0	0	0	9	9	9	5
	SegMoy	10	16	22	9	24	40	54	16
	Seg Var	0	1	1	0	15	17	18	8
	Pelt	6	9	13	11	28	37	45	28
	Pett	9	12	16	12	11	14	21	12
	Br	1	2	3	2	2	3	4	4
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	11	15	18	13.18	24	35	46	7.75
	Bar	5	10	12	8.87	26	31	37	13.82
	GLR	7	10	12	12.72	21	24	26	12.84
	MW	7	10	15	14.38	15	18	27	10.2
	Moo	7	10	13	8.66	27	39	49	9.83
	LP	10	13	18	9.65	27	35	42	7.74
	KS	8	11	15	14.88	17	25	33	14.01
	CVM	6	7	13	13.06	15	21	32	9.99

TABLE A.44 – Résultats des détections de ruptures sur 100 simulations de Protéines, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$  ; méthode avec classifications



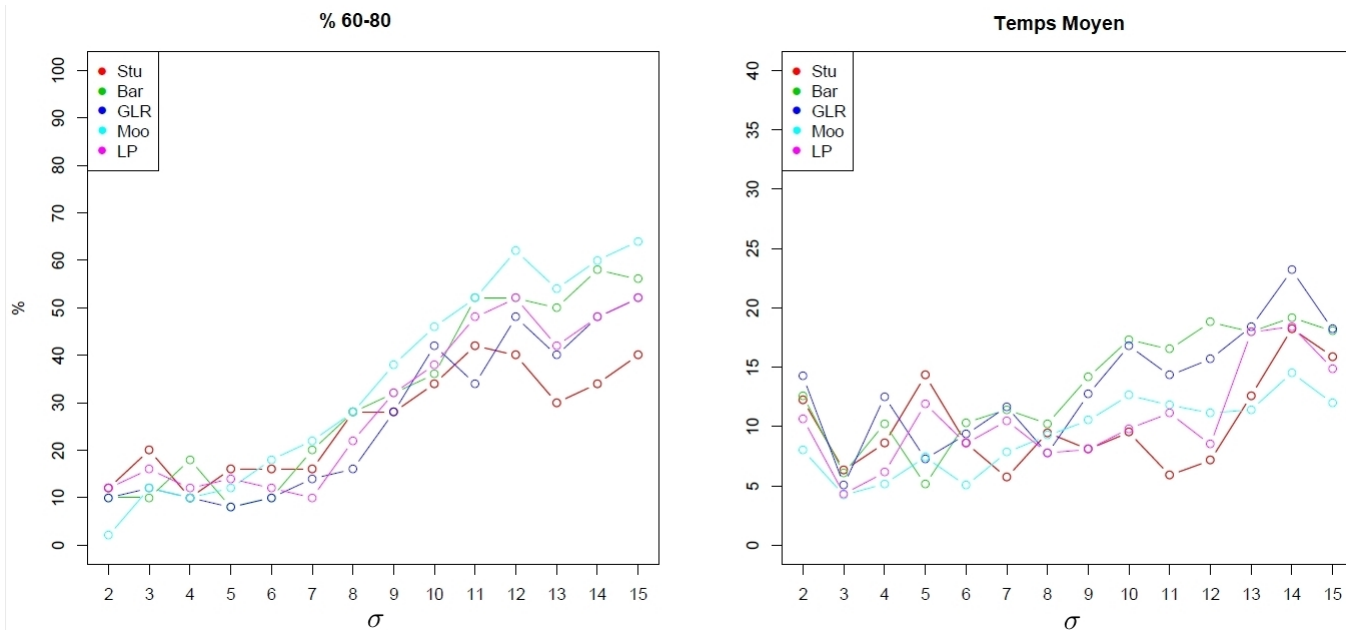


FIGURE A.20 – Pourcentages de simulations de Protéines avec ajout d’un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 5 à 20 ; Méthode avec classifications

		$\mathcal{N}(0, 5T; 1)$			
	Méthode	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	47	71	79	86
Ligne	BinVar	7	18	39	40
	SegMoy	23	66	90	97
	SegVar	4	19	57	76
	Pelt	38	61	69	86
	Pett	61	70	74	80
	Br	41	70	81	87
	Méthode	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	46	68	71	14.62
	Bar	16	22	28	34.09
	GLR	50	76	80	17.1
	MW	48	68	71	17
	Moo	14	41	75	7.96
	LP	41	78	83	10.71
	KS	51	76	80	12.52
	CVM	47	68	75	14.59

TABLE A.45 – Résultats des détections de ruptures sur 100 simulations de Protéines, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec classifications

Méthode	0	1	2	3
StuR	32	48	18	2
BarR	44,5	45,5	9,5	0,5
GLRR	43	44	12	1
MWR	37,5	48,5	13	1
MooR	69	25	6	0
LPR	48,5	44	7	0,5
KSR	39	49,5	11	0,5
CVMR	38	48	13,5	0,5

TABLE A.46 – Pourcentages de simulations où les méthodes de détection en ligne (avec classifications) ont trouvé 0,1,2,... ruptures sur 200 simulations de protéines de 190 individus

