

# Course Overview

Jeongyoun Ahn

KAIST

# Statistics, the Grammar of Science (Karl Pearson)

Life and most aspects of it is inherently random variable. "in the sense of repeated trials"

Ch 2 ~ Ch 6 → Uncertainty from the nature. (cannot be reduced)

Ch 7 ~ → Uncertainty from not enough information (= data)

Statistics allows us to make informed decisions in the face of uncertainty.

Statistics, as a field, lays ground rules for

- ▶ how data should be collected. ✱ →
- ▶ how decision can be made based on the data.

# Population and Samples

**Population** is the entire collection of objects on which an investigation is focused.

A **census** measures every member of the population. — impossible

A **sample** is any subset of the population. Our decisions are only as good as our sample.

A **variable** is a characteristic of interest for the objects in a population.

# Descriptive vs. Inferential Statistics


**Descriptive** statistics uses graphical or numerical methods to describe the sample.

**Inferential** statistics draws inference from the sample about the population.

- ▶ **Frequentist** approach is the focus of this course. We interpret probability as the long-run chance of an outcome's occurring in repeated trials.
- ▶ **Bayesian** approach is a popular alternative to the classical inference.

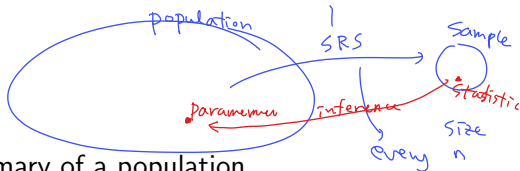
# Statistical Procedure

1. Set the goal - What do we want to show?
2. Collect data (experimental design) - What kind and how much data need to be collected? *What pop. do your data represent?*
3. Describe the data - summarize and describe the prominent features of data. (e.g., histogram, scatter plot, mean, variance, etc.)
4. Analyze the data (inferential statistics)
  - ▶ estimation, prediction, test, decision
  - ▶ generalize from a sample to a population
5. Conclusion based on the goal - how to assess the strength of the conclusion?

Statistics  a field of study  
plural form of statistic

Simple random  
? Sampling

# Parameter and Statistic



**Parameter:** a numerical summary of a population

- ▶ Population mean: average of a numerical measure
- ▶ Population proportion: fraction having a particular characteristic

**Statistic:** a numerical summary of a sample

- ▶ Sample mean
- ▶ Sample proportion

Inferences depend on the sample being representative of the population.

ideally  
SRS

# Classification of a variable I

- ▶ Qualitative
  - measurement is a set of unordered categories.
- ▶ Quantitative — *Value matters*
  - values of the variable differ in magnitude
- ▶ Ordinal
  - values are categories but with natural ordering

# Classification of a variable II

## ▶ Discrete

- takes finite (countable) number of values

e.g. education levels  
- blood types  
- no. books read

## ▶ Continuous

- can take any value within an interval, infinite possibilities

e.g. weight, . . . .

## \* Remark

- all categorical variables are discrete
- quantitative variables could be discrete or continuous
- sometimes it depends on a situation

age



# Descriptive study of data



Want to estimate population distribution using sample distribution.

- ▶ tabular/graphical representation
  - quantitative: frequency table, dot diagram, histogram, line diagram, stem-and-leaf display
  - categorical: contingency table, pie chart, bar chart

# Descriptive study of data



## ► numerical representation

- center: mean, median, percentiles, trimmed mean, Winsorized mean
- 10% Quantiles

- variation: variance, standard deviation, range, interquartile range

mean = 50  
Stdev = 10

## - empirical rule:

If data are bell-shape,

mean $\pm 1$ stdev	contains	68%
2	"	95%
3	"	99.7%

Quantiles

$Q_1$  - 25<sup>th</sup> percentile

$Q_2$  = median

$Q_3$  = 75<sup>th</sup> percentile

$$IQR = Q_3 - Q_1$$

# Descriptive study of data

Chebyshev

$\bar{X} \pm 2S$  will

contain 75%

Whisker  
► Boxplot

- Graphical display of five-number summary
- location, variation, skewness, outliers

