# MATH 450: Mathematical statistics

Oct 8th, 2020

Lecture 12: Information

MATH 450: Mathematical statistics

| Week 2 · · · · ·  | Chapter 6: Statistics and Sampling<br>Distributions |
|-------------------|---|
| Week 4 · · · · ·  | Chapter 7: Point Estimation                         |
| Week 7 · · · · ·  | Chapter 8: Confidence Intervals                     |
| Week 10 · · · · • | Chapter 9: Test of Hypothesis                       |
| Week 11 · · · · · | Chapter 10: Two-sample inference                    |
| Week 13 · · · · · | Regression  |

▲日 ▶ ▲圖 ▶ ▲ 画 ▶ ▲ 画 ▶ →

æ

# Chapter 7: Overview

7.1 Point estimate

- unbiased estimator
- mean squared error
- 7.2 Methods of point estimation
  - method of moments
  - method of maximum likelihood.
- 7.3 Sufficient statistic
- 7.4 Information and Efficiency
  - Large sample properties of the maximum likelihood estimator

## Sufficient statistic

MATH 450: Mathematical statistics

・ロト ・回 ト ・ ヨト ・ ヨト …

æ

### Definition

A statistic  $T = t(X_1, ..., X_n)$  is said to be sufficient for making inferences about a parameter  $\theta$  if the joint distribution of  $X_1, X_2, ..., X_n$  given that T = t does not depend upon  $\theta$  for every possible value t of the statistic T.

T is sufficient for  $\theta$  if and only if nonnegative functions g and h can be found such that

$$f(x_1, x_2, \ldots, x_n; \theta) = g(t(x_1, x_2, \ldots, x_n), \theta) \cdot h(x_1, x_2, \ldots, x_n)$$

i.e. the joint density can be factored into a product such that one factor, h does not depend on  $\theta$ ; and the other factor, which does depend on  $\theta$ , depends on x only through t(x).

### Definition

The *m* statistics  $T_1 = t_1(X_1, \ldots, X_n)$ ,  $T_2 = t_2(X_1, \ldots, X_n)$ , ...,  $T_m = t_m(X_1, \ldots, X_n)$  are said to be jointly sufficient for the parameters  $\theta_1, \theta_2, \ldots, \theta_k$  if the joint distribution of  $X_1, \ldots, X_n$  given that

$$T_1=t_1, T_2=t_2,\ldots, T_m=t_m$$

does not depend upon  $\theta_1, \theta_2, \ldots, \theta_k$  for every possible value  $t_1, t_2, \ldots, t_m$  of the statistics.

 $T_1, T_2, \ldots, T_m$  are sufficient for  $\theta_1, \theta_2, \ldots, \theta_k$  if and only if nonnegative functions g and h can be found such that

$$f(x_1, x_2, \ldots, x_n; \theta_1, \theta_2, \ldots, \theta_k) = g(t_1, t_2, \ldots, t_m, \theta_1, \theta_2, \ldots, \theta_k)$$
$$\cdot h(x_1, x_2, \ldots, x_n)$$

• Let  $X_1, X_2, ..., X_n$  be a random sample from  $\mathcal{N}(\mu, \sigma^2)$ 

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

Prove that

$$T_1 = X_1 + \ldots + X_n, \qquad T_2 = X_1^2 + X_2^2 + \ldots + X_n^2$$

are jointly sufficient for the two parameters  $\mu$  and  $\sigma.$ 

★ ∃ ► < ∃ ►</p>

• Let  $X_1, X_2, ..., X_n$  be a random sample from a Gamma distribution

$$f_X(x) = \frac{1}{\Gamma(\alpha)\beta^{\alpha}} x^{\alpha-1} e^{-x/\beta}$$

where  $\alpha, \beta$  is unknown.

Prove that

$$T_1 = X_1 + \ldots + X_n, \qquad T_2 = \prod_{i=1}^n X_i$$

are jointly sufficient for the two parameters  $\alpha$  and  $\beta$ .

• • = • • = •

## Information

MATH 450: Mathematical statistics

## Definition

The Fisher information  $I(\theta)$  in a single observation from a pmf or pdf  $f(x; \theta)$  is the variance of the random variable  $U = \frac{\partial \ln f(X, \theta)}{\partial \theta}$ , which is

$$I(\theta) = Var \left[ \frac{\partial \ln f(X, \theta)}{\partial \theta} \right]$$

Note: We always have E[U] = 0

# Fisher information

We have

$$\sum_{x} f(x,\theta) = 1 \quad \forall \theta$$

Thus

$$E[U] = E\left[\frac{\partial \ln f(X,\theta)}{\partial \theta}\right]$$
$$= \sum_{x} \frac{\partial \ln f(x,\theta)}{\partial \theta} f(x,\theta)$$
$$= \sum_{x} \frac{\partial f(x,\theta)}{\partial \theta} = 0$$

문 문 문

## Problem

Let X be distributed by

$$\begin{array}{c|cc} x & 0 & 1 \\ \hline f(x,\theta) & 1-\theta & \theta \end{array}$$

Compute  $I(X, \theta)$ .

