### **Basic Statistics**

### Sudipta Das

Assistant Professor,

Department of Computer Science,
Ramakrishna Mission Vivekananda Educational & Research Institute

# Outline I

- Statistical Inference
  - Introduction
  - Point Estimation



Chapter 8: Statistical Inference

### Statistical Inference I

- The main objective in any statistical enquiry is the properties of one or more population.
  - However, the population(s) is (are) usually unknown to us, and we simply have a sample from the population (or, a sample from each of the given populations)

### Statistical Inference II

Statistical Inference:-

Given the properties of the sample (or, of the samples), to infer about those of the population(s) is the problem of statistical inference

- It is analogous to the inductive logic, the only difference being that the induction is achieved under probabilistic framework
  - Probability comes due to random sampling
- It is a process of going over from the known sample to unknown population.

### Statistical Inference III

### Statistical set-up of the problem of inference

- Let  $(X_1, X_2, ..., X_n)$  be a random sample of size n drawn from a population (discrete/continuous) with p.m.f/p.d.f  $f(\underline{x}; \underline{\theta}) = f_{\underline{\theta}}(\underline{x})$ , where  $\underline{\theta}$  is the unknown parameter(s) of interest.
  - Our problem is to infer about  $\underline{\theta}$
- Let  $\Theta$  be the set of all possible values of  $\theta$ 
  - Θ is called the parameter space
- Note:
  - In the problem of statistical inference,  $\Theta$  is known, although  $\theta$  is unknown.
  - Example 1:  $X_1, X_2, \dots, X_n \sim \text{Bernoulli}(p)$ 
    - $\theta = p$  is unknown,  $\Theta = [0, 1]$  is known
  - Example 2:  $X_1, X_2, \dots, X_n \sim \text{Normal}(\mu, \sigma)$ 
    - $\underline{\theta} = [\mu, \sigma]'$  is unknown,  $\Theta = (-\infty, \infty) \times (0, \infty)$  is known
    - $\theta = \mu$  is unknown,  $\Theta = (-\infty, \infty)$  is known
    - $\theta = \sigma$  is unknown,  $\Theta = (0, \infty)$  is known

## Statistical Inference IV

- Statistical Inference
  - Estimation
    - i Point Estimation
    - ii Interval Estimation
  - 4 Hypothesis-testing



## Statistical Inference V

#### 1 Estimation:-

Here, we have **no idea** about the true value of  $\theta$  and the problem is **to estimate** the likely value of  $\theta$  on the basis of the random sample  $(X_1, X_2, \dots, X_n)$  drawn from the population

## Statistical Inference VI

#### i Point Estimation:-

Here, we estimate  $\theta$  by a **single** value (i.e., by a point)

- Let  $T = T(X_1, X_2, ..., X_n)$  be a statistic which is used to estimate the parameter  $\theta$ , is called an **estimator** of  $\theta$
- For the observed sample  $(X_1 = x_1, X_2 = x_2, ..., X_n = x_n)$ , the observed value of the estimator, namely,

$$t = T(x_1, x_2, \ldots, x_n)$$

is called an estimate of  $\theta$ 

## Statistical Inference VII

#### ii Interval Estimation:-

Here, we estimate  $\theta$  by an **interval** of values

• Let  $T_1 = T_1(X_1, X_2, \dots, X_n)$  and  $T_2 = T_2(X_1, X_2, \dots, X_n)$  be two statistics such that

$$P[T_1 \le \theta \le T_2] = 1 - \alpha,$$

where  $\alpha$  is a pre-assigned small quantity. Usually, we take  $\alpha=0.05$  or 0.01 etc.

• If  $\alpha=0.05$ , then  $P[T_1 \le \theta \le T_2]=0.95$ . Hence, the observed values of  $[T_1,T_2]=[t_1,t_2]$ , say, is called a 95% confidence interval of  $\theta$ 

## Statistical Inference VIII

### 2 Hypothesis-testing:-

Here we have some idea about the true value of  $\theta$ , in the form of a hypothesis, say,  $\theta = \theta_0$ ,

• Our problem is **to judge or test** the validity/ feasibility/ tenability of the given hypothesis  $\theta=\theta_0$  on the basis of random sample of the population

Chapter 8a: Point Estimation

### Point Estimation I

- In this case, we estimate  $\theta$  by a **single** value (i.e., by a point)
  - Let  $T = T(X_1, X_2, ..., X_n)$  be a statistic which is used to estimate the parameter  $\theta$ , is called an **estimator** of  $\theta$
  - For the observed sample  $(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$ , the observed value of the estimator, namely,

$$t = T(x_1, x_2, \ldots, x_n)$$

is called an **estimate** of  $\theta$ .

### Point Estimation II

- Methods of calculating point estimates
  - Method of moments
  - Method of maximum likelihood

## Point Estimation III

Method of moments:

It consists in equating the first few moments of the population  $(\mu'_k = E[X^k])$  with the corresponding moments of the sample

$$\left(m'_{k} = \frac{1}{n} \sum_{i}^{n} x_{i}\right), \text{ i.e.,}$$

$$\mu'_{k} = m'_{k}$$

### Point Estimation IV

- The method of moments procedure:
   Suppose there are *I* parameters to be estimated, say θ = (θ<sub>1</sub>,...,θ<sub>I</sub>)<sup>*I*</sup>.
  - Find *I* population moments,  $\mu'_k$ , for k = 1, 2, ..., I.
    - $\mu'_{k}$  will contain one or more parameters  $\theta_{1}, \ldots, \theta_{l}$ .
  - Find the corresponding *I* sample moments,  $m'_k$ , for k = 1, 2, ..., I.
    - The number of sample moments should equal the number of parameters to be estimated.
  - From the system of equations,  $\mu'_k = m'_k$ , for k = 1, 2, ..., I, solve for the parameter  $\theta = (\theta_1, ..., \theta_I)'$ ;
    - This will be a moment estimator of  $\hat{\theta}$

