

Polycopié pour les UE *Statistique Inférentielle* et
Tests.

Vincent Brault




Introduction

Ce cours repose principalement sur des notes de cours de Gérard Biau (Université Pierre et Marie Curie), Anatoli Juditsky (Université Grenoble Alpes), Jean-François Le Gall (Université Paris-Saclay) et Marie-Anne Poursat (Université Paris-Saclay).

Par défaut, les corrections des exercices ne sont pas fournies : soit votre raisonnement est juste du début jusqu'à la fin et, dans ce cas, votre résultat est bon ; soit vous hésitez sur l'une des justifications et même si le résultat est bon, il ne vaut rien. En revanche, lorsque des astuces existent, elles sont mises en fin de chapitre.

Enfin, merci aux étudiants qui me font des retours sur mes coquilles ou fautes d'inattention.

Pour les lecteurices qui souhaiteraient plus d'exercices ou une vision différente de ce qui est présenté dans ce manuscrit, voici quelques références :

- *Mathématiques et statistique pour les sciences de la nature* de Biau et al. (2010) pour des exercices appliqués (Chapitres 5 et 6 principalement).
- *Probabilités, analyse des données et Statistique* de Saporta (2006) pour des exercices théoriques (Chapitre 13 et 14 principalement).
- *Statistique inférentielle : cours et exercices corrigés* de Fourdrinier (2002) pour des exercices théoriques (Chapitres 1, 4, 5 et 6 principalement).
- Pour des exercices en , il y a les livres où François Husson est co-auteur¹ où vous retrouverez aussi des jeux de données sur son site. En particulier :
 - *Statistique avec R, 3ème édition revue et augmentée* de Cornillon et al. (2012).
 - *R pour la statistique et la science des données* de Husson et al. (2018).
 - *Statistiques générales pour utilisateurs. 2 - Exercices et corrigés* de Husson et Pagès (2013).

1. Retrouvez la liste sur sa page web <https://husson.github.io/books.html>

Planning

Le planning mis ici est à titre indicatif et sera régulièrement mis à jour en fonction de l'avancée.

Cours	Date	Estimation
Stat Inf 1	26-09-2025	Cours statistiques descriptives et modélisation.
Stat Inf 2	01-10-2025	TD/TP statistiques descriptives et modélisation.
Stat Inf 3	03-10-2025	Cours sur les estimateurs.
Stat Inf 4	08-10-2025	TD/TP sur les estimateurs.
Stat Inf 5	10-10-2025	Cours sur les propriétés d'un estimateur.
Test 1	13-10-2025	Cours sur l'introduction aux tests.
Stat Inf 6	15-10-2025	TD/TP sur les propriétés d'un estimateur.
Stat Inf 7	22-10-2025	TD/TP sur les intervalles de confiance et l'introduction aux tests.
Stat Inf 8	24-10-2025	Examen Stat Inf.
Test 2	05-11-2025	TD/TP
Test 3	07-11-2025	Cours
Test 4	12-11-2025	TD/TP
Test 5	14-11-2025	Cours
Test 6	19-11-2025	TD/TP
Test 7	26-11-2025	TD/TP
Test 8	28-11-2025	Examen Test

Table des matières

I	UE Statistique inférentielle	6
1	Statistique descriptive	7
1.1	Avant-propos	7
1.2	Rappels/Pré-requis	8
1.3	Généralité sur la démarche scientifique	8
1.4	Individus et variables statistiques	10
1.4.1	Individus statistiques	10
1.4.2	Variables statistiques	12
1.5	Statistiques descriptives	14
1.5.1	Statistiques descriptives univariées	14
1.5.2	Statistiques descriptives bivariées	21
1.5.3	Statistiques descriptives multivariées	23
2	Modélisation	24
2.1	Objectifs	24
2.2	Introduction	24
2.3	Modèle statistique	25
2.4	Solutions des exercices	28
3	Estimation	30
3.1	Objectifs	30
3.2	Rappels/Pré-requis	30
3.3	Introduction	30
3.4	Consistance d'un estimateur	31
3.5	Construction d'un estimateur	32
3.5.1	Méthode des moments	32
3.5.2	Méthode du maximum de vraisemblance	33
3.6	Qualité d'un estimateur	36
3.6.1	Loi asymptotique	36
3.6.2	Estimation sans biais	41
3.6.3	Estimation optimale	42
3.7	Solutions des exercices	48
4	Intervalle et région de confiance	51
4.1	Objectifs	51
4.2	Introduction	51
4.3	Premières constructions	52
4.4	Intervalle de confiance de niveaux obtenus par des inégalités de probabilités	53
4.4.1	Cas de variances uniformément bornée	53
4.4.2	Cas de lois à supports inclus dans un compact donné	54
4.5	Intervalle de confiance asymptotique	55
4.5.1	Cas de variances uniformément bornées	55
4.5.2	Estimation consistante de la variance	55
4.5.3	Stabilisation de la variance	56

II	UE Test	57
5	Introduction aux tests	58
5.1	Objectifs	58
5.2	Formalisme et démarche expérimentale	58
5.2.1	Mesure de la qualité d'un test	59
5.2.2	Dissymétrie des rôles des hypothèses \mathcal{H}_0 et \mathcal{H}_1	60
5.2.3	Démarche de construction et mise en œuvre pratique	60
5.2.4	Qualités d'un test	61
5.3	Un outil important : p -valeur	61
5.4	Test du rapport de vraisemblance	62
5.5	Démarche générale	64
5.6	Solutions des exercices	64
6	Tests sur les variables quantitatives	66
6.1	Objectifs	66
6.2	Test d'une espérance	66
6.2.1	Cas gaussien avec la variance connue	66
6.2.2	Cas gaussien avec la variance inconnue	68
6.2.3	Test de la normalité	72
6.2.4	Cas non gaussien	79
6.2.5	Quel(s) test(s) utilisé(s) ?	79
6.3	Test de comparaison des égalités d'espérances de deux échantillons	80
6.3.1	Cas d'échantillons appariés	80
6.3.2	Test de comparaison des égalités des variances de deux échantillons	82
6.3.3	Cas d'échantillons indépendants	84
6.3.4	Quel(s) test(s) utilisé(s) ?	88
6.3.5	Généralisation à plus de deux échantillons dans le cas gaussien	89
6.4	Test de corrélation	91
6.5	Solutions des exercices	93
7	Tests sur les variables qualitatives	96
7.1	Objectifs	96
7.2	Test d'une probabilité pour une variable de Bernoulli	96
7.2.1	Méthode exacte	96
7.2.2	Méthode approchée	97
7.3	Test de comparaison de deux probabilités pour deux échantillons appariés	98
7.4	Test de comparaison de deux probabilités pour deux échantillons indépendants	99
7.5	Généralisation à plus de deux échantillons appariés	100
7.6	Test de comparaison de plus de deux probabilités pour des échantillons indépendants	100
8	Tests non-paramétriques	102
8.1	Test de Kolmogorov-Smirnov	102
9	Autres tests	104
9.1	Botanique des tests	104
9.1.1	Tests paramétriques classiques	104
9.1.2	Tests d'adéquation d'une loi de probabilité à des données	108

Première partie

UE Statistique inférentielle

Chapitre 1

Statistique descriptive

"Quels que soient les progrès des connaissances humaines, il y aura toujours place pour l'ignorance et par suite pour le hasard et la probabilité."

Émile Borel

1.1 Avant-propos

Le but du chapitre de *Statistique Descriptive* est de comprendre les fondements de la *Statistique*. La statistique est une branche des mathématiques qui étudie les statistiques c'est-à-dire les résultats de calculs statistiques réalisés à partir des jeux de données.

Il faut faire attention à différencier le domaine des *probabilités* où nous nous donnons un modèle dans lequel nous effectuons des calculs et la *statistique* où nous observons des données et nous tentons de les expliquer en faisant éventuellement appel à un modèle probabiliste.



Attention au piège

En tant que branche des mathématiques, la statistique est une science rigoureuse qui impose un certain nombre de règles et d'obligations. Il n'est donc pas possible de faire tout ce que nous souhaitons juste parce que ce serait *joli* et que *cela montrerait mieux que nous avons raison*.

Nous verrons que la statistique est souvent mal utilisée ce qui la décrédibilise parfois aux yeux des personnes qui n'ont pas eu de formations poussées en statistique (comme la citation de Winston Churchill en témoin). Le rôle des statisticien-ne-s est de s'assurer que les conclusions obtenues à partir de jeux de données proviennent d'études rigoureuses.

Pour conclure cette introduction, nous vous invitons à faire vôtre le serment d'Hippocrate du *Data Scientist*¹ :

1. **Intégrité scientifique et rigueur** : J'exploiterai les données avec toute la rigueur requise et en conformité avec les meilleurs standards de ma profession.
2. **Transparence** : J'informerai de manière compréhensible et précise toutes les parties prenantes sur les finalités, les modalités et les implications potentielles de mon utilisation des données.
3. **Équité** : Je veillerai à toujours m'assurer que des individus ou des groupes ne soient pas discriminés par rapport à des critères illégaux ou illégitimes, de manière directe ou indirecte, sur la base de mes travaux sur les données.
4. **Respect** : J'exercerai mon activité professionnelle en respectant la vie privée et la dignité des personnes dans toutes leurs dimensions.
5. **Responsabilité et indépendance** : J'assumerai mes responsabilités en cas de manquement ou de conflit d'intérêt et je donnerai l'alerte si des actes illégaux liés à des données sont constatés.

1. un complément peut être obtenu sur la page suivante : <https://hippocrate.tech/>

1.2 Rappels/Pré-requis

Ce chapitre est essentiellement un rappel de notions généralement vues en lycée ou dans les cours de statistique² (comme dans la formation STID par exemple).

1.3 Généralité sur la démarche scientifique

Une étude statistique a pour but de répondre à une **question** dans un domaine d'application particulier à partir d'un **jeu de données**.

Une étude statistique se déroule en plusieurs étapes :

1. La réflexion sur le **protocole** à suivre pour le recueil des données (plan d'expériences, plan de sondage, élaboration d'un questionnaire ...);
2. Le **recueil** et le **codage** des données (fichier informatique, base de données);
3. L'**exploration** des données (statistique descriptive, analyse des données, fouille de données, ...), sans chercher à les modéliser;
4. Éventuellement, **prétraitement des données** (recodage, agrégation, transformations, création de nouvelles données...);
5. Si les données sont issues d'un échantillon, **modélisation statistique** : statistique inférentielle, appel à un modèle probabiliste;
6. **Prévision** et/ou **prise de décision** (réponse à la question initiale).

Chacune de ces étapes est importante pour que les résultats de l'étude répondent correctement à la question posée initialement.

La (mauvaise) utilisation de la statistique à travers les âges

Lors de la pandémie du virus appelé *covid19*, une molécule, appelée *hydroxychloroquine*, a été mise en avant médiatiquement : essentiellement deux camps se sont opposés ; l'un disant que ce médicament soignait la maladie et l'autre qu'il était inefficace voire toxique. Nous allons reprendre deux études présentées par chacun des camps pour montrer que le protocole ne permet pas de conclure ni dans un sens, ni dans l'autre.

La première étude est celle faite par Gautret et al. (2020) sur l'effet jugé positif de la molécule ; elle était appelée dans les médias comme *la première étude du professeur Raoult*. Dans cette étude, plusieurs patients ont été suivis durant 6 jours dans plusieurs établissements : l'Institut Hospitalier Universitaire Infection Méditerranée à Marseille où ils ont pris un traitement chloroquine+azythromicine et un groupe n'ayant eu aucun traitement provenant de différents établissements (Nice, Avignon et Besançon). Les différents reproches peuvent être regroupés par thématiques :

- Le premier problème fut l'évaluation de l'état des patients. En effet, dans chacun des groupes, les patients ont été testés tous les jours et, pour un même patient, il arrivait que les résultats changent entre positif et négatif chaque jour. Les tests n'étaient donc pas fiables à 100%.
- Le deuxième reproche fut le faible nombre de personnes impliquées dans l'essai clinique (à savoir 26 patients traités) et la diversité des profils. En effet, comme il y avait quasiment unicité des profils, il n'était pas possible de différencier clairement l'effet du traitement par rapport aux autres co-variables potentielles (comme l'âge). Au moment où ce polycopié est rédigé, une estimation est de 2,3% de morts. Ceci veut dire que sur 26 individus, on peut s'attendre à n'avoir que 0,6 mort ce qui est en dessous d'un mort. En conclusion, s'il n'y a pas de mort dans l'étude, ceci peut être indépendant du traitement.
- Le plus gros reproche fut l'exclusion de 6 des 26 patients pour *non-suivis*. Ceci fut problématique car si 2 patients ont choisi d'arrêter d'eux-mêmes (donc nous ne pouvons critiquer ce choix), pour les 4 autres, ils ont été sortis car 3 sont partis en soin intensif et 1 est décédé. Or, écarter les patients qui ne répondent pas positivement au traitement ne peut qu'augmenter le taux de réussite du reste de la cohorte.

2. En particulier, ce polycopié est une synthèse de celui distribué aux étudiants de STID (voir Brault (2021b) disponible à l'adresse https://www-ljk.imag.fr/membres/Vincent.Brault/Cours/Cours_1A.pdf) et celui donné aux étudiants du M1 de Mathématiques Générales (voir Brault (2021a) disponible à l'adresse <https://www-ljk.imag.fr/membres/Vincent.Brault/Cours/CoursStatistique.pdf>).



- Le quatrième reproche fut également le suivi du groupe *placebo* (c'est-à-dire le groupe qui n'avait pas pris le traitement) car il a été fait de manière plus légère (les charges virales n'ont pas été faites quotidiennement par exemple) dans des CHU différents et par des équipes différentes ; il existe donc un biais non négligeable.

La deuxième étude, dite *du Lancet* du nom du journal qui la publia, regroupait cette fois 96 032 patients atteints du covid19 et traités avec de la chloroquine ou de l'hydroxychloroquine et révélait un important sur-risque de mortalité. Après avoir été étudiées de plus près, il s'est finalement avéré que les données étaient certainement fausses et/ou corrigées arbitrairement. Nous mettons ici quelques éléments qui ont permis de mettre en évidence le problème :

- Le taux de mortalité était très supérieur à tout ce qui avait été enregistré jusqu'alors. Même si la conclusion de l'étude était la sur-mortalité du traitement, des conclusions aussi flagrantes auraient mérité une meilleure critique de la part des auteurs.
- Le nombre de morts attribués pour l'Australie était supérieur (73 patients morts dans 5 hôpitaux) que toutes les données remontées par le pays entier (68 morts en tout). De plus, un hôpital fut surpris de se retrouver dans l'étude alors qu'il n'avait jamais transmis ses données.
- Le descriptif des patients (la prise en charge autre que le traitement par exemple) était très flou. Il n'était pas possible de savoir si des co-variables auraient expliqué certaines morts.
- Les doses données ont également surpris puisque certains patients d'Amérique auraient eu des doses de 600mg alors qu'aucun hôpital n'avait donné de doses supérieures à 500mg.

Après toute cette polémique, l'article fut finalement retiré car la société qui avait fourni les données n'était pas capable de prouver leur exactitude. Nous renvoyons vers le site de Libération les statisticien-ne· intéressé-e-s : https://www.liberation.fr/checknews/2020/06/02/pourquoi-l-etude-du-lancet-sur-l-hydroxychloroquine-est-elle-sous-le-feu-des-critiques_1789844

Nous voyons bien par ces deux exemples contraires l'importance d'un protocole correctement établi. L'une de ces deux affirmations est peut-être juste, toutefois, sans une étude rigoureuse, aucune des deux n'est acceptée par la totalité de la communauté scientifique.



Attention au piège

Il arrive parfois que les statisticien-ne-s ne soient associé-e-s qu'à partir de l'étape 3. Ceci est souvent trop tard car les données récoltées (parfois des plans à plusieurs millions d'euros) ne permettent pas de répondre à la question posée.

La réflexion sur le protocole doit également être l'occasion de considérer les erreurs potentielles. Par exemple, si nous récupérons des données de manière manuscrite, il y a un risque d'erreur humaine ou de mauvaise lecture de ce qui est marqué. A l'opposé, mettre des garde fous sur des recueils virtuels peut empêcher d'avoir certains résultats. Les étapes d'exploration et de prétraitement des données sont également importantes pour repérer les erreurs potentielles.

La (mauvaise) utilisation de la statistique à travers les âges

Un exemple classique d'erreur qui a encore la vie dure dans notre inconscient collectif est le fait que les épinards soient remplis de fer :

- En 1870, le biochimiste allemand nommé E. von Wolf découvre que les épinards contiennent environ 2,7mg de fer pour 100g. L'histoire raconte que, quand sa secrétaire recopia la valeur, elle oublia la virgule ce qui multiplia par 10 la teneur.
- En 1881, un chercheur nommé Gustav von Bunge réévalue la teneur en fer mais se trompe entre le poids des épinards frais et des épinards déshydratés. Comme les épinards sont constitués à 90% d'eau, l'erreur était à nouveau de 1 pour 10.
- Dans les années 1930 à 1937, la communauté scientifique réévalua cette teneur et découvrit des deux erreurs mais elle ne réussit pas à convaincre le public. Il fallut attendre l'article de Hamblin (1981) de 1981, soit plus d'un siècle après, pour que les médias reprennent l'information.

Nous conseillons le site *Science & fourchette* recensant les nombreuses erreurs commises au fil des ans sur les épinards : <http://sciencefourchette.com/2014/04/11/popeye-est-une-supercherie/> dont nous avons tiré une partie des informations.



HISTORY

1.4 Indivus et variables statistiques

Les deux éléments importants à maîtriser dans une étude statistique sont les individus et les variables statistiques.

1.4.1 Indivus statistiques

Débutons par introduire deux définitions.

Définitions 1 (Population et individus statistiques)

La **population statistique** est l'ensemble concerné par une étude statistique. L'**individu statistique** est un élément unique de cette population.

La population statistique doit être correctement choisie en fonction de la question posée.

Exemple fil rouge

Dans l'étude fil rouge provenant de Micheaux et al. (2011) que nous allons utiliser tout au long de cette partie, nous analyserons le comportement alimentaire des personnes âgées de la région de Bordeaux en 2000. Il faut alors correctement caractériser ce qu'est *une personne âgée* (plus de 60 ans? Plus de 80 ans? Ou alors plus de 40 ans?) et la *région de Bordeaux* (est-ce la ville? L'agglomération proche? Le département qui contient la ville de Bordeaux?).



Attention au piège

Les constatations faites durant l'étude portent uniquement sur la population étudiée (en tenant compte de ses spécificités).

Par exemple, si nous étudions la survie des passagers du Titanic, la population est l'ensemble des passagers et tous les passagers. Il ne faut donc pas étudier uniquement les passagers de première classe ou ceux d'un autre paquebot et extrapoler les résultats.



La (mauvaise) utilisation de la statistique à travers les âges

Dans un tweet du 28 juin 2021, le compte @CNEWS publiait les conclusions d'une étude sur le fait que 40% des nouveaux cas de Covid en Israël étaient vaccinés (voir la gauche de la figure 1.1) avec une introduction laissant penser que la *couverture vaccinale importante* n'empêchait une large propagation du virus. Or, comme le rappelle le compte @sc_cath, malgré 87% de la population vaccinée, nous ne retrouvons que 40% dans les nouveaux cas. Ceci est dû à une plus faible contamination des personnes vaccinées (environ 10%) que des personnes non vaccinées (environ 100%).

En poussant le raisonnement à l'extrême limite, si 100% des israéliens étaient vaccinés alors le moindre cas serait forcément un vacciné et le titre serait alors *100% des nouveaux cas sont vaccinés*. La question est ici de savoir si le fait de vacciner largement la population diminue globalement le nombre de nouveaux cas ou pas par rapport à une population non vaccinée.

Remarque

Les termes *population* et *individus* font naturellement penser à des êtres humains. Néanmoins, l'individu peut également être un groupement d'êtres humains (par exemple si nous étudions le comportement de plusieurs villes, ce sera la ville l'individu statistique) ou des objets (par exemple un capteur de température).

Exemple

Si nous étudions les personnages de Star Wars, ce sont des personnes fictives qui peuvent également être des robots par exemple.



FIGURE 1.1 – Exemple d'étude où la conclusion (tweet de gauche) fut reprise sans recul sur le fait que la population était largement vaccinée et donc, sans avoir pris en compte les probabilités conditionnelles (rappelées par le tweet de droite).



Attention au piège

Le fait qu'il y ait un être humain dans une étude ne veut pas dire que l'individu statistique sera un être humain. Par exemple, si nous observons le suivi d'un traitement donné à des patients et que nous faisons des prélèvements à différents moments, l'individu statistique sera alors le couplage *patient + date du prélèvement*.

Point méthode

Un individu statistique doit être **unique**. Si vous avez un doute sur la qualification de votre individu statistique, demandez-vous si un même individu se retrouve plusieurs fois. Dans l'exemple précédent, nous voyons que si nous prenons uniquement le patient alors il sera associé à plusieurs prélèvements.

Exemple

Si nous étudions les vidéos de la plateforme *Youtube*, la caractérisation de l'unicité se fait par l'identifiant de la vidéo. En effet, une même vidéo (par exemple un clip de musique) peut être déposée par plusieurs youtubeurs et chacune de ces vidéos est pourtant unique.

Une fois la population statistique définie, le mieux serait de récupérer les informations pour tous les individus présents. En pratique, c'est bien sûr compliqué pour des questions de temps et d'argent. Nous introduisons alors de nouvelles notions :

Définitions 2 ((Sous-)échantillon)

Un **échantillon (statistique)** est un sous-ensemble de la population statistique. On appelle **sous-échantillon** une partie de cet échantillon.

Remarque

Pour plus d'information sur la façon de recueillir un échantillon, nous vous recommandons le cours (optionnel) sur les techniques de sondage du M2 SSD.

Attention au piège

Il est important de s'assurer que ce sous-ensemble soit vraiment représentatif de la population. Par exemple, si nous étudions la proportion de cancers en France dans la population et que l'échantillon ne contient que des hommes, nous aurons une sur-représentation des cancers des testicules et une sous-représentation des cancers du sein.

1.4.2 Variables statistiques

Une fois les individus correctement définis, nous devons chercher quelles seront les informations importantes à analyser pour répondre à la question.

Définition 3 (Variable statistique)

Nous appelons **variable statistique** toute information recueillie sur les individus statistiques.

Exemple

Dans le cadre d'une enquête sur la consommation des personnes âgées, nous pouvons leur demander leur âge, leur poids, leurs habitudes alimentaires (boivent-elles du café ou pas ? Si oui, combien de tasses par jour ?), la ville dans laquelle elles vivent...

Les variables statistiques peuvent être classées en deux types contenant chacun deux sous-catégories.

Définitions 4 (Types de variables)

Le premier type est celui des **variables qualitatives** qui contient toutes les variables pour lesquelles aucune opération mathématique n'est possible. Nous distinguons deux sous-catégories :

- Les variables qualitatives **nominales** lorsqu'il n'y a pas d'ordre sur ses modalités.
- Les variables qualitatives **ordinales** lorsqu'un ordre accepté par tout le monde peut être fait.

Le deuxième type est celui des **variables quantitatives** qui contient toutes les variables pour lesquelles il est possible de faire des opérations. Nous distinguons deux sous-catégories :

- Les variables quantitatives **discrètes** lorsqu'il y a un faible nombre de valeurs possibles, relevant essentiellement d'un comptage.
- Les variables quantitatives **continues** dans les autres cas.

Exemples

- Genre d'une personne : variable qualitative nominale car il n'est pas possible d'additionner *un homme + une femme* et qu'il n'y a pas d'ordre naturel et accepté par tous entre les hommes et les femmes.
- La situation familiale (célibataire, marié-e, veuf-ve, en couple...) : variable qualitative nominale car il n'est pas possible d'additionner *célibataire + en couple* par exemple et qu'il n'y a pas d'ordre naturel et accepté par tous.
- "Êtes-vous *totalelement d'accord, plutôt d'accord, plutôt pas d'accord* ou *pas du tout d'accord* sur le fait que le professeur de statistique est vraiment pédagogue ?" est une question renvoyant une réponse qui est une variable qualitative ordinale puisqu'un ordre naturel est présent mais on ne peut pas faire d'opérations arithmétiques sur les modalités de réponse.
- Le nombre d'enfants d'une personne : variable quantitative discrète puisque l'essentiel de la population se concentre sur un faible nombre d'enfants et qu'on peut calculer le *nombre d'enfants moyens par personne*.

- La pression atmosphérique : variable quantitative continue puisque nous pouvons calculer une pression moyenne et que, si on a un appareil suffisamment précis, il est possible de relever un ensemble continu de valeurs.

Généralement, il y a deux types d'erreurs dans la classification des variables.



Attention au piège (Variables qualitatives)

La première difficulté rencontrée est pour différencier une variable qualitative ordinaire d'une variable qualitative nominale. Il est important de se souvenir que *l'ordre doit être accepté par tout le monde*. Par exemple, nous pourrions demander la couleur des cheveux de plusieurs personnes et proposer un ordre suivant l'intensité de la couleur. Mais est-ce que tout le monde serait d'accord sur cet ordre ? Par exemple, est-ce que la couleur *roux* est plus intense ou moins intense que *châtain* ou *blond* ?



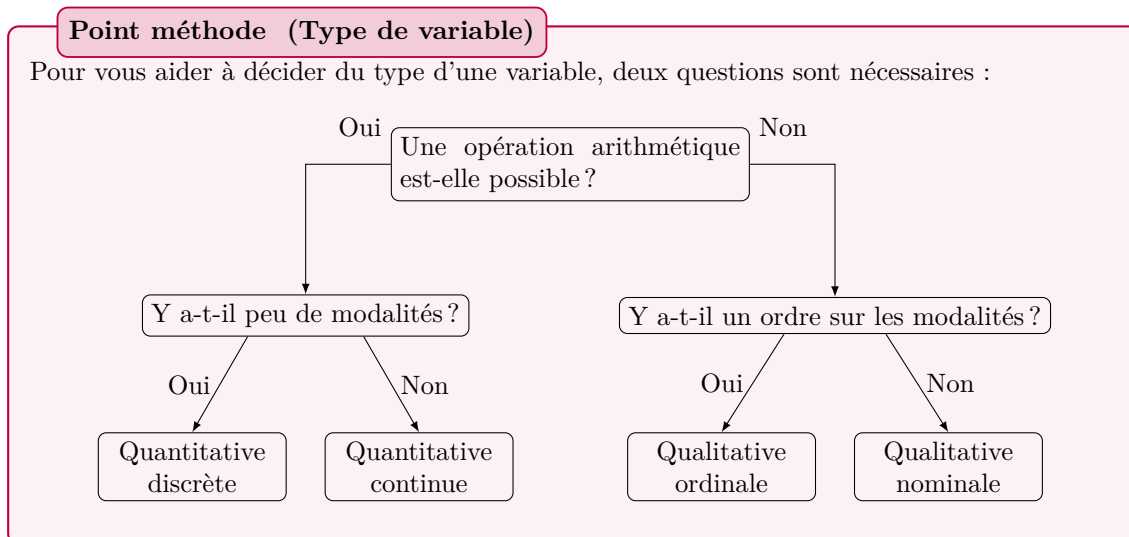
Attention au piège (Variables quantitatives)

La frontière entre *discret* et *continu* est quelque fois floue et il est possible qu'une même donnée puisse être considérée comme continue ou discrète suivant la situation.

Par exemple, si nous regardons les poids de personnes avec un outil suffisamment précis, nous pourrions avoir autant de poids que d'individus (donc une variable quantitative continue). À l'opposé, si nous regardons les poids de personnes ayant entre 50 et 55 kilogrammes avec un outil ne permettant d'avoir que des valeurs entières, nous n'aurons que 6 valeurs (donc une variable quantitative discrète).

Un autre exemple peut être le nombre de **j'aime** sur un réseau social (twitter, facebook, instagram...) : c'est un nombre entier et, généralement, les utilisateurs ont peu de notifications (donc une variable quantitative discrète) mais pour les influenceur·se·s, ce nombre peut être très grand et surtout très différent d'une publication à une autre.

Pour vous aider, voici un arbre de décision :



Indépendamment de leur type, il faut également différencier les variables qui nous permettront de répondre à la question.

Définitions 5 (Variables explicatives et à expliquer)

Nous appelons **variable à expliquer** les variables que nous cherchons à décrire, expliqué ou prédire afin de répondre à la question posée en début d'étude. Nous appelons **variable explicative** les variables utilisées pour expliquer, décrire ou prédire les variables à expliquer.

Exemple

⋮ Dans le cas des passagers du Titanic, nous cherchons à expliquer la variable *survie* en fonction de différentes variables explicatives (âge, prix du billet, port d'embarquement...).

La (mauvaise) utilisation de la statistique à travers les âges

Avec l'avènement des algorithmes automatiques et du *Big Data*, de nouvelles questions se posent quant au choix et à l'utilisation de certaines variables explicatives. Par exemple, nous nous apercevons que, mal utilisés, certains algorithmes reproduisent, voire accentuent, les inégalités déjà présentes. En 2015, des chercheurs et chercheuses ont montré que des profils masculins avaient plus de chances de se voir proposer des offres à hauts salaires que des profils féminins. Ceci est dû au fait que jusqu'à présent, les hommes ont des salaires plus élevés et, comme les concepteurs n'ont pas corrigé ce biais, les algorithmes ont cru que c'était la norme. Nous recommandons la vidéo *Algocratie : L'inégalité programmée - #DATAGUEULE 84* de la chaîne *Data Gueule* proposée par Goetz et al. (2014).

1.5 Statistiques descriptives

Dans la démarche énoncée dans la partie 1.3, la statistique descriptive prend sa place dans la partie d'exploration et de **prétraitement des données**. Ce sera à partir des résultats présentés que nous pourrons proposer une modélisation statistique (qui sera abordée dans tout le reste de la formation) qui mènera à une prévision et/ou à une prise de décision. Dans ce cadre, trois grandes catégories existent :

- Univarié.
- Bivarié.
- Multivarié.

La statistique descriptive désigne un ensemble de techniques dont le but est de

- explorer, découvrir l'information contenue dans les données,
- les représenter graphiquement,
- détecter des premières tendances.

A chacun de ces buts correspond une ou plusieurs techniques (voir le tableau 1.1).

TABLE 1.1 – Mise en relation des objectifs (à gauche) avec les outils possibles (au milieu). La dernière colonne correspond au cours dans lequel nous verrons ces notions.

Objectif	Techniques	Cours
Explorer, découvrir les données, résumer l'information	tableaux statistiques résumés statistiques	Univarié
Représenter graphiquement	graphiques	Univarié
Détecter des tendances	indicateurs de liaison	Bivarié et multivarié

Il est important de ne pas utiliser une technique pour un autre objectif. De même, nous verrons par la suite que chaque outil correspond à un ou plusieurs objectifs précis : leur choix est donc primordial.

La suite de cette partie consiste à présenter ces outils suivants les catégories (univariée, bivariée puis multivariée).

1.5.1 Statistiques descriptives univariées

Pour la description, les statistiques descriptives univariées permettent d'observer les variables indépendamment des autres. Nous pouvons les décomposer en plusieurs grandes catégories :

- Tableau
- Résumés statistiques
 - ◆ de position
 - ◆ de dispersion
 - ◆ de forme



- Graphique

TABLE 1.2: Tableau récapitulatif des outils et de quand les appliquer. Le symbole 🚩 signifie qu'il faut faire attention pour l'utiliser (par exemple, que sur les variables positives) ou qu'il y a une différence vis-à-vis de la définition usuelle (comme pour les fractiles dans le cas discret).

Nom	Qualitative		Quantitative	
	Nominale	Ordinale	Discrète	Continue
Tableaux				
Tri à plat	✓	✓	✓	
Résumés statistiques de position				
Centile			🚩	✓
Classe modale				✓
Décile			🚩	✓
Fractile			🚩	✓
Médiane			🚩	✓
Maximum			✓	✓
Minimum			✓	✓
Mode	✓	✓	✓	
Moyenne			✓	✓
Quartile			🚩	✓
Résumés statistiques de dispersion				
Écart-type			✓	✓
Étendue			✓	✓
Longueur interquartile			✓	✓
Rapport interdécile			🚩	🚩
Variance			✓	✓
Résumés statistiques de forme				
Kurtosis			✓	✓
Skewness			✓	✓
Graphiques				
Boxplot ou boîte à moustaches			✓	✓
Diagramme circulaire	🚩	🚩		
Diagramme de Pareto	✓			
Diagramme empilé	✓	✓		
Diagramme en barres ou tuyaux d'orgues	✓	✓		
Diagramme en bâtons			✓	
Fonction de répartition empirique			✓	✓
Histogramme				✓
Polygone régulier				✓

Dans le tableau 1.2, tous les outils ont été récapitulés. Le symbole 🚩 signifie qu'il y a une particularité vis-à-vis de l'utilisation ou de la définition classique. Notamment, il est important de noter :

- Dans le cas discret, les fractiles (médiane, quartiles, déciles et centiles notamment) doivent appartenir à l'ensemble de définition (par exemple, un nombre entier si on parle du nombre d'enfants dans un foyer). Les logiciels utilisent par défaut la version continue qui peut faire des pondérations si le fractile se trouve entre deux valeurs distinctes.
- Le rapport interdécile ne s'utilise que sur des variables strictement positives.
- Le diagramme circulaire ne s'utilise que quand il y a un nombre très réduit (typiquement quatre ou moins) de modalités effectives (c'est-à-dire de modalités prises par au moins un individu).

Une autre façon de résumer l'information est de faire un diagramme en mettant les outils dans la zone qui correspond aux types de variables sur lesquelles les outils s'appliquent (voir la figure 1.2).

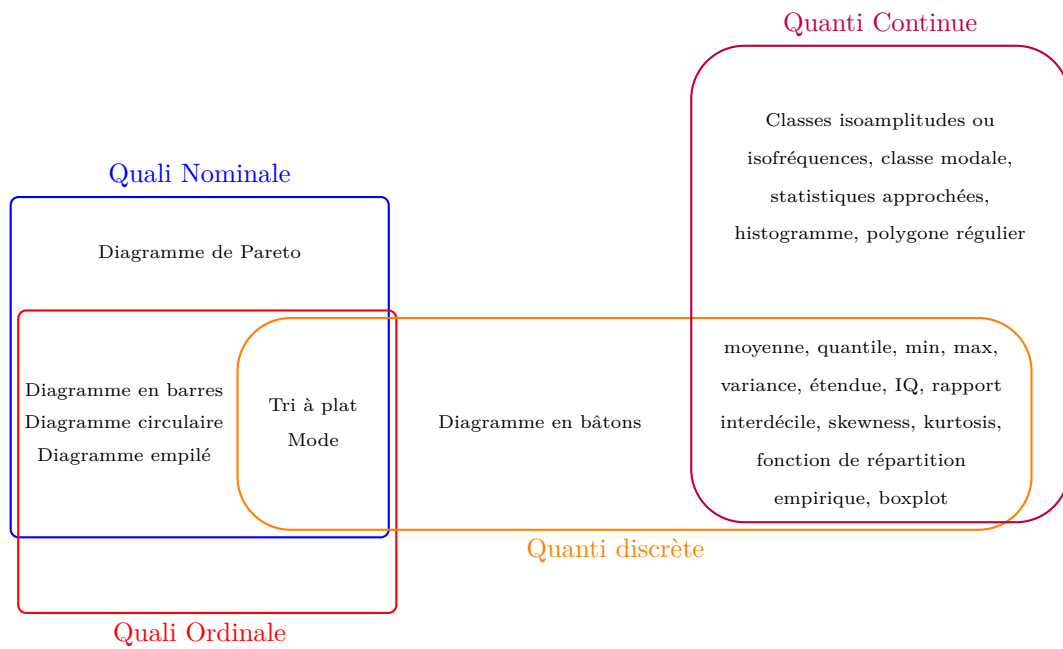


FIGURE 1.2 – Diagramme des outils en fonctions des types de variables.

Tableaux

Pour présenter le tri à plat, nous avons besoin de quelques définitions.

Définitions 6 (Effectifs et fréquences)

Étant donnée une modalité a_k de la variable qualitative \mathbf{X} , nous appelons **effectif de la modalité**, noté N_k , le nombre d'individus prenant la modalité a_k . Nous appelons **fréquence de la modalité**, notée f_k , la proportion d'individus prenant la modalité a_k parmi ceux qui ont donné une réponse. Nous avons donc les relations suivantes :

$$f_k = \frac{N_k}{N} \text{ et } N_k = N f_k$$

sous la condition que $N > 0$.

Dans le cas de variables qualitatives ordonnées et de quantitatives discrètes, nous introduisons le principe de fréquence cumulée :

Définition 7 (Fréquence cumulée pour variable qualitative ordonnée)

Étant donnée une modalité a_k de la variable qualitative **ordonnée** \mathbf{X} , nous appelons **fréquence cumulée de la modalité**, notée F_k , la proportion d'individus prenant la modalité a_k ou une modalité inférieure. Nous avons donc les relations suivantes :

$$F_k = \sum_{\ell=1}^k f_{\ell} = \frac{1}{N} \sum_{\ell=1}^k N_{\ell}.$$

Une autre façon de l'écrire est :

$$F_k = \sum_{\ell \leq k} f_{\ell} = \frac{1}{N} \sum_{\ell \leq k} N_{\ell}.$$

Avec ces notations, nous pouvons introduire le tableau de tri à plat.

Définition 8 (Tri à plat)

Un **tri à plat** est un tableau récapitulatif toutes les informations précédentes sous la forme suivante :

Qualitative nominale			Qualitative ordinale ou quantitative discrète			
Modalités	Effectifs	Fréquences	Modalités	Effectifs	Fréquences	Cumulées
a_1	N_1	f_1	a_1	N_1	f_1	$F_1 = f_1$
a_2	N_2	f_2	a_2	N_2	f_2	F_2
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
a_k	N_k	f_k	a_k	N_k	f_k	F_k
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
a_K	N_K	f_K	a_K	N_K	f_K	$F_K = 1$
Total	N	1	Total	N	1	

Si la variable contient des données manquantes (c'est-à-dire si $N < n$), nous précisons le pourcentage de répondants lorsque nous présenterons le tableau (généralement, nous le mettons dans la légende du tableau ; voir le chapitre ??).

Résumés statistiques de dispersion

Dans cette partie, nous présentons quelques résumés statistiques de dispersion.

Définition 9 (Étendue)

L'**étendue** (*range* en anglais) de la distribution x_1, \dots, x_n , notée W , est définie comme l'écart (positif) entre la plus grande et la plus petite valeur :

$$W = \max_{1 \leq i \leq n} x_i - \min_{1 \leq i \leq n} x_i.$$

Nous introduisons maintenant la longueur de l'intervalle inter-quartile.

Définition 10 (Intervalle interquartile)

L'**intervalle interquartile** de la distribution x_1, \dots, x_n est l'intervalle de bornes Q_1 et Q_3 ; c'est-à-dire $[Q_1; Q_3]$. La **longueur de cet intervalle** est un indicatif de dispersion, noté IQ , valant :

$$IQ = Q_3 - Q_1.$$

**Attention au piège**

L'erreur la plus commune faite par les étudiants est de dire que *l'intervalle interquartile vaut 5* par exemple : il y a confusion entre l'intervalle (c'est-à-dire l'objet $[Q_1; Q_3]$) et sa longueur.

Point méthode (Vérifications)

Comme nous avons les relations suivantes

$$\min_{1 \leq i \leq n} x_i \leq Q_1 \leq Q_3 \leq \max_{1 \leq i \leq n} x_i$$

alors nous devons vérifier les deux points suivants :

1. $IQ \geq 0$ (comme pour W).
2. $IQ \leq W$.

Résumés statistiques de forme

Les résumés statistiques de forme sont basés sur la forme de la distribution d'une gaussienne.

Définition 11 (Coefficient d'asymétrie)

Le **coefficient d'asymétrie** (ou *skewness* en anglais) mesure la dissymétrie de la distribution. Sa formule vaut :

$$G_1(x) = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\sigma_x^3}$$

où \bar{x} est la moyenne empirique des valeurs des observations de la variable X et σ_x est son écart-type.

Remarque

Plus la distribution sera symétrique, plus le coefficient sera proche de 0.

Définition 12 (Coefficient d'aplatissement)

Le **coefficient d'aplatissement** (ou *kurtosis*) s'utilise dans le cadre d'une distribution symétrique (donc si le précédent coefficient est proche de 0) et mesure la répartition des poids. Sa formule vaut :

$$G_2(x) = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{\sigma_x^4} - 3.$$

Remarque

Il sera proche de 0 si la distribution ressemble à une loi gaussienne et négatif si les points sont répartis uniformément.

Graphiques

Dans cette partie, nous présentons les graphiques.

Boxplot ou boîte à moustaches : La boîte à moustaches est un diagramme généralement présenté horizontalement dans les cours de lycée. Nous faisons le choix de le représenter verticalement comme le font de nombreux logiciels.

Définition 13 (Boxplot ou boîte à moustaches (variante))

La variante du **boxplot** ou **boîte à moustaches** est la suivante :

1. Nous traçons 3 traits horizontaux de mêmes longueurs placés aux niveaux de chacun des 3 quartiles centraux : Q_1 , médiane et Q_3 .
2. Deux traits relient verticalement les extrémités du premier (Q_1) et du troisième quartiles (Q_3). Ce trait est censé passer par les extrémités du trait horizontal symbolisant la médiane.
3. Nous calculons la longueur interquartile IQ .
4. Nous regardons si le maximum est plus grand que $Q_3 + 1,5IQ$:
 - Si la réponse est oui :
 - (a) Nous traçons un trait horizontal à l'ordonnée $Q_3 + 1,5IQ$ (plus petit que les autres traits horizontaux).
 - (b) puis nous relierons le centre de ce trait avec le centre du trait symbolisant le troisième quartile Q_3 .
 - (c) Nous mettons une croix rouge à l'emplacement du maximum et dans l'alignement des centres des traits horizontaux.

(d) (Facultatif) Enfin, nous mettons des croix bleues pour chaque modalité effective (c'est-à-dire ayant un effectif non nul) entre $Q_3 + 1,5IQ$ et le maximum (également dans l'alignement des centres des traits horizontaux).

- Si la réponse est non (et que le maximum est donc plus petit que $Q_3 + 1,5IQ$), nous procédons de la sorte :

(a) Nous traçons le trait horizontal correspondant au maximum.

(b) Nous relierons ce trait au trait symbolisant le troisième quartile.

5. Nous faisons la procédure symétrique avec le bas du graphique en comparant le minimum avec la valeur $Q_1 - 1,5IQ$.

Dans le cas où le maximum serait plus grand que $Q_3 + 1,5IQ$ et/ou que le minimum serait plus petit que $Q_1 - 1,5IQ$, nous devons préciser combien d'individus se trouvent dans les intervalles

$\left[\min_{1 \leq i \leq n} x_i; Q_1 - 1,5IQ \right]$ et $\left[Q_3 + 1,5IQ; \max_{1 \leq i \leq n} x_i \right]$; nous appelons ces points des **outsiders**.

Remarque

Il existe une autre variante consistant à ne pas dépasser $Q_1 - 3IQ$ et $Q_3 - 3IQ$: si le maximum (resp. le minimum) est plus grand que $Q_3 - 3IQ$ (resp. plus petit que $Q_1 - 3IQ$), nous mettons la croix rouge au niveau de $Q_3 + 3IQ$ (resp. $Q_1 - 3IQ$) et nous ne mettons pas de croix bleues au delà de cette limite.

Point logiciel



Dans le logiciel **Excel** , il n'existe pas de fonctions pour en construire simplement ; nous verrons en TP une procédure pour le faire à partir de nuage de points. Dans le langage **R** , nous pouvons utiliser la fonction `boxplot` ou le package `ggplot2` par exemple.

Diagramme circulaire : Ce diagramme est connu aussi sous le nom de *camembert*.

Définition 14 (Diagramme circulaire)

Le **diagramme circulaire** consiste à découper un disque en plusieurs zones de telle sorte que :

- Chaque zone est délimitée par deux rayons et une partie du cercle.
- Chaque angle (et donc chaque zone) est proportionnel à la fréquence de la modalité associée.



Attention au piège

Bien que très visuel, le diagramme circulaire rend les comparaisons entre fréquences des modalités plus difficiles à faire car nous ne sommes pas vraiment habitués à comparer des angles (par opposition à la comparaison des fréquences dans le cas des diagrammes en tuyaux d'orgue). En particulier, il est fortement déconseillé d'utiliser ce diagramme si :

- Les fréquences sont assez proches.
- Il y a beaucoup de modalités.

Point méthode (Calcul de l'angle)

Comme un cercle fait au total 360° et que chaque région est proportionnelle à la fréquence de la modalité, nous calculons l'angle de chaque région par la formule $360^\circ \times f_k$ où f_k est la fréquence de la modalité a_k .

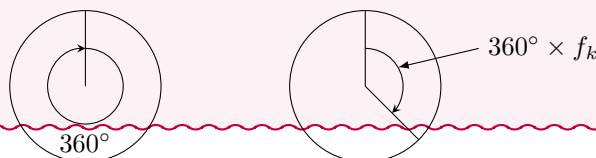


Diagramme en barres et diagramme de Pareto : Pour présenter ces diagrammes, nous commençons par celui en barres.

Définition 15 (Diagramme en tuyaux d'orgue ou en barres)

Le **diagramme en tuyaux d'orgue** ou **en barres** consiste à représenter les fréquences des modalités d'une variable qualitative à l'aide de barres rectangulaires verticales parallèles de telle sorte que :

- Toutes les barres ont la même largeur à la base.
- La hauteur de chaque barre représente la valeur de la fréquence (l'axe des ordonnées est gradué suivant les valeurs des fréquences).
- Les barres sont toutes espacées de la même distance (elles ne se collent pas).
- Le nom de chaque modalité doit apparaître en dessous de chaque barre.
- Dans le cas de variables qualitatives ordinales, il faut conserver l'ordre naturel des modalités.

Le diagramme de Pareto est un diagramme en tuyaux d'orgue particulier.

Définition 16 (Diagramme de Pareto)

Le **diagramme de Pareto** est un diagramme en tuyaux d'orgue où les modalités sont réorganisées par ordre décroissant de fréquences.



Attention au piège

Comme nous intervenons sur l'ordre des modalités, le diagramme de Pareto ne s'utilise jamais pour des variables qualitatives ordonnées.

Diagramme empilé : Le diagramme empilé permet une représentation rapide des fréquences cumulées.

Définition 17 (Diagramme empilé)

Le **diagramme empilé** consiste à représenter les fréquences des modalités à l'aide de rectangles empilés les uns au-dessus des autres de telle sorte que :

- Tous les rectangles ont la même largeur.
- La hauteur de chaque rectangle correspond à la fréquence de chaque modalité.
- Les rectangles sont empilés les uns au-dessus des autres de telle sorte que la hauteur totale fasse 100%.

