

1 Bayesian Inference

1.1 Introduction and Terminology

In Bayesian analysis, θ is treated as a random variable with a **prior** density encapsulating the beliefs about θ before the data are collected. Denoting by $\pi_{\theta}(\theta)$ the prior density, by Bayes Theorem,

$$\pi_{\theta|\underline{x}}(\theta|\underline{x}) = \frac{f_{\underline{x}|\theta}(\underline{x}|\theta)\pi_{\theta}(\theta)}{f_{\underline{x}}(\underline{x})} = \frac{f_{\underline{x}|\theta}(\underline{x}|\theta)\pi_{\theta}(\theta)}{\int f_{\underline{x}|\theta}(\underline{x}|\theta)\pi_{\theta}(\theta)d\theta} \quad (1)$$

The denominator in (1), $f_{\underline{x}}(\underline{x})$ is sometimes termed the **marginal likelihood** of the data \underline{x} , and does not depend on θ . The term $\pi_{\theta|\underline{x}}(\theta|\underline{x})$ in (1) is the **posterior distribution** of θ given \underline{x} , which encapsulates the beliefs about θ in light of the data. Note that for the posterior calculation in (1)

$$\pi_{\theta|\underline{x}}(\theta|\underline{x}) \propto f_{\underline{x}|\theta}(\underline{x}|\theta)\pi_{\theta}(\theta)$$

ignoring terms that do not involve θ ; $f_{\underline{x}}(\underline{x})$ for fixed \underline{x} acts as a normalizing constant.

1.2 A Variance Lemma

The following inequality proves that the posterior variance is larger than the prior variance, and gives a general version of the result for vector parameters. For two $k \times k$ matrices A and B , write $A \geq B$ if $A - B$ is **non-negative definite**, that is, if

$$\underline{x}^T(A - B)\underline{x} \geq 0 \quad \text{for all } \underline{x} \in \mathbb{R}^k$$

Lemma 1.1 For any two vector random variables \underline{X} ($k_1 \times 1$) and \underline{Y} ($k_2 \times 1$) having some joint probability structure,

$$\text{Var}_{f_{\underline{X}}}[\underline{X}] \geq E_{f_{\underline{Y}}}[\text{Var}_{f_{\underline{X}|\underline{Y}}}[\underline{X}|\underline{Y} = \underline{y}]]$$

where both sides are $k_1 \times k_1$ matrices.

Proof. Denote the marginal and conditional expectations of \underline{X} by

$$\underline{\mu}^X = E_{f_{\underline{X}}}[\underline{X}] \quad \underline{\mu}^{X|\underline{y}} = E_{f_{\underline{X}|\underline{Y}}}[\underline{X}|\underline{Y} = \underline{y}]$$

and then define

$$\underline{\mu}^{X|\underline{Y}} = E_{f_{\underline{X}|\underline{Y}}}[\underline{X}|\underline{Y}]$$

as the **random variable** formed by conditioning on $\underline{Y} = \underline{y}$ as \underline{y} varies according to $f_{\underline{Y}}$. We have that

$$\begin{aligned} \text{Var}_{f_{\underline{X}}}[\underline{X}] &= E_{f_{\underline{X}}}[(\underline{X} - \underline{\mu}^X)(\underline{X} - \underline{\mu}^X)^T] = E_{f_{\underline{Y}}} [E_{f_{\underline{X}|\underline{Y}}}[(\underline{X} - \underline{\mu}^X)(\underline{X} - \underline{\mu}^X)^T]|\underline{Y} = \underline{y}] \\ &= E_{f_{\underline{Y}}} [E_{f_{\underline{X}|\underline{Y}}}[(\underline{X} - \underline{\mu}^{X|\underline{y}} + \underline{\mu}^{X|\underline{y}} - \underline{\mu}^X)(\underline{X} - \underline{\mu}^{X|\underline{y}} + \underline{\mu}^{X|\underline{y}} - \underline{\mu}^X)^T]|\underline{Y} = \underline{y}] \\ &= E_{f_{\underline{Y}}} [E_{f_{\underline{X}|\underline{Y}}}[(\underline{X} - \underline{\mu}^{X|\underline{y}})(\underline{X} - \underline{\mu}^{X|\underline{y}})^T]|\underline{Y} = \underline{y}] + \\ &\quad 2E_{f_{\underline{Y}}} [E_{f_{\underline{X}|\underline{Y}}}[(\underline{X} - \underline{\mu}^{X|\underline{y}})(\underline{\mu}^{X|\underline{y}} - \underline{\mu}^X)^T]|\underline{Y} = \underline{y}] \\ &\quad + E_{f_{\underline{Y}}} [E_{f_{\underline{X}|\underline{Y}}}[(\underline{\mu}^{X|\underline{y}} - \underline{\mu}^X)(\underline{\mu}^{X|\underline{y}} - \underline{\mu}^X)^T]|\underline{Y} = \underline{y}] \end{aligned}$$

