

# Introduction à la statistique bayésienne

Notes de cours

Cours DSXS 315

ISAE-SUPAERO 2018–2019

Xavier Gendre et Florian Simatos

21 novembre 2018

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Généralités . . . . .	5
1.2	Exemples historiques . . . . .	5
1.3	Exemple fil rouge : le modèle gaussien à variance connue . . . . .	6
<b>2</b>	<b>Motivations et justifications</b>	<b>8</b>
2.1	Théorème de De Finetti . . . . .	8
2.2	Justification axiomatique . . . . .	8
2.3	Théorie de la décision . . . . .	10
2.3.1	Critères fréquentistes . . . . .	10
2.3.2	Approche bayésienne . . . . .	14
2.3.3	Admissibilité et estimateurs bayésiens . . . . .	16
2.3.4	Minimaxité et estimateurs bayésiens . . . . .	16
2.3.5	Estimateurs ponctuels classiques et lien avec la théorie de la décision . . . . .	18
<b>3</b>	<b>Quelques notions de théorie de l'information</b>	<b>22</b>
3.1	Entropie . . . . .	22
3.2	Divergence de Kullback–Leibler et entropie croisée . . . . .	22
3.3	Score et information de Fisher . . . . .	23
<b>4</b>	<b>Choix de la distribution a priori</b>	<b>25</b>
4.1	Approximations paramétriques . . . . .	25
4.2	Maximum d'entropie . . . . .	25
4.3	Lois conjuguées et familles exponentielles . . . . .	26
4.3.1	Résultats généraux . . . . .	26
4.3.2	Modèle gaussien à variance connue . . . . .	29
4.3.3	Modèle gaussien à moyenne connue . . . . .	29
4.3.4	Modèle gaussien à moyenne et précision inconnues . . . . .	30
4.4	Lois a priori impropres . . . . .	32
4.5	Zellner . . . . .	34
4.6	Lois a priori non-informatives . . . . .	34
4.6.1	Lois invariantes et a priori de Laplace . . . . .	34
4.6.2	La loi a priori de Jeffreys . . . . .	34
<b>5</b>	<b>Comportement asymptotique des estimateurs bayésiens</b>	<b>37</b>
5.1	Régularité et différentiabilité en moyenne quadratique . . . . .	37
5.2	Comportement asymptotique de l'estimateur du maximum de vraisemblance . . . . .	37
5.2.1	Théorème et classe de Glivenko–Cantelli . . . . .	38
5.2.2	Retour sur l'estimateur du maximum de vraisemblance . . . . .	39
5.3	Comportements asymptotiques des estimateurs bayésiens . . . . .	41
5.3.1	Comportement asymptotique du MAP . . . . .	41
5.3.2	Comportement asymptotique de la densité a posteriori . . . . .	42
5.3.3	Application du Théorème de Bernstein-von Mises à la loi de Jeffreys : a priori de référence . . . . .	44
5.3.4	Moyenne a posteriori . . . . .	45
5.4	Inégalité de van Trees . . . . .	46

<b>6</b>	<b>Les tests d'hypothèses bayésiens</b>	<b>48</b>
6.1	Cadre fréquentiste . . . . .	48
6.2	Cadre bayésien . . . . .	48
6.2.1	Estimateurs bayésiens . . . . .	48
6.2.2	Le facteur de Bayes . . . . .	50
6.2.3	Hypothèses nulles simples . . . . .	50
6.2.4	Loi a priori impropres . . . . .	51
6.2.5	Régions de confiance . . . . .	52
<b>A</b>	<b>Tableau de quelques lois absolument continues</b>	<b>54</b>

*Bien que la connaissance débute avec l'expérimentation, il ne s'ensuit pas que la connaissance soit entièrement déduite de l'expérimentation.*

Kant

*On dit souvent qu'il faut expérimenter sans idée préconçue. Cela n'est pas possible ; non seulement ce serait rendre toute expérience stérile, mais on le voudrait qu'on ne le pourrait pas. Chacun porte en soi sa conception du monde dont il ne peut se défaire si aisément.*

Poincaré

# 1 Introduction

## 1.1 Généralités

On se place dans le cadre de l'estimation paramétrique classique, i.e., on a un modèle paramétrique  $\{f_\theta : \theta \in \Theta\}$  avec  $f_\theta$  une densité de probabilité (par défaut, sur  $\mathbb{R}$  et  $\Theta \subset \mathbb{R}$ ). On a un échantillon  $x = (x_1, \dots, x_n)$  supposé, sauf mention explicite du contraire, i.i.d. et tirée selon une “vraie” densité  $f_{\theta_0}$ . Le but est de faire une “bonne” estimation de  $\theta_0$  à partir de l'échantillon  $x$ .

Dans le cadre fréquentiste classique, que vous avez par exemple vu en tronc commun de première année, cette estimation prend la valeur d'une variable aléatoire  $\hat{\theta} \in \Theta$  qui vit dans le même espace que  $\theta_0$ . L'estimateur le plus connu, et sur lequel nous reviendrons régulièrement dans le cours, est l'estimateur du maximum de vraisemblance qui, comme son nom l'indique, sélectionne le paramètre  $\theta$  qui maximise la vraisemblance de l'échantillon :

$$\hat{\theta}_n^{\text{MV}} \in \arg \max_{\theta \in \Theta} \prod_{i=1}^n f_\theta(x_i).$$

On remarque ainsi en particulier que dans ce cadre fréquentiste, aucune estimation n'est proposée en l'absence d'observation. En pratique en revanche, face à un problème d'estimation et même en l'absence d'observation on est amenés à proposer un estimateur. C'est ce que l'on appelle la loi a priori en statistique bayésienne, qui sera  $\pi(\theta)$ , et qui représente le degré de confiance, l'état des connaissances que l'on a sur le paramètre à estimer. L'approche bayésienne consiste à mettre à jour cette loi a priori en fonction des observations : on obtient ainsi la loi a posteriori, notée  $\pi(\theta | x)$ , qui est la loi de  $\theta$  conditionnée par les observations. On calcule cette loi à l'aide de la formule de Bayes, que l'on note

$$\pi(\theta | x) \propto \pi(\theta)f(x | \theta)$$

où, dans le contexte bayésien, on note  $f(x | \theta) = f_\theta(x)$  pour mettre en avant que  $f_\theta$  est la loi d'échantillonnage conditionnée par le paramètre  $\theta$ , qui est effectivement rendu aléatoire via la loi a priori  $\pi$ . Une différence majeure entre l'approche fréquentiste et bayésienne est donc qu'en statistique bayésienne, l'estimation du paramètre n'est pas un nombre (ou un vecteur), mais une distribution de probabilités sur l'espace des paramètres : c'est la loi a posteriori, qui résulte de la loi a priori et des observations. Nous verrons qu'il y a en fait d'autres motivations ou justifications de cette approche, et que le choix de la loi a priori est un des enjeux majeurs de la statistique bayésienne. Pour illustrer ces concepts on commence par considérer deux exemples historiques.

## 1.2 Exemples historiques

### Exemple 1.1.

---

Le problème suivant a été posé par Bayes en 1763. Une boule de billard  $W$  roule sur une ligne de longueur un, avec une probabilité uniforme de s'arrêter n'importe où. Supposons qu'elle s'arrête en  $p$ . Une deuxième boule  $O$  roule alors  $n$  fois dans les mêmes conditions, et on note  $X$  le nombre de fois que la boule  $O$  s'arrête à la gauche

de  $W$ . Avant d'avoir lancé  $O$ , quelle estimation de  $p$  pouvez-vous faire ? Comment cette estimation est-elle modifiée après avoir lancé  $O$  et Connaisant  $X$  ?

On pourrait adopter un raisonnement fréquentiste en raisonnant conditionnellement à  $p$  et considérer  $X/n$  comme estimateur, qui converge par le théorème central limite. Ce raisonnement, justifié, ignore le fait que  $p$  elle-même est aléatoire et même suit une loi uniforme. On peut essayer d'exploiter cette information pour calculer la nouvelle loi de  $p$  mise à jour avec les observations : on verra que cette loi s'appelle la loi a posteriori. Le théorème de Bayes nous donne

$$\mathbb{P}(a < p < b \mid X = x) = \frac{\int_a^b \binom{n}{x} p^x (1-p)^{n-x} dp}{\int_0^1 \binom{n}{x} p^x (1-p)^{n-x} dp} \propto \int_a^b p^x (1-p)^{n-x} dp$$

et on reconnaît donc la densité d'une loi beta de paramètre  $x + 1$  et  $n - x + 1$ .

Fin exemple 1.1.

### Exemple 1.2.

Le problème suivant a été posé par Laplace en 1773. Une urne contient un nombre  $n$  de boules noires ou blanches. Si la première boule sortie de l'urne est blanche, quelle est la probabilité que la proportion  $p$  de boules blanches soit  $p_0$  ?

Ici la situation est un peu différente de l'exemple précédent, puisqu'on pose explicitement une question en terme de probabilité concernant la proportion de boules blanches : là on sort forcément du cadre fréquentiste habituel, et on voit qu'il faut nécessairement faire une hypothèse de base sur la loi de  $p$ . Par exemple, Laplace a supposé que  $p$  était uniformément réparti dans  $\{2/n, \dots, (n-1)/n\}$  et la règle de Bayes nous donne alors comme loi a posteriori

$$\mathbb{P}(p = p_0 \mid \text{données}) = \frac{2np_0}{n(n-1) - 2}$$

ce qui représente la loi biaisée.

Fin exemple 1.2.

## 1.3 Exemple fil rouge : le modèle gaussien à variance connue

Dans tout le cours, on considèrera un exemple dit fil rouge qui servira à illustrer les notions introduites : il s'agit du modèle d'échantillonnage gaussien à variance  $\sigma^2$  connue, avec une loi a priori elle aussi gaussienne de paramètre  $(\mu_0, \sigma_0^2)$  :

$$\pi(\theta) = \frac{1}{(2\pi\sigma_0^2)^{1/2}} e^{-(\theta - \mu_0)^2 / (2\sigma_0^2)} \quad (1.1)$$

et

$$f(x \mid \theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{k=1}^n (x_k - \theta)^2\right). \quad (1.2)$$

Lorsque l'on sera amené à manipuler des variables gaussiennes, on utilisera les deux résultats suivants.

**Lemme 1.1.** *Si  $X$  est un vecteur gaussien de moyenne  $\mu$  et de matrice de variance-covariance  $\Sigma$ , alors  $MX + A$  est aussi un vecteur gaussien de moyenne  $M\mu + A$  et de matrice de variance-covariance  $M\Sigma M^T$ .*

**Théorème 1.2.** Soit  $X = (X_1, X_2)$  un vecteur gaussien tel que  $X_1$  est absolument continu. Alors  $X_2$  conditionnellement à  $X_1$  est un vecteur gaussien de moyenne

$$\mathbb{E}(X_2 | X_1) = \mathbb{E}(X_2) + \text{Cov}(X_2, X_1)\text{Var}(X_1)^{-1}(X_1 - \mathbb{E}(X_1))$$

et de matrice de variance-covariance

$$\Sigma_{X_1} = \text{Var}(X_2) - \text{Cov}(X_2, X_1)\text{Var}(X_1)^{-1}\text{Cov}(X_1, X_2).$$

Pour  $x$  un vecteur on notera  $|x| = \sum |x_i|$  sa norme  $L_1$  et  $\|x\|$  sa norme  $L_2$ , i.e.,  $\|x\|^2 = \sum x_i^2$ . On déduit en particulier du second résultat la loi a posteriori.

**Lemme 1.3.** Dans l'exemple fil rouge, la loi a posteriori est la loi

$$\mathcal{N}\left(p\bar{x} + q\mu_0, \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)^{-1}\right) \quad \text{avec } p = 1 - q = \frac{\sigma_0^2}{\sigma^2/n + \sigma_0^2}. \quad (1.3)$$

*Démonstration.* On fait d'abord le calcul en dimension 1. On a  $\theta \sim \mathcal{N}(\mu_0, \sigma_0^2)$  et  $x | \theta \sim \mathcal{N}(\theta, \sigma^2)$  ce que l'on peut écrire sous la forme  $x = \theta + \sigma\varepsilon$  et  $\theta = \theta_0 + \sigma_0\delta$  avec  $\varepsilon, \delta \sim \mathcal{N}(0, 1)$ . Donc  $(x, \theta)$  est un vecteur gaussien et  $x$  admet une densité, si bien que l'on trouve que  $\theta | x$  suit une loi gaussienne de moyenne

$$\mu_0 + \text{Cov}(\theta, x)\text{Var}(x)^{-1}(x - \mu_0)$$

et de variance

$$\sigma_0^2 - \text{Cov}(\theta, x)^2\text{Var}(x)^{-1}.$$

On calcule  $\text{Var}(x) = \sigma^2 + \sigma_0^2$  et  $\text{Cov}(\theta, x) = \sigma_0^2$  ce qui donne le résultat pour  $n = 1$ . Pour  $n \geq 1$ , on se ramène au cas  $n = 1$  : en effet, la vraisemblance est donnée par

$$\begin{aligned} f(x | \theta) &\propto_{\theta} \exp\left(-\frac{1}{2\sigma^2} \sum_i (x_i - \theta)^2\right) \\ &\propto_{\theta} \exp\left(-\frac{1}{2\sigma^2} \sum_i (\theta^2 - 2\theta x_i)\right) \\ &\propto_{\theta} \exp\left(-\frac{n\theta^2 - 2n\theta\bar{x}}{2\sigma^2}\right) \\ &\propto_{\theta} f(\bar{x} | \theta) \end{aligned}$$

et donc  $\pi(\theta | x) \propto \pi(\theta)f(x | \theta) \propto \pi(\theta)f(\bar{x} | \theta) \propto \pi(\theta | \bar{x})$  ce qui donne le résultat puisque  $\bar{x} | \theta \sim \mathcal{N}(\theta, \sigma^2/n)$ .  $\square$

Ce résultat permet dès à présent d'illustrer plusieurs idées que l'on reverra pendant le cours :

1. la moyenne a posteriori est une combinaison linéaire de l'information a priori via  $\mu_0$  et de l'information apportée par les données via le maximum de vraisemblance  $\bar{x}$ . En outre, lorsque la taille de l'échantillon augmente, l'effet de la loi a priori s'estompe et l'estimation se base principalement sur les données ;
2. si l'on centre et que l'on met à l'échelle, la loi a posteriori suit une loi gaussienne.

## 2 Motivations et justifications

### 2.1 Théorème de De Finetti

Dans l'approche bayésienne, les observations  $x$  ne sont pas i.i.d., mais ne le sont que conditionnellement au paramètre  $\theta$  vu comme un paramètre aléatoire. En fait, le théorème de De Finetti montre que cette situation est caractéristique de suites dites échangeables.

**Définition 2.1.** Une suite  $x = (x_1, x_2, \dots)$  est dite échangeable si pour tout  $n$  et pour toute permutation  $\sigma : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ ,  $(x_1, \dots, x_n)$  et  $(x_{\sigma(1)}, \dots, x_{\sigma(n)})$  ont la même loi.

Evidemment, une suite i.i.d. est échangeable, mais la réciproque n'est pas vraie. Un contre-exemple simple est donnée par la suite  $(x_0 + x_k, k \geq 1)$  avec les  $(x_k, k \geq 0)$  i.i.d.. Néanmoins, le théorème de De Finetti montre qu'une suite échangeable n'est qu'un mélange de suites i.i.d.

**Théorème 2.2** (De Finetti). *Une suite de variables aléatoires  $(x_1, x_2, \dots)$  de densité  $f$  est échangeable si et seulement si il existe une vraisemblance  $f(x | \theta)$  et une mesure de probabilités  $P$  sur  $\Theta$  telles que pour tout  $n$ ,*

$$f(x_1, \dots, x_n) = \int f(x_1 | \theta) \cdots f(x_n | \theta) P(d\theta).$$

L'exemple suivant illustre les relations entre le théorème de De Finetti et l'approche bayésienne.

**Exemple 2.3.** \_\_\_\_\_

Historiquement les moteurs de recherche ont utilisé des modèles de "sacs de mots" pour modéliser ces documents, dans lesquels l'ordre des mots ne compte pas. Si l'on voit un mot et qu'il est français, alors on s'attend à ce que le reste du document soit en français. Cela montre que les mots ne sont pas i.i.d., mais qu'ils le sont possiblement conditionnellement à un paramètre, par exemple ici la langue.

\_\_\_\_\_ **Fin exemple 2.3.**

### 2.2 Justification axiomatique

**Principe de conditionnement.** *Si deux expériences  $E_1$  et  $E_2$  sur le paramètre  $\theta$  sont possibles et si on choisit une de ces expériences au hasard avec probabilité  $p$ , alors l'inférence sur  $\theta$  ne doit dépendre que de l'expérience choisie.*

**Exemple 2.4.** \_\_\_\_\_

Dans un laboratoire de recherche, une quantité physique  $\theta$  doit être mesurée par un appareil efficace, mais très souvent utilisé, qui donne une mesure  $x_1 \sim N(\theta, 1)$ , avec une probabilité  $p = 0,5$ , ou grâce à un autre appareil, moins précis mais plus disponible, qui donne  $x_2 \sim N(\theta, 10)$ . L'appareil a été choisi au hasard selon la disponibilité de l'appareil le plus précis. L'inférence sur  $\theta$  ne devrait donc pas dépendre du fait que le second appareil aurait pu être choisi. Cependant, un intervalle de confiance au seuil 5% prenant en compte cette sélection, soit donc moyennant entre toutes les expériences possibles, est de demi-longueur 5,19, tandis que l'intervalle associé à  $E_1$  est de demi-longueur 0,62.

En effet, si  $a$  est la demi-longueur, on a en moyennisant sur les expériences

$$\mathbb{P}(\theta - a \leq x \leq \theta + a) = \frac{1}{2} \left[ \mathbb{P}\left(|N| \leq \frac{a}{\sigma_1}\right) + \mathbb{P}\left(|N| \leq \frac{a}{\sigma_2}\right) \right] = 1 - \alpha$$

alors qu'en ne considérant que  $E_1$ ,

$$\mathbb{P}(\theta - a \leq x \leq \theta + a) = \mathbb{P}\left(|N| \leq \frac{a}{\sigma_1}\right) = 1 - \alpha$$

---

**Fin exemple 2.4.**

**Principe de vraisemblance.** *L'information apportée par une observation de  $x$  sur  $\theta$  est entièrement contenu dans la fonction de vraisemblance  $f(\theta | x)$ .*

**Exemple 2.5.**

Soit  $\theta$  le biais d'une pièce de monnaie. On a un échantillon consistant de 9 fois face et 3 fois pile, obtenu par l'une des deux expériences suivantes :

1. on lance la pièce 12 fois et on regarde le nombre de pile ;
2. on la lance la pièce un nombre de fois suffisant pour voir 3 fois pile.

Dans les deux cas, on peut vérifier que la vraisemblance  $f(x | \theta) \propto \theta^9(1 - \theta)^3$ , et donc le principe de vraisemblance suggère que l'inférence devrait être la même. Néanmoins, si on teste l'hypothèse  $H_0 : \theta = 1/2$  contre  $H_1 : \theta > 1/2$  avec un risque de première espèce fixé, alors la réponse dépend de l'expérience considérée.

---

**Fin exemple 2.5.**

**Définition 2.3.**  $T$  est une statistique exhaustive si, alors que  $x \sim f(x | \theta)$ , la loi de  $x$  conditionné à  $T(x)$  ne dépend pas de  $\theta$ .

En d'autres termes, une statistique est exhaustive si elle contient toute l'information apportée sur  $x$  par  $\theta$ . L'importance théorique des statistiques exhaustives découle en grande partie du théorème de Rao–Blackwell.

**Théorème 2.4** (Théorème de Rao–Blackwell). *Si  $T$  est une statistique exhaustive, alors pour tout estimateur  $\hat{\theta}$  l'erreur quadratique moyenne est améliorée en considérant comme estimateur l'espérance de  $\hat{\theta}$  conditionnellement à  $T$ , i.e.,*

$$\mathbb{E}_\theta \left[ \left( \theta - \mathbb{E}_\theta(\hat{\theta} | T) \right)^2 \right] \leq \mathbb{E}_\theta \left[ \left( \theta - \hat{\theta} \right)^2 \right].$$

*Démonstration.* En considérant  $\hat{\theta} - \theta$  on peut considérer sans perte de généralité que  $\theta = 0$ . On a

$$\mathbb{E}_\theta \left[ \hat{\theta} \mathbb{E}_\theta(\hat{\theta} | T) \right] = \mathbb{E}_\theta \left[ \mathbb{E}_\theta(\hat{\theta} | T)^2 \right]$$

et donc

$$0 \leq \mathbb{E}_\theta \left[ \left( \hat{\theta} - \mathbb{E}_\theta(\hat{\theta} | T) \right)^2 \right] = \mathbb{E}_\theta \left( \hat{\theta}^2 \right) - \mathbb{E}_\theta \left[ \mathbb{E}_\theta(\hat{\theta} | T)^2 \right]$$

ce qui prouve le résultat. □

**Principe d'exhaustivité.** *Deux observations  $x$  et  $y$  donnant la même valeur d'une statistique exhaustive  $T$  doivent conduire à la même inférence sur  $\theta$ .*

**Théorème 2.5** (Birnbaum). *Le principe de vraisemblance est équivalent à la conjonction des principes d'exhaustivité et de conditionnement.*

L'approche bayésienne satisfait ces trois principes, contrairement à l'approche fréquentiste, comme illustré sur les exemples ci-dessus. En effet, dans le cadre bayésien toute inférence est faite à partir de la loi a posteriori, et satisfait donc le principe de vraisemblance.

