Lecture A3. Statistics Review

Sim, Min Kyu, Ph.D., mksim@seoultech.ac.kr

서울과학기술대학교 데이터사이언스학과

## Population and Sample

- A population set (모집단) is the entire group that you want to draw conclusions about.
- A sample set (표본집단) is the subset of population that you have an access to collect data from.
- The size of the sample is always less than the total size of the population.
- It is a researcher's primary concern to draw conclusion on the population set, by studying the behavior from the sample set.

$$\hat{m} \quad \hat{v} \quad \hat{a}$$

## Population statistics

- Population
  - Suppose that you are interested in Korean male's hand length. Let $X$ be a distribution [random var.] of population set (entire Korean male's hand length).
  - Let $\mu$ be the mean of $X$ and $\sigma^2$ be the variance of $X$.
  - That is, $\mu = \mathbb{E}X$ and $\sigma^2 = \mathbb{E}[(X - \mathbb{E}X)^2]$. (분산공식)
  - These *population statistics* are what we are after, specifically, *population mean* and *population variance*. (모분산)
  - Since these are what we aim to estimate, we often call them as *true values*, specifically *true mean* and *true variance*.
- Sample
  - In order to estimate $\mu$ and $\sigma^2$ you collect $n$ samples of Korean male's hand length.
  - Typically, these collected samples are denoted as $X_1, X_2, ..., X_n$, or $\{X_i, 1 \leq i \leq n\}$. sample set.

  $X$는 모두 서로다른 값을가짐.

  $X_1, X_2 \cdots X_n \longrightarrow \mu, \sigma^2.$

# Sample statistics

$\{X_i, 1 \leq i \leq N\} \rightarrow f(X_i) \approx \mu$

분포를 추어내고

비갈따이나.

sample mean

population mean

- Estimation
  - You want to draw conclusions on the *population mean* ($\mu$) and *population variance* ($\sigma^2$) by studying the sample $\{X_i, 1 \leq i \leq n\}$.
  - From the sample, we compute some value that should be similar to population statistics.
- Sample Mean
  - It is known that $\sum_{i=1}^{n} X_i / n$ is similar value to the population mean. ← $\mu$
  - This quantity is typically notated as $\overline{X}$, i.e., $\overline{X} = \sum_{i=1}^{n} X_i / n$.
  - This quantity is called as *sample mean* for obvious reason.
  - Sample mean is obtained by taking an arithmetic average of all samples.
- Sample Variance
  - It is known that $\frac{\sum (X_i - \overline{X})^2}{n-1}$ is similar value to the population variance. $s^2$ $\sigma^2$
  - This quantity is typically notated as $s^2$, i.e. $s^2 = \frac{\sum (X_i - \overline{X})^2}{n-1}$.
  - This quantity is called as *sample variance* for obvious reason.
  - Sample variance is obtained by 1) summing up squared deviations of all samples and 2) divide it by $n-1$.

- Summary

| | Mean | Variance |
|---|---|---|
| Population | $\mu = \mathbb{E}X$ | $\sigma^2 = \mathbb{E}[(X - \mathbb{E}X)^2]$ |
| Sample | $\overline{X} = \frac{\sum_{i=1}^n X_i}{n}$ | $s^2 = \frac{\sum (X_i - \overline{X})^2}{n-1}$ |

Population *Statistics*

Sample *Statistics*

$\hat{\mu}$ : $\mu$의 unbiased estimator ($ex)$ $\overline{X}$)

$\hat{\sigma}^2$ : $\sigma^2$의 " " ($ex)$ $s^2$, ·····)

# Estimation

- Remind that it is mentioned that 'Sample mean is *believed to be a similar value* to the population mean'.
- Like such, we call the process of 'Finding sample statistics that is *believed to be a similar value* to the population statistics.' as **estimation**. 통계량 추정.

  (true statistics)
- For true mean $\mu$, there may be various estimation efforts that aims to find similar value to the $\mu$. We call these *similar value to the true value*, as an **estimator**.
- Again, *estimator* is not a true value, but an estimation effort. To distinguish between the *true value* and *estimator*. Notation of 'hat', or $\hat{\cdot}$ is typically used. For example, $\hat{\mu}$ indicates an estimator for $\mu$, and $\hat{\sigma}^2$ indicates an estimator for $\sigma^2$.

- Sample mean serves as *an estimator* for the true mean.
- Sample variance serves as *an estimator* for the true variance.

## Desired properties of estimators

$100 \text{ of } X_1 \ldots X_{10}$

$M?$

A: $\hat{M} = \dfrac{X_1 + \cdots + X_{10}}{10}$

B: $\hat{M} = \dfrac{X_{(5)} + X_{(6)}}{2}$

- Is $\dfrac{\sum_{i=1}^{n} X_i}{n}$ a good estimator for the true mean? What it means by *good*?

