

# Chapter 5:

## Probabilistic Features of the Distributions of Certain Sample Statistics

5.1 Introduction

5.2 sampling Distribution

5.3 Distribution of the Sample Mean  $\bar{X}$

5.4 Distribution of the Difference Between Two Sample Means  $\bar{X}_1 - \bar{X}_2$

5.5 Distribution of the Sample Proportion  $\hat{P}$

5.6 Distributions of the difference between two sample proportions  $\hat{P}_1 - \hat{P}_2$

# Introduction

In this Chapter we will discuss the probability distributions of some statistics.

As we mention earlier, a statistic is a measure computed from the random sample. As the sample values vary from sample to sample, the value of the statistic varies accordingly.

A statistic is a random variable; it has a probability distribution, a mean and a variance.

	Sample	Population
Size	$n$	$N$
Mean	$\bar{X}$	$\mu$
Variance	$S^2$	$\sigma^2$
Standard deviation	$S$	$\sigma$
Proportion	$\hat{P}$	$P$
Difference between two means	$\bar{X}_1 - \bar{X}_2$	$\mu_1 - \mu_2$
Difference between two proportions	$\hat{P}_1 - \hat{P}_2$	$P_1 - P_2$

**Statistics**

**Parameter**

## 5.2 sampling Distribution :

The probability distribution of a statistic is called the sampling distribution of that statistic.

The sampling distribution of the statistic is used to make statistical inference about the unknown parameter.

## 5.3 Distribution of the Sample Mean ( $\bar{X}$ ) :

Suppose that we have a population with mean  $\mu$  and variance  $\sigma^2$ .

Suppose that  $x_1, x_2, \dots, x_n$  is a random sample of size ( $n$ ) selected randomly from this population. We know that the sample mean is:

$$\bar{X} = \frac{\sum_{i=1}^n x_n}{n}$$

Suppose that we select several random samples of size  $n = 5$  :

	1 <sup>st</sup> sample	2 <sup>nd</sup> sample	3 <sup>rd</sup> sample	...	Last sample
Sample value	28	31	14	...	17
	30	20	31	...	32
	34	31	25	...	29
	34	40	27	...	31
	17	28	32	...	30
Sample Mean $\bar{x}$	28.4	29.9	25.8	...	27.8

- The value of the sample mean  $\bar{X}$  varies from random sample to another.
- The value of  $\bar{X}$  is random and it depends on the random sample.
- The sample mean  $\bar{X}$  is a random variable.
- The probability distribution of  $\bar{X}$  is called the sampling distribution of the sample mean  $\bar{X}$ .

### •Questions :

- What is the sampling distribution of the sample mean  $\bar{X}$  ?
- What is the mean of the sample mean  $\bar{X}$  ?
- What is the variance of the sample mean  $\bar{X}$  ?

### Result (1) : (Mean & variance of $\bar{X}$ )

If  $x_1, x_2, \dots, x_n$  is a random sample of size  $(n)$  from any distribution with mean  $\mu$  and variance  $\sigma^2$  , then:

1. The mean of  $\bar{X}$  is :  $\mu_{\bar{X}} = \mu$
2. The variance of  $\bar{X}$  is :  $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$
3. The standard deviation of  $\bar{X}$  is called the standard error and is defined by :  $\sigma_{\bar{X}} = \sqrt{\sigma_{\bar{X}}^2} = \frac{\sigma}{\sqrt{n}}$

### Result (2): (Sampling from normal population)

If  $x_1, x_2, \dots, x_n$  is a random sample of size  $(n)$  from a normal population with mean  $\mu$  and variance  $\sigma^2$ , that is  $\text{Normal}(\mu, \sigma^2)$ , then the sample mean  $\bar{X}$  has a normal distribution with mean  $\mu$  and variance  $\frac{\sigma^2}{n}$ , that is :

1.  $\bar{X} \sim \text{Normal}(\mu, \frac{\sigma^2}{n})$
2.  $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim \text{Normal}(0, 1)$

We use this result when sampling from **normal distribution with known variance  $\sigma^2$** .

