

$$P(C_j) \rightsquigarrow \theta \mapsto \Pi(\theta) \quad \text{prior distribution}$$

$$P(C_j|B) \rightsquigarrow \theta \mapsto \Pi_{\text{post}}(\theta|Y=y) \quad \text{posterior distribution}$$

$$P(B|C_j) \rightsquigarrow y \mapsto f_Y(y|\theta)$$

$$\sum_k P(B|C_k) P(C_k) \rightsquigarrow y \mapsto \int f_Y(y|\theta) \Pi(\theta) d\theta$$

Suppose $Y = \sum_{j=1}^N X_j \sim \text{binomial}(N, p)$

$$\Pi(p) \sim \text{beta}(\alpha, \beta) \quad \theta \in [0, 1]$$

fixed

$$\Rightarrow \Pi_{\text{post}}(p|y) = \frac{\text{binomial}(y|N, p) \text{beta}(p|\alpha, \beta)}{\int_0^1 \text{binomial}(y|N, s) \text{beta}(s|\alpha, \beta) ds} \quad \textcircled{1}$$

$$\textcircled{1} = \binom{N}{y} p^y (1-p)^{N-y} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1}$$

$$= \binom{N}{y} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{y+\alpha-1} (1-p)^{N-y+\beta-1}$$

$$\textcircled{2} = \int \textcircled{1} ds = \binom{N}{y} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(y+\alpha)\Gamma(N-y+\beta)}{\Gamma(N+\alpha+\beta)}$$

$$\Rightarrow \Pi_{\text{post}}(p|y) = \frac{\textcircled{1}}{\textcircled{2}} = \frac{\Gamma(N+\alpha+\beta)}{\Gamma(y+\alpha)\Gamma(N-y+\beta)} p^{y+\alpha-1} (1-p)^{N-y+\beta-1}$$

$$= \text{beta}(y+\alpha, N-y+\beta)(p)$$

By taking expectation of the prior and posterior distributions

$$E_{\text{prior}}(p) = \frac{\alpha}{\alpha+\beta}$$

$$E_{\text{post}}(p) = \frac{y+\alpha}{N+\alpha+\beta} = \frac{N}{\alpha+\beta+N} \frac{y}{N} + \frac{\alpha+\beta}{\alpha+\beta+N} \frac{\alpha}{\alpha+\beta} \rightarrow \text{estimation on } p \text{ without any prior idea}$$

The posterior expectation is a linear combination

of prior expectation and experimental expectation.

Remark: As N becoming large, the prior expectation is becoming less and less important

Remark: The linear combination is not an accident;

\rightsquigarrow notion of **conjugate family** of a distribution.

IV.2 Evaluating estimators

Different estimators for θ can give you different values for θ .

Which one should we choose?

2.1 Mean square error

Framework $\underline{X} = (X_1, \dots, X_N)$ with $X_j = X \sim f(\cdot | \theta)$.

Let $W = W(\underline{X})$ be an estimator (\equiv statistic) for θ .

Def. The mean square error (MSE) of W is defined by

$$E((W - \theta)^2) \text{ still a function of } \theta$$

⚠ We made the choice of $s \mapsto s^2$. Later there will be a generation.

Observe that

$$\begin{aligned} E((W - \theta)^2) &= E((W - E(W) + E(W) - \theta)^2) \\ &= E((W - E(W))^2) + E(2(W - E(W))(E(W) - \theta)) + E((E(W) - \theta)^2) \\ &= \text{Var}_\theta(W) - (E_\theta(W) - \theta)^2 \end{aligned}$$

Def. If W is an estimator for θ , we set

$$\text{Bias}_\theta(W) := E_\theta(W) - \theta$$

If $\text{Bias}_\theta(W) = 0$ we say that W is unbiased.

$$\Rightarrow E_\theta((W - \theta)^2) = \text{Var}_\theta(W) + (\text{Bias}_\theta(W))^2$$

Remark: In Section II.1, we saw that if $E(X) = \mu$ and $\text{Var}(X) = \sigma^2$ then

$$\begin{aligned} E(\bar{X}) &= \mu \quad \text{and} \quad E(S^2) = \sigma^2 \\ \uparrow & \quad \quad \quad \uparrow \\ \text{sample mean} & \quad \quad \quad \text{sample variance} \end{aligned}$$

$\Rightarrow E(\bar{X})$ and $E(S^2)$ are unbiased estimators for μ and σ^2 resp.

