

Chapter 8

Some complements to Bayesian estimation

8.1 Choice of a prior law in the Bayesian estimation

One of the main difficulties in the application of Bayesian theory in practice is the choice or the attribution of the direct probabilities $f(x|\theta)$ and $\pi(\theta)$. In general, $f(x|\theta)$ is obtained via an appropriate model relating the observable quantity X to the parameters θ and is well accepted. The choice or the attribution of the prior $\pi(\theta)$ has been, and still is, the main subject of discussion and controversy between the Bayesian and orthodox statisticians.

Here, I will try to give a brief summary of different approaches and different tools that can be used to attribute a prior probability distribution. There are mainly four tools:

- use of some invariance principles
- use of maximum entropy (ME) principle
- use of conjugate and reference priors
- use of other information criteria

8.1.1 Invariance principles

Définition 1 [Group invariance] A probability model $f(x|\theta)$ is said to be invariant (or closed) under the action of a group of transformations \mathcal{G} if, for every $g \in \mathcal{G}$, there exists a unique $\theta^* = \bar{g}(\theta) \in \mathcal{T}$ such that $y = g(x)$ is distributed according to $f(y|\theta^*)$.

Exemple 1 Any probability density function in the form $f(x|\theta) = f(x - \theta)$ is invariant under the *translation* group

$$\mathcal{G} : \{g_c(x) : g_c(x) = x + c, \quad c \in \mathbf{R}\} \tag{8.1}$$

This can be verified as follows

$$x \sim f(x - \theta) \longrightarrow y = x + c \sim f(y - \theta^*) \quad \text{with} \quad \theta^* = \theta + c$$

Example 2 Any probability density function in the form $f(x|\theta) = \frac{1}{\theta}f(\frac{x}{\theta})$ is invariant under the *multiplicative* or *scale* transformation group

$$\mathcal{G} : \{g_s(x) : g_s(x) = s x, \quad s > 0\} \quad (8.2)$$

This can be verified as follows

$$x \sim \frac{1}{\theta}f\left(\frac{x}{\theta}\right) \longrightarrow y = s x \sim \frac{1}{\theta^*}f\left(\frac{y}{\theta^*}\right) \quad \text{with} \quad \theta^* = s \theta$$

Example 3 Any probability density function in the form $f(x|\theta_1, \theta_2) = \frac{1}{\theta_2}f(\frac{x-\theta_1}{\theta_2})$ is invariant under the *affine* transformation group

$$\mathcal{G} : \{g_{a,b}(x) : g_{a,b}(x) = a x + b, \quad a > 0, b \in \mathbf{R}\} \quad (8.3)$$

This can be verified as follows

$$x \sim \frac{1}{\theta_2}f\left(\frac{x-\theta_1}{\theta_2}\right) \longrightarrow y = a x + b \sim \frac{1}{\theta_2^*}f\left(\frac{y-\theta_1^*}{\theta_2^*}\right) \quad \text{with} \quad \theta_2^* = a \theta_2, \quad \theta_1^* = a \theta_1 + b.$$

Example 4 Any multi variable probability density function in the form $f(\mathbf{x}|\boldsymbol{\theta}) = f(\mathbf{x}-\boldsymbol{\theta})$ is invariant under the translation group

$$\mathcal{G} : \{g_{\mathbf{c}}(\mathbf{x}) : g_{\mathbf{c}}(\mathbf{x}) = \mathbf{x} - \mathbf{c}, \quad \mathbf{c} \in \mathbf{R}^n\} \quad (8.4)$$

Example 5 Any multi variable probability density function in the form $f(\mathbf{x}) = f(\|\mathbf{x}\|)$ is invariant under the orthogonal transformation group

$$\mathcal{G} : \left\{g_{\mathbf{A}}(\mathbf{x}) : g_{\mathbf{A}}(\mathbf{x}) = \mathbf{A} \mathbf{x}, \quad \mathbf{A}^t \mathbf{A} = \mathbf{A} \mathbf{A}^t = \mathbf{I}\right\} \quad (8.5)$$

Example 6 Any multi variable probability density function in the form $f(\mathbf{x}|\theta) = \frac{1}{\theta}f\left(\frac{\|\mathbf{x}\|}{\theta}\right)$ is invariant under the following transformation group

$$\mathcal{G} : \left\{g_{\mathbf{A},s}(\mathbf{x}) : g_{\mathbf{A},s}(\mathbf{x}) = s \mathbf{A} \mathbf{x}, \quad \mathbf{A}^t \mathbf{A} = \mathbf{A} \mathbf{A}^t = \mathbf{I}, \quad s > 0.\right\} \quad (8.6)$$

This can be verified as follows

$$\mathbf{x} \sim \frac{1}{\theta}f\left(\frac{\|\mathbf{x}\|}{\theta}\right) \longrightarrow \mathbf{y} = s \mathbf{A} \mathbf{x} \sim \frac{1}{\theta^*}f\left(\frac{\|\mathbf{y}\|}{\theta^*}\right) \quad \text{with} \quad \theta^* = s \theta.$$