### Result (3): (Central Limit Theorem: Sampling from Non-normal population)

Suppose that  $x_1, x_2, \dots, x_n$  is a random sample of size  $(n)$  from a non-normal population with mean  $\mu$  and variance  $\sigma^2$ . if the sample size  $n$  is large ( $n \geq 30$ ), then the sample mean  $\bar{X}$  has approximately a normal distribution with mean  $\mu$  and variance  $\frac{\sigma^2}{n}$ , that is :

1.  $\bar{X} \approx \text{Normal}(\mu, \frac{\sigma^2}{n})$  (approximately)
2.  $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \approx \text{Normal}(0, 1)$  (approximately)



## Notes :

- “ $\approx$ ” means “approximately distributed”.
- We use this result when sampling from non-normal distribution with known variance  $\sigma^2$  and with large sample size ( $n \geq 30$ ).

## Result (4): (used when $\sigma^2$ is unknown + normal distribution)

If  $x_1, x_2, \dots, x_n$  is a random sample of size ( $n$ ) from a normal distribution with mean  $\mu$  and variance  $\sigma^2$ , that is  $\text{Normal}(\mu, \sigma^2)$ , then the statistic :

$$1. \quad \bar{X} \sim t(n - 1)$$

$$2. \quad T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t(n-1)$$

Notation : degrees of freedom =  $df = v = n - 1$

We use this result when sampling from normal distribution with Unknown variance  $\sigma^2$  and with small sample size ( $n < 30$ ).

## Application:

### Example: (Sampling distribution of the sample mean)

Suppose that the time duration of a minor surgery is approximately normally distributed with mean equal to 800 seconds and a standard deviation of 40 seconds. Find the probability that a random sample of 16 surgeries will have average time duration of less than 775 seconds.

### Solution:

$X$  = the duration of the surgery

$$\mu = 800, \sigma = 40, \sigma^2 = 40^2 = 1600, n = 16$$

$$X \sim \text{Normal}(\mu = 800, \sigma^2 = 1600)$$

Calculating mean, variance, and standard error (standard deviation) of the sample mean  $\bar{X}$  :

The mean of  $\bar{X}$  :  $\mu_{\bar{X}} = \mu = 800$

The variance of  $\bar{X}$ :  $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n} = \frac{1600}{16} = 100$

The standard deviation of  $\bar{X}$ :  $\sigma_{\bar{X}} = \sqrt{\sigma_{\bar{X}}^2} = \frac{\sigma}{\sqrt{n}} = \frac{40}{\sqrt{16}} = 10$

The sampling distribution of  $\bar{X}$  is :

mean  $\mu_{\bar{X}} = 800$  and variance  $\sigma_{\bar{X}}^2 = 100$  . That is :

$$\bar{X} \sim \text{Normal}(800, 100)$$

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{\bar{X} - 800}{10} \sim \text{Normal}(0, 1)$$

**Result (2):**  
normal distribution with  
known variance  $\sigma^2$

The probability that a random sample of 16 surgeries will have average time duration of less than 775 seconds equals to :

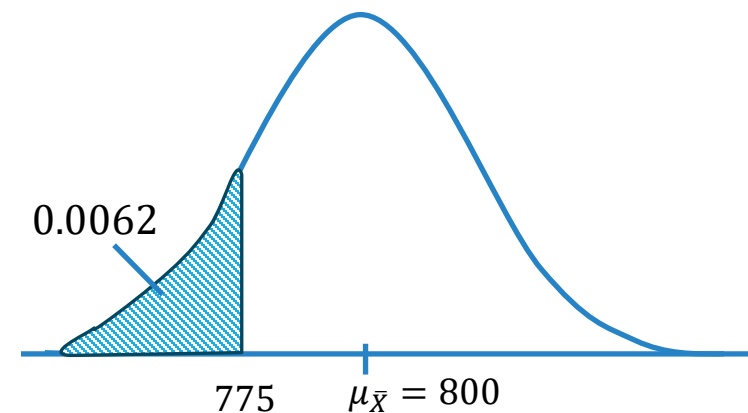
$$P(\bar{X} < 775) = P\left(\frac{\bar{X} - 800}{10} < \frac{775 - 800}{10}\right) = P(Z < -2.5) = 0.00621$$

