

# Statistical Decision Theory and Bayesian Analysis

## CH 1 BASIC CONCEPTS

### 1.1 Introduction



Example 1: A drugmaker cares about

- $\theta_1$ : % of people for whom the drug is effective
- $\theta_2$ : % of market.

Classical statistics use only sample information.

Non-sample information



- Loss (Statistics are pessimistic)
- Prior

E.g. A lady claims to tell whether tea or milk is poured into the cup first



A music expert claims to distinguish Haydn score from Mozart score

A drunker claims to predict the outcome of a fair win.

In 10 trials, they are all wrong.

Frequentist's hypothesis test

$$H_0: \theta = 0.5 \text{ (The person is guessing)}$$

$H_0$  is rejected with one tailed significance level of  $2^{-10}$

$$\text{N.B. } p\text{-value}(\mathcal{D}) \triangleq P(f(\hat{\mathcal{D}}) \geq f(\mathcal{D}) \mid \hat{\mathcal{D}} \sim H_0)$$

(See Machine Learning: A Probabilistic Perspective)  
P.213

here:  $f$  is the summarization of the outcome



## §1.2 Basic Elements

### Definitions:

$\theta$ : state of nature / parameter

that affects the decision process

$\Theta$ : parameter space

$a$ : action (decision)  $\leftarrow$  Think of an estimator  $\hat{\theta}$

$A$ : all actions

Loss function:  $L(\theta, a) : \Theta \times A \rightarrow \mathbb{R}$

A particular action  $a_0$  with true state of nature  $\theta_0$ ,

results in loss  $L(a_0, \theta_0)$ .

Samples  $X = (X_1, \dots, X_n)$   $X_i$  iid.

Sample space  $\mathcal{X}$ : possible outcomes of  $X_i$

Reminders:  $A$  is not related to  $\mathcal{A}$

- Let  $A$  be a set of events.  $A \subseteq \mathcal{X}$

$$P_\theta(A) = \int_A f(x|\theta) dx$$

or  $P_\theta(A) = \sum_{x \in A} f(x|\theta).$

- $E_\theta[h(X)] = \int_{\mathcal{X}} h(x) \cdot f(x|\theta) dx$

or  $E_\theta[h(X)] = \sum_{x \in \mathcal{X}} h(x) \cdot f(x|\theta).$

LOSS, Risk, etc

Definition: If  $\pi^*(\theta)$  is the believed probability distribution of  $\theta$  at the time of decision making. the Bayesian expected loss of an action  $a$  is

$$P(\pi^*, a) = E^{\pi^*} L(\theta, a) = \int_{\Theta} L(\theta, a) dF^{\pi^*}(\theta)$$

Definition A (nonrandomized) decision rule  $\delta(x)$  is a function from  $\mathcal{X}$  into  $\mathcal{A}$ . If the data sample is  $x$ , then the action is  $\delta(x)$

