

Chapitre 1

Introduction

1.1 Rôle de l'analyse statistique

L'objectif d'une méthode statistique est de tirer, au vu d'observations d'un phénomène, une *inférence* au sujet de la loi générant ces observations, afin, soit d'analyser le passé, soit de prévoir le futur.

La *statistique*, au même titre que la *physique*, est donc à envisager comme *interprétation* du monde et non comme explication de celui-ci. Lorsque la modélisation déterministe d'un phénomène est trop complexe, ou lorsque les observations sont incertain on fait appel à la modélisation probabiliste qui permet de dépasser le stade descriptif des approches déterministes.

L'inférence s'accompagne donc d'une modélisation probabiliste du phénomène observé. Deux approches sont alors en compétition, modélisation *paramétrique* et modélisation *non paramétrique*. Dans ce travail nous nous limiterons à la modélisation *paramétrique* où, on représente la distribution des observations par une loi $p(x|\boldsymbol{\theta})$ où seul le paramètre $\boldsymbol{\theta}$ de dimension finie est inconnue.

Cette perspective nous semble plus pragmatique, dans la mesure où elle considère qu'un nombre fini d'observations ne peut permettre d'estimer efficacement qu'un nombre fini de paramètres. L'objet de l'inférence statistique est alors de déterminer un ensemble fini de paramètres $\boldsymbol{\theta} = \{\theta_1, \dots, \theta_k\}$ à partir d'un nombre fini d'observations $\boldsymbol{x} = \{x_1, \dots, x_n\}$ d'une quantité X dont on a modélisé sa génération par la loi $p(x|\boldsymbol{\theta})$.

Définition 1 [Modèle paramétrique] Un modèle statistique paramétrique consiste en l'observation d'une variable X , de loi $p(x|\boldsymbol{\theta})$, où seul le paramètre $\boldsymbol{\theta}$, appartenant à un sous-ensemble Θ d'un espace vectoriel de dimension finie, n'est pas connue.

D'une manière plus générale, une fois le modèle construit, on cherche à établir une inférence sur $\boldsymbol{\theta}$, c'est-à-dire à utiliser les observations \boldsymbol{x} afin d'évaluer $\boldsymbol{\theta}$, en vue d'une décision liée à ces paramètres. L'inférence peut ne concerner qu'une partie des composantes de $\boldsymbol{\theta}$. Elle peut se formuler de manière précise comme "Quelle est la valeur de θ_1 ? ", ou vague comme " θ_1 est-il positive ? ", ou encore " θ_2 est-il entre a et $a + \Delta a$? ". Dans le premier cas on parle de *estimation* et dans le deuxième cas on parle de *test d'hypothèse*.

Comparée à la théorie des probabilités, la démarche statistique est fondamentalement une démarche d'*inversion* qui est de remonter des "effets" (observations) aux "causes" (paramètres).

1.2 Approche statistique classique

Dans l'approche classique de statistique cette inversion est flagrante dans la notion de *vraisemblance*. En effet, on écrit

$$l(\boldsymbol{\theta}|\mathbf{x}) = \prod_{i=1}^n p(x_i|\boldsymbol{\theta})$$

en considérant $l(\boldsymbol{\theta}|\mathbf{x})$ comme une fonction de $\boldsymbol{\theta}$, on la normalise (quand cela est possible) pour en faire une fonction densité sur $\boldsymbol{\theta}$ et on l'utilise soit pour estimation soit pour le test d'hypothèse sur $\boldsymbol{\theta}$. Par exemple, en estimation on cherche la valeur $\hat{\boldsymbol{\theta}}_{\text{MV}}$ qui maximise $l(\boldsymbol{\theta}|\mathbf{x})$, c'est l'*estimation au sens du maximum de vraisemblance* :

$$\hat{\boldsymbol{\theta}}_{\text{MV}} = \arg \max_{\boldsymbol{\theta} \in \Theta} \{l(\boldsymbol{\theta}|\mathbf{x})\}$$

En test d'hypothèse on utilise le rapport de deux vraisemblances

$$r = \frac{l(\boldsymbol{\theta} = \mathbf{a}|\mathbf{x})}{l(\boldsymbol{\theta} \neq \mathbf{a}|\mathbf{x})}$$

pour décider entre deux hypothèses $H_0 : \boldsymbol{\theta} = \mathbf{a}$ et $H_1 : \boldsymbol{\theta} \neq \mathbf{a}$.

