## CH3   PRIOR INFORMATION AND SUBJECTIVE PROBABILITY

Determing prior: discrete variables $\Big\{$ Scoring / Betting

### § 3.2 Determining Prior Density

- Histogram approach
- The relative likelihood approach
- Matching a given functional form
  - Estimating prior moments $\mu, \sigma^2$
    - ☺ The tail of a density can have a drastic effect on its moments   Eg. $\int_b^\infty \theta \cdot (k\theta^{-2}) d\theta = \infty$
  - Estimating fractiles
  - Equivalent sample size / device of imaginary results

For normal distribution, the posterior with normal prior : $N(\mu, \sigma^2)$

$$\left(\frac{\sigma^2}{\sigma^2 + 1/n}\right) \bar{x} + \left(\frac{1/n}{\sigma^2 + 1/n}\right) \mu$$

$\sigma^2 = 1/n^*$  equivalent to have $1/\sigma^2$ samples of mean $\mu$.

- ☺ Useful only when certain specific functional forms
- ☹ Tend to considerably underestimate the amount of information carried by a sample of size $n$.

  - CDF Determination  (CDF: cumulative distribution function)
    1° Subjectively determine several $\alpha$-fractiles, $z(\alpha)$.
    2° Plot the points $(\alpha, z(\alpha))$ and sketch a smooth curve joining them

§ 3.3  Noninformative Priors

No (or minimal) prior information available + Compelling Bayesian analysis

$$\Downarrow$$

Noninformative Prior

<u>Example</u> :  Discrete variables: uniform

<u>Example</u> :  $\Theta = (-\infty, \infty)$  uniform $\Rightarrow$  $\pi(\theta) = \overset{\tfrac{1}{\|}}{c} > 0.$

$\int \pi(\theta)\, d\theta = \infty$  may or may not cause problems

Severe (though unjustified) criticism:

Lack of Invariance under Transformation

<u>Example</u>:  Let  $\eta = \exp\{\theta\}.$

$\pi^*(\eta) = \eta^{-1}\,\pi(\log \eta)$

Let $y = g(x)$, $x = h(y)$
$(h = g^{-1})$
$f_Y(y) = |h'(y)|\, f_X(h(y))$

<u>Example</u>  (Noninformative priors for location problems)

Suppose $\mathcal{X}$ and $\Theta$ are subsets of $\mathbb{R}^p$. and density of $X$ is of the form $f(x-\theta)$

$$\left(\text{E.g.}\quad x - \theta \sim N(\theta, \underset{\uparrow}{\Sigma})\right)$$
$\Sigma$ fixed.

$c \in \mathbb{R}^p$. fixed

Imagine that, instead of $X$, we observe $Y = X + c$.

Then $Y$ has density $f(y - \eta)$

$\Rightarrow$ The $(X, \theta)$ and $(Y, \eta)$ problems are identical in structure

♪ Noninformative priors in general settings, please see textbooks pp. 87–8

Let $\pi_1$ and $\pi_2$ denote the noninformative priors for $\theta$ and $\eta$

Invariant
noninformative :  $\qquad P^{\pi_1}(\theta \in A) = P^{\pi_2}(\eta \in A) \qquad \forall A$ in $\mathbb{R}^p$
assumption

Then $\qquad P^{\pi_2}(\eta \in A) = P^{\pi_1}(\theta + c \in A) = P^{\pi_1}(\theta \in A-c)$

Combining the above equations:

$$P^{\pi_1}(\theta \in A) = P^{\pi_1}(\theta \in A-c)$$

$$\int_A \pi(\theta)\,d\theta = \int_{A-c} \pi(\theta)\,d\theta = \int_A \pi(\theta-c)\,d\theta$$

*Unnecessary for*
∫ *intuitive thinking!*
∫ *Insufficient for*
*mathematical proof*

It can be shown that

$$\pi(\theta) = \pi(\theta-c)$$

Let $\theta = c$, $\qquad\qquad \pi(c) = \pi(0)$

**Example** (Noninformative priors for scale problems)

$$y = \frac{x}{\sigma} \qquad\qquad f_Y(y) = \sigma^{-1} f_X\left(\frac{x}{\sigma}\right)$$