The second expectation is zero, as in the interior expectation

$$E_{f_{\underline{X}|\underline{Y}}}[(\underline{X} - \underline{\mu}^{X|y})(\underline{\mu}^{X|y} - \underline{\mu}^X)^\top | \underline{Y} = \underline{y}] = E_{f_{\underline{X}|\underline{Y}}}[(\underline{X} - \underline{\mu}^{X|y}) | \underline{Y} = \underline{y}](\underline{\mu}^{X|y} - \underline{\mu}^X)^\top = 0.$$

Therefore

$$\begin{aligned} \text{Var}_{f_{\underline{X}}}[\underline{X}] &= E_{f_{\underline{Y}}}[\text{Var}_{f_{\underline{X}|\underline{Y}}}[\underline{X} | \underline{Y} = \underline{y}]] + E_{f_{\underline{Y}}}[\text{E}_{f_{\underline{X}|\underline{Y}}}[(\underline{\mu}^{X|y} - \underline{\mu}^X)(\underline{\mu}^{X|y} - \underline{\mu}^X)^\top] | \underline{Y} = \underline{y}]] \\ &= E_{f_{\underline{Y}}}[\text{Var}_{f_{\underline{X}|\underline{Y}}}[\underline{X} | \underline{Y} = \underline{y}]] + \text{Var}_{f_{\underline{Y}}}[\text{E}_{f_{\underline{X}|\underline{Y}}}[\underline{\mu}^{X|y} | \underline{Y} = \underline{y}]] \\ &= E_{f_{\underline{Y}}}[\text{Var}_{f_{\underline{X}|\underline{Y}}}[\underline{X} | \underline{Y} = \underline{y}]] + \text{Var}_{f_{\underline{Y}}}[\underline{\mu}^{X|Y}] \end{aligned}$$

Thus, using this iterated expectation argument, and denoting by

$$\text{Var}_{f_{\underline{X}|\underline{Y}}}[\underline{X} | \underline{Y}]$$

the random variable formed by constructing $\text{Var}_{f_{\underline{X}|\underline{Y}}}[\underline{X} | \underline{Y} = \underline{y}]$ as \underline{y} varies according to $f_{\underline{Y}}$, we have

$$\begin{aligned} \text{Var}_{f_{\underline{X}}}[\underline{X}] &= E_{f_{\underline{Y}}}[\text{Var}_{f_{\underline{X}|\underline{Y}}}[\underline{X} | \underline{Y}]] + \text{Var}_{f_{\underline{Y}}}[\text{E}_{f_{\underline{X}|\underline{Y}}}[\underline{X} | \underline{Y}]] \\ &\geq E_{f_{\underline{Y}}}[\text{Var}_{f_{\underline{X}|\underline{Y}}}[\underline{X} | \underline{Y}]] \end{aligned}$$

as the second term is non-negative definite. ■

Corollary : The variance of the posterior distribution, $\text{Var}_{\pi_{\underline{\theta}|\underline{X}}}[\underline{\theta} | \underline{X} = \underline{x}]$ satisfies

$$\text{Var}_{\pi_{\underline{\theta}}}[\underline{\theta}] \geq E_{f_{\underline{X}}}[\text{Var}_{\pi_{\underline{\theta}|\underline{X}}}[\underline{\theta} | \underline{X}]]$$

that is, the prior variance is at least as big as expected posterior variance. ■

1.3 Bayesian Updating

The Bayesian calculation in (1) acts sequentially, that is, for data \underline{x}_1

$$\pi_{\underline{\theta}|\underline{X}}(\underline{\theta} | \underline{x}_1) = \frac{f_{\underline{X}|\underline{\theta}}(\underline{x}_1 | \underline{\theta})\pi_{\underline{\theta}}(\underline{\theta})}{f_{\underline{X}}(\underline{x}_1)} = \frac{f_{\underline{X}|\underline{\theta}}(\underline{x}_1 | \underline{\theta})\pi_{\underline{\theta}}(\underline{\theta})}{\int f_{\underline{X}|\underline{\theta}}(\underline{x}_1 | \underline{\theta})\pi_{\underline{\theta}}(\underline{\theta}) d\underline{\theta}} \quad (2)$$

contains the information about $\underline{\theta}$ in light of the data \underline{x}_1 and prior assumptions. If new (independent and identically distributed to \underline{x}_1) data \underline{x}_2 becomes available, then the posterior for $\underline{\theta}$ in light of the combined data $(\underline{x}_1, \underline{x}_2)$ is

$$\pi_{\underline{\theta}|\underline{X}}(\underline{\theta} | \underline{x}_1, \underline{x}_2) = \frac{f_{\underline{X}|\underline{\theta}}(\underline{x}_1, \underline{x}_2 | \underline{\theta})\pi_{\underline{\theta}}(\underline{\theta})}{f_{\underline{X}}(\underline{x}_1, \underline{x}_2)} = \frac{f_{\underline{X}|\underline{\theta}}(\underline{x}_1, \underline{x}_2 | \underline{\theta})\pi_{\underline{\theta}}(\underline{\theta})}{\int f_{\underline{X}|\underline{\theta}}(\underline{x}_1, \underline{x}_2 | \underline{\theta})\pi_{\underline{\theta}}(\underline{\theta}) d\underline{\theta}}.$$

