

## Ch7.5 Maximum likelihood estimators (MLEs)

$\theta$  fixed, not random

Maximum likelihood estimation relies on likelihood functions while avoiding the use of prior distributions and loss functions.

Given a likelihood function, it determines the value of  $\theta$  that maximizes it.

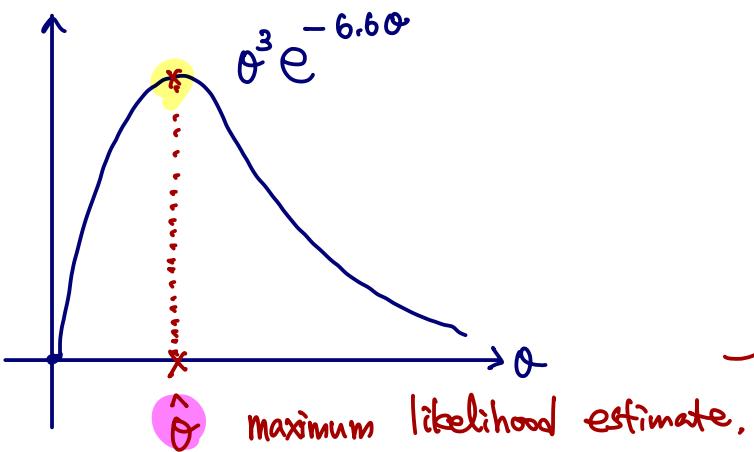
Let  $X_1, \dots, X_n$  be a random sample from a distribution with  $f(x|\theta)$  where  $\theta$  is unknown. For fixed  $x = (x_1, \dots, x_n)$ , the joint probability distribution  $f_n(x|\theta)$  is called the likelihood function when it is regarded as a function of  $\theta$  given the observed data  $x_1, \dots, x_n$ .

It quantifies how likely different values of  $\theta$  are given the data.

\* Aim to find a value of  $\theta$  that is most consistent with data.

**EG. [Lifetimes of Electronic Components]** Lifetimes  $X_1, \dots, X_n$  of electronic components are modeled as iid exponential random variables with parameter  $\theta$ . The observed data are  $x_1 = 3, x_2 = 1.5, x_3 = 2.1$ .

$$f_3(x_1, x_2, x_3 | \theta) \stackrel{\text{iid}}{=} \prod_{i=1}^3 f(x_i | \theta) \stackrel{\text{iid}}{=} \prod_{i=1}^3 \theta e^{-\theta x_i} = \theta^3 e^{-\theta \sum_{i=1}^3 x_i} = \theta^3 e^{-6.6\theta}$$



With  $\hat{\theta}$ , can explain data best in terms of likelihood under the statistical model "iid exponential".

**Definition.** For each  $x$ , let  $\delta(x) \in \Omega$  be such that  $\delta(x)$  maximizes the likelihood function  $f_n(x|\theta)$ .

← MLE.

$\hat{\theta} = \delta(x)$  is called the maximum likelihood estimate of  $\theta$ .

$\hat{\theta} = \delta(X)$  is called the maximum likelihood estimator of  $\theta$ .

$f^* \leftarrow$  Bayes estimators.

To simplify computations (especially when dealing with products of probabilities), the logarithm of the likelihood function is often used.

*natural log.*

**Definition.** Call  $L(\theta) = \log f_n(x|\theta)$  the log likelihood function.

**Fact:** The log function is monotonically increasing on  $(0, \infty)$   
 value of  $\theta$  that maximizes  $L(\theta)$   
 !!

Value of  $\theta$  that maximizes  $f_n(x|\theta)$

Find  $\theta$  that maximizes  $L(\theta) = \hat{\theta}$  MLE.

**EG. [Lifetimes of Electronic Components]**

$$f_3(d_1, d_2, d_3 | \theta) = \theta^3 e^{-6.6\theta}.$$

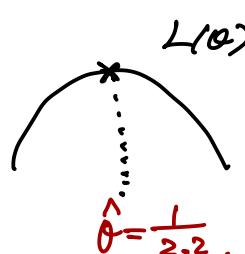
$$\Rightarrow L(\theta) = \log f_3(d_1, d_2, d_3 | \theta) = \log \theta^3 e^{-6.6\theta}.$$

$$= \log \theta^3 + \log e^{-6.6\theta} = 3 \log \theta - 6.6\theta.$$

Find  $\theta$  that maximizes  $L(\theta) = 3 \log \theta - 6.6\theta$  !!

