

- 1 Population, random samples, statistics and sampling distribution
- 2 Law of large numbers and CLT
- 3 Parameter Estimation: Point Estimation

X

Population, random samples, statistics and sampling distribution

In statistics, a **population** is a set of objects or a certain kind of experiment that generates certain outcomes. A specific property of these objects is analyzed statistically.

Examples:

Population	Property
Undergraduate students in NTU	CGPA
Stars in the universe	Luminosity
Chess players in Singapore	Elo rating
Rolling a dice repeatedly	outcomes of rolls

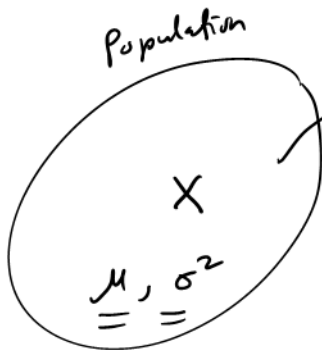
- Instead of the whole population, often only a **random subset** is selected (easier, more efficient) for measurements of the property of interest.

$n = \text{sample size}$

- These measurements x_1, x_2, \dots, x_n (also called **observations/data**) can be modelled by random variables X_1, X_2, \dots, X_n (called **random sample**), which are assumed to be **i.i.d (identically independently distributed)**,
- The distribution of the random variables X_i is called **population distribution**. ($\mathbb{E}[X_i]$ is called the **population mean**; $\text{Var}[X_i]$ is called the **population variance**).

...[continued]

- n is called the **sample size**.
- x_1, \dots, x_n can be viewed as realizations of i.i.d random variables X_1, \dots, X_n .



random sample.

$$\frac{X_1, X_2, \dots, X_n}{\text{i.i.d}}$$



estimate parameters
like μ, σ^2 , median, .
etc.

statistic $T(X_1, X_2, \dots, X_n)$

Example 1

- Population: Undergraduate students at NTU
- Property: CGPA
- Population Distribution: $N(\mu, \sigma^2)$
- Random sample: n randomly chosen NTU students X_1, \dots, X_n
- Observation/Data: $x_1, \dots, x_n \in [0, 5]$
- Statistical model: $X_1, \dots, X_n \text{ i.i.d } \sim N(\mu, \sigma^2)$.

e.g. If μ is unknown, we may estimate μ using the sample mean $\frac{X_1 + X_2 + \dots + X_n}{n}$.

Example 2

- Population: Tossing a fair coin 10 times
- Property: Number of heads among the 10 tosses.
- Population Distribution: $\text{Binomial}(10, 0.5)$
- Random sample: n repetitions of 10 tosses.
- Observation/Data: $x_1, \dots, x_n \in \{0, 1, \dots, 10\}$ ✓
- Statistical model: $X_1, \dots, X_n \text{ i.i.d } \sim \text{Binomial}(\underline{10}, \underline{0.5})$.

Let X_1, \dots, X_n be a random sample.

- A real valued function $T(X_1, \dots, X_n)$ is called a **statistic**.
- The distribution of a **statistic** is called a **sampling distribution**.

of T .

Example 3

Let X_1, \dots, X_n be a random sample. Some examples of statistics.

- $T_1 = \sum_{i=1}^n X_i^2$
- $T_2 = \min\{X_1, \dots, X_n\}$
- $T_3 = X_1$

$$T_4 = \frac{X_1 + X_2 + \dots + X_n}{n} \quad \checkmark$$

$$T_5 = \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}}$$

Let X_1, \dots, X_n be an i.i.d random sample.

- **Population distribution:** distribution of X_i
- **Sampling distribution:** distribution of a statistic based on X_1, \dots, X_n

\downarrow
 $T(X_1, \dots, X_n)$
 \downarrow
estimate of some
parameter like mean, variance,
median etc.

$X_i \sim \text{Population distribution.}$

Example 4

Let X_1, \dots, X_n be a random sample.

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

Note that \bar{X} and S^2 are statistics. Their distributions are examples of sampling distributions.

