

01/20/22

- ① in-person instruction from Feb/01
- ② office hours will be remote
- ③ HW#1 solution will be posted

normal example w/ a proper prior for unknown  $\theta$  &  $\sigma^2$

## † Maximum Likelihood Estimator (MLE)

- The maximum likelihood estimation approach is a way to implement the likelihood principle.

$\ell(\hat{\theta}^{\text{MLE}} \mid x)$  is at least as great as  $\ell(\theta \mid x)$  for every  $\theta \in \Theta$ .

- A lot more in STAT 205(B)!

## † Conditionality Principle

- **Def** If two experiments on the parameter  $\theta$ ,  $\mathcal{E}_1$  and  $\mathcal{E}_2$ , are available and if one of these two experiments is selected with probability  $p$ , the resulting inference on  $\theta$  should only depend on the selected experiment.
- Any inference should depend only on the outcome observed and not on any other outcome we might have observed.
- This sharply contrasts with a frequentist approach (which cares long-run behavior over repeated experiments). e.g. unbiasedness, significance levels, and power of tests, etc., violate the conditionality principle.

## Example 14 (JB-p25)

- Suppose a substance to be analyzed can be sent either to a laboratory in New York or a laboratory in CA. The two labs seem equally good, so a fair coin is flipped to choose between them, with “heads” denoting that the lab in NY will be chosen. The coin is flipped and comes up tails, so the CA lab is used.
- After a while, the experimental results come back and a conclusion and report must be developed. Should this conclusion take into account the fact that the coin could have been heads and hence that the experiment in NY might have been performed instead?
- Common sense and conditional view point say **NO**, but the frequentist approach calls for averaging over all possible data, even the possible NY data.

- Example 1.3.7:** In research laboratory, a physical quantity  $\theta$  can be measured by a precise but often busy machine, which provides a measurement,  $x_1 \sim N(\theta, 0.1)$ , with probability  $p = 0.5$ , or through a less precise but always available machine, which gives  $x_2 \sim N(\theta, 10)$ . The machine being selected at random, depending on the availability of the more precise machine, the inference on  $\theta$  when it has been selected should not depend on that fact that the alternative machine *could have been selected*.

$$x_1 \sim N(\theta, 0.1) \quad \text{w/p } 0.5$$

$$x_2 \sim N(\theta, 10) \quad \text{w/p } 0.5$$

conditionality principle

Suppose we know the precise machine was used,

$$x_1 \pm 1.96 \cdot \sqrt{0.1}$$

Otherwise,  $x_2 \pm 1.96 \times \sqrt{10}$ .

$$x=0$$

$$\longrightarrow (-5.19, 5.19)$$

$$\begin{aligned}
 & \text{Find a st. } \underbrace{P_r(X < a)}_{= 0.025} \\
 & = \underbrace{P_r(P)}_{= 0.5} \cdot \underbrace{P_r(X < a | P)}_{N(\theta, 0.1)} + \underbrace{(1 - P_r(P))}_{= 0.5} \cdot \underbrace{P_r(X < a | NP)}_{N(\theta, 10)} \\
 & \Rightarrow a = -5.19
 \end{aligned}$$

- **Theorem 1.3.8** (Birnbaum, 1962) The Likelihood Principle is equivalent to the conjunction of the Sufficiency and the Conditionality Principles.

*See page 18 for the proof.*

- Read Sec 1.3.4 for more about the likelihood principle.

# STAT 206B

## Chapter 2: Decision-Theoretic Foundations

Winter 2022

## † Statistical Decision Theory

- CR Chapter 2 & JB Chapters 1 & 2.
- Decision theory deals the problem of making decisions
- Statistical decision theory: Making decisions in the presence of statistical knowledge (statistical knowledge explains some of the uncertainties involved in the decision problem)



## JB Example 1 (page 5)

- Consider the situation of a drug company deciding whether or not to market a new pain killer. Two of the many factors affecting its decision are
  - ★★ the proportion of people for which the drug will prove effective ( $\theta_1$ )
  - ★★ the proportion of market the drug will capture ( $\theta_2$ )
- *Examples of decision problems:* estimate  $\theta_1$  &  $\theta_2$ , decide whether or not to market the drug, how much to market, what price to charge, etc.
- $\theta_1$  and  $\theta_2$  are unknown  $\Rightarrow$  conduct experiments to obtain statistical information about them.
- This is a problem of statistical decision theory!

## JB Example 1 (page 5) (contd)

- Consider the problem of estimating  $\theta_2$  (the proportion of market the drug will capture).
  - ★★ Let me use  $\theta$ , not  $\theta_2$  from now on.
  - ⇒ Parameter space:  $\theta \in \Theta = [0, 1]$ .
- Goal: Estimating  $\theta$   $\Leftrightarrow$  Choosing a number from interval  $[0, 1]$ 
  - ⇒ Your decision  $d$  will be a number  $\in \mathcal{D} = [0, 1]$ .
- Decisions are more commonly called “actions”.

