

# Decision Theory and Bayesian Inference

It's a unifying framework for theory of statistics including estimation and testing.

## (1) Definition:

Suppose "a" is an action,  $a \in A$ , the action space.

$\Rightarrow$  The choice of action "a" depends on:

- i) Observations of R.V.: Data  $\vec{X}$   
where  $\vec{X} \in S$ , sample space.
- ii)  $d$ , the decision function, i.e.  
 $d: S \rightarrow A$ .  $d(\vec{X}) = a$

$\Rightarrow$  The prob dist of  $\vec{X}$  depends on the parameter  $\theta$ , called state of nature.

$\theta \in \Lambda$ , the space of parameters

$\Rightarrow$  Then we can define a loss function  $l(\theta, a)$

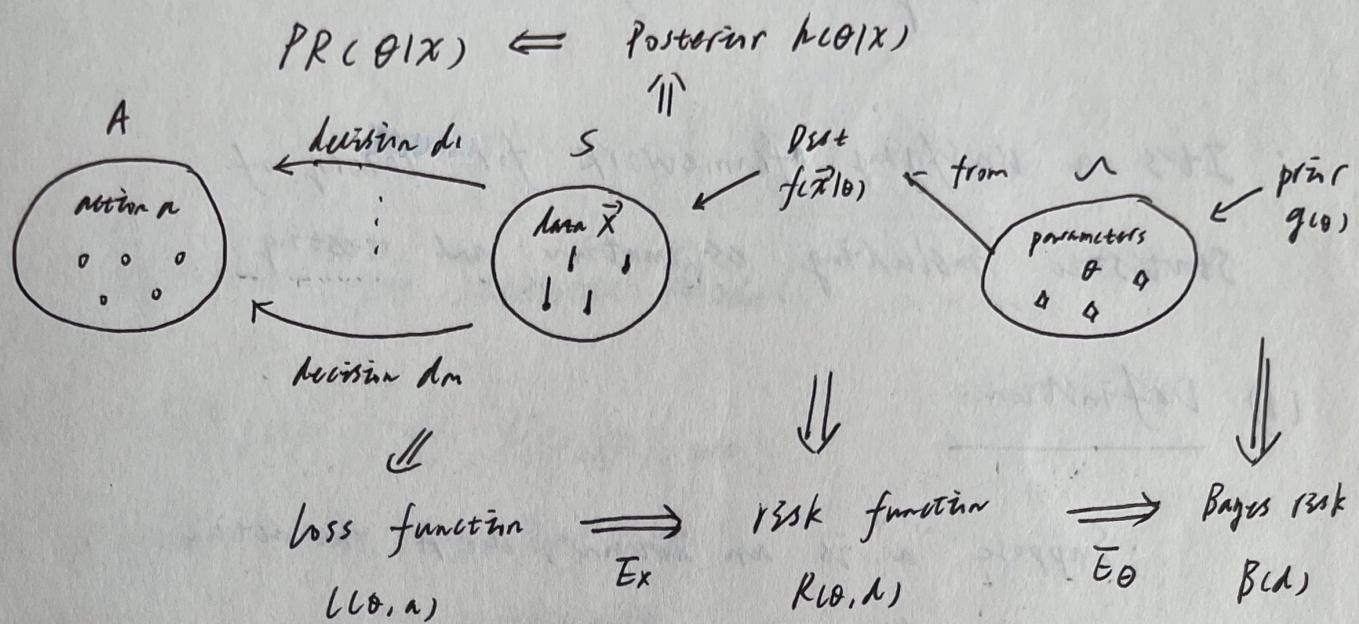
on  $\Lambda \times A$ . Since  $a = d(\vec{X}) \therefore l(\theta, a) = l(\theta, d(\vec{X}))$

The expected loss of  $d(\vec{X})$  is the risk function:

$$R(\theta, d) = E_{\vec{X}}(l(\theta, d(\vec{X}))) \quad (\text{It depends on } \theta)$$

$\Rightarrow$  Our aim is to minimize  $R(\theta, \lambda)$  by

choosing a good decision  $d(\vec{x})$ !



e.g., estimate  $V(\theta)$ , where  $X_k \sim f(x|\theta)$ , i.i.d. get data  $\vec{X}$ .