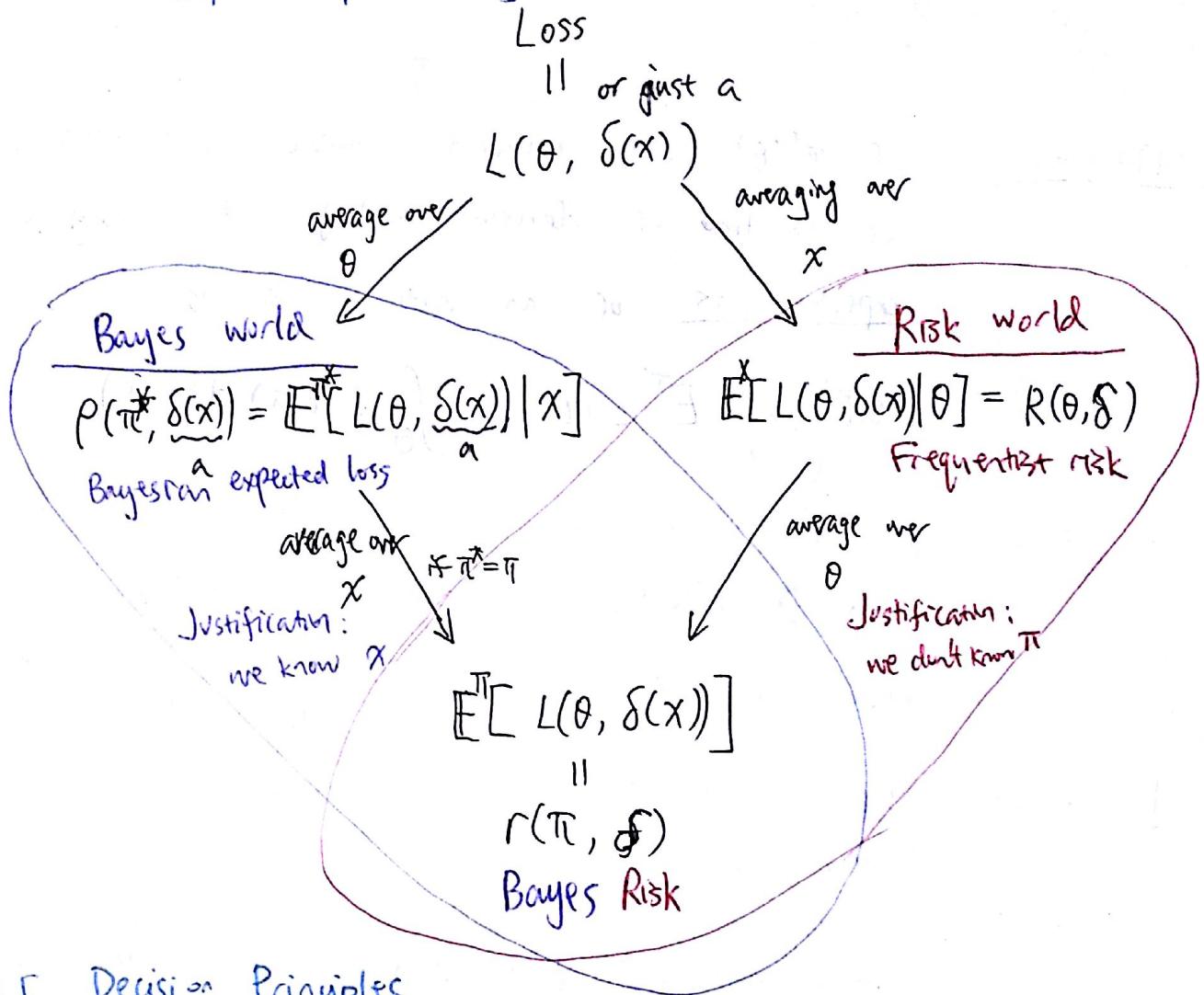
Definition The frequentist's risk function of a decision rule  $\delta(x)$  is

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

For no-data problem,  $R(\theta, \delta) = L(\theta, \delta)$

Definition The Bayes risk of a decision rule  $\delta$ , with respect to a prior distribution  $\pi$  on  $\Theta$  is

$$r(\pi, \delta) = E^{\pi} [R(\theta, \delta)]$$



### § 1.5 Decision Principles

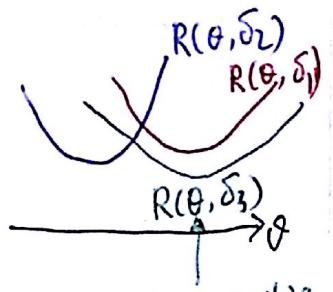
Which decision to use?

Bayesian 1°  $x$  known Bayesian expected loss

2°  $x$  unknown Bayes risk  
(or no-data)

Frequentist. Additional criterion needed.

- Admissibility
- Min max
- Invariance
- (Resort to  $\pi$ ). Bayes Risk.

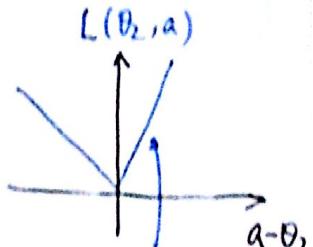


Example 1. Drugmaker wants to estimate the proportion of the market  $\theta_2$ . Obviously,  $\Omega_2 = \{\theta_2 : 0 \leq \theta_2 \leq 1\} = [0, 1]$

An action is an estimate of  $\theta_2$ .  $A = [0, 1]$

The company determines that the loss function is

$$L(\theta_2, a) = \begin{cases} \theta_2 - a & \text{if } \theta_2 - a \geq 0 \\ 2(a - \theta_2) & \text{if } \theta_2 - a \leq 0 \end{cases}$$