## PSA

		N(6,8; 1)				N(50,62; 1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	21	29	33	26	9	19	25	11
Ligne	BinVar	23	27	28	13	25	40	51	35
	SegMoy	15	21	22	7	12	21	23	6
	SegVar	32	36	37	17	33	51	66	47
	Pelt	78	78	78	74	56	56	56	72
	Pett	96	96	96	96	23	29	39	29
	Br	96	96	96	94	4	7	9	7
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	53	53	53	8.03	69	69	69	6.36
	Bar	4	4	4	24.61	10	10	10	23.79
	GLR	9	9	9	23.34	9	9	9	23.43
	MW	95	95	95	8.58	100	100	100	6.31
	Moo	31	31	31	17.86	92	92	92	6.49
	LP	93	93	93	10.46	99	99	99	4.41
	KS	98	98	98	5.8	96	96	96	5.48
	CVM	96	96	96	7.47	100	100	100	5.88

TABLE A.47 – Résultats des détections de ruptures sur 100 simulations de PSA, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec classifications

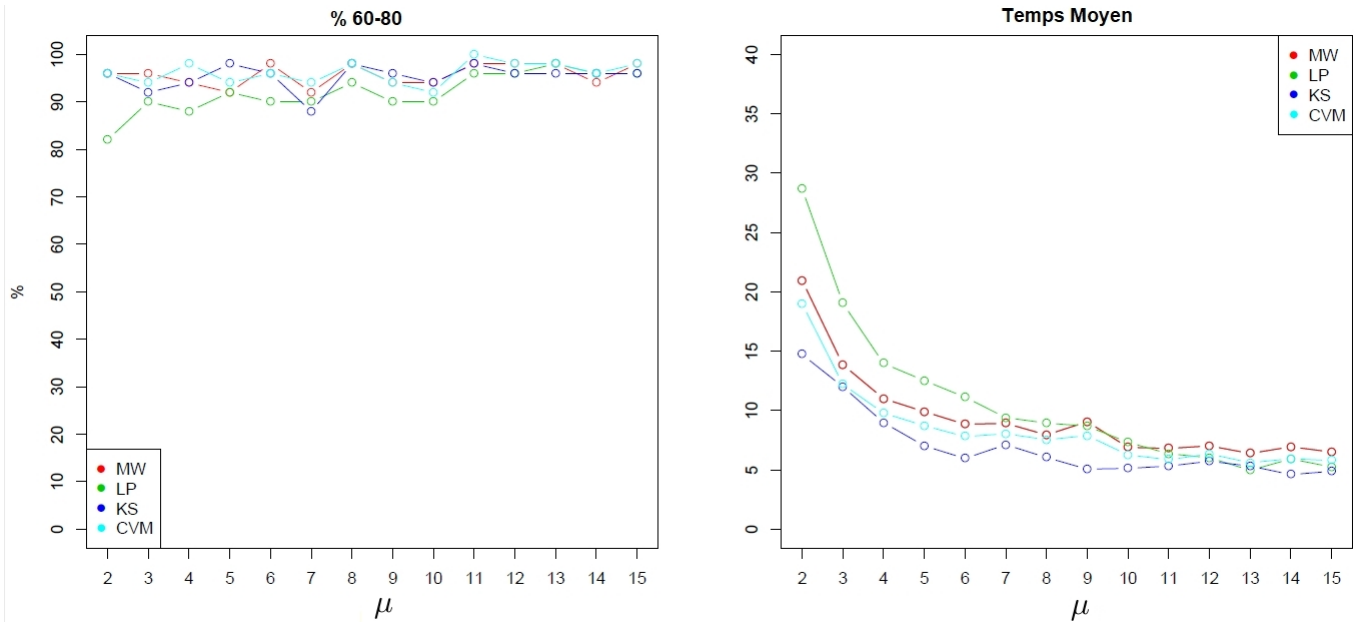


FIGURE A.21 – Pourcentages de simulations de PSA avec ajout d'un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 2 à 15 ; Méthode avec classifications

		$N(0; 6,8)$				$N(0; 50,62)$			
Méthode		%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	19	24	28	23	10	12	16	21
Ligne	BinVar	46	51	54	28	20	24	31	22
	SegMoy	7	12	17	15	3	5	10	11
	SegVar	44	48	53	35	14	19	27	24
	Pelt	79	83	83	77	51	52	55	50
	Pett	13	16	17	18	21	26	35	23
	Br	4	5	6	4	8	10	12	15
Méthode		%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	48	56	65	8.28	56	56	56	8.16
	Bar	7	8	9	24.34	6	6	6	24.82
	GLR	10	11	12	23.55	7	7	7	25.12
	MW	25	35	38	5.08	52	63	69	5.69
	Moo	76	79	80	11.14	84	84	84	7.88
	LP	87	90	92	9.37	99	99	99	5.01
	KS	80	85	88	15.26	96	97	97	9.66
	CVM	50	58	64	10.8	95	99	99	12.62

TABLE A.48 – Résultats des détections de ruptures sur 100 simulations de PSA, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$ ; méthode avec classifications