Theorem 5 (Random sample from Normal distribution)

Let X_1, \dots, X_n be observations of a random sample of size n from the normal distribution $N(\mu, \sigma^2)$. Then the sample mean \bar{X} and sample variance S^2 are independent, and ✓

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

Recall: $Z_i \sim N(0,1)$ i.i.d.
 $Z_i^2 \sim \chi^2(1)$

$$\sum_{i=1}^n Z_i^2 \sim \chi^2(n).$$

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\frac{(n-1)S^2}{\sigma^2} = \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2$$

$$= \sum_{i=1}^n z_i^2$$

$$\sim \chi^2(n-1)$$

$$z_i \sim N(0, 1).$$

$$\sum_{i=1}^n z_i = 0$$

$$\begin{aligned}
 \sum \left(\frac{x_i - \bar{x}}{s} \right) &= \frac{1}{s} \sum_{i=1}^n (x_i - \bar{x}) \\
 &= \frac{1}{s} \left(\sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} \right) \\
 &= \frac{1}{s} \left(n\bar{x} - n\bar{x} \right) \\
 &= 0.
 \end{aligned}$$

$$\bar{x} = \frac{\sum x_i}{n}$$

$$\underline{z_1 + z_2 + z_3 = 0}$$

$$\underbrace{z_1 + z_2 + \dots}_{\quad} + z_n = 0.$$

$$X_i \sim N(\mu, \sigma^2) \text{ i.e. } M_{X_i}(t) = e^{\mu t + \sigma^2 t^2 / 2}.$$

Theorem 6 (Random sample from Normal distribution)

Let X_1, \dots, X_n i.i.d $\sim N(\mu, \sigma^2)$. The sampling distribution of the sample mean \bar{X} is given by

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

This implies that the **standardized sample mean** $\frac{(\bar{X} - \mu)}{\sigma/\sqrt{n}} \sim N(0, 1)$.

Recall:

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

$$\text{MGF of } X = M_X(t) = e^{\mu t + \frac{\sigma^2 t^2}{2}}$$

$$\bullet M_{aX}(t) = M_X(at)$$

$$\bullet M_{X+Y}(t) = M_X(t) \cdot M_Y(t) \quad \checkmark$$

$$M_{\bar{X}}(t) = M_{\frac{1}{n}(X_1 + \dots + X_n)}(t)$$

$$= M_{X_1 + \dots + X_n}\left(\frac{1}{n}t\right) \quad \text{by } \star.$$

$$= \underbrace{M_{X_1}\left(\frac{1}{n}t\right)}_{\text{by } \star} \underbrace{M_{X_2}\left(\frac{1}{n}t\right)}_{\text{by } \star} \dots \underbrace{M_{X_n}\left(\frac{1}{n}t\right)}_{\text{by } \star} \quad \checkmark$$

$$= \left(e^{\mu \frac{t}{n} + \sigma^2 \frac{t^2}{2n^2}} \right) \left(e^{\mu \frac{t}{n} + \sigma^2 \frac{t^2}{2n^2}} \right) \dots \left(e^{\mu \frac{t}{n} + \sigma^2 \frac{t^2}{2n^2}} \right) \quad \text{by } \star \star$$

$$= e^{\left(\mu \frac{t}{n} + \sigma^2 \frac{t^2}{2n^2} \right) n} = e^{\mu t + \frac{\sigma^2}{n} \frac{t^2}{2}}.$$

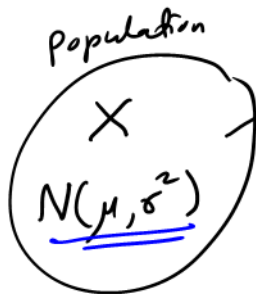
$$\text{So } \bar{X} \sim N(\mu, \frac{\sigma^2}{n}).$$

or equivalently,

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

Summary:

X_1, X_2, \dots, X_n



$$\bullet \bar{X} = \frac{X_1 + \dots + X_n}{n}$$

$$\bullet S^2 = \frac{\sum (X_i - \bar{X})^2}{n-1}$$

Now know: $\left\{ \begin{array}{l} \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1) \quad \checkmark \\ \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1) \quad \checkmark \end{array} \right.$

Q: what if population distribution
is not normal ?
is unknown ?

Remark: For increasing sample size n , the variance $\frac{\sigma^2}{n}$ tends to 0, and so the distribution of the sample mean \bar{X} tends to the distribution of the constant μ . It turns out that this is true even if the random sample is not from a normal distribution!