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) = N(800, 100)$$

$$\mu_{\bar{X}} = 800$$

$$\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n} = 100$$

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} = 10$$



### Example :

If the mean and standard deviation of serum iron values for healthy men are 120 and 15 microgram/100ml, respectively, what is the probability that a random sample of size 50 normal men will yield a mean between 115 and 125 microgram/100ml ?

### Solution :

$X$  = the serum iron value

$$\mu = 120, \sigma = 15, \sigma^2 = 15^2 = 225, n = 50$$

$$X \sim \text{Normal} (\mu = 120, \sigma^2 = 225)$$

**Calculating mean, variance, and standard error (standard deviation) of the sample mean  $\bar{X}$  :**

The mean of  $\bar{X}$  :  $\mu_{\bar{X}} = \mu = 120$

The variance of  $\bar{X}$  :  $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n} = \frac{225}{50} = 4.5$

The standard deviation of  $\bar{X}$  :  $\sigma_{\bar{X}} = \sqrt{\sigma_{\bar{X}}^2} = \frac{\sigma}{\sqrt{n}} = \frac{15}{\sqrt{50}} = 2.1213$

Using the central limit theorem  $\bar{X}$  has a normal distribution

with mean  $\mu_{\bar{X}} = 120$  and variance  $\sigma_{\bar{X}}^2 = 4.5$

That is :

$$\bar{X} \sim \text{Normal}(120, 4.5)$$

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{\bar{X} - 120}{2.12} \sim \text{Normal}(0, 1)$$

**Result (3):**  
non-normal distribution with  
known variance  $\sigma^2$  and with  
large sample size ( $n \geq 30$ ).

The probability that a random sample of size 50 men will yield a mean between 115 and 125 microgram/100ml equals to :

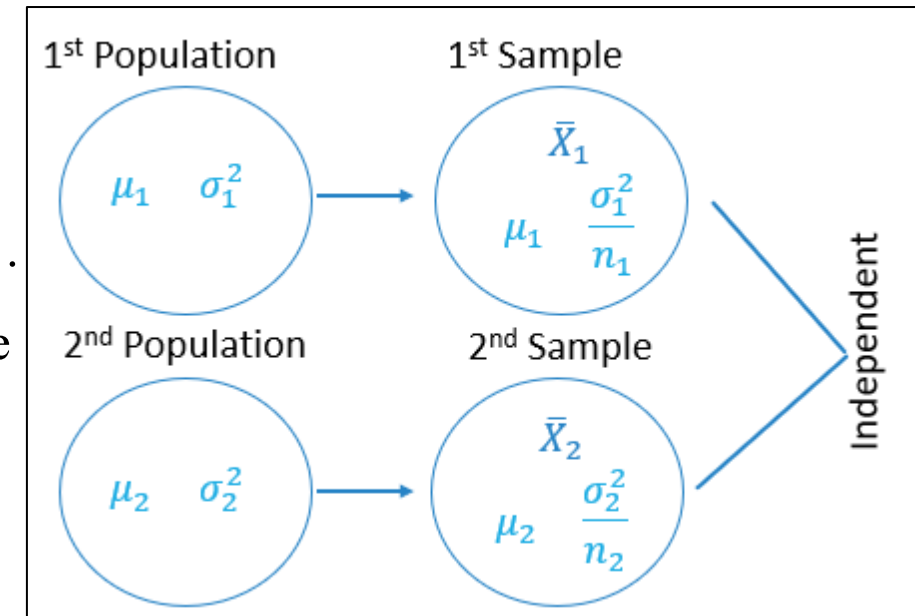
$$\begin{aligned} P(115 < \bar{X} < 125) &= P\left(\frac{115-120}{2.1213} < \frac{\bar{X}-120}{2.12} < \frac{125-120}{2.1213}\right) \\ &= P(-2.36 < Z < 2.36) \\ &= P(Z < 2.36) - P(Z < -2.36) \\ &= 0.99086 - 0.00914 \\ &= 0.98172 \end{aligned}$$

