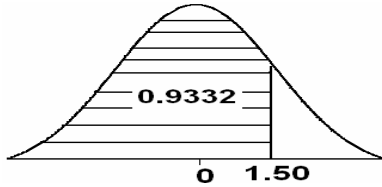(iv) $P(Z = a) = 0$ for every $a$.

**Example:**

Suppose that $Z \sim N(0,1)$

(1) $P(Z \le 1.50) = 0.9332$
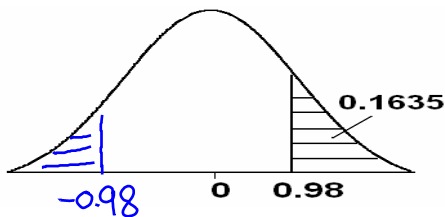


| Z | 0.00 | 0.01 | … |
|---|---|---|---|
| : | ⇓ | | |
| 1.50 ⇒ | 0.9332 | | |
| : | | | |

(2)

$P(Z \ge 0.98) = 1 - P(Z \le 0.98)$

$\qquad = 1 - 0.8365$

$\qquad = 0.1635$



| Z | 0.00 | … | 0.08 |
|---|---|---|---|
| : | : | : | ⇓ |
| : | | … | … | ⇓ |
| 0.90⇒ | ⇒ | ⇒ | 0.8365 |

or

P(Z>0.98)=P(Z< -0.98)= 0.1635

(3)

$P(-1.33 \le Z \le 2.42) =$

$P(Z \le 2.42) - P(Z \le -1.33)$

$\qquad = 0.9922 - 0.0918$

$\qquad = 0.9004$



| Z | … | | −0.03 |
|---|---|---|---|
| : | : | | ⇓ |
| −1.30 | ⇒ | | 0.0918 |
| : | | | |
| | | | |

(4) $P(Z \le 0) = P(Z \ge 0) = 0.5$

**Notation:**

$P(Z \le Z_A) = A$

0.5 | 0.5

0

$$P(Z < Z_A) = A$$

$Z \sim N(0,1)$

$P(Z<Z_A)$

A | 1- A

$Z_A$

For example:

$$P(Z < Z_{0.025}) = 0.025$$

0.025 | 0.975

$Z_{0.025}$

$$P(Z < Z_{0.90}) = 0.90$$

0.90 | 0.10

$Z_{0.90}$

**Result:**

Since the pdf of Z~N(0,1) is symmetric about 0, we have:

$$Z_A = - Z_{1-A}$$

For example:  
$$Z_{0.35} = - Z_{1-0.35} = - Z_{0.65}$$
$$Z_{0.86} = - Z_{1-0.86} = - Z_{0.14}$$

Z(0.975)= 1.96  
Z(0.025)= - Z(0.975)= -1.96



**Example:**

Suppose that $Z \sim N(0,1)$.

If $P(Z \le a) = 0.9505\,3$

Then $a = 1.65$

| Z | … | 0.05 | … |
|---|---|------|---|
| : | | ⇑ | |
| 1.60 | ⇐ | 0.9505 3 | |
| : | | | |



$P(Z < a) = 0.9505$

$P(Z < Z_{0.9505}) = 0.9505$

$a = Z_{0.9505}$

0.9505 | 0.0495

$a = Z_{0.9505} = 1.65$

a

**Example:**

Suppose that Z~N(0,1). Find the value of $k$ such that P(Z≤$k$)= 0.0207.

| Z | … | −0.04 | |
|---|---|---|---|
| : | : | ⇑ | |
| | | ⇑ | |
| −2.0 | ⇐⇐ | 0.0207 | |
| : | | | |

**Solution:**

.$k = -2.04$

Notice that $k = Z_{0.0207} = -2.04$



**Example:**

If $Z \sim N(0,1)$, then:

$Z_{0.90} = 1.285$   $Z_{0.90} = (Z_{0.89973} + Z_{0.90147})/2 = (1.28 + 1.29)/2 = 1.285$

$Z_{0.95} = 1.645$

$Z_{0.975} = 1.96$   $Z_{.95} = (Z_{0.94950} + Z_{0.95053})/2 = (1.64 + 1.65)/2 = 1.645$

$Z_{0.99} = 2.325$   $Z_{0.99} = (Z_{0.98983} + Z_{0.99010})/2 = (2.32 + 2.33)/2 = 2.325$



Using the result:  $Z_A = - Z_{1-A}$

$Z_{0.10} = - Z_{0.90} = - 1.285$

$Z_{0.05} = - Z_{0.95} = - 1.645$

$Z_{0.025} = - Z_{0.975} = - 1.96$

$Z_{0.01} = - Z_{0.99} = - 2.325$

## Calculating Probabilities of Normal $(\mu, \sigma^2)$:

- Recall the result:

$X \sim \text{Normal} (\mu, \sigma^2) \iff Z = \dfrac{X - \mu}{\sigma} \sim \text{Normal} (0,1)$

- $X \le a \iff \dfrac{X-\mu}{\sigma} \le \dfrac{a-\mu}{\sigma} \iff Z \le \dfrac{a-\mu}{\sigma}$

1. $P(X \le a) = P\left(Z \le \dfrac{a-\mu}{\sigma}\right) =$ From the table.

2. $P(X \ge a) = 1 - P(X \le a) = 1 - P\left(Z \le \dfrac{a-\mu}{\sigma}\right)$

3. $P(a \le X \le b) = P(X \le b) - P(X \le a)$

$$= P\left(Z \le \dfrac{b-\mu}{\sigma}\right) - P\left(Z \le \dfrac{a-\mu}{\sigma}\right)$$

4. $P(X = a) = 0$, for every $a$.

## 4.7 Normal Distribution Application:

**Example**

Suppose that the hemoglobin levels of healthy adult males are approximately normally distributed with a mean of 16 and a variance of 0.81.

(a) Find that probability that a randomly chosen healthy adult male has a hemoglobin level less than 14.

(b) What is the percentage of healthy adult males who have hemoglobin level less than 14?

(c) In a population of 10,000 healthy adult males, how many would you expect to have hemoglobin level less than 14?

**Solution:**

$X =$ hemoglobin level for healthy adults males

Mean: $\mu = 16$

Variance: $\sigma^2 = 0.81$

Standard deviation: $\sigma = 0.9$

$X \sim$ Normal (16, 0.81)

(a) The probability that a randomly chosen healthy adult male has hemoglobin level less than 14 is $P(X \le 14)$.

$$P(X \le 14) = P\left(Z \le \frac{14 - \mu}{\sigma}\right)$$

$$= P\left(Z \le \frac{14 - 16}{0.9}\right)$$

$$= P(Z \le -2.22)$$

$$= 0.0132 \text{[1]}$$



(b) The percentage of healthy adult males who have hemoglobin level less than 14 is:

$$P(X \le 14) \times 100\% = 0.0132 \times 100\% = 1.32\%$$

(c) In a population of 10000 healthy adult males, we would expect that the number of males with hemoglobin level less than 14 to be:

$$P(X \le 14) \times 10000 = 0.0132 \times 10000 = 132 \text{ males}$$

**Example:**
Suppose that the birth weight of Saudi babies has a normal distribution with mean $\mu=3.4$ and standard deviation $\sigma=0.35$.
(a) Find the probability that a randomly chosen Saudi baby has a birth weight between 3.0 and 4.0 kg.
(b) What is the percentage of Saudi babies who have a birth weight between 3.0 and 4.0 kg?
(c) In a population of 100000 Saudi babies, how many would you expect to have birth weight between 3.0 and 4.0 kg?