Diagramme en bâtons : Le diagramme en bâtons est spécifique aux variables quantitatives discrètes.

Définition 18 (Diagramme en bâtons)

Le **diagramme en bâtons** consiste à représenter les fréquences des modalités d'une variable quantitative discrète à l'aide de segments verticaux de coordonnées $(a_k, 0)$ et (a_k, f_k) :

- Le segment est parallèle à l'axe des ordonnées (ou perpendiculaire à l'axe des abscisses suivant le point de vue).
- Il se situe en abscisse au niveau de la modalité a_k .
- Il va de 0 à f_k .

Éventuellement, on peut rajouter un point en haut de chaque segment.

Fonction de répartition empirique : La fonction de répartition empirique est une fonction dont le graphique est peu utilisé alors qu'il contient un très grand nombre d'informations.

Définition 19 (Fonction de répartition empirique)

La **fonction de répartition empirique** est une fonction définie pour tout $t \in \mathbb{R}$ par :

$$\widehat{F}(t) = \sum_{\substack{k=1 \\ \text{tel que } a_k \leq t}}^K f_k = \sum_{k=1}^K f_k \mathbb{1}_{\{a_k \leq t\}}$$

où $\mathbb{1}_{\{a_k \leq t\}}$ vaut 1 si et seulement si $a_k \leq t$ et 0 sinon.

Proposition 1 (Propriété de la fonction de répartition empirique)

Les propriétés de la fonction de répartition empirique sont les suivantes :

- Elle est croissante et constante par morceaux (nos disons aussi qu'elle est *en escalier croissant*).
- Chaque saut se fait au moment d'une modalité avec une fréquence non nulle.
- La hauteur de chaque saut correspond à la fréquence de la modalité associée.
- Elle vaut 0 avant la première modalité et 1 après la dernière modalité.
- Elle est continue à droite ; c'est-à-dire que, au moment de la modalité a_k la fonction vaut toujours F_k qui est la valeur du plateau du-dessus.

1.5.2 Statistiques descriptives bivariées

Le principe des croisements est de regarder les corrélations possibles entre deux variables. Ainsi, il est intéressant comment une variable quantitative se comporte vis-à-vis d'une autre variable quantitative.



Attention au piège

Attention, tous les croisements ne se représentent pas (voir le tableau 1.6). Par exemple, il est possible de regarder les salaires de personnes en fonction de leur sexe (une variable quantitative en fonction d'une variable qualitative). L'autre sens est plus compliqué à représenter. Dans le cas de la description, cela n'a pas d'importance car seule une corrélation est recherchée (pas une causalité).

Nous voyons dans le tableau 1.6 qu'il y a un indicateur particulier pour le cas d'un croisement entre deux ordinales. Si l'une des deux variables n'est pas une ordinale (par exemple, le croisement d'une nominale avec une ordinale), il n'existe plus d'indicateurs particuliers : il faut se contenter de considérer que l'ordinale est une qualitative (presque) comme une autre.

TABLE 1.3 – Résumé des possibilités d'expliquer une variable (colonnes) en fonction d'une variable (lignes) : ✓ s'il existe un indicateur de liaison, ✗ s'il n'y a pas de théorie et 🤝 si c'est possible mais qu'il faut prendre un indicateur plus général (dans le cas où il faut se contenter qu'une ordinale est une qualitative).

à expliquer en fonction de	Qualitative	Quantitative	Ordinale
Qualitative	✓	✓	🤝
Quantitative	✗	✓	✗
Ordinale	🤝	🤝	✓

Comme précédemment, il y a trois groupes d'outils qui peuvent être utilisés :

- Tableau (voir le tableau 1.4)

- Graphique (voir le tableau 1.5)
- Indicateurs de liaisons (voir le tableau 1.6).

TABLE 1.4 – Tableaux qui peuvent être calculés pour chaque croisement.

à expliquer en fonction de	Qualitative	Quantitative	Ordinale
Qualitative	Contingence, distribution conjointe, distribution conditionnelle, effectifs théoriques	Tri à plat par sous populations	
Quantitative		Tableau des variables centrées	
Ordinale			Tableau des concordances

TABLE 1.5 – Graphiques qui peuvent être calculés pour chaque croisement.

à expliquer en fonction de	Qualitative	Quantitative
Qualitative	Diagrammes empilés, Diagrammes en barres	Boxplots
Quantitative		Nuages de points

TABLE 1.6 – Indicateurs de liaisons qui peuvent être calculés pour chaque croisement.

à expliquer en fonction de	Qualitative	Quantitative	Ordinale
Qualitative	Khi2 χ_n^2 , Phi2 ϕ^2 , V^2 de Cramer	Variances intra et inter, rapport de corrélation η^2	
Quantitative		Covariance, coefficient de corrélation linéaire (de Pearson)	
Ordinale			Tau de Kendall

Il est important de souvenir que les indicateurs du tableau 1.6 ne permettent qu'une description : il faudra confirmer par un test ensuite (voir le cours de *Compléments en tests statistiques*).



Attention au piège

Attention, une corrélation ne veut pas dire qu'il y a une causalité. Par exemple, il est observé une corrélation entre les coups de soleil et une augmentation dans les ventes de glaces. Même si cette corrélation peut être intéressante mais il ne faut pas voir une causalité : par exemple, en tant que vendeur de glaces, je peux me dire qu'il faut que j'ouvre mon échoppe si je vois qu'il y a beaucoup de coups de soleil, cela ne veut pas dire que les coups de soleil impliquent que je vais vendre plus de glaces.

En particulier, il arrive que certaines co-variables n'aient pas été prises en compte dans le cadre de l'étude et viennent fausser les résultats. Nous renvoyons par exemple les statisticien-ne-s intéressé-e-s à la vidéo *Chocolat, corrélation et moustaches de chats*³ de la chaîne Youtube *la statistique expliquée à mon chat*⁴ proposée par Uyttendaele et al. (2016).

3. Disponible à l'url suivante : <https://www.youtube.com/watch?v=aOX0pIwBCvw>

4. L'url de la chaîne est la suivante : https://www.youtube.com/channel/UCWty1tzwZW_ZNSp5GVGteaA



La (mauvaise) utilisation de la statistique à travers les âges

Pour sensibiliser à la différence entre corrélation et causalité, *les décodeurs* du journal *Le monde* proposent un générateur de corrélation aléatoire :

https://www.lemonde.fr/les-decodeurs/article/2019/01/02/correlation-ou-causalite-brillez-en-societe-avec-notre-generateur-aleatoire-de-comparaisons-absurdes_5404286_4355770.html

N'hésitez pas à vous amuser avec.

1.5.3 Statistiques descriptives multivariées

Les statistiques descriptives multivariées portent sur des techniques impliquant plus de deux variables. Par exemple, l'*Analyse en Composantes Principales* (ou ACP) ou la *Classification Ascendante Hiérarchique* (CAH) peuvent être citées. Toutes ces méthodes sont abordées dans le cours d'*analyse de données* ou encore dans des livres tels que Husson et al. (2016) ou sur le site *Wikistat* de Besse.

Chapitre 2

Modélisation

"Le loto, c'est un impôt sur les gens qui ne comprennent pas les statistiques."

2.1 Objectifs

L'objectif de ce court chapitre est de comprendre les fondements de la statistique, la philosophie et le principe de modélisation.

2.2 Introduction

Commençons par quelques exemples d'application de la statistique :

Exemple

Adapté d'un exemple tiré du livre de Robert (2006). Considérant la proportion proportion de naissances masculines à Paris, Pierre-Simon de Laplace s'intéresse à savoir s'il y a plus de naissances d'hommes que de femmes. Pour cela, il va recenser 251 527 naissances masculines et 241 945 naissances féminines en 1785.

Une naissance peut être considérée comme le résultat d'une variable X_i de Bernoulli de paramètre $\theta \in [0; 1]$:

$$X_i = \begin{cases} \text{homme} & \text{avec probabilité } \theta, \\ \text{femme} & \text{avec probabilité } 1 - \theta. \end{cases}$$

X_1, \dots, X_n est alors l'échantillon des $n = 493\,472$ résultats de naissances observées. La question qu'on peut alors se poser est de savoir si $\theta = 1/2$, $\theta < 1/2$ ou $\theta > 1/2$.

Exemple

Adapté d'une histoire sur la *cryptographie automatique* de Ycart (2017). Une ancienne technique de cryptographie consiste à associer chaque lettre de l'alphabet à une autre de façon unique : nous obtenons ainsi une bijection entre un texte lisible et un texte codé. Pour décoder, il suffit de faire la transformation inverse. Par exemple, avec le tableau de correspondance suivant :

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
K	R	D	U	C	W	M	G	B	E	F	P	N	H	I	Y	O	S	X	J	A	V	L	T	Q	Z

nous pouvons décrypter le texte :

KYYSCHUSC PK XJKJBXJBOAC CXJ SCYIXKHJ¹

Une technique de *cryptanalyse* pourrait être d'essayer toutes les combinaisons possibles mais cela représente $26! \approx 4 \times 10^{26}$ possibilités pour la langue française (à la raison d'une possibilité toutes les secondes, cela prendrait un peu moins de 10^{10} fois l'âge de l'univers). Une autre idée a été proposée dès le 9^{ème} siècle par Al-Kindi basée sur une méthode statistique : en fait, il s'est aperçu que toutes les lettres ne revenaient pas de la même façon. Par exemple, dans la langue française, le *e* est celui qui revient le plus souvent (avec une fréquence de 12,10% pour le corpus de wikipédia en 2008 par exemple

si nous séparons les e sans accent de ceux avec accent ²).

L'idée derrière est de dire que chaque lettre est le résultat d'une variable aléatoire X_i prenant ses valeurs dans $\{A, B, \dots, Z\}$ avec une probabilité (p_A, p_B, \dots, p_Z) . En calculant les fréquences dans le texte, nous pouvons essayer de retrouver le codage.

Notons toutefois que ce système demande des textes codés relativement longs. Pour accélérer l'estimation, il est conseillé de coupler avec les probabilités conditionnelles d'une lettre sachant une autre ; par exemple, la probabilité d'avoir deux w successifs est quasiment nulle.

En particulier, la procédure proposée par Alan Turing (1912-1954) pour décrypter Enigma fonctionna car les allemands envoyés régulièrement les bulletins météorologiques avec une chaîne de caractère régulière et assez longue à chaque fois (voir par exemple Lehning (2006)).

Exemple

Dans le film de présentation ³, nous nous étions intéressés à la hauteur d'un fleuve en fonction des pluies tombées dans les montagnes. Pour cela, nous supposons que chaque observation est le résultat d'une variable aléatoire Y_i qui dépend d'une co-variable x_i (la quantité de pluie tombée dans les montagnes) de façon polynomiale ; c'est-à-dire qu'il existe un polynôme P_m de degrés m inconnu et de coefficients inconnus tels que pour tout $i \in \mathbb{N}$:

$$Y_i = P_m(x_i) + \varepsilon_i \text{ avec } \varepsilon_i \stackrel{\text{i.i.d}}{\sim} \mathcal{N}(0, \sigma^2)$$

avec σ également inconnu. La question est alors d'estimer au mieux chacun des paramètres notamment le degré m du polynôme car la théorie de la sélection de modèle nous permet de voir que le polynôme qui collerait le mieux aux données serait un polynôme de degrés n (passant par toutes les observations) mais ce serait également le plus éloigné du vrai polynôme ⁴. Ce phénomène est appelé du *surapprentissage*.

2.3 Modèle statistique

Commençons par quelques définitions.

Définitions 20 (Modèle statistique)

Un **modèle statistique** est la donnée d'un espace mesurée (E, \mathcal{E}) et d'une famille $(\mathbb{P}_\theta)_{\theta \in \Theta}$ de mesures de probabilité d'indexation Θ . Le modèle associé est noté $(E, \mathcal{E}, \mathbb{P}_\theta, \theta \in \Theta)$.

Quand il existe $d \in \mathbb{N}^*$ tel que $\Theta \subset \mathbb{R}^d$, nous disons que le modèle est **paramétrique**. Sinon, il est dit **non-paramétrique**.

Remarque

θ est inconnu mais son ensemble d'appartenance Θ est connu.

Exemple

Pour la loi de Bernoulli, nous savons que $\theta \in [0; 1]$.

Attention au piège

Le plus souvent, l'indexation se fait sur une seule famille de probabilités. Dans de très rares cas, nous pouvons imaginer une concurrence entre plusieurs familles de lois. Par exemple, sur l'ensemble fini \mathbb{N} , nous pouvons mettre en concurrence des lois de Poisson de paramètres $\lambda \in \mathbb{R}^+$ et des lois géométriques de paramètres $p \in [0; 1]$.

1. La solution étant **APPRENDRE LA STATISTIQUE EST REPOSANT** .

2. Voir par exemple https://fr.wikipedia.org/wiki/Fr%C3%A9quence_d%27apparition_des_lignes_en_fran%C3%A7ais

3. disponible ici : <https://www.youtube.com/watch?v=wQ-1Zj9sB1M>

4. Une vidéo pour illustrer ce problème est disponible ici : <https://www.youtube.com/watch?v=YuP1d0IbWfQ>

La démarche statistique consiste à donner de l'information sur le paramètre θ en s'appuyant sur la notion d'observation ; on parle alors d'**inférence** ou de **statistique inférentielle**.

Définition 21 (Observation)

Une **observation** X est une variable aléatoire à valeur dans E et dont la loi appartient à la famille $(\mathbb{P}_\theta)_{\theta \in \Theta}$.

Point méthode (Hypothèse fondamentale)

L'hypothèse fondamentale de la statistique est de supposer qu'une observation X suit une loi \mathbb{P}_X et qu'il existe un paramètre $\theta \in \Theta$ tel que $\mathbb{P}_X = \mathbb{P}_\theta$.

Remarque

Dans le langage courant, nous confondons souvent l'observation X qui est une variable aléatoire et sa réalisation x qui est fixée.

Définition 22 (n -échantillon)

Étant donné $n \in \mathbb{N}^*$, un **n -échantillon de loi** \mathbb{P}_X est la donnée de n variables aléatoires X_1, \dots, X_n indépendantes et identiquement distribuées selon la loi \mathbb{P}_θ .

Remarque

Lorsque l'observation X a la structure d'un n -échantillon de loi \mathbb{P}_θ alors $\mathbb{P}_X = (\mathbb{P}_\theta)^{\otimes n}$.

Exercices 2.1

Reprenre les exercices en début de cours et dire si les modèles sont paramétriques ou non-paramétriques. Dans le cas de modèles paramétriques, donnez l'espace des paramètres. Dans tous les cas, donner l'espace des observations.

Exemple

Une entreprise de téléphonie pense que le nombre d'appels journaliers passés par chaque client suit une loi de Poisson et souhaite tester cette hypothèse pour pouvoir fixer les tarifs. Pour ce faire, elle va choisir aléatoirement des clients et récupérer le nombre d'appels. Formellement, les observations X_1, \dots, X_n sont supposées indépendantes et de même loi et l'entreprise se demande donc si celle-ci suit une loi de Poisson et, dans ce cas, quel est le paramètre ?

Remarque

Il peut arriver que l'exigence d'indépendance soit affaiblie voir totalement abandonnée. C'est par exemple le cas dans la plupart des analyses de série chronologique (Voir par exemple le cours de *Données dépendantes 1 : séries temporelles*).

Point méthode (Quelle loi choisir ?)

La problématique du choix de la loi dans le cas de modèles non-paramétriques peut être simplifiée en regardant l'ensemble de définition. Voici quelques exemples classiques :

- Si l'ensemble est dénombrable :
 - S'il n'a que deux éléments : loi de Bernoulli.
 - Si c'est $\{0, \dots, n\}$: loi binomiale, loi uniforme, loi hypergéométrique.
 - Si c'est \mathbb{N} ou \mathbb{N}^* : loi de Poisson, loi géométrique, loi binomiale négative.

- Si l'ensemble est non dénombrable :
 - Si c'est l'intervalle $[0, 1]$: loi uniforme, loi bêta.
 - Si c'est un intervalle borné : loi uniforme.
 - Si c'est \mathbb{R}^+ : loi exponentielle, loi gamma, loi inverse gamma.
 - Si c'est $[a, +\infty[$ avec $a > 0$: loi de Pareto
 - Si c'est \mathbb{R} : loi gaussienne ou normale, loi de Cauchy, loi de Laplace.

Une autre solution consiste à utiliser les liens entre les lois. Pour cela, nous renvoyons à la figure 2.1.

Définition 23 (Identifiabilité)

Étant donné un modèle paramétré par θ , nous dirons que le modèle est **identifiable** si l'application $\theta \mapsto \mathbb{P}_\theta$ est injective.

Contre-exemple



Le contre-exemple le plus simple est celui de la loi gaussienne centrée de variance inconnue $\mathcal{N}(0, \theta^2)$:

- Si $\Theta = \mathbb{R}_+^*$, le modèle est identifiable.
- Si $\Theta = \mathbb{R}^*$, le modèle n'est pas identifiable car pour tout $\theta \neq 0$, les paramètres θ et $-\theta$ fournissent la même loi.

Exercices 2.2



Pour les lois mentionnées dans le point méthode précédent, donner l'espace de paramètre Θ de telle sorte que les lois correspondent à des modèles identifiables. Pour rappel, nous mettons les densités :

Nom	Abréviation	Support	Densité $f(x)$ ou $\mathbb{P}(X = x)$
Bernoulli	$\mathcal{B}(p)$	$\{0, 1\}$	$p^x(1-p)^{1-x}$
Bêta	$\mathcal{Be}(a, b)$	$]0, 1[$	$\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1}(1-x)^{b-1}$
Binomiale	$\mathcal{Bin}(n; p)$	$\{0, 1, \dots, n\}$	$\binom{n}{x} p^x (1-p)^{n-x}$
Binomiale négative	$\mathcal{NBin}(r; p)$	\mathbb{N}	$\binom{x+n-1}{x} p^n (1-p)^x$
Cauchy	$\mathcal{Cau}(\mu, a)$	\mathbb{R}	$\frac{\pi}{1 + (\frac{x-\mu}{a})^2}$
Exponentielle	$\mathcal{E}(\lambda)$	\mathbb{R}^+	$\lambda e^{-\lambda x}$
Gamma	$\mathcal{Ga}(\lambda)$	\mathbb{R}_+^*	$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$
Gaussienne ou normale	$\mathcal{N}(\mu, \sigma^2)$	\mathbb{R}	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$
Géométrique	$\mathcal{Géo}(p)$	\mathbb{N}^*	$(1-p)^{x-1} p$
Hypergéométrique	$\mathcal{H}(n, p, A)$	$\{0, 1, \dots, n\}$	$\frac{\binom{pA}{x} \binom{(1-p)A}{n-x}}{\binom{A}{n}}$
Inverse gamma	$\mathcal{IGa}(\lambda)$	\mathbb{R}_+^*	$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{-\alpha-1} e^{-\frac{\beta}{x}}$
Laplace	$\mathcal{Laplace}(\mu, b)$	\mathbb{R}	$\frac{1}{2b} e^{-\frac{ x-\mu }{b}}$
Pareto	$\mathcal{Pa}(a, k)$	$[a, +\infty[$	$\frac{ka^k}{x^{k+1}}$
Poisson	$\mathcal{P}(\lambda)$	\mathbb{N}	$e^{-\lambda} \frac{\lambda^x}{x!}$
Uniforme	$\mathcal{U}(\{0, 1, \dots, n\})$	$\{0, 1, \dots, n\}$	$\frac{1}{n+1}$
Uniforme	$\mathcal{U}([a, b])$	$[a, b]$	$\frac{1}{b-a}$

Hypothèse

Dans la suite de ce cours, nous supposons généralement avoir un modèle statistique identifiable $(E, \mathcal{E}, \mathbb{P}_\theta, \theta \in \Theta)$. De plus, nous noterons souvent $\mathbf{X} = (X_1, \dots, X_n)$ pour n -échantillon.

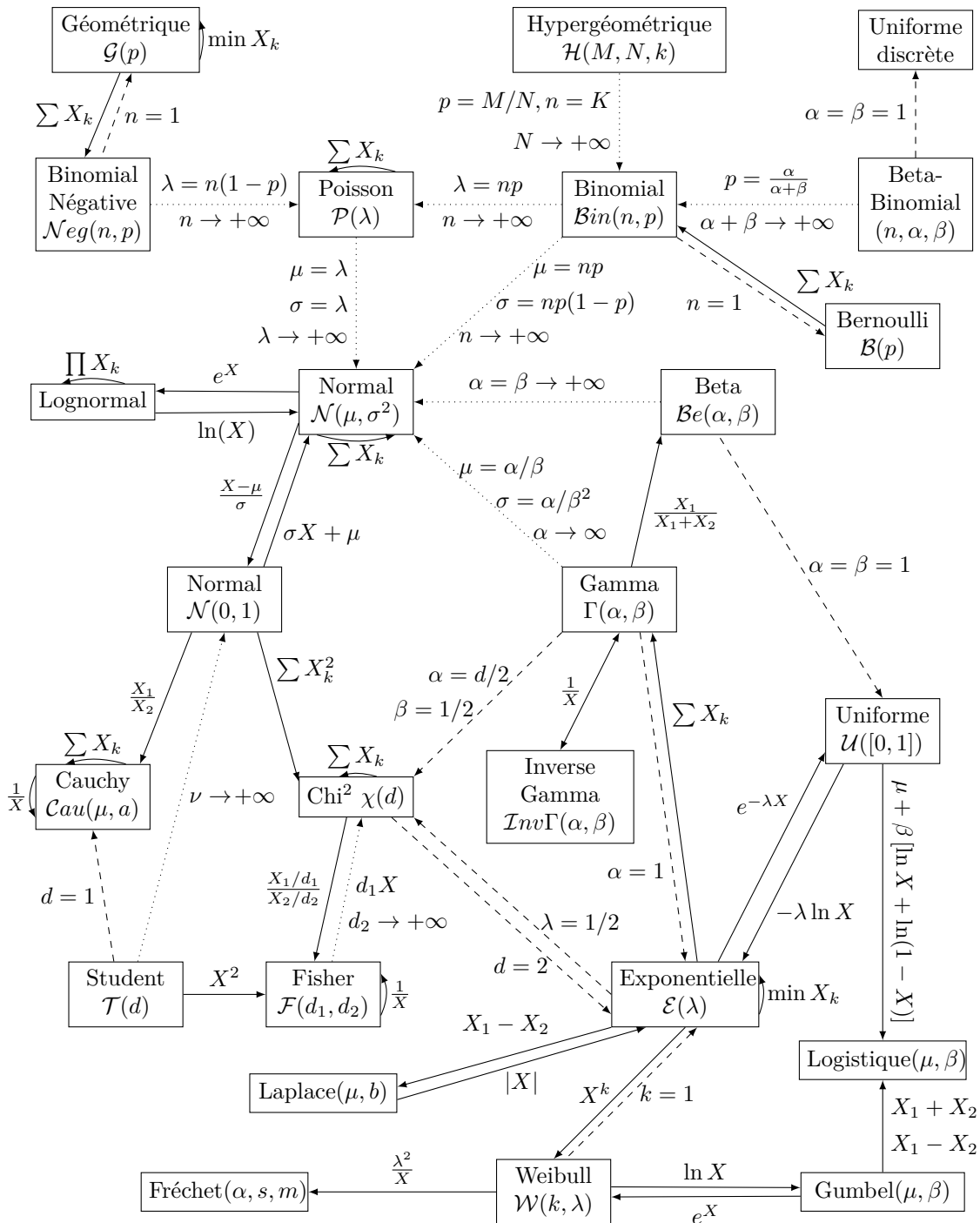


FIGURE 2.1 – Graphe de liens entre différentes lois : les flèches en traits pleins représentent une transformation de la loi du vecteur \mathbf{X} (début de la flèche) vers la variable $Y = f(\mathbf{X})$ (bout de la flèche), les flèches en pointillés longs indiquent un cas particulier et les flèches en pointillés courts un cas limite ou une approximation. Ce graphe est adapté du travail de Leemis (1986).

2.4 Solutions des exercices

Exercices 2.1

--- Pour la première, c'est une loi de Bernoulli (donc paramétrique), la deuxième peut être supposée paramétrique (puisque nous sommes dans un simplexe de \mathbb{R}^d) et la dernière n'est pas paramétrique si nous

ne connaissons pas le degré.



Exercices 2.2

Nom	Abréviation	Support	Densité $f(x)$ ou $\mathbb{P}(X = x)$	
Bernoulli	$\mathcal{B}(p)$	$\{0, 1\}$	$p^x(1-p)^{1-x}$	$p \in [0, 1]$
Bêta	$\mathcal{B}e(a, b)$	$]0, 1[$	$\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1}(1-x)^{b-1}$	$(a, b) \in \mathbb{R}_+^* \times \mathbb{R}_+^*$
Binomiale	$\mathcal{B}in(n; p)$	$\{0, 1, \dots, n\}$	$\binom{n}{x} p^x (1-p)^{n-x}$	$(n, p) \in \mathbb{N}^* \times [0; 1]$
Binomiale négative	$\mathcal{N}Bin(n; p)$	\mathbb{N}	$\binom{x+n-1}{x} p^n (1-p)^x$	$(n, p) \in \mathbb{N}^* \times [0; 1]$
Cauchy	$\mathcal{C}au(\mu, a)$	\mathbb{R}	$\frac{\pi}{1 + (\frac{x-\mu}{a})^2}$	$(\mu, a) \in \mathbb{R} \times \mathbb{R}_+^*$
Exponentielle	$\mathcal{E}(\lambda)$	\mathbb{R}^+	$\lambda e^{-\lambda x}$	$\lambda \in \mathbb{R}_+^*$
Gamma	$\mathcal{G}a(\lambda)$	\mathbb{R}_+^*	$\frac{\beta^\alpha}{\Gamma \alpha} x^{\alpha-1} e^{-\beta x}$	$(\alpha, \beta) \in \mathbb{R}_+^* \times \mathbb{R}_+^*$
Gaussienne ou normale	$\mathcal{N}(\mu, \sigma^2)$	\mathbb{R}	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$	$(\mu, \sigma) \in \mathbb{R} \times \mathbb{R}_+^*$
Géométrique	$\mathcal{G}éo(p)$	\mathbb{N}^*	$(1-p)^{x-1} p$	$p \in]0, 1]$
Hypergéométrique	$\mathcal{H}(n, p, A)$	$\{0, 1, \dots, n\}$	$\frac{\binom{pA}{x} \binom{(1-p)A}{n-x}}{\binom{A}{n}}$	$A \in \mathbb{N}$ $n \in \{0, \dots, A\}$ $p \in \{p \in [0, 1] \mid pA \text{ et } (1-p)A \text{ entiers}\}$
Inverse gamma	$\mathcal{I}G\alpha(\alpha, \beta)$	\mathbb{R}_+^*	$\frac{\beta^\alpha}{\Gamma \alpha} x^{-\alpha-1} e^{-\frac{\beta}{x}}$	$(\alpha, \beta) \in \mathbb{R}_+^* \times \mathbb{R}_+^*$
Laplace	Laplace(μ, b)	\mathbb{R}	$\frac{1}{2b} e^{-\frac{ x-\mu }{b}}$	$(\mu, b) \in \mathbb{R} \times \mathbb{R}_+^*$
Pareto	$\mathcal{P}a(a, k)$	$[a, +\infty[$	$\frac{ka^k}{x^{k+1}}$	$(a, k) \in \mathbb{R}_+^* \times \mathbb{R}_+^*$
Poisson	$\mathcal{P}(\lambda)$	\mathbb{N}	$e^{-\lambda} \frac{\lambda^x}{x!}$	\mathbb{R}^+
Uniforme	$\mathcal{U}(\{0, 1, \dots, n\})$	$\{0, 1, \dots, n\}$	$\frac{1}{n+1}$	$n \in \mathbb{N}^*$
Uniforme	$\mathcal{U}([a, b])$	$[a, b]$	$\frac{1}{b-a}$	$(a, b) \in \{(a, b) \in \mathbb{R}^2 \mid a < b\}$

Chapitre 3

Estimation

"Vous [les statisticien.ne.s] n'êtes jamais contents : avant, nous n'avions pas assez de données et les résultats obtenus n'étaient pas souvent concluants ; maintenant, nous en avons trop et nous ne pouvons pas avoir de résultats dans des temps corrects."

Discussion avec un biologiste.

3.1 Objectifs

Le but de ce chapitre est de comprendre le principe d'estimation et ce qui caractérise un *bon* estimateur. Nous allons donc introduire un certain nombre de notions et présenterons les estimateurs les plus connus.

3.2 Rappels/Pré-requis

Dans ce chapitre, nous aurons besoin des notions suivantes.

Probabilités : Différentes convergences, loi forte des grands nombres

3.3 Introduction

Le but de ce chapitre est d'estimer au mieux une quantité dépendante du paramètre inconnu θ . Pour bien différencier les outils que nous utilisons, nous allons commencer par un certain nombre de notations.

Définitions 24 (Paramètre d'intérêt)

Étant donnée une famille de probabilité $(\mathbb{P}_\theta)_{\theta \in \Theta}$ et un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$, nous notons θ^* le paramètre *inconnu* de la loi de X ; c'est-à-dire que $\mathbb{P}_X = \mathbb{P}_{\theta^*}$.

La **quantité à estimer** ou **paramètre d'intérêt** est notée $g(\theta^*)$ avec $g : \Theta \rightarrow \mathbb{R}^d$ où $d \in \mathbb{N}^*$.

Remarque

⊥ Souvent, la quantité à estimer est simplement θ^* c'est-à-dire pour g valant l'identité.

Maintenant, nous allons parler de l'estimation :

Définition 25 (Estimateur)

Étant donné un paramètre d'intérêt $g(\theta^*)$, nous appelons un **estimateur** de $g(\theta^*)$ toute application, notée \hat{g} , mesurable en l'observation X et indépendante de θ^* .

Remarque

⊥ Tout estimateur est donc de la forme $\hat{g} = h(X)$ avec h une application mesurable $E \rightarrow \mathbb{R}^d$.

Exemple

Par exemple, lorsque nous cherchons à estimer θ^* , nous proposons un estimateur noté $\hat{\theta}$. Cet estimateur est donc une fonction des observations.

**Attention au piège**

La notion d'indépendance de θ^* peut paraître étrange mais il arrive parfois que des auteurs proposent un estimateur directement lié à ce dernier. Par exemple, un estimateur de la forme $2\theta^*$ ne fonctionne pas.

Remarque

On peut trouver toute sorte d'estimateurs. Un exemple naïf serait l'estimateur qui renvoie toujours la valeur 0 par exemple. Il faut donc réfléchir à obtenir un *bon* estimateur et en quel sens ?

3.4 Consistance d'un estimateur

La première propriété que nous voulons pour un estimateur est que son estimation soit relativement proche de la vraie valeur.

Définition 26 (Consistance d'un estimateur)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ et un estimateur $\hat{g}_n = h_n(X_1, \dots, X_n)$, nous dirons que l'estimateur est

- **consistant** si pour tout θ^* de Θ , \hat{g}_n converge en \mathbb{P}_{θ^*} -probabilité vers $g(\theta^*)$ lorsque n tend vers $+\infty$. C'est-à-dire que :

$$\forall \varepsilon > 0, \mathbb{P}(\|g(\theta^*) - \hat{g}_n\|_{+\infty} > \varepsilon) \xrightarrow{n \rightarrow +\infty} 0$$

où $\|\cdot\|_{+\infty}$ est la **norme infinie** sur \mathbb{R}^d c'est-à-dire définie par :

$$\forall \mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\|_{+\infty} = \max_{1 \leq j \leq d} |x_j|.$$

- **fortement consistant** si pour tout θ^* de Θ l'estimateur \hat{g}_n converge \mathbb{P}_{θ^*} -presque sûrement vers $g(\theta^*)$ lorsque n tend vers $+\infty$. C'est-à-dire que :

$$\mathbb{P}\left(\lim_{n \rightarrow +\infty} \hat{g}_n = g(\theta^*)\right) = 1.$$

**Attention au piège**

En dimension finie, toutes les normes sont équivalentes donc nous pouvons choisir celle que nous préférons. En dimension infinie, il faut préciser par rapport à quelle norme nous avons la consistance.

Remarque

Par abus de langage, on parle souvent de consistance d'un estimateur là où on devrait plutôt parler de la consistance de la suite d'estimateurs.

Exemple

Dans l'exemple de la proportion de naissances masculines dans Paris, il est naturel d'estimer le paramètre θ^* par la moyenne empirique

$$\hat{\theta}_n(X_1, \dots, X_n) = \bar{\mathbf{X}}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

Alors, par la loi forte des grands nombres, nous savons que c'est un estimateur fortement consistant du paramètre.

D'autres propriétés peuvent être utilisées pour caractériser un bon estimateur, nous les verrons un peu plus tard.

3.5 Construction d'un estimateur

Dans cette section, nous présenterons deux types d'estimateurs : celui des moments et celui du maximum de vraisemblance.

3.5.1 Méthode des moments

Commençons par un exercice :

Exercices 3.1

Nous avons lancé un dé 1000 fois et nous avons obtenus 237 fois un chiffre inférieur ou égal à 3. Pensez-vous que le dé est équilibré ? Pourquoi ?

Modéliser l'expérience : comment d'écrire la situation ? Quelle loi est derrière ? Quel estimateur pourrions-nous proposer ? Quelle est sa loi exacte ? Sa loi asymptotique ? Que nous dit la loi forte des grands nombres ?

Le principe de la méthode des moments repose simplement sur la loi forte des grands nombres. Pour ce faire, nous supposons posséder un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi \mathbb{P}_{θ^*} et une fonction $\phi : E \rightarrow \mathbb{R}$ mesurable telle que $\phi(X_1)$ possède un moment d'ordre 1.

Pour estimer la valeur $g(\theta^*) = \mathbb{E}_{\theta^*}[\phi(X_1)]$, nous nous appuyons sur la moyenne empirique :

$$\hat{g}_n = \frac{1}{n} \sum_{i=1}^n \phi(X_i)$$

qui, d'après la loi forte des grands nombres, converge presque sûrement vers $g(\theta^*)$. Par conséquent, la suite d'estimateurs est fortement consistante.

Remarque

Cette méthode fonctionne dès que les paramètres d'intérêt peuvent s'exprimer comme des fonctions des moments de la loi des observations X_i .

Exemple

Nous estimons par exemple l'espérance $\mu_1 = \mathbb{E}[X]$ par la **moyenne empirique** $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$.

Si les variables X_i sont à valeurs réelles et possèdent un moment d'ordre 2 alors nous pouvons estimer $\mu_2 = \mathbb{E}[X^2]$ par $\hat{\mu}_{2,n} = \frac{1}{n} \sum_{i=1}^n X_i^2$.

Définition 27 (Méthode des moments (MM))

Étant donné un n -échantillon (X_1, \dots, X_n) de \mathbb{R}^d , q fonctions ϕ_1, \dots, ϕ_q allant respectivement de \mathbb{R}^d dans $\mathbb{R}^{n_1}, \dots, \mathbb{R}^{n_q}$ ayant chacune un moment d'ordre 1 ; c'est-à-dire que $\mathbb{E}_{\theta^*}[\|\phi(X_1)\|]$ est finie et une fonction Ψ allant de $\prod_{j=1}^q \mathbb{R}^{n_j}$ dans \mathbb{R}^m , la **méthode des moments** ou (MM) consiste à estimer $g(\theta^*) = \Psi(g_1(\theta^*), \dots, g_q(\theta^*))$ où $g_j(\theta^*) = \mathbb{E}_{\theta^*}[\phi_j(\theta^*)]$ par

$$\hat{g}_n = \Psi\left(\frac{1}{n} \sum_{i=1}^n \phi_1(X_i), \dots, \frac{1}{n} \sum_{i=1}^n \phi_q(X_i)\right).$$

Exemple

Étant donné un n -échantillon (X_1, \dots, X_n) de loi admettant un moment d'ordre 2, la variance $\sigma^2 = \mathbb{V}_{\theta^*} [X_1] = \mathbb{E}_{\theta^*} [X_1^2] - \mathbb{E}_{\theta^*} [X_1]^2$ est estimée par

$$\widehat{s}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2$$

appelée la **variance empirique**.

Exemple

Étant donné un n -échantillon (X_1, \dots, X_n) de loi de Poisson $\mathcal{P}(\theta^*)$ de paramètre θ^* alors nous avons deux estimateurs des moments possibles pour θ^* :

- Comme $\mathbb{E}_{\theta^*} [X_1] = \theta^*$, nous pouvons prendre la moyenne empirique.
- Comme $\mathbb{V}_{\theta^*} [X_1] = \theta^*$, nous pouvons prendre la variance empirique.

Exemple

Étant donné un n -échantillon (X_1, \dots, X_n) de loi $\mathcal{N}(\mu, \sigma^2)$ alors :

- Si σ^2 est connu, la méthode des moments suggère d'estimer μ par la moyenne empirique.
- Si μ est connue, la méthode des moments suggère d'estimer σ^2 par

$$\widehat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2.$$

- Si μ et σ^2 sont tous les deux inconnus alors la méthode des moments suggère d'estimer μ par la moyenne empirique et σ^2 par la variance empirique.

3.5.2 Méthode du maximum de vraisemblance**Attention au piège**

Dans cette partie, nous supposons que toutes les lois étudiées sont dominées presque sûrement par une mesure μ . Typiquement, nous aurons :

- la mesure de Lebesgue pour les lois sur \mathbb{R}^d .
- la mesure de comptage pour les cas de variables à valeurs dans un espace dénombrable.

Avant d'introduire la méthode, nous avons besoin d'une notation et d'une définition.

Notation

Pour tout paramètre $\theta \in \Theta$, nous notons f_θ la densité de \mathbb{P}_θ par rapport à la mesure dominante μ . En particulier, dans le cas discret, nous avons pour tout $x \in S$, $f_\theta(x) = \mathbb{P}_\theta(X = x) = \mathbb{P}(X = x; \theta)$.

Certains auteurs notent la densité par $p_\theta(\cdot)$ ou $p(\cdot; \theta)$.

Remarque

L'idée de la méthode du maximum de vraisemblance est de partir du constat que si nous avons une observation x et une famille de loi $(\mathbb{P}_\theta)_{\theta \in \Theta}$ alors nous pouvons calculer pour chaque θ la probabilité d'avoir eu l'observation. Par exemple, dans le cas discret, cela revient à calculer $\mathbb{P}_\theta(X = x) = \mathbb{P}(X = x; \theta)$. Une fois ces calculs faits, le θ le plus *vraisemblable* est simplement celui qui correspond à la plus forte probabilité d'avoir eu l'observation.

Définition 28 (Vraisemblance)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de densité f_{θ^*} , nous appelons **vraisemblance du n -échantillon \mathbf{X}** , notée $V_{\mathbf{X}}$, la fonction :

$$\begin{aligned} V_{\mathbf{X}} : \Theta &\rightarrow \mathbb{R} \\ \theta &\mapsto V_{\mathbf{X}}(\theta) = f_{\theta}(\mathbf{X}) \end{aligned}$$

Définition 29 (Méthode du maximum de vraisemblance (MMV))

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de densité f_{θ^*} , la **méthode du maximum de vraisemblance (MMV)** consiste à estimer $g(\theta^*) = \theta^*$ par l'élément de Θ maximisant la vraisemblance $V_{\mathbf{X}}$.

Si un tel point $\hat{\theta}_n$ existe, on l'appelle **l'estimateur du maximum de vraisemblance**.

Attention au piège

Nous pouvons faire trois remarques importantes :

- L'estimateur du maximum de vraisemblance n'existe pas toujours : c'est-à-dire qu'il n'y a pas toujours un point qui maximise la vraisemblance.
- S'il existe, il n'est pas obligatoirement unique : c'est-à-dire qu'il peut y avoir plusieurs points qui maximisent la vraisemblance.
- S'il existe et qu'il est unique, il n'est pas obligatoirement calculable : nous pouvons juste étudier la vraisemblance pour voir qu'elle doit admettre un maximum sans pouvoir trouver une formule explicite.

Proposition 2

Dans le cas d'un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi commune \mathbb{P}_{θ^*} , la vraisemblance prend la forme produit :

$$V_{\mathbf{X}}(\theta) = \prod_{i=1}^n f_{\theta}(X_i).$$

La conséquence est que nous préférons souvent maximiser la **log-vraisemblance** :

$$\log V_{\mathbf{X}}(\theta) = \sum_{i=1}^n \log f_{\theta}(X_i).$$

Preuve

Le produit vient de l'indépendance des variables $(X_i)_{i \in \{1, \dots, n\}}$ et la formule du fait que nous avons une simple densité.

Remarque

L'inconvénient du maximum de vraisemblance est que nous n'estimons que le paramètre θ^* mais pas des fonctions de celui-ci. En revanche, nous verrons par la suite que certains estimateurs du maximum de vraisemblance ont des meilleures propriétés que l'estimateur des moments.

Un calcul assez simple permet d'obtenir les estimateurs suivants :

**Exercices 3.2**

Dans le cas d'un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi de Bernoulli $\mathcal{B}(\theta^*)$ ou de Poisson $\mathcal{P}(\theta^*)$, l'estimateur du maximum de vraisemblance est le même que celui des moments.

De même pour les modèles gaussiens $\mathcal{N}(\mu, \sigma^2)$ que nous connaissons μ ou σ^2 voir aucun des deux.

Un exemple assez intéressant d'estimateur du maximum de vraisemblance est le suivant :

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de même loi $\mathcal{U}([0, \theta^*])$ avec $\theta^* \in \mathbb{R}_+^*$, nous avons :

$$\begin{aligned} V_{\mathbf{X}}(\theta) &= \prod_{i=1}^n \left(\frac{1}{\theta} \mathbb{1}_{[0, \theta]}(X_i) \right) \\ &= \frac{1}{\theta^n} \prod_{i=1}^n \mathbb{1}_{[0, \theta]}(X_i) \\ &= \theta^{-n} \mathbb{1}_{[X_{(n)}, +\infty]}(\theta) \text{ où } X_{(n)} = \max_{1 \leq i \leq n} X_i. \end{aligned}$$

Nous voyons que pour tout $n \in \mathbb{N}^*$, la fonction $\theta \mapsto \theta^{-n}$ est décroissante et positive mais l'indicatrice $\theta \mapsto \mathbb{1}_{[X_{(n)}, +\infty]}(\theta)$ est nulle pour tout θ strictement plus petit que $X_{(n)}$. La valeur maximale de la fonction est donc atteinte pour $\hat{\theta}_n = X_{(n)}$. Notons, que c'est bien n qui est fixé et θ qui tend vers l'infini (si on commence à regarder dans l'autre sens, on pourrait croire que $0 < \theta < 1$ pose problème).

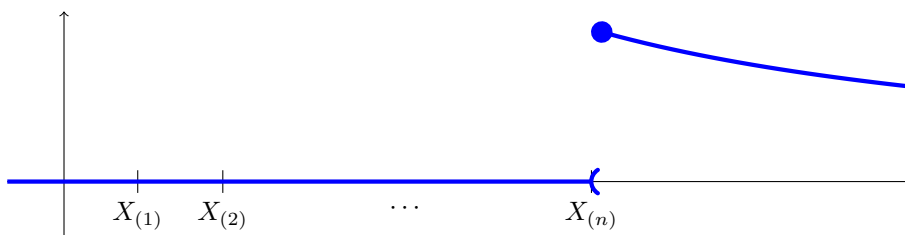


FIGURE 3.1 – Représentation en bleu de la fonction $\theta \mapsto \theta^{-n} \mathbb{1}_{[X_{(n)}, +\infty]}(\theta)$.

Enfin, voici un contre-exemple pour se rappeler qu'il n'est pas toujours possible de calculer le maximum de vraisemblance :

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de même loi $\mathcal{Cau}(\theta^*, 1)$ c'est-à-dire de densité

$$f_{\theta}(x) = \frac{1}{\pi} \times \frac{1}{1 + (x - \theta)^2}.$$

Le calcul de la log-vraisemblance est alors :

$$\log V_{\mathbf{X}}(\theta) = -n \log \pi - \sum_{i=1}^n \log [1 + (X_i - \theta)^2].$$

Cette fonction est définie sur \mathbb{R} , dérivable en θ , tend vers $-\infty$ lorsque θ tend vers $\pm\infty$ donc elle admet au moins un maximum qui est un point où la dérivée s'annule. En revanche, le calcul n'est pas explicite car la dérivée est une somme de fractions qui, une fois réduite au même dénominateur, est la somme de polynôme de degrés $2n - 1$. Il y a donc beaucoup de racines qui n'ont pas forcément de formules explicites (exceptées dans certains rares cas).

3.6 Qualité d'un estimateur

Dans cette partie, nous discutons des différentes propriétés que peut avoir un estimateur. Nous avons déjà vu l'indispensable consistance d'un estimateur qui permet de garantir que ce dernier sera *suffisamment proche* lorsque nous aurons suffisamment d'observations mais nous pouvons être intéressés par la vitesse à laquelle il va converger ou encore par la qualité non-asymptotique de son estimation.

3.6.1 Loi asymptotique

Comme nous allons parler de convergence en loi, nous commençons par un petit rappel et un théorème que nous utiliserons.

Définition 30 (Convergence en loi)

Étant donnée une suite $(X_n)_{n \in \mathbb{N}^*}$ de variables aléatoires d'un même espace de probabilité $(E, \mathcal{E}, \mathbb{P})$, nous disons que la suite $(X_n)_{n \in \mathbb{N}^*}$ **converge en loi** vers la variable X si, pour toute fonction φ continue bornée sur E à valeurs dans \mathbb{R} , nous avons :

$$\lim_{n \rightarrow +\infty} \mathbb{E}[\varphi(X_n)] = \mathbb{E}[\varphi(X)].$$

Nous notons alors :

$$X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X.$$

Remarque

Dans le cas de \mathbb{R}^d , nous pouvons montrer qu'il suffit de prendre des fonctions dans un espace dont l'adhérence englobe les fonctions continues à support compact sur \mathbb{R}^d (voir par exemple le polycopié de Le Gall (2006), proposition 10.3.3). En particulier, nous pourrions être intéressés par prendre des fonctions lipschitzienne (voir la preuve du lemme 5 de Slutsky).

Dans notre cas, nous allons souvent utiliser la proposition suivante :

Proposition 3 (Cas de variables réelles)

Étant données une suite $(X_n)_{n \in \mathbb{N}^*}$ de variables aléatoires d'un même espace de probabilité $(E, \mathcal{E}, \mathbb{P})$ et une variable X sur le même espace, nous avons l'équivalence des propositions suivantes :

- (1) $(X_n)_{n \in \mathbb{N}^*}$ converge en loi vers X ;
- (2) Pour tout $x \in \mathbb{R}$ où F_X est continue (c'est-à-dire que $\mathbb{P}(X = x) = 0$), nous avons la convergence simple :

$$\lim_{n \rightarrow +\infty} F_{X_n}(x) = F_X(x)$$

où $F_X(x) = \mathbb{P}(X \leq x)$ pour tout $x \in \mathbb{R}$.