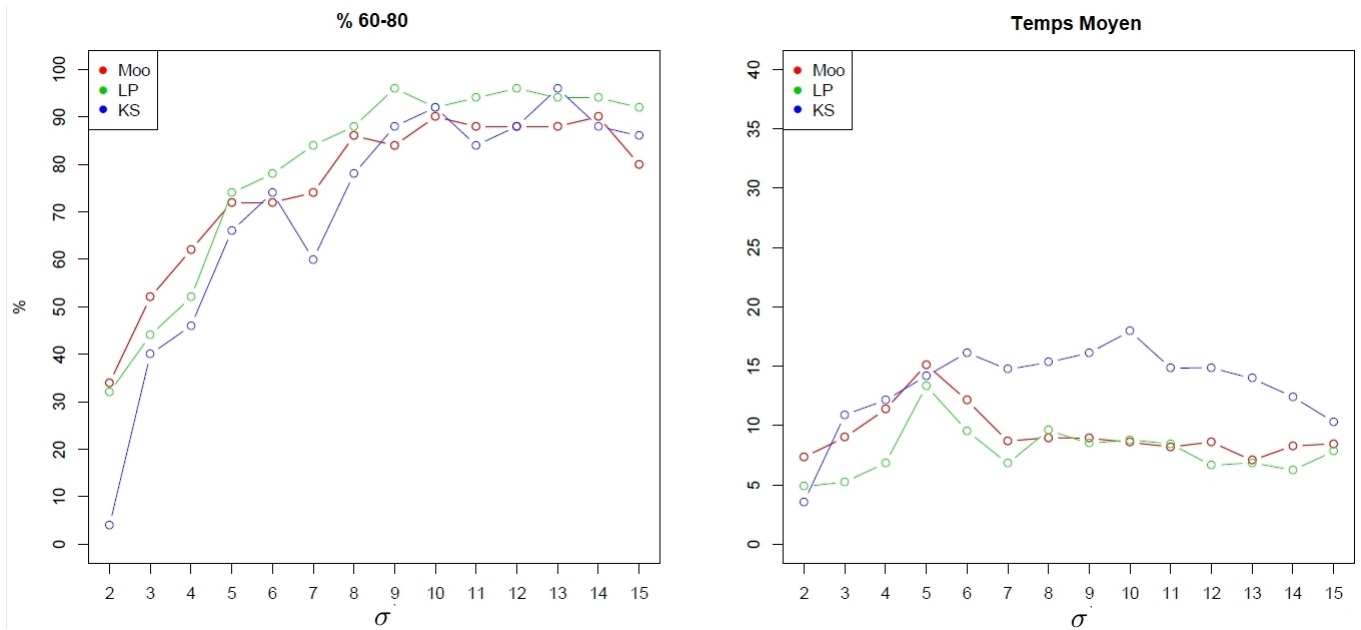


FIGURE A.22 – Pourcentages de simulations de PSA avec ajout d'un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne données ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 2 à 15 ; Méthode avec classifications

Méthode		$N(0, 5T; 1)$			
		%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	15	22	27	23
Ligne	BinVar	15	50	77	69
	SegMoy	12	21	24	16
	SegVar	12	44	76	64
	Pelt	65	74	76	69
	Pett	79	80	87	78
	Br	62	77	77	69
	Méthode		%60-80	%60-90	%60-100
En Ligne	Stu	56	61	61	9.04
	Bar	2	2	2	22.4
	GLR	7	8	8	21.41
	MW	96	96	96	11.98
	Moo	21	52	70	12.41
	LP	93	95	95	11.64
	KS	93	93	93	9.74
	CVM	97	97	97	10.9

TABLE A.49 – Résultats des détections de ruptures sur 100 simulations de PSA, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec classifications

Méthode	0
StuR	100
BarR	100
GLRR	100
MWR	100
MooR	100
LPR	100
KSR	100
CVMR	100

TABLE A.50 – Pourcentages de simulations où les méthodes de détection en ligne (avec classifications) ont trouvé 0,1,2,... ruptures sur 200 simulations de PSA de 20 individus

## Urée

		N(6,9; 1)				N(34,24; 1)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110_130
Hors	BinMoy	16	19	30	12	51	51	51	76
Ligne	BinVar	8	9	13	3	18	20	23	16
	SegMoy	24	29	40	9	59	59	60	80
	SegVar	12	14	21	11	25	30	36	25
	Pelt	37	43	46	27	72	72	72	90
	Pett	21	29	38	27	65	65	66	82
	Br	34	35	38	32	67	67	67	87
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	74	76	78	19.45	93	93	93	7.28
	Bar	16	21	23	28.1	31	31	31	22.28
	GLR	35	40	40	26.42	50	50	50	17.6
	MW	80	80	80	17.38	89	89	89	9.67
	Moo	17	23	29	27.2	93	93	93	8.19
	LP	80	81	82	20.69	93	93	93	7.38
	KS	83	83	83	13.45	89	89	89	8.19
	CVM	82	82	82	15.59	91	91	91	9.5

TABLE A.51 – Résultats des détections de ruptures sur 100 simulations de Urée, VISTA 2, avec ajout d'un bruit  $\mathcal{N}(\mu, 1)$  à  $i = 70$ ; méthode avec classifications