**Solution:**
X = birth weight of Saudi babies
Mean: $\mu = 3.4$
Standard deviation: $\sigma = 0.35$
Variance: $\sigma^2 = (0.35)^2 = 0.1225$
X ~ Normal $(3.4, 0.1225)$
(a) The probability that a randomly chosen Saudi baby has a birth weight between 3.0 and 4.0 kg is $P(3.0 < X < 4.0)$

$$P(3.0 < X < 4.0) = P(X \leq 4.0) - P(X \leq 3.0)$$

$$= P\left(Z \leq \frac{4.0 - \mu}{\sigma}\right) - P\left(Z \leq \frac{3.0 - \mu}{\sigma}\right)$$

$$= P\left(Z \leq \frac{4.0 - 3.4}{0.35}\right) - P\left(Z \leq \frac{3.0 - 3.4}{0.35}\right)$$

$$= P(Z \leq 1.71) - P(Z \leq -1.14)$$

$$= 0.9564 - 0.1271 = 0.8293$$



(b) The percentage of Saudi babies who have a birth weight between 3.0 and 4.0 kg is

P(3.0<X<4.0) × 100% = 0.8293× 100% = 82.93%

(c) In a population of 100,000 Saudi babies, we would expect that the number of babies with birth weight between 3.0 and 4.0 kg to be:

P(3.0<X<4.0) × 100000 = 0.8293× 100000 = 82930 babies

## Standard Normal Table
Areas Under the Standard Normal Curve

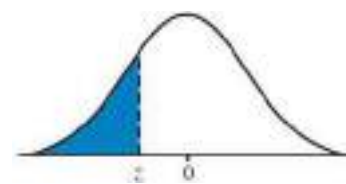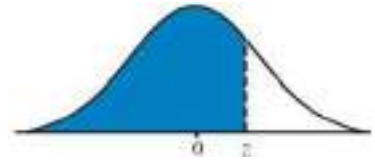| z | -0.09 | -0.08 | -0.07 | -0.06 | -0.05 | -0.04 | -0.03 | -0.02 | -0.01 | -0.00 | z |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **-3.50** | 0.00017 | 0.00017 | 0.00018 | 0.00019 | 0.00019 | 0.00020 | 0.00021 | 0.00022 | 0.00022 | 0.00023 | **-3.50** |
| **-3.40** | 0.00024 | 0.00025 | 0.00026 | 0.00027 | 0.00028 | 0.00029 | 0.00030 | 0.00031 | 0.00032 | 0.00034 | **-3.40** |
| **-3.30** | 0.00035 | 0.00036 | 0.00038 | 0.00039 | 0.00040 | 0.00042 | 0.00043 | 0.00045 | 0.00047 | 0.00048 | **-3.30** |
| **-3.20** | 0.00050 | 0.00052 | 0.00054 | 0.00056 | 0.00058 | 0.00060 | 0.00062 | 0.00064 | 0.00066 | 0.00069 | **-3.20** |
| **-3.10** | 0.00071 | 0.00074 | 0.00076 | 0.00079 | 0.00082 | 0.00084 | 0.00087 | 0.00090 | 0.00094 | 0.00097 | **-3.10** |
| **-3.00** | 0.00100 | 0.00104 | 0.00107 | 0.00111 | 0.00114 | 0.00118 | 0.00122 | 0.00126 | 0.00131 | 0.00135 | **-3.00** |
| **-2.90** | 0.00139 | 0.00144 | 0.00149 | 0.00154 | 0.00159 | 0.00164 | 0.00169 | 0.00175 | 0.00181 | 0.00187 | **-2.90** |
| **-2.80** | 0.00193 | 0.00199 | 0.00205 | 0.00212 | 0.00219 | 0.00226 | 0.00233 | 0.00240 | 0.00248 | 0.00256 | **-2.80** |
| **-2.70** | 0.00264 | 0.00272 | 0.00280 | 0.00289 | 0.00298 | 0.00307 | 0.00317 | 0.00326 | 0.00336 | 0.00347 | **-2.70** |
| **-2.60** | 0.00357 | 0.00368 | 0.00379 | 0.00391 | 0.00402 | 0.00415 | 0.00427 | 0.00440 | 0.00453 | 0.00466 | **-2.60** |
| **-2.50** | 0.00480 | 0.00494 | 0.00508 | 0.00523 | 0.00539 | 0.00554 | 0.00570 | 0.00587 | 0.00604 | 0.00621 | **-2.50** |
| **-2.40** | 0.00639 | 0.00657 | 0.00676 | 0.00695 | 0.00714 | 0.00734 | 0.00755 | 0.00776 | 0.00798 | 0.00820 | **-2.40** |
| **-2.30** | 0.00842 | 0.00866 | 0.00889 | 0.00914 | 0.00939 | 0.00964 | 0.00990 | 0.01017 | 0.01044 | 0.01072 | **-2.30** |
| **-2.20** | 0.01101 | 0.01130 | 0.01160 | 0.01191 | 0.01222 | 0.01255 | 0.01287 | 0.01321 | 0.01355 | 0.01390 | **-2.20** |
| **-2.10** | 0.01426 | 0.01463 | 0.01500 | 0.01539 | 0.01578 | 0.01618 | 0.01659 | 0.01700 | 0.01743 | 0.01786 | **-2.10** |
| **-2.00** | 0.01831 | 0.01876 | 0.01923 | 0.01970 | 0.02018 | 0.02068 | 0.02118 | 0.02169 | 0.02222 | 0.02275 | **-2.00** |
| **-1.90** | 0.02330 | 0.02385 | 0.02442 | 0.02500 | 0.02559 | 0.02619 | 0.02680 | 0.02743 | 0.02807 | 0.02872 | **-1.90** |
| **-1.80** | 0.02938 | 0.03005 | 0.03074 | 0.03144 | 0.03216 | 0.03288 | 0.03362 | 0.03438 | 0.03515 | 0.03593 | **-1.80** |
| **-1.70** | 0.03673 | 0.03754 | 0.03836 | 0.03920 | 0.04006 | 0.04093 | 0.04182 | 0.04272 | 0.04363 | 0.04457 | **-1.70** |
| **-1.60** | 0.04551 | 0.04648 | 0.04746 | 0.04846 | 0.04947 | 0.05050 | 0.05155 | 0.05262 | 0.05370 | 0.05480 | **-1.60** |
| **-1.50** | 0.05592 | 0.05705 | 0.05821 | 0.05938 | 0.06057 | 0.06178 | 0.06301 | 0.06426 | 0.06552 | 0.06681 | **-1.50** |
| **-1.40** | 0.06811 | 0.06944 | 0.07078 | 0.07215 | 0.07353 | 0.07493 | 0.07636 | 0.07780 | 0.07927 | 0.08076 | **-1.40** |
| **-1.30** | 0.08226 | 0.08379 | 0.08534 | 0.08691 | 0.08851 | 0.09012 | 0.09176 | 0.09342 | 0.09510 | 0.09680 | **-1.30** |
| **-1.20** | 0.09853 | 0.10027 | 0.10204 | 0.10383 | 0.10565 | 0.10749 | 0.10935 | 0.11123 | 0.11314 | 0.11507 | **-1.20** |
| **-1.10** | 0.11702 | 0.11900 | 0.12100 | 0.12302 | 0.12507 | 0.12714 | 0.12924 | 0.13136 | 0.13350 | 0.13567 | **-1.10** |
| **-1.00** | 0.13786 | 0.14007 | 0.14231 | 0.14457 | 0.14686 | 0.14917 | 0.15151 | 0.15386 | 0.15625 | 0.15866 | **-1.00** |
| **-0.90** | 0.16109 | 0.16354 | 0.16602 | 0.16853 | 0.17106 | 0.17361 | 0.17619 | 0.17879 | 0.18141 | 0.18406 | **-0.90** |
| **-0.80** | 0.18673 | 0.18943 | 0.19215 | 0.19489 | 0.19766 | 0.20045 | 0.20327 | 0.20611 | 0.20897 | 0.21186 | **-0.80** |
| **-0.70** | 0.21476 | 0.21770 | 0.22065 | 0.22363 | 0.22663 | 0.22965 | 0.23270 | 0.23576 | 0.23885 | 0.24196 | **-0.70** |
| **-0.60** | 0.24510 | 0.24825 | 0.25143 | 0.25463 | 0.25785 | 0.26109 | 0.26435 | 0.26763 | 0.27093 | 0.27425 | **-0.60** |
| **-0.50** | 0.27760 | 0.28096 | 0.28434 | 0.28774 | 0.29116 | 0.29460 | 0.29806 | 0.30153 | 0.30503 | 0.30854 | **-0.50** |
| **-0.40** | 0.31207 | 0.31561 | 0.31918 | 0.32276 | 0.32636 | 0.32997 | 0.33360 | 0.33724 | 0.3409 | 0.34458 | **-0.40** |
| **-0.30** | 0.34827 | 0.35197 | 0.35569 | 0.35942 | 0.36317 | 0.36693 | 0.37070 | 0.37448 | 0.37828 | 0.38209 | **-0.30** |
| **-0.20** | 0.38591 | 0.38974 | 0.39358 | 0.39743 | 0.40129 | 0.40517 | 0.40905 | 0.41294 | 0.41683 | 0.42074 | **-0.20** |
| **-0.10** | 0.42465 | 0.42858 | 0.43251 | 0.43644 | 0.44038 | 0.44433 | 0.44828 | 0.45224 | 0.45620 | 0.46017 | **-0.10** |
| **-0.00** | 0.46414 | 0.46812 | 0.47210 | 0.47608 | 0.48006 | 0.48405 | 0.48803 | 0.49202 | 0.49601 | 0.50000 | **-0.00** |

