

# UNIT 1.4 SAMPLING DISTRIBUTION OF THE MEAN

Prof. Ryung Kim  
ryung.kim@einstein.yu.edu

1

## UNIT 1.4 (PnG p.4)

2

- *This unit investigates the properties of the sample mean (i.e. average) when repeated samples are drawn from a population, thus introducing an important concept known as the central limit theorem. This theorem provides a foundation for quantifying the uncertainty associated with the inferences being made.*

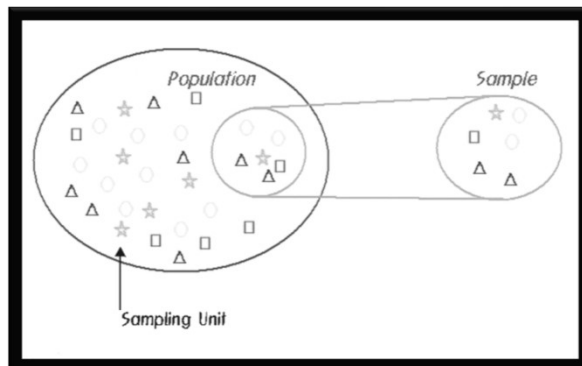
2

# 1. SAMPLING DISTRIBUTIONS

Prof. Ryung Kim  
ryung.kim@einstein.yu.edu

3

## What is Statistics?



### Population and Sample

- Population: the complete collection of all elements to be studied. All subjects to be studied are included.
- Sample: a subcollection of members selected from a population.

Often, it is impossible to investigate characteristics of whole population. So we investigate a part (sample) of population and make inference about population using the sample characteristics.

4

Image by Lee J, SNU

4

## Example

5

- To investigate the mean daily calorie intake among adults in Worcester, MA, an agency called randomly selected 2000 adults in Worcester who are older than 20 for phone interview.
  - ▣ Population
  - ▣ Sample

5

## Population and Sample

6

- A parameter is a number that describes the population.
  - ▣ e.g. Population mean
  - ▣ e.g. True mean daily calorie intake of all adults in Worcester
- A statistic is a number that describes a sample.
  - ▣ e.g. Sample mean
  - ▣ e.g. Average daily calorie intake of the sample of 2,000 people

6

## Sampling variability and Sampling distribution

7

- The value of a statistic, such as the sample mean  $\bar{x}$ , depends on the particular values included in the sample, and varies from sample to sample. This variability of a statistic is called *sampling variability*.

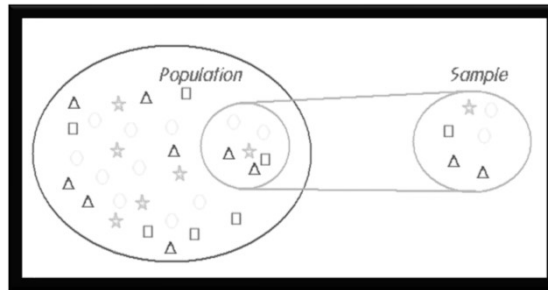


Image by Prof. Lee J, SNU

7

## Sampling variability and Sampling distribution

8

- If we repeatedly choose samples from a population, value of sample statistic will vary.
- The *sampling distribution* of the distribution of sample statistic value (e.g. sample mean or sample proportion) of all possible samples of the same size  $n$  taken from the same population.
- In other words, sampling distribution is a table describing the sampling variability.

8

## Sampling distribution

9

- **This is one of the most important concept in this course.**

9

### Difference sources of uncertainty: Variability often called 'Error'

- **Sampling error (non-systemic error; random error)**
  - Variability due to sampling
  - We can quantify the uncertainty.
- **Systematic error**
  - Measurement (instrumentation) error or observer Error
  - We cannot quantify the uncertainty.