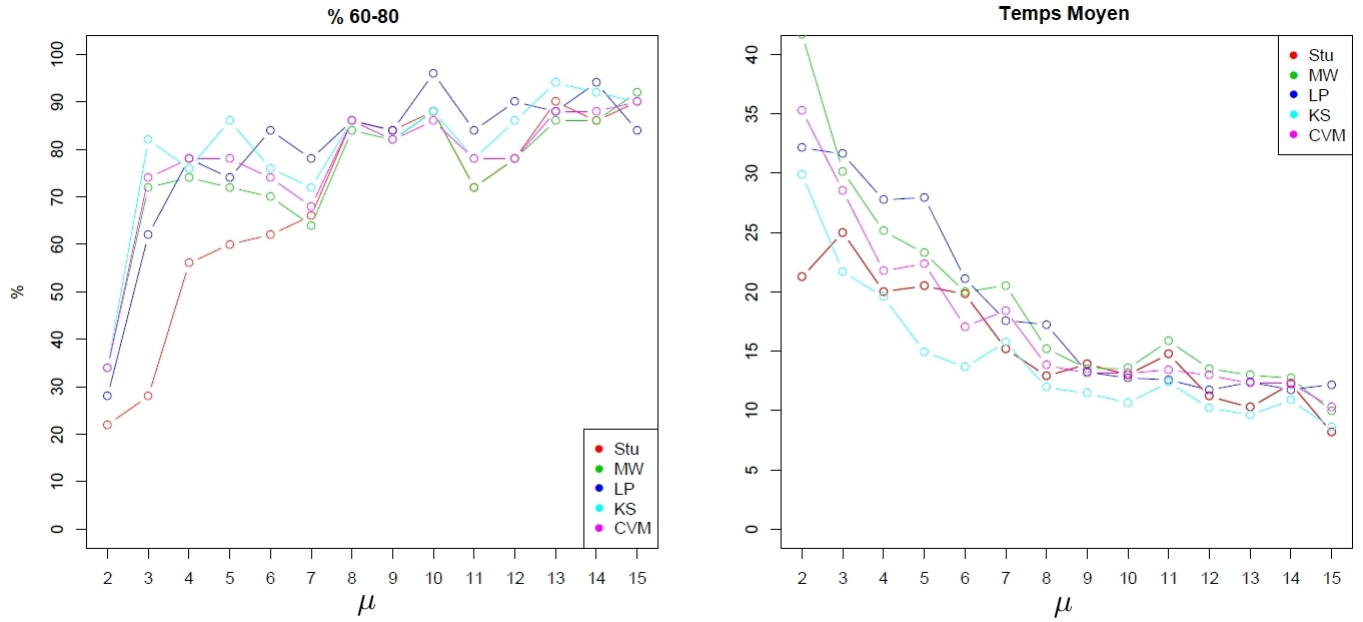


FIGURE A.23 – Pourcentages de simulations d’Urée avec ajout d’un bruit  $\mathcal{N}(\mu; 1)$  de taille 50 à  $i = 70$ , où les méthodes en ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\mu$  allant de 2 à 15 ; Méthode avec classifications

		N(0 ; 6,9)				N (0 ; 34,24)			
	Méthode	%60-80	%60-90	%60-100	%110-130	%60-80	%60-90	%60-100	%110_130
Hors	BinMoy	36	46	52	29	30	43	50	25
Ligne	BinVar	20	22	24	9	70	71	71	76
	SegMoy	33	47	56	30	31	47	57	25
	Seg Var	19	22	25	11	75	76	76	78
	Pelt	44	57	65	34	61	72	74	62
	Pett	23	28	36	14	20	28	35	13
	Br	27	34	41	22	16	25	29	17
	Méthode	%60-80	%60-90	%60-100	Temps Moyen	%60-80	%60-90	%60-100	Temps Moyen
En Ligne	Stu	34	48	57	18.76	79	79	79	11.66
	Bar	17	24	26	28.71	38	38	38	22.57
	GLR	22	22	25	30	40	40	40	22.93
	MW	30	41	48	19.42	73	77	81	16.98
	Moo	70	80	86	18.76	92	92	92	8.96
	LP	66	77	83	17.64	84	84	84	11.1
	KS	38	45	49	21.7	78	79	80	16.4
	CVM	35	46	51	19.35	72	76	77	17.51

TABLE A.52 – Résultats des détections de ruptures sur 100 simulations de Urée, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0, \sigma)$  à  $i = 70$  ; méthode avec classifications