## Point Estimation V

### Point Estimation (method of moment)

To Estimate	Notation	Point
Mean	X	$\frac{1}{n}\sum_{i=1}^{n}X_{i}$
Proportion	ĝ	$\frac{1}{n}\sum_{i=1}^{n}I_{(X_i=1)}$
Variance	$S_n^2$	$\frac{1}{n}\sum_{i=1}^n(X_i-\bar{X})^2$

### Point Estimation VI

• Method of maximum likelihood: It consists in choosing as estimator of  $\underline{\theta}$  that statistic, which when substituted for  $\theta$ , maximizes the likelihood function

$$L = f_{\underline{X}}(\underline{x},\underline{\theta}) = f_{X_1,\ldots,X_n}(x_1,\ldots,x_n,\theta_1,\ldots,\theta_l).$$

- Procedure to find maximum likelihood estimate (mle):
  - Define the likelihood function,  $L(\theta)$ .
  - Often it is easier to take the natural logarithm (ln) of  $L(\theta)$ .
  - When applicable, differentiate  $I(\theta) = InL(\theta)$  with respect to  $\theta$ , and then equate the derivative to zero.
  - Solve for the parameter  $\theta$ , and we will obtain  $\hat{\theta}$ .
  - Check whether it is a maximizer or global maximizer.

### Point Estimation VII

### Some Desirable Properties of Point Estimators

- Unbiased
- Sufficiency
- Consistency
- Efficiency

## Unbiased Estimators I

#### **Unbiased Estimators**

• A point estimator  $\hat{\theta}$  is called an unbiased estimator of the parameter  $\theta$  if

$$E(\hat{\theta}) = \theta$$

for all possible values of  $\theta$ .

- Otherwise  $\hat{\theta}$  is said to be biased.
  - Furthermore, the bias of  $\hat{\theta}$  is given by

$$B(\hat{\theta}) = E[\hat{\theta}] - \theta.$$

# Unbiased Estimators II

#### Theorems:

- The sample mean,  $\bar{X}\left(=\frac{1}{n}\sum_{i=1}^n X_i\right)$  is an unbiased estimator of the population mean  $\mu$ .
  - Sketch of proof:  $E[\bar{X}] = E\left(\frac{1}{n}\sum_{i=1}^{n}X_i\right) = \mu$

## Unbiased Estimators III

- The statistic,  $S_n^2 \left( = \frac{1}{n} \sum_{i=1}^n (X_i \bar{X})^2 \right)$  is not an unbiased estimator of the population variance  $\sigma^2$ .
  - Sketch of proof:

$$E\left[S_{n}^{2}\right] = E\left(\frac{1}{n}\sum_{i=1}^{n}(X_{i} - \bar{X})^{2}\right) = \frac{1}{n}E\left(\sum_{i=1}^{n}\left((X_{i} - \mu) - (\bar{X} - \mu)\right)^{2}\right)$$

$$= \frac{1}{n}\left(\sum_{i=1}^{n}E(X_{i} - \mu)^{2} - nE(\bar{X} - \mu)^{2}\right) = \frac{1}{n}\left(n\sigma^{2} - n\frac{\sigma^{2}}{n}\right) = \left(\frac{n-1}{n}\right)\sigma^{2}$$

- The sample variance,  $S^2\left(=\frac{1}{n-1}\sum_{i=1}^n(X_i-\bar{X})^2\right)$  is an unbiased estimator of the population variance  $\sigma^2$ .
  - Sketch of proof:  $E\left[S^2\right] = E\left[\frac{n}{n-1}S_n^2\right] = \frac{n}{n-1}E\left[S_n^2\right] = \sigma^2$

## **Unbiased Estimators IV**

### Mean Square Error of an Estimator

• The mean square error of the estimator  $\hat{\theta}$ , denoted by  $MSE(\hat{\theta})$ , is defined as

$$MSE(\hat{\theta}) = E(\hat{\theta} - \theta)^2$$
.

- Theorem:  $MSE(\hat{\theta}) = Var(\hat{\theta}) + B^2(\hat{\theta})$ .
  - Sketch of proof:

$$\begin{split} \mathit{MSE}(\hat{\theta}) &= E(\hat{\theta} - \theta)^2 = E\left((\hat{\theta} - E(\hat{\theta})) + (E(\hat{\theta}) - \theta)\right)^2 \\ &= E\left(\hat{\theta} - E(\hat{\theta})^2 + \left(E(\hat{\theta}) - \theta\right)^2 + 2\left(E(\hat{\theta}) - \theta\right)\left(E(\hat{\theta}) - E(\hat{\theta})\right) \\ &= \mathit{Var}(\hat{\theta}) + B^2(\hat{\theta}). \end{split}$$

## Unbiased Estimators V

• If  $\hat{\theta}$  is an unbiased estimator of  $\theta$ , then

$$B^2(\hat{\theta}) = 0$$
 and  $MSE(\hat{\theta}) = Var(\hat{\theta})$ .

- Minimum variance unbiased estimator (MVUE) of  $\theta$  :
  - The unbiased estimator  $\hat{\theta}$  that minimizes the mean square error is called the *MVUE* of  $\theta$ .

## Sufficient Estimators I

#### Sufficient Estimators

- A statistic U is a sufficient statistic for a parameter  $\theta$  if U contains all the information available in the data about the value of  $\theta$ .
- If U(X) is a sufficient statistic for a parameter θ, then any inference about θ should depend on the sample X only through the value U(X).
  - That is, if  $\mathbf{x}$  and  $\mathbf{y}$  are two sample points such that  $U(\mathbf{x}) = U(\mathbf{y})$ , then the inference about  $\theta$  should be the same whether  $\mathbf{X} = \mathbf{x}$  or  $\mathbf{X} = \mathbf{y}$  is observed.
- Example
  - the sample mean  $(\bar{X})$  may contain all the relevant information about the parameter  $\mu$ , and in that case  $U(\mathbf{X}) = \bar{X}$  is called a sufficient statistic for  $\mu$ .