From these examples we see also that any invariance transformation group \mathcal{G} on $x \in \mathcal{X}$ induces a corresponding transformation group $\bar{\mathcal{G}}$ on $\theta \in \mathcal{T}$. For example for the translation invariance \mathcal{G} on $x \in \mathcal{X}$ induces the following translation group on $\theta \in \mathcal{T}$

$$\bar{\mathcal{G}} : \{\bar{g}_c(\theta) : \bar{g}_c(\theta) = \theta + c, \quad c \in \mathbf{R}\} \quad (8.7)$$

and the scale invariance \mathcal{G} on $x \in \mathcal{X}$ induces the following translation group on $\theta \in \mathcal{T}$

$$\bar{\mathcal{G}} : \{\bar{g}_s(\theta) : \bar{g}_s(\theta) = s \theta, \quad s > 0\} \quad (8.8)$$

We just see that for an invariant family of $f(x|\theta)$ we have a corresponding invariant family of prior laws $\pi(\theta)$. To be complete, we have also to consider the cost function to be able to define the Bayesian estimate.

Définition 2 [Invariant cost functions] Assume a probability model $f(x|\theta)$ is invariant under the action of the group of transformations \mathcal{G} . Then the cost function $c[\theta, \hat{\theta}]$ is said to be invariant under the group of transformations $\tilde{\mathcal{G}}$ if, for every $g \in \mathcal{G}$ and $\hat{\theta} \in \mathcal{T}$, there exists a unique $\hat{\theta}^* = \tilde{g}(\hat{\theta}) \in \mathcal{T}$ with $\tilde{g} \in \tilde{\mathcal{G}}$ such that

$$c[\theta, \hat{\theta}] = c[\tilde{g}(\theta), \hat{\theta}^*] \quad \text{for every } \theta \in \mathcal{T}.$$

Définition 3 [Invariant estimate] For an invariant probability model $f(x|\theta)$ under the group of transformation Gc and an invariant cost function $c[\theta, \hat{\theta}]$ under the corresponding group of transformation $\tilde{\mathcal{G}}$, an estimate $\hat{\theta}$ is said to be invariant or *equivariant* if

$$\hat{\theta}(g(x)) = \tilde{g}(\hat{\theta}(x))$$

Exemple 7 Estimation of θ from the data coming from any model of the kind $f(x|\theta) = f(x - \theta)$ with a quadratic cost function $c[\theta, \hat{\theta}] = (\theta - \hat{\theta})^2$ is equivariant and we have

$$\mathcal{G} = \bar{\mathcal{G}} = \tilde{\mathcal{G}} = \{g_c(x) : g_c(x) = x - c, \quad c \in \mathbb{R}\}$$

Exemple 8 Estimation of θ from the data coming from any model of the kind $f(x|\theta) = \frac{1}{\theta} f(\frac{x}{\theta})$ with the entropy cost function

$$c[\theta, \hat{\theta}] = \frac{\theta}{\hat{\theta}} - \ln\left(\frac{\theta}{\hat{\theta}}\right) - 1$$

is equivariant and we have

$$\begin{aligned} \mathcal{G} &= \{g_s(x) : g_s(x) = s x, \quad s > 0\} \\ \bar{\mathcal{G}} = \tilde{\mathcal{G}} &= \{g_s(\theta) : g_s(\theta) = s \theta, \quad s > 0\} \end{aligned}$$

Proposition 1 [Invariant Bayesian estimate] Suppose that a probability model $f(x|\theta)$ is invariant under the group of transformations \mathcal{G} and that there exists a probability distribution $\pi^*(\theta)$ on \mathcal{T} which is invariant under the group of transformations $\bar{\mathcal{G}}$, *i.e.*,

$$\pi^*(\bar{g}(A)) = \pi^*(A)$$

for any measurable set $A \in \mathcal{T}$. Then the Bayes estimator associated with π^* , noted $\hat{\theta}^*$ minimizes

$$\int R(\theta, \hat{\theta}) \pi^*(\theta) d\theta = \int R(\theta, \bar{g}(\hat{\theta})) \pi^*(\theta) d\theta = \int \mathbb{E} \left[c \left[\theta, \bar{g}(\hat{\theta}(X)) \right] \right] \pi^*(\theta) d\theta \quad \text{over } \hat{\theta}.$$

If this Bayes estimator is unique, it satisfies

$$\hat{\theta}^*(x) = \tilde{g}^{-1}(\hat{\theta}^*(g(x)))$$

Therefore, a Bayes estimator associated with an invariant prior and a strictly convex invariant cost function is almost equivariant.