## † Action and Action Space (Decision and Decision Space)

- $d \in \mathcal{D}$ :  $d$  denotes an action (decision) and  $\mathcal{D}$  the set of all possible actions under consideration (action space, decision space).

e.g. Problem of estimating  $\theta$ :

$$\underline{\mathcal{D}} = \Theta \text{ and } d \in \Theta = \mathcal{D}.$$

e.g. Testing problem:

$$\mathcal{D} = \{accept, reject\}.$$

- In the JB example,  $\Theta = \mathcal{D}$  (true for an estimation problem, but not necessarily for other problems).

## † Loss Function

$\theta$ : true state

$L(\theta, d)$

$d$ : decision

$d = \delta(x)$

- Consider the standard estimation problem, i.e.,  $\mathcal{D} = \Theta$ .
- Def 2.1.1** A loss function is any function  $L$  from  $\Theta \times \mathcal{D}$  in  $[0, +\infty)$ .

★★ The loss function evaluates the penalty (or error)  $L(\theta, d)$  associated with the decision (action)  $d$  when the parameter takes the value  $\theta$  for all  $(\theta, d) \in \Theta \times \mathcal{D}$ .

- Utility( $U$ ) and Loss ( $L$ )

gain

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

★★ Read CR Section 2.2 and JB Chapter 2 for details.

- Will discuss usual loss functions (Section 2.5).

• **Example 2.1.2:** Consider the problem of estimating the mean  $\theta$  of a normal vector,  $x \mid \theta \sim N_n(\theta, \Sigma)$ , where  $\Sigma$  is a known diagonal matrix with diagonal elements  $\sigma_i^2$ ,  $i = 1, \dots, n$ .

★★  $\mathcal{D} = \Theta = \mathbb{R}^n$

★★ Consider

$$L(\theta, \delta) = \sum_{i=1}^n \left( \frac{\delta_i - \theta_i}{\sigma_i} \right)^2,$$

where  $\delta_i$ : an estimator of  $\theta_i$  (the  $i$ -th component of  $\theta$ ).

\*  $L$  takes its minimum at 0, e.g.,  $L(t) = t^2$  i.e., the global estimation error is the sum of the squared componentwise errors.

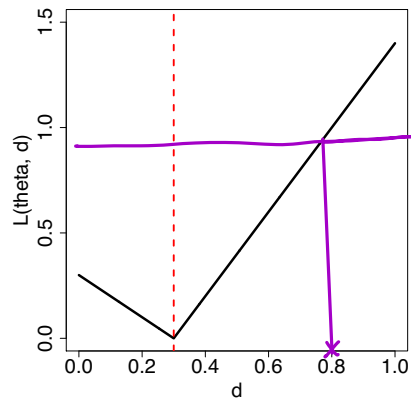
\*  $L(\theta, \delta)$  prevents the overall loss from being heavily affected by components with a large variance.

## JB Example 1 (page 5) (contd)

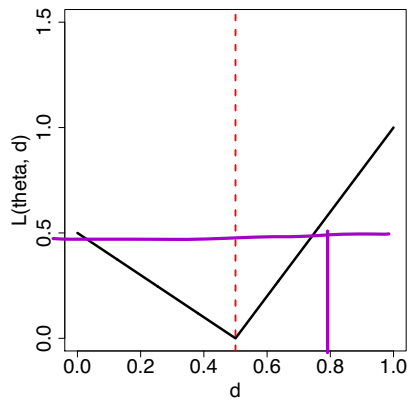
- The company thinks that an overestimation of demand (and hence overproduction of the drug) is twice as costly as an underestimate of demand and that otherwise the loss is linear in the error.
- The company might determine the loss function to be

$$L(\theta, d) = \begin{cases} |\theta - d| & \text{if } \theta \geq d \text{ (underestimation),} \\ 2|\theta - d| & \text{if } \theta < d \text{ (overestimation),} \end{cases}$$

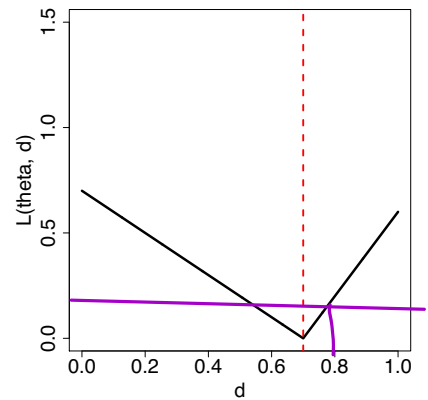
## JB Example 1 (page 5) (contd)