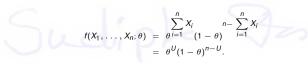
### Sufficient Estimators II

#### Formal definition

- Let  $X_1, ..., X_n$  be a random sample from a probability distribution with unknown parameter  $\theta$ .
  - Then, the statistic  $U = g(X_1, \ldots, X_n)$  is said to be **sufficient statistic** for  $\theta$  if  $f_{X_1, \ldots, X_n}(x_1, \ldots, x_n | U = u)$  does not depend on  $\theta$  for any value of u.
  - An estimator of  $\theta$  that is a function of a sufficient statistic for  $\theta$  is said to be a **sufficient estimator** of  $\theta$ .

# Sufficient Estimators III

- Example: Let  $X_1, \ldots, X_n$  be i.i.d. Bernoulli random variables with parameter  $\theta$ . Then  $U = \sum_{i=1}^n X_i$  is sufficient for  $\theta$ .
- Sketch of proof:



Since,  $U \sim Bin(n, \theta)$ ,

$$f(U;\theta) = {}^{n}C_{U}\theta^{U}(1-\theta)^{n-U}.$$

Thus,

$$f(x_1,\ldots,x_n|U=u)=\frac{f(x_1,\ldots,x_n,u)}{f_U(u)}=\left\{\begin{array}{cc}\frac{1}{n_{Cu}}&\text{; if }u=\sum x_i\\0&\text{; otherwise.}\end{array}\right.$$

## Sufficient Estimators IV

Neyman-Fisher factorization theorem to spot a sufficient statistic.

• Theorem: -Let U be a statistic based on the random sample  $X_1, \ldots, X_n$ . Then, U is a sufficient statistic for  $\theta$  if and only if the joint p.d.f or p.m.f.,  $f(x_1, \ldots, x_n; \theta)$  can be factored into two non-negative functions, i.e.,

$$f(x_1,\ldots,x_n;\theta)=g(u,\theta)h(x_1,\ldots,x_n), \text{ for all } x_1,\ldots,x_n,$$

#### where

- $g(u, \theta)$  is a function only of u and  $\theta$
- and  $h(x_1, ..., x_n)$  is a function of only  $x_1, ..., x_n$  and not of  $\theta$ .

## Sufficient Estimators V

- Sketch of proof: discrete case
  - Sufficient => Factorization

$$f(x_1,\ldots,x_n;\theta) = P_{\theta}(X_1 = x_1,\ldots,X_n = x_n, U = u)$$

$$= P_{\theta}(X_1 = x_1,\ldots,X_n = x_n|U = u)P_{\theta}(U = u)$$

$$\stackrel{\text{suff.}}{=} h(x_1,\ldots,x_n)g(u,\theta).$$

Factorization => Sufficient

$$\begin{split} P_{\theta}(X_1 = x_1, \dots, X_n = x_n | U = u) &= \frac{P_{\theta}(X_1 = x_1, \dots, X_n = x_n, U = u)}{P_{\theta}(U = u)} \\ &= \begin{cases} \frac{P_{\theta}(X_1 = x_1, \dots, X_n = x_n, U = u)}{P_{\theta}(U = u)} & \text{if } (x_1, \dots, x_n) \in A_u \\ 0 & \text{if } (x_1, \dots, x_n) \notin A_u, \end{cases} \end{split}$$

where  $A_U$  is the set of all  $(x_1, \ldots, x_n)$  such that U maps it into u.

## Sufficient Estimators VI

• When  $(x_1, x_2, \ldots, x_n) \in A_u$ ,

$$P_{\theta}(X_{1} = x_{1}, \dots, X_{n} = x_{n} | U = u) = \frac{P_{\theta}(X_{1} = x_{1}, \dots, X_{n} = x_{n}, U = u)}{P_{\theta}(U = u)}$$

$$= \frac{P_{\theta}(X_{1} = x_{1}, \dots, X_{n} = x_{n})}{P_{\theta}(U = u)}$$

$$= \frac{f(x_{1}, \dots, x_{n}, \theta)}{\sum_{(x_{1}, \dots, x_{n}) \in A_{u}} f(x_{1}, \dots, x_{n}, \theta)}$$

face = =

$$\stackrel{\text{fact}}{=} \frac{g(u,\theta)h(x_1,\ldots,x_n)}{\sum_{(x_1,\ldots,x_n)\in A_U} g(u,\theta)h(x_1,\ldots,x_n)}$$

$$= \frac{h(x_1,\ldots,x_n)}{\sum_{(x_1,\ldots,x_n)\in A_U} h(x_1,\ldots,x_n)} \perp \theta$$

• When  $(x_1, x_2, \ldots, x_n) \notin A_u$ ,

$$P_{\theta}(X_1 = x_1, \dots, X_n = x_n | U = u) = \frac{P_{\theta}(X_1 = x_1, \dots, X_n = x_n, U = u)}{P_{\theta}(U = u)}$$
  
=  $0 \perp \theta$ 

## Sufficient Estimators VII

- Procedure to verify Sufficiency
  - **①** Obtain the joint pdf or pmf  $f_{\theta}(x_1, \ldots, x_n)$ .
  - 2 If necessary, rewrite the joint pdf or pmf in terms of the given statistic and parameter so that one can use the factorization theorem.
  - Define the functions g and h, in such a way that g is a function of the statistic and parameter only and h is a function of the observations only.
  - If step 3 is possible, then the statistic is sufficient. Otherwise, it is not sufficient.

### Sufficient Estimators VIII

- Example: Let  $X_1, \ldots, X_n$  denote a random sample from a geometric population with parameter p. Show that  $\bar{X}$  is sufficient for p.
  - Sketch of proof:

$$f(x_1, \dots, x_n) = \prod_{i=1}^n p(1-p)^{x_i-1}$$

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

$$= p^n(1-p)^{-n+n\bar{x}}$$

$$= g(\bar{x}, p) h(x_1, \dots, x_n)$$

### Sufficient Estimators IX

• Joint Sufficiency: Two statistics  $U_1$  and  $U_2$  are said to be jointly sufficient for the parameters  $\theta_1$  and  $\theta_2$  if the conditional distribution of  $X_1, \ldots, X_n$  given  $U_1$  and  $U_2$  does not depend on  $\theta_1$  or  $\theta_2$ . In general, the statistic  $U = (U_1, \ldots, U_n)$  is jointly sufficient for  $\theta = (\theta_1, \ldots, \theta_n)$  if the conditional distribution of  $X_1, \ldots, X_n$  given U is free of  $\theta$ .