We choose  $L(\theta, \lambda(\vec{x})) = [V(\theta) - \lambda(\vec{x})]^2$ : quadratic loss func

## (2) Minimization:

- Differentials of  $R(\theta, \lambda)$  depend on unknown  $\theta$ .
- minimizing  $R(\theta, \lambda)$ 
  - i)  $R(\theta, \lambda)$  depend on unknown  $\theta$ .
  - ii) For different  $\theta_1, \theta_2$ . It may happen:  
 $R(\theta_1, \lambda_1) > R(\theta_2, \lambda_2)$  how to choose?  
 $R(\theta_2, \lambda_1) < R(\theta_2, \lambda_2)$

### ① Def of

#### minimax rule:

- Consider the worst case:  $\sup_{\theta \in \Theta} R(\theta, \lambda)$

$d^*$  is the minimax rule. If  $\sup_{\theta \in \Theta} R(\theta, d^*)$

$= \inf_{\lambda} \sup_{\theta \in \Theta} R(\theta, \lambda)$ ,  $d^*$  may not exist!

Remark: It's very conservative to consider the worst case which isn't likely to occur.

## ② Bayesian Rule:

We assume  $\theta \in \Theta$ , is random, with a prior dist. Then the Bayesian risk of decision function  $d$  is  $B(d) = E_{\theta} [R(\theta, d)]$

Def: Bayesian rule is a decision func.  $d^{**}$  attain the  $\min_d B(d)$

Remark: It can be interpreted as average of risk with weight form.

## $\Rightarrow$ Posterior Analysis:

### A method for finding Bayesian Rule:

Suppose  $g(\theta)$  prior dist of  $\theta$ ,  $f_{X|\theta}(x)$  condition of  $X$ .

$$\Rightarrow f_{X,\theta} = g(\theta) f_{X|\theta} \Rightarrow \text{Sum/Integration} = f_X$$

We obtain:  $h_{\theta|x} = f_{X,\theta} / f_X$ . posterior dist of  $\theta$ .

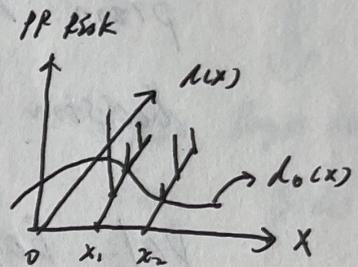
Def: Posterior risk:  $E_{\theta|x} [l(\theta, d(x))] = PR(\theta|x)$

Remark: The observed data  $X=x$  updates the p-risk!

Thm. If  $h_0(x)$  is a decision func. minimizes the posterior risk for each  $x$ . Then  $h_0$  is Bayesian Rule.

$$\begin{aligned}
 \text{Pf: } B(\lambda) &= E_{\theta} (R(\theta, d)) = E_{\theta} (E_{x| \theta} (l(\theta, d(x)) | \theta)) \\
 &= \int \left[ \int l(\theta, d(x)) f_{x|\theta}(x) dx \right] g_{\theta}(\theta) d\theta \\
 &= \int \left[ \int l(\theta, d(x)) h_{\theta|x}(x) d\theta \right] f_x(x) dx \\
 &= \int E_{\theta|x} (l(\theta, d(x)) f_x(x) dx)
 \end{aligned}$$

Then  $B(\lambda)$  is minimized!



$\Rightarrow$  Algorithm:

- 1°) Calculate  $h(\theta|x)$ , for each  $x$ .
- 2°) Calculate  $E_{\theta|x} (l(\theta, d(x)))$  for each  $x$ .
- 3°) Find  $d(x)$  that minimizes every  $PR(\theta|x_0)$ , fixed  $x_0$ .

(3) Application of Decision Theory:

Estimation

- ①  $\left\{ \begin{array}{l} \text{Action space } A \rightarrow \text{Parameter Space } \mathcal{N}. \\ \text{Decision Func. } d(x) \rightarrow \text{estimator of } \theta. \\ l(\theta, d(x)) = [\theta - d(x)]^2 \text{ or } |\theta - d(x)| \end{array} \right.$

$$\text{Thm. i) } E_{\theta|x}((\theta - \hat{\theta})^2 | x) = \text{Var}_{\theta|x}(\theta | x) + [E_{\theta|x}(\theta | x) - \hat{\theta}]^2$$

Then  $\hat{\theta} = E_{\theta|x}(\theta | x)$  is the best predictor.

ii)  $E_{\theta|x}(|\theta - \hat{\theta}| | x)$  has the best predictor = median.

$$\text{Pf: For ii) } E_{\theta|x}(|\theta - \hat{\theta}| | x) = \int |\theta - \hat{\theta}| h_{\theta|x}(\theta) d\theta$$

$$= \int_{\hat{\theta}}^{\infty} (\theta - \hat{\theta}) h_{\theta|x} d\theta + \int_{-\infty}^{\hat{\theta}} (\hat{\theta} - \theta) h_{\theta|x} d\theta \stackrel{a}{=} f(\hat{\theta})$$

$$\frac{\partial f}{\partial \hat{\theta}} = 0 \Rightarrow -\int_{\hat{\theta}}^{\infty} h_{\theta|x} \cos d\theta + \int_{-\infty}^{\hat{\theta}} h_{\theta|x} \cos d\theta = 0$$

$\therefore \hat{\theta} = \text{median of } h_{\theta|x}$ !