On utilise $l(\boldsymbol{\theta}|\mathbf{x})$ comme si elle était une fonction densité de probabilité de $\boldsymbol{\theta}$ conditionnellement aux observations \mathbf{x} . Cette inversion est purement formelle, alors que dans l'approche bayésienne, comme nous le verrons un peu plus tard, cette inversion se fait d'une manière plus satisfaisante par la règle de Bayes.

Exemple 1 [Détection d'un signal constant bruité] Soit $\{x_i = \theta + b_i, i = 1, \dots, n\}$ les n -échantillons d'un signal constant noyé dans un bruit blanc gaussien, *i.e.*; $b_i \sim \mathbf{N}(b_i|0, 1)$ ou d'une manière équivalente $x_i \sim \mathbf{N}(x_i|\theta, 1)$. Il s'agit alors d'abord de détecter l'existence d'un signal et s'il existe d'estimer sa valeur θ . Pour cela on note que

$$\begin{aligned} l(\theta|\mathbf{x}) &= \prod_{i=1}^n \mathbf{N}(x_i|\theta, 1) = c_1 \exp \left[-\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2 \right] \\ &= c_2 \exp \left[-\frac{1}{2} [s^2 + n(\bar{x} - \theta)^2] \right] \end{aligned}$$

avec

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \text{et} \quad s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Ainsi, la détection est équivalente au test d'hypothèse $H_0 : \theta \neq 0$ contre $H_1 : \theta = 0$. On définit alors le rapport de vraisemblance r

$$r = \frac{l(\theta \neq 0|\mathbf{x})}{l(\theta = 0|\mathbf{x})} = \frac{c \exp \left[-\frac{1}{2} [s^2 + n(\bar{x} - \theta)^2] \right]}{c \exp \left[-\frac{1}{2} [s^2 + n\bar{x}^2] \right]} = \exp \left[-\frac{1}{2} [n(\bar{x} - \theta)^2 - n\bar{x}^2] \right],$$

et un seuil r_0 et on décide pour H_0 si $r > r_0$, *i.e.*; si

$$(\bar{x} - \theta)^2 - \bar{x}^2 < \frac{2}{n} \log r_0.$$

Notons que le choix du seuil r_0 est arbitraire. Un choix usuel est $r_0 = 1$ ce qui conduit à la condition

$$(\bar{x} - \theta)^2 < \bar{x}^2,$$

pour acceptation de l'hypothèse H_0 . Pour estimer θ par maximum de vraisemblance on obtient

$$\hat{\boldsymbol{\theta}}_{\text{MV}} = \arg \max_{\theta} \{l(\theta|\mathbf{x})\} = \bar{x}.$$

Exemple 2 [Estimation de l'amplitude d'un signal sinusoïdale bruité] Soit $\{x_i = \theta \sin(\omega t_i) + b_i, i = 1, \dots, n\}$ les n -échantillons d'un signal sinusoïdale noyé dans un bruit blanc gaussien, *i.e.*; $b_i \sim \mathbf{N}(b_i|0, 1)$ ou d'une manière équivalente $x_i \sim \mathbf{N}(x_i|\theta \sin(\omega t_i), 1)$. Il s'agit alors d'estimer θ par la méthode du maximum de vraisemblance. Pour cela on note

$$\begin{aligned} l(\theta|\mathbf{x}) &= \prod_{i=1}^n \mathbf{N}(x_i|\theta \sin(\omega t_i), 1) \\ &= c \exp \left[-\frac{1}{2} \sum_{i=1}^n [x_i - \theta \sin(\omega t_i)]^2 \right] \end{aligned}$$

et ainsi, on obtient

$$\hat{\theta}_{\text{MV}} = \arg \max_{\theta} \{l(\theta|\mathbf{x})\} = \arg \min_{\theta} \left\{ \sum_{i=1}^n [x_i - \theta \sin(\omega t_i)]^2 \right\}.$$

On retrouve la notion de l'estimation au sens des moindres carrés.