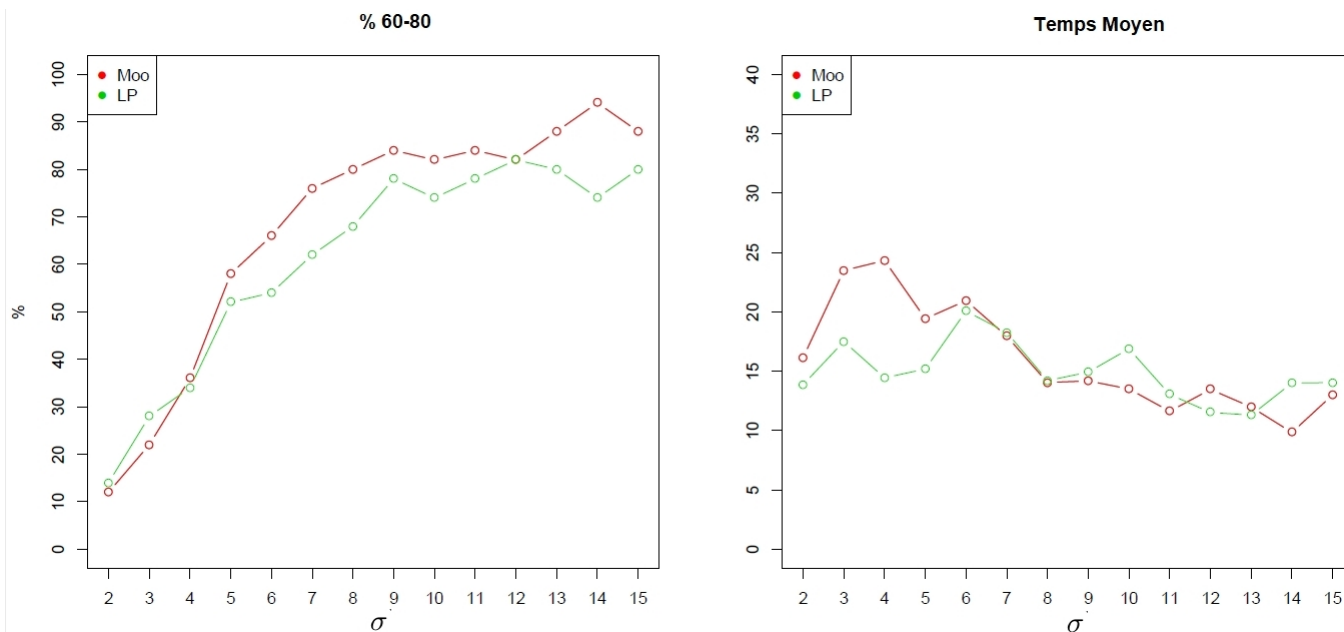


FIGURE A.24 – Pourcentages de simulations d’Urée avec ajout d’un bruit  $\mathcal{N}(0, \sigma)$  de taille 50 à  $i = 70$ , où les méthodes en ligne données ont trouvé la rupture entre  $i = 60$  et  $i = 80$  et temps de détection moyen ;  $\sigma$  allant de 2 à 15 ; Méthode avec classifications

		$\mathcal{N}(0, 5T; 1)$			
Méthode		%60-80	%60-90	%60-100	%110-130
Hors	BinMoy	1	18	65	63
Ligne	BinVar	0	2	12	5
	SegMoy	2	9	47	68
	SegVar	0	6	23	11
	Pelt	9	39	74	65
	Pett	1	25	70	35
	Br	2	16	63	72
	Méthode		%60-80	%60-90	%60-100
En Ligne	Stu	56	82	84	13
	Bar	14	21	23	32.15
	GLR	21	38	39	25.6
	MW	73	74	74	19.39
	Moo	3	28	65	15.45
	LP	64	80	81	18.89
	KS	76	78	78	15.68
	CVM	76	77	77	18.48

TABLE A.53 – Résultats des détections de ruptures sur 100 simulations de Urée, VISTA 2, avec ajout d’un bruit  $\mathcal{N}(0.5T, 1)$  à  $i = 70$  ; méthode avec classifications

Méthode	0	1	2	3
StuR	12	59,5	28	0,5
BarR	0	7,5	79,5	13
GLRR	0	14,5	78	7,5
MWR	41	46,5	12	0,5
MooR	57	40	3	0
LPR	50,5	41,5	8	0
KSR	39,5	43,5	15,5	1,5
CVMR	38,5	47	14	0,5

TABLE A.54 – Pourcentages de simulations où les méthodes de détection en ligne (avec classifications) ont trouvé 0,1,2,... ruptures sur 200 simulations de Urée de 190 individus

## A.6 Étude de l'impact de l'incertitude de mesure sur les valeurs vraies estimées des paramètres biologiques

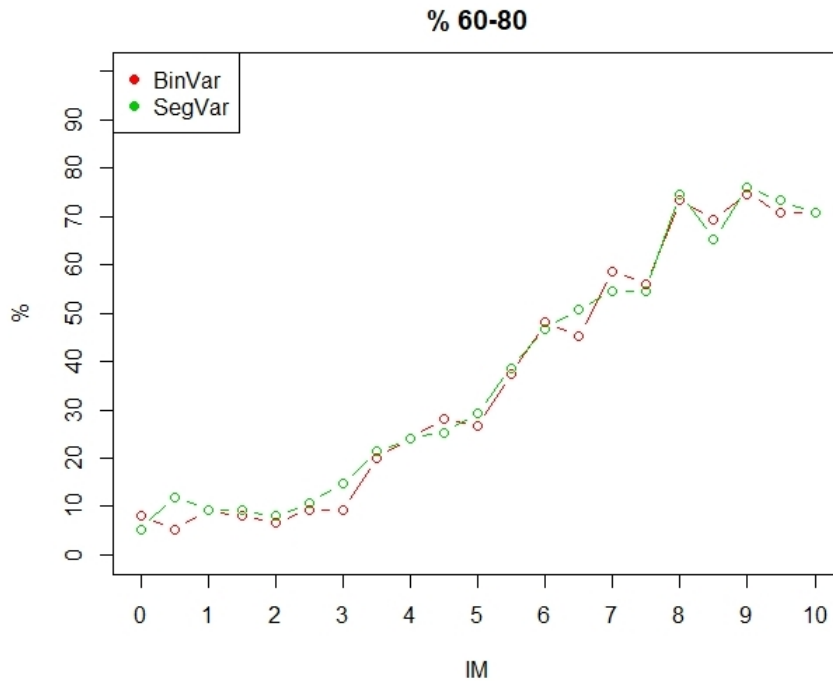


FIGURE A.25 – Pourcentages de simulations de chlore réalisées sur les fonctions de densité des valeurs vraies avec ajout d'un bruit  $\mathcal{N}(m, \%IM \times m)$  de taille 50 à  $i = 70$ , où les méthodes hors ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$ .