Preuve

Une preuve complète peut être trouvée dans le polycopié de Le Gall (2006), proposition 10.3.2 et sa conséquence.

Dans un premier temps, nous allons nous concentrer sur l'estimation de $g(\theta^*) = \theta^*$. Le principe de normalité asymptotique va plus loin que la consistance puisque nous nous intéressons à la vitesse de convergence de $\hat{\theta}_n$ vers θ^* .

Définition 31 (Vitesse de convergence)

Étant donné un estimateur \hat{g}_n du paramètre $g(\theta^*)$ et une suite $(a_n)_{n \in \mathbb{N}^*}$ tendant vers $+\infty$, nous disons que $(a_n)_{n \in \mathbb{N}^*}$ est l'ordre de la **vitesse de convergence** de l'estimateur \hat{g}_n s'il



existe une variable aléatoire Y de loi non triviale indépendante de \hat{g}_n telle que

$$a_n (\hat{g}_n - g(\theta^*)) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} Y.$$

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi commune admettant un moment d'ordre 2, d'espérance μ^* et de matrice de variance covariance Σ^* alors, d'après le théorème de limite centrale, nous avons :

$$\sqrt{n} (\bar{\mathbf{X}}_n - \mu^*) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \Sigma^*).$$

Sa vitesse de convergence est donc \sqrt{n} .



Exercices 3.3

Dans le cas d'un n -échantillon suivant une loi de Poisson $\mathcal{P}(\theta^*)$, nous avons déjà qu'un estimateur possible est également $\hat{\theta}_n = \bar{\mathbf{X}}_n$ et nous aurons dans ce cas :

$$\sqrt{n} (\bar{\mathbf{X}}_n - \theta^*) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \theta^*).$$

Certains auteurs utilisent parfois la notion d'équivalence asymptotique :

$$\bar{\mathbf{X}}_n \underset{+\infty}{\overset{\mathcal{L}}{\sim}} \mathcal{N}\left(\theta^*, \frac{\theta^*}{n}\right).$$

Définition 32 (Asymptotiquement normal)

Un estimateur \hat{g}_n est dit **asymptotiquement normal** s'il converge vers une loi asymptotique normale à une vitesse de l'ordre de \sqrt{n} . Donc il faut vérifier les trois conditions suivantes :

- La vitesse de convergence est en \sqrt{n} .
- La convergence a lieu en loi.
- La loi limite est une loi normale non triviale.

Remarque

┃ Dans les cas classiques et notamment pour les estimateurs des moments, ce sera bien souvent le cas.

Contre-exemple

Si X_1, \dots, X_n est un n -échantillon suivant une loi $\mathcal{U}([0, \theta^*])$ alors, nous avons vu que l'estimateur du maximum de vraisemblance était $\hat{\theta}_n = \max\{X_1, \dots, X_n\} = X_{(n)}$ et calculons la loi de $n(\theta^* - \hat{\theta}_n)$. Comme, \mathbb{P}_{θ^*} -presque sûrement, θ^* est plus grand que toutes les variables X_i , nous avons que $n(\theta^* - \hat{\theta}_n)$ est \mathbb{P}_{θ^*} -presque sûrement positif. Donc, pour tout $x \in \mathbb{R}^+$, nous avons :

$$\begin{aligned} F_{n(\theta^* - \hat{\theta}_n)}(x) &= \mathbb{P}\left(n(\theta^* - \hat{\theta}_n) \leq x\right) \\ &= \mathbb{P}\left(\theta^* - \hat{\theta}_n \leq \frac{x}{n}\right) \\ &= \mathbb{P}\left(\theta^* - \frac{x}{n} \leq \hat{\theta}_n\right) \\ &= 1 - \mathbb{P}\left(\theta^* - \frac{x}{n} > \hat{\theta}_n\right) \\ &= 1 - \mathbb{P}\left(\max\{X_1, \dots, X_n\} < \theta^* - \frac{x}{n}\right) \end{aligned}$$

$$\begin{aligned}
&= 1 - \mathbb{P}\left(\left(X_1 < \theta^* - \frac{x}{n}\right) \text{ et } \left(X_2 < \theta^* - \frac{x}{n}\right) \text{ et } \dots \text{ et } \left(X_n < \theta^* - \frac{x}{n}\right)\right) \\
&= 1 - \prod_{i=1}^n \mathbb{P}\left(X_i < \theta^* - \frac{x}{n}\right) \\
&= 1 - \left[\prod_{i=1}^n \left(\frac{\theta^* - \frac{x}{n}}{\theta^*}\right)\right] \mathbb{1}_{\{\theta^* - \frac{x}{n} > 0\}} \\
&= 1 - \left(1 - \frac{x}{n\theta^*}\right)^n \mathbb{1}_{]-\infty; n\theta^*]}(x) \\
&\xrightarrow[n \rightarrow +\infty]{} 1 - e^{-\frac{1}{\theta^*}x}
\end{aligned}$$

et nous reconnaissons la fonction de répartition de la loi exponentielle de paramètre $1/\theta^*$. Par conséquent, la loi asymptotique de $n(\theta^* - \hat{\theta}_n)$ est $\mathcal{E}\left(\frac{1}{\theta^*}\right)$.

Remarque

Dans le contre-exemple précédent, nous voyons que la vitesse est en n par opposition à l'estimateur des moments qui converge à la vitesse \sqrt{n} .

Dans cette partie, la proposition la plus importante est celle de la méthode Delta :

Théorème 4 (Méthode Delta)

Étant données une suite $(U_n)_{n \in \mathbb{N}^*}$ de vecteurs aléatoires de \mathbb{R}^d avec $d \in \mathbb{N}^*$, une suite déterministe de réels $(a_n)_{n \in \mathbb{N}^*}$ et une application $\ell : \mathbb{R}^d \rightarrow \mathbb{R}^p$ avec $p \in \mathbb{N}^*$ telles que :

- a_n tend vers $+\infty$.
- Il existe $U^* \in \mathbb{R}^d$ un vecteur déterministe et $V \in \mathbb{R}^d$ un vecteur aléatoire tels que

$$a_n (U_n - U^*) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} V.$$

- ℓ est différentiable en U^* de différentielle notée $D\ell(U^*) \in \mathcal{M}_{p \times d}(\mathbb{R})$.

Alors, nous avons :

$$a_n (\ell(U_n) - \ell(U^*)) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} D\ell(U^*) V.$$

Pour la démonstration de ce théorème, nous avons besoin de rappeler le lemme de Slutsky.

Lemme 5 (Slutsky)

Étant données $(X_n)_{n \in \mathbb{N}^*}$ et $(Y_n)_{n \in \mathbb{N}^*}$ deux suites de variables aléatoires prenant leurs valeurs respectivement dans \mathbb{R}^d et \mathbb{R}^p avec $d, p \in \mathbb{N}^*$ éventuellement différents telles que X_n converge en loi vers un vecteur aléatoire X et Y_n converge en probabilité vers un vecteur déterministe C alors :

$$(X_n, Y_n) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} (X, C).$$

Remarque

Le plus souvent, nous appliquons à cette convergence jointe une fonction g ce qui permet de déduire :

$$g(X_n, Y_n) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} g(X, C).$$

Les fonctions les plus utilisées sont par exemple la somme (si $d = p$), le produit (si $p = 1$) ou la division (si $C \neq 0$ et $p = 1$).

**Preuve**

Pour démontrer le lemme de Slutsky, nous prenons une fonction $\varphi : \mathbb{R}^d \times \mathbb{R}^p \rightarrow \mathbb{R}$ lipschitzienne et bornée de constante de lipschitz L et M un majorant de la valeur absolue de la fonction. Commençons par utiliser l'inégalité triangulaire pour montrer que :

$$|\mathbb{E}[\varphi(X_n, Y_n)] - \mathbb{E}[\varphi(X, C)]| \leq \underbrace{|\mathbb{E}[\varphi(X_n, Y_n)] - \mathbb{E}[\varphi(X_n, C)]|}_{=(1)} + \underbrace{|\mathbb{E}[\varphi(X_n, C)] - \mathbb{E}[\varphi(X, C)]|}_{=(2)}.$$

Comme la fonction $x \mapsto \varphi(x, C)$ est bornée et que la suite $(X_n)_{n \in \mathbb{N}^*}$ tend en loi vers X alors, par le théorème du porte-manteau, le second terme (2) tend vers 0.

Pour le premier terme (1), nous prenons $\varepsilon > 0$ et nous utilisons le fait que la fonction soit lipschitzienne :

$$\begin{aligned} (1) &= |\mathbb{E}[\varphi(X_n, Y_n)] - \mathbb{E}[\varphi(X_n, C)]| \\ &\leq \mathbb{E}[|\varphi(X_n, Y_n) - \varphi(X_n, C)|] \\ &\leq \mathbb{E}[|\varphi(X_n, Y_n) - \varphi(X_n, C)| \mid \|Y_n - C\| > \varepsilon] \mathbb{P}(\|Y_n - C\| > \varepsilon) \\ &\quad + \mathbb{E}[|\varphi(X_n, Y_n) - \varphi(X_n, C)| \mid \|Y_n - C\| \leq \varepsilon] \mathbb{P}(\|Y_n - C\| \leq \varepsilon) \\ &\leq \mathbb{E}[|\varphi(X_n, Y_n)| + |\varphi(X_n, C)| \mid \|Y_n - C\| > \varepsilon] \mathbb{P}(\|Y_n - C\| > \varepsilon) \\ &\quad + \mathbb{E}[L \|Y_n - C\| \mid \|Y_n - C\| \leq \varepsilon] \mathbb{P}(\|Y_n - C\| \leq \varepsilon) \\ &\leq \mathbb{E}[2M \mid \|Y_n - C\| > \varepsilon] \mathbb{P}(\|Y_n - C\| > \varepsilon) \\ &\quad + \mathbb{E}[L\varepsilon \mid \|Y_n - C\| \leq \varepsilon] \mathbb{P}(\|Y_n - C\| \leq \varepsilon) \\ &\leq 2M\mathbb{P}(\|Y_n - C\| > \varepsilon) + L\varepsilon\mathbb{P}(\|Y_n - C\| \leq \varepsilon) \end{aligned}$$

Comme $(Y_n)_{n \in \mathbb{N}^*}$ converge en probabilité vers C , nous avons pour tout $\varepsilon > 0$:

$$\limsup_{n \rightarrow +\infty} |\mathbb{E}[\varphi(X_n, Y_n)] - \mathbb{E}[\varphi(X_n, C)]| \leq L\varepsilon$$

ce qui permet de conclure.

À l'aide de ce lemme, nous pouvons maintenant démontrer la méthode Delta :

Preuve

Avant de commencer, rappelons que la définition de la différentiabilité est la suivante :

$$\forall M > 0, \forall \varepsilon > 0, \exists \delta > 0,$$

$$\forall U' \in \mathbb{R}^d, \|U' - U^*\| \leq \delta \Rightarrow \|\ell(U') - \ell(U^*) - D\ell(U^*)(U' - U^*)\| \leq \frac{M}{\varepsilon} \|U' - U^*\|.$$

Commençons par décomposer la variable étudiée pour introduire la notion de différentiabilité :

$$a_n(\ell(U_n) - \ell(U^*)) = \underbrace{a_n(\ell(U_n) - \ell(U^*)) - D\ell(U^*)a_n(U_n - U^*)}_{=(1)} + \underbrace{D\ell(U^*)a_n(U_n - U^*)}_{=(2)}$$

et montrons que le premier terme converge vers 0 et en probabilité. Pour cela, nous fixons $\varepsilon > 0$ et $M > 0$ et nous choisissons $\delta > 0$ tel que

$$\begin{aligned} \mathbb{P}(\|(1)\| > \varepsilon) &= \mathbb{P}(\|a_n(\ell(U_n) - \ell(U^*)) - D\ell(U^*)a_n(U_n - U^*)\| > \varepsilon) \\ &= \mathbb{P}(a_n \|\ell(U_n) - \ell(U^*) - D\ell(U^*)(U_n - U^*)\| > \varepsilon) \\ &= \mathbb{P}(a_n \|\ell(U_n) - \ell(U^*) - D\ell(U^*)(U_n - U^*)\| > \varepsilon \mid \|U_n - U^*\| \leq \delta) \mathbb{P}(\|U_n - U^*\| \leq \delta) \\ &\quad + \mathbb{P}(a_n \|\ell(U_n) - \ell(U^*) - D\ell(U^*)(U_n - U^*)\| > \varepsilon \mid \|U_n - U^*\| > \delta) \mathbb{P}(\|U_n - U^*\| > \delta) \\ &\leq \mathbb{P}(a_n \|U_n - U^*\| > M \mid \|U_n - U^*\| \leq \delta) + \mathbb{P}(a_n \|U_n - U^*\| > a_n \delta) \\ &\leq \mathbb{P}(a_n \|U_n - U^*\| > M) + \mathbb{P}(a_n \|U_n - U^*\| > a_n \delta) \\ &\leq 2\mathbb{P}(a_n \|U_n - U^*\| > \min(M, a_n \delta)). \end{aligned}$$

Ceci étant vrai pour tout $n \in \mathbb{N}^*$ et comme a_n tend vers l'infini, nous avons pour n assez grand :

$$\begin{aligned} \limsup_{n \rightarrow +\infty} \mathbb{P}(\|(1)\| > \varepsilon) &\leq \limsup_{n \rightarrow +\infty} 2\mathbb{P}(a_n \|U_n - U^*\| > M) \\ &\leq 2\mathbb{P}(\|V\| \geq M) \end{aligned}$$

car $a_n(U_n - U^*)$ converge en loi vers V et par le lemme de Porte-Manteau. Or, ceci est vrai pour tout M donc nous pouvons faire tendre M vers $+\infty$ et nous avons donc la convergence en probabilité vers 0.

Enfin, par les hypothèses, nous savons que la deuxième partie (2) converge en loi vers $D\ell(U^*)V$ donc, en utilisant le lemme de Slutsky sur la somme, nous avons le résultat.

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ suivant une loi de Poisson $\mathcal{P}(\theta^*)$, nous avons vu qu'un estimateur possible de θ^* est la variance empirique $\hat{\theta}_{n,2} = \overline{X_n^2} - \overline{X_n}^2$. Par la méthode Delta, nous montrons que :

$$\sqrt{n}(\hat{\theta}_{n,2} - \theta^*) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 2\theta^{*2} + \theta^*).$$

Exercices 3.4

Montrer le résultat de l'exemple.

Corollaire 6 (Méthode Delta et estimateur des moments)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de vecteurs de \mathbb{R}^d , de loi commune admettant un moment d'ordre 2, d'espérance μ et de matrice de variance-covariance Σ et ℓ une fonction différentiable en μ de différentielle notée $D\ell(\mu)$ alors

$$\sqrt{n}(\ell(\overline{\mathbf{X}}_n) - \ell(\mu)) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, D\ell(\mu)\Sigma D\ell(\mu)^T).$$

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ suivant une loi exponentielle $\mathcal{E}(\theta^*)$, nous cherchons estimer la variance $1/\theta^{*2}$ par l'estimateur des moments

$$\hat{s}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2 = \overline{X^2} - \overline{X}^2 = \ell(\overline{X^2}, \overline{X}) \text{ avec } \ell(x, y) = x - y^2.$$

Or, nous savons que, pour tout $m \in \mathbb{N}^*$, nous avons $\mathbb{E}_{\theta^*}[X^m] = \frac{m!}{\theta^{*m}}$ et ainsi :

$$\begin{aligned} \mathbb{V}_{\theta^*}[X_1] &= \mathbb{E}_{\theta^*}[X_1^2] - \mathbb{E}_{\theta^*}[X_1]^2 \\ &= \frac{2!}{\theta^{*2}} - \left(\frac{1!}{\theta^{*1}} \right)^2 \\ &= \frac{1}{\theta^{*2}}, \\ \text{Cov}_{\theta^*}(X_1, X_1^2) &= \mathbb{E}_{\theta^*}[X_1^3] - \mathbb{E}_{\theta^*}[X_1]\mathbb{E}_{\theta^*}[X_1^2] \\ &= \frac{3!}{\theta^{*3}} - \frac{1!}{\theta^{*1}} \times \frac{2!}{\theta^{*2}} \\ &= \frac{4}{\theta^{*3}}, \\ \mathbb{V}_{\theta^*}[X_1^2] &= \mathbb{E}_{\theta^*}[X_1^4] - \mathbb{E}_{\theta^*}[X_1^2]^2 \end{aligned}$$

$$\begin{aligned}
&= \frac{4!}{\theta^{*4}} - \left(\frac{2!}{\theta^{*2}}\right)^2 \\
&= \frac{20}{\theta^{*4}}.
\end{aligned}$$

Par conséquent, nous avons par le théorème de limite centrale :

$$\sqrt{n} \left(\left(\frac{\overline{X^2}}{\overline{X}} \right) - \left(\frac{\frac{2}{\theta^{*2}}}{\frac{1}{\theta^*}} \right) \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \frac{20}{\theta^{*4}} & \frac{4}{\theta^{*3}} \\ \frac{4}{\theta^{*3}} & \frac{1}{\theta^{*2}} \end{pmatrix} \right).$$

Pour appliquer la méthode Delta, nous devons calculer la différentielle :

$$D\ell(x, y) = (1, -2y).$$

De plus, nous avons :

$$\begin{aligned}
D\ell \left(\frac{2}{\theta^{*2}}, \frac{1}{\theta^*} \right) \Sigma D\ell \left(\frac{2}{\theta^{*2}}, \frac{1}{\theta^*} \right)^T &= \begin{pmatrix} 1 & -\frac{2}{\theta^*} \end{pmatrix} \begin{pmatrix} \frac{20}{\theta^{*4}} & \frac{4}{\theta^{*3}} \\ \frac{4}{\theta^{*3}} & \frac{1}{\theta^{*2}} \end{pmatrix} \begin{pmatrix} 1 \\ -\frac{2}{\theta^*} \end{pmatrix} \\
&= \begin{pmatrix} 1 & -\frac{2}{\theta^*} \end{pmatrix} \begin{pmatrix} \frac{12}{\theta^{*4}} \\ \frac{2}{\theta^{*3}} \end{pmatrix} \\
&= \frac{8}{\theta^{*4}}.
\end{aligned}$$

Par la méthode Delta, nous avons donc :

$$\sqrt{n} \left(\hat{s}_n^2 - \frac{1}{\theta^{*2}} \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(0, \frac{8}{\theta^{*4}} \right)$$

c'est-à-dire que l'estimateur est asymptotiquement normal.

3.6.2 Estimation sans biais

Le principe d'estimation est de vérifier si on peut *espérer* avoir une estimation parfaite en moyenne.

Définitions 33 (Biais)

Étant donné une observation X et un estimateur $\hat{g} = h(X)$ de $g(\theta^*)$ ayant un moment d'ordre 1; c'est-à-dire tel que $\mathbb{E}_{\theta^*} [\|\hat{g}\|] < +\infty$, pour tout $\theta^* \in \Theta$, nous appelons **biais** la fonction $b : \Theta \rightarrow \mathbb{R}^d$ définie pour tout $\theta^* \in \Theta$ par :

$$b(\theta^*) = \mathbb{E}_{\theta^*} [\hat{g}] - g(\theta^*).$$

L'estimateur \hat{g} sera dit **sans biais** si $b(\theta^*) = 0$ pour tout $\theta^* \in \Theta$; ou encore si l'espérance $\mathbb{E}_{\theta^*} [\hat{g}] = g(\theta^*)$.

Exemple

Lorsqu'un moment d'ordre 1 existe, la moyenne empirique est toujours un estimateur sans biais; en revanche, celui de la variance est biaisé.

Astuce algorithmique

Généralement, nous trouvons l'estimateur non biaisé de la variance dans les logiciels par défaut. De plus en plus, nous avons aussi l'estimateur classique.

Exemple

Dans le cas d'un n -échantillon de loi $\mathcal{U}([0, \theta^*])$, l'estimateur des moments est sans biais mais celui du maximum de vraisemblance est biaisé (il sous-estime la vraie valeur). Néanmoins, nous préférons le second car il converge plus vite vers la bonne valeur.

Contre-exemple

Dans le cas d'un modèle $\mathcal{N}(\theta^*, 1)$ où $\theta^* \in \mathbb{R}$, l'estimateur $\hat{g} = 0$ est nul pour $\theta^* = 0$ mais différent de zéro dès que $\theta^* \neq 0$, il est donc biaisé puisqu'il ne vérifie pas la condition pour tout $\theta^* \in \mathbb{R}$.

Remarque

Nous pouvons faire trois remarques :

1. Le caractère non-biaisé d'un estimateur doit être vérifié pour tout $\theta^* \in \Theta$; c'est-à-dire que, quelque soit la quantité à estimer, le biais doit être nul (pas juste pour des valeurs).
2. Le critère est non-asymptotique.
3. Le critère peut être discutable : à un moment, le biais était jugé très important (c'est pour cela que beaucoup de logiciels proposent des versions non biaisées par défaut de certains estimateurs). Néanmoins, il vaut toujours mieux un estimateur qui converge plus vite vers la bonne valeur qu'un estimateur non biaisé. De même, un estimateur non biaisé qui possède une forte variance a parfois une plus forte probabilité de fournir une estimation très loin de la quantité à estimer qu'un paramètre biaisé avec une faible variance.

Par rapport à la troisième remarque, nous préférons nous contenter de la version asymptotique.

Définition 34 (Asymptotiquement sans biais)

Étant donnée une observation X et un estimateur $\hat{g} = h(X)$ de $g(\theta^*)$ ayant un moment d'ordre 1 ; c'est-à-dire tel que $\mathbb{E}_{\theta^*} [\|\hat{g}\|] < +\infty$, pour tout $\theta^* \in \Theta$, l'estimateur \hat{g}_n sera dit **asymptotiquement sans biais** si pour tout $\theta^* \in \Theta$:

$$\lim_{n \rightarrow +\infty} \mathbb{E}_{\theta^*} [\hat{g}_n] = g(\theta^*).$$

Point méthode

Pour le biais, il suffit donc de procéder en deux étapes :

1. Nous calculons $\mathbb{E}_{\theta^*} [\hat{g}_n]$:
 - S'il vaut $g(\theta^*)$ pour tout $n \in \mathbb{N}^*$, l'estimateur est sans biais donc asymptotiquement sans biais aussi.
2. Sinon, nous regardons la limite :
 - Si elle vaut $g(\theta^*)$, l'estimateur est asymptotiquement sans biais aussi.
 - Sinon, il peut y avoir des valeurs de $\theta^* \in \Theta$ pour lesquelles, l'estimateur n'est pas consistant.

3.6.3 Estimation optimale

Dans cette partie, nous allons parler du risque d'un estimateur qui représente à quel point un estimateur peut s'écarter ou non de la valeur qu'il estime et nous essayerons de minimiser ce risque.

Risque d'un estimateur

Dans la partie précédente, nous nous sommes intéressé(e)s au fait que les estimations soient, en espérance, égale à la vraie valeur (statistique de position). Toutefois, nous pouvons être en espérance proche mais avec une forte dispersion (statistique de dispersion). Dans cette partie, nous étudions ainsi la variabilité des estimations et nous commençons par introduire le risque.

Définition 35 (Risque d'un estimateur)

Étant donné un estimateur \hat{g} de $g(\theta^*)$, nous associons son **risque (quadratique)** défini comme l'application suivante :

$$\begin{aligned} R : \Theta &\rightarrow [0, +\infty] \\ \theta^* &\mapsto \mathbb{E}_\theta \left[(\hat{g} - g(\theta))^2 \right] \end{aligned}$$

Remarque

Si plusieurs estimateurs sont étudiés en même temps, il est préférable d'indexer la fonction R par le nom de l'estimateur.

Un estimateur sera d'autant meilleur que son risque est faible sur Θ .

Proposition 7

Étant donné un estimateur \hat{g} de $g(\theta^*)$, nous avons la relation suivante :

$$R(\theta) = (\mathbb{E}_\theta [\hat{g}] - g(\theta))^2 + \mathbb{E}_\theta \left[(\hat{g} - \mathbb{E}_\theta [\hat{g}])^2 \right] = b(\theta)^2 + \mathbb{V}_\theta [\hat{g}].$$

Preuve

La démonstration consiste simplement à introduire l'espérance :

$$\begin{aligned} R(\theta) &= \mathbb{E}_\theta \left[(\hat{g} - g(\theta))^2 \right] \\ &= \mathbb{E}_\theta \left[(\hat{g} - \mathbb{E}_\theta [\hat{g}] + \mathbb{E}_\theta [\hat{g}] - g(\theta^*))^2 \right] \\ &= \mathbb{E}_\theta \left[(\hat{g} - \mathbb{E}_\theta [\hat{g}])^2 \right] + 2\mathbb{E}_\theta \left[(\hat{g} - \mathbb{E}_\theta [\hat{g}]) (\mathbb{E}_\theta [\hat{g}] - g(\theta^*)) \right] + \mathbb{E}_\theta \left[(\mathbb{E}_\theta [\hat{g}] - g(\theta^*))^2 \right] \\ &= \mathbb{V}_\theta [\hat{g}] + 2\mathbb{E}_\theta [\hat{g} - \mathbb{E}_\theta [\hat{g}]] (\mathbb{E}_\theta [\hat{g}] - g(\theta^*)) + b(\theta)^2 \\ &= \mathbb{V}_\theta [\hat{g}] + 2 \underbrace{(\mathbb{E}_\theta [\hat{g}] - \mathbb{E}_\theta [\hat{g}])}_{=0} (\mathbb{E}_\theta [\hat{g}] - g(\theta^*)) + b(\theta)^2 \\ &= b(\theta)^2 + \mathbb{V}_\theta [\hat{g}] \end{aligned}$$

Corollaire 8

Étant donné un estimateur \hat{g} sans biais de $g(\theta^*)$, le risque est égal à la variance de l'estimateur.

Dans le cadre des estimateurs non biaisés, nous nous intéressons donc à ceux qui minimisent uniformément la variance.

Définition 36 (Estimateur UMVU)

Étant donnée une collection \mathcal{G} d'estimateurs non biaisés de $g(\theta^*)$, c'est-à-dire un ensemble d'applications mesurable en l'observation X à valeurs dans l'espace d'arrivée de $g(\theta^*)$, nous appelons **estimateur UMVU** pour (**Uniform Minimum Variance Unbiased**) tout estimateur \hat{g} qui minimise la variance de façon uniforme c'est-à-dire vérifiant :

$$\forall \theta \in \Theta, \forall \hat{h} \in \mathcal{G}, \mathbb{V}_{\theta^*} [\hat{g}] \leq \mathbb{V}_{\theta^*} [\hat{h}].$$

Information de Fisher

Dans cette partie et les suivantes, nous avons besoin de restreindre l'ensemble des lois étudiées.

Définition 37 (Modèle paramétrique régulier)

Étant donné un modèle paramétrique $(f_\theta)_{\theta \in \Theta}$, nous dirons qu'il est **régulier** s'il vérifie les trois conditions suivantes :

1. Θ est un ouvert de \mathbb{R}^p .
2. Le support des densités f_θ , c'est-à-dire l'ensemble $\{x \in \mathbb{R}^d \mid f_\theta(x) > 0\}$, est indépendant de θ et la fonction $\theta \mapsto f_\theta$ est dérivable sur Θ sauf éventuellement pour un nombre fini de points.
3. Pour toute application g mesurable en X dont l'espérance du module est finie sous la loi f_θ pour tout θ , c'est-à-dire que $\mathbb{E}_\theta [\|g(X)\|] < +\infty$, nous supposons que l'application $\theta \mapsto \int_{\mathcal{X}} g(x) f_\theta(x) dx$ est dérivable et que sa dérivée vaut :

$$\frac{\partial}{\partial \theta} \int_{\mathcal{X}} g(x) f_\theta(x) dx = \int_{\mathcal{X}} g(x) \frac{\partial}{\partial \theta} f_\theta(x) dx$$

où $\frac{\partial}{\partial \theta}$ représente le gradient (donc si $p = 1$, c'est juste la dérivée).

Remarque

Pour les deux premiers points, nous pouvons faire des remarques :

1. Un exemple classique d'ouverts de \mathbb{R}^p est le pavé $]a_1, b_1[\times]a_2, b_2[\times \dots \times]a_p, b_p[$ avec $(a_1, b_1, a_2, \dots, b_p) \in \mathbb{R}^{2p}$.
2. La loi uniforme $\mathcal{U}([0, \theta])$ ne vérifie pas ces conditions. Ceci ne veut pas dire que les résultats que nous démontrerons par la suite ne s'appliquent pas forcément ; par contre, les démonstrations seront plus compliquées.

Définition 38 (Score)

Étant donné un modèle paramétrique régulier $(f_\theta)_{\theta \in \Theta}$ et un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$, nous appelons **score**, noté $U(\cdot)$, la fonction $U : \theta \mapsto \mathbb{R}$ définie par

$$U(\theta) = \frac{\partial}{\partial \theta} \log V_{\mathbf{X}}(\theta) = \frac{\partial}{\partial \theta} \left(\sum_{i=1}^n \log f_\theta(X_i) \right).$$

Propriétés 9

Dans un modèle paramétrique régulier, nous avons pour tout $\theta \in \Theta$:

$$\mathbb{E}_\theta [U(\theta)] = 0.$$

Preuve

Pour montrer cela, il suffit de faire le calcul. Dans le cas continu, nous avons :

$$\begin{aligned} \mathbb{E}_\theta [U(\theta)] &= \mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \left(\sum_{i=1}^n \log f_\theta(X_i) \right) \right] \\ &= \sum_{i=1}^n \mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \log f_\theta(X_i) \right] \end{aligned}$$



$$\begin{aligned}
&= n \int_{\mathcal{X}} \frac{\partial}{\partial \theta} \log f_{\theta}(x) f_{\theta}(x) dx \\
&= \sum_{i=1}^n \int_{\mathcal{X}} \frac{\frac{\partial}{\partial \theta} f_{\theta}(x)}{f_{\theta}(x)} f_{\theta}(x) dx \\
&= \sum_{i=1}^n \int_{\mathcal{X}} \frac{\partial}{\partial \theta} f_{\theta}(x) dx \\
&= \sum_{i=1}^n \frac{\partial}{\partial \theta} \underbrace{\int_{\mathcal{X}} f_{\theta}(x) dx}_{=1} \\
&= 0.
\end{aligned}$$

Remarque

L'estimateur du maximum de vraisemblance peut être vu comme la valeur annulant le score.

Définition 39 (Information de Fisher)

Étant donné un modèle paramétrique régulier $(f_{\theta})_{\theta \in \Theta}$ et un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$, nous appelons **information de Fisher**, notée $I_n(\cdot)$, la matrice symétrique de taille $p \times p$ définie pour tout $\theta \in \Theta$ par :

$$I_n(\theta) = \mathbb{E}_{\theta^*} [U(\theta)U(\theta)^T] = \left(\mathbb{E} \left[\frac{\partial \log V_{\mathbf{X}}(\theta)}{\partial \theta_i} \frac{\partial \log V_{\mathbf{X}}(\theta)}{\partial \theta_j} \right] \right)_{1 \leq i, j \leq p}.$$

Bornes de Cramer-Rao

Dans cette partie, nous faisons des hypothèses supplémentaires sur l'ensemble de définitions et sur la régularité des fonctions :

Hypothèse

Nous supposons que nous sommes dans un modèle paramétrique régulier $(f_{\theta})_{\theta \in \Theta}$ et que $\Theta \subset \mathbb{R}$.

Corollaire 10 (Cas unidimensionnel)

Dans le cas où $\theta \in \mathbb{R}$, nous avons :

$$I_n(\theta) = nI_1(\theta) = n \int \frac{\left[\frac{\partial}{\partial \theta} f_{\theta}(x) \right]^2}{f_{\theta}(x)} \mathbb{1}_{\{f_{\theta}(x) > 0\}} dx.$$

Preuve

Le fait que $I_n(\theta) = nI_1(\theta)$ vient de l'indépendance des observations et la deuxième partie est simplement le calcul.

Remarque

Ce corollaire permet de ne se concentrer que sur les propriétés de $I_1(\theta)$ que nous trouvons souvent simplement sous la notation $I(\theta)$.

Corollaire 11 (Information de Fisher et variance)

Étant donné un modèle paramétrique régulier $(f_\theta)_{\theta \in \Theta}$ unidimensionnel, nous avons :

$$I_1(\theta) = \mathbb{V}_\theta \left[\frac{\partial}{\partial \theta} \log f_\theta(X) \right].$$

Preuve

Pour cela, nous faisons le calcul de la variance :

$$\begin{aligned} \mathbb{V}_\theta \left[\frac{\partial}{\partial \theta} \log f_\theta(X) \right] &= \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \log f_\theta(X) \right)^2 \right] - \mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \log f_\theta(X) \right]^2 \\ &= I_1(\theta) - \underbrace{\mathbb{E}_\theta [U(\theta)]^2}_{=0}. \end{aligned}$$

Proposition 12 (Information de Fisher et dérivée seconde)

Étant donné un modèle paramétrique régulier $(f_\theta)_{\theta \in \Theta}$ unidimensionnel, si nous pouvons dériver deux fois sous le signe somme alors nous avons :

$$I_1(\theta) = -\mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(X) \right].$$

Preuve

Il suffit de faire le calcul :

$$\begin{aligned} \mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(X) \right] &= \int_{\{x \in \mathcal{X} | f_\theta(x) > 0\}} \left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(x) \right] f_\theta(x) dx \\ &= \int_{\{x \in \mathcal{X} | f_\theta(x) > 0\}} \left[\frac{\partial}{\partial \theta} \frac{\frac{\partial}{\partial \theta} f_\theta(x)}{f_\theta(x)} \right] f_\theta(x) dx \\ &= \int_{\{x \in \mathcal{X} | f_\theta(x) > 0\}} \left[\frac{\frac{\partial^2}{\partial \theta^2} f_\theta(x)}{f_\theta(x)} - \frac{\left[\frac{\partial}{\partial \theta} f_\theta(x) \right]^2}{f_\theta^2(x)} \right] f_\theta(x) dx \\ &= \int_{\{x \in \mathcal{X} | f_\theta(x) > 0\}} \frac{\partial^2}{\partial \theta^2} f_\theta(x) dx - \int_{\{x \in \mathcal{X} | f_\theta(x) > 0\}} \frac{\left[\frac{\partial}{\partial \theta} f_\theta(x) \right]^2}{f_\theta(x)} dx \\ &= \frac{\partial^2}{\partial \theta^2} \int_{\{x \in \mathcal{X} | f_\theta(x) > 0\}} f_\theta(x) dx - I_1(\theta) \\ &= -I_1(\theta) \end{aligned}$$

Exercices 3.5

Montrer de deux façons différentes que pour un n -échantillon de loi de Poisson, nous avons $I_n(\theta) = n/\theta$.

Théorème 13 (Borne de Cramer-Rao)

Étant donné un estimateur $\widehat{g} = h(X)$ sans biais de $g(\theta)$ avec g une fonction dérivable sur Θ , si le modèle est régulier d'information de Fisher strictement positive sur Θ et si $\mathbb{E}_\theta [h^2(X)]$ est

bornée alors nous avons :

$$\forall \theta \in \Theta, \forall_{\theta} [\hat{g}] \geq \frac{(g'(\theta))^2}{I_1(\theta)}.$$



Preuve

Comme \hat{g} est supposé sans biais, nous savons que $\mathbb{E}_{\theta} [\hat{g}] = g(\theta)$ donc nous avons :

$$\begin{aligned} g'(\theta) &= \frac{\partial}{\partial \theta} g(\theta) \\ &= \frac{\partial}{\partial \theta} \mathbb{E}_{\theta} [\hat{g}] \\ &= \frac{\partial}{\partial \theta} \mathbb{E}_{\theta} [h(X)] \\ &= \frac{\partial}{\partial \theta} \int_{\{x \in \mathcal{X} | f_{\theta}(x) > 0\}} h(x) f_{\theta}(x) dx \\ &= \int_{\{x \in \mathcal{X} | f_{\theta}(x) > 0\}} h(x) \frac{\partial}{\partial \theta} f_{\theta}(x) dx \\ &= \int_{\{x \in \mathcal{X} | f_{\theta}(x) > 0\}} h(x) \frac{\partial}{\partial \theta} f_{\theta}(x) dx - \underbrace{\mathbb{E}_{\theta} [h(X)]}_{\text{bornée}} \underbrace{\int_{\{x \in \mathcal{X} | f_{\theta}(x) > 0\}} \frac{\partial}{\partial \theta} f_{\theta}(x) dx}_{=0} \\ &= \int_{\{x \in \mathcal{X} | f_{\theta}(x) > 0\}} [h(x) - \mathbb{E}_{\theta} [h(X)]] \frac{\partial}{\partial \theta} f_{\theta}(x) dx \\ &= \int_{\{x \in \mathcal{X} | f_{\theta}(x) > 0\}} [h(x) - \mathbb{E}_{\theta} [h(X)]] \sqrt{f_{\theta}(x)} \frac{\partial f_{\theta}(x)}{f_{\theta}(x)} \sqrt{f_{\theta}(x)} dx \\ &\leq \sqrt{\int_{\{x \in \mathcal{X} | f_{\theta}(x) > 0\}} [h(x) - \mathbb{E}_{\theta} [h(X)]]^2 f_{\theta}(x) dx} \sqrt{\int_{\{x \in \mathcal{X} | f_{\theta}(x) > 0\}} \frac{[\frac{\partial}{\partial \theta} f_{\theta}(x)]^2}{f_{\theta}^2(x)} f_{\theta}(x) dx} \\ &\quad \text{par Cauchy-Schwartz} \\ &\leq \sqrt{\mathbb{V}_{\theta} [h(X)]} \sqrt{\int_{\{x \in \mathcal{X} | f_{\theta}(x) > 0\}} \frac{[\frac{\partial}{\partial \theta} f_{\theta}(x)]^2}{f_{\theta}(x)} dx} \\ &\leq \sqrt{\mathbb{V}_{\theta} [\hat{g}]} \sqrt{I_1(\theta)}. \end{aligned}$$

En passant au carré et en divisant par $I_1(\theta)$ (qui est strictement positif), nous avons le résultat.

Définition 40 (Borne de Cramer-Rao)

Dans le cas où la quantité d'intérêt est $g(\theta)$ avec g dérivable, la **borne de Cramer-Rao** est la quantité :

$$\frac{(g'(\theta))^2}{I_1(\theta)}.$$

Corollaire 14 (Borne de Cramer-Rao et estimateur UMVU)

Si un estimateur sans biais \hat{g} vérifie pour tout $\theta \in \Theta$ que sa variance est égale à la borne de Cramer-Rao alors c'est un estimateur UMVU.

Corollaire 15 (Borne de Cramer-Rao pour un n -échantillon)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ d'information de Fisher I_n et sous les conditions du théorème précédent, nous avons :

$$\forall \theta \in \Theta, \mathbb{V}_\theta [h(X_1, \dots, X_n)] \geq \frac{(g'(\theta))^2}{nI_1(\theta)}.$$

Exemple

Étant donné $\mathbf{X} = (X_1, \dots, X_n)$ un n -échantillon de loi de Poisson $\mathcal{P}(\theta^*)$ alors nous montrons que $\hat{\theta}_{n,1} = \bar{\mathbf{X}}_n$ est un UMVU.

Point méthode

Pour un estimateur *UMVU*, il faut procéder en plusieurs étapes :

1. Vérifier s'il est non biaisé (sinon, cela s'arrête là).
2. Calculer la dérivée première ou la dérivée seconde en θ de la fonction densité f_θ (suivant si nous pouvons dériver qu'une seule ou deux fois sous le signe somme).
3. Calculer l'information de Fisher pour une observation.
4. Calculer la variance de l'estimateur.
5. Vérifier si borne de Cramer-Rao est atteinte ou non.

3.7 Solutions des exercices**Exercices 3.1**

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ suivant une loi de Poisson $\mathcal{P}(\theta^*)$, nous avons vu qu'un estimateur possible de θ^* est la variance empirique $\hat{\theta}_{n,2} = \bar{X}_n^2 - \bar{\mathbf{X}}_n^2 = \ell(\bar{X}_n^2, \bar{\mathbf{X}}_n)$ où $\ell(x, y) = x - y^2$ et nous souhaitons appliquer la méthode Delta.

Pour cela, nous prenons le vecteur $U_n = \begin{pmatrix} X_n^2 \\ \bar{\mathbf{X}}_n \end{pmatrix}$ dont nous devons calculer la loi asymptotique. Or, si nous prenons $Y_i = \begin{pmatrix} X_i^2 \\ X_i \end{pmatrix}$, nous voyons que le théorème de la limite centrale pourra nous aider à trouver la loi puisque nous avons des observations (Y_1, \dots, Y_n) indépendantes et de même loi. Commençons par calculer le moment d'ordre 1 de Y_i , pour cela, nous rappelons les moments suivants :

$$\begin{aligned} \mathbb{E}[X_i] &= \theta^*, \quad \mathbb{E}[X_i^2] = \theta^{*2} + \theta^*, \\ \mathbb{E}[X_i^3] &= \theta^{*3} + 3\theta^{*2} + \theta^* \quad \text{et} \quad \mathbb{E}[X_i^4] = \theta^{*4} + 6\theta^{*3} + 7\theta^{*2} + \theta^*. \end{aligned}$$

Ceci nous permet de calculer les moments d'ordre 1 :

$$\begin{aligned} \mathbb{E}[Y_i] &= \begin{pmatrix} \mathbb{E}[X_1^2] \\ \mathbb{E}[X_1] \end{pmatrix} \\ &= \begin{pmatrix} \theta^{*2} + \theta^* \\ \theta^* \end{pmatrix}. \end{aligned}$$

Pour la matrice de variance-covariance, nous avons :

$$\begin{aligned} \mathbb{V}[X_i^2] &= \mathbb{E}[X_i^4] - \mathbb{E}[X_i^2]^2 \\ &= \theta^{*4} + 6\theta^{*3} + 7\theta^{*2} + \theta^* - (\theta^{*2} + \theta^*)^2 \end{aligned}$$



$$\begin{aligned}
&= \theta^{*4} + 6\theta^{*3} + 7\theta^{*2} + \theta^* - \theta^{*4} - 2\theta^{*3} - \theta^{*2} \\
&= 4\theta^{*3} + 6\theta^{*2} + \theta^*, \\
\mathbb{V}[X_i] &= \theta^*, \\
\text{Cov}(X_i^2, X_i) &= \mathbb{E}[X_i^3] - \mathbb{E}[X_i^2] \mathbb{E}[X_i] \\
&= \theta^{*3} + 3\theta^{*2} + \theta^* - \theta^{*3} - \theta^{*2} \\
&= 2\theta^{*2} + \theta^*.
\end{aligned}$$

Au final, par le théorème de la limite centrale, nous avons :

$$\sqrt{n} \left(U_n - \begin{pmatrix} \theta^{*2} + \theta^* \\ \theta^* \end{pmatrix} \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 4\theta^{*3} + 6\theta^{*2} + \theta^* & 2\theta^{*2} + \theta^* \\ 2\theta^{*2} + \theta^* & \theta^* \end{pmatrix} \right).$$

Maintenant, calculons la dérivée :

$$Dl \begin{pmatrix} \theta^{*2} + \theta^* & \theta^* \end{pmatrix} = \begin{pmatrix} 1 \\ -2\theta^* \end{pmatrix}.$$

Il reste donc à calculer la variance :

$$\begin{aligned}
&(1 \quad -2\theta^*) \begin{pmatrix} 4\theta^{*3} + 6\theta^{*2} + \theta^* & 2\theta^{*2} + \theta^* \\ 2\theta^{*2} + \theta^* & \theta^* \end{pmatrix} \begin{pmatrix} 1 \\ -2\theta^* \end{pmatrix} \\
&= (4\theta^{*3} + 6\theta^{*2} + \theta^* - 4\theta^{*3} - 2\theta^{*2} \quad 2\theta^{*2} + \theta^* - 2\theta^{*2}) \begin{pmatrix} 1 \\ -2\theta^* \end{pmatrix} \\
&= (4\theta^{*2} + \theta^* \quad \theta^*) \begin{pmatrix} 1 \\ -2\theta^* \end{pmatrix} \\
&= 4\theta^{*2} + \theta^* - 2\theta^{*2} \\
&= 2\theta^{*2} + \theta^*
\end{aligned}$$

Donc, au final, nous avons :

$$\sqrt{n} (\hat{\theta}_{n,2} - \theta^*) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} (0, 2\theta^{*2} + \theta^*).$$

Exercices 3.2

Nous pouvons au choix utiliser l'un ou l'autre des corollaires. Dans les deux cas, nous commençons par calculer le logarithme de la densité :

$$\begin{aligned}
\log f_\theta(x) &= \log \left(e^{-\theta} \frac{\theta^x}{x!} \right) \\
&= -\theta + x \log \theta - \log x!.
\end{aligned}$$

$$\begin{aligned}
I_1(\theta) &= \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \log f_\theta(X) \right)^2 \right] \\
&= \mathbb{E}_\theta \left[\left(-1 + \frac{X}{\theta} \right)^2 \right] \\
&= \mathbb{V}_\theta \left[\frac{X}{\theta} \right] \text{ car } \frac{X}{\theta} \sim \mathcal{E}(1), \\
&= \frac{1}{\theta^2}
\end{aligned}$$

$$= \frac{1}{\theta}.$$

Autre solution, nous calculons la dérivée seconde :

$$\begin{aligned} I_1(\theta) &= -\mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(X) \right] \\ &= -\mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \left(-1 + \frac{X}{\theta} \right) \right] \\ &= -\mathbb{E}_\theta \left[-\frac{X}{\theta^2} \right] \\ &= \frac{\mathbb{E}_\theta [X]}{\theta^2} \\ &= \frac{\theta}{\theta^2} \\ &= \frac{1}{\theta}. \end{aligned}$$

Le résultat final vient du fait que $I_n(\theta) = nI_1(\theta)$.

Chapitre 4

Intervalle et région de confiance

"Capitaine, je suis sûr à 80% que c'est le meurtrier."

Adrien Monk, personnage éponyme de la série *Monk*.

4.1 Objectifs

Le but de ce chapitre est de comprendre qu'un risque est toujours possible et qu'il faut donc en tenir compte.

4.2 Introduction

Ce chapitre commence par un certain nombre de définitions.

Définition 41 (Région de confiance)

Étant donné $\alpha \in]0, 1[$, une **région de confiance de $g(\theta^*)$ au niveau $1 - \alpha$** est un ensemble $\hat{\mathcal{C}}$ construit mesurablement par rapport à l'observation X et tel que pour tout $\theta^* \in \Theta$, nous ayons :

$$\mathbb{P}_{\theta^*} \left(g(\theta^*) \in \hat{\mathcal{C}} \right) \geq 1 - \alpha.$$

Lorsque l'inégalité devient une égalité, nous disons que le **niveau de confiance est exactement égal à $1 - \alpha$** .