1.3 Approche bayésienne

L'idée de base dans cette approche est dans le fait que l'on suppose connaître une information *a priori* supplémentaire sur les paramètres θ que l'on peut la traduire sous la forme d'une loi *a priori* $p(\theta)$. L'objectif étant donc d'utiliser cette information supplémentaire. Sachant que l'information contenue dans les observations \mathbf{x} est contenue dans $p(\mathbf{x}|\theta)$ et l'information *a priori* sur θ dans $p(\theta)$, on peut utiliser la règle de Bayes pour combiner ces deux types d'informations en définissant la loi *a posteriori*

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta) p(\theta)}{\int p(\mathbf{x}|\theta) p(\theta) d\theta}$$

qui contiendra donc toute information sur θ .

On remarque que l'inversion de cause à effet est ici beaucoup plus naturelle. Elle se fait d'une manière cohérente, car l'état de connaissance *a priori* sur θ traduite par la loi *a priori* $p(\theta)$ est transformé, après les observations \mathbf{x} , en état de connaissance *a posteriori* par la loi *a posteriori* $p(\theta|\mathbf{x})$.

Remarquons aussi que si les observations x_i sont indépendantes on a

$$p(\mathbf{x}|\theta) = \prod_{i=1}^n p(x_i|\theta) = l(\theta|\mathbf{x}),$$

et la loi *a posteriori* peut s'écrire

$$p(\theta|\mathbf{x}) \propto l(\theta|\mathbf{x}) p(\theta).$$

Définition 2 [Modèle bayésien] On appelle *modèle bayésien* la donnée d'un modèle paramétrique, $p(\mathbf{x}|\theta)$, et d'une loi *a priori* $p(\theta)$ sur les paramètres.

Étant donnée la loi des observations $p(\mathbf{x}|\theta)$ et la loi *a priori* $p(\theta)$ on peut construire

(a) la loi jointe de (θ, \mathbf{x}) ,

$$\phi(\theta, \mathbf{x}) = p(\mathbf{x}|\theta) p(\theta);$$

(b) la loi marginale de \mathbf{x} ,

$$m(\mathbf{x}) = \int \phi(\theta, \mathbf{x}) d\theta = \int p(\mathbf{x}|\theta) p(\theta) d\theta;$$

et

(c) la loi *a posteriori* de θ , obtenue par la formule de Bayes,

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{\int p(\mathbf{x}|\theta)p(\theta) d\theta} = \frac{p(\mathbf{x}|\theta)p(\theta)}{m(\mathbf{x})}.$$

Exemple 3 Si on suppose que $x \sim \mathbf{Bin}(x|\theta, n)$ et $\theta \sim \mathbf{Beta}(\theta|\alpha, \beta)$, alors on a

$$p(x|\theta) = \mathbf{Bin}(x|\theta, n) = C_n^x \theta^x (1-\theta)^{n-x},$$

$$p(\theta) = \mathbf{Beta}(\theta|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

et on peut calculer la loi jointe

$$\phi(\theta, x) = \frac{C_n^x}{B(\alpha, \beta)} \theta^{\alpha+x-1} (1-\theta)^{\beta+n-x-1},$$

la loi marginale

$$m(x) = \frac{C_n^x}{B(\alpha, \beta)} B(\alpha+x, n+\beta-x)$$

$$= C_n^x \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha+x)\Gamma(n+\beta-x)}{\Gamma(\alpha+\beta+n)},$$

et la loi *a posteriori*

$$p(\theta|x) = \frac{1}{B(\alpha+x, \beta+n-x)} \theta^{\alpha+x-1} (1-\theta)^{\beta+n-x-1} = \mathbf{Beta}(\alpha+x, n+\beta-x).$$

Le tableau 1.3 montre la relation entre ces différentes lois pour quelque cas classiques des lois usuelles.