② Def: i) d, d<sub>2</sub> two decision functions  $\in A$ .

says  $d_1$  dominates  $d_2$ :  $R(d_1, h) \leq R(d_2, h)$ ,  $\forall h \in \Omega$

ii) d. a decision func. d is admissible. if d is not strictly dominated by any other decision func.

Thm. i) If n is discrete.  $d^*$  is Bayesian rule w.r.t.

prior pmf  $g(\theta)$ . where  $g(\theta) > 0, \forall \theta$ .

ii) If n is dense.  $d^*$  is Bayesian rule w.r.t

prior plf  $g(\theta)$ .  $g(\theta) > 0, \forall \theta$ .  $R(\theta, d)$  is conti.

of  $\theta$ .  $\forall d$ .

Then  $d^*$  is admissible.

Remark: It claims the relation between Bayesian rule and admissible.

Pf: If  $\exists \theta'$  st.  $R(\theta, \theta^*) \geq R(\theta, \theta')$ ,  $\forall \theta$ .

$\exists \theta_0 \text{ s.t. } \theta \in (\theta_0 - h, \theta_0 + h) \Rightarrow R(\theta, \theta^*) > R(\theta, \theta') + \varepsilon$ .

Then check:  $B(d^*) - B(d) > 0$ , contradiction!

#### (4) Bayesian View for prob.:

##### (Personal Opinion)

① Bayesian prob is personal. It varies from person to person, embodying the beliefs of person. (Subjective)

② Bayesian rule describes the prob. evolves with experience.

⇒ Difference between "B"

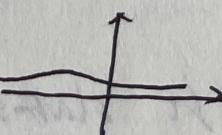
and "F" Approach:

##### ① Point Estimation:

$\left\{ \begin{array}{l} F: \theta \text{ is fixed, unknown, not random} \rightarrow \text{likelihood} \rightarrow \max_{\theta} L(\theta|x) \\ \text{Find MLE.} \end{array} \right.$

$\left\{ \begin{array}{l} B: \theta \text{ has prior dist } g(\theta) \rightarrow \text{posterior} \rightarrow \text{centering:} \\ \text{dist } h(\theta|x) \quad \text{Find E}(\theta|x) \end{array} \right.$

Remark: From  $p(\theta|x) = \frac{f(x|\theta)g(\theta)}{\int f(x|\theta)g(\theta)d\theta}$   $\therefore p(\theta|x) \propto f(x|\theta)g(\theta)$

$\Rightarrow$  If  $\theta$  is uniform:  almost const.

Then  $g(\theta)$  has little effect.  $p(\theta|x) \propto f(x|\theta) = L(\theta|x)$

### ③ Interval Estimation:

Frequentist:  $P(\theta \in [\theta_L(\vec{X}), \theta_U(\vec{X})] | \theta) = 1-\alpha$ .  $\vec{X}$  is random,  $\theta$  fixed.  
 If  $\vec{x}$  is observed data, then  $P(\theta \in [\theta_L(\vec{x}), \theta_U(\vec{x})]) = 0$  or 1.  
 Bayesian:  $P(\theta \in [\theta_L(\vec{x}), \theta_U(\vec{x})] | \vec{X} = \vec{x}) = 1-\beta$ .  $\theta$  is random.  
 $\vec{x}$  is the observed data, fixed.

### ④ Testing:

Frequentist: prob of Type I, II, error.  
 Bayesian: After observing data, compare posterior prob.

### (5) Bayesian Inference

#### for Normal Dist:

Suppose  $X|M \sim N(\mu, \sigma^2)$ ,  $M \sim N(M_0, \sigma_0^2)$ ,  $\sigma$  is known.

$\Rightarrow$  Posterior dist of  $M$  is  $N(M_1, \sigma_1^2)$ , where

$$M_1 = \frac{f_0}{f_0 + f} M_0 + \frac{f}{f_0 + f} x \cdot \frac{1}{\sigma^2} = \frac{1}{f_0 + f}$$

$$f_0 = \frac{1}{\sigma_0^2}, \quad f = \frac{1}{\sigma^2}.$$

Pf: Since the const. is for normalization.