$$\Rightarrow \frac{dL(\theta)}{d\theta} = \frac{3}{\theta} - 6.6 \underset{\text{set}}{=} 0 \Rightarrow \theta = \frac{3}{6.6} = \frac{1}{2.2}$$

Since  $\left. \frac{d^2 L(\theta)}{d\theta^2} \right|_{\theta=\frac{1}{2.2}} = -\frac{3}{\theta^2} \Bigg|_{\theta=\frac{1}{2.2}} = -\frac{3}{\frac{1}{2.2^2}} < 0$ ,



$$\hat{\theta} = \frac{1}{2.2}.$$

**EG.** Given  $X_1, \dots, X_n | \theta \sim \text{iid Bernoulli } \theta$ , where  $0 \leq \theta \leq 1$ , find the MLE  $\hat{\theta}$ .

$$f_n(d_1, \dots, d_n | \theta) = \theta^{\sum_{i=1}^n d_i} (1-\theta)^{n - \sum_{i=1}^n d_i}.$$

$$\Rightarrow L(\theta) = \log \theta^{\sum d_i} \cdot (1-\theta)^{n - \sum d_i}$$

$$= \log \theta^{\sum d_i} + \log (1-\theta)^{n - \sum d_i}$$

$$= \sum d_i \cdot \log \theta + (n - \sum d_i) \log (1-\theta).$$

1)  $\sum d_i \notin \{0, n\}$  :  $\frac{dL(\theta)}{d\theta} = \frac{\sum d_i}{\theta} - \frac{n - \sum d_i}{1-\theta} \stackrel{\text{set}}{=} 0$

$\times \theta(1-\theta)$

$$\Rightarrow (1-\theta) \sum d_i - \theta(n - \sum d_i) = 0$$

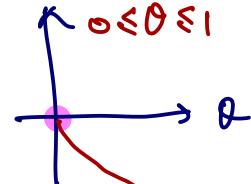
$$\Rightarrow \sum d_i - \cancel{\sum d_i \theta} - n\theta + \cancel{\sum d_i \theta} = 0$$

$$\Rightarrow \sum d_i = n\theta \Rightarrow \theta = \frac{\sum d_i}{n} = \bar{x}_n.$$

$0 \leq \theta \leq 1$

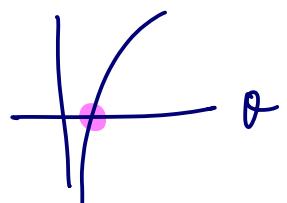
2)  $\sum d_i = 0$  ;  $L(\theta) = n \log(1-\theta)$  is decreasing

$$\hat{\theta} = 0 (= \bar{x}_n)$$



3)  $\sum d_i = n$  ;  $L(\theta) = n \log \theta$  is increasing

$$\hat{\theta} = 1 (= \bar{x}_n)$$



$$\Rightarrow \hat{\theta} = \bar{x}_n.$$

$X_1, X_2, \dots, X_n | \theta \sim \text{iid} \sim f(x| \theta)$  unknown, fixed.

Aim to construct an estimator for  $\theta$ .  
a function of  $x_1, \dots, x_n$

In MLE,  $\hat{\theta}$  that maximizes  $f_n(x_1, \dots, x_n | \theta)$   
likelihood function.  
a function of  $\theta$  given  $x_1, \dots, x_n$ .

$\hat{\theta}$  is most consistent w/ data under the likelihood  
function considered.

↖ unknown.

EG. Given  $X_1, \dots, X_n | \mu, \sigma^2 \sim \text{iid normal } \theta = (\mu, \sigma^2)$ , where  $-\infty < \mu < \infty, \sigma^2 > 0$ , find the MLE  $\hat{\theta} = (\hat{\mu}, \hat{\sigma}^2)$ .

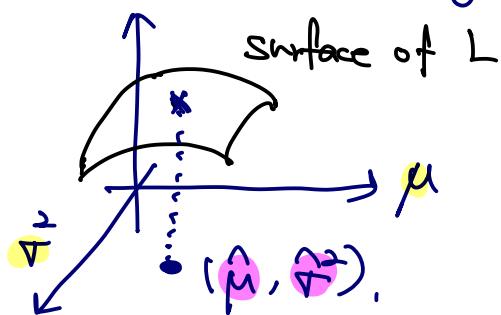
$$f_n(x_1, \dots, x_n | \mu, \sigma^2) \stackrel{\text{iid}}{=} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$= \frac{1}{\sqrt{2\pi}^n} \times \frac{1}{\sqrt{\sigma^2}^n} \times e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2} = (\sigma^2)^{-\frac{n}{2}}$$

