

Statistical

Inference.

Q. What is the average height of a BSDS
student? population.

We observe the heights of the current batch
of students. sample,

Goal: To make inference about the parameter
based on the sample at hand.

To connect the gap between the sample and
the population, we make use of statistical
models.

Assume: The height of a BSDS student
is H . $H \sim N(\mu, \sigma^2)$.

The observations we have

$$H_1, H_2, \dots, H_{62} \stackrel{\text{iid}}{\sim} H$$

$$H_1, H_2, \dots, H_{62} \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2).$$

Statistical
modeling

Average height of a BSDS student

$$\mathbb{E}(H) = \mu$$

We want to estimate μ based on H_1, \dots, H_{62} .

Eg. What is the average lifetime of a Phillips LED bulb?

parameters of interest

population

If Phillips gives a one-year warranty on its bulbs, what proportion of bulbs need to be replaced?

parameters of interest

Fix a time period T (e.g. 1000 hours).

Let 500 bulbs run for T hours.

If a bulb is fused, we know when it got fused.

If it is not fused, then we only know that it lasted for T hours.

The Lifetime of a Phillips LED bulb is L

$L \sim \text{Exponential} (\text{mean} = \lambda)$.

We observe

$$\tilde{L}_1, \dots, \tilde{L}_{500}$$

where

$$\tilde{L}_i = \begin{cases} L_i & \text{if } L_i \leq T \\ T & \text{if } L_i > T \end{cases} \quad \text{Truncated at } T.$$

$\tilde{L}_1, \dots, \tilde{L}_{500}$ are iid.

They have the same distribution as L truncated at T .

The distribution of \tilde{L} involves λ and T

From this, we can get to λ .

In statistical inference,

we have a population. (an infinite one).

We have a sample from the population
(finite).

We postulate certain (statistical) model
on the population.

Eg. $H \sim N(\mu, \sigma^2)$, $L \sim \text{Exp}(\text{mean} = \lambda)$

We also make certain assumptions on the
sample collection mechanism.

Eg. $H_1, \dots, H_{50} \stackrel{\text{iid}}{\sim} H$, $\tilde{L}_1, \dots, \tilde{L}_{50} \stackrel{\text{iid}}{\sim} \tilde{L} = \begin{cases} L & \text{if } L \leq T \\ T & \text{o.w.} \end{cases}$

Our task is to make inference about the unknown parameters of the population based on the sample.

Parametric model: When we assume the form of the population distribution to be known (except for the unknown parameters).

Nonparametric model: We do not assume the form of the distribution to be known.

Eg. $X \sim f$ where f is continuous.

$X \sim f$, and $E(X)$ exists
 $E(X^n)$ exists.

We will focus on parametric models.

Throughout the course, we will assume we have observations

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f_{\theta}$$

where f_{θ} is a known parametric form.

The form of f_{θ} is known upto the parameter θ .
We know f_{θ} completely as soon as we know θ .

Eg. $f_{\theta} = N(\mu, \sigma^2)$, $\theta = (\mu, \sigma^2)$

$$f_{(\mu, \sigma^2)} = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}, x \in \mathbb{R}$$

Eg. $f_\theta = \text{Exponential}(\text{mean} = \lambda)$, $\theta = \lambda$

$$f_\lambda = \frac{1}{\lambda} e^{-\frac{x}{\lambda}}, \quad x > 0.$$

Eg. $f_\theta = \text{Binomial}(n, p)$ $\theta = p$

$$f_p = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n$$

We have a parametric model

The diagram shows a horizontal axis with two 'X' marks. Above the axis, the letters 'iid' are written under a wavy line. Below the axis, there is a bracket spanning both 'X' marks, with the symbol f_θ written above it. To the right of the bracket, the text 'theta is the unknown parameter.' is written.

We have a parametric model

$X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta$. θ is the unknown parameter.

We want to make inference about θ .

- Estimation: Estimate the value of θ

based on X_1, \dots, X_n

want a single value
(point estimation)

want an interval where
 θ lies with high probability
(interval estimation)

- Testing: We want to check whether θ

prescribes to some preconceived values that we have in mind,

Eg. X_1, \dots, X_n $\stackrel{iid}{\sim}$ Binomial (S, p)

We may want to estimate p .

$$\frac{\bar{X}}{S}$$

$$\left(\frac{\bar{X}}{S} - c, \frac{\bar{X}}{S} + c \right)$$

We may want to test whether the value of p equals 0.1

$$H_0: p = 0.1 \quad vs \quad H_1: p \neq 0.1$$

$$H_0: p < 0.1 \quad vs. \quad H_1: p > 0.1$$

In our statistical model, the values of the parameters can vary within certain sets.

These are known as the parameter space.

Eg. $f_\theta = N(\mu, \sigma^2)$. $\mu \in \mathbb{R}$
 $\sigma > 0$

Eg. $F_\theta = \text{Exp}(\text{mean}=1)$, $\lambda > 0$.

Eg. $F_\theta = \text{Bin}(5, p)$, $p \in [0, 1]$

In this course, we will focus on optimal
inferential procedures.

Note: Optimality of the procedure comes as a
result of our modeling assumption.

The same procedure which is optimal under one
setup may not be good under a different setup.

Grading

Mid - term - 30

End - term - 50

2x class tests - 20

100

References:

- Statistical Inference - Casella & Berger
- Introduction to the Theory of Statistics - Mood, Graybill & Boes
- Introduction to Mathematical Statistics - Hogg, McKean, Craig
- Theory of Point Estimation - Lehmann & Casella
- Testing Statistical Hypotheses - Lehmann & Romano

$X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta$, θ unknown parameter
 $\theta \in \mathbb{H}$ \leftarrow parameter space.

We want to do inference on θ .

(a) To estimate θ

(b) Testing hypotheses involving θ .

X_i 's take values in $\mathcal{X} \leftarrow$ sample space.

Eg. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$, $p \in [0,1]$

$X_i \in \{0,1\} \leftarrow \mathcal{X}$

\mathbb{H}

To estimate p .

To test. $p = \frac{1}{2}$. vs. $p \neq \frac{1}{2}$.

$p > \frac{1}{2}$ / $p < \frac{1}{2}$.

Eg. $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, $\mu \in \mathbb{R}$, $\sigma^2 > 0$

$$\underline{\theta} = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_+$$
$$(-\infty, \infty) \times (0, \infty)$$

$$X_i \in \mathbb{R}$$

$$X_i \in (-\infty, \infty) \subset \mathcal{X}.$$

(H)

To estimate $\underline{\theta}$: μ, σ^2 .

To test $\underline{\theta} = \underline{\theta}_0 \left(\begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix} = \begin{pmatrix} \mu_0 \\ \sigma_0^2 \end{pmatrix} \right)$ vs. $\underline{\theta} \neq \underline{\theta}_0$.

To test $\mu = \mu_0$ vs $\mu \neq \mu_0$

In this case, σ^2 is a nuisance parameter.

$$X_1, \dots, X_n \stackrel{iid}{\sim} f_{\theta}, \quad \theta \in \mathbb{H}$$

The first step is to summarize the data.
data reduction.

Eg. $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ ← sample mean

A reduction from \mathbb{R}^n to \mathbb{R} .

Eg. $s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ ← sample variance

A reduction from \mathbb{R}^n to \mathbb{R}_+

Eg. The minimum or the maximum values.

$$X_{(1)} = \min_{i=1, \dots, n} X_i, \quad X_{(n)} = \max_{i=1, \dots, n} X_i$$

Eg. the middlemost value / median

$$\frac{x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n}{2}+1\right)}}{2} \quad \text{if } n \text{ even}$$

$$x_{\left(\frac{n+1}{2}\right)} \quad \text{if } n \text{ odd.}$$

Eg. Skewness, Kurtosis. . - - -

How to choose the summary?

Statistic: A statistic $T = T(\underline{x})$ is a measurable function of the sample.

— Something that can be computed based on the sample
(it should only be based on the sample)

Measurable function: A random function $T(\underline{X})$ is measurable if we can assign probabilities to it.

In particular, $\mathbb{P}_\theta(T(\underline{X}) \in O)$ is defined for every open set O .

We can also define $\mathbb{P}_\theta(T(\underline{X}) \in (a, b))$

$\mathbb{P}_\theta(T(\underline{X}) \in (-\infty, t])$ is well defined

$= \mathbb{P}_\theta(T(\underline{X}) \leq t) \leftarrow \text{cdf of } T(\underline{X})$

$[a, b], (a, b], [a, b]$
 $(-\infty, b], (-\infty, b)$
 $(a, \infty), [a, \infty)$

- Usually, most functions are measurable. But not all.
- All the functions that we will encounter in this course are ^{measurable}.

x_1, \dots, x_n : sample



$T = T(\underline{x}) = T(x_1, \dots, x_n)$: summary statistic.
(can be a vector)

Starting from the sample we can compute the statistic.

But, usually, not the other way round.

Which summary to use?

Depends on the modelling assumption.

Modelling assumption: $X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta$

We write $p_\theta(\underline{x}) = p_\theta(x_1, \dots, x_n)$: joint dist of the sample

$$= \prod_{i=1}^n f_\theta(x_i)$$

Idea: Choose the summary statistic in such a way that no information on θ is lost in the process.

- The distribution of the sample given the value of the summary statistic does not involve θ .
- The conditional distribution of \underline{X} given $T(\underline{X})$ is free of θ .

Def: We say that a statistic $T(\underline{x})$ is sufficient for θ if the conditional distribution of \underline{x} given $T(\underline{x})$ is free of θ .

- Knowing the statistic is sufficient to get all the information about θ .

Eg. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$, $p \in [0, 1]$

$$f_{\theta}(\underline{x}) = \prod_{i=1}^n f_{\theta}(x_i) \quad | \quad f_{\theta}(x_i) = \begin{cases} p & \text{if } x_i = 1 \\ 1-p & \text{if } x_i = 0 \end{cases}$$

$$= \prod_{i=1}^n \left\{ p^{x_i} (1-p)^{1-x_i} \right\}$$

$$= p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i}$$

$$T(\underline{x}) = \sum_{i=1}^n x_i \sim \text{Binomial}(n, p)$$

$$P_\theta(\underline{x} = \underline{x} | T(\underline{x}) = t) = \frac{P_\theta(\{\underline{x} = \underline{x}\} \cap \{T(\underline{x}) = t\})}{P_\theta(T(\underline{x}) = t)}$$

$$P_\theta(T(\underline{x}) = t) = \begin{cases} \binom{n}{t} p^t (1-p)^{n-t}, & t = 0, 1, \dots, n. \\ 0 & \text{otherwise.} \end{cases}$$

$$P_\theta(\{\underline{x} = \underline{x}\} \cap \{T(\underline{x}) = t\}) = \begin{cases} 0 & \text{if } T(\underline{x}) \neq t \\ P_\theta(\underline{x} = \underline{x}) & \text{if } T(\underline{x}) = t \end{cases}$$

$$P_\theta(\underline{x} = \underline{x}) = P_\theta(\underline{x}) = p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i}$$

Numerator = $\begin{cases} 0 & \text{if } T(\underline{x}) \neq t \\ p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i} & \text{if } T(\underline{x}) = t \end{cases}$

$| T(\underline{x}) = \sum_{i=1}^n x_i$

$$\text{Numerator} = \begin{cases} 0 & \text{if } T(\underline{x}) \neq t \\ b^t (1-b)^{n-t} & \text{if } T(\underline{x}) = t \end{cases}$$

$$P_{\theta}(\underline{X} = \underline{x} \mid T(\underline{X}) = t) = \frac{\begin{cases} 0 & \text{if } T(\underline{x}) \neq t \\ b^t (1-b)^{n-t} & \text{if } T(\underline{x}) = t \end{cases}}{\binom{n}{t} b^t (1-b)^{n-t}}$$

$$\text{So, } T(\underline{X}) = \sum_{i=1}^n X_i \text{ is } \begin{cases} 0 & \text{if } T(\underline{x}) \neq t \\ \frac{1}{\binom{n}{t}} & \text{if } T(\underline{x}) = t \end{cases}$$

a sufficient statistic for free of b .
 b . Knowing the total no. of success is enough for b .

Eg. $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, 1), \mu \in \mathbb{R}$.

$$f_{\theta}(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2}}, x \in \mathbb{R}$$

$$\begin{aligned} p_{\theta}(\underline{x}) &= \prod_{i=1}^n f_{\theta}(x_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_i - \mu)^2} \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_i^2 + \mu^2 - 2\mu x_i)} \\ &= \left(\frac{1}{\sqrt{2\pi}} \right)^n e^{-\frac{1}{2} \sum_{i=1}^n x_i^2 - \frac{n\mu^2}{2} + \mu \sum_{i=1}^n x_i} \end{aligned}$$

$$T(\underline{x}) = \sum_{i=1}^n x_i \sim N(n\mu, n)$$

pdf of $T(\underline{x})$

$$q_{\theta}(t) = \frac{1}{\sqrt{2\pi} n} e^{-\frac{1}{2n} (t - n\mu)^2}$$

$$= \frac{1}{\sqrt{2\pi} n} e^{-\frac{1}{2n} (t^2 + n^2\mu^2 - 2n\mu t)}$$

$$= \frac{1}{\sqrt{2\pi} n} e^{-\frac{t^2}{2n} - \frac{n\mu^2}{2} + nt}$$

$$= \frac{1}{\sqrt{2\pi} n} e^{-\frac{t^2}{2n}}$$

The pdf of \underline{x} given $T(\underline{x})$

$$\tilde{p}_{\theta}(\underline{x} = \underline{x} | T(\underline{x}) = t) = \begin{cases} 0 & \text{if } T(\underline{x}) \neq t \\ \frac{p_{\theta}(\underline{x})}{q_{\theta}(t)} & \text{if } T(\underline{x}) = t \end{cases}$$

$$= \begin{cases} 0 & \text{if } T(\underline{x}) \neq t \\ \frac{\left(\frac{1}{\sqrt{2\pi}}\right)^n e^{-\frac{1}{2} \sum_{i=1}^n x_i^2 - \frac{n\mu^2}{2}}}{\sqrt{2\pi n}} & \text{if } T(\underline{x}) = t \end{cases}$$

free of μ

$T(\underline{x}) = \sum_{i=1}^n x_i$ is sufficient for μ

The same conclusion holds if the variance σ^2 is known (not necessarily 1)

$$T(\underline{x}) = \sum_{i=1}^n x_i$$

Result: A statistic $T = T(\underline{x})$ is sufficient for θ if the ratio $\frac{p_\theta(\underline{x})}{q_\theta(T(\underline{x}))}$ is free of θ

where $p_\theta(\underline{x})$ is the "joint" pdf/pmf of the sample \underline{X} and $q_\theta(t)$ is the pdf/pmf of the statistic $T(\underline{x})$.

We need to check this for all $\underline{x} \in \mathcal{X}$.

