

01/25/22

- HW #2 - Q1 : mixture - skip until we cover Ex 3.4.1

- Q10 : $f(x = \theta - 1 \mid \theta) = f(x = \theta + 1 \mid \theta) = \frac{1}{2}$

x_1, x_2

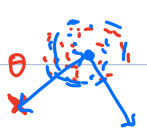
$\pi(\theta)$: consider an arbitrary distribution over integers.

- Lecture Capture - as soon as possible

JB Example 4(p10) Assume $X \mid \theta \sim N(\theta, 1)$. The goal is estimating θ under **squared-error loss**, $L(\theta, d) = (\theta - d)^2$. Consider the decision rule $\delta_c(x) = \text{c}x$. (SEL)

- Find $R(\theta, \delta_c)$.

$$\begin{aligned}
 R(\theta, \delta_c) &= E_{\theta}(L(\theta, \delta_c)) \\
 &= E_{\theta}((\theta - \delta_c(x))^2) \\
 &= E_{\theta} \left[\left\{ \theta - E_{\theta}(\delta_c(x)) - \delta_c(x) \right\}^2 \right] \\
 &= E_{\theta} \left[\left\{ \theta - E_{\theta}(\delta_c(x)) \right\}^2 \right. \\
 &\quad \left. + 2 \left\{ \theta - E_{\theta}(\delta_c(x)) \right\} \left\{ E_{\theta}(\delta_c(x)) - \delta_c(x) \right\} \right. \\
 &\quad \left. + \left\{ E_{\theta}(\delta_c(x)) - \delta_c(x) \right\}^2 \right]
 \end{aligned}$$



$$= \underbrace{\{\theta - E_\theta(\delta_c(x))\}^2}_{= \text{bias}} + E_\theta \left[\{E_\theta(\delta_c(x)) - \delta_c(x)\}^2 \right]$$

$$\underbrace{\delta_c(x)}_{\text{Estimator (decision rule)}} = \{\theta - E_\theta(\delta_c(x))\}^2 + \underbrace{\text{Var}(\delta_c(x))}$$

Estimator (decision rule) $\delta_c(x) = cX$

$$E_\theta(\delta_c(x)) = E_\theta(cX) = c\theta$$

$$\begin{aligned} \Rightarrow R(\theta, \delta_c) &= (\theta - c\theta)^2 + \text{Var}(cX) \\ &= (1-c)^2 \theta^2 + c^2 \underbrace{\text{Var}(X)}_{=1} \\ &= (1-c)^2 \theta^2 + c^2 \end{aligned}$$

eg)	$c = 1/2$	$\Rightarrow R(\theta, \delta_{1/2}) = \frac{1}{4}\theta^2 + \frac{1}{4} \leftarrow$
$\delta_1(x) = x$	$c = 1$	$\Rightarrow R(\theta, \delta_1) = 1 \leftarrow$
$E_\theta(\delta_2(x)) = E_\theta(x) = \theta$	$c = 2$	$\Rightarrow R(\theta, \delta_2) = \theta^2 + 4 \leftarrow$

