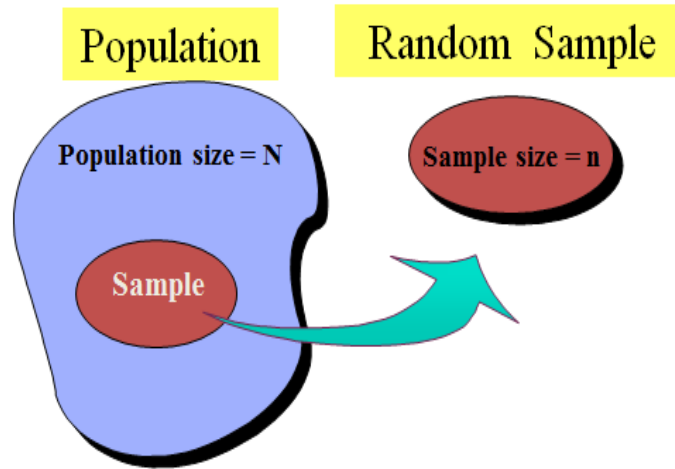


Summary---Topic 4: Sampling Distributions

Terminology



- Sample: a subset of population
- Statistics: characteristics / measures describe random sample
- A statistic is a random variable (A parameter is a constant, not a random variable)
- Sampling distribution: probability distribution of a statistic

Terminology

Population Distribution of Random Variable, X

Population mean,

$$\mu = \frac{\sum_{i=1}^N x_i}{N},$$

or= Expected value, $E[X] = \sum_i x_i P(x_i)$

(**Sample mean**, $\bar{X} = \frac{\sum_{i=1}^n x_i}{n}$, is an estimator of μ)

Population standard deviation, σ :

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

or= $\sqrt{\sum_i (x_i - E[X])^2 P(x_i)}$

(**Sample standard deviation,**

$S = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n-1}}$, is an estimator of σ)

Sampling Distribution of Sample Means

Mean of sample means,

$$\mu_{\bar{X}} = \sum_i \bar{x}_i P(\bar{x}_i) = \mu$$

Standard deviation of sample means,

$$\sigma_{\bar{X}} = \sqrt{\sum_i (\bar{x}_i - \mu_{\bar{X}})^2 P(\bar{x}_i)} = \sigma / \sqrt{n}$$

Central Limit Theorem(CLT)

- When sample size is large enough, the **sample means distribution** is approximately normal regardless of the population distribution.

$$\bar{X} \sim N(\mu_{\bar{X}}, \sigma_{\bar{X}}^2) \sim N(\mu, (\sigma/\sqrt{n})^2)$$

- Standardization of sample means distribution:

$$Z = \frac{\bar{X} - \mu_{\bar{X}}}{\sigma_{\bar{X}}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

- Rule of thumb: when sample size is larger than or equal to 30 (i.e. $n \geq 30$), we can apply the Central Limit Theorem.

Case	Population Distribution	Sample Size	Sampling Distribution of Means	Assumption Requirement
1	Normal	Whatever	Normal	Nil
2	Unknown / not normal	$n \geq 30$	By the CLT, the sample means distribution is approximately normal.	Nil
3	Unknown	$n < 30$	Not normal	Assume population distribution is normal so that the sample means also follow normal distribution.
4	Known, and not normal	$n < 30$	Not normal	Cannot make assumption as population distribution is known.

Exercises and Solutions

Q1. Given a normal distribution with $\mu = 100$ and $\sigma = 12$, if you select a sample of $n = 36$, what is the probability that \bar{X} is

- a) **Less than 95?**
- b) Between 95 and 97.5?
- c) Above 102.2?
- d) There is a 65% chance that \bar{X} is above what value?

*Mean of sample means, $\mu_{\bar{X}} = \sum_i \bar{x}_i P(\bar{x}_i) = \mu$

*Standard deviation of sample means, $\sigma_{\bar{X}} = \sqrt{\sum_i (\bar{x}_i - \mu_{\bar{X}})^2 P(\bar{x}_i)} = \sigma/\sqrt{n}$

Solution:

According to the question, we know that $X \sim N(100, 12^2)$.

