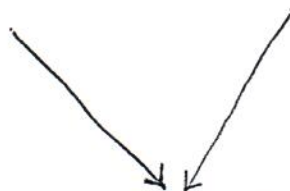# Statistical Decision Theory & Inference

Statistical Experiment(s)

$\downarrow$

Statistical Data          Decision Theory

State of nature
$\theta$ : parameter

Classical Statistics $\longrightarrow$ uses only sample information
collected through experiments;

but usually doesn't use extra evidence
or some other extra information. But
decision theory uses all kinds of informations.

Decision theory $\longrightarrow$ uses prior information
about $\theta$.

$\longrightarrow$ Consequence of using/taking
a particular decision

$\downarrow$

Technically known as
reward / loss ( in statistics)
utility ( in Economics)

The basic difference between the decision theory and the classical statistics is the loss function. In decision theory we are always guided by the loss (or reward) when we estimate $\theta$ by $\hat{\theta}$. More generally, we can view the whole statistics subject from a decision theoretic point of justifying suitable loss functions.
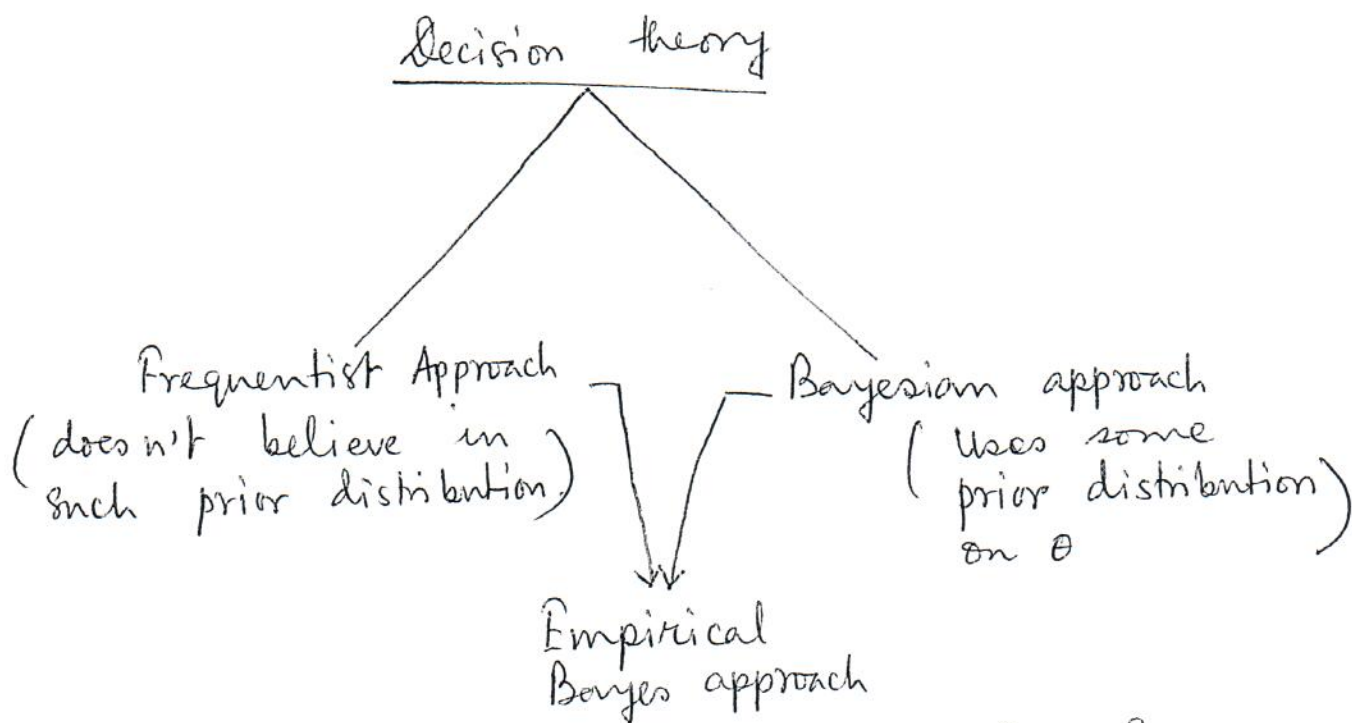
Suppose $X \sim N(\theta, 1)$, $\theta \in \mathbb{R}$, unknown. In classical statistics, UMVUE, MLE is $\hat{\theta} = X$. A decision theorist would ask why only $\hat{\theta} = X$? Why not try with $\hat{\theta}_c = cX$, $c \in (0, 1)$?