But note also that

$$\pi_{\underline{\theta}|\underline{X}}(\underline{\theta} | \underline{x}_1, \underline{x}_2) = \frac{f_{\underline{X}|\underline{\theta}}(\underline{x}_2 | \underline{\theta})\pi_{\underline{\theta}|\underline{X}}(\underline{\theta} | \underline{x}_1)}{f_{\underline{X}}(\underline{x}_2 | \underline{x}_1)}$$

where $\pi_{\underline{\theta}|\underline{X}}(\underline{\theta} | \underline{x}_1)$ is the posterior for $\underline{\theta}$ from (2), and

$$f_{\underline{X}}(\underline{x}_2 | \underline{x}_1) = \frac{f_{\underline{X}}(\underline{x}_1, \underline{x}_2)}{f_{\underline{X}}(\underline{x}_1)} = \frac{\int f_{\underline{X}|\underline{\theta}}(\underline{x}_1 | \underline{\theta})\pi_{\underline{\theta}}(\underline{\theta}) d\underline{\theta}}{\int f_{\underline{X}|\underline{\theta}}(\underline{x}_1, \underline{x}_2 | \underline{\theta})\pi_{\underline{\theta}}(\underline{\theta}) d\underline{\theta}} \quad (3)$$

1.4 Sufficiency

Direct from (1), if $\underline{T}(X)$ is a sufficient statistic for θ , it follows that

$$\pi_{\theta|\underline{X}}(\theta|x) = \frac{f_{\underline{X}|\theta}(x|\theta)\pi_{\theta}(\theta)}{f_{\underline{X}}(x)} = \frac{g(\underline{T}(x), \theta)h(x)\pi_{\theta}(\theta)}{f_{\underline{X}}(x)} = \left[\frac{h(x)}{f_{\underline{X}}(x)} \right] g(\underline{T}(x), \theta)\pi_{\theta}(\theta)$$

where $f_{\underline{X}|\theta}(x|\theta) = g(\underline{T}(x), \theta)h(x)$ by the Neyman factorization result. Thus the posterior distribution of θ only depends on the data through $\underline{T}(x)$.

Lemma 1.2 If $\underline{T}(X)$ is a sufficient statistic for θ (in the classical sense) then

$$\pi_{\theta|\underline{X}}(\theta|x) = \pi_{\theta|\underline{X}}(\theta|\underline{T}(x))$$

for all prior specifications $\pi_{\theta}(\theta)$.

Proof. By definition

$$f_{\underline{X}|\theta}(x|\theta) = f_{\underline{X}, \underline{T}(X)|\theta}(x, t|\theta)$$

if $t = \underline{T}(x)$, and zero otherwise. Thus, by sufficiency,

$$f_{\underline{X}|\theta}(x|\theta) = f_{\underline{X}|\underline{T}(X)}(x|t)f_{\underline{T}|\theta}(t|\theta)$$

and hence

$$\begin{aligned} \pi_{\theta|\underline{X}}(\theta|x) \propto f_{\underline{X}|\theta}(x|\theta)\pi_{\theta}(\theta) &= f_{\underline{X}|\underline{T}(X)}(x|t)f_{\underline{T}(X)|\theta}(t|\theta)\pi_{\theta}(\theta) \\ &\propto f_{\underline{T}(X)|\theta}(t|\theta)\pi_{\theta}(\theta) \\ &\propto \pi_{\theta|\underline{T}(X)}(\theta|t) \end{aligned}$$

with the constant of proportionality equal to one, as both sides must integrate to one. ■
Statistic $\underline{T}(X)$ is sufficient for θ in the Bayesian sense if

$$\pi_{\theta|\underline{X}}(\theta|x) \propto f_{\underline{T}(X)|\theta}(t|\theta)\pi_{\theta}(\theta)$$

Lemma 1.3 Statistic $\underline{T}(X)$ is sufficient in the Bayesian sense if and only if it is sufficient in the Classical sense.