eg) $X|\theta \sim N(\theta, 1)$, assume $\theta \sim N(0, \tau^2)$

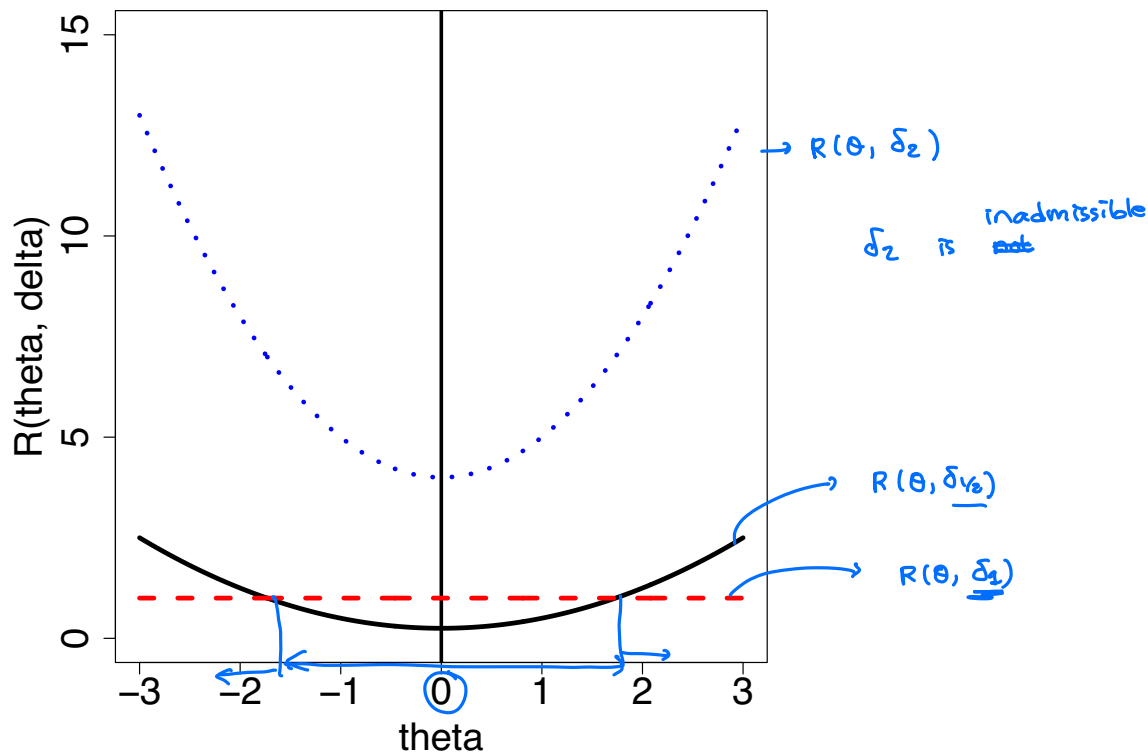
$$\Rightarrow \theta|X \sim N \left(\underbrace{\left(\frac{1}{1} + \frac{1}{\tau^2} \right)^{-1} \left(\frac{X}{1} + \frac{0}{\tau^2} \right)}_{= \left(\frac{\tau^2}{1+\tau^2} \right) X}, \left(\frac{1}{1} + \frac{1}{\tau^2} \right)^{-1} \right)$$

Under the squared error loss function,

$$d = E^\theta(\theta|X) = \underbrace{\left(\frac{\tau^2}{1+\tau^2} \right)}_{\neq 1} X$$

$$\begin{aligned} &E_\theta \left[\underbrace{\{\theta - E_\theta(\delta_c(x))\}}_{=0} \underbrace{\{E_\theta(\delta_c(x)) - \delta_c(x)\}}_{=0} \right] \\ &= \underbrace{\{\theta - E_\theta(\delta_c(x))\}}_{=0} \times E_\theta \{E_\theta(\delta_c(x)) - \delta_c(x)\} \\ &= \underbrace{\{E_\theta(\delta_c(x)) - E_\theta(\delta_c(x))\}}_{=0} = 0 \end{aligned}$$

JB Example 4(p10) (contd) Plot of $R(\theta, \delta_c)$



- Difficulties associated with using $R(\theta, \delta)$.
 - ★★ For each $\theta \in \Theta$, $R(\theta, \delta)$ is the expected loss based on an average over the random $X \in \mathcal{X}$.
 - \Rightarrow *long-run* performance of $\delta(x)$ and **not** directly for the given observation x .
 - ★★ A function of $\theta \in \Theta$ & θ is unknown.
 - \Rightarrow The frequentist approach $R(\theta, \delta)$ does not induce a total ordering on the set of procedures.

† How can Frequentists choose δ ?