$$\Rightarrow L(\mu, \sigma^2) = \log \frac{1}{\sqrt{2\pi}^n} + \log \frac{1}{\sqrt{\sigma^2}^n} + \log e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2}$$

$$= C - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

Note that  $L$  is smooth and has the maximum point  
(verified using 2nd order derivatives and Jacobian)



$$\textcircled{1} \quad \frac{\partial L(\mu, \sigma^2)}{\partial \mu} = 0, \quad \textcircled{2} \quad \frac{\partial L(\mu, \sigma^2)}{\partial \sigma^2} = 0.$$

$$\textcircled{1} \quad \frac{\partial L(\mu, \sigma^2)}{\partial \mu} = 0 - 0 + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu) \stackrel{\text{set}}{=} 0 \Rightarrow \sum (x_i - \mu) = 0$$

$$\Rightarrow \sum x_i - n\mu = 0 \Rightarrow \mu = \frac{\sum x_i}{n} = \bar{x}_n$$

$$\textcircled{2} \quad \frac{\partial L(\mu, \sigma^2)}{\partial \sigma^2} = 0 - \frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum (x_i - \mu)^2 \stackrel{\text{set}}{=} 0$$

$$- n\sigma^2 + \sum (x_i - \mu)^2 = 0 \Rightarrow \sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2.$$

$$\Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \bar{x}_n)^2$$

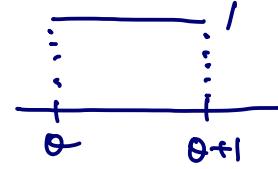
sample variance

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

Note MLEs aren't always unique. MLEs don't always exist;

Page 423  $X_1, \dots, X_n | \theta \sim \text{iid continuous uniform on } [\theta, \theta + 1]$

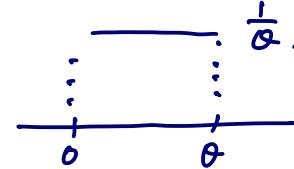
$$\begin{aligned}
f(x|\theta) &= 1 \times I(\theta \leq x \leq \theta+1) \\
\Rightarrow f_n(x_1, \dots, x_n | \theta) &= \prod_{i=1}^n 1 \times I(x_i \leq \theta+1) \\
&= I(0 \leq x_1, \dots, x_n \leq \theta+1) \\
&= I(\theta \leq \min\{x_1, \dots, x_n\} \leq \max\{x_1, \dots, x_n\} \leq \theta+1).
\end{aligned}$$



The likelihood function can be maximized at any values btw  $\max\{x_1, \dots, x_n\} - 1$  and  $\min\{x_1, \dots, x_n\}$ .  
 $\Rightarrow$  MLEs are not unique.

Page 422  $X_1, \dots, X_n | \theta \sim \text{iid continuous uniform on } (0, \theta)$

$$\begin{aligned}
f(x|\theta) &= \frac{1}{\theta} \times I(0 < x < \theta) \\
\Rightarrow f_n(x_1, \dots, x_n | \theta) &= \prod_{i=1}^n \frac{1}{\theta} I(0 < x_i < \theta) \\
&= \frac{1}{\theta^n} I(0 < x_1, \dots, x_n < \theta) \\
&= \frac{1}{\theta^n} I(0 < \min\{x_1, \dots, x_n\} \leq \max\{x_1, \dots, x_n\} < \theta)
\end{aligned}$$



Pick  $\theta$  as small as possible subject to \*

There exists no such  $\theta$  because

$\theta$  cannot be chosen equal to  $\max\{x_1, \dots, x_n\}$ .

$\Rightarrow$  MLEs of  $\theta$  do not exist!

## The MLEs for many of the Ch5 special distributions

$$X_1, \dots, X_n | \theta \sim \text{iid } f(x|\theta)$$

↙ unknown parameter.

distribution, $f(x \theta)$	MLE $\hat{\theta}$
Bernoulli $\theta$	$\hat{\theta} = \bar{X}_n$ ← proportion of successes. (lecture note).
Poisson $\theta$	$\hat{\theta} = \bar{X}_n$ #9 in Ch7.t
exponential $\theta$	$\hat{\theta} = \frac{1}{\bar{X}_n}$ #2
uniform $[a, \theta]$	$\hat{\theta} = \max \{X_1, \dots, X_n\}$ , #11.
uniform $[\theta, b]$	$\hat{\theta} = \min \{X_1, \dots, X_n\}$
normal $\mu, \sigma^2$ , both unknown	$\hat{\mu} = \bar{X}_n$ $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ lecture note.
normal $\theta, \sigma^2$ , mean unknown	$\hat{\theta} = \bar{X}_n$
normal $\mu, \theta$ , variance unknown	$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{\mu})^2$