## Standard Normal Table (continued)
Areas Under the Standard Normal Curve

| z | 0.00 | 0.01 | 0.02 | 0.03 | 0.04 | 0.05 | 0.06 | 0.07 | 0.08 | 0.09 | z |
|------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|------|
| 0.00 | 0.50000 | 0.50399 | 0.50798 | 0.51197 | 0.51595 | 0.51994 | 0.52392 | 0.52790 | 0.53188 | 0.53586 | 0.00 |
| 0.10 | 0.53983 | 0.54380 | 0.54776 | 0.55172 | 0.55567 | 0.55962 | 0.56356 | 0.56749 | 0.57142 | 0.57535 | 0.10 |
| 0.20 | 0.57926 | 0.58317 | 0.58706 | 0.59095 | 0.59483 | 0.59871 | 0.60257 | 0.60642 | 0.61026 | 0.61409 | 0.20 |
| 0.30 | 0.61791 | 0.62172 | 0.62552 | 0.62930 | 0.63307 | 0.63683 | 0.64058 | 0.64431 | 0.64803 | 0.65173 | 0.30 |
| 0.40 | 0.65542 | 0.65910 | 0.66276 | 0.66640 | 0.67003 | 0.67364 | 0.67724 | 0.68082 | 0.68439 | 0.68793 | 0.40 |
| 0.50 | 0.69146 | 0.69497 | 0.69847 | 0.70194 | 0.70540 | 0.70884 | 0.71226 | 0.71566 | 0.71904 | 0.72240 | 0.50 |
| 0.60 | 0.72575 | 0.72907 | 0.73237 | 0.73565 | 0.73891 | 0.74215 | 0.74537 | 0.74857 | 0.75175 | 0.75490 | 0.60 |
| 0.70 | 0.75804 | 0.76115 | 0.76424 | 0.76730 | 0.77035 | 0.77337 | 0.77637 | 0.77935 | 0.78230 | 0.78524 | 0.70 |
| 0.80 | 0.78814 | 0.79103 | 0.79389 | 0.79673 | 0.79955 | 0.80234 | 0.80511 | 0.80785 | 0.81057 | 0.81327 | 0.80 |
| 0.90 | 0.81594 | 0.81859 | 0.82121 | 0.82381 | 0.82639 | 0.82894 | 0.83147 | 0.83398 | 0.83646 | 0.83891 | 0.90 |
| 1.00 | 0.84134 | 0.84375 | 0.84614 | 0.84849 | 0.85083 | 0.85314 | 0.85543 | 0.85769 | 0.85993 | 0.86214 | 1.00 |
| 1.10 | 0.86433 | 0.86650 | 0.86864 | 0.87076 | 0.87286 | 0.87493 | 0.87698 | 0.87900 | 0.88100 | 0.88298 | 1.10 |
| 1.20 | 0.88493 | 0.88686 | 0.88877 | 0.89065 | 0.89251 | 0.89435 | 0.89617 | 0.89796 | 0.89973 | 0.90147 | 1.20 |
| 1.30 | 0.90320 | 0.90490 | 0.90658 | 0.90824 | 0.90988 | 0.91149 | 0.91309 | 0.91466 | 0.91621 | 0.91774 | 1.30 |
| 1.40 | 0.91924 | 0.92073 | 0.92220 | 0.92364 | 0.92507 | 0.92647 | 0.92785 | 0.92922 | 0.93056 | 0.93189 | 1.40 |
| 1.50 | 0.93319 | 0.93448 | 0.93574 | 0.93699 | 0.93822 | 0.93943 | 0.94062 | 0.94179 | 0.94295 | 0.94408 | 1.50 |
| 1.60 | 0.94520 | 0.94630 | 0.94738 | 0.94845 | 0.94950 | 0.95053 | 0.95154 | 0.95254 | 0.95352 | 0.95449 | 1.60 |
| 1.70 | 0.95543 | 0.95637 | 0.95728 | 0.95818 | 0.95907 | 0.95994 | 0.96080 | 0.96164 | 0.96246 | 0.96327 | 1.70 |
| 1.80 | 0.96407 | 0.96485 | 0.96562 | 0.96638 | 0.96712 | 0.96784 | 0.96856 | 0.96926 | 0.96995 | 0.97062 | 1.80 |
| 1.90 | 0.97128 | 0.97193 | 0.97257 | 0.97320 | 0.97381 | 0.97441 | 0.97500 | 0.97558 | 0.97615 | 0.97670 | 1.90 |
| 2.00 | 0.97725 | 0.97778 | 0.97831 | 0.97882 | 0.97932 | 0.97982 | 0.98030 | 0.98077 | 0.98124 | 0.98169 | 2.00 |
| 2.10 | 0.98214 | 0.98257 | 0.98300 | 0.98341 | 0.98382 | 0.98422 | 0.98461 | 0.98500 | 0.98537 | 0.98574 | 2.10 |
| 2.20 | 0.98610 | 0.98645 | 0.98679 | 0.98713 | 0.98745 | 0.98778 | 0.98809 | 0.98840 | 0.98870 | 0.98899 | 2.20 |
| 2.30 | 0.98928 | 0.98956 | 0.98983 | 0.99010 | 0.99036 | 0.99061 | 0.99086 | 0.99111 | 0.99134 | 0.99158 | 2.30 |
| 2.40 | 0.99180 | 0.99202 | 0.99224 | 0.99245 | 0.99266 | 0.99286 | 0.99305 | 0.99324 | 0.99343 | 0.99361 | 2.40 |
| 2.50 | 0.99379 | 0.99396 | 0.99413 | 0.99430 | 0.99446 | 0.99461 | 0.99477 | 0.99492 | 0.99506 | 0.99520 | 2.50 |
| 2.60 | 0.99534 | 0.99547 | 0.99560 | 0.99573 | 0.99585 | 0.99598 | 0.99609 | 0.99621 | 0.99632 | 0.99643 | 2.60 |
| 2.70 | 0.99653 | 0.99664 | 0.99674 | 0.99683 | 0.99693 | 0.99702 | 0.99711 | 0.99720 | 0.99728 | 0.99736 | 2.70 |
| 2.80 | 0.99744 | 0.99752 | 0.99760 | 0.99767 | 0.99774 | 0.99781 | 0.99788 | 0.99795 | 0.99801 | 0.99807 | 2.80 |
| 2.90 | 0.99813 | 0.99819 | 0.99825 | 0.99831 | 0.99836 | 0.99841 | 0.99846 | 0.99851 | 0.99856 | 0.99861 | 2.90 |
| 3.00 | 0.99865 | 0.99869 | 0.99874 | 0.99878 | 0.99882 | 0.99886 | 0.99889 | 0.99893 | 0.99896 | 0.9990 | 3.00 |
| 3.10 | 0.99903 | 0.99906 | 0.99910 | 0.99913 | 0.99916 | 0.99918 | 0.99921 | 0.99924 | 0.99926 | 0.99929 | 3.10 |
| 3.20 | 0.99931 | 0.99934 | 0.99936 | 0.99938 | 0.99940 | 0.99942 | 0.99944 | 0.99946 | 0.99948 | 0.99950 | 3.20 |
| 3.30 | 0.99952 | 0.99953 | 0.99955 | 0.99957 | 0.99958 | 0.99960 | 0.99961 | 0.99962 | 0.99964 | 0.99965 | 3.30 |
| 3.40 | 0.99966 | 0.99968 | 0.99969 | 0.99970 | 0.99971 | 0.99972 | 0.99973 | 0.99974 | 0.99975 | 0.99976 | 3.40 |
| 3.50 | 0.99977 | 0.99978 | 0.99978 | 0.99979 | 0.99980 | 0.99981 | 0.99981 | 0.99982 | 0.99983 | 0.99983 | 3.50 |