## 5.4 Distribution of the Difference Between Two Sample Means ( $\bar{X}_1 - \bar{X}_2$ ) :

Suppose that we have two populations:

- 1-st population with mean  $\mu_1$  and variance  $\sigma_1^2$  .
- 2-nd population with mean  $\mu_2$  and variance  $\sigma_2^2$  .
- We are interested in comparing  $m_1$  and  $m_2$ , or equivalently, making inferences about the difference between the means ( $\mu_1 - \mu_2$ ) .
- We independently select a random sample of size  $n_1$  from the 1-st population and another random sample of size  $n_2$  from the 2-nd population .

- Let  $\bar{X}_1$  and  $S_1^2$  be the sample mean and the sample variance of the 1-st sample .
- Let  $\bar{X}_2$  and  $S_2^2$  be the sample mean and the sample variance of the 2-nd sample
- The sampling distribution of  $\bar{X}_1 - \bar{X}_2$  is used to make inferences about  $\mu_1 - \mu_2$  .



## The sampling distribution of $\bar{X}_1 - \bar{X}_2$ :

**Result (1) :** The mean , variance and the standard deviation of  $\bar{X}_1 - \bar{X}_2$  are :

1. The mean of  $\bar{X}_1 - \bar{X}_2$  is :  $\mu_{\bar{X}_1 - \bar{X}_2} = \mu_1 - \mu_2$
2. The variance of  $\bar{X}_1 - \bar{X}_2$  is :  $\sigma_{\bar{X}_1 - \bar{X}_2}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$
3. The standard deviation of  $\bar{X}_1 - \bar{X}_2$  is called the standard error and is defined by :  $\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\sigma_{\bar{X}_1 - \bar{X}_2}^2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

• **Note:** Square roots distribute over multiplication or division, but not addition or subtraction.

$$\sqrt{a + b} \neq \sqrt{a} + \sqrt{b}$$

• **In general:**  $Z = \frac{\text{value} - \text{Mean}}{\text{Standard deviation}}$

## Result (2) :

If the two random samples were selected from normal distributions ( or non-normal distributions with large sample sizes ) with known variances  $\sigma_1^2$  and  $\sigma_2^2$  , then the difference between the sample means  $\bar{X}_1 - \bar{X}_2$  has a normal distribution with mean  $\mu_1 - \mu_2$  and variance  $\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$  , that is :

$$1. \quad \bar{X}_1 - \bar{X}_2 \sim \text{Normal} (\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$$

$$2. \quad Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim \text{Normal} (0,1)$$

We use this result when sampling :

1-from normal distribution with known variance  $\sigma_1^2, \sigma_2^2$  .

2-from non-normal distribution with known variance  $\sigma_1^2, \sigma_2^2$  ,and with large sample size ( $n_1, n_2 \geq 30$  ).



### Example:

Suppose it has been established that for a certain type of client (type A) the average length of a home visit by a public health nurse is 45 minutes with standard deviation of 15 minutes, and that for second type (type B) of client the average home visit is 30 minutes long with standard deviation of 20 minutes. If a nurse randomly visits 35 clients from the first type and 40 clients from the second type, what is the probability that the average length of home visit of first type will be greater than the average length of home visit of second type by 20 or more minutes ?