## Ch7.6 Properties of MLEs

1. invariance property

$$\text{failure rate} = \frac{1}{\text{orange}}$$

**Q.**  $X_1, \dots, X_n | \theta \sim \text{iid exponential } \theta$ . What is the MLE of  $\theta^2$ ?

$$\hat{\theta} = \text{MLE of } \theta = \frac{1}{\bar{X}_n}$$

Can we say  $\left(\frac{1}{\bar{X}_n}\right)^2$  is an MLE of  $\theta^2$ ? Yes by invariance thm.

**Theorem (Invariance).** Let  $\hat{\theta}$  be an MLE of  $\theta$ , and let  $g(\theta)$  be a function of  $\theta$ , then an MLE of  $g(\theta)$  is  $g(\hat{\theta})$ .

$$P(X=x) = \frac{e^{-\theta} \cdot \theta^x}{x!} \quad x=0,$$

**EG.**  $X_1, \dots, X_n | \theta \sim \text{iid Poisson } \theta$ . What is the MLE of  $P(X_1 = 0 | \theta) = e^{-\theta}$ ?

$$\hat{\theta} = \bar{X}_n$$

$$\text{MLE of } e^{-\theta} = g(\theta) \Rightarrow e^{-\hat{\theta}} = e^{-\bar{X}_n}.$$

**EG.**  $X_1, \dots, X_n | \mu, \sigma^2 \sim \text{iid normal } \mu, \sigma^2$ . What the an MLE of  $\sigma$  and  $E[X_1^2]$ ?

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

$$\text{Var}[X_i] + (E[X_i])^2$$

$$* \quad \tau = \sqrt{\hat{\sigma}^2} \Rightarrow \hat{\tau} = \sqrt{\hat{\sigma}^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2} \quad \hat{\sigma}^2 + \mu^2$$

$$* \quad \sigma^2 + \mu^2 \Rightarrow \hat{\sigma}^2 + (\hat{\mu})^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 + (\bar{X}_n)^2.$$

2. Asymptotic behaviors behavior of MLEs when  $n$  is large.  
 $(n \rightarrow \infty)$

Fact: Under some regularity conditions (often hold in practice), MLEs are

1. consistent for true  $\theta_0$ , i.e.,  $\hat{\theta} \xrightarrow{P} \theta_0$ , as  $n \rightarrow \infty$ .

2. asymptotically normal, i.e.,  $\hat{\theta} \approx \text{normal } \sim \theta_0, \frac{1}{n I(\theta_0)}$ ,  
 where  $I(\theta_0)$  is Fisher Information (in Ch8).

### Method of moments (MoM)

It is a simple technique for estimating parameters. It match empirical  
 (sample) moments to the population moments of a distribution.  
 (data.) (theoretical values)

Assume that  $X_1, \dots, X_n$  form a random sample from a distribution indexed by a  $k$ -dimensional parameter  $\theta = (\theta_1, \dots, \theta_k)$ . ex) normal  $\mu, \sigma^2$  2-dimensional  
 $\theta_1 \quad \theta_2$

### Definition.

Define the  $j^{th}$  sample moments by  $m_j = \frac{1}{n} \sum_{i=1}^n X_i^j$ , for  $j = 1, 2, \dots$

Define the  $j^{th}$  population moments by  $\mu_j = E[X_i^j]$ , for  $j = 1, 2, \dots$

The method of moments estimator, denoted by  $\tilde{\theta} = (\tilde{\theta}_1, \dots, \tilde{\theta}_k)$  for  $\theta$ , is the solution to the system of  $k$  equations  $\nwarrow k$  unknown parameters.

$$\begin{cases} m_1 = \mu_1 \\ m_2 = \mu_2 \\ \vdots \\ m_k = \mu_k \end{cases} \quad \text{Solve for } \theta = (\theta_1, \dots, \theta_k)$$

EG  $X_1, \dots, X_n | \mu, \sigma^2 \sim \text{iid normal } \mu, \sigma^2$       two unknown parameters.