When we are trying with $\hat{\theta}$ to guess $\theta$, the error is $|\hat{\theta} - \theta|$. It may happen that for every sq. unit of deviation, we lose \$5. i.e., total loss, $L$ is

$$L = 5(\hat{\theta} - \theta)^2.$$

As a decision theorist, you should try to 'minimize' your loss $5(\hat{\theta} - \theta)^2$ choosing $\hat{\theta}$ suitably.

Well, we can view the use of UMVUE also from a decision theoretic point. If our loss is defined as $L = |bias| = |(E(\hat{\theta}) - \theta)|$, then there is no argument against using $X = \hat{\theta}$,

Broadly speaking, the decision theory is divided in two main schools.

$$\underline{\text{Decision theory}}$$

Frequentist Approach
(doesn't believe in such prior distribution)

Bayesian approach
(Uses some prior distribution on $\theta$)

Empirical Bayes approach

## An example where Bayesian Approach differs from all other methods.

Keep on tossing a coin. $X = \#$ heads.

$\theta = P(H)$, $\quad 1-\theta = P(T)$, $\quad 0 \le \theta \le 1$.

After 12 tosses, you get 9 Heads & 3 tails.

Q. Guess $\theta$, the probability of head.

Note: Do you really know the probability distri. of the data $(9-H, 3-T)$? The probability distribution depends heavily on the mechanism of the experiment.

Observe carefully, that we get 9-H, 3T because,

either (i) The experiment was stopped right after 12th toss

or (ii) the experiment was stopped right after 3rd T.

(i) $\implies$ $P(9\text{-}H, 3\text{-}T) = \binom{12}{9} \theta^9 (1-\theta)^3$, $X \sim Bin(12, \theta)$

(ii) $\implies$ $P(9\text{-}H, 3\text{-}T) = \binom{11}{2} \theta^9, (1-\theta)^3$, $X \sim Neg.Bin(12, 3, \theta)$

The Bayesians treat both the models (i) & (ii) in the same way and draws same conclusion regarding $\theta$; i.e., they are more interested in the part containing $\theta$, which is same for both the models. Where as frequentists or the classical statisticians treat both those two models differently.

Using classical statistics, if we want to test $H_0: \theta = \frac{1}{2}$ vs. $H_1: \theta > \frac{1}{2}$

model (i) $\implies$ ~~reject~~ Accept $H_0$ at 5% level

where as model (ii) $\implies$ reject $H_0$ at 5% level.

But if we apply (we will see) Bayesian techniques, we get the same answer from both the modes.

$X \sim NB(n, r, \theta)$ : # Successes before $r^{\underline{th}}$ failure

$= P(X=x) = \binom{n-1}{r-1} \theta^x (1-\theta)^r$ , $x = 0, 1, 2, \cdots\cdots$