# Sufficient Estimators X

 Factorization criteria for Joint Sufficiency Theorem: -

The two statistics  $U_1$  and  $U_2$  are jointly sufficient for  $\theta_1$  and  $\theta_2$  if and only if the likelihood function can be factored into two non-negative functions,

$$f(x_1,...,x_n;\theta_1,\theta_2) = g(u_1,u_2;\theta_1,\theta_2)h(x_1,...,x_n)$$

where  $g(u_1, u_2; \theta_1, \theta_2)$  is only a function of  $u_1, u_2; \theta_1$  and  $\theta_2$ , and  $h(x_1, \dots, x_n)$  is free of  $\theta_1$  or  $\theta_2$ 

# Sufficient Estimators XI

- Example: Let  $X_1, \ldots, X_n$  be a random sample from  $N(\mu, \sigma^2)$ .
  - If  $\mu$  is unknown and  $\sigma^2 = \sigma_0^2$  is known, then  $U_1 = \bar{X}$  is a sufficient statistic for  $\mu$ .

$$L = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma_0}\right)^2} = \underbrace{\left(2\pi\sigma_0^2\right)^{-n/2} e^{-\frac{1}{2\sigma_0^2} \left[\sum_{i=1}^{n} x_i^2\right]}}_{h(x_1, \dots, x_n)} \times \underbrace{e^{-\frac{1}{2\sigma_0^2} \left[-2n\mu\bar{x} + n\mu^2\right]}}_{g(u_1, \theta_1)}$$

## Sufficient Estimators XII

② If  $\mu = \mu_0$  is known and  $\sigma^2$  is unknown, then  $U_1 = \sum_{i=1}^n (X_i - \mu_0)^2$  is a sufficient statistic for  $\sigma^2$ .

$$L = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left(\frac{x_i - \mu_0}{\sigma}\right)^2} = \underbrace{(2\pi)^{-n/2}}_{h(x_1, \dots, x_n)} \times \underbrace{\sigma^{-n} e^{-\sum_{i=1}^{n} (x_i - \mu_0)^2}}_{g(u_1, \theta_1)}$$

## Sufficient Estimators XIII

If  $\mu$  and  $\sigma^2$  are both unknown, then  $U_1 = \sum_{i=1}^n X_i$  and  $U_2 = \sum_{i=1}^n X_i^2$  are jointly sufficient for  $\mu$  and  $\sigma^2$ .

$$L = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left(\frac{x_{i} - \mu}{\sigma}\right)^{2}} = \underbrace{(2\pi)^{-n/2}}_{h(x_{1}, \dots, x_{n})} \times \underbrace{\sigma^{-n} e^{-\frac{1}{2\sigma^{2}} \left[\sum_{i=1}^{n} x_{i}^{2} - 2\mu \sum_{i=1}^{n} x_{i} + n\mu^{2}\right]}}_{g(u_{1}, \theta_{1})}$$

### Sufficient Estimators XIV

- Theorem: -
  - If U is a sufficient statistic for  $\theta$ , then the maximum likelihood estimator of  $\theta$ , if unique, is a function of U.
    - Sketch of proof.

$$f(x_1,...,x_n;\theta)\stackrel{suff}{=} g(u,\theta)h(x_1,...,x_n).$$

Thus, the joint pdf/pmf depends on  $\theta$  only through the statistic U. To maximize L we need to maximize  $g(U, \theta)$ .

## Sufficient Estimators XV

• Theorem: -

Let  $X_1, \ldots, X_n$  be a random sample from a population with pdf or pmf of the exponential form

$$f(x_1,\ldots,x_n;\theta) = \left\{ \begin{array}{cc} exp[k(x)c(\theta) + S(x) + d(\theta)] & , x \in B \\ 0 & , x \notin B, \end{array} \right.$$

where B does not depend on the parameter  $\theta$ . The statistic

$$U = \sum_{i=1}^{n} k(X_i)$$
 is sufficient for  $\theta$ .

Sketch of proof.

$$f(x_1,...,x_n;\theta) = exp\left[c(\theta)\sum_{i=1}^n k(x_i) + \sum_{i=1}^n S(x_i) + nd(\theta)\right]$$

$$= exp\left[c(\theta)\sum_{i=1}^n k(x_i) + nd(\theta)\right] exp\left[\sum_{i=1}^n S(x_i)\right].$$

### Sufficient Estimators XVI

#### Some Observations on Sufficiency

- All the statistics need not be sufficient, in addition
- Any function of a sufficient statistic needs not to be sufficient, however
- Any one-to-one function of a sufficient statistic is also sufficient statistic

### Sufficient Estimators XVII

- RAO-BLACKWELL Theorem: -Let  $X_1, \ldots, X_n$  be a random sample with joint pmf or pdf  $f(x_1, \ldots, x_n; \theta)$  and let  $U = (U_1, \ldots, U_n)$  be jointly sufficient for  $\theta = (\theta_1, \ldots, \theta_n)$ . If T is any unbiased estimator of  $k(\theta)$ , and if  $T^* = E(T|U)$ , then:
  - $T^*$  is an unbiased estimator of  $k(\theta)$ . Sketch of proof:  $-ET^* = E(E(T|U)) = E(T) = k(\theta)$ . Hence,  $T^*$  is an unbiased estimator of  $k(\theta)$ .
  - T\* is a function of U, and does not depend on θ.
     Sketch of proof: Because U is sufficient for θ, the conditional distribution of any statistic (hence, for T), given

U, does not depend on  $\theta$ .

•  $Var(T^*) \leq Var(T)$  for every  $\theta$ , and  $Var(T^*) < Var(T)$  for some  $\theta$  unless  $T^* = T$  with probability 1. Sketch of proof:  $-Var(T) = E(Var(T|U)) + Var(E(T|U)) = E(Var(T|U)) + Var(T^*)$ . Because Var(T|U) > 0 for all U, it follows that E(Var(T|U)) > 0. Hence,  $Var(T^*) < Var(T)$ .