### 2.3 Théorie de la décision

La dernière justification que nous présenterons provient de la théorie de la décision. En pratique, les observations sont là pour aider un décideur à prendre une décision, vendre un stock, etc. Une décision est donc représentée par une fonction  $\delta : \mathcal{O} \rightarrow D$  avec  $D$  l'ensemble des actions possibles, et l'on appellera une telle fonction une règle de décision (la décision étant plus l'action). Par simplicité, on se restreint dans ce cours au cas où la décision à prendre est une estimation, i.e.,  $D = \mathcal{O}$ , et on notera donc  $\mathcal{D} = \{\delta : \mathcal{O} \rightarrow \Theta \text{ mesurable}\}$  l'ensemble des règles de décision, i.e., estimateurs, possibles.

**Exemple 2.6.** \_\_\_\_\_

#### TODO

Pour fixer les idées : exemple de décision qui n'est pas une estimation.

\_\_\_\_\_ **Fin exemple 2.6.**

La théorie de la décision repose sur deux éléments :

- une **fonction de coût** ou **fonction de perte**  $L : \Theta \times D \rightarrow \mathbb{R}$  qui à une valeur de paramètre  $\theta \in \Theta$  et une décision  $d \in D$  associe un coût  $L(\theta, d) \in \mathbb{R}$  ;
- une **règle de décision**  $\delta : \mathcal{O} \rightarrow D$  qui à des observations  $x \in \mathcal{O}$  associe une décision/observation  $\delta(x) \in D = \Theta$ .

Les fonctions de coût les plus classiques correspondent à des distances  $L_p$ , par exemple la *perte quadratique*  $L(\theta, d) = (\theta - d)^2$  ou encore la *perte absolue*  $L(\theta, d) = |\theta - d|$ . Dans le cas où  $\Theta$  est fini ou dénombrable, on peut aussi considérer la fonction de coût binaire  $L(\theta, d) = \mathbf{1}(\theta \neq d)$ , i.e.,  $L(\theta, d) = 0$  si  $\theta = d$  et 1 sinon.

#### 2.3.1 Critères fréquentistes

Dans le cadre de la théorie de la décision, la fonction de coût  $L$  est supposée fixée une fois pour toute, et on s'intéresse à ce qui se passe lorsque le paramètre  $\theta \in \Theta$  varie ainsi que la règle de décision  $\delta$ . Ainsi, on considère le **risque** ou **risque fréquentiste**  $R_f(\theta, \delta)$  associé au paramètre  $\theta$  et à la règle de décision  $\delta$  défini par le coût moyen lorsque le paramètre vaut  $\theta$  :

$$R_f(\theta, \delta) = \mathbb{E}_\theta [L(\theta, \delta(\cdot))] = \int L(\theta, \delta(x))f(x | \theta)dx.$$

**Exemple fil rouge 2.7.** \_\_\_\_\_

Pour l'exemple fil rouge, on prendra une fonction de coût quadratique :  $L(\theta, \delta) = (\theta - \delta)^2$  et une fonction linéaire en les observations comme règle de décision/estimateur,

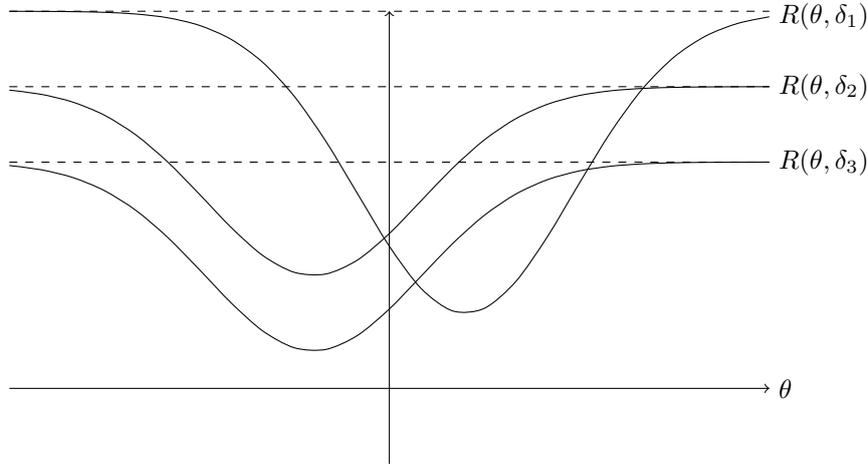


FIGURE 1 – Un exemple de trois risques fréquentistes :  $\delta_3$  est uniformément meilleur que  $\delta_2$ , on dit qu'elle lui est préférable et donc  $\delta_2$  est inadmissible. Par contre, il n'est a priori pas évident de comparer  $\delta_1$  à  $\delta_2$  et  $\delta_3$ .

i.e.,  $\delta = a^T x$  pour un certain vecteur  $a$ . Si  $a = (1/n, \dots, 1/n)$  on a alors  $\delta = \bar{x}$ , la moyenne empirique, et par ailleurs, l'estimateur est non biaisé si et seulement si  $A := a^T \mathbf{1}_n = \sum a_i = 1$ . On a alors

$$\begin{aligned}
 R_f(\theta, \delta) &= \mathbb{E}_\theta \left[ \left( \theta - a^T X \right)^2 \right] \\
 &= \mathbb{E}_\theta \left[ \left( \theta - \theta a^T \mathbf{1}_n - (a^T X - \theta a^T \mathbf{1}_n) \right)^2 \right] \\
 &= (1 - A)^2 \theta^2 + \text{Var}(a^T X) \\
 &= (1 - A)^2 \theta^2 + \|a\|^2 \sigma^2.
 \end{aligned}$$

Si  $\delta$  est non biaisé, on a donc un risque indépendant de  $\theta$  et qui augmente avec la variance  $\sigma^2$  du modèle d'échantillonnage. Pour  $a = n^{-1} \mathbf{1}_n$  on a un risque qui diminue en  $1/n$ , ce qui correspond au risque minimal dans cette classe d'estimateurs.

Fin exemple fil rouge 2.7.

On étudie maintenant la performance fréquentiste d'une règle de décision via deux critères : l'admissibilité et la minimaxité.

**Définition 2.6.** On dit qu'une règle de décision  $\delta_1$  est **préférable** à une autre règle de décision  $\delta_2$  si les deux conditions suivantes sont satisfaites :

1. le risque associé à  $\delta_1$  n'est jamais pire que celui associé à  $\delta_2$ , i.e.,  $R_f(\theta, \delta_1) \leq R_f(\theta, \delta_2)$  pour tout  $\theta \in \Theta$  ;
2. le risque associé à  $\delta_1$  est strictement meilleur que celui associé à  $\delta_2$  pour au moins une valeur de paramètre  $\theta$ , i.e., il existe  $\theta_0$  tel que  $R_f(\theta_0, \delta_1) < R_f(\theta_0, \delta_2)$

On dit qu'une règle de décision est **admissible** si aucune autre règle ne lui est préférable.

**Remarque 2.1.**

Pour montrer qu'une règle  $\delta$  est admissible, il suffit donc de montrer que si  $R_f(\theta, \delta') \leq R_f(\theta, \delta)$  pour tout  $\theta$ , alors  $\delta = \delta'$ .

**Fin remarque 2.1.**

Par exemple, sur la Figure 1  $\delta_3$  n'est pas admissible puisque  $\delta_3$  lui est préférable. Le premier exemple suivant montre que l'admissibilité n'induit pas un ordre total. Le deuxième exemple montre que la notion d'inadmissibilité n'est pas très forte, notamment du fait que l'on ne peut pas toujours comparer deux règles. Néanmoins, si une règle n'est pas admissible alors il existe une règle dont le risque est moindre quelle que soit la vraie valeur du paramètre  $\theta$  et on a donc tout intérêt à prendre cette autre règle.

**Exemple fil rouge 2.8.**

Si on reprend le risque intégré  $R_f(\theta, \delta) = (1-A)^2\theta^2 + \|a\|^2\sigma^2$  obtenu dans l'exemple 2.7, on voit qu'en fonction des valeurs de  $A$  et  $\|a\|$  les règles de décision peuvent être ou non comparables. Si on se restreint aux règles de décision avec  $A = 1$ , i.e., les règles sans biais, alors on peut comparer deux règles  $a_1$  et  $a_2$  et  $a_1$  sera préférable à  $a_2$  si  $\|a_1\| \leq \|a_2\|$ .

**Fin exemple fil rouge 2.8.**

**Exemple 2.9.**

Soit  $\hat{\theta} = \theta_0$  pour un certain  $\theta_0 \in \Theta$  fixé, et on considère le coût quadratique  $L(\theta, d) = (\theta - d)^2$ . Soit  $T$  un autre estimateur avec  $R_f(\theta, T) \leq R_f(\theta, \hat{\theta})$  pour tout  $\theta$ . En particulier, pour  $\theta = \theta_0$  on a  $R_f(\theta_0, T) \leq R_f(\theta_0, \hat{\theta}) = 0$  et donc  $\mathbb{E}_{\theta_0}((\theta_0 - T)^2) = 0$ , i.e.,  $T = \theta_0$   $\mathbb{P}_{\theta_0}$ -presque sûrement. Ainsi,  $\hat{\theta}$  est admissible.

**Fin exemple 2.9.**

En général, si on restreint la classe d'estimateurs considérée, on peut établir des résultats d'admissibilité, par exemple dans la classe des estimateurs non biaisés comme illustré dans l'exemple 2.8. Néanmoins, ça n'est pas toujours le cas comme le montre l'exemple suivant, proposé par Stein en 1955.

**Exemple 2.10.**

On cherche à estimer un paramètre  $\theta \in \mathbb{R}^p$  à l'aide d'une observation  $X \sim \mathcal{N}(\theta, 1)$ . L'estimateur "ordinaire" est simplement donné par  $\hat{\theta} = X$  et présente un risque  $R(\theta, \hat{\theta}) = p$ . En revanche, l'estimateur

$$\hat{\theta}' = \left(1 - \frac{p-2}{\|X\|^2}\right) X$$

présente lui un risque

$$R(\hat{\theta}', \theta) = p - (p-2)\mathbb{E}\left(\frac{1}{\|X\|}\right)$$

qui est strictement meilleur dès lors que  $p \geq 3$ . Ce résultat est particulièrement surprenant, puisqu'il dit que pour estimer les moyennes des  $X_i$  qui sont pourtant **indépendants**, il est préférable d'utiliser un estimateur qui les combine! Cf. par exemple [5] pour plus de détails.

**Fin exemple 2.10.**

Un moyen un peu "brutal" de comparer le risque de deux règles de décision est de simplement comparer leur risque maximal : cela conduit à la notion de minimaxité. Par exemple, sur la Figure 1, on a  $\sup_{\theta \in \Theta} R_f(\theta, \delta_1) \geq \sup_{\theta \in \Theta} R_f(\theta, \delta_2)$  ce qui pourra amener à considérer  $\delta_2$  plutôt que  $\delta_1$ . Dans ce cas,  $\delta_2$  est appelé estimateur minimax.

**Définition 2.7.** On appelle **risque minimax** associé à la fonction de coût  $L$  la valeur

$$\bar{R} = \inf_{\delta \in \mathcal{D}} \sup_{\theta \in \Theta} R_f(\theta, \delta)$$

et **estimateur minimax** tout estimateur  $\delta_0$  avec

$$\sup_{\theta \in \Theta} R_f(\theta, \delta_0) = \bar{R}.$$

L'interprétation de  $\bar{R}$  est la suivante : je choisis une règle de décision, et la Nature choisit le pire paramètre. Je cherche donc à choisir le meilleur paramètre avec ces règles.

**Exemple fil rouge 2.11.**

Pour le risque  $R_f(\theta, \delta) = (1 - A)^2 \theta^2 + \|a\|^2 \sigma^2$ , le risque minimax est infini si l'estimateur est biaisé et que  $\Theta$  n'est pas borné. Si par exemple  $\Theta = [0, \theta_0]$ , alors on a  $\sup_{\theta} R_f(\theta, \delta) = R_f(\theta_0, \delta)$  et tout  $a$  qui minimise  $(1 - A)^2 \theta_0^2 + \|a\|^2 \sigma^2$  est minimax.

**Fin exemple fil rouge 2.11.**

Sous certaines hypothèses générales, l'existence d'un estimateur minimax est assurée par le théorème suivant.

**Théorème 2.8** ([4, Théorème 2.20]). *Si  $\mathcal{D} \subset \mathbb{R}^k$  est convexe et compact et si  $L(\theta, d)$  est continue et convexe en tant que fonction de  $d$  pour chaque valeur de  $\theta$ , alors il existe un estimateur minimax.*

Le critère minimax est extrêmement conservatif puisqu'il ne considère que les valeurs extrêmes du risque. Le risque intégré introduit ci-dessous pallie ce problème en prenant en compte tous les risques possibles, mais en pondérant les valeurs du paramètre. La terminologie en lien avec Bayes (risque de Bayes, règle de décision bayésienne, estimateur bayésien) sera justifiée dans la prochaine section.

**Définition 2.9.** Pour toute règle de décision  $\delta \in \mathcal{D}$ , le **risque de Bayes** ou **risque intégré** de  $\delta$  par rapport à une loi a priori  $\pi$  est défini par

$$R_B(\delta) = \int R_f(\theta, \delta) \pi(\theta) d\theta = \int L(\theta, \delta(x)) f(x | \theta) \pi(\theta) d\theta dx = \mathbb{E}[L(\theta, \delta(\cdot))]$$

Le **risque de Bayes selon  $\pi$**  est la plus petite valeur possible du risque de Bayes :

$$\underline{R} = \inf_{\delta \in \mathcal{D}} R_B(\delta).$$

Une règle de décision  $\delta_\pi$  est une **règle de décision bayésienne** pour l'a priori  $\pi$  si

$$R_B(\delta_\pi) = \inf_{\delta \in \mathcal{D}} R_B(\delta).$$

Si  $\delta_\pi$  est une règle bayésienne, l'estimateur  $\delta_\pi(x)$  associé est appelé **estimateur bayésien**.

**Exemple fil rouge 2.12.**

Le risque de Bayes est donné par

$$R_B(\delta) = \mathbb{E}((1 - A)^2 \theta^2 + \|a\|^2 \sigma^2)$$

où  $\theta \sim \mathcal{N}(\mu_0, \sigma_0^2)$ , ce qui donne

$$R_B(\delta) = (1 - A)^2 \mathbb{E}(\theta^2) + \|a\|^2 \sigma^2 = (1 - A)^2 (\sigma_0^2 + \mu_0^2) + \|a\|^2 \sigma^2.$$

Si l'on se restreint à la classe des estimateurs non biaisés, les estimateurs bayésiens et minimax coïncident donc puisqu'ils correspondent tous les deux à minimiser  $\|a\|^2$  sous la contrainte  $a^T \mathbf{1}_n = 1$ . Dans le cas général, le risque minimax est infini alors que le risque de Bayes reste fini.

Fin exemple fil rouge 2.12.

### Exemple 2.13.

On considère  $\pi(\theta) \sim \mathcal{E}(1)$ ,  $f(x | \theta) \sim \mathcal{P}(\theta)$  et les estimateurs linéaires  $\delta_c(x) = cx$ . On considère d'abord le coût quadratique  $L(\theta, d) = (\theta - d)^2$  : on a alors

$$R_f(\delta_c, \theta) = \mathbb{E}[(cx - \theta)^2 | \theta] = c^2 \text{Var}(x | \theta) + (c - 1)^2 \theta^2 = c^2 \theta + (c - 1)^2 \theta^2.$$

Alors  $\delta_\gamma$  est préférable à  $\delta_c$  si et seulement si

$$\begin{aligned} \forall \theta > 0, \gamma^2 \theta + (\gamma - 1)^2 \theta^2 &< c^2 \theta + (c - 1)^2 \theta^2 \\ \iff \forall \theta > 0, \gamma^2 + (\gamma - 1)^2 \theta &< c^2 + (c - 1)^2 \theta \\ \iff \gamma < c \text{ et } (\gamma - 1)^2 &< (c - 1)^2 \end{aligned}$$

Ainsi,  $\delta_c$  pour  $c > 1$  n'est pas admissible, puisque tout  $\delta_\gamma$  avec  $\gamma \in (c_*, c)$  avec  $c_* < c$  et  $(1 - c_*)^2 = (1 - c)^2$  lui est préférable (et on vérifie directement que  $\gamma = 1$  lui est préférable, puisque pour  $c > 1$  on a  $R_f(\theta, \delta_c) \geq \theta = R_f(\theta, \delta_1)$ ). Par contre, les arguments ci-dessus montrent aussi que  $\delta_c$  avec  $c \leq 1$  est admissible :  $c^2$  représente la pente à l'origine et  $(c - 1)^2$  le taux d'accroissement pour  $\theta$  grand, et donc pour  $c \geq c' \leq 1$  on a  $c^2 \leq (c')^2$  mais  $(1 - c)^2 \geq (1 - c')^2$  et donc les courbes se croisent nécessairement. Le meilleur estimateur, qui minimise le risque intégré

$$R_B(\delta_c) = c^2 + 2(1 - c)^2,$$

est obtenu pour  $c$  solution de  $2c + 4(c - 1) = 0$ , i.e.,  $c = 2/3$ . Si on prend la fonction de coût  $L'(\theta, d) = (1 - d/\theta)^2 = \theta^{-2} L(\theta, d)$  on obtient

$$R'_f(\delta_c, \theta) = c^2 \theta^{-1} + (c - 1)^2.$$

Le raisonnement pour l'admissibilité reste le même, par contre le risque intégré vaut maintenant  $+\infty$ .

Fin exemple 2.13.

### 2.3.2 Approche bayésienne

*Un fréquentiste est quelqu'un content de regarder d'autres données qu'il aurait pu avoir mais qu'il n'a pas eu.*

Michael Jordan

La volonté de minimiser le risque fréquentiste est un peu paradoxal, puisqu'elle repose sur un comportement moyen des données, et non pas les données à disposition. A l'inverse, l'approche bayésienne cherche à minimiser le risque a posteriori.

**Définition 2.10.** Le risque a posteriori  $\varrho(\delta, \pi | x)$  est le risque intégré par rapport à la loi a posteriori :

$$\varrho(\delta, \pi | x) = \mathbb{E}[L(\theta, \delta(x)) | x] = \int L(\theta, \delta(x)) \pi(\theta | x) d\theta.$$

Dans l'approche bayésienne, on va donc chercher à minimiser le risque a posteriori.

**Exemple fil rouge 2.14.**

Dans l'exemple fil rouge, on a alors

$$\varrho(\delta, \pi | x) = \mathbb{E} \left[ \left( \theta - a^T x \right)^2 | x \right] = \left( \mathbb{E}(\theta | x) - a^T x \right)^2 + \text{Var}(\theta | x)$$

et puisque

$$\theta \sim \mathcal{N} \left( p\bar{x} + q\mu_0, \left( \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1} \right)$$

par (1.3), on obtient

$$\varrho(\delta, \pi | x) = \left( p\bar{x} + q\mu_0 - a^T x \right)^2 + \left( \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1}.$$

En particulier, on voit que si l'on veut que le risque a posteriori disparaisse lorsque  $n \rightarrow \infty$ , il faut choisir une suite d'estimateurs  $a_n$  telle que  $a_n^T x \rightarrow \mu_0$ , par exemple  $a = n^{-1} \mathbf{1}_n$  qui mène au risque a posteriori

$$\varrho(\delta, \pi | x) = \left( \frac{\sigma^2}{\sigma^2 + n\sigma_0^2} \right)^2 (\mu_0 - \bar{x})^2 + \left( \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1}. \quad (2.1)$$

**Fin exemple fil rouge 2.14.**

De manière assez élégante, cette approche rejoint l'approche fréquentiste via le risque de Bayes ou risque intégré, qui consiste à intégrer le risque fréquentiste par rapport à la loi a priori. Le résultat suivant justifie cette terminologie, et montre notamment qu'un estimateur bayésien correspond à l'intégration de la fonction de coût par rapport à la loi a posteriori.

**Théorème 2.11.**  $\delta_\pi$  est une règle de décision bayésienne pour l'a priori  $\pi$  si et seulement si  $\delta_\pi(x)$  minimise le risque a posteriori pour toute observation  $x \in \mathcal{O}$  on a

$$\delta_\pi(x) \in \arg \min_{d \in \Theta} \mathbb{E} [L(\theta, d) | x].$$

**Remarque 2.2.**

En fait, on a rigoureusement  $\mathbb{P}(\delta_\pi(x) \in \arg \min \dots) = 1$ , mais ici comme dans le reste du cours on passera sous silence ces considérations issues de la théorie de la mesure.

