

Lecture 8 . Bayes Estimators and Average Risk Optimality.

9th November 2020.

* Explore an alternative approach to achieve optimality.

* X (data). $P = \{P_\theta : \theta \in \Omega\}$, $L(\theta, \delta)$ loss function, $R(\theta, \delta)$

$$E_{\theta}^{\Lambda}(L(\theta, \delta))$$

* Average Risk Optimality.

We need to introduce a measure Λ over the parameter space Ω .

This measure Λ can be viewed as an assignment of weights to each of the parameter values $\theta \in \Omega$ ^p ~~a priori~~. [i.e. before any data is observed].

\times fixed unknown constant

(H)

Given a measure Λ , our objective is to find an estimator δ_{Λ}

which minimises the average risk, which is given by

$$r(\Lambda, \delta) = \int R(\theta, \delta) d\Lambda(\theta) = E_{(H)}(R(\theta, \delta))$$

If Λ is a probability distribution on Ω , we call Λ the prior distribution. Correspondingly, the estimator δ_{Λ} , if exists, is called the Bayes estimator with respect to Λ , and the minimised average risk is called the Bayes risk.

$$\begin{aligned} r(\Lambda, \delta) &= E_{(x, \Theta)}(L(\Theta, \delta(x))) \\ &= E_{\Theta}(E_x(L(\Theta, \delta(x)) | \Theta)) \\ &= E_{\Theta}(R(\Theta, \delta)). \end{aligned}$$

↓ data.

We shall pay attention to $E(L(\Theta, \delta(x)) | X=x)$ the conditional risk at (almost) every value of X . Notice that the expectation here is taken with respect to the cond. distribution of Θ given X . $(\Theta | X=x)$.

[THEOREM]. Suppose $\Theta \sim \Lambda$ and $X | \Theta = \theta \sim P_\theta$. If

(a) there exists δ_0 , an estimator of $g(\theta)$ with finite risk for all θ , and

(b) there exists a value $\delta_{\Lambda}(x)$ that minimises

$$\rightarrow E(L(\Theta, \delta_{\Lambda}(x)) | X=x) \text{ for almost every } x,$$

then δ_A is a Bayes estimator with respect to Π .

Note: The almost sure statement (*) is defined with respect to the marginal distribution of X , which is given by

$$\underline{P(X \in A) = \int P_\theta(X \in A) d\Pi(\theta)}.$$

Proof. Under the assumptions of the theorem (a) & (b), for other estimator δ' , say, and for almost every x ,

$$E(L(\Theta, \delta_A(x)) | X=x) \leq E(L(\Theta, \delta'(x)) | X=x).$$

Taking expectation over X , we obtain

$$E(L(\Theta, \delta_A(X))) \leq E(L(\Theta, \delta'(X)))$$

for all δ' .

(Example). If we consider the squared loss function $L(\theta, d) = (\theta - d)^2$, to find the Bayes estimator, we need to minimise

$$E((g(\Theta) - \delta(X))^2 | X=x)$$

and in this case, the Bayes estimator is $\delta_A(x) = \underline{E(g(\Theta)|X)}$, the posterior mean of $g(\Theta)$ given $X=x$.

$$= E(fg(\Theta) - E(g(\Theta)|X) + E(g(\Theta)|X) - \delta(x))^2 | X=x$$

$$= E(fg(\Theta) - E(g(\Theta)|X))^2 | X=x$$

$$+ E(E(g(\Theta)|X) - \delta(x))^2 | X=x$$

↑
Posterior mean $E(\Theta|X=x)$

$$\text{posterior} = \frac{\text{joint}}{\text{marginal}} = \frac{\text{prior} \times \text{likelihood}}{\text{marginal.}}$$

$$\text{Bayes' Theorem} \rightarrow p(\theta|x) = \frac{p(\theta, x)}{\int p(\theta', x) d\theta'} = \frac{p(x|\theta) \pi(\theta)}{\int p(x|\theta') \pi(\theta') d\theta'}$$

posterior \propto prior \times likelihood

(Example) Suppose $X \sim \text{Binomial}(n, \theta)$ given $\Theta = \theta$ and that Θ has a prior distribution Beta (α, β) . The prior density is given by

$$\pi(\theta; \alpha, \beta) = \underbrace{\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}}_{\text{prior}} \theta^{\alpha-1} (1-\theta)^{\beta-1} I(0 < \theta < 1).$$

non-informative Conjugate prior
(共轭先验)
flat prior.

