

Estimation, and Decision Theory

PUBH 8442: Bayes Decision Theory and Data Analysis

Eric F. Lock
UMN Division of Biostatistics, SPH
elock@umn.edu

02/01/2021

- ▶ Full posterior distribution $p(\theta | \mathbf{y})$ is nice...
 - ▶ But often wish to simplify inference and conclusions
- ▶ Possible *point estimates* for θ :

- ▶ Posterior mode:

$$\hat{\theta} = \operatorname{argmax}_{\theta} p(\theta | \mathbf{y})$$

- ▶ Posterior expectation:

$$\hat{\theta} = E_{\theta | \mathbf{y}} \theta = \int \theta p(\theta | \mathbf{y}) d\theta$$

- ▶ Posterior median:

$$P(\theta \leq \hat{\theta} | \mathbf{y}) = \int_{-\infty}^{\hat{\theta}} p(\theta | \mathbf{y}) d\theta = \frac{1}{2}$$

- ▶ Posterior mode
 - ▶ Easy to compute because only need to work with numerator

$$p(\theta | \mathbf{y}) \propto p(\theta)p(\mathbf{y} | \theta)$$

- ▶ Sometimes called generalized maximum likelihood estimation
- ▶ Posterior expectation uses full posterior
- ▶ Posterior median is more robust to outliers in \mathbf{y} & posterior tails
- ▶ How to choose which estimate is optimal for a given application?

Loss function

- ▶ A *loss function* $l(\text{truth}, a)$ gives the loss incurred for an action a given the (usually unknown) truth.
 - ▶ Here “loss” is abstract - could be loss to society, loss in terms of model accuracy, etc.
 - ▶ want to minimize loss
- ▶ If the truth is given by parameter θ , the *posterior risk* is

$$\rho(p_{\theta}, \mathcal{A}) \quad E_{\theta | \mathbf{y}} l(\theta | a) = \int l(\theta | a) p(\theta | \mathbf{y}) d\theta$$

- ▶ Averaging the loss over the posterior for θ
- ▶ For point estimation, our action a is given by an estimate $\hat{\theta}$:
 $l(\theta, \hat{\theta})$

Point estimate loss

- Squared error loss is commonly used

$$l(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$$

- Posterior risk for squared error loss is minimized by posterior expectation:

$$\rho(\mathbb{P}_\theta, \hat{\theta})$$

$$E_{\theta|y} \theta = \operatorname{argmin}_{\hat{\theta}} E_{\theta|y} (\theta - \hat{\theta})^2$$

$$E(\theta - \hat{\theta})^2 = E(\underbrace{\theta - E\theta}_{\text{zero mean}} + \underbrace{E\theta - \hat{\theta}}_{\text{constant}})^2$$

$$= E[(\theta - E\theta)^2 + 2(\theta - E\theta)(E\theta - \hat{\theta}) + (E\theta - \hat{\theta})^2]$$

$$= E(\theta - E\theta)^2 + 0 + (E\theta - \hat{\theta})^2 \geq E(\theta - E\theta)^2 = \rho(\mathbb{P}_\theta, E\theta)$$

- Posterior risk for absolute loss $l(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$ is minimized by the posterior median
 - Homework

Decision rules

- ▶ Denote the space of allowable actions \mathcal{A} ($a \in \mathcal{A}$)
- ▶ Denote sample space (possible data observations) \mathcal{Y} ($\mathbf{y} \in \mathcal{Y}$)
- ▶ A *decision rule* $d \in \mathcal{D} : \mathcal{Y} \rightarrow \mathcal{A}$ is a rule for determining an action based on data.
- ▶ Formal decision-theoretic framework:
 - ▶ *prior distribution*: $p(\theta), \theta \in \Theta$
 - ▶ *sampling distribution*: $p(\mathbf{y} | \theta)$
 - ▶ *allowable actions*: $a \in \mathcal{A}$
 - ▶ *decision rules*: $d \in \mathcal{D} : \mathcal{Y} \rightarrow \mathcal{A}$
 - ▶ *loss function*: $l(\theta, a)$

Normal-normal point estimate

- ▶ Recall for y_1, \dots, y_n iid $\text{Normal}(\mu, \sigma^2)$ and $p(\mu) = \text{Normal}(\mu_0, \tau^2)$:

$$p(\mu | \mathbf{y}) = \text{Normal} \left(\frac{\sigma^2 \mu_0 + n\tau^2 \bar{y}}{\sigma^2 + n\tau^2}, \frac{\sigma^2 \tau^2}{\sigma^2 + n\tau^2} \right)$$

- ▶ Consider decision for point estimate of μ :
 - ▶ *prior distribution*: $p(\mu) = \text{Normal}(\mu_0, \tau^2)$, $\mu \in \mathbb{R}$
 - ▶ *sampling distribution*: y_1, \dots, y_n iid $\text{Normal}(\mu, \sigma^2)$, $\mathbf{y} \in \mathbb{R}^n$
 - ▶ *allowable actions*: $a \in \mathcal{A} = \mathbb{R}$
 - ▶ *decision rule*: $d(\mathbf{y}) = \frac{\sigma^2 \mu_0 + n\tau^2 \bar{y}}{\sigma^2 + n\tau^2}$
 - ▶ *loss function*: $l(\mu, a) = (\mu - a)^2$ (or $l(\mu, a) = |\mu - a|$)

Example: Coke bottles (cont.)