$\sigma$ : a scale parameter $\quad$ *Eg*. $\sigma \sim N(0, \beta^2)$

Imagine that, instead of observing $\cancel{\varnothing} X$ we observe the random

variable $\cancel{\text{prex}}$ $(c>0)$ $\qquad\qquad$ Note that $X \sim \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right)$
$\qquad Y = cX$
Let $\eta = c \cdot \sigma$. $\qquad Y \sim \eta^{-1} f\left(\frac{y}{\eta}\right) \qquad\qquad Y = cX \sim \frac{1}{c\sigma} f\left(\frac{x}{c\sigma}\right)$

If $\mathcal{X} = \mathbb{R}^1$ of $\mathcal{X} = (0,\infty)$, then
$$(X, \sigma) \text{ is equivalent to } (Y, \eta)$$

Denote $\pi_1$ and $\pi_2$ the prior of $\sigma$ and $\eta$

$$P^{\pi_1}(\sigma \in A) = P^{\pi_2}(\eta \in A)$$

Since $\eta = c\sigma$

$$P^{\pi_2}(\eta \in A) = P^{\pi_1}(\sigma \in c^{-1}A) \qquad c^{-1}A = \{c^{-1}z : z \in A\}$$

Combining the above equations

$$P^{\pi}(\sigma \in A) = P^{\pi}(\sigma \in c^{-1}A)$$

Thus

$$\int_A \pi(\sigma)\, d\sigma = \int_{c^{-1}A} \pi(\sigma)\, d\sigma = \int_A \pi(c^{-1}\sigma)\, c^{-1} d\sigma$$

Choosing $\sigma = c$ in $\quad \pi(\sigma) = c^{-1}\pi(c^{-1}\sigma)$

$$\pi(c) = c^{-1}\pi(1).$$

Note that $\int_0^\infty \sigma^{-1}\, d\sigma = \infty$, $\pi$ is an improper prior.

## Example (The "Table Entry" Problem)

Observation: The frequencies of the integer $1..9$ being the first significant digit of the table entries are $\log\dfrac{\log(1 + i^{-1})}{\log 10}$

Explanation by "non informative priors"

$$\tilde{\pi}(\sigma) = \sigma^{-1}$$

Normalize on $(1, 10)$

$$\pi(\sigma) = \frac{\sigma^{-1}}{\log 10}$$

The probability of $i$ being the first significant digit

$$P_i = \int_i^{i+1} [\sigma \log 10]^{-1}\, d\sigma = \frac{\log(1 + i^{-1})}{\log 10}$$

May be coincidence, but intriguing

§ 3.4 Maximum Entropy Priors

See Adam L Berger et al., A Maximum Entropy Approach to Natural Language Processing. 1996.

Note: I remember that one or a few formulas in the above paper are wrong, when solving the Lagrangian.

§ 3.5 Using the Marginal Distribution to Determine the Prior

- **Definition**: The joint density of $X$ and $\theta$ is

$$h(x, \theta) = f(x|\theta) \pi(\theta)$$

The **marginal density**

$$m(x|\pi) = \int_\theta f(x|\theta) \, dF^\pi(\theta) = \begin{cases} \int_\theta f(x|\theta) \pi(\theta) \, d(\theta) & \text{(continuous)} \\ \sum_\theta f(x|\theta) \pi(\theta) \end{cases}$$

- Information about $m$ : $\begin{cases} \text{subjective knowledge} \\ \text{data itself (empirical Bayes)} \end{cases}$

- We consider also restricted classes of priors denoted as $\Gamma$

1° Priors of a given functional form
$$\Gamma = \{\pi : \pi(\theta) = g(\theta|\lambda), \ \lambda \in \Lambda\}$$

2° Priors of a given structural form
<u>Eg</u> $\theta_i$ independent: β $\Gamma = \{\pi : \pi(\theta) = \prod_{i=1}^{p} \pi_i(\theta_i)\}$

3° Priors close to an elicited prior
$$\Gamma = \{\pi : \pi(\theta) = (1-\varepsilon) \pi_0(\theta) + \varepsilon q(\theta), \ q \in \mathcal{Q}\}$$
<span>↗ elicited prior</span>     <span>↑ class of contamination</span>

- The ML-II approach to prior selection

Definition: Suppose $\Gamma$ is a class of priors under consideration, and that $\hat{\pi} \in \Gamma$ satisfies (for the observed data $x$).

$$m(x|\hat{\pi}) = \sup_{\pi \in \Gamma} m(x|\pi)$$

Then $\hat{\pi}$ will be called the <u>type II maximum likelihood</u> prior or <u>ML-II prior</u>.