Also  $Var(T^*) = Var(T)$  iff Var(T|U) = 0 or T is a function of U, in which case  $T^* = E(T|U) = T$ .

### Sufficient Estimators XVIII

#### More Observations

- If one is searching for an unbiased estimator with minimal variance, it has to be restricted to functions of a sufficient statistics.
- If  $k(\theta) = \theta$ , and T is an unbiased estimator of  $\theta$ , then  $T^* = E(T|U)$  will typically give the MVUE of  $\theta$ .

### Sufficient Estimators XIX

 Minimal sufficient statistic:
 A sufficient statistic T(X) is called a minimal sufficient statistic if for any other sufficient statistic T'(X),

$$T(X) = g(T'(X)),$$

i.e., T(X) is a function of T'(X).

• Intuitively, a minimal sufficient statistic most efficiently captures all possible information about the parameter  $\theta$ .

## Sufficient Estimators XX

- Lehmann and Scheffe method to find a minimal sufficient statistic
  - Let  $X_1, ..., X_n$  be a random sample with pdf or pmf f(x) that depends on a parameter  $\theta$ .
  - Let  $(x_1, ..., x_n)$  and  $(y_1, ..., y_n)$  be two different sets of values of  $(X_1, ..., X_n)$ .
    - Let  $\frac{L(\theta;x_1,\ldots,x_n)}{L(\theta;y_1,\ldots,y_n)}$  be the ratio of the likelihoods evaluated at these two points.
  - Suppose it is possible to find a function  $g(x_1, ..., x_n)$  such that this ratio will be free of the unknown parameter  $\theta$  if and only if  $g(x_1, ..., x_n) = g(y_1, ..., y_n)$ ,
    - in other words

$$\frac{L(\theta;x_1,\ldots,x_n)}{L(\theta;y_1,\ldots,y_n)} \text{ is independent of } \theta \Leftrightarrow g(x_1,\ldots,x_n) = g(y_1,\ldots,y_n).$$

• If such a function g can be found, then  $g(X_1, \ldots, X_n)$  is a minimal sufficient statistic for  $\theta$ .

## Sufficient Estimators XXI

- Example: Let (X<sub>1</sub>,..., X<sub>n</sub>) be a random sample from the Bernoulli(p), where
   p is unknown. Find a minimal sufficient statistic for p.
  - Solution: The ratio of the likelihoods is

$$\frac{L(x_1,\ldots,x_n)}{L(y_1,\ldots,y_n)} = \frac{\rho(x_1,\ldots,x_n)}{\rho(y_1,\ldots,y_n)} = \frac{\rho^{\sum x_i}(1-\rho)^{n-\sum x_i}}{\rho^{\sum y_i}(1-\rho)^{n-\sum y_i}} = \left(\frac{\rho}{1-\rho}\right)^{\sum x_i-\sum y_i}$$

This ratio is to be independent of p, if and only if

$$g(x_1,\ldots,x_n) = \sum_{i=1}^n x_i = \sum_{i=1}^n y_i = g(y_1,\ldots,y_n).$$

Therefore,  $g(X_1, \dots, X_n) = \sum_{i=1}^n X_i$  is a minimal sufficient statistic for p.

### Consistent Estimators I

#### Consistent Estimators

 A statistic is a consistent estimator if its value becomes closer to the value of the true parameter which is being estimated, as the sample size becomes larger.

## Consistent Estimators II

#### Formal definition

• The estimator  $\hat{\theta}_n$  is said to be a consistent estimator of  $\theta$  if, for any  $\epsilon > 0$ .

$$\lim_{n\to\infty} P\left[|\hat{\theta}_n - \theta| \le \epsilon\right] = 1$$

or equivalently,

$$\lim_{n\to\infty} P\left[|\hat{\theta}_n - \theta| > \epsilon\right] = 0,$$

i.e.  $\hat{\theta}_n$  converges in probability to  $\theta$ .

### Consistent Estimators III

- Example: Let  $X_1, \ldots, X_n$  be a random sample with true mean  $\mu$  and finite variance,  $\sigma^2$ . Then, the sample mean  $\bar{X}$  is a consistent estimator of the population mean  $\mu$ .
  - Sketch of proof: -Note that, for any positive r.v. X (Markov Inequality)

$$E[X] = \int_0^\epsilon x f_X(x) dx + \int_\epsilon^\infty x f_X(x) dx \ge \int_0^\epsilon x f_X(x) dx + \epsilon \int_\epsilon^\infty f_X(x) dx \ge \epsilon P(X \ge \epsilon).$$

Hence, (Chebyshev's Inequality)

$$P\left[|\bar{X} - \mu| \geq \epsilon\right] = P\left[(\bar{X} - E[\bar{X}])^2 \geq \epsilon^2\right] \leq \frac{E(\bar{X} - E[\bar{X}])^2}{\epsilon^2} = \frac{Var(\bar{X})}{\epsilon^2},$$

Thus,

$$\lim_{n\to\infty} P\left[|\bar{X}-\mu| \geq \epsilon\right] \leq \lim_{n\to\infty} \frac{\mathit{Var}(\bar{X})}{\epsilon^2} = \lim_{n\to\infty} \frac{\sigma^2}{n\epsilon^2} = 0$$

#### Consistent Estimators IV

#### Test for Consistency (Sufficient conditions)

• Theorem: - An unbiased estimator  $\hat{\theta}_n$  of  $\theta$  is a consistent estimator for  $\theta$  if

$$\lim_{n\to\infty} Var(\hat{\theta}_n) = 0.$$

- $\bullet \quad \text{Sketch of proof:} \lim_{n \to \infty} P\left[ |\hat{\theta}_n \theta| \geq \epsilon \right] \stackrel{U.E.}{=} \lim_{n \to \infty} P\left[ |\hat{\theta}_n E(\hat{\theta}_n)| \geq \epsilon \right] \stackrel{C.I.}{\leq} \lim_{n \to \infty} \frac{\text{Var}(\hat{\theta}_n)}{\epsilon^2} = 0$
- An unbiased estimator  $\hat{\theta}_n$  is consistent if  $Var(\hat{\theta}_n) \to 0$  as  $n \to \infty$ .
- Theorem: An estimator  $\hat{\theta}_n$ , with finite variance, is a consistent estimator for  $\theta$  if

$$\lim_{n\to\infty} E(\hat{\theta}_n - \theta)^2 = 0.$$

- $\bullet \quad \text{Sketch of proof: } -\lim_{\substack{n \to \infty}} P\left[|\hat{\theta}_n \theta| \geq \epsilon\right] \overset{C.I.}{\leq} \lim_{\substack{n \to \infty}} \frac{E(\hat{\theta}_n \theta)^2}{\epsilon^2} = 0$
- A biased estimator  $\hat{\theta}_n$  is consistent if both  $Var(\hat{\theta}_n) \to 0$  and  $B(\hat{\theta}_n) \to 0$  as  $n \to \infty$ .