- Recall:
 - Coke bottles are filled with calibration $\text{Normal}(12, 0.01)$
 - Given machine with calibration μ , bottles filled with $\text{Normal}(\mu, 0.05)$
 - For $n = 5$ and $\bar{y} = 11.88$ oz, $p(\mu | \mathbf{y}) = \text{Normal}(11.94, 0.005)$

- Action is to estimate $\hat{\mu} = 11.94$



- Posterior risk under squared error loss is the posterior variance, 0.005

$$\begin{aligned} & E_{\mu|y} L(\mu, \hat{\mu}) \\ &= \int (\mu - \hat{\mu})^2 p(\mu|y) d\mu \\ &= \text{Var}_{\mu|y}(\mu) \end{aligned}$$

- ▶ The *frequentist risk* of a decision rule d is

$$\begin{aligned}R(\theta, d) &= E_{\mathbf{y} | \theta} l(\theta, d(\mathbf{y})) \\ &= \int l(\theta, d(\mathbf{y})) p(\mathbf{y} | \theta) d\mathbf{y}\end{aligned}$$

- ▶ Loss averaged over \mathbf{y} , given θ .
 - ▶ Note: does not depend on prior $p(\theta)$
- ▶ A rule d is *inadmissible* if $\exists d^*$ with

$$R(\theta, d^*) \leq R(\theta, d) \quad \forall \theta \in \Theta$$

and $<$ for some $\theta \in \Theta$

- ▶ Implies another rule is universally “better”
- ▶ A rule that is not inadmissible is *admissible*
- ▶ Admissible rules are not necessarily good

Example: Coke bottles (cont.)

- ▶ The (poor) rule $d(\mathbf{y}) = 5$ is admissible because it is unbeatable when $\mu = 5$! $R(\mu, 5) = (\mu - 5)^2$
- ▶ The rule $d(\mathbf{y}) = \bar{y}$ has frequentist risk

$$\begin{aligned} R(\mu, \bar{y}) &= \text{Var}_{\mathbf{y} | \mu} \bar{y} \\ E_{\mathbf{y} | \mu} (\mu - \bar{y})^2 &= \frac{\sigma^2}{n} \\ &= 0.05/5 = 0.01 \end{aligned}$$

- ▶ Does not depend on μ

Example: Coke bottles (cont.)

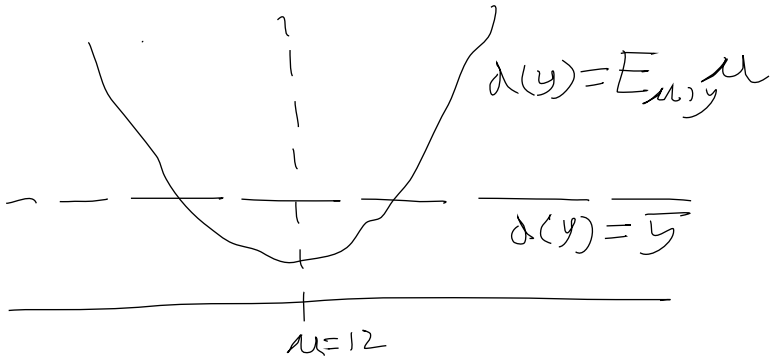
- Note: posterior mean has form $B\mu_0 + (1 - B)\bar{y}$, where

$$B = \frac{\sigma^2}{\sigma^2 + n\tau^2}.$$

- For coke example with $n = 5$, $B = 0.5$
- The posterior mean rule $d(\mathbf{y}) = E_{\mu|\mathbf{y}} \mu$ has frequentist risk

$$\begin{aligned} R(\mu, E_{\mu|\mathbf{y}} \mu) &= B^2(\mu - \mu_0)^2 + (1 - B)^2 \text{Var}_{\mathbf{y}|\mu} \bar{y} \\ &= 0.25(\mu - 12)^2 + 0.0025 \end{aligned}$$

$$\begin{aligned} &\rightarrow E_{\mathbf{y}|\mu} (\mu - B\mu_0 - (1-B)\bar{y})^2 \\ &= E (B(\mu - \mu_0) + (1-B)(\mu - \bar{y}))^2 \\ &= B^2(\mu - \mu_0)^2 + (1-B)^2 E (\mu - \bar{y})^2 \\ &\quad + 0 \end{aligned}$$