L'introduction d'une loi *a priori* sur les paramètres θ était véritablement révolutionnaire, et aujourd'hui encore, divise les statisticiens. Les orthodoxes considèrent que les paramètres sont *inconnus* mais pas *aléatoire*, et que seule une grandeur aléatoire a droit à avoir une loi de probabilité. C'est cette notion de *quantité ou variable aléatoire* qui est considérée comme une *réalité* qui les conduit à considérer une loi de probabilité comme la représentation de la *réalité* et non pas une représentation de notre état de connaissance sur cette quantité.

Alors que les bayésiens considèrent que les paramètres et les observations, sont deux quantités liées, l'une observable mais incertain (mesurable avec une précision finie), l'autre non observable, mais sur laquelle nous avons une information très incomplète. On utilise alors la théorie des probabilités pour modéliser cette incertitude (manque de précision) des observations par l'attribution des lois $p(\mathbf{x}|\theta)$ et $p(\theta)$, la règle de Bayes pour le calcul de $p(\theta|\mathbf{x})$ et la Statistiques pour tirer une inférence et des conclusions (décision : estimation ou test) sur θ à partir de cette loi *a posteriori*.

Exemple 4 [Estimation de l'amplitude d'un signal sinusoïdale bruité] Prenons le cas de l'exemple 2 où $\{x_i = \theta \sin(\omega t_i) + b_i, i = 1, \dots, n\}$ sont les n -échantillons d'un signal sinusoïdale noyé dans un bruit blanc gaussien, *i.e.*; $x_i \sim \mathbf{N}(x_i|\theta \sin(\omega t_i), 1)$ et supposons que nous savons que l'amplitude du signal en moyenne est μ et sa variation autour de μ est en moyenne égale à λ . Il n'est pas alors irrationnelle d'attribuer une loi *a priori* gaussienne $p(\theta) = \mathbf{N}(\theta|\mu, \lambda)$ à θ pour traduire cette connaissance. (Nous verrons plus tard que ce choix correspond à une loi à entropie maximale sous les contraintes de la connaissance de

Relation entre la loi des observations, la loi *a priori*, la loi marginale et la loi *a posteriori*

loi des observation $f(x \theta)$	loi <i>a priori</i> $\pi(\theta)$	loi marginale $m(x) = \int f(x \theta) \pi(\theta) d\theta$	loi <i>a posteriori</i> $\pi(\theta x) = \frac{f(x \theta) \pi(\theta)}{m(x)}$
Variables discrètes			
Binomiale Bin ($x n, \theta$)	Bêta Bet ($\theta \alpha, \beta$)	Binomiale-Bêta BinBet ($x \alpha, \beta, n$)	Bêta Bet ($\theta \alpha + x, \beta + n - x$)
Binomiale négative NegBin ($x n, \theta$)	Bêta Bet ($\theta \alpha, \beta$)	Binomiale-Bêta négative NegBinBet ($x \alpha, \beta, \theta$)	Bêta Bet ($\theta \alpha + n, \beta + x$)
Poisson Pn ($x \theta$)	Gamma Gam ($\theta \alpha, \beta$)	Poisson-Gamma PnGam ($x \alpha, \beta, 1$)	Gamma Gam ($\theta \alpha + x, \beta + 1$)
Variables continues			
Gamma Gam ($x \nu, \theta$)	Gamma Gam ($\theta \alpha, \beta$)	Gamma-Gamma GamGam ($x \alpha, \beta, \nu$)	Gamma Gam ($\theta \alpha + \nu, \beta + x$)
Exponentielle Ex ($x \theta$)	Gamma Gam ($\theta \alpha, \beta$)	Pareto Par ($x \alpha, \beta$)	Gamma Gam ($\theta \alpha + 1, \beta + x$)
Bêta Bet ($x \alpha, \theta$)	Exponentielle Ex ($\theta \lambda$)	? $?(x \alpha, \lambda)$	Exponentielle Ex ($\theta \lambda - \log(1 - x)$)
Normale N ($x \theta, \sigma^2$)	Normale N ($\theta \mu, \tau^2$)	Normale N ($x \mu + \theta, \tau^2$)?	Normale N ($\mu \frac{\mu\sigma^2 + \tau^2 x}{\sigma^2 + \tau^2}, \frac{\sigma^2 \tau^2}{\sigma^2 + \tau^2}$)
Normale N ($x \mu, \lambda\theta$)	Gamma Gam ($\theta \frac{\alpha}{2}, \frac{\alpha}{2}$)	Student (t) St ($x \mu, \lambda, \alpha$)	Gamma Gam ($\theta \frac{\alpha+1}{2}, \frac{\alpha}{2} + \frac{1}{2}(\mu - x)^2$)
Normale N ($x \mu, 1/\theta$)	Gamma Gam ($\theta \alpha, \beta$)	Student (t)?? St ($x \alpha, \beta$)	Gamma Gam ($\theta \alpha + \frac{1}{2}, \beta + \frac{1}{2}(\mu - x)^2$)