Obviously, the model density is $f(x|\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$, in which case the posterior distribution of Θ given X is

$$\pi(\theta|x) \propto \left\{ \binom{n}{x} \theta^x (1-\theta)^{n-x} \cdot \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \right\}^{k_{\text{new}}}$$

$$\frac{\text{likelihood}}{\int \text{NUMERATOR}(\theta) d\theta} \rightarrow \text{normalising constant.}$$

If $X \sim \text{Beta}(\alpha, \beta)$

$$E(X) = \frac{\alpha}{\alpha + \beta}.$$

$$\propto \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\sim \text{Beta}(x+\alpha, n-x+\beta), \quad \frac{x+\alpha}{x+\alpha+n-\beta}$$

$$\text{meaning that the posterior mean of } \theta | X = \frac{x+\alpha}{n+\alpha+\beta}.$$

Remark: The posterior mean can be rewritten as:

$$\rightarrow \frac{x+\alpha}{n+\alpha+\beta} = \underbrace{\frac{n}{n+\alpha+\beta} \left(\frac{x}{n} \right)}_{\omega} + \underbrace{\frac{\alpha}{n+\alpha+\beta} \left(\frac{\alpha}{\alpha+\beta} \right)}_{1-\omega}$$

Contribution from the prior.

Shrink the estimate towards $\alpha/(\alpha+\beta)$

"weighted average of the sample mean \bar{x}_n and the prior mean $\alpha/(\alpha+\beta)$ "

As $(n \rightarrow \infty)$ $E(\theta | X) \rightarrow \bar{x}_n$. (Let the data "speak for themselves") empirical evidence/observations...

(Example). (Normal Mean Estimation)

Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, \sigma^2)$, with σ^2 known. Furthermore,

let $\theta \sim N(\mu, b^2)$ where μ and b^2 are two fixed hyperparameters.

Then the posterior dist. of $\theta | X$ is

$$\pi(\theta | x) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(X_i - \theta)^2}{2\sigma^2}\right\} \cdot \frac{1}{\sqrt{2\pi b^2}} \exp\left\{-\frac{(\theta - \mu)^2}{2b^2}\right\}$$

likelihood prior

$$\propto \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \theta)^2 - \frac{1}{2b^2} (\theta - \mu)^2\right\}$$

$$= \exp\left\{\frac{1}{\sigma^2} \sum_{i=1}^n x_i \theta - \frac{n\theta^2}{2\sigma^2} - \frac{1}{2b^2} \theta^2 + \frac{\mu}{b^2} \theta\right\}$$

$$(a+b)^2 = a^2 + 2ab + b^2$$

$$\propto \exp\left\{-\frac{1}{2} \left(\frac{n}{\sigma^2} + \frac{1}{b^2} \right) \theta^2 + \left(\frac{n\bar{x}}{\sigma^2} + \frac{\mu}{b^2} \right) \theta\right\} \dots \exp\left(-\frac{1}{2\sigma^2} (\theta - \bar{x})^2\right)$$

The posterior distribution of θ given X is normal. with $(\tilde{\mu}, \tilde{\sigma}^2)$ as

$$\exp\left\{-\frac{1}{2\tilde{\sigma}^2} (\theta - \tilde{\mu})^2\right\}, \text{ where } \begin{cases} \tilde{\mu} = \frac{n\bar{x}/\sigma^2 + \mu/b^2}{n/\sigma^2 + 1/b^2} \\ \tilde{\sigma}^2 = \frac{1}{n/\sigma^2 + 1/b^2} \end{cases}.$$