Actually, invariant probability distributions are rare. The following are some examples:

Exemple 9 If $\pi(\theta)$ is invariant under the translation group \mathcal{G}_c , it satisfies $\pi(\theta) = \pi(\theta + c)$ for every θ and for every c , which implies that $\pi(\theta) = \pi(0)$ uniformly on \mathbb{R} and this leads to the Lebesgue measure as an invariant measure.

Exemple 10 If $\theta > 0$ and $\pi(\theta)$ is invariant under the scale group \mathcal{G}_s , it satisfies $\pi(\theta) = s \pi(s\theta)$ for every $\theta > 0$ and for every $s > 0$, which implies that $\pi(\theta) = 1/\theta$.

Note that in both cases the invariant laws are improper.

8.2 Conjugate priors

The conjugate prior concept is tightly related to the sufficient statistic and exponential families.

Définition 4 [Sufficient statistics] When $X \sim P_\theta(x)$, a function $h(X)$ is said to be a sufficient statistic for $\{P_\theta(x), \theta \in \mathcal{T}\}$ if the distribution of X conditioned on $h(X)$ does not depend on θ for $\theta \in \mathcal{T}$.

Définition 5 [Minimal sufficiency] A function $h(X)$ is said to be minimal sufficient for $\{P_\theta(x), \theta \in \mathcal{T}\}$ if it is a function of every other sufficient statistic for $P_\theta(x)$.

A minimal sufficient statistic contains the whole information brought by the observation $X = x$ about θ .

Proposition 2 [Factorization theorem] Suppose that $\{P_\theta(x), \theta \in \mathcal{T}\}$ has a corresponding family of densities $\{p_\theta(x), \theta \in \mathcal{T}\}$. A statistic T is sufficient for θ if and only if there exist functions g_θ and h such that

$$p_\theta(x) = g_\theta(T(x)) h(x) \quad (8.9)$$

for all $x \in \Gamma$ and $\theta \in \mathcal{T}$.

Exemple 11 If $X \sim \mathcal{N}(\theta, 1)$ then $T(x) = x$ can be chosen as a sufficient statistic.

Exemple 12 If $\{X_1, X_2, \dots, X_n\}$ are i.i.d. and $X_i \sim \mathcal{N}(\theta, 1)$ then

$$\begin{aligned} f(\mathbf{x}|\theta) &= (2\pi)^{-n/2} \exp \left[-\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2 \right] \\ &= \exp \left[-\frac{1}{2} \sum_{i=1}^n x_i^2 \right] (2\pi)^{-n/2} \exp \left[-\frac{n}{2} \theta^2 \right] \exp \left[\theta \sum_{i=1}^n x_i \right] \end{aligned}$$

and we have $T(\mathbf{x}) = \sum_{i=1}^n x_i$.

Note that, in this case, we need to know n and $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. Note also that we can write

$$f(\mathbf{x}|\theta) = a(\mathbf{x}) g(\theta) \exp[\theta T(\mathbf{x})]$$

where

$$g(\theta) = (2\pi)^{-n/2} \exp \left[-\frac{n}{2} \theta^2 \right] \quad \text{and} \quad a(\mathbf{x}) = \exp \left[-\frac{1}{2} \sum_{i=1}^n x_i^2 \right]$$

Exemple 13 If $X \sim \mathcal{N}(0, \theta)$ then $T(x) = x^2$ can be chosen as a sufficient statistic.

Exemple 14 If $X \sim \mathcal{N}(\theta_1, \theta_2)$ then $T_1(x) = x^2$ and $T_2(x) = x$ can be chosen as a set of sufficient statistics.

Exemple 15 If $\{X_1, X_2, \dots, X_n\}$ are i.i.d. and $X_i \sim \mathcal{N}(\theta_1, \theta_2)$ then

$$\begin{aligned} f(\mathbf{x}|\theta_1, \theta_2) &= (2\pi)^{-n/2} \theta_2^{-1/2} \exp \left[-\frac{1}{2\theta_2} \sum_{i=1}^n (x_i - \theta_1)^2 \right] \\ &= (2\pi)^{-n/2} \theta_2^{-1/2} \exp \left[-\frac{n\theta_1^2}{2\theta_2} \right] \exp \left[-\frac{1}{2\theta_2} \sum_{i=1}^n x_i^2 + \frac{\theta_1}{\theta_2} \sum_{i=1}^n x_i \right] \end{aligned}$$

and we have $T_1(\mathbf{x}) = \sum_{i=1}^n x_i$ and $T_2(\mathbf{x}) = \sum_{i=1}^n x_i^2$.