**Fin remarque 2.2.**

*Démonstration du Théorème 2.11.* Par le théorème de l'espérance totale, on a

$$R_B(\delta) = \mathbb{E} [L(\theta, \delta(\cdot))] = \mathbb{E} [\varrho(\delta, \pi | x)] = \mathbb{E} [\mathbb{E} [L(\theta, \delta(x)) | x]].$$

Si  $\delta_\pi(x)$  minimise  $\mathbb{E} [L(\theta, d) | x]$  en  $d$  pour chaque  $x$ , alors  $R_B(\delta_\pi) \leq R_B(\delta)$  pour toute règle de décision  $\delta$  au vu de l'égalité précédente et  $\delta_\pi$  est donc une règle de décision bayésienne. Réciproquement, soit  $\delta_\pi$  une règle de décision bayésienne et  $\delta \in \mathcal{D}$  une autre règle de décision avec  $\delta(x) \in \arg \min_{d \in \Theta} \mathbb{E} [L(\theta, d) | x]$ . Alors

$$0 \leq R_B(\delta) - R_B(\delta_\pi) = \mathbb{E} [\mathbb{E} [L(\theta, \delta(x)) | x] - \mathbb{E} [L(\theta, \delta_\pi(x)) | x]].$$

Par définition de  $\delta_\pi$  on a  $\mathbb{E} [L(\theta, \delta(x)) | x] - \mathbb{E} [L(\theta, \delta_\pi(x)) | x] \leq 0$  ce qui implique que  $\mathbb{E} [L(\theta, \delta(x)) | x] = \mathbb{E} [L(\theta, \delta_\pi(x)) | x]$  pour presque tout  $x$  et donc  $\delta_\pi(x) \in \arg \min_d \mathbb{E} [L(\theta, d) | x]$ .  $\square$

### Exemple fil rouge 2.15.

Pour un risque quadratique on a par définition

$$L(\theta, d | x) = \int (\theta - d)^2 f(\theta | x) d\theta = \text{Var}(\theta | x) + [\mathbb{E}(\theta | x) - d]^2$$

et donc

$$\delta_\pi(x) = \mathbb{E}(\theta | x) = p\bar{x} + q\mu_0$$

en utilisant (1.3) pour la dernière égalité.

Fin exemple fil rouge 2.15.

### 2.3.3 Admissibilité et estimateurs bayésiens

**Proposition 2.12.** Une règle de décision bayésienne  $\delta_\pi$  est admissible si au moins une des conditions suivantes est satisfaite :

- $\delta_\pi$  est l'unique estimateur bayésien ;
- $\pi$  admet une densité strictement positive sur  $\Theta$  et  $\theta \mapsto R_f(\theta, \delta)$  est continue pour tout  $\delta \in \mathcal{D}$  ;
- $\Theta$  est fini et  $\pi(\{\theta\}) > 0$  pour tout  $\theta \in \Theta$ .

*Démonstration de a).* Soit  $\delta$  une règle de décision telle que  $R_f(\theta, \delta) \leq R_f(\theta, \delta_\pi)$  pour tout  $\theta \in \Theta$ . Puisque  $R_B(\delta) = \int R_f(\theta, \delta) \pi(\theta) d\theta$ , en intégrant sous  $\pi$  on obtient  $R_B(\delta) \leq R_B(\delta_\pi) = \inf_{\delta'} R(\Pi, \delta')$ , donc  $\delta'$  est aussi une règle de Bayes pour  $\pi$ . Donc  $\delta' = \delta$  par hypothèse, et donc  $R_f(\theta, \delta) = R_f(\theta, \delta')$  ce qui montre que  $\delta$  est admissible.  $\square$

*Démonstration de b).* Soit  $\delta$  une règle de décision telle que  $R_f(\theta, \delta) \leq R_f(\theta, \delta_\pi)$  pour tout  $\theta \in \Theta$ . Alors  $\mathbb{P}(R_f(\theta, \delta) = R_f(\theta, \delta_\pi)) = 1$  puisque

$$0 \leq \int_{\theta} (R_f(\theta, \delta_\pi) - R_f(\theta, \delta)) \pi(\theta) d\theta = R_B(\delta_\pi) - R_B(\delta) \leq 0$$

où la dernière inégalité provient de la définition de  $\delta_\pi$ . On a donc que  $R_f(\theta, \delta_\pi) \pi(\theta) = R_f(\theta, \delta) \pi(\theta)$  presque partout, et les hypothèses sur  $\pi$  et  $R_f$  impliquent  $R_f(\theta, \delta_\pi) = R_f(\theta, \delta)$  pour tout  $\theta$ .  $\square$

*Démonstration de c).* Le raisonnement est le même que précédemment.  $\square$

### 2.3.4 Minimaxité et estimateurs bayésiens

**Proposition 2.13.** Le risque de Bayes est toujours plus petit que le risque minimax, i.e.,  $\underline{R} \leq \bar{R}$ .

*Démonstration.* Pour toute règle de décision  $\delta \in \mathcal{D}$  on a

$$\underline{R} \leq R_B(\delta) = \int R_f(\theta, \delta) \pi(\theta) d\theta \leq \sup_{\theta} R_f(\theta, \delta)$$

si bien que  $\underline{R} \leq \sup_{\theta} R_f(\theta, \delta)$  pour toute règle de décision  $\delta$ , et donc  $\underline{R} \leq \inf_{\delta} \sup_{\theta} R_f(\theta, \delta) = \bar{R}$ .  $\square$

**Définition 2.14.** On dit que le problème d'estimation admet une valeur si  $\bar{R} = \underline{R}$ .

**Théorème 2.15.** Soit  $\delta_\pi$  une règle de décision bayésienne pour une distribution a priori  $\pi$ . Si

$$R_B(\delta_\pi) = \sup_{\theta \in \Theta} R_f(\theta, \delta_\pi) \quad (2.2)$$

alors :

1.  $\delta_\pi$  est une règle minimax et le problème d'estimation admet une valeur ;
2.  $\underline{R}(\pi) = \sup_{\pi'} \underline{R}(\pi')$ , i.e., on dit que  $\pi$  est une loi la moins favorable ;
3. toute règle minimax est une règle bayésienne pour  $\pi$ . En particulier, si  $\delta_\pi$  est l'unique règle de Bayes pour  $\pi$ , alors  $\delta_\pi$  est l'unique règle minimax.

*Démonstration.* Pour toute règle de décision  $\delta$  on a

$$\sup_{\theta} R_f(\theta, \delta_\pi) = R_B(\delta_\pi) \leq R_B(\delta) \leq \sup_{\theta} R_f(\theta, \delta)$$

ce qui montre bien que  $\delta_\pi$  est minimax. De plus, si  $\delta$  est une autre règle minimax, alors on a égalité dans les inégalités ci-dessous et donc on a en particulier  $R_B(\delta) = R_B(\delta_\pi)$  et donc  $\delta$  est une règle bayésienne.

Pour démontrer le deuxième point, on note que pour toute loi a priori  $\pi'$

$$\underline{R}(\pi') = \inf_{\delta} R_B(\delta, \pi') \leq \inf_{\delta} \sup_{\theta} R_f(\theta, \delta) = \underline{R}(\pi)$$

où la dernière égalité vient du premier point. □

La condition 2.2 peut paraître assez forte, et effectivement elle l'est. Elle est satisfaite en particulier si  $R_f(\theta, \delta_\pi)$  ne dépend pas de  $\theta$  ce qui est par exemple le cas dans l'exemple suivant.

**Exemple fil rouge 2.16.**

Comme on l'a déjà vu dans l'exemple 2.7, dans le cas non biaisé le risque fréquentiste ne dépend pas de  $\theta$  et la règle de décision bayésienne est donc minimax, ce que l'on avait déjà vu dans l'exemple 2.12.

**Fin exemple fil rouge 2.16.**

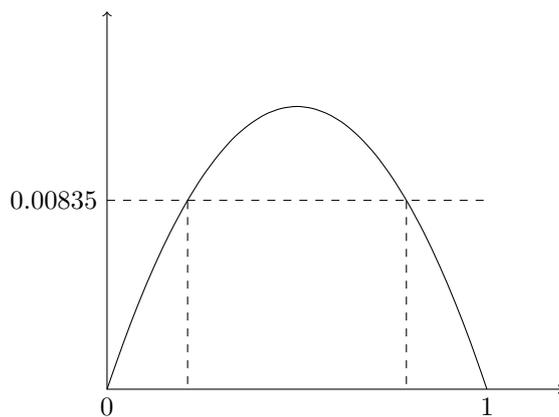


FIGURE 2 – Comparaison du risque minimax et du risque bayésien dans l'exemple 2.17.

---

**Exemple 2.17.**

On reprend le problème d'estimation du paramètre d'une loi de Bernoulli, i.e., le modèle paramétrique est  $f(x | \theta) = \theta^{|x|}(1 - \theta)^{1-|x|}$  et on considère un a priori qui suit une loi  $\mathcal{B}e(\alpha, \beta)$ . La loi a posteriori est donc une loi  $\mathcal{B}e(\alpha + |x|, \beta + n - |x|)$ , de moyenne

$$\frac{\alpha + |x|}{\alpha + \beta + n}.$$

Considérant une fonction coût quadratique, l'estimateur bayésien  $\hat{\theta}$  est donc donné par la moyenne a posteriori, et donc sous  $\mathbb{P}_\theta$ , on a

$$\hat{\theta} = \frac{\alpha + |x|}{\alpha + \beta + n}.$$

et en utilisant la décomposition biais-variance, le risque fréquentiste est donné par

$$\begin{aligned} R_f(\theta, \hat{\theta}) &= \left(\theta - \mathbb{E}_\theta(\hat{\theta})\right)^2 + \text{Var}_\theta(\hat{\theta}) \\ &= \left(\theta - \frac{\alpha + n\theta}{\alpha + \beta + n}\right)^2 + \frac{n\theta(1-\theta)}{(\alpha + \beta + n)^2} \\ &= \frac{\{(\alpha + \beta)^2 - n\}\theta^2 + \{n - 2\alpha(\alpha + \beta)\}\theta + \alpha^2}{(\alpha + \beta + n)^2} \end{aligned}$$

Cela ne dépend pas de  $\theta$  si  $(\alpha + \beta)^2 = n$  et  $2\alpha(\alpha + \beta) = n$ , ce qui se résout par  $\alpha = \beta = \sqrt{n}/2$ , et ce qui donne alors

$$\hat{\theta} = \frac{\sqrt{n}/2 + |x|}{\sqrt{n} + n} = \frac{1}{2 + 2\sqrt{n}} + \frac{n}{\sqrt{n} + n} \bar{x}_n$$

et le risque correspondant vaut

$$R_f(\theta, \hat{\theta}) = \frac{n}{4(n + \sqrt{n})^2}.$$

Comme ce risque ne dépend pas de  $\theta$ , le Théorème 2.15 implique que  $\hat{\theta}$  est minimax. La figure 2 compare le coût minimax  $n/(4(n + \sqrt{n})^2)$  au risque obtenu pour l'estimateur du maximum de vraisemblance obtenu pour la moyenne empirique  $\hat{\eta} = n^{-1}\text{Binom}(n, \theta)$  qui vaut

$$\mathbb{E}_\theta((\theta - \hat{\eta})) = \frac{\theta(1-\theta)}{n}.$$

La théorie nous garantit  $\hat{\theta}$  présente un risque optimal, mais qu'il est très conservatif puisque pour une large plage de valeurs de  $\theta$ , la moyenne empirique offre un risque plus faible.

---

**Fin exemple 2.17.**

### 2.3.5 Estimateurs ponctuels classiques et lien avec la théorie de la décision

Etant donné la distribution a posteriori  $\pi(\theta | x)$ , de nombreux choix sont possibles si l'on souhaite faire une estimation ponctuelle : le mode, la moyenne, la médiane, etc. En fait, chaque estimateur provient d'une fonction de coût différente.

**Proposition 2.16.** *L'estimateur de Bayes est donné par :*

- le mode  $\arg \max_{\theta \in \Theta} \pi(\theta | x)$  dans le cas d'une fonction de coût binaire :  $L(\theta, d) = \mathbb{1}(d \neq \theta)$  ;

- l'espérance conditionnelle  $\mathbb{E}(\theta \mid x)$  dans le cas d'une fonction de coût quadratique :  $L(\theta, d) = (\theta - d)^2$  ;
- une médiane de  $\pi(\theta \mid x)$  dans le cas d'une fonction de coût  $L_1$  :  $L(\theta, d) = |\theta - d|$ .

Une fonction de coût binaire n'est bien définie que si  $\Theta$  est fini, sinon dans le cas continu on aurait toujours un risque fréquentiste qui vaut 1. Par contre, le mode est toujours bien défini.

**Définition 2.17.** On dit que  $\delta$  est un estimateur du maximum a posteriori (MAP) si  $\delta(x) \in \arg \max_{\theta \in \Theta} \pi(\theta \mid x)$  pour tout  $x \in \mathcal{O}$ .

**Lemme 2.18.** On a  $\arg \min_{m \in \mathbb{R}} \mathbb{E}((X - m)^2) = \mathbb{E}(X)$ .

*Démonstration.* On a  $\mathbb{E}((X - m)^2) = \mathbb{E}(X^2) - 2m\mathbb{E}(X) + m^2$ . □

**Lemme 2.19.** Soit  $X$  à densité : alors

$$\arg \min_{m \in \mathbb{R}} \mathbb{E}(|X - m|) = \left\{ x : \mathbb{P}(X < x) = \frac{1}{2} \right\}.$$

*Démonstration.* On a

$$\begin{aligned} \mathbb{E}(|X - m|) &= \int_{-\infty}^m (m - x)f(x)dx + \int_m^{\infty} (x - m)f(x)dx \\ &= m\mathbb{P}(X < m) - \int_{-\infty}^m xf(x)dx + \int_m^{\infty} xf(x)dx - m\mathbb{P}(X > m). \end{aligned}$$

Si on dérive par rapport à  $m$ , on obtient donc

$$\frac{d}{dm} \mathbb{E}(|X - m|) = \mathbb{P}(X < m) + mf(m) - mf(m) - mf(m) + mf(m) - \mathbb{P}(X > m)$$

ce qui donne le résultat. □

*Démonstration pour  $L(\theta, d) = \mathbf{1}(\theta \neq d)$ .* Dans ce cas, le risque fréquentiste est

$$R_f(\theta, \delta) = \mathbb{E}_\theta(L(\theta, \delta)) = \mathbb{P}_\theta(\theta \neq \delta);$$

le risque intégré est

$$R_B(\delta) = \mathbb{E}(L(\theta, \delta)) = \mathbb{P}(\theta \neq \delta);$$

et donc les estimateurs bayésiens maximisent  $d \mapsto \mathbb{P}(\theta = d \mid x) = \pi(d \mid x)$ . □

*Démonstration pour  $L(\theta, d) = (\theta - d)^2$ .* Dans ce cas, le risque fréquentiste est

$$R_f(\theta, \delta) = \mathbb{E}_\theta(L(\theta, \delta)) = \mathbb{E}_\theta \left[ (\theta - \delta)^2 \right];$$

le risque intégré est

$$R_B(\theta, \delta) = \mathbb{E}(L(\theta, \delta)) = \mathbb{E} \left[ (\theta - \delta)^2 \right];$$

et donc les estimateurs bayésiens maximisent  $d \mapsto \mathbb{E}[(\theta - d)^2 \mid x]$  qui est l'espérance conditionnelle par le Lemme 2.18 (ou par définition de l'espérance conditionnelle comme projection  $L_2$ ). □

*Démonstration pour  $L(\theta, d) = |\theta - d|$ .* Dans ce cas, le risque fréquentiste est

$$R_f(\theta, \delta) = \mathbb{E}_\theta(L(\theta, \delta)) = \mathbb{E}_\theta |\theta - \delta|;$$

le risque intégré est

$$R_B(\theta, \delta) = \mathbb{E}(L(\theta, \delta)) = \mathbb{E} |\theta - \delta|;$$

et donc les estimateurs bayésiens maximisent  $d \mapsto \mathbb{E}[|\theta - d| \mid x]$  qui est une médiane par le Lemme 2.19.  $\square$

---

**Exemple fil rouge 2.18.**

Dans le cas gaussien, tous ces estimateurs coïncident.

**Fin exemple fil rouge 2.18.**

---

**Exemple 2.19.**

On reprend l'exemple original de Bayes (exemple 1.1) qui correspond à l'estimation du paramètre d'une loi de Bernoulli dans un modèle d'échantillonnage binomial, i.e.,  $f(x \mid \theta) \propto \theta^{|x|}(1 - \theta)^{n - |x|}$ . Lorsque la loi a priori est  $\mathcal{B}e(\alpha, \beta)$ , on a vu dans l'exemple 2.17 que la loi a posteriori était la loi  $\mathcal{B}e(\alpha + |x|, n + \beta - |x|)$ . On a déjà calculé la moyenne a posteriori, qui vaut

$$\mathbb{E}(\theta \mid x) = \frac{\alpha + |x|}{n + \alpha + \beta} = \frac{\alpha}{n + \alpha + \beta} + \frac{n}{n + \alpha + \beta} \bar{x}_n \quad (2.3)$$

qui est donc obtenue comme dans le cas gaussien comme une pondération du maximum de vraisemblance et la loi a priori. En considérant le logarithme de la densité a posteriori, on voit que le MAP est obtenu en résolvant

$$\frac{\alpha + |x| - 1}{\theta} - \frac{n - |x| + \beta - 1}{1 - \theta} = 0,$$

i.e.,

$$\hat{\theta}^{\text{MAP}} = \frac{\alpha - 1}{n + \alpha + \beta - 2} + \frac{n}{n + \alpha + \beta - 2} \bar{x}_n$$

qui coïncide avec la moyenne a posteriori mais pour une loi a priori  $\mathcal{B}e(\alpha - 1, \beta - 1)$ . La médiane a posteriori n'a pas de forme explicite, elle fait intervenir la fonction beta incomplète régularisée.

**Fin exemple 2.19.**

---

**Exemple 2.20.**

On considère un modèle gaussien à variance  $\sigma^2$  connue :

$$f(x \mid \theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\|x - \theta\|^2 / (2\sigma^2)}.$$

Pour un a priori gaussien standard  $\pi(\theta) \propto e^{-(\theta - \mu_0)^2 / (2\sigma_0^2)}$ , la loi a posteriori est alors donnée par

$$\pi(\theta \mid x) \propto e^{-\|x - \theta\|^2 / 2 - (\theta - \mu_0)^2 / (2\sigma_0^2)} \propto e^{-(n+1)\theta^2 / 2 + \theta n \bar{x}} \propto e^{-(n+1)\theta^2 / 2 + \theta n \bar{x}}$$

La loi a posteriori est donc encore une loi gaussienne, de moyenne  $n\bar{x} / (n + 1)$  et de variance  $1 / (n + 1)$ . Pour une loi a priori exponentielle  $e^{-x} \mathbb{1}(x \geq 0)$ , on a une loi a posteriori donnée par

$$\pi(\theta \mid x) \propto e^{-\|x - \theta\|^2 / 2 - \theta} \propto e^{-n\theta^2 / 2 + (n\bar{x} - 1)\theta} \propto e^{-n(\theta - (\bar{x}_n - 1/n))^2 / 2}$$

et on reconnaît donc une loi gaussienne de moyenne  $\bar{x}_n - 1/n$  et de variance  $1/n$ . Dans tous ces cas le MAP, la moyenne a posteriori et la médiane sont donc égales. Si maintenant on prend comme a priori la loi de Cauchy  $\pi(\theta) \propto 1/(\theta^2 + 1)$ , on obtient

$$\pi(\theta | x) \propto e^{-\|x-\theta\|^2/2}/(\theta^2 + 1)$$

et là par contre, la loi a posteriori n'est plus gaussienne. Le MAP est alors donné par la résolution de l'équation

$$-\frac{n}{2\sigma^2}(\theta - \bar{x}) - \frac{2\theta}{\theta^2 + 1} = 0$$

et la moyenne a posteriori par

$$\frac{\int \frac{\theta}{\theta^2+1} e^{-\|x-\theta\|^2/2} d\theta}{\int \frac{1}{(\theta^2+1)} e^{-\|x-\theta\|^2/2} d\theta}$$

---

**Fin exemple 2.20.**

### 3 Quelques notions de théorie de l'information

#### 3.1 Entropie

**Définition 3.1.** Etant donné une distribution  $f$ , on définit son entropie différentielle

$$H(f) = - \int f(x) \log f(x) dx \in [-\infty, +\infty].$$

**Remarque 3.3.**

On remarque que l'entropie différentielle peut être négative, par exemple en considérant une loi uniforme, et peut aussi valoir  $\pm\infty$ . Pour une loi discrète, l'entropie est mieux définie car

$$h(p) = - \sum_x p(x) \log p(x) \in [0, \infty).$$

L'entropie est une mesure de la quantité d'information contenue dans une loi de probabilité : une entropie (discrète) faible correspond à une loi de probabilité très resserrée et donc avec beaucoup d'information. En effet, on a par exemple  $h(p) = 0$  si et seulement si  $p$  est une mesure de Dirac. L'entropie différentielle s'approche par l'entropie discrète via une constante additive. En effet, si  $f$  est une densité, on peut considérer la loi discrétisée

$$p(x) = \frac{1}{2\varepsilon} \int_{x-\varepsilon}^{x+\varepsilon} f$$

et alors

**TODO**

Compléter la remarque

**Fin remarque 3.3.**

#### 3.2 Divergence de Kullback–Leibler et entropie croisée

**Définition 3.2.** La divergence de Kullback–Leibler entre deux densités de probabilité  $f$  et  $g$  est donnée par

$$D(f, g) = \int f(x) \log \left( \frac{f(x)}{g(x)} \right) dx.$$

Par convention, on a  $0/0 = 0$  et  $D(f, g) = +\infty$  si l'ensemble  $\{x : f(x) \neq 0, g(x) = 0\}$  n'est pas de mesure nulle.

