

Probability and Random Process (SWE3026)

Statistical Inference

JinYeong Bak

jy.bak@skku.edu

College of Computing, SKKU

Objectives

We will:

- 1) Learn statistical inference**
- 2) Understand two different statistical inference methods**

Statistical Inference

Definition

A collection of methods that deal with drawing conclusions from data that are prone to random variation

Example

We want to predict the outcome of an election. But we cannot poll the entire population, so we will choose a random sample from the population.



Statistical Inference

Frequentist (classical) inference

The unknown quantity which we want to estimate from data is assumed to be a fixed (deterministic, non-random) quantity,

It is to be estimated by the observed data

Bayesian inference

The unknown quantity is assumed to be a random variable,

We have some initial guess about the distribution of the quantity and we update the distribution using Bayes Rule

Random Sampling

Definition

The collection of random variables X_1, X_2, \dots, X_n is said to be a random sample of size n if they are independent and identically distributed (i.i.d)

Random Sampling

Example

We want to identify the height distribution of people. So we define n random variables X_1, X_2, \dots, X_n where X_i is the height of the i -th sampled person with replacement. (X_i is the height of the i -th person that is chosen uniformly and independently from the population.)

We can estimate the average height in the population as follows:

$$\hat{\Theta} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

: point estimator for the average height in the population

Random Sampling

Properties

- The X_i 's are independent
- $F_{X_1}(x) = F_{X_2}(x) = \cdots = F_{X_n}(x) = F_X(x)$
- $E[X_i] = E[X] = \mu < \infty$
- $0 < Var(X_i) = Var(X) = \sigma^2 < \infty$

Sample Mean

Definition

$$\bar{X} = \frac{X_1 + X_2 + \cdots + X_n}{n}$$

Properties

- $E[\bar{X}] = \mu$
- $Var(\bar{X}) = \frac{\sigma^2}{n}$
- $\lim_{n \rightarrow \infty} P(|\bar{X} - \mu| \geq \epsilon) = 0$
- $\lim_{n \rightarrow \infty} P(Z_n \leq x) = \Phi(x)$ where $Z_n = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$

Good Estimator?

We want to know an unknown parameter θ (i.e., expected value of a R.V)

To estimate θ , we need some random samples X_1, X_2, \dots, X_n

Point estimator $\hat{\Theta} = h(X_1, X_2, \dots, X_n)$

One of the estimators for expected value of a R.V is sample mean

$$\hat{\Theta} = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

Here, how can we know a good estimator?

A good estimator $\hat{\Theta}$ is able to give us values that are close to the real value of θ

Bias

Definition

The bias of a estimator $\hat{\Theta}$:

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

Definition

Unbiased estimator:

$$B(\hat{\Theta}) = 0$$

Example

Let X_1, X_2, \dots, X_n be a random sample. Show that the sample mean

$$\hat{\theta} = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

is an unbiased estimator of $\theta = E[X_i]$

Mean Squared Error

How to ensure that an estimator is a good estimator?

Definition

Mean Squared Error (MSE) of a estimator $\hat{\Theta}$:

$$MSE(\hat{\Theta}) = E[(\hat{\Theta} - \theta)^2]$$

Example

Let X_1, X_2, \dots, X_n be a random sample. Its mean is $E[X_i] = \theta$, and variance $Var(X_i) = \sigma^2$. Show that $\hat{\Theta}_2 = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$ is better estimator than $\hat{\Theta}_1 = X_1$

Mean Squared Error

Definition

$$MSE(\hat{\Theta}) = E[(\hat{\Theta} - \theta)^2]$$

Properties

- $MSE(\hat{\Theta}) = Var(\hat{\Theta}) + B(\hat{\Theta})^2$
- Consistent estimator if $\lim_{n \rightarrow \infty} P(|\hat{\Theta}_n - \theta| \geq \epsilon) = 0$, for all $\epsilon > 0$

Sample Variance

Definition

$$s^2 = \frac{1}{n-1} \sum_{k=1} (X_k - \bar{X})^2 = \frac{1}{n-1} \left(\sum_{k=1} X_k^2 - n\bar{X}^2 \right)$$

Example

Let X_1, X_2, \dots, X_n be a random sample. Find the bias of sample variance.

$$s^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2 = \frac{1}{n-1} \left(\sum_{k=1}^n X_k^2 - n\bar{X}^2 \right)$$

Likelihood Estimation

Definition

Let X_1, X_2, \dots, X_n be a random sample from a distribution with a parameter θ . Suppose that we have observed $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$.