$\qquad\qquad\qquad\qquad\qquad\qquad n = x + r$

$E(x) = \dfrac{rp}{q} = r\dfrac{\theta}{(1-\theta)}$

(ii) $\quad n = 12$
$\qquad r = 3 \qquad E(x) = 3\dfrac{\theta}{(1-\theta)}$

$\Longrightarrow \quad 3\left(\dfrac{\hat{\theta}}{1-\theta}\right) = 9$

$\Longrightarrow \quad \dfrac{\hat{\theta}}{1-\hat{\theta}} = 3$

Where as

$X \sim B(n, \theta)$

$E(x) = n\theta$

(i) $\qquad E(x) = 12\theta$

$\Longrightarrow \dfrac{x}{12} = \hat{\theta}$

$\Longleftrightarrow \hat{\theta} = 3/4$

is unbiased

$\Longrightarrow \quad \hat{\theta} = 3 - 3\hat{\theta}$

$\Longrightarrow \quad 4\hat{\theta} = 3$

$\Longleftrightarrow \quad \hat{\theta} = 3/4$

NOT unbiased

---

(i) Test $\quad H_0 : \theta = 1/2 \qquad vs \qquad H_1 : \theta \neq 1/2$

$\Lambda = \ln \dfrac{\sup\limits_{H_0} L(\theta|x)}{\sup\limits_{H_0 \cup H_1} L(\theta|x)} = \ln \dfrac{\left(\frac{1}{2}\right)^{12}}{\left(\frac{x}{n}\right)^x \left(1-\frac{x}{n}\right)^{n-x}}$

$\left(\dfrac{\hat{\theta}}{1-\theta}\right) = \dfrac{\frac{x}{n}}{1-\frac{x}{n}} = \left(\dfrac{x}{n-x}\right)$

If $\quad \Lambda > c \qquad$ then accept $H_0$
$\qquad\quad \leq c \qquad$ then rej $H_0$

i.e,

The basic concepts of decision theory was derived from the mathematical Game theory which in the simplest situation consists of two players — Player-I and player-II. In game theory both the players are taking actions against another & the actions are called strategies. Each player has its own set of strategies.

Let, the set of strategies available to P-I be $\Theta$, while $\mathcal{A}$ will stand for the set of possible strategies for P-II. If P-I chooses strategy $\theta \in \Theta$ and P-II chooses strategy $a \in \mathcal{A}$, there is a payoff $L(\theta, a)$ for P-I. This is called a two-person <u>zero sum</u> game. It can be assumed that $L(\theta, a) \geqslant K$ (for some finite constant K)

<u>Example</u> : Two contestants simultaneously put up either one or two fingers. P-I wins if the sum of the digits showing is odd and if it is even, P-II wins. The winner receives in $ the sum of the digits showing, this being paid by the loser.

So, $\Theta = \{1, 2\}$, $\mathcal{A} = \{1, 2\}$

| P-I $\Theta$ \ P-II $\mathcal{A}$ | 1 | 2 |
|---|---|---|
| 1 | −2 | 3 |
| 2 | 3 | −4 |

In decision theory, the role of player-I is played by the nature and P-II is the statistician. The set of all strategies for the nature is $\Theta$, called the <u>parameter space</u> and the $A$, the set of statistician's strategies is called the action space. Here, when nature chooses $\theta \in \Theta$ and the statistician chooses $a \in A$, the statistician loses $L(\theta, a)$. Without loss of generality, assume $L(\theta, a) \geq 0 \; \forall a \in A$, and $\theta \in \Theta$.

So, $\quad L : \Theta \times A \longrightarrow \mathbb{R}^+$.

Hence, we identify a statistical decision theory by the triplet $(\Theta, A, L)$.

<u>Major differences between the Game theory and the decision theory.</u>

(1) The game theory, both the players are rational (assumed) and both are trying to maximize their gain (or minimizing loss) simultaneous. A common approach is the 'minimax' approach where each player is guarding against the worst situation.

$$P-II \longrightarrow \quad \min_{a \in A} \; \max_{\theta \in \Theta} L(\theta, a) = \overline{V}$$

$$P-I \longrightarrow \quad \max_{\theta \in \Theta} \; \min_{a \in A} L(\theta, a) = \underline{V}$$

Usually $\quad \underline{V} \leq \overline{V}$

If there exists $a_0$ and $\theta_0$ $\quad \exists$

$$\max_{\theta \in \Theta} L(\theta, a_0) = \min_a \max_\theta L(\theta, a)$$

$$= \max_\theta \min_a L(\theta, a) = \min_a L(\theta_0, a)$$

i.e, $\quad \underline{V} = \overline{V} = V$ we say that the game (say)

is determined $\quad$ and $\quad$ the value of the game

is $V$.

In decision theory, the nature is not 'rational'. It chooses $\theta \in \Theta$ arbitrarily, or it doesn't plan to maximize $L(\theta, a)$. Its the statistician's turn to minimize $L(\theta, a)$. In this sense, it is a 'kind of' one sided game (though the minimax principle can be applied).