- An additional principle must be introduced to select a specific rule for use.
e.g. δ_1 is preferred to δ_3 under some concept of optimality.
- Some important frequentist decision principles (CR 2.4 + a lot in JB)
 - ★★ Bayes risk principle
 - ★★ minimax
 - ★★ admissibility
 - ★★ restricted classes: e.g. we only consider unbiased estimators.
- Bayes estimators are *often optimal* for the frequentist concepts of optimality.

† The Bayes Risk Principle

- The frequentist risk of a decision rule $\delta(x)$ is a function of θ .

$$R(\theta, \delta) = \int_{\mathcal{X}} L(\theta, \delta(x)) f(x | \theta) dx.$$

- The *integrated risk* (also called Bayes Risk) is the frequentist risk averaged over Θ according to their prior $\pi(\theta)$.

$$\begin{aligned} \underbrace{r(\pi, \delta)}_{\text{Bayes Risk}} &= E^{\pi} [R(\theta, \delta)] \\ &= \int_{\Theta} \underbrace{\int_{\mathcal{X}} L(\theta, d) f(x | \theta) dx}_{= R(\theta, \delta)} \pi(\theta) d\theta. \end{aligned}$$

- $r(\pi, \delta)$ is a real number associated with estimator δ .

⇒ Induces a total ordering on the set of estimators, so allows for the direct comparison of estimators.

† Any connection between $r(\pi, \delta)$ and $\rho(\pi, \delta | x)$?

⇒ They lead to the same decision.

- **Th 2.3.2** An estimator minimizing the integrated risk $r(\pi, \delta)$ can be obtained by selecting, for every $x \in \mathcal{X}$, the value $\delta(x)$ which minimizes the posterior expected loss, $\rho(\pi, \delta | x)$, since

$$r(\pi, \delta) = \int_{\mathcal{X}} \rho(\theta, \delta(x) | x) m(x) dx.$$

$$\begin{aligned} \underline{r(\pi, \delta)} &= E^{\pi} (R(\theta, \delta)) \\ &= \int_{(\mathcal{H})} R(\theta, \delta) \pi(\theta) d\theta \\ &= \int_{(\mathcal{H})} \int_{\mathcal{X}} L(\theta, \delta) \underbrace{f(x|\theta) \pi(\theta)}_{= \pi(\theta|x) \cdot m(x)} dx d\theta \\ &= \int_{\mathcal{X}} \underbrace{\int_{(\mathcal{H})} L(\theta, \delta) \pi(\theta|x) d\theta}_{= \underline{\rho(\pi, \delta | x)}} m(x) dx \end{aligned}$$

- Def 2.3.3

- ★★ **A Bayes estimator** associated with a prior distribution π and a loss function L is any estimator δ^π , which minimizes $r(\pi, \delta)$.
 - ★★ For every $x \in \mathcal{X}$, it is given by $\delta^\pi(x)$ (**a Bayes action**), argument of $\min_d \rho(\pi, d \mid x)$.
 - ★★ The value $r(\pi) = r(\pi, \delta^\pi)$ is then called **the Bayes risk**.
- **JB Def 9, p160** If π is an improper prior, but $\delta^\pi(x)$ is an action which minimizes $\rho(\pi, d \mid x)$ for each x with $m(x) > 0$, then δ^π is called a **generalized Bayes rule**.

† Minimaxity: Minimize the expected loss in the least favorable case (\Leftrightarrow protect against the worst possible state of nature, conservative!)