- There are many criteria for *good* estimator such as $\mathbb{E}\hat{M} = M$
  - *unbiased* estimator - Expected value of estimator must be same as true value.
  - *consistent* estimator - As the number of sample increases, the estimator converges to the true value.
  - *maximum-likelihood (ML) estimator* - The probability that the estimator is exactly equal to true value is maximal.

- For mathematical expression, let's notate the true statistics we are after as $\theta$, and the estimator as $\hat{\theta}$. Then,
  - $\hat{\theta}$ is an *unbiased* estimator if $\mathbb{E}\hat{\theta} = \theta$.
  - $\hat{\theta}$ is a *consistent estimator* if $\hat{\theta} \to \theta$ as $n \to \infty$.
  - $\hat{\theta}$ is a *maximum-likelihood (ML)* estimator if $\hat{\theta} = argmax_x \mathbb{P}(\theta = x)$.

$$\hat{\theta} = \underset{x}{argmax} \ \mathbb{P}(\theta = x)$$

- It is known that *sample mean*
  - $\frac{\sum_{i=1}^{n} X_i}{n}$ is an *unbiased*, *consistent*, and *maximum-likelihood* estimator for the (true mean)
  - $\frac{\sum (X_i - \overline{X})^2}{n-1}$ ← *sample variance* is an *unbiased* and *consistent* estimator for the true variance, but it is not a *maximum-likelihood* estimator.
  - $\frac{\sum (X_i - \overline{X})^2}{n}$ is a *consistent* and *maximum-likelihood* estimator for the true variance, but it is not an *unbiased* estimator. In other words, it is *biased* estimator.

$\mu. \ \hat{\mu}, \ \overline{X}$

$6^2. \ \hat{6}^2, \ s^2$

$\overline{X} = \frac{\sum\limits_{i=1}^{n} X_i}{n}$

$s^2 = \frac{\sum\limits_{i=1}^{n} (X_i - \overline{X})^2}{n-1}$

- $\overline{X} \overset{?}{=} \hat{\mu} = ?$ 서로같음..?
- $s^2 \overset{?}{=} \hat{6}^2 = ?$ 서로같음?
- $\frac{\sum\limits_{i=1}^{n} (X_i - \overline{X})^2}{n} \overset{?}{=} \hat{6}^2 = ?$ 서로같음?

$\frac{\sum (X_i - \overline{X})^2}{n-1}$ vs $\frac{\sum (X_i - \overline{X})^2}{n}$

- Normal variable $X \sim N(\mu, \sigma^2)$

pdf

$\sigma^2$
large

$\mu$

pdf

$\sigma^2$
small

$\mu$

# Central limit theorem (CLT)

## Theorem 1

*For a random variable $X$, whatever the distribution of $X$ is, its sample mean $\overline{X}$ follows a normal distribution as long as the number of samples $n$ is larger than 30. That is*

$$\overline{X} \sim N(\mu, \sigma^2/n)$$

*(handwritten annotations: "estimator", "true value", $\frac{\sigma^2}{n} \to 0$ as $n \to \infty$)*

- It is intriguing that the population distribution may not be a normal distribution, but the sample mean from the population will always follow a normal distribution as long as the number of sample is larger than 30.
- It is also intriguing that the uncertainty of closeness between the estimator and true value is nicely quantified with the variance $\sigma^2/n$.
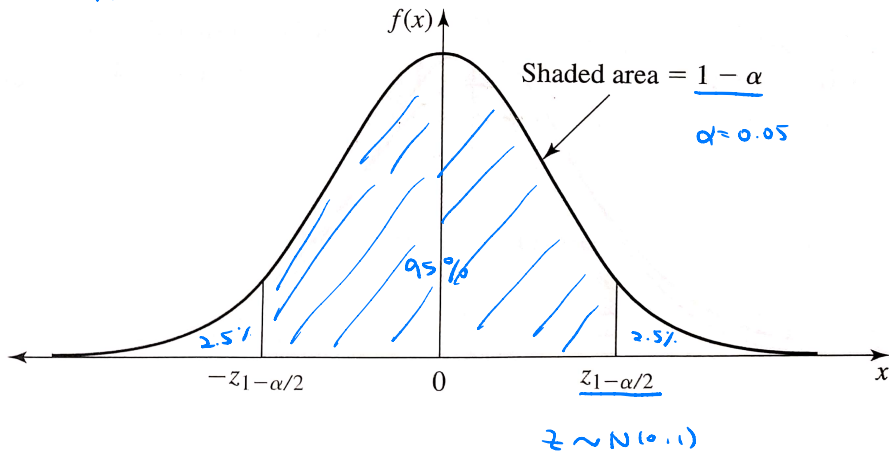
### Exercise 1