① we speak of sampling distribution
of statistic $T(X_1, \dots, X_n)$ ^{random sample}.

if $T(X_1, \dots, X_n) = \text{sample mean} = \bar{X}$

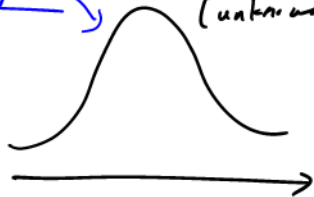
then sampling distribution of \bar{X} = distribution
of \bar{X}

(2)

random sample
sample size = n .



sampling distribution
(unknown).



$$R_1 = (x_1^1, x_2^1, \dots, x_n^1) \rightarrow \bar{x}^1$$

$$R_2 = (x_1^2, x_2^2, \dots, x_n^2) \rightarrow \bar{x}^2$$

...

$$R_m = (x_1^m, x_2^m, \dots, x_n^m) \rightarrow \bar{x}^m$$

\bar{x}
(depends
on n)

Law of large numbers and CLT

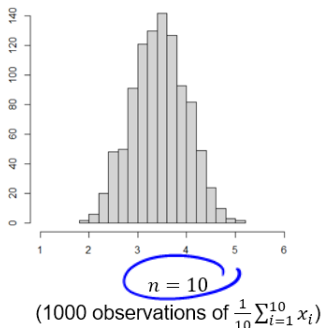
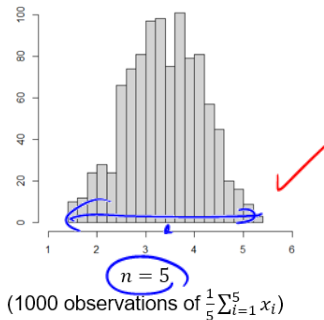
$$\begin{aligned} \text{True mean } \mu &= 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} \\ &\quad + 3 \cdot \frac{1}{6} + 4 \cdot \frac{1}{6} \\ &\quad + 5 \cdot \frac{1}{6} + 6 \cdot \frac{1}{6} \\ &= 3.5 \end{aligned}$$

An experiment:

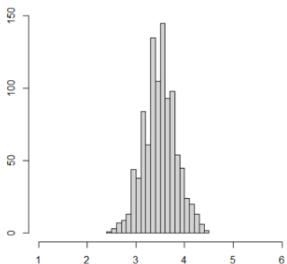
- Roll a fair dice n times.
- Compute average $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ where $x_i \in \{1, 2, \dots, 6\}$ is the outcome of the i th roll.
- Repeat this 1000 times to get 1000 observations for \bar{X} .
- Plot a histogram of these 1000 observations to visualize the distribution of the average.

generate sampling distribution of \bar{x} .

Distribution of \bar{X} (average result of rolling dice n times)

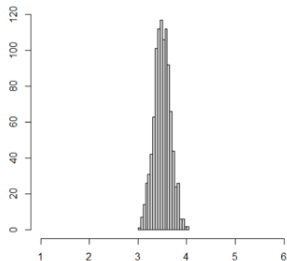


Distribution of \bar{X} (average result of rolling dice n times)



$n = 25$

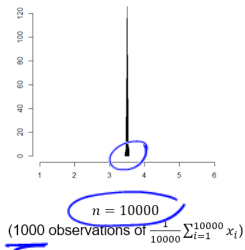
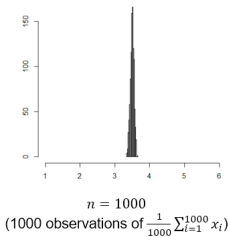
(1000 observations of $\frac{1}{25} \sum_{i=1}^{25} x_i$)



$n = 100$

(1000 observations of $\frac{1}{100} \sum_{i=1}^{100} x_i$)

Distribution of \bar{X} (average result of rolling dice n times)



From this experiment, when n increases, the probability that \bar{X} is close to the population mean $\mathbb{E}[X_i] = 3.5$ is getting higher. This fact is formalized by the **Law of Large Numbers**.

Theorem 7 (Law of Large Numbers)

Let X_1, \dots, X_n be i.i.d such that $\mu = \mathbb{E}[X_i]$ exists. Let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$.
Then

$$\mathbb{P}(|\bar{X} - \mu| < \epsilon) \rightarrow 1, \text{ as } n \rightarrow \infty,$$

for all $\epsilon > 0$.