## <span style="color:red">CHAPTER 5: Probabilistic Features of the Distributions of Certain Sample Statistics</span>

## 5.1 Introduction:

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

As we mention earlier, a statistic is measure computed form 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.

## 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:
## (Sampling Distribution of the Sample Mean $\overline{X}$ ):

Suppose that we have a population with mean $\mu$ and variance $\sigma^2$. Suppose that $X_1, X_2, ..., X_n$ is a random sample of size ($n$) selected randomly from this population. We know that the sample mean is:

$$\overline{X} = \frac{\sum_{i=1}^{n} X_i}{n}.$$

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

| | 1st sample | 2nd sample | 3rd sample | … | Last sample |
|---|---|---|---|---|---|
| Sample values | 28 30 34 34 17 | 31 20 31 40 28 | 14 31 25 27 32 | . . . . | 17 32 29 31 30 |
| Sample mean $\overline{x}$ | 28.4 | 29.9 | 25.8 | … | 27.8 |

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

**Some Results about Sampling Distribution of $\overline{X}$:**

**Result (1): (mean & variance of $\overline{X}$)**

If $X_1, X_2, \ldots, X_n$ is a random sample of size $n$ from any distribution with mean $\mu$ and variance $\sigma^2$; then:

1. The mean of $\overline{X}$ is: $\mu_{\overline{X}} = \mu$.

2. The variance of $\overline{X}$ is: $\sigma_{\overline{X}}^2 = \dfrac{\sigma^2}{n}$.

3. The Standard deviation of $\overline{X}$ is call the standard error and is defined by: $\sigma_{\overline{X}} = \sqrt{\sigma_{\overline{X}}^2} = \dfrac{\sigma}{\sqrt{n}}$.

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

If $X_1, X_2, \ldots, X_n$ is a random sample of size $n$ from a normal population with mean $\mu$ and variance $\sigma^2$; that is Normal$(\mu, \sigma^2)$, then the sample mean has a normal distribution with mean $\mu$ and variance $\sigma^2/n$, that is:

1. $\overline{X} \sim$ Normal $\left( \mu, \dfrac{\sigma^2}{n} \right)$.

2. $Z = \dfrac{\overline{X} - \mu}{\sigma/\sqrt{n}} \sim$ Normal $(0,1)$.

We use this result when sampling from <mark>normal distribution</mark> with <mark>known variance $\sigma^2$</mark>.

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

Suppose that $X_1, X_2, \ldots, X_n$ is a random sample of size $n$ from non-normal population with mean $\mu$ and variance $\sigma^2$. If the sample size $n$ is large $(n \geq 30)$, then the sample mean has <mark>approximately</mark> a normal distribution with mean $\mu$ and variance $\sigma^2/n$, that is

$$1. \quad \overline{X} \approx \text{Normal}\left(\mu, \frac{\sigma^2}{n}\right) \qquad \text{(approximately)}$$

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

Note: "$\approx$" means "approximately distributed".
We use this result when sampling from <mark>non-normal distribution</mark> with <mark>known variance $\sigma^2$</mark> and with <mark>large sample size.</mark>

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

If $X_1, X_2, \ldots, X_n$ is a random sample of size $n$ from a normal distribution with mean $\mu$ and unknown variance $\sigma^2$; that is $\text{Normal}(\mu, \sigma^2)$, then the statistic:

$$T = \frac{\overline{X} - \mu}{S/\sqrt{n}}$$

has a t- distribution with $(n-1)$ degrees of freedom, where S is the sample standard deviation given by:

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

We write:

$$T = \frac{\overline{X} - \mu}{S/\sqrt{n}} \sim t(n-1)$$

<mark>Notation:</mark> degrees of freedom = df = $\nu$

**The t-Distribution:** (Section 6.3. pp 172-174)

- Student's t distribution.
- t-distribution is a distribution of a continuous random variable.