Proof. For the if, see the previous Lemma. For the only if, suppose

$$\pi_{\theta|\underline{X}}(\theta|x) = f_{\underline{T}(X)|\theta}(t|\theta)\pi_{\theta}(\theta)h(x).$$

By (1),

$$\pi_{\theta|\underline{X}}(\theta|x) = \frac{f_{\underline{X}|\theta}(x|\theta)\pi_{\theta}(\theta)}{f_{\underline{X}}(x)} \implies \frac{f_{\underline{X}|\theta}(x|\theta)}{f_{\underline{X}}(x)} = \frac{\pi_{\theta|\underline{X}}(\theta|x)}{\pi_{\theta}(\theta)} = f_{\underline{T}(X)|\theta}(t|\theta)h(x).$$

Hence

$$f_{\underline{X}|\theta}(x|\theta) = f_{\underline{T}(X)|\theta}(t|\theta)h(x)f_{\underline{X}}(x) = g(t, \theta)h^*(x)$$

where $g(t, \theta) = f_{\underline{T}(X)|\theta}(t|\theta)$ and $h^*(x) = h(x)f_{\underline{X}}(x)$. Thus $\underline{T}(X)$ is sufficient in the classical sense. ■

2 Construction of Prior Distributions

In the Bayesian formulation, the prior density $\pi_{\underline{\theta}}(\underline{\theta})$ plays an important role. There are several methods via which the prior can be specified quantitatively; from historical or training data; by subjective assessment, similar to the subjective assessment of probabilities in elementary probability theory; by matching to a desired functional form; or in a **non-informative** or **vague** specification, where the prior probability is spread evenly across the parameter space.

2.1 Conjugate Priors

For some models, a **conjugate prior** can be chosen; this prior combines with the likelihood in such a way to give an analytically tractable posterior calculation. Consider a class of distributions \mathcal{F} indexed by parameter $\underline{\theta}$

$$\mathcal{F} = \left\{ f_{X|\underline{\theta}}(x|\underline{\theta}) : \underline{\theta} \in \Theta \right\}$$

A class \mathcal{P} of prior distributions for $\underline{\theta}$ is a conjugate family for \mathcal{F} if the posterior distribution for $\underline{\theta}$ resulting from data \underline{x} is an element of \mathcal{P} for all $f_{X|\underline{\theta}} \in \mathcal{F}$, $\pi_{\underline{\theta}} \in \mathcal{P}$ and $\underline{x} \in \mathcal{X}$.

Example 2.1 Suppose that $f_{X|\underline{\theta}}(x|\underline{\theta})$ is an Exponential Family distribution

$$f_{X|\underline{\theta}}(x|\underline{\theta}) = h(x)c(\underline{\theta}) \exp \left\{ \sum_{j=1}^k t_j(x)w_j(\underline{\theta}) \right\}$$

so that for a random sample of size n

$$L(\underline{\theta}|\underline{x}) = h(\underline{x})\{c(\underline{\theta})\}^n \exp \left\{ \sum_{j=1}^k T_j(\underline{x})w_j(\underline{\theta}) \right\} \quad (4)$$

where $\underline{T}(\underline{x}) = (T_1(\underline{x}), \dots, T_k(\underline{x}))^\top$ and

$$T_j(\underline{x}) = \sum_{i=1}^n t_j(x_i).$$

Suppose that

$$\pi_{\underline{\theta}}(\underline{\theta}) = d(\alpha, \underline{\beta})\{c(\underline{\theta})\}^\alpha \exp \left\{ \sum_{j=1}^k \beta_j w_j(\underline{\theta}) \right\} \quad (5)$$

where α and $\underline{\beta} = (\beta_1, \dots, \beta_k)^\top$ are **hyperparameters**. Combining equations (4) and (5) yields the posterior distribution up to proportionality as

$$\begin{aligned} \pi_{\underline{\theta}|\underline{x}}(\underline{\theta}|\underline{x}) &\propto \{c(\underline{\theta})\}^{\alpha+n} \exp \left\{ \sum_{j=1}^k [\beta_j + T_j(\underline{x})]w_j(\underline{\theta}) \right\} \\ &= \{c(\underline{\theta})\}^{\alpha^*} \exp \left\{ \sum_{j=1}^k \beta_j^* w_j(\underline{\theta}) \right\} \end{aligned}$$

The normalizing constant can be deduced to be $d(\alpha + n, \underline{\beta} + \underline{T}(\underline{x}))$, and hence the posterior distribution has the same functional form as the prior, but with parameters updated to

$$\alpha^* = \alpha + n \quad \underline{\beta}^* = (\beta_1^*, \dots, \beta_k^*)^\top = (\beta_1 + T_1(\underline{x}), \dots, \beta_k + T_k(\underline{x}))^\top.$$

2.2 Ignorance Priors

A non-informative prior expresses **prior ignorance** about the parameter of interest.

- If $\Theta = \{\underline{\theta}_1, \dots, \underline{\theta}_k\}$ (that is, $\underline{\theta}$ is known to take one of a finite number of possible values). Then a non-informative prior places equal probability on each value, that is,

$$\pi_{\underline{\theta}}(\underline{\theta}) = \frac{1}{k} \quad \underline{\theta} \in \Theta.$$

- If Θ is a **bounded region**, then a natural non-informative prior is **constant** on Θ .
- If the parameter space Θ is uncountable and unbounded, however, a non-informative prior specification is more difficult to construct. A naive prior specification would be to set $\pi_{\underline{\theta}}(\underline{\theta})$ to be a constant, although this prior does not give a valid probability measure as it does not integrate to 1 over Θ .