In other words, for increasing sample size, the distribution of the sample mean \bar{X} tends to the distribution of the constant μ .

- In practice, we often encounter i.i.d random samples which are not normally distributed.
- The population distribution may even be totally unknown. ✓
- In this situation, the exact distribution of \bar{X} cannot be determined.
- For large samples, however, the Central Limit Theorem provides an approximation to the distribution of \bar{X} . ✓

Theorem 8 (Central Limit Theorem (CLT))

Let X_1, \dots, X_n i.i.d with $\mathbb{E}[X_i] = \mu$ and $\text{Var}[X_i] = \sigma^2 < \infty$. Then

$$\mathbb{P}\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq x\right) \rightarrow \Phi(x) \text{ for } n \rightarrow \infty.$$

Here, $\Phi(x)$ is the CDF of standard normal. ✓

This means for large n , the standardized sample mean $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ approximately has a standard normal distribution.

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

n large

The CLT is often used to **approximate** probabilities of sum of i.i.d:

$$\begin{aligned}\mathbb{P}(a \leq \sum_{i=1}^n X_i \leq b) &= \mathbb{P}\left(\frac{a}{n} \leq \bar{X} \leq \frac{b}{n}\right) \\&= \mathbb{P}\left(\frac{a - n\mu}{n} \leq \bar{X} - \mu \leq \frac{b - n\mu}{n}\right) \\&= \mathbb{P}\left(\frac{a - n\mu}{\sqrt{n}} \leq \sqrt{n}(\bar{X} - \mu) \leq \frac{b - n\mu}{\sqrt{n}}\right) \\&= \mathbb{P}\left(\frac{a - n\mu}{\sigma\sqrt{n}} \leq \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \leq \frac{b - n\mu}{\sigma\sqrt{n}}\right) \\&\approx \Phi\left(\frac{b - n\mu}{\sigma\sqrt{n}}\right) - \Phi\left(\frac{a - n\mu}{\sigma\sqrt{n}}\right),\end{aligned}$$

Handwritten notes: A blue arrow points from the expression $\sqrt{n}(\bar{X} - \mu)$ to $\sim N(0,1)$. A blue box is drawn around $\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}$. A blue checkmark is at the bottom right.

by CLT when n is large.

Example 9

X_1, \dots, X_{100} i.i.d $\sim \text{Bernoulli}(0.8)$. Approximate
 $\mathbb{P}(70 \leq X_1 + \dots + X_{100} \leq 90)$.

$n = 100$

Solution. $\mathbb{P}(70 \leq \sum_{i=1}^{100} X_i \leq 90)$

$$= \mathbb{P}\left(\frac{70}{100} \leq \frac{1}{100} \sum_{i=1}^{100} X_i \leq \frac{90}{100}\right)$$

$$= \mathbb{P}\left(\frac{0.7 - \mu}{\sigma/\sqrt{n}} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq \frac{0.9 - \mu}{\sigma/\sqrt{n}}\right)$$

$$\approx \mathbb{P}\left(\frac{0.7 - 0.8}{\sqrt{0.8(0.2)}/\sqrt{100}} \leq Z \leq \frac{0.9 - 0.8}{\sqrt{0.8(0.2)}/\sqrt{100}}\right) \quad (\text{by CLT})$$

$$\approx P(-2.5 \leq z \leq 2.5)$$

$$= \Phi(2.5) - \Phi(-2.5)$$



$$n = 20.$$

Example 10

Let X_1, X_2, \dots, X_{20} be a i.i.d random sample of size 20 from the uniform distribution $U(0, 1)$. Let $Y = X_1 + X_2 + \dots + X_{20}$. Use CLT to approximate the following probabilities.