Overestimate  
incurs more loss  
( $a > \theta_2$ )

Experiment: a survey among  $n$  people, among whom  $X$  people would like to buy drug. ( $X \sim \mathcal{B}(n, \theta_2)$ )

$$f(x|\theta_2) = \binom{n}{x} \theta_2^x (1-\theta_2)^{n-x}$$

↑  
Binomial

Prior: From previous experience, the drugmaker believes  $\theta_2 \sim U[0.1, 0.2]$

Thus,

$$\pi(\theta_2) = 10 I_{(0.1, 0.2)}(\theta_2)$$

Bayesian expected loss Assume no experiment is conducted

$$\text{i.e. } \pi^*(\theta_2) = \pi(\theta_2) = 10 I_{(0.1, 0.2)}(\theta_2)$$

$$\rho(\pi, a) = \int_0^1 L(\theta_2, a) \pi(\theta_2) d\theta_2$$

$$\begin{aligned} a \text{ is what you} &= \int_0^a 2(a - \theta_2) \cdot 10 \cdot I_{(0.1, 0.2)}(\theta_2) d\theta_2 + \int_a^1 (\theta_2 - a) \cdot 10 I_{(0.1, 0.2)}(\theta_2) d\theta_2 \\ \text{decide to do} \\ \text{when making} \\ \text{decisions} &= \begin{cases} 0.15 - a & \text{if } a \leq 0.1 \\ 15a^2 - 4a + 0.3 & \text{if } 0.1 \leq a \leq 0.2 \\ 2a - 1.3 & \text{if } a \geq 0.2 \end{cases} \end{aligned}$$

Conditional Bayes Principle

$$\text{To minimize } \rho(\pi, a), \quad a^* = \frac{2}{15} \quad \text{and} \quad \min \rho(\pi, a) = \frac{1}{30}$$

Example 3. An investor must decide whether to buy a rather risky ZZZ bonds.

Loss matrix:

	$a_1$	$a_2$
invest	-500	-300
not invest	1000	-300
$\theta = p(\text{win})$		0.9
$p(\text{lose}) = 1 - \theta$		0.1

$\downarrow \max_{\text{win}}$        $\downarrow \max_{\text{lose}}$

Prior: The investor determines according to his/her experience that

$$\pi(\theta) = 0.9 \quad \text{i.e., } p(\text{win}) = 0.9 \quad p(\text{lose}) = 0.1$$

### Bayesian expected loss

$$\begin{aligned} p(\pi, a_1) &= E^\pi [L(\theta, a_1)] = L(\theta_1, a_1) \pi(\theta_1) + L(\theta_2, a_1) \pi(\theta_2) \\ &= 0.9 \cdot (-500) + 0.1 \cdot 1000 \\ &= -350 \end{aligned}$$

$$\begin{aligned} p(\pi, a_2) &= E^\pi [L(\theta, a_2)] = L(\theta_1, a_2) \pi(\theta_1) + L(\theta_2, a_2) \pi(\theta_2) \\ &= -300 \cdot 0.9 - 300 \cdot 0.1 \\ &= -300 \end{aligned}$$

### Conditional Bayesian Decision Principle

To minimize  $p(\pi, a)$ , with respect to  $a$ , we choose  $a = a_1$ ,  $\min_a p = -350$

### Minimax Principle

$$\sup_\theta L(\theta, a_1) = \max \{-500, 1000\} = 1000$$

$$\sup_\theta L(\theta, a_2) = \max \{-300, -300\} = -300$$

To minimize  $\sup_\theta L(\theta, a)$ , we choose  $a = a_2$

$$\inf_a \sup_\theta L(\theta, a) = -300$$

Example 4:  $X \sim N(\theta, 1)$

action: to estimate  $\theta$ .

$$\text{Assume } L(\theta, a) = (\theta - a)^2$$

Consider the decision rule

$$\delta_c(x) = cx$$

Frequentist risk

$$R(\theta, \delta_c) = E_\theta^X [L(\theta, \delta_c(x))]$$

$$= E_\theta^X [(\theta - cx)^2]$$

$$= E_\theta^X [(c(\theta - x) + (1-c)\theta)^2]$$

$$= E_\theta^X [c^2(\theta - x)^2] + 2 \cdot E_\theta^X [c(1-c) \cdot \theta \cdot (\theta - x)] + E_\theta^X [(1-c)^2 \theta^2]$$