Since $\mu_{\bar{X}} = 100$, and $\sigma_{\bar{X}} = \frac{12}{\sqrt{36}} = 2$,

and X is normally distributed, we conclude that the sampling distribution of the sample mean $\bar{X} \sim N(100, 2^2)$.

To compute non-standard normal probabilities: 1. Do standardization; 2. Check the Standard Normal Table.

$$\text{a) } P(\bar{X} < 95) = P\left(\frac{\bar{X}-100}{2} < \frac{95-100}{2}\right) = P(Z < -2.5) = 0.0062$$

Exercises and Solutions

Q1. Given a normal distribution with $\mu = 100$ and $\sigma = 12$, if you select a sample of $n = 36$, what is the probability that \bar{X} is

- a) Less than 95?
- b) Between 95 and 97.5?**
- c) Above 102.2?**
- d) **There is a 65% chance that \bar{X} is above what value?**

Solution: $\bar{X} \sim N(100, 2^2)$.

$$\begin{aligned} \text{b) } P(95 < \bar{X} < 97.5) &= P\left(\frac{95-100}{2} < \frac{\bar{X}-100}{2} < \frac{97.5-100}{2}\right) \\ &= P(-2.5 < Z < -1.25) = 0.1056 - 0.0062 = 0.0994 \end{aligned}$$

$$\begin{aligned} \text{c) } P(\bar{X} > 102.2) &= P\left(\frac{\bar{X}-100}{2} > \frac{102.2-100}{2}\right) = P(Z > 1.1) \\ &= 1 - P(Z < 1.1) = 1 - 0.8643 = 0.1357 \end{aligned}$$

d) We need find the value of “ a ” such that $P(\bar{X} > a) = 0.65 \rightarrow P(\bar{X} < a) = 0.35$

$$\rightarrow P\left(\frac{\bar{X}-100}{2} < \frac{a-100}{2}\right) = 0.35 \rightarrow P\left(Z < \frac{a-100}{2}\right) = 0.35 \rightarrow \frac{a-100}{2} = -0.39 \rightarrow a = 99.22$$

Q2. The diameter of a brand of Ping-Pong balls is normally distributed, with a mean of 1.30 inches and a standard deviation of 0.05 inch. If you select a random sample of 25 Ping-Pong balls,

- a) **What is the sampling distribution of the mean?**
- b) **What is the probability that the sample mean is less than 1.28 inches?**
- c) What is the probability that the sample mean is between 1.31 and 1.33 inches?
- d) The probability is 60% that the sample mean will be between what two values, symmetrically distributed around the population mean?

*Mean of sample means, $\mu_{\bar{X}} = \sum_i \bar{x}_i P(\bar{x}_i) = \mu$

*Standard deviation of sample means, $\sigma_{\bar{X}} = \sqrt{\sum_i (\bar{x}_i - \mu_{\bar{X}})^2 P(\bar{x}_i)} = \sigma / \sqrt{n}$

Solution:

a) Let X be the diameter of a brand of Ping-Pong balls, $X \sim N(1.3, 0.05^2)$.

Since $\mu_{\bar{X}} = 1.3$, and $\sigma_{\bar{X}} = \frac{0.05}{\sqrt{25}} = 0.01$, and X is normally distributed,

we conclude that the sample means is normally distributed, $\bar{X} \sim N(1.3, 0.01^2)$.

b)
$$P(\bar{X} < 1.28) = P\left(\frac{\bar{X} - 1.3}{0.01} < \frac{1.28 - 1.3}{0.01}\right) = P(Z < -2) = 0.0228$$

Q2. The diameter of a brand of Ping-Pong balls is normally distributed, with a mean of 1.30 inches and a standard deviation of 0.05 inch. If you select a random sample of 25 Ping-Pong balls,

c) What is the probability that the sample mean is between 1.31 and 1.33 inches?

d) The probability is 60% that the sample mean will be between what two values, symmetrically distributed around the population mean?