Result 2:
- Recall that, if $X_1$, $X_2$, ..., $X_n$ is a random sample of size $n$ from a normal distribution with mean $\mu$ and variance $\sigma^2$, i.e. $N(\mu, \sigma^2)$, then

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

We can apply this result only when $\sigma^2$ is known!

- If $\sigma^2$ is unknown, we replace the population variance $\sigma^2$ with the sample variance $S^2 = \dfrac{\sum\limits_{i=1}^{n}(X_i - \overline{X})^2}{n-1}$ to have the following statistic

$$T = \frac{\overline{X} - \mu}{S/\sqrt{n}}$$

**Recall:**

If $X_1$, $X_2$, ..., $X_n$ is a random sample of size $n$ from a normal distribution with mean $\mu$ and variance $\sigma^2$ is unknown, i.e. $N(\mu, \sigma^2)$, then the statistic:
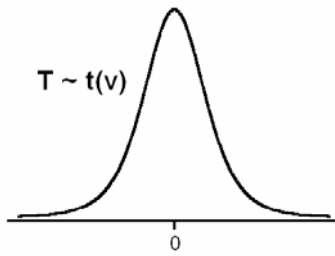
$$T = \frac{\overline{X} - \mu}{S/\sqrt{n}}$$

has a t-distribution with $(n-1)$ degrees of freedom ($df = \ = n-1$), and we write T~ t(ν) or T~ t(n–1).
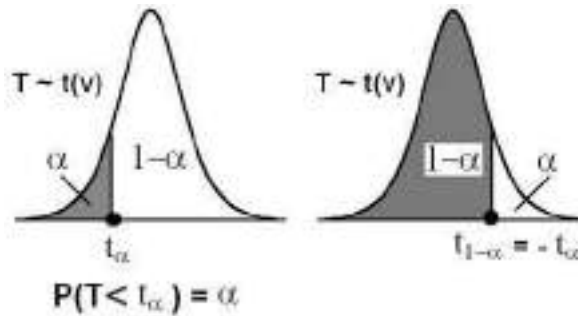
**Note:**

- t-distribution is a continuous distribution.
- The value of t random variable range from $-\infty$ to $+\infty$ (that is, $-\infty < t < \infty$).
- The mean of t distribution is 0.
- It is symmetric about the mean 0.
- The shape of t-distribution is similar to the shape of the standard normal distribution.
- t-distribution $\rightarrow$ Standard normal distribution as $n \rightarrow \infty$.

  i.e. If (n) go to infinity , the t distribution approximately normal distribution

88

$$T \sim t(v)$$

**Notation: (t $_\alpha$)**



$$T \sim t(v) \qquad T \sim t(v)$$

$$\alpha \quad 1-\alpha \qquad 1-\alpha \quad \alpha$$

$$t_\alpha \qquad t_{1-\alpha} = -t_\alpha$$

$$P(T < t_\alpha) = \alpha$$

- $t_\alpha$ = The t-value under which we find an area equal to $\alpha$
  = The t-value that leaves an area of $\alpha$ to the left.
- The value $t_\alpha$ satisfies: $P(T < t_\alpha) = \alpha$.
- Since the curve of the pdf of T~ t(v) is symmetric about 0, we have

$$t_{1-\alpha} = -t_\alpha$$

For example: $\quad t_{0.1} = -t_{1-0.1} = -t_{0.9}$

$$t_{0.975} = -t_{1-0.975} = -t_{0.025}$$

- Values of $t_\alpha$ are tabulated in a special table for several values of $\alpha$ and several values of degrees of freedom. (Table E, appendix p. A-40 in the textbook).

**Example:**
Find the t-value with $v=14$ (df) that leaves an area of:
(a)   0.95 to the left.
(b)   0.95 to the right.
**Solution:**
$v = 14$   (df);  T~ t(14)
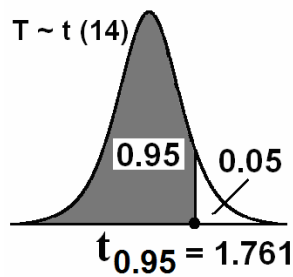(a) The t-value that leaves an area of 0.95 to the left is
$t_{0.95} = 1.761$.

$$t_{0.95} = 1.761 \qquad t_{0.95} = 1.761$$

(b) The t-value that leaves an area of 0.95 to the right is
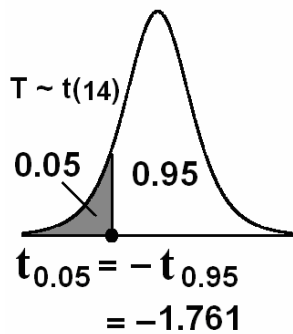$$t_{0.05} = -t_{1-0.05} = -t_{0.95} = -1.761$$



$$t_{0.05} = -t_{0.95} \qquad t_{0.05} = -1.761$$
$$= -1.761$$

**Note:** Some t-tables contain values of $\alpha$ that are greater than or equal to 0.90. When we search for small values of $\alpha$ in these tables, we may use the fact that:
$$t_{1-\alpha} = -t_{\alpha}$$

**Example:**
For $v = 10$ degrees of freedom (df), find $t_{0.93}$ and $t_{0.07}$.

**Solution:**
$t_{0.93} = (1.372+1.812)/2 = 1.592$ (from the table)
$t_{0.07} = -t_{1-0.07} = -t_{0.93} = -1.592$ (using the rule: $t_{1-\alpha} = -t_{\alpha}$)



$$t_{0.93} = \frac{1.372+1.812}{2}$$
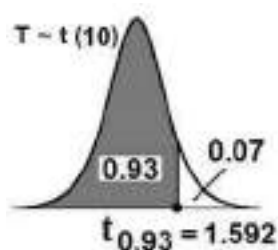$$= 1.592$$

90

# The t-Distribution

**Find :**

The t-value that leaves an area of 0.975 to the <mark>left</mark> (use $v = 12$) is

$t_{0.975} = 2.179$

The t-value that leaves an area of 0.90 to the <mark>right</mark> (use $v = 16$) is

$t_{0.10} = -t_{1-0.10} = -t_{0.90} = -1.337$

The t-value that leaves an area of 0.025 to the <mark>right</mark> $(use\ v = 8)$ is

$t_{0.975} = 2.306$

The t-value that leaves an area of 0.025 to the <mark>left</mark> $(use\ v = 8)$ is

$t_{0.025} = -t_{1-0.025} = -t_{0.975} = -2.306$

The t-value that leaves an area of 0.93 to the <mark>left</mark> $(use\ v = 10)$ is

$t_{0.93} = \dfrac{t_{0.90} + t_{0.95}}{2} = \dfrac{1.372 + 1.812}{2} = 1.592$

The t-value that leaves an area of 0.07 to the <mark>left</mark> $(use\ v = 10)$ is