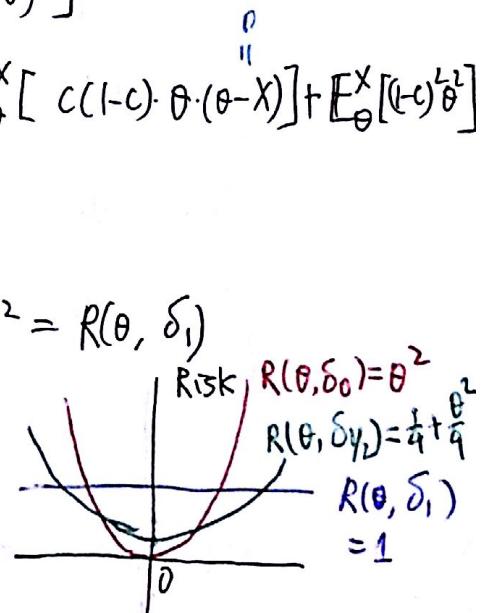
$$= c^2 + (1-c)^2 \theta^2$$

For  $c > 1$

$$R(\theta, \delta_c) > c^2 + (1-c)^2 \theta^2 > c^2 = R(\theta, \delta_1)$$

Thus  $\delta_c$  ( $c > 1$ ) inadmissible

For  $0 \leq c \leq 1$  : noncomparable



Definition A decision rule  $\delta_1$  is R-better

then a decision rule  $\delta_2$  if  $R(\theta, \delta_1) \leq R(\theta, \delta_2), \forall \theta \in \mathbb{R}$

Definition A decision rule is admissible if there exists no R-better decision rule. A decision rule is inadmissible if there does exist an R-better decision rule.

Bayes riskSuppose  $\pi(\theta) = N(\theta | 0, \tau^2)$ 

$$\begin{aligned} r(\pi, \delta_c) &= E^\pi[R(\theta, \delta_c)] \\ &= E^\pi[C^2 + (1-c)^2 \theta^2] \\ &= c^2 + (1-c)^2 \tau^2 \end{aligned}$$

The Bayes Risk Principle A decision rule  $\delta_1$  is preferred to a rule  $\delta_2$  if

$$r(\pi, \delta_1) < r(\pi, \delta_2)$$

A decision rule minimizes  $r(\pi, \delta)$  is optimal,  
if it is called a Bayes rule, denoted at  $\delta^\pi$ .

$r(\pi) = r(\pi, \delta^\pi)$  is called Bayes risk for  $\pi$

In the aforementioned example

$$f(\pi) = r(\pi),$$

$$r(\pi, \delta_c) = c^2 + (1-c)^2 \tau^2$$

$$\frac{\partial}{\partial c} r(\pi, \delta_c) = 2c + 2(1-c)(-1)\tau^2 \stackrel{\Delta}{=} 0$$

$$\Rightarrow c^* = \frac{\tau^2}{1+\tau^2}$$

$$\text{Thus } r(\pi) = r(\pi, \delta_{c^*}) = \left(\frac{\tau^2}{1+\tau^2}\right)^2 + \left(\frac{1}{1+\tau^2}\right)^2 \cdot \tau^2 = \frac{\tau^2}{1+\tau^2}$$

## § 1.6. Foundations

SDT/BA-CH-1-9

### § 1.6.1. Misuse of CLT

#### Misuse of Classical Inference Procedures

- p-value is harmful!

My informal note on the risk function of frequentist hypothesis test:

Example:  $X_1, \dots, X_n \sim N(\theta, 1)$   $R(\theta, \delta) = \begin{cases} \alpha & \text{if } \theta_0 < \delta \\ \text{type II Err. if } \theta_0 \geq \delta \end{cases}$   
Conduct a size  $\alpha = 0.05$  test. The solution is not unique.

$$H_0: \theta = 0 \Leftrightarrow H_1: \theta \neq 0$$



The classical test: reject  $H_0$  if  $\sqrt{n} |\bar{x}| > 1.96$

What about  $n = 10^{24}$ ?