Result: (Factorization Theorem) Let $p_{\theta}(\underline{x})$ be the joint
pdf/pmf of the sample \underline{X} . A statistic $T = T(\underline{X})$
is sufficient for θ if and only if we can
factorize $p_{\theta}(\cdot)$ as

$$p_{\theta}(\underline{x}) = \underbrace{g_{\theta}(T(\underline{x}))}_{\text{involves } \theta} \underbrace{h(\underline{x})}_{\text{free of } \theta} \quad \forall \underline{x} \in \mathcal{X}.$$

only depends on
the statistic T

Eg. $X_1, \dots, X_n \stackrel{\text{i.i.d}}{\sim} \text{Bernoulli}(p)$

$$p_{\theta}(\underline{x}) = p^{\sum x_i} (1-p)^{n - \sum x_i} = \left(\frac{p}{1-p} \right)^{\sum x_i} (1-p)^n$$

$$p_\theta(\underline{x}) = \begin{cases} p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}, & x_i \in \{0,1\}, i=1, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

$$= g_\theta\left(\sum_{i=1}^n x_i\right) h(\underline{x})$$

where $g_\theta(t) = \left(\frac{p}{1-p}\right)^t (1-p)^n, t=0, 1, \dots, n$

$$h(\underline{x}) = \begin{cases} 1 & \text{if } x_i \in \{0,1\}, i=1, \dots, n \\ 0 & \text{o.w.} \end{cases}$$

$$T(\underline{x}) = \sum_{i=1}^n x_i \text{ is sufficient for } p.$$

$X_1, \dots, X_n \leftarrow$ sample

$p_\theta(\underline{x}), \theta \in \mathbb{H} \leftarrow$ statistical model

$T = T(\underline{X})$ is sufficient for θ if the conditional dist. of \underline{X} given $T(\underline{X})$ is free of θ .

Result: (Factorization Theorem) A statistic $T = T(\underline{X})$ is sufficient for θ if and only if we can factorize

$$p_\theta(\underline{x}) = g_\theta(T(\underline{x})) h(\underline{x}) \quad \forall \underline{x} \in \mathcal{X}, \theta \in \mathbb{H}$$

Proof: If part : We have to show :

if $p_\theta(\underline{x}) = g_\theta(T(\underline{x})) h(\underline{x}) \quad \forall \underline{x} \in \mathcal{X}$
Then $T(\underline{x})$ is sufficient for θ .

We have to show : $IP_\theta(\underline{X} = \underline{x} | T(\underline{X}) = t)$ is free of θ

$$IP_\theta(\underline{X} = \underline{x} | T(\underline{X}) = t) = \begin{cases} 0 & \text{if } T(\underline{x}) \neq t \\ \frac{IP_\theta(\underline{X} = \underline{x})}{IP_\theta(T(\underline{X}) = t)} & \text{if } T(\underline{x}) = t \end{cases}$$

$$\frac{IP_\theta(\underline{X} = \underline{x})}{IP_\theta(T(\underline{X}) = t)} = \frac{p_\theta(\underline{x})}{\sum_{y \in \mathcal{Y}: T(y)=t} p_\theta(y)}$$

$$\underline{y} \in \mathcal{Y} : T(\underline{y}) = t$$

$$\frac{P_{\theta}(\underline{x} = \underline{x})}{P_{\theta}(T(\underline{x}) = T(\underline{x}))} = \frac{p_{\theta}(\underline{x})}{\sum_{\underline{y} \in \mathcal{X}: T(\underline{y}) = T(\underline{x})} p_{\theta}(\underline{y})}$$

$$= \frac{g_{\theta}(T(\underline{x})) h(\underline{x})}{\sum_{\underline{y} \in \mathcal{X}: T(\underline{y}) = T(\underline{x})} g_{\theta}(T(\underline{y})) h(\underline{y})}$$

$$= \frac{g_{\theta}(T(\underline{x}))}{g_{\theta}(T(\underline{x})) \sum_{\underline{y} \in \mathcal{X}: T(\underline{y}) = T(\underline{x})} h(\underline{y})} \quad \leftarrow \text{free of } \theta$$

$\therefore T(\underline{x})$ is sufficient for θ .

Note: The derivation shows us

$$\text{IP}_\theta(\tau(\underline{x}) = t) = \sum_{\underline{y} \in \mathcal{X} : \tau(\underline{y}) = t} p_\theta(\underline{y})$$

$$= \sum_{\underline{y} \in \mathcal{X} : \tau(\underline{y}) = t} g_\theta(\tau(\underline{y})) h(\underline{y})$$

$$= g_\theta(t) \sum_{\underline{y} \in \mathcal{X} : \tau(\underline{y}) = t} h(\underline{y}).$$

Only if part: $T(\underline{x})$ is sufficient for θ

To show: $p_\theta(\underline{x}) = g_\theta(T(\underline{x})) h(\underline{n}) \quad \forall \underline{x} \in \mathcal{X}$

$T(\underline{x})$ is sufficient $\Rightarrow P_\theta(\underline{x} = \underline{x} | T(\underline{x}) = t)$ is free of θ

$$p_\theta(\underline{x}) = P_\theta(\underline{x} = \underline{x})$$

$$= \sum P_\theta(\underline{x} = \underline{x} | T(\underline{x}) = t) \cdot P_\theta(T(\underline{x}) = t)$$

$t : t = T(\underline{x}) \text{ for some } \underline{x} \in \mathcal{X}$

ie $t \in \text{range of } T$

$$= \sum_{t \text{ in range of } T} P_\theta(\underline{x} = \underline{x} | T(\underline{x}) = t) \cdot P_\theta(T(\underline{x}) = t)$$

$$P_{\theta}(\underline{X} = \underline{x} \mid T(\underline{X}) = t) = 0 \quad \text{if } T(\underline{x}) \neq t$$

$$\sum_{t \text{ in range } T} P_{\theta}(\underline{X} = \underline{x} \mid T(\underline{X}) = t) / P_{\theta}(T(\underline{X}) = t)$$

$$= \sum_{t : t = T(\underline{x})} P_{\theta}(\underline{X} = \underline{x} \mid T(\underline{X}) = t) / P_{\theta}(T(\underline{X}) = t).$$

$$= P_{\theta}(T(\underline{X}) = T(\underline{x})) / P_{\theta}(\underline{X} = \underline{x} \mid T(\underline{X}) = T(\underline{x}))$$

is free of θ $h(\underline{x})$

a function of θ and x only via $T(x)$

$$g_{\theta}(T(x)) \Rightarrow p_{\theta}(\underline{x}) = g_{\theta}(T(\underline{x})) h(\underline{x}).$$

Some observations

- Suppose we have a statistic $T = T(\underline{x})$

$$T: \mathcal{X} \rightarrow T(\mathcal{X})$$

For $t \in T(\mathcal{X})$, define $A_t = \left\{ \underline{x} \in \mathcal{X} : T(\underline{x}) = t \right\}$

\sim is an equivalence relation if

- $x \sim x$,
- $x \sim y, y \sim x$,
- $x \sim y, y \sim z, z \sim x$
- $T(\underline{x}) = T(\underline{x})$
- $T(\underline{x}) = T(\underline{y}) \Rightarrow T(\underline{y}) = T(\underline{x})$
- $T(\underline{x}) = T(\underline{y}), T(\underline{y}) = T(\underline{z}) \Rightarrow T(\underline{x}) = T(\underline{z})$
- $x \sim y$ if $T(\underline{x}) = T(\underline{y})$

So, \sim defined as $x \sim y$ if $T(x) = T(y)$
is an equivalent relation.

$$A_t = \left\{ \underline{x} : T(\underline{x}) = t \right\}$$

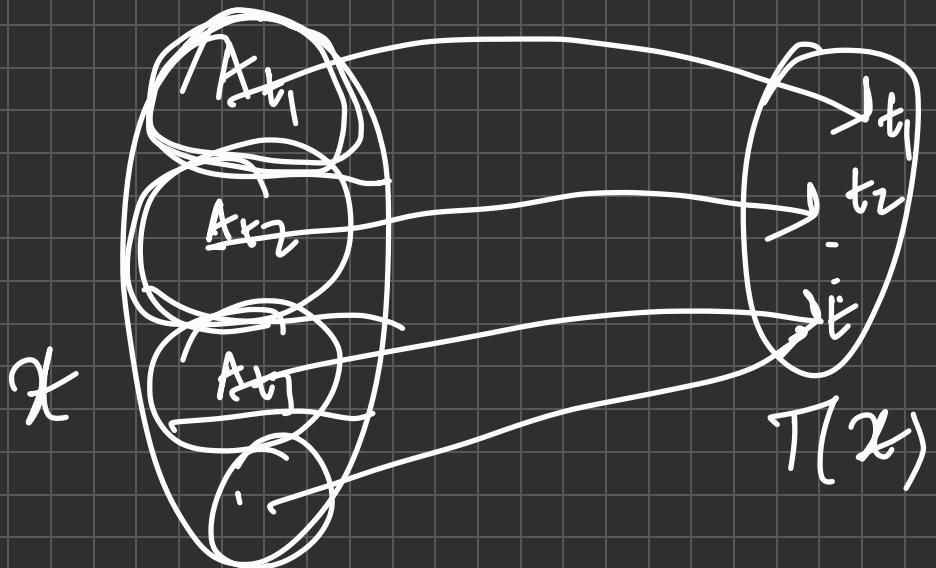
$\underline{x} \in A_{T(x)}$ \Rightarrow every $\underline{x} \in A_t$ for some $t \in T(\underline{x})$

$\underline{x} \sim \underline{y}$ and $\underline{x} \in A_t$, then $\underline{y} \in A_t$

Every \underline{x} can belong to only one A_t for some $t \in T(\underline{x})$

$\underline{x} \in A_t, A_s$ $t \neq s$.

The space X is partitioned by the sets A_t ,
 $t \in T(\underline{x})$.
These are known as equivalence classes.



$$X_1 - X_n \rightarrow X_1 + \dots + X_n$$

A statistic $T(x)$ is sufficient if it contains all the information about θ .

A sufficient statistic is also called a sufficient partition.

• $P_\theta(\underline{X} = \underline{x} \mid T(\underline{X}) = T(\underline{x}))$ is free of θ .

This is known to us.

We can generate observations from this distribution

Ex.

For the Bernoulli case

$$P_\theta(\underline{X} = \underline{x} \mid T(\underline{X}) = t) = \begin{cases} 0 & \text{if } \sum x_i \neq t \\ \frac{1}{\binom{n}{\sum x_i}} & \text{if } \sum x_i = t \end{cases}$$

$$P_\theta(\underline{X} = \underline{x} \mid T(\underline{X}) = T(\underline{x})) = \frac{1}{\binom{n}{\sum x_i}}, \quad x_i \in \{0, 1\}, \quad i=1, \dots, n$$

Generati $\underline{\tilde{X}}' \sim IP_{\theta} (\underline{\tilde{X}} = \underline{x} \mid T(\underline{\tilde{X}}) = T(\underline{x}))$

$$IP_{\theta} (\underline{\tilde{X}}' = \underline{x}) = \sum_{T(\underline{x})} IP_{\theta} (\underline{\tilde{X}}' = \underline{\tilde{x}} \mid T(\underline{\tilde{X}}') = T(\underline{x})) \cdot IP_{\theta} (T(\underline{\tilde{X}}') = T(\underline{x}))$$

$$= \sum_{T(\underline{x})} IP_{\theta} (\underline{\tilde{X}} = \underline{x} \mid T(\underline{\tilde{X}}) = T(\underline{x})) \cdot IP_{\theta} (T(\underline{\tilde{X}}) = T(\underline{x}))$$

$$= IP_{\theta} (\underline{\tilde{X}} = \underline{x})$$

In The Bernoulli example

$$P_{\theta}(\underline{x}' = \underline{x}' \mid T(\underline{x}) = t) = \begin{cases} 0 & \text{if } \sum_{i=1}^n x_i' \neq t \\ \frac{1}{\binom{n}{\sum_{i=1}^n x_i'}} & \text{if } \sum_{i=1}^n x_i' = t \end{cases}$$

From \underline{x} , we get $T(\underline{x}) = t_0$

$$\frac{1}{\binom{n}{\sum x_i'}}, \text{ if } \sum x_i' = t_0$$

\underline{x} ← Person 1 has access to it
 \sim can compute

Person 2 only has access to $T(\underline{x})$

Once the sample is observed. $\underline{X} = \underline{x} \rightarrow T(\underline{x})$

Person 1 has access to entire \underline{x}

Person 2 only gets to know $T(\underline{x})$

We have a statistical model $f_\theta(\cdot)$

$X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta(\cdot)$

$P_\theta(\underline{X} = \underline{x} \mid T(\underline{x}) = t)$ is free of θ
← completely known

Once we know t , $\Pr_{\theta}(\underline{X} = \underline{x} \mid T(\underline{X}) = t)$ is a distribution of \underline{X} . and it is completely known.

So, we can generate a random sample \underline{Y} from this distribution.

That is, $\Pr_{\theta}(\underline{Y} = \underline{y} \mid t) = \Pr_{\theta}(\underline{X} = \underline{y} \mid T(\underline{X}) = t)$

↑ conditional distribution depending on
the value of t .

Unconditional distribution of \underline{Y} .

$$\begin{aligned}
 \text{IP}(\underline{Y} = \underline{y}) &= \sum_t \underbrace{\text{IP}_t(\underline{Y} = \underline{y})}_{t} \text{IP}_\theta(t) \\
 &= \sum_t \text{IP}_\theta(\underline{X} = \underline{y} \mid T(\underline{X}) = t) \text{IP}_\theta(T(\underline{X}) = t) \\
 &= \text{IP}_\theta(\underline{X} = \underline{y}).
 \end{aligned}$$

So, the unconditional dist. of \underline{Y} is same as
 The unconditional dist. of \underline{X} .

Eg. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$

We have seen $T(\underline{x}) = \sum_{i=1}^n X_i$ is sufficient for θ .

\uparrow # of 1's in the sample

$$P_\theta(\underline{X} = \underline{y} \mid T(\underline{X}) = t) = \begin{cases} 0 & \text{if } \sum_{i=1}^n y_i \neq t \\ \frac{1}{\binom{n}{t}} & \text{if } \sum_{i=1}^n y_i = t \end{cases}$$

Once we observe $\underline{x} = (x_1, \dots, x_n)$

we can compute $t = \sum_{i=1}^n x_i$.

We know $P_\theta(\underline{X} = \underline{x} \mid T(\underline{X}) = t)$

Generate $\underline{y} = (y_1 \dots y_n)$ from this distribution -

$$P_t(y_1 \dots y_n) = \begin{cases} 0 & \text{if } \sum_{i=1}^n y_i \neq t \\ \frac{1}{\binom{n}{t}} & \text{if } \sum_{i=1}^n y_i = t \end{cases}$$

choose t places at random (out of the n places)

Put 1 in those t places.

Put 0 in the rest

This gives us $y_1 \dots y_n$

In a particular problem, There are multiple sufficient statistics.

Eg. $X_1 - X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$

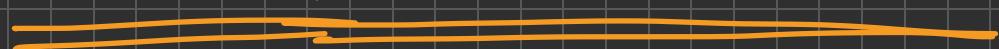
$(X_1 - X_n)$ is sufficient

$(X_{(1)}, \dots, X_{(n)})$ is sufficient

$\sum_{i=1}^n X_i$ is sufficient.