Nous avons immédiatement la généralisation pour un n -échantillon :

Définition 42 (Région de confiance asymptotique)

Étant donné $\alpha \in]0, 1[$, une **région de confiance asymptotique de $g(\theta^*)$ au niveau $1 - \alpha$** est une suite d'ensembles $(\hat{\mathcal{C}}_n)_{n \in \mathbb{N}}$, chacun construit mesurablement par rapport au n -échantillon (X_1, \dots, X_n) et telle que pour tout $\theta^* \in \Theta$, nous ayons :

$$\liminf_{n \rightarrow +\infty} \mathbb{P}_{\theta^*} \left(g(\theta^*) \in \hat{\mathcal{C}}_n \right) \geq 1 - \alpha.$$

Remarque

Usuellement, nous prenons $\alpha = 5\%$ ou $\alpha = 1\%$.

Remarque

Étant donné un niveau de confiance fixé, une région de confiance est d'autant meilleure que son aire est petite.

Définition 43 (Intervalle de confiance)

Si $g(\theta^*) \in \mathbb{R}$, nous parlons plutôt d'**intervalle de confiance**.

**Attention au piège**

Une région de confiance est une variable aléatoire et nous calculons donc la probabilité que cette variable aléatoire englobe le paramètre d'intérêt qui lui est fixé. En particulier, il ne faut pas confondre la région de confiance qui est aléatoire et son estimée (sa réalisation).

Si je prends une région de confiance $\hat{C}_n = [\hat{A}_n, \hat{B}_n]$ telle que :

$$\mathbb{P}_{\theta^*} (\theta^* \in [\hat{A}_n, \hat{B}_n]) = 95\%$$

et que je calcule une réalisation de celle-ci, par exemple $[0.33, 0.75]$ alors on ne peut plus dire que $\mathbb{P}_{\theta^*} (\theta^* \in [0.33, 0.75]) = 95\%$ puisque tout est fixé :

- Soit θ^* est dans l'intervalle et la probabilité vaut 1.
- Soit il ne l'est pas et la probabilité est nulle.

La notion de région de confiance est de dire que si nous répétons l'expérience un grand nombre de fois alors le paramètre d'intérêt se trouvera dans la région de confiance estimée en moyenne 95% du temps.

4.3 Premières constructions

La première façon de construire un intervalle de confiance est d'utiliser les quantiles :

Définition 44 (Quantile d'ordre p)

Pour tout $\beta \in]0; 1[$, nous appelons **quantile d'ordre β** d'une variable aléatoire X de fonction de répartition F_X est définie par

$$q_\beta = \inf \{q \in \mathbb{R} \mid F_X(q) \geq \beta\}.$$

Il est parfois noté $F^{-1}(\beta)$.

En particulier, nous avons :

Définitions 45 (Quantiles particuliers)

- La **médiane** est le quantile d'ordre 50%.
- Les **quartiles** sont les quantiles d'ordre 25*k*% avec $k \in \{1, 2, 3\}$.
- Les **déciles** sont les quantiles d'ordre 10*k*% avec $k \in \{1, \dots, 9\}$.
- Les **centiles** sont les quantiles d'ordre *k*% avec $k \in \{1, \dots, 99\}$.

Remarque

Si la fonction est inversible (voir figure 4.1 (a)), le quantile est directement l'inversion de cette dernière. Si la fonction possède un plateau horizontal (voir figure 4.1 (b)), il faut prendre au début du plateau (la fonction quantile aura donc un saut). Si la fonction possède un saut (voir figure 4.1 (c)), il faut prendre la valeur correspondante à ce saut (la fonction quantile aura donc un plateau horizontal).

Exemple

⋮ Pour la loi $\mathcal{N}(0, 1)$, deux quantiles classiques (à connaître) sont : $q_{0,975} \approx 1,96$ et $q_{0,95} \approx 1,645$.

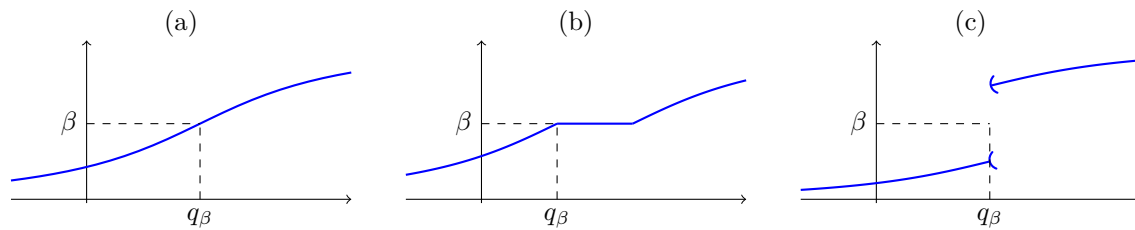


FIGURE 4.1 – Calcul de différents quantiles en fonction du type de fonction de répartition : inversible (a), avec un plateau (b) ou avec un saut (c).



Exercices 4.1

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de même loi $\mathcal{N}(\theta^*, 1)$. Proposer un intervalle de confiance de niveau exactement 95% de la forme $] -\infty, a]$, $[a, +\infty[$ et un **intervalle bilatère** c'est-à-dire de la forme $[a, b]$. Lequel préférez-vous ?

Point méthode (Méthode du pivot)

Pour calculer un intervalle de confiance, nous procédons souvent à l'aide de la **méthode dite du pivot** :

1. Choix d'un estimateur.
2. Calcul de sa loi en fonction de θ^* .
3. Transformation de l'estimateur pour obtenir une statistique dont la loi ne dépend plus de θ^* .
4. Détermination des quantiles de la loi pour estimer la région de confiance adéquat.

4.4 Intervalle de confiance de niveaux obtenus par des inégalités de probabilités

Dans cette partie, nous présentons des intervalles de confiance obtenus à partir d'inégalités de probabilités classiques.

4.4.1 Cas de variances uniformément bornée

Proposition 16 (Inégalité de Bienaymé-Tchebychev)

Soit une variable aléatoire Y admettant un moment d'ordre 2 alors nous avons :

$$\forall \varepsilon > 0, \mathbb{P}(|Y - \mathbb{E}[Y]| \geq \varepsilon) \leq \frac{\mathbb{V}[Y]}{\varepsilon^2}.$$

Ce théorème est une conséquence de l'inégalité de Markov :

Lemme 17 (Inégalité de Markov)

Soit Z une variable aléatoire réelle presque sûrement positive ou nulle alors nous avons :

$$\forall a > 0, \mathbb{P}(Z \geq a) \leq \frac{\mathbb{E}[Z]}{a}.$$

**Preuve**

Pour tout $a > 0$, nous avons presque sûrement $Z \geq a \mathbb{1}_{\{Z \geq a\}}$ et nous avons donc :

$$\mathbb{E}[Z] \geq \mathbb{E}[a \mathbb{1}_{\{Z \geq a\}}] = a \mathbb{P}(Z \geq a) \Leftrightarrow \mathbb{P}(Z \geq a) \leq \frac{\mathbb{E}[Z]}{a}.$$

La preuve de la proposition 16 est alors :

**Preuve**

Comme la fonction $x \mapsto x^2$ est bijective sur \mathbb{R}^+ à valeurs dans \mathbb{R}^+ , nous avons par le lemme 17 pour tout $\varepsilon > 0$:

$$\mathbb{P}(|Y - \mathbb{E}[Y]| \geq \varepsilon) = \mathbb{P}\left((Y - \mathbb{E}[Y])^2 \geq \varepsilon^2\right)$$

et nous pouvons appliquer l'inégalité de Markov avec $Z = (Y - \mathbb{E}[Y])^2$ et $a = \varepsilon^2$.

Exemple

Étant donné un n -échantillon de loi $\mathcal{B}(\theta^*)$, l'estimateur du maximum de vraisemblance est la moyenne $\bar{\mathbf{X}}_n$ et nous avons pour tout $\varepsilon > 0$:

$$\mathbb{P}_{\theta^*}(|\bar{\mathbf{X}}_n - \theta^*| \geq \varepsilon) \leq \frac{\theta^*(1 - \theta^*)}{n\varepsilon^2} \leq \frac{1}{4n\varepsilon^2} = \alpha.$$

Donc, un intervalle de confiance bilatéral de niveau $1 - \alpha$ est :

$$IC_{1-\alpha}(\theta^*) = \left[\bar{\mathbf{X}}_n - \frac{1}{2\sqrt{n\alpha}}, \bar{\mathbf{X}}_n + \frac{1}{2\sqrt{n\alpha}} \right].$$

Pour $\alpha = 5\%$ et $n = 100$, la longueur est d'environ 0.45.

4.4.2 Cas de lois à supports inclus dans un compact donné

Lemme 18 (Inégalité de Hoeffding)

Étant donnés, (Y_1, \dots, Y_n) des variables aléatoires réelles indépendantes et (a_1, \dots, a_n) et (b_1, \dots, b_n) des n -uplets de réels tels que pour tout $i \in \{1, \dots, n\}$:

- $\mathbb{E}[Y_i] = 0$,
- $a_i \leq Y_i \leq b_i$ presque sûrement.

Alors, pour tout $t > 0$, nous avons :

$$\mathbb{P}\left(\sum_{i=1}^n Y_i > t\right) \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

Corollaire 19

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi commune à support inclus dans $[a, b]$ admettant un moment d'ordre 1, nous avons :

$$\forall \varepsilon > 0, \mathbb{P}_{\theta^*}(|\bar{\mathbf{X}}_n - \mathbb{E}_{\theta^*}[X_1]| \geq \varepsilon) \leq 2 \exp\left(-\frac{2n\varepsilon^2}{(b-a)^2}\right).$$

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi $\mathcal{B}(\theta^*)$ alors nous avons presque sûrement $a = 0 \leq X_i \leq b = 1$ et donc pour tout $\varepsilon > 0$:

$$\mathbb{P}_{\theta^*} (|\bar{\mathbf{X}}_n - \theta^*| \geq \varepsilon) \leq 2e^{-2n\varepsilon^2}.$$

Le choix de $\varepsilon = \sqrt{1/(2n) \log(2/\alpha)}$ conduit à :

$$IC_{1-\alpha}(\theta^*) = \left[\bar{\mathbf{X}}_n - \sqrt{\frac{1}{2n} \log \frac{2}{\alpha}}, \bar{\mathbf{X}}_n + \sqrt{\frac{1}{2n} \log \frac{2}{\alpha}} \right].$$

Pour $\alpha = 5\%$ et $n = 100$, la longueur est d'environ 0.27.

Remarque

De manière générale, comme la variance d'une variable aléatoire encadrée par a et b est plus petite que $(b-a)^2/4$, comparer la précision des deux méthodes revient à comparer $1/\sqrt{2\alpha}$ et $\sqrt{\log(2/\alpha)}$. Pour les petites valeurs de α , c'est l'intervalle fourni par l'inégalité de Hoeffding qui propose un meilleur résultat au prix d'une hypothèse plus forte.

4.5 Intervalle de confiance asymptotique

Le théorème de la limite centrale permet d'avoir à chaque fois un intervalle de confiance de niveau asymptotique $1 - \alpha$ pour l'espérance d'une loi puisque nous avons :

$$\sqrt{n} (\bar{\mathbf{X}}_n - \mathbb{E}_{\theta^*} [X_1]) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \mathbb{V}_{\theta^*} [X_1]).$$

Nous pouvons alors proposer l'intervalle :

$$\left[\bar{\mathbf{X}}_n \pm q_{1-\alpha/2} \sqrt{\frac{\mathbb{V}_{\theta^*} [X_1]}{n}} \right]$$

comme intervalle de confiance (asymptotique) à condition de connaître la variance. Or, cette variance dépend quasiment tout le temps de θ^* .

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi $\mathcal{B}(\theta^*)$ alors nous avons :

$$\sqrt{n} (\bar{\mathbf{X}}_n - \theta^*) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \theta^*(1 - \theta^*))$$

et il faudrait isoler θ^* mais c'est compliqué.

4.5.1 Cas de variances uniformément bornées

Dans l'exemple précédent, si nous admettons connaître la variance $\theta^*(1 - \theta^*)$ alors nous pouvons proposer :

$$IC_{1-\alpha}(\theta^*) = \left[\bar{\mathbf{X}}_n \pm q_{1-\alpha/2} \sqrt{\frac{\theta^*(1 - \theta^*)}{n}} \right] \subset \left[\bar{\mathbf{X}}_n \pm \frac{q_{1-\alpha/2}}{2\sqrt{n}} \right].$$

Pour $\alpha = 5\%$ et $n = 100$, nous trouvons une longueur de 0.196.

4.5.2 Estimation consistante de la variance

S'il existe une suite d'estimateurs $(\hat{\sigma}_n^2)_{n \in \mathbb{N}^*}$ de la variance consistante alors d'après le lemme de Slutsky, nous avons pour tout $\theta^* \in \Theta$:

$$\sqrt{\frac{n}{\hat{\sigma}_n^2}} (\bar{\mathbf{X}}_n - \mathbb{E}_{\theta^*} [X_1]) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

D'où l'intervalle de confiance de niveau asymptotiquement exact $1 - \alpha$:

$$IC_{1-\alpha}(\theta^*) = \left[\bar{\mathbf{X}}_n \pm q_{1-\alpha/2} \sqrt{\frac{\widehat{\sigma}_n^2}{n}} \right].$$

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi $\mathcal{B}(\theta^*)$ alors nous avons :

$$IC_{1-\alpha}(\theta^*) = \left[\bar{\mathbf{X}}_n \pm q_{1-\alpha/2} \sqrt{\frac{\bar{\mathbf{X}}_n (1 - \bar{\mathbf{X}}_n)}{n}} \right].$$

4.5.3 Stabilisation de la variance

À l'aide de la méthode Delta, nous avons que pour toute fonction φ continument dérivable sur $\Theta \subset \mathbb{R}$, nous avons :

$$\sqrt{n} (\varphi(\bar{\mathbf{X}}_n) - \varphi(\mathbb{E}_{\theta^*} [X_1])) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(0, \varphi'^2 \mathbb{V}_{\theta^*} [X_1] \right).$$

Donc, si prenons φ telle que $\varphi'(\theta^*) = 1/\sqrt{\mathbb{V}_{\theta^*} [X_1]}$ alors nous aura la convergence en loi suivante :

$$\sqrt{n} (\varphi(\bar{\mathbf{X}}_n) - \varphi(\mathbb{E}_{\theta^*} [X_1])) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1)$$

et nous pouvons en tirer l'intervalle de confiance suivant :

$$IC_{1-\alpha}(\theta^*) = \varphi^{-1} \left(\left[\varphi(\bar{\mathbf{X}}_n) \pm \frac{q_{1-\alpha/2}}{\sqrt{n}} \right] \right)$$

c'est-à-dire l'image inverse de l'intervalle.

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi $\mathcal{B}(\theta^*)$, nous devons résoudre $\varphi'(\theta^*) = 1/\sqrt{\theta^*(1-\theta^*)}$. Or, nous savons que la dérivée de la fonction arcsin est :

$$\arcsin'(x) = \frac{1}{\sqrt{1-x^2}}$$

alors en prenant $\varphi(x) = 2 \arcsin \sqrt{x}$, nous avons :

$$IC_{1-\alpha}(\theta^*) = \left[\sin^2 \left(\arcsin \sqrt{\bar{\mathbf{X}}_n} \pm \frac{q_{1-\alpha/2}}{\sqrt{n}} \right) \right].$$

Deuxième partie

UE Test

Chapitre 5

Introduction aux tests

"Dans la justice, il y a deux visions :

- En France, nous sommes présumé innocent jusqu'à ce que notre culpabilité soit établie ;
- Aux États-Unis, nous sommes présumé coupable jusqu'à ce que notre innocence soit établie."

5.1 Objectifs

La notion de test est assez complexe car elle repose sur une asymétrie des deux hypothèses mises en concurrence.

5.2 Formalisme et démarche expérimentale

Dans le cadre du modèle paramétrique $(E, \mathcal{E}, \mathbb{P}_\theta, \theta \in \Theta)$, nous nous donnons deux sous-ensembles Θ_0 et Θ_1 **disjoints** et inclus dans Θ (nous n'imposons pas qu'ils forment une partition). À la vue d'une observation, nous souhaitons décider si $\theta^* \in \Theta_0$ ou pas ; si ce n'est pas le cas, nous considérerons que $\theta^* \in \Theta_1$.

Définitions 46 (Hypothèses)

Étant donnés deux sous-ensembles Θ_0 et Θ_1 **disjoints** et inclus dans Θ , nous définissons respectivement :

\mathcal{H}_0 : l'**hypothèse nulle** considère que $\theta^* \in \Theta_0$.

\mathcal{H}_1 : l'**hypothèse alternative** considère que $\theta^* \in \Theta_1$.

Pour tout $j \in \{0, 1\}$, nous dirons que \mathcal{H}_j est **simple** si Θ_j est réduit à un singleton et **composite** sinon.

Exemple

⋮ Dans le cas d'une pièce, nous nous intéressons souvent à savoir si la probabilité d'avoir *pile* est égale à celle d'avoir *face* ($\Theta_0 = \{1/2\}$) ou pas ($\Theta_1 = [0, 1] \setminus \{1/2\}$).
⋮

À l'aide de ces notations, nous pouvons introduire la définition d'un test :

Définitions 47 (Test)

Nous appelons **test de l'hypothèse \mathcal{H}_0 contre l'hypothèse \mathcal{H}_1** toute fonction $\phi(X)$ à valeur dans $\{0, 1\}$ avec ϕ mesurable et pouvant dépendre de Θ_0 et Θ_1 :

- Lorsque $\phi(X) = 0$, nous disons que nous **conservons** l'hypothèse \mathcal{H}_0 .
- Lorsque $\phi(X) = 1$, nous disons que nous **rejetons** l'hypothèse \mathcal{H}_0 et nous **acceptons** l'hypothèse \mathcal{H}_1 .

Remarque

Un test peut toujours s'écrire $\phi(X) = \mathbb{1}_{\{X \in \mathcal{R}\}}$ où $\mathcal{R} \in \mathcal{E}$ ou encore $\phi(X) = \mathbb{1}_{\{h(X) \in \mathcal{R}'\}}$ avec h mesurable et $\mathcal{R}' \in \mathcal{B}(\mathbb{R})$ la tribu des boréliens de \mathbb{R} .

Définitions 48 (Statistique de test)

Dans le cas où $\phi(X) = \mathbb{1}_{\{h(X) \in \mathcal{R}'\}}$ avec h mesurable et $\mathcal{R}' \in \mathcal{B}(\mathbb{R})$ la tribu des boréliens de \mathbb{R} , nous appelons \mathcal{R}' la **région de rejet** et $h(X)$ la **statistique de test**. Le complémentaire de la région de rejet est parfois appelée **région d'acceptation**.

5.2.1 Mesure de la qualité d'un test

Pour mesurer la qualité d'un test, nous nous intéressons d'une part à la probabilité d'accepter à tort et d'autre part de rejeter à tort.

Définitions 49 (Risques de première et seconde espèce)

Les **risques de première** et **de seconde espèce** du test $\phi(X)$ sont définis comme les fonctions $\underline{\alpha}$ sur Θ_0 et $\underline{\beta}$ sur Θ_1 par :

$$\begin{aligned} \underline{\alpha} : \Theta_0 &\rightarrow [0, 1] & \text{et} & \quad \underline{\beta} : \Theta_1 \rightarrow [0, 1] \\ \theta^* &\mapsto \mathbb{P}_{\theta^*}(\phi(X) = 1) & & \quad \theta^* \mapsto \mathbb{P}_{\theta^*}(\phi(X) = 0) \end{aligned}$$

et nous pouvons introduire la **puissance d'un test** comme la fonction $1 - \underline{\beta}$ et sa **taille** comme le réel α^* défini par

$$\alpha^* = \sup_{\theta^* \in \Theta_0} \underline{\alpha}(\theta^*) = \sup_{\theta^* \in \Theta_0} \mathbb{P}_{\theta^*}(\phi(X) = 1).$$

Enfin, nous disons que le test est de **niveau** α si sa taille α^* est inférieure ou égale à α .

Comme pour les intervalles de confiances, la notion de test s'étend également au cas asymptotique

Définition 50 (Version asymptotique)

Étant donnée une suite $(\phi_n)_{n \in \mathbb{N}}$ de tests et en notant α_n^* la taille de ϕ_n , nous disons que la suite de tests est de **niveau asymptotique** α si

$$\limsup_{n \rightarrow +\infty} \alpha_n^* \leq \alpha.$$

Point méthode (Représentation sous forme de tableau)

Nous pouvons résumer les informations précédentes de la façon suivante :

θ^*	Θ_0	Θ_1
$\phi(X)$		
0	Conservation correcte	Risque seconde espèce : $\underline{\beta}$
1	Risque première espèce : $\underline{\alpha}$	Puissance

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ iid de loi $\mathcal{N}(\theta^*, 1)$ avec $\theta \in \mathbb{R}$, nous voulons tester :

$$\begin{cases} \mathcal{H}_0 : \theta^* \geq 1, \\ \mathcal{H}_1 : \theta^* < 1, \end{cases} \quad \text{soit} \quad \begin{cases} \Theta_0 = [1, +\infty[, \\ \Theta_1 =]-\infty, 1]. \end{cases}$$

L'intuition nous conduit à choisir $\phi(X) = \mathbb{1}_{\{\bar{X}_n < 1\}}$: le test est fondé sur la statistique \bar{X}_n et la région de rejet $]-\infty, 1[$.

On voit que pour tout $\theta^* \geq 1$, $\alpha(\theta^*) = \mathbb{P}_{\theta^*}(\bar{X}_n < 1) = F_{\mathcal{N}(0,1)}(\sqrt{n}(1 - \theta^*))$ car $\bar{X}_n \sim \mathcal{N}(\theta^*, 1/n)$ et pour tout $\theta^* < 1$, $\beta(\theta^*) = 1 - F_{\mathcal{N}(0,1)}(\sqrt{n}(1 - \theta^*))$.

En revanche, la taille du test vaut $\alpha^* = 1/2$ ce qui veut dire qu'on ne fait pas mieux qu'un lancer de pile ou face.

5.2.2 Dissymétrie des rôles des hypothèses \mathcal{H}_0 et \mathcal{H}_1

Si nous voulons minimiser l'erreur de première espère, il suffirait de toujours conserver l'hypothèse \mathcal{H}_0 ; à l'opposé, si nous souhaitons minimiser l'erreur de seconde espèce, il suffit de rejeter à chaque fois l'hypothèse \mathcal{H}_0 au profit de l'hypothèse \mathcal{H}_1 . Nous voyons donc que minimiser les deux erreurs en même temps est compliquée (voir impossible sauf dans certains cas). Il a donc été choisi de dissymétriser le problème en faisant jouer un rôle particulier à \mathcal{H}_0 et en contrôlant exclusivement le risque de première espèce. On prendra par exemple pour \mathcal{H}_0 :

- Une hypothèse communément admise : par exemple, il n'y a pas de traces de vie autre que sur Terre dans notre système solaire.
- Une hypothèse de prudence ou conservatrice : le nouveau médicament qui est en train d'être testé ne soigne pas mieux que l'ancien.
- Ou plus simplement la seule hypothèse que l'on peut facilement formuler.

Dans cette vision le statisticien impose le niveau α (par exemple 5%) et acceptera de rejeter à tort son hypothèse dans 5% des cas.

5.2.3 Démarche de construction et mise en œuvre pratique

Exemple

Étant donné un n -échantillon X_1, \dots, X_n de loi normale $\mathcal{N}(\theta^*, 1)$ et regardons le test suivant :

$$\begin{cases} \mathcal{H}_0 : \theta^* \geq 1, \\ \mathcal{H}_1 : \theta^* < 1. \end{cases}$$

La statistique naturelle est toujours \bar{X}_n et nous choisissons, à la vue de \mathcal{H}_1 , une région de rejet de la forme $]-\infty, k_\alpha[$ c'est-à-dire une statistique de la forme $\phi(X) = \mathbb{1}_{\{\bar{X}_n < k_\alpha\}}$. La question revient donc au choix de k_α .

Pour ce faire, nous calculons la taille du test :

$$\begin{aligned} \alpha^* &= \sup_{\theta^* \in \Theta_0} \mathbb{P}_{\theta^*}(\bar{X}_n < k_\alpha) \\ &= \sup_{\theta^* \geq 1} \mathbb{P}_{\theta^*}(\sqrt{n}(\bar{X}_n - \theta^*) < \sqrt{n}(k_\alpha - \theta^*)) \\ &= \sup_{\theta^* \geq 1} F[\sqrt{n}(k_\alpha - \theta^*)] \text{ où } F \text{ est la fonction de répartition d'une gaussienne centrée réduite,} \\ &= F[\sqrt{n}(k_\alpha - 1)] \text{ car } F \text{ est croissante.} \end{aligned}$$

Donc, pour que le niveau du test soit exactement α , il suffit que :

$$\begin{aligned} \alpha = \alpha^* &\Leftrightarrow \alpha = F[\sqrt{n}(k_\alpha - 1)] \\ &\Leftrightarrow q_\alpha = \sqrt{n}(k_\alpha - 1) \text{ où } q_\alpha \text{ est le quantile d'une gaussienne centrée réduite,} \\ &\Leftrightarrow \frac{q_\alpha}{\sqrt{n}} + 1 = k_\alpha. \end{aligned}$$

Donc, le test est de la forme $\phi(X) = \mathbb{1}_{\{\bar{X}_n < \frac{q_\alpha}{\sqrt{n}} + 1\}}$ et a un niveau exactement de α .

5.2.4 Qualités d'un test

Nous recensons un certain nombre de propriétés éventuelles.

Définitions 51

Un test est dit **sans biais** lorsque sa fonction puissance vérifie $1 - \underline{\beta} > \alpha$ sur Θ_1 .

Une suite de tests est dit **consistante** s'ils sont tous de niveaux $\bar{\alpha}$ et si, pour tout $\theta^* \in \Theta_1$, nous avons :

$$\underline{\beta}_n(\theta^*) \xrightarrow{n \rightarrow +\infty} 0.$$

La **robustesse** exprime l'absence de sensibilité d'un test, éventuellement asymptotique, au modèle (où à la forme de l'observation). Entre d'autres termes, si la modélisation du phénomène est mauvaise, est-ce que nous obtenons les mêmes conclusions ?

Comme pour les estimateurs, nous pouvons chercher des tests ayant des propriétés plus intéressantes que les autres.

Définition 52 (Test uniformément plus puissant)

Pour deux tests $\phi(X)$ et $\phi'(X)$ de niveau α , nous disons que $\phi(X)$ est **uniformément plus puissant (UPP)** que $\phi'(X)$ si leurs fonctions puissances $1 - \underline{\beta}$ et $1 - \underline{\beta}'$ vérifient $1 - \underline{\beta} \geq 1 - \underline{\beta}'$; c'est-à-dire que pour tout $\theta^* \in \Theta_1$, nous avons :

$$\mathbb{P}_{\theta^*}(\phi(X) = 1) \geq \mathbb{P}_{\theta^*}(\phi'(X) = 1).$$

Nous disons que $\phi(X)$ est $UPP(\alpha)$

5.3 Un outil important : p -valeur

Le principe de la p -valeur (ou p -value en anglais) est expliquée sur la figure 5.1 : plutôt que de demander la valeur α visée par l'utilisateur, la plupart des logiciels renvoie la valeur estimée à partir de la moyenne ; il ne reste plus qu'à l'utilisateur de comparer cette valeur à la cible.

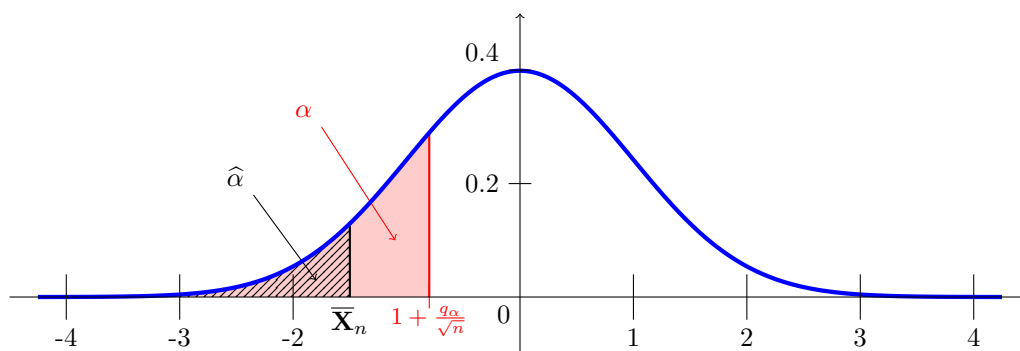


FIGURE 5.1 – Représentation schématique de la p -valeur (zone hachurée) par rapport à un niveau α visé (zone rouge).

Définition 53 (p -valeur)

Supposons avoir construit une famille de test $\phi_\alpha(X)$, chacun de niveau α , pour $\alpha \in]0, 1[$. La p -valeur associée à cette famille et à l'observation X est le réel défini par :

$$\hat{\alpha}(X) = \sup \{ \alpha \in]0, 1[\mid \phi_\alpha(X) = 0 \}.$$

Remarque

En d'autres termes, $\hat{\alpha}(X)$ est le plus grand niveau autorisant l'acceptation de \mathcal{H}_0 et forme donc, en un certain sens, un indice de crédibilité de \mathcal{H}_0 à la vue de X . La p -valeur mesure l'adéquation entre \mathcal{H}_0 et les observations : plus $\hat{\alpha}(X)$ est faible et plus nos informations contredisent \mathcal{H}_0 .

Le cas typique est celui où la famille $(\phi_\alpha(X))_{\alpha \in]0,1[}$ est p.s. croissante en α , c'est-à-dire que $\phi_\alpha(X)$ décroît vers 0 lorsque α décroît vers 0 (ce qui est une condition naturelle). Ceci équivaut à dire que la famille des régions de rejet est croissante au sens de l'inclusion. Dans ce cas, on a également :

$$\hat{\alpha}(X) = \inf \{ \alpha \in]0,1[\mid \phi_\alpha(X) = 0 \}.$$

En outre, le test ϕ_α de niveau α conserve \mathcal{H}_0 à la vue de X pour $\alpha < \hat{\alpha}(X)$ et rejette $\alpha > \hat{\alpha}(X)$.

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi $\mathcal{N}(\theta^*, 1)$ et si nous voulons tester :

$$\begin{cases} \mathcal{H}_0 : \theta^* \geq 1, \\ \mathcal{H}_1 : \theta^* < 1. \end{cases}$$

Nous avons vu que la région de rejet est $\bar{X}_n < k_\alpha = 1 + \frac{q_\alpha}{\sqrt{n}}$.

La p -valeur est caractérisée par le cas limite $k_{\hat{\alpha}_n} = \bar{X}_n$ soit $\hat{\alpha}_n(\mathbf{X}) = F(\sqrt{n}(\bar{X}_n - 1))$.

Avec $\bar{X}_4 = 0.5$, nous obtenons $\hat{\alpha}(\mathbf{X}) \approx 15,9\%$ et nous conservons donc l'hypothèse \mathcal{H}_0 avec un niveau de 10% ou même de 5%.

5.4 Test du rapport de vraisemblance

Il s'agit d'une méthode générale pour construire des tests $UPP(\alpha)$ sous certaines hypothèses.

Définitions 54 (Statistique du rapport de vraisemblance)

Étant donnés deux ensembles Θ_0 et Θ_1 représentant les deux hypothèses, nous introduisons la **statistique du rapport de vraisemblance** notée h par :

$$h(X) = \frac{\sup_{\theta \in \Theta_1} V_X(\theta)}{\sup_{\theta \in \Theta_0} V_X(\theta)}$$

où nous rappelons que $V_X(\theta) = f_\theta(X)$ est la vraisemblance du modèle. La zone de rejet naturelle est donc de la forme $]k_\alpha, +\infty[$ où k_α est à fixer suivant le contexte et le **test du rapport de vraisemblance** s'écrit donc $\phi(X) = \mathbb{1}_{\{h(X) > k_\alpha\}}$.

Remarque

Si la valeur $\sup_{\theta \in \Theta_0} V_X(\theta)$ est nulle, nous supposons que le test vaut 1 puisque le modèle 0 est rejeté. Toutefois, si le numérateur et le dénominateur sont nuls, il faudra remettre en cause la pertinence des hypothèses.

Exemple

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ iid de loi $\mathcal{N}(\theta^*, 1)$ avec $\theta^* \in \mathbb{R}$, nous voulons tester :

$$\begin{cases} \mathcal{H}_0 : \theta^* \geq 1, \\ \mathcal{H}_1 : \theta^* < 1. \end{cases}$$

La vraisemblance s'écrit alors :

$$V_X(\theta) = \left(\frac{1}{\sqrt{2\pi}} \right)^n \exp \left[-\frac{1}{2} \sum_{i=1}^n (X_i - \theta)^2 \right]$$

est donc strictement croissante sur $] -\infty, \bar{\mathbf{X}}_n]$ et strictement décroissante sur $[\bar{\mathbf{X}}_n, +\infty[$. Ainsi, nous avons le rapport de vraisemblance suivant :

$$h(\mathbf{X}) = \begin{cases} \frac{V_{\mathbf{X}}(\bar{\mathbf{X}}_n)}{V_{\mathbf{X}}(\theta_1)} & \text{si } \bar{\mathbf{X}}_n < 1, \\ \frac{V_{\mathbf{X}}(\theta_1)}{V_{\mathbf{X}}(\bar{\mathbf{X}}_n)} & \text{si } \bar{\mathbf{X}}_n \geq 1. \end{cases}$$

Or, pour tout θ_0 et θ_1 de \mathbb{R} , nous avons :

$$\frac{V_{\mathbf{X}}(\theta_1)}{V_{\mathbf{X}}(\theta_0)} = \exp\left(\frac{n}{2} [\theta_0^2 - \theta_1^2 + 2\bar{\mathbf{X}}_n (\theta_1 - \theta_0)]\right)$$

de sorte que le rapport de vraisemblance s'écrit $h(X) = \exp\left[\frac{n}{2} (\bar{\mathbf{X}}_n - 1)^2 (2\mathbb{1}_{\{\bar{\mathbf{X}}_n < 1\}} - 1)\right]$. Or, il s'agit d'une fonction strictement décroissante en $\bar{\mathbf{X}}_n$ donc le test du rapport de vraisemblance peut se réécrire de la forme $\mathbb{1}_{\{\bar{\mathbf{X}}_n < k_\alpha\}}$.

Théorème 20 (Lemme de Neyman-Pearson)

Étant donné un modèle paramétré par Θ et deux points θ_0 et θ_1 de Θ distincts, nous voulons tester :

$$\begin{cases} \mathcal{H}_0 : \theta^* = \theta_0, \\ \mathcal{H}_1 : \theta^* = \theta_1. \end{cases}$$

Étant donné $\alpha \in]0, 1[$, si nous supposons qu'il existe un seul k_α tel que le test de rapport de vraisemblance $\phi(X) = \mathbb{1}_{\{h(X; \theta_0, \theta_1) > k_\alpha\}}$ avec

$$h(X; \theta_0, \theta_1) = \frac{\sup_{\theta \in \Theta_1} V_X(\theta)}{\sup_{\theta \in \Theta_0} V_X(\theta)} = \frac{f_{\theta_1}(X)}{f_{\theta_0}(X)}$$

vérifie $\mathbb{P}_{\theta_0}(\phi(X) = 1) = \alpha$ alors ce test est *UPP*(α).



Preuve

Nous considérons qu'il existe un autre test Ψ de niveau α tel que $\mathbb{E}_{\theta_0}[\Psi(X)] = \mathbb{P}_{\theta_0}(\Psi(X) = 1) \leq \alpha$. Nous avons donc par construction $\mathbb{E}_{\theta_0}[\phi(X) - \Psi(X)] \geq 0$. Nous voulons donc montrer que le test de rapport de vraisemblance est au moins aussi puissant c'est-à-dire que

$$\mathbb{E}_{\theta_1}[\phi(X) - \Psi(X)] = \mathbb{P}_{\theta_1}(\phi(X) = 1) - \mathbb{P}_{\theta_1}(\Psi(X) = 1) \geq 0.$$

Pour ce faire, nous allons faire un changement de densité en remarquant que partout où la densité θ_0 est non nulle, nous avons

$$f_{\theta_1}(x) = \frac{f_{\theta_1}(x)}{f_{\theta_0}(x)} f_{\theta_0}(x).$$

Nous avons donc :

$$\begin{aligned} & \mathbb{E}_{\theta_1}[\phi(X) - \Psi(X)] - k_\alpha \mathbb{E}_{\theta_0}[\phi(X) - \Psi(X)] \\ &= \mathbb{E}_{\theta_0} \left[(\phi(X) - \Psi(X)) \frac{f_{\theta_1}(X)}{f_{\theta_0}(X)} \right] + \mathbb{E}_{\theta_1} \left[(\phi(X) - \Psi(X)) \mathbb{1}_{\{f_{\theta_0}(X)=0\}} \right] - k_\alpha \mathbb{E}_{\theta_0}[\phi(X) - \Psi(X)] \\ &= \mathbb{E}_{\theta_0} \left[(\phi(X) - \Psi(X)) \left(\frac{f_{\theta_1}(X)}{f_{\theta_0}(X)} - k_\alpha \right) \right] + \mathbb{E}_{\theta_1} \left[(\phi(X) - \Psi(X)) \mathbb{1}_{\{f_{\theta_0}(X)=0\}} \right] \text{ car } k_\alpha \text{ fini si } \alpha \in]0, 1[. \end{aligned}$$

Or, $\phi(X) = \mathbb{1}_{\{h(X; \theta_0, \theta_1) > k_\alpha\}}$ donc si elle vaut 0 nous avons :

$$\mathbb{E}_{\theta_0} \left[\underbrace{-\Psi(X)}_{>0} \underbrace{\left(\frac{f_{\theta_1}(X)}{f_{\theta_0}(X)} - k_\alpha \right)}_{<0} \right] > 0,$$

et si elle vaut 1, l'espérance est également positive par construction. De plus, pour les valeurs x telles que $f_{\theta_0}(x) = 0$, l'hypothèse nulle ne peut pas être choisie donc comme k_α est fini, la statistique vaut 1. Au final, comme $k_\alpha > 0$, nous avons donc :

$$\mathbb{E}_{\theta_1} [\phi(X) - \Psi(X)] \geq k_\alpha \mathbb{E}_{\theta_0} [\phi(X) - \Psi(X)] \geq 0.$$

5.5 Démarche générale

S'il ne faut retenir qu'une seule chose de ce cours, c'est la démarche générale. Toute la suite consistera à l'appliquer dans des cas particuliers :

Point méthode

Démarche générale des tests La démarche générale pour faire un test est la suivante :

1. Bien comprendre la problématique pour proposer une modélisation adéquate.
2. Définir l'hypothèse \mathcal{H}_0 qui sera la plus conservatrice.
3. Définir l'hypothèse \mathcal{H}_1 alternative.
4. Définir la statistique de test et la région de rejet. Pour cela, il sera nécessaire de :
 - Proposer un estimateur du paramètre d'intérêt.
 - Définir la loi sous l'hypothèse \mathcal{H}_0 , éventuellement asymptotique, de cet estimateur.
 - Étant donné $\alpha \in [0; 1]$, proposer une règle de décision de niveau de confiance α , éventuellement asymptotique.
5. Étudier la puissance du test en déterminant la probabilité de rejeter à raison l'hypothèse \mathcal{H}_0 :
 - Soit de façon théorique si des calculs sont possibles.
 - Soit de façon empirique à l'aide de simulations.

Exercices 5.1

Une entreprise d'appels téléphoniques possède 10 personnes dont le contrat consiste à gérer une centaine d'appels par jour. Les syndicats affirment qu'ils ne sont pas assez nombreux pour traiter toutes les demandes journalières et qu'il faudrait embaucher une personne de plus. L'entreprise vous embauche afin de savoir s'il est effectivement nécessaire de recruter quelqu'un en plus ou pas et, pour cela, elle vous fournit le nombre d'appels journaliers des n derniers jours. Quel test proposez-vous ?

Indications : La loi de Poisson est souvent utilisée pour les comptages. De plus, si X_1, \dots, X_n sont des variables aléatoires indépendantes suivant des lois de Poisson de paramètres $\lambda_1, \dots, \lambda_n$ alors la somme $\sum_{i=1}^n X_i$ est une variable aléatoire de Poisson de paramètre $\sum_{i=1}^n \lambda_i$.

5.6 Solutions des exercices

Exercices 5.1

1. Nous supposons que le nombre d'appels journaliers suit une loi de Poisson de paramètre θ^* .
2. L'hypothèse la plus conservatrice pour l'entreprise est que les 10 personnes suffisent. Comme chaque personne peut faire en moyenne 100 appels, cela représente 1000 appels par jour donc l'hypothèse \mathcal{H}_0 est :

$$\mathcal{H}_0 : \theta^* = 1000$$

ou, autrement dit, $\Theta_0 = \{1000\}$.

3. L'hypothèse proposée par les syndicats est qu'il manque une personne donc qu'il y a plutôt 1100 appels quotidiens. L'hypothèse \mathcal{H}_1 est donc :

$$\mathcal{H}_1 : \theta^* = 1100$$

ou, autrement dit, $\Theta_1 = \{1100\}$.

4. Comme nous sommes dans un cas où les deux hypothèses sont simples, nous savons par le lemme de Neyman-Pearson que le test du rapport de vraisemblance est *uniformément plus puissant*. Nous commençons donc par calculer la vraisemblance étant donnée un paramètre $\theta \in \mathbb{R}_+^*$ quelconque :

$$\begin{aligned} V_{\mathbf{X}}(\theta) &= \prod_{i=1}^n \left(e^{-\theta} \frac{\theta^{X_i}}{X_i!} \right) \\ &= e^{\sum_{i=1}^n (-\theta)} \frac{\theta^{\sum_{i=1}^n X_i}}{\prod_{i=1}^n (X_i!)} \\ &= e^{-n\theta} \frac{\theta^{\sum_{i=1}^n X_i}}{\prod_{i=1}^n (X_i!)}. \end{aligned}$$

Ainsi, la statistique de test est :

$$\begin{aligned} h(X; 1000, 1100) &= \frac{\sup_{\theta \in \Theta_1} V_X(\theta)}{\sup_{\theta \in \Theta_0} V_X(\theta)} \\ &= \frac{V_{\mathbf{X}}(1100)}{V_{\mathbf{X}}(1000)} \\ &= \frac{e^{-1100n} \frac{1100^{\sum_{i=1}^n X_i}}{\prod_{i=1}^n (X_i!)}}{e^{-1000n} \frac{1000^{\sum_{i=1}^n X_i}}{\prod_{i=1}^n (X_i!)}} \\ &= e^{-1100n+1000n} \left(\frac{1100}{1000} \right)^{\sum_{i=1}^n X_i} \\ &= e^{-100n} \left(\frac{11}{10} \right)^{\sum_{i=1}^n X_i}. \end{aligned}$$

La statistique de test est donc du type :

$$\phi(\mathbf{x}) = \mathbb{1}_{\{h(X; 1000, 1100) > k_\alpha\}} \quad (5.1)$$

où k_α est choisi tel que le niveau de confiance soit de α . Or, la fonction $x \mapsto e^{-100n} \left(\frac{11}{10} \right)^x$ est croissante en x donc l'équation (5.1) revient à trouver une valeur k_α^{somme} telle que :

$$\phi(\mathbf{x}) = \mathbb{1}_{\{\sum_{i=1}^n X_i > k_\alpha^{\text{somme}}\}} \quad (5.2)$$

ou une valeur $k_\alpha^{\text{moyenne}}$ telle que :

$$\phi(\mathbf{x}) = \mathbb{1}_{\{\bar{\mathbf{x}}_n > k_\alpha^{\text{moyenne}}\}}. \quad (5.3)$$

Deux solutions s'offre à nous :

- Comme, sous \mathcal{H}_0 , les X_i sont un n -échantillon de loi $\mathcal{P}(1000)$ alors $\sum_{i=1}^n X_i$ suit une loi de Poisson de paramètre $\sum_{i=1}^n 1000 = 1000n$ et nous pouvons prendre pour k_α^{somme} le quantile d'ordre $1 - \alpha$ de la loi de Poisson $\mathcal{P}(1000n)$. Dans ce cas, le test (5.2) sera quasiment de niveau exact α (*quasiment* car la loi étant entière, il y a une petite fluctuation. Nous verrons cela en TP).
- Sinon, sous \mathcal{H}_0 , $\bar{\mathbf{X}}_n$ suit asymptotiquement une loi normale de moyenne 1000 et de variance $1000/n$, nous pouvons recentrer et normaliser les termes et avoir le test (5.3) de niveau asymptotique α .

Chapitre 6

Tests sur les variables quantitatives

6.1 Objectifs

Comprendre d'où viennent les tests relatifs aux variables quantitatives.

6.2 Test d'une espérance

Dans le cadre de variables quantitatives, les premières statistiques étudiées sont souvent la position et, en particulier, l'espérance. Ainsi, nous disposons d'un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de variables aléatoires suivant une loi d'espérance μ^* et de variance σ^2 ; tous les deux inconnus. Dans notre cas, nous cherchons à savoir si la valeur de l'espérance vaut une valeur μ_0 particulière ou si elle vaut autre chose :

$$\begin{cases} \mathcal{H}_0 : \mu^* = \mu_0, \\ \mathcal{H}_1 : \mu^* \neq \mu_0. \end{cases}$$

Remarque

Suivant le contexte, nous pouvons avoir une information supplémentaire sur la forme de μ^* dans l'hypothèse alternative par exemple :

$$\mathcal{H}_1 : \mu^* < \mu_0 \text{ ou } \mathcal{H}_1 : \mu^* > \mu_0.$$

Dans ce cas, nous adapterons les zones des prédictions.

Pour estimer μ^* , nous pouvons prendre l'estimateur des moments $\bar{\mathbf{X}}_n$ (qui est, pour beaucoup de lois classiques, l'estimateur du maximum de vraisemblance également).