Solution:

$$\begin{aligned} \text{c) } P(1.31 < \bar{X} < 1.33) &= P\left(\frac{1.31-1.3}{0.01} < \frac{\bar{X}-1.3}{0.01} < \frac{1.33-1.3}{0.01}\right) = P(1 < Z < 3) \\ &= P(Z < 3) - P(Z < 1) = 0.99865 - 0.8413 = 0.15735 \end{aligned}$$

$$\text{d) } 60\% \text{ of the values are between two values } \rightarrow P(a < \bar{X} < b) = 0.6$$

$$\text{symmetrically distributed around the mean} \rightarrow \begin{cases} P(\bar{X} < a) = 0.2 \\ P(\bar{X} < b) = 0.8 \end{cases}$$

Do the standardization, we have

$$\begin{cases} P\left(\frac{\bar{X}-1.3}{0.01} < \frac{a-1.3}{0.01}\right) = 0.2 \\ P\left(\frac{\bar{X}-1.3}{0.01} < \frac{b-1.3}{0.01}\right) = 0.8 \end{cases} \quad \rightarrow \quad \begin{cases} P\left(Z < \frac{a-1.3}{0.01}\right) = 0.2 \\ P\left(Z < \frac{b-1.3}{0.01}\right) = 0.8 \end{cases} \quad \rightarrow \quad \begin{cases} \frac{a-1.3}{0.01} = -0.84 \\ \frac{b-1.3}{0.01} = 0.84 \end{cases}$$

$$\rightarrow a = 1.2916, b = 1.3084$$

Q3. Time spent using e-mail per session is normally distributed with $\mu = 8$ minutes and $\sigma = 2$ minutes. If you select a random sample of 16 sessions,

- a) **What is the probability that the sample mean is between 7.8 and 8.2 minutes?**
- b) What is the probability that the sample mean is between 7.5 and 8 minutes?
- c) If you select a random sample of 100 sessions, what is the probability that the sample means is between 7.8 and 8.2 minutes?
- d) Explain the difference in the results of (a) and (c).

*Mean of sample means, $\mu_{\bar{X}} = \sum_i \bar{x}_i P(\bar{x}_i) = \mu$

*Standard deviation of sample means, $\sigma_{\bar{X}} = \sqrt{\sum_i (\bar{x}_i - \mu_{\bar{X}})^2 P(\bar{x}_i)} = \sigma / \sqrt{n}$

Solution:

Let X be the time spent using e-mail per session, $X \sim N(8, 2^2)$.

When $n = 16$, since $\mu_{\bar{X}} = 8$, and $\sigma_{\bar{X}} = \frac{2}{\sqrt{16}} = 0.5$, and X is normally distributed,

we conclude that the sample means is normally distributed, $\bar{X} \sim N(8, 0.5^2)$.

$$\begin{aligned} \text{a) } P(7.8 < \bar{X} < 8.2) &= P\left(\frac{7.8-8}{0.5} < \frac{\bar{X}-8}{0.5} < \frac{8.2-8}{0.5}\right) = P(-0.4 < Z < 0.4) \\ &= P(Z < 0.4) - P(Z < -0.4) = 0.6554 - 0.3446 = 0.3108 \end{aligned}$$

Q3. Time spent using e-mail per session is normally distributed with $\mu = 8$ minutes and $\sigma = 2$ minutes. If you select a random sample of 16 sessions,

- b) What is the probability that the sample mean is between 7.5 and 8 minutes?**
- c) If you select a random sample of 100 sessions, what is the probability that the sample means is between 7.8 and 8.2 minutes?**
- d) Explain the difference in the results of (a) and (c).

Solution:

We conclude that, $\bar{X} \sim N(8, 0.5^2)$.

$$\begin{aligned} \text{b) } P(7.5 < \bar{X} < 8) &= P\left(\frac{7.5-8}{0.5} < \frac{\bar{X}-8}{0.5} < \frac{8-8}{0.5}\right) = P(-1 < Z < 0) \\ &= P(Z < 0) - P(Z < -1) = 0.5 - 0.1587 = 0.3413 \end{aligned}$$

c) When $n = 100$, since $\mu_{\bar{X}} = 8$, and $\sigma_{\bar{X}} = \frac{2}{\sqrt{100}} = 0.2$, and X is normally distributed, we conclude that the sample means is normally distributed, $\bar{X} \sim N(8, 0.2^2)$.