10

## 2. THE CENTRAL LIMIT THEOREM

Prof. Ryung Kim  
ryung.kim@einstein.yu.edu

11

### Distribution of a sample mean

12

If  $X_1, X_2, \dots, X_n$  is an independent random sample with size  $n$  from a common distribution:

- The mean of the sample:

$$\mu_{\bar{x}} = \mu$$

- The standard deviation of sample mean (often called the *standard error*):

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

12

## Distribution of a sample mean

13

1. If the original population is itself normally distributed, then the sample means will be normally distributed for any sample size ( $n$ ).
2. For samples of size ( $n$ ) larger than 30, the distribution of the sample means can be approximated reasonably well by a normal distribution. The approximation gets better as the sample size  $n$  becomes larger. (“Central Limit Theorem”)

13

## Formal Central Limit Theorem

14

Let  $X_1, X_2, \dots, X_n$  be an independent random sample of size  $n$  from a distribution (which may or may not be normal) with mean  $\mu$  and standard deviation  $\sigma$ .

1. If sample size ( $n$ ) is large enough, the distribution of sample mean  $\bar{X}$  is approximately a normal distribution.
2. And the mean and standard deviation of the sample means are  $\mu$  and  $\sigma/\sqrt{n}$ .

14

### 3. APPLICATIONS OF THE DISTRIBUTION OF SAMPLE MEANS

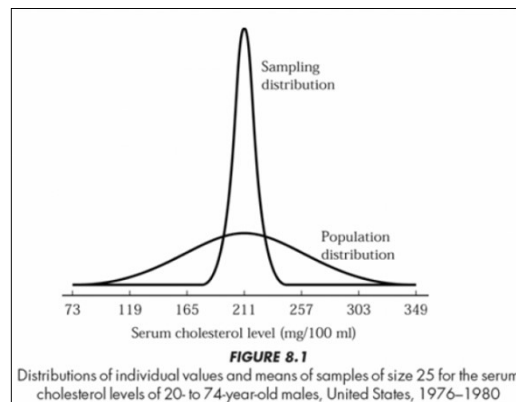
Prof. Ryung Kim  
ryung.kim@einstein.yu.edu

15

### Example – serum cholesterol levels

16

- Serum cholesterol levels for all 20~74 years old males in the U.S. have the mean 211 mg/100ml and the standard deviation 46 mg/100ml. (Population distribution is approximately normal.)
- If we select repeated samples of size 25 from the population, what proportion of the samples will have a mean value of 230 mg/ml or above?  
(Distribution of sample mean)



PnG

16



## Example – cont.

17

- The sample mean approximately follows normal distribution with mean 230 mg/ml and standard deviation  $46/\sqrt{25}$  (mg/ml).

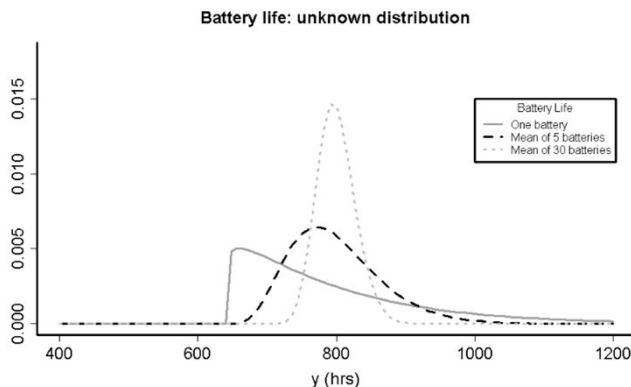
$$Z = \frac{230-211}{\frac{46}{\sqrt{25}}} = 2.065 \quad P(Z > 2.07) = 0.0195$$

17

## Example – Life of batteries

18

- The life of a certain brand battery has mean 800 hours and a standard deviation of 150 hours. The distribution of the battery life is known to be skewed (hence, not normal). What is the probability that the mean life time of 30 batteries is less than 700 hours? (→ Central Limit Theorem)



$$P\left(Z < \frac{700 - 800}{150/\sqrt{30}}\right) \\ = P(Z < -3.651) \\ = 0.0339$$

'Very unlikely'

18

- Reference

- Principles of Biostatistics (Pagano and Gauvreau)
- *Elementary Statistics* by Triola, 10<sup>th</sup> edition.
- *Applied Statistics for Engineers and Scientists* by Petrucci et al.

- Acknowledgements

- Prof. Jayson Wilbur, WPI
- Prof. Balgobin Nandram, WPI
- Prof. Lee Jaeyong, Seoul National University
- Some slides provided by Pearson Education, Inc Publishing as Pearson Addison-Wesley