- *Is $\overline{X}$ an unbiased estimator for $\mu$, why or why not?*
- *Is $\overline{X}$ a consistent estimator for $\mu$, why or why not?*
- *Is $\overline{X}$ a ML estimator for $\mu$, why or why not?*

- ~~Questions~~

# Normal variable's quantile

random



$f(x)$

Shaded area = $1 - \alpha$

$\alpha = 0.05$

$95\%$

$2.5\%$        $2.5\%$

$-z_{1-\alpha/2}$        $0$        $z_{1-\alpha/2}$        $x$

$z \sim N(0, 1)$

$z_{1-0.05/2} = z_{0.975} = 1.96$

- From $\overline{X} \sim N(\mu, \sigma^2/n)$, we can use normal distribution's property to say:

$\mu$?
$x_1, x_2 \cdots x_n$
$\overline{X} \sim N(\mu, \sigma^2/n)$

$$\mathbb{P}[\mu - 1.96\frac{\sigma}{\sqrt{n}} \leq \overline{X} \leq \mu + 1.96\frac{\sigma}{\sqrt{n}}] = 0.95$$

① $\overline{X}$ oriented CI

- Two issues with the above confidence interval.
  1. The above expression is a confidence interval for the estimator $\overline{X}$ not for the true value $\mu$.
  2. We do not know the true value $\sigma$.

$z \sim N(0,1)$
$-z \sim N(0,1)$

- To tackle the first issue, the following effort is made.

$$\overline{X} \sim N(\mu, \sigma^2/n) \quad \Rightarrow \quad \frac{\overline{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1) = Z$$

② 0.95
$\mathbb{P}[\overline{X} - 1.96\frac{\sigma}{\sqrt{n}} \leq \mu \leq \overline{X} + 1.96\frac{\sigma}{\sqrt{n}}]$

"$\mu$ oriented CI"

$$\Rightarrow \quad \frac{\mu - \overline{X}}{\sigma/\sqrt{n}} \sim Z$$

$$\Rightarrow \quad \mu \sim N(\overline{X}, \sigma^2/n)$$

- From the last expression, $\mu \sim N(\overline{X}, \sigma^2/n)$, we still have the second issue of not knowing $\sigma$. We must replace $\sigma$ with $s$.

- In replacing $\sigma$ with $s$, it is known that $\boxed{\dfrac{\mu - \overline{X}}{\sigma/\sqrt{n}} \sim Z}$ becomes

$$\boxed{\frac{\mu - \overline{X}}{s/\sqrt{n}} \sim t_{n-1}}$$

$\sigma^2 = \mathbb{E}X^2 - (\mathbb{E}X)^2$

$s^2 = \dfrac{\Sigma(x_i - \overline{x})^2}{n-1}$

Student t-dist.

- Now we are ready to state the confidence interval for $\mu$ as following.

③

$$\mathbb{P}[\overline{X} - t_{0.975, n-1} \cdot s/\sqrt{n} \leq \mu \leq \overline{X} + t_{0.975, n-1} \cdot s/\sqrt{n}] = 0.95$$

- To get the some sence of what $t_{0.975, n-1}$ might be depending on $n$,
  - If $n = 30$, $\mathbb{P}[\overline{X} - 2.045 \cdot s/\sqrt{30} \leq \mu \leq \overline{X} + 2.045 \cdot s/\sqrt{30}] = 0.95$
  - If $n = 60$, $\mathbb{P}[\overline{X} - 2.000 \cdot s/\sqrt{60} \leq \mu \leq \overline{X} + 2.000 \cdot s/\sqrt{60}] = 0.95$
  - If $n = 120$, $\mathbb{P}[\overline{X} - 1.980 \cdot s/\sqrt{120} \leq \mu \leq \overline{X} + 1.980 \cdot s/\sqrt{120}] = 0.95$
  - If $n$ is bigger, $\mathbb{P}[\overline{X} - 1.960 \cdot s/\sqrt{n} \leq \mu \leq \overline{X} + 1.960 \cdot s/\sqrt{n}] = 0.95$
- For the most applications in this course, $n$ is so big enough that we are generally just fine using 1.96.

## Normal dist. vs $t$ dist.



$f(x)$

Standard normal distribution

$t$ distribution with 4 df

$n-1$

0

$x$

## Exercise 2

*You randomly sample 1,600 Korean male and measured their hand length. The sample mean is 20cm and the sample standard deviation is 2cm. What is the 95% confidence interval for Korean male's hand length?*

"Man can learn nothing unless he proceeds from the known to the unknown. - Claude Bernard"