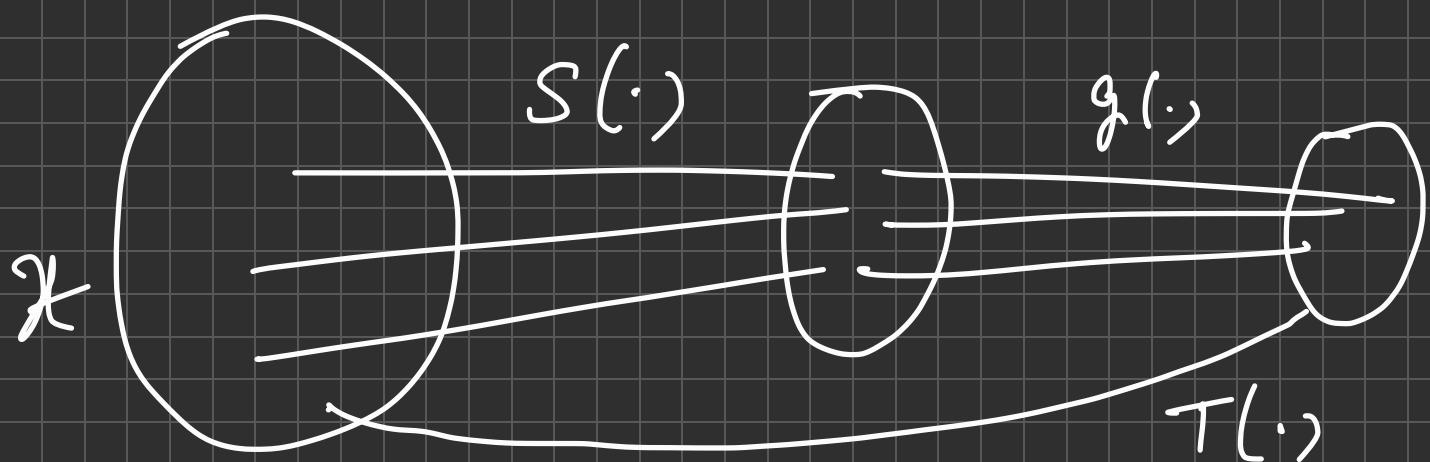
} which one
to choose?

We should go for the one which gives us
the most possible reduction.



Def: A statistic $T(\underline{x})$ is said to be minimal sufficient for a parameter θ if it is sufficient and for any sufficient statistic $S(\underline{x})$ of θ , $T(\underline{x})$ can be written as a function of $S(\underline{x})$.

- $T(\underline{x})$ is sufficient.
- If $S(\underline{x})$ is sufficient, then there exists g such that $T(\underline{x}) = g(S(\underline{x})). \forall \underline{x} \in \mathcal{X}$.



A minimal sufficient statistic is a possible reduction to every sufficient statistic.

Minimal sufficient statistics is not unique.

Results
Two statistics $T(\underline{x})$ and $S(\underline{x})$ are both minimal sufficient for θ if and only if there exist functions $g(\cdot)$ and $h(\cdot)$ such that $T(\underline{x}) = g(S(\underline{x}))$ and $S(\underline{x}) = h(T(\underline{x}))$.

Remark: Let $p_\theta(\underline{x})$ be the joint pdf / pmf of the sample. If $\frac{p_\theta(\underline{x})}{p_\theta(\underline{y})}$ is free of θ if and only if $T(\underline{x}) = T(\underline{y})$, then $T(\underline{x})$ is a minimal sufficient statistic for θ .

Eg. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$

$$\frac{p_\theta(\underline{x})}{p_\theta(\underline{y})} = \frac{\theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i}}{\theta^{\sum_{i=1}^n y_i} (1-\theta)^{n - \sum_{i=1}^n y_i}} = \theta^{\sum_{i=1}^n x_i - \sum_{i=1}^n y_i} (1-\theta)^{n - \sum_{i=1}^n x_i + \sum_{i=1}^n y_i}$$

is free of θ if and only if $\sum_{i=1}^n x_i = \sum_{i=1}^n y_i$.

So, $T(\underline{x}) = \sum_{i=1}^n x_i$ is minimal sufficient for θ .

Eg. $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, 1)$

$$p_\theta(\underline{x}) = \left(\frac{1}{\sqrt{2\pi}} \right)^n e^{-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2}, \quad x_i \in \mathbb{R}, \quad i=1, \dots, n$$

$$= \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \left\{ \sum_{i=1}^n x_i^2 + n\mu^2 - 2\mu \sum_{i=1}^n x_i \right\}}$$

$$\frac{p_\theta(\underline{x})}{p_\theta(\underline{y})} = e^{-\frac{1}{2} \left\{ \sum_{i=1}^n x_i^2 - \sum_{i=1}^n y_i^2 - 2\mu \left(\sum_{i=1}^n x_i - \sum_{i=1}^n y_i \right) \right\}}$$

This is free of μ if and only if $\sum_{i=1}^n x_i = \sum_{i=1}^n y_i$.

So, $\sum_{i=1}^n x_i$ is a minimal sufficient statistic for μ .

Ex. $X_1, \dots, X_n \stackrel{iid}{\sim} U(0, \theta), \theta > 0$

$$p_\theta(x) = \begin{cases} \frac{1}{\theta^n} & \text{if } 0 < x_i < \theta \quad \forall i \\ 0 & \text{o.w.} \end{cases}$$

$$= \begin{cases} \frac{1}{\theta^n} & \text{if } x_{(1)} > 0 \text{ and } x_{(n)} < \theta \\ 0 & \text{o.w.} \end{cases}$$

$$= \frac{1}{\theta^n} \mathbb{1}\{x_{(1)} > 0\} \mathbb{1}\{x_{(n)} < \theta\}$$

$$p_\theta(y) = \frac{1}{\theta^n} \mathbb{1}\{y_{(1)} > 0\} \mathbb{1}\{y_{(n)} < \theta\}$$

$\frac{p_\theta(x)}{p_\theta(y)}$ is free of θ if and only if $\mathbb{1}\{x_{(n)} < \theta\} = \mathbb{1}\{y_{(n)} < \theta\}$

This is true iff $x_{(n)} = y_{(n)}$. So, $x_{(n)}$ is minimal suff. for θ .

Sufficiency: Gives us a reduction of the data without sacrificing any information on the relevant parameter.

Minimal sufficiency: Gives us the "best" possible sufficient reduction

Result: If $\frac{p_\theta(x)}{p_\theta(y)}$ is free of θ if and only if $T(x) = T(y)$, then $T(X)$ is a minimal sufficient statistic for θ .

Proof: We are given $\frac{p_\theta(\underline{x})}{p_\theta(\underline{y})}$ is free of $\theta \Leftrightarrow T(\underline{x}) = T(\underline{y})$

We want to show: $T(\underline{x})$ is a minimal suff. statistic for θ .

- $T(\underline{x})$ is sufficient
- For any suff. statistic $S(\underline{x})$, we can find a fun. $g(\cdot)$ s.t. $T(\underline{x}) = g(S(\underline{x})) \quad \forall \underline{x} \in \mathcal{X}$.
- $T(\underline{x})$ is sufficient iff $p_\theta(\underline{x}) = g_\theta(T(\underline{x})) h(\underline{x}) \quad \forall \underline{x} \in \mathcal{X}$

Define $\mathcal{T} = \{ t : t = T(\underline{x}) \text{ for some } \underline{x} \in \mathcal{X} \}$

← image of \mathcal{X} under $T(\cdot)$

$A_t = \{ \underline{x} \in \mathcal{X} : T(\underline{x}) = t \}$ ← level set.

Now, for a fixed $t \in T$, fix an element $x_t \in A_t$.

Take $\underline{x} \in \mathcal{X}$. If $T(\underline{x}) = t$, then $x_{T(\underline{x})} = x_t \in A_t$

$\xrightarrow{\quad}$

$\Rightarrow \underline{x} \in A_t$.

$x_{T(\underline{x})}$ is the fixed element in A_t

\underline{x} and $x_{T(\underline{x})}$ belong to the same set A_t .

$$T(\underline{x}) = T(x_{T(\underline{x})}).$$

By the condition of the theorem $\frac{p_\theta(\underline{x})}{p_\theta(x_{T(\underline{x})})} = h(\underline{x})$
is free of θ .

$$p_\theta(\underline{x}) = p_\theta(x_{T(\underline{x})}) h(\underline{x}) \quad \forall \underline{x}$$

Def'm $g_\theta(b) = p_\theta(x_t)$, $t \in \mathcal{T}$.

$g_\theta(\cdot)$ is well-defined since x_t is a fixed element for every $t \in \mathcal{T}$.

With this, we set $p_\theta(\underline{x}) = g_\theta(T(\underline{x})) \tilde{h}(\underline{x})$
Hence

$\Rightarrow T(\underline{x})$ is sufficient for θ .

Next, we will show that $T(\underline{x})$ is minimal.

Let $S(\underline{x})$ be a sufficient statistic.

\Rightarrow we can find \tilde{g} and \tilde{h} st.

$p_\theta(\underline{x}) = \tilde{g}_\theta(S(\underline{x})) \tilde{h}(\underline{x}) \quad \forall \underline{x} \in \mathcal{X}$.

$$\frac{p_\theta(x)}{p_\theta(y)} = \frac{\tilde{g}_\theta(s(x))}{\tilde{g}_\theta(s(y))} \frac{\tilde{h}(x)}{\tilde{h}(y)}$$

If $s(x) = s(y)$, then $\frac{p_\theta(x)}{p_\theta(y)} = \frac{\tilde{h}(x)}{\tilde{h}(y)} \leftarrow \text{free of } \theta$

By the condition of the theorem, $T(x) = T(y)$

So, $s(x) = s(y) \Rightarrow T(x) = T(y)$.

Thus, There must exist a function $g(\cdot)$ st.

$T(x) = g(s(x)) \quad \forall x \in \mathcal{X}$.

So, $T(\cdot)$ is minimal.

Result: Let f and g be two functions such that

$f(x) = f(y) \Rightarrow g(x) = g(y)$. Then, there exists a function h such that $g(x) = h(f(x)) \forall x$

Proof: Take $z = f(x)$.

Define $h(z) = g(x)$

Suppose $z = f(x)$ and $z = f(y) \Rightarrow f(x) = f(y)$

\Rightarrow the value of the function h $\Rightarrow g(x) = g(y)$.

does not depend on the choice of x

So, h is a valid function.

And, by definition, $g(x) = h(f(x))$. (Proved).

One-parameter exponential family

A pmf / pdf $f_\theta(\cdot)$ belongs to the (one-parameter) exponential family if

$$f_\theta(x) = e^{a(\theta) z(x) - b(\theta)} c(x) \quad \forall x \in \mathcal{S} \quad \theta \in \mathcal{H}$$

\mathcal{S} is the support of the distribution (pmf / pdf)

\mathcal{H} is the parameter space.

The support \mathcal{S} must not involve the parameter θ .

Eg. Uniform $(0, \theta)$, $f_\theta(x) = \begin{cases} \frac{1}{\theta}, & \text{if } x \in (0, \theta) \\ 0 & \text{ow.} \end{cases}$

$\mathcal{S} = (0, \theta) \leftarrow$ depends on θ . Not in exponential family.

$$\begin{aligned}
 p_{\theta}(\underline{x}) &= \prod_{i=1}^n f_{\theta}(x_i) \\
 &= \prod_{i=1}^n \left\{ e^{a(\theta) \tau(x_i) - b(\theta)} c(x_i) \right\} \\
 &= e^{a(\theta) \sum_{i=1}^n \tau(x_i) - n b(\theta)} \prod_{i=1}^n c(x_i)
 \end{aligned}$$

$$= g_{\theta}(\tau(\underline{x})) h(\underline{x})$$

$$\begin{aligned}
 g_{\theta}(t) &= e^{a(\theta) t - n b(\theta)}, \quad h(\underline{x}) = \prod_{i=1}^n c(x_i) \\
 \tau(\underline{x}) &= \sum_{i=1}^n \tau(x_i)
 \end{aligned}$$

$$\Rightarrow \tau(\underline{x}) = \sum_{i=1}^n \tau(x_i) \text{ is sufficient for } \theta.$$

$$\text{In fact, } \frac{p_{\theta}(\underline{x})}{p_{\theta}(\underline{y})} = e^{a(\theta) \left\{ \sum_{i=1}^n T(x_i) - \sum_{i=1}^n T(y_i) \right\} \frac{\prod_{i=1}^n C(x_i)}{\prod_{i=1}^n C(y_i)}}$$

This is free of θ iff

$$\sum_{i=1}^n T(x_i) = \sum_{i=1}^n T(y_i)$$

\Updownarrow

$$T(\underline{x}) = T(\underline{y})$$

So, $T(\underline{x}) = \sum_{i=1}^n T(x_i)$ is a minimal sufficient statistic for θ .

Multi parameters exponential family

$\underline{\theta} = (\theta_1, \dots, \theta_r) \leftarrow r \text{ parameters}$.

A pdf / pmf belongs to the multi parameter exponential family if

$$f_{\underline{\theta}}(x) = e^{\sum_{j=1}^k a_j(\theta) \tau_j(x) - b(\theta)} c(x), \quad x \in \mathcal{S}, \quad \underline{\theta} \in \mathcal{H}.$$

Eg. $N(\mu, \sigma^2)$. $\underline{\theta} = (\mu, \sigma^2)$, $\mathcal{S} = \mathbb{R}$, $\mathcal{H} = \mathbb{R} \times (0, \infty)$

$$f_{\mu, \sigma^2}(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}, \quad x \in \mathbb{R}$$

$$= \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} e^{-\frac{1}{2\sigma^2}(x^2 + \mu^2 - 2\mu x)}$$

$$= \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} e^{-\frac{x^2}{2\sigma^2} - \frac{\mu^2}{2\sigma^2} + \frac{\mu x}{\sigma^2}}$$

$$= e^{-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{\mu^2}{2\sigma^2} - \log \sigma}$$

$$= e$$

$$a_1(\theta) \tau_1(x) + a_2(\theta) \tau_2(x) - b(\theta)$$

$$\frac{1}{\sqrt{2\pi}}$$

$$c(x)$$

$$= e$$

$$a_1(\theta) = -\frac{1}{2\sigma^2}, \quad a_2(\theta) = \frac{\mu}{\sigma^2}$$

$$\tau_1(x) = x^2, \quad \tau_2(x) = x$$

$$b(\theta) = \frac{\mu^2}{2\sigma^2} + \log \sigma$$

$$c(x) = \frac{1}{\sqrt{2\pi}} \quad \forall x \in \mathbb{R}$$

f_{μ, σ^2} belongs
to the multiparameter
exponential family

$$p_{\theta}(\underline{x}) = \prod_{i=1}^n \left\{ e^{\sum_{j=1}^k a_j(\theta) z_j(x_i)} - b(\theta) c(x_i) \right\}$$

$$= \prod_{i=1}^n \left\{ \left[\sum_{j=1}^k a_j(\theta) z_j(x_i) - b(\theta) \right] + \frac{b(\theta)}{c(x_i)} \right\}$$

$$= e^{\sum_{j=1}^k \left\{ a_j(\theta) \sum_{i=1}^n z_j(x_i) \right\} - nb(\theta)} \prod_{i=1}^n c(x_i)$$

$$= e^{\sum_{j=1}^k \left\{ a_j(\theta) \sum_{i=1}^n z_j(x_i) \right\}}$$

$$\Rightarrow T(\underline{x}) = \left(\sum_{i=1}^n z_1(x_i), \dots, \sum_{i=1}^n z_k(x_i) \right) \text{ is sufficient for } \underline{\theta}$$

$$\frac{p_{\theta}(x)}{p_{\theta}(y)} = e^{\sum_{j=1}^K \left\{ a_j(\theta) \sum_{i=1}^n \zeta_j(x_i) \right\} - \sum_{j=1}^K \left\{ a_j(\theta) \sum_{i=1}^n \zeta_j(y_i) \right\}}$$

$$= e^{\sum_{j=1}^K a_j(\theta) \left\{ \sum_{i=1}^n \zeta_j(x_i) - \sum_{i=1}^n \zeta_j(y_i) \right\}}$$

This ratio is free of θ if and only if

$$\sum_{i=1}^n \zeta_j(x_i) = \sum_{i=1}^n \zeta_j(y_i) \quad \forall j=1, \dots, K$$

$$\Rightarrow T(\underline{x}) = \left(\sum_{i=1}^n \zeta_1(x_i), \dots, \sum_{i=1}^n \zeta_K(x_i) \right)$$

minimal sufficient statistic for θ .