6.2.1 Cas gaussien avec la variance connue

Nous commençons par supposer que nous connaissons la variance σ^2 . Dans ce cas, nous avons la loi suivante :

Proposition 21 (Loi de la variable $\bar{\mathbf{X}}_n$ dans le cas gaussien avec variance σ^2 connue)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi gaussienne $\mathcal{N}(\mu^*, \sigma^2)$ avec σ^2 connue alors la loi de $\bar{\mathbf{X}}_n$ suit une loi gaussienne $\mathcal{N}(\mu^*, \sigma^2/n)$:

$$\bar{\mathbf{X}}_n = \frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}(\mu^*, \sigma^2/n).$$

Preuve

Pour montrer cette proposition, nous avons besoin d'un résultat sur les vecteurs gaussiens :



Lemme 22 (Linéarité des vecteurs gaussiens)

Étant donnée $\mathbf{X} = (X_1, \dots, X_d)$ un vecteur gaussien d'espérance m et de matrice de variance covariance Σ alors pour toute matrice $A \in \mathcal{M}_{k \times d}(\mathbb{R})$ et pour tout vecteur $b \in \mathbb{R}^k$, le vecteur $A\mathbf{X} + b$ est gaussien de loi $\mathcal{N}(Am + b, A\Sigma A^T)$.

Démonstration. Ce lemme est une conséquence de la définition d'un vecteur gaussien. Pour montrer cette propriété, nous avons besoin de la forme de la fonction caractéristique et des propriétés de cette dernière. Pour une démonstration plus complète, nous encourageons les lectrices et lecteurs intéressés à lire le polycopié de Le Gall (2006)¹ ou la proposition 16 du polycopié de Brault (2021a)². \square

Le lemme 22 a pour conséquence que toute combinaison linéaire de variables gaussiennes suit une loi gaussienne (éventuellement triviale). Ainsi, nous savons que $\bar{\mathbf{X}}_n$ suit une loi normale et il ne reste qu'à connaître les paramètres. Par linéarité de l'espérance, nous avons :

$$\mathbb{E}[\bar{\mathbf{X}}_n] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] = \frac{1}{n} \sum_{i=1}^n \mu^* = \frac{n\mu^*}{n} = \mu^*.$$

Et pour la variance, comme les variables sont indépendantes, nous avons :

$$\begin{aligned} \mathbb{V}[\bar{\mathbf{X}}_n] &= \mathbb{V}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \\ &= \frac{1}{n^2} \mathbb{V}\left[\sum_{i=1}^n X_i\right] \quad \text{car la variance est quadratique,} \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{V}[X_i] \quad \text{car les observations sont indépendantes,} \\ &= \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}. \end{aligned}$$

Ainsi, nous obtenons une loi gaussienne de moyenne μ^* et de variance σ^2/n .

Ainsi, si nous connaissons la variance dans le cas de loi gaussienne, nous avons le point méthode suivant :

Point méthode (Test d'espérance ; cas gaussien, variance connue)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi gaussienne $\mathcal{N}(\mu^*, \sigma^2)$ avec σ^2 connue alors la statistique de test est

$$T_n = \sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu_0}{\sigma}$$

qui a pour loi une gaussienne centrée réduite sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau exact $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\mu^* \neq \mu_0$	$\mu^* < \mu_0$	$\mu^* > \mu_0$
Θ_1	$\mathbb{R} \setminus \{\mu_0\}$	$] -\infty, \mu_0[$	$] \mu_0; +\infty[$
Règle	$ T_n > q_{1-\alpha/2}$	$T_n < q_\alpha$	$T_n > q_{1-\alpha}$

où q_α est le quantile d'ordre α d'une loi normale centrée réduite ; c'est-à-dire que si X suit une loi gaussienne centrée réduite $\mathcal{N}(0, 1)$ alors q_α est tel que :

$$\mathbb{P}(X \leq q_\alpha) = \alpha.$$

1. Disponible à l'adresse <https://www.imo.universite-paris-saclay.fr/~jean-francois.le-gall/IPPA2.pdf>

2. Disponible à l'adresse https://membres-ljk.imag.fr/Vincent.Brault/Cours/CoursStatistique_IF.pdf

**Preuve**

À partir de la proposition 21 qui dit que $\bar{\mathbf{X}}_n$ suit une loi gaussienne d'espérance μ_0 et de variance σ^2/n et du lemme 22, nous montrons que, sous l'hypothèse \mathcal{H}_0 qui suppose que $\mu^* = \mu_0$, T_n suit une loi gaussienne dont il faut calculer l'espérance et la variance :

$$\mathbb{E}[T_n] = \mathbb{E}\left[\sqrt{n}\frac{\bar{\mathbf{X}}_n - \mu_0}{\sigma}\right] = \sqrt{n}\frac{1}{\sigma} [\mathbb{E}[\bar{\mathbf{X}}_n] - \mu_0] = \sqrt{n}\frac{1}{\sigma} [\mu_0 - \mu_0] = 0.$$

Pour la variance, nous avons :

$$\begin{aligned} \mathbb{V}[T_n] &= \mathbb{V}\left[\sqrt{n}\frac{\bar{\mathbf{X}}_n - \mu_0}{\sigma}\right] \\ &= \frac{n}{\sigma^2} [\mathbb{V}[\bar{\mathbf{X}}_n] - \mathbb{V}[\mu_0]] \quad \text{car la variance est quadratique,} \\ &= \frac{n}{\sigma^2} [\sigma^2/n - 0] \quad \text{car la variance d'une constante est nulle,} \\ &= \frac{n\sigma^2}{n\sigma^2} = 1. \end{aligned}$$

Et la loi de T_n est une loi normale centrée réduite.

**Exercices 6.1**

Nous laissons en exercice le calcul de la puissance pour chaque règle de décision. En particulier, dans le cas où l'hypothèse \mathcal{H}_1 vaut $\mu^* < \mu_0$ ou $\mu^* > \mu_0$, montrer que la suite de tests est consistante.

6.2.2 Cas gaussien avec la variance inconnue

Dans le cas où la variance n'est pas connue, nous pouvons l'estimer. Dans ce cas, nous avons deux estimateurs possibles à savoir l'estimateur du maximum de vraisemblance ou basé sur les moments d'ordre 2

$$\hat{\sigma}_n^2 = \overline{X_n^2} - \bar{\mathbf{X}}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i\right)^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{\mathbf{X}}_n)^2$$

ou sa version non biaisée :

$$S_n^2 = \frac{n}{n-1} \hat{\sigma}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{\mathbf{X}}_n)^2.$$

**Exercices 6.2**

Montrer que $\hat{\sigma}_n^2$ est un estimateur biaisé de σ^2 alors que S_n^2 est sans biais.

Ainsi, nous allons étudier la statistique :

$$T_n = \sqrt{n}\frac{\bar{\mathbf{X}}_n - \mu^*}{S_n}$$

dont nous devons trouver la loi :

Proposition 23 (Loi de la variable S_n)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de la loi gaussienne $\mathcal{N}(\mu^*, \sigma^2)$ avec σ^2 inconnue alors la variable $(n-1)S_n^2/\sigma^2$ suit une loi du $\chi^2(n-1)$; c'est-à-dire une loi du χ^2 à $n-1$ degrés de libertés.

$$(n-1)S_n^2/\sigma^2 = \sum_{i=1}^n \frac{(X_i - \bar{\mathbf{X}}_n)^2}{\sigma^2} \sim \chi^2(n-1).$$



Preuve

La preuve se base sur un théorème que nous admettrons car il utilise des notions de projections (donc d'espérances conditionnelles) qui ne seront pas abordées (les lectrices intéressé·e·s peuvent aller regarder le polycopié de Le Gall (2006) pour avoir la démonstration). Nous mettons toutefois l'énoncé ici.

Lemme 24 (Théorème de Cochran)

Étant donné un vecteur gaussien \mathbf{X} de \mathbb{R}^d de loi $\mathcal{N}(m, \sigma^2 \mathbb{I}_d)$ et F_1, \dots, F_P des sous-espaces vectoriels de \mathbb{R}^d , orthogonaux deux à deux et de somme \mathbb{R}^d dont nous notons P_{F_p} les matrices de projection orthogonale sur F_p et n_p la dimension de F_p alors nous avons :

- les vecteurs aléatoires $P_{F_1} \mathbf{X}, \dots, P_{F_P} \mathbf{X}$ sont deux à deux indépendants et de lois respectives $\mathcal{N}(P_{F_1} m, \sigma^2 P_{F_1}), \dots, \mathcal{N}(P_{F_P} m, \sigma^2 P_{F_P})$.
- les variables aléatoires $\frac{\|P_{F_1}(\mathbf{X}-m)\|^2}{\sigma^2}, \dots, \frac{\|P_{F_P}(\mathbf{X}-m)\|^2}{\sigma^2}$ sont deux à deux indépendantes et sont de lois respectives $\chi^2(n_1), \dots, \chi^2(n_P)$.

Pour démontrer la proposition 23, nous commençons par définir pour tout $i \in \{1, \dots, n\}$, la variable :

$$Z_i = \frac{X_i - \bar{\mathbf{X}}_n}{\sigma}.$$

Calculons l'espérance et la matrice de variance-covariance de la variable $\mathbf{Z} = (Z_1, \dots, Z_n)$. Sous l'hypothèse \mathcal{H}_0 , nous avons de plus pour tout $i \in \{1, \dots, n\}$:

$$\begin{aligned} \mathbb{E}[Z_i] &= \mathbb{E}\left[\frac{X_i - \bar{\mathbf{X}}_n}{\sigma}\right] \\ &= \frac{\mathbb{E}[X_i] - \mathbb{E}[\bar{\mathbf{X}}_n]}{\sigma} \\ &= \frac{\mu^* - \mu^*}{\sigma} \\ &= 0. \end{aligned}$$

De plus, pour tout $i \in \{1, \dots, n\}$, nous avons :

$$\begin{aligned} \mathbb{V}[Z_i] &= \mathbb{V}\left[\frac{X_i - \bar{\mathbf{X}}_n}{\sigma}\right] \\ &= \frac{1}{\sigma^2} \mathbb{V}[X_i - \bar{\mathbf{X}}_n] \\ &= \frac{1}{\sigma^2} \{ \mathbb{V}[X_i] + \mathbb{V}[\bar{\mathbf{X}}_n] - 2\text{Cov}(X_i, \bar{\mathbf{X}}_n) \} \\ &= \frac{1}{\sigma^2} \left\{ \sigma^2 + \frac{\sigma^2}{n} - 2\text{Cov}\left(X_i, \frac{1}{n} \sum_{i'=1}^n X_{i'}\right) \right\} \\ &= \frac{n+1}{n} - \frac{2}{n\sigma^2} \sum_{i'=1}^n \text{Cov}(X_i, X_{i'}) \\ &= \frac{n+1}{n} - \frac{2}{n\sigma^2} \left[\sum_{\substack{i'=1 \\ i' \neq i}}^n \underbrace{\text{Cov}(X_i, X_{i'})}_{=0} + \text{Cov}(X_i, X_i) \right] \\ &= \frac{n+1}{n} - \frac{2}{n\sigma^2} \mathbb{V}[X_i] \\ &= \frac{n-1}{n}. \end{aligned}$$

Donc la loi commune des Z_i est centrée et matrice de variance covariance :

$$\Gamma_n = \mathbb{I}_n - \begin{pmatrix} 1/\sqrt{n} \\ \vdots \\ 1/\sqrt{n} \end{pmatrix} (1/\sqrt{n} \quad \dots \quad 1/\sqrt{n}).$$

La loi de la variable $\frac{Z_1 + \dots + Z_n}{\sqrt{n}}$ est une loi normale centrée de covariance Γ_n . Or, cette loi est le projeté d'un vecteur aléatoire de \mathbb{R}^K suivant une loi normale centrée réduite sur l'hyperplan orthogonal au vecteur colonne $(1/\sqrt{n}, \dots, 1/\sqrt{n})^T$.

De plus, nous voyons que la variable $(n-1)S_n^2/\sigma^2$ est la norme au carré de cette variable donc, d'après le théorème de Cochran, la loi asymptotique est une loi du χ^2 à $n-1$ degrés de liberté.

Ainsi, nous obtenons le corollaire suivant :

Corollaire 25 (Loi de la statistique)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de la loi gaussienne $\mathcal{N}(\mu^*, \sigma^2)$ avec σ^2 inconnue alors la loi de

$$T_n = \sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{S_n}$$

suit une loi de Student centrée à $n-1$ degrés de liberté :

$$T_n \sim \mathcal{T}(n-1).$$

Preuve

Avant de faire la démonstration, rappelons la définition d'une loi de Student :

Définition 55 (Loi de Student)

Soient Z une variable de loi gaussienne centrée réduite et U une variable aléatoire indépendante de Z et distribuée suivant une loi du χ^2 à k degrés de liberté alors la variable :

$$\frac{Z}{\sqrt{U/k}}$$

suit une loi de Student centrée à k degrés de liberté $\mathcal{T}(k)$.

Nous avons :

$$\begin{aligned} \sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{\sqrt{S_n^2}} &= \frac{\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{\sigma}}{\frac{\sqrt{S_n^2}}{\sigma}} \\ &= \frac{\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{\sigma}}{\sqrt{\frac{(n-1)S_n^2}{(n-1)\sigma^2}}} \\ &= \frac{Z}{\sqrt{\frac{Y}{n-1}}} \end{aligned}$$

avec Z qui suit une loi normale centrée réduite et Y qui suit une loi du χ^2 à $n-1$ degrés de liberté indépendante d'après le théorème 24 de Cochran donc, par définition de la loi de Student, nous avons le résultat.

Ainsi, nous avons :

Point méthode (Test d'espérance ; cas gaussien, variance inconnue)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi gaussienne $\mathcal{N}(\mu^*, \sigma^2)$ avec σ^2 inconnue alors la statistique de test est

$$T_n = \sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu_0}{S_n}$$

qui a pour loi une loi de Student centrée à $n-1$ degrés de liberté sous l'hypothèse \mathcal{H}_0 et la


règle de rejet de niveau exact $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\mu^* \neq \mu_0$	$\mu^* < \mu_0$	$\mu^* > \mu_0$
Θ_1	$\mathbb{R} \setminus \{\mu_0\}$	$] -\infty, \mu_0[$	$] \mu_0; +\infty[$
Règle	$ T_n > s_{1-\alpha/2}^{(n-1)}$	$T_n < s_{\alpha}^{(n-1)}$	$T_n > s_{1-\alpha}^{(n-1)}$

où $s_{\alpha}^{(n-1)}$ est le quantile d'ordre α d'une loi de Student centrée à $n - 1$ degrés de liberté; c'est-à-dire que si S suit une loi de Student centrée à $n - 1$ degrés de liberté $\mathcal{T}(n - 1)$ alors $s_{\alpha}^{(n-1)}$ est tel que :

$$\mathbb{P}(S \leq s_{\alpha}^{(n-1)}) = \alpha.$$

Codage en R

La fonction  permettant de faire le test bilatéral est `t.test`. Par exemple :

```
1 ## Simulation d'un jeu de données
2 X <- rnorm(100)
3 ## Test sous H_0 : mu=0 -> acceptation avec 95% de chances
4 t.test(X,mu=0)
5 ## Test sous H_0 : mu=1 -> rejet normalement
6 t.test(X,mu=1)
```

Remarque

Nous pouvons aussi utiliser $\hat{\sigma}_n$ à la place de S_n dans le calcul de la statistique de T_n mais, dans ce cas, le test est de niveau asymptotique seulement.

Point méthode (Test d'espérance ; cas gaussien, variance inconnue avec $\hat{\sigma}_n$)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi gaussienne $\mathcal{N}(\mu^*, \sigma^2)$ avec σ^2 inconnue alors la statistique de test est

$$T_n = \sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu_0}{\hat{\sigma}_n}$$

qui a pour loi **asymptotique** une loi de Student centrée à $n - 1$ degrés de liberté sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau **asymptotique** exact $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\mu^* \neq \mu_0$	$\mu^* < \mu_0$	$\mu^* > \mu_0$
Θ_1	$\mathbb{R} \setminus \{\mu_0\}$	$] -\infty, \mu_0[$	$] \mu_0; +\infty[$
Règle	$ T_n > s_{1-\alpha/2}^{(n-1)}$	$T_n < s_{\alpha}^{(n-1)}$	$T_n > s_{1-\alpha}^{(n-1)}$

où $s_{\alpha}^{(n-1)}$ est le quantile d'ordre α d'une loi de Student centrée à $n - 1$ degrés de liberté; c'est-à-dire que si S suit une loi de Student centrée à $n - 1$ degrés de liberté $\mathcal{T}(n - 1)$ alors $s_{\alpha}^{(n-1)}$ est tel que :

$$\mathbb{P}(S \leq s_{\alpha}^{(n-1)}) = \alpha.$$

Preuve

Par la même démonstration que pour la proposition 23, nous montrons que la variable $n\hat{\sigma}_n^2/\sigma^2$ suit une loi du $\chi^2(n - 1)$; c'est-à-dire une loi du χ^2 à $n - 1$ degrés de liberté. Ainsi, nous avons :

$$\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{\sqrt{\hat{\sigma}_n^2}} = \frac{\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{\sigma}}{\frac{\sqrt{\hat{\sigma}_n^2}}{\sigma}}$$



$$\begin{aligned}
&= \frac{\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{\sigma}}{\frac{1}{\sqrt{n}} \sqrt{\frac{n \hat{\sigma}_n^2}{\sigma^2}}} \\
&= \frac{1}{\sqrt{\frac{n}{n-1}}} \times \frac{\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{\sigma}}{\sqrt{\frac{n \hat{\sigma}_n^2}{\sigma^2 (n-1)}}} \\
&= \underbrace{\frac{1}{\sqrt{\frac{n-1}{n}}}}_{\xrightarrow[n \rightarrow +\infty]{} 1}} \times \frac{Z}{\sqrt{\frac{Y}{n-1}}}
\end{aligned}$$

avec Z qui suit une loi normale centrée réduite et Y qui suit une loi du χ^2 à $n - 1$ degrés de liberté indépendante d'après le théorème 24 de Cochran et comme asymptotiquement $n/(n - 1)$ tend vers 1 alors, asymptotiquement et sous l'hypothèse que $\mu^* = \mu_0$, T_n suit une loi de Student centrée à $n - 1$ degrés de liberté.

6.2.3 Test de la normalité

Comme vous pouvez le voir dans les sections 6.2.1 et 6.2.2, nous avons supposé que les variables étaient gaussiennes. Dans la section 6.2.4 suivante, nous verrons comment faire lorsque les variables ne semblent pas gaussiennes. Toutefois, la question qui se pose alors est de savoir comment tester si une variable est gaussienne ou non. Pour cela, il existe différents tests (voir le chapitre 8), nous présentons les principes de deux d'entre eux. En particulier, nous faisons le test suivant :

$$\begin{cases} \mathcal{H}_0 : \mathbf{X} \text{ suit une loi normale,} \\ \mathcal{H}_1 : \mathbf{X} \text{ ne suit pas une loi normale.} \end{cases} \quad (6.1)$$

Remarque

Remarquons que le test n'empêche pas les corrélations entre les observations.

Test de Shapiro-Wilk (voir Shapiro et Wilk (1965)) : Le test de Shapiro-Wilk repose sur la même idée que celui des graphiques *QQ-plot* c'est-à-dire au fait d'ordonner les valeurs de la plus petite à la plus grande et de les comparer aux quantiles d'une loi gaussienne centrée réduite. Concrètement, étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$, nous organisons les variables de la plus petite à la plus grande $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ et nous les comparons aux espérances des quantiles d'une loi normale ayant les mêmes paramètres. Autrement dit, si nous notons $\mathbf{m} = (m_1, \dots, m_n)$ les espérances des statistiques d'ordre d'un n -échantillon suivant une loi normale et $\hat{\Sigma}_n$ la matrice de variance covariance estimée des observations :

$$\forall (i, j) \in \{1, \dots, n\}^2, \quad (\hat{\Sigma}_n)_{i,j} = \text{Cov}_n(X_i, X_j),$$

nous comparons les observations aux constantes $\mathbf{a} = (a_1, \dots, a_n)$ définies par :

$$\mathbf{a} = \frac{\mathbf{m}^T \hat{\Sigma}_n^{-1}}{\left[\mathbf{m}^T \hat{\Sigma}_n^{-1} \hat{\Sigma}_n^{-1} \mathbf{m} \right]^{1/2}}.$$

Ainsi, nous avons la méthode suivante :

Point méthode (Test de normalité ; méthode de Shapiro et Wilk (1965))

Étant données des variables aléatoires $\mathbf{X} = (X_1, \dots, X_n)$ de la loi gaussienne multivariée $\mathcal{N}_n(\boldsymbol{\mu}^*, \Sigma_n^2)$ alors la statistique


$$W_n = \frac{\left[\sum_{i=1}^n a_i X_{(i)} \right]^2}{\sum_{i=1}^n (X_i - \bar{\mathbf{X}}_n)^2}$$

suit une loi qui ne porte pas de nom mais dont les valeurs peuvent être calculées à l'aide de simulations de Monte-Carlo. La règle de décision du test (6.1) est alors :

$$W_n > w_{1-\alpha}^n$$

où $w_{1-\alpha}^n$ est le quantile d'ordre α estimé par simulation de Monte-Carlo (voir par exemple Royston (1982)).

Codage en R

La fonction  permettant de faire le test de normalité en utilisant la méthode de Shapiro-Wilk est `shapiro.test`. Par exemple :

```
1 ## Acceptation avec 95% de chances
2 shapiro.test(rnorm(100))
3 ## Rejet normalement
4 shapiro.test(rexp(100))
```

Remarque

Comme nous l'avons dit précédemment, ce test est à mettre en lien avec les graphiques *QQ*-plot. Sur la figure 6.1, nous avons représenté les *QQ*-plot pour 4 types de distributions différentes. Les deux du haut sont des distributions gaussiennes (indépendantes à gauche et corrélées à droite) ; les *p*-valeurs ne sont pas trop petites donc on va avoir tendance à conserver, à raison, l'hypothèse de normalité. Les deux distributions sont respectivement une loi exponentielle $\mathcal{E}(1)$ et une loi Cauchy $\mathcal{Cau}(0,1)$; dans ces deux cas, l'hypothèse \mathcal{H}_0 est largement rejetée, à raison, au profit de l'hypothèse \mathcal{H}_1 .

Codage en R

Le code utilisé pour créer la figure 6.1 est le suivant :

```
1 #####
2 # Packages
3 #####
4 library(ggpubr) # Pour qqplot
5 library(grid) # Pour affichage multiple
6
7 #####
8 # Préparation pour représenter les 4 graphiques
9 #####
10 grid.newpage()
11 pushViewport(viewport(layout = grid.layout(2,2)))
12
13 #####
14 # Simulation gaussienne iid
15 ###
16 n<-100
17 dat<-data.frame(x=rnorm(n))
18 ## Affichage
19 q<-ggqqplot(dat, x = "x",
20             ggtheme = theme_pubclean()+
21             ggtitle("Variables gaussiennes iid", subtitle = paste0("p-value=", signif(shapiro.test(dat$x)$p.value),2))
22             print(q, vp=viewport(layout.pos.row = 1, layout.pos.col = 1))
23
24 #####
25 # Simulation gaussienne corrélée
26 #####
27 ## Triangulaire inférieure
28 P<-matrix(runif(n*n)*2-1,n,n)
29 P[abs(P)<0.3]<-0
```

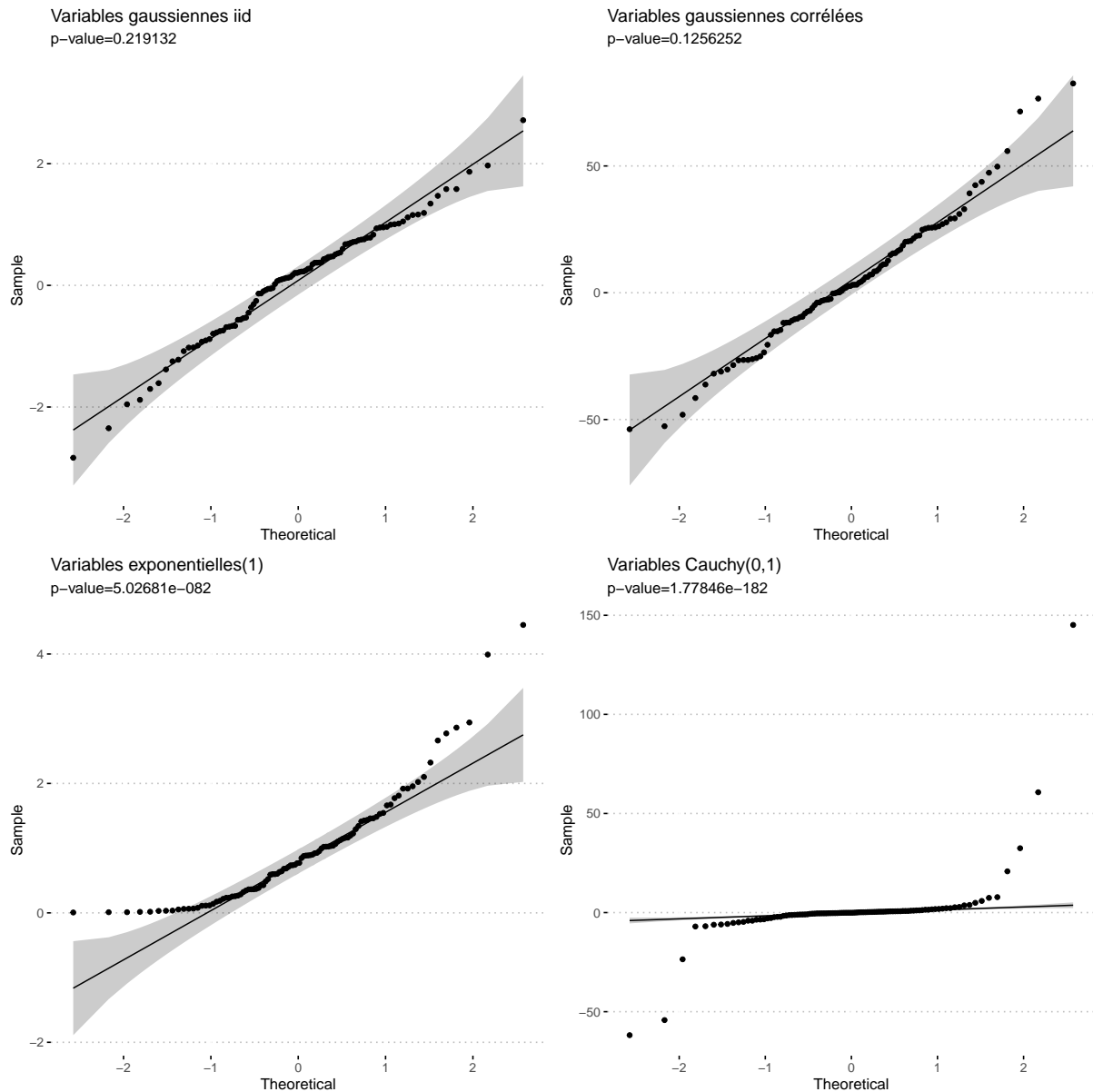


FIGURE 6.1 – QQ -plot de plusieurs 100-échantillons suivant différentes lois : loi gaussienne indépendante (en haut à gauche), loi gaussienne corrélée (en haut à droite), loi exponentielle $\mathcal{E}(1)$ (en bas à gauche) et loi Cauchy $\mathcal{Cau}(0, 1)$ (en bas à droite). La p -valeur obtenue par le test de Shapiro-Wilk a été ajouté à chaque fois dans le sous titre.

```

31 for (i in 1:(n-1)){
32   for (j in (i+1):n){
33     P[i,j]=0
34   }
35 }
36 ## Diagonale forte pour éviter les valeurs nulles
37 diag(P)<-runif(n)+0.5
38 Sigma<-P%*%t(P)
39 ## Vérification de la matrice
40 sum(abs(Sigma-t(Sigma))) ## Symétrique
41 min(eigen(Sigma)$values) ## Définie positive
42 ## Simulation
43 dat<-data.frame(x=Sigma%*%rnorm(n)+1)
44 ## Affichage

```

```

45 q<-ggqqplot(dat, x = "x",
46           ggtheme = theme_pubclean()+
47           ggtitle("Variables gaussiennes corrélées",subtitle = paste0("p-value=",signif(shapiro.test(dat$x)$p.value),2))
48 print(q,vp=viewport(layout.pos.row = 1, layout.pos.col = 2))
49
50 #####
51 # Simulation loi exponentielle
52 #####
53 ## Simulation
54 dat<-data.frame(x=rexp(n))
55 ## Affichage
56 q<-ggqqplot(dat, x = "x",
57           ggtheme = theme_pubclean()+
58           ggtitle("Variables exponentielles(1)",subtitle = paste0("p-value=",signif(shapiro.test(dat$x)$p.value),2))
59 print(q,vp=viewport(layout.pos.row = 2, layout.pos.col = 1))
60
61 #####
62 # Simulation loi Cauchy
63 #####
64 ## Simulation
65 dat<-data.frame(x=rcauchy(n))
66 ## Affichage
67 q<-ggqqplot(dat, x = "x",
68           ggtheme = theme_pubclean()+
69           ggtitle("Variables Cauchy(0,1)",subtitle = paste0("p-value=",signif(shapiro.test(dat$x)$p.value),2))
70 print(q,vp=viewport(layout.pos.row = 2, layout.pos.col = 2))

```

Test de Lilliefors : Le test de Lilliefors est une adaptation du test de Kolmogorov-Smirnov qui est développé dans la section 8.1 sur les tests non paramétriques. Avant ça, nous avons besoin de redéfinir la fonction de répartition empirique (dont la définition a été donnée en section 19) :

Définition 56 (Fonction de répartition empirique)

Étant donné un n -échantillon \mathbf{X} de variables réelles, nous appelons **fonction de répartition empirique** la fonction définie sur \mathbb{R} par :

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \leq x\}}.$$

Ainsi, nous avons la méthode suivante :

Point méthode (Test de normalité ; méthode de Lilliefors (1967))

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de la loi gaussienne $\mathcal{N}(\mu^*, \sigma^2)$ alors nous faisons la procédure suivante :

1. Calcul de $\bar{\mathbf{X}}_n$ et de $\hat{\sigma}_n^2$ les estimateurs de l'espérance et de la variance de \mathbf{X} .
2. Calcul du maximum de la dispersion entre la fonction de répartition empirique F_n de l'échantillon \mathbf{X} et la fonction de répartition $F_{\mathcal{N}(\bar{\mathbf{X}}_n, \hat{\sigma}_n^2)}$ d'une loi gaussienne ayant les paramètres estimés en 1 :


$$L_n = \sqrt{n} \sup_{x \in \mathbb{R}} \left| F_{\mathcal{N}(\bar{\mathbf{X}}_n, \hat{\sigma}_n^2)}(x) - F_n(x) \right|$$

alors la statistique L_n suit une distribution de Lilliefors qui se calcule aussi par une méthode de Monte-Carlo. La règle de décision du test (6.1) est alors :

$$L_n > \ell_{1-\alpha}^n$$

où $\ell_{1-\alpha}^n$ est le quantile d'ordre α d'une distribution de Lilliefors estimé par simulation de Monte-Carlo.

Codage en R

La fonction  permettant de faire le test de normalité en utilisant la méthode de Lilliefors est `lillie.test` du package `nortest`. Par exemple :

```
1 ## Chargement du package
2 library(nortest)
3 ## Acceptation avec 95% de chances
4 lillie.test(rnorm(100))
5 ## Rejet normalement
6 lillie.test(rexp(100))
```

Remarque

Comme nous l'avons dit précédemment, ce test est basé sur les fonctions de répartition. Sur la figure 6.2, nous avons représenté les fonctions de répartition empirique et les fonctions de répartition théoriques d'une loi gaussienne $\mathcal{N}(\bar{\mathbf{X}}_n, \hat{\sigma}_n^2)$ pour 4 types de distributions différentes. Les deux du haut sont des distributions gaussiennes (indépendantes à gauche et corrélées à droite) ; les p -valeurs ne sont pas trop petites donc on va avoir tendance à conserver, à raison, l'hypothèse de normalité. Les deux distributions sont respectivement une loi exponentielle $\mathcal{E}(1)$ et une loi Cauchy $\mathcal{Cau}(0, 1)$; dans ces deux cas, l'hypothèse \mathcal{H}_0 est largement rejetée, à raison, au profit de l'hypothèse \mathcal{H}_1 .

Codage en R

Le code utilisé pour créer la figure 6.1 est le suivant :

```
1 #####
2 # Packages
3 #####
4 library(ggplot2) # Pour qqplot
5 library(grid) # Pour affichage multiple
6 library(nortest) # Pour le test
7
8 #####
9 # Préparation pour représenter les 4 graphiques
10 #####
11 grid.newpage()
12 pushViewport(viewport(layout = grid.layout(2,2)))
13
14 ###
15 # Simulation gaussienne iid
16 ###
17 n <- 100
```

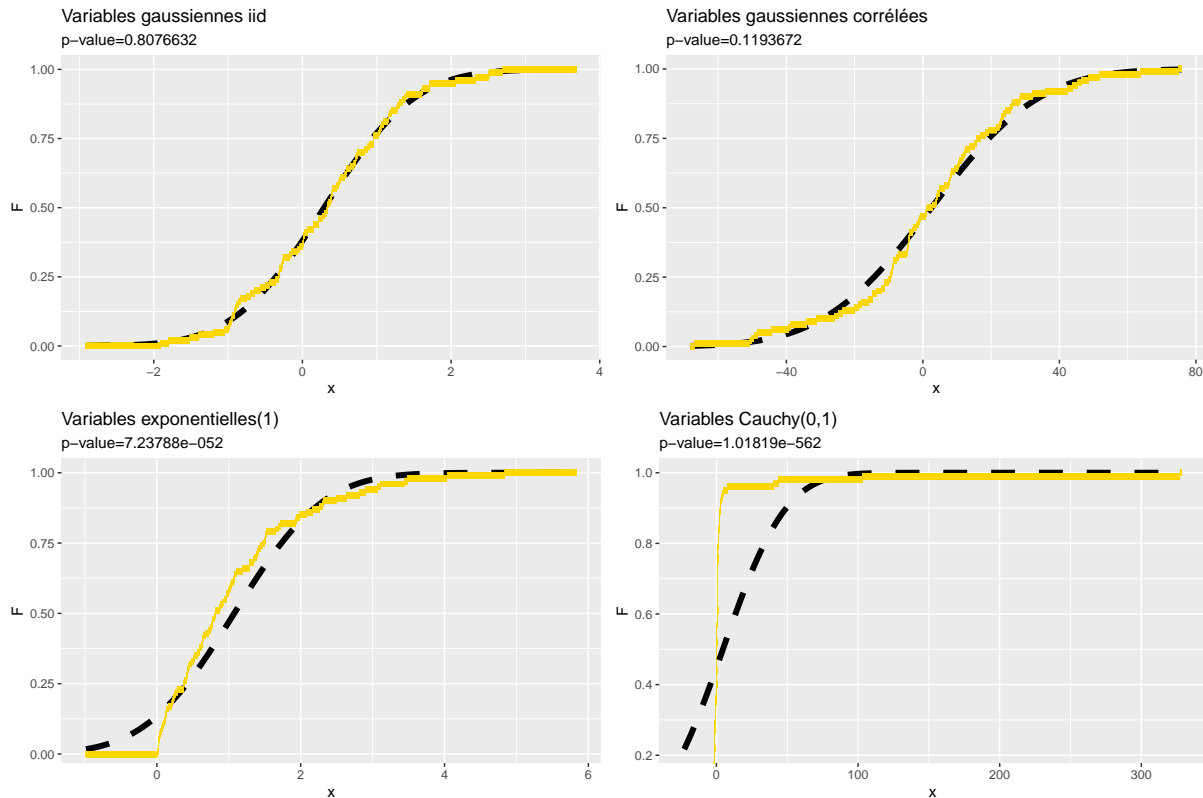


FIGURE 6.2 – Fonctions de répartition empirique (en or) et les fonctions de répartition théoriques (en noir) d'une loi gaussienne $\mathcal{N}(\bar{X}_n, \hat{\sigma}_n^2)$ de plusieurs 100-échantillons suivant différentes lois : loi gaussienne indépendante (en haut à gauche), loi gaussienne corrélée (en haut à droite), loi exponentielle $\mathcal{E}(1)$ (en bas à gauche) et loi Cauchy $\mathcal{Cau}(0, 1)$ (en bas à droite). La p -valeur obtenue par le test de Shapiro-Wilk a été ajoutée à chaque fois dans le sous titre.

```

18 X<-rnorm(n)
19 ### Fonction de répartition de référence
20 x_ref<-seq(min(X)-1,max(X)+1,by=0.1)
21 dat<-data.frame(x=x_ref,F=pnorm(x_ref,mean = mean(X),sd=sd(X)))
22 ## Affichage
23 q<-ggplot(dat, aes(x = x,y=F))+geom_line(lwd=2,lty="dashed")+
24   ggtitle("Variables gaussiennes iid",subtitle = paste0("p-value=",signif(lillie.test(X)$p.value),2))
25 ## Ajout fonction de répartition empirique
26 X_ord<-c(x_ref[1],sort(X),x_ref[length(x_ref)])
27 for (it in 1:(length(X_ord)-1)){
28   q<-q+geom_segment(x = X_ord[it],xend=X_ord[it+1],
29     y=(it-1)/n,yend=(it-1)/n,col="gold",alpha=0.5,lwd=2)
30 }
31 ## Affichage
32 print(q,vp=viewport(layout.pos.row = 1, layout.pos.col = 1))
33
34 #####
35 # Simulation gaussienne corrélée
36 #####
37 ## Triangulaire inférieure
38 P<-matrix(runif(n*n)*2-1,n,n)
39 P[abs(P)<0.3]<-0
40 for (i in 1:(n-1)){
41   for (j in (i+1):n){
42     P[i,j]=0
43   }
44 }

```

```

45 ## Diagonale forte pour éviter les valeurs nulles
46 diag(P)<-runif(n)+0.5
47 Sigma<-P%*%t(P)
48 ## Vérification de la matrice
49 sum(abs(Sigma-t(Sigma))) ## Symétrique
50 min(eigen(Sigma)$values) ## Définie positive
51 ## Simulation
52 X<-Sigma%*%rnorm(n)+1
53 ### Fonction de répartition de référence
54 x_ref<-seq(min(X)-1,max(X)+1,by=0.1)
55 dat<-data.frame(x=x_ref,F=pnorm(x_ref,mean = mean(X),sd=sd(X)))
56 ## Affichage
57 q<-ggplot(dat, aes(x = x,y=F))+geom_line(lwd=2,lty="dashed")+
58   ggtitle("Variables gaussiennes corrélées",subtitle = paste0("p-value=",signif(lillie.test(X)$p.value),2))
59 ## Ajout fonction de répartition empirique
60 X_ord<-c(x_ref[1],sort(X),x_ref[length(x_ref)])
61 for (it in 1:(length(X_ord)-1)){
62   q<-q+geom_segment(x = X_ord[it],xend=X_ord[it+1],
63     y=(it-1)/n,yend=(it-1)/n,col="gold",alpha=0.5,lwd=2)
64 }
65 ## Affichage
66 print(q,viewport(layout.pos.row = 1, layout.pos.col = 2))
67
68 #####
69 # Simulation loi exponentielle
70 #####
71 ## Simulation
72 X<-rexp(n)
73 ### Fonction de répartition de référence
74 x_ref<-seq(min(X)-1,max(X)+1,by=0.1)
75 dat<-data.frame(x=x_ref,F=pnorm(x_ref,mean = mean(X),sd=sd(X)))
76 ## Affichage
77 q<-ggplot(dat, aes(x = x,y=F))+geom_line(lwd=2,lty="dashed")+
78   ggtitle("Variables exponentielles(1)",subtitle = paste0("p-value=",signif(lillie.test(X)$p.value),2))
79 ## Ajout fonction de répartition empirique
80 X_ord<-c(x_ref[1],sort(X),x_ref[length(x_ref)])
81 for (it in 1:(length(X_ord)-1)){
82   q<-q+geom_segment(x = X_ord[it],xend=X_ord[it+1],
83     y=(it-1)/n,yend=(it-1)/n,col="gold",alpha=0.5,lwd=2)
84 }
85 ## Affichage
86 print(q,viewport(layout.pos.row = 2, layout.pos.col = 1))
87
88 #####
89 # Simulation loi Cauchy
90 #####
91 ## Simulation
92 X<-rcauchy(n)
93 ### Fonction de répartition de référence
94 x_ref<-seq(min(X)-1,max(X)+1,by=0.1)
95 dat<-data.frame(x=x_ref,F=pnorm(x_ref,mean = mean(X),sd=sd(X)))
96 ## Affichage
97 q<-ggplot(dat, aes(x = x,y=F))+geom_line(lwd=2,lty="dashed")+
98   ggtitle("Variables Cauchy(0,1)",subtitle = paste0("p-value=",signif(lillie.test(X)$p.value),2))
99 ## Ajout fonction de répartition empirique
100 X_ord<-c(x_ref[1],sort(X),x_ref[length(x_ref)])
101 for (it in 1:(length(X_ord)-1)){
102   q<-q+geom_segment(x = X_ord[it],xend=X_ord[it+1],
103     y=(it-1)/n,yend=(it-1)/n,col="gold",alpha=0.5,lwd=2)
104 }
105 ## Affichage
106 print(q,viewport(layout.pos.row = 2, layout.pos.col = 2))

```

6.2.4 Cas non gaussien

Enfin, si le test de normalité rejette l'hypothèse que la distribution soit gaussienne au profit d'une distribution gaussienne mais que nous admettons que la distribution ait un moment d'ordre 2 alors nous pouvons réutiliser la statistique de la section 6.2.2 mais avec une loi asymptotique :

Point méthode (Test d'espérance ; cas non gaussien, variance inconnue avec S_n)

Étant donné un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi admettant un moment d'ordre 2 alors la statistique de test est

$$T_n = \sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu_0}{S_n}$$

qui a pour loi **asymptotique** une loi de Student centrée à $n - 1$ degrés de liberté sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau **asymptotique** exact $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\mu^* \neq \mu_0$	$\mu^* < \mu_0$	$\mu^* > \mu_0$
Θ_1	$\mathbb{R} \setminus \{\mu_0\}$	$] -\infty, \mu_0[$	$] \mu_0; +\infty[$
Règle	$ T_n > s_{1-\alpha/2}^{(n-1)}$	$T_n < s_{\alpha}^{(n-1)}$	$T_n > s_{1-\alpha}^{(n-1)}$

où $s_{\alpha}^{(n-1)}$ est le quantile d'ordre α d'une loi de Student centrée à $n - 1$ degrés de liberté ; c'est-à-dire que si S suit une loi gde Student centrée à $n - 1$ degrés de liberté $\mathcal{T}(n - 1)$ alors $s_{\alpha}^{(n-1)}$ est tel que :

$$\mathbb{P}(S \leq s_{\alpha}^{(n-1)}) = \alpha.$$



Preuve

D'après le théorème de la limite centrale et comme la distribution admet un moment d'ordre 2, la loi de $\bar{\mathbf{X}}_n$ converge vers une loi gaussienne d'espérance μ^* et de variance inconnue que nous noterons σ^2 . Par la même démonstration que pour la proposition 23, nous montrons que la variable $n\hat{\sigma}_n^2/\sigma^2$ suit une loi asymptotique du $\chi^2(n - 1)$; c'est-à-dire une loi du χ^2 à $n - 1$ degrés de liberté. Ainsi, nous avons :

$$\begin{aligned} \sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{\sqrt{\hat{\sigma}_n^2}} &= \frac{\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{\sigma}}{\frac{\sqrt{\hat{\sigma}_n^2}}{\sigma}} \\ &= \frac{\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu^*}{\sigma}}{\sqrt{\frac{(n-1)\hat{\sigma}_n^2}{\sigma^2(n-1)}}} \\ &= \times \frac{Z}{\sqrt{\frac{Y}{n-1}}} \end{aligned}$$

avec Z qui suit une loi asymptotique normale centrée réduite et Y qui suit une loi asymptotique du χ^2 à $n - 1$ degrés de liberté indépendante d'après le théorème 24 de Cochran alors, asymptotiquement et sous l'hypothèse que $\mu^* = \mu_0$, T_n suit une loi de Student centrée à $n - 1$ degrés de liberté.

6.2.5 Quel(s) test(s) utilisé(s) ?

Pour faire le bilan de cette partie, nous avons créé un arbre de décision pour savoir quel(s) test(s) utilisé(s) disponible sur la figure 6.3.

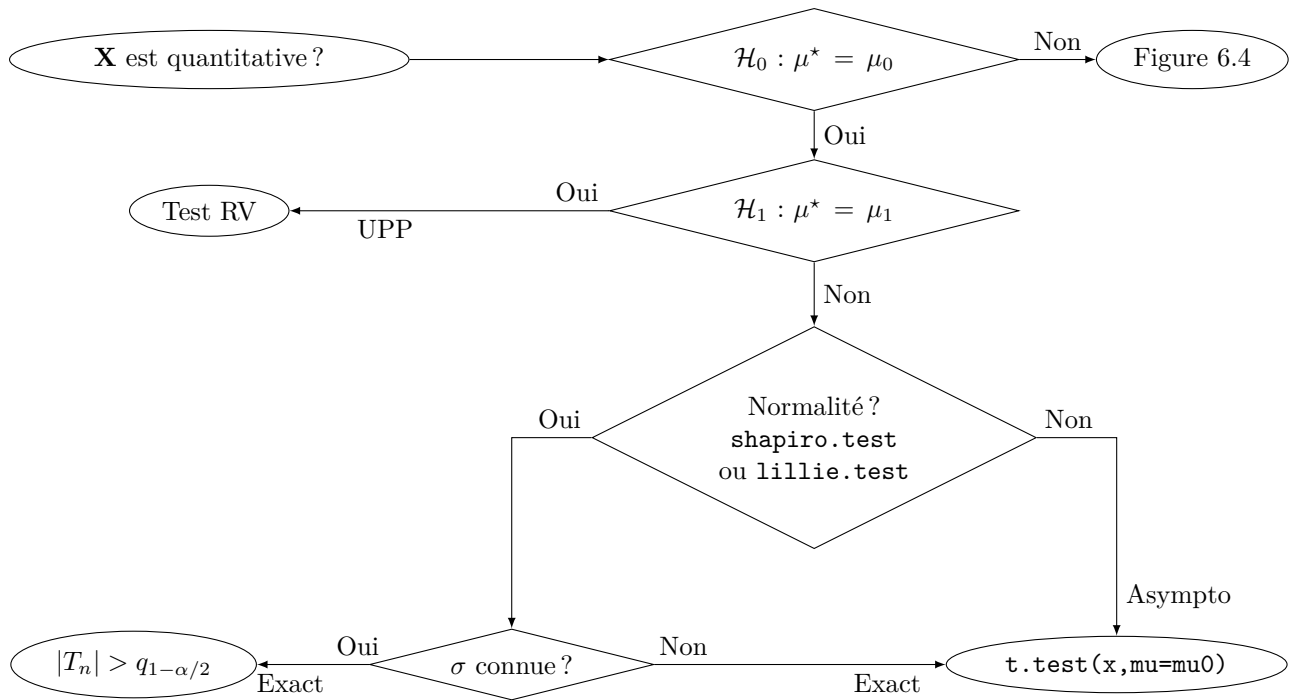


FIGURE 6.3 – Représentation de l’arbre de décision pour savoir quel test utilisé pour tester si l’espérance est égale à une valeur fixée en amont.