Note also that we can write

$$f(\mathbf{x}|\theta) = a(\mathbf{x}) g(\theta_1, \theta_2) \exp \left[\frac{\theta_1}{\theta_2} T_1(\mathbf{x}) - \frac{1}{2\theta_2} T_2(\mathbf{x}) \right]$$

where

$$g(\theta_1, \theta_2) = (2\pi)^{-n/2} \theta_2^{-1/2} \exp \left[-\frac{n\theta_1^2}{2\theta_2} \right] \quad \text{and} \quad a(\mathbf{x}) = 1.$$

In this case, $\frac{\theta_1}{\theta_2}$ and $\frac{-1}{2\theta_2}$ are called canonical parametrization. It is also usual to use n , $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and $\bar{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2$ as the sufficient statistics.

Exemple 16 If $X \sim \mathbf{Gam}(\alpha, \theta)$ then $T(x) = x$ can be chosen as a sufficient statistic.

Exemple 17 If $X \sim \mathbf{Gam}(\theta, \beta)$ then $T(x) = \ln x$ can be chosen as a sufficient statistic.

Exemple 18 If $X \sim \mathbf{Gam}(\theta_1, \theta_2)$ then $T_1(x) = \ln x$ and $T_2(x) = x$ can be chosen as a set of sufficient statistics.

Exemple 19 If $\{X_1, X_2, \dots, X_n\}$ are i.i.d. and $X_i \sim \mathbf{Gam}(\theta_1, \theta_2)$ then it is easy to show that $T_1(\mathbf{x}) = \sum_{i=1}^n \ln x_i$ and $T_2(\mathbf{x}) = \sum_{i=1}^n x_i$.

Définition 6 [Exponential family] A class of distributions $\{P_\theta(\mathbf{x}), \theta \in \mathcal{T}\}$ is said to be an exponential family if there exist: $a(\mathbf{x})$ a function of Γ on \mathbb{R} , $g(\theta)$ a function of \mathcal{T} on \mathbb{R}^+ , $\phi_k(\theta)$ functions of \mathcal{T} on \mathbb{R} , and $h_k(\mathbf{x})$ functions of Γ on \mathbb{R} such that

$$\begin{aligned} p_\theta(\mathbf{x}) = p(\mathbf{x}|\theta) &= a(\mathbf{x}) g(\theta) \exp \left[\sum_{k=1}^K \phi_k(\theta) h_k(\mathbf{x}) \right] \\ &= a(\mathbf{x}) g(\theta) \exp \left[\boldsymbol{\phi}^t(\theta) \mathbf{h}(\mathbf{x}) \right] \end{aligned}$$

for all $\theta \in \mathcal{T}$ and $x \in \Gamma$. This family is entirely determined by $a(\mathbf{x})$, $g(\theta)$, and $\{\phi_k(\theta), h_k(\mathbf{x}), k = 1, \dots, K\}$ and is noted $\mathbf{Exfn}(\mathbf{x}|a, g, \boldsymbol{\phi}, \mathbf{h})$

Particular cases:

- When $a(\mathbf{x}) = 1$ and $g(\theta) = \exp[-b(\theta)]$ we have

$$p(\mathbf{x}|\theta) = \exp \left[\boldsymbol{\phi}^t(\theta) \mathbf{h}(\mathbf{x}) - b(\theta) \right]$$

and is noted $\mathbf{CExf}(\mathbf{x}|b, \boldsymbol{\phi}, \mathbf{h})$.

- Natural exponential family:

When $a(\mathbf{x}) = 1$, $g(\boldsymbol{\theta}) = \exp[-b(\boldsymbol{\theta})]$, $\mathbf{h}(\mathbf{x}) = \mathbf{x}$ and $\boldsymbol{\phi}(\boldsymbol{\theta}) = \boldsymbol{\theta}$ we have

$$p(\mathbf{x}|\boldsymbol{\theta}) = \exp[\boldsymbol{\theta}^t \mathbf{x} - b(\boldsymbol{\theta})] \mathbf{E}\mathbf{x}\mathbf{f}(\mathbf{x}|b).$$

and is noted $\mathbf{N}\mathbf{E}\mathbf{x}\mathbf{f}(\mathbf{x}|b)$.

- Scalar random variable with a vector parameter:

$$\begin{aligned} p(x|\boldsymbol{\theta}) &= \mathbf{E}\mathbf{x}\mathbf{f}(x|a, g, \boldsymbol{\phi}, \mathbf{h}) \\ &= a(x)g(\boldsymbol{\theta}) \exp\left[\sum_{k=1}^K \phi_k(\boldsymbol{\theta})h_k(x)\right] \\ &= a(x)g(\boldsymbol{\theta}) \exp[\boldsymbol{\phi}^t(\boldsymbol{\theta})\mathbf{h}(x)] \end{aligned}$$

and is noted $\mathbf{E}\mathbf{x}\mathbf{f}\mathbf{k}(x|a, g, \boldsymbol{\phi}, \mathbf{h})$.

- Scalar random variable with a scalar parameter:

$$p(x|\theta) = \mathbf{E}\mathbf{x}\mathbf{f}(x|a, g, \phi, h) = a(x)g(\theta) \exp[\phi(\theta)h(x)]$$

and is noted $\mathbf{E}\mathbf{x}\mathbf{f}(\mathbf{x}|a, g, \phi, h)$.