- **The Minimax Principle.** JB p18 δ_1 is preferred to δ_2 if

$$\sup_{\theta} R(\theta, \delta_1) < \sup_{\theta} R(\theta, \delta_2).$$

- **Def 2.4.3** The **minimax risk** associated with a loss function L is the value

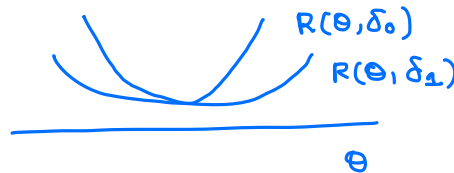
$$\bar{R} = \inf_{\delta \in \mathcal{D}} \sup_{\theta} R(\theta, \delta) = \inf_{\delta \in \mathcal{D}} \sup_{\theta} E_{\theta} \{L(\theta, \delta(x))\},$$

and a **minimax estimator** is any (possibly randomized) estimator δ_0 such that

$$\sup_{\theta} R(\theta, \delta_0) = \bar{R}$$

JB Example 4 (contd) Assume $X \mid \theta \sim N(\theta, 1)$. The goal is estimating θ under squared-error loss, $L(\theta, a) = (\theta - a)^2$. Consider the decision rule $\delta_c(x) = cx$. Find the minimax rule.

† Admissibility



- **Def 2.4.19** An estimator δ_0 is inadmissible if there exists an estimator δ_1 which dominates δ_0 , that is, such that, for every θ

$$R(\theta, \delta_0) \geq R(\theta, \delta)$$

and, for at least one value θ_0 of the parameter,

$$R(\theta_0, \delta_0) > R(\theta_0, \delta).$$

Otherwise, δ_0 is said to be admissible.

- What is the underlying idea as a criterion?

Inadmissible estimators should not be considered at all since they can be uniformly improved!

JB Example 4 (contd) Assume $X \mid \theta \sim N(\theta, 1)$. The goal is estimating θ under squared-error loss, $L(\theta, a) = (\theta - a)^2$. Consider the decision rule $\delta_c(x) = cx$ with $c > 1$. Is the rule admissible?

* Admissibility is related (stronger than minmax) to the Bayesian paradigm.

- Admissibility is automatically satisfied by most Bayes estimators.
- **Prop 2.4.22** If a prior distribution π is strictly positive on Θ , with finite Bayes risk and the risk function, $R(\theta, \delta)$, is a continuous function of θ for every δ , the Bayes estimator δ^π is admissible.
- Want to learn more? *Read CR 2 and JB 4.8*

Example 2.4.6 (Stein Phenomenon) Suppose a p -dimensional vector, $\mathbf{X} \sim N_p(\boldsymbol{\theta}, I_p)$ and consider the problem of estimating $\boldsymbol{\theta}$ (a p -dim vector). Assume the quadratic loss function

$$L(\boldsymbol{\theta}, \boldsymbol{\delta}) = (\boldsymbol{\theta} - \boldsymbol{\delta})'(\boldsymbol{\theta} - \boldsymbol{\delta}) = \|\boldsymbol{\theta} - \boldsymbol{\delta}\|^2 = \sum_{i=1}^p (\theta_i - \delta_i)^2$$

$\delta_c(x) = cx$
 $\delta_1(x) = x$

- The maximum likelihood estimator $\delta_1(\mathbf{X}) = \mathbf{X}$

★★ The least squares estimator in standard regression setting

★★ For $p = \underline{1}$ or 2, it is admissible and the unique minimax estimator.

$$\delta_{JS} = \left(1 - \frac{(p-2)}{\|\mathbf{X}\|^2} \right) \mathbf{X}$$

$$R(\boldsymbol{\theta}, \delta_{JS}) \leq R(\boldsymbol{\theta}, \delta_1) \text{ for all } \boldsymbol{\theta}$$

$$R(\boldsymbol{\theta}, \delta_{JS}) < R(\boldsymbol{\theta}, \delta_1) \text{ for some } \boldsymbol{\theta}$$

\Rightarrow dominates $\delta_1 = \mathbf{X}$ for $p \geq 3$

positive part JS

$$\delta_c^+ = \begin{cases} \left(1 - \frac{c}{\|\mathbf{X}\|^2} \right) \mathbf{X} & \text{if } \|\mathbf{X}\|^2 > c \\ 0 & \text{o.w.} \end{cases}$$