6.3 Test de comparaison des égalités d’espérances de deux échantillons

Dans cette partie, nous supposons que nous avons au moins deux échantillons $\mathbf{X} = (X_1, \dots, X_{n_{\mathbf{X}}})$ et $\mathbf{Y} = (Y_1, \dots, Y_{n_{\mathbf{Y}}})$ dont les lois associées possèdent une espérance $\mu_{\mathbf{X}}^*$ et $\mu_{\mathbf{Y}}^*$ et nous regardons si ces espérances sont égales ou pas. Dans ce cas, nous étudions donc le test paramétrique suivant :

$$\begin{cases} \mathcal{H}_0 : \mu_{\mathbf{X}}^* = \mu_{\mathbf{Y}}^*, \\ \mathcal{H}_1 : \mu_{\mathbf{X}}^* \neq \mu_{\mathbf{Y}}^*. \end{cases} \quad (6.2)$$

Hypothèse (Test paramétrique et espérance)

Pour enfoncer le clou, nous insistons sur le fait que nous ne pouvons comparer les espérances qu’à condition que les deux lois possèdent une espérance.

6.3.1 Cas d’échantillons appariés

Nous commençons par définir ce que sont les échantillons appariés :

Définition 57 (Échantillons appariés)

Nous disons que deux échantillons \mathbf{X} et \mathbf{Y} sont **appariés** si pour tout $i \in \{1, \dots, n\}$, X_i et Y_i appartiennent au même individu. Dans ce cas, nous regarderons la loi de la paire de variables $(X_i, Y_i)_{1 \leq i \leq n}$.

Nous avons le corollaire suivant :

Corollaire 26

Si nous supposons que les échantillons \mathbf{X} et \mathbf{Y} sont appariés alors $n_{\mathbf{X}} = n_{\mathbf{Y}} = n$.

Pour ce faire, nous regardons pour tout $i \in \{1, \dots, n\}$, nous regardons la différence $D_i = X_i - Y_i$ de chaque paire et le test (6.2) devient alors le test suivant :

$$\begin{cases} \mathcal{H}_0 : \mu_{\mathbf{D}}^* = 0, \\ \mathcal{H}_1 : \mu_{\mathbf{D}}^* \neq 0. \end{cases}$$

Ainsi, nous pouvons reprendre les tests présentés dans la section 6.2 c'est-à-dire en adaptant la statistique T_n présentée dans le corollaire 25 :

Définition 58 (Test d'égalité d'espérance ; statistique pour échantillons appariés)

Pour tester deux échantillons appariés (\mathbf{X}, \mathbf{Y}) , nous utilisons la statistique de test suivante :

$$T_n = \sqrt{n} \frac{\bar{\mathbf{D}}_n}{S_n^{(\mathbf{D})}}$$

avec $\bar{\mathbf{D}}_n$ la moyenne empirique du n -échantillon $\mathbf{D} = (D_1, \dots, D_n)$ et $S_n^{(\mathbf{D})} = \frac{1}{n-1} \sum_{i=1}^n (D_i - \bar{\mathbf{D}}_n)^2$ l'estimateur non biaisée de la variance.

Cas gaussien : Dans le cas où les deux échantillons sont gaussiens, par le lemme 22, l'échantillon \mathbf{D} est donc gaussien et nous avons la loi suivante :

Point méthode (Test d'égalité des espérances ; cas gaussien et échantillons appariés)

Étant donnés deux n -échantillons appariés $(\mathbf{X}, \mathbf{Y}) = [(X_1, Y_1), \dots, (X_n, Y_n)]$ de loi gaussienne $\mathcal{N} \left[\begin{pmatrix} \mu_{\mathbf{X}}^* \\ \mu_{\mathbf{Y}}^* \end{pmatrix}, \Sigma_2 \right]$ avec Σ_2 la matrice de variance covariance alors la statistique de test est

$$T_n = \sqrt{n} \frac{\bar{\mathbf{D}}_n}{S_n^{(\mathbf{D})}}$$

qui a pour loi une loi de Student centrée à $n - 1$ degrés de liberté sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau exact $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\mu_{\mathbf{X}}^* \neq \mu_{\mathbf{Y}}^*$	$\mu_{\mathbf{X}}^* < \mu_{\mathbf{Y}}^*$	$\mu_{\mathbf{X}}^* > \mu_{\mathbf{Y}}^*$
Θ_1	$\{(x, y) \in \mathbb{R}^2 \mid x \neq y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x < y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x > y\}$
Règle	$ T_n > s_{1-\alpha/2}^{(n-1)}$	$T_n < s_{\alpha}^{(n-1)}$	$T_n > s_{1-\alpha}^{(n-1)}$

où $s_{\alpha}^{(n-1)}$ est le quantile d'ordre α d'une loi de Student centrée à $n - 1$ degrés de liberté ; c'est-à-dire que si S suit une loi de Student centrée à $n - 1$ degrés de liberté $\mathcal{T}(n - 1)$ alors $s_{\alpha}^{(n-1)}$ est tel que :

$$\mathbb{P}(S \leq s_{\alpha}^{(n-1)}) = \alpha.$$



Codage en R

La fonction **R** permettant de faire le test bilatéral est `t.test`. Par exemple :

```
1 ## Simulation d'un jeu de données
2 X <- rnorm(100)
3 Y <- rnorm(100)
4 ## Test sous H_0 : mu=0 -> acceptation avec 95% de chances
5 t.test(X,Y,paired=T)
6 ## Test sous H_0 : mu=1 -> rejet normalement
7 t.test(X,Y+1,paired=T)
8 ## Simulation d'un jeu de données avec variances différentes
9 X <- rnorm(100)*2
10 Y <- rnorm(100)
11 ## Test sous H_0 : mu=0 -> acceptation avec 95% de chances
```

```

12 t.test(X,Y,paired=T)
13 ## Test sous H_0 : mu=1 -> rejet normalement
14 t.test(X,Y+1,paired=T)

```

Cas non gaussien : Dans le cas où les deux échantillons ne sont pas gaussiens, nous supposons toutefois que les lois admettent un moment d'ordre 2 et nous utilisons à nouveau le théorème de la limite centrale :

Point méthode (Test d'égalité des espérances ; cas non gaussien et échantillons appariés)

Étant donnés deux n -échantillons appariés $(\mathbf{X}, \mathbf{Y}) = [(X_1, Y_1), \dots, (X_n, Y_n)]$ de loi quelconque admettant un moment d'ordre 2 alors la statistique de test est

$$T_n = \sqrt{n} \frac{\bar{\mathbf{D}}_n}{S_n^{(\mathbf{D})^2}}$$


qui a pour loi **asymptotique** une loi de Student centrée à $n - 1$ degrés de liberté sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau **asymptotique** $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\mu_{\mathbf{X}}^* \neq \mu_{\mathbf{Y}}^*$	$\mu_{\mathbf{X}}^* < \mu_{\mathbf{Y}}^*$	$\mu_{\mathbf{X}}^* > \mu_{\mathbf{Y}}^*$
Θ_1	$\{(x, y) \in \mathbb{R}^2 \mid x \neq y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x < y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x > y\}$
Règle	$ T_n > s_{1-\alpha/2}^{(n-1)}$	$T_n < s_{\alpha}^{(n-1)}$	$T_n > s_{1-\alpha}^{(n-1)}$

où $s_{\alpha}^{(n-1)}$ est le quantile d'ordre α d'une loi de Student centrée à $n - 1$ degrés de liberté ; c'est-à-dire que si S suit une loi de Student centrée à $n - 1$ degrés de liberté $\mathcal{T}(n - 1)$ alors $s_{\alpha}^{(n-1)}$ est tel que :

$$\mathbb{P}(S \leq s_{\alpha}^{(n-1)}) = \alpha.$$

Codage en R

La fonction  permettant de faire le test bilatéral est `t.test`. Par exemple :

```

1 ## Simulation d'un jeu de données
2 X <- rexp(100)
3 Y <- rexp(100)
4 ## Test sous H_0 : mu=0 -> acceptation avec 95% de chances
5 t.test(X,Y,paired=T)
6 ## Test sous H_0 : mu=1 -> rejet normalement
7 t.test(X,exp(100,1/2),paired=T)

```

6.3.2 Test de comparaison des égalités des variances de deux échantillons

Dans le cas de deux échantillons indépendants, nous commençons par supposer que les lois ont un moment d'ordre 2 et nous notons $\sigma_{\mathbf{X}}^{*2}$ et $\sigma_{\mathbf{Y}}^{*2}$ les variances associées à chaque échantillon. La première question à se poser est de savoir si les variances sont égales ou non. Pour ce faire, nous regardons donc le test :

$$\begin{cases} \mathcal{H}_0 : \sigma_{\mathbf{X}}^{*2} = \sigma_{\mathbf{Y}}^{*2}, \\ \mathcal{H}_1 : \sigma_{\mathbf{X}}^{*2} \neq \sigma_{\mathbf{Y}}^{*2}. \end{cases} \quad (6.3)$$

Pour ce faire, nous allons utiliser les deux estimateurs non biaisés $S_{n_{\mathbf{X}}}^{(\mathbf{X})^2} = \frac{1}{n_{\mathbf{X}}-1} \sum_{i=1}^{n_{\mathbf{X}}} (X_i - \bar{\mathbf{X}}_{n_{\mathbf{X}}})^2$ et $S_{n_{\mathbf{Y}}}^{(\mathbf{Y})^2} = \frac{1}{n_{\mathbf{Y}}-1} \sum_{i=1}^{n_{\mathbf{Y}}} (Y_i - \bar{\mathbf{Y}}_{n_{\mathbf{Y}}})^2$. Comme les deux valeurs sont (strictement) positives (sauf si toutes les observations sont égales), nous étudions le rapport des deux estimations :

$$F_n = \frac{S_{n_{\mathbf{X}}}^{(\mathbf{X})^2}}{S_{n_{\mathbf{Y}}}^{(\mathbf{Y})^2}}.$$

Cas gaussien :

Si les observations sont gaussiennes, nous savons par le théorème 24 de Cochran que les deux statistiques $(n_{\mathbf{X}} - 1) S_{n_{\mathbf{X}}}^{(\mathbf{X})^2}$ et $(n_{\mathbf{Y}} - 1) S_{n_{\mathbf{Y}}}^{(\mathbf{Y})^2}$ suivent des lois du χ^2 à $n_{\mathbf{X}} - 1$ et $n_{\mathbf{Y}} - 1$ degrés de libertés respectivement et nous avons :

Définitions 59 (Test d'égalité des variances ; statistique pour échantillons indépendants)

Pour tester deux échantillons indépendants $\mathbf{X} = (X_1, \dots, X_{n_{\mathbf{X}}})$ et $\mathbf{Y} = (Y_1, \dots, Y_{n_{\mathbf{Y}}})$, nous utilisons la statistique de test suivante :

$$F_n = \frac{S_{n_{\mathbf{X}}}^{(\mathbf{X})^2}}{S_{n_{\mathbf{Y}}}^{(\mathbf{Y})^2}}.$$

avec $S_{n_{\mathbf{X}}}^{(\mathbf{X})^2} = \frac{1}{n_{\mathbf{X}} - 1} \sum_{i=1}^{n_{\mathbf{X}}} (X_i - \bar{X}_{n_{\mathbf{X}}})^2$ et $S_{n_{\mathbf{Y}}}^{(\mathbf{Y})^2} = \frac{1}{n_{\mathbf{Y}} - 1} \sum_{i=1}^{n_{\mathbf{Y}}} (Y_i - \bar{Y}_{n_{\mathbf{Y}}})^2$ des estimateurs non biaisés respectifs de la variance pour chaque échantillon.

Et nous obtenons

Point méthode (Test d'égalité des variances ; cas gaussien)

Étant donnés deux n -échantillons indépendants $\mathbf{X} = (X_1, \dots, X_{n_{\mathbf{X}}})$ et $\mathbf{Y} = (Y_1, \dots, Y_{n_{\mathbf{Y}}})$ de lois gaussiennes $\mathcal{N}(\mu_{\mathbf{X}}^*, \sigma_{\mathbf{X}}^{*2})$ et $\mathcal{N}(\mu_{\mathbf{Y}}^*, \sigma_{\mathbf{Y}}^{*2})$ alors la statistique de test est

$$F_n = \frac{S_{n_{\mathbf{X}}}^{(\mathbf{X})^2}}{S_{n_{\mathbf{Y}}}^{(\mathbf{Y})^2}}$$


qui a pour loi une loi de Fisher de paramètres $n_{\mathbf{X}} - 1$ et $n_{\mathbf{Y}} - 1$ sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau exact $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\sigma_{\mathbf{X}}^{*2} \neq \sigma_{\mathbf{Y}}^{*2}$	$\sigma_{\mathbf{X}}^{*2} < \sigma_{\mathbf{Y}}^{*2}$	$\sigma_{\mathbf{X}}^{*2} > \sigma_{\mathbf{Y}}^{*2}$
Θ_1	$\{(x, y) \in \mathbb{R}_+^2 \mid x \neq y\}$	$\{(x, y) \in \mathbb{R}_+^2 \mid x < y\}$	$\{(x, y) \in \mathbb{R}_+^2 \mid x > y\}$
Règle	$F_n < f_{\alpha/2}^{(n_{\mathbf{X}}-1, n_{\mathbf{Y}}-1)}$ ou $F_n > f_{1-\alpha/2}^{(n_{\mathbf{X}}-1, n_{\mathbf{Y}}-1)}$	$F_n < f_{\alpha}^{(n_{\mathbf{X}}-1, n_{\mathbf{Y}}-1)}$	$F_n > f_{1-\alpha}^{(n_{\mathbf{X}}-1, n_{\mathbf{Y}}-1)}$

où $f_{\alpha}^{(n_{\mathbf{X}}-1, n_{\mathbf{Y}}-1)}$ est le quantile d'ordre α d'une loi de Fisher de paramètres $n_{\mathbf{X}} - 1$ et $n_{\mathbf{Y}} - 1$; c'est-à-dire que si F suit une loi de Fisher de paramètres $n_{\mathbf{X}} - 1$ et $n_{\mathbf{Y}} - 1$ $\mathcal{F}(n_{\mathbf{X}} - 1, n_{\mathbf{Y}} - 1)$ alors $f_{\alpha}^{(n_{\mathbf{X}}-1, n_{\mathbf{Y}}-1)}$ est tel que :

$$\mathbb{P}(F \leq f_{\alpha}^{(n_{\mathbf{X}}-1, n_{\mathbf{Y}}-1)}) = \alpha.$$

Codage en R

La fonction  permettant de faire le test bilatéral est `t.test`. Par exemple :

```
1 ## Simulation d'un jeu de données
2 X <- rnorm(100)
3 Y <- rnorm(200)+1
4 ## Test sous H_0 : sigma_X=sigma_Y -> acceptation avec 95% de chances
5 var.test(X,Y)
6 ## Test sous H_0 : sigma_X=sigma_Y -> rejet normalement
7 var.test(X,Y*2)
```

Cas non gaussien :

Dans ce cas, nous utilisons une nouvelle fois le théorème de la limite centrale :

Point méthode (Test d'égalité des variances ; cas non gaussien)

Étant donnés deux n -échantillons indépendants $\mathbf{X} = (X_1, \dots, X_{n_X})$ et $\mathbf{Y} = (Y_1, \dots, Y_{n_Y})$ de lois admettant un moment d'ordre deux alors la statistique de test est

$$F_n = \frac{S_{n_X}^{(\mathbf{X})^2}}{S_{n_Y}^{(\mathbf{Y})^2}}$$


qui a pour loi **asymptotique** une loi de Fisher de paramètres $n_X - 1$ et $n_Y - 1$ sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau **asymptotique** $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\sigma_{\mathbf{X}}^{*2} \neq \sigma_{\mathbf{Y}}^{*2}$	$\sigma_{\mathbf{X}}^{*2} < \sigma_{\mathbf{Y}}^{*2}$	$\sigma_{\mathbf{X}}^{*2} > \sigma_{\mathbf{Y}}^{*2}$
Θ_1	$\{(x, y) \in \mathbb{R}_+^2 \mid x \neq y\}$	$\{(x, y) \in \mathbb{R}_+^2 \mid x < y\}$	$\{(x, y) \in \mathbb{R}_+^2 \mid x > y\}$
Règle	$F_n < f_{\alpha/2}^{(n_X-1, n_Y-1)}$ ou $F_n > f_{1-\alpha/2}^{(n_X-1, n_Y-1)}$	$F_n < f_{\alpha}^{(n_X-1, n_Y-1)}$	$F_n > f_{1-\alpha}^{(n_X-1, n_Y-1)}$

où $f_{\alpha}^{(n_X-1, n_Y-1)}$ est le quantile d'ordre α d'une loi de Fisher de paramètres $n_X - 1$ et $n_Y - 1$; c'est-à-dire que si F suit une loi de Fisher de paramètres $n_X - 1$ et $n_Y - 1$ $\mathcal{F}(n_X - 1, n_Y - 1)$ alors $f_{\alpha}^{(n_X-1, n_Y-1)}$ est tel que :

$$\mathbb{P}(F \leq f_{\alpha}^{(n_X-1, n_Y-1)}) = \alpha.$$

Codage en R

La fonction  permettant de faire le test bilatéral est `t.test`. Par exemple :

```
1 ## Simulation d'un jeu de données
2 X <- rpois(100,2)
3 Y <- rexp(200,1/sqrt(2))
4 ## Test sous H_0 : sigma_X=sigma_Y -> acceptation avec 95% de chances
5 var.test(X,Y)
6 ## Test sous H_0 : sigma_X=sigma_Y -> rejet normalement
7 var.test(X, rexp(200,1))
```

6.3.3 Cas d'échantillons indépendants

Dans le cas de deux échantillons indépendants, nous commençons par supposer que les lois ont un moment d'ordre 2 et nous notons $\sigma_{\mathbf{X}}^{*2}$ et $\sigma_{\mathbf{Y}}^{*2}$ les variances associées à chaque échantillon. La première question à se poser est de savoir si les variances sont égales ou non. Pour ce faire, nous renvoyons à la section 6.3.2.

Cas des variances égales

Dans le cas où les variances $\sigma_{\mathbf{X}}^{*2}$ et $\sigma_{\mathbf{Y}}^{*2}$ sont égales alors nous regarderons la statistique suivante :

Définitions 60 (Test d'égalité d'espérance ; statistique pour échantillons indépendants)

Pour tester deux échantillons indépendants $\mathbf{X} = (X_1, \dots, X_{n_X})$ et $\mathbf{Y} = (Y_1, \dots, Y_{n_Y})$, nous utilisons la statistique de test suivante :

$$Z_{n_X, n_Y} = \frac{\bar{\mathbf{X}}_{n_X} - \bar{\mathbf{Y}}_{n_Y}}{S_{n_X, n_Y}^2}$$

avec $\bar{\mathbf{X}}_{n_X}$ et $\bar{\mathbf{Y}}_{n_Y}$ les moyennes empiriques respectives et S_{n_X, n_Y}^2 est un estimateur non biaisé



de la variance commune

$$S_{n_X, n_Y} = \sqrt{\frac{(n_X - 1) S_{n_X}^{(X)^2} + (n_Y - 1) S_{n_Y}^{(Y)^2}}{n_X + n_Y - 2}}$$

avec $S_{n_X}^{(X)^2} = \frac{1}{n_X - 1} \sum_{i=1}^{n_X} (X_i - \bar{X}_{n_X})^2$ et $S_{n_Y}^{(Y)^2} = \frac{1}{n_Y - 1} \sum_{i=1}^{n_Y} (Y_i - \bar{Y}_{n_Y})^2$ des estimateurs non biaisés respectifs de la variance pour chaque échantillon.

Cas gaussien : Dans le cas où les deux échantillons sont gaussiens, nous avons la loi suivante :

Point méthode (Test d'égalité des espérances ; cas gaussien et échantillons indépendants)

Étant donnés deux n -échantillons indépendants $\mathbf{X} = (X_1, \dots, X_{n_X})$ et $\mathbf{Y} = (Y_1, \dots, Y_{n_Y})$ de lois gaussiennes $\mathcal{N}(\mu_X^*, \sigma_X^{*2})$ et $\mathcal{N}(\mu_Y^*, \sigma_Y^{*2})$ avec $\sigma_X^{*2} = \sigma_Y^{*2}$ alors la statistique de test est

$$Z_{n_X, n_Y} = \frac{\bar{X}_{n_X} - \bar{Y}_{n_Y}}{S_{n_X, n_Y}}$$

qui a pour loi une loi de Student centrée à $n_X + n_Y - 2$ degrés de liberté sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau exact $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\mu_X^* \neq \mu_Y^*$	$\mu_X^* < \mu_Y^*$	$\mu_X^* > \mu_Y^*$
Θ_1	$\{(x, y) \in \mathbb{R}^2 \mid x \neq y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x < y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x > y\}$
Règle	$ Z_n > s_{1-\alpha/2}^{(n_X+n_Y-2)}$	$Z_n < s_{\alpha}^{(n_X+n_Y-2)}$	$Z_n > s_{1-\alpha}^{(n_X+n_Y-2)}$

où $s_{\alpha}^{(n_X+n_Y-2)}$ est le quantile d'ordre α d'une loi de Student centrée à $n_X + n_Y - 2$ degrés de liberté ; c'est-à-dire que si S suit une loi de Student centrée à $n_X + n_Y - 2$ degrés de liberté $\mathcal{T}(n_X + n_Y - 2)$ alors $s_{\alpha}^{(n_X+n_Y-2)}$ est tel que :

$$\mathbb{P}(S \leq s_{\alpha}^{(n_X+n_Y-2)}) = \alpha.$$



Codage en R

La fonction permettant de faire le test bilatéral est `t.test`. Par exemple :

```
1 ## Simulation d'un jeu de données
2 X <- rnorm(100)
3 Y <- rnorm(200)
4 ## Test sous H_0 : mu=0 -> acceptation avec 95% de chances
5 t.test(X, Y, var.equal=T)
6 ## Test sous H_0 : mu=1 -> rejet normalement
7 t.test(X, Y+1, var.equal=T)
```

Cas non gaussien : Dans le cas où les deux échantillons ne sont pas gaussiens, nous supposons toutefois que les lois admettent un moment d'ordre 2 et nous utilisons à nouveau le théorème de la limite centrale :

Point méthode (Test d'égalité des espérances ; cas non gaussien et échantillons indépendants)

Étant donnés deux n -échantillons indépendants $\mathbf{X} = (X_1, \dots, X_{n_X})$ et $\mathbf{Y} = (Y_1, \dots, Y_{n_Y})$ de loi quelconque admettant un moment d'ordre 2 alors la statistique de test est

$$Z_{n_X, n_Y} = \frac{\bar{X}_{n_X} - \bar{Y}_{n_Y}}{S_{n_X, n_Y}}$$

qui a pour loi une loi **asymptotique** de Student centrée à $n_{\mathbf{X}} + n_{\mathbf{Y}} - 2$ degrés de liberté sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau **asymptotique** $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\mu_{\mathbf{X}}^* \neq \mu_{\mathbf{Y}}^*$	$\mu_{\mathbf{X}}^* < \mu_{\mathbf{Y}}^*$	$\mu_{\mathbf{X}}^* > \mu_{\mathbf{Y}}^*$
Θ_1	$\{(x, y) \in \mathbb{R}^2 \mid x \neq y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x < y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x > y\}$
Règle	$ Z_n > s_{1-\alpha/2}^{(n_{\mathbf{X}}+n_{\mathbf{Y}}-2)}$	$Z_n < s_{\alpha}^{(n_{\mathbf{X}}+n_{\mathbf{Y}}-2)}$	$Z_n > s_{1-\alpha}^{(n_{\mathbf{X}}+n_{\mathbf{Y}}-2)}$

où $s_{\alpha}^{(n_{\mathbf{X}}+n_{\mathbf{Y}}-2)}$ est le quantile d'ordre α d'une loi de Student centrée à $n_{\mathbf{X}} + n_{\mathbf{Y}} - 2$ degrés de liberté; c'est-à-dire que si S suit une loi de Student centrée à $n_{\mathbf{X}} + n_{\mathbf{Y}} - 2$ degrés de liberté $\mathcal{T}(n_{\mathbf{X}} + n_{\mathbf{Y}} - 2)$ alors $s_{\alpha}^{(n_{\mathbf{X}}+n_{\mathbf{Y}}-2)}$ est tel que :

$$\mathbb{P}\left(S \leq s_{\alpha}^{(n_{\mathbf{X}}+n_{\mathbf{Y}}-2)}\right) = \alpha.$$

Codage en R

La fonction **R** permettant de faire le test bilatéral est **t.test**. Par exemple :

```
1 ## Simulation d'un jeu de données
2 X <- rchisq(100,25)
3 Y <- rchisq(200,25)
4 ## Test sous H_0 : mu=0 -> acceptation avec 95% de chances
5 t.test(X,Y,var.equal=T)
6 ## Test sous H_0 : mu=1 -> rejet normalement
7 t.test(X,Y+5,var.equal=T)
```

Cas des variances inégales

Dans le cas où les variances $\sigma_{\mathbf{X}}^{*2}$ et $\sigma_{\mathbf{Y}}^{*2}$ sont différentes alors nous regarderons la statistique suivante :

Définitions 61 (Test d'égalité d'espérance ; statistique pour échantillons indépendants)

Pour tester deux échantillons indépendants $\mathbf{X} = (X_1, \dots, X_{n_{\mathbf{X}}})$ et $\mathbf{Y} = (Y_1, \dots, Y_{n_{\mathbf{Y}}})$, nous utilisons la statistique de test suivante :

$$Z_{n_{\mathbf{X}}, n_{\mathbf{Y}}} = \frac{\bar{\mathbf{X}}_{n_{\mathbf{X}}} - \bar{\mathbf{Y}}_{n_{\mathbf{Y}}}}{\sqrt{\frac{S_{n_{\mathbf{X}}}^{(\mathbf{X})}}{n_{\mathbf{X}}} + \frac{S_{n_{\mathbf{Y}}}^{(\mathbf{Y})}}{n_{\mathbf{Y}}}}$$

avec $\bar{\mathbf{X}}_{n_{\mathbf{X}}}$ et $\bar{\mathbf{Y}}_{n_{\mathbf{Y}}}$ les moyennes empiriques respectives et $S_{n_{\mathbf{X}}}^{(\mathbf{X})2} = \frac{1}{n_{\mathbf{X}}-1} \sum_{i=1}^{n_{\mathbf{X}}} (X_i - \bar{\mathbf{X}}_{n_{\mathbf{X}}})^2$ et $S_{n_{\mathbf{Y}}}^{(\mathbf{Y})2} = \frac{1}{n_{\mathbf{Y}}-1} \sum_{i=1}^{n_{\mathbf{Y}}} (Y_i - \bar{\mathbf{Y}}_{n_{\mathbf{Y}}})^2$ des estimateurs non biaisés respectifs de la variance pour chaque échantillon.

Cas gaussien : Dans le cas où les deux échantillons sont gaussiens, nous avons la loi suivante :

Point méthode (Test de Welch d'égalité des espérances)

Étant donnés deux n -échantillons indépendants $\mathbf{X} = (X_1, \dots, X_{n_{\mathbf{X}}})$ et $\mathbf{Y} = (Y_1, \dots, Y_{n_{\mathbf{Y}}})$ de lois gaussiennes $\mathcal{N}(\mu_{\mathbf{X}}^*, \sigma_{\mathbf{X}}^{*2})$ et $\mathcal{N}(\mu_{\mathbf{Y}}^*, \sigma_{\mathbf{Y}}^{*2})$ avec $\sigma_{\mathbf{X}}^{*2} \neq \sigma_{\mathbf{Y}}^{*2}$ alors la statistique de test est

$$Z_{n_{\mathbf{X}}, n_{\mathbf{Y}}} = \frac{\bar{\mathbf{X}}_{n_{\mathbf{X}}} - \bar{\mathbf{Y}}_{n_{\mathbf{Y}}}}{\sqrt{\frac{S_{n_{\mathbf{X}}}^{(\mathbf{X})}}{n_{\mathbf{X}}} + \frac{S_{n_{\mathbf{Y}}}^{(\mathbf{Y})}}{n_{\mathbf{Y}}}}$$

qui a pour loi une loi de Student centrée à ν degrés de liberté sous l'hypothèse \mathcal{H}_0 avec ν la

valeur entière la plus proche de :

$$\frac{\left[\frac{S_{n_X}^{(X)}}{n_X} + \frac{S_{n_Y}^{(Y)}}{n_Y} \right]^2}{\frac{S_{n_X}^{(X)2}}{(n_X-1)n_X^2} + \frac{S_{n_Y}^{(Y)2}}{(n_Y-1)n_Y^2}}$$

et la règle de rejet de niveau approché $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$\mu_X^* \neq \mu_Y^*$	$\mu_X^* < \mu_Y^*$	$\mu_X^* > \mu_Y^*$
Θ_1	$\{(x, y) \in \mathbb{R}^2 \mid x \neq y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x < y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x > y\}$
Règle	$ Z_n > s_{1-\alpha/2}^{(\nu)}$	$Z_n < s_{\alpha}^{(\nu)}$	$Z_n > s_{1-\alpha}^{(\nu)}$

où $s_{\alpha}^{(\nu)}$ est le quantile d'ordre α d'une loi de Student centrée à ν degrés de liberté ; c'est-à-dire que si S suit une loi de Student centrée à ν degrés de liberté $\mathcal{T}(\nu)$ alors $s_{\alpha}^{(\nu)}$ est tel que :

$$\mathbb{P}(S \leq s_{\alpha}^{(\nu)}) = \alpha.$$

Ce test est connu sous le nom de **test de Welch**.

Codage en R

La fonction `t.test` permettant de faire le test bilatéral est `t.test`. Par exemple :

```
1 ## Simulation d'un jeu de données
2 X <- rnorm(100)
3 Y <- rnorm(200)*2
4 ## Test sous H_0 : mu=0 -> acceptation avec 95% de chances
5 t.test(X,Y,var.equal=F)
6 ## Test sous H_0 : mu=1 -> rejet normalement
7 t.test(X,Y+1,var.equal=F)
```

Cas non gaussien : Dans le cas où les deux échantillons ne sont pas gaussiens, nous supposons toutefois que les lois admettent un moment d'ordre 2 et nous utilisons à nouveau le théorème de la limite centrale :

Point méthode (Test d'égalité des espérances ; cas non gaussien et échantillons indépendants)

Étant donnés deux n -échantillons indépendants $\mathbf{X} = (X_1, \dots, X_{n_X})$ et $\mathbf{Y} = (Y_1, \dots, Y_{n_Y})$ de loi quelconque admettant des moments d'ordre 2 différents alors la statistique de test est

$$Z_{n_X, n_Y} = \frac{\bar{\mathbf{X}}_{n_X} - \bar{\mathbf{Y}}_{n_Y}}{\sqrt{\frac{S_{n_X}^{(X)}}{n_X} + \frac{S_{n_Y}^{(Y)}}{n_Y}}}$$

qui a pour loi une loi **asymptotique** de Student centrée à ν degrés de liberté sous l'hypothèse \mathcal{H}_0 avec ν la valeur entière la plus proche de :

$$\frac{\left[\frac{S_{n_X}^{(X)}}{n_X} + \frac{S_{n_Y}^{(Y)}}{n_Y} \right]^2}{\frac{S_{n_X}^{(X)2}}{(n_X-1)n_X^2} + \frac{S_{n_Y}^{(Y)2}}{(n_Y-1)n_Y^2}}$$

et la règle de rejet de niveau **asymptotique** approché $\alpha \in]0; 1[$ est donc :


\mathcal{H}_1	$\mu_X^* \neq \mu_Y^*$	$\mu_X^* < \mu_Y^*$	$\mu_X^* > \mu_Y^*$
Θ_1	$\{(x, y) \in \mathbb{R}^2 \mid x \neq y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x < y\}$	$\{(x, y) \in \mathbb{R}^2 \mid x > y\}$
Règle	$ Z_n > s_{1-\alpha/2}^{(\nu)}$	$Z_n < s_{\alpha}^{(\nu)}$	$Z_n > s_{1-\alpha}^{(\nu)}$

où $s_{\alpha}^{(\nu)}$ est le quantile d'ordre α d'une loi de Student centrée à ν degrés de liberté ; c'est-à-dire que si S suit une loi de Student centrée à ν degrés de liberté $\mathcal{T}(\nu)$ alors $s_{\alpha}^{(n_{\mathbf{X}}+n_{\mathbf{Y}}-2)}$ est tel que :

$$\mathbb{P}\left(S \leq s_{\alpha}^{(\nu)}\right) = \alpha.$$



Codage en R

La fonction  permettant de faire le test bilatéral est `t.test`. Par exemple :

```

1 ## Simulation d'un jeu de données
2 X <- rexp(100)
3 Y <- rchisq(200,1)
4 ## Test sous H_0 : mu=0 -> acceptation avec 95% de chances
5 t.test(X,Y,var.equal=F)
6 ## Test sous H_0 : mu=1 -> rejet normalement
7 t.test(X, rexp(200,1/2), var.equal=F)

```

6.3.4 Quel(s) test(s) utilisé(s) ?

Pour faire le bilan de cette partie, nous avons créé un arbre de décision pour savoir quel(s) test(s) utilisé(s) disponible sur la figure 6.4.

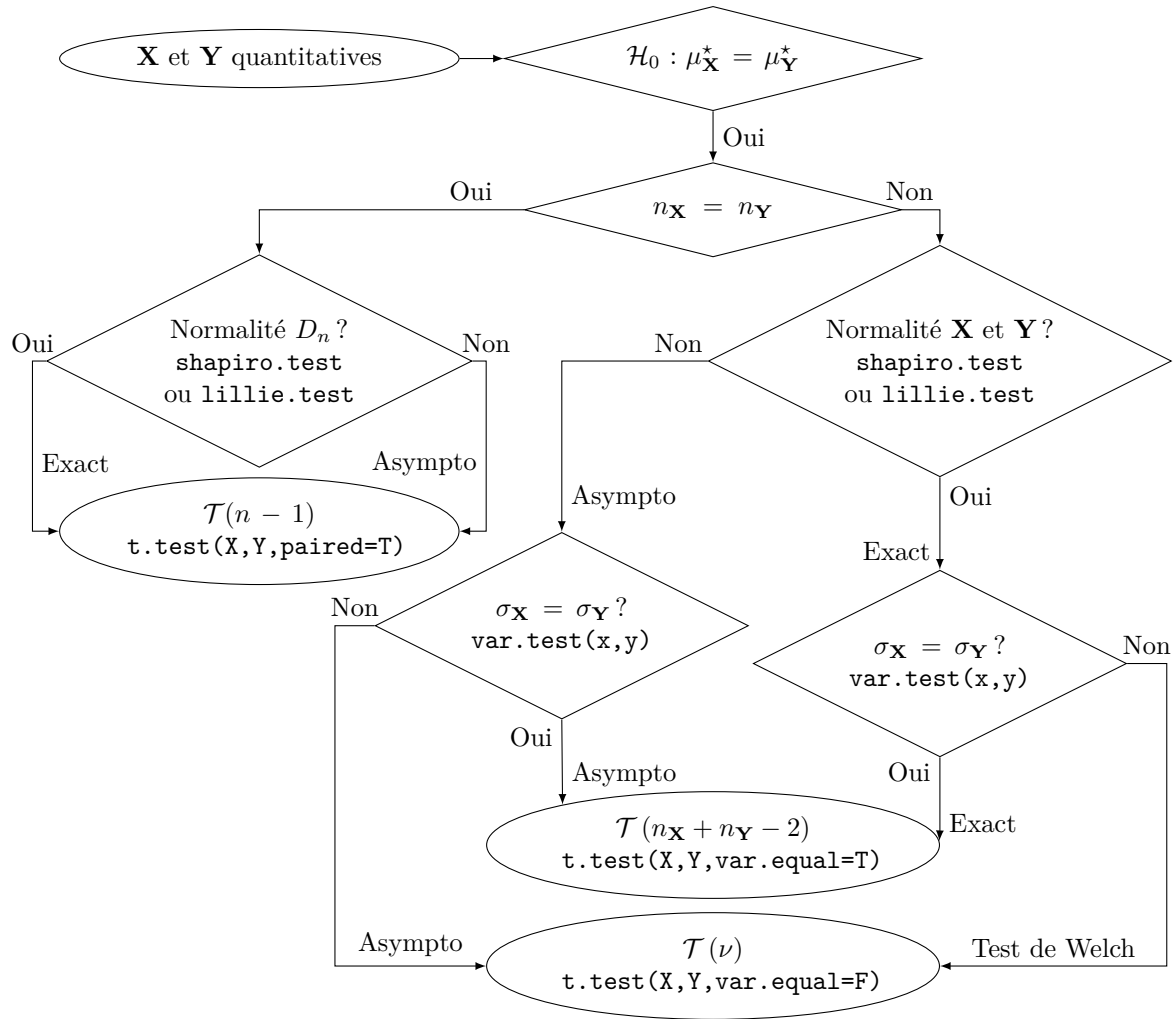


FIGURE 6.4 – Représentation de l’arbre de décision pour savoir quel test utilisé dans le cas de deux échantillons où nous cherchons à savoir si leurs espérances sont égales.

6.3.5 Généralisation à plus de deux échantillons dans le cas gaussien

Ces tests peuvent être généralisés à plus de deux échantillons. Nous mettons ici les grandes idées.

Comparaison de variances, échantillons indépendants

Nous supposons posséder K échantillons $(X_{1,1}, \dots, X_{n_1,1})$, $(X_{1,2}, \dots, X_{n_2,2})$ et ainsi de suite jusqu’à $(X_{1,K}, \dots, X_{n_K,K})$ admettant des moments d’ordre deux. Nous testons donc les variances $\sigma_1^*, \dots, \sigma_K^*$:

$$\begin{cases} \mathcal{H}_0 : \sigma_1^* = \sigma_2^* = \dots = \sigma_K^*, \\ \mathcal{H}_1 : \exists k \neq k', \sigma_k^* \neq \sigma_{k'}^*. \end{cases}$$

On introduit les estimateurs non biaisés de variance et nous avons :

Définitions 62 (Test d’égalité des variances ; statistique pour plusieurs échantillons indépendants)

Pour tester plusieurs échantillons indépendants $(X_{1,1}, \dots, X_{n_1,1})$, $(X_{1,2}, \dots, X_{n_2,2})$ et ainsi de suite jusqu’à $(X_{1,K}, \dots, X_{n_K,K})$ admettant des moments d’ordre deux, nous utilisons la statistique de test suivante :

$$F_n = \frac{(n - k) \ln \left[\frac{1}{n-K} \sum_{k=1}^K (n_k - 1) S_k^2 \right] - \sum_{k=1}^K (n_k - 1) \ln [S_k^2]}{1 + \frac{1}{3(K-1)} \left[\sum_{k=1}^K \left(\frac{1}{n_k-1} \right) - \frac{1}{n-K} \right]}$$

où $n = \sum_{k=1}^K n_k$ et pour tout $k \in \{1, \dots, K\}$, $S_k^2 = \frac{1}{n_k} \sum_{i=1}^{n_k} (X_{i,k} - \bar{X}_{\cdot,k})^2$ avec $\bar{X}_{\cdot,k}$ est la moyenne de l'échantillon k .

Et nous obtenons

Point méthode (Test de Barlett)

Étant donnés plusieurs échantillons indépendants $(X_{1,1}, \dots, X_{n_1,1})$, $(X_{1,2}, \dots, X_{n_2,2})$ et ainsi de suite jusqu'à $(X_{1,K}, \dots, X_{n_K,K})$ admettant des moments d'ordre deux alors la statistique de test est

$$F_n = \frac{(n - k) \ln \left[\frac{1}{n-K} \sum_{k=1}^K (n_k - 1) S_k^2 \right] - \sum_{k=1}^K (n_k - 1) \ln [S_k^2]}{1 + \frac{1}{3(K-1)} \left[\sum_{k=1}^K \left(\frac{1}{n_k-1} \right) - \frac{1}{n-K} \right]}$$


qui a pour loi **asymptotique** une loi du $\chi^2(K-1)$ à $K-1$ degrés de liberté sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau exact $\alpha \in]0; 1[$ est donc :

$$\begin{array}{l|l} \mathcal{H}_1 & \exists k \neq k', \sigma_k^* \neq \sigma_{k'}^* \\ \Theta_1 & \mathbb{R}_+^K \setminus \left\{ (x_1, x_2, \dots, x_K) \in \mathbb{R}_+^K \mid x_1 = x_2 = \dots = x_K \right\} \\ \text{Règle} & F_n > s_{1-\alpha}^{(K-1)} \end{array}$$

où $s_{1-\alpha}^{(K-1)}$ est le quantile d'ordre α d'une loi du $\chi^2(K-1)$ à $K-1$ degrés de liberté ; c'est-à-dire que si S suit une loi du $\chi^2(K-1)$ à $K-1$ degrés de liberté alors $s_{1-\alpha}^{(K-1)}$ est tel que :

$$\mathbb{P} \left(S \leq s_{1-\alpha}^{(K-1)} \right) = \alpha.$$

Codage en R

La fonction  permettant de faire le test bilatéral est `bartlett.test`. Par exemple :

```
1 ## Simulation d'un jeu de données
2 X<-data.frame(X=c(rnorm(100),rnorm(200)+1,rnorm(300)-1),
3               groupe=factor(c(rep(1,100),rep(2,200),rep(3,300))))
4
5 ## Test sous H_0 : sigma_X=sigma_Y -> acceptation avec 95% de chances
6 bartlett.test(X$X,X$groupe)
7 ## Test sous H_0 : sigma_X=sigma_Y -> rejet normalement
8 X<-data.frame(X=c(rnorm(100),2*rnorm(200)+1,rnorm(300)-1),
9               groupe=factor(c(rep(1,100),rep(2,200),rep(3,300))))
10 bartlett.test(X$X,X$groupe)
```

Comparaison d'espérances, échantillons appariés

Si les échantillons sont appariés alors nous avons une matrice $\mathbf{X} = (X_{i,k})$ où chaque colonne $X_{\cdot,k}$ possède sa propre loi d'espérance μ_k^* . Nous testons donc :

$$\begin{cases} \mathcal{H}_0 : \mu_1^* = \mu_2^* = \dots = \mu_K^*, \\ \mathcal{H}_1 : \exists k \neq k', \mu_k^* \neq \mu_{k'}^*. \end{cases}$$

Ce test revient à tester l'effet du facteur expérience dans un modèle d'ANOVA 2 sans interaction (voir le cours de modélisation linéaire). Nous avons :



Définitions 63 (Test d'égalité des espérances ; statistique pour plusieurs échantillons appariés)

Pour tester plusieurs échantillons appariés $\mathbf{X} = (X_{i,k})$, nous utilisons la statistique de test suivante :

$$F_n = \frac{n \sum_{k=1}^K (\bar{X}_{\cdot,k} - \bar{X}_{\cdot,\cdot})^2}{\frac{1}{n-1} \sum_{k=1}^K \sum_{i=1}^n (X_{i,k} - \bar{X}_{\cdot,k} - \bar{X}_{i,\cdot} + \bar{X}_{\cdot,\cdot})^2}$$

où $\bar{X}_{\cdot,k}$ et $\bar{X}_{i,\cdot}$ sont respectivement les moyennes par colonne et par ligne et $\bar{X}_{\cdot,\cdot}$ est la moyenne globale.