(2) $\quad$ There is another constraint for the statistician. Before, taking any decision, he performs a statistical experiment whose outcome is random, $X$ and depending on $X$, he chooses an optimal action.

More over, ~~the di~~ the outcome $X$ has a probability distribution $P_\theta$, $\theta \in \Theta$ which depends on nature's choice $\theta$ ; $\quad$ i.e., the

gets some information about the state of the nature. So, the statistical decision theory, strictly speaking, is not a fair game.

Let $\mathcal{X}$ be the sample space i.e., the outcome $X$ of the stat. exp. takes values in $\mathcal{X}$. Instead of using $a$, to denote the action of the statistician, we use $\delta(X)$ which depends on $\theta$. So our loss function is

$$L(\,:\,\delta(X), \theta) \geqslant 0$$

and let $\mathscr{D} = \{$ set of all possible $\delta$'s $\}$. $\delta(X)$ is called a '<u>rule</u>' (decision rule).

<u>Basic assumption</u> : $L(\delta(X), \theta)$ is convex in $(\delta(X) - \theta)$.

Our goal is to choose $\delta(X)$ such that $L(\delta(X), \theta)$ is "small". But $L(\delta(X), \theta)$ depends on $X$, so the loss is again a "random quantity". So, we try to minimize average loss, called risk, defined as

$$R(\delta, \theta) = E_\theta\, L(\delta(X), \theta).$$

So, to minimize $R(\delta, \theta)$, it's not a question of choosing a particular value $\delta(X)$, its the question of choosing the structure or functional form $\delta$ of the decision rule.

<u>Loss functions</u> ;   (Two general loss functions)

1. <u>Squared error loss ( quadratic loss)</u>

The loss function $\bcancel{L(\theta, \delta)} = \bcancel{(\theta-\delta)}$

$$L(\delta(x), \theta) = (\delta(x) - \theta)^2 \quad \text{is called}$$

squared error loss.

Why should we use squared error loss ?
(SEL)

(i)   SEL was first introduced in estimation problems
when unbiased estimators of $\theta$ being considered,
since, $R(\delta, \theta) = E(\delta(x) - \theta)^2$ would then be the
variance; otherwise, its simply MSE.

(ii) Another reason for the popularity of SEL is
due to the relationship to classical least
squares theory which in turn connects the
maximum likelihood estimation method under normality

(iii) Finally, the SEL makes the calculations much
easier.

---

In the multivariate setup when we estimate
$\theta = (\theta_1, \ldots, \theta_p)'$, by $\underset{\sim}{\delta}(\underset{\sim}{x}) = (\delta_1(\underset{\sim}{x}), \ldots, \delta_p(\underset{\sim}{x}))$,
the natural generalization is

$$L(\underset{\sim}{\delta}(\underset{\sim}{x}), \underset{\sim}{\theta}) = \| \underset{\sim}{\delta}(\underset{\sim}{x}) - \underset{\sim}{\theta} \|^2$$
$$= (\underset{\sim}{\delta}(\underset{\sim}{x}) - \underset{\sim}{\theta})'(\underset{\sim}{\delta}(\underset{\sim}{x}) - \underset{\sim}{\theta}).$$

(iv) Another justification for the squared error loss is :—
Let the loss be a fn. of the deviation ; i.e.,

$$L(\delta, \theta) = g(\delta - \theta) \quad \text{(by Taylor's expansion series)}$$
$$= g(0) + (\delta - \theta) g'(0) + \tfrac{1}{2}(\delta - \theta)^2 g''(0)$$

(neglecting the higher order terms of $(\delta - \theta)$ )

$$\implies L(\delta, \theta) \approx K_1 \left((\delta - \theta) + K_2\right)^2 + K_3$$

for suitable constants $K_1, K_2, K_3$

So, basically $L(\delta, \theta) \propto (\delta + c - \theta)^2$

if we redefine $\delta^* = \delta + c$

then minimizing $L(\delta, \theta)$ w.r.t. $\delta$ is same as minimizing $L^*(\delta^*, \theta) = (\delta^* - \theta)^2$ which brings us

to SEL again.
( See Page – 60 , Berger )