$t_{0.07} = -t_{0.93} = -\left(\dfrac{t_{0.90} + t_{0.95}}{2}\right) = -1.592$

$P(T < K) = 0.90 \quad , df = 10$
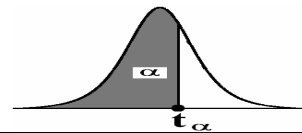
$K = 1.372$

$P(T \geq K) = 0.95 \quad , df = 15$

$K = -1.753$

$P(T < 2.110) = ? \quad (df = 17)$

$P(T < 2.110) = 0.975$

$P(T \leq 2.718) = ? \quad (df = 11) \quad P(T \leq 2.718) = 0.99$

*Critical Values of the t-distribution ($t_\alpha$)*



| ν=df | $t_{0.90}$ | $t_{0.95}$ | $t_{0.975}$ | $t_{0.99}$ | $t_{0.995}$ |
|---|---|---|---|---|---|
| 1 | 3.078 | 6.314 | 12.706 | 31.821 | 63.657 |
| 2 | 1.886 | 2.920 | 4.303 | 6.965 | 9.925 |
| 3 | 1.638 | 2.353 | 3.182 | 4.541 | 5.841 |
| 4 | 1.533 | 2.132 | 2.776 | 3.747 | 4.604 |
| 5 | 1.476 | 2.015 | 2.571 | 3.365 | 4.032 |
| 6 | 1.440 | 1.943 | 2.447 | 3.143 | 3.707 |
| 7 | 1.415 | 1.895 | 2.365 | 2.998 | 3.499 |
| 8 | 1.397 | 1.860 | 2.306 | 2.896 | 3.355 |
| 9 | 1.383 | 1.833 | 2.262 | 2.821 | 3.250 |
| 10 | 1.372 | 1.812 | 2.228 | 2.764 | 3.169 |
| 11 | 1.363 | 1.796 | 2.201 | 2.718 | 3.106 |
| 12 | 1.356 | 1.782 | 2.179 | 2.681 | 3.055 |
| 13 | 1.350 | 1.771 | 2.160 | 2.650 | 3.012 |
| 14 | 1.345 | 1.761 | 2.145 | 2.624 | 2.977 |
| 15 | 1.341 | 1.753 | 2.131 | 2.602 | 2.947 |
| 16 | 1.337 | 1.746 | 2.120 | 2.583 | 2.921 |
| 17 | 1.333 | 1.740 | 2.110 | 2.567 | 2.898 |
| 18 | 1.330 | 1.734 | 2.101 | 2.552 | 2.878 |
| 19 | 1.328 | 1.729 | 2.093 | 2.539 | 2.861 |
| 20 | 1.325 | 1.725 | 2.086 | 2.528 | 2.845 |
| 21 | 1.323 | 1.721 | 2.080 | 2.518 | 2.831 |
| 22 | 1.321 | 1.717 | 2.074 | 2.508 | 2.819 |
| 23 | 1.319 | 1.714 | 2.069 | 2.500 | 2.807 |
| 24 | 1.318 | 1.711 | 2.064 | 2.492 | 2.797 |
| 25 | 1.316 | 1.708 | 2.060 | 2.485 | 2.787 |
| 26 | 1.315 | 1.706 | 2.056 | 2.479 | 2.779 |
| 27 | 1.314 | 1.703 | 2.052 | 2.473 | 2.771 |
| 28 | 1.313 | 1.701 | 2.048 | 2.467 | 2.763 |
| 29 | 1.311 | 1.699 | 2.045 | 2.462 | 2.756 |
| 30 | 1.310 | 1.697 | 2.042 | 2.457 | 2.750 |
| 35 | 1.3062 | 1.6896 | 2.0301 | 2.4377 | 2.7238 |
| 40 | 1.3030 | 1.6840 | 2.0210 | 2.4230 | 2.7040 |
| 45 | 1.3006 | 1.6794 | 2.0141 | 2.4121 | 2.6896 |
| 50 | 1.2987 | 1.6759 | 2.0086 | 2.4033 | 2.6778 |
| 60 | 1.2958 | 1.6706 | 2.0003 | 2.3901 | 2.6603 |
| 70 | 1.2938 | 1.6669 | 1.9944 | 2.3808 | 2.6479 |
| 80 | 1.2922 | 1.6641 | 1.9901 | 2.3739 | 2.6387 |
| 90 | 1.2910 | 1.6620 | 1.9867 | 2.3685 | 2.6316 |
| 100 | 1.2901 | 1.6602 | 1.9840 | 2.3642 | 2.6259 |
| 120 | 1.2886 | 1.6577 | 1.9799 | 2.3578 | 2.6174 |
| 140 | 1.2876 | 1.6558 | 1.9771 | 2.3533 | 2.6114 |
| 160 | 1.2869 | 1.6544 | 1.9749 | 2.3499 | 2.6069 |
| 180 | 1.2863 | 1.6534 | 1.9732 | 2.3472 | 2.6034 |
| 200 | 1.2858 | 1.6525 | 1.9719 | 2.3451 | 2.6006 |
| ∞ | 1.282 | 1.645 | 1.960 | 2.326 | 2.576 |

**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 , $\sigma^2 = 1600$
X~N(800, 1600)
Sample size: $n$=16
Calculating mean, variance, and standard error (standard deviation) of the sample mean $\bar{X}$ :

Mean of $\bar{X}$ : $\qquad \mu_{\bar{X}} = \mu = 800$

Variance of $\bar{X}$ : $\qquad \sigma_{\bar{X}}^2 = \dfrac{\sigma^2}{n} = \dfrac{1600}{16} = 100$

Standard error (standard deviation) of $\bar{X}$ : $\sigma_{\bar{X}} = \dfrac{\sigma}{\sqrt{n}} = \dfrac{40}{\sqrt{16}} = 10$

~~Using result 2~~
~~Using the central limit theorem,~~ $\bar{X}$ has a normal distribution with mean $\mu_{\bar{X}} = 800$ and variance $\sigma_{\bar{X}}^2 = 100$ , that is:
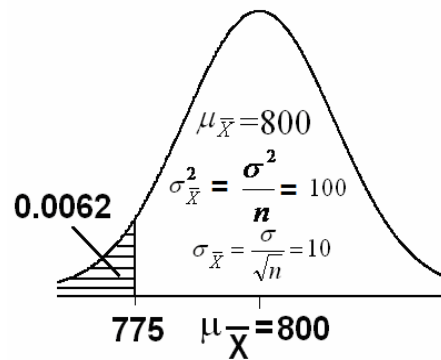
$$\bar{X} \sim N(\mu, \dfrac{\sigma^2}{n}) = N(800,100)$$

$$\Leftrightarrow Z = \dfrac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \dfrac{\bar{X} - 800}{10} \sim N(0,1)$$