- Simple scalar exponential family:

$$p(x|\theta) = \theta \exp[-\theta x] = \exp[-\theta x + \ln \theta], \quad x \geq 0, \quad \theta \geq 0.$$

Définition 7 [Conjugate distributions] A family \mathcal{F} of probability distributions $\pi(\boldsymbol{\theta})$ on \mathcal{T} is said to be conjugate (or closed under sampling) if, for every $\pi(\boldsymbol{\theta}) \in \mathcal{F}$, the posterior distribution $\pi(\boldsymbol{\theta}|\mathbf{x})$ also belongs to \mathcal{F} .

The main argument for the development of the conjugate priors is the following: When the observation of a variable X with a probability law $f(x|\theta)$ modifies the prior $\pi(\theta)$ to a posterior $\pi(\theta|x)$, the information conveyed by x about θ is obviously limited, therefore it should not lead to a modification of the whole structure of $\pi(\theta)$, but only of its parameters.

Définition 8 [Conjugate priors] Assume that $f(\mathbf{x}|\boldsymbol{\theta}) = l(\boldsymbol{\theta}|\mathbf{x}) = l(\boldsymbol{\theta}|\mathbf{t}(\mathbf{x}))$ where $\mathbf{t} = \{n, \mathbf{s}\} = \{n, s_1, \dots, s_k\}$ is a vector of dimension $k+1$ and is sufficient statistic for $f(\mathbf{x}|\boldsymbol{\theta})$. Then, if there exists a vector $\{\tau_0, \boldsymbol{\tau}\} = \{\tau_0, \tau_1, \dots, \tau_k\}$ such that

$$\pi(\boldsymbol{\theta}|\boldsymbol{\tau}) = \frac{f(\mathbf{s} = (\tau_1, \dots, \tau_k)|\boldsymbol{\theta}, n = \tau_0)}{\int f(\mathbf{s} = (\tau_1, \dots, \tau_k)|\boldsymbol{\theta}', n = \tau_0) d\boldsymbol{\theta}'}$$

exists and defines a family \mathcal{F} of distributions for $\boldsymbol{\theta} \in \mathcal{T}$, then the posterior $\pi(\boldsymbol{\theta}|\mathbf{x}, \boldsymbol{\tau})$ will remain in the same family \mathcal{F} . The prior distribution $\pi(\boldsymbol{\theta}|\boldsymbol{\tau})$ is then a conjugate prior for the sampling distribution $f(\mathbf{x}|\boldsymbol{\theta})$.

Proposition 3 [Sufficient statistics for the exponential family] For a set of n i.i.d. samples $\{x_1, \dots, x_n\}$ of a random variable $X \sim \mathbf{Exfn}(x|a, g, \boldsymbol{\theta}, \mathbf{h})$ we have

$$\begin{aligned} f(\mathbf{x}|\boldsymbol{\theta}) &= \prod_{j=1}^n f(x_j|\boldsymbol{\theta}) = [g(\boldsymbol{\theta})]^n \left(\prod_{j=1}^n a(x_j) \right) \exp \left[\sum_{k=1}^K \phi_k(\boldsymbol{\theta}) \sum_{j=1}^n h_k(x_j) \right] \\ &= g^n(\boldsymbol{\theta}) a(\mathbf{x}) \exp \left[\boldsymbol{\phi}^t(\boldsymbol{\theta}) \sum_{j=1}^n \mathbf{h}(x_j) \right], \end{aligned}$$

where $a(\mathbf{x}) = \prod_{j=1}^n a(x_j)$. Then, using the factorization theorem it is easy to see that

$$\mathbf{t} = \left\{ n, \sum_{j=1}^n h_1(x_j), \dots, \sum_{j=1}^n h_K(x_j) \right\}$$

is a sufficient statistic for $\boldsymbol{\theta}$.

Proposition 4 [Conjugate priors of the Exponential family] A conjugate prior family for the exponential family

$$f(\mathbf{x}|\boldsymbol{\theta}) = a(\mathbf{x}) g(\boldsymbol{\theta}) \exp \left[\sum_{k=1}^K \phi_k(\boldsymbol{\theta}) h_k(\mathbf{x}) \right]$$

is given by

$$\pi(\boldsymbol{\theta}|\tau_0, \boldsymbol{\tau}) = z(\boldsymbol{\tau}) [g(\boldsymbol{\theta})]^{\tau_0} \exp \left[\sum_{k=1}^K \tau_k \phi_k(\boldsymbol{\theta}) \right]$$

