

## Chapter 5: Special topics

---

Yonghyun Kwon

Department of Mathematics, Korea Military Academy

# Sampling distribution

---

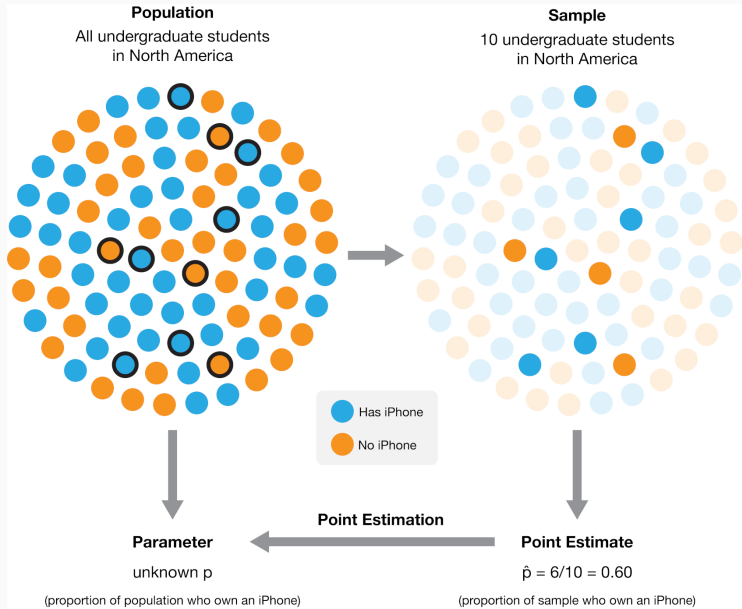
**Example:** Estimating the average lifetime of company light bulbs.

- *Population:* The entire set of objects under study.
- *Parameter:* A numerical value describing a characteristic of the population (e.g., population mean, population variance, population proportion).
- *Sample:* A subset drawn from the population.

## Key Idea

Because testing the entire population may be impossible (cost, time), we use samples to estimate population characteristics.





Suppose a population follows some probability distribution  $F$ .

- A **random sample** of size  $n$  is a collection of  $n$  independent random variables

$$X_1, X_2, \dots, X_n \sim F$$

drawn from the same distribution.

**Example:** Selecting 10 bulbs independently and recording their lifetimes. Then  $X_1, \dots, X_{10}$  are **independent and identically distributed (i.i.d.)** random variables with the same distribution as the population lifetime.



# Statistic and Sampling Distribution

- A **statistic** is a function of the sample that does not involve unknown parameters.
- The probability distribution of a statistic is called a **sampling distribution**.
- Since a statistic is computed from a random sample, it is itself a random variable.

## Examples of statistic:

(Sample mean) :  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

(Sample variance) :  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$



## Practice:

Let  $X$  be a discrete random variable taking values 0, 1, 2 with probabilities 0.3, 0.6, and 0.1. Taking a random sample of size  $n = 2$ , find the distribution of  $\bar{X} = (X_1 + X_2)/2$ .

$X$	0	1	2
$f(X)$	0.3	0.6	0.1

Population distribution of  $X$

$X_1 \backslash X_2$	0	1	2
0	0.09	0.18	0.03
1	0.18	0.36	0.06
2	0.03	0.06	0.01

Joint distribution of  $(X_1, X_2)$

$\bar{X}$	0	0.5	1	1.5	2
$P(\bar{X})$	0.09	0.36	0.42	0.12	0.01

Sampling distribution of  $\bar{X}$



# Expectation and Variance of the Sample Mean

## Expectation and Variance of the Sample Mean

Let  $X_1, X_2, \dots, X_n$  be a random sample from a population with mean  $\mu$  and variance  $\sigma^2$ . Then

$$E[\bar{X}] = \mu, \quad \text{Var}(\bar{X}) = \frac{\sigma^2}{n},$$

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

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

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n}$$

## Sample Mean from a Normal Population

### Recall: a linear combination of two RVs

If  $X \sim N(\mu_X, \sigma_X^2)$ ,  $Y \sim N(\mu_Y, \sigma_Y^2)$  are independent,  
 $aX + bY \sim N(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2)$

If  $X_1, \dots, X_n$  is a random sample from  $N(\mu, \sigma^2)$ ,

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right),$$

$$E(\bar{X}) = \mu$$

and standardizing  $\bar{X}$  gives

$$SD(\bar{X}) = \frac{\sigma}{\sqrt{n}}$$

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1).$$

$$Z = \frac{\bar{X} - \text{mean}}{SE} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$



## Practice

The fuel efficiency (km/L) of a new car model follows  $X \sim N(18, 2^2)$ . Suppose we randomly select  $n = 4$  cars and record their fuel efficiency  $X_1, \dots, X_4$ . Compute  $P(\bar{X} \leq 17)$ .