If $\Gamma = \{\pi : \pi(\theta) = g(\theta|\lambda), \lambda \in \Lambda\}$

then $\sup_{\pi \in \Gamma} m(x|\pi) = \sup_{\lambda \in \Lambda} m(x|g(\theta|\lambda))$

Example: Let $X \sim N(\theta, \sigma_f^2)$

$\theta \sim N(\mu_\pi, \sigma_\pi^2)$

Then $m(x) = N(x_i| \mu_\pi, \sigma_\pi^2 + \sigma_f^2)$ $(\forall i)$

$$m(x|\pi) = \prod_{i=1}^{p} m_0(x_i|\pi_0)$$

$$= \prod_{i=1}^{p} \frac{1}{[2\pi(\sigma_\pi^2 + \sigma_f^2)]^{1/2}} \exp\left\{-\frac{(x_i - \mu_\pi)^2}{2(\sigma_\pi^2 + \sigma_f^2)}\right\}$$

$$= [2\pi(\sigma_\pi^2 + \sigma_f^2)]^{-p/2} \exp\left\{-\frac{\sum_{i=1}^{p}(x_i - \mu_\pi)^2}{2(\sigma_\pi^2 + \sigma_f^2)}\right\}$$

$$= [2\pi(\sigma_\pi^2 + \sigma_f^2)]^{-p/2} \exp\left\{-\frac{ps^2}{2(\sigma_\pi^2 + \sigma_f^2)}\right\} \exp\left\{\frac{-p(\bar{x}-\mu_\pi)^2}{2(\sigma_\pi^2 + \sigma_f^2)}\right\}$$

where $\bar{x} = \sum_{i=1}^{p} x_i / p.$   $s^2 = \sum_{i=1}^{p}(x_i - \bar{x})^2 / p$

To maximize $m(x|\pi)$ with respect to $\mu_\pi$ and $\sigma_\pi^2$. we first observe that $\mu_\pi$ has to be $\bar{x}$

Then we optimize, with respect to $\sigma_\pi^2$,

$$\psi(\sigma_\pi^2) = \left[ 2\pi( \sigma_\pi^2 + \sigma_f^2) \right]^{-p/2} \exp\left\{ \frac{-ps^2}{2(\sigma_\pi^2 + \sigma_f^2)} \right\}$$

We instead optimize $\lg \psi(\sigma_\pi^2)$. Indeed.

$$\frac{d}{d\sigma_\pi^2} \lg \psi(\sigma_\pi^2) = \frac{-p/2}{(\sigma_\pi^2 + \sigma_f^2)} + \frac{ps^2}{2(\sigma_\pi^2 + \sigma_f^2)^2} \overset{\triangle}{=} 0$$

$$\Rightarrow \sigma_\pi^2 = s^2 - \sigma_f^2.$$

Also we observe that $\sigma_\pi^2 \geq 0$.

Hence $\sigma_\pi^2 = \max\{0, s^2 - \sigma_f^2\}$

In conclusion ML-II prior is

$$\hat{\pi}_0 = \mathcal{N}(\hat{\mu}_\pi, \hat{\sigma}_\pi^2) \quad \text{where} \quad \hat{\mu}_\pi = \bar{x} \quad \text{and} \quad \hat{\sigma}_\pi^2 = \max\{0, s^2 - \sigma_f^2\}$$

**<u>Example</u>** For any $\pi$ in the $\varepsilon$-contamination class

$$\Gamma = \left\{ \pi: \pi(\theta) = (1-\varepsilon)\pi_0(\theta) + \varepsilon q(\theta), \quad q \in \mathcal{Q} \right\}$$

$$m(x|\pi) = \int f(x|\theta) \overbrace{\left[ (1-\varepsilon)\pi_0(\theta) + \varepsilon q(\theta) \right]}^{\pi(\theta)} d\theta$$

$$= (1-\varepsilon)\, m(x|\pi_0) + \varepsilon \cdot m(x|q)$$

$\rightarrowtail$ The ML-II prior can be found by maximizing $m(x|q)$
over $q \in \mathcal{Q}$, and using the maximizing $\hat{q}$ in the expression for $\pi$.