Another Example:

Example:  $X \sim N(\theta, 9)$ .

Kepler's data

$$H_0: \theta \leq 0 \Leftrightarrow H_1: \theta > 0.$$

A sample of 9 observations resulting in  $\bar{x} = 1$ . has  $\alpha = 0.16$ .

Moderately convincing.

Collecting more data?



### § 1.6.2. The Frequentist Perspective (frequentists' Justification)



The original idea was from J. Neyman and E. Pearson.

A frequentist intends to

- Find out a method independent of  $\theta$  or its prior knowledge, i.e., get a value  $\bar{R}$  by  $S(x)$  and  $L(\theta, S(x))$ , which satisfies that
- the expectation of  $L$  is  $\bar{R}$  while Repeating  $S(x)$ .

Example. For dealing with standard univariate normal theory problems consider the usual 95% confidence rule for the unknown mean  $\theta$ .

$$\delta(X) = (\bar{x} - ts, \bar{x} + ts)$$

where  $\bar{x}$  is the sample mean and  $s$  is the sample std.

$t$  is the appropriate percentile from relevant  $t$  distribution.

Suppose we measure the performance of  $\delta$  by loss

$$L(\theta, \delta(X)) = 1 - I_{\delta(X)}(\theta) = \begin{cases} 0, & \text{if } \theta \in \delta(X) \\ 1, & \text{if } \theta \notin \delta(X) \end{cases}$$

Then the frequentist risk is,

$$\begin{aligned} R((\theta, \sigma), \delta) &= \mathbb{E}_{\theta, \sigma}^X [L(\theta, \delta(X))] \\ &= P_{\theta, \sigma} \{ \delta(X) \text{ does not contain } \theta \} \\ &= 0.05 \quad (\text{At } \theta, \sigma) \end{aligned}$$

The universal measure of performance is, if we have different problems  $\theta_i, \sigma_i$  with data  $X^{(i)}$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N L(\theta_i, \delta(X^{(i)})) = 0.05 \equiv \bar{R}$$

GREAT!

However, in general  $R(\theta, \delta)$   $\stackrel{\text{?}}{}$  has not eliminate dependence of  $\theta$ ,  $\stackrel{\text{?}}{}$  is repeated over  $X$  rather than use of  $\delta$ .

Bound of  $R$  for certain  $\theta$  is NOT satisfactory.

Example: Consider a simple null hypothesis test

$$H_0: \theta = \theta_0 \leftrightarrow H_1: \theta = \theta_1$$

Let us define loss to be "0-1" loss

Then the risk function

$$R(\theta_0, \delta) = \alpha_0 \quad (\text{Type I error, FP})$$

$$R(\theta_1, \delta) = \alpha_1 \quad (\text{Type II error, FN})$$

We claim only 1% of true null hypotheses will be rejected.

However, we can say nothing ~~that the null hypothesis~~ about how often our errs when rejecting.

		Accept	reject
Actual	$\theta_0$	1	0
	$\theta_1$	0	?

Let  $\alpha_1 = 0.99$  and  $p(\theta_0) = p(\theta_1) = 0.5$

Then  $P(\text{Error} | \text{Reject}) = 0.5$



"real" error, because we do not know  $\theta$

### 1.6.3 The Conditional Perspective.

Example. Let  $X = (X_1, X_2)$ .  $X_i$  iid.

$$p_\theta(X_i = \theta - 1) = p_\theta(X_i = \theta + 1) = \frac{1}{2}$$

Frequentist:  $\delta(X) = \begin{cases} \frac{1}{2}(X_1 + X_2) & \text{if } X_1 \neq X_2 \\ X_1 - 1 & \text{if } X_1 = X_2 \end{cases}$

Loss: "0-1" loss half chance  $X_1 = X_2$

$$R(\theta, \delta) = \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

↑ half chance  $X_1 \neq X_2$       ↑ half chance  $\delta(X) = \theta$

if  $X_1 = X_2$

Nonsensical to report  $\delta$  risk even though  $X$  is provided.