Now if $X \sim n(\mu, \sigma^2)$ then $E((\bar{X} - \mu)^2) = \text{Var}(\bar{X}) = \frac{\sigma^2}{N}$

$$\text{and } E((S^2 - \sigma^2)^2) = \frac{2}{N-1} \sigma^4$$

Question: Can we find better with a smaller MSE estimator for μ and σ^2 ?

Remark: for $X \sim n(\mu, \sigma^2)$ consider

$$\tilde{S}^2 := \frac{1}{N} \sum_{j=1}^N (X_j - \bar{X})^2 = \frac{N-1}{N} S^2$$

$$\Rightarrow E(\tilde{S}^2) = \frac{N-1}{N} \sigma^2 \quad \Rightarrow \tilde{S}^2 \text{ is biased}$$

$$\text{But } E((\tilde{S}^2 - \sigma^2)^2) = \frac{2N-1}{N^2} \sigma^4 < \frac{2}{N-1} \sigma^4$$

$\Rightarrow \tilde{S}^2$ has a lower MSE, despite being biased.

~ Finding a better estimator is not a simple question.

We shall consider only unbiased estimators.

Def. An estimator W^* for θ parameter with value unknown is the **best unbiased estimator** for θ if $\text{Var}_\theta(W^*) \leq \text{Var}_\theta(W)$ for any estimator W for θ for any θ value of parameter with W^* and W unbiased.

Thm. (7.3.19) The best unbiased estimator is unique if it exists.

But it does not tell you how to find it.

Question: Can we get $\text{Var}_\theta(W) = 0$?

Thm. (7.3.9 + 7.3.10) (Cramer-Rao inequality) (based on Cauchy-Schwartz inequality)

$$\text{If } \frac{d}{d\theta} E_\theta(W) \equiv \frac{d}{d\theta} \int_{\mathbb{R}^N} W(x) f_x(x|\theta) dx = \int_{\mathbb{R}^N} \frac{\partial}{\partial \theta} [W(x) f_x(x|\theta)] dx$$

and $\text{Var}_\theta(W) < \infty$ then $\rightarrow = \theta$

$$\text{Var}_\theta(W) \geq \frac{\left(\frac{d}{d\theta} E_\theta(W) \right)^2}{N E_\theta \left[\left(\frac{\partial}{\partial \theta} \ln f_x(\cdot|\theta) \right)^2 \right]} = 1$$

information number
or Fisher information

If $\text{Var}_\theta(W)$ saturates satisfies equality in this inequality,

then W is the best unbiased estimator for any value of θ .

⚠ Otherwise for two W the $\text{Var}_\theta(W)$ may cross.

Remark: The following relation holds:

$$\text{If } \frac{d}{d\theta} E \left(\frac{\partial}{\partial \theta} \ln f(\cdot|\theta) \right) = \int \frac{\partial}{\partial \theta} \left[\frac{\partial}{\partial \theta} \ln f(\cdot|\theta) \right] f(x|\theta) dx \text{ then}$$

$$E_\theta \left[\left(\frac{\partial}{\partial \theta} \ln f(\cdot|\theta) \right)^2 \right] = -E_\theta \left(\frac{\partial^2}{\partial \theta^2} \ln f(\cdot|\theta) \right)$$

2.2 Loss function optimality

So far, $\text{MSE} = E((W-\theta)^2) = E(\mathcal{L}(W, \theta))$ with $\mathcal{L}(s, \theta) = (s-\theta)^2$

but other function $\mathcal{L} = \text{loss function}$ could be chosen. For example:

$$\cdot \mathcal{L}(s, \theta) = |s - \theta|$$

$$\cdot \mathcal{L}(s, \theta) = \frac{(s-\theta)^2}{|s|+1}$$

$$\cdot \mathcal{L}(s, \theta) = \begin{cases} (s-\theta)^2 & \text{for } s < \theta \\ 10(s-\theta)^2 & \text{for } s \geq \theta \end{cases}$$

$$\cdot \mathcal{L}(s, \theta) = \frac{s}{\theta} - 1 - \ln \left(\frac{s}{\theta} \right) \dots$$

Then the **risk function** is defined by $R(\theta, W) := E(\mathcal{L}(W, \theta))$

and we want to have a value of R close to 0.

Remark: the nice feature $E((W-\theta)^2) = \text{Var}(W) + \text{Bias}(W)^2$

will not be possible in general.

\Rightarrow minimizing $R(\theta, W)$ can be complicated for other function \mathcal{L} .

\Rightarrow It is a complicated question.