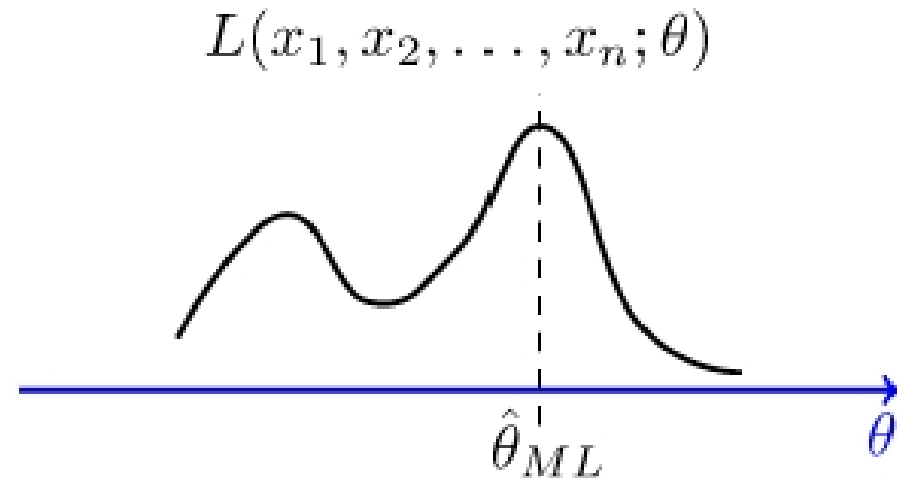
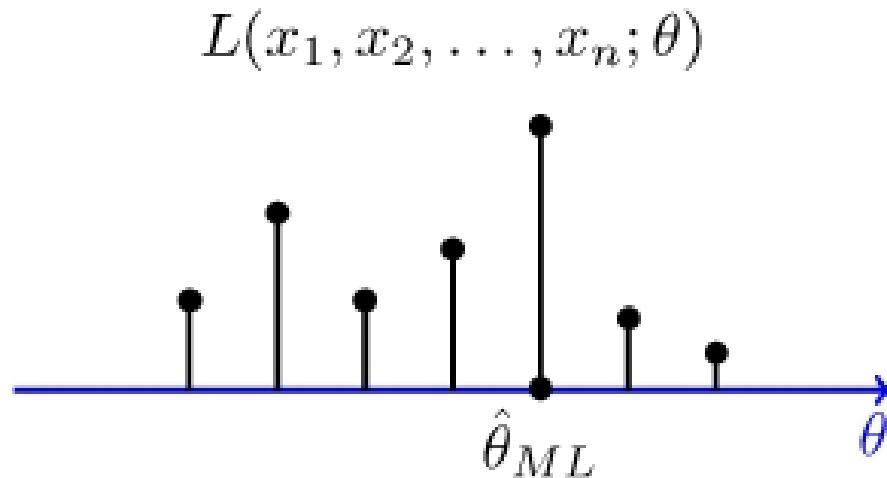
The likelihood function

- $L(x_1, x_2, \dots, x_n; \theta) = P_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta)$ if X_i s are discrete
- $L(x_1, x_2, \dots, x_n; \theta) = f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta)$ if X_i s are continuous

Maximum Likelihood Estimation

Definition

A systematic way of parameter estimation that finds the parameter value that maximizes the likelihood function



Example

I have a bag that contains 3 balls. Each ball is either red or blue, but I have no information in addition to this. Thus, the number of blue balls, call it θ , might be 0, 1, 2, or 3. I am allowed to choose 4 balls at random from the bag with replacement. We define the random variables X_1, X_2, X_3 , and X_4 as follows

$$X_i = \begin{cases} 1 & \text{if the } i\text{th chosen ball is blue} \\ 0 & \text{if the } i\text{th chosen ball is red} \end{cases}$$

Note that $X_i \sim \text{Bernoulli}(\frac{\theta}{3})$. After doing my experiment, I observe the following values of X_i s

$$x_1 = 1, x_2 = 0, x_3 = 1, x_4 = 1$$

1. Find the probability of the observed sample
2. For which value of θ is the probability of the observed sample is the largest?

I have a bag that contains 3 balls. Each ball is either red or blue, but I have no information in addition to this. Thus, the number of blue balls, call it θ , might be 0, 1, 2, or 3. I am allowed to choose 4 balls at random from the bag with replacement. We define the random variables X_1, X_2, X_3 , and X_4 as follows

$$X_i = \begin{cases} 1 & \text{if the } i\text{th chosen ball is blue} \\ 0 & \text{if the } i\text{th chosen ball is red} \end{cases}$$

Note that $X_i \sim \text{Bernoulli}(\frac{\theta}{3})$. After doing my experiment, I observe the following values of X_i s

$$x_1 = 1, x_2 = 0, x_3 = 1, x_4 = 1$$

1. Find the probability of the observed sample
2. For which value of θ is the probability of the observed sample is the largest?

$$\begin{aligned} P_{X_1} P_{X_2} P_{X_3} P_{X_4} &= P_{X_1}(1) \cdot P_{X_2}(0) \cdot P_{X_3}(1) \cdot P_{X_4}(1) \\ &= \left(\frac{\theta}{3}\right) \left(1 - \frac{\theta}{3}\right) \cdot \left(\frac{\theta}{3}\right) \cdot \left(\frac{\theta}{3}\right) \\ &= \left(\frac{\theta}{3}\right)^3 \cdot \left(1 - \frac{\theta}{3}\right) \end{aligned}$$

Example

For the following random samples, find the maximum likelihood estimate of θ :

$X_i \sim \text{Exponential}(\theta)$ and we have observed

$$(x_1, x_2, x_3, x_4) = (1.23, 3.32, 1.98, 2.12)$$