$$\begin{aligned} P(7.8 < \bar{X} < 8.2) &= P\left(\frac{7.8-8}{0.2} < \frac{\bar{X}-8}{0.2} < \frac{8.2-8}{0.2}\right) = P(-1 < Z < 1) \\ &= P(Z < 1) - P(Z < -1) = 0.8413 - 0.1587 = 0.6826 \end{aligned}$$

Q3. Time spent using e-mail per session is normally distributed with $\mu = 8$ minutes and $\sigma = 2$ minutes. If you select a random sample of 16 sessions,

d) Explain the difference in the results of (a) and (c).

Solution:

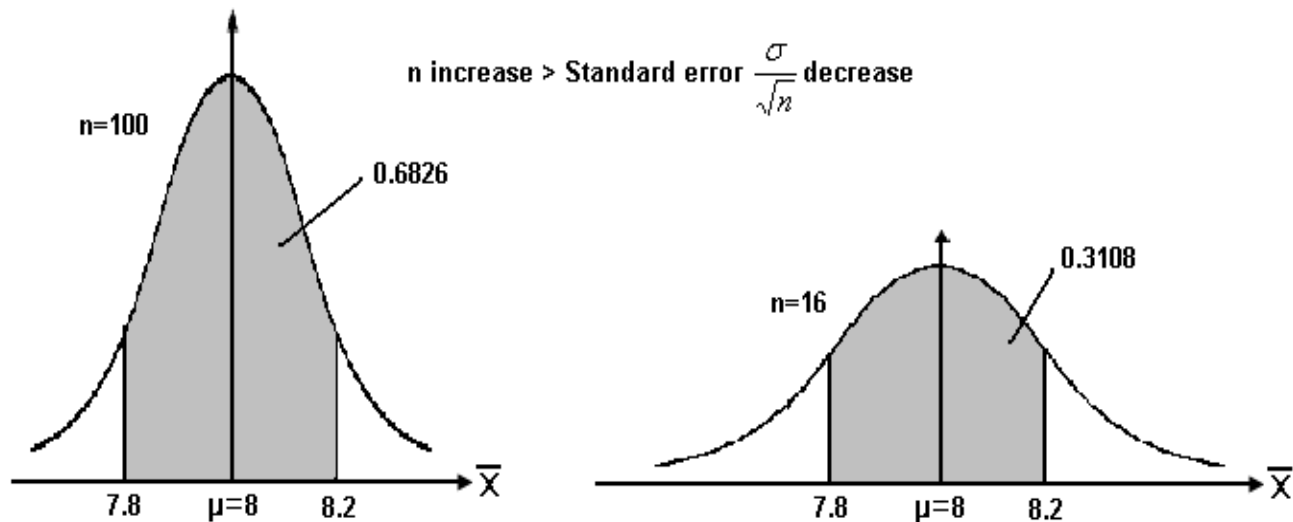
c) When $n = 100$,

$$\bar{X} \sim N(8, 0.2^2);$$

$$P(7.8 < \bar{X} < 8.2) = 0.6826$$

a) When $n = 16$, $\bar{X} \sim N(8, 0.5^2)$;

$$P(7.8 < \bar{X} < 8.2) = 0.3108$$



With the sample size increasing from $n = 16$ to $n = 100$, **more sample means will be closer to the population mean**. The standard error of the sample mean of size 100 is much smaller than that of size 16, so the likelihood(probability) that the sample mean will fall within 0.2 minutes of the population mean is much higher for samples of size 100 (probability = 0.6826) than for samples of size 16 (probability = 0.3108).

Q4. In a recent survey concerning the age (to the nearest year) and weight (to the nearest 10 lb) of first-year university students, the following probability distribution was obtained:

Age	Weight				
	100	110	120	130	140
19	0.02	0.09	0.09	0.01	0.02
20	0.06	0.15	0.2	0.05	0.03
21	0.02	0.06	0.11	0.04	0.05

A sample of 36 first-year students is taken. Find the approximate chance that their total weight is at most 4350 lb.