Eg. $N(\mu, \sigma^2)$

$$T_1(x) = x^2, \quad T_2(x) = x$$

$\left(\sum_{i=1}^n T_1(x_i), \sum_{i=1}^n T_2(x_i) \right)$ is minimal suff. fr
 (μ, σ^2)

$$= \left(\sum_{i=1}^n x_i^2, \sum_{i=1}^n x_i \right)$$

In many situations, the order statistics is the best possible reduction that we can get.

Order statistics

For a random sample X_1, \dots, X_n , the order statistics are $X_{(1)} \leq \dots \leq X_{(n)}$

$$X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta.$$

$$p_\theta(x) = \prod_{i=1}^n f_\theta(x_i) = \prod_{i=1}^n f_\theta(X_{(i)})$$

So, $(X_{(1)}, \dots, X_{(n)})$ is sufficient for θ .

Order statistics are sufficient even under the nonparametric regime.

Suppose that $X_1, \dots, X_n \stackrel{iid}{\sim} f$, where f is a pdf.

So, X_i 's are continuous random variables.

Therefore, $P(X_i = X_j) = 0$

So, the order statistics are distinct (with prob-1)

$$X_{(1)} < X_{(2)} < \dots < X_{(n)}$$

$$P(\underline{x} = \underline{x} \mid X_{(1)} = y_1, \dots, X_{(n)} = y_n) = \begin{cases} 0 & \text{if the ordered values of } n \text{ do not match } y_1, \dots, y_n \\ \frac{1}{n!} & \text{otherwise.} \end{cases}$$

$\Rightarrow (X_{(1)}, \dots, X_{(n)})$ is sufficient for f

Order statistics

Data: X_1, \dots, X_n

Ordered values: $X_{(1)} < X_{(2)} < \dots < X_{(n)}$

If X_i 's are continuous, Then the probability of two random variables taking the same value is 0
In that case $P(X_i$'s are all distinct) = 1.

Order statistics are sufficient if the underlying distribution is continuous.

The joint distribution / pdf of the order statistics

$$f_{X_{(1)}, \dots, X_{(n)}}(y_1, \dots, y_n) = \begin{cases} n! f(y_1) \cdots f(y_n), & \text{if } y_1 < y_2 < \dots < y_n \\ 0 & \text{otherwise.} \end{cases}$$

Marginal distribution of the order statistics:

Pdf of $X_{(j)}$.

$U(0, \theta)$ — $X_{(n)}$ is minimal suff. for θ .

$U(\theta, 1)$ — $X_{(1)}$ is minimal suff. for θ

If n is odd; Then $\widehat{X}_{\left(\frac{n+1}{2}\right)}$ is the median

$X_{\left(\frac{n}{4}\right)} - X_{\left(\frac{n}{4}\right)}$ = inter-quartile range.

Pdf of $X_{(1)}$:

The cdf of $X_{(1)}$

$$F_{(1)}(y) = \Pr(X_{(1)} \leq y) = 1 - \Pr(X_{(1)} > y)$$

$$\{X_{(1)} > y\} \Leftrightarrow \{X_1 > y, X_2 > y, \dots, X_n > y\}$$

$$\therefore \Pr(X_{(1)} > y) = \Pr(X_1 > y, X_2 > y, \dots, X_n > y)$$

$$(X_1, \dots, X_n \text{ are indep}) \therefore \Pr(X_1 > y) \Pr(X_2 > y) \dots \Pr(X_n > y)$$

$$(X_1, \dots, X_n \text{ are identically distributed}) = \left\{ \Pr(X_1 > y) \right\}^n$$
$$= \left\{ 1 - F_X(y) \right\}^n \quad (F_X(y) = \Pr(X_1 \leq y))$$

$$\begin{aligned}
 F_{(1)}(y) &= \Pr(X_{(1)} \leq y) = 1 - \Pr(X_{(1)} > y) \\
 &= 1 - \left\{ 1 - \Pr(X_1 \leq y) \right\}^n \\
 &= \left\{ 1 - F_X(y) \right\}^n
 \end{aligned}$$

If $F_X(\cdot)$ is differentiable, then $F_{(1)}(\cdot)$ is differentiable

In that case, the pdf of $X_{(1)}$ is

$$\begin{aligned}
 f_{(1)}(y) &= -n \left\{ 1 - F_X(y) \right\}^{n-1} (-f_X(y)) \\
 &= n \left\{ 1 - F_X(y) \right\}^{n-1} f_X(y).
 \end{aligned}$$

$X_{(1)}$ has the same range as X_1 .

The cdf of $X_{(n)}$

$$F_{(n)}(y) = \Pr(X_{(n)} \leq y)$$

$$\{X_{(n)} \leq y\} \Leftrightarrow \{X_1 \leq y, \dots, X_n \leq y\}$$

$$\therefore \Pr(X_{(n)} \leq y) = \Pr(X_1 \leq y, \dots, X_n \leq y)$$

$$= \Pr(X_1 \leq y) \cdots \Pr(X_n \leq y) \quad (\text{indep}).$$

$$= \left\{ \Pr(X_1 \leq y) \right\}^n$$

$$= \left\{ F_X(y) \right\}^n$$

$$\therefore F_{(n)}(y) = \left\{ F_X(y) \right\}^n$$

The pdf of $X_{(n)}$

$$f_{(n)}(y) = n \{F_X(y)\}^{n-1} f_X(y)$$

The pdf's of $X_{(1)}$ & $X_{(n)}$ are

$$f_{(1)}(y) = n \{1 - F_X(y)\}^{n-1} f_X(y)$$

$$f_{(n)}(y) = n \{F_X(y)\}^{n-1} f_X(y)$$

Both $X_{(1)}$ and $X_{(n)}$ have the same range as X_1
(marginally).

$X_{(j)}$ For general j :

$$\bigcirc X_1 \bigcirc - \bigcirc X_j$$

The cdf of $X_{(j)}$

$$F_{(j)}(y) = \text{IP}(X_{(j)} \leq y).$$

Define $W_i^0 = \begin{cases} 1 & \text{if } X_i \leq y \\ 0 & \text{if } X_i > y \end{cases} = \mathbb{I}\{X_i \leq y\}$

$\{X_{(j)} \leq y\} (=)$ at least j many of the W_i 's
are 1.

$$\sum_{i=1}^n W_i \geq j$$

$$w_i = \mathbb{1} \{x_i \leq y\}$$

$$\{x_{(i)} \leq y\} \Leftrightarrow \sum_{j=1}^n w_i \geq j$$

w_i 's are Bernoulli

w_i 's are independent.

$$\Pr(w_i = 1) = \Pr(x_i \leq y) = F_X(y) \leftarrow \text{true if } i \\ (\because x_i \text{'s are identical})$$

$\therefore w_1, \dots, w_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(F_X(y))$

$$\Rightarrow \sum_{i=1}^n w_i \sim \text{Binomial}(n, F_X(y))$$

$$\{X_{(j)} \leq y\} \quad (=) \quad \sum_{i=1}^n w_i \geq j$$

$$\sum_{i=1}^n w_i \sim \text{Binomial}(n, F_X(y))$$

$$\therefore \text{IP}\left(X_{(j)} \leq y\right) = \text{IP}\left(\text{Bin}(n, F_X(y)) \geq j\right)$$

$$= \sum_{k=j}^n \text{IP}\left(\text{Bin}(n, F_X(y)) = k\right)$$

$$= \sum_{k=j}^n \binom{n}{k} \{F_X(y)\}^k \left\{1 - F_X(y)\right\}^{n-k}$$

$$=: F_{(j)}(y) .$$

If $F_X(\cdot)$ is differentiable, then the pdf of $X_{(j)}$ is

$$\sum_{k=j}^n \binom{n}{k} \left[\left\{ F_X(y) \right\}^k \left\{ 1 - F_X(y) \right\}^{n-k} \left\{ -f_X(y) \right\} \right. \\ \left. + k \left\{ F_X(y) \right\}^{k-1} f_X(y) \left\{ 1 - F_X(y) \right\}^{n-k} \right]$$

$$= \sum_{k=j}^n \binom{n}{k} \left\{ F_X(y) \right\}^{k-1} \left\{ 1 - F_X(y) \right\}^{n-k-1} f_X(y) \\ \left[-(n-k) F_X(y) + k \left\{ 1 - F_X(y) \right\} \right]$$

$$= \sum_{k=j}^n \left\{ F_X(y) \right\}^{k-1} \left\{ 1 - F_X(y) \right\}^{n-k-1} f_X(y) \\ \left[\binom{n}{k} k \left\{ 1 - F_X(y) \right\} - \binom{n}{k} (n-k) F_X(y) \right]$$

$$\sum_{k=j}^n \left\{ F_X(y) \right\}^{k-1} \left\{ 1 - F_X(y) \right\}^{n-k-1} f_X(y) \\ \left[\binom{n}{k} k \left\{ 1 - F_X(y) \right\} - \binom{n}{k-1} (n-k) F_X(y) \right]$$

$$\binom{n}{k} k = \frac{n!}{k!(n-k)!} \quad k = \frac{n!}{(k-1)!(n-k)!} = \frac{n!}{(k-1)!(n-k+1)!} \\ = (n-k+1) \binom{n}{k-1}$$

$$= f_X(y) \sum_{k=j}^n \left\{ F_X(y)^{k-1} \left\{ 1 - F_X(y) \right\}^{n-k} \binom{n}{k-1} (n-k+1) \right. \\ \left. - F_X(y)^k \left\{ 1 - F_X(y) \right\}^{n-k-1} \binom{n}{k} (n-k) \right\}$$

$$a_k = F_X(y)^{k-1} \left\{ 1 - F_X(y) \right\}^{n-k} \binom{n}{k-1} (n-k+1) \\ = f_X(y) \sum_{k=j}^n (a_k - a_{k+1})$$

$$\sum_{k=j}^n (a_k - a_{k+1}) = a_j - a_{j+1} + a_{j+1} - a_{j+2} + \dots + a_h - a_{h+1}$$

$$= a_j - a_{h+1}$$

$$a_j - a_{h+1} = \left\{ F_X(y) \right\}^{j-1} \left\{ 1 - F_X(y) \right\}^{h-j} \binom{n}{j-1} \binom{h-j}{h-j}$$

$$- \left\{ F_X(y) \right\}^h \left\{ 1 - F_X(y) \right\}^{-1} \binom{n}{h} \underbrace{\binom{h-h-1+1}{h-1+1}}_0$$

$$= \left\{ F_X(y) \right\}^{j-1} \left\{ 1 - F_X(y) \right\}^{h-j} \frac{n!}{(j-1)! (h-j)!}$$

i.e. The bdf of $X_{(j)}$

$$f_{(j)}(y) = \frac{n!}{(j-1)! (h-j)!} \left\{ F_X(y) \right\}^{j-1} \left\{ 1 - F_X(y) \right\}^{h-j} f_X(y)$$

So, The pdf of $X_{(j)}$ is

$$f_{(j)}(y) = \frac{n!}{(j-1)! (n-j)!} \{F_X(y)\}^{j-1} \{1-F_X(y)\}^{n-j} f_X(y).$$

Eg. Let $X_1 - X_n \stackrel{\text{i.i.d}}{\sim} U(0, 1)$
Then, $f_X(y) = \begin{cases} 1 & \text{if } 0 < y < 1 \\ 0 & \text{otherwise.} \end{cases}$

$$F_X(y) = \begin{cases} 0 & \text{if } y < 0 \\ y & \text{if } 0 \leq y < 1 \\ 1 & \text{if } y \geq 1 \end{cases}$$

$$\therefore f_{(j)}(y) = \frac{n!}{(j-1)! (n-j)!} y^{j-1} (1-y)^{n-j}, \quad 0 < y < 1$$

Joint distribution of two order statistics

The j.t. pdf of $X_{(i)}$ and $X_{(j)}$ is where $i < j$

$$f_{(i,j)}(y, z) = \frac{(i-1)! (j-i-1)! (n-j)!}{(i-1)! (j-i-1)! (n-j)!} \left\{ F_X(y) \right\}^{i-1} \left\{ F_X(z) - F_X(y) \right\}^{j-i-1} \left\{ 1 - F_X(z) \right\}^{n-j}, \quad y < z$$

$$f_X(y) f_X(z)$$

$$f_{(i,j)}(y, z) = \frac{n!}{(i-1)! (j-i-1)! (n-j)!} f_X(y) f_X(z)$$

$$\left\{ F_X(y) \right\}^{i-1} \left\{ F_X(z) - F_X(y) \right\}^{j-i-1} \left\{ 1 - F_X(z) \right\}^{n-j}, \quad y < z$$

Auxiliary statistic

A statistic $S = S(\underline{X})$ is said to be auxiliary if the distribution of S is free of the unknown parameter θ .

Eg. $X_1, X_2 \stackrel{\text{iid}}{\sim} \text{Normal}(\theta, 1)$, $\theta \in \mathbb{R}$.

$X_1 - X_2 \sim \text{Normal}(0, 2) \leftarrow$ free of θ .

$X_1 - X_2$ is auxiliary.

Eg. $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} N(\theta, 1)$, $\theta \in \mathbb{R}$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad S^2 = \frac{1}{n-1} \sum_{i=1}^{n-1} (X_i - \bar{X})^2$$

$\bar{X} \sim N\left(\theta, \frac{1}{n}\right)$, $(n-1)S^2 \sim \chi_{n-1}^2$, \bar{X}, S^2 indep.
 S^2 is auxiliary. \uparrow free of θ

Eg. X_1, \dots, X_n iid Uniform $(0, \theta)$, $\theta > 0$.

$X_{(n)}$ is minimal sufficient.

$\frac{X_{(1)}}{X_{(n)}}$ has distribution free of θ .

$\frac{X_{(i)}}{X_{(j)}}$ is ancillary.

$\frac{X_{(i)}}{X_{(j)}}$, $1 \leq i < j \leq n$ is ancillary.

Eg Location family of distributions

X_1, \dots, X_n iid f_θ , $f_\theta(x) = f(x - \theta)$ $\forall x$ $\forall \theta$

where f is a known pdf (does not involve any unknown parameter)

In this case, the cdf satisfies

$$F_\theta(x) = \int_{-\infty}^x f_\theta(t) dt = \int_{-\infty}^x f(t-\theta) dt$$
$$= \int_{-\infty}^{x-\theta} f(z) dz = F(x-\theta)$$

↑ cdf corr. to f .

