

# AMS 572 Data Analysis I

## Basic Concepts of Inference

Pei-Fen Kuan

Applied Math and Stats, Stony Brook University

# Sampling

- ▶ A function of statistics is the provision of techniques for making inductive inferences, i.e, generalizations beyond the actual data in hand
- ▶ The goal of inductive inference is to find out something about a target population by examining a sample of it. (Read Chapter 3 of sampling designs).
- ▶ Inductive inference is accomplished by constructing a model that describes the origin of the data and a model for data collection (sampling)

**Definition:** The random variables  $X_1, \dots, X_n$  are called a *random sample* of size  $n$  from the population  $f$  if  $X_1, \dots, X_n$  are independent and identically distributed (iid) with marginal pdf or pmf  $f$ .

Note: The definition implies that the joint pdf or pmf of  $X_1, \dots, X_n$  is

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

# Statistics and sample moments

- ▶ A *statistic* is a function of observable random variables. The probability distribution of a statistic is called its sampling distribution
- ▶ Sample mean:

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

- ▶ Sample variance:

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2, n > 1$$

# Sampling from normal distribution

Let  $X_1, \dots, X_n$  be a random sample from  $N(\mu, \sigma^2)$ . Then  $\bar{X}_n$  has  $N(\mu, \sigma^2/n)$  distribution.

# Statistical Inference

- ▶ Statistical inference: Using statistics and probability theory to draw conclusions about parameters
- ▶ Two modes of inference:
  - ▶ **Estimation**: attempt to estimate value of parameter(s) and quantify uncertainty about these estimate(s)
  - ▶ **Hypothesis testing**: posit certain values for parameters and test whether the observed data are consistent with the hypothesis

- ▶ *Estimand*: parameter of interest we are trying to estimate; a constant; Eg  $\mu$
- ▶ *Estimator*: the statistic used to estimate the estimand; a random variable; Eg  $\bar{X}$
- ▶ *Estimate*: a realization of an estimator from an observed data set; Eg  $\bar{x} = 36.3$

# Point estimation:

## Methods of Maximum Likelihood and Moment Estimator

Definition: Let  $X_1, \dots, X_n$  be a sample with joint p.d.f or p.m.f  $f(\mathbf{x}|\theta)$ . Given that  $\mathbf{X} = \mathbf{x}$  is observed, the function of  $\theta$  defined by

$$L(\theta|\mathbf{x}) = f(\mathbf{x}|\theta)$$

is called the likelihood function.



# Point estimation: Method of Maximum Likelihood

Definition: For a given observed sample  $\mathbf{x}$ , let  $\theta(\hat{\mathbf{x}})$  be a value in the parameter space at which  $L(\theta|\mathbf{x})$  attains its maximum.  $\theta(\hat{\mathbf{x}})$  is called the maximum likelihood estimator (MLE) of  $\theta$ .

Example: Suppose  $X_i \sim N(\mu, \sigma^2)$  for  $i = 1, \dots, n$ . Derive the MLE for  $\mu$  and  $\sigma^2$ .

## Point estimation: Method of Moments

Method of moment method is one of the oldest method, and is simple to use. This method almost always yields some sort of estimate (MME).

Work by equating the sample moments to the population moments:

$$E(X) = \frac{1}{n} \sum_{i=1}^n X_i$$

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

$$\vdots$$

$$E(X^k) = \frac{1}{n} \sum_{i=1}^n X_i^k$$

# Point estimation: Method of Moments

Example: Suppose  $X_i \sim N(\mu, \sigma^2)$  for  $i = 1, \dots, n$ . Derive the MME for  $\mu$  and  $\sigma^2$ .

# Methods of Evaluating Estimators



Precise; Not Accurate



Accurate; Precise



Not Accurate; Not Precise



Accurate; Not Precise

[blackswanfarming.com](http://blackswanfarming.com)

# Methods of Evaluating Estimators

Definition: The bias of an estimator  $\hat{\theta}$  of  $\theta$  is

$$\text{bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$$

(measures accuracy)

Example: Are the MLEs of  $\mu$  and  $\sigma^2$  for normal distribution unbiased? If not, find an unbiased estimator for these parameters.

# Methods of Evaluating Estimators

Definition: The mean squared error (MSE) of an estimator  $\hat{\theta}$  of  $\theta$  is

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

(measures precision and accuracy)

In addition,  $\text{MSE}(\hat{\theta}) = \text{Var}(\hat{\theta}) + \text{bias}(\hat{\theta})^2$

# Methods of Evaluating Estimators

Example 6.4 of Tamhane and Dunlop:

Compare the MSE of these estimators for  $\sigma^2$ ,

$$S_1^2 = \sum_{i=1}^n (X_i - \bar{X})^2 / (n - 1) \text{ and } S_2^2 = \sum_{i=1}^n (X_i - \bar{X})^2 / n.$$