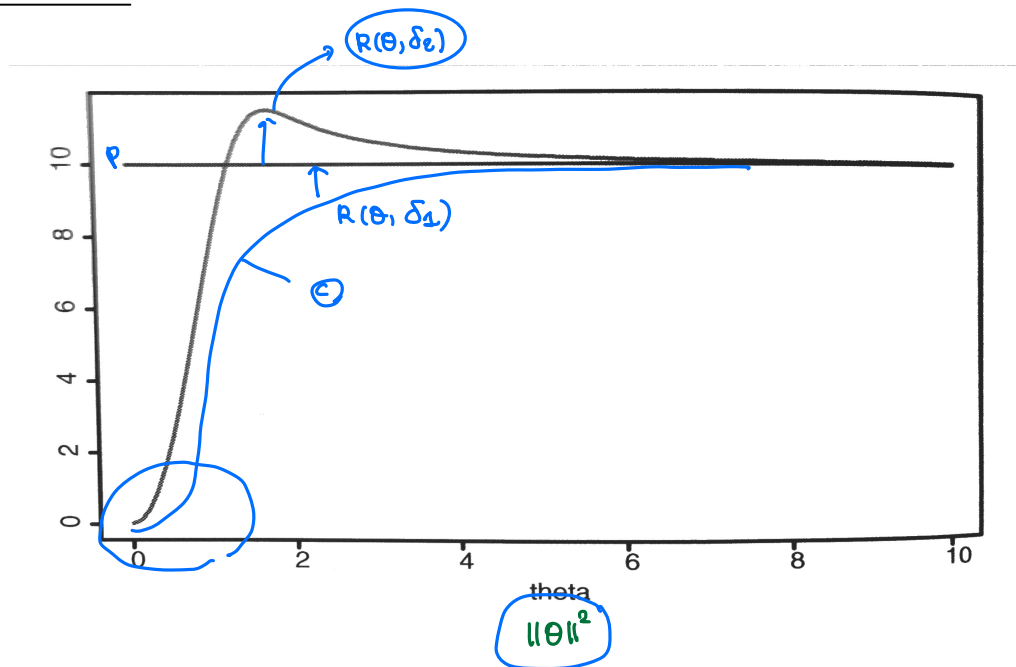
$$p-2 \leq c \leq 2 \cdot (p-2) \leq 2p-4 < \underline{\underline{2p-1}}$$

Example 2.4.6 (contd)

- Consider the positive part James-Stein estimator,

$$\delta_2(\mathbf{X}) = \begin{cases} \left(1 - \frac{2p-1}{\|\mathbf{X}\|^2}\right) \mathbf{X} & \text{if } \|\mathbf{X}\|^2 \geq 2p-1, \\ \mathbf{0} & \text{ow.} \end{cases}$$

♣ Figure 2.4.1 Comparison of the risks of δ_1 and δ_2 for $p = 10$



Example 2.4.6 (contd)

- ★★ δ_2 cannot be minimax.
- ★★ δ_2 is definitely superior on some (the most interesting) part of the parameter space.
- ★★ “The Stein effect”: allows to borrow information from the other components, even when they are independent and deal with totally different estimation problems.
- ★★ Sometimes the minimax rule is not useful! (or sometimes may not exist)
- ★★ Following James and Stein, extensive research on this has been done – *Shrinkage estimators*

† Usual loss functions (CR Section 2.5)

- Quadratic loss $\leftarrow E(\theta(x))$

- Absolute loss \leftarrow posterior median

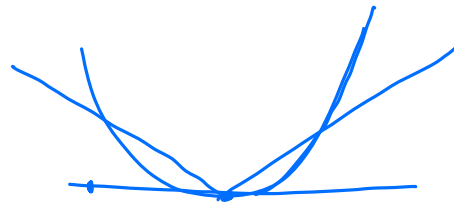
- 0-1 loss \leftarrow posterior mode

- Intrinsic loss – entropy distance (Kullback-Leibler divergence), Hellinger loss...