Hence, the posterior mean of $\theta | X$ is $\frac{n\bar{x}/\sigma^2 + \mu/b^2}{n/\sigma^2 + 1/b^2}$

data.

$$\frac{n/\sigma^2}{\bar{x}} + \frac{1/b^2}{\mu} \text{ prior mean } [N(\mu, \infty)]$$

$$\begin{array}{c}
 \text{Bayes estimator} \\
 \delta_n \xrightarrow{\sim} \mu \text{ if we} \\
 \text{adopt the squared} \\
 \text{loss function...}
 \end{array}
 \left| \begin{array}{c}
 \frac{n/\sigma^2 + 1/b^2}{n/\sigma^2 + 1/b^2} \cdot \frac{n/\sigma^2 + 1/b^2}{n/\sigma^2 + 1/b^2} \\
 1 \text{ as } n \rightarrow \infty \\
 0 \text{ as } n \rightarrow \infty
 \end{array} \right| \quad \begin{array}{l}
 \text{prior: } n(\mu; \sigma) \\
 \text{Suppose } b^2 \rightarrow \infty \dots \\
 \delta_n \xrightarrow{n(\mu; \sigma)} \bar{x} + \frac{1/n}{n/\sigma^2 + 1/n} \\
 = \underline{\bar{x}}
 \end{array}$$

(Example). Assume that we consider $L(\theta, d) = \omega(\theta) \{d - g(\theta)\}^2$, where $\omega(\theta) \geq 0$, which can be interpreted as a weight function. Our goal is to find the corresponding Bayes estimator, which minimizes $E(\omega(\theta) \{g(\theta) - d\}^2 | X=x)$ with respect to d . $(*)$

$(*)$ can be rewritten as

$$d^2 E(\omega(\theta) | X=x) - 2d E(\omega(\theta) g(\theta) | X=x) + E(\omega(\theta) g(\theta)^2 | X=x) \quad (*)$$

Taking derivative of $(*)$ w.r.t. d , we obtain

$$\begin{aligned} & 2d^* E(\omega(\theta) | X=x) - 2E(\omega(\theta) g(\theta) | X=x) = 0 \\ \Rightarrow \quad & \delta(x) = d^* = \frac{E(\omega(\theta) g(\theta) | X=x)}{E(\omega(\theta) | X=x)}. \end{aligned}$$

In particular, if $\omega(\cdot) \equiv 1$, $\delta_n(x)$ (with $\omega(\cdot) \equiv 1$) = $E(g(\theta) | X=x)$.

[THEOREM]. (TPE 4.2.3). If δ is unbiased for $g(\theta)$ with $r(1, \delta) < \infty$

Can \bar{X}_n be
a Bayes estimator?
and $E(g(\theta)^2) < \infty$, then δ is not Bayes under the
Squared loss function unless its average risk is zero,
i.e. $E_{(\theta, \bar{X})} (\{\delta(x) - g(\theta)\}^2) = 0$.

Proof. let δ be an unbiased estimator under the squared loss function. Then, we know that δ is the posterior mean,

$$\text{i.e. } \delta(x) = E(g(\theta) | X=x) \text{ a.s..} \quad *$$

Thus, we have

$$\begin{aligned} E(\delta(x) g(\theta)) &= E(E(\delta(x) g(\theta) | X)) = E(\delta(x) E(g(\theta) | X)) \\ &\stackrel{*}{=} E(\delta^2(x)). \end{aligned} \quad (1)$$

Also,

$$\begin{aligned} E(\delta(x) g(\theta)) &= E(E(\delta(x) g(\theta) | \theta)) = E(g(\theta) E(\delta(x) | \theta)) \\ &\stackrel{\text{unbiasedness of } \delta}{=} E(g^2(\theta)). \end{aligned} \quad (2)$$