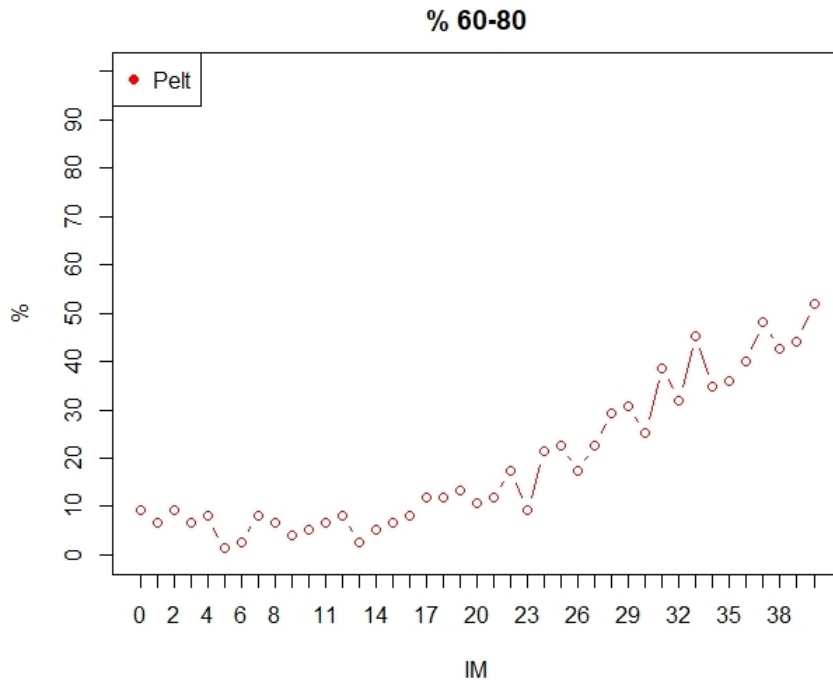


FIGURE A.26 – Pourcentages de simulations de créatinine réalisées sur les fonctions de densité des valeurs vraies avec ajout d'un bruit  $\mathcal{N}(m, \%IM \times m)$  de taille 50 à  $i = 70$ , où les méthodes hors ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$ .

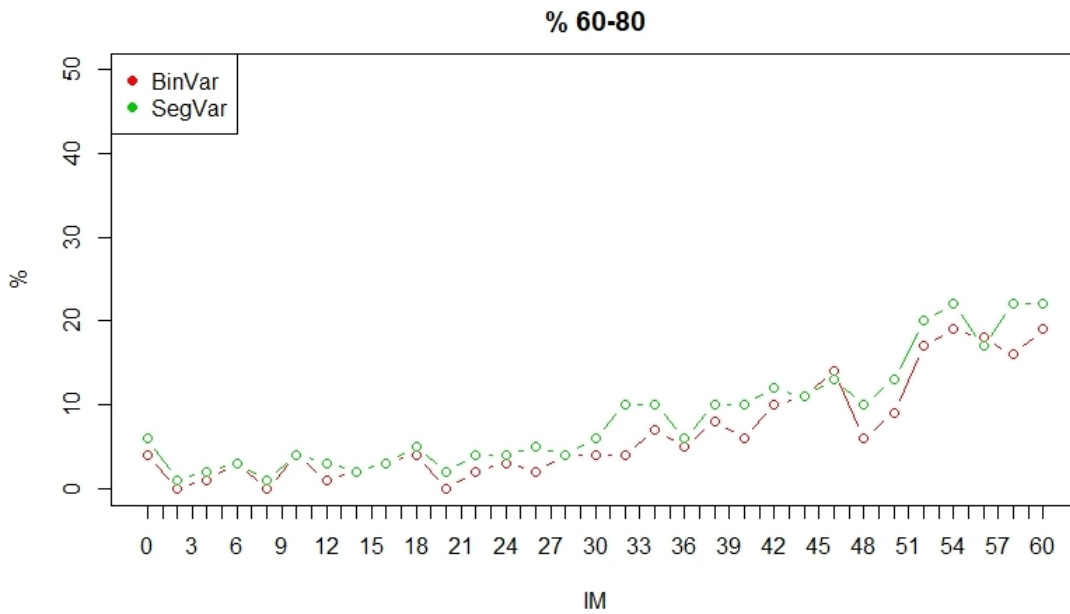


FIGURE A.27 – Pourcentages de simulations de ferritine réalisées sur les fonctions de densité des valeurs vraies avec ajout d'un bruit  $\mathcal{N}(m, \%IM \times m)$  de taille 50 à  $i = 70$ , où les méthodes hors ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$ .

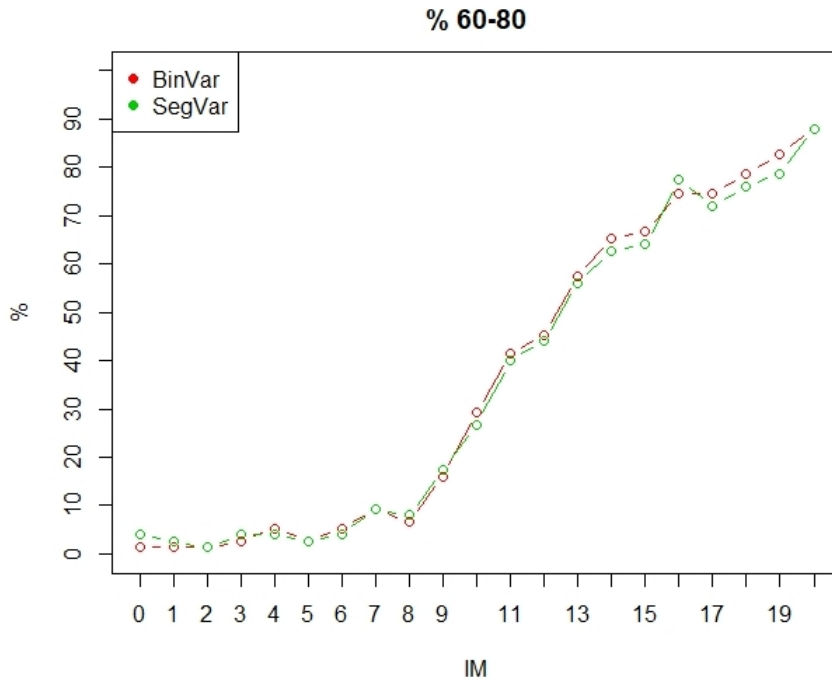


FIGURE A.28 – Pourcentages de simulations de protéines réalisées sur les fonctions de densité des valeurs vraies avec ajout d’un bruit  $\mathcal{N}(m, \%IM \times m)$  de taille 50 à  $i = 70$ , où les méthodes hors ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$ .

