# Decision Theory and
# Bayesian Inference

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

(1) Definition:

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

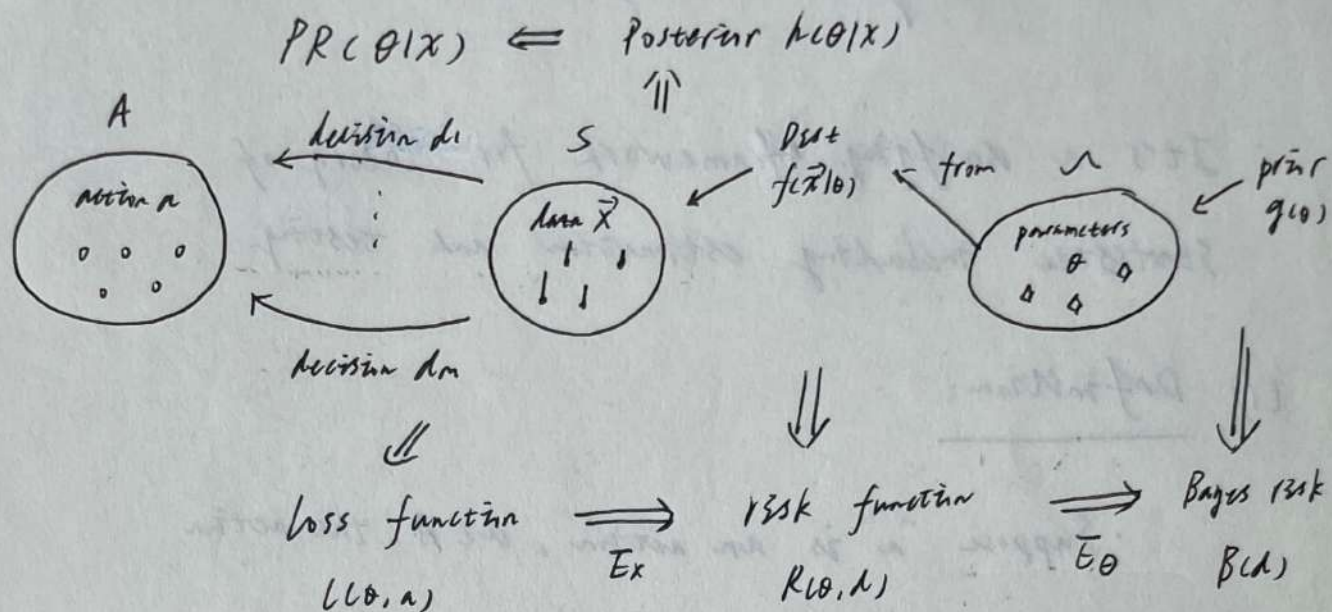$\Rightarrow$ The choice of action "a" depends on:

    i) Observations of r.v. : Data $X$ where $X \in S$. 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_X( l(\theta, d(\vec{X}))) $ (it depends on $\theta$)

$\Rightarrow$ Our aim is to minimize $R(\theta, d)$ by
choosing a good decision $d(\vec{x})$!

$$PR(\theta|x) \Longleftarrow \text{Posterior } h(\theta|x)$$



$$\underset{E_{x}}{\text{loss function}} \Longrightarrow \underset{R(\theta, d)}{\text{risk function}} \Longrightarrow \underset{\bar{E}_{\theta}}{\text{Bayes risk}}$$
$$L(\theta, a) \qquad\qquad R(\theta, d) \qquad\qquad B(d)$$

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

we choose $L(\theta, d(\vec{x})) = [V(\theta) - d(\vec{x})]^2$. quadratic loss func

## (2) Minimization:

· Difficulties of $\qquad$ i) $R(\theta, d)$ depend on unknown $\theta$.

minimizing $R(\theta, d)$ $\qquad$ ii) For different $\theta_1, \theta_2$. It may happen:
$$R(\theta_1, d_1) > R(\theta_1, d_2)$$ how to choose?
$$R(\theta_2, d_1) < R(\theta_2, d_2)$$

① Def of

minimax rule:

· Consider the worst case: $\sup\limits_{\theta \in \Lambda} R(\theta, d)$

$d^*$ is the minimax rule. If $\sup\limits_{\theta \in \Lambda} R(\theta, d^*)$

$= \inf\limits_{d} \sup\limits_{\theta \in \Lambda} R(\theta, d)$, $d^*$ may not exist!