**Proposition 3.3.**  $D(f, g) \geq 0$  avec égalité si et seulement si  $f = g$  ( $p.p.$ ).

*Démonstration.* C'est l'inégalité de Jensen : pour toute variable aléatoire  $X$  et fonction  $\varphi$  convexe, on a sous réserve d'intégrabilité  $\mathbb{E}(\varphi(X)) \geq \varphi(\mathbb{E}(X))$  avec égalité si et seulement si  $\varphi$  est linéaire ou  $X$  est une constante. Donc si  $X \sim f$  on a  $D(f, g) = \mathbb{E}(-\log(g(X)/f(X)))$  et puisque  $-\log$  est convexe, l'inégalité de Jensen donne  $D(f, g) \leq -\log \mathbb{E}(g(X)/f(X)) = 0$  avec égalité, puisque  $\log$  n'est pas linéaire, si et seulement si  $f(X)/g(X)$  est une constante, i.e., si  $f = g$  ( $f$ - $p.p.$ ).  $\square$

**Définition 3.4.** L'entropie croisée entre  $f$  et  $g$  est définie par

$$CE(f, g) = - \int f(x) \log g(x) dx.$$

Ainsi,  $\text{CE}(f, g) = D(f, g) + H(f)$  et donc minimiser l'entropie croisée revient à minimiser la divergence de Kullback–Leibler, ce qui par le résultat précédent est obtenu pour  $g = f$ , i.e.,

$$f = \arg \min_g \text{CE}(f, g) = \arg \min_g D(f, g). \quad (3.1)$$

### 3.3 Score et information de Fisher

**Définition 3.5.** Soit  $\{f_\theta : \theta \in \Theta\}$  un modèle paramétrique avec  $\theta \mapsto f_\theta(x)$  dérivable pour tout  $x$  et dans  $L_2(f_\theta)$ . On définit le score de ce modèle paramétrique par

$$\ell(x | \theta) = \partial_\theta \log f(x | \theta) = \frac{\partial_\theta f(x | \theta)}{f(x | \theta)}$$

et l'information de Fisher par

$$I(\theta) = \mathbb{E}_\theta(\partial_\theta \log f_\theta(x)^2) = \int \partial_\theta \log f_\theta(x)^2 f_\theta(x) dx.$$

**Lemme 3.6.** Lorsque le modèle est régulier, on a  $I(\theta) = -\mathbb{E}_\theta \partial_\theta^2 \log f_\theta(x_1) = \mathbb{E}(\partial_\theta \log f_\theta(x)^2)$ .

*Démonstration.* En effet,

$$\begin{aligned} -\mathbb{E}_\theta \partial_\theta^2 \log f_\theta(x) &= -\int \partial_\theta^2 \log f_\theta(x) f_\theta(x) dx \\ &= -\int \left( \frac{\partial_\theta^2 f_\theta(x)}{f_\theta(x)} - \left( \frac{\partial_\theta f_\theta(x)}{f_\theta(x)} \right)^2 \right) f_\theta(x) dx \\ &= -\int \partial_\theta^2 f_\theta(x) dx + I(\theta) \end{aligned}$$

et  $\int \partial_\theta^2 f_\theta(x) dx = \partial_\theta^2 \int f_\theta(x) dx = 0$ . □

Le score peut être interprété comme une mesure de la vitesse à laquelle la densité change lorsque l'on change légèrement le paramètre  $\theta$ . L'information de Fisher, qui n'est autre que la moyenne du score au carré, permet donc d'obtenir une version moyennisée de cette mesure. Donc si l'information de Fisher est grande, cela veut dire que la distribution va changer rapidement lorsque l'on bouge le paramètre, et que donc la distribution avec le paramètre  $\theta_0$  va être "bien différente" et donc peut être "bien distinguée" des distributions avec des paramètres relativement éloignés de  $\theta_0$ . Cela veut dire que l'on devrait être capable d'estimer  $\theta_0$  correctement à partir des données. D'un autre côté, si l'information de Fisher est petite, cela veut dire que la distribution est très similaire à d'autres distributions qui ont des paramètres distincts et donc qu'il sera plus difficile de discriminer entre les deux, et donc l'estimation sera moins bonne.

Dans la suite on aura besoin du résultat technique suivant, qui explicite l'information de Fisher lorsque l'on change de variable.

**Proposition 3.7.**  $I(\theta) = I(h(\theta))(h'(\theta))^2$ .

*Démonstration.* Par définition,  $I(h(\theta))$  est l'information de Fisher associé au modèle d'échantillonnage  $g(x | \eta) = f(x | h^{-1}(\eta))$  avec  $\eta \in h(\Theta)$ . Le vecteur du score associé est

$$\tilde{\ell}(x | \eta) = \ell(x | h^{-1}(\eta))(h^{-1})'(\eta) = \frac{\ell(x | h^{-1}(\eta))}{h' \circ h^{-1}(\eta)}$$

et donc pour  $\eta = h(\theta)$ , on a

$$\begin{aligned} I(h(\theta)) &= \int [\tilde{\ell}(x | \eta)]^2 g(x | \eta) dx \\ &= \int \left( \frac{\ell(x | h^{-1}(\eta))}{h' \circ h^{-1}(\eta)} \right)^2 f(x | h^{-1}(\eta)) dx \\ &= \frac{1}{h'(\theta)^2} I(\theta) \end{aligned}$$

ce qui prouve le résultat. □

## 4 Choix de la distribution a priori

### 4.1 Approximations paramétriques

Consiste à restreindre  $\pi$  à une certaine famille paramétrique et à calibrer les paramètres en fonction des informations disponibles type moyenne/médiane. Une autre manière de “calibrer” est de rajouter une couche d’aléa et supposer que les hyperparamètres sont eux-mêmes tirés selon une distribution de probabilité : on parle alors de modèles hiérarchiques, qui ne seront pas abordés dans ce cours, cf. par exemple [4, Chapitre 10].

### 4.2 Maximum d’entropie

Il arrive que l’on ait de l’information sur la loi a priori sous la forme de moments  $\int g(\theta)\pi(\theta)d\theta$ . Dans ce cas, on peut utiliser la méthode du maximum d’entropie, qui repose sur le principe d’incertitude de Laplace, qui postule d’utiliser la loi uniforme lorsqu’aucune information n’est disponible. Puisqu’une loi avec une entropie contient peu d’information, le résultat suivant fournit une justification de ce principe.

**Lemme 4.1.** *Soit  $A$  un ensemble de mesure de Lebesgue  $m(A)$  finie. Alors la mesure uniforme sur  $A$  est l’unique solution du problème d’optimisation*

$$\begin{aligned} & \arg \max H(f) \\ & \text{t.q. } f \in \mathcal{P} \end{aligned}$$

*Démonstration.* On a

$$H(f) = \text{CE}(f, u) - D(f, u) \leq \text{CE}(f, u) = - \int_A f(x) \log 1/\lambda(A) dx = \log \lambda(A)$$

et puisque  $\log \lambda(A) = H(u)$  on obtient bien le résultat.  $\square$

Ainsi, lorsque de l’information sur la loi a priori est connue sous la forme de moments, il est naturel de considérer le problème suivant.

$$\begin{aligned} & \arg \max H(f) \\ & \text{t.q. } f \in \mathcal{P} \\ & \text{et } \int g_k(x)f(x)dx = 0, k = 1, \dots, m \end{aligned}$$

Si on considérait ce problème sur un espace d’états fini, on pourrait utiliser la méthode des multiplicateurs de Lagrange et l’on obtiendrait que les solutions sont données par (cf. Remarque 4.4 ci-dessous)

$$\pi^*(x) \propto \exp \left( \sum_{k=1}^m \lambda_k g_k(x) \right).$$

On admettra que ce résultat reste vrai en dimension infinie. Néanmoins, un problème majeur de cette approche dans le continu est qu’elle n’est pas stable par reparamétrisation : si  $\theta$  suit la loi uniforme sur  $[0, 1]$ , alors  $\theta^2$  est de densité  $1/(2\sqrt{t})$  et ne suit en particulier pas une loi uniforme. Ce problème n’est pas présent en discret, mais d’une certaine manière parce que l’espace d’état change.

On peut aussi prendre une mesure de référence  $f_0$  et considérer plutôt le problème

$$\begin{aligned} & \arg \min D(f||f_0) \\ & \text{t.q. } f \in \mathcal{P} \\ & \text{et } \int g_k(x)f(x)dx = 0, k = 1, \dots, m \end{aligned}$$

dont la solution est donnée par

$$\pi^*(x) \propto \exp\left(\sum_{k=1}^m \lambda_k g_k(x)\right) f_0(x).$$

Dans ce cas,  $\pi_0$  est la solution d'entropie minimale et on peut donc prendre pour celle-ci une distribution non-informative, cf. Section 4.6.

---

**Remarque 4.4.**

Dans le cas discret, on peut prouver ce résultat (qui généralise le résultat précédent en prenant  $f_0$  la loi uniforme) à l'aide des multiplicateurs de Lagrange. En effet, on cherche  $p = (p_1, \dots, p_n)$  et

$$L(p, \lambda) = \sum_{i=1}^n p_i \log(p_i/p_0(i)) + \sum_{k=1}^m \sum_{i=1}^n \lambda_k g_k(i) p_i$$

et donc

$$\partial_{p_i} L = \log(p_i/p_0(i)) + 1 + \sum_{k=1}^m \lambda_k g_k(i)$$

et donc

$$\partial_{p_i} L = 0 \iff p_i = p_0(i) \exp\left(-1 - \sum_{k=1}^m \lambda_k g_k(i)\right)$$

---

**Fin remarque 4.4.**

Un dernier problème est dû au fait que l'on ne peut pas toujours normaliser la distribution obtenue, par exemple si on met seulement un moment d'ordre un on a  $\pi^*(x) \propto e^{\lambda x}$  qui n'est pas intégrable sur  $\mathbb{R}$ . Par contre, si on impose deux moments alors  $\pi^*(x) \propto e^{\lambda_1 x + \lambda_2 x^2}$  qui est bien intégrable pour  $\lambda_2 < 0$ .

## 4.3 Lois conjuguées et familles exponentielles

### 4.3.1 Résultats généraux

**Définition 4.2.** La famille  $\{f(\cdot | \theta)\}$  est dite exponentielle de dimension  $k$  si elle peut s'écrire

$$f(x | \theta) = C(\theta)h(x) \exp(R(\theta) \cdot T(x))$$

avec  $R$  et  $T$  à valeurs dans  $\mathbb{R}^k$ , et  $\cdot$  le produit scalaire. Si  $R = T = \text{Id}$  on parle de famille naturelle.

A noter que  $C(\theta)$  peut être vue comme une constante de normalisation. Par ailleurs, un changement de variable de  $x$  en  $z = T(x)$  et une reparamétrisation  $\eta = R(\theta)$  amène à la forme naturelle.

---

**Exemple 4.21.**

La loi exponentielle  $f(x | \theta) = \theta e^{-\theta x} \mathbf{1}(x > 0)$  constitue une famille exponentielle avec  $C(\theta) = \theta$ ,  $h(x) = \mathbf{1}(x > 0)$ ,  $R(\theta) = -\theta$  et  $T(x) = x$ . Ainsi, elle est presque naturelle.

---

**Fin exemple 4.21.**

---

**Exemple 4.22.**

Le modèle gaussien à variance connue  $f(x | \theta) = (2\pi\sigma^2)^{-1/2} e^{-(x-\mu)^2/(2\sigma^2)}$  constitue une famille exponentielle avec  $C(\theta) = (2\pi\sigma^2)^{-1/2} e^{-\theta^2/(2\sigma^2)}$ ,  $h(x) = e^{-x^2/(2\sigma^2)}$ ,  $T(x) = x$  et  $R(\theta) = -\theta^2/(2\sigma^2)$ . Ainsi, on obtient la forme naturelle en paramétrant par  $\eta = -\theta^2/(2\sigma^2)$ .

---

**Fin exemple 4.22.**

---

**Exemple 4.23.**

Le modèle gaussien à moyenne connue  $f(x | \theta) = (2\pi\theta)^{-1/2} e^{-(x-\mu)^2/(2\theta)}$  constitue une famille exponentielle avec  $C(\theta) = (2\pi\theta)^{-1/2}$ ,  $h(x) = 1$ ,  $T(x) = -(x - \mu)^2/2$  et  $R(\theta) = 1/\theta$ . Ainsi, on obtient la forme naturelle en paramétrant par  $\eta = 1/\theta$  et  $z = (x - \mu)^2/2$ .

---

**Fin exemple 4.23.**

---

**Exemple 4.24.**

Le modèle gaussien  $f(x | \mu, \tau) = \tau(2\pi)^{-1/2} e^{-(x-\mu)^2/(2\theta)}$  constitue une famille exponentielle avec  $C(\theta) = (2\pi\theta)^{-1/2}$ ,  $h(x) = 1$ ,  $T(x) = -(x - \mu)^2/2$  et  $R(\theta) = 1/\theta$ . Ainsi, on obtient la forme naturelle en paramétrant par  $\eta = 1/\theta$  et  $z = (x - \mu)^2/2$ .

---

**Fin exemple 4.24.**

---

**Exemple 4.25.**

Plusieurs lois sont exponentielles :

- La loi exponentielle :  $h(x) = \mathbb{1}(x > 0)$ ,  $C(\theta) = \theta$ ,  $R(\theta) = -\theta$  et  $T(x) = x$ ;
- la loi normale :  $h(x) = 1$ ,  $C(\theta) = e^{\mu^2/(2\sigma^2)}/\sqrt{2\pi\sigma^2}$ ,  $T(x) = (x, x^2)$  et  $R(\theta) = (\mu/\sigma^2, -1/(2\sigma^2))$ ;
- la loi de Poisson :  $h(x) = 1/x!$ ,  $C(\theta) = e^{-\theta}$ ,  $T(x) = x$  et  $R(\theta) = \log \theta$ ;
- la loi de Dirichlet, et donc la loi beta et donc la loi uniforme : la densité est  $\frac{\prod_{i=1}^n \Gamma(\theta_i)}{\Gamma(|\theta|)} \prod_{i=1}^n x_i^{\theta_i-1}$ , ce qui correspond à  $T(x) = (\log x_i)$ ,  $R(\theta) = \theta - \mathbf{1}_n$

---

**Fin exemple 4.25.**

---

**Exemple 4.26.**

Plusieurs lois sont exponentielles :

- La loi exponentielle :  $h(x) = \mathbb{1}(x > 0)$ ,  $C(\theta) = \theta$ ,  $R(\theta) = -\theta$  et  $T(x) = x$ ;
- la loi normale :  $h(x) = 1$ ,  $C(\theta) = e^{\mu^2/(2\sigma^2)}/\sqrt{2\pi\sigma^2}$ ,  $T(x) = (x, x^2)$  et  $R(\theta) = (\mu/\sigma^2, -1/(2\sigma^2))$ ;
- la loi de Poisson :  $h(x) = 1/x!$ ,  $C(\theta) = e^{-\theta}$ ,  $T(x) = x$  et  $R(\theta) = \log \theta$ ;
- la loi de Dirichlet, et donc la loi beta et donc la loi uniforme : la densité est  $\frac{\prod_{i=1}^n \Gamma(\theta_i)}{\Gamma(|\theta|)} \prod_{i=1}^n x_i^{\theta_i-1}$ , ce qui correspond à  $T(x) = (\log x_i)$ ,  $R(\theta) = \theta - \mathbf{1}_n$

---

**Fin exemple 4.26.**

---

**Remarque 4.5.**

$T$  est une statistique exhaustive si, alors que  $x \sim f(x | \theta)$ , la loi de  $x$  conditionné à  $T(x)$  ne dépend pas de  $\theta$ . Le théorème de factorisation garantit que si  $T$  est une statistique exhaustive, alors  $f(x | \theta)$  peut s'écrire  $f(x | \theta) = g(T(x) | \theta)h(x | T(x))$  avec  $g$  la densité de  $T$ . Le Lemme de Pitman-Koopman dit par ailleurs que si une famille de lois  $f(\cdot | \theta)$  à support constant est telle que, à partir d'une taille d'échantillon

suffisamment grande, il existe une statistique exhaustive de taille fixe, alors la famille est exponentielle.

**Fin remarque 4.5.**

**Proposition 4.3.** Soit  $f(x | \theta) = h(x)e^{\theta \cdot x - \psi(\theta)}$  la loi générique d'une famille exponentielle. Une famille conjuguée pour  $f(x | \theta)$  est donnée par

$$\pi(\theta | \mu, \lambda) \propto e^{\theta \cdot \mu - \lambda \psi(\theta)}.$$

La loi a posteriori correspondante est  $\pi(\theta | \mu + x, \lambda + 1)$ .

*Démonstration.* On a

$$\pi(\theta | x) \propto e^{\theta \cdot x - \psi(\theta) + \theta \cdot \mu - \lambda \psi(\theta)} \propto \pi(\theta | \mu + x, \lambda + 1).$$

□

Il existe des familles conjuguées en dehors des familles exponentielles comme le montre le résultat suivant.

**Lemme 4.4.** La loi de Pareto  $\mathcal{Pa}(\alpha, \theta)$  est la loi sur  $\mathbb{R}$  de densité donnée par  $\alpha \theta^\alpha x^{-\alpha-1} \mathbf{1}(x \geq \theta)$ . La famille de lois a priori  $\pi(\theta | a, b) \propto \theta^{a+1} \mathbf{1}(0 \leq \theta \leq b)$  est conjuguée par les lois de Pareto.

*Démonstration.* Si  $\pi = \pi(\cdot | a, b)$  et  $f(x | \theta) \sim \mathcal{Pa}(\alpha, \theta)$ , alors

$$\pi(\theta | x) \propto \theta^{a+1} \mathbf{1}(\theta \leq b) \times \theta^\alpha \mathbf{1}(\theta \leq x) \propto \pi(\theta | a + \alpha, \min(b, x)).$$

□

**Lemme 4.5.** Soit  $\mathcal{F}$  la famille conjuguée naturelle d'une famille exponentielle. Alors l'ensemble des mélanges de  $N$  lois conjuguées

$$\hat{\mathcal{F}}_n = \left\{ \sum_{i=1}^N \omega_i \pi(\theta | \mu_i, \lambda_i) : \omega \geq 0, |\omega| = 1 \right\}$$

est aussi une famille conjuguée. De plus, si

$$\pi(\theta) = \sum_{i=1}^N \omega_i \pi(\theta | \mu_i, \lambda_i)$$

alors la loi a posteriori est un mélange

$$\pi(\theta | x) = \sum_{i=1}^N \omega'_i(x) \pi(\theta | \mu_i + x, \lambda_i + 1)$$

avec

$$\omega'_i(x) = \frac{\omega_i K(\mu_i, \lambda_i) / K(\mu_i + x, \lambda_i + 1)}{\sum_{j=1}^N \omega_j K(\mu_j, \lambda_j) / K(\mu_j + x, \lambda_j + 1)}$$

avec  $K$  la constante de normalisation de  $\pi(\theta | \lambda, \mu)$ , i.e.,  $\pi(\theta | \lambda, \mu) = K(\mu, \lambda) e^{\theta \cdot \mu - \lambda \psi(\theta)}$ .

*Démonstration.* Soit

$$\begin{aligned} f_i(x) &= \int \pi(\theta \mid \lambda_i, \mu_i) f(x \mid \theta) d\theta \\ &= K(\lambda_i, \mu_i) \int e^{\theta \mu_i - \lambda_i \psi(\theta)} h(x) e^{\theta \mu - \lambda \psi(\theta)} d\theta \\ &= K(\lambda_i, \mu_i) / K(\lambda_i + 1, \mu_i + x) \end{aligned}$$

de telle sorte que

$$\pi(\theta \mid \lambda_i, \mu_i) f(x \mid \theta) = f_i(x) \pi_i(\theta \mid \lambda_i + 1, \mu_i + x).$$

On a donc

$$\pi(\theta \mid x) \propto \pi(\theta) f(x \mid \theta) = \sum_i \omega_i \pi(\theta \mid \lambda_i, \mu_i) f(x \mid \theta) = \sum_i \omega_i f_i(x) \pi_i(\theta \mid \lambda_i + 1, \mu_i + x)$$

ce qui donne le résultat avec  $\omega'_i(x) \propto \omega_i f_i(x)$ .  $\square$

L'intérêt de ce résultat est que les combinaisons linéaires sont denses dans l'ensemble des distributions : en théorie, on peut donc approcher n'importe quelle loi a priori par un mélange de lois conjuguées, cf. par exemple [4, Théorème 3.24]. On étudie maintenant plus en détail le cas gaussien.

#### 4.3.2 Modèle gaussien à variance connue

Il s'agit de l'exemple du fil rouge. En fait, le Lemme 1.3 nous dit que pour le modèle paramétrique

$$f(x \mid \theta) \propto_x \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2\right),$$

la famille de lois normales forme une famille conjuguée.

#### 4.3.3 Modèle gaussien à moyenne connue

On considère maintenant le cas où la moyenne  $\mu$  est connue, et l'on cherche à estimer la variance  $\theta = \sigma^2$ . Le modèle paramétrique est alors

$$f(x \mid \theta) = \frac{1}{\sqrt{2\pi\theta^n}} \exp\left(-\frac{1}{2\theta} \sum_{i=1}^n (x_i - \mu)^2\right).$$

**Définition 4.6.** La loi Gamma de paramètre  $\alpha, \beta > 0$ , notée  $\mathcal{G}(\alpha, \beta)$ , est la loi sur  $\mathbb{R}$  de densité

$$g(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \mathbf{1}(x > 0).$$

Pour  $\alpha = 1$  on retrouve donc la loi exponentielle.

**Définition 4.7.** La loi Inverse Gamma de paramètre  $\alpha, \beta > 0$ , notée  $\mathcal{IG}(\alpha, \beta)$ , est la loi de  $1/Z$  avec  $Z \sim \mathcal{G}(\alpha, \beta)$ .