$F_\theta(x) = F(x-\theta)$, F is completely known

Eg. Normal $(\theta, 1)$

$$f_\theta(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\theta)^2}{2}}, x \in \mathbb{R}$$

$$= f(x-\theta), \text{ where } f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}, z \in \mathbb{R}$$

Eg. Cauchy $(\theta, 1)$

$$f_{\theta}(x) = \frac{1}{\pi} \frac{1}{1 + (x-\theta)^2}, \quad x \in \mathbb{R}$$

$$= f(x-\theta), \quad f(z) = \frac{1}{\pi} \frac{1}{1+z^2}, \quad z \in \mathbb{R}$$

Eg. Laplace $(\theta, 1)$ is a member of the location family.

If X has p.d.f $f_{\theta}(x) = f(x-\theta)$

Then $X-\theta$ has p.d.f $f(x)$

$\Rightarrow X-\theta = Y$ has a known distribution

$$X = \theta + Y$$

For location family:

$$X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta(x) = f(x-\theta)$$

$X_1 - \theta, \dots, X_n - \theta \stackrel{iid}{\sim} f(x) \leftarrow$ completely known.
random variables, but not observable.

$$X_{(n)} - X_{(1)} = \max_i X_i - \min_i X_i$$

$$= \max_i (X_i - \theta) + \theta - \left\{ \min_i (X_i - \theta) + \theta \right\}$$

$$= \max_i (X_i - \theta) + \cancel{\theta} - \min_i (X_i - \theta) - \cancel{\theta}$$

$$= \max_i Y_i - \min_i Y_i, \quad Y_i = X_i - \theta \\ = Y_{(n)} - Y_{(1)}$$

$Y_1, \dots, Y_n \stackrel{iid}{\sim} f(\cdot) \leftarrow$ completely known $\Rightarrow Y_{(n)} - Y_{(1)}$ has
dist. free of θ .

$\therefore X_{(h)} - X_{(1)} = Y_{(h)} - Y_{(1)}$ is a ancillary.

The range $R = X_{(h)} - X_{(1)}$ for a location family $f(x-\theta)$ is a ancillary.

Scale family $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f_\theta(x), \theta > 0$.

f_θ is said to belong to the scale family if $f_\theta(x) = \frac{1}{\theta} f\left(\frac{x}{\theta}\right)$, $\forall x \neq 0$, where $f(\cdot)$ is completely known.

Auxiliary statistic A statistic $S = S(\underline{x})$ is auxiliary if its distribution is free of the unknown parameter θ .

Sufficient statistics The part of the sample that contains all the information about the unknown parameter θ .

Minimal sufficient statistics Give us the best possible reduction.

Maybe auxiliary statistics & minimal suff. stat. are unrelated ← NOT True.

Eg. Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Uniform}(\theta, \theta+1)$, $\theta \in \mathbb{R}$.

A minimal suff. statistic for θ is $(X_{(1)}, X_{(n)})$.

$\Rightarrow (X_{(n)} - X_{(1)}, X_{(1)} + X_{(n)})$ is also minimal sufficient.

$\text{Uniform}(\theta, \theta+1) = \text{Uniform}(0, 1) + \theta$.

$$f_\theta(y) = \begin{cases} 1 & \text{if } \theta < y < \theta+1 \\ 0 & \text{o.w.} \end{cases} \quad | \quad f(y) = \begin{cases} 1 & \text{if } 0 < y < 1 \\ 0 & \text{o.w.} \end{cases}$$

$$= \begin{cases} 1 & \text{if } 0 < y - \theta < 1 \\ 0 & \text{o.w.} \end{cases}$$

$= f(y - \theta)$ \leftarrow a distribution from the location family.

$\therefore X_{(n)} - X_{(1)}$ is ancillary.

The distribution of $X_{(n)} - X_{(1)}$ can be computed directly.

$$f_R(r) = n(n-1) r^{n-2}, \quad 0 < r < 1 \quad] \text{ for } \theta.$$

$$R = X_{(n)} - X_{(1)}.$$

$X_{(n)} - X_{(1)}$ is a ancillary.

$(X_{(n)} - X_{(1)}, X_{(1)} + X_{(n)})$ is minimal sufficient.

They are not independent.

When is it true that a minimal sufficient statistic is unrelated to an ancillary statistic?

Complete statistic

Defn.: A statistic $T = T(\underline{x})$ is said to be complete if

$$\text{IE}_\theta \{ g(T) \} = 0 \quad \forall \theta \Rightarrow \text{IP}_\theta (g(T) = 0) = 1 \quad \forall \theta.$$

$g(\cdot)$ is a real-valued function

$g(T) = 0$ with pr. 1

Eg. X_1, \dots, X_n $\stackrel{iid}{\sim}$ Bernoulli(θ), $0 < \theta < 1$

$$T = \sum_{i=1}^n X_i \sim \text{Binomial}(n, \theta)$$

Take any function $g(\cdot)$

$$E_\theta \{ g(T) \} = \sum_{t=0}^n g(t) \text{IP}_\theta (T=t)$$

$$= \sum_{t=0}^n g(t) \binom{n}{t} \theta^t (1-\theta)^{n-t}$$

$$E_\theta \{ g(T) \} = 0 \quad \forall \theta \in (0, 1)$$

$$\Leftrightarrow \sum_{t=0}^n g(t) \binom{n}{t} \theta^t (1-\theta)^{n-t} = 0 \quad \forall \theta \in (0, 1)$$

$$(\Rightarrow) (\text{---}) \sum_{t=0}^n g(t) \binom{n}{t} \left(\frac{\theta}{1-\theta} \right)^t = 0 \quad \forall \theta \in (0, 1)$$

$$(\Rightarrow) \sum_{t=0}^n g(t) \binom{n}{t} n^t = 0 \quad \forall n \in (0, \infty) \quad \begin{cases} n = \frac{\theta}{1-\theta} \\ \theta \in (0, 1) \Rightarrow n \in (0, \infty) \end{cases}$$

↑
does not involve n .

a polynomial in n of degree n

$$(\Rightarrow) g(t) \binom{n}{t} = 0 \quad \forall t = 0, 1, \dots, n$$

$$(\Rightarrow) g(t) = 0 \quad \forall t = 0, 1, \dots, n \quad [\because \binom{n}{t} \neq 0] \quad \begin{cases} \text{since this} \\ \text{is a} \\ \text{polynomial} \\ \text{of degree } n \\ \text{which is} \\ \text{identically 0} \end{cases}$$

$$\therefore E_\theta \{ g(T) \} = 0 \Rightarrow g(T) = 0 \text{ w.p. 1.}$$

$$\therefore T = \sum_{i=1}^n X_i \text{ is complete.}$$

Ex. X_1, \dots, X_n iid Uniform $(0, \theta)$, $\theta > 0$ [$\theta \in (0, \infty)$]

$T = X_{(n)}$ is minimal sufficient for θ .

$$f_{(n)}(y) = n \left(\frac{y}{\theta}\right)^{n-1} \frac{1}{\theta}, \quad 0 < y < \theta$$
$$= n \frac{y^{n-1}}{\theta^n}, \quad 0 < y < \theta$$

Take a function g .

$$E_\theta[g(T)] = \int g(t) f_T(t) dt$$

$$= \int_0^\theta g(t) n \frac{t^{n-1}}{\theta^n} dt$$

$$= \frac{n}{\theta^n} \int_0^\theta g(t) t^{n-1} dt$$

$$E_\theta [g(T)] = 0 \quad \forall \theta \in (0, \infty)$$

$$\Leftrightarrow \int_0^\theta g(t) t^{n-1} dt = 0 \quad \forall \theta \in (0, \infty)$$

The LHS as a function of θ is differentiable
and its derivative is

$$g(\theta) \theta^{n-1}$$

The derivative of the RHS w.r.t. θ is 0

$$\therefore g(\theta) \theta^{n-1} = 0 \quad \forall \theta \in (0, \infty) \left[\begin{array}{l} \text{since LHS} \\ \text{two derivatives} \\ \text{should match} \end{array} \right]$$

$$\Rightarrow g(\theta) = 0 \quad \forall \theta \in (0, \infty)$$

$$\Rightarrow \mathbb{P}_\theta [g(T) = 0] = 1 \quad \therefore T = X_{(n)} \text{ is complete.}$$

Def: A complete statistic that is sufficient is called a complete sufficient statistic.

Result: A complete sufficient statistic is minimal sufficient.

Pf: Let $T = T(\underline{x})$ be complete & sufficient.

Let $S = S(\underline{x})$ be a sufficient statistic

We will show that $T = g(S)$ for some $g(\cdot)$.

Define $\psi(T) = T - \overline{IE}_{\theta}(T|S)$ S is sufficient.

$$\overline{|}$$

does not depend
on θ ($\because S$ is
sufficient)

$X|S$ is free of θ .

$T|S$ is free of θ .

$$\begin{aligned} \mathbb{E}_\theta[\psi(T)] &= \mathbb{E}_\theta[T - \mathbb{E}_\theta(T|S)] \\ &= \mathbb{E}_\theta(T) - \mathbb{E}_\theta \mathbb{E}_\theta(T|S) \\ &= \mathbb{E}_\theta(T) - \mathbb{E}_\theta(T) = 0 \quad \forall \theta \end{aligned}$$

T is complete $\Rightarrow \psi(T) = 0$ w.p 1.

$$\Rightarrow T - \mathbb{E}_\theta(T|S) = 0 \quad \text{w.p 1}$$

$$\Rightarrow T = \mathbb{E}_\theta(T|S) \quad \text{w.p 1}$$

$$\mathbb{P}_\theta \left[T = \mathbb{E}_\theta(T|S) \right] = 1$$

a function of S and does not involve θ .

$$\Rightarrow \mathbb{P}_\theta \left[T = g(S) \right] = 1$$

$\therefore T$ can be written as a function of S .

$\Rightarrow T$ is a minimal sufficient statistic.

If T is complete and S is sufficient, then we can always write T as a function of S . But, it is not always true that T is sufficient. However, if T is sufficient, then it must be minimal sufficient.

Rem11: A complete sufficient statistic is independent of any ancillary statistic. (Basu's Theorem)

Pf: Let T be complete sufficient.

S be ancillary

We will show:

$$IP_{\theta} [S \in A | T] = IP_{\theta} (S \in A) \text{ for every } A.$$

$IP_{\theta} (S \in A) \leftarrow$ is free of θ

$IP_{\theta} (S \in A | T) \leftarrow$ free $\nabla \theta$ (since T is suff.)

Define $p_A = IP_{\theta} (S \in A)$, $q_A(T) = IP_{\theta} (S \in A | T)$

$$\begin{aligned} \text{IP}(A) &= \int_A f(x) dx = \int_{-\infty}^{\infty} \underbrace{1}_{\mathbb{1}} [x \in A] f(x) dx \\ &= \mathbb{E} \{ \mathbb{1}_{(X \in A)} \} \\ &= \mathbb{E} \{ \mathbb{1}_A \} \end{aligned}$$

$$\text{IP}(A) = \mathbb{E}_X (\mathbb{1}_A)$$

$$\text{IP}(A|B) = \mathbb{E} (\mathbb{1}_A | B)$$

$$\mathbb{E} \text{IP}(A|B) = \mathbb{E} \mathbb{E}(\mathbb{1}_A | B) = \mathbb{E} (\mathbb{1}_A) = \text{IP}(A)$$

$$\therefore p_A = \mathbb{P}_\theta(S \in A) = \mathbb{E}_\theta[\mathbb{P}_\theta(S \in A | T)]$$

$$= \mathbb{E}_\theta[q_A(T)] \quad \forall \theta$$

$$\therefore \mathbb{E}_\theta[q_A(T) - p_A] = 0 \quad \forall \theta$$

T is complete $\Rightarrow \mathbb{P}_\theta\{q_A(T) = p_A\} = 1 \quad \forall \theta$

$$\therefore \text{wp } \mathbb{P}$$

$$\mathbb{P}_\theta(S \in A | T) = \mathbb{P}_\theta(S \in A) \quad \forall A$$

The conditional dist of S given T equals
The unconditional dist of S

$\Rightarrow T \perp \! f \! \perp S$ are independent.

Completeness : T is said to be complete if

$$E_{\theta} \{ g(T) \} = 0 \quad \forall \theta \in \Theta \text{ implies } P_{\theta} \{ g(T) = 0 \} = 1.$$

Result: If T is complete and sufficient, then

T is minimal sufficient.

Result: (Basu's Theorem) If T is complete sufficient
and S is ancillary, then T and S are independent.

Eg. Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$, $\theta \in \mathbb{R}$.

$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ is minimal sufficient for θ .

\bar{X} is also complete

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

$$S^2 = \frac{1}{2n(n-1)} \sum_{1 \leq i \neq j \leq n} (x_i - x_j)^2$$

$$\frac{1}{2n(n-1)} \sum_{1 \leq i \neq j \leq n} (x_i - \bar{x} - x_j + \bar{x})^2$$

$$(x_i - \bar{x})^2 + (x_j - \bar{x})^2 - 2(x_i - \bar{x})(x_j - \bar{x})$$

$$\sum_{i,j=1}^n (x_i - \bar{x})^2 = n \sum_{i=1}^n (x_i - \bar{x})^2$$

equal.

$$\sum_{i,j=1}^n (x_i - \bar{x})^2 = n \sum_{j=1}^n (x_j - \bar{x})^2$$

$$\sum_{i,j=1}^n (x_i - \bar{x})(x_j - \bar{x}) = \sum_{i=1}^n (x_i - \bar{x}) \sum_{j=1}^n (x_j - \bar{x}) = 0$$

$$\frac{1}{2n(n-1)} \sum_{1 \leq i \neq j \leq n} (x_i - x_j)^2 = \frac{1}{2n(n-1)} \left\{ 2n \sum_{i=1}^n (x_i - \bar{x})^2 \right\}$$

$$= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = S^2$$

$$X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$$

$$\underbrace{X_1 - \theta}, \dots, \underbrace{X_n - \theta} \stackrel{iid}{\sim} N(0, 1)$$

$$Y_1, \dots, Y_n \stackrel{iid}{\sim} N(0, 1)$$

$$\begin{aligned} y_i &= x_i - \theta \\ x_i &= y_i + \theta \end{aligned}$$

$$x_i - x_j = y_i + \theta - (y_j + \theta) = y_i - y_j$$

$$S^2 = \frac{1}{2n(n-1)} \sum_{i \neq j} (x_i - x_j)^2 = \frac{1}{2n(n-1)} \sum_{i \neq j} (y_i - y_j)^2$$

$\therefore S^2$ is auxiliary

\bar{X} is complete sufficient.

S^2 is ancillary.

So, by Basu's Theorem \bar{X} and S^2 are independent.

Eg. Let X_1, \dots, X_n iid $\text{Unif}(0, \theta)$, $\theta > 0$

$X_{(n)}$ is complete sufficient.