Conditionalist: if  $x_1 \neq x_2$   $a = \frac{1}{2}(x_1 + x_2)$

$$\rho(\pi^*, a) = 0$$

$$\text{if } x_1 = x_L \quad \pi^*(\theta = x_L) = \frac{1}{2}$$

$$\pi^*(\theta = x_1 + 1) = \frac{1}{2}$$

Take action  $a = x_1 - 1$

$$\rho(\pi^*, a) = 0.5$$

A Conditionalist knows whether one is in the 100% or 50% case.

Example:

		reject H <sub>0</sub>	accept H <sub>0</sub>
		1	2
f(x θ=0)	1	0.005	0.005
	2	0.99	
f(x θ=1)	0.0051	0.9849	0.01

Let  $\theta = 0, 1$ .

$x = 1, 2$  or  $3$

$f(x|\theta)$  as left.

The classical most powerful test of  $H_0: \theta = 0 \leftrightarrow H_1: \theta = 1$  at level  $\alpha = 0.01$ :

Reject  $H_0$  (conclude  $H_1$ ) if  $x = 1$  or  $2$

Accept  $H_0$  if  $x = 3$

Type I error  $\alpha_0 = 0.01$

Type II error  $\alpha_1 = 0.01$

Standard frequentist: upon observing  $x=1$ , report the decision is  $H_1$ , with error probability

Conditionalist:

when observing  $x=1$ , likelihood ratio

$$\frac{f(x|\theta=0)}{f(x|\theta=1)} = \frac{0.005}{0.0051} \approx 1$$



Example: Cox 1988.

Suppose a substance to be analyzed can be sent either to a laboratory in New York or a laboratory in California. The two labs seem equally good, so a fair coin is flipped to choose between them, with "heads" denoting that the lab in New York will be chosen. The coin is flipped and comes up tails, so the California lab is used. After a while, the experimental result 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 the experiment in New York might

have been performed instead? Common sense (and the conditional viewpoint) views no, that only the experiment actually performed is relevant, but frequentist reasoning would call for averaging over all possible data, even the possible New York data.

### § 1.6.4 The Likelihood Principle

Definition For observed data,  $x$ , the function  $l(\theta) = f(x|\theta)$ , considered as a function of  $\theta$ , is called the likelihood function.

The Likelihood Principle In making inferences or decisions about  $\theta$  after  $x$  is observed, all relevant experimental information is contained in the likelihood function for observed  $x$ . Furthermore two likelihood functions contain the same information about  $\theta$  if they are proportional to each other (as functions of  $\theta$ )

Example:  $p(\text{Head})$  of a coin is  $\theta$

9 heads and 3 tails

Let  $X = \# \text{ heads}$

How is the experiment conducted?

case #1: The experiment consisted of a predetermined 12 flips

$$l_1(\theta) = f_1(x|\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} = 220 \cdot \theta^9 (1-\theta)^3$$

Notes

case #2: The experiment consisted of flipping the coin

1° By "proportion", we impose constraint on the same parameter.

$$l_2(\theta) = f_2(x|\theta) = \binom{n+x-1}{x} \theta^9 (1-\theta)^3 = 55 \cdot \theta^9 (1-\theta)^3$$

2° Likelihood principle says

Classical hypothesis test  
about experimental information

$$H_0: \theta = \frac{1}{2} \leftrightarrow H_1: \theta > \frac{1}{2}$$

$$\alpha_1 = P_{Y_2}(X \geq 9) = f_1(9|\theta=\frac{1}{2}) + f_1(10|\theta=\frac{1}{2}) + f_1(11|\theta=\frac{1}{2}) + f_1(12|\theta=\frac{1}{2}) \\ = 0.075$$

$$\alpha_2 = P_{Y_2}(X \geq 9) = f_2(9|\theta=\frac{1}{2}) + f_2(10|\theta=\frac{1}{2}) + \dots \\ = 0.0325$$

Example. Let  $\mathcal{X} = \{1, 2, 3\}$ ,  $\mathcal{H} = \{0, 1\}$ . Two experiments  $E_1, E_2$  with observing data  $X_1$  and  $X_2$ , respectively.