Sample moment

$$m_1 = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n$$

$$m_2 = \frac{1}{n} \sum_{i=1}^n X_i^2$$

population moment

$$\mu_1 = E[X_i] = \mu$$

$$\mu_2 = E[X_i^2] = \text{Var}[X_i] + (E[X_i])^2 = \tau^2 + \mu^2$$

$$\Rightarrow \begin{cases} m_1 = \mu_1 \\ m_2 = \mu_2 \end{cases} \Rightarrow \begin{cases} \bar{X}_n = \mu \\ \frac{1}{n} \sum_{i=1}^n X_i^2 = \tau^2 + \mu^2 \end{cases} \Rightarrow \text{Solve for } \mu \text{ and } \tau^2 !$$

MoM for  $\mu$

$$\Rightarrow \hat{\mu} = \bar{X}_n \text{ and } \hat{\tau}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - (\bar{X}_n)^2 \stackrel{\oplus}{=} \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

" $\hat{\mu}$  (MLE for  $\mu$ )

" $\hat{\tau}^2$  (MLE for  $\tau^2$ )

$$\textcircled{*} \quad \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n} \sum_{i=1}^n (X_i^2 - 2X_i \bar{X}_n + \bar{X}_n^2) = \frac{1}{n} \left\{ \sum X_i^2 - 2\bar{X}_n \sum X_i + n\bar{X}_n^2 \right\}$$

$$= \frac{1}{n} \sum X_i^2 - 2(\bar{X}_n)^2 + \bar{X}_n^2 = \frac{1}{n} \sum X_i^2 - (\bar{X}_n)^2.$$

### Note.

- In some cases, MoM = MLE (see the previous example)
- When MLE  $\neq$  MoM, generally MLE is considered "better" because MLEs are asymptotically efficient! ( $n \rightarrow \infty$ ) (small variance)
  - ⇒ Roughly speaking, MLEs are more precise when  $n$  is large.
- But MoM is typically easy to compute. ↪ when you need a quick estimate!

$$\text{EG } X_1, \dots, X_n | \alpha \sim \text{iid gamma } \alpha, 1 \quad \xrightarrow{\hspace{1cm}} f(x|\alpha) = \frac{1}{\Gamma(\alpha)} \alpha^{\alpha-1} e^{-\alpha}, \alpha > 0$$

$$\text{MoM: } m_1 = \bar{x}_n \quad \text{and} \quad \mu_1 = E[X_1] = \frac{\alpha}{1} = \alpha.$$

$$\Rightarrow m_1 = \mu_1 \Rightarrow \hat{\alpha} = \bar{x}_n \leftarrow \text{easy to derive!}$$

$$\text{MLE: } f_n(x_1, \dots, x_n | \alpha) = \prod_{i=1}^n \frac{1}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-\alpha x_i} = \frac{1}{\Gamma(\alpha)^n} (\pi x_i)^{\alpha-1} e^{-\alpha \sum x_i}$$

$$\begin{aligned} \Rightarrow L(\alpha) &= \log \frac{1}{\Gamma(\alpha)^n} + \log (\pi x_i)^{\alpha-1} + \log e^{-\alpha \sum x_i} \\ &= -n \log \Gamma(\alpha) + (\alpha-1) \sum_{i=1}^n \log x_i - \sum_{i=1}^n \alpha x_i. \end{aligned}$$

$$\Rightarrow \frac{dL(\alpha)}{d\alpha} = -n \times \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} + \sum_{i=1}^n \log x_i \stackrel{\text{set}}{=} 0.$$

$$\Rightarrow \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} - \frac{1}{n} \sum_{i=1}^n \log x_i = 0.$$

Need to solve for  $\alpha$  to get MLE  $\hat{\alpha}$ !

No closed form solution!

$\Rightarrow$  Numerical method is needed!







## Ch7.7 Sufficient statistics

What if Bayes estimators, MLE, and MoM are unavailable?

R.A. Fisher introduces the concept of sufficient statistics, which provides a method for finding a good estimator for  $\theta$ .

$\hookrightarrow$  function of  $X_1, \dots, X_n$ .

**Concept:** We seek statistics  $T(X_1, \dots, X_n)$  that are **sufficient** for finding a good estimator of  $\theta$  in the sense that

1. it summarizes all the information contained in the data, and
2. knowing the individual values of  $X_1, \dots, X_n$  gives no additional information on  $\theta$ .

EG  $X_1, \dots, X_n | \mu \sim \text{iid normal } \mu, 1$ .

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

knowing  $X_1, \dots, X_n$  gives no additional information to estimate  $\mu$ .