★★ What are the Bayesian estimators $\delta^\pi(x)$ under the classical loss functions?

† The quadratic loss

$$L(\theta, d) = (\theta - d)^2$$



- most common evaluation criterion – simplicity
- penalize large deviations too heavily

- **Prop 2.5.1** The Bayes estimator δ^π associated with the prior distribution π and with the quadratic loss $L(\theta, d) = (\theta - d)^2$, is the posterior expectation,

$$\delta^\pi(x) = \mathbb{E}^\pi(\theta \mid x) = \frac{\int_{\Theta} \theta f(x \mid \theta) \pi(\theta) d\theta}{\int_{\Theta} f(x \mid \theta) \pi(\theta) d\theta}.$$

- **Cor 2.5.2** The Bayes estimator δ^π associated with π and with the weighted quadratic loss $L(\theta, d) = \underbrace{w(\theta)}_{\geq 0}(\theta - d)^2$, where $w(\theta)$ is a nonnegative function, is

$$\delta^\pi(x) = \frac{\mathbb{E}^\pi(w(\theta) \cdot \theta \mid x)}{\mathbb{E}^\pi(w(\theta) \mid x)}$$

- **Cor 2.5.3** When $\Theta \in \mathbb{R}^p$, the Bayes estimator δ^π associated with the prior distribution π and with the quadratic loss $L(\theta, d) = (\theta - d)^t Q (\theta - d)$, is the posterior mean $\delta^\pi(x) = \mathbb{E}^{\theta \mid x}(\theta \mid x)$, for every positive -definite symmetric $p \times p$ matrix Q . $\mathbb{E}^\pi(\theta \mid x)$

† The absolute error loss

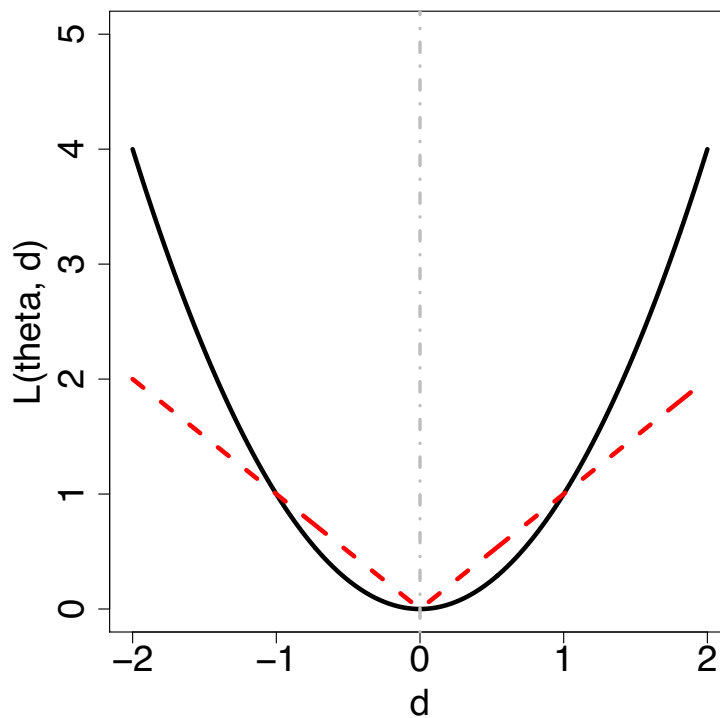
$$L(\theta, d) = |\theta - d|$$

Or multilinear function (more general)

$$L(\theta, d) = \begin{cases} k_1(\theta - d) & \text{if } \theta - d \geq 0, \\ k_2(d - \theta) & \text{if } \theta - d < 0, \end{cases}$$

- slow down the progression of the quadratic loss for large errors and has a robustifying effect.
- k_1 and k_2 reflect the relative importance of underestimation or overestimation.
- $k_1 = k_2 \Rightarrow$ the absolute error loss

♣ squared error loss vs absolute error loss



- **Prop 2.5.5** A Bayes estimator associated with the prior distribution π and the multilinear loss is a $k_1/(k_1 + k_2)$ fractile of $\pi(\theta \mid x)$.
- If $k_1 = k_2$, the Bayes estimator is the posterior median.