- Since  $X_i \sim N(18, 2^2)$ ,  $E(\bar{X}) = 18$ ,  $Var(\bar{X}) = \frac{\sigma^2}{n} = \frac{4}{4} = 1$

$$\bar{X} \sim N(18, 1)$$

- Standardizing:

$$\begin{aligned} P(\bar{X} \leq 17) &= P\left(\frac{\bar{X} - 18}{1} \leq \frac{17 - 18}{1}\right) \\ &= P(Z \leq -1) = 0.1587 \end{aligned}$$



# Central Limit Theorem (CLT)

Even if the population distribution is *not normal*:  $\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim N(0,1)$

## Central Limit Theorem (CLT)

For a random sample  $X_1, \dots, X_n$  from a population with mean  $\mu$  and variance  $\sigma^2 < \infty$ ,

$$\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \sim N(0, 1)$$

- When  $n$  is large,  $\bar{X}$  is approximately normal, regardless of the population distribution.



## Practice

The monthly salary (in 10,000 KRW) of new 4-year college graduates has mean 266 and standard deviation 15. Suppose we take a random sample of  $n = 81$  graduates. What is the probability that their average salary exceeds 270?

$$E(\bar{X}) = 266, \text{Var}(\bar{X}) = \frac{\sigma^2}{n} = \frac{15^2}{81}$$

Since  $n = 81$  is large enough, we can apply the CLT.

$$\bar{X} \sim N\left(266, \left(\frac{15}{9}\right)^2\right)$$

$$P(\bar{X} \geq 270) = P\left(Z \geq \frac{270 - 266}{5/3}\right) = P(Z \geq 2.4) \approx 0.0082.$$



## Sample Variance Distribution

- To make inference about population variance  $\sigma^2$ , we use the sample variance

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

- Its distribution is related to the *chi-squared distribution*.



## Chi-squared Distribution

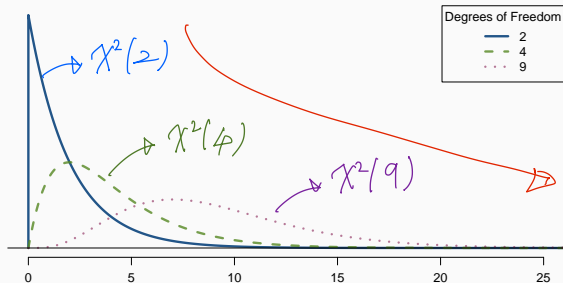
If  $Z_1, Z_2, \dots, Z_n \sim N(0, 1)$  independently,  $V := Z_1^2 + Z_2^2, \dots + Z_n^2$  follows chi-square distribution with *degrees of freedom*  $n$ .

$$V \sim \chi^2(n),$$

- For small  $n$ , the distribution is *skewed to the right*.
- As  $n$  increases, it becomes more symmetric and approaches normal.

# Chi-Square Distribution

Which of the following is false about chi-square distribution?



As the df increases,

- (a) the center of the  $\chi^2$  distribution increases as well.
- (b) the variability of the  $\chi^2$  distribution increases as well.
- (c) the shape of the  $\chi^2$  distribution becomes more skewed.

## Additivity of the Chi-Squared Distribution

If  $V_1 \sim \chi^2(n_1)$  and  $V_2 \sim \chi^2(n_2)$  independently, then

$$V_1 + V_2 \sim \chi^2(n_1 + n_2).$$

## Relationship to Normal distribution

If  $Z \sim N(0, 1)$ , then

$$Z^2 \sim \chi^2(1).$$

## Distribution of the Sample Variance

Let  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$  independently. Then:

- $\bar{X}$  and  $S^2$  are independent.
- The sample variance satisfies

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1).$$

---

*Note:* This is fundamental for building confidence intervals and hypothesis tests about  $\sigma^2$ .