**Solution :**

Type	Mean (Average)	Standard deviation	Variance	Sample size
<b>A (First type)</b>	$\mu_1=45$	$\sigma_1=15$	$\sigma_1^2=225$	$n_1=35$ (large)
<b>B (Second type)</b>	$\mu_2=30$	$\sigma_2=20$	$\sigma_2^2=400$	$n_2=40$ (large)

**The mean, the variance and the standard deviation of  $\bar{X}_1 - \bar{X}_2$  are:**

1. The mean of  $\bar{X}_1 - \bar{X}_2$  is :  $\mu_{\bar{X}_1 - \bar{X}_2} = \mu_1 - \mu_2 = 45 - 30 = 15$
2. The variance of  $\bar{X}_1 - \bar{X}_2$  is :  $\sigma_{\bar{X}_1 - \bar{X}_2}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} = \frac{225}{35} + \frac{400}{40} = 16.4286$
3. The standard deviation of  $\bar{X}_1 - \bar{X}_2$  is called the standard error and is defined by :  $\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} = \sqrt{16.4286} = 4.0532$

The sampling distribution of  $\bar{X}_1 - \bar{X}_2$  is :

$$\bar{X}_1 - \bar{X}_2 \sim \text{Normal} (15 , 16.4286)$$

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - 15}{\sqrt{16.4286}} \sim \text{Normal} (0,1)$$

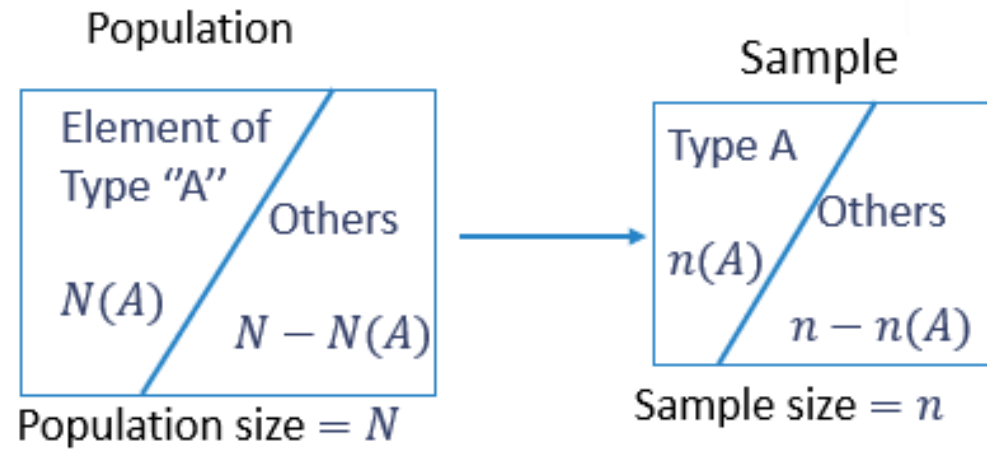
**Result :**  
non-normal distribution with  
known variance and with  
large sample size .

the probability that the average length of home visit of first type will be greater than the average length of home visit of second type by 20 or more minutes equals to :

$$\begin{aligned} P(\bar{X}_1 > \bar{X}_2 + 20) &= \\ &= P(\bar{X}_1 - \bar{X}_2 > 20) \\ &= P\left( Z > \frac{20-15}{4.0532} \right) \\ &= P( Z > 1.23) \\ &= 1 - P(Z < 1.23) \\ &= 1 - 0.89065 \\ &= 0.10935 \end{aligned}$$

## 5.5 Distribution of the Sample Proportion ( $\hat{P}$ ):

The sampling distribution of  $\hat{P}$  is used to make inferences about  $p$ .



- **For the population:**

$N(A)$  = number of elements in the population with a specified characteristic "A"

$N$  = total number of elements in the population (population size)

The population proportion is :

$$P = \frac{N(A)}{N} \quad (P \text{ is a parameter})$$

- **For the sample:**

$n(A)$  = number of elements in the sample with the same characteristic "A"

$n$  = sample size

The sample proportion is :

$$\hat{P} = \frac{n(A)}{n} \quad (\hat{P} \text{ is a statistic})$$

### Result (1):

- The mean of the sample proportion  $\hat{P}$  is the population proportion (P) , that is:  $\mu_{\hat{P}} = P$
- The variance of the sample proportion  $\hat{P}$  is :  $\sigma_{\hat{P}}^2 = \frac{P(1-P)}{n} = \frac{pq}{n}$  (where  $q=1-p$ )