We just care the ratio (about  $M$ ):

$$\begin{aligned}
 h(M|x) &\propto f(x|M) g(M) \propto e^{-\frac{1}{2\sigma^2}(x-M)^2 - \frac{1}{2\sigma_0^2}(M-M_0)^2} \\
 &= e^{-\frac{1}{2\sigma^2} \left[ \frac{1}{\sigma^2} + \frac{1}{\sigma_0^2} \right] M^2 - 2M \left( \frac{x}{\sigma^2} + \frac{M_0}{\sigma_0^2} \right) + \frac{x^2}{\sigma^2} + \frac{M_0^2}{\sigma_0^2}} = a \\
 &\propto e^{-\frac{1}{2\sigma^2} (M - \frac{b}{n})^2} = b
 \end{aligned}$$

Cor. For  $n$  samples  $\vec{X} = (x_1 \dots x_n)$ , i.i.d.

$$M|\vec{X} \sim N(\frac{s_0}{n\beta + s_0} M_0 + \frac{n\bar{x}}{n\beta + s_0}, \frac{1}{n\beta + s_0})$$

Remark: i) Note that  $\frac{s_0}{n\beta + s_0} + \frac{n\bar{x}}{n\beta + s_0} = 1$ . We mix up the prior and the data to generate the posterior. (Weighted Average!)

ii)  $\frac{1}{n\beta + s_0} < \sigma_v^2$ . Which means the ~~post~~ of  $M|\vec{X}$  is more concentrated. It carries more information (informative)

iii) If  $n$  is large enough. Then the data will dominate the prior dist!

iv) For objectiveness, the prior needs to be vague, noninformative!

### (b) Bayesian Inference

for Binomial Dist:

$$\cdot X|p \sim \text{Bin}(n, p), p \sim \text{Beta}(a, b)$$

$$\Rightarrow p|\vec{x} \sim \text{Beta}(a+x, b+n-x)$$

Similarly,  $M_{\text{post}} = \frac{a+b}{n+b+a} \cdot \frac{a}{a+b} + \frac{n}{n+b+a} \cdot \bar{x}$

If  $n \rightarrow \infty$ ,  $M_{\text{post}} \rightarrow \bar{x}$ !

Remark: Define:

$$\begin{cases} \mathcal{G} = \text{family of prior dist. func for } \Theta \\ \mathcal{H} = \text{family of conditional dist. func for } X \end{cases}$$

$\Rightarrow \mathcal{G}$  is called conjugate prior to  $\mathcal{H}$ :

If the posterior of  $\mathcal{G}$  under  $\mathcal{H}$  also belongs to  $\mathcal{G}$ :

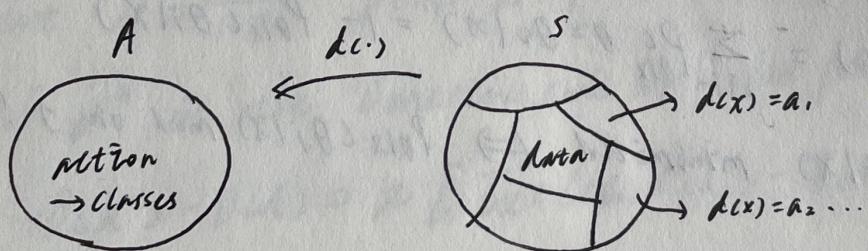
e.g.  $\mathcal{G}$ : Normal dist.  $\mathcal{H}$ : Normal dist  $\rightarrow \mathcal{G}/\mathcal{H} = \mathcal{G}$

$\mathcal{H}$ : Binomial dist  $\mathcal{G}$ : Beta dist  $\rightarrow \mathcal{G}/\mathcal{H} = \mathcal{G}$

## (7) Application of

### Decision Theory:

#### ① Classification:



Let  $A = \{a_1: \in \text{Class } \theta_1, a_2: \in \text{Class } \theta_2, \dots, a_m: \in \text{Class } \theta_m\}$ .

Then  $d(\cdot)$  is function to classify datas in  $S$ .

$\Rightarrow$  For parameter  $\theta_i$ , let  $A = \{\theta_i\}^m$ .

$\pi_i = P(\theta = \theta_i)$ , s.t.  $\sum_i^n \pi_i = 1$ .  $f(x|\theta_i)$  is known.  $1 \leq i \leq m$

If  $l_{ij}$  is the loss in classifying class  $i$  to class  $j$ .