En utilisant la formule de changement de variable

$$f_{\varphi(X)}(y) = f_X(\varphi^{-1}(y)) |(\varphi^{-1})'(y)|$$

avec  $\varphi(x) = \varphi^{-1}(x) = 1/x$ , on obtient donc que  $\mathcal{IG}(\alpha, \beta)$  est la loi de densité

$$g(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{-\alpha-1} e^{-\beta/x} \mathbb{1}(x > 0)$$

**Lemme 4.8.** *Les lois inverse gamma forment une famille de lois conjuguées pour le modèle paramétrique gaussien à moyenne connue.*

*Démonstration.* Si  $\pi \sim \mathcal{IG}(\alpha, \beta)$ , on a alors

$$\pi(\theta | x) \propto \frac{1}{\theta^{n/2+\alpha+1}} \exp\left(-\frac{1}{2\theta} \sum_{i=1}^n (x_i - \mu)^2 - \beta/\theta\right)$$

et l'on reconnaît la loi  $\mathcal{IG}(\alpha + n/2, \beta + \sum (x_i - \mu)^2/2)$  □

**Lemme 4.9.** *Les lois gamma forment une famille de lois conjuguées pour le modèle paramétrique gaussien à moyenne connue et paramétrée par la précision  $\tau = 1/\sigma^2$ .*

#### 4.3.4 Modèle gaussien à moyenne et précision inconnues

On considère maintenant le modèle à moyenne et précision inconnues, i.e.,  $\theta = (\mu, \tau)$  avec  $\tau = 1/\sigma^2$  et

$$f(x | \mu, \tau) \propto \tau^{n/2} \exp\left(-\frac{\tau}{2} \sum_{i=1}^n (x_i - \mu)^2\right).$$

Une idée naturelle est de considérer comme loi a priori le produit des deux lois conjuguées précédentes, i.e.,

$$\pi(\mu, \tau) \propto \exp\left(-\frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) \times \sigma^{\alpha+1} \exp(-\beta\tau).$$

Néanmoins, cela mènerait à la loi a posteriori

$$\pi(\mu, \tau | x) \propto \tau^{n/2+\alpha+1} \exp\left(-\frac{\tau}{2} \sum_{i=1}^n (x_i - \mu)^2 - \frac{(\mu - \mu_0)^2}{2\sigma_0^2} - \beta\tau\right)$$

Néanmoins, à cause du terme “croisé”

$$\tau \sum_{i=1}^n (x_i - \mu)^2$$

on voit que cette densité a posteriori ne se factorise pas, i.e., ne peut pas s'écrire comme le produit d'une loi normale et d'une loi gamma. Néanmoins, on va montrer dans le résultat suivant que, si au lieu de prendre  $\mu$  qui suit une loi normale indépendante de  $\tau$ ,  $\mu$  suit une loi normale *conditionnellement* à  $\tau$ , alors on a bien une famille de lois conjuguées.

**Lemme 4.10.** Si  $\tau \sim \mathcal{G}(\alpha, \beta)$  et  $\mu \mid \tau \sim \mathcal{N}(\mu_0, n_0\tau)$ , alors

$$\mu \mid \tau, x \sim \mathcal{N}\left(\frac{n}{n+n_0}\bar{x} + \frac{n_0}{n+n_0}\mu_0, \frac{1}{(n+n_0)\tau}\right)$$

et

$$\tau \mid x \sim \mathcal{G}\left(\alpha + \frac{n}{2}, \beta + \frac{1}{2}\sum_{i=1}^n (x_i - \bar{x})^2 + \frac{nn_0}{2(n+n_0)}(\bar{x} - \mu_0)^2\right)$$

*Démonstration.* D'après le Lemme 1.3, si la variance  $\sigma^2$  est fixe,

$$\mu \sim \mathcal{N}(\mu_0, \sigma_0^2) \text{ et } x \mid \mu \sim \mathcal{N}(\mu, \sigma^2),$$

alors

$$\mu \mid x \sim \mathcal{N}\left(\frac{\sigma_0^2}{\sigma_0^2 + \sigma^2/n}\bar{x} + \frac{\sigma^2/n}{\sigma_0^2 + \sigma^2/n}\mu_0, \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)^{-1}\right).$$

Dans le cas qui nous intéresse, même si  $\sigma^2$  est aléatoire, on peut appliquer ce résultat conditionnellement à  $\tau$ , et puisque  $\sigma^2 = 1/\tau$  et  $\sigma_0^2 = 1/(n_0\tau)$ , on obtient

$$\mu \mid x, \tau \sim \mathcal{N}\left(\frac{1/(n_0\tau)}{1/(n\tau) + 1/(n_0\tau)}\bar{x} + \frac{1/(n\tau)}{1/(n\tau) + 1/(n_0\tau)}\mu_0, (n_0\tau + n\tau)^{-1}\right)$$

ce qui donne bien le résultat. On regarde maintenant la loi de  $\tau$  conditionnellement à  $x$ . Pour cela on regarde la loi jointe de  $(\mu, \tau)$  puis on extrait la marginale :

$$\begin{aligned} \pi(\mu, \tau \mid x) &\propto \pi(\tau)\pi(\mu \mid \tau)f(x \mid \mu, \tau) \\ &\propto \tau^{\alpha-1}e^{-\beta\tau} \times (n_0\tau)^{1/2}e^{-n_0\tau(\mu-\mu_0)^2/2} \times \tau^{n/2}e^{-\frac{\tau}{2}\sum_{i=1}^n (x_i - \mu)^2} \\ &\propto \tau^{n/2+\alpha-1/2} \exp\left(-\beta\tau - \frac{n_0\tau}{2}(\mu - \mu_0)^2 - \frac{\tau}{2}\sum_{i=1}^n (x_i - \mu)^2\right) \\ &\propto \tau^{n/2+\alpha-1/2} \exp\left(-\beta\tau - \frac{\tau}{2}\sum_{i=1}^n (x_i - \bar{x})^2\right) \\ &\quad \times \exp\left(-\frac{n_0\tau}{2}(\mu - \mu_0)^2 - \frac{n\tau}{2}(\bar{x} - \mu)^2\right) \end{aligned}$$

On calcule

$$\begin{aligned} n_0(\mu - \mu_0)^2 + n(\bar{x} - \mu)^2 &= (n+n_0)\left(\mu^2 - \frac{2\mu}{n+n_0}(n_0\mu_0 + n\bar{x})\right) + n_0\mu_0^2 + n\bar{x}^2 \\ &= (n+n_0)\left(\mu - \frac{1}{n+n_0}(n_0\mu_0 + n\bar{x})\right)^2 \\ &\quad - \frac{1}{n+n_0}(n_0\mu_0 + n\bar{x})^2 + n_0\mu_0^2 + n\bar{x}^2 \end{aligned}$$

et

$$-\frac{1}{n+n_0}(n_0\mu_0 + n\bar{x})^2 + n_0\mu_0^2 + n\bar{x}^2 = \frac{nn_0}{n+n_0}(\mu_0^2 + \bar{x}^2 - 2\mu_0\bar{x}) = \frac{nn_0}{n+n_0}(\mu_0 - \bar{x})^2$$

et donc

$$\begin{aligned} & \int \exp\left(-\frac{\tau}{2} (n_0(\mu - \mu_0)^2 + n(\bar{x} - \mu)^2)\right) d\mu \\ & \propto \exp\left(-\frac{\tau}{2} \frac{nn_0}{n+n_0} (\mu_0 - \bar{x})^2\right) \int e^{-\tau(n+n_0)\mu^2/2} d\mu \\ & \propto \tau^{-1/2} \exp\left(-\frac{\tau}{2} \frac{nn_0}{n+n_0} (\mu_0 - \bar{x})^2\right) \end{aligned}$$

et donc finalement,

$$\pi(\tau | x) \propto \tau^{n/2+\alpha-1} \exp\left(-\beta\tau - \frac{\tau}{2} \sum_{i=1}^n (x_i - \bar{x})^2 - \frac{\tau}{2} \frac{nn_0}{n+n_0} (\mu_0 - \bar{x})^2\right)$$

On reconnaît bien la densité de la loi Gamma attendue.  $\square$

#### 4.4 Loix a priori impropres

Les lois a priori impropres apparaissent naturellement lorsque l'on considère des modèles paramétriques invariants, par exemple par translation. C'est le cas notamment de l'exemple fil rouge, pour lequel

$$f(x | \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\theta)^2/(2\sigma^2)} = f(x - \theta | 0)$$

Dans ce cas, il est naturel de chercher une loi a priori  $\pi$  qui respecte cette invariance, i.e., telle que  $\pi(\theta - a) = \pi(\theta)$  pour tout  $a \in \mathbb{R}$ . La solution de cette équation est bien évidemment la fonction constante  $\pi(\theta) = \pi(0)$ , qui présente le problème de ne pas être intégrable (sauf pour  $\pi \equiv 0$ ) : on dit que  $\pi$  est une loi a priori **impropre**. Ainsi, une loi a priori impropre ne peut pas être interprétée comme une distribution de probabilité.

Néanmoins, le fait que  $\pi$  ne soit pas intégrable ne présente pas forcément de problème puisque cela n'empêche pas forcément la distribution a priori d'être bien définie. En effet, pour que la mesure a posteriori  $\pi(\theta | x)$  puisse être proportionnelle à  $\pi(\theta)f(x | \theta)$ , tout ce qu'il faut est que cette fonction soit intégrable, i.e.,  $\int \pi(\theta)f(x | \theta)d\theta < \infty$ . Si c'est le cas, on peut alors effectivement définir

$$\pi(\theta | x) = \frac{\pi(\theta)f(x | \theta)}{\int \pi(\theta)f(x | \theta)}.$$

C'est bien le cas dans l'exemple du fil rouge, puisque l'on a

$$\int \pi(\theta)f(x | \theta)d\theta = \pi(0) \int \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\theta)^2/(2\sigma^2)} d\theta = \pi(0)$$

ce qui donne donc

$$\pi(\theta | x) \propto f(x | \theta).$$

On remarque en particulier que la valeur de la constante n'a pas d'importance.

Il existe de nombreuses raisons de considérer des lois impropres (cf. [4, Chapitre 1.5]), parmi lesquelles :

1. dans le cas précédent, la mesure de Lebesgue est une généralisation naturelle de la loi uniforme, justifiée comme on l'a vu par la méthode du maximum d'entropie et défendue par Laplace dans son principe de la raison insuffisante. Dit autrement, en l'absence de toute information il est naturel de vouloir considérer la mesure uniforme sur  $\mathbb{R}$ , même si celle-ci n'est pas intégrable ;
2. dans de nombreux cas, les mesures impropres apparaissent comme limite de mesures propres. Elles peuvent donc être interprétées comme un cas extrême où la précision de l'information a priori a complètement disparu. Dans l'exemple du fil rouge, on peut ainsi voir la mesure de Lebesgue comme limite d'un loi normale centrée et de variance  $\sigma^2 \rightarrow \infty$  ;
3. comme on le verra plus loin, elles apparaissent naturellement dans le cadre des lois non-informatives telles que la loi de Jeffreys.

Un autre exemple de loi a priori impropre qui apparaît par des considérations d'invariance est la mesure  $\pi(\sigma) = 1/\sigma$  sur  $\mathbb{R}_+$ , motivée par des lois qui satisfont une relation d'invariance d'échelle

$$f(x | \theta) = \frac{1}{\theta} f\left(\frac{x}{\theta} | 1\right). \quad (4.1)$$

C'est par exemple le cas du modèle gaussien centré à variance inconnue, où

$$f(x | \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-x^2/(2\sigma^2)} = \frac{1}{\sigma} f\left(\frac{x}{\sigma} | 1\right).$$

Comme dans le cas de l'invariance par translation, la relation d'invariance d'échelle (4.1) suggère de prendre une loi a priori qui satisfait l'équation fonctionnelle

$$\pi(\theta) = \frac{1}{A} \pi\left(\frac{\theta}{A}\right)$$

dont la solution est donnée par  $\pi(\theta) = A/\theta$ .

---

**Exemple 4.27.**

Dans le cas du modèle gaussien centré, cela donne alors

$$\pi(\theta | x) \propto \frac{1}{\theta} \times \frac{1}{\theta^{n/2}} e^{-\|x\|^2/(2\theta)}.$$

Le symbole de proportionnalité est alors bien définie puisque pour tout  $n \geq 1$  on a bien

$$\int \frac{1}{\theta^{n/2+1}} e^{-\|x\|^2/(2\theta)} d\theta = \int \eta^{n/2-1} e^{-\eta\|x\|^2/2} d\eta < \infty.$$

---

**Fin exemple 4.27.**

Dans l'exemple ci-dessus, dans le cas  $n = 1$  l'intégrabilité n'est en fait obtenue que pour  $x \neq 0$ , mais comme cet évènement est de mesure nulle il semble raisonnable de le négliger. Dans le cas discret l'exemple suivant montre que la situation peut être plus subtile.

---

**Exemple 4.28.**

[4, Exemple 1.27] Soit une observation binomiale  $x \sim \text{Bin}(n, p)$ . Quelques auteurs contestent le choix de Laplace de la loi uniforme sur  $[0, 1]$  comme distribution a priori

automatique, car celle-ci apparaît comme étant biaisée contre les valeurs extrêmes 0 et 1. Ils proposent de considérer plutôt l'a priori de Haldane (1931)

$$\pi^*(p) \propto \frac{1}{p(1-p)}.$$

Dans ce cas, la distribution a posteriori est bien définie si la loi marginale

$$f(x) = \int_0^1 \frac{1}{p(1-p)} \binom{n}{x} p^x (1-p)^{n-x} dp$$

est finie, ce qui est le cas pour  $x \neq 0, n$ . La difficulté en 0 peut être résolue en notant que  $\pi^*$  apparaît comme limite de lois beta dénormalisées :

$$\pi_{\alpha,\beta}(p) = p^{\alpha-1} (1-p)^{\beta-1}$$

lorsque  $\alpha$  et  $\beta$  tendent vers 0. Ces distributions donnent comme loi a posteriori  $\mathcal{B}e(\alpha + x, \beta + n - x)$ , malgré l'absence de facteur normalisant, puisque le choix de ce tte constante n'a pas d'impact. La distribution a posteriori  $\pi_{\alpha,\beta}(p | x)$  a pour espérance

$$\delta_{\alpha,\beta}^\pi(x) = \frac{x + \alpha}{\alpha + \beta + n}$$

qui tend vers  $x/n$  quand  $\alpha$  et  $\beta$  tendent vers 0. Si la moyenne a posteriori est la quantité d'intérêt, nous pouvons alors étendre la procédure inférentielle aux cas  $x = 0$  et  $x = n$  en considérant également  $x/n$  comme un estimateur bayésien (uniquement) formel.

---

**Fin exemple 4.28.**

## 4.5 Zellner

## 4.6 Lois a priori non-informatives

### 4.6.1 Lois invariantes et a priori de Laplace

Lorsqu'aucune information n'est disponible sur la loi a priori, il est raisonnable de la choisir à partir de la seule information disponible, à savoir le modèle d'échantillonnage  $\{f(\cdot | \theta) : \theta \in \Theta\}$  : on parle dans ce cas de lois a priori *non-informatives*. Les lois impropres découlant de principes d'invariance rencontrées dans la section 4.4 rentrent elles aussi dans le cadre des lois non-informatives, puisqu'elles n'ont été choisies qu'à partir d'information disponible sur le modèle d'échantillonnage.

La loi uniforme préconisée par Laplace, ou la mesure de Lebesgue dans le cas d'un espace non compact, peut elle aussi être considérée comme une loi non-informative. Néanmoins, une de ses limitations majeures, déjà évoquées dans la section 4.2 sur le maximum d'entropie

### 4.6.2 La loi a priori de Jeffreys

$I(\theta)$  mesure la capacité du modèle à discriminer entre  $\theta$  et  $\theta \pm d\theta$  via la pente moyenne de  $\log f(x | \theta)$ . Favoriser les valeurs de  $\theta$  pour lesquelles  $I(\theta)$  est plus grande équivaut à minimiser l'influence de la loi a priori et est donc aussi non informatif que possible. On a donc envie de prendre comme loi a priori une loi  $\pi(\theta) \propto \varphi(I(\theta))$  avec  $\varphi$  croissante. On souhaite aussi satisfaire l'invariance par reparamétrisation. Si  $\pi$  était la loi d'une variable aléatoire  $X$ , i.e.,  $\pi = f_X$ , on

voudrait alors  $\pi = f_{h(X)}$  avec  $h$  difféomorphisme. Or le théorème de changement de variable nous assure que

$$f_{h(X)}(y) = f_X(h^{-1}(y)) |(h^{-1})'(y)| = \frac{f_X(h^{-1}(y))}{|h'(h^{-1}(y))|}$$

et donc la contrainte d'invariance par reparamétrisation impose

$$\pi(h(x)) = \pi(x)/h'(x)$$

et donc

$$\varphi(I(h(x))) = \varphi(I(x))/h'(x) = \varphi(I(x)/h'(x)^2)$$

d'après la Proposition 3.7, ce qui impose  $\varphi(x) = \sqrt{x}$ .

**Définition 4.11.** L'a priori de Jeffreys est défini par

$$\pi^*(\theta) \propto \sqrt{I(\theta)}.$$

En fait, la loi de Jeffreys est fréquemment impropre et donc la constante qui apparaît n'est pas importante.

**Exemple fil rouge 4.29.** 

---

On a dans ce cas

$$f(x | \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\theta)^2/(2\sigma^2)}$$

et donc le score est défini par

$$\ell(\theta; x) = \partial_\theta \log f(x | \theta) = -\frac{x - \theta}{\sigma^2}$$

et donc l'information de Fisher par

$$I(\theta) = \int \ell(\theta, x)^2 f(x | \theta) dx = \frac{1}{\sigma^2}$$

qui est une constante : on retrouve donc bien la loi impropre suggérée par l'invariance d'échelle.

**Fin exemple fil rouge 4.29.**

**Exemple 4.30.** 

---

On considère maintenant le cas du modèle gaussien centré à variance inconnue :

$$f(x | \theta) = \frac{1}{\sqrt{2\pi\theta}} e^{-x^2/(2\theta)}.$$

Alors

$$\ell(\theta; x) = \partial_\theta \log f(x | \theta) = -\frac{1}{\theta} + \frac{x^2}{2\theta^2}$$

et donc

$$I(\theta) = \mathbb{E}_\theta \left[ \left( -\frac{1}{\theta} + \frac{x^2}{2\theta^2} \right)^2 \right] = \mathbb{E}_1 \left[ \left( -\frac{1}{\theta} + \frac{x^2}{2\theta} \right)^2 \right] = \frac{1}{\theta^2} \mathbb{E}_1 [(x-1)^2] = \frac{1}{2\theta^2}$$

ce qui donne  $\pi(\theta) \propto 1/\sigma$ , encore une fois, comme suggéré par la relation d'invariance.

**Fin exemple 4.30.**

**Exemple 4.31.** 

---

**TODO**

Modèle binomial

**Fin exemple 4.31.**

Si l'on veut pouvoir regarder le modèle gaussien avec moyenne et variance inconnues, il faut pouvoir généraliser l'information de Fisher et l'a priori de Jeffreys en plus grande dimension. Lorsque  $\theta \in \mathbb{R}^d$  avec  $d \geq 1$ , l'information de Fisher  $I(\theta)$  est une matrice, définie comme

$$- \left( \mathbb{E}_\theta \left( \frac{\partial^2 \log f(x | \theta)}{\partial \theta_i \partial \theta_j} \right) \right)_{1 \leq i, j \leq d}$$

L'a priori de Jeffreys est alors donnée par  $\sqrt{\det(I(\theta))}$ .

**Exemple 4.32.**

Dans le cas du modèle gaussien avec moyenne et variance inconnues, on a alors  $\theta = (\mu, \sigma^2)$  et

$$f(x | \theta) = \frac{1}{\sqrt{2\pi\theta_2}} e^{-(x-\theta_1)^2/(2\theta_2)}$$

et donc

$$\log f(x | \theta) = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \theta_2 - \frac{(x - \theta_1)^2}{2\theta_2}$$

et donc

$$\begin{cases} \frac{\partial^2 \log f(x | \theta)}{\partial \theta_1^2} = -\frac{1}{\theta_2}, \\ \frac{\partial^2 \log f(x | \theta)}{\partial \theta_2^2} = \frac{1}{2\theta_2^2} - \frac{(x - \theta_1)^2}{\theta_2^3}, \\ \frac{\partial^2 \log f(x | \theta)}{\partial \theta_1 \partial \theta_2} = -\frac{\theta_1 - x}{\theta_2^2}, \end{cases}$$

et donc

$$I(\theta) = - \begin{pmatrix} -1/\theta_2 & \mathbb{E}_\theta(x - \theta_1)/\theta_2^2 \\ \mathbb{E}_\theta(x - \theta_1)/\theta_2^2 & 1/(2\theta_2^2) - \mathbb{E}_\theta(x - \theta_1)^2/\theta_2^2 \end{pmatrix} = \begin{pmatrix} 1/\theta_2 & 0 \\ 0 & 1/(2\theta_2^2) \end{pmatrix}$$

et donc l'a priori de Jeffreys est donné par  $\sqrt{1/\theta_2^3} = 1/\sigma^{3/2}$ .

**Fin exemple 4.32.**

## 5 Comportement asymptotique des estimateurs bayésiens

On s'intéresse au comportement des estimateurs bayésiens dans le régime asymptotique  $n \rightarrow \infty$ .

### 5.1 Régularité et différentiabilité en moyenne quadratique

On va examiner des conditions techniques sous lesquelles les estimateurs convergent. Ces conditions impliquent des conditions de régularité que l'on énonce ici dans le cas à densité sur  $\mathbb{R}$ , i.e., on considère un modèle paramétrique  $\{f_\theta, \theta \in \Theta\}$  avec  $f_\theta$  densité sur  $\mathbb{R}$ . On généralise au cas vectoriel en remplaçant principalement les dérivées par des gradients.