The associated posterior law is

$$\pi(\boldsymbol{\theta}|\mathbf{x}, \tau_0, \boldsymbol{\tau}) \propto [g(\boldsymbol{\theta})]^{n+\tau_0} a(\mathbf{x}) z(\boldsymbol{\tau}) \exp \left[\sum_{k=1}^K \left(\tau_k + \sum_{j=1}^n h_k(x_j) \right) \phi_k(\boldsymbol{\theta}) \right].$$

We can rewrite this in a more compact way:

If

$$f(\mathbf{x}|\boldsymbol{\theta}) = \mathbf{Exfn}(\mathbf{x}|a(\mathbf{x}), g(\boldsymbol{\theta}), \boldsymbol{\phi}, \mathbf{h}),$$

then a conjugate prior family is

$$\pi(\boldsymbol{\theta}|\boldsymbol{\tau}) = \mathbf{Exfn}(\boldsymbol{\theta}|g^{\tau_0}, z(\boldsymbol{\tau}), \boldsymbol{\tau}, \boldsymbol{\phi}),$$

and the associated posterior law is

$$\pi(\boldsymbol{\theta}|\mathbf{x}, \boldsymbol{\tau}) = \mathbf{Exfn}(\boldsymbol{\theta}|g^{n+\tau_0}, a(\mathbf{x}) z(\boldsymbol{\tau}), \boldsymbol{\tau}', \boldsymbol{\phi})$$

where

$$\tau'_k = \tau_k + \sum_{j=1}^n h_k(x_j)$$

or

$$\boldsymbol{\tau}' = \boldsymbol{\tau} + \bar{\mathbf{h}}, \quad \text{with} \quad \bar{h}_k = \sum_{j=1}^n h_k(x_j).$$

Définition 9 [Conjugate priors of natural exponential family] If

$$f(\mathbf{x}|\boldsymbol{\theta}) = a(\mathbf{x}) \exp \left[\boldsymbol{\theta}^t \mathbf{x} - b(\boldsymbol{\theta}) \right]$$

Then a conjugate prior family is

$$\pi(\boldsymbol{\theta}|\boldsymbol{\tau}_0) = g(\boldsymbol{\theta}) \exp \left[\boldsymbol{\tau}_0^t \boldsymbol{\theta} - d(\boldsymbol{\tau}_0) \right]$$

and the corresponding posterior is

$$\pi(\boldsymbol{\theta}|\mathbf{x}, \boldsymbol{\tau}_0) = g(\boldsymbol{\theta}) \exp \left[\boldsymbol{\tau}_n^t \boldsymbol{\theta} - d(\boldsymbol{\tau}_n) \right] \quad \text{with} \quad \boldsymbol{\tau}_n = \boldsymbol{\tau}_0 + \bar{\mathbf{x}}$$

where

$$\bar{\mathbf{x}}_n = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$$

A slightly more general notation which gives some more explicit properties of the conjugate priors of the natural exponential family is the following:

If

$$f(\mathbf{x}|\boldsymbol{\theta}) = a(\mathbf{x}) \exp \left[\boldsymbol{\theta}^t \mathbf{x} - b(\boldsymbol{\theta}) \right]$$

Then a conjugate prior family is

$$\pi(\boldsymbol{\theta}|\alpha_0, \boldsymbol{\tau}_0) = g(\alpha_0, \boldsymbol{\tau}_0) \exp \left[\alpha_0 \boldsymbol{\tau}_0^t \boldsymbol{\theta} - \alpha_0 b(\boldsymbol{\tau}_0) \right]$$

The posterior is

$$\pi(\boldsymbol{\theta}|\alpha_0, \boldsymbol{\tau}_0, \mathbf{x}) = g(\alpha, \boldsymbol{\tau}) \exp \left[\alpha \boldsymbol{\tau}^t \boldsymbol{\theta} - \alpha b(\boldsymbol{\tau}) \right]$$

with

$$\alpha = \alpha_0 + n \quad \text{and} \quad \boldsymbol{\tau} = \frac{\alpha_0 \boldsymbol{\tau}_0 + n \bar{\mathbf{x}}}{(\alpha_0 + n)}$$

and we have the following properties:

$$\mathbb{E}[\mathbf{X}|\boldsymbol{\theta}] = \mathbb{E}[\bar{\mathbf{X}}|\boldsymbol{\theta}] = \nabla b(\boldsymbol{\theta})$$

$$\mathbb{E}[\nabla b(\boldsymbol{\Theta})|\alpha_0, \boldsymbol{\tau}_0] = \boldsymbol{\tau}_0$$

$$\mathbb{E}[\nabla b(\boldsymbol{\theta})|\alpha_0, \boldsymbol{\tau}_0, \mathbf{x}] = \frac{n \bar{\mathbf{x}} + \alpha_0 \boldsymbol{\tau}_0}{\alpha_0 + n} = \pi \bar{\mathbf{x}}_n + (1 - \pi) \boldsymbol{\tau}_0, \quad \text{with} \quad \pi = \frac{n}{\alpha_0 + n}$$