TAB. 1.1 – Dans ce tableau, lorsque la loi des observations dépend de plusieurs paramètres on suppose que seul un paramètre noté θ est inconnu. Par exemple dans la septième ligne $f(x|\theta) = \mathbf{N}(x|\theta, \sigma^2)$ on suppose que σ^2 est connu alors que dans la huitième ligne $f(x|\theta) = \mathbf{N}(x|\mu, \lambda\theta)$ on suppose que μ et λ sont connus.

la moyenne et de la variance.) Avec ce choix et en utilisant les mêmes notations que celles de l'exemple 2 on a

$$\begin{aligned} p(\mathbf{x}|\theta) &= \prod_{i=1}^n \mathbf{N}(x_i|\theta \sin(\omega t_i), 1) \\ &= c_1 \exp \left[-\frac{1}{2} \sum_{i=1}^n [x_i - \theta \sin(\omega t_i)]^2 \right], \\ p(\theta) &= c_2 \exp \left[-\frac{1}{2\lambda} (\theta - \mu)^2 \right] \end{aligned}$$

et

$$\begin{aligned} p(\theta|\mathbf{x}) &= c_3 p(\mathbf{x}|\theta) p(\theta) \\ &= c \exp \left[-\frac{1}{2} \left[\sum_{i=1}^n [x_i - \theta \sin(\omega t_i)]^2 + \frac{1}{\lambda} (\theta - \mu)^2 \right] \right]. \end{aligned}$$

On remarque que la loi *a posteriori* est une loi gaussienne, par conséquent elle est entièrement définie par sa moyenne et sa variance. Dans ce cas, les estimateurs MAP ou la moyenne *a posteriori* sont équivalents. On a

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} \{p(\theta|\mathbf{x})\} = \arg \min_{\theta} \{J(\theta)\},$$

avec

$$J(\theta) = \sum_{i=1}^n [x_i - \theta \sin(\omega t_i)]^2 + \frac{1}{\lambda} (\theta - \mu)^2.$$

On retrouve la notion de l'estimation au sens des moindres carrés.

On a beaucoup critiqué (et cela continue toujours) le choix de la loi *a priori* $\pi(x)$, mais il en va de même pour celui des observations $p(x|\theta)$. Ce qui est étonnant, car ils admettent que X , une quantité physique, suive une loi, par exemple gaussienne de moyen θ , mais ils n'admettent pas que θ qui représente une autre quantité physique qui peut être de la même nature que X (car c'est sa moyenne par exemple), puisse suivre une autre loi gaussienne avec une moyenne μ !

1.4 Vraisemblance et statistique exhaustive

L'analyse statistique (bayésienne ou non bayésienne) est régie par un certain nombre de notions, définitions, règles et principes que j'essaie de les énumérer. Notons que nous connaissons déjà les notions de

- la loi de probabilité $p(x|\theta)$ d'une grandeur X ;
- les observations $\mathbf{x} = \{x_1, \dots, x_n\}$ que l'on appelle parfois un n -échantillon ;
- la fonction de vraisemblance $l(\theta|\mathbf{x}) = \prod_{i=1}^n p(x_i|\theta)$;
- l'espace des observations \mathcal{X} ;
- l'espace des paramètres Θ ;
- la loi jointe des observations $p(\mathbf{x}|\theta)$;
- la loi jointe des observations et des paramètres $p(\mathbf{x}, \theta)$;
- la loi *a priori* des paramètres $p(\theta)$; et
- la loi *a posteriori* des paramètres $p(\theta|\mathbf{x})$;