Solution:

Since $\alpha = 1 - 0.02 - 0.09 - \dots - 0.05 = 0.2$, we can summarize the distribution of weight, denoted as X:

Weight x_i	100	110	120	130	140
$P(x_i)$	0.1	0.3	0.4	0.1	0.1

The population mean: $\mu = E[X] = \sum_i x_i P(x_i) = 100 \times 0.1 + \dots + 140 \times 0.1 = 118$

The population std:

$$\begin{aligned}\sigma &= \sqrt{\sum_i (x_i - E[X])^2 P(x_i)} = \sqrt{(100 - 118)^2 \times 0.1 + \dots + (140 - 118)^2 \times 0.1} \\ &= \sqrt{116} = 10.77\end{aligned}$$

Q4. A sample of 36 first-year students is taken. Find the approximate chance that their total weight is at most 4350 lb.

Solution:

Weight x_i	100	110	120	130	140
$P(x_i)$	0.1	0.3	0.4	0.1	0.1

The population mean: $\mu = 118$

The population std: $\sigma = 10.77$

When $n = 36$, we obtain that $\mu_{\bar{X}} = 118$, and $\sigma_{\bar{X}} = \frac{10.77}{\sqrt{36}} = 1.795$.

Although X is not normally distributed, by Central Limit Theorem, the sampling distribution of mean is normal due to $n > 30$.

Therefore, we conclude that the sample means is normally distributed, $\bar{X} \sim N(118, 1.795^2)$.

$$\begin{aligned} P(\text{total weight} \leq 4350) &= P\left(\bar{X} \leq \frac{4350}{36}\right) = P\left(\frac{\bar{X}-118}{1.795} \leq \frac{4350/36-118}{1.795}\right) \\ &= P(Z \leq 1.58) = 0.9429 \end{aligned}$$

Q5*. At the CityU Computer Service Centre, the loading time for e-Portal page on Internet Explorer is normally distributed with mean 3 seconds.

A random sample of 5 computers is drawn. What is the chance that their total loading time is at least 15 seconds?

*Mean of sample means, $\mu_{\bar{X}} = \sum_i \bar{x}_i P(\bar{x}_i) = \mu$

*Standard deviation of sample means, $\sigma_{\bar{X}} = \sqrt{\sum_i (\bar{x}_i - \mu_{\bar{X}})^2 P(\bar{x}_i)} = \sigma/\sqrt{n}$

Solution:

Let X be the loading time, then $X \sim N(3, \sigma^2)$.

Since $\mu_{\bar{X}} = 3$, and $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{5}}$, and X is normally distributed, we get that the average loading time $\bar{X} \sim N(3, (\sigma/\sqrt{5})^2)$.

$$\begin{aligned} P(\text{total loading time} \geq 15) &= P\left(\bar{X} \geq \frac{15}{5}\right) = P(\bar{X} \geq 3) \\ &= P\left(\frac{\bar{X}-3}{\sigma/\sqrt{5}} \geq \frac{3-3}{\sigma/\sqrt{5}}\right) = P(Z \geq 0) = 0.5 \end{aligned}$$

Q6. Suppose there is a population with population size $N = 3$. The variable of interest is the salary (X) of individuals. The values of X are 18, 20 and 22 (in thousand dollars).

a) Find the mean (μ) and standard deviation (σ) for the population distribution.

In the process of developing sampling distribution, all possible samples (taken with replacement) of size $n = 2$ are obtained. The sample mean (\bar{X}) is considered as the sample statistic.

- b) What are the possible values of this sample mean random variable? Develop the probability distribution of the sample mean.
- c) Show that the sample statistic \bar{X} is an unbiased estimator of μ .
- d) Denote the standard deviation of \bar{X} , verify the following relationship: $\sigma_{\bar{X}} = \sigma/\sqrt{n}$.
- e) Does the sampling distribution of \bar{X} follows a Normal Distribution? Explain.