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

$$P(\bar{X} < 775) = P\left( \dfrac{\bar{X} - \mu}{\sigma/\sqrt{n}} < \dfrac{775 - \mu}{\sigma/\sqrt{n}} \right) = P\left( \dfrac{\bar{X} - 800}{10} < \dfrac{775 - 800}{10} \right)$$

$$= P\left( Z < \dfrac{775 - 800}{10} \right) = P(Z < -2.50) = 0.0062$$

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



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

$\sigma_{\overline{X}}^2 = \dfrac{\sigma^2}{n} = 100$

0.0062

$\sigma_{\overline{X}} = \dfrac{\sigma}{\sqrt{n}} = 10$

775    $\mu_{\overline{X}} = 800$

**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

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

$X \approx N(120, 225)$

Sample size: $n=50$

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

Mean of $\overline{X}$ :      $\mu_{\overline{X}} = \mu = 120$

Variance of $\overline{X}$ :    $\sigma_{\overline{X}}^2 = \dfrac{\sigma^2}{n} = \dfrac{225}{50} = 4.5$

Standard error (standard deviation) of $\overline{X}$ : $\sigma_{\overline{X}} = \dfrac{\sigma}{\sqrt{n}} = \dfrac{15}{\sqrt{50}} = 2.12$

Using the central limit theorem, $\overline{X}$ has a normal distribution with mean $\mu_{\overline{X}} = 120$ and variance $\sigma_{\overline{X}}^2 = 4.5$, that is:

$$\overline{X} \sim N(\mu, \frac{\sigma^2}{n}) = N(120, 4.5)$$

$$\Leftrightarrow Z = \frac{\overline{X} - \mu}{\sigma/\sqrt{n}} = \frac{\overline{X} - 120}{2.12} \sim N(0,1)$$

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

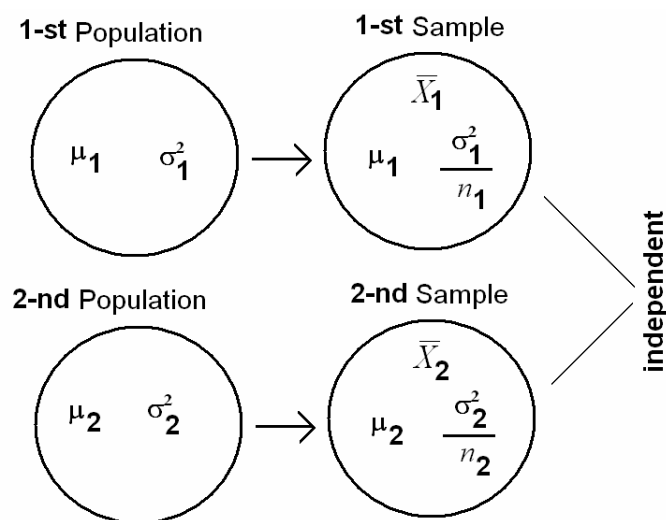$$P(115 < \overline{X} < 125) = P\left( \frac{115 - \mu}{\sigma/\sqrt{n}} < \frac{\overline{X} - \mu}{\sigma/\sqrt{n}} < \frac{125 - \mu}{\sigma/\sqrt{n}} \right)$$

$$= P\left(\frac{115-120}{2.12} < \frac{\overline{X} - \mu}{\sigma/\sqrt{n}} < \frac{125-120}{2.12}\right) = P\left(-2.36 < Z < 2.36\right)$$

$$= P\left(Z < 2.36\right) - P\left(Z < -2.36\right)$$

$$= 0.9909 - 0.0091$$

$$= 0.9818$$

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

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

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

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

In general: Z= (value - Mean)/ Standard deviation

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

### Result:

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

Mean of $\overline{X}_1 - \overline{X}_2$ is: $\qquad \mu_{\overline{X}_1 - \overline{X}_2} = \mu_1 - \mu_2$

Variance of $\overline{X}_1 - \overline{X}_2$ is: $\qquad \sigma^2_{\overline{X}_1 - \overline{X}_2} = \dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}$

Standard error (standard) deviation of $\overline{X}_1 - \overline{X}_2$ is:

$$\sigma_{\overline{X}_1 - \overline{X}_2} = \sqrt{\sigma^2_{\overline{X}_1 - \overline{X}_2}} = \sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}$$

### Result:

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 $(\overline{X}_1 - \overline{X}_2)$ has a normal distribution with mean $(\mu_1 - \mu_2)$ and variance $((\sigma_1^2 / n_1) + (\sigma_2^2 / n_2))$, that is:

- $\overline{X}_1 - \overline{X}_2 \sim N\left(\mu_1 - \mu_2 , \dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}\right)$

- $Z = \dfrac{(\overline{X}_1 - \overline{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}} \sim N(0,1)$

### Application:

### 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:**

$$\overline{X}_1 > \overline{X}_2 + 20$$

For the first type:

$\mu_1 = 45$

$\sigma_1 = 15$

$\sigma_1^2 = 225$

$n_1 = 35$ is large

For the second type:

$\mu_2 = 30$

$\sigma_2 = 20$

$\sigma_2^2 = 400$

$n_2 = 40$ is large

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

Mean of $\overline{X}_1 - \overline{X}_2$ is:

$$\mu_{\overline{X}_1-\overline{X}_2} = \mu_1 - \mu_2 = 45 - 30 = 15$$

Variance of $\overline{X}_1 - \overline{X}_2$ is:

$$\sigma_{\overline{X}_1-\overline{X}_2}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} = \frac{225}{35} + \frac{400}{40} = 16.4286$$

Standard error (standard) deviation of $\overline{X}_1 - \overline{X}_2$ is:

$$\sigma_{\overline{X}_1-\overline{X}_2} = \sqrt{\sigma_{\overline{X}_1-\overline{X}_2}^2} = \sqrt{16.4286} = 4.0532$$

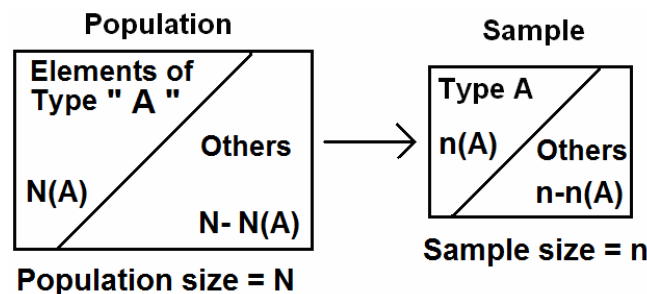The sampling distribution of $\overline{X}_1 - \overline{X}_2$ is:

$$\overline{X}_1 - \overline{X}_2 \sim N(15, 16.4286)$$

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

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 is:

$$P(\bar{X}_1 - \bar{X}_2 > 20) = P\left( \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}} > \frac{20 - (\mu_1 - \mu_2)}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}} \right)$$

$$= P\left( Z > \frac{20 - 15}{4.0532} \right) = P(Z>1.23) = 1 - P(Z<1.23)$$

$$= 1 - 0.8907$$

$$= 0.1093$$

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



- 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

   $$p = \frac{N(A)}{N}$$     (p 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

   $$\hat{p} = \frac{n(A)}{n}$$     ( $\hat{p}$ is a statistic)

- The sampling distribution of $\hat{p}$ is used to make inferences

97

about p.