† The 0-1 loss: the penalty associated with an estimator δ is 0 if the answer is correct and 1 otherwise.

Example 2.5.6 Consider the test of $H_0 : \theta \in \Theta_0$ vs $H_1 : \theta \notin \Theta_0$. Then $\mathcal{D} = \{0, 1\}$, where 1 stands for acceptance of H_0 and 0 for rejection. For the 0-1 loss,

if $\theta \in \Theta_0$ & $d=1$, 0
 H_0 is true

$$\rightarrow L(\theta, d) = \begin{cases} 1 - d & \text{if } \theta \in \Theta_0, \\ d & \text{otherwise.} \end{cases}$$

• posterior expected loss

case 1: $d=1$

$$\begin{aligned} p(\pi, 1 | x) &= E^\pi(L(\theta, 1) | x) \\ &= \int_{\Theta_0} L(\theta, 1) \pi(\theta | x) d\theta \end{aligned}$$

$$= \underbrace{\int_{\mathbb{H}_0} \overset{=0}{L(\theta, 1)} \pi(\theta | x) d\theta}_{=0} + \int_{\mathbb{H}_0^c} \underbrace{L(\theta, 1)}_{=1} \pi(\theta | x) d\theta$$

$$= \int_{\mathbb{H}_0^c} \pi(\theta | x) d\theta$$

$$= \underline{P_r(\theta \in \mathbb{H}_0^c | x)}$$

Case 2 : $d=0$

$$p(\pi, 0 | x) = \underline{P_r(\theta \in \mathbb{H}_0 | x)}$$

① change $d=1$ if $P_r(\theta \in \mathbb{H}_0^c | x) \leq P_r(\theta \in \mathbb{H}_0 | x)$
 $d=0$ o.w

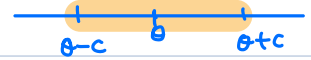
- **Prop 2.5.7** The Bayes estimator associated with the prior distribution π and with the 0-1 loss is

$$\delta^\pi(x) = \begin{cases} 1 & \text{if } \Pr(\theta \in \Theta_0 \mid x) > \Pr(\theta \notin \Theta_0 \mid x), \\ 0 & \text{otherwise,} \end{cases}$$

i.e., $\delta^\pi(x)$ is equal to 1 if and only if $\Pr(\theta \in \Theta \mid x) > 1/2$.

Consider the 0-1 loss function for the estimation problem.

$$L(\theta, d) = \begin{cases} 0 & \text{if } |\theta - d| \leq c \\ 1 & \text{o.w.} \end{cases}$$

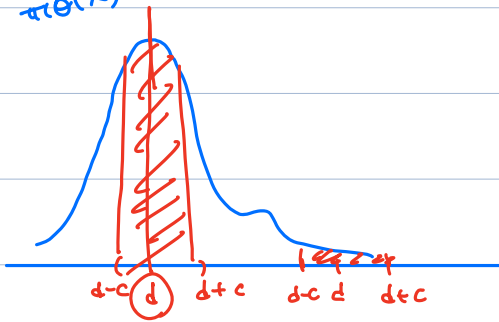


$$p(\pi, d | x) = \int_{\mathbb{R}} L(\theta, d) \pi(\theta | x) d\theta$$

$$= \Pr(\theta \notin (d-c, d+c) | x)$$



$\pi(\theta | x)$



$c \rightarrow 0$, $\underline{d} =$ posterior mode

$=$ MAP
(maximum a posteriori estimator)