A prior distribution $\pi_{\underline{\theta}}(\underline{\theta})$ for parameter $\underline{\theta}$ is termed **improper** if it does not integrate to 1.

Even for improper priors, (1) can be used to compute the posterior density, which itself will often not be improper (that is, the posterior is **proper**) However, if $\underline{\phi} = \underline{g}(\underline{\theta})$ is a transformation of $\underline{\theta}$ that may of inferential interest, then by elementary transformation results, including the Jacobian of the transform $J(\underline{\theta} \rightarrow \underline{\phi})$, it follows that

$$\pi_{\underline{\theta}}(\underline{\theta}) = c \quad \implies \quad \pi_{\underline{\phi}}(\underline{\phi}) = c \times J(\underline{\theta} \rightarrow \underline{\phi})$$

which may **not** be constant, and hence a **non-uniform** prior on $\underline{\phi}$ results. This is perhaps unsatisfactory, and so the following procedure may be preferable.

2.3 Jeffreys' Prior

Consider the prior $\pi_{\underline{\theta}}(\underline{\theta})$ for parameter $\underline{\theta}$ in probability model $f_{X|\underline{\theta}}(x|\underline{\theta})$ determined by

$$\pi_{\underline{\theta}}(\underline{\theta}) \propto |I(\underline{\theta})|^{1/2}$$

where $I(\underline{\theta})$ is Fisher's Information, and $|I(\underline{\theta})|$ indicates the absolute value of the determinant of $I(\underline{\theta})$. Recall that

$$I(\underline{\theta}) = E_{f_{X|\underline{\theta}}} \left[\underline{S}(X; \underline{\theta}) \underline{S}(X; \underline{\theta})^T \right] = -E_{f_{X|\underline{\theta}}} [\Psi(X; \underline{\theta})]$$

and also that $\underline{S}(X; \underline{\theta})$ is the $k \times 1$ vector score function with j th element

$$S_j(X; \underline{\theta}) = \frac{\partial}{\partial \theta_j} \log f_{X|\underline{\theta}}(x|\underline{\theta}) \quad j = 1, \dots, k$$

and $\Psi(X; \underline{\theta})$ is the $k \times k$ matrix of second partial derivatives with (j, l) th element

$$\frac{\partial^2}{\partial \theta_j \partial \theta_l} \log f_{X|\underline{\theta}}(x|\underline{\theta})$$

Example 2.2 *Binomial*(m, θ). We have

$$\begin{aligned}\log f_{X|\theta}(x|\theta) &= \log \binom{m}{x} + x \log \theta + (m-x) \log(1-\theta) \\ S(x; \theta) &= \frac{x}{\theta} - \frac{(m-x)}{(1-\theta)} \\ \Psi(x; \theta) &= -\frac{x}{\theta^2} - \frac{(m-x)}{(1-\theta)^2}\end{aligned}$$

so therefore

$$I(\theta) = -E_{f_{X|\theta}} \left[-\frac{X}{\theta^2} - \frac{(m-X)}{(1-\theta)^2} \right] = \frac{m\theta}{\theta^2} + \frac{m(1-\theta)}{(1-\theta)^2} = \frac{m}{\theta(1-\theta)}$$

and hence

$$\pi_{\theta}(\theta) \propto |I(\theta)|^{1/2} = \{\theta(1-\theta)\}^{-1/2}$$

Lemma 2.1 Jeffreys' prior is invariant under 1-1 transformations, that is, if $\underline{\phi} = \underline{\phi}(\underline{\theta})$, then the prior for $\underline{\phi}$ obtained by reparameterization from $\underline{\theta}$ to $\underline{\phi}$ in the prior for $\underline{\theta}$, is precisely Jeffreys' prior for $\underline{\phi}$.

Proof. Let $\underline{\phi} = \underline{\phi}(\underline{\theta})$ be a 1-1 transformation. Denote by $\ell_{\theta}(x|\theta)$ and $\ell_{\phi}(x|\underline{\phi})$ the log pdfs in the two parameterizations. Then by the rules of partial differentiation

$$\frac{\partial \ell_{\phi}}{\partial \phi_j} = \sum_{l=1}^k \frac{\partial \ell_{\theta}}{\partial \theta_l} \frac{\partial \theta_l}{\partial \phi_j} \quad j = 1, \dots, k$$

so that

$$\underline{S}(X; \underline{\phi}) = \Lambda(\underline{\theta}, \underline{\phi}) \underline{S}(X; \underline{\theta})$$