If $\mathcal{Q}$ is the class of anything.

$$m(x|q) = \int f(x|\theta) q(\theta) d\theta$$

To maximize this, we choose ~~a's prior~~ $q$ to be a one-point
distribution centered at $\hat{\theta}$ ($\hat{\theta}$ is the ML of $\theta$ given data).

Then
$$\hat{\pi} = (1-\varepsilon)\pi_0 + \varepsilon\langle\hat{\theta}\rangle$$

- The Moment Approach $\begin{cases} \text{"given functional form" of } \Gamma \\ \text{relate prior moments to moments of marginals} \end{cases}$

<u>Lemma</u> Let $\mu_f(\theta)$, $\sigma_f^2(\theta)$ be conditional mean and variance of $X$

(conditioned on $\theta$, i.e., $f(x|\theta)$). Let $\mu_m$ and $\sigma_m^2$ denote

the marginal mean and variance of $X$ (w.r.t. $m(x)$)

Assuming these quantities exist, then

$$\mu_m = \mathbb{E}^\pi[\mu_f(\theta)]$$

$$\sigma_m^2 = \mathbb{E}^\pi[\sigma_f^2(\theta)] + \mathbb{E}^\pi[(\mu_f(\theta)-\mu_m)^2]$$

<u>Proof.</u> (the continuous case)

$$\mu_m = \mathbb{E}^m[X] = \int_{\mathcal{X}} x \cdot m(x) \, dx$$

$$= \int_{\mathcal{X}} x \cdot \int_\Theta f(x|\theta)\pi(\theta) \, d\theta \, dx$$

$$= \int_\Theta \pi(\theta) \underbrace{\int_{\mathcal{X}} x f(x|\theta) \, dx}_{\mu_f(\theta)} \, d\theta$$

$$= \int_\Theta \pi(\theta) \cdot \mu_f(\theta) \, d\theta$$

$$\sigma_m^2 = \mathbb{E}^m[(X-\mu_m)^2] = \int_{\mathcal{X}} (x-\mu_m)^2 \int_\Theta f(x|\theta)\pi(\theta) \, d\theta \, dx$$

$$= \int_\Theta \pi(\theta) * \int_{\mathcal{X}} (x-\mu_m)^2 f(x|\theta) \, dx \, d\theta$$

$$= \mathbb{E}^\pi[\mathbb{E}_\theta^f[X-\mu_m]^2]$$

$$= \mathbb{E}^\pi[\mathbb{E}_\theta^f[(x-\mu_f(\theta)) + (\mu_f(\theta)-\mu_m)]^2]$$

$$= \mathbb{E}^\pi\Big[\underbrace{\mathbb{E}_\theta^f[x-\mu_f(\theta)]^2}_{\sigma_f^2(\theta)} +$$

$$2\,\mathbb{E}_\theta^f[(x-\mu_f(\theta))\cdot(\mu_f(\theta)-\mu_m)] +$$

$$\mathbb{E}_\theta^f[\mu_f(\theta)-\mu_m]^2$$

$$\Big]$$

$$= \mathbb{E}^\pi[\sigma_f^2(\theta)] + \mathbb{E}^\pi[\mu_f(\theta)-\mu_m]^2$$

Gotto

**Corollary** i) If $\mu_f(\theta) = \theta$, then $\mu_m = \mu_\pi$, where $\mu_\pi = E^\pi[\theta]$, prior mean

ii) If, in addition, $\sigma_f^2(\theta) = \sigma_f^2$ is a constant independent of $\theta$,

then $\sigma_m^2 = \sigma_f^2 + \sigma_\pi^2$, where $\sigma_\pi^2$ is the prior variance.