**Definition.** Let  $X_1, \dots, X_n | \theta$  be iid random variables with  $f(x|\theta)$  and  $T = T(X_1, \dots, X_n)$  be a statistic. If the conditional distribution of  $X_1, \dots, X_n | T, \theta$  depends only on  $T$  but not on  $\theta$ , then  $T$  is a sufficient statistic.

dist. of  $X_1, \dots, X_n \sim$  depends on  $\theta$ .      dist. of  $X_1, \dots, X_n | T \sim$  independent of  $\theta$ .

EG  $X_1, \dots, X_n | \theta \sim \text{iid Bernoulli } \theta$ .  $\leftarrow$  prob. of success.

Claim:  $T(X_1, \dots, X_n) = \sum_{i=1}^n X_i$  is sufficient for  $\theta$ !  
↑ # of successes in n trials.

$$P(X_1=d_1, \dots, X_n=d_n | T=t) = \frac{P(X_1=d_1, \dots, X_n=d_n, T=t)}{P(T=t), T \sim \text{Binomial } n, \theta}$$

$$= \begin{cases} \frac{\theta^t (1-\theta)^{n-t}}{\binom{n}{t} \theta^t (1-\theta)^{n-t}} & \sum d_i = t \\ 0 & \sum d_i \neq t \end{cases}, \quad \text{is independent of } \theta.$$

## "Sufficient statistics"

unknown.

$$x_1, \dots, x_n | \theta \sim f(\cdot | \theta)$$

Consider  $T(x_1, \dots, x_n)$  to estimate  $\theta$ .

- \*  $T$  is considered "sufficient" for estimating  $\theta$ , if it captures all the information about  $\theta$ , (no need to examine  $x_1, \dots, x_n$  to estimate  $\theta$ ).
- \* Use Factorization thm. to find a sufficient statistic

$$f_n(x_1, \dots, x_n | \theta) = u(x_1, \dots, x_n) \times v[T(x_1, \dots, x_n), \theta]$$

$\Rightarrow T(x_1, \dots, x_n)$  is sufficient for  $\theta$ .

Is there a simpler method for finding sufficient statistics?

Fisher - Neyman

Theorem (Factorization Criterion). Let  $X_1, \dots, X_n | \theta$  be iid random variables with  $f(x | \theta)$  and  $T = T(X_1, \dots, X_n)$  be a statistic. Then  $T$  is **sufficient** for  $\theta$  if and only if

$f_n(d_1, \dots, d_n | \theta)$  can be factored as follows  
likelihood function

$$f_n(d_1, \dots, d_n | \theta) = u(d_1, \dots, d_n) \times v[T(d_1, \dots, d_n), \theta],$$

*depends on data*      *depends on a function of data and  $\theta$ .*

where  $u$  and  $v$  are nonnegative functions.

EG  $X_1, \dots, X_n | \theta \sim$  iid Poisson  $\theta$

$$f_n(d_1, \dots, d_n | \theta) = \prod_{i=1}^n \frac{e^{-\theta} \cdot \theta^{d_i}}{d_i!} = \frac{e^{-n\theta} \cdot \theta^{\sum d_i}}{\prod d_i!}$$

$$= \frac{1}{\prod d_i!} \times e^{-n\theta} \cdot \theta^{\sum d_i} \Rightarrow \text{By Factorization thm.}$$

*u(d<sub>1</sub>, ..., d<sub>n</sub>)*       $v[\sum d_i, \theta]$

$T(X_1, \dots, X_n) = \sum X_i$   
is sufficient for  $\theta$ .

EG  $X_1, \dots, X_n | \theta \sim$  iid uniform  $[0, \theta]$

$$f_n(d_1, \dots, d_n | \theta) = \prod_{i=1}^n \frac{1}{\theta} \times \mathbb{1}(0 \leq d_i \leq \theta) = \frac{1}{\theta^n} \times \mathbb{1}(0 \leq d_1, \dots, d_n \leq \theta)$$

$$= \mathbb{1}(0 \leq d_1, \dots, d_n) \times \frac{1}{\theta^n} \mathbb{1}(\max\{d_1, \dots, d_n\} \leq \theta)$$

*u(d<sub>1</sub>, ..., d<sub>n</sub>)*       $v[\max\{d_1, \dots, d_n\}, \theta]$