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

③ Bayesian Rule:

' We assume $\theta \in \textcircled{H}$. 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 $\min_d \{B(d)\}$

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

$\Rightarrow$ Posterior Analyses:

A method for finding Bayesian Rule:

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

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

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

Def: Posterior risk: $E_{\theta|X}(L(\theta, d(x))) = \int R(\theta|x)$

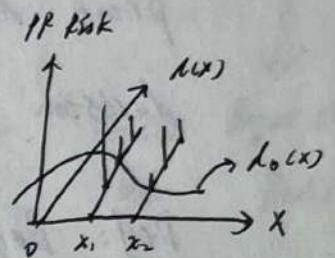Remark: The observed data $X=x$ updates the p-risk!

**Thm.** If $d_0(x)$ is a desicison func. minimizes the posterior risk for each $x$. Then $d_0$ is Bayesian Rule.

**Pf:**
$$B(d) = E_\theta ( R(\theta, d)) = E_\theta ( E_x ( l(\theta, d(x) | \theta)))$$

$$= \int [ \int l(\theta, d(x)) f_{x|\theta}(x) \, dx ] g_\theta(\theta) \, d\theta$$

$$=. \int [ \int l(\theta, d(x)) h_{\theta|x}(\theta) \, d\theta ] f_x(x) \, dx$$

$$= \int E_{\theta|x} ( l(\theta, d(x)) f_x(x) \, dx$$



Then $B(d)$ is minimized !

$\Rightarrow$ **Algoisthm:**

1°) Caculate $h(\theta|x)$. for each $x$.

2°) Caculate $E_{\theta|x} ( l(\theta, d(x)))$ for each $x$.

3°) Find $d(x) \in A$ minimizes every $PR(\theta|x)$. fixed $x_0$.

(3) **Application of Decision Theory:**

**Estimation**

① $\begin{cases} \text{Action space } A \longrightarrow \text{Parameter Space } \Omega. \\ \text{Decision Func. } d(x) \longrightarrow \text{estimator of } \theta. \\ l(\theta, d(x)) = [\theta - d(x)]^2 \text{ or } |\theta - d(x)| \end{cases}$

__Thm.__ i) $E_{\theta|x}((\theta - \hat{\theta})^2 | x) = Var_{\theta|x}(\theta | x) + [E_{\theta|x}(\theta|x) - \hat{\theta}]^2$

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.

__Pf:__ For ii) $E_{\theta|x}(|\theta - \hat{\theta}| | x) = \int |\theta - \hat{\theta}| h_{\theta|x} c_{\theta} d\theta$

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

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

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

③ __Def:__ i) $d_1, d_2$ two decision Functions $\in A$.
says $d_1$ __dominates__ $d_2$ : $R(\theta, d_1) \leq R(\theta, d_2)$. $\forall \theta \in$ ①

ii) $d$ a decision Func. $d$ is __admissable__. If $d$ is not strictly dominated by any other decision Func.

__Thm.__ i) If $\Lambda$ is discrete. $d^*$ is Bayesian rule w.r.t. prior pmf $g(\theta)$. where $g(\theta) > 0$. $\forall \theta$.

ii) If $\Lambda$ is dense. $d^*$ is Bayesian rule w.r.t prior pdf $g(\theta)$. $g(\theta) > 0$. $\forall \theta$. $R(\theta, d)$ is conti. of $\theta$. $\forall d$.

Then $d^*$ is admissible.

__Remark:__ It claims the relation between Bayesian Rule and admissible.

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

$\exists \theta_0 h$. $\theta \in (\theta_0 - h, \theta_0 + h) \Rightarrow R(\theta, \delta^*) > R(\theta, \delta') + \varepsilon$.

Then check: $B(\delta^*) - B(\delta') > 0$. contradict!

## (4) Bayesian View for prob. :

<u>( Personal Opinion )</u>

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.