$\rightsquigarrow$ $\mu_m$, $\sigma_m^2$ can usually be estimated by ML-II or subjective experience,

We can then solve the prior.

**Example**     Let $X \sim N(\theta, 1)$.     $\Gamma = \{ N(\mu_\pi, \sigma_\pi^2) \}$

If we know, either by subjective experience or type II ML,

that prective density yield $\mu_m = 1$, $\sigma_m^2 = 3$. of X

Using corollary, & we have $\mu_m = \mu_\pi$, $\sigma_m^2 = 1 + \sigma_\pi^2$

$\phantom{Using corollary, & we have} \underset{1}{\|} \phantom{\mu_m = \mu_\pi,} \underset{3}{\|}$

Thus, we conclude $\pi = N(1, 2)$

- The Distance Approach to Prior Selection

$\begin{cases} \Gamma \text{ not a "given functional form"} \\ \text{considerable information available about } m. \end{cases}$

$\Rightarrow 1^0$ estimate $m$.

$2^0$ use the integral relationship $m(x) = \int_\Theta f(x|\theta) dF^\pi(\theta)$ to estimate $\pi$

I.e., seek an estimate of $\pi$, say $\hat{\pi}$,

yielding $m_{\hat{\pi}}(x) = \int_\Theta f(x|\theta) dF^{\hat{\pi}}(\theta)$

is close to $\hat{m}(x)$.

By "close," we minimize $KL(\hat{m} \| m_{\hat{\pi}})$, given by

$KL(\hat{m} \| m_{\hat{\pi}}) = E^{\hat{m}}\left[ \log \frac{\hat{m}(X)}{m_{\hat{\pi}}(X)} \right] = \begin{cases} \int_x \hat{m}(x) \log\left[ \frac{\hat{m}(x)}{m_{\hat{\pi}}(x)} \right] dx & \text{(continuous)} \\ \sum_x \hat{m}(x) \log\left[ \frac{\hat{m}(x)}{m_{\hat{\pi}}(x)} \right] & \text{(discrete)} \end{cases}$

not related to $\hat{\pi}$

$$KL(\hat{m}, m_{\hat{\pi}}) = \mathbb{E}^{\hat{m}}\left[\lg \frac{\hat{m}(X)}{m_{\hat{\pi}}(X)}\right] = \mathbb{E}^{\hat{m}}\left[\lg \hat{m}(X)\right] - \mathbb{E}^{\hat{m}}\left[\lg m_{\hat{\pi}}(X)\right]$$

Minimizing $\quad KL(\hat{m} \| m_{\hat{\pi}}) \quad \Longleftrightarrow \quad$ maximizing $\mathbb{E}^{\hat{m}}\left[\lg m_{\hat{\pi}}(X)\right]$

CASE: $\quad \textcircled{H} = \{\theta_1, \cdots \theta_K\}$ finite.

Let $\quad p_i = \hat{\pi}(\theta_i)$ .

Then $\quad m_{\hat{\pi}}(x) = \sum_{i=1}^{k} f(x|\theta_i) p_i$

The problems becomes to maximize

$$\mathbb{E}^{\hat{m}}\left[\lg \left(\sum_{i=1}^{K} f(x|\theta_i) p_i\right]\right] = \sum_{j=1}^{n} q_j \frac{1}{n} \lg \left(\sum_{i=1}^{k} f(x_j|\theta_i) p_i\right)$$

CASE: $\quad \textcircled{H}$ continuous , the problem becomes very difficult.

$\textcircled{:(}$

- **Hierarchical Priors** 

stage 1: $\quad \Gamma = \{\pi_1(\theta|\lambda): \pi_1$ is of a given functional form, $\lambda \in \Lambda\}$

stage 2: $\quad \pi_2(\lambda) \quad$ on hyperparameter $\lambda$

- more robust
- usually use noninformative prior
- more stages are rarely used
- Hierarchical prior is a convenient representation

$$\pi(\theta) = \int_{\Lambda} \pi_1(\theta|\lambda) \, dF^{\pi_2}(\lambda)$$

is the standard prior distribution.

for criticism