$\Rightarrow T = \max\{X_1, \dots, X_n\}$  is sufficient for  $\theta$

EG  $X_1, \dots, X_n | \theta \sim \text{iid exponential } \theta$

$$f(\mathbf{x}|\theta) = \theta e^{-\theta x_i}$$

$$f_n(x_1, \dots, x_n | \theta) = \prod_{i=1}^n \theta e^{-\theta x_i} = \theta^n e^{-\theta \sum_{i=1}^n x_i}$$

$$= \frac{1}{\theta^n} \times \theta^n e^{-\theta \sum_{i=1}^n x_i}$$

$$\therefore [x_1, \dots, x_n]$$

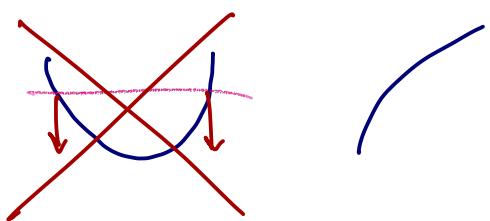
By factorization thm.,  $T(X_1, \dots, X_n) = \sum_{i=1}^n X_i$  is sufficient for  $\theta$ .

Is  $h(T)$  also sufficient for  $\theta$  when  $T$  is sufficient for  $\theta$ ?

**Theorem** Let  $T$  be a sufficient statistic for  $\theta$ , and  $h$  be a 1-1 function. Then,  $h(T)$  is a sufficient statistic for  $\theta$ .

EG  $X_1, \dots, X_n | \theta \sim \text{iid exponential } \theta$

$$\text{rate} = \frac{1}{\text{average lifetime}}$$



Is  $\hat{\theta} = \text{MLE of } \theta = \frac{1}{\bar{X}_n} = \frac{n}{\sum_{i=1}^n x_i}$  sufficient for  $\theta$ ?

1.  $\sum X_i$  is sufficient for  $\theta$ .

2.  $h(t) = \frac{n}{t}$ ,  $t > 0$ , is 1-1



$\Rightarrow h(\sum X_i) = \frac{n}{\sum X_i} = \hat{\theta}$  is sufficient for  $\theta$ .



$$\text{EG } X_1, \dots, X_n | \mu, \sigma^2 \text{ iid normal } (\mu, \sigma^2) \rightarrow f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$f_n(x_1, \dots, x_n | \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$= \frac{1}{\sqrt{2\pi}^n} \times \frac{1}{\sqrt{\sigma^2}^n} e^{-\frac{1}{2\sigma^2} \times \sum_{i=1}^n (x_i-\mu)^2}$$

$$= \frac{1}{\sqrt{2\pi}^n} \times \frac{1}{\sqrt{\sigma^2}^n} e^{-\frac{1}{2\sigma^2} \sum x_i^2 + \frac{2\mu}{2\sigma^2} \sum x_i - \frac{n\mu^2}{2\sigma^2}}$$

" " " "  $\left[ \sum x_i, \sum x_i^2, \mu, \sigma^2 \right]$

1.  $\Rightarrow \left( \sum_{i=1}^n x_i, \sum_{i=1}^n x_i^2 \right)$  is jointly sufficient for  $(\mu, \sigma^2)$ .

IS  $\left( \hat{\mu} = \bar{x}_n, \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \right)$  jointly sufficient  
 for  $(\mu, \sigma^2)$  ?

" " " "  $\frac{\sum x_i}{n}, \frac{\sum x_i^2}{n} - \frac{(\sum x_i)^2}{n^2}$  // work

2.  $h(x|y) = \left( \frac{\mathbf{x}}{n}, \frac{\mathbf{y}}{n} - \frac{\mathbf{x}^2}{n^2} \right)$  is 1-1,

$\Rightarrow (\hat{\mu}, \hat{\sigma}^2) = h(\sum x_i, \sum x_i^2)$  is jointly sufficient  
 for  $(\mu, \sigma^2)$ .

Is there a set of jointly sufficient statistics that exists across all problems?

**Definition.** Let  $X_1, \dots, X_n$  be a random sample. Order the sample

$$\min\{x_1, \dots, x_n\} = Y_1 \leq Y_2 \leq \dots \leq Y_{n-1} \leq Y_n = \max\{x_1, \dots, x_n\}.$$

$Y_1, \dots, Y_n$  are called the order statistics.