Définition 3 [Statistique exhaustive] Pour un n -échantillon $\mathbf{x} = \{x_1, \dots, x_n\}$ avec $x_i \sim p(x_i|\boldsymbol{\theta})$, une statistique $\mathbf{t}(\mathbf{x})$, fonctions des observation \mathbf{x} à valeur dans Θ , est dite *exhaustive* si la loi $p(\mathbf{x}|\boldsymbol{\theta})$ ne dépende de \mathbf{x} qu'à travers de $\mathbf{t}(\mathbf{x})$, i.e.;

$$p(\mathbf{x}|\boldsymbol{\theta}) = f(\boldsymbol{\theta}, \mathbf{t}(\mathbf{x})).$$

La connaissance de $\boldsymbol{\theta}$ et $\mathbf{t}(\mathbf{x})$ défini entièrement la loi $p(\mathbf{x}|\boldsymbol{\theta})$. Notons aussi l'équivalence entre cette condition et les conditions suivantes dans l'approche non bayésienne:

$$l(\boldsymbol{\theta}|\mathbf{x}) = l(\boldsymbol{\theta}|\mathbf{t}(\mathbf{x}))$$

et

$$l(\boldsymbol{\theta}|\mathbf{x}, \mathbf{t}(\mathbf{x})) = l(\boldsymbol{\theta}|\mathbf{t}(\mathbf{x})).$$

Exemple 5 Si on suppose que les observations x_i suivent

$$x_i \sim \mathbf{N}(x_i|\mu, \sigma^2),$$

alors on a

$$\begin{aligned} p(\mathbf{x}|\mu, \sigma^2) &= \prod_{i=1}^n \mathbf{N}(x_i|\mu, \sigma^2) \\ &= c\sigma^{-n} \exp\left[-\frac{1}{2\sigma^2}[s^2 + n(\bar{x} - \mu)^2]\right] \\ &= f(\mu, \sigma^2, \bar{x}, s^2), \end{aligned}$$

avec

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \text{et} \quad s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Ainsi, si on note $\boldsymbol{\theta} = [\mu, \sigma^2]$ et $\mathbf{t}(\mathbf{x}) = [\bar{x}, s^2]$ on a

$$p(\mathbf{x}|\boldsymbol{\theta}) = f(\boldsymbol{\theta}, \mathbf{t}(\mathbf{x})),$$

ou en terme de la fonction de vraisemblance on a

$$l(\boldsymbol{\theta}|\mathbf{x}) = l(\boldsymbol{\theta}|\mathbf{t}(\mathbf{x})),$$

ce qui signifie que $\mathbf{t}(\mathbf{x}) = [\bar{x}, s^2]$ est une statistique suffisante pour ces observations, et la connaissance de \mathbf{t} définit entièrement la loi des observations.

La statistique $\mathbf{t}(\mathbf{x})$ contient donc "tout" l'information apportée par \mathbf{x} sur $\boldsymbol{\theta}$, et sous certaines conditions de régularité, la densité de \mathbf{x} se factorise et s'écrit

$$p(\mathbf{x}|\boldsymbol{\theta}) = g(\mathbf{t}(\mathbf{x})|\boldsymbol{\theta}) h(\mathbf{x}|\mathbf{t}(\mathbf{x})),$$

où g est la densité de $\mathbf{t}(\mathbf{x})$.

Définition 4 [Statistique ancillaire] Pour un n -échantillon $\mathbf{x} = \{x_1, \dots, x_n\}$ avec $x_i \sim p(x_i|\boldsymbol{\theta})$, une statistique $\mathbf{a}(\mathbf{x})$ est dite ancillaire par rapport aux paramètres $\boldsymbol{\theta}$ si la distribution de \mathbf{x} conditionnellement à $[\mathbf{a}, \boldsymbol{\theta}]$ ne dépend pas à \mathbf{a} , i.e.;

$$p(\mathbf{x}|\mathbf{a}(\mathbf{x}), \boldsymbol{\theta}) = p(\mathbf{x}|\boldsymbol{\theta}).$$