Observe that.

$$\begin{aligned} E(\{\delta(x) - g(\theta)\}^2) &= E(\delta^2(x)) - 2E(\delta(x) g(\theta)) + E(g^2(\theta)) \\ &= E(\delta^2(x)) - E(\delta(x) g(\theta)) \\ &\quad + E(g^2(\theta)) - E(\delta(x) g(\theta)) \end{aligned}$$

$$\begin{aligned}
 &= E(\delta^2(x)) - E(\delta^2(x)) \\
 &\quad + E(g^2(\theta)) - E(g^2(\theta)) \quad (\text{due to } ① \text{ and } ②) \\
 &= 0.
 \end{aligned}$$

Thus, we have that $E(f\delta(x) - g(\theta))^2 = 0$, i.e. the average risk is zero. The claim is thus proved. \square

(Example). Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\theta, \sigma^2)$, with σ^2 known. Is \bar{X} Bayes under the squared loss function for some choice of the prior distribution?

Observe that $E(\bar{X} | \theta) = \theta$, hence \bar{X} is unbiased for θ . The corresponding average risk under the squared loss function is given by

$$E_{(\bar{X}, \theta)}(\{\bar{X} - \theta\}^2) = \frac{\sigma^2}{n} \neq 0.$$

$\Rightarrow \bar{X}$ is not Bayes estimator under any prior distribution... \square

guaranteed admissibility \leftarrow [THEOREM] (TPE 6.2.4) A unique Bayes estimator (a.s. for all P_θ) is admissible.

An estimator is admissible if it is not uniformly dominated by some other estimator. δ is said to be inadmissible if and only if there exists δ' such that

$$\left\{ \begin{array}{l} R(\theta, \delta') \leq R(\theta, \delta) \quad \forall \theta \in \Omega \text{ and} \\ R(\theta, \delta') < R(\theta, \delta) \quad \text{for some } \theta \in \Omega. \end{array} \right\}$$

Proof. Suppose δ_1 is Bayes for Λ , and for some δ' , $R(\theta, \delta') \leq R(\theta, \delta_1)$ for all $\theta \in \Omega$. If we take expectations w.r.t Λ , the inequality above is preserved and we can write

$$\int_{\theta \in \Omega} R(\theta, \delta') d\Lambda(\theta) \leq \int_{\theta \in \Omega} R(\theta, \delta_1) d\Lambda(\theta)$$

This implies that δ' is also Bayes because δ' has less (or equal) risk than δ_1 which minimizes the average risk. Hence $\delta' = \delta_1$ with probability one for all P_θ .

[Question: When is a Bayes estimator unique?]

[THEOREM]. Let Ω be the marginal distribution of X , that is

$$\Omega(E) = \int P(X \in E | \theta) d\Lambda(\theta).$$

TPE 1-7.27
uniqueness of
the minimizer of
a strictly convex

Then, under a strictly convex loss function, δ_1 is unique (a.s. for all P_θ) if

- (proved by contradiction)
- * (a) $r(\Delta, \delta_\Delta)$ is finite and $\not\rightarrow$ finiteness for Chapman
 - * (b) $P_\theta < \underline{\alpha}$ (absolute continuity)
 - $P_\theta(X \in A)$ is a continuous function of θ for all meas. sets of A .
- Ω is defined on the support Δ .
- Benefits of Bayes:
- { ① Admissible
 - ② Incorporate prior information \rightarrow frequentist domain knowledge
 - ③ ...

Next lecture: Minimax Estimator.

$$\underbrace{\int R(\theta, \delta) d\Delta(\theta)}_{\text{average.}} X \rightarrow \boxed{\sup_{\theta \in \Omega} R(\theta, \delta)}$$

maximum
(most conservative)

Worst-case scenario/optimality

Testing of Statistical Hypotheses, (UMP)