## Distribution of the Sample Proportion

For a Bernoulli population with success probability  $p$ :

$$X_i \sim B(1, p), \quad i = 1, \dots, n, \text{ independently,}$$

the sample proportion is

$$\sum_{i=1}^n X_i = X \sim B(n, p)$$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i. \quad \hat{p} = \frac{X}{n} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}$$

Since  $\sum_{i=1}^n X_i \sim B(n, p)$ ,  $E(\hat{p}) = p$  and  $\text{Var}(\hat{p}) = \frac{p(1-p)}{n}$ .

By the CLT, when  $n$  is large,

$$\hat{p} \sim N\left(p, \frac{p(1-p)}{n}\right), \text{ or } \frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \sim N(0, 1).$$



## Practice

In a city, Party A's support rate is  $p = 0.6$ . A random sample of  $n = 30$  voters is surveyed. What is the probability that the sample proportion  $\hat{p}$  is less than 0.5?

The expectation and variance of the sample proportion  $\hat{p}$  are

$$E(\hat{p}) = 0.6 \quad \text{Var}(\hat{p}) = \frac{p(1-p)}{n} = \frac{0.6 \times 0.4}{30} = 0.008$$

Since  $np = 18 \geq 10$  and  $n(1-p) = 12 \geq 10$ , we can apply CLT:

$$\hat{p} \sim \mathcal{N}(0.6, 0.008)$$

Therefore,

$$P(\hat{p} \leq 0.5) = P\left(\frac{\hat{p} - 0.6}{\sqrt{0.008}} \leq \frac{0.5 - 0.6}{\sqrt{0.008}}\right) = P(Z \leq -1.12) \approx 0.1314.$$



## Point estimation

---

# Point Estimation

- *Point estimation* is the process of using sample information to estimate of an unknown population parameter.
- Statistics such as the sample mean, sample variance, and sample proportion are used as point estimators.
- *Estimator and Estimate:*
  - A statistic  $\hat{\theta}(X_1, \dots, X_n)$  used to estimate a parameter  $\theta$  is called an *estimator*.
  - The computed value  $\hat{\theta}(x_1, \dots, x_n)$  from observed data  $(x_1, \dots, x_n)$  is called the *estimate*.
- In general, an estimator of the parameter  $\theta$  is denoted by  $\hat{\theta}$ .



## Example: Estimator and Estimate

Suppose  $X_1, \dots, X_n$  is a random sample from a population with mean  $\mu$  and variance  $\sigma^2$ .

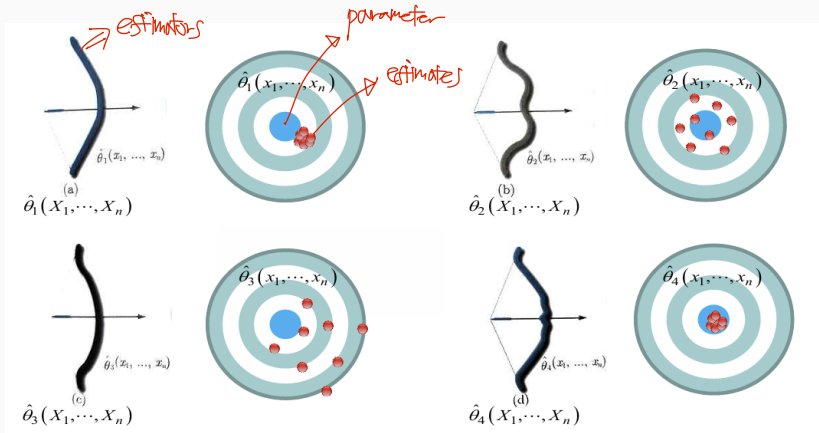
	Mean	Variance	Proportion
Parameter	$\mu$	$\sigma^2$	$P$
Estimator	$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$	$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$	$\hat{P}$
Estimate	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$	$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$	$\hat{p}$

Plugging in observed data  $x_1, \dots, x_n$  gives the estimate.



# Properties of Estimators

Multiple estimators can be used to estimate the same parameter.