**Définition 5.1.** Le modèle paramétrique  $\{f_\theta, \theta \in \Theta\}$  est dit régulier si :

1.  $\Theta$  est un ouvert de  $\mathbb{R}^k$  ;
2. pour tout  $x$ ,  $\theta \mapsto f_\theta(x)$  est continûment dérivable, de dérivée  $\partial_\theta f_\theta$  ;
3. la fonction du score  $\ell_\theta = \partial_\theta f_\theta / f_\theta$  satisfait  $\int \ell_\theta(x)^2 f_\theta(x) dx < \infty$  ;
4. l'information de Fisher  $I(\theta) = \int \ell_\theta(x)^2 f_\theta(x) dx$  est strictement positive et continue en  $\theta$ .

Lorsque le modèle est régulier dans le sens ci-dessus, alors il est aussi différentiable en moyenne quadratique, cf. par exemple [7, Lemme 7.6].

**Définition 5.2.** Un modèle paramétrique  $\{f_\theta, \theta \in \Theta\}$  est dit différentiable en moyenne quadratique s'il existe une fonction mesurable  $\dot{\ell}_\theta$  telle que, lorsque  $\theta \rightarrow \theta_0$ ,

$$\int \left[ \sqrt{f_\theta(x)} - \sqrt{f_{\theta_0}(x)} - \frac{1}{2}(\theta - \theta_0)\dot{\ell}_{\theta_0}(x)\sqrt{f_{\theta_0}(x)} \right]^2 dx = o((\theta - \theta_0)^2).$$

Un des intérêts de la notion plus générale de différentiabilité en moyenne quadratique est de s'affranchir de l'hypothèse de continuité de  $\theta \mapsto f_\theta(x)$ , qui n'est par exemple pas satisfaite dès lors que les lois n'ont pas toutes le même support (cf. par exemple la loi uniforme, qui n'est d'ailleurs pas non plus différentiable en moyenne quadratique, cf. [7, Exemple 7.9]). A noter par ailleurs que si le modèle est régulier, alors on a en fait  $\dot{\ell}_\theta = \ell_\theta$  puisque

$$\partial_\theta \sqrt{f_\theta} = \frac{\partial_\theta f_\theta}{2\sqrt{f_\theta}} = \frac{\partial_\theta f_\theta}{2f_\theta} \sqrt{f_\theta}.$$

Mais, encore une fois, la différentiabilité quadratique ne nécessite pas la différentiabilité de  $f_\theta$ .

### 5.2 Comportement asymptotique de l'estimateur du maximum de vraisemblance

L'estimateur du maximum de vraisemblance joue un rôle important dans l'étude du comportement asymptotique des estimateurs bayésiens, et on commence donc pas son étude. On définit

$$L_n(\theta) = \frac{1}{n} \log f_\theta(x)$$

la log-vraisemblance pondérée. La loi des grands nombres donne alors le résultat suivant.

**Lemme 5.3.** *Si les observations sont i.i.d. tirées selon la vraie loi  $f_{\theta_0}$ , alors pour tout  $\theta \in \Theta$  on a  $L_n(\theta) \xrightarrow{\text{P.S.}} L(\theta) := \mathbb{E}_{\theta_0}(\log f_{\theta}(x_1)) = -\text{CE}(f_{\theta_0}, f_{\theta})$ .*

Ainsi, le résultat d'optimisation (3.1) donne  $\theta_0 = \arg \max_{\theta} L(\theta)$  et puisque  $\hat{\theta}^{\text{MV}} = \arg \max L_n$ , la convergence  $\hat{\theta}^{\text{MV}} \rightarrow \theta_0$  est naturelle. Néanmoins, la convergence ponctuelle  $L_n(\theta) \xrightarrow{\text{P.S.}} L(\theta)$  n'est pas suffisante pour assurer que le maximiseur de  $L_n$  converge vers le maximiseur de  $L$ . Il faut typiquement une convergence fonctionnelle plus forte et une unicité du maximiseur. On a par exemple le résultat suivant.

**Théorème 5.4** ([7, Théorème 5.7]). *Soit  $M_n$  des fonctions aléatoires et  $M$  une fonction déterministe. Si pour tout  $\varepsilon > 0$  on a*

$$\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{\mathbb{P}} 0 \quad \text{et} \quad \sup_{\theta: d(\theta, \theta_0) \geq \varepsilon} M(\theta) < M(\theta_0) \quad (5.1)$$

alors  $\hat{\theta}_n^{\text{MV}} \xrightarrow{\mathbb{P}} \theta_0$  avec  $\hat{\theta}_n^{\text{MV}} = \arg \max M_n$ .

*Démonstration.* Puisque  $\|M_n - M\|_{\infty} \xrightarrow{\mathbb{P}} 0$  on a en particulier  $M_n(\theta_0) \xrightarrow{\mathbb{P}} M(\theta_0)$  et donc, puisque  $\theta_0$  maximise  $M$ , on a

$$0 \leq M(\theta_0) - M(\hat{\theta}_n^{\text{MV}}) \leq M_n(\theta_0) - M(\hat{\theta}_n^{\text{MV}}) + o_{\mathbb{P}}(1).$$

Puisque  $\hat{\theta}_n^{\text{MV}} \in \arg \max M_n$  cela donne

$$0 \leq M(\theta_0) - M(\hat{\theta}_n) \leq M_n(\hat{\theta}_n^{\text{MV}}) - M(\hat{\theta}_n^{\text{MV}}) + o_{\mathbb{P}}(1) \leq \|M_n - M\|_{\infty} + o_{\mathbb{P}}(1)$$

et donc  $M(\theta_0) - M(\hat{\theta}_n) \xrightarrow{\mathbb{P}} 0$ . Par ailleurs, la deuxième condition implique que si  $M(\theta)$  est proche de  $M(\theta_0)$ , alors nécessairement  $\theta$  est proche de  $\theta_0$ . Ainsi, la convergence  $M(\hat{\theta}_n) \xrightarrow{\mathbb{P}} M(\theta_0)$  implique que  $\hat{\theta}_n \xrightarrow{\mathbb{P}} \theta_0$ .  $\square$

**Remarque 5.6.**

Sous les hypothèses (5.1), on a en fait  $\hat{\theta}_n \rightarrow \theta_0$  pour tout estimateur  $\hat{\theta}_n$  qui satisfait  $\hat{\theta}_n \geq M_n(\theta_0) - o_{\mathbb{P}}(1)$ .

**Fin remarque 5.6.**

La deuxième hypothèse du théorème ci-dessus implique en particulier que  $M$  admet un unique maximiseur, par ailleurs bien séparé. Concernant l'estimateur du maximum de vraisemblance, cette hypothèse est vérifiée si le modèle est identifiable par la Proposition 3.3 et il s'agit donc de vérifier la première condition dans (5.1). Pour comprendre cette convergence uniforme, on fait un détour via le théorème de Glivenko–Cantelli.

**5.2.1 Théorème et classe de Glivenko–Cantelli**

Si on définit  $\mathcal{F} = \{f_{\theta} : \theta \in \Theta\}$ ,  $\mathbb{P}_n f = \frac{1}{n} \sum_{k=1}^n f(x_k)$  la moyenne sous la mesure empirique et  $\mathbb{P} f = \mathbb{E}_{\theta_0} f(x)$  la moyenne sous la vraie valeur du paramètre, on peut donc réécrire la condition de convergence uniforme qui apparaît dans (5.1)

$$\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{\mathbb{P}} 0$$

en

$$\sup \{ |\mathbb{P}_n f - \mathbb{P} f| : f \in \mathcal{F} \} \xrightarrow{\mathbb{P}} 0. \quad (5.2)$$

Cela mène à la définition générale suivante.

**Définition 5.5.** Une classe de fonctions mesurables  $\mathcal{F}$  est dite  $\mathbb{P}$ -Glivenko-Cantelli si la convergence (5.2) est satisfaite.

Lorsque  $\mathcal{F}$  est l'ensemble des fonctions indicatrices d'intervalles de la forme  $(-\infty, t]$  on obtient le théorème de Glivenko–Cantelli que certains d'entre vous ont déjà vu puisqu'il est à la base du test de Kolmogorov–Smirnov. En effet, pour  $\mathcal{F} = \{x \mapsto \mathbb{1}(x \leq t) : t \in \mathbb{R}\}$  on a

$$\sup \{ |\mathbb{P}_n f - \mathbb{P} f| : f \in \mathcal{F} \} = \sup_{t \in \mathbb{R}} |F_n(t) - F(t)|$$

avec  $F_n(t) = \frac{1}{n} \sum_{k=1}^n \mathbb{1}(x_k \leq t)$  la fonction de répartition empirique et  $F(t) = \mathbb{P}_{\theta_0}(x \leq t)$  la vraie fonction de répartition.

**Théorème 5.6** (Théorème de Glivenko–Cantelli). *Si  $x_1, x_2, \dots$  sont i.i.d. de fonction de répartition, alors  $\|F_n - F\|_\infty \xrightarrow{\mathbb{P}} 0$ .*

*Démonstration.* Pour  $\varepsilon > 0$  donné, il existe une partition  $-\infty = t_0 < t_1 < \dots < t_k = \infty$  telle que  $F(t_i) - F(t_{i-1}) < \varepsilon$ . Pour  $t_{i-1} \leq t < t_i$  on a donc

$$F_n(t_{i-1}) - F(t_{i-1}) - \varepsilon \leq F_n(t) - F(t) \leq F_n(t_i) - F(t_i) + \varepsilon$$

et donc

$$\limsup_n \|F_n - F\|_\infty \leq \varepsilon$$

puisque  $F_n(t_i) \xrightarrow{\text{P.s.}} F(t_i)$  et  $F_n(t_i) \xrightarrow{\text{P.s.}} F(t_i)$  pour chaque  $i$  par la loi forte des grands nombres.  $\square$

L'idée de la preuve ci-dessus est de trouver des “enveloppes” qui entourent la limite attendue  $F$ .

**Définition 5.7.** Un crochet  $[\ell, u]$  est l'ensemble des fonctions  $\{f : \ell \leq f \leq u\}$ . Un crochet  $[\ell, u]$  est un  $\varepsilon$ -crochet si  $\mathbb{P}(u - \ell) \leq \varepsilon$ . Le nombre crochétant  $N(\varepsilon, \mathcal{F})$  est le nombre minimal de  $\varepsilon$ -crochets nécessaires pour couvrir  $\mathcal{F}$ .

**Théorème 5.8** ([7, Théorème 19.4]). *Si  $N(\varepsilon, \mathcal{F}) < \infty$  pour tout  $\varepsilon > 0$ , alors  $\mathcal{F}$  est  $\mathbb{P}$ -Glivenko–Cantelli.*

### 5.2.2 Retour sur l'estimateur du maximum de vraisemblance

On utilise maintenant les résultats ci-dessus pour établir des résultats de convergence de l'estimateur du maximum de vraisemblance.

**Théorème 5.9** (Consistance du MV, [7, Exemple 19.8]). *Si  $\Theta$  est compact,  $\sup_{\theta} f_{\theta}$  est intégrable et  $\theta \mapsto f_{\theta}(x)$  est continue pour chaque  $x$ , alors  $\{f_{\theta} : \theta \in \Theta\}$  est  $\mathbb{P}$ -Glivenko–Cantelli, i.e.,  $L_n \rightarrow L$  uniformément en probabilités.*

*En particulier, si les hypothèses ci-dessus sont satisfaites,  $L$  est continue et le modèle est identifiable alors  $\hat{\theta}_n^{\text{MV}} \xrightarrow{\mathbb{P}} \theta_0$ .*

*Démonstration.* On montre que la première partie des hypothèses implique que le nombre crochetant est fini. Pour  $B \subset \Theta$  une boule ouverte, on définit  $f_B(x) = \inf_{\theta \in B} f_\theta(x)$  et  $f^B(x) = \sup_{\theta \in B} f_\theta(x)$ . Si  $B_m$  est une suite de boules avec un centre commun  $\theta$  et un rayon décroissant vers 0, alors on a  $f^{B_m} - f_{B_m} \downarrow 0$  par la continuité supposée, pour chaque  $x$  et donc dans  $L_1$  par convergence dominée (en utilisant l'hypothèse  $\sup_\theta f_\theta$  intégrable). Ainsi, pour  $\varepsilon > 0$  fixé, pour chaque  $\theta$  il existe une boule ouverte  $B$  centrée sur  $\theta$  telle que le crochet  $[f_B, f^B]$  a une taille au plus  $\varepsilon$ . Par la compacité de  $\Theta$ , on peut extraire des  $B$  que l'on vient de construire un recouvrement fini. Les crochets correspondant recouvrent  $\mathcal{F}$  et donc le nombre crochetant est fini.

Ainsi, pour avoir la convergence de l'estimateur du maximum de vraisemblance il ne reste qu'à montrer la deuxième condition dans (5.1), ce qui suit directement de la continuité et du fait que  $\theta_0$  est l'unique maximiseur de la log-vraisemblance lorsque le modèle est identifiable.  $\square$

### Exemple fil rouge 5.33.

Pour  $f_\theta(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\theta)^2/(2\sigma^2)}$ ,  $\theta \mapsto f_\theta(x)$  est bien intégrable pour tout  $x \in \mathbb{R}$ . Par contre,  $\sup_{\theta \in \mathbb{R}} f_\theta(x) = (2\pi\sigma^2)^{-1/2}$  n'est pas intégrable mais si on se restreint à  $\Theta = [\theta_*, \theta_0]$  on a alors

$$\sup_{\theta_* \leq \theta \leq \theta_0} f_\theta(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \begin{cases} 1 & \text{si } \theta_* \leq x \leq \theta_0, \\ e^{-(x-\theta_0)^2/(2\sigma^2)} & \text{si } \theta_* \leq x < \theta_0, \end{cases}$$

avec  $\theta_0 = \theta_*$  ou  $\theta_0 = \theta_0$  en fonction de  $x$ . Ainsi, en se restreignant à un compact on a bien que  $\sup_\theta f_\theta$  est intégrable.

**Fin exemple fil rouge 5.33.**

On rappelle maintenant un résultat de normalité asymptotique. Dans le reste de ces notes,  $N$  désigne une variable aléatoire normale standard.

**Théorème 5.10** (Normalité asymptotique du MV, [7, Théorème 5.39]). *Si les conditions suivantes sont satisfaites :*

- le modèle est identifiable et différentiable en moyenne quadratique ;
- il existe une fonction mesurable  $\ell$  avec  $\mathbb{P}\ell^2 < \infty$  telle que pour tous  $\theta_1, \theta_2$  dans un voisinage de  $\theta_0$  :

$$|\log f_{\theta_1}(x) - \log f_{\theta_2}(x)| \leq \ell(x) \|\theta_1 - \theta_2\| ;$$

- $I(\theta_0) > 0$  ;
- $\hat{\theta}_n^{\text{MV}}$  est consistant ;

alors  $\sqrt{n}(\hat{\theta}_n^{\text{MV}} - \theta_0) \xrightarrow{L} I(\theta)^{-1/2} N$ .

*Eléments de démonstration.* On présente l'idée de la preuve lorsque  $f_\theta$  est deux fois continûment dérivable, et on note  $L'_n$  et  $L''_n$  et  $L'$  et  $L''$  les dérivées première et seconde de  $L_n$  et  $L$ , respectivement. Puisque  $L'_n(\hat{\theta}^{\text{MV}}) = 0$ , on obtient par l'expansion de Taylor

$$0 = L'_n(\hat{\theta}^{\text{MV}}) = L'_n(\theta_0) + (\hat{\theta}^{\text{MV}} - \theta_0)L''_n(\theta_1)$$

pour un certain  $\theta_1$  entre  $\theta_0$  et  $\hat{\theta}_n^{\text{MV}}$ , et donc

$$\sqrt{n}(\theta_0 - \hat{\theta}) = \frac{\sqrt{n}L'_n(\theta_0)}{L''_n(\theta_1)}$$

On a

$$\sqrt{n}L'_n(\theta_0) = \frac{1}{\sqrt{n}} \sum_{k=1}^n \partial_{\theta} \log f_{\theta}(x_i) |_{\theta=\theta_0}$$

et puisque  $\mathbb{E}_{\theta_0}(\partial_{\theta} \log f_{\theta}(x_1) |_{\theta=\theta_0}) = L'(\theta_0) = 0$  puisque  $\theta_0$  maximise  $L$ , on obtient  $\sqrt{n}L'_n(\theta) \xrightarrow{L} \sqrt{\mathbb{E}_{\theta_0}(\partial_{\theta} \log f_{\theta_0}(x)^2)}N = I(\theta_0)^{1/2}N$ . Par ailleurs,

$$L''_n(\theta_1) = \frac{1}{n} \sum_{i=1}^n \partial_{\theta}^2 \log f_{\theta_1}(x_i) \rightarrow \mathbb{E}_{\theta_0} \partial_{\theta}^2 \log f_{\theta}(x) = -I(\theta_0)$$

en utilisant pour la convergence le fait que  $\hat{\theta}_n^{\text{MV}} \xrightarrow{\mathbb{P}} \theta_0$ , ce qui donne le résultat.  $\square$

### 5.3 Comportements asymptotiques des estimateurs bayésiens

Ibragimov et Has'minskiï [3, Théorème II.2.1] présentent un résultat très général garantissant qu'une suite d'estimateurs bayésiens est consistante et asymptotiquement normale. Nous nous restreindrons ici à une discussion sur l'estimateur du maximum a posteriori, sur la normalité asymptotique de la densité a posteriori (théorème de Bernstein-von Mises) et à une discussion rapide sur la moyenne a posteriori.

#### 5.3.1 Comportement asymptotique du MAP

Par définition, on a

$$\hat{\theta}_n^{\text{MAP}} \in \arg \max_{\theta} M_n(\theta) \text{ avec } M_n(\theta) = L_n(\theta) + \frac{1}{n} \log \pi(\theta).$$

On a donc  $M_n(\theta) \xrightarrow{\text{P.S.}} L(\theta)$  mais en général, même si  $L_n \xrightarrow{\text{P.S.}} L$  uniformément, à cause du terme  $\frac{1}{n} \log \pi(\theta)$  on n'aura pas convergence uniforme de  $M_n$  vers  $L$  (sauf si  $\sup |\log \pi| < \infty$ , i.e.,  $\sup \pi < \infty$  et  $\inf \{\pi(\theta) : \theta \text{ t.q. } \pi(\theta) > 0\} > 0$ ). Pour éviter la deuxième condition qui est très forte, on se ramène à la maximisation de la vraisemblance (et non de la log-vraisemblance) ce qui nécessite une hypothèse supplémentaire que  $L$  est bornée. Ces hypothèses doivent probablement pouvoir être allégées mais je n'ai pas trouvé de référence<sup>1</sup>.

#### TODO

Trouver de meilleures hypothèses, par exemple  $\theta \in L^1$  ?

**Théorème 5.11** (Consistance du MAP). *Si les hypothèses du Théorème 5.9 sont satisfaites, que  $L$  et  $\pi$  sont bornées et que  $\pi(\theta_0) > 0$ , alors  $\hat{\theta}_n^{\text{MAP}} \xrightarrow{\mathbb{P}} \theta_0$ .*

1. Par exemple, en utilisant la remarque (5.6), on pourrait s'affranchir de l'hypothèse  $L$  bornée si l'on pouvait montrer que  $\frac{1}{n} \log \pi(\hat{\theta}_n^{\text{MAP}}) \rightarrow 0$ , ce qui est consistant avec l'hypothèse  $\pi(\theta_0) > 0$  et le comportement attendu  $\hat{\theta}_n^{\text{MV}}$ .

*Démonstration.* Sous les hypothèses du Théorème 5.9, on a  $\|L_n - L\|_\infty \xrightarrow{\mathbb{P}} 0$  et donc  $\|e^{L_n} - e^L\|_\infty \xrightarrow{\mathbb{P}} 0$  aussi, puisque

$$\left| e^{L_n(\theta)} - e^{L(\theta)} \right| \leq e^{L(\theta)} \left( e^{\|L_n - L\|_\infty} - 1 \right),$$

et donc  $M_n = (\pi(\theta))^{1/n} e^{L_n(\theta)} \xrightarrow{\mathbb{P}} e^{L(\theta)}$  uniformément. Puisque  $\hat{\theta}^{\text{MAP}} = \arg \max M_n$  on peut donc appliquer le Théorème 5.4.  $\square$

### TODO

Préciser les hypothèses et la conclusion.

**Théorème 5.12** (Normalité asymptotique du MAP). *Sous certaines hypothèses, on a  $\sqrt{n}(\hat{\theta}^{\text{MAP}} - \theta_0)$  converge vers une loi normale.*

### Exemple 5.34.

On considère l'exemple de Bayes avec une loi a priori  $\pi \sim \mathcal{Be}(\alpha, \beta)$  et un modèle d'échantillonnage  $\mathcal{Bin}(n, p)$ , si bien que la loi a posteriori est  $\mathcal{Be}(\alpha + |x|, n + \beta - |x|)$ . On a calculé le MAP dans l'exemple 2.19 :

$$\hat{\theta}_n^{\text{MAP}} = \frac{\alpha - 1}{n + \alpha + \beta - 2} + \frac{n}{n + \alpha + \beta - 2} \bar{x}$$

sur lequel on voit bien directement la consistance et la normalité asymptotique (et seulement asymptotique cette fois, contrairement à l'exemple du fil rouge).

**Fin exemple 5.34.**

### 5.3.2 Comportement asymptotique de la densité a posteriori

Pour discuter de la consistance de la densité a posteriori on commence par reprendre le modèle de Bernoulli.