$\frac{X_{(1)}}{X_{(n)}}$ is ancillary

\therefore Basu's Theorem tells us that $\frac{X_{(1)}}{X_{(n)}}$ and $X_{(n)}$ are independent

Let X_1, \dots, X_n $\stackrel{iid}{\sim} N(\theta, 1)$, $T = \bar{X} \sim N\left(\theta, \frac{1}{n}\right)$

$$\begin{aligned} E_\theta \{ g(T) \} &= \frac{1}{\sqrt{2\pi \frac{1}{n}}} \int_{-\infty}^{\infty} g(t) e^{-\frac{1}{2\frac{1}{n}}(t-\theta)^2} dt \\ &= \sqrt{\frac{n}{2\pi}} \int_{-\infty}^{\infty} g(t) e^{-\frac{n(t-\theta)^2}{2}} dt = 0 \quad \forall \theta \in \mathbb{R} \end{aligned}$$

Integral transforms: A mapping between function classes, defined in terms of integrals.

Let f be a continuous function.

Define

$$\psi(t) = \int_0^\infty e^{-tx} f(x) dx, \quad t \in \mathbb{R}$$

\leftarrow one-sided Laplace transform.

$$\psi(t) = \int_{-\infty}^\infty e^{-tx} f(x) dx, \quad t \in \mathbb{R}$$

\leftarrow two-sided Laplace transform

$$\psi(t) = \int_0^\infty x^{t-1} f(x) dx, \quad t > 0$$

\leftarrow Mellin's transform

$$\psi(t) = \int_0^\infty \frac{1}{t+x} f(x) dx, \quad t > 0$$

\leftarrow Stieltjes transform

$$\psi(t) = \int_{-\infty}^\infty e^{itx} f(x) dx, \quad t \in \mathbb{R}$$

\leftarrow Fourier transform

All these transforms have the property that
if $\psi_f(t) = \psi_g(t)$ $\forall t$ then $f(x) = g(x) \forall x$

$$f \mapsto \psi_f, \quad g \mapsto \psi_g$$

$$\psi_f \equiv \psi_g \Rightarrow f \equiv g$$

In particular, if $\psi_f \equiv 0$ ($\psi_f(t) = 0 \forall t$)
 $\lim_{t \rightarrow \infty} f \equiv 0$ ($f(x) = 0 \forall x$)

$$E_\theta(g(T)) = \int_{-\infty}^{\infty} g(t) f(t) dt \leftarrow c \psi_g(\theta)$$

$$\psi_g(\theta) = 0 \quad \forall \theta \Rightarrow g(t) = 0 \quad \forall t.$$

Back to the normal example:

$$E_\theta \{ g(T) \} = \sqrt{\frac{n}{2\pi}} \int_{-\infty}^{\infty} g(t) e^{-\frac{n(t-\theta)^2}{2}} dt$$

$$= \sqrt{\frac{n}{2\pi}} \int_{-\infty}^{\infty} g(t) e^{-\frac{nt^2}{2}} e^{-\frac{n\theta^2}{2}} e^{nt\theta} dt$$

$$= \sqrt{\frac{n}{2\pi}} e^{-\frac{n\theta^2}{2}} \int_{-\infty}^{\infty} g(t) e^{-\frac{nt^2}{2}} e^{nt\theta} dt$$

$$E_\theta \{ g(T) \} = 0 \quad \forall \theta$$

$$\Rightarrow \int_{-\infty}^{\infty} g(t) e^{-\frac{nt^2}{2}} e^{nt\theta} dt = 0 \quad \forall \theta$$

By uniqueness of integral transforms,

$$g(t) e^{-\frac{nt^2}{2}} = 0 \quad \forall t \in \mathbb{R}$$

$$e^{-\frac{nt^2}{2}} > 0 \quad \forall t \Rightarrow g(t) = 0 \quad \forall t \in \mathbb{R}.$$

$$\therefore P_{\theta} \{ g(T) = 0 \} = 1. \quad \forall \theta.$$

$\therefore T = \bar{x}$ is complete.

Complete sufficient statistic for exponential family

The pdf or pmf is given by

$$e^{\sum_{j=1}^K a_j(\theta) \tau_j(x) - b(\theta)} c(x), \quad x \in \mathcal{X}$$

↑ free of θ

If X_1, \dots, X_n are iid with the above pdf/pmf,

then $\left(\sum_{i=1}^n \tau_1(X_i), \dots, \sum_{i=1}^n \tau_K(X_i) \right)$ is minimal sufficient for θ .

Result: Let X_1, \dots, X_n be iid with pdf/pmf

$$e^{\sum_{j=1}^K a_j(\theta) \tau_j(x) - b(\theta)} c(x), \quad x \in \mathcal{X}$$

Then, $T = \left(\sum_{i=1}^n \tau_1(X_i), \dots, \sum_{i=1}^n \tau_K(X_i) \right)$ is complete sufficient for θ if $\{(a_1(\theta), \dots, a_K(\theta)) : \theta \in \mathbb{G}\}$ contains a K -dimensional open set.

Eg. $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2)$. $\therefore \mu \in \mathbb{R}, \sigma > 0$

$$a_1(\mu, \sigma^2) = \mu, \quad a_2(\mu, \sigma^2) = \frac{1}{\sigma^2} \quad \textcircled{A} = \mathbb{R} \times (0, \infty)$$

$$\left\{ (a_1(\mu, \sigma^2), a_2(\mu, \sigma^2)) : (\mu, \sigma^2) \in \textcircled{A} \right\}$$

$$= \left\{ \left(\mu, \frac{1}{\sigma^2} \right) : \mu \in \mathbb{R}, \sigma > 0 \right\}$$

$$= \left\{ (\nu, \kappa) : \nu \in \mathbb{R}, \kappa > 0 \right\}, \quad \nu = \mu, \kappa = \frac{1}{\sigma^2}$$

two-dimensional
contains a open subset of \mathbb{R}^2 .

$\therefore \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2 \right)$ is complete sufficient for (μ, σ^2)

Eg. X_1, \dots, X_n iid $N(\theta, \theta)$, $\theta > 0$

$\sum_{i=1}^n X_i^2$ is minimal sufficient.

$\therefore (\sum X_i, \sum X_i^2)$ is not minimal sufficient.

$\therefore (\sum X_i, \sum X_i^2)$ is not complete.

Result: Let T be a complete sufficient statistic.

Then, any 1-1 function of T is also complete sufficient.

Pf: Suppose S is a 1-1 function of T .

$\exists h(\cdot)$ s.t. $S = h(T)$

$$E_\theta \{ g(S) \} = 0 \quad \forall \theta \Rightarrow P_\theta \{ g(S) = 0 \} = 1 \quad \forall \theta$$

$$E_{\theta} \{ g(S) \} = E_{\theta} \{ g(h(T)) \} = 0 \quad \forall \theta$$

$$\Leftrightarrow E_{\theta} \{ g \circ h(T) \} = 0 \quad \forall \theta$$

$$\Rightarrow P_{\theta} \{ g \circ h(T) = 0 \} = 1 \quad \forall \theta$$

$$= P_{\theta} \{ g(S) = 0 \} = 1 \quad \forall \theta$$

$\Rightarrow S$ is also complete sufficient.

Sufficiency principle: In a statistical inference problem,

if there is a (minimal) sufficient statistic, Then our inference should depend on the sample only via the (minimal) sufficient statistic.

That is, if $T = T(\underline{X})$ is (minimal) sufficient.

For two sets of observations \underline{x} and \underline{y} , we have

$T(\underline{x}) = T(\underline{y})$. Then, our inference should be the same irrespective of whether we observe $\underline{X} = \underline{x}$ or $\underline{X} = \underline{y}$.

Likelihood.

$$x_1, \dots, x_n \stackrel{\text{iid}}{\sim} f_\theta$$

$$p_\theta(\underline{x}) = \prod_{i=1}^n f_\theta(x_i) = g_\theta(T(\underline{x})) h(\underline{x})$$

$$p_\theta(\underline{y}) = g_\theta(T(\underline{y})) h(\underline{y})$$

$$p_\theta(\underline{x}) = \frac{h(\underline{x})}{h(\underline{y})} \frac{g_\theta(T(\underline{x}))}{g_\theta(T(\underline{y}))} p_\theta(\underline{y})$$

Def: Let X_1, \dots, X_n have joint dist (pdf/pdf) $p_\theta(\cdot)$ which is indexed by $\theta \in \mathbb{H}$. The likelihood function $L(\theta)$ is defined as

$$L(\theta) = L(\theta | \underline{x}) = p_\theta(\underline{x}) \quad \forall \theta \in \mathbb{H}$$

Eg: $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$

$$p_\theta(\underline{x}) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{n-x_i}.$$

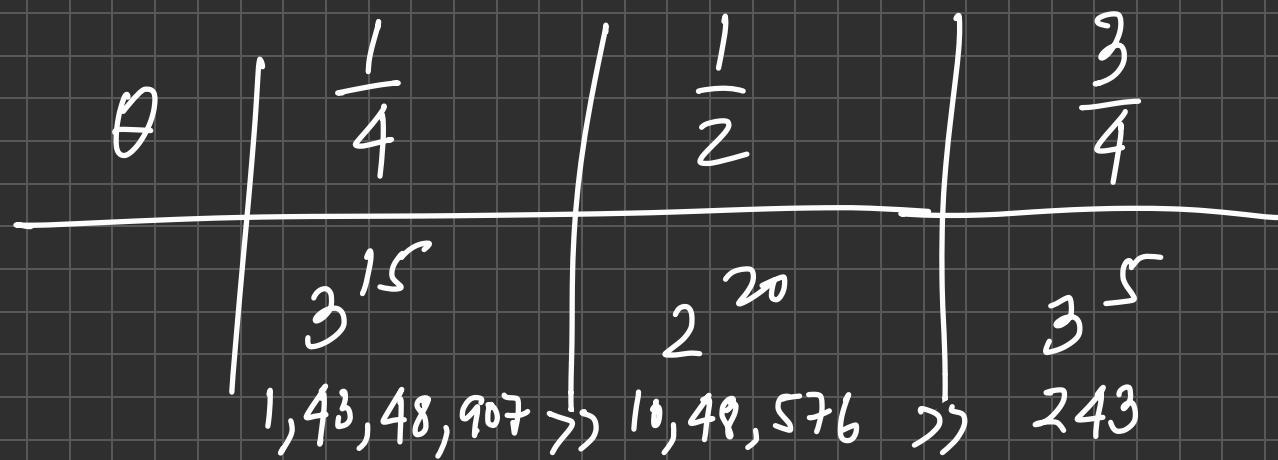
$\theta = \left\{ \frac{1}{4}, \frac{1}{2}, \frac{3}{4} \right\}$. we observe 5 successes out of 20 draws.

$$n = 20, \sum_{i=1}^n x_i = 5$$

$$\text{If } \theta = \frac{1}{4}, \text{ then } L(\theta) = \left(\frac{1}{4}\right)^5 \left(\frac{3}{4}\right)^{15} = \frac{3^{15}}{4^{20}}.$$

$$\theta = \frac{1}{2}, \quad L(\theta) = \left(\frac{1}{2}\right)^5 \left(\frac{1}{2}\right)^{15} = \frac{2^5 \cdot 2^{15}}{4^{20}} = \frac{2^{20}}{4^{20}}.$$

$$\theta = \frac{3}{4}, \quad L(\theta) = \left(\frac{3}{4}\right)^5 \left(\frac{1}{4}\right)^{15} = \frac{3^5}{4^{20}}.$$



The likelihood: Let x_1, \dots, x_n be a random sample with
jt. pdf / pmf $p_\theta(\underline{x})$, $\theta \in \mathbb{H}$. Then The likelihood
of θ is defined as

$$L(\theta) = L(\theta | \underline{x}) = p_\theta(\underline{x}) \quad \forall \theta \in \mathbb{H}$$

\uparrow $\forall \underline{x} \in \mathcal{X}$.

$L(\theta | \underline{x})$ viewed as a function of θ .
 $L(\theta | \underline{x})$ tells us how likely it is that the true
parameter is θ when we observe \underline{x} as our sample.

Likelihood principle: Let $L(\theta | \underline{x})$ be the likelihood function. Suppose for two sets of observations \underline{x} and \underline{y} , we have $L(\theta | \underline{x}) = C(\underline{x}, \underline{y}) L(\theta | \underline{y})$. Then, the inference on θ based on whether we observe \underline{x} or \underline{y} should be the same.

Eg. X_1, \dots, X_n $\stackrel{\text{iid}}{\sim}$ Bernoulli(θ) $0 < \theta < 1$

$$\frac{L(\theta | \underline{x})}{L(\theta | \underline{y})} = \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i}, \quad 0 \leq \sum_{i=1}^n x_i \leq n$$

Let \underline{y} be another observation.

$$\frac{L(\theta | \underline{y})}{L(\theta | \underline{x})} = \theta^{\sum_{i=1}^n y_i} (1-\theta)^{n - \sum_{i=1}^n y_i}$$

$$\frac{L(\theta | \underline{x})}{L(\theta | \underline{y})} = \theta^{\sum_{i=1}^n x_i - \sum_{i=1}^n y_i} (1-\theta)^{\sum_{i=1}^n y_i - \sum_{i=1}^n x_i}$$

\uparrow from θ if $\sum_{i=1}^n x_i = \sum_{i=1}^n y_i$

$$L(\theta | \underline{x}) = L(\theta | \underline{y}) \quad \text{if} \quad \sum_{i=1}^n x_i = \sum_{i=1}^n y_i;$$

So, inference on θ based on two samples \underline{x} and \underline{y} should be the same if $\sum_{i=1}^n x_i = \sum_{i=1}^n y_i$.

- For two samples if the no. of successes are the same, Then our inference about θ should also match.

$$L(\theta | \underline{x}) = p_\theta(\underline{x})$$

$$\frac{L(\theta | \underline{x})}{L(\theta | \underline{y})} = \frac{p_\theta(\underline{x})}{p_\theta(\underline{y})} \leftarrow \text{is free of } \theta \text{ iff } T(\underline{x}) \sim T(\underline{y})$$

$\Leftrightarrow T(\underline{X})$ is min. suff.

\therefore If $T(\underline{x})$ is minimal sufficient for θ

and $T(\underline{x}) = T(\underline{y})$, Then

$$L(\theta | \underline{x}) = C(\underline{x}, \underline{y}) L(\theta | \underline{y})$$

So, inference on θ based on \underline{x} and that based on
 \underline{y} should be the same.

Therefore, inference on θ based on \underline{x} should go
via $T(\underline{x})$.

That is, our inference based on the sample only,
depends on the minimal sufficient statistic.

Point Estimation

$$X_1, \dots, X_n \stackrel{iid}{\sim} f_{\theta}, \quad \theta \in \mathbb{H}$$

We want to estimate θ .

A point estimator of θ is any statistic $T(\underline{x})$ that takes values in \mathbb{H} .

With the observed sample \underline{x}_1 , we get the estimate $T(\underline{x}_1)$.

- Which T to use? How to estimate?
- What can we say about the qualities of T ?

How to estimate? How do we construct estimators?

We will discuss about :

- Method of moments
- Maximum likelihood estimator

Method of moments estimators

$$x_1, \dots, x_n \stackrel{\text{i.i.d}}{\sim} f_{\theta} \quad \theta = (\theta_1, \dots, \theta_K)$$
$$\mu_j' = E_{\theta}(x^j), \quad m_j' = \frac{1}{n} \sum_{i=1}^n x_i^j$$

$\uparrow a_j'(\theta)$: depends on θ

$\uparrow b_j$: free of θ
can be observed.