Desirable properties: *Unbiasedness, Efficiency, Consistency*

## Unbiasedness

An estimator  $\hat{\theta}$  of parameter  $\theta$  is *unbiased* if

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



Its *bias* is defined as  $B(\hat{\theta}) = E(\hat{\theta}) - \theta = 0$

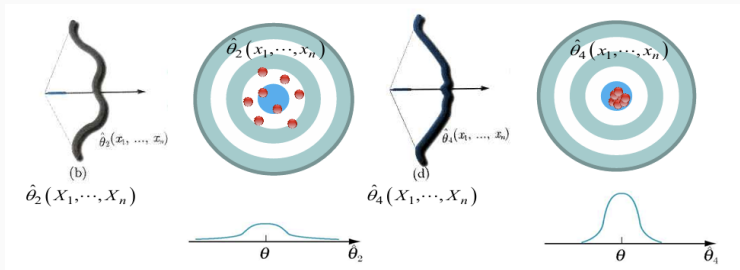
Suppose  $X_1, \dots, X_n$  is a random sample from a population with mean  $\mu$  and variance  $\sigma^2$ . Show that the sample mean  $\bar{X}$  is an unbiased estimator of  $\mu$ .

From slide 6,

$$\begin{aligned} E(\bar{X}) &= E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E(X_i) \quad \parallel \mu \\ &= \frac{1}{n} \sum_{i=1}^n \mu = \mu \end{aligned}$$

# Efficiency

Among unbiased estimators, the one with smaller variance is preferred.



Given two estimators  $\hat{\theta}_1$  and  $\hat{\theta}_2$  of the same parameter  $\theta$ , if

$$\text{Var}(\hat{\theta}_1) < \text{Var}(\hat{\theta}_2),$$

then  $\hat{\theta}_1$  is said to be more *efficient*.

## Example: Efficiency

Suppose  $X_1, X_2, X_3$  a random sample from a population with mean  $\mu$  and variance  $\sigma^2$ . Consider two estimators of  $\mu$ :

$$T_1 = \frac{X_1 + 2X_2 + X_3}{4}, \quad T_2 = \frac{X_1 + X_2 + X_3}{3}.$$

Which estimator is more efficient?

$$\text{Var}(T_1) = \frac{1}{16} \underbrace{\text{Var}(X_1)}_{=\sigma^2} + \frac{4}{16} \underbrace{\text{Var}(X_2)}_{=\sigma^2} + \frac{1}{16} \underbrace{\text{Var}(X_3)}_{=\sigma^2} = \frac{3}{8} \sigma^2$$

(independence)  
⇓

$$\text{Var}(T_2) = \frac{1}{9} (\text{Var}(X_1) + \text{Var}(X_2) + \text{Var}(X_3)) = \frac{1}{3} \sigma^2.$$

Since  $\frac{3}{8} \sigma^2 > \frac{1}{3} \sigma^2$ ,  $T_2$  is more efficient than  $T_1$



## Consistency

An estimator  $\hat{\theta}$  of parameter  $\theta$  is *consistent* if

$$\lim_{n \rightarrow \infty} E(\hat{\theta}) = \theta \text{ and } \lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}) = 0$$

As the sample size increases, a consistent estimator *converges* in probability to the true parameter.

## Example: Sample Mean

### Properties of the Sample Mean

If  $X_1, \dots, X_n$  is a random sample from a population with mean  $\mu$  and variance  $\sigma^2$ , then the sample mean  $\bar{X}$  is both an **unbiased** and **consistent** estimator of  $\mu$ .

Unbiasedness follows directly from slide 21. Also,

$$E(\bar{X}) = \mu, \quad \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

gives

$$\lim_{n \rightarrow \infty} E(\bar{X}) = \lim_{n \rightarrow \infty} \mu = \mu, \quad \lim_{n \rightarrow \infty} \text{Var}(\bar{X}) = \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} = 0$$

Hence  $\bar{X}$  is consistent for  $\mu$ .



## Example: Sample Variance

### Properties of the Sample Variance

If  $X_1, \dots, X_n$  is a random sample from a population with variance  $\sigma^2$ , then the sample variance  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$  is both an *unbiased* and *consistent* estimator of  $\sigma^2$ .