Hint:

• If 
$$x = 1$$
, then  $f(x, \theta) = \theta$ . Thus

$$u(x) = \frac{\partial \ln f(x,\theta)}{\partial \theta} = \frac{1}{\theta}$$

• How about x = 0?

돈 돈 돈

# Example

## Problem

Let X be distributed by

$$\begin{array}{c|c} x & 0 & 1 \\ \hline f(x,\theta) & 1-\theta & \theta \end{array}$$

Compute  $I(X, \theta)$ .

We have

$$Var[U] = E[U^2] - (E[U])^2 = E[U^2]$$
$$= \sum_{x=0,1} U^2(x)f(x,\theta)$$
$$= \frac{1}{(1-\theta)^2} \cdot (1-\theta) + \frac{1}{\theta^2} \cdot \theta$$

문 문 문

Assume a random sample  $X_1, X_2, ..., X_n$  from the distribution with pmf or pdf  $f(x, \theta)$  such that the set of possible values does not depend on  $\theta$ . If the statistic  $T = t(X_1, X_2, ..., X_n)$  is an unbiased estimator for the parameter  $\theta$ , then

$$Var(T) \geq rac{1}{n \cdot I( heta)}$$

Recall that E[U] = 0 and  $E[T] = \theta$  (since T is an unbiased estimator of  $\theta$ ) we have

$$Cov(T, U) = E[TU] - E[U] \cdot E[T]$$
$$= \sum_{x} t(x) \frac{\partial \ln f(x, \theta)}{\partial \theta} f(x, \theta)$$
$$= \sum_{x} t(x) \frac{\partial f(x, \theta)}{\partial \theta} \frac{1}{f(x, \theta)} f(x, \theta)$$
$$= \frac{\partial}{\partial \theta} \left( \sum_{x} t(x) f(x, \theta) \right) = 1$$

▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ …

æ

## The Cauchy–Schwarz inequality shows that

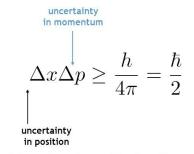
$$Cov(T, U) \leq \sqrt{Var(T) \cdot Var(U)}$$

which implies

$$Var(T) \geq rac{1}{I( heta)}.$$

∃ ► < ∃ ►</p>

# Heisenberg's Uncertainty Principle



The more accurately you know the position (i.e., the smaller  $\Delta x$  is), the less accurately you know the momentum (i.e., the larger  $\Delta p$  is); and vice versa

MATH 450: Mathematical statistics

| 4 同 ト 4 ヨ ト 4 ヨ ト

Let  $T = t(X_1, X_2, ..., X_n)$  is an unbiased estimator for the parameter  $\theta$ , the ratio of the lower bound to the variance of T is its efficiency

$$Efficiency = \frac{1}{nI(\theta)V(T)} \le 1$$

T is said to be an efficient estimator if T achieves the Cramer–Rao lower bound (i.e., the efficiency is 1).

Note: An efficient estimator is a minimum variance unbiased (MVUE) estimator.

Given a random sample  $X_1, X_2, ..., X_n$  from the distribution with pmf or pdf  $f(x, \theta)$  such that the set of possible values does not depend on  $\theta$ . Then for large n the maximum likelihood estimator  $\hat{\theta}$ has approximately a normal distribution with mean  $\theta$  and variance  $\frac{1}{n \cdot l(\theta)}$ . More precisely, the limiting distribution of  $\sqrt{n}(\hat{\theta} - \theta)$  is normal

with mean 0 and variance  $1/I(\theta)$ .

# Minimum variance unbiased estimator (MVUE)

MATH 450: Mathematical statistics

э

### Definition

Among all estimators of  $\theta$  that are unbiased, choose the one that has minimum variance. The resulting  $\hat{\theta}$  is called the minimum variance unbiased estimator (MVUE) of  $\theta$ .

Recall:

- Mean squared error = variance of estimator +  $(bias)^2$
- unbiased estimator  $\Rightarrow$  bias =0

 $\Rightarrow$  MVUE has minimum mean squared error among unbiased estimators

Question: Let  $X_1, \ldots, X_n$  be a random sample from a distribution with mean  $\mu$ . What is the best estimator of  $\mu$ ?

Answer: It depends.

- Normal distribution ightarrow sample mean  $ar{X}$
- Cauchy distribution ightarrow sample median  $ilde{X}$
- $\bullet \ {\sf Uniform} \ {\sf distribution} \ {\rightarrow}$

$$\hat{X}_e = \frac{\text{largest number} + \text{smaller number}}{2}$$

• In all cases, 10% trimmed mean performs pretty well

Let  $X_1, \ldots, X_n$  be a random sample from a normal distribution with mean  $\mu$ . Then the estimator  $\hat{\mu} = \bar{X}$  is the MVUE for  $\mu$ .

MATH 450: Mathematical statistics