Conjugate priors

Observation law $p(x \theta)$	Prior law $p(\theta \tau)$	Posterior law $p(\theta x, \tau) \propto p(\theta \tau)p(x \theta)$
Discrete variables		
Binomial $\mathbf{Bin}(x n, \theta)$	Beta $\mathbf{Bet}(\theta \alpha, \beta)$	Beta $\mathbf{Bet}(\theta \alpha + x, \beta + n - x)$
Negative Binomial $\mathbf{NegBin}(x n, \theta)$	Beta $\mathbf{Bet}(\theta \alpha, \beta)$	Beta $\mathbf{Bet}(\theta \alpha + n, \beta + x)$
Multinomial $\mathbf{M}_k(x \theta_1, \dots, \theta_k)$	Dirichlet $\mathbf{Di}_k(\theta \alpha_1, \dots, \alpha_k)$	Dirichlet $\mathbf{Di}_k(\theta \alpha_1 + x_1, \dots, \alpha_k + x_k)$
Poisson $\mathbf{Pn}(x \theta)$	Gamma $\mathbf{Gam}(\theta \alpha, \beta)$	Gamma $\mathbf{Gam}(\theta \alpha + x, \beta + 1)$
Gamma $\mathbf{Gam}(x \nu, \theta)$	Gamma $\mathbf{Gam}(\theta \alpha, \beta)$	Gamma $\mathbf{Gam}(\theta \alpha + \nu, \beta + x)$
Beta $\mathbf{Bet}(x \alpha, \theta)$	Exponential $\mathbf{Ex}(\theta \lambda)$	Exponential $\mathbf{Ex}(\theta \lambda - \log(1 - x))$
Normal $\mathbf{N}(x \theta, \sigma^2)$	Normal $\mathbf{N}(\theta \mu, \tau^2)$	Normal $\mathbf{N}\left(\mu \frac{\mu\sigma^2 + \tau^2 x}{\sigma^2 + \tau^2}, \frac{\sigma^2 \tau^2}{\sigma^2 + \tau^2}\right)$
Continuous variables		
Normal $\mathbf{N}(x \mu, 1/\theta)$	Gamma $\mathbf{Gam}(\theta \alpha, \beta)$	Gamma $\mathbf{Gam}\left(\theta \alpha + \frac{1}{2}, \beta + \frac{1}{2}(\mu - x)^2\right)$
Normal $\mathbf{N}(x \theta, \theta^2)$	Generalized inverse Normal $\mathbf{INg}(\theta \alpha, \mu, \sigma) \propto$ $ \theta ^{-\alpha} \exp\left[-\frac{1}{2\sigma^2} \left(\frac{1}{\theta} - \mu\right)^2\right]$	Generalized inverse Normal $\mathbf{INg}(\theta \alpha_n, \mu_n, \sigma_n)$

Table 8.1: Relation between the sampling distributions, their associated conjugate priors and their corresponding posteriors

8.3 Non informative priors based on Fisher information

Another notion of information related to the maximum likelihood estimation is the Fisher information. In this section, first we give some definitions and results related to this notion and we see how this is used to define non informative priors.

Proposition 5 [Information Inequality] Let $\hat{\theta}$ be an estimate of the parameter θ in a family $\{P_\theta; \theta \in \mathcal{T}\}$ and assume that the following conditions hold:

1. The family $\{P_\theta; \theta \in \mathcal{T}\}$ has a corresponding family of densities $\{p_\theta(x); \theta \in \mathcal{T}\}$, all with the same support.
2. $p_\theta(x)$ is differentiable for all $\theta \in \mathcal{T}$ and all x in its support.
3. The integral

$$g(\theta) = \int_{\Gamma} h(x) p_\theta(x) \mu(\mathrm{d}x)$$

exists and is differentiable for $\theta \in \mathcal{T}$, for $h(x) = \hat{\theta}(x)$ and for $h(x) = 1$ and

$$\frac{\partial g(\theta)}{\partial \theta} = \int_{\Gamma} h(x) \frac{\partial p_\theta(x)}{\partial \theta} \mu(\mathrm{d}x)$$

Then

$$\mathrm{Var}_\theta[\hat{\theta}(X)] \geq \frac{\left[\frac{\partial}{\partial \theta} \mathrm{E}_\theta \left\{ \hat{\theta}(X) \right\} \right]^2}{I_\theta} \quad (8.10)$$

where

$$I_\theta \stackrel{\text{def}}{=} \mathrm{E}_\theta \left\{ \left[\frac{\partial}{\partial \theta} \ln p_\theta(X) \right]^2 \right\} \quad (8.11)$$

Furthermore, if $\frac{\partial^2}{\partial \theta^2} p_\theta(x)$ exists for all $\theta \in \mathcal{T}$ and all x in the support of $p_\theta(x)$, and if

$$\int \frac{\partial^2}{\partial \theta^2} p_\theta(x) \mu(\mathrm{d}x) = \frac{\partial^2}{\partial \theta^2} \int p_\theta(x) \mu(\mathrm{d}x)$$

then I_θ can be computed via

$$I_\theta = -\mathrm{E}_\theta \left\{ \frac{\partial^2}{\partial \theta^2} \ln p_\theta(X) \right\} \quad (8.12)$$

The quantity defined in (8.11) is known as *Fisher's information* for estimating θ from X , and (8.10) is called the *information inequality*.