The standard error (standard deviation) of the sample proportion  $\hat{P}$  is :  $\sigma_{\hat{P}} = \sqrt{\frac{pq}{n}}$

### Result (2):

- For large sample size (  $n \geq 30, np > 5, nq > 5$  ), the sample proportion  $\hat{P}$  has approximately a normal distribution with  $\mu_{\hat{P}} = P$  and  $\sigma_{\hat{P}}^2 = \frac{pq}{n}$  , that is:

$$\hat{P} \approx \text{Normal} \left( P, \frac{pq}{n} \right) \quad (\text{approximately})$$

$$Z = \frac{\hat{P} - P}{\sqrt{\frac{pq}{n}}} \sim \text{Normal} (0,1)$$

### Example:

Suppose that 45% of the patients visiting a certain clinic are females. If a sample of 35 patients was selected at random, find the probability that:

1. The proportion of females in the sample will be greater than 0.4.
2. The proportion of females in the sample will be between 0.4 and 0.5.

### Solution:

$n = 35$  ( large )

$p =$  The population proportion of females  $= 45\% = \frac{45}{100} = 0.45$  ( $q = 1 - p = 0.55$ )

$\hat{P} =$  The sample proportion (proportion of females in the sample)

- The mean of the sample proportion  $\hat{P}$  is :  $P=0.45$
- The variance of the sample proportion  $\hat{P}$  is :  $\frac{pq}{n} = \frac{(0.45)(0.55)}{35} = 0.0071$
- The standard error (standard deviation) of the sample proportion  $\hat{P}$  is :  $\sqrt{\frac{pq}{n}} = \sqrt{0.0071} = 0.084$

- $n = 35 \geq 30$ ,  $np = 35(0.45) = 15.75 > 5$ ,  $nq = 35(0.55) = 19.25 > 5$

The distribution of  $\hat{P}$  is:  $\hat{P} \approx \text{Normal } (P=0.45, \frac{pq}{n} = 0.0071)$