---

2. ## 0 - 1 Loss

In 0-1 loss, the decision space $\mathcal{A}$ consists of two elements ; i.e., $\mathcal{A} = \{\delta_0, \delta_1\}$. We want to guess $\theta$, $\theta \in \Theta$. More over $\Theta = \Theta_0 \cup \Theta_1$ (disjoint)
The loss is described as follows –

$$\mathcal{X} = \mathcal{X}_0 \cup \mathcal{X}_1$$

we take $\delta_0$ if $X \in \mathcal{X}_0$ & $\delta_1$ if $X \in \mathcal{X}_1$.

$$L(\delta_0, \theta) = \begin{cases} 0 & \text{if} \quad \theta \in \Theta_0 \\ 1 & \text{if} \quad \theta \in \Theta_1 \end{cases}$$

$$L(\delta_1, \theta) = \begin{cases} 0 & \text{if} \quad \theta \in \Theta_1 \\ 1 & \text{if} \quad \theta \in \Theta_0 \end{cases}$$

So, the risk fn. is given as

$$R(\delta_0, \theta) = E\, L(\delta_0, \theta) = P(x \in \mathcal{X}_0 \mid \theta \in \Theta_1)$$
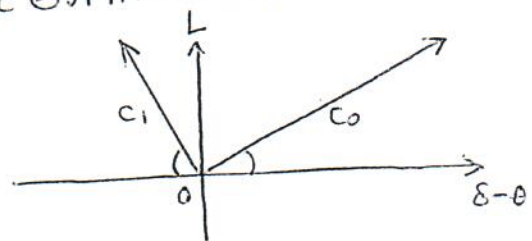$$= \text{Type-II error} = 1 - \text{Power}.$$

$$R(\delta_1, \theta) = E\, L(\delta_1, \theta) = P(x \in \mathcal{X}_1 \mid \theta \in \Theta_0)$$
$$= \text{Type-I error} = \text{size/level}.$$

This is infact the classical hypothesis testing.
(See Ferguson, 198 )

---

Apart from the squared error loss, in point estimation, another common loss is 'Linear loss' given as —

$$L(\delta, \theta) = \begin{cases} c_0 (\delta - \theta) & \text{if} \quad \delta \geqslant \theta \\ c_1 (\theta - \delta) & \text{if} \quad \delta < \theta \end{cases}$$

The constants $c_0$ & $c_1$ can be chosen to reflect the relative importance of over estimation or underestimation.



If $c_0 = c_1 = c$,

it is absolute error loss $\propto |\delta - \theta|$.

There are several other loss functions, like entropy losses, invariant squared error losses etc. which will be discussed when necessary.

---

Criticisms: The use of loss functions is often criticised for

(i) being inappropriate for the inference problem

(ii) Even if one derives a loss function from the economic utility function, it may be extremely difficult to handle.

(iii) Loss functions are 'non robust', in the sense that, ~~an~~ a "good" estimator for the loss $L_1$ may not be good for another loss $L_2$.

---

To tackle (i), we should concentrate on most intuitive loss functions; however we have to admit our limitations regarding (ii), But we can do something for (iii). We can look for estimators which perform more or less "well" under various loss functions and are called robust estimators — a major field of research.

# Admissibility & Completeness

__Def 1.__ (a) A decision rule $\delta_1$ is said to be __as good as__ a rule $\delta_2$ if $R(\delta_1, \theta) \leq R(\delta_2, \theta)$ $\forall \theta \in \circledR$.
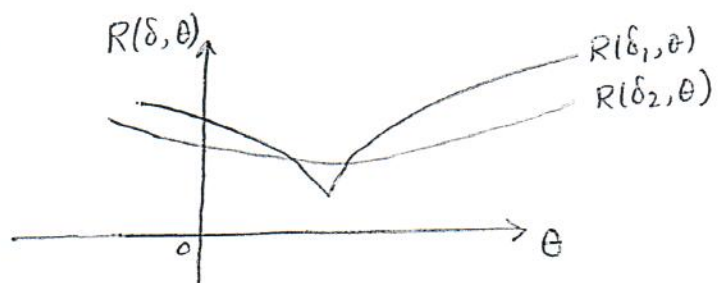
(b) A rule $\delta_1$ is said to be __better than__ a rule $\delta_2$ if $R(\delta_1, \theta) \leq R(\delta_2, \theta)$ $\forall \theta \in \circledR$ and $R(\delta_1, \theta) < R(\delta_2, \theta)$ for some $\theta_1$