Then, we can find Bayesian Rule:

$$h(\theta_3 | x) = \frac{\pi_3 f(x|\theta_3)}{\sum_i \pi_i f(x|\theta_i)} = P_{\text{C}}(\theta_3 = \theta_3 | X=x), \text{ let } h(x) = a_j$$

We obtain  $PR(j|x) = \sum h_{ij} h(\theta_i|x)$ . Posterior risk

$\Rightarrow$  Choose  $j$  to minimize  $PR(j|x)$ , for  $1 \leq j \leq m$

Then  $d: x \rightarrow a_j$  is a Bayesian rule.

e.g. (0-1 loss)

$$l_{ij} = \begin{cases} 0, & i=j \quad (\text{Because the classification} \\ & \text{is "qualitative", not "quantitative"}) \\ 1, & i \neq j \end{cases}$$

$$\begin{aligned} R(\xi, d) &= E_{x \sim P(x)} l(\theta_3, d(x)) = \sum_j l_{ij} P_{\theta_3}(d(x)=j) \\ &= \sum_{j \neq i} P_{\theta_3}(d(x)=j) = 1 - P_{\theta_3}(d(x)=i) \end{aligned}$$

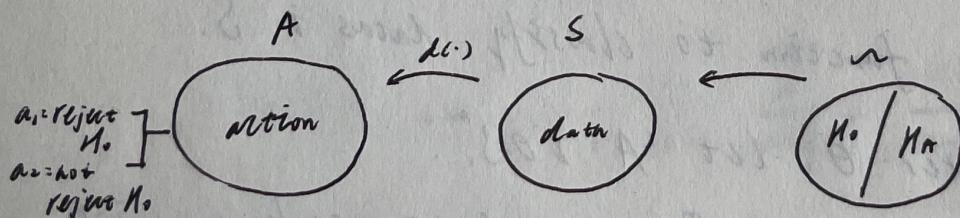
(From frequentist approach, minimax etc.)

$\Rightarrow$  For Bayesian approach:

$$PR(j|x) = \sum_{i \neq j} P_{\theta_3}(\theta = \theta_i | x) = 1 - P_{\theta_3}(\theta_j | x)$$

$\therefore PR(j|x)$  minimized  $\Leftrightarrow P_{\theta_3}(\theta_j | x) \max \text{ on } j$ !

② Hypothesis Testing:



Then both type I and type II errors are misclassification errors.

$\Rightarrow$  Bayesian approach:

$p(H_0) = \lambda$ ,  $p(H_A) = 1 - \lambda$ , then we obtain:

$$\frac{p(H_0|x)}{p(H_A|x)} = \frac{p(H_0)}{p(H_A)} \cdot \frac{p(x|H_0)}{p(x|H_A)} > 1. \text{ (LHS} = \frac{p(H_0|x)}{p(H_A|x)}$$

under 0-1 loss)

$$\Leftrightarrow \frac{p(x|H_0)}{p(x|H_A)} > \frac{1-\lambda}{\lambda} = c. \text{ which means:}$$

assign  $x \xrightarrow{d} H_0$  when  $\frac{p(x|H_0)}{p(x|H_A)} > c$ , get

$d(x)$  is Bayesian rule under 0-1 loss.

Remark: For protecting  $H_0$ , we can assign large " $\lambda$ " to  $p(H_0)$ .

Alternative proof of N-P Lemma:

$d^*$ : test accept  $H_0 \Leftrightarrow \frac{f_0(x)}{f_1(x)} > c$ , with level  $\alpha^*$ .

Then  $d^*$  is the most power test of level  $\alpha \leq \alpha^*$

Pf: let  $c^* = \frac{\alpha}{1-\alpha}$ ,  $p(H_0) = \lambda$ ,  $p(H_A) = 1 - \lambda$ .

$\therefore d^*$  is the Bayesian rule with this prior, with 0-1 loss

$$\therefore \beta(d^*) - \beta(d) = \lambda [E_{x \sim L(H_0, d^*(x))} - E_{x \sim L(H_0, d(x))}]$$

$$+ (1-\lambda) [E_{x \sim L(H_A, d^*(x))} - E_{x \sim L(H_A, d(x))}]$$

$$= \lambda(\alpha^* - \alpha) + (1-\lambda) [E_{x \sim L(H_A, d^*(x))} - E_{x \sim L(H_A, d(x))}] \leq 0$$

$$\therefore E_{x \sim L(H_A, d^*(x))} \leq E_{x \sim L(H_A, d(x))}$$

$$\text{i.e. } \beta_{d^*} \geq \beta_d.$$