where $\Lambda(\underline{\theta}, \underline{\phi})$ is the $k \times k$ matrix with (j, l) th element

$$\frac{\partial \theta_l}{\partial \phi_j}$$

In fact, $\Lambda(\underline{\theta}, \underline{\phi})$ is just the Jacobian of the transformation from $\underline{\theta}$ to $\underline{\phi}$, $J(\underline{\theta} \rightarrow \underline{\phi})$. Hence

$$I_{\phi}(\underline{\phi}) = \Lambda(\underline{\theta}, \underline{\phi}) I_{\theta}(\underline{\theta}) \Lambda(\underline{\theta}, \underline{\phi})^{\top}$$

and so

$$|I_{\phi}(\underline{\phi})| = |\Lambda(\underline{\theta}, \underline{\phi}) I_{\theta}(\underline{\theta}) \Lambda(\underline{\theta}, \underline{\phi})^{\top}| = |\Lambda(\underline{\theta}, \underline{\phi})|^2 |I_{\theta}(\underline{\theta})|$$

and

$$|I_{\phi}(\underline{\phi})|^{1/2} = |\Lambda(\underline{\theta}, \underline{\phi})| |I_{\theta}(\underline{\theta})|^{1/2}.$$

Thus

$$\pi_{\underline{\phi}}(\underline{\phi}) \propto |I_{\phi}(\underline{\phi})|^{1/2} = |\Lambda(\underline{\theta}, \underline{\phi})| |I_{\theta}(\underline{\theta})|^{1/2} = |\Lambda(\underline{\theta}, \underline{\phi})| \pi_{\underline{\theta}}(\underline{\theta})$$

and Jeffreys' prior for $\underline{\phi}$ is identical to the one that would be obtained by constructing Jeffreys' prior for $\underline{\theta}$ and reparameterizing to $\underline{\phi}$. ■

Example 2.3 *Binomial*(m, θ). Suppose that $\phi = \theta/(1 - \theta)$ (so that $\theta = \phi/(1 + \phi)$). Then

$$\begin{aligned}\log f_{X|\theta}(x|\phi) &= \log \binom{m}{x} + x \log \phi + m \log(1 + \phi) \\ S(x; \phi) &= \frac{x}{\phi} - \frac{m}{(1 + \phi)} \\ \Psi(x; \phi) &= -\frac{x}{\phi^2} + \frac{m}{(1 + \phi)^2}\end{aligned}$$

so therefore

$$I(\phi) = -E_{f_{X|\phi}} \left[-\frac{X}{\phi^2} + \frac{m}{(1 + \phi)^2} \right] = \frac{m\phi}{(1 + \phi)\phi^2} - \frac{m}{(1 + \phi)^2} = \frac{m}{\phi(1 + \phi)^2}$$

and hence

$$\pi_\phi(\phi) \propto |I(\phi)|^{1/2} = \{\phi(1 + \phi)^2\}^{-1/2}.$$

Now, recall that Jeffreys' prior for θ takes the form

$$\pi_\theta(\theta) \propto \{\theta(1 - \theta)\}^{-1/2}$$

The Jacobian of the transformation from θ to ϕ is $(1 + \phi)^2$, and thus using the univariate transformation theorem

$$\pi_\phi(\phi) \propto \{\phi/(1 + \phi)^2\}^{-1/2}(1 + \phi)^2 = \{\phi(1 + \phi)^2\}^{-1/2}$$

matching the result found above.

2.4 Location and Scale Parameters

Consider the one-dimensional case: θ is a **location parameter** if

$$f_{X|\theta}(x|\theta) = f_X(x - \theta).$$

θ is a **scale parameter** if

$$f_{X|\theta}(x|\theta) = \frac{1}{\theta} f_X\left(\frac{x}{\theta}\right)$$

For a location parameter, for a non-informative prior, it is required to have, for set $A \subset \Theta$

$$\int_A \pi_\theta(\theta) d\theta = \int_{A_c} \pi_\theta(\theta) d\theta$$

where $A_c = \{\theta : \theta - c \in A\}$ for scalar c . Therefore, for all c ,

$$\int_{A_c} \pi_\theta(\theta) d\theta = \int_A \pi_\theta(\theta - c) d\theta \implies \pi_\theta(\theta) = \pi_\theta(\theta - c) \implies \pi_\theta(\theta) = \text{constant}$$

For a scale parameter, it is required to have, for set $A \subset \Theta$

$$\int_A \pi_\theta(\theta) d\theta = \int_{A_c} \pi_\theta(\theta) d\theta$$

where now $A_c = \{\theta : c\theta \in A\}$ for scalar c . But, therefore, for all c ,

$$\int_{A_c} \pi_\theta(\theta) d\theta = \int_A \pi_\theta(c\theta) d\theta \implies \pi_\theta(\theta) = c\pi_\theta(c\theta) \implies \pi_\theta(\theta) \propto \frac{1}{\theta}$$

3 Bayesian Inference and Decision Making

The key components of a **decision problem** are as follows;

- a **decision** δ is to be made, and the decision is selected from some set \mathcal{D} of alternatives.
- a true **state of nature**, v , lying in set Υ
- a **loss function**, $\mathcal{L}(\delta, v)$, for decision δ and state v , which records the loss (or penalty) incurred when the true state of nature is v and the decision made is δ .