(Notation : $\delta_1 \succ \delta_2$)

(c) A rule $\delta_1$ is said to be __equivalent__ to a rule $\delta_2$ if $R(\delta_1, \theta) = R(\delta_2, \theta)$ $\forall \theta \in \circledR$.
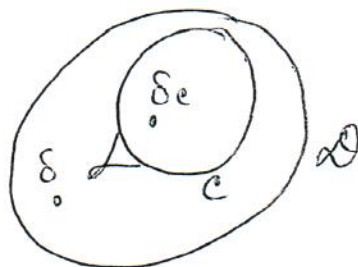
A above definitions give a partial ordering of the space $\mathscr{D}$ of decision rules.

__Def - 2.__ A rule $\delta$ is said to be __admissible__ if there exists no rule better than $\delta$.

A rule is said to be __inadmissible__ if it is not admissible.

Admissibility is an optimum property, although in a very weak sense. Conversely, we could never feel very proud about a rule if it is inadmissible.

<u>Def-3</u>. (a) A class $C$ of decision rules, $C \subset \mathscr{D}$, is said to be complete, if, given any $\delta \in \mathscr{D}$, $\delta \notin C$, $\exists$ a rule $\delta_c \in C \ni \delta_c$ is better than $\delta$.



(b) A class $C_*$ of decision rules is said to be essentially complete, if, given any rule $\delta \in \mathscr{D}$, $\delta \notin C$, $\exists$ a rule $\delta_* \in C_* \ni \delta_*$ is as good as $\delta$

<u>Lemma 1</u>. If $C$ is a complete class and $A$ denotes the class of all admissible rules, then $A \subseteq C$.

<u>Pf</u>:- (EX)

<u>Lemma 2</u>. If $C_*$ is an essentially complete class and there exists an admissible rule $\delta \notin C_*$, then $\exists \ \delta' \in C_*$ which is equivalent to $\delta$.

<u>Pf</u>:- (EX)

<u>Def-4</u>. A class $C_0$ of decision rules is said to be <u>minimal complete</u> if $C_0$ is complete and if no proper subclass of $C$ is complete; i.e., $C_0$ is the smallest complete class.

Similarly we can define minimal essentially comple class.

Theorem 1. If a minimal comple class exists, it consists of exactly the admissible rules.

Pf:− ( See Ferguson, p−56).

## Minimaxity

Def 5. (a) A decision rule $\delta_0$ is said to be minimax

if $$\sup_{\theta \in \Theta} R(\delta_0, \theta) = \inf_{\delta \in \mathscr{D}} \sup_{\theta \in \Theta} R(\delta, \theta).$$

The value on the right hand side is called the minimax value.

(b) Let $\epsilon > 0$. A decision rule $\delta_0$ is said to

be $\epsilon -$ minimax if

$$\sup_{\theta \in \Theta} R(\delta_0, \theta) \leq \inf_{\delta \in \mathscr{D}} \sup_{\theta \in \Theta} R(\delta, \theta) + \epsilon$$

## Geometric interpretation:

First we define the risk set $S$

$$S = \{ (\ldots, R(\delta, \theta), \ldots ) \neq \theta \in \Theta \mid \delta \in \mathscr{D} \}$$

If $\Theta$ is finite, say $\Theta = \{ \theta_1, \theta_2 \}$

then, $$S = \{ (R(\delta, \theta_1), R(\delta, \theta_2)) \mid \delta \in \mathscr{D} \}$$

for rule $\delta$, $$\sup_{\theta \in \Theta} R(\delta, \theta) = R(\delta, \theta_1) \vee R(\delta, \theta_2)$$