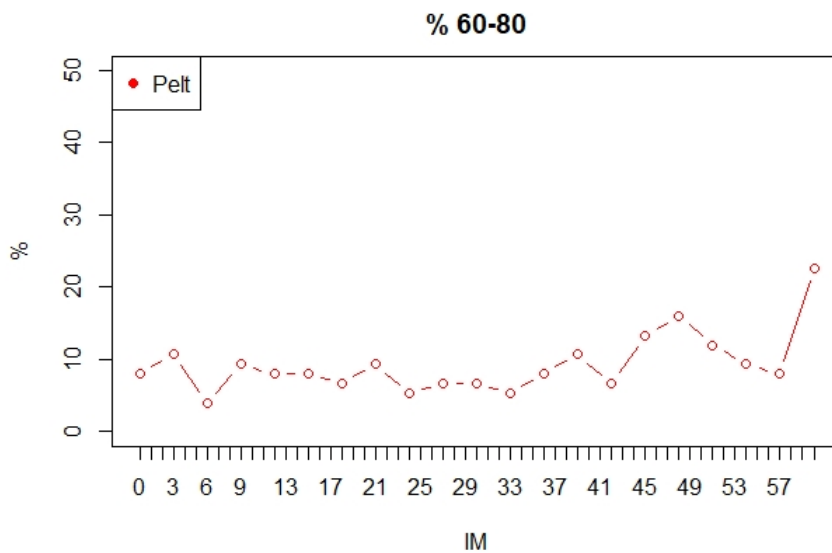


FIGURE A.29 – Pourcentages de simulations de PSA réalisées sur les fonctions de densité des valeurs vraies avec ajout d’un bruit  $\mathcal{N}(m, \%IM \times m)$  de taille 50 à  $i = 70$ , où les méthodes hors ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$ .

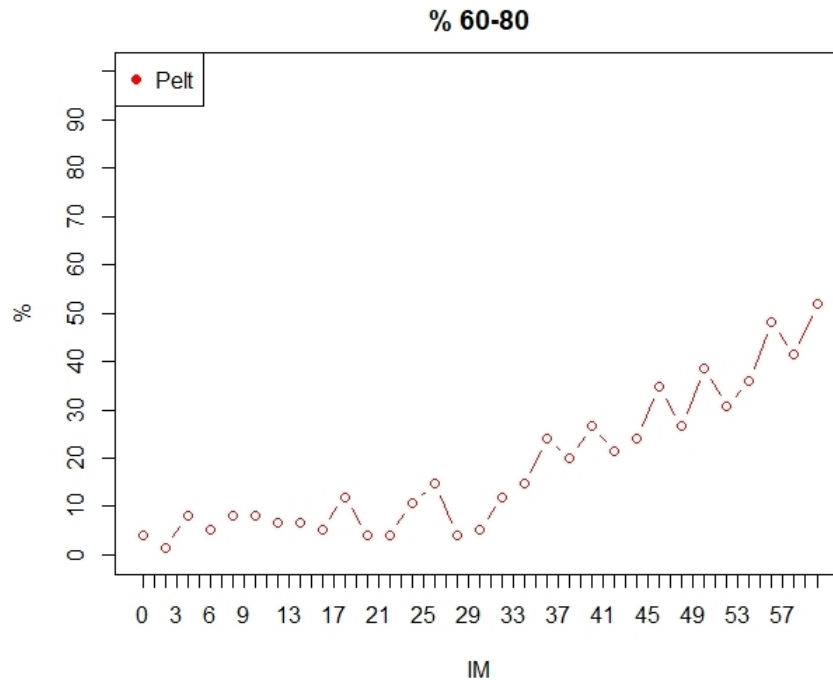


FIGURE A.30 – Pourcentages de simulations d’urée réalisées sur les fonctions de densité des valeurs vraies avec ajout d’un bruit  $\mathcal{N}(m, \%IM \times m)$  de taille 50 à  $i = 70$ , où les méthodes hors ligne ont trouvé la rupture entre  $i = 60$  et  $i = 80$ .

# Bibliographie

- [1] Norme iso 5725-6 :1994, exactitude (justesse et fidélité) des résultats et méthodes de mesure. "<https://www.iso.org/fr/standard/11837.html>", 1994.
- [2] Norme nf en iso 15189 :2012, laboratoire de biologie médicales - exigences concernant la qualité et la compétence. "<https://www.iso.org/fr/standard/56115.html>", 2012.
- [3] Loi 2013 - 442 du 30 mai 2013, réforme de la biologie médicale. "<https://www.legifrance.gouv.fr/affichTexte.do?cidTexte=JORFTEXT000027478077>", 30 mai 2013.
- [4] J. Bai and P. Perron. Computation and analysis of multiple structural change models. *Journal of applied econometrics*, 18 :1–22, 2002.
- [5] Biomnis. Ferritine. "<https://www.eurofins-biomnis.com/referentiel/liendoc/precis/FERRITINE.pdf>", 2012.
- [6] Coffrac. Guide technique d'accréditation pour l'évaluation des incertitudes de mesure en biologie medicale, sh gta 14. "<https://tools.cofrac.fr/documentation/SH-GTA-14>", 2011.
- [7] Coffrac. Guide technique d'accréditation : contrôle de qualité en biologie médicale sh gta 06. "<https://tools.cofrac.fr/documentation/SH-GTA-06>", 2012.
- [8] Julia EYCHENNE. *Budgets éruptifs et origine des paroxysmes explosifs andésitiques en système ouvert : l'éruption d'août 2006 du Tungurahua en Equateur*. PhD thesis, Université Blaise Pascal, 2012.
- [9] J. Fan. On the optimal rates of convergence for nonparametric deconvolution problems. *The Annals of Statistics*, 19 :1257–1272, 1991.
- [10] C. Fraley and A.E. Raftery. Model-based methods of classification : Using the mclust software in chemometrics. *Journal of Statistical Software*, 18, 2007.
- [11] C. Fraley, A.E. Raftery, and L. Scrucca. Package mclust. "<https://cran.r-project.org/web/packages/mclust/mclust.pdf>", July 2019.