- ▶ A *minimax rule* d satisfies

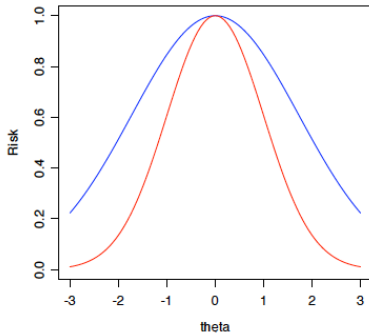
$$\sup_{\theta \in \Theta} R(\theta, d) \leq \sup_{\theta \in \Theta} R(\theta, d^*), \quad \forall d^* \in \mathcal{D}.$$

- ▶ Chose d to minimize maximum risk.
 - ▶ Prepare for “the worst case scenario”
- ▶ May not be unique, admissible, or favorable.

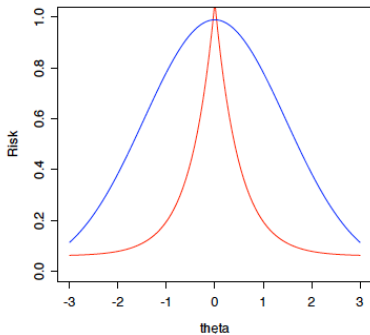
Minimax rules



Inadmissible minimax rule



Counterintuitive minimax rule



Credit: Victor Panaretos

- ▶ The *Bayes risk* for given decision rule, for prior p_θ , is

$$r(p_\theta, d) = E_\theta E_{\mathbf{y}|\theta} l(\theta, d(\mathbf{y})) = \int R(\theta, d) p(\theta) d\theta$$

- ▶ Equivalently,

$$r(p_\theta, d) = E_{\mathbf{y}} E_{\theta|\mathbf{y}} l(\theta, d(\mathbf{y})) = \int \rho(p_\theta, d(\mathbf{y})) p(\mathbf{y}) d\mathbf{y}$$

where $p(\mathbf{y})$ is the marginal distribution of \mathbf{y} .

- ▶ Also called “preposterior risk”
 - ▶ The expected risk before obtaining any data

Risks galore!

- ▶ *Loss function* $l(\theta, d(\mathbf{y}))$
 - ▶ Function of θ and \mathbf{y}
- ▶ *Posterior risk* $\rho(p_\theta, d(\mathbf{y}))$
 - ▶ Function of \mathbf{y} , averaged over θ
- ▶ *Frequentist risk* $R(\theta, d)$
 - ▶ Function of θ , averaged over \mathbf{y}
- ▶ *Bayes risk* $r(p_\theta, d)$
 - ▶ Averaged over both \mathbf{y} and θ

Bayes decision rules

- A *Bayes decision rule* minimizes Bayes risk:

$$\operatorname{argmin}_{d \in \mathcal{D}} r(p_\theta, d)$$

- Bayes rules are generally admissible
 - If a Bayes rule is unique, it is admissible.

If d is inadmissible,

$$\exists d^* \text{ s.t. } R(\theta, d^*) \leq R(\theta, d) \quad \forall \theta$$

$$\rightarrow R(\theta, d^*) P(\theta) \leq R(\theta, d) P(\theta)$$

$$\rightarrow \int R(\theta, d^*) P(\theta) \leq \int R(\theta, d) P(\theta)$$

$$\rightarrow r(p_\theta, d^*) \leq r(p_\theta, d), \quad \therefore d \text{ not a unique BR}$$

Normal-normal model

- ▶ For the normal-normal model with $l(\mu, \hat{\mu}) = (\mu - \hat{\mu})^2$, we've shown

$$\operatorname{argmin}_{d \in \mathcal{D}} \rho(p_\theta, d(\mathbf{y}))$$

is given by the posterior mean for any \mathbf{y}

- ▶ So, it is a Bayes decision rule.
- ▶ The Bayes risk is given by

$$r(p_\theta, d) = \frac{\sigma^2 \tau^2}{\sigma^2 + n\tau^2}$$

$\int \rho(\mu_n, d(\mathbf{y})) p(\mathbf{y}) d\mathbf{y}$

$\int \frac{\sigma^2 \tau^2}{\sigma^2 + n\tau^2} p(\mathbf{y}) d\mathbf{y} =$

Example: Coke bottles (cont.)

- ▶ For the Coke bottling example,

$$r(p_\theta, d) = 0.005$$

- ▶ This is expected loss before checking any bottles.
- ▶ Recall expected risk after collecting bottles (posterior risk) was also 0.005
 - ▶ Equivalent in this case, because posterior risk does not depend on \mathbf{y} .
- ▶ Bayes risk of $d(\mathbf{y}) = \bar{y}$ is 0.01 – twice as large.

$$\int R(\mu, d) P(\mu) = \frac{\sigma^2}{n} \underbrace{\int p(\mu) d\mu}_1 = 0.01$$