	$X_1$				$X_2$		
	1	2	3		1	2	3
$f_1(X_1 0)$	0.90	0.05	0.05		0.26	0.73	0.01
$f_1(X_1 1)$	0.09	0.055	0.855		0.026	0.803	0.171

Likelihood Principle says:  $E_1$  and  $E_2$  are equivalent given  $X$  because  $l_1(\theta) \propto l_2(\theta)$

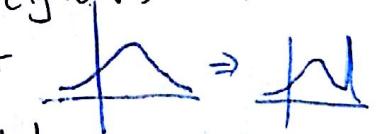
Frequentists: accept  $H_0$  if  $X_1=1$  (or  $X_2=1$ )  
reject  $H_0$  if  $X_1=2, 3$  (or  $X_2=2, 3$ )

Type I error      Experiment I: 0.1      Experiment II: 0.74

Type II error      Experiment I: 0.09      Experiment II: 0.026

Example:



- Engineer draws a random sample of electron tubes and measure the plate voltages under certain conditions with a very accurate voltmeter
- A statistician examines the samples which seem to be Gaussian distributed, varying from 75 to 99,  $m_{\text{est}}=87$ ,  $s_{\text{std}}=4$
- Later the statistician visited the engineer's lab, and found the voltmeter reads 100 at most  
  
Hence a new analysis is needed  

- The engineer says he has a high-range voltmeter and he would have used it if the voltage > 100
- The statistician is happy
- The engineer says actually the high-range voltmeter is broken
- The statistician says a new analysis is needed.

## The Weak Conditionality Principle



Suppose one can perform either of two experiments  $E_1$  or  $E_2$ , both pertaining to  $\theta$ , and that the actual experiment conducted is the mixed experiment of first choosing  $J=1$  or  $2$  with probability  $\frac{1}{2}$  each (independent of  $\theta$ ), and then performing experiment  $E_J$ . Then the actual information about  $\theta$  obtained from the overall mixed experiment should depend only on the experiment  $E_J$  that is actually performed.

Sufficiency + Weak Conditionality Principle  $\Rightarrow$  Likelihood Principle

Frequentists and Conditionals agree on sufficiency. Then frequentists are in trouble.

### § 1.7 Sufficient Statistics.

#### Definition

Let  $X$  be a random variable whose distribution depends on the unknown parameter  $\theta$ , but is otherwise unknown. A function  $T$  of  $X$  is said to be sufficient statistic for  $\theta$  if the  $f(x|T=t)$  is independent of  $\theta$ .

#### Fisher's factorization Theorem

$$T(X) \text{ is sufficient for } \theta \Leftrightarrow f_{\theta}(x) = h(x) g_{\theta}[T(x)].$$

#### Definition Partition induced by $T$

$$X_t = \{x \in X : T(x)=t\}$$

Argue: a sufficient statistic contain all information concerning  $\theta$  from samples.

Given  $T(x) = t$ , denote  $f_t(x)$  as a distribution over  $\mathcal{X}$ .

The density  $f_t(x)$  does not depend on  $\theta$  because  $f(X | T=t) \perp\!\!\!\perp \theta$

$\Rightarrow f_t(x)$  is explicitly calculable from  $f_\theta(x|\theta)$  for whatever  $\theta$ )

Random variable  $X$  can be thought of as



1° first arising from the random generation of  $T$ .

2° followed by the random choice of  $x$  from  $\mathcal{X}_t$ .

We would now argue sufficiency is sensible from a frequentist's perspective.

Theorem: Assume that  $T$  is a sufficient statistic for  $\theta$ ,

let  $\delta_0^*(x, \cdot)$  be any randomized rule in  $\mathcal{D}^*$ .

Then (subject to measurability conditions) there

exists a randomized rule  $\delta_1^*(t, \cdot)$  depending only

on  $T(x)$ , which is R-equivalent to  $\delta_0^*$

$$\delta^*: \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$$

( $\delta^*(x, a)$  is the probability that we take action  $a$  given data  $x$ )

Proof: If  $A \subseteq \mathcal{A}$  and  $t \in \mathcal{T}$ , we design  $\delta_1^*$  base on  $\delta_0^*$  as

$$\delta_1^*(t, A) = E^{X|t} [\delta_0^*(X, A)]$$

i.e.  $\delta_1^*(t, \cdot)$  is formed by averaging  $\delta_0^*$  over  $\mathcal{X}_t$

$$\begin{aligned} L(\theta, \delta_1^*(t, \cdot)) &= \mathbb{E}^{\delta_1^*(t, \cdot)} [L(\theta, a)] \\ &= \mathbb{E}^{X|t} \mathbb{E}^{\delta_0^*(X, \cdot)} [L(\theta, a)] \end{aligned}$$