**Theorem.** The order statistics are sufficient in any random samples.

$$\{Y_1, \dots, Y_n\} = \{X_1, \dots, X_n\}.$$

**EG**  $X_1, X_2, X_3 | \mu, \sigma^2$  iid normal  $\mu, \sigma^2$

$$x_1 = 1.1, x_2 = 1.6, x_3 = 0.3$$

$$1. \quad \sum_{i=1}^3 x_i = 3, \quad \sum_{i=1}^3 x_i^2 = 3.96$$

$$2. \quad \hat{\mu} = \bar{x}_3 = 1, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^3 (x_i - \bar{x}_3)^2 = .43$$

$$3. \quad y_1 = .3, \quad y_2 = 1.1, \quad y_3 = 1.6$$

} three sets  
of sufficient  
statistics  
for  $(\mu, \sigma^2)$

To be useful, a sufficient statistic should be a simpler function of the data, capturing all the relevant information in the most concise form possible. In other words, we seek the simplest sufficient statistic! such as 1 and 2.



## Ch7.9 Improving an estimator

In this section, we show how to improve upon an estimator using sufficient statistics.

Given: A random sample  $X_1, \dots, X_n | \theta$  from  $f(x_1 | \theta)$ .  
Unknown parameter.

Goal: Improve your estimator  $\delta = \delta(X) = \delta(X_1, \dots, X_n)$ .

$\delta(X_1, \dots, X_n)$ : a function of  $X_1, \dots, X_n$ , not necessarily sufficient statistics!

We first need to

1. pick a performance metric (ex) MSE and then
2. use vocabulary from decision theory.

**Definition (Mean Squared Error).** At given  $\theta$ , the MSE of estimator  $\delta$  is

$$R(\theta, \delta) = E\left[\left(\delta(X) - \theta\right)^2\right]$$

$\nwarrow$  risk of  $\delta$

if  $R(\theta, \delta)$  is small across  $\theta$ !

### Definition.

1. Estimator  $\delta_1$  is at least as good as  $\delta_2$  if

$$R(\theta, \delta_1) \leq R(\theta, \delta_2), \quad \forall \theta \in \Omega.$$

2. Estimator  $\delta_1$  is admissible or dominates  $\delta_2$  if

- $\delta_1$  is at least as good as  $\delta_2$ , and  $\delta_1$  is strictly better than  $\delta_2$  for some  $\theta$ .
- $\exists \theta \in \Omega$  s.t.  $R(\theta, \delta_1) < R(\theta, \delta_2)$

We want to use admissible estimators!

**Theorem (Rao-Blackwell).** Let  $T$  be sufficient for  $\theta$ ,  $\delta(X)$  be an estimator, and

$$\delta_0(T) = E[\delta(X) | T] : \text{a function of } T.$$

Then,  $\delta_0$  is at least as good as  $\delta$ . If  $\delta$  is not a function of a sufficient statistic, then  $\delta_0$  dominates  $\delta$ .

EG  $X_1, \dots, X_n | p \sim \text{iid Bernoulli}(p)$ . Let  $\delta(x) = x_1$ .

Sufficient statistic for  $p$ :  $T = \sum_{i=1}^n X_i$  (Ch7.7)

Then, by R-B thm.,  $\delta_0 = E[X_1 | T]$  is at least as good as  $\delta = X_1$ , and  $\delta_0$  dominates  $\delta$  because  $X_1$  is not a function of  $T = \sum X_i$

$$\begin{aligned}\delta_0 &= E[X_1 | T] = \underbrace{1 \times P(X_1=1 | T)}_{=} + 0 \times \cancel{P(X_1=0 | T)}, \\ &= \frac{P(X_1=1 \cap T=t)}{P(T=t)} = \frac{P(X_1=1 \cap \sum_{i=1}^n X_i = t)}{P(\sum_{i=1}^n X_i = t)} \\ &= \frac{P(X_1=1 \cap \sum_{i=2}^n X_i = t-1)}{P(\sum_{i=1}^n X_i = t)} = \frac{P(X_1=1) \cancel{P(\sum_{i=2}^n X_i = t-1)}}{P(\sum_{i=1}^n X_i = t)} \\ &\quad \text{Since } X_1, \sum_{i=2}^n X_i \text{ are independent. } \sum_{i=2}^n X_i \sim \text{Binomial}(n-1, p) \\ &= \frac{P\left(\binom{n-1}{t-1} p^{t-1} (1-p)^{(n-1)-(t-1)}\right)}{\binom{n}{t} p^t (1-p)^{n-t}} \\ &= \frac{\binom{n-1}{t-1}}{\binom{n}{t}} = \frac{(n-1)!}{(t-1)! (n-t)!} \\ &= \frac{(n-1)!}{n!} \times \frac{t!}{(t-1)!} = \frac{t}{n} \Rightarrow \delta_0 = \bar{X}_n.\end{aligned}$$