Solution:

a) The population mean: $\mu = \frac{\sum_{i=1}^N x_i}{N} = \frac{18+20+22}{3} = 20$

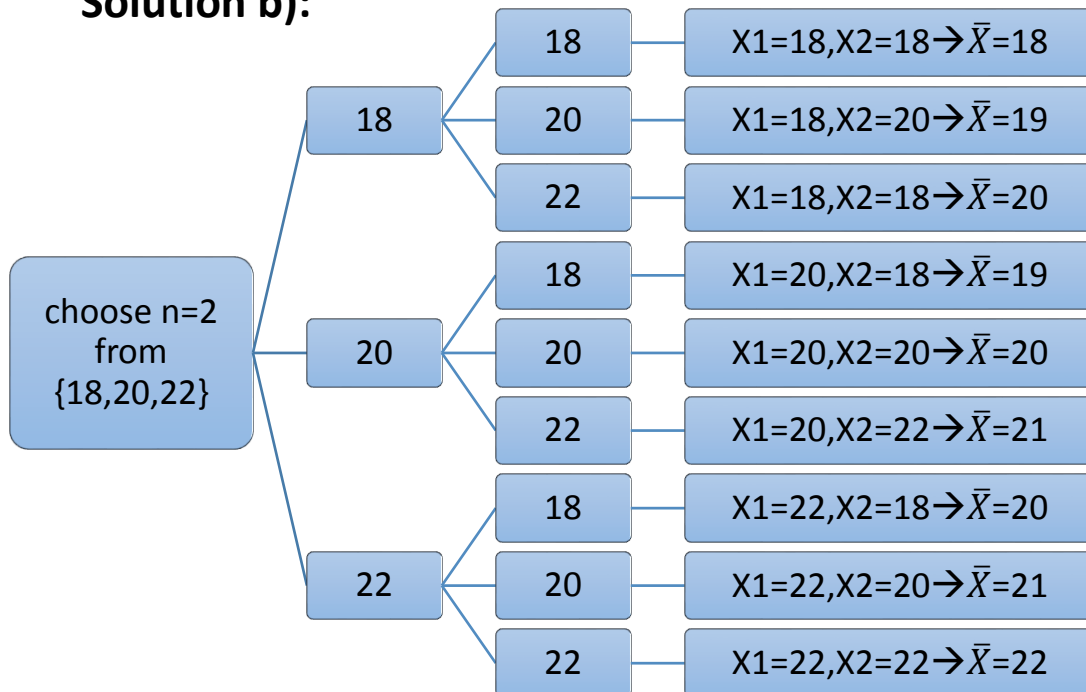
The population std: $\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}} = \sqrt{\frac{(18-20)^2 + (20-20)^2 + (22-20)^2}{3}} = 1.63$

Q6. Suppose there is a population with population size $N = 3$. The variable of interest is the salary (X) of individuals. The values of X are 18, 20 and 22 (in thousand dollars).

In the process of developing sampling distribution, all possible samples (taken with replacement) of size $n = 2$ are obtained. The sample mean (\bar{X}) is considered as the sample statistic.

b) What are the possible values of this sample mean random variable? Develop the probability distribution of the sample mean.

Solution b):



\bar{x}_i	Frequency	$P(\bar{x}_i)$
18	1	1/9
19	2	2/9
20	3	3/9
21	2	2/9
22	1	1/9

Q6. Suppose there is a population with population size $N = 3$. The variable of interest is the salary (X) of individuals. The values of X are 18, 20 and 22 (in thousand dollars).

In the process of developing sampling distribution, all possible samples (taken with replacement) of size $n = 2$ are obtained. The sample mean (\bar{X}) is considered as the sample statistic.

- c) Show that the sample statistic \bar{X} is an unbiased estimator of μ .
- d) Denote the standard deviation of \bar{X} , verify the following relationship: $\sigma_{\bar{X}} = \sigma/\sqrt{n}$.
- e) Does the sampling distribution of \bar{X} follows a Normal Distribution? Explain.