### Exemple 5.35.

On continue l'exemple 5.34, dans lequel

$$\begin{aligned} \pi(\theta | x) &\propto \theta^{|x| + \alpha - 1} (1 - \theta)^{n - |x| + \beta - 1} \\ &= \exp((|x| + \alpha - 1) \log \theta + (n - |x| + \beta - 1) \log(1 - \theta)) \\ &= \exp((n + \alpha + \beta - 2)(p \log \theta + (1 - p) \log(1 - \theta))) \end{aligned}$$

avec  $p = (|x| + \alpha - 1)/(n + \alpha + \beta - 2)$ . On a donc

$$\pi(\theta | x) \propto \exp(-(n + \alpha + \beta - 2) \text{CE}(\mathcal{Ber}(p), \mathcal{Ber}(\theta))).$$

Ainsi, on a

$$\frac{\pi(\theta | x)}{\pi(\theta' | x)} = \exp[-(n + \alpha + \beta - 2) (\text{CE}(\mathcal{Ber}(p), \mathcal{Ber}(\theta)) - \text{CE}(\mathcal{Ber}(p), \mathcal{Ber}(\theta')))]$$

et puisque l'entropie croisée  $\text{CE}(\mathcal{Ber}(p), \mathcal{Ber}(\theta))$  est minimisée en  $\theta = p$ , on en déduit que

$$\mathbb{P}(p - \varepsilon \leq \theta \leq p + \varepsilon | x) \rightarrow 1.$$

Puisque  $p \rightarrow \theta_0$  presque sûrement par la loi forte des grands nombres, on en déduit donc que  $\theta \rightarrow \theta_0$  dans un certain sens.

**Fin exemple 5.35.**

La consistance est valable en tout généralité dès lors que les observations sont à valeurs dans un espace euclidien.

**Théorème 5.13** (Consistance de la densité a posteriori, [7, Théorème 10.10]).  
*Si le modèle est identifiable, alors pour toute loi a priori  $\pi$  sur  $\Theta$  la suite des lois a posteriori est consistante pour  $\pi$ -presque tout  $\theta$ , i.e., pour  $\pi$ -presque tout  $\theta$  et toute fonction continue bornée  $f$  on a*

$$\int f(\eta)\pi(\eta | x)d\eta \xrightarrow{\mathbb{P}} f(\theta).$$

La limitation de ce résultat est que la convergence n'est garantie que pour  $\pi$ -presque tout  $\theta$ . On considère maintenant un exemple plus simple.

L'exemple suivant est tiré de [1]. Dans cet exemple, on ne suppose pas que la vraie distribution appartient au modèle paramétrique considéré : dans ce cas, la loi a posteriori se concentre autour du paramètre qui minimise la distance de Kullback–Leibler.

**Théorème 5.14.** *Si  $\Theta$  est compact et  $A$  est un voisinage de  $\theta_0$  avec  $\pi(A) > 0$ , alors  $\pi(A | x) \xrightarrow{\text{p.s.}} 1$ , où  $\theta_0 = \arg \min D(f_{\theta_0}, f_\theta)$ .*

*Démonstration.* On se ramène au cas discret en recouvrant  $A$  d'un nombre fini (puisque  $A$  est compact) de boules dont une seule contient  $\theta_0$ . Il suffit de montrer que  $\mathbb{P}(\theta = \theta_0 | x) \rightarrow 1$  lorsque  $\Theta$  est fini et  $\pi(\theta_0) = \mathbb{P}(\theta = \theta_0) > 0$ . On a

$$\log \left( \frac{\pi(\theta | x)}{\pi(\theta_0 | x)} \right) = \log \left( \frac{\pi(\theta)}{\pi(\theta_0)} \right) + \sum_{i=1}^n \log \left( \frac{f(x_i | \theta)}{f(x_i | \theta_0)} \right)$$

et donc

$$\frac{1}{n} \log \left( \frac{\pi(\theta | x)}{\pi(\theta_0 | x)} \right) \xrightarrow{\text{p.s.}} -D(f_{\theta_0}, f_\theta).$$

Puisque  $D(f_{\theta_0}, f_\theta) < 0$  pour  $\theta \neq \theta_0$ , cela implique que  $\pi(\theta | x)/\pi(\theta_0 | x) \xrightarrow{\text{p.s.}} 0$  pour  $\theta \neq \theta_0$  et par suite, puisque  $\Theta$  est fini, que  $\pi(\theta_0 | x) \xrightarrow{\text{p.s.}} 1$ .  $\square$

Concernant la normalité asymptotique, il s'agit du Théorème de Bernstein–von Mises. Encore une fois, l'intuition est assez claire : on a

$$\begin{aligned} f(x | \theta) &= f(x | \hat{\theta}^{\text{MV}}) \exp \left( -n(L_n(\theta) - L_n(\hat{\theta}^{\text{MV}})) \right) \\ &\approx f(x | \hat{\theta}^{\text{MV}}) \exp \left( -\frac{1}{2}n(\theta - \hat{\theta}^{\text{MV}})^2 L_n''(\hat{\theta}^{\text{MV}}) \right) \\ &= f(x | \hat{\theta}^{\text{MV}}) \exp \left( -\frac{1}{2}n(\theta - \hat{\theta}^{\text{MV}})^2 I(\theta_0) \right) \end{aligned}$$

et donc

$$\pi(\theta | x) = \pi(\theta) \frac{f(x | \theta)}{f(x)} \propto \pi(\theta) \exp \left( -\frac{1}{2}n(\theta - \hat{\theta}^{\text{MV}})^2 I(\theta_0) \right)$$

Or si  $\bar{f}(\theta | x)$  est la densité de  $\sqrt{n}(\theta - \hat{\theta}_n^{\text{MV}})$  conditionnellement à  $x^2$ , on a

$$\bar{f}(\theta | x) = \frac{1}{\sqrt{n}} \pi \left( \hat{\theta}^{\text{MV}} + \frac{\theta}{\sqrt{n}} | x \right) \propto \pi \left( \hat{\theta}^{\text{MV}} + \frac{\theta}{\sqrt{n}} \right) \exp \left( -\frac{1}{2}\theta^2 I(\theta_0) \right)$$

---

2.  $\hat{\theta}^{\text{MV}}$  étant une fonction mesurable de  $x$ , on a n'a pas besoin de préciser la structure de dépendance entre  $\theta$  et  $\hat{\theta}_n^{\text{MV}}$ .

Cela suggère donc que  $\sqrt{n}(\theta - \hat{\theta}^{\text{MV}})$  suit asymptotiquement une loi normale, ce qui est bien le cas. Il existe plusieurs versions de ce résultat appelé théorème de Bernstein-von Mises.

**Théorème 5.15** ([7, Théorème 10.1]). *On suppose que :*

- l'échantillon  $x = (x_1, \dots, x_n)$  est constitué de  $n$  observations i.i.d. de loi  $f_\theta$  ;
- le modèle paramétrique  $\{f_\theta, \theta \in \Theta\}$  est différentiable en moyenne quadratique en  $\theta_0$  avec  $I(\theta_0) > 0$  ;
- $\pi$  est absolument continue dans un voisinage de  $\theta_0$  avec  $\pi(\theta_0) > 0$  ;
- pour tout  $\varepsilon > 0$  il existe une suite de tests  $\phi_n$  tels que  $\mathbb{P}_{\theta_0} \phi_n \rightarrow 0$  et  $\sup_{|\theta - \theta_0| \geq \varepsilon} \mathbb{P}_n(1 - \phi_n) \rightarrow 0$  pour tout  $\varepsilon > 0$ .

Alors la loi a posteriori de  $\sqrt{n}(\theta - \hat{\theta}_n^{\text{MV}})$  sachant  $x$  converge en variation totale en probabilité vers une loi gaussienne centrée de variance  $I(\theta)^{-1}$ , i.e.,

$$\mathbb{P}_{\theta_0} \left( \int \left| \frac{1}{\sqrt{n}} \pi \left( \hat{\theta}^{\text{MV}} + \frac{t}{\sqrt{n}} \mid x \right) - \frac{I(\theta_0)^{1/2}}{\sqrt{2\pi}} e^{-t^2 I(\theta_0)/2} \right| dt \geq \varepsilon \right) \xrightarrow{n \rightarrow \infty} 0.$$

Dans ce résultat, on compare le paramètre à un estimateur du maximum de vraisemblance, et donc sous des hypothèses garantissant la consistance de l'estimateur du maximum de vraisemblance on obtient que la loi de  $\theta$  conditionnelle aux observations suit une loi gaussienne. On peut aussi avoir des résultats de convergence uniforme sur les compacts de la densité a posteriori. Par exemple, sous certaines hypothèses techniques Schervish [6, Théorème 7.89] obtient le résultat suivant :

$$\mathbb{P}_{\theta_0} \left( \sup_{t \in B} |\pi'(t \mid x) - \phi(t)| \geq \varepsilon \right) \xrightarrow{n \rightarrow \infty} 0$$

pour tout  $\varepsilon > 0$  et tout compact  $B$ , où  $\phi$  est la densité de la loi normale standard et  $\pi'$  la densité de  $\sqrt{n}(\theta - \hat{\theta}_n^{\text{MV}})$ .

Lorsque qu'une loi converge, il est naturel de s'attendre à ce que la moyenne associée converge aussi. Ainsi, sous des hypothèses supplémentaires, un corollaire du théorème de Bernstein-von Mises est que la moyenne a posteriori est un estimateur consistant et asymptotiquement normal.

### 5.3.3 Application du Théorème de Bernstein-von Mises à la loi de Jeffreys : a priori de référence

Pour une loi a priori  $\pi$  on définit

$$G(\pi) = \mathbb{E}[D(\pi(\cdot \mid x) \mid \pi)] = \mathbb{E}\{\mathbb{E}[D(\pi(\cdot \mid x) \mid \pi) \mid \theta]\}$$

Or, le Théorème 5.15 de Bernstein-von Mises nous assure que sous  $\mathbb{P}_\theta$ ,

$$\pi(\theta' \mid x) = \sqrt{n} \pi'(\sqrt{n}(\theta' - \hat{\theta}) \mid x) \approx \sqrt{\frac{nI(\theta)}{2\pi}} e^{-nI(\theta)(\theta' - \hat{\theta})^2/2}$$

et donc en utilisant  $\hat{\theta} \approx \theta$

$$\begin{aligned}
\mathbb{E} [D(\pi(\cdot | x) || \pi) | \theta] &\approx \int \sqrt{\frac{nI(\theta)}{2\pi}} e^{-nI(\theta)(\theta' - \theta)^2/2} \log \left( \frac{\sqrt{\frac{nI(\theta')}{2\pi}} e^{-nI(\theta')(\theta' - \theta)^2/2}}{\pi(\theta')} \right) d\theta' \\
&= \int \sqrt{\frac{I(\theta)}{2\pi}} e^{-I(\theta)u^2/2} \log \left( \frac{\sqrt{\frac{nI(\theta)}{2\pi}} e^{-I(\theta)u^2/2}}{\pi(\theta + u/\sqrt{n})} \right) du \\
&\approx \int \sqrt{\frac{I(\theta)}{2\pi}} e^{-I(\theta)u^2/2} \log \left( \frac{\sqrt{\frac{nI(\theta)}{2\pi}} e^{-I(\theta)u^2/2}}{\pi(\theta)} \right) du \\
&= \frac{1}{2} \log \left( \frac{n}{2\pi} \right) + \log \left( \frac{\sqrt{I(\theta)}}{\pi(\theta)} \right) du - \frac{I(\theta)}{2} \int u^2 \sqrt{\frac{I(\theta)}{2\pi}} e^{-I(\theta)u^2/2} du \\
&= \frac{1}{2} \log \left( \frac{n}{2\pi} \right) + \log \left( \frac{\sqrt{I(\theta)}}{\pi(\theta)} \right) du - \frac{I(\theta)}{2} \int u^2 \sqrt{\frac{I(\theta)}{2\pi}} e^{-I(\theta)u^2/2} du \\
&= \frac{1}{2} \log \left( \frac{n}{2\pi} \right) + \log \left( \frac{\sqrt{I(\theta)}}{\pi(\theta)} \right) du - \frac{1}{2}
\end{aligned}$$

et donc

$$\begin{aligned}
G(\pi) &\approx \mathbb{E} \left[ \frac{1}{2} \log \left( \frac{n}{2\pi} \right) + \log \left( \frac{\sqrt{I(\theta)}}{\pi(\theta)} \right) - \frac{1}{2} \right] \\
&= \frac{1}{2} \log \left( \frac{n}{2\pi} \right) + \int \log \left( \frac{\sqrt{I(\theta)}}{\pi(\theta)} \right) \pi(\theta) d\theta - \frac{1}{2}
\end{aligned}$$

et donc si l'on cherche à minimiser cela, on trouve l'a priori de Jeffreys.

### 5.3.4 Moyenne a posteriori

On s'intéresse maintenant à l'estimateur bayésien  $\hat{\theta}_n = \mathbb{E}(\theta | x_n)$  donné par la moyenne a posteriori. Puisque

$$\hat{\theta}_n = \int \eta \pi(\eta | x) d\eta,$$

on est presque dans le cadre d'application du Théorème 5.13, sauf qu'on voudrait l'appliquer à  $f(\eta) = \eta$  qui n'est pas bornée. En fait, le cœur de la preuve du Théorème 5.13 consiste à prouver que  $\theta$  est mesurable par rapport à  $x_\infty = (x_1, x_2, \dots)$ , i.e., qu'il existe une fonction mesurable  $h$  telle que

$$\theta = h(x_\infty).$$

Une fois ce résultat admis, on peut alors prouver la consistance de la moyenne a posteriori grâce à un argument de martingale. En effet, la suite  $(\hat{\theta}_n, n \geq 1)$  vérifie une propriété bien particulière, à savoir

$$\mathbb{E}(\hat{\theta}_{n+1} | x_n) = \mathbb{E}[\mathbb{E}(\theta | x_{n+1}) | x_n] = \mathbb{E}(\theta | x_n) = \hat{\theta}_n.$$

La deuxième égalité vient de la propriété de la tour, à savoir

$$\mathbb{E}[\mathbb{E}(X | Y, Z) | Y] = \mathbb{E}(X | Y).$$

On a donc l'égalité remarquable  $\mathbb{E}(\hat{\theta}_{n+1} | x_n) = \hat{\theta}_n$  qui est caractéristique d'une structure de martingale. Dans le cas présent, où  $\hat{\theta}_n = \mathbb{E}(\theta | x_n)$ , on a le résultat fondamental suivant.

**Théorème 5.16.** *Si  $\mathbb{E}(|\theta|) < \infty$  et le modèle est identifiable, alors  $\hat{\theta}_n$  converge presque sûrement et dans  $L^1$  vers  $\mathbb{E}(\theta | x_1, x_2, \dots)$ .*

Puisque  $\theta$  est mesurable par rapport à  $x_\infty$  on obtient bien la consistance de  $\hat{\theta}_n$ . Concernant la normalité asymptotique on se contentera à nouveau de calculs heuristiques. On a

$$\hat{\theta} = \int \theta \pi(\theta | x) d\theta$$

et donc le changement de variable  $\eta = \sqrt{n}(\theta - \hat{\theta}^{\text{MV}})$  donne

$$\hat{\theta} = \frac{1}{\sqrt{n}} \int \left( \hat{\theta}^{\text{MV}} + \frac{\eta}{\sqrt{n}} \right) \pi \left( \hat{\theta}^{\text{MV}} + \frac{\eta}{\sqrt{n}} | x \right) d\eta$$

et donc

$$\sqrt{n}(\hat{\theta} - \hat{\theta}^{\text{MV}}) = \int \frac{\eta}{\sqrt{n}} \pi \left( \hat{\theta}^{\text{MV}} + \frac{\eta}{\sqrt{n}} | x \right) d\eta.$$

Le Théorème de Bernstein–von Mises suggère donc que  $\sqrt{n}(\hat{\theta} - \hat{\theta}^{\text{MV}}) \xrightarrow{L} 0$ .

#### TODO

Préciser les hypothèses.

**Théorème 5.17.**  $\sqrt{n}(\hat{\theta} - \hat{\theta}^{\text{MV}}) \xrightarrow{L} 0$  et donc  $\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{L} \mathcal{N}(0, 1/I(\theta_0))$  si l'estimateur du maximum de vraisemblance converge.

## 5.4 Inégalité de van Trees

On présente dans cette section l'inégalité de van Trees, qui peut servir à borner inférieurement le risque minimax (puisque  $\max_{\theta} R_f(\delta, \theta) \geq \int R_f(\delta, \theta) \pi(\theta) d\theta$ ). Cf. [2, 8, 9, 10] pour des généralisations (notamment les bornes de Weiss–Weinstein) et le lien avec l'inégalité de Fréchet–Darmois–Cramer–Rao.

**Théorème 5.18.** *Si le modèle est régulier, que  $I(\theta) > 0$  et que  $\hat{g}$  est un estimateur de carré intégrable et sans biais de  $g(\theta)$ , alors pour tout  $\theta \in \Theta$  on a  $\text{Var}_{\theta} \hat{g} \geq \frac{g'(\theta)^2}{nI(\theta)}$ .*

*Démonstration.* Puisque  $\theta = \mathbb{E}_{\theta} \hat{g} = \int g(x) f_{\theta}(x) dx$  on obtient en dérivant

$$\begin{aligned} 1 &= \int g(x) \partial_{\theta} f_{\theta}(x) dx \\ &= \int (g(x) - \theta) \partial_{\theta} f_{\theta}(x) dx \\ &\leq \left\{ \int (g(x) - \theta)^2 f_{\theta}(x) dx \right\} \times \left\{ \int \frac{\partial_{\theta} f_{\theta}(x)^2}{f_{\theta}(x)} dx \right\} \end{aligned}$$

en utilisant  $\int \partial_\theta f_\theta(x) dx = 0$  pour la deuxième égalité et Cauchy–Schwarz dans  $L_2(f_\theta)$  pour la dernière inégalité, ce qui donne le résultat.  $\square$

**Théorème 5.19** (Inégalité de van Trees). *On suppose que le modèle est différentiable en moyenne quadratique en tout  $\theta$ , et que  $\theta \mapsto f(x | \theta)$  est  $C^1$  pour tout  $x$ . On suppose que  $\pi$  est dérivable sur  $[a, b]$ , nulle aux bords, et on note*

$$J = \mathbb{E}[\log \pi(\theta)'^2] = \int \frac{\pi'(\theta)^2}{\pi(\theta)} d\theta.$$

On suppose en outre que  $g$  est  $C^1$  sur  $[a, b]$  et telle que  $\mathbb{E}(|g'(\theta)|) < \infty$ . Alors pour tout estimateur  $T$  on a

$$\mathbb{E}[(T - g(\theta))^2] \geq \frac{[\mathbb{E}(g'(\theta))]^2}{\mathbb{E}(I(\theta)) + J}$$

avec  $I(\theta) = \mathbb{E}[\log f(x | \theta)'^2 | \theta]$  l'information de Fisher.

*Démonstration.* On a

$$\mathbb{E}(g'(\theta)) = \int g'(\theta)\pi(\theta)d\theta = \int \left\{ \int g'(\theta)\pi(\theta)f(x | \theta)d\theta \right\} dx.$$

Par intégration par parties,

$$\begin{aligned} \int g'(\theta)\pi(\theta)f(x | \theta)d\theta &= [(g(\theta) - T(x))\pi(\theta)f(x | \theta)]_a^b \\ &\quad - \int (g(\theta) - T(x)) \partial_\theta(\pi(\theta)f(x | \theta))d\theta \end{aligned}$$

Par hypothèse, les termes de bord sont nuls et donc

$$\begin{aligned} \mathbb{E}(g'(\theta)) &= - \int (g(\theta) - T(x)) (\pi'(\theta)f(x | \theta) + \pi(\theta)\partial_\theta f(x | \theta)) d\theta dx \\ &= - \int (g(\theta) - T(x)) \left( \frac{\pi'(\theta)}{\pi(\theta)} + \frac{\partial_\theta f(x | \theta)}{f(x | \theta)} \right) \pi(\theta)f(x | \theta) d\theta dx \end{aligned}$$

puis par Cauchy–Schwarz,

$$\begin{aligned} [\mathbb{E}(g'(\theta))]^2 &\leq \left\{ \int (g(\theta) - T(x))^2 \pi(\theta)f(x | \theta) d\theta dx \right\} \\ &\quad \times \left\{ \int \left( \frac{\pi'(\theta)}{\pi(\theta)} + \frac{\partial_\theta f(x | \theta)}{f(x | \theta)} \right)^2 \pi(\theta)f(x | \theta) d\theta dx \right\} \end{aligned}$$

Le premier terme donne le terme  $\mathbb{E}[(T - g(\theta))^2]$ . Quant au second terme, on obtient en développant le carré

$$\int \left( \frac{\pi'(\theta)}{\pi(\theta)} + \frac{\partial_\theta f(x | \theta)}{f(x | \theta)} \right)^2 \pi(\theta)f(x | \theta) d\theta dx = J + I(\theta) + \int \pi'(\theta)\partial_\theta f(x | \theta) d\theta dx$$

et puisque  $\int \partial_\theta f(x | \theta) dx = 0$  on obtient bien le résultat voulu.  $\square$

## 6 Les tests d'hypothèses bayésiens

Le cadre théorique des tests d'hypothèse est le suivant : on se donne deux sous-ensembles disjoints  $\Theta_0, \Theta_1 \subset \Theta$  de l'ensemble des paramètres, et l'on souhaite décider si le vrai paramètre appartient à  $\Theta_0$  ou à  $\Theta_1$ . Formellement, il s'agit donc essentiellement de faire de l'inférence sur la fonction  $\mathbf{1}(\theta \in \Theta_0)$ , ou bien, en terme de théorie de la décision, de prendre une décision à valeurs dans  $\{0, 1\}$ , mais ce cadre inférentiel présente certaines particularités propres qui motivent son étude séparée.