Broadly, we aim to select the decision to minimize the loss. In the presence of uncertainty, we minimize the **expected loss**.

In an estimation context, the decision is the **estimate** of the parameter, and the true state of nature is the true value of the parameter. The minimization of expected loss proceeds in slightly different fashion dependent on whether data are observed;

- (a) With prior $\pi_{\theta}(\theta)$ and no data, the **expected loss** (or **Bayes loss**) is defined as

$$E_{\pi_{\theta}} [\mathcal{L}(\delta, \theta)] = \int \mathcal{L}(\delta, \theta) \pi_{\theta}(\theta) d\theta$$

The optimal Bayesian decision is

$$\hat{\delta}_B = \underset{\delta \in \mathcal{D}}{\operatorname{argmin}} E_{\theta} [\mathcal{L}(\delta, \theta)]$$

that is, the decision that minimizes the Bayes loss.

- (b) If data are available, the optimal decision will intuitively become a function of the data. Suppose now that the decision in light of the data is denoted $\delta(\underline{x})$ (a function from \mathcal{X} to \mathcal{D} , and the associated loss is $\mathcal{L}(\delta(\underline{x}), \theta)$)

The **risk** associated with decision $\delta(\underline{X})$ is the expected loss associated with $\delta(\underline{X})$, with the expectation taken over the distribution of \underline{X} given θ

$$R_{\delta}(\theta) = E_{f_{\underline{X}|\theta}} [\mathcal{L}(\delta(\underline{X}), \theta)] = \int_{\Theta} \mathcal{L}(\delta(\underline{x}), \theta) f_{\underline{X}|\theta}(\underline{x}|\theta) d\underline{x}$$

The **Bayes risk** is the expected risk $R_{\delta}(\theta)$ associated with $\delta(\underline{X})$, with the expectation taken over the prior distribution of θ

$$\begin{aligned} R(\delta) &= E_{\pi_{\theta}} [R_{\delta}(\theta)] = E_{\pi_{\theta}} \left[E_{f_{\underline{X}|\theta}} [\mathcal{L}(\delta(\underline{X}), \theta)] \right] = \int_{\Theta} \left\{ \int_{\mathcal{X}} \mathcal{L}(\delta(\underline{x}), \theta) f_{\underline{X}|\theta}(\underline{x}|\theta) d\underline{x} \right\} \pi_{\theta}(\theta) d\theta \\ &= \int_{\Theta} \int_{\mathcal{X}} \mathcal{L}(\delta(\underline{x}), \theta) f_{\underline{X}}(\underline{x}) \pi_{\theta|\underline{X}}(\theta|\underline{x}) d\underline{x} d\theta \\ &= \int_{\mathcal{X}} \left\{ \int_{\Theta} \mathcal{L}(\delta(\underline{x}), \theta) \pi_{\theta|\underline{X}}(\theta|\underline{x}) d\theta \right\} f_{\underline{X}}(\underline{x}) d\underline{x} \end{aligned}$$

where, by equation (1), $f_{\underline{X}|\theta}(\underline{x}|\theta) \pi_{\theta}(\theta) = f_{\underline{X}}(\underline{x}) \pi_{\theta|\underline{X}}(\theta|\underline{x})$. With prior $\pi_{\theta}(\theta)$ and fixed data \underline{x} the optimal Bayesian decision, termed the **Bayes rule** is $\hat{\delta}_B = \underset{\delta \in \mathcal{D}}{\operatorname{argmin}} R(\delta)$ so that

$$\hat{\delta}_B = \underset{\delta \in \mathcal{D}}{\operatorname{argmin}} \int_{\mathcal{X}} \left\{ \int_{\Theta} \mathcal{L}(\delta(\underline{x}), \theta) \pi_{\theta|\underline{X}}(\theta|\underline{x}) d\theta \right\} f_{\underline{X}}(\underline{x}) d\underline{x} = \underset{\delta \in \mathcal{D}}{\operatorname{argmin}} \int_{\Theta} \mathcal{L}(\delta(\underline{x}), \theta) \pi_{\theta|\underline{X}}(\theta|\underline{x}) d\theta$$

as only the inner integral depends on δ . That is, the decision that minimizes the Bayes risk minimizes **posterior expected loss** in making decision δ , with expectation taken with respect to the posterior distribution $\pi_{\theta|\underline{X}}(\theta|\underline{x})$.