$\Rightarrow$ <u>Difference between "B"</u>

<u>and "F" Approach</u> :

① <u>Point Estimation</u>:

$\begin{cases} F: \theta \text{ is fixed. unknown. not random} \longrightarrow \text{likelihood} \longrightarrow \text{maximal:} \\ \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad L(\theta|x) \qquad\qquad \text{Find MLE.} \\ \\ B: \theta \text{ has prior Asst } g(\theta) \longrightarrow \text{Posterior} \longrightarrow \text{Centering:} \\ \qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \text{Asst} \qquad\qquad \text{Find } E(\theta|x) \\ \qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad h(\theta|x) \end{cases}$

<u>Remark</u>: From $h(\theta|x) = \dfrac{f(x|\theta) g(\theta)}{\int f(x|\theta) g(\theta) d\theta}$ $\therefore h(\theta|x) \propto f(x|\theta) g(\theta)$

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

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

② Interval Estimation:

$\begin{cases}\text{Frequentist}: P(\theta \in [\theta_L(\vec{X}), \theta_U(\vec{X})] \mid \theta) = 1-\alpha. \ \vec{X} \text{ is random}, \theta \text{ fixed}. \\ \quad \text{If } \vec{x} \text{ is observed data, then } P(\theta \in [\theta_L(\vec{x}), \theta_U(\vec{x})]) = 0 \text{ or } 1. \\ \text{Bayesian}: P(\theta \in [\theta_L(\vec{x}), \theta_U(\vec{x})] \mid \vec{X} = \vec{x}) = 1-\alpha. \ \theta \text{ is random}.\end{cases}$

$\vec{x}$ is the observed data. fixed.

③ Testing:

$\begin{cases}\text{Frequentist}: \text{prob of Type I. II. error.} \\ \text{Bayesian}: \text{After observing data, compute posterior prob.}\end{cases}$

⑤ Bayesian Inference

for Normal Dist:

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

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

$$M_1 = \frac{\xi_0}{\xi + \xi_0} M_0 + \frac{\xi}{\xi + \xi_0} X. \quad \sigma_1^2 = \frac{1}{\xi + \xi_0}$$

$$\xi_0 = \frac{1}{\sigma_0^2}, \quad \xi = \frac{1}{\sigma^2}.$$

Pf: Since the const. is for normalization.
We just care the ratio (about $M$):

$$h(M|x) \propto f(x|M) g(M) \propto e^{-\frac{1}{2\sigma^2}(X-M)^2 - \frac{1}{2\sigma_0^2}(M-M_0)^2}$$

$$= e^{-\frac{1}{2}\left(M^2 \underbrace{[\frac{1}{\sigma^2} + \frac{1}{\sigma_0^2}]}_{=a} - 2M \underbrace{(\frac{x}{\sigma^2} + \frac{M_0}{\sigma_0^2})}_{=b} + \frac{x^2}{\sigma^2} + \frac{M_0^2}{\sigma_0^2}\right)}$$

$$\propto e^{-\frac{1}{2a^{-1}}(M - \frac{b}{a})^2}$$

**Cor.** For $n$ samples $\vec{X} = (X_1 \cdots X_n)$, i.i.d.

$$\mu \mid \vec{X} \sim N(\frac{s_0}{ns + s_0} \mu_0 + \frac{ns}{ns + s_0} \bar{x}, \frac{1}{ns + s_0})$$

**Rmk:** i) Note that $\frac{s_0}{ns + s_0} + \frac{ns}{ns + s_1} = 1$. We mix up the prior and the data to generate the posterior. (Weighted Average!)

ii) $\frac{1}{ns + s_0} < \sigma^2$. Which means the dist of $\mu \mid \vec{x}$ 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, non-informative!

(6) Bayesian Inference for Binomial Dist:

. $X \mid p \sim B(n.p)$. $p \sim Beta(a,b)$.

$\Rightarrow p \mid x \sim Beta(a+x, n+b-x)$

Similary. $\mu_{post} = \frac{a+b}{a+b+n} \cdot \frac{a}{a+b} + \frac{n}{a+b+n} \cdot \bar{x}$

If $n \to \infty$. $\mu_{post} \to \bar{x}$!