## Consistent Estimators V

- Example: Let  $X_1, ..., X_n$  be a random sample from  $N(\mu, \sigma^2)$  population.
  - Then the sample variance  $S^2$  is a consistent estimator for  $\sigma^2$ .
    - Sketch of proof: -

$$(n-1)\frac{S^2}{\sigma^2} \sim \chi^2(n-1) \Rightarrow Var\left[(n-1)\frac{S^2}{\sigma^2}\right] = 2(n-1) \Rightarrow Var(S^2) = \frac{2\sigma^4}{n-1}.$$
  
 $S^2$  is u.e. of  $\sigma^2$  and  $\lim_{n\to\infty} Var(S^2) = 0.$ 

- The maximum likelihood estimators  $\bar{X}$  and  $S_n^2$  for  $\mu$  and  $\sigma^2$ , respectively, are consistent estimators for  $\mu$  and  $\sigma^2$ .
  - Sketch of proof: -

$$\bar{X}$$
 is u.e. of  $\mu$  and  $\lim_{n\to\infty} Var(\bar{X})=0$ , thus  $\bar{X}$  is consistent estimator of  $\mu$ .

$$\begin{aligned} & \textit{Var}(S_n^2) = \textit{Var}\left(\frac{n-1}{n}S^2\right) = \frac{2(n-1)\sigma^4}{n^2}, \text{ thus } \lim_{n \to \infty} \textit{Var}(S_n^2) = 0 \text{ and } \\ & \textit{B}(S_n^2) = \frac{n-1}{n}\sigma^2 - \sigma^2 = -\frac{1}{n}\sigma^2, \text{ thus } \lim_{n \to \infty} \textit{B}(S_n^2) = 0 \end{aligned}$$

## Efficiency I

#### Efficiency

 It is a relative comparison between variances of two unbiased/biased estimators

# Efficiency II

#### Formal Definitions

• If  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are two unbiased estimators for  $\theta$ , the efficiency of  $\hat{\theta}_1$  relative to  $\hat{\theta}_2$  is the ratio

$$e(\hat{\theta}_1, \hat{\theta}_2) = \frac{Var(\hat{\theta}_2)}{Var(\hat{\theta}_1)}.$$

- If  $Var(\hat{\theta}_1) < Var(\hat{\theta}_2)$ , or equivalently,  $e(\hat{\theta}_1, \hat{\theta}_2) > 1$ , then,  $\hat{\theta}_1$  is relatively more efficient than  $\hat{\theta}_2$ .
- If  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are two biased estimators for  $\theta$ , the efficiency of  $\hat{\theta}_1$  relative to  $\hat{\theta}_2$  is the ratio

$$e(\hat{\theta}_1, \hat{\theta}_2) = \frac{E(\hat{\theta}_2 - \theta)^2}{E(\hat{\theta}_1 - \theta)^2} = \frac{MSE(\hat{\theta}_2)}{MSE(\hat{\theta}_1)}.$$

• If  $MSE(\hat{\theta}_1) < MSE(\hat{\theta}_2)$ , or equivalently,  $e(\hat{\theta}_1, \hat{\theta}_2) > 1$ , then,  $\hat{\theta}_1$  is relatively more efficient than  $\hat{\theta}_2$ .

52 / 64

## Efficiency III

- Example:- Let  $X_1, \ldots, X_n, n \ge 2$  be a random sample from a normal population with a true mean  $\mu$  and variance  $\sigma^2$ . Consider the following two estimators of  $\sigma^2$ :  $\theta_1 = S^2$ , and  $\theta_2 = S^2_n$ . Find  $e(\theta_1, \theta_2)$ .
  - Sketch of proof: -

$$\begin{split} \mathit{MSE}(\hat{\theta}_1) &= 0 + \frac{2\sigma^4}{n-1} = \frac{2\sigma^4}{n-1} \\ \mathit{MSE}(\hat{\theta}_2) &= \left( -\frac{\sigma^2}{n} \right)^2 + \frac{2(n-1)\sigma^4}{n^2} = \frac{(2n-1)\sigma^4}{n^2} \\ e(\hat{\theta}_1, \hat{\theta}_2) &= \frac{(n-1)(2n-1)}{2n^2} < 1, \text{ for } n \geq 2. \end{split}$$

Hence,  $S_n^2$  is relatively more efficient than  $S^2$ .

# Efficiency IV

#### Uniformly Minimum Variance Unbiased Estimator

• An unbiased estimator  $\hat{\theta}_0$ , is said to be a uniformly minimum variance unbiased estimator (UMVUE) for the parameter  $\theta$  if, for any other unbiased estimator  $\hat{\theta}$ ,

$$Var(\hat{\theta}_0) \leq Var(\hat{\theta}),$$

for all possible values of  $\theta$ .