Notons aussi l'équivalence entre cette condition et les conditions suivantes dans l'approche non bayésienne:

$$l(\boldsymbol{\theta}|\mathbf{x}, \mathbf{a}(\mathbf{x})) = l(\boldsymbol{\theta}|\mathbf{x})$$

et

$$l(\boldsymbol{\theta}|\mathbf{t}(\mathbf{x}), \mathbf{a}(\mathbf{x})) = l(\boldsymbol{\theta}|\mathbf{t}(\mathbf{x})).$$

Exemple 6 Dans l'exemple précédent, si on définit $a_1(\mathbf{x}) = \min\{x_1, \dots, x_n\}$, $a_2(\mathbf{x}) = \max\{x_1, \dots, x_n\}$ et $\mathbf{a} = [a_1, a_2]$ on a

$$p(\mathbf{x}|\mu, \sigma^2, a_1(\mathbf{x}), a_2(\mathbf{x})) = p(\mathbf{x}|\mu, \sigma^2)$$

ou encore

$$p(\mathbf{x}|\boldsymbol{\theta}, \mathbf{a}) = p(\mathbf{x}|\boldsymbol{\theta}),$$

autrement dit, la connaissance de \mathbf{a} n'apporte rien par rapport à la connaissance de $\boldsymbol{\theta}$.

Principe d'Exhaustivité : Deux observations pour lesquelles une statistique exhaustive prend la même valeur doivent conduire à la même inférence sur $\boldsymbol{\theta}$.

Principe de Vraisemblance : Toute l'information sur $\boldsymbol{\theta}$ tirée des observations \mathbf{x} est contenue dans la vraisemblance $l(\boldsymbol{\theta}|\mathbf{x})$, et si \mathbf{x}_1 et \mathbf{x}_2 sont tels qu'il existe une constante c telle que, pour tout $\boldsymbol{\theta}$,

$$l(\boldsymbol{\theta}|\mathbf{x}_1) = cl(\boldsymbol{\theta}|\mathbf{x}_2),$$

ils apportent la même information sur $\boldsymbol{\theta}$ et doivent conduire à la même inférence.

1.5 Lien entre les deux approches

Pensant à l'aspect d'inversion de la statistique, il est tentant de considérer, sous réserve de l'intégrabilité, la fonction de vraisemblance $l(\boldsymbol{\theta}|\mathbf{x})$ comme une loi de probabilité sur $\boldsymbol{\theta}$ dont l'estimateur du maximum de vraisemblance serait la mode. En effet, celle-ci est équivalente à la loi *a posteriori* $p(\boldsymbol{\theta}|\mathbf{x})$ lorsqu'on choisit une loi *a priori* uniforme pour $p(\boldsymbol{\theta})$. Une telle position était, par exemple, celle de Laplace qui considérait que l'absence d'information *a priori* justifiait le choix de la loi uniforme. De même, Fisher, en introduisant l'analyse fiduciaire, voulait mettre en œuvre le principe de vraisemblance sans passer par une approche bayésienne. Cette position, défendable lorsque $\boldsymbol{\theta}$ est un paramètre de position, conduit cependant à des paradoxes et des contradictions qui montrent clairement la nécessité de la théorie bayésienne plus élaborée incluant les notions des lois *a priori non informatives*, des lois *conjuguées* et des lois *a priori de références*.

L'approche à privilégier pour l'inférence bayésienne est celle passant par la loi *a posteriori*. En effet, travaillant conditionnellement aux observations, cette approche suit d'une manière cohérente l'idée d'inversion des causes aux effets, tout en restant fidèle au principe de vraisemblance. En fait, la loi *a posteriori* représente l'actualisation de l'information *a priori*, $p(\boldsymbol{\theta})$, au vu de l'information contenue dans les observations \mathbf{x} , au travers de la vraisemblance $l(\boldsymbol{\theta}|\mathbf{x})$.