Et nous obtenons

Point méthode (Test d'égalité des espérances ; cas gaussien)

Étant donnés plusieurs échantillons appariés $\mathbf{X} = (X_{i,k})$ de lois gaussiennes alors la statistique de test est

$$F_n = \frac{n \sum_{k=1}^K (\bar{X}_{\cdot,k} - \bar{X}_{\cdot,\cdot})^2}{\frac{1}{n-1} \sum_{k=1}^K \sum_{i=1}^n (X_{i,k} - \bar{X}_{\cdot,k} - \bar{X}_{i,\cdot} + \bar{X}_{\cdot,\cdot})^2}$$

qui a pour loi une loi de Fisher de paramètres $K - 1$ et $(n - 1)(K - 1)$ sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau exact $\alpha \in]0; 1[$ est donc :

$$\begin{array}{l|l} \mathcal{H}_1 & \exists k \neq k', \mu_k^* \neq \mu_{k'}^* \\ \Theta_1 & \mathbb{R}^K \setminus \{ (x_1, x_2, \dots, x_K) \in \mathbb{R}^K \mid x_1 = x_2 = \dots = x_K \} \\ \text{Règle} & F_n < f_{\alpha/2}^{(K-1, (n-1)(K-1))} \text{ ou } F_n > f_{1-\alpha/2}^{(K-1, (n-1)(K-1))} \end{array}$$

où $f_{\alpha}^{(K-1, (n-1)(K-1))}$ est le quantile d'ordre α d'une loi de Fisher de paramètres $K - 1$ et $(n - 1)(K - 1)$; c'est-à-dire que si F suit une loi de Fisher de paramètres $n_{\mathbf{X}} - 1$ et $n_{\mathbf{Y}} - 1$ $\mathcal{F}(n_{\mathbf{X}} - 1, n_{\mathbf{Y}} - 1)$ alors $f_{\alpha}^{(K-1, (n-1)(K-1))}$ est tel que :

$$\mathbb{P} \left(F \leq f_{\alpha}^{(K-1, (n-1)(K-1))} \right) = \alpha.$$

**Codage en R**

Sous condition que les espérances soient égales (voir le test de Barrett), la fonction **R** permettant de faire le test bilatéral est `aov`. Par exemple :

```
1 ## Simulation d'un jeu de données
2 X<-data.frame(X=c(rnorm(100),rnorm(200),rnorm(300)),
3               groupe=factor(c(rep(1,100),rep(2,200),rep(3,300))))
4 ## Test sous H_0 : sigma_X=sigma_Y -> acceptation avec 95% de chances
5 summary(aov(X$X~X$groupe))
6 ## Test sous H_0 : sigma_X=sigma_Y -> rejet normalement
7 X<-data.frame(X=c(rnorm(100),rnorm(200)+1,rnorm(300)-1),
8               groupe=factor(c(rep(1,100),rep(2,200),rep(3,300))))
9 summary(aov(X$X~X$groupe))
```

6.4 Test de corrélation

Enfin, le dernier test des variables quantitatives est celui des corrélations. Dans ce cadre, nous supposons avoir des échantillons appariés $(X_i, Y_i)_{1 \leq i \leq n}$ qui proviennent d'une loi binormale d'espérance $(\mu_{\mathbf{X}}^*, \mu_{\mathbf{Y}}^*)$ et de matrice de variance covariance :

$$\begin{pmatrix} \sigma_{\mathbf{X}}^{*2} & \sigma_{\mathbf{X}}^* \sigma_{\mathbf{Y}}^* \rho \\ \sigma_{\mathbf{X}}^* \sigma_{\mathbf{Y}}^* \rho & \sigma_{\mathbf{Y}}^{*2} \end{pmatrix}$$

avec $\rho \in \mathbb{R}$ appelé **coefficient de corrélation** (voir la section 1.5.2). L'estimateur naturel est le coefficient de corrélation empirique :

Définition 64 (Coefficient de corrélation empirique)

Le coefficient de corrélation empirique $\hat{\rho}_n$ est défini par :

$$\hat{\rho}_n = \frac{\text{Cov}_n(\mathbf{X}, \mathbf{Y})}{\hat{\sigma}_n^{(\mathbf{X})} \hat{\sigma}_n^{(\mathbf{Y})}} = \frac{\sum_{i=1}^n (X_i - \bar{X}_n) (Y_i - \bar{Y}_n)}{\sqrt{\sum_{i=1}^n (X_i - \bar{X}_n)^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y}_n)^2}}.$$

Ainsi, la question que nous nous posons est de savoir si $\rho = 0$ ou pas :

$$\begin{cases} \mathcal{H}_0 : \rho = 0, \\ \mathcal{H}_1 : \rho \neq 0. \end{cases}$$

Nous avons :

Définitions 65 (Test de corrélation)

Pour tester la corrélation de deux échantillons appariés $(X_i, Y_i)_{1 \leq i \leq n}$, nous utilisons la statistique de test suivante :

$$T_n = \sqrt{n-2} \times \frac{\hat{\rho}_n}{\sqrt{1 - \hat{\rho}_n^2}}.$$

Et nous obtenons

Point méthode (Test de corrélation)

Étant donnés deux échantillons appariés $(X_i, Y_i)_{1 \leq i \leq n}$ de loi binormale alors la statistique de test est

$$T_n = \sqrt{n-2} \times \frac{\hat{\rho}_n}{\sqrt{1 - \hat{\rho}_n^2}}$$

qui a pour loi une loi de Student $\mathcal{T}(n-2)$ à $n-2$ degrés de libertés sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau exact $\alpha \in]0; 1[$ est donc :

$$\begin{array}{c|c} \mathcal{H}_1 & \rho \neq 0 \\ \Theta_1 & \mathbb{R}^* \\ \hline \text{Règle} & |T_n| > t_{1-\alpha/2}^{(n-2)} \end{array}$$

où $t_\alpha^{(n-2)}$ est le quantile d'ordre α d'une loi de Student $\mathcal{T}(n-2)$ à $n-2$ degrés de libertés ; c'est-à-dire que si T suit une loi de Student $\mathcal{T}(n-2)$ à $n-2$ degrés de libertés alors $t_\alpha^{(n-2)}$ est tel que :

$$\mathbb{P}(T \leq t_\alpha^{(n-2)}) = \alpha.$$

Codage en R

La fonction  permettant de faire le test bilatéral est `cor.test`. Par exemple :

```
1 ## Simulation d'un jeu de données
2 X<-rnorm(100)
3 Y<-rnorm(100)
4 ## Test sous H_0 : corrélation -> acceptation avec 95% de chances
5 cor.test(X,Y,method="pearson")
6 ## Test sous H_0 : corrélation -> rejet normalement
7 cor.test(X,X+0.5*Y,method="pearson")
```



6.5 Solutions des exercices



Exercices 6.1

Pour la puissance, nous supposons que μ^* est différent de μ_0 donc, toujours d'après la proposition 21 $\bar{\mathbf{X}}_n$ est de loi gaussienne d'espérance $\mu^* \neq \mu_0$ et de variance σ^2/n . Ainsi, nous avons que T_n suit une loi gaussienne dont il faut calculer l'espérance et la variance :

$$\mathbb{E}[T_n] = \mathbb{E}\left[\sqrt{n}\frac{\bar{\mathbf{X}}_n - \mu_0}{\sigma}\right] = \sqrt{n}\frac{1}{\sigma} [\mathbb{E}[\bar{\mathbf{X}}_n] - \mu_0] = \sqrt{n}\frac{1}{\sigma} [\mu^* - \mu_0].$$

Et par la même démonstration que pour le point méthode, la variance est égale à 1. Ainsi, sous l'hypothèse que $\mu^* \neq \mu_0$, nous avons :

$$T_n \sim \mathcal{N}\left(\frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0], 1\right)$$

donc nous avons par exemple dans le cas où \mathcal{H}_1 est $\mu^* \neq \mu_0$:

$$\begin{aligned} \underline{\beta}(\mu^*) &= \mathbb{P}(|T_n| \leq q_{1-\alpha/2}) \\ &= \mathbb{P}(-q_{1-\alpha/2} < T_n \leq q_{1-\alpha/2}) \\ &= 1 - \mathbb{P}(T_n \leq -q_{1-\alpha/2}) + \mathbb{P}(T_n \leq q_{1-\alpha/2}) - [\mathbb{P}(T_n \leq q_{1-\alpha/2}) - \mathbb{P}(T_n \leq -q_{1-\alpha/2})] \\ &= \mathbb{P}(T_n \leq q_{1-\alpha/2}) - \mathbb{P}(T_n < -q_{1-\alpha/2}) \\ &= \mathbb{P}(T_n \leq q_{1-\alpha/2}) - \mathbb{P}(T_n \leq -q_{1-\alpha/2}) \text{ car } T_n \text{ est continue,} \\ &= \mathbb{P}\left(T_n - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0] \leq q_{1-\alpha/2} - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0]\right) \\ &\quad - \mathbb{P}\left(T_n - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0] \leq -q_{1-\alpha/2} - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0]\right) \\ &= F_{\mathcal{N}(0,1)}\left(q_{1-\alpha/2} - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0]\right) - F_{\mathcal{N}(0,1)}\left(-q_{1-\alpha/2} - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0]\right) \end{aligned}$$

où $F_{\mathcal{N}(0,1)}$ est la fonction de répartition d'une loi gaussienne centrée réduite. En particulier, nous voyons que plus μ^* est loin de μ_0 plus l'erreur de type deux va tendre vers 0.

Par la même démonstration, dans le cas où $\mu^* > \mu_0$, nous montrons que le risque de seconde espèce $\underline{\beta}$ vaut :

$$\begin{aligned} \underline{\beta}(\mu^*) &= \mathbb{P}(\Phi(X) = 0) \\ &= \mathbb{P}(T_n \leq q_{1-\alpha}) \\ &= \mathbb{P}\left(T_n - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0] \leq q_{1-\alpha} - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0]\right) \\ &= F_{\mathcal{N}(0,1)}\left(q_{1-\alpha} - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0]\right) \end{aligned}$$

et nous voyons que pour tout $\mu^* > \mu_0$, nous avons :

$$\begin{aligned} \mu^* > \mu_0 &\Leftrightarrow \mu^* - \mu_0 > 0 \\ &\Leftrightarrow \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0] > 0 \\ &\Leftrightarrow -\frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0] < 0 \\ &\Leftrightarrow q_{1-\alpha} - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0] < q_{1-\alpha} \\ &\Leftrightarrow F_{\mathcal{N}(0,1)}\left(q_{1-\alpha} - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0]\right) < F_{\mathcal{N}(0,1)}(q_{1-\alpha}) \text{ car } F_{\mathcal{N}(0,1)} \text{ est croissante,} \\ &\Leftrightarrow \underline{\beta}(\mu^*) < 1 - \alpha \end{aligned}$$

$$\begin{aligned} \Leftrightarrow -\underline{\beta}(\mu^*) &> \alpha - 1 \\ \Leftrightarrow 1 - \underline{\beta}(\mu^*) &> \alpha \end{aligned}$$

et la suite de tests est bien sans biais. De plus, nous avons pour tout $\mu^* > \mu_0$:

$$\begin{aligned} \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0] \xrightarrow[n \rightarrow +\infty]{} +\infty &\Rightarrow -\frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0] \xrightarrow[n \rightarrow +\infty]{} -\infty \\ &\Rightarrow q_{1-\alpha} - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0] \xrightarrow[n \rightarrow +\infty]{} -\infty \\ &\Rightarrow F_{\mathcal{N}(0,1)} \left(q_{1-\alpha} - \frac{\sqrt{n}}{\sigma} [\mu^* - \mu_0] \right) \xrightarrow[n \rightarrow +\infty]{} 0 \end{aligned}$$

et la suite de tests est bien consistante (puisque sans biais et avec un risque de seconde espèce qui tend vers 0 lorsque n tend vers l'infini).



Exercices 6.2

Nous commençons par remarquer que pour tout $i \in \{1, \dots, n\}$, $X_i - \bar{\mathbf{X}}_n$ est une loi gaussienne par le lemme 22 dont nous devons calculer la moyenne et la variance :

$$\mathbb{E}[X_i - \bar{\mathbf{X}}_n] = \mathbb{E}[X_i] - \mathbb{E}[\bar{\mathbf{X}}_n] = \mu^* - \mu^* = 0.$$

Et pour la variance, il faut faire attention au fait que $\bar{\mathbf{X}}_n$ n'est pas indépendant de X_i puisque la somme inclut cette variable. Il faut donc la traiter à part dans le calcul :

$$\begin{aligned} \mathbb{V}[X_i - \bar{\mathbf{X}}_n] &= \mathbb{V}[X_i] + \mathbb{V}[\bar{\mathbf{X}}_n] - 2\text{Cov}(X_i, \bar{\mathbf{X}}_n) \\ &= \sigma^2 + \frac{\sigma^2}{n} - 2\text{Cov}\left(X_i, \frac{1}{n} \sum_{i'=1}^n X_{i'}\right) \\ &= \frac{(n+1)\sigma^2}{n} - \frac{2}{n} \sum_{i'=1}^n \text{Cov}(X_i, X_{i'}) \\ &= \frac{(n+1)\sigma^2}{n} - \frac{2}{n} \left[\sum_{\substack{i'=1 \\ i' \neq i}}^n \underbrace{\text{Cov}(X_i, X_{i'})}_{=0} + \text{Cov}(X_i, X_i) \right] \\ &= \frac{(n+1)\sigma^2}{n} - \frac{2}{n} \mathbb{V}[X_i] \\ &= \frac{(n-1)\sigma^2}{n}. \end{aligned}$$

Or, comme les variables $Z_i = X_i - \bar{\mathbf{X}}_n$ sont centrées alors pour tout $i \in \{1, \dots, n\}$, nous avons :

$$\mathbb{E}[(X_i - \bar{\mathbf{X}}_n)^2] = \mathbb{E}[Z_i^2] = \mathbb{V}[Z_i] = \frac{(n-1)\sigma^2}{n}.$$

Pour vérifier le biais, nous devons calculer l'espérance de S_n^2 :

$$\begin{aligned} \mathbb{E}[\hat{\sigma}_n^2] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{\mathbf{X}}_n)^2\right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(X_i - \bar{\mathbf{X}}_n)^2] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[Z_i^2] \end{aligned}$$

$$\begin{aligned} &= \frac{1}{n} \sum_{i=1}^n \frac{(n-1)\sigma^2}{n} \\ &= \frac{n(n-1)\sigma^2}{n \cdot n} \\ &= \frac{(n-1)\sigma^2}{n} \end{aligned}$$

et l'estimateur est biaisé alors que nous avons :

$$\mathbb{E}[S_n^2] = \mathbb{E}\left[\frac{n}{n-1}\widehat{\sigma}_n^2\right] = \frac{n}{n-1}\mathbb{E}[\widehat{\sigma}_n^2] = \frac{n}{n-1} \times \frac{(n-1)\sigma^2}{n} = \sigma^2$$

qui est non biaisé.

Chapitre 7

Tests sur les variables qualitatives

7.1 Objectifs

Dans ce chapitre, nous nous concentrons sur les variables qualitatives. Dans ce cadre, nous ne pouvons regarder que les proportions de chaque modalité.

7.2 Test d'une probabilité pour une variable de Bernoulli

La variable de Bernoulli permet de caractériser la présence ou l'absence d'une particularité dans la population. Par exemple, nous pouvons nous intéresser à la présence du caractère A par opposition à un autre caractère. Une utilisation documentée de cette question est celle de Pierre-Simon de Laplace sur la proportion de naissances masculines à Paris (présentée en section 2.2). Dans ce cadre, nous supposons avoir un n -échantillon de loi de Bernoulli $\mathcal{B}(p^*)$ et nous nous demandons si elle est égale à une certaine valeur $p_0 \in]0; 1[$:

$$\begin{cases} \mathcal{H}_0 : p^* = p_0, \\ \mathcal{H}_1 : p^* \neq p_0. \end{cases}$$

L'estimateur usuel de p^* , que ce soit celui des moments ou du maximum de vraisemblance, est $\hat{p}_n = \bar{\mathbf{X}}_n$. Deux solutions s'offrent alors à nous pour résoudre ce problème.

7.2.1 Méthode exacte

La première solution est d'utiliser la loi de $n\bar{\mathbf{X}}_n$ qui est une loi binomiale de paramètre n et p^* . Ainsi, nous avons :

Point méthode (Test d'une probabilité ; niveau exact)

Étant donné un échantillon $\mathbf{X} = (X_i)_{1 \leq i \leq n}$ de loi de Bernoulli $\mathcal{B}(p^*)$ alors la statistique de test est

$$n\bar{\mathbf{X}}_n$$

qui a pour loi une loi binomiale $\text{Bin}(n, p_0)$ sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau exact $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$p^* \neq p_0$	$p^* < p_0$	$p^* > p_0$
Θ_1	$]0, p_0[\cup]p_0, 1[$	$]0, p_0[$	$]p_0, 1[$
Règle	$n\bar{\mathbf{X}}_n < b_{\alpha/2}^{(n, p_0)}$ ou $n\bar{\mathbf{X}}_n > b_{1-\alpha/2}^{(n, p_0)}$	$n\bar{\mathbf{X}}_n < b_{\alpha}^{(n, p_0)}$	$n\bar{\mathbf{X}}_n > b_{1-\alpha}^{(n, p_0)}$

où $b_{\alpha}^{(n, p_0)}$ est le quantile d'ordre α d'une loi binomiale de paramètres n et p_0 ; c'est-à-dire que si B suit une loi binomiale de paramètres n et p_0 alors $b_{\alpha}^{(n, p_0)}$ est tel que :

$$\mathbb{P}(B \leq b_{\alpha}^{(n, p_0)}) = \alpha.$$

**Codage en R**

La fonction permettant de faire le test bilatéral est `binom.test`. Par exemple :

```

1 ## Simulation d'un jeu de données
2 p_0<-0.5
3 n<-100
4 X <- rbinom(n,1,p_0)
5 ## Test sous H_0 : p_0=0.5 -> acceptation avec 95% de chances
6 binom.test(c(sum(X),n-sum(X)),p = p_0)
7 ## Test sous H_0 : p_0=0.75 -> rejet normalement
8 binom.test(c(sum(X),n-sum(X)),p = 0.75)

```

7.2.2 Méthode approchée

La deuxième solution est d'utiliser le théorème de la limite centrale en sachant que nous connaissons la variance théorique sous la loi de \mathcal{H}_0 . Ainsi, nous avons :

Définition 66 (Test d'une probabilité ; statistique approchée)

Pour tester un échantillon $\mathbf{X} = (X_i)_{1 \leq i \leq n}$ de loi de Bernoulli $\mathcal{B}(p^*)$ et vérifier si $p^* = p_0$, nous utilisons la statistique de test suivante :

$$S_n = \sqrt{n} \frac{\bar{\mathbf{X}}_n - p_0}{\sqrt{p_0(1-p_0)}}$$

avec $\bar{\mathbf{X}}_n$ la moyenne empirique du n -échantillon.

Alors nous avons :

Point méthode (Test d'une probabilité ; niveau approché)

Étant donné un échantillon $\mathbf{X} = (X_i)_{1 \leq i \leq n}$ de loi de Bernoulli $\mathcal{B}(p^*)$ alors la statistique de test est

$$S_n = \sqrt{n} \frac{\bar{\mathbf{X}}_n - p_0}{\sqrt{p_0(1-p_0)}}$$

qui a pour loi **asymptotique** une loi normale $\mathcal{N}(0, 1)$ centrée réduite sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau **asymptotique** $\alpha \in]0; 1[$ est donc :

\mathcal{H}_1	$p^* \neq p_0$	$p^* < p_0$	$p^* > p_0$
Θ_1	$]0, p_0[\cup]p_0, 1[$	$]0, p_0[$	$]p_0, 1[$
Règle	$ S_n > q_{1-\alpha/2}$	$S_n < q_\alpha$	$S_n > q_{1-\alpha}$

où q_α est le quantile d'ordre α d'une loi normale centrée réduite ; c'est-à-dire que si X suit une loi gaussienne centrée réduite $\mathcal{N}(0, 1)$ alors q_α est tel que :

$$\mathbb{P}(X \leq q_\alpha) = \alpha.$$

**Codage en R**

La fonction permettant de faire le test bilatéral est `prop.test`. Par exemple :

```

1 ## Simulation d'un jeu de données
2 p_0<-0.5
3 n<-100
4 X <- rbinom(n,1,p_0)
5 ## Test sous H_0 : p_0=0.5 -> acceptation avec 95% de chances
6 prop.test(sum(X),n,p = p_0)
7 ## Test sous H_0 : p_0=0.75 -> rejet normalement
8 prop.test(sum(X),n,p = 0.75)

```

⊥

7.3 Test de comparaison de deux probabilités pour deux échantillons appariés

Dans cette section, nous supposons avoir un n -échantillon apparié $(X_i, Y_i)_{1 \leq i \leq n}$ prenant leurs valeurs dans l'ensemble $\{0, 1\}^2$ suivant des lois de Bernoulli $\mathcal{B}(p_X^*)$ et $\mathcal{B}(p_Y^*)$ et symbolisant, par exemple, la présence d'un caractère dans deux expériences. Par exemple, si $X_i = 0$ et $Y_i = 0$, cela signifie que l'individu i n'a pas eu la caractéristique dans aucune des expériences. La question que nous pouvons nous poser est l'égalité des probabilités :

$$\begin{cases} \mathcal{H}_0 : p_X^* = p_Y^*, \\ \mathcal{H}_1 : p_X^* \neq p_Y^*. \end{cases} \quad (7.1)$$

Nous avons donc 4 possibilités pour les paires (X_i, Y_i) et nous comptons le nombre d'individus de chaque paire :

		Y	
		0	1
X	0	$n_{0,0}$	$n_{0,1}$
	1	$n_{1,0}$	$n_{1,1}$

En particulier, nous voyons que pour estimer p_X^* et p_Y^* , nous pouvons compter le nombre moyen de fois où X (resp. Y) vaut 1 :

$$\hat{p}_{X,n} = \frac{n_{1,1} + n_{1,0}}{n} \quad \text{et} \quad \hat{p}_{Y,n} = \frac{n_{1,1} + n_{0,1}}{n}.$$

En particulier, nous voyons que la différence entre les deux estimateurs se fait uniquement sur les valeurs de $n_{1,0}$ et $n_{0,1}$. Ainsi, le test (7.1) est équivalent au test :

$$\begin{cases} \mathcal{H}_0 : p_{0,1}^{*(X,Y)} = p_{1,0}^{*(X,Y)}, \\ \mathcal{H}_1 : p_{0,1}^{*(X,Y)} \neq p_{1,0}^{*(X,Y)} \end{cases}$$

où $p_{0,1}^{*(X,Y)}$ et $p_{1,0}^{*(X,Y)}$ sont les probabilités croisées. En théorie, si ces probabilités sont égales, cela signifie que $n_{1,0} \approx n_{0,1}$ et donc, si nous notons $n_{1/0} = n_{1,0} + n_{0,1}$, elles valent chacune la moitié. Autrement dit, nous pouvons résumer par le tableau suivant :

		(X, Y)		Total
		(1, 0)	(0, 1)	
Observé	$n_{1,0}$	$n_{0,1}$	$n_{1/0}$	
Théorique	$\frac{n_{1/0}}{2}$	$\frac{n_{1/0}}{2}$	$n_{1/0}$	

L'idée est donc de regarder la différence entre les effectifs observés et les théoriques. Nous obtenons ainsi :

Point méthode (Test de Mc Nemar)

Étant donné deux échantillons appariés (X, Y) ayant chacune une loi de Bernoulli $\mathcal{B}(p_X^*)$ et $\mathcal{B}(p_Y^*)$ alors la statistique de test est

$$mcN_n = \frac{(n_{1,0} - \frac{n_{1/0}}{2})^2}{\frac{n_{1/0}}{2}} = \frac{(n_{1,0} - n_{0,1})^2}{n_{1,0} + n_{0,1}}$$

qui a pour loi **asymptotique** une loi du Chi2 $\chi^2(1)$ à 1 degrés de liberté sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau **asymptotique** $\alpha \in]0; 1[$ est donc :

$$\begin{array}{l|l} \mathcal{H}_1 & p_X^* \neq p_Y^* \\ \Theta_1 &]0; 1[^2 \setminus \{(x, y) \in]0; 1[^2 \mid x = y\} \\ \text{Règle} & mcN_n > s_{1-\alpha}^{(1)} \end{array}$$

où $s_\alpha^{(1)}$ est le quantile d'ordre α d'une loi du Chi2 $\chi^2(1)$ à 1 degrés de liberté ; c'est-à-dire que si S suit une loi du Chi2 $\chi^2(1)$ à 1 degrés de liberté alors $s_\alpha^{(1)}$ est tel que :

$$\mathbb{P}(S \leq s_\alpha^{(1)}) = \alpha.$$



Codage en R

La fonction `R` permettant de faire le test bilatéral est `mcnemar.test`. Par exemple :

```

1 ## Simulation d'un jeu de données
2 p_0<-0.5
3 n<-100
4 tab<-table(rbinom(n,1,p_0),rbinom(n,1,p_0))
5 ## Test sous H_0 : p_X=p_Y -> acceptation avec 95% de chances
6 mcnemar.test(tab)
7 ## Test sous H_0 : p_X= p_Y -> rejet normalement
8 mcnemar.test(table(rbinom(n,1,p_0),rbinom(n,1,0.9)))

```

7.4 Test de comparaison de deux probabilités pour deux échantillons indépendants

Dans le cas de deux échantillons \mathbf{X} et \mathbf{Y} indépendants de tailles $n_{\mathbf{X}}$ et $n_{\mathbf{Y}}$, le fait de comparer les deux probabilités $p_{\mathbf{X}}^*$ et $p_{\mathbf{Y}}^*$ comme dans le test (7.1) revient à tester la différence des deux probabilités :

$$\begin{cases} \mathcal{H}_0 : p_{\mathbf{X}}^* - p_{\mathbf{Y}}^* = 0, \\ \mathcal{H}_1 : p_{\mathbf{X}}^* - p_{\mathbf{Y}}^* \neq 0. \end{cases}$$

Or, par analogie, les deux estimateurs sont :

$$\hat{p}_{\mathbf{X},n} = \frac{1}{n_{\mathbf{X}}} \sum_{i=1}^{n_{\mathbf{X}}} X_i = \frac{\mathbf{X}_+}{n_{\mathbf{X}}} \text{ et } \hat{p}_{\mathbf{Y},n} = \frac{\mathbf{Y}_+}{n_{\mathbf{Y}}},$$

nous voyons que si les probabilités sont égales, c'est à dire qu'elles valent p^* , alors un estimateur est :

$$\hat{p}_{n,\mathbf{X},\mathbf{Y}} = \frac{\mathbf{X}_+ + \mathbf{Y}_+}{n_{\mathbf{X}} + n_{\mathbf{Y}}}.$$

A partir de ces notations, nous avons donc :

Point méthode (Test d'égalité des proportions ; échantillons indépendants)

Étant donné deux échantillons indépendants \mathbf{X} et \mathbf{Y} de tailles $n_{\mathbf{X}}$ et $n_{\mathbf{Y}}$ ayant chacune une loi de Bernoulli $\mathcal{B}(p_{\mathbf{X}}^*)$ et $\mathcal{B}(p_{\mathbf{Y}}^*)$ alors la statistique de test est

$$S_{n_{\mathbf{X}},n_{\mathbf{Y}}} = \frac{\hat{p}_{\mathbf{X},n} - \hat{p}_{\mathbf{Y},n}}{\sqrt{\hat{p}_{n,\mathbf{X},\mathbf{Y}}(1 - \hat{p}_{n,\mathbf{X},\mathbf{Y}}) \left(\frac{1}{n_{\mathbf{X}}} + \frac{1}{n_{\mathbf{Y}}} \right)}}$$

qui a pour loi **asymptotique** une loi normale $\mathcal{N}(0, 1)$ centrée réduite sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau **asymptotique** $\alpha \in]0, 1[$ est donc :

\mathcal{H}_1	$p_{\mathbf{X}}^* \neq p_{\mathbf{Y}}^*$	$p_{\mathbf{X}}^* < p_{\mathbf{Y}}^*$	$p_{\mathbf{X}}^* > p_{\mathbf{Y}}^*$
Θ_1	$]0, p_0[\cup]p_0, 1[$	$]0, p_0[$	$]p_0, 1[$
Règle	$ S_{n_{\mathbf{X}},n_{\mathbf{Y}}} > q_{1-\alpha/2}$	$S_{n_{\mathbf{X}},n_{\mathbf{Y}}} < q_{\alpha}$	$S_{n_{\mathbf{X}},n_{\mathbf{Y}}} > q_{1-\alpha}$

où q_{α} est le quantile d'ordre α d'une loi normale centrée réduite ; c'est-à-dire que si X suit une loi gaussienne centrée réduite $\mathcal{N}(0, 1)$ alors q_{α} est tel que :

$$\mathbb{P}(X \leq q_{\alpha}) = \alpha.$$

Codage en R

La fonction  permettant de faire le test bilatéral est `prop.test`. Par exemple :

```

1 ## Simulation d'un jeu de données
2 nX<-100
3 nY<-200
4 X<-rbinom(nX,1,0.5)
5 Y<-rbinom(nY,1,0.5)
6 ## Test sous H_0 : p_X=p_Y -> acceptation avec 95% de chances
7 prop.test(matrix(c(sum(X),sum(Y),nX,nY),2,2))
8 ## Test sous H_0 : p_X= p_Y -> rejet normalement
9 prop.test(matrix(c(sum(X),sum(rbinom(nY,1,0.75)),nX,nY),2,2))

```


7.5 Généralisation à plus de deux échantillons appariés

Dans le cas de K échantillons appariés $\mathbf{X} = (X_{i,k})$ dont nous cherchons à savoir si toutes les proportions sont égales (c'est-à-dire que, qu'importe l'expérience, nous retrouvons la même proportion), nous avons le test :

$$\begin{cases} \mathcal{H}_0 : p_1^* = p_2^* = \dots = p_K^*, \\ \mathcal{H}_1 : \exists k \neq k', p_k^* \neq p_{k'}^*. \end{cases}$$

Dans ce cas, le test utilisé est une généralisation du test de Mc Nemar appelé test de Cochran-Mantel-Haenszel.

Codage en R

La fonction  permettant de faire le test bilatéral est `mantelhaen.test`. Par exemple :

```

1 ## Simulation d'un jeu de données
2 X<-rbinom(100,1,0.5)
3 Y<-rbinom(100,1,0.5)
4 Z<-rbinom(100,1,0.5)
5 ## Test sous H_0 : p_X=p_Y -> acceptation avec 95% de chances
6 mantelhaen.test(table(X,Y,Z))
7 ## Test sous H_0 : p_X= p_Y -> rejet normalement
8 mantelhaen.test(table(X,rbinom(100,1,0.1),rbinom(100,1,0.9)))

```

7.6 Test de comparaison de plus de deux probabilités pour des échantillons indépendants

Enfin, nous nous intéressons au cas où un individu peut prendre une caractéristique A_i avec $i \in \{1, \dots, I\}$ dans une population M_j où $j \in \{1, \dots, J\}$ et nous nous intéressons au fait que la probabilité $p_{i,j}^*$ d'avoir cette caractéristique dans la population M_j soit invariable qu'importe la population. Autrement dit, pour tout couple (j, j') de $\{1, \dots, J\}^2$, $p_{i,j}^*$ est égale à $p_{i,j'}^*$. De plus, nous ajoutons l'hypothèse que cette propriété est vraie qu'importe la caractéristique. Au final, nous avons donc :

$$\begin{cases} \mathcal{H}_0 : \forall i \in \{1, \dots, I\}, \forall (j, j') \in \{1, \dots, J\}^2, p_{i,j}^* = p_{i,j'}^*, \\ \mathcal{H}_1 : \exists i \in \{1, \dots, I\} \text{ et } \exists (j, j') \in \{1, \dots, J\}^2, p_{i,j}^* \neq p_{i,j'}^*. \end{cases}$$

Comme précédemment, nous regardons les effectifs théoriques si nous avons l'indépendance et les comparons aux effectifs observés. Pour cela, nous avons besoin des notations suivantes :

		Population				Total
		M_1	M_2	\dots	M_J	
Caractéristique	A_1	$n_{1,1}$	$n_{1,2}$	\dots	$n_{1,J}$	$n_{1,\cdot} = \sum_{j=1}^J n_{1,j}$
	A_2	$n_{2,1}$	$n_{2,2}$	\dots	$n_{2,J}$	$n_{2,\cdot} = \sum_{j=1}^J n_{2,j}$
	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
	A_I	$n_{I,1}$	$n_{I,2}$	\dots	$n_{I,J}$	$n_{I,\cdot} = \sum_{j=1}^J n_{I,j}$
Total	$n_{\cdot,1} = \sum_{i=1}^I n_{i,1}$	$n_{\cdot,2} = \sum_{i=1}^I n_{i,2}$	\dots	$n_{\cdot,J} = \sum_{i=1}^I n_{i,J}$	$n = \sum_{i=1}^I \sum_{j=1}^J n_{i,j}$	

et nous avons donc :

Point méthode (Test de Chi2 d'indépendance de Pearson)

Étant donné K échantillons indépendants ayant chacune une loi de Bernoulli $\mathcal{B}(p_k^*)$ alors la statistique de test est

$$\chi_n^2 = \sum_{i=1}^I \sum_{j=1}^J \frac{(n_{i,j} - c_{i,j})^2}{c_{i,j}} \text{ avec } c_{i,j} = \frac{n_{i,\cdot} \cdot n_{\cdot,j}}{n}$$

qui a pour loi **asymptotique** une loi du Chi2 $\chi^2 [(I-1)(J-1)]$ à $(I-1)(J-1)$ degrés de liberté sous l'hypothèse \mathcal{H}_0 et la règle de rejet de niveau **asymptotique** $\alpha \in]0; 1[$ est donc :

$$\frac{\mathcal{H}_1}{\text{Règle}} \parallel \frac{\exists i \in \{1, \dots, I\} \text{ et } \exists (j, j') \in \{1, \dots, J\}^2, p_{i,j}^* \neq p_{i,j'}^*}{\chi_n^2 > s_{1-\alpha}^{[(I-1)(J-1)]}}$$

où $s_{\alpha}^{[(I-1)(J-1)]}$ est le quantile d'ordre α d'une loi du Chi2 $\chi^2(1)$ à $(I-1)(J-1)$ degrés de liberté; c'est-à-dire que si S suit une loi du Chi2 $\chi^2 [(I-1)(J-1)]$ à $(I-1)(J-1)$ degrés de liberté alors $s_{\alpha}^{[(I-1)(J-1)]}$ est tel que :

$$\mathbb{P}(S \leq s_{\alpha}^{[(I-1)(J-1)]}) = \alpha.$$

Codage en R



La fonction permettant de faire le test unilatéral est `chisq.test`.

Chapitre 8

Tests non-paramétriques

8.1 Test de Kolmogorov-Smirnov

Le deuxième test d'adéquation est celui de Kolmogorov-Smirnov et se base cette fois sur la fonction de répartition. Le principe est d'estimer cette dernière.

Définition 67 (Fonction de répartition empirique)

Étant donné un n -échantillon \mathbf{X} de variables réelles, nous appelons **fonction de répartition empirique** la fonction définie sur \mathbb{R} par :

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \leq x\}}.$$

Le principe du test est de calculer le plus grand écart entre la fonction de répartition recherchée et la fonction de répartition empirique.

Proposition 27 (Test de Kolmogorov-Smirnov)

Étant donné un n -échantillon \mathbf{X} de variable réelle et une loi de référence de fonction de répartition $F^{(0)}$, alors nous définissons la statistique de test

$$D_n = \sqrt{n} \sup_{x \in \mathbb{R}} |F^{(0)}(x) - F_n(x)|$$

et, sous l'hypothèse \mathcal{H}_0 , nous avons :

$$D_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \sup_{x \in \mathbb{R}} |B_{F^{(0)}(x)}|$$

où $(B_t)_{t \in [0;1]}$ est un pont brownien ; c'est-à-dire un mouvement brownien conditionné par le fait que $B_0 = B_1 = 0$.

En particulier, dans le cas d'une fonction de répartition continue, nous avons pour toute constante c strictement positive :

$$\mathbb{P}(D_n > c) \xrightarrow[n \rightarrow +\infty]{} \alpha(c) = 2 \sum_{r=1}^{+\infty} (-1)^{r-1} e^{-2c^2 r^2}.$$

Remarque

Nous remarquons qu'il existe une formule explicite pour calculer la p -valeur même si nous ne devons nous contenter d'approximations numériques dans la pratique. En particulier, $\alpha(1,36) \approx 0,05$. Sur la figure 9.1, nous avons représenté la fonction $c \mapsto \alpha(c)$.

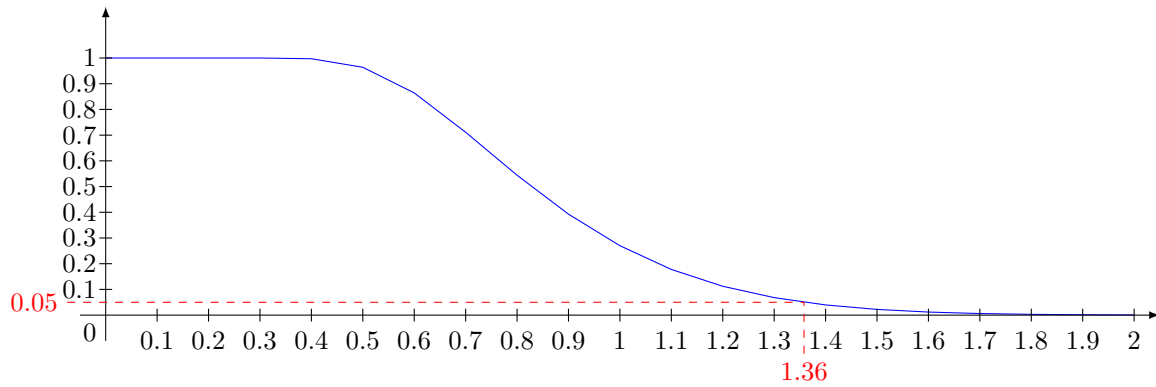


FIGURE 8.1 – Représentation de la fonction $c \mapsto \alpha(c)$ qui représente la probabilité $\mathbb{P}\left(\sup_{t \in [0;1]} |B_t| \geq c\right)$ (donc la p -valeur). La ligne hachée rouge permet de retrouver la valeur correspondante à une probabilité de 0.05.

Corollaire 28 (Décision)

À partir de la proposition 36, nous pouvons rejeter l'hypothèse nulle si la statistique D_n dépasse le quantile d'ordre $1 - \alpha$ de la loi présentée dans la proposition.

Chapitre 9

Autres tests

9.1 Botanique des tests

Dans cette partie, nous mettons quelques tests classiques.

9.1.1 Tests paramétriques classiques

Dans cette section, nous mettons des tests classiques.

Test de proportion

Dans cette partie, nous testons la proportion d'une quantité (par exemple, le nombre d'étudiants du M1 de Maths Générale aimant la statistique) par rapport à une quantité fixée. Comme cette question peut être ramenée à une loi de Bernoulli, nous supposons avoir un n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$ de loi de Bernoulli $\mathcal{B}(p^*)$, une valeur p_0 de référence et nous regardons l'un des tests suivants :

$$\begin{cases} \mathcal{H}_0 : p^* = p_0, \\ \mathcal{H}_1 : p^* \neq p_0 \text{ ou } p^* < p_0 \text{ ou } p^* > p_0 \text{ (choix à fixer suivant le contexte)}. \end{cases}$$

Notons que nous n'abordons pas le cas $\mathcal{H}_1 : p^* = p_1$ car le théorème 20 montre que le test de vraisemblance est uniformément plus puissant.

Pour la suite, nous présentons le cas où l'hypothèse alternative est $p^* \neq p_0$, la démarche étant similaire pour les autres. Pour estimer le paramètre p^* , nous utilisons la moyenne empirique qui est l'estimateur du maximum de vraisemblance et celui des moments basé sur l'espérance. Par le théorème de la limite centrale, nous avons, sous l'hypothèse \mathcal{H}_0 , :

$$\sqrt{n} (\bar{\mathbf{X}}_n - p_0) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, p_0(1 - p_0)).$$

Ainsi, la statistique la plus logique est :

$$T_n = \sqrt{n} \frac{\bar{\mathbf{X}}_n - p_0}{\sqrt{p_0(1 - p_0)}}$$

qui suit asymptotiquement une loi normale centrée réduite.

Corollaire 29 (Décision)

À partir des propriétés que nous venons de voir, nous pouvons rejeter l'hypothèse nulle si la statistique T_n dépasse le quantile d'ordre $1 - \alpha/2$ ou est plus petite que le quantile d'ordre $\alpha/2$ de la loi gaussienne centrée réduite.

Test des moyennes

Dans cette partie, nous nous plaçons dans le cas d'un n -échantillon avec un moment d'ordre 2 et nous cherchons à savoir si son moment d'ordre 1 μ^* est égal à une valeur connue ou pas. Ainsi, nous regardons l'un des tests suivants :

$$\begin{cases} \mathcal{H}_0 : \mu^* = \mu_0, \\ \mathcal{H}_1 : \mu^* \neq \mu_0 \text{ ou } \mu^* < \mu_0 \text{ ou } \mu^* > \mu_0 \text{ (choix à fixer suivant le contexte)}. \end{cases}$$

Comme pour le test des proportions, une estimation de la moyenne se fait grâce à la moyenne empirique $\bar{\mathbf{X}}_n$ et, sous l'hypothèse, \mathcal{H}_0 , nous avons :

$$\sqrt{n} (\bar{\mathbf{X}}_n - \mu_0) \underset{n \rightarrow +\infty}{\xrightarrow{\mathcal{L}}} \mathcal{N}(0, \mathbb{V}[X_1]).$$

Si nous connaissons la variance, nous pouvons reprendre la méthode précédente. Nous nous plaçons donc dans le cas où la variance est inconnue.

Hypothèse

Dans la suite de cette partie, nous supposons que la variance est inconnue.

Comme nous ne connaissons pas la moyenne, nous l'estimons par l'estimateur empirique :

$$S_n^2 = \overline{X_n^2} - \bar{\mathbf{X}}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{\mathbf{X}}_n)^2.$$

Nous commençons par calculer la loi de cet estimateur.

Proposition 30 (Loi de S_n^2)

Étant donné un n -échantillon \mathbf{X} admettant un moment d'ordre 2 alors nS_n^2/σ (où σ est l'écart-type commun) suit une loi du $\chi^2(n-1)$; c'est-à-dire une loi du χ^2 à $n-1$ degrés de liberté.

Preuve

La preuve se base sur un théorème que nous admettrons car il utilise des notions de projections (donc d'espérances conditionnelles) qui ne seront abordées que durant le deuxième semestre. Nous mettons toutefois l'énoncé ici.

Lemme 31 (Théorème de Cochran)

Étant donné un vecteur gaussien \mathbf{X} de \mathbb{R}^d de loi $\mathcal{N}(m, \sigma^2 \mathbb{I}_d)$ et F_1, \dots, F_P des sous-espaces vectoriels de \mathbb{R}^d , orthogonaux deux à deux et de somme \mathbb{R}^d dont nous notons P_{F_p} les matrices de projection orthogonale sur F_p et n_p la dimension de F_p alors nous avons :

- les vecteurs aléatoires $P_{F_1}\mathbf{X}, \dots, P_{F_P}\mathbf{X}$ sont deux à deux indépendants et de lois respectives $\mathcal{N}(P_{F_1}m, \sigma^2 P_{F_1}), \dots, \mathcal{N}(P_{F_P}m, \sigma^2 P_{F_P})$.
- les variables aléatoires $\frac{\|P_{F_1}(\mathbf{X}-m)\|^2}{\sigma^2}, \dots, \frac{\|P_{F_P}(\mathbf{X}-m)\|^2}{\sigma^2}$ sont deux à deux indépendantes et sont de lois respectives $\chi^2(n_1), \dots, \chi^2(n_P)$.

Pour démontrer la proposition 30, nous commençons par définir pour tout $i \in \{1, \dots, n\}$, la variable :

$$Z_i = \frac{X_i - \bar{\mathbf{X}}_n}{\sigma}.$$

Calculons l'espérance et la matrice de variance-covariance de la variable $\mathbf{Z} = (Z_1, \dots, Z_n)$. Sous l'hypothèse \mathcal{H}_0 , nous avons de plus pour tout $i \in \{1, \dots, n\}$:

$$\begin{aligned} \mathbb{E}[Z_i] &= \mathbb{E}\left[\frac{X_i - \bar{\mathbf{X}}_n}{\sigma}\right] \\ &= \frac{\mathbb{E}[X_i] - \mathbb{E}[\bar{\mathbf{X}}_n]}{\sigma} \\ &= \frac{\mu^* - \mu^*}{\sigma} \\ &= 0. \end{aligned}$$



De plus, pour tout $i \in \{1, \dots, n\}$, nous avons :

$$\begin{aligned}
 \mathbb{V} \left[\frac{X_i - \bar{\mathbf{X}}_n}{\sigma} \right] &= \frac{1}{\sigma^2} \mathbb{V} \left[X_i \frac{1}{n} \sum_{j=1}^n 1 - \frac{1}{n} \sum_{j=1}^n X_j \right] \\
 &= \frac{1}{n^2 \sigma^2} \mathbb{V} \left[\sum_{j=1}^n (X_i - X_j) \right] \\
 &= \frac{1}{n^2 \sigma^2} \mathbb{V} \left[\sum_{\substack{j=1 \\ j \neq i}}^n (X_i - X_j) \right] \\
 &= \frac{1}{n^2 \sigma^2} \left(\sum_{\substack{j=1 \\ j \neq i}}^n \mathbb{V}[X_i - X_j] + \sum_{\substack{j_1=1 \\ j_1 \neq i}}^n \sum_{\substack{j_2=1 \\ j_2 \neq j_1 \text{ et } j_2 \neq i}}^n \text{Cov}(X_i - X_{j_1}, X_i - X_{j_2}) \right) \\
 &= \frac{1}{n^2 \sigma^2} \left[\sum_{\substack{j=1 \\ j \neq i}}^n \left(\mathbb{V}[X_i] + \mathbb{V}[X_j] - \underbrace{2 \text{Cov}(X_i, X_j)}_{=0} \right) \right. \\
 &\quad \left. + \sum_{\substack{j_1=1 \\ j_1 \neq i}}^n \sum_{\substack{j_2=1 \\ j_2 \neq j_1 \text{ et } j_2 \neq i}}^n \left(\mathbb{V}[X_i] - \underbrace{\text{Cov}(X_{j_1}, X_i)}_{=0} - \underbrace{\text{Cov}(X_i, X_{j_2})}_{=0} + \underbrace{\text{Cov}(X_{j_1}, X_{j_2})}_{=0} \right) \right] \\
 &= \frac{1}{n^2 \sigma^2} \left(\sum_{\substack{j=1 \\ j \neq i}}^n 2\sigma^2 + \sum_{\substack{j_1=1 \\ j_1 \neq i}}^n \sum_{\substack{j_2=1 \\ j_2 \neq j_1 \text{ et } j_2 \neq i}}^n \sigma^2 \right) \\
 &= \frac{1}{n^2 \sigma^2} [2(n-1)\sigma^2 + (n-1)(n-2)\sigma^2] \\
 &= \frac{n(n-1)\sigma^2}{n^2 \sigma^2} \\
 &= 1 - \frac{1}{n}.
 \end{aligned}$$

Enfin, pour tout $1 \leq i_1 \neq i_2 \leq n$, nous avons :

$$\begin{aligned}
 &\text{Cov} \left(\frac{X_{i_1} - \bar{\mathbf{X}}_n}{\sigma}, \frac{X_{i_2} - \bar{\mathbf{X}}_n}{\sigma} \right) \\
 &= \frac{1}{\sigma^2} \text{Cov}(X_{i_1} - \bar{\mathbf{X}}_n, X_{i_2} - \bar{\mathbf{X}}_n) \\
 &= \frac{1}{\sigma^2} \left[\underbrace{\text{Cov}(X_{i_1}, X_{i_2})}_{=0} - \text{Cov}(X_{i_1}, \bar{\mathbf{X}}_n) - \text{Cov}(\bar{\mathbf{X}}_n, X_{i_2}) + \text{Cov}(\bar{\mathbf{X}}_n, \bar{\mathbf{X}}_n) \right] \\
 &= \frac{1}{\sigma^2} \left[-\frac{1}{n} \sum_{j=1}^n \text{Cov}(X_{i_1}, X_j) - \frac{1}{n} \sum_{j=1}^n \text{Cov}(X_j, X_{i_2}) + \mathbb{V}[\bar{\mathbf{X}}_n] \right] \\
 &= \frac{1}{\sigma^2} \left(-\frac{1}{n} \mathbb{V}[X_{i_1}] - \frac{1}{n} \mathbb{V}[X_{i_2}] + \frac{1}{n^2} \mathbb{V} \left[\sum_{j=1}^n X_j \right] \right) \\
 &= \frac{1}{n\sigma^2} \left[-2\sigma^2 + \frac{1}{n} \left(\sum_{j=1}^n \mathbb{V}[X_j] + \sum_{\substack{j_1=1 \\ j_2=1 \\ j_2 \neq j_1}}^n \sum_{\substack{j_2=1 \\ j_2 \neq j_1}}^n \underbrace{\text{Cov}(X_{j_1}, X_{j_2})}_{=0} \right) \right]
 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n\sigma^2} \left(-2\sigma^2 + \frac{n\sigma^2}{n} \right) \\
&= \frac{-\sigma^2}{n\sigma^2} \\
&= \frac{-1}{n}
\end{aligned}$$

Donc la loi commune des Z_i est centrée et matrice de variance covariance :

$$\Gamma_n = \mathbb{I}_n - \begin{pmatrix} 1/\sqrt{n} \\ \vdots \\ 1/\sqrt{n} \end{pmatrix} (1/\sqrt{n} \quad \cdots \quad 1/\sqrt{n}).$$

Par le théorème multivarié de la limite centrale, la loi de la variable $\frac{Z_1 + \dots + Z_n}{n}$ converge vers loi normale centrée de covariance Γ_n . Or, cette loi est le projeté d'un vecteur aléatoire de \mathbb{R}^K suivant une loi normale centrée réduite sur l'hyperplan orthogonal au vecteur colonne $(1/\sqrt{n}, \dots, 1/\sqrt{n})^T$.