To check  $S^2$  is unbiased for  $\sigma^2$ , we show that  $E[S^2] = \sigma^2$ :

$$\begin{aligned}(n-1)E[S^2] &= E\left[\sum_{i=1}^n (X_i - \bar{X})^2\right] \\ &= \left(\sum_{i=1}^n E[(X_i - \mu)^2] - nE[(\bar{X} - \mu)^2]\right) \\ &= \left(n\sigma^2 - n \cdot \frac{\sigma^2}{n}\right) = (n-1)\sigma^2,\end{aligned}\tag{1}$$

and (1) holds because

$$\sum_{i=1}^n (X_i - \mu)^2 = \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^n (\bar{X} - \mu)^2 + 2(\bar{X} - \mu) \underbrace{\sum_{i=1}^n (X_i - \bar{X})}_0.$$



## Example: Sample Proportion

### Properties of the Sample Proportion

If  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$  independently, the sample proportion

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$$

$\hat{p} = \frac{X}{n}$       $X \sim B(n, p)$   
 $= \frac{1}{n} \sum_{i=1}^n X_i$

is both an *unbiased* and *consistent* estimator of  $p$ .

$p(1-p)$   
 $\frac{p(1-p)}{n}$

Recall that

$$E(\hat{p}) = p \quad \text{Var}(\hat{p}) = \frac{p(1-p)}{n}$$

from slide 16. Thus,  $\hat{p}$  is consistent for  $p$ . Also,

$$\lim_{n \rightarrow \infty} E(\hat{p}) = \lim_{n \rightarrow \infty} p = p, \quad \lim_{n \rightarrow \infty} \text{Var}(\hat{p}) = \lim_{n \rightarrow \infty} \frac{p(1-p)}{n} = 0$$

Hence  $\hat{p}$  is consistent for  $p$ .



## Exercises

1. An employee at a call center takes between 0.5 and 6.5 minutes to respond to a phone call, and this response time is assumed to be uniformly distributed. If the response times of 60 independent calls handled by this employee are observed, approximate the probability that the average response time is at least 4 minutes.



## Exercises

2. When Team A and Team B play a game in a sport, probabilities that Team A wins or loses are 0.6 and 0.4, respectively. Each outcome gives 3 points for a win and 0 points for a loss. Let  $X_1, X_2, \dots, X_n$  denote the points earned by Team A in  $n$  independent games (assume independence across games).

- (1) Find the joint probability mass function of  $X_1$  and  $X_2$ .
- (2) Let  $\bar{X}_2$  be the average score over 2 games. Find the probability mass function of  $\bar{X}_2$ .
- (3) Find the expectation  $E[\bar{X}_2]$  and variance  $Var(\bar{X}_2)$ .
- (4) Suppose the two teams play 54 games. Approximate the probability that Team A's average score is at least 2.



3. Suppose we obtain a random sample  $X_1, X_2, \dots, X_n$  from a population with probability density function

$$f(x) = \begin{cases} \frac{1}{3\theta}, & -\theta < x < 2\theta, (\theta > 0), \\ 0, & \text{otherwise.} \end{cases}$$

- (1) Find the constant  $k$  such that the estimator  $k\bar{X}$  is an unbiased estimator for  $\theta$ , where  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ .
- (2) For the value of  $k$  obtained in (1), show that  $k\bar{X}$  is a consistent estimator for  $\theta$ .

4. Suppose  $X_1, X_2, \dots, X_n$  are obtained from a Bernoulli distribution  $B(1, p)$ .

(1) Consider the estimators

$$\hat{p}_1 = \frac{1}{n} \sum_{i=1}^n X_i, \quad \hat{p}_2 = \frac{1}{2}(X_1 + X_2).$$

Determine whether each estimator is unbiased for  $p$ .

- (2) Between  $\hat{p}_1$  and  $\hat{p}_2$ , determine which estimator is more efficient (assume  $n > 2$ ).
- (3) Show that the estimator  $\hat{p}_1$  is a consistent estimator for  $p$ .

## Exercises

5. Suppose we obtain a random sample  $X_1, X_2, \dots, X_n$  from the population with density

$$f(x) = \begin{cases} \frac{1}{\theta + 1} e^{-x/(\theta+1)}, & x > 0, \theta > -1, \\ 0, & \text{otherwise.} \end{cases}$$

- (1) Find the expectation  $E(X_1)$  and variance  $\text{Var}(X_1)$ .
- (2) Consider the estimators

$$T_1 = X_1 - 1, \quad T_2 = \bar{X} - 1 = \left( \frac{1}{n} \sum_{i=1}^n X_i \right) - 1, \quad (n \geq 2).$$

Determine whether each is unbiased for  $\theta$ .

- (3) Between  $T_1$  and  $T_2$ , which estimator is more efficient?
- (4) Determine whether the more efficient estimator in (3) is consistent for  $\theta$ .