Disposant ainsi d'une distribution de probabilité sur $\boldsymbol{\theta}$, le champ de l'inférence est beaucoup plus vaste que dans le cadre classique qui se contentait de $l(\boldsymbol{\theta}|\mathbf{x})$. On peut calculer moyenne, médiane, modes ou des régions de confiance sous la forme de régions de plus forte densité *a posteriori*. On peut aussi, en définissant un espace de décisions et d'une fonction de coût définir tout autres estimateurs classiques, ainsi que des estimateurs de la précision des estimateurs employés. On peut également définir rigoureusement la probabilité *a posteriori* d'une hypothèse $H_0 : \boldsymbol{\theta} \in \Theta_0$, c'est-à-dire $P(\boldsymbol{\theta} \in \Theta_0|\mathbf{x})$, ou comparer deux hypothèses $H_0 : \boldsymbol{\theta} \in \Theta_0$ et $H_1 : \boldsymbol{\theta} \in \Theta_1$ en comparant leurs probabilités respectives $\pi_0(H_0) = P(\boldsymbol{\theta} \in \Theta_0|\mathbf{x})$ et $\pi_1(H_1) = P(\boldsymbol{\theta} \in \Theta_1|\mathbf{x})$.

À ce stade, on peut dire que l'approche statistique classique où toute l'inférence est basée sur la vraisemblance (normalisée et utilisée comme une fonction densité) est un cas particulier de l'approche bayésienne avec une loi *a priori* uniforme. En effet, si $p(\boldsymbol{\theta}) = c$, la loi *a posteriori* $p(\boldsymbol{\theta}|\mathbf{x}) \propto l(\boldsymbol{\theta}|\mathbf{x})$ et, d'après le principe de vraisemblance toute inférence

tirée de ces deux approches seront équivalentes. Cependant, seule l'approche bayésienne permet d'introduire des informations complémentaires sur θ sous forme d'une loi *a priori* non-uniforme.

Une fois la loi *a posteriori* $p(\theta|\mathbf{x})$ calculée, il reste encore à savoir l'utiliser correctement. Deux avis alors se rencontrent.

Le premier consiste à transmettre cette loi à l'utilisateur comme telle qui en fera ce que lui semble bon. Par exemple, quand il s'agit d'un seul paramètre il peut tout simplement tracer la courbe et en déduire

- la mode :

$$M[\theta] = \sup_{\theta \in \Theta} p(\theta|\mathbf{x}),$$

- la moyenne :

$$E[\theta] = \int_{\theta \in \Theta} \theta p(\theta|\mathbf{x}) d\theta,$$

- la variance :

$$\text{Var}[\theta] = \int_{\theta \in \Theta} [\theta - E[\theta]]^2 p(\theta|\mathbf{x}) d\theta,$$

etc.

Dans le cas d'un vecteur des paramètres θ il est plus difficile de résumer la loi *a posteriori*. On peut cependant calculer ces mêmes quantités en utilisant les lois marginales $p(\theta_i|\mathbf{x})$:

- la mode :

$$M[\theta_i] = \sup_{\theta_i \in \Theta_i} p(\theta_i|\mathbf{x}),$$

- la moyenne :

$$E[\theta_i] = \int_{\theta_i \in \Theta_i} \theta_i p(\theta_i|\mathbf{x}) d\theta_i,$$

- la variance :

$$\text{Var}[\theta_i] = \int_{\theta_i \in \Theta_i} [\theta_i - E[\theta_i]]^2 p(\theta_i|\mathbf{x}) d\theta_i,$$

mais aussi

- la covariance :

$$\text{Cov}[\theta_i, \theta_j] = \iint_{\theta \in \Theta} [\theta_i - E[\theta_i]] [\theta_j - E[\theta_j]] p(\theta|\mathbf{x}) d\theta.$$

Notons aussi que la mode $M[\theta_i]$ n'est pas forcément égale à $M[\theta]_i$ où :

$$M[\theta] = \sup_{\theta \in \Theta} p(\theta|\mathbf{x})$$

Le deuxième consiste à définir un estimateur $\delta(\theta)$ suivant un critère. En effet, dans l'approche théorie de la décision, on définit une fonction *coût* $C(\theta, \delta)$ (ou une fonction *utilité*) qui mesure, d'une certaine manière, la conformité de la décision, et choisit comme le meilleur estimateur celui qui minimiserait le coût moyen. Le choix ou la détermination de la fonction coût ne peut se faire qu'une manière subjective.

Dans les chapitres qui vont suivre nous détaillerons un peu plus ces deux approches.