It can be shown that

$$\text{MSE}(S_1^2) = \frac{2\sigma^4}{n - 1}$$

and

$$\text{MSE}(S_2^2) = \frac{2n - 1}{n^2} \sigma^4$$

Thus,  $\text{MSE}(S_1^2) > \text{MSE}(S_2^2)$

- ▶ On average  $S_2^2$  will be closer to  $\sigma^2$  if MSE is used as criterion
- ▶ On average  $S_2^2$  underestimates  $\sigma^2$



# Confidence Interval

Definition: Let  $X_1, \dots, X_n$  be random variables with joint p.d.f or p.m.f  $f(\mathbf{x}|\theta)$ . The random interval  $(T_1(\mathbf{X}), T_2(\mathbf{X}))$  is called a  $100(1 - \alpha)\%$  confidence interval for  $\theta$  if

$$P(T_1(\mathbf{X}) < \theta < T_2(\mathbf{X})) = 1 - \alpha$$

where  $0 < \alpha < 1$ .

- ▶  $T_1(\mathbf{X})$  and  $T_2(\mathbf{X})$  are called the lower and upper confidence limits, respectively.
- ▶  $1 - \alpha$  is called the confidence coefficient.

# Confidence Interval Interpretation

- ▶ If we draw 100 different random samples, on average  $100(1 - \alpha)\%$  of them will contain  $\theta$
- ▶ How can one decrease the width of confidence interval?

We will study the pivotal quantity method for deriving confidence interval in Chapter 7 lecture notes.

# Hypothesis testing

1. Set up a hypothesis
  2. Collect data
  3. Infer from the data whether hypothesis is plausible
- Examples:
- Will folic acid supplementation reduce the risk of stroke?
  - Is the return the same in project 1 and project 2?

# Hypothesis testing

- ▶ Null hypothesis  $H_0$
- ▶ Example of null hypothesis:  
The incidence of stroke will be the same in those taking folic acid supplements and those not taking folic acid supplements  
Investing in project 1 will yield the same return as investing in project 2

# Null and Alternative

- ▶ In a test of a hypothesis, we are testing whether some population parameter has a particular value
- ▶ For example,

$$H_0 : \theta = \theta_0$$

where  $\theta_0$  is a known constant

- ▶ The **alternative hypothesis** is complement of null hypothesis

$$H_a : \theta \neq \theta_0$$

# Test statistic

- ▶ Once the data are collected, we will compute a *test statistic* related to  $\theta$ , say  $S(\hat{\theta})$
- ▶  $S(\hat{\theta})$  is a random variable, since it is computed from a sample
- ▶  $S(\hat{\theta})$  will have a particular probability distribution under the assumption  $H_0$ , say  $F_0[S(\hat{\theta})]$

# Test statistic

- ▶ Under  $F_0$ , we compute the probability that we would observe  $S(\hat{\theta})$  or a value more extreme than  $S(\hat{\theta})$  if the null  $H_0$  was true
- ▶ If this probability is large, the data are consistent with  $H_0$
- ▶ If this probability is small,  $H_0$  is probably not true

# Interpretation

- ▶ Usually if the probability is small, we conclude  $H_0$  is not true; i.e., we “reject”  $H_0$
- ▶ If the probability is large, we have not proven  $H_0$ . We say that “we failed to reject  $H_0$ ”
- ▶ We can never prove  $H_0$  is true!



# Significance Level

- ▶ How do we decide if the probability is too small?
- ▶ Prior to seeing the data, we select a value  $\alpha$  such that: if the computed probability is less than or equal to  $\alpha$ , we reject  $H_0$
- ▶  $\alpha$  is known as *significance level*

# Critical Region and Value

- ▶ We have a statistic  $S(\hat{\theta})$  with distribution  $F_0$  under the null hypothesis
- ▶ We specify  $\alpha$  and under  $F_0$  determine a *critical region* or *rejection region*  $C_\alpha$  such that

$$\Pr[S(\hat{\theta}) \in C_\alpha | H_0] = \alpha$$

- ▶ Values at the boundaries of  $C_\alpha$  are called *critical values*

# Critical Region and Value

- ▶ From the data we compute  $S(\hat{\theta})$
- ▶ If  $S(\hat{\theta}) \in C_\alpha$ , we reject  $H_0$
- ▶ If  $S(\hat{\theta}) \notin C_\alpha$ , the data are consistent with  $H_0$  and we do not reject  $H_0$

# Type I error, Type II error and Power

		Truth	
		$H_0$	$H_a$
Decision	Do not reject		Type II error
	Reject	Type I error	Power

# Tests of Hypotheses: Seven Steps

1. Design study
2. Establish null hypothesis
3. Determine test statistic to be employed
4. Choose significance level  $\alpha$  and establish  $C_\alpha$
5. Carry out study and collect data
6. Compute statistic from data
7. If statistic is in  $C_\alpha$ , reject  $H_0$

We will study step by step of setting up a hypothesis test in Chapter 7 slides.

Next lecture on R and SAS