Solution: Recall that $\mu = 20, \sigma = 1.63$.

$$c) \mu_{\bar{X}} = \sum_i \bar{x}_i P(\bar{x}_i) = 18 \times \frac{1}{9} + \dots + 22 \times \frac{1}{9} = 20 = \mu.$$

$$\begin{aligned} d) \sigma_{\bar{X}} &= \sqrt{\sum_i (\bar{x}_i - \mu_{\bar{X}})^2 P(\bar{x}_i)} \\ &= \sqrt{(18 - 20)^2 \times \frac{1}{9} + \dots + (22 - 20)^2 \times \frac{1}{9}} = \sqrt{\frac{4}{3}} = 1.15 \\ &= \sigma/\sqrt{2}. \end{aligned}$$

\bar{x}_i	$P(\bar{x}_i)$
18	1/9
19	2/9
20	3/9
21	2/9
22	1/9

- e) The sampling distribution of \bar{X} does not follow a Normal Distribution.
 - population distribution of X is not normal
 - $n < 30$, cannot apply Central Limit Theorem

Q7. To investigate the length of time working for an employer, researchers at the CityU sampled 344 business students and asked them a question: Over the course of your lifetime, what is the maximum number of years you expect to work for any one employer? The resulting sample had sample mean $\bar{X} = 19.1$ years and sample standard deviation $S = 6$ years. Assume the sample of students was randomly selected from the 5800 undergraduate students in CityU.

- a) What are reasonable estimators of population mean and population standard deviation?**
- b) What is the sampling distribution of \bar{X} ? Why?**

Solution:

- a) Sample mean and sample standard deviation are reasonable estimators of population mean and population standard deviation.
- b) According to Central Limit Theorem, with large sample size ($n = 344$), sample mean \bar{X} follows a normal distribution approximately, with mean 19.1 and standard deviation $\frac{6}{\sqrt{344}} = 0.3235$.

Q7. To investigate the length of time working for an employer, researchers at the CityU sampled 344 business students and asked them a question: Over the course of your lifetime, what is the maximum number of years you expect to work for any one employer? The resulting sample had sample mean $\bar{X} = 19.1$ years and sample standard deviation $S = 6$ years. Assume the sample of students was randomly selected from the 5800 undergraduate students in CityU.

- c) If the population mean was 18.5 years, what is $P(\bar{X} \geq 19.1 \text{ years})$?**
- d) If the population mean was 19.5, what is $P(\bar{X} = 19.1 \text{ years})$?**
- e) If $P(\bar{X} \geq 19.1 \text{ years}) = 0.5$, what is the population mean?**

Solution:

c) $\bar{X} \sim N(18.5, 0.3235^2)$

$$\begin{aligned} P(\bar{X} \geq 19.1 \text{ years}) &= P\left(\frac{\bar{X} - 18.5}{0.3235} \geq \frac{19.1 - 18.5}{0.3235}\right) = P(Z \geq 1.8547) \\ &= 1 - P(Z < 1.8547) = 1 - 0.9678 = 0.0322 \end{aligned}$$

d) $\bar{X} \sim N(19.5, 0.3235^2)$, $P(\bar{X} = 19.1 \text{ years}) = 0$

e) We know that $P(\bar{X} \geq 19.1) = 0.5 \rightarrow P(\bar{X} < 19.1) = 0.5$

$$\rightarrow P\left(\frac{\bar{X} - \mu}{0.3235} < \frac{19.1 - \mu}{0.3235}\right) = 0.5 \rightarrow \frac{19.1 - \mu}{0.3235} = 0 \rightarrow \mu = 19.1 \text{ years}$$

Q7. To investigate the length of time working for an employer, researchers at the CityU sampled 344 business students and asked them a question: Over the course of your lifetime, what is the maximum number of years you expect to work for any one employer? The resulting sample had sample mean $\bar{X} = 19.1$ years and sample standard deviation $S = 6$ years. Assume the sample of students was randomly selected from the 5800 undergraduate students in CityU.

f) If $P(\bar{X} \geq 19.1 \text{ years}) = 0.2$, without calculation, can you tell that the population mean is greater or less than 19.1 years? Explain.

Solution:

f) Yes. The population mean will be less than 19.1 years.

Because $\bar{X} \sim N(\mu, 0.3235^2)$ and for a normal distribution mean is equal to median.

However, now we observe $P(\bar{X} \geq 19.1 \text{ years}) = 0.2 < 0.5 \rightarrow \text{median} < 19.1$

Therefore, the population mean should be less than 19.1 years.