(a)  $\theta = \underline{0.3}$



(b)  $\theta = \underline{0.5}$



(c)  $\theta = \underline{0.7}$

## JB Example 1 (page 5) (contd)

- Conduct a sample survey to obtain sample information about  $\theta$  would be to.
- For example, assume  $n$  people are interviewed, and the number  $x$  who would buy the drug is observed. A reasonable choice for such  $x$  might be  $x \sim \text{Bin}(n, \theta)$ ,

$$\underline{f(x \mid \theta)} = \binom{n}{x} \theta^x (1 - \theta)^{n-x}.$$

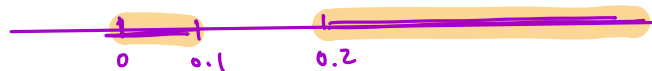
- $\mathcal{X}$ : sample space (the set of all possible outcomes),  $x$ : a particular realization,  $x \in \mathcal{X}$ .



## JB Example 1 (page 5) (contd)

- There could be considerable prior information about  $\theta$ , arising from previous introductions of new similar drugs into the market.
- Suppose that new drugs tended to capture between  $1/10$  and  $1/5$  of the market, with all values between  $1/10$  and  $1/5$  being equally likely. That is,

$$\pi(\theta) = 10, \text{ for } \theta \in (0.1, 0.2).$$



## † A fundamental basis of Bayesian Decision Theory

- Statistical inference should start with the rigorous determination of three factors;

✓  
★★ the distribution family for the observations (sampling distribution),  $f(x | \theta)$  for  $x \in \mathcal{X}$

★★ the prior distribution for the parameter  $\pi(\theta)$ ,  $\theta \in \Theta$

✓  
★★ the loss association with the decisions,  $L(\theta, \delta) \in [0, +\infty)$

$$\delta \in \mathcal{D}$$

$$\delta(x) = d$$

$$\delta(x) \\ x \longrightarrow d$$

† Decision Rule ( $\delta$ )

$$\delta(x) = \bar{x} = \frac{\sum_{i=1}^n x_i}{n} \in \mathcal{D} = \mathcal{H}$$

- **JB Def 2** (p9): A (nonrandomized) decision rule  $\delta(x)$  is a function from  $\mathcal{X}$  into  $\mathcal{D}$ , i.e., the allocation of a decision to each outcome  $x \sim f(x | \theta)$  from a random experiment.
- If  $x$  is the observed value (assumed to follow  $f(x | \theta)$ ), then  $\delta(x)$  is the action that will be taken.
- In estimation problems, decision rule  $\delta$ , from  $\mathcal{X}$  to  $\mathcal{D}$  is usually called estimator (while the *value*  $\delta(x)$  is called *estimate* of  $\theta$ ).
- **JB Example 1 (page 5)**:  $\delta(x) = x/n$  (sample proportion): this does not incorporate the loss function or prior information

## † Bayesian Approach to Decision Theory

- Minimize the expected loss of a decision  $d$  for the believed distribution of  $\theta$  at the time of decision making, i.e.,  $\pi(\theta \mid x)$ .
- The *posterior expected loss* of decision  $d$ , when the posterior distribution is  $\pi(\theta \mid x)$ ,

$$\begin{aligned}\rho(\pi, d \mid x) &= E^{\pi} [L(\theta, d) \mid x] \\ &= \int_{\Theta} L(\theta, d) \pi(\theta \mid x) d\theta.\end{aligned}$$

$\Rightarrow \rho(\pi, d \mid x)$  averages the error (loss) according to the posterior distribution of  $\theta$ , conditionally on the observed data.

- A *Bayes decision*,  $\delta^{\pi}(x)$  is any decision  $d \in \mathcal{D}$  which minimizes  $\rho(\pi, d \mid x)$ .

**JB Example 4**(p10) Assume  $X | \theta \sim N(\theta, 1)$ . The goal is estimating  $\theta$ .

★★  $d \in \mathcal{D} = \mathbb{R}$ .

★★ sampling distribution:  $N(\theta, 1)$

★★ prior distribution:  $N(\mu, \tau^2)$

$$\Rightarrow N\left(\underbrace{\left(\frac{1}{1} + \frac{1}{\tau^2}\right)^{-1} \left(\frac{x}{1} + \frac{\mu}{\tau^2}\right)}_{= \mu_1}, \underbrace{\left(\frac{1}{1} + \frac{1}{\tau^2}\right)^{-1}}_{= \tau_1^2}\right)$$

★★ loss function: squared-error loss,  $L(\theta, d) = (\theta - d)^2$

Find the posterior expected loss for any  $d \in \mathcal{D}$ . ( $\mu_1 - d$ ) ·  $E(\theta - \mu_1 | x)$