# Efficiency V

## Cramer-Rao Inequality:

(Lower bound for the variance of any estimator)

• Theorem: - Let  $X_1, \ldots, X_n$  be a sample with pdf  $f(\mathbf{x}|\theta)$ , and let  $W(\mathbf{X}) = W(X_1, \ldots, X_n)$  be any estimator satisfying

$$\frac{d}{d\theta} E_{\theta}[W(\mathbf{X})] = \int_{\mathcal{X}} \frac{\delta}{\delta \theta} [W(\mathbf{x}) f(\mathbf{x}|\theta)] d\mathbf{x}$$

and

$$Var_{\theta}[W(\mathbf{X})] < \infty$$
,

then

$$Var_{ heta}[W(\mathbf{X})] \geq rac{\left(rac{d}{d heta}E_{ heta}[W(\mathbf{X})]
ight)^2}{E_{ heta}\left[\left(rac{\delta}{\delta heta}\ln f(\mathbf{x}| heta)
ight)^2
ight]}.$$

# Efficiency VI

Sketch of proof: -Note that

$$\begin{split} \frac{d}{d\theta} E_{\theta}[W(\mathbf{X})] &= \int_{\mathcal{X}} \frac{\delta}{\delta \theta} [W(\mathbf{x}) f(\mathbf{x} | \theta)] d\mathbf{x} \\ &= \int_{\mathcal{X}} W(\mathbf{x}) \left[ \frac{\delta}{\delta \theta} f(\mathbf{x} | \theta) \right] d\mathbf{x} \\ &= \int_{\mathcal{X}} W(\mathbf{x}) \left[ \frac{\delta f(\mathbf{x} | \theta)}{\delta \theta} \frac{1}{f(\mathbf{x} | \theta)} \right] f(\mathbf{x} | \theta) d\mathbf{x} \\ &= E_{\theta} \left[ W(\mathbf{X}) \frac{\delta}{\delta \theta} \ln f(\mathbf{x} | \theta) \right] \\ &= Cov_{\theta} \left( W(\mathbf{X}), \frac{\delta}{\delta \theta} \ln f(\mathbf{x} | \theta) \right), \end{split}$$

since 
$$E_{\theta} \left[ \frac{\delta}{\delta \theta} \ln f(\mathbf{x}|\theta) \right] = \int_{\mathcal{X}} \left[ \frac{\delta}{\delta \theta} \ln f(\mathbf{x}|\theta) \right] f(\mathbf{x}|\theta) d\mathbf{x} = \int_{\mathcal{X}} \left[ \frac{\delta}{\delta \theta} f(\mathbf{x}|\theta) \right] d\mathbf{x} = \frac{d}{d\theta} \left[ \int_{\mathcal{X}} f(\mathbf{x}|\theta) d\mathbf{x} \right] = 0.$$

Also

$$Var_{\theta}\left(\frac{\delta}{\delta\theta}\ln f(\mathbf{x}|\theta)\right) = E_{\theta}\left[\left(\frac{\delta}{\delta\theta}\ln f(\mathbf{x}|\theta)\right)^{2}\right].$$

Now,

$$\frac{\left[\textit{Cov}_{\theta}\left(\textit{W}(\mathbf{X}), \frac{\delta}{\delta\theta} \ln f(\mathbf{x}|\theta)\right)\right]^{2}}{\textit{Var}_{\theta}\left[\textit{W}(\mathbf{X})\right] \textit{Var}_{\theta}\left[\frac{\delta}{\delta\theta} \ln f(\mathbf{x}|\theta)\right]} \leq 1 \Rightarrow \frac{\left[\frac{d}{d\theta} \textit{E}_{\theta}[\textit{W}(\mathbf{X})]\right]^{2}}{\textit{E}_{\theta}\left[\left(\frac{\delta}{\delta\theta} \ln f(\mathbf{x}|\theta)\right)^{2}\right]} \leq \textit{Var}_{\theta}\left[\textit{W}(\mathbf{X})\right]$$

# Efficiency VII

Corollary: -

Let  $X_1, \ldots, X_n$  be an iid random sample from a population with pdf or pmf  $f_{\theta}(x)$  that depends on a parameter  $\theta$ . If  $\hat{\theta} = W(\mathbf{x})$  is an unbiased estimator of  $\psi(\theta)$ , then

$$extstyle extstyle ext$$

## Efficiency VIII

Sketch of proof: - Note that

$$\begin{split} E\left[\left(\frac{\delta}{\delta\theta}\ln f(\mathbf{x}|\theta)\right)^2\right] &= E\left[\left(\frac{\delta}{\delta\theta}\ln\prod_{i=1}^n f_\theta(x_i)\right)^2\right] = E\left[\left(\sum_{i=1}^n \frac{\delta}{\delta\theta}\ln f_\theta(x_i)\right)^2\right] \\ &= \sum_{i=1}^n E\left[\left(\frac{\delta}{\delta\theta}\ln f_\theta(x_i)\right)^2\right] + \sum_{i\neq j} E\left[\left(\frac{\delta}{\delta\theta}\ln f_\theta(x_i)\right)\left(\frac{\delta}{\delta\theta}\ln f_\theta(x_j)\right)\right] \\ &= nE\left[\left(\frac{\delta}{\delta\theta}\ln f_\theta(x)\right)^2\right] + \sum_{i\neq j} E\left[\left(\frac{\delta}{\delta\theta}\ln f_\theta(x_i)\right)\right] E\left[\left(\frac{\delta}{\delta\theta}\ln f_\theta(x_j)\right)\right] \\ &= nE\left[\left(\frac{\delta}{\delta\theta}\ln f_\theta(x)\right)^2\right] \end{split}$$

and

$$\left[\frac{d}{d\theta}E_{\theta}[\textit{W}(\mathbf{x})]\right]^2 = \left[\frac{d}{d\theta}E_{\theta}[\hat{\theta}]\right]^2 = \left[\frac{d}{d\theta}\psi(\theta)\right]^2.$$

Hence,

$$\textit{Var}(\hat{\theta}) \geq \frac{\left[\frac{d}{d\theta} E_{\theta}[\hat{\theta}]\right]^{2}}{E_{\theta}\left[\left(\frac{\delta}{\delta \theta} \ln f(\mathbf{x}|\theta)\right)^{2}\right]} = \frac{\left[\frac{d}{d\theta} \psi(\theta)\right]^{2}}{nE\left[\left(\frac{\delta}{\delta \theta} \ln f_{\theta}(\mathbf{x})\right)^{2}\right]}.$$

# Efficiency IX

Corollary: -

Let  $X_1, \ldots, X_n$  be an iid random sample from a population with pdf or pmf  $f_{\theta}(x)$  that depends on a parameter  $\theta$ . If  $\hat{\theta}$  is an unbiased estimator of  $\theta$ , then

$$Var(\hat{ heta}) \geq rac{1}{nE\left[\left(rac{\delta}{\delta heta} \ln f_{ heta}(x)
ight)^{2}
ight]}.$$

# Efficiency X

#### Efficient Estimator

• If  $\hat{\theta}$  is an unbiased estimator of  $\theta$  and if

$$Var(\hat{ heta}) = rac{1}{nE\left[\left(rac{\delta}{\delta heta}\ln f_{ heta}(x)
ight)^{2}
ight]}$$

then  $\hat{\theta}$  is a uniformly minimum variance unbiased estimator (UMVUE) of  $\theta$ .