- [12] R. Killick. Package changepoint. "<https://cran.r-project.org/web/packages/changepoint/changepoint.pdf>", October 2016.
- [13] R. Killick and I.A. Eckley. changepoint : An r package for changepoint analysis. *Journal of Statistical Software*, 58, June 2014.
- [14] R. Killick, P. Fearnhead, and I.A. Eckley. Optimal detection of changepoints with a linear computational cost. *Journal of American Statistical Association*, 107 :1590–1598, October 2012.
- [15] Sebastien Leibrandt and Jean-Luc Lepennec. Towards fast and routine analyses of volcanic ash morphometry for eruption surveillance applications. *Journal of Volcanology and Geothermal Research*, (297) :11–27, 2015.
- [16] CHU Liège. Albumine sang, référentiel des examens. "[https://www.chu.ulg.ac.be/jcms/c\\_498407/fr/albumine-sang](https://www.chu.ulg.ac.be/jcms/c_498407/fr/albumine-sang)".
- [17] CHU Liège. Chlorures référentiel des examens. "[http://www.chu.ulg.ac.be/jcms/c\\_498669/fr/chlorures-sang](http://www.chu.ulg.ac.be/jcms/c_498669/fr/chlorures-sang)".
- [18] CHU Liège. Créatinine référentiel des examens. "[http://www.chu.ulg.ac.be/jcms/c\\_498839/fr/creatinine-sang](http://www.chu.ulg.ac.be/jcms/c_498839/fr/creatinine-sang)".
- [19] CHU Liège. Protéines totales référentiel des examens. "[http://www.chu.ulg.ac.be/jcms/c\\_496968/proteines-totales-sang](http://www.chu.ulg.ac.be/jcms/c_496968/proteines-totales-sang)".
- [20] CHU Liège. Psa total référentiel des examens. "[http://www.chu.ulg.ac.be/jcms/c\\_351722/fr/psa-total](http://www.chu.ulg.ac.be/jcms/c_351722/fr/psa-total)".
- [21] CHU Liège. Urée référentiel des examens. "[http://www.chu.ulg.ac.be/jcms/c\\_498724/fr/uree-sang](http://www.chu.ulg.ac.be/jcms/c_498724/fr/uree-sang)".
- [22] European Federation of Clinical Chemistry and Laboratory Medicine (EFLM). Eflm biological variation database. "<https://biologicalvariation.eu/>", 2019.
- [23] A.N. Pettitt. A non-parametric approach to the change point problem. *Journal of the Royal Statistical Society Series C*, 28 :126–135, 1979.
- [24] T. Pohlert. Package trend. "<https://cran.r-project.org/web/packages/trend/trend.pdf>", July 2018.
- [25] C. Ricos and al. Desirable specifications for total error, imprecision, and bias, derived from intra- and inter-individual biologic variation. "<https://www.westgard.com/biodatabase1.htm>", 2014.

- [26] G.J. Ross. Package cpm. "<https://cran.r-project.org/web/packages/cpm/cpm.pdf>", July 2015.
- [27] G.J. Ross. Parametric and nonparametric sequential change detection in r : The cpm package. *Journal of Statistical Software*, 66(3) :1–20, 2015.
- [28] L. Rosyl. *Procédure des CIQ*. CHU Clermont-Ferrand, 2014.
- [29] L. Scrucca, M. Fop, B. Murphy, and A.E. Raftery. mclust 5 : Clustering, classification and density estimation using gaussian finite mixture models. *The R journal*, August 2016.
- [30] X-F. Wang and B. Wang. Deconvolution estimation in measurement error models : The r package decon. *Journal of Statistical Software*, 39(10), 2011.
- [31] X-F. Wang and B. Wang. Package decon. "<https://cran.r-project.org/web/packages/decon/decon.pdf>", February 2015.
- [32] J. Westgard, Barry PL, Hunt MR, and Groth T. A multi-rule shewhart chart for control in clinical chemistry. *Clin Chem*, pages 493–501, 1981.
- [33] QC Westgard. Desirable specifications for total error, imprecision, and bias, derived from intra- and inter-individual biologic variation. "<https://www.westgard.com/biodatabase1.htm>".
- [34] A. Zeileis, C. Kleiber, W. Krämer, and K. Hornik. Testing and dating of structural changes in practice. *Computational Statistics and Data Analysis*, 44 :109–123, 2003.
- [35] A. Zeileis, F. Leisch, K. Hornik, c. Kleiber, B. Hanser, and E. Merkle. Package strucchange. "<https://cran.r-project.org/web/packages/strucchange/strucchange.pdf>", 2019.