For the particular case in which $\hat{\theta}$ is unbiased $\mathrm{E}_\theta \left\{ \hat{\theta}(X) \right\} = \theta$, the information inequality becomes

$$\mathrm{Var}_\theta[\hat{\theta}(X)] \geq \frac{1}{I_\theta} \quad (8.13)$$

Expression $\frac{1}{I_\theta}$ is known as the *Cramer-Rao lower bound* (CRLB).

Exemple 20 [The information Inequality for exponential families] Assume that \mathcal{T} is open and p_θ is given by

$$p_\theta(x) = a(x) g(\theta) \exp[g(\theta) h(x)]$$

Then it can be shown that

$$I_\theta \stackrel{\text{def}}{=} \mathbb{E}_\theta \left\{ \left[\frac{\partial}{\partial \theta} \ln p_\theta(X) \right]^2 \right\} = |g'(\theta)|^2 \text{Var}_\theta(h(X)) \quad (8.14)$$

and

$$\frac{\partial}{\partial \theta} \mathbb{E}_\theta \{h(X)\} = g'(\theta) \text{Var}_\theta(h(X)) \quad (8.15)$$

and thus, if we choose $\hat{\theta}(x) = h(x)$ we obtain the lower bound in the information inequality (8.10)

$$\text{Var}_\theta[\hat{\theta}(X)] = \frac{\left[\frac{\partial}{\partial \theta} \mathbb{E}_\theta \{ \hat{\theta}(X) \} \right]^2}{I_\theta} \quad (8.16)$$

Définition 10 [Non informative priors] Assume $X \sim f(x|\theta) = p_\theta(x)$ and assume that

$$I_\theta \stackrel{\text{def}}{=} \mathbb{E}_\theta \left\{ \left[\frac{\partial}{\partial \theta} \ln p_\theta(X) \right]^2 \right\} = -\mathbb{E}_\theta \left\{ \frac{\partial^2}{\partial \theta^2} \ln p_\theta(X) \right\} \quad (8.17)$$

Then, a non informative prior $\pi(\theta)$ is defined as

$$\pi(\theta) \propto I_\theta^{1/2} \quad (8.18)$$

Définition 11 [Non informative priors, case of vector parameters]

Assume $X \sim f(x|\boldsymbol{\theta}) = p_{\boldsymbol{\theta}}(x)$ and assume that

$$I_{ij}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\mathbb{E}_\theta \left\{ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln p_{\boldsymbol{\theta}}(X) \right\} \quad (8.19)$$

Then, a non informative prior $\pi(\boldsymbol{\theta})$ is defined as

$$\pi(\boldsymbol{\theta}) \propto |\mathbf{I}(\boldsymbol{\theta})|^{1/2} \quad (8.20)$$

where $\mathbf{I}(\boldsymbol{\theta})$ is the Fisher information matrix with the elements $I_{ij}(\boldsymbol{\theta})$.

Exemple 21 If

$$f(\mathbf{x}|\boldsymbol{\theta}) = a(\mathbf{x}) \exp[\boldsymbol{\theta}^t \mathbf{x} - b(\boldsymbol{\theta})]$$

then

$$\mathbf{I}(\boldsymbol{\theta}) = \nabla \nabla^t b(\boldsymbol{\theta})$$

and

$$\pi(\boldsymbol{\theta}) \propto |\mathbf{I}(\boldsymbol{\theta})|^{1/2} = \left| \prod_{i=1}^n \frac{\partial^2 \theta_i}{\partial b(\boldsymbol{\theta})^2} \right|^{1/2}$$

Exemple 22 If

$$f(\mathbf{x}|\boldsymbol{\theta}) = \mathcal{N}(\mu, \sigma^2), \quad \boldsymbol{\theta} = (\mu, \sigma^2)$$

then

$$I(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}} \left\{ \begin{pmatrix} \frac{1}{\sigma^2} & \frac{2(X-\mu)}{\sigma^3} \\ \frac{2(X-\mu)}{\sigma^3} & \frac{3(X-\mu)^2}{\sigma^4} - \frac{1}{\sigma^2} \end{pmatrix} \right\} = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{2}{\sigma^2} \end{pmatrix}$$

and

$$\pi(\boldsymbol{\theta}) = \pi(\mu, \sigma^2) \propto \frac{1}{\sigma^4}$$