Therefore

$$\begin{aligned} R(\theta, \delta_1^*) &= \mathbb{E}^T [L(\theta, \delta_1^*(T, \cdot))] \\ &= \mathbb{E}^T \mathbb{E}^{X|T} \mathbb{E}^{\delta_0^*(X, \cdot)} [L(\theta, a)] \\ &= \mathbb{E}^X \mathbb{E}^{\delta_0^*(X, \cdot)} [L(\theta, a)] \quad \mathbb{E}^X [h(x)] = \mathbb{E}^T \mathbb{E}^{X|T} [h(X)] \\ &= \mathbb{E}^X [L(\theta, \delta_0^*(X, \cdot))] \\ &= R(\theta, \delta_0^*) \quad \# \end{aligned}$$

[by substituting the design of  $\delta_1^*$ , and then changing the order of expectation.]

The result implies that it is only necessary to consider rules based on a sufficient statistic.

The Likelihood Principle immediately implies sufficiency

$$l(\theta) = f(x|\theta) = h(x) g(T(x)|\theta) \propto g(T(x)|\theta)$$

### § 1.7 Convexity

Theorem: Assume that

1.  $\mathcal{A} \subseteq \mathbb{R}^m$  convex
2.  $\forall \theta \in \mathbb{R}$ ,  $L(\theta, a)$  is convex w.r.t  $a$
3. Let  $\delta^*$  be a randomized in  $\mathcal{D}^*$  for which  $\mathbb{E}^{\delta^*(x, \cdot)} [\|a\|] < \infty \quad \forall x \in \mathbb{R}^m$

Then.  $\forall x$  and  $\theta$  the nonrandomized decision rule

$$\delta(x) = \mathbb{E}^{\delta^*(x, \cdot)} [a]$$

has  $L(\theta, \delta(x)) \leq L(\theta, \delta^*(x, \cdot))$

Proof: 1°  $\delta(x) \in \mathcal{A}$

2°  $L(\theta, \delta(x)) = L(\theta, \mathbb{E}^{\delta^*(x, \cdot)} [a]) \leq \mathbb{E}^{\delta^*(x, \cdot)} [L(\theta, a)] = L(\theta, \delta^*(x, a)) \quad \#$

Jensen's Inequality:  $g(\mathbb{E}x) \leq \mathbb{E}[g(x)]$  if  $g$  convex.

Theorem (Rao-Blackwell) Assume

1°  $A \subseteq \mathbb{R}^m$  convex

2°  $\forall \theta \in \Theta, L(\theta, a)$  convex in  $a$

3° Let  $T(x)$  be a sufficient statistics, and  $\delta_0 \in \mathcal{D}$

Then a nonrandomized decision rule  $\delta_1$ , ~~non~~, defined by

$$\delta_1(t) = \mathbb{E}^{X|t} [\delta_0(X)]$$

is R-equivalent to or R-better than  $\delta_0$ , provided that the expectation exists.

↓ Jensen's inequality

$$L(\theta, \delta_1(t)) = L(\theta, \mathbb{E}^{X|t} [\delta_0(X)]) \leq \mathbb{E}^{X|t} [L(\theta, \delta_0(X))]$$

$$\text{Then, } R(\theta, \delta_1) = \mathbb{E}^T [L(\theta, \delta_1(T))]$$

$$\leq \mathbb{E}^T \mathbb{E}^{X|T} L(\theta, \delta_0(X))$$

$$= \mathbb{E}^X [L(\theta, \delta_0(X))]$$

$$= R(\theta, \delta_0)$$

Theorem p16

Another Proof:

$$\text{Define } \delta_1^*(t, A) = \mathbb{E}^{X|t} [I_A(\delta_0(X))] \iff \delta_0$$

$$\text{Then } \delta_1(t) = \mathbb{E}^{\delta_1^*(t, \cdot)} [a] = \mathbb{E}^{X|t} \mathbb{E}^{\mu_X} [a] = \mathbb{E}^{X|t} [\delta_0(X)]$$

$$\text{Let } \mu_X(A) = I_A(\delta_0(X))$$

is R-equivalent or R-better than  $\delta_0$

Theorem p17

#