De plus, nous voyons que la variable nS_n^2/σ est la norme au carré de cette variable donc, d'après le théorème de Cochran, la loi asymptotique est une loi du χ^2 à $n - 1$ degrés de libertés.

Corollaire 32 (Loi de la statistique)

Étant donné un n -échantillon \mathbf{X} admettant un moment d'ordre 2 alors

$$T_n = \sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu_0}{\sqrt{S_n^2}}$$

suit asymptotiquement une loi de Student centrée à $n - 1$ degrés de liberté.



Preuve

Nous avons :

$$\begin{aligned}
\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu_0}{\sqrt{S_n^2}} &= \frac{\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu_0}{\sigma}}{\frac{\sqrt{S_n^2}}{\sigma}} \\
&= \frac{\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu_0}{\sigma}}{\frac{1}{\sqrt{n}} \sqrt{\frac{nS_n^2}{\sigma^2}}} \\
&= \frac{1}{\sqrt{\frac{n}{n-1}}} \times \frac{\sqrt{n} \frac{\bar{\mathbf{X}}_n - \mu_0}{\sigma}}{\sqrt{\frac{nS_n^2}{\sigma^2(n-1)}}} \\
&= \underbrace{\frac{1}{\sqrt{\frac{n-1}{n}}}}_{\rightarrow 1 \text{ as } n \rightarrow +\infty} \times \frac{Z}{\sqrt{\frac{Y}{n-1}}}
\end{aligned}$$

avec Z qui suit asymptotiquement une loi normale centrée réduite et Y qui suit une loi du χ^2 à $n - 1$ degrés de liberté donc, par définition de la loi de Student, nous avons le résultat.

Corollaire 33 (Décision)

À partir des propriétés que nous venons de voir, nous pouvons rejeter l'hypothèse nulle si la statistique T_n dépasse le quantile d'ordre $1 - \alpha/2$ ou est plus petite que le quantile d'ordre $\alpha/2$ de la loi Student centrée de $n - 1$ degrés de libertés.

9.1.2 Tests d'adéquation d'une loi de probabilité à des données

Dans cette partie, nous intéressons au cas où le n -échantillon \mathbf{X} suit une loi de référence de densité f_0 . Le test effectué est donc le suivant :

$$\begin{cases} \mathcal{H}_0 : & \text{La loi de } \mathbf{X} \text{ a pour densité } f_0, \\ \mathcal{H}_1 : & \text{La loi de } \mathbf{X} \text{ n'a pas pour densité } f_0. \end{cases}$$

Nous présentons ici deux tests qui permettent de répondre à cette question.

Test du χ^2 d'adéquation

Dans cette partie, nous présentons essentiellement le cas où les valeurs de \mathbf{X} sont dans un ensemble fini $I = \{a_1, \dots, a_K\}$.

Définitions 68 (Loi multinomiale)

Étant donnée un n -échantillon \mathbf{X} prenant ses valeurs dans un ensemble fini $I = \{a_1, \dots, a_K\}$, nous notons p_k la probabilité que la variable X_1 soit égale à a_k pour tout k compris entre 1 et K :

$$p_k = \mathbb{P}(X_1 = a_k).$$

Pour tout $k \in \{1, \dots, K\}$, nous notons N_k la variable aléatoire qui compte le nombre de variables X_i qui ont pris la modalité a_k :

$$N_k = \sum_{i=1}^n \mathbb{1}_{\{X_i = a_k\}}.$$

Nous disons alors que la variable $\mathbf{N} = (N_1, \dots, N_K)$ suit une **loi multinomiale de paramètres n et (p_1, \dots, p_K)** , notée $\mathcal{M}(n; p_1, \dots, p_K)$, dont la densité est :

$$\mathbb{P}(N_1 = n_1, \dots, N_K = n_k) = \frac{n!}{\prod_{k=1}^K n_k!} \left(\prod_{k=1}^K p_k \right) \mathbb{1}_{\{\sum_{k=1}^K n_k = n\}}.$$

Remarque

L'ensemble des valeurs possibles de la variable \mathbf{N} est donc l'ensemble des valeurs discrètes d'un simplexe :

$$\left\{ (n_1, \dots, n_K) \in \{0, \dots, n\}^K \mid \sum_{k=1}^K n_k = n \right\}.$$

Comme une loi discrète est entièrement caractérisée par les probabilités de valoir chaque modalité, le test revient à savoir si :

$$\begin{cases} \mathcal{H}_0 : & \forall k \in \{1, \dots, K\}, p_k^* = p_k^{(0)}, \\ \mathcal{H}_1 : & \exists k \in \{1, \dots, K\}, p_k^* \neq p_k^{(0)}. \end{cases}$$

Proposition 34 (Test du χ^2 d'adéquation)

Étant donné un n -échantillon de loi définie sur un ensemble discret de probabilité p_k^* va valoir la modalité a_k et une loi de référence de probabilité $p_k^{(0)}$ de valoir la modalité a_k , alors nous définissons la statistique de test

$$T_n = \sum_{k=1}^K \frac{(N_k - np_k^{(0)})^2}{np_k^{(0)}}$$

et, sous l'hypothèse \mathcal{H}_0 , nous avons :

$$T_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi^2(K-1)$$

c'est-à-dire une loi du χ^2 à $K - 1$ degrés de liberté.



Preuve

Pour démontrer la proposition 34, nous commençons par définir pour tout $i \in \{1, \dots, n\}$, la variable :

$$Z_i = \left(\frac{\mathbb{1}_{\{X_i=a_1\}} - p_1^{(0)}}{\sqrt{p_1^{(0)}}}, \dots, \frac{\mathbb{1}_{\{X_i=a_K\}} - p_K^{(0)}}{\sqrt{p_K^{(0)}}} \right).$$

Calculons l'espérance et la matrice de variance-covariance de la variable $\mathbf{Z} = (Z_1, \dots, Z_n)$. Sous l'hypothèse \mathcal{H}_0 , nous avons de plus pour tout $i \in \{1, \dots, n\}$:

$$\begin{aligned} \mathbb{E}[Z_i] &= \mathbb{E} \left[\left(\frac{\mathbb{1}_{\{X_i=a_1\}} - p_1^{(0)}}{\sqrt{p_1^{(0)}}}, \dots, \frac{\mathbb{1}_{\{X_i=a_K\}} - p_K^{(0)}}{\sqrt{p_K^{(0)}}} \right) \right] \\ &= \left(\mathbb{E} \left[\frac{\mathbb{1}_{\{X_i=a_1\}} - p_1^{(0)}}{\sqrt{p_1^{(0)}}} \right], \dots, \mathbb{E} \left[\frac{\mathbb{1}_{\{X_i=a_K\}} - p_K^{(0)}}{\sqrt{p_K^{(0)}}} \right] \right) \\ &= \left(\frac{\mathbb{E}[\mathbb{1}_{\{X_i=a_1\}}] - p_1^{(0)}}{\sqrt{p_1^{(0)}}}, \dots, \frac{\mathbb{E}[\mathbb{1}_{\{X_i=a_K\}}] - p_K^{(0)}}{\sqrt{p_K^{(0)}}} \right) \\ &= (0, \dots, 0). \end{aligned}$$

De plus, pour tout $k \in \{1, \dots, K\}$, nous avons :

$$\begin{aligned} \mathbb{V} \left[\frac{\mathbb{1}_{\{X_i=a_k\}} - p_k^{(0)}}{\sqrt{p_k^{(0)}}} \right] &= \frac{\mathbb{V}[\mathbb{1}_{\{X_i=a_k\}}]}{p_k^{(0)}} \\ &= \frac{p_k^{(0)}(1 - p_k^{(0)})}{p_k^{(0)}} \\ &= 1 - p_k^{(0)}. \end{aligned}$$

Enfin, pour tout $1 \leq k \neq \ell \leq K$, nous avons :

$$\begin{aligned} \text{Cov} \left(\frac{\mathbb{1}_{\{X_i=a_k\}} - p_k^{(0)}}{\sqrt{p_k^{(0)}}}, \frac{\mathbb{1}_{\{X_i=a_\ell\}} - p_\ell^{(0)}}{\sqrt{p_\ell^{(0)}}} \right) &= \frac{1}{\sqrt{p_k^{(0)} p_\ell^{(0)}}} \text{Cov} \left(\mathbb{1}_{\{X_i=a_k\}} - p_k^{(0)}, \mathbb{1}_{\{X_i=a_\ell\}} - p_\ell^{(0)} \right) \\ &= \frac{1}{\sqrt{p_k^{(0)} p_\ell^{(0)}}} \left(\mathbb{E} \left[\left(\mathbb{1}_{\{X_i=a_k\}} - p_k^{(0)} \right) \left(\mathbb{1}_{\{X_i=a_\ell\}} - p_\ell^{(0)} \right) \right] - \underbrace{\mathbb{E} \left[\mathbb{1}_{\{X_i=a_k\}} - p_k^{(0)} \right]}_{=0} \underbrace{\mathbb{E} \left[\mathbb{1}_{\{X_i=a_\ell\}} - p_\ell^{(0)} \right]}_{=0} \right) \\ &= \frac{1}{\sqrt{p_k^{(0)} p_\ell^{(0)}}} \mathbb{E} \left[\mathbb{1}_{\{X_i=a_k\}} \mathbb{1}_{\{X_i=a_\ell\}} - p_k^{(0)} \mathbb{1}_{\{X_i=a_\ell\}} - \mathbb{1}_{\{X_i=a_k\}} p_\ell^{(0)} + p_k^{(0)} p_\ell^{(0)} \right] \\ &= \frac{1}{\sqrt{p_k^{(0)} p_\ell^{(0)}}} \left(\underbrace{\mathbb{E} \left[\mathbb{1}_{\{X_i=a_k\}} \mathbb{1}_{\{X_i=a_\ell\}} \right]}_{=0 \text{ car } k \neq \ell} - p_k^{(0)} \mathbb{E} \left[\mathbb{1}_{\{X_i=a_\ell\}} \right] - \mathbb{E} \left[\mathbb{1}_{\{X_i=a_k\}} \right] p_\ell^{(0)} + p_k^{(0)} p_\ell^{(0)} \right) \\ &= \frac{-p_k^{(0)} p_\ell^{(0)}}{\sqrt{p_k^{(0)} p_\ell^{(0)}}} \end{aligned}$$

$$= -\sqrt{p_k^{(0)} p_\ell^{(0)}}.$$

Donc la loi commune des Z_i est centrée et matrice de variance covariance :

$$\Gamma = \mathbb{I}_K - \begin{pmatrix} \sqrt{p_1^{(0)}} \\ \vdots \\ \sqrt{p_K^{(0)}} \end{pmatrix} \begin{pmatrix} \sqrt{p_1^{(0)}} & \cdots & \sqrt{p_K^{(0)}} \end{pmatrix}.$$

Par le théorème multivarié de la limite centrale, la loi de la variable $\frac{Z_1 + \dots + Z_n}{n}$ converge vers la loi normale centrée de covariance Γ . Or, cette loi est le projeté d'un vecteur aléatoire de \mathbb{R}^K suivant une loi normale centrée réduite sur l'hyperplan orthogonal au vecteur colonne $\left(\sqrt{p_1^{(0)}}, \dots, \sqrt{p_K^{(0)}}\right)^T$.

De plus, nous voyons que la variable T_n est la norme au carré de cette variable donc, d'après le théorème de Cochran, la loi asymptotique est une loi du χ^2 à $K - 1$ degrés de liberté.



Attention au piège

Quand nous appliquons le théorème, l'erreur la plus classique est de penser que c'est une loi du χ^2 à K degrés de liberté (au lieu de $K - 1$). Une aide mémoire possible est de se souvenir que la dernière coordonnée du vecteur est totalement définie une fois les $K - 1$ premières sont fixées (nous le voyons notamment dans la matrice Γ de la démonstration qui est la matrice identité moins la matrice de rang 1).

Corollaire 35 (Décision)

À partir de la proposition 34, nous pouvons rejeter l'hypothèse nulle si la statistique T_n dépasse le quantile d'ordre $1 - \alpha$ de la loi du χ^2 à $K - 1$ degrés de liberté.

Remarque

Il existe des versions de ce test pour le cas de données continues en découpant en classes et en calculant les probabilités d'appartenir à une classe (voir par exemple Tenenhaus (2007), chapitre 2, section 3, sous-section 5). L'inconvénient est alors que le résultat dépend du choix de la classe.

La (mauvaise) utilisation de la statistique à travers les âges

Dans leurs articles controversés, Reinhart et Rogoff (2010) étudient l'influence de la dette publique sur la croissance d'un pays. En plus d'une erreur Excel, qui a largement été documentée (voir par exemple la vidéo Youtube de Louapre ou son blog lou (2020)), les auteurs choisissent des classes qu'on peut juger arbitraires (notamment la borne critique de 90% du PIB) qui a des conséquences sur beaucoup de politiques d'austérité (dont la France).

Test de Kolmogorov-Smirnov

Le deuxième test d'adéquation est celui de Kolmogorov-Smirnov et se base cette fois sur la fonction de répartition. Le principe est d'estimer cette dernière.

Définition 69 (Fonction de répartition empirique)

Étant donné un n -échantillon \mathbf{X} de variables réelles, nous appelons **fonction de répartition empirique** la fonction définie sur \mathbb{R} par :

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \leq x\}}.$$



HISTORY

Le principe du test est de calculer le plus grand écart entre la fonction de répartition recherchée et la fonction de répartition empirique.

Proposition 36 (Test de Kolmogorov-Smirnov)

Étant donné un n -échantillon \mathbf{X} de variable réelle et une loi de référence de fonction de répartition $F^{(0)}$, alors nous définissons la statistique de test

$$D_n = \sqrt{n} \sup_{x \in \mathbb{R}} |F^{(0)}(x) - F_n(x)|$$

et, sous l'hypothèse \mathcal{H}_0 , nous avons :

$$D_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \sup_{x \in \mathbb{R}} |B_{F^{(0)}(x)}|$$

où $(B_t)_{t \in [0;1]}$ est un pont brownien ; c'est-à-dire un mouvement brownien conditionné par le fait que $B_0 = B_1 = 0$.

En particulier, dans le cas d'une fonction de répartition continue, nous avons pour toute constante c strictement positive :

$$\mathbb{P}(D_n > c) \xrightarrow[n \rightarrow +\infty]{} \alpha(c) = 2 \sum_{r=1}^{+\infty} (-1)^{r-1} e^{-2c^2 r^2}.$$

Remarque

Nous remarquons qu'il existe une formule explicite pour calculer la p -valeur même si nous ne devons nous contenter d'approximations numériques dans la pratique. En particulier, $\alpha(1,36) \approx 0,05$. Sur la figure 9.1, nous avons représenté la fonction $c \mapsto \alpha(c)$.

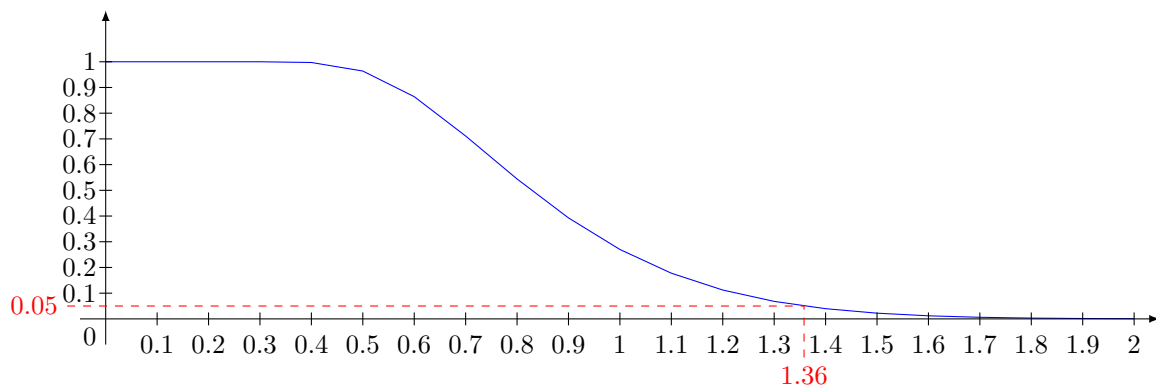


FIGURE 9.1 – Représentation de la fonction $c \mapsto \alpha(c)$ qui représente la probabilité $\mathbb{P}\left(\sup_{t \in [0;1]} |B_t| \geq c\right)$ (donc la p -valeur). La ligne hachée rouge permet de retrouver la valeur correspondante à une probabilité de 0.05.

Corollaire 37 (Décision)

À partir de la proposition 36, nous pouvons rejeter l'hypothèse nulle si la statistique D_n dépasse le quantile d'ordre $1 - \alpha$ de la loi présentée dans la proposition.

Bibliographie

- Les politiques d'austérité : à cause d'une erreur excel? <https://scienceetonnante.com/2020/04/17/austerite-excel/>, 04 2020. Accessed : 2020-11-26.
- P. Besse. Analyse en composantes principales (acp). Wikistat Toulouse. URL <http://www.math.univ-toulouse.fr/~besse/Wikistat/pdf/st-m-explo-acp>.
- G. Biau, J. Droniou, et M. Herzlich. *Mathématiques et statistique pour les sciences de la nature*. EDP Sciences, 2010.
- V. Brault. Polycopié pour l'ue statistique. *M1 Mathématiques Générales*, 2021a. URL <https://www-ljk.imag.fr/membres/Vincent.Brault/Cours/CoursStatistique.pdf>.
- V. Brault. Polycopié des cours dispensés dans le cadre de la formation stid grenoble. *Département STID de Grenoble*, 2021b. URL https://www-ljk.imag.fr/membres/Vincent.Brault/Cours/Cours_1A.pdf.
- P.-A. Cornillon, A. Guyader, F. Husson, N. Jégou, J. Josse, M. Kloareg, E. Matzner-Løber, et L. Rouvière. *Statistique avec R*. 2012.
- D. Fourdrinier. *Statistique inférentielle : cours et exercices corrigés*. Dunod, 2002.
- P. Gautret, J.-C. Lagier, P. Parola, L. Meddeb, M. Mailhe, B. Doudier, J. Courjon, V. Giordanengo, V. E. Vieira, H. T. Dupont, et al. Hydroxychloroquine and azithromycin as a treatment of covid-19 : results of an open-label non-randomized clinical trial. *International journal of antimicrobial agents*, page 105949, 2020.
- J. Goetz, S. Lapoix, et H. Poulain. *Data gueule*, 2014. URL <https://www.youtube.com/user/datagueule/>.
- T. Hamblin. Fake. *British medical journal (Clinical research ed.)*, 283(6307) :1671, 1981.
- F. Husson et J. Pagès. *Statistiques générales pour utilisateurs : 2, Exercices et corrigés*. Presses universitaires de Rennes, 2013.
- F. Husson, S. Lê, et J. Pagès. *Analyse de données avec R*. Presses universitaires de Rennes, 2016.
- F. Husson, P.-A. Cornillon, A. Guyader, N. Jégou, J. Josse, N. Klutchnikoff, E. Le Pennec, E. Matzner-Løber, L. Rouvière, et B. Thieurmél. *R pour la statistique et la science des données*. Presses universitaires de Rennes, 2018.
- J.-F. Le Gall. Intégration, probabilités et processus aléatoires. *Ecole Normale Supérieure de Paris*, 2006. URL <https://www.imo.universite-paris-saclay.fr/~jflgall/IPPA2.pdf>.
- L. M. Leemis. Relationships among common univariate distributions. *The American Statistician*, 40(2) : 143–146, 1986.
- H. Lehning. *Cryptographie et codes secrets. l'art de cacher*. 26, 2006.
- H. W. Lilliefors. On the kolmogorov-smirnov test for normality with mean and variance unknown. *Journal of the American statistical Association*, 62(318) :399–402, 1967.
- D. Louapre. Les politiques d'austérité : à cause d'une erreur excel? URL https://www.youtube.com/watch?v=yeX_Zs7zztY&ab_channel=ScienceEtonnante.
- P. L. Micheaux, R. Drouilhet, et B. Liquet. *Le logiciel R*. 2011.

- C. M. Reinhart et K. S. Rogoff. Growth in a time of debt. *American economic review*, 100(2) :573–78, 2010.
- C. Robert. *Le choix bayésien : Principes et pratique*. Springer Science & Business Media, 2006.
- J. P. Royston. Algorithm as 181 : the w test for normality. *Applied Statistics*, pages 176–180, 1982.
- G. Saporta. *Probabilités, analyse des données et statistique*. Editions technip, 2006.
- S. S. Shapiro et M. B. Wilk. An analysis of variance test for normality (complete samples). *Biometrika*, 52(3-4) :591–611, 1965.
- M. Tenenhaus. *Statistique : méthodes pour décrire, expliquer et prévoir*, volume 680. Dunod Paris, France :, 2007.
- N. Uyttendaele, L. Maugeri, et G. Grisi. La statistique expliquée à mon chat, 2016. URL https://www.youtube.com/channel/UCWty1tzwZW_ZNSp5GVGteaA/featured.
- B. Ycart. Histoires de mathématiques. <https://hist-math.u-ga.fr/>, 2017.


Index

- acceptation
 - Acceptation
 - Région, 59
- Acceptation
 - de l'hypothèse \mathcal{H}_1 , 58
- Alternatif
 - Hypothèse alternative, 58
- Aplatissement
 - Coefficient, 18
- Apparié
 - Échantillons appariés, 80
- Arcsinus, 56
- Asymétrie
 - Coefficient, 18
- Asymptotiquement normal, 37
- Asymptotiquement sans biais, 42

- Barlett
 - Test de Barlett, 90
- Barre
 - Diagramme en barres, 15, 20
- Bâton
 - Diagramme en bâtons, 15, 20
- Bernoulli
 - Densité, 27
 - Loi, 26, 27
- Bêta
 - Densité, 27
 - Loi, 27
- Biais, 41
 - Asymptotiquement sans, 42
 - Sans, 41
 - Test sans biais, 61
- Bienaymé-Tchebychev
 - Inégalité de Bienaymé-Tchebychev, 53
- Bilatère
 - Intervalle, 53
- Binomial
 - Densité, 27
 - Densité de la binomiale négative, 27
 - Loi, 26, 27
- binomiale
 - Loi binomiale négative, 26, 27
- Binomiale négative
 - Densité, 27
 - Loi, 26, 27
- Bivarié
 - Statistique descriptive bivariée, 21
- Boite
 - à moustaches, 15
- Boîte
 - à moustaches, 18
- Borne
 - de Cramer-Rao, 47
- Boxplot, 15, 18


- Camembert, 19
- Cauchy
 - Densité, 27
 - Loi, 27
- Centile, 15, 52
- Chat
 - Statistique expliquée à mon chat, 22
- χ^2
 - Test, 108
- Chocolat
 - corrélation et moustaches de chats, 22
- Circulaire
 - Diagramme circulaire, 15, 19
- Classe
 - modale, 15
- Cochran
 - Test de Cochran-Mantel-Haenszel, 100
 - Théorème, 69, 105
- Codage, 8
- Coefficient
 - d'aplatissement, 18
 - d'asymétrie, 18
 - de corrélation empirique, 92
 - de corrélation linéaire (de Pearson), 22
- Comparaison
 - Test de comparaison des égalités d'espérances de deux échantillons, 80
 - Test de comparaison des égalités d'espérances de deux échantillons appariés, 80
 - Test de comparaison des égalités d'espérances de deux échantillons gaussiens appariés, 81
 - Test de comparaison des égalités d'espérances de deux échantillons gaussiens indépendants, 85, 86
 - Test de comparaison des égalités d'espérances de deux échantillons indépendants, 84
 - Test de comparaison des égalités d'espérances de deux échantillons non gaussiens appariés, 82

- Test de comparaison des égalités d'espérances de deux échantillons non gaussiens indépendants, 85, 87
- Test de comparaison des égalités des variances de deux échantillons, 82
- Composite
 - Hypothèse composite, 58
- Confiance
 - intervalle, 52
 - Région de confiance asymptotique de $g(\theta^*)$ au niveau $1 - \alpha$, 51
 - Région de confiance de $g(\theta^*)$ au niveau $1 - \alpha$, 51
- Conservation
 - de l'hypothèse \mathcal{H}_0 , 58
- Consistance
 - Suite de tests consistante, 61
- Consistant, 31
 - Fortement, 31
- Continu
 - Variable, 12
- Convergence
 - en loi, 36
 - en probabilité, 31
 - presque sûre, 31
 - Vitesse, 36
- Corrélation
 - Chocolat, corrélation et moustaches de chats, 22
 - Coefficient de corrélation empirique, 92
 - Test de corrélation, 92
- Covariance, 22
 - Matrice de variance covariance, 72
- Cramer
 - Borne de Cramer-Rao, 47
- Cryptanalyse, 24
- Data
 - Scientist, 7
 - Scientist : serment, 7
- Décile, 15, 52
- Décision
 - Prise de décision, 8, 14
- Delta
 - Méthode, 38, 40
- Démarche
 - statistique, 26
- Densité, 33
 - Bernoulli, 27
 - Binomiale, 27
 - Binomiale négative, 27
 - Bêta, 27
 - de \mathbb{P}_θ , 33
 - de Cauchy, 27
 - de Laplace, 27
 - de Pareto, 27
 - Exponentielle, 27
 - Gamma, 27
 - gaussienne, 27
- Géométrique, 27
- Hypergéométrique, 27
- Inverse gamma, 27
- normale, 27
- Poisson, 27
- Uniforme, 27
- Descriptif
 - Statistique descriptive bivariée, 21
 - Statistique descriptive multivariée, 23
 - Statistique descriptive univariée, 14
- Diagramme
 - circulaire, 15, 19
 - de Pareto, 15, 20
 - empilé, 15, 20
 - en barres, 15, 20
 - en bâtons, 15, 20
 - en tuyaux d'orgue, 20
 - en tuyaux d'orgues, 15
- Discret
 - Variable, 12
- Dispersion
 - Résumé statistique de dispersion, 14, 15
- Distribution
 - de Lilliefors, 76
- Données
 - Explorations des données, 8, 14
 - Jeu de données, 8
 - Pré-traitement des données, 8, 14
- Écart-type, 15
- Échantillon
 - n -échantillon de loi \mathbb{P} , 26
 - apparié, 80
 - Vraisemblance du n -échantillon, 34
- Echantillon
 - statistique, 11
- Effectif
 - de la modalité a_k , 16
 - Modalité effective, 15
- Égalité
 - Test de comparaison des égalités d'espérances de deux échantillons, 80
 - Test de comparaison des égalités d'espérances de deux échantillons appariés, 80
 - Test de comparaison des égalités d'espérances de deux échantillons gaussiens appariés, 81
 - Test de comparaison des égalités d'espérances de deux échantillons gaussiens indépendants, 85, 86
 - Test de comparaison des égalités d'espérances de deux échantillons indépendants, 84
 - Test de comparaison des égalités d'espérances de deux échantillons non gaussiens appariés, 82
 - Test de comparaison des égalités d'espérances de deux échantillons non gaussiens indépendants, 85, 87
 - Test de comparaison des égalités des variances de deux échantillons, 82

- Empilé
 - Diagramme empilé, 15, 20
- Empirique
 - Coefficient de corrélation empirique, 92
 - Fonction de répartition empirique, 15, 21, 75, 102, 110
 - Moyenne, 32
 - Variance, 33
- Équité
 - Serment d'Hippocrate du Data Scientist, 7
- Espèce
 - Risque de première espèce, 59
 - Risque de seconde espèce, 59
- Espérance
 - Test d'une espérance, 66
 - Cas gaussien variance connue, 66, 67
 - Cas gaussien variance inconnue, 68, 70
 - Cas gaussien variance inconnue avec $\hat{\sigma}_n$, 71
 - Cas non gaussien variance inconnue avec S_n , 79
 - Test d'égalité des espérances
 - Cas gaussien et échantillons appariés, 81
 - Cas gaussien et échantillons indépendants, 85, 86
 - Cas gaussien, multiples échantillons appariés, 91
 - Cas non gaussien et échantillons appariés, 82, 85, 87
 - Test de comparaison des égalités d'espérances de deux échantillons, 80
 - Test de comparaison des égalités d'espérances de deux échantillons appariés, 80
 - Test de comparaison des égalités d'espérances de deux échantillons gaussiens appariés, 81
 - Test de comparaison des égalités d'espérances de deux échantillons gaussiens indépendants, 85, 86
 - Test de comparaison des égalités d'espérances de deux échantillons indépendants, 84
 - Test de comparaison des égalités d'espérances de deux échantillons non gaussiens appariés, 82
 - Test de comparaison des égalités d'espérances de deux échantillons non gaussiens indépendants, 85, 87
- Estimateur, 30
 - consistant, 31
 - de la variance, 68
 - des moments, 40
 - du maximum de vraisemblance, 34
 - UMVU, 43
 - Uniform Minimum Variance Unbiased (UMVU), Haenszel
 - 43
 - Vitesse de convergence, 36
- Étendue, 15, 17
- Excel 
 - Boxplot, 19
 - Boîte à moustaches, 19
- Explicatif
 - Variable, 13
- Expliquer
 - Variable à expliquer, 13
- Exploration
 - des données, 8, 14
- Exponentielle
 - Densité, 27
 - Loi, 27
- Fisher
 - Information, 45
- Fonction
 - de répartition empirique, 15, 21, 75, 102, 110
- Fondamental
 - Hypothèse, 26
- Forme
 - Résumé statistique de forme, 14, 15
- Fortement
 - consistant, 31
- Fractile, 15
- Fréquence
 - cumulée de la modalité a_k , 16
 - de la modalité a_k , 16
- Gamma
 - Densité, 27
 - Densité inverse gamma, 27
 - Loi, 27
 - Loi inverse gamma, 27
- Gaussien
 - Densité, 27
 - Linéarité des vecteurs gaussiens, 67
 - Loi, 27
 - Test de comparaison des égalités d'espérances de deux échantillons gaussiens appariés, 81
 - Test de comparaison des égalités d'espérances de deux échantillons gaussiens indépendants, 85, 86
 - Test de comparaison des égalités d'espérances de deux échantillons non gaussiens appariés, 82
 - Test de comparaison des égalités d'espérances de deux échantillons non gaussiens indépendants, 85, 87
- Géométrique
 - Densité, 27
 - Loi, 26, 27
- graphique, 15, 22
- Graphique, 15
 - QQ-plot, 72
- Haenszel
 - Test de Cochran-Mantel-Haenszel, 100
- Hippocrate
 - Serment d'Hippocrate du Data Scientist, 7
- Histogramme, 15
- Hoeffding
 - Inégalité de Hoeffding, 54

- Hypergéométrie
 - Densité, 27
 - Loi, 26, 27
- Hypothèse
 - alternative, 58
 - composite, 58
 - fondamentale, 26
 - nulle, 58
 - simple, 58
 - Test, 58
- Identifiabilité, 27
- Identifiable
 - Modèle, 27
- Indépendance
 - et responsabilité (Serment d'Hippocrate du Data Scientist), 7
- Indicateur
 - de liaisons, 22
- Individu
 - statistique, 10
- Inégalité
 - de Bienaymé-Tchebychev, 53
 - de Hoeffding, 54
 - de Markov, 53
- Inférence, 26
- Inférentiel
 - Statistique inférentielle, 26
- Information
 - de Fisher, 45
- Intégrité
 - et rigueur (Serment d'Hippocrate du Data Scientist), 7
- Interdécile
 - Rapport interdécile, 15
- Intérêt
 - Paramètre, 30
- Interquartile
 - Intervalle, 17
 - Longueur interquartile, 15
 - Longueur de l'intervalle, 17
- Intervalle
 - bilatère, 53
 - de confiance, 52
 - interquartile, 17
 - interquartile (longueur), 17
- Inverse gamma
 - Densité, 27
 - Loi, 27
- Jeu
 - de données, 8
- χ^2
 - Test, 108
- Kolmogorov
 - Test de Kolmogorov-Smirnov, 75, 102, 110
- Kurtosis, 15, 18
- Laplace
 - Densité, 27
 - Loi, 27
 - Pierre-Simon, 24
- Lemme
 - de Neyman-Pearson, 63
- Liaison
 - Indicateur de liaisons, 22
- Lilliefors
 - Distribution de Lilliefors, 76
 - Test de Lilliefors, 75
- Linéarité
 - des vecteurs gaussiens, 67
- Log-vraisemblance, 34
- Loi
 - binomiale, 26, 27
 - binomiale négative, 26, 27
 - bêta, 27
 - Convergence en loi, 36
 - de Bernoulli, 26, 27
 - de Cauchy, 27
 - de Laplace, 27
 - de Pareto, 27
 - de Poisson, 26, 27
 - de Student, 70
 - exponentielle, 27
 - gamma, 27
 - gaussienne, 27
 - géométrique, 26, 27
 - hypergéométrique, 26, 27
 - inverse gamma, 27
 - multinomiale de paramètres n et (p_1, \dots, p_K) , 108
 - normale, 27
 - uniforme, 26, 27
- Longueur
 - de l'intervalle interquartile, 17
 - interquartile, 15
- Mantel
 - Test de Cochran-Mantel-Haenszel, 100
- Markov
 - Inégalité de Markov, 53
- Matrice
 - de variance covariance, 72
- Maximum, 15
 - Estimateur du maximum de vraisemblance, 34
 - Méthode du maximum de vraisemblance, 34
- Médiane, 15, 52
- Méthode
 - Delta, 38, 40
 - des moments, 32
 - du maximum de vraisemblance, 34
 - du pivot, 53
- Minimum, 15
- Modalité
 - Effectif de la modalité a_k , 16
 - Effective, 15
 - Fréquence cumulée de la modalité a_k , 16

- Fréquence de la modalité a_k , 16
- Mode, 15
 - Classe modale, 15
- Modèle
 - identifiable, 27
 - non-paramétrique, 25
 - paramétrique, 25
 - paramétrique régulier, 44
 - statistique, 25
- Modélisation
 - statistique, 8, 14
- Moment
 - Estimateur, 40
 - Méthode, 32
- Moustache
 - Boîte à moustaches, 15
 - Boîte, 18
- Moyenne, 15
 - empirique, 32
- Multinomial
 - Loi multinomiale de paramètres n et (p_1, \dots, p_K) , Pearson, 108
- Bultivarié
 - Statistique descriptive multivariée, 23
- n -échantillon
 - de loi \mathbb{P} , 26
 - vraisemblance, 34
- Nemar
 - Test de Mc Nemar, 98
- Neyman
 - Lemme de Neyman-Pearson, 63
- Niveau
 - asymptotique d'un test, 59
 - d'un test, 59
 - de confiance exactement égale à $1 - \alpha$, 51
 - Région de confiance asymptotique de $g(\theta^*)$ au niveau $1 - \alpha$, 51
 - Région de confiance de $g(\theta^*)$ au niveau $1 - \alpha$, 51
- Nominal
 - Variable, 12
- Non-paramétrique, 25
 - Modèle, 25
- Normal
 - asymptotiquement, 37
 - Densité, 27
 - Loi, 27
- Normalité
 - Test de Lilliefors, 75
 - Test de normalité, 72, 76
 - Test de Shapiro-Wilk, 72
- Norme
 - infinie, 31
- Nul
 - Hypothèse nulle, 58
- Observation, 26
- Ordinal
 - Variable, 12
- Orgue
 - Diagramme en tuyaux d'orgue, 20
 - Diagramme en tuyaux d'orgues, 15
- Outsider, 19
- Paramètre
 - d'intérêt, 30
- Paramétrique, 25
 - Modèle, 25
 - Modèle non-paramétrique, 25
 - Modèle régulier, 44
 - Non-paramétrique, 25
- Pareto
 - Densité, 27
 - Diagramme de Pareto, 15, 20
 - Loi, 27
- Pearson
 - Coefficient de corrélation linéaire (de Pearson), 22
 - Lemme de Neyman-Pearson, 63
 - Test de Chi2 d'indépendance de Pearson, 101
- Pivot
 - Méthode, 53
- Plat
 - Tri à plat, 15, 17
- Poisson
 - Densité, 27
 - Loi, 26, 27
- Polygone
 - régulier, 15
- Population
 - statistique, 10
- Postion
 - Résumé statistique de position, 14, 15
- Premier
 - Risque de première espèce, 59
- Pré-traitement
 - des données, 8, 14
- Prévision, 8, 14
- Probabilité
 - Test d'une probabilité; niveau approché, 97
 - Test d'une probabilité; niveau exact, 96
 - Test d'égalité des proportions; échantillons indépendants, 99
- Probabilités, 7
- Protocole, 8
- Puissance
 - d'un test, 59
 - Test uniformément plus puissant, 61
- p -value, 61
- p -valeur, 61
- Quadratique
 - Risque, 43
- Qualitatif
 - Variable, 12
- Quantile

- d'ordre β , 52
- d'une loi binomiale, 96
- d'une loi de Fisher de paramètres $n_{\mathbf{X}} - 1$ et $n_{\mathbf{Y}} - 1$, 83, 84
- d'une loi de Fisher de paramètres $K - 1$ et $(n - 1)(K - 1)$, 91
- d'une loi de Student centrée à $n - 1$ degrés de liberté, 71, 79, 81, 82, 85-88
- d'une loi de Student à $n - 2$ degrés de libertés, 92
- d'une loi du Chi2 $\chi^2(1)$ à 1 degrés de liberté, 98
- d'une loi du Chi2 $\chi^2[(I - 1)(J - 1)]$ à $(I - 1)(J - 1)$ degrés de liberté, 101
- d'une loi normale centrée réduite, 67, 97, 99
- Quantitatif
 - Variable, 12
- Quantité
 - à estimer, 30
- Quartile, 15, 52
- Question, 8
- R** 
 - Boxplot, 19
 - Boîte à moustaches, 19
- Rao
 - Borne de Cramer-Rao, 47
- Rapport
 - de vraisemblance, 62
 - interdécile, 15
 - Statistique du rapport de vraisemblance, 62
 - Test du rapport de vraisemblance, 62
- Recueil, 8
- Région
 - d'acceptation, 59
 - de confiance asymptotique de $g(\theta^*)$ au niveau $1 - \alpha$, 51
 - de confiance de $g(\theta^*)$ au niveau $1 - \alpha$, 51
 - de rejet, 59
- Régulier
 - Modèle paramétrique régulier, 44
 - Polygone régulier, 15
- Rejet
 - de l'hypothèse \mathcal{H}_0 , 58
 - Région, 59
- Répartition
 - Fonction de répartition empirique, 15, 21, 75, 102, 110
- Respect
 - Serment d'Hippocrate du Data Scientist, 7
- Responsabilité
 - et indépendance (Serment d'Hippocrate du Data Scientist), 7
- Résumé
 - statistique de dispersion, 14, 15
 - statistique de forme, 14, 15
 - statistique de position, 14, 15
- Rigueur
 - et intégrité (Serment d'Hippocrate du Data Scientist), 7
- Risque, 43
 - de première espèce, 59
 - de seconde espèce, 59
- Sans
 - Asymptotiquement sans biais, 42
 - biais, 41
- Score, 44
- Seconde
 - Risque de seconde espèce, 59
- Serment
 - d'Hippocrate du Data Scientist, 7
- Shapiro
 - Test de Shapiro-Wilk, 72
- Simple
 - Hypothèse simple, 58
- Skewness, 15, 18
- Slutsky
 - Lemme, 38
- Smirnov
 - Test de Kolmogorov-Smirnov, 75, 102, 110
- Sous-échantillon
 - statistique, 11
- Statistique, 7
 - de test, 59
 - de vraisemblance, 62
 - descriptive, 7, 14
 - descriptive bivariée, 21
 - descriptive multivariée, 23
 - descriptive univariée, 14
 - Démarche, 26
 - Echantillon, 11
 - expliqué à mon chat, 22
 - Individu, 10
 - inférentielle, 26
 - Modèle, 25
 - Modélisation, 8, 14
 - Population, 10
 - Sous-échantillon, 11
 - Variable, 12
- Student
 - Loi, 70
- Suite
 - de tests consistante, 61
- Surapprentissage, 25
- tableau
 - Tableau, 14, 21
- Tableau, 15
- Taille
 - d'un test, 59
- Test, 58
 - d'une espérance, 66
 - Cas gaussien variance connue, 66, 67
 - Cas gaussien variance inconnue, 68, 70
 - Cas gaussien variance inconnue avec $\hat{\sigma}_n$, 71

- Cas non gaussien variance inconnue avec S_n , 79
- d'une probabilité; niveau approché, 97
- d'une probabilité; niveau exact, 96
- d'égalité des espérances
 - Cas gaussien et échantillons appariés, 81
 - Cas gaussien et échantillons indépendants, 85, 86
 - Cas gaussien, multiples échantillons appariés, 91
 - Cas non gaussien et échantillons appariés, 82
 - Cas non gaussien et échantillons indépendants, 85, 87
- d'égalité des proportions; échantillons indépendants, 99
- d'égalité des variances
 - Cas gaussien, 83
 - Cas non gaussien, 84
- de Barlett, 90
- de Chi2 d'indépendance de Pearson, 101
- de Cochran-Mantel-Haenszel, 100
- de comparaison des égalités d'espérances de deux échantillons, 80
- de comparaison des égalités d'espérances de deux échantillons appariés, 80
- de comparaison des égalités d'espérances de deux échantillons gaussiens appariés, 81
- de comparaison des égalités d'espérances de deux échantillons gaussiens indépendants, 85, 86
- de comparaison des égalités d'espérances de deux échantillons indépendants, 84
- de comparaison des égalités d'espérances de deux échantillons non gaussiens appariés, 82
- de comparaison des égalités d'espérances de deux échantillons non gaussiens indépendants, 85, 87
- de comparaison des égalités des variances de deux échantillons, 82
- de corrélation, 92
- de Kolmogorov-Smirnov, 75, 102, 110
- de l'hypothèse \mathcal{H}_0 contre l'hypothèse \mathcal{H}_1 , 58
- de Lilliefors, 75
- de Mc Nemar, 98
- de normalité, 72, 76
- de Shapiro-Wilk, 72
- de vraisemblance, 62
- de Welch, 86
- du χ^2 , 108
- Niveau, 59
- niveau asymptotique, 59
- Puissance, 59
- sans biais, 61
- Statistique, 59
- Suite de tests consistante, 61
- Taille, 59
- uniformément plus puissant, 61
- Théorème
 - de Cochran, 69, 105
- Transparence
 - Serment d'Hippocrate du Data Scientist, 7
- Tri
 - à plat, 15, 17
- Tuyau
 - Diagramme en tuyaux d'orgue, 20
 - Diagramme en tuyaux d'orgues, 15
- UMVU
 - Estimateur, 43
- Uniform
 - Minimum Variance Unbiased (UMVU), 43
- Uniforme
 - Densité, 27
 - Loi, 26, 27
- Uniformément
 - Test uniformément plus puissant, 61
- Univarié
 - Statistique descriptive univariée, 14
- Valeur
 - p-value*, 61
 - p-valeur*, 61
- Variable
 - explicative, 13
 - qualitative, 12
 - qualitative nominale, 12
 - qualitative ordinale, 12
 - quantitative, 12
 - quantitative continue, 12
 - quantitative discrète, 12
 - statistique, 12
 - à expliquer, 13
- Variance, 15
 - empirique, 33
 - Matrice de variance covariance, 72
 - Test d'égalité des variances
 - Cas gaussien, 83
 - Cas non gaussien, 84
 - Test de comparaison des égalités des variances de deux échantillons, 82
- Vecteur
 - Linéarité des vecteurs gaussiens, 67
- Vitesse
 - de convergence, 36
- Vraisemblance, 34
 - du n -échantillon \mathbf{X} , 34
 - Estimateur du maximum, 34
 - Log-vraisemblance, 34
 - Méthode du maximum, 34
 - Rapport de vraisemblance, 62
 - Statistique du rapport de vraisemblance, 62
 - Test du rapport de vraisemblance, 62
- Welch
 - Test de Welch, 86
- Wilk
 - Test de Shapiro-Wilk, 72

Accronymes

- MMV : méthode du maximum de vraisemblance, 34
 MM : méthode des moments, 32
 UMVU : Uniform Minimum Variance Unbiased, 43

Acronyme

- $UPP(\alpha)$: test uniformément plus puissant, 61

Notations

- $D\ell(\cdot)$: différentielle de l'application ℓ , 38
 $F^{-1}(\beta)$: quantile d'ordre β , 52
 F_k : fréquence cumulée de la modalité a_k , 16
 IQ : longueur de l'intervalle interquartile, 17
 $I_n(\cdot)$: information de Fisher, 45
 N_k : effectif de la modalité a_k , 16
 $R(\cdot)$: risque quadratique, 43
 $V_{\mathbf{X}}$: vraisemblance du n -échantillon \mathbf{X} , 34
 W : étendue d'une distribution, 17
 $X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X$: convergence en loi de la suite $(X_n)_{n \in \mathbb{N}^+}$ vers la variable X , 36
 $X_{(n)}$: valeur maximale du n -échantillon $\mathbf{X} = (X_1, \dots, X_n)$, 35
 Θ_0 : Sous-ensemble de Θ lié à l'hypothèse nulle, 58
 Θ_1 : Sous-ensemble de Θ lié à l'hypothèse alternative, 58
 \bar{X}_n : moyenne empirique, 32
 α^* : taille d'un test, 59
 $\underline{\alpha}$: Risque de première espèce, 59
 $\underline{\beta}$: Risque de seconde espèce, 59
 \widehat{g} : estimateur de $g(\theta^*)$, 30
 \mathcal{H}_0 : Hypothèse nulle, 58
 \mathcal{H}_1 : Hypothèse alternative, 58
 $\mathcal{M}(n; p_1, \dots, p_K)$: Loi multinomiale de paramètres n et (p_1, \dots, p_K) , 108
 $\|\cdot\|_{+\infty}$: norme infinie, 31
 $\widehat{\theta}_n$: estimateur du maximum de vraisemblance, 34
 θ^* : paramètre *inconnu* d'intérêt, 30
 \widehat{s}_n^2 : variance empirique, 33
 $b(\cdot)$: biais, 41
 f_k : fréquence de la modalité a_k , 16
 f_θ : densité de \mathbb{P}_θ , 33
 $g(\theta^*)$: quantité à estimer, 30
 $p(\cdot \cdot \cdot; \theta)$: densité de \mathbb{P}_θ , 33
 $p_\theta(\cdot)$: densité de \mathbb{P}_θ , 33
 q_β : quantile d'ordre β , 52