Find  $d^{opt}$  =  $\arg \min_{d \in \mathcal{D}} p(\pi, d | x)$

$$E((X - E(X))^2) = \text{Var}(X)$$

$$p(\pi, d | x) = E^\pi(L(\theta, d) | x) = E^\pi((\theta - d)^2 | x)$$

$$= E^\pi((\theta - \mu_1 - d + \mu_1)^2 | x)$$

$$= E^\pi((\theta - \mu_1)^2 + 2(\theta - \mu_1)(\underline{\mu_1 - d}) + (\mu_1 - d)^2 | x)$$

$$= \underbrace{E^\pi((\theta - \mu_1)^2 | x)}_{\text{Var}(\theta | x)} + E((\mu_1 - d)^2 | x)$$

$$= \sigma_1^2 + \underbrace{E((\mu_1 - d)^2 | x)}_{z_0} = (\mu_1 - d)^2$$

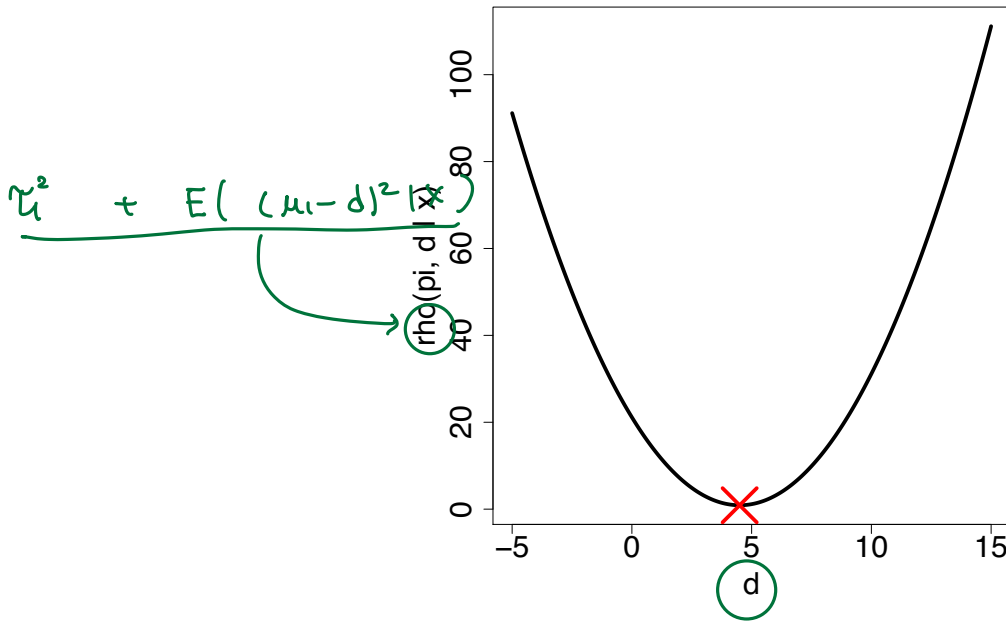
$d = \mu_1 \rightarrow$  Bayes decision is  $d = \mu_1$

**JB Example 4**(p10) (contd) Suppose  $x = 5$  is observed. Assume  $\mu = 0$  and  $\tau^2 = 9$ .

$$\delta(x) = \left( \frac{1}{1} + \frac{1}{\tau^2} \right)^{-1} \left( \frac{x}{1} + \frac{\mu}{\tau^2} \right)$$

$$= \mu_1$$

$$\underline{\delta(5) = 4.5}$$



- Find  $\delta^\pi(x)$ .

† Frequentist Risk (Average Loss, CR Section 2.3 & JB Section 1.3)

- In the frequentist paradigm, the long run performance of  $\delta(x)$  by varying  $x \in \mathcal{X}$  is the key.
- **JB Def 3 (p9) & CR p61:** The *frequentist risk (or average risk)* of a decision rule  $\delta(x)$  is defined by

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

⇒ The error (loss) is averaged over the different values of  $x$  proportionally to the density  $f(x | \theta)$ .

- Suppose we have multiple estimators and want to compare them (or even want to select the best estimator). *How?*



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

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

$$\begin{aligned}
 R(\theta, \delta_c) &= E_{\theta} ( L(\theta, \delta_c) ) \\
 &= E_{\theta} ( (\theta - c \cdot x)^2 ) \\
 &\quad \quad \quad \text{+ c}\theta \\
 &= E_{\theta} ( (c(\theta - x) + (1-c)\theta)^2 ) \\
 &= E_{\theta} ( \underbrace{c^2(\theta - x)^2} + 2c\underbrace{(\theta - x)} + (1-c)^2\theta^2 ) \\
 &= c^2 \cdot 1 + 0 + (1-c)^2\theta^2
 \end{aligned}$$