Results for Different Loss Functions in the One Parameter case

(I) Under **squared-error loss**

$$\mathcal{L}(\delta(\underline{x}), \theta) = (\delta(\underline{x}) - \theta)^2$$

the Bayes rule for estimating θ is

$$\hat{\delta}_B(\underline{x}) = \hat{\theta}_B(\underline{x}) = \mathbb{E}_{\pi_{\theta|\underline{X}}} [\theta | \underline{X} = \underline{x}] = \int \theta \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta$$

that is, the **posterior expectation**. The expected posterior loss is

$$\int \mathcal{L}(\delta(\underline{x}), \theta) \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta = \int (\delta(\underline{x}) - \theta)^2 \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta$$

which needs to be minimized with respect to $\delta(\underline{x})$. Write $\delta = \delta(\underline{x})$. Then

$$\frac{d}{d\delta} \left\{ \int (\delta - \theta)^2 \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta \right\} = \int \frac{d}{d\delta} \left\{ (\delta - \theta)^2 \right\} \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta = \int 2(\delta - \theta) \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta$$

and equating this to zero gives

$$\int (\delta - \theta) \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta = 0 \quad \implies \quad \delta = \int \theta \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta = \mathbb{E}_{\pi_{\theta|\underline{X}}} [\theta | \underline{X} = \underline{x}]$$

and hence the optimal δ is the posterior expectation as stated.

(II) Under **absolute error loss**

$$\mathcal{L}(\delta(\underline{x}), \theta) = |\delta(\underline{x}) - \theta|$$

the Bayes rule for estimating θ is the solution of

$$\int_{-\infty}^{\delta(\underline{x})} \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta = \frac{1}{2}$$

that is, the **posterior median**. The expected posterior loss is

$$\int \mathcal{L}(\delta(\underline{x}), \theta) \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta = \int |\delta(\underline{x}) - \theta| \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta$$

which needs to be minimized with respect to $\delta(\underline{x})$. Set $\delta = \delta(\underline{x})$. Then

$$\int |\delta - \theta| \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta = \int_{-\infty}^{\delta} (\delta - \theta) \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta + \int_{\delta}^{\infty} (\theta - \delta) \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta \quad (6)$$

Differentiating with respect to δ the first term using the product rule yields

$$\begin{aligned} \frac{d}{d\delta} \left\{ \int_{-\infty}^{\delta} (\delta - \theta) \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta \right\} &= \frac{d}{d\delta} \left\{ \delta \int_{-\infty}^{\delta} \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta - \int_{-\infty}^{\delta} \theta \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta \right\} \\ &= \delta \pi_{\theta|\underline{X}}(\delta | \underline{x}) + \int_{-\infty}^{\delta} \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta - \delta \pi_{\theta|\underline{X}}(\delta | \underline{x}). \end{aligned}$$

Similarly

$$\frac{d}{d\delta} \left\{ \int_{\delta}^{\infty} (\theta - \delta) \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta \right\} = -\delta \pi_{\theta|\underline{X}}(\delta | \underline{x}) - \int_{\delta}^{\infty} \pi_{\theta|\underline{X}}(\theta | \underline{x}) d\theta + \delta \pi_{\theta|\underline{X}}(\delta | \underline{x})$$

Thus, equating the derivative of equation (6) to zero yields

$$\int_{-\infty}^{\delta} \pi_{\theta|\underline{X}}(\theta|\underline{x}) d\theta - \int_{\delta}^{\infty} \pi_{\theta|\underline{X}}(\theta|\underline{x}) d\theta = 0$$

so that

$$\int_{-\infty}^{\delta} \pi_{\theta|\underline{X}}(\theta|\underline{x}) d\theta = \int_{\delta}^{\infty} \pi_{\theta|\underline{X}}(\theta|\underline{x}) d\theta = \frac{1}{2}$$

and hence the optimal δ is the posterior median.

(III) Under **zero-one loss**

$$\mathcal{L}(\delta(\underline{x}), \theta) = \begin{cases} 0 & \delta(\underline{x}) = \theta \\ 1 & \delta(\underline{x}) \neq \theta \end{cases}$$

the Bayes rule for estimating θ is

$$\hat{\delta}_B(\underline{x}) = \hat{\theta}_B(\underline{x}) = \operatorname{argmax}_{\theta \in \Theta} \pi_{\theta|\underline{X}}(\theta|\underline{x})$$

that is, the **posterior mode**. The expected posterior loss is

$$\int \mathcal{L}(\delta(\underline{x}), \theta) \pi_{\theta|\underline{X}}(\theta|\underline{x}) d\theta = \int_{\Theta \setminus \delta(\underline{x})} \pi_{\theta|\underline{X}}(\theta|\underline{x}) d\theta$$

which needs to be minimized with respect to $\delta(\underline{x})$. The optimal $\delta(\underline{x})$ is thus the posterior mode as stated