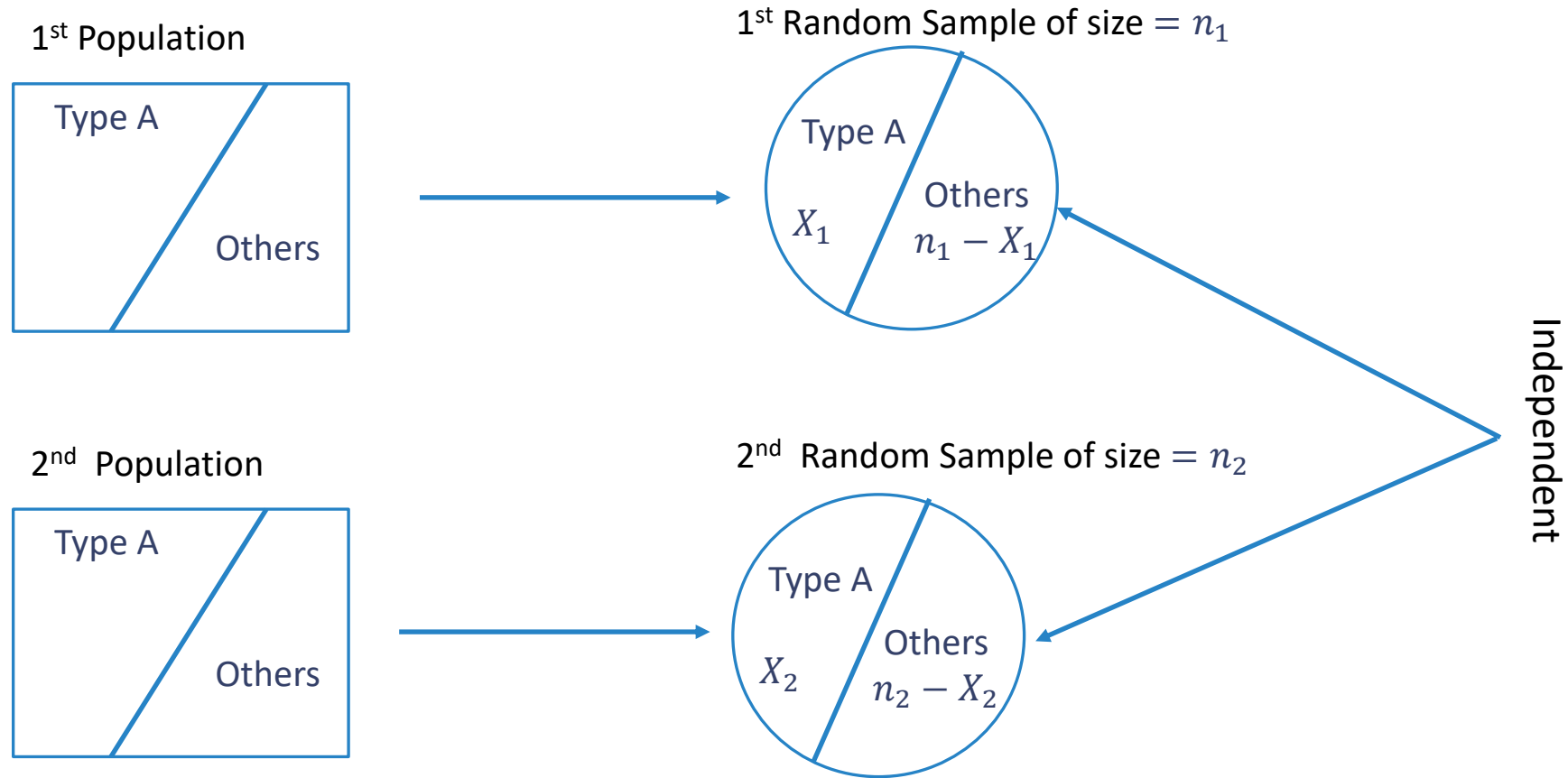
- The probability that the proportion of females in the sample will be greater than 0.4 :

$$\begin{aligned} P(\hat{P} > 0.4) &= P(Z > \frac{0.4-0.45}{\sqrt{0.0071}}) \\ &= P(Z > -0.59) = 1 - P(Z < -0.59) \\ &= 1 - 0.2776 \\ &= 0.7224 \end{aligned}$$

- The probability that the proportion of females in the sample will be between 0.4 and 0.5 :

$$\begin{aligned} P(0.4 < \hat{P} < 0.5) &= P(\frac{0.4-0.45}{\sqrt{0.0071}} < Z < \frac{0.5-0.45}{\sqrt{0.0071}}) \\ &= P(-0.59 < Z < 0.59) \\ &= P(Z < 0.59) - P(Z < -0.59) \\ &= 0.7224 - 0.2776 \\ &= 0.4448 \end{aligned}$$

## 5.6 Distributions of the difference between two sample proportions ( $\hat{P}_1 - \hat{P}_2$ ):



Suppose that we have two populations:

- $P_1$  = proportion of elements of type (A) in the 1-st population.
- $P_2$  = proportion of elements of type (A) in the 2-nd population.
- We are interested in comparing  $P_1$  and  $P_2$  , or equivalently, making inferences about  $P_1 - P_2$ .
- We independently select a random sample of size  $n_1$  from the 1-st population and another random sample of size  $n_2$  from the 2-nd population:
- Let  $X_1$  = no. of elements of type (A) in the 1-st sample.
- Let  $X_2$  = no. of elements of type (A) in the 2-nd sample.
- $\hat{P}_1 = \frac{X_1}{n_1}$  = sample proportion of the 1-st sample .
- $\hat{P}_2 = \frac{X_2}{n_2}$  = sample proportion of the 2-nd sample .
- The sampling distribution of  $\hat{P}_1 - \hat{P}_2$  is used to make inferences about  $P_1 - P_2$ .