MOM estimators are obtained by solving:

$$\begin{aligned} a_1(\theta) &= b_1 \\ a_2(\theta) &= b_2 \\ &\vdots \\ a_K(\theta) &= b_K \end{aligned} \quad \left. \right\}$$

Another version :

$$\mu_1 = \mathbb{E}_\theta(X), \quad m_1 = \frac{1}{n} \sum_{i=1}^n x_i \leftarrow b_1,$$

$$\mu_j = \mathbb{E}_\theta((X - \mu_1)^j) \quad j \geq 2,$$

$$m_j = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^j$$

The MLE estimators $\hat{\theta}_1, \dots, \hat{\theta}_K$ are obtained by solving :

$$\begin{aligned} a_1(\theta) &= b_1 \\ \tilde{a}_2(\theta) &= \tilde{b}_2 \\ &\vdots \\ \tilde{a}_K(\theta) &= \tilde{b}_K \end{aligned} \quad \left. \right\}$$

$$\mu'_j = \mathbb{E}_\theta (x^j) = \mathbb{E}_\theta \left((\underbrace{x - \mu_1 + \mu_1}_{Y})^j \right)$$

$$= \mathbb{E}_\theta \left\{ \sum_{k=0}^j \binom{j}{k} (x - \mu_1)^{j-k} \mu_1^k \right\}$$

$$= \sum_{k=0}^j \binom{j}{k} \mu_1^k \underbrace{\mathbb{E}_\theta ((x_1 - \mu_1)^{j-k})}_{\mu_{j-k}}$$

$$= \sum_{k=0}^j \binom{j}{k} \mu_1^k \mu_{j-k}$$

$$\mu_j = \mathbb{E}_\theta ((x - \mu_1)^j) = \mathbb{E}_\theta \left\{ \sum_{k=0}^j \binom{j}{k} x^k (-\mu_1)^{j-k} \right\}$$

$$= \sum_{k=0}^j (-1)^{j-k} \binom{j}{k} \mu_1^{(j-k)} \mu_k'$$

- we can get the central moments from the raw moments by using linear combinations.
We can also get the raw moments from the central moments using linear combinations.
- . This is true for both the population moments and the sample moments.

Therefore, the MOM estimators based on the two approaches are the same.

We use one or the other based on which one is more convenient to us.

Maximum likelihood estimator:

$$x_1, \dots, x_n \stackrel{\text{iid}}{\sim} f_{\theta}, \quad \theta \in \mathbb{H}$$

$$L(\theta | \underline{x}) = \prod_{i=1}^n f_{\theta}(x_i), \quad \theta \in \mathbb{H}$$

The maximum likelihood estimator (MLE) of θ is the maximizer of the likelihood.

$$\hat{\theta} = \hat{\theta}(\underline{x}) = \underset{\theta \in \mathbb{H}}{\operatorname{argmax}} L(\theta | \underline{x}) \quad \underline{x} \in \mathcal{X}.$$

$\hat{\theta}(\underline{x})$ is the MLE.

Eg. X_1, \dots, X_n iid Bernoulli(θ), $0 \leq \theta \leq 1$.

$$L(\theta | \underline{x}) = \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i}, \quad 0 \leq \theta \leq 1$$

log-likelihood:

$$\ell(\theta | \underline{x}) = \log L(\theta | \underline{x}) = \sum_{i=1}^n x_i \log \theta + (n - \sum_{i=1}^n x_i) \log(1-\theta)$$

$$\frac{\partial}{\partial \theta} \ell(\theta | \underline{x}) = \sum_{i=1}^n x_i \cdot \frac{1}{\theta} - (n - \sum_{i=1}^n x_i) \cdot \frac{1}{1-\theta}$$

$$\frac{\partial^2}{\partial \theta^2} \ell(\theta | \underline{x}) = - \frac{\sum_{i=1}^n x_i}{\theta^2} - \frac{(n - \sum_{i=1}^n x_i)}{(1-\theta)^2} < 0 \quad \forall \theta \in (0, 1)$$

$$\begin{aligned} \frac{\partial}{\partial \theta} \ell(\theta | \underline{x}) = 0 \quad (\Rightarrow) \quad \frac{\sum_{i=1}^n x_i}{\theta} &= \frac{n - \sum_{i=1}^n x_i}{1-\theta} \quad (\Leftarrow) \quad \frac{1-\theta}{\theta} = \frac{n - \sum_{i=1}^n x_i}{\sum_{i=1}^n x_i} \\ \frac{1}{\theta} &= \frac{n}{\sum_{i=1}^n x_i} \quad \therefore \theta = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}. \quad \hat{\theta}_{MLE} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i. \end{aligned}$$

Eg. X_1, \dots, X_n iid Uniform $(0, \theta)$, $\theta > 0$

$$L(\theta | x) = \begin{cases} \left(\frac{1}{\theta}\right)^n, & 0 \leq x_1, \dots, x_n \leq \theta \\ 0, & \text{o.w.} \end{cases}$$

$$= \frac{1}{\theta^n} \mathbb{1}\{x_{(n)} \leq \theta\} \mathbb{1}\{x_{(1)} > 0\}$$

$\theta \mapsto \frac{1}{\theta^n}$ is decreasing.

$\therefore L(\theta | x)$ is max. at the smallest value of θ

which is $x_{(n)}$.

$$\therefore \hat{\theta}_{MLE} = X_{(n)}.$$

Eg. X_1, \dots, X_n iid Bernoulli(θ), $0 \leq \theta \leq 1$

θ : prob. of success in one trial.

θ^2 : prob of getting two successes in two trials

$$\hat{\theta}_{MLE} = \bar{X}$$

What is the MLE of θ^2 ? \bar{X}^2 . Since MLE's are invariant.

$$X_1, \dots, X_n \sim p_{\theta}(\cdot), \quad \theta \in \mathbb{H}$$

j.t. dist (pdf/pmf) of the random sample

likelihood function is

$$L(\theta) = L(\theta | \underline{x}) = p_{\theta}(\underline{x}), \quad \theta \in \mathbb{H}, \quad \underline{x} \in \mathcal{X}$$

Maximum likelihood estimator (MLE) is

$$\hat{\theta}_{MLE} = \underset{\theta \in \mathbb{H}}{\operatorname{argmax}} L(\theta | \underline{x})$$

$$\hat{\theta}(x) = \underset{\theta \in \mathbb{H}}{\operatorname{argmax}} L(\theta | x), \quad \forall x \in \mathcal{X}$$

$$\hat{\theta}_{MLE} = \hat{\theta}(x).$$

log-likelihood

$$l(\theta) = l(\theta | \underline{x}) = \log L(\theta | \underline{x}), \quad \theta \in \mathbb{R}, \quad \underline{x} \in \mathcal{X}.$$

Since $t \mapsto \log(t)$ is strictly increasing

$$\operatorname{argmax}_{\theta \in \mathbb{R}} L(\theta | \underline{x}) = \operatorname{argmax}_{\theta \in \mathbb{R}} l(\theta | \underline{x})$$

If $l(\theta) = l(\theta | \underline{x})$ is twice differentiable,

Then $\hat{\theta}_{\text{MLE}}$ is found by solving

$$l'(\theta) = 0 \leftarrow \text{likelihood equation.}$$

and verifying $l''(\hat{\theta}) < 0$

$\hat{\theta}$ is obtained as a solution to the likelihood equation.

The parameter space \mathbb{H} plays an important role in the determination of MLE.

Eg. $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta)$

$$\hat{\theta}_{\text{MLE}} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad (\text{True if } \theta \in [0, 1])$$

Suppose $\mathbb{H} = \left(-\frac{1}{4}, \frac{1}{2} \right) . \frac{1}{4} < \theta < \frac{1}{2}$

$$\hat{\theta}_{\text{MLE}} = \operatorname{argmax}_{\theta \in \mathbb{H}} L(\theta | \underline{x})$$

Eg. $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, 1), \mu \in \mathbb{R} \quad \mathbb{H} = \mathbb{R}$

$$\hat{\mu}_{\text{MLE}} = \bar{X}$$

If $\mu > 0$.

$$\hat{\theta}_{\text{MLE}} = \begin{cases} \bar{X} & \text{if } \bar{X} > 0 \\ 0 & \text{if } \bar{X} \leq 0 \end{cases} \quad (\text{Technically undefined})$$

MLE of different parameterization

Eg. $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta)$, $\theta \in [0, 1]$

θ ← prob. of success in a single trial

θ^2 ← prob. of two successes in two trials

We are interested in θ^2 .

What is the MLE of θ^2 ? (\bar{X}^2)

$$L(\theta) = L(\theta | \underline{x}) = \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i}, \theta \in [0, 1]$$

$$\left(t = \sum_{i=1}^n x_i \right) = \theta^t (1-\theta)^{n-t}, \theta \in [0, 1]$$

$$\eta = \theta^2 \quad (\Rightarrow) \quad \theta = \sqrt{\eta}, \quad \theta \in [0, 1] \Leftrightarrow \eta \in [0, 1]$$

↑ parameter of interest.

$$\tilde{L}(\eta) = [L(\theta)]_{\theta=\sqrt{\eta}} = [\theta^t (1-\theta)^{n-t}]_{\theta=\sqrt{\eta}}$$

$$= (\sqrt{\eta})^t (1-\sqrt{\eta})^{n-t}, \quad \eta \in [0, 1]$$

$$\tilde{\ell}(\eta) = \frac{t}{2} \log \eta + (n-t) \log (1-\sqrt{\eta})$$

$$\begin{aligned} \tilde{\ell}'(\eta) &= \frac{t}{2\eta} + (n-t) \frac{-\frac{1}{2\sqrt{\eta}}}{1-\sqrt{\eta}} \\ &= \frac{t}{2\eta} - \frac{n-t}{2\sqrt{\eta}(1-\sqrt{\eta})} \end{aligned}$$

$$\tilde{\ell}''(\eta) = -\frac{t}{2\eta^2} + \frac{(n-t) \frac{1}{2\sqrt{\eta}} - 1}{2\eta(1-\sqrt{\eta})^2}$$

$$\tilde{\ell}''(\bar{\eta}) < 0 \quad (\text{check this}).$$

$$\therefore \hat{\theta}_{MLE} = \bar{x}^2$$

$$\begin{aligned} \tilde{\ell}'(\eta) &= 0 \\ \Rightarrow \frac{t}{2\eta} &= \frac{n-t}{2\sqrt{\eta}(1-\sqrt{\eta})} \\ \Rightarrow \frac{t}{n-t} &= \frac{\sqrt{\eta}}{1-\sqrt{\eta}} \end{aligned}$$

$$\Rightarrow \frac{t}{n} = \sqrt{\eta}$$

$$\Rightarrow \eta = \left(\frac{t}{n}\right)^2 = \bar{x}^2$$

Now, take an arbitrary likelihood $L(\theta) = L(\theta | \underline{x})$
 $\theta \in \Theta$

Suppose g is a 1-1 function

We are interested in $\eta = g(\theta)$, $\eta \in g(\Theta)$
 $= \{g(\theta) : \theta \in \Theta\}$

$$\eta = g(\theta), g \text{ is 1-1} (\Rightarrow) \theta = g^{-1}(\eta)$$

Now, $\tilde{L}(\eta) = [L(\theta)]_{\theta = g^{-1}(\eta)} = L(g^{-1}(\eta)), \eta \in g(\Theta)$

$$\hat{\eta}(\underline{x}) = \underset{\eta \in g(\Theta)}{\operatorname{argmax}} \tilde{L}(\eta | \underline{x})$$

$$= \underset{\eta \in g(\Theta)}{\operatorname{argmax}} L(g^{-1}(\eta))$$

$L(\theta)$ is max. if $\theta = \hat{\theta}(\underline{x}) \Rightarrow \hat{L}(\eta)$ is max if
 $g^{-1}(\eta) = \hat{\theta}(\underline{x})$

$$\hat{\eta}(\underline{x}) = \underset{\eta \in g(\mathbb{H})}{\operatorname{argmax}} \tilde{L}(\eta | \underline{x}) = \underset{\eta \in g(\mathbb{H})}{\operatorname{argmax}} L(g^{-1}(\eta) | \underline{x})$$

$L(\theta | \underline{x})$ is max. if $\theta = \hat{\theta}(\underline{x}) \Rightarrow \tilde{L}(\eta | \underline{x})$ is max. if $g^{-1}(\eta) = \hat{\theta}(\underline{x})$

($\Rightarrow \tilde{L}(\eta | \underline{x})$ is max. if $\eta = g(\hat{\theta}(\underline{x}))$)

$$\therefore \hat{\eta}_{MLE} = g(\hat{\theta}_{MLE})$$

This property holds even if g is not 1-1.

Result: Suppose that in a particular statistical problem $\hat{\theta}_{MLE}$ is the MLE of a parameter θ . If $\eta = g(\theta)$ is another parameter, then the MLE of η is $\hat{\eta}_{MLE} = g(\hat{\theta}_{MLE})$.

In many cases, closed-form expression for MLE does not exist.

Eg. Suppose X_1, \dots, X_n ^{iid} ~ Cauchy($\theta, 1$). $\theta \in \mathbb{R}$

$$f_\theta(x) = \frac{1}{\pi} \frac{1}{1 + (x - \theta)^2}, \quad x \in \mathbb{R}$$

$$L(\theta) = L(\theta | x) = \prod_{i=1}^n f_\theta(x_i) = \frac{1}{\pi^n} \prod_{i=1}^n \frac{1}{1 + (x_i - \theta)^2}$$

We know that a unique maximizer of $L(\cdot)$ exists.

A closed-form expression for the MLE does not exist.

In such cases, we use some kind of optimization technique to find the MLE.

Suppose $\ell(\theta|x) = \log L(\theta|x)$, $\theta \in \mathbb{H}$

Let ℓ' and ℓ'' exist

So, $\hat{\theta}_{MLE}$ is obtained by solving $\ell'(\theta) = 0$.

However, it cannot be solved analytically.

Newton-Raphson method : To find the root of $f(x)=0$

Start with x_0 and define $x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$

$$x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$$

:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

till convergence

The root is taken to be

The convergent point.

In our case, $f = \ell'$, $f' = \ell''$.

So, our iterative procedure would look like :

- Start with $\theta_0 \leftarrow \text{initial guess}$
- Define $\theta_{n+1} = \theta_n - \frac{\ell'(\theta_n)}{\ell''(\theta_n)}$, $n=0, 1, 2, \dots$
- continue till convergence.
- $\hat{\theta}_{MLE}$ is the convergent point

$$l(\theta) = l(\theta | x) = \log L(\theta | x) = \log \left(\prod_{i=1}^n f_\theta(x_i) \right)$$

$$= \sum_{i=1}^n \log f_\theta(x_i)$$

$$l'(\theta) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f_\theta(x_i)$$

$$l''(\theta) = \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f_\theta(x_i)$$

If often true:

$$\mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \log f_\theta(x) \right] = 0$$

$$\mathbb{E}_\theta \left[\left\{ \frac{\partial}{\partial \theta} \log f_\theta(x) \right\}^2 \right]$$

$$= - \mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(x) \right]$$

$I(\theta) = \text{Fisher's information}$

If, instead of using $\ell''(\theta)$, we use

$E_{\theta}[\ell''(\theta | \underline{x})]$ in Newton-Raphson iterations,

then the resulting algorithm is known as Fisher's Scoring.