Sometimes  $\hat{\theta}$  is also referred to as an efficient estimator.

# Efficiency XI

• Result:- If the function  $f(\cdot)$  is sufficiently smooth, specifically if  $\frac{d}{d\theta} E_{\theta} \left( \frac{\delta}{\delta \theta} \ln f_{\theta}(x) \right) = \int \frac{\delta}{\delta \theta} \left[ \left( \frac{\delta}{\delta \theta} \ln f_{\theta}(x) \right) f_{\theta}(x) \right] dx$ , then

$$E_{\theta}\left[\left(\frac{\delta}{\delta\theta}\ln f_{\theta}(x)\right)^{2}\right] = -E_{\theta}\left(\frac{\delta^{2}}{\delta\theta^{2}}\ln f_{\theta}(x)\right) = Var\left[\frac{\delta}{\delta\theta}\ln f_{\theta}(x)\right]$$

and for an unbiased estimator  $\hat{\theta}$  for  $\theta$  the Cramer-Rao inequality can be rewritten as

$$Var(\hat{\theta}) \geq \frac{1}{-nE\left(rac{\delta^2}{\delta\theta^2}\ln f_{\theta}(x)
ight)} = \frac{1}{nVar\left[rac{\delta}{\delta\theta}\ln f_{\theta}(x)
ight]}.$$

## Efficiency XII

Sketch of proof: -

$$\begin{split} E_{\theta} \left( \frac{\delta^2}{\delta \theta^2} \ln f_{\theta}(x) \right) &= E_{\theta} \left[ \frac{\delta}{\delta \theta} \left( \frac{\delta}{\delta \theta} \ln f_{\theta}(x) \right) \right] = E_{\theta} \left[ \frac{\delta}{\delta \theta} \left( \frac{\frac{\delta}{\delta \theta} f_{\theta}(x)}{f_{\theta}(x)} \right) \right] \\ &= E_{\theta} \left[ \left( \frac{\frac{\delta^2}{\delta \theta^2} f_{\theta}(x)}{f_{\theta}(x)} \right) - \left( \frac{\frac{\delta}{\delta \theta} f_{\theta}(x)}{f_{\theta}(x)} \right)^2 \right] \\ &= \int \frac{\delta^2}{\delta \theta^2} f_{\theta}(x) dx - E_{\theta} \left[ \left( \frac{\delta}{\delta \theta} \ln f_{\theta}(x) \right)^2 \right] \\ &= \frac{\delta}{\delta \theta} \int \frac{\delta}{\delta \theta} f_{\theta}(x) dx - E_{\theta} \left[ \left( \frac{\delta}{\delta \theta} \ln f_{\theta}(x) \right)^2 \right] \\ &= \frac{d}{d \theta} E_{\theta} \left[ \frac{\delta}{\delta \theta} \ln f_{\theta}(x) \right] - E_{\theta} \left[ \left( \frac{\delta}{\delta \theta} \ln f_{\theta}(x) \right)^2 \right] \\ &= -E_{\theta} \left[ \left( \frac{\delta}{\delta \theta} \ln f_{\theta}(x) \right)^2 \right] \end{split}$$

Hence,

$$\mathit{Var}(\hat{\theta}) \geq \frac{1}{n \mathsf{E}_{\theta} \left[ \left( \frac{\delta}{\delta \theta} \ln f_{\theta}(x) \right)^2 \right]} = \frac{1}{-n \mathsf{E}_{\theta} \left( \frac{\delta^2}{\delta \theta^2} \ln f_{\theta}(x) \right)} = \frac{1}{n \mathit{Var}_{\theta} \left[ \frac{\delta}{\delta \theta} \ln f_{\theta}(x) \right]}.$$

# Efficiency XIII

- Example: Let  $X_1, \ldots, X_n$  be a random sample from an  $N(\mu, \sigma^2)$  population. Then  $\hat{X}$  is an efficient estimator for  $\mu$ .
  - Sketch of proof:  $I(x, \mu) = \ln f(x, \mu) = c \frac{(x-\mu)^2}{2\sigma^2}$ . Thus,

$$\frac{\delta}{\delta\mu}I(x,\mu) = \frac{x-\mu}{\sigma^2} \text{ and } \frac{\delta^2}{\delta\mu^2}I(x,\mu) = -\frac{1}{\sigma^2}.$$

Hence

$$\frac{1}{n E\left[\left(\frac{\delta}{\delta \theta} \ln f_{\theta}(x)\right)^{2}\right]} = \frac{1}{n E\left[\left(\frac{\delta}{\delta \mu} l(x, \mu)\right)^{2}\right]} = \frac{1}{n E\left[-\frac{\delta^{2}}{\delta^{2} \mu} l(x, \mu)\right]} = \frac{\sigma^{2}}{n} = Var(\bar{X})$$

# Efficiency XIV

#### Note: -

- For a given problem UMVUE may not exist.
- Even when an UMVUE exists, it is not necessary that it have a variance equal to the Cramer—Rao lower bound.
- The term  $I(\theta) = E\left[\left(\frac{\delta}{\delta \theta} \ln f_{\theta}(x)\right)^{2}\right]$  is called the Fisher information.
- It can be shown that the Fisher information in a sample of size n, denoted by  $I_n(\theta)$ , is n times the Fisher information in one observation. That is.

$$I_n(\theta) = nI(\theta).$$