### 6.1 Cadre fréquentiste

Dans le cadre fréquentiste, **tous les raisonnements et calculs sont conditionnés par une hypothèse**. On se pose donc des questions du genre : *Si l'hypothèse nulle ou alternative est vraie, quelle est la probabilité de tel ou tel événement ?* Par exemple, le risque de première espèce est la probabilité de rejeter l'hypothèse nulle en supposant qu'elle est vraie. Ainsi, un risque de première espèce faible garantit qu'on ne rejettera pas l'hypothèse nulle à tort. Le risque de deuxième espèce est quant à lui la probabilité de rejeter l'hypothèse alternative en supposant qu'elle est vraie. Ces deux risques sont antagoniques, i.e., avoir un risque de première espèce faible se traduit en général par un risque de deuxième espèce élevé. L'approche classique veut que l'on fixe le risque de premier espèce, puis que l'on choisisse le test qui, à ce risque de première espèce fixé, présente le plus faible risque de deuxième espèce. Cette approche induit donc une asymétrie très forte entre les hypothèses nulle et alternative, qui ne jouent pas des rôles interchangeables. Ainsi, intervertir les deux rôles peut mener à des conclusions différentes. Avec cette approche, où le risque de première espèce est privilégié, il est cohérent de choisir une hypothèse nulle "conservative", i.e., telle que les conséquences de l'accepter à tort sont limitées (puisque le risque de la rejeter est faible).

Par ailleurs, les tests d'hypothèse fréquentistes sont souvent utilisés, ou tout du moins interprétés, à tort, et l'on est souvent enclins de parler de probabilité d'une hypothèse, par exemple d'interpréter un risque comme la probabilité de l'hypothèse : dans le cadre fréquentiste, cette interprétation n'a aucun sens. Néanmoins, dans le cadre bayésien cela devient tout à fait normal.

### 6.2 Cadre bayésien

#### 6.2.1 Estimateurs bayésiens

On retrouve cela à l'aide de la théorie de la décision. Neyman et Pearson ont proposé le coût 0-1 :

$$L(\theta, d) = \begin{cases} 0 & \text{si } d = \mathbf{1}(\theta \in \Theta_0), \\ 1 & \text{sinon.} \end{cases}$$

Une approche plus générale consiste à pénaliser différemment les erreurs de type I et II, ce qui correspond à la fonction de coût suivante :

$$L(\theta, d) = \begin{cases} 0 & \text{si } d = \mathbf{1} (\theta \in \Theta_0), \\ a_0 & \text{si } \theta \in \Theta_0 \text{ et } d = 0, \\ a_1 & \text{si } \theta \in \Theta_1 \text{ et } d = 1. \end{cases}$$

L'estimateur de Bayes associé est alors donné par le résultat suivant.

**Proposition 6.1.** *Sous le coût  $a_0$ - $a_1$  ci-dessus, l'estimateur de Bayes associé à la loi a priori  $\pi$  est*

$$\delta_\pi(x) = \begin{cases} 1 & \text{si } \mathbb{P}(\theta \in \Theta_0 | x) > \frac{a_1}{a_0 + a_1}, \\ 0 & \text{sinon.} \end{cases}$$

*Démonstration.* On utilise le Théorème 2.11 : on a

$$\begin{aligned} \mathbb{E}[L(\theta, d) | x] &= a_0 \mathbf{1}(d = 0) \mathbb{P}(\theta \in \Theta_0 | x) + a_1 \mathbf{1}(d = 1) \mathbb{P}(\theta \in \Theta_1 | x) \\ &= a_0(1 - d) \mathbb{P}(\theta \in \Theta_0 | x) + a_1 d \mathbb{P}(\theta \in \Theta_1 | x) \end{aligned}$$

et donc

$$\arg \min_{d \in \{0,1\}} \mathbb{E}[L(\theta, d) | x] = \begin{cases} 1 & \text{si } a_0 \mathbb{P}(\theta \in \Theta_0 | x) \geq a_1 \mathbb{P}(\theta \in \Theta_1 | x) \\ 0 & \text{sinon} \end{cases}$$

□

On voit donc que la décision est uniquement basée sur la probabilité a posteriori que l'hypothèse soit vraie, ce qui est naturel dans le cadre bayésien. Par ailleurs, la décision ne dépend que de la fonction de coût via le ratio  $a_0/a_1$ , contrairement au cas fréquentiste qui, en plus de la fonction de coût, nécessite de fixer un seuil  $\alpha$ . On remarque que plus  $a_0/a_1$  est grand, i.e., plus une réponse incorrecte est pénalisée sous  $H_0$  relativement à  $H_1$ , plus la probabilité a posteriori de  $H_0$  doit être petite pour être rejetée.

**Exemple fil rouge 6.36.**

Pour tester  $H_0 : \theta < 0$  dans l'exemple du fil rouge, on calcule à l'aide du Lemme 1.3

$$\mathbb{P}(\theta < 0 | x) = \mathbb{P}(\tau N + \mu(\bar{x}) < 0 | x) = \phi(-\mu(\bar{x})/\tau)$$

avec  $\mu(s) = ps + q\mu_0$  et  $\tau^2 = 1/\sigma_0^2 + n/\sigma^2$ , et  $H_0$  est donc acceptée lorsque  $-\mu(\bar{x}) > z_{a_0, a_1} \tau$  avec  $\phi(z_{a_0, a_1}) = a_1/(a_0 + a_1)$ , i.e.,

$$\bar{x} < -\frac{\sigma^2}{n\sigma_0^2} \mu_0 - \left(1 + \frac{\sigma^2}{n\sigma_0^2}\right) \tau z.$$

**Fin exemple fil rouge 6.36.**

### 6.2.2 Le facteur de Bayes

**Définition 6.2.** Le facteur de Bayes est le rapport des probabilités a posteriori des hypothèses nulle et alternative sur le rapport des probabilités a priori de ces mêmes hypothèses, soit

$$B_{01}^{\pi}(x) = \frac{\mathbb{P}(\theta \in \Theta_0 | x)}{\mathbb{P}(\theta \in \Theta_1 | x)} \bigg/ \frac{\pi(\theta \in \Theta_0)}{\pi(\theta \in \Theta_1)}$$

Ce rapport évalue la modification de la vraisemblance de l'ensemble  $\Theta_0$  par rapport à celle de l'ensemble  $\Theta_0$  due à l'observation et peut se comparer naturellement à 1, bien qu'une échelle de comparaison exacte doive être fondée sur une fonction de coût. Dans le cas particulier d'hypothèses simples où  $\Theta_0 = \{\theta_0\}$  et  $\Theta_1 = \{\theta_1\}$ , le facteur de Bayes se simplifie et devient le rapport de vraisemblance classique

$$B_{01}^{\pi}(x) = \frac{f(x | \theta_0)}{f(x | \theta_1)}.$$

De manière plus générale, ce rapport peut être perçu comme un rapport de vraisemblance bayésien, car, si  $\pi_i(\theta) \propto \mathbb{1}(\theta \in \Theta_i) \pi(\theta)$  est la loi a priori sous  $H_i$ , alors  $B_{01}^{\pi}(x)$  peut s'écrire

$$\begin{aligned} B_{01}^{\pi}(x) &= \frac{\int_{\Theta_0} f(x | \theta) \pi(\theta) d\theta}{\int_{\Theta_1} f(x | \theta) \pi(\theta) d\theta} \bigg/ \frac{\int_{\Theta_0} \pi(\theta) d\theta}{\int_{\Theta_1} \pi(\theta) d\theta} \\ &= \frac{\int_{\Theta_0} f(x | \theta) \pi(\theta) d\theta}{\int_{\Theta_0} \pi(\theta) d\theta} \bigg/ \frac{\int_{\Theta_1} f(x | \theta) \pi(\theta) d\theta}{\int_{\Theta_1} \pi(\theta) d\theta} \\ &= \frac{\int_{\Theta_0} f(x | \theta_0) \pi_0(\theta) d\theta}{\int_{\Theta_1} f(x | \theta_1) \pi_1(\theta) d\theta} = \frac{m_0(x)}{m_1(x)} \end{aligned}$$

ce qui revient donc à remplacer les vraisemblances par des marginales sous les deux hypothèses.

**Définition 6.3.** L'échelle de Jeffreys est la suivante :

1. si  $\log_{10}(B_{10}^{\pi})$  varie entre 0 et 0,5, la certitude que  $H_0$  est fautive est **faible** ;
2. si elle varie entre 0,5 et 1, cette certitude est **substantielle** ;
3. si elle est entre 1 et 2, elle est **forte** ;
4. si elle est au-dessus de 2, elle est **décisive**,

avec la même échelle en faveur de  $H_0$  pour les valeurs négatives.

### 6.2.3 Hypothèses nulles simples

Une hypothèse nulle simple du genre  $\Theta_0 = \{\theta_0\}$  peut ne pas avoir de sens, par exemple il ne semble pas y avoir de sens à se demander si la probabilité qu'il pleuve demain vaut 0,7163891256... Dans certains cas néanmoins, cela peut avoir un sens, par exemple dans le cas discret, ou dans le cas de sélection de modèle si l'on teste la nullité d'un paramètre, ou encore, en astrophysique, tester si l'univers est en expansion, s'il se contracte ou s'il est stable revient à tester si la constante de Hubble est plus grande, plus petite ou égale à une valeur spécifique  $h_0$ .

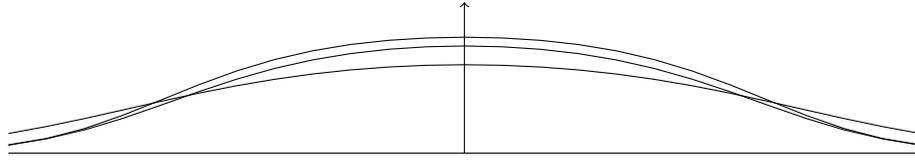


FIGURE 3 – Evolution de la probabilité a posteriori de  $H_0$  en fonction de  $z = x/\rho$  pour  $\varrho_0 = 1/2$  et différentes valeurs de  $\alpha = \sigma_0^2/\sigma^2$ .

Dans le cas à densité, la probabilité a posteriori d'une telle hypothèse est toujours nulle et il faut donc modifier la démarche. Pour cela, on procède à une modification de la loi a priori, qui est une démarche qui s'applique de manière plus générale dès qu'on teste un ensemble de mesure nulle. Dans le cas présent, on se donne  $\varrho_0$  la probabilité a priori que  $\theta = \theta_0$  et  $\pi_1$  la densité alternative, si bien que la loi a priori s'écrit

$$\pi(\theta) = \varrho_0 \mathbf{1}(\theta = \theta_0) + (1 - \varrho_0) \pi_1(\theta)$$

et la probabilité a posteriori de  $H_0$  est donnée par

$$\pi(\Theta_0 | x) = \frac{f(x | \theta_0) \varrho_0}{\int f(x | \theta) \pi(\theta) d\theta} = \frac{f(x | \theta_0) \varrho_0}{f(x | \theta_0) \varrho_0 + (1 - \varrho_0) \int_{\Theta_1} f(x | \theta) \pi_1(\theta) d\theta}$$

et le facteur de Bayes par

$$B_{01}^{\pi}(x) = \frac{f(x | \theta_0)}{m_1(x)}.$$

C'est ce que l'on obtient lorsque l'on considère  $\Theta_0 = \{\theta : |\theta - \theta_0| \leq \varepsilon\}$  et que l'on fait tendre  $\varepsilon$  vers 0.

### Exemple fil rouge 6.37.

Si dans l'exemple fil rouge on cherche maintenant à tester  $H_0 : \theta = 0$ , il semble raisonnable de prendre  $\pi_1 = \mathcal{N}(0, \sigma_0^2)$ . Alors

$$\frac{m_1(x)}{f(x | 0)} = \frac{\sigma}{\sqrt{\sigma^2 + \sigma_0^2}} \frac{e^{-x^2/(2(\sigma^2 + \sigma_0^2))}}{e^{-x^2/(2\sigma^2)}} = \frac{1}{\sqrt{1 + \alpha}} \exp\left(\frac{\alpha z^2}{2(1 + \alpha)}\right)$$

avec  $\alpha = \sigma_0^2/\sigma^2$  et  $z = x/\sigma$ , ce qui donne pour la probabilité

$$\pi(\theta = 0 | x) = \left[1 + \frac{1 - \varrho_0}{\varrho_0} \frac{1}{\sqrt{1 + \alpha}} \exp\left(\frac{\alpha z^2}{2(1 + \alpha)}\right)\right]^{-1}.$$

La Figure 3 montre l'évolution de la probabilité de  $H_0$  en fonction de  $z = x/\rho$  pour  $\varrho_0 = 1/2$  et différentes valeurs de  $\alpha = \sigma_0^2/\sigma^2$ . Paradoxalement, pour certaines valeurs de  $z$ , augmenter la variance a priori augmente la probabilité d'accepter  $H_0$ .

**Fin exemple fil rouge 6.37.**

## 6.2.4 Loi a priori impropres

Si l'on a ardemment défendu l'utilisation des lois a priori impropres pour l'inférence, ça n'est plus aussi clair dans le cadre des tests d'hypothèses. D'une

part, le simple fait de tester des hypothèses, et pas n'importe lesquelles, est en contradiction avec l'absence d'information supposée à la base des lois a priori impropres. D'autre part, l'utilisation de lois impropres dans le cadre des tests d'hypothèses mènent à des contradictions apparemment superficielles : par exemple, les résultats ne sont plus invariants par multiplication de la loi a priori par une constante.

Nous illustrons ce point sur l'exemple du fil rouge. On souhaite tester  $H_0 : \theta = 0$  contre  $H_1 : \theta \neq 0$ . Si nous utilisons la loi a priori impropre  $\pi_1(\theta) = 1$ , la loi a priori est alors

$$\pi(\theta) = \frac{1}{2} \mathbb{1}(\theta = 0) + \frac{1}{2}$$

et la probabilité a posteriori de  $H_0$  est

$$\pi(\theta = 0 | x) = \frac{e^{-x^2/2}}{e^{-x^2/2} + \int e^{-(x-\theta)^2/2} d\theta} = \frac{1}{1 + \sqrt{2\pi}e^{x^2/2}}$$

Ainsi, la probabilité a posteriori de  $H_0$  est bornée par  $1/(1 + \sqrt{2\pi}) \approx 0,285$ . Ceci implique que la loi a posteriori est plutôt biaisée contre  $H_0$ , même dans le cas le plus favorable.

En outre, une difficulté conceptuelle est que dans le cadre inférentiel précédent, les lois a priori impropres avaient été justifiées comme limite de lois propres, par exemple la mesure de Lebesgue comme limite de la densité normale avec une variance infinie. Néanmoins, dans le cadre des tests d'hypothèse cela n'est plus le cas. Par exemple, on a calculé dans l'Exemple 6.37 une probabilité a posteriori égale à

$$\pi(\theta = 0 | x) = \left[ 1 + \frac{1 - \varrho_0}{\varrho_0} \frac{1}{\sqrt{1 + \alpha}} \exp\left(\frac{\alpha z^2}{2(1 + \alpha)}\right) \right]^{-1}.$$

Lorsque la variance a priori  $\sigma_0^2$  tend vers  $+\infty$ , i.e., que  $\alpha \rightarrow \infty$ , cette probabilité tend vers 1 qui est à la fois inutile et différent de la réponse obtenue en considérant directement la loi a priori impropre. Ce comportement est une manifestation du paradoxe de Jeffreys-Lindley.

### 6.2.5 Régions de confiance

**Définition 6.4.** Une région de crédibilité de niveau  $\alpha \in ]0, 1[$  est un sous-ensemble  $R_\alpha = R_\alpha(x)$  de  $\Theta$  tel que  $\mathbb{P}(\theta \in R_\alpha | x) \geq 1 - \alpha$ . Une région Highest Posterior Density (HPD) de niveau  $\alpha$  est une région de crédibilité de niveau  $\alpha$  de la forme

$$\{\theta : \pi(\theta | x) \geq h\}$$

où  $h$  est le plus grand seuil tel que  $\mathbb{P}(\theta \in R | x) \geq 1 - \alpha$ .

Une fois de plus, le fait que, dans la formulation bayésienne,  $\theta$  ait une probabilité donnée d'appartenir à une région fixée  $R$  est plus attrayant que l'interprétation fréquentiste d'une région aléatoire ayant une probabilité donnée de contenir le paramètre inconnu  $\theta$ .

Il est assez clair d'après sa définition qu'une région HPD est de volume minimal parmi les régions  $\alpha$ -crédible. Une région HPD n'est en général pas connexe, mais on a le résultat suivant.

**Proposition 6.5.** *Si  $\pi(\cdot | x)$  est continue et n'a qu'un seul maximum local, alors toute région HPD est connexe.*

Si la densité a posteriori admet un unique maximum local dans  $R$ , les régions HPD sont des intervalles dont les bornes sont des quantiles de la loi a posteriori. Même si ces quantiles ont des formes explicites, comme pour certaines familles de lois conjuguées, calculer analytiquement des quantiles optimaux n'est pas forcément simple. On peut néanmoins obtenir une approximation des régions HPD en utilisant le comportement asymptotiquement gaussien de la densité a posteriori.

**Proposition 6.6.** *Sous les hypothèses du théorème de Bernstein-von Mises, l'intervalle*

$$I_\alpha = \left[ \hat{\theta}^{\text{MAP}} - q_{1-\alpha/2} \frac{1}{\sqrt{I(\hat{\theta}^{\text{MAP}})n}}, \hat{\theta}^{\text{MAP}} + q_{1-\alpha/2} \frac{1}{\sqrt{I(\hat{\theta}^{\text{MAP}})n}} \right]$$

est un intervalle de crédibilité de niveau  $\alpha$  asymptotiquement en probabilités :  $\mathbb{P}(\theta \in I_\alpha | x) \xrightarrow{\mathbb{P}} 1 - \alpha$ .

On ne fournit qu'une ébauche de preuve : si  $\Pi(\cdot | x)$  est la fonction de répartition associée à  $\pi(\cdot | x)$ , on a

$$\begin{aligned} \mathbb{P}(\theta \in I_\alpha | x) &= \mathbb{P}\left(\theta \leq \hat{\theta}^{\text{MAP}} + \frac{q}{\sqrt{In}} \mid x\right) - \mathbb{P}\left(\theta \leq \hat{\theta}^{\text{MAP}} - \frac{q}{\sqrt{In}} \mid x\right) \\ &= \mathbb{P}\left(\sqrt{n}(\theta - \hat{\theta}^{\text{MAP}}) \leq \frac{q}{\sqrt{I}} \mid x\right) - \mathbb{P}\left(\sqrt{n}(\theta - \hat{\theta}^{\text{MAP}}) \leq -\frac{q}{\sqrt{I}} \mid x\right) \end{aligned}$$

et donc le théorème de Bernstein von Mises donne

$$\mathbb{P}(\theta \in I_\alpha | x) \rightarrow \mathbb{P}(N \leq q | x) - \mathbb{P}(N \leq -q | x) = 1 - \alpha.$$

## A Tableau de quelques lois absolument continues

Loi	Paramètres	Densité	Moyenne
Loi gaussienne	$\sigma > 0, \mu \in \mathbb{R}$	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	$\mu$
Loi exponentielle	$\lambda > 0$	$\lambda e^{-\lambda x} \mathbf{1}(x > 0)$	$\frac{1}{\lambda}$
Loi exponentielle symétrique	$\lambda > 0$	$\frac{\lambda}{2} e^{-\lambda x }$	0
Loi de Cauchy	—	$\frac{1}{\pi(x^2 + 1)}$	Non définie
Loi Gamma	$\alpha, \beta > 0$	$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \mathbf{1}(x > 0)$	$\frac{\alpha}{\beta}$
Loi Beta	$\alpha, \beta > 0$	$\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)} x^{\alpha-1} (1-x)^{\beta-1} \mathbf{1}(0 < x < 1)$	$\frac{\alpha}{\alpha + \beta}$

TABLE 1 – Densités des lois continues usuelles

## Références

- [1] Andrew Gelman, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin. *Bayesian data analysis*. Texts in Statistical Science Series. CRC Press, Boca Raton, FL, third edition, 2014.
- [2] Richard D. Gill and Boris Y. Levit. Applications of the Van Trees inequality: a Bayesian Cramér-Rao bound. *Bernoulli*, 1(1-2):59–79, 1995.
- [3] I. A. Ibragimov and R. Z. Has'minskiĭ. *Statistical estimation*, volume 16 of *Applications of Mathematics*. Springer-Verlag, New York-Berlin, 1981. Asymptotic theory, Translated from the Russian by Samuel Kotz.
- [4] Christian P. Robert. *Le choix bayésien*. Statistique et probabilités appliquées. Springer, Paris, second edition, 2006.
- [5] R. J. Samworth. Steins's paradox. *Eureka*, 62:38–41, 2012.
- [6] Mark J. Schervish. *Theory of statistics*. Springer Series in Statistics. Springer-Verlag, New York, 1995.
- [7] A. W. van der Vaart. *Asymptotic statistics*, volume 3 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 1998.
- [8] E. Weinstein and A. J. Weiss. Lower bounds on the mean square estimation error. *Proceedings of the IEEE*, 73(9):1433–1434, Sept 1985.
- [9] Ehud Weinstein and Anthony J. Weiss. A general class of lower bounds in parameter estimation. *IEEE Trans. Inform. Theory*, 34(2):338–342, 1988.
- [10] Anthony J. Weiss and Ehud Weinstein. A lower bound on the mean-square error in random parameter estimation. *IEEE Trans. Inform. Theory*, 31(5):680–682, 1985.