So, the iterates become

$$\cdot \theta_{n+1} = \theta_n - \frac{\ell'(\theta_n)}{E_{\theta}[\ell''(\theta_n | \underline{x})]}$$

$$= \theta_n + \frac{\ell'(\theta_n)}{n I(\theta_n)}, \text{ where } I(\theta) \text{ is the Fisher information.}$$

MLE with multiple parameters

$$L(\underline{\theta}) = L(\underline{\theta} | \underline{x}) = p_{\underline{\theta}}(\underline{x}), \quad \underline{\theta} \in \mathbb{H}, \quad \forall \underline{x} \in \mathcal{X}$$

$$\hat{\underline{\theta}}(\underline{x}) = \underset{\underline{\theta} \in \mathbb{H}}{\operatorname{argmax}} L(\underline{\theta} | \underline{x})$$

$$\hat{\underline{\theta}}_{\text{MLE}} = \hat{\underline{\theta}}(\underline{x})$$

$$l(\underline{\theta} | \underline{x}) = \log L(\underline{\theta} | \underline{x}), \quad \underline{\theta} \in \mathbb{H}$$

If $l(\underline{\theta})$ is twice-differentiable (all second order partial derivatives exist)

Then we solve

$$\nabla l(\underline{\theta}) = 0 \quad \leftarrow \text{likelihood equations}$$

$\nabla^2 l(\underline{\theta})$ is negative-definite.

In case a closed form does not exist, we can use a similar algorithm to Newton-Raphson.

$$\underline{\theta}_0$$

$$\underline{\theta}_{n+1} = \underline{\theta}_n - (\nabla^2 \ell(\underline{\theta}_n))^{-1} \nabla \ell(\underline{\theta}_n), \quad n=0, 1, \dots$$

till convergence.

Fisher Scoring.

$$\underline{\theta}_{n+1} = \underline{\theta}_n - (n I(\underline{\theta}_n))^\top \nabla \ell(\underline{\theta}_n)$$

$I(\theta)$ is the Fisher's information matrix.

How to judge an estimator?

We have a model $X_1, \dots, X_n \sim f_\theta$, $\theta \in \mathbb{H}$

θ is the parameter of interest.

$T = T(X)$ is our estimator.

$|T - \theta|^2$ ← gives the departure of our estimator
from the truth
↑
a random variable

Take expectation:

$$MSE_\theta(T) = E_\theta (T - \theta)^2 \leftarrow \text{mean squared error.}$$

$$\begin{aligned} MSE_\theta(T) &= E_\theta (T - \theta)^2 = E_\theta (T - E_\theta(T) + (E_\theta(T) - \theta))^2 \\ &= E_\theta \left\{ (T - E_\theta(T))^2 \right\} + E_\theta \left\{ (E_\theta(T) - \theta)^2 \right\} + 2 E_\theta \left\{ (T - E_\theta(T)) (E_\theta(T) - \theta) \right\} \end{aligned}$$

$$\text{MSE}_\theta(T) = \mathbb{E}_\theta \left\{ (T - \mathbb{E}_\theta(T))^2 + [\mathbb{E}_\theta(T) - \theta]^2 \right\}$$

$$= \text{var}_\theta(T) + \text{Bias}_\theta^2(T) \quad \text{Bias}_\theta(T) \text{ or } B_\theta(T)$$

L bias variance decomposition

We would like to have the estimator which has the minimum MSE.

More precisely, we would like to have T such that

$$\text{MSE}_\theta(T) \leq \text{MSE}_\theta(S) \quad \forall \theta \in \Theta \text{ and all estimator } S.$$

No such estimator exists.

Take any $\theta_0 \in \Theta$

Define $S = \theta_0$ with form 1. \leftarrow The statistic S always takes the value. irrespective of the sample

Then, $MSE_{\theta_0}(S) = 0$

and $MSE_{\theta_0}(T) > 0$ unless $T = \theta_0$ with prob 1.

So, setting an estimator which has lower MSE than any other estimator for all values of the parameters is not possible.

To set something meaningful, we need to restrict our class of estimators.

In particular, we focus on estimators which are unbiased. That is, $Bias_{\theta}(T) = 0 \forall \theta$

Def: An estimator T of θ is said to be unbiased if $E_{\theta}(T) = \theta \quad \forall \theta \in \Theta$.

Def: An estimator T is said to be best unbiased estimator or uniformly minimum variance unbiased estimator (UMVUE) of θ if it has the least variance among all unbiased estimators of θ .

That is : • $IE_\theta(T) = \theta \quad \forall \theta \in \Theta$

• $\forall S$ st. $IE_\theta(S) = \theta \quad \forall \theta \in \Theta$

$$\text{var}_\theta(T) \leq \text{var}_\theta(S) \quad \forall \theta \in \Theta$$

This also entails that $MSE_\theta(T) \leq MSE_\theta(S) \quad \forall \theta \in \Theta$.

$p_{\theta}(\underline{x})$ ← joint pmf/pdf of the ^{random} sample (X_1, \dots, X_n)

X_1, \dots, X_n ← random sample

$$p_{\theta}(\underline{x}) = P_{\theta} [X_1 = x_1, \dots, X_n = x_n] \leftarrow X_1, \dots, X_n \text{ discrete}$$

= joint pdf of (X_1, \dots, X_n) evaluated at (x_1, \dots, x_n)

Take any estimator $T = T(\underline{x})$. Then,

$$E_{\theta} [T(\underline{x})] = \int_{\mathcal{X}} T(\underline{x}) p_{\theta}(\underline{x}) d\underline{x}, \quad \begin{array}{l} \mathcal{X} \text{ is the} \\ \text{support of } p_{\theta}(\cdot) \end{array}$$

\mathcal{X} can depend on θ .

If the LHS is differentiable w.r.t. θ , then

$$\frac{d}{d\theta} E_{\theta} [T(\underline{x})] = \frac{d}{d\theta} \int_{\mathcal{X}} T(\underline{x}) p_{\theta}(\underline{x}) d\underline{x}$$

$$\frac{d}{d\theta} E_{\theta}[T(\underline{x})] = \frac{d}{d\theta} \int_{\mathcal{X}} T(\underline{x}_i) p_{\theta}(\underline{x}_i) d\underline{x}_i$$

$$= \int_{\mathcal{X}} \frac{\partial}{\partial \theta} \left\{ T(\underline{x}_i) p_{\theta}(\underline{x}_i) \right\} d\underline{x}_i \quad \begin{matrix} \text{(assuming this} \\ \text{change of differentiation} \\ \text{and integral is)} \\ \text{permissible} \end{matrix}$$

$$= \int_{\mathcal{X}} T(\underline{x}_i) \frac{\partial}{\partial \theta} p_{\theta}(\underline{x}_i) d\underline{x}_i$$

$$= \int_{\mathcal{X}} T(\underline{x}_i) \frac{\frac{\partial}{\partial \theta} p_{\theta}(\underline{x}_i)}{p_{\theta}(\underline{x}_i)} p_{\theta}(\underline{x}_i) d\underline{x}_i$$

$$= \int_{\mathcal{X}} T(\underline{x}_i) \frac{\partial}{\partial \theta} \ln p_{\theta}(\underline{x}_i) p_{\theta}(\underline{x}_i) d\underline{x}_i$$

$$= E_{\theta} \left[T(\underline{x}) \frac{\partial}{\partial \theta} \ln p_{\theta}(\underline{x}) \right]$$

For any estimator $T = T(\underline{x})$,

$$\frac{d}{d\theta} \mathbb{E}_{\theta} [T(\underline{x})] = \mathbb{E}_{\theta} \left[T(\underline{x}) \frac{\partial}{\partial \theta} \log p_{\theta}(\underline{x}) \right]$$

$$\mathbb{E}_{\theta} \left[\frac{\partial}{\partial \theta} \log p_{\theta}(\underline{x}) \right] = \int_{\mathcal{X}} \frac{\partial}{\partial \theta} \log p_{\theta}(x_i) p_{\theta}(x_i) dx_i$$

$$= \int_{\mathcal{X}} \frac{\frac{\partial}{\partial \theta} p_{\theta}(x_i)}{p_{\theta}(x_i)} p_{\theta}(x_i) dx_i$$

$$= \int_{\mathcal{X}} \frac{1}{\partial \theta} p_{\theta}(x_i) dx_i$$

$$= \frac{d}{d\theta} \int_{\mathcal{X}} p_{\theta}(x_i) dx_i$$

$$= \frac{d}{d\theta} | = 0$$

[assuming the interchange is formal]

Now,

$$\begin{aligned} \underset{\theta}{\text{cov}} \left(T(\underline{x}), \frac{\partial}{\partial \theta} \log p_{\theta}(\underline{x}) \right) &= \mathbb{E}_{\theta} \left[T(\underline{x}) \frac{\partial}{\partial \theta} \log p_{\theta}(\underline{x}) \right] - \\ &\quad \mathbb{E}_{\theta} [T(\underline{x})] \mathbb{E}_{\theta} \left[\frac{\partial}{\partial \theta} \log p_{\theta}(\underline{x}) \right] \\ &= \mathbb{E}_{\theta} \left[T(\underline{x}) \frac{\partial}{\partial \theta} \log p_{\theta}(\underline{x}) \right] \\ &= \frac{d}{d\theta} \mathbb{E}_{\theta} [T(\underline{x})] \end{aligned}$$

For two random variables Y and Z

$$\{ \text{cov}(Y, Z) \}^2 \leq \text{var}(Y) \text{var}(Z)$$

Cauchy-Schwarz inequality.

$$\Rightarrow \text{var}(Y) \geq \frac{\{ \text{cov}(Y, Z) \}^2}{\text{var}(Z)}$$

$$\begin{aligned} \therefore \text{var}_{\theta}(T(\underline{x})) &\geq \frac{\left\{ \text{cov}_{\theta}\left(T(\underline{x}), \frac{\partial}{\partial \theta} \ln p_{\theta}(\underline{x})\right) \right\}^2}{\text{var}_{\theta}\left(\frac{\partial}{\partial \theta} \ln p_{\theta}(\underline{x})\right)} \\ &= \frac{\left\{ \frac{d}{d\theta} \mathbb{E}_{\theta}[T(\underline{x})] \right\}^2}{\mathbb{E}_{\theta}\left[\left\{ \frac{\partial}{\partial \theta} \ln p_{\theta}(\underline{x}) \right\}^2\right]}, \quad \left[\because \mathbb{E}_{\theta}\left[\frac{\partial}{\partial \theta} \ln p_{\theta}(\underline{x})\right] = 0 \right] \end{aligned}$$

So, the variance of any estimator must satisfy this bound. This bound is known as the Cramér-Rao lower bound (for the variance of an estimator).

If T is an unbiased estimator of θ , then

$$\text{var}_{\theta}(T(\underline{x})) \geq \frac{1}{\mathbb{E}_{\theta}\left[\left\{ \frac{\partial}{\partial \theta} \ln p_{\theta}(\underline{x}) \right\}^2\right]}.$$

$$\mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \ln p_\theta(\underline{x}) \right] = 0 \quad \forall \theta$$

$$(\Rightarrow) \int \frac{\partial}{\partial \theta} \ln p_\theta(\underline{x}) p_\theta(\underline{x}) d\underline{x} = 0 \quad \forall \theta$$

$$\Rightarrow \frac{d}{d\theta} \int \frac{\partial}{\partial \theta} \ln p_\theta(\underline{x}) p_\theta(\underline{x}) d\underline{x} \rightarrow$$

$$\Rightarrow \int \frac{\partial}{\partial \theta} \left\{ \frac{\partial}{\partial \theta} \ln p_\theta(\underline{x}) p_\theta(\underline{x}) \right\} d\underline{x} = 0$$

$$\Rightarrow \int \left\{ \frac{\partial^2}{\partial \theta^2} \ln p_\theta(\underline{x}) p_\theta(\underline{x}) + \frac{\partial}{\partial \theta} \ln p_\theta(\underline{x}) \frac{\partial}{\partial \theta} p_\theta(\underline{x}) \right\} d\underline{x} = 0$$

$$\Rightarrow \int \left\{ \frac{\partial^2}{\partial \theta^2} \ln p_\theta(\underline{x}) p_\theta(\underline{x}) + \left\{ \frac{\partial}{\partial \theta} \ln p_\theta(\underline{x}) \right\}^2 p_\theta(\underline{x}) \right\} d\underline{x} = 0$$

$$\int_{\mathcal{X}} \frac{\partial^2}{\partial \theta^2} \log p_\theta(\underline{x}) p_\theta(\underline{x}) d\underline{x} + \int_{\mathcal{X}} \left\{ \frac{\partial}{\partial \theta} \log p_\theta(\underline{x}) \right\}^2 p_\theta(\underline{x}) d\underline{x} = 0$$

$$\Rightarrow \int_{\mathcal{X}} \left\{ \frac{\partial}{\partial \theta} \log p_\theta(\underline{x}) \right\}^2 p_\theta(\underline{x}) d\underline{x} = - \int_{\mathcal{X}} \frac{\partial^2}{\partial \theta^2} \log p_\theta(\underline{x}) p_\theta(\underline{x}) d\underline{x}$$

$$\Rightarrow \mathbb{E}_\theta \left[\left| \frac{\partial}{\partial \theta} \log p_\theta(\underline{x}) \right|^2 \right] = - \mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \log p_\theta(\underline{x}) \right]$$

Then,

$$\text{Var}_\theta (\tau(\underline{x})) \geq - \frac{1}{\mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \log p_\theta(\underline{x}) \right]}.$$

When $X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta$

$$p_\theta(\underline{x}) = \prod_{i=1}^n f_\theta(x_i)$$

$$\log p_\theta(\underline{x}) = \sum_{i=1}^n \log f_\theta(x_i)$$

$$\frac{\partial}{\partial \theta} \log p_\theta(\underline{x}) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f_\theta(x_i)$$

$$\frac{\partial^2}{\partial \theta^2} \log p_\theta(\underline{x}) = \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f_\theta(x_i)$$

$$\therefore -E_\theta \left[\frac{\partial^2}{\partial \theta^2} \log p_\theta(\underline{x}) \right] = - \sum_{i=1}^n E_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(x_i) \right]$$

$$\therefore -E_\theta \left[\frac{\partial^2}{\partial \theta^2} \log p_\theta(\underline{x}) \right] = -n E_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(x_i) \right] \\ = n I(\theta)$$

$I(\theta) = -E_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \log f_{\theta}(x_1) \right]$ is the Fisher information.

We can also start with

$$\frac{\partial}{\partial \theta} \log p_{\theta}(x) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f_{\theta}(x_i) \text{ and show that}$$

$$E_{\theta} \left[\left\{ \frac{\partial}{\partial \theta} \log p_{\theta}(x) \right\}^2 \right] = n E_{\theta} \left[\left\{ \frac{\partial}{\partial \theta} \log f_{\theta}(x_1) \right\}^2 \right] \\ = n I(\theta).$$

$$\therefore \text{var}_{\theta} (T(x)) \geq \underbrace{\frac{d}{d\theta} E_{\theta} [T(x)]^2}_{n I(\theta)}$$

If T is unbiased for θ , Then $\frac{d}{d\theta} E_{\theta} [T(x)] = 1$

If T is unbiased for $a(\theta)$, Then $\frac{d}{d\theta} E_{\theta} [T(x)] = a'(\theta)$.