Remark: Define:

$$\begin{cases} G = \text{family of prior dist } g(\theta) \text{ for } \textcircled{\theta} \\ H = \text{family of conditional dist } f(x|\theta) \end{cases}$$
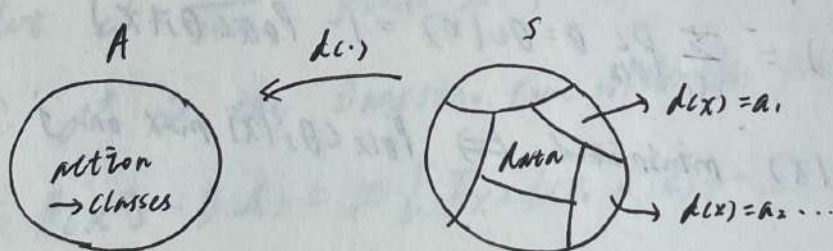
$\Rightarrow$ $G$ is called conjugate prior to $H$:
if the posterior of $G$ under $H$ also
belongs to $G$:

e.g. $G$: Normal dist. $H$: Normal dist $\longrightarrow G|H = G$

$H$: Binomial dist $G$: Beta dist $\longrightarrow G|H = G$

## (7) Application of

## Decision Theory:

### ① Classification:



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

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

$\Rightarrow$ For parameter $\theta$, let $A = \{\theta_i\}_1^m$.

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

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

Then, we can find Bayesian Rule:

$$p(\theta_i \mid x) = \frac{z_i f(x \mid \theta_i)}{\sum z_k f(x \mid \theta_k)} = P(\theta = \theta_i \mid X = x). \text{ Let } d(x) = a_j$$

We obtain $PR(j \mid x) = \sum l_{ij} p(\theta_i \mid x)$. Posterior risk

$\Rightarrow$ Choose $j$ to minimize $PR(j \mid x)$. for $\forall 1 \leq i \leq m$

Then $d: x \longmapsto a_j$ is a Bayesian Rule.

e.g. (0-1 loss)

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

$$R(i, d) = E_{x} \, l(\theta_i, d(x)) = \sum_j l_{ij} \, P_{\theta_i}(d(x) = j)$$

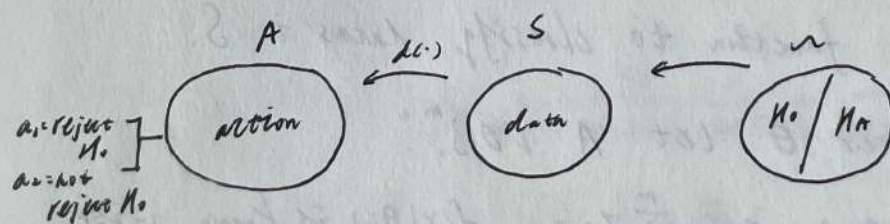$$= \sum_{j \neq i} P_{\theta_i}(d(x) = j) = 1 - P_{\theta_i}(d(x) = i)$$

(From frequentist approach, minimax it )

$\Rightarrow$ For Bayesian approach:

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

$\therefore PR(j \mid x)$ minimized $\Leftrightarrow P_{\theta \mid x}(\theta_j \mid x)$ max on $j$ !

② Hypothesis Testing:



$a_1$: reject $H_0$ ⎤
                  ⎦  (action)  ←  $d(\cdot)$  (data)  ←  $\sim$  ( $H_0 / H_n$ )
$a_2$: not reject $H_0$

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

$\Rightarrow$ **Bayesian approach:**

$p(H_0) = z$ . $p(H_A) = 1-z$ . then we obtain:

$$\frac{p(H_0|x)}{p(H_A|x)} = \frac{p(H_0)}{p(H_A)} \frac{p(x|H_0)}{p(x|H_A)} > 1 . \left(\text{LRS} = \frac{p_R(H_0|x)}{p_R(H_A|x)}\right.$$

$$\text{under } 0\text{-}1 \text{ loss}\Big)$$

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

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

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

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

**Alternative proof of N-P Lemma:**

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

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

Pf: let $c^1 = \frac{z}{1-z}$ . $p(H_0) = z$ . $p(H_A) = 1-z$ .

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

$$\therefore \beta(d^*) - \beta(d) = z\left[ E_x(l(H_0, d^*(x))) - E_x(l(H_0, d(x))) \right]$$

$$+ (1-z)\left[ E_x(l(H_A, d^*(x))) - E_x(l(H_A, d(x))) \right]$$

$$= z(\alpha^* - \alpha) + (1-z)\left[ E_x(l(H_A, d^*(x))) - E_x(l(H_A, d(x))) \right] \leq 0$$

$$\therefore E_x(l(H_A, d^*(x))) \leq E_x(l(H_A, d(x)))$$

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