(a) $\mathbb{P}(Y \leq 9.1)$; ✓

(b) $\mathbb{P}(8.5 \leq Y \leq 11.7)$. ✓

Solution. Note that $\mathbb{E}[X_i] = 1/2$ and $\text{Var}[X_i] = 1/12$ for $i = 1, \dots, 20$.

$$\mu = \frac{1}{2}$$

$$\frac{1}{3} - \frac{1}{4} = \left[\frac{x^3}{3} \right]_0^1 - \frac{1}{4}$$

$$\sigma^2 = \frac{1}{12} = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \int_0^1 x^2 dx - \left(\frac{1}{2}\right)^2$$

$$P(\bar{Y} \leq 9.1) = P\left(\sum_{i=1}^{20} X_i \leq 9.1\right)$$

$$= P\left(\frac{1}{20} \sum_{i=1}^{20} X_i \leq \frac{9.1}{20}\right)$$

$$= P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq \frac{\frac{9.1}{20} - \frac{1}{2}}{\frac{1}{12}/\sqrt{20}}\right)$$

$$= P\left(Z \leq \frac{\frac{9.1}{20} - \frac{1}{2}}{\frac{1}{12}/\sqrt{20}}\right)$$

$$\mathbb{P}(Y \leq 9.1) \approx \Phi\left(\frac{9.1 - 20(1/2)}{\sqrt{1/12}\sqrt{20}}\right) = \Phi(-0.6971) \stackrel{=}{=} 0.2429.$$

$$\begin{aligned} \mathbb{P}(8.5 \leq Y \leq 11.7) &\approx \Phi\left(\frac{11.7 - 20(1/2)}{\sqrt{1/12}\sqrt{20}}\right) - \Phi\left(\frac{8.5 - 20(1/2)}{\sqrt{1/12}\sqrt{20}}\right) \\ &= \Phi(-1.162) - \Phi(1.317) = 0.7835. \end{aligned}$$



Recall: MGF $X \sim \text{Poisson}(\lambda)$

$$M_X(t) = e^{\lambda(e^t - 1)}$$

Example 11

Explain how a Poisson distribution with mean $\lambda = 20$ can be approximated with the use of a normal distribution.

Let $Y \sim \text{Poisson}(\lambda = 20)$.

$\overset{11}{\text{Poisson}(1) + \dots + \text{Poisson}(1)}$
20 times.

Let $Y_i \sim \text{Poisson}(\lambda=1)$

$$\begin{aligned} M_{Y_1 + \dots + Y_{20}}(t) &= M_{Y_1}(t) \dots M_{Y_{20}}(t) \\ &= e^{1 \cdot (e^t - 1)} \dots e^{1 \cdot (e^t - 1)} \\ &= e^{(e^t - 1) \times 20} \\ &= M_Y(t) \quad Y \sim \text{Poisson}(20) \end{aligned}$$

By CLT:

$$Y_i \sim \text{Poisson}(\lambda=1)$$

$$Y = Y_1 + \dots + Y_{20}.$$

$$\frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \approx N(0, 1).$$

$$\frac{\frac{Y}{20} - 1}{\frac{1}{\sqrt{20}}} \approx N(0, 1) \quad \begin{array}{l} \mu = E[Y_i] = 1 \\ \sigma^2 = \text{Var}[Y_i] = 1 \end{array}$$
$$\boxed{\sqrt{20} \left(\frac{Y}{20} - 1 \right) \approx N(0, 1)}$$

e.g $\sqrt{n_0} \left(\frac{Y}{n_0} - 1 \right) \sim N(0, 1)$

$\hookrightarrow P(Y \leq a) = P\left(\frac{Y}{n_0} \leq \frac{a}{n_0}\right)$
 $= P\left(\frac{Y}{n_0} - 1 \leq \frac{a}{n_0} - 1\right)$
 $= P\left(\underbrace{\sqrt{n_0} \left(\frac{Y}{n_0} - 1\right)} \leq \sqrt{n_0} \left(\frac{a}{n_0} - 1\right)\right)$
 $\approx \underline{\underline{\Phi}}\left(\sqrt{n_0} \left(\frac{a}{n_0} - 1\right)\right) . \checkmark$

Take-away :

$$\textcircled{1} \quad \left. \begin{array}{l} \bar{X} \sim N(\mu, \frac{\sigma^2}{n}) \\ \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1) \end{array} \right\} \begin{array}{l} \text{assume} \\ \text{population} \\ \text{is } N(\mu, \sigma^2) \\ \text{i.e. } X_i \sim N(\mu, \sigma^2) \end{array}$$

$$\textcircled{2} \quad \text{CLT : } \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \quad n \text{ large}$$

use CLT to
approximate sum of i.i.d

$$\begin{aligned} \mu &= E[X_i] \\ \sigma^2 &= \text{Var}[X_i] \end{aligned}$$