## The sampling distribution of $\hat{P}_1 - \hat{P}_2$ :

### Result (1 ):

The mean, the variance and the standard error (standard deviation) of  $\hat{P}_1 - \hat{P}_2$  are :

1.The mean of  $\hat{P}_1 - \hat{P}_2$  :  $\mu_{\hat{P}_1 - \hat{P}_2} = P_1 - P_2$

2. The variance of  $\hat{P}_1 - \hat{P}_2$  is :  $\sigma_{\hat{P}_1 - \hat{P}_2}^2 = \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}$

3. The standard error (standard deviation) of  $\hat{P}_1 - \hat{P}_2$  is :  $\sigma_{\hat{P}_1 - \hat{P}_2} = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$

Where :  $q_1 = 1 - p_1$  and  $q_2 = 1 - p_2$  .

### Result (2) :

For large samples sizes ( $n_1 \geq 30, n_2 \geq 30, n_1 p_1 > 5, n_1 q_1 > 5, n_2 p_2 > 5, n_2 q_2 > 5$ ) , we have :

$$\hat{P}_1 - \hat{P}_2 \approx \text{Normal} \left( P_1 - P_2, \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2} \right) \quad (\text{approximately})$$

$$Z = \frac{(\hat{P}_1 - \hat{P}_2) - (P_1 - P_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} \sim \text{Normal} (0,1)$$

**Example :**

Suppose that 40% of Non-Saudi residents have medical insurance and 30% of Saudi residents have medical insurance in a certain city. We have randomly and independently selected a sample of 130 Non-Saudi residents and another sample of 120 Saudi residents. What is the probability that the difference between the sample proportions, will be between 0.05 and 0.2 ?

$P_1$  = population proportion of non-Saudi with medical insurance

$P_2$  = population proportion of Saudi with medical insurance

$\hat{P}_1$  = sample proportion of non-Saudi with medical insurance

$\hat{P}_2$  = sample proportion of Saudi with medical insurance

Type	proportion	q=1-p	Sample size
Non-Saudi residents	$P_1=0.40$	$q_1=0.60$	$n_1=130$ (large)
Saudi residents	$P_2=0.30$	$q_2=0.70$	$n_2=120$ (large)

The mean of  $\hat{P}_1 - \hat{P}_2$  :  $\mu_{\hat{P}_1 - \hat{P}_2} = P_1 - P_2 = 0.4 - 0.3 = 0.1$

The variance of  $\hat{P}_1 - \hat{P}_2$  :  $\sigma_{\hat{P}_1 - \hat{P}_2}^2 = \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2} = \frac{(0.4)(0.6)}{130} + \frac{(0.3)(0.7)}{120} = 0.0036$

The standard error (standard deviation) of  $\hat{P}_1 - \hat{P}_2$  :  $\sigma_{\hat{P}_1 - \hat{P}_2} = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}} = \sqrt{0.0036} = 0.0599$

The distribution of  $\hat{P}_1 - \hat{P}_2$  is:  $(n_1 \geq 30, n_2 \geq 30, n_1 p_1 > 5, n_1 q_1 > 5, n_2 p_2 > 5, n_2 q_2 > 5)$

$$\hat{P}_1 - \hat{P}_2 \approx \text{Normal}(0.1, 0.0036)$$

The probability that the difference between the sample proportions, will be between 0.05 and 0.2 is :

$$\begin{aligned} P(0.05 < \hat{P}_1 - \hat{P}_2 < 0.2) &= P\left( \frac{0.05 - (P_1 - P_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} < Z < \frac{0.2 - (P_1 - P_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} \right) \\ &= P\left( \frac{0.05 - 0.1}{0.0599} < Z < \frac{0.2 - 0.1}{0.0599} \right) \\ &= P(-0.83 < Z < 1.67) \\ &= P(Z < 1.67) - P(Z < -0.83) \\ &= 0.95254 - 0.20327 = 0.74927 \end{aligned}$$