**Result:**

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:

$$\sigma_{\hat{p}}^2 = \frac{p(1-p)}{n} = \frac{pq}{n}. \qquad \text{(where q=1 -p)}$$

The standard error (standard deviation) of the sample proportion ($\hat{p}$) is:

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}} = \sqrt{\frac{pq}{n}}$$

**Result:**

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

$$\hat{p} \sim N\left(p, \frac{pq}{n}\right) \qquad \text{(approximately)}$$

$$Z = \frac{\hat{p} - p}{\sqrt{\frac{pq}{n}}} \sim N(0,1) \qquad \text{(approximately)}$$

**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 = $\frac{45}{100} = 0.45$

- $\hat{p}$ = 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:

$$\frac{p(1-p)}{n} = \frac{pq}{n} = \frac{0.45(1-0.45)}{35} = 0.0071.$$

- The standard error (standard deviation) of the sample proportion ($\hat{p}$) is:

$$\sqrt{\frac{p(1-p)}{n}} = \sqrt{0.0071} = 0.084$$

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

1. The probability that the sample proportion of females ($\hat{p}$) will be greater than 0.4 is:

$$P(\hat{p} > 0.4) = 1 - P(\hat{p} < 0.4) = 1 - P\left(\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} < \frac{0.4 - p}{\sqrt{\frac{p(1-p)}{n}}}\right)$$

$$= 1 - P\left(Z < \frac{0.4 - 0.45}{\sqrt{\frac{0.45(1-0.45)}{35}}}\right) = 1 - P(Z < -0.59)$$

$$= 1 - 0.2776 = 0.7224$$

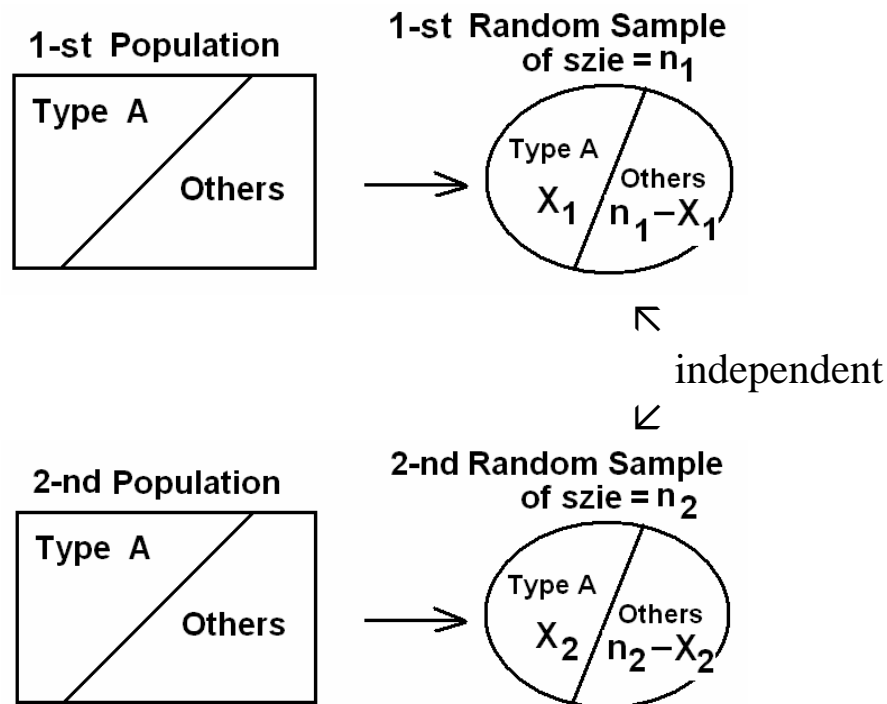2. The probability that the sample proportion of females ($\hat{p}$) will be between 0.4 and 0.5 is:

$$P(0.4 < \hat{p} < 0.5) = P(\hat{p} < 0.5) - P(\hat{p} < 0.4)$$

$$= P\left(\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} < \frac{0.5 - p}{\sqrt{\frac{p(1-p)}{n}}}\right) - 0.2776$$

$$= P\left(Z < \frac{0.5 - 0.45}{\sqrt{\frac{0.45(1-0.45)}{35}}}\right) - 0.2776$$

$$= P(Z < 0.59) - 0.2776$$
$$= 0.7224 - 0.2776$$
$$= 0.4448$$

## 5.6 Distribution 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.
- $\hat{p}_1 = \dfrac{X_1}{n_1}$ = sample proportion of the 1-st sample

- $\hat{p}_2 = \dfrac{X_2}{n_2}$ = 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:**
The mean, the variance and the standard error (standard deviation) of $\hat{p}_1 - \hat{p}_2$ are:

- Mean of $\hat{p}_1 - \hat{p}_2$ is:

$$\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2$$

- Variance of $\hat{p}_1 - \hat{p}_2$ is:

$$\sigma^2_{\hat{p}_1 - \hat{p}_2} = \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}$$

- Standard error (standard deviation) of $\hat{p}_1 - \hat{p}_2$ is:

$$\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$$

- $q_1 = 1 - p_1$ and $q_2 = 1 - p_2$

**Result:**
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 that $\hat{p}_1 - \hat{p}_2$ has approximately normal distribution with mean $\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2$ and variance $\sigma^2_{\hat{p}_1 - \hat{p}_2} = \dfrac{p_1 q_1}{n_1} + \dfrac{p_2 q_2}{n_2}$, that is:

$$\hat{p}_1 - \hat{p}_2 \sim N\left(p_1 - p_2, \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}\right) \quad \text{(Approximately)}$$

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\dfrac{p_1 q_1}{n_1} + \dfrac{p_2 q_2}{n_2}}} \sim N(0,1) \quad \text{(Approximately)}$$

**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, $\hat{p}_1 - \hat{p}_2$, will be between 0.05 and 0.2?

**Solution:**

$p_1$ = population proportion of non-Saudi with medical insurance.
$p_2$ = population proportion of Saudi with medical insurance.
$\hat{p}_1$ = sample proportion of non-Saudis with medical insurance.
$\hat{p}_2$ = sample proportion of Saudis with medical insurance.

q1=0.6   $p_1 = 0.4$        $n_1$=130  **> 30**
q2=0.7   $p_2 = 0.3$        $n_2$=120  **>30**

$$\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2 = 0.4 - 0.3 = 0.1$$

$$\sigma^2_{\hat{p}_1 - \hat{p}_2} = \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2} = \frac{(0.4)(0.6)}{130} + \frac{(0.3)(0.7)}{120} = 0.0036$$

$$\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}} = \sqrt{0.0036} = 0.06$$

The probability that the difference between the sample proportions, $\hat{p}_1 - \hat{p}_2$, will be between 0.05 and 0.2 is:

$$P(0.05 < \hat{p}_1 - \hat{p}_2 < 0.2) = P(\hat{p}_1 - \hat{p}_2 < 0.2) - P(\hat{p}_1 - \hat{p}_2 < 0.05)$$

$$= P\left( \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} < \frac{0.2 - (p_1 - p_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} \right)$$

$$-\,\mathrm{P}\left(\frac{(\hat{p}_1-\hat{p}_2)-(p_1-p_2)}{\sqrt{\dfrac{p_1\,q_1}{n_1}+\dfrac{p_2\,q_2}{n_2}}}<\frac{0.05-(p_1-p_2)}{\sqrt{\dfrac{p_1\,q_1}{n_1}+\dfrac{p_2\,q_2}{n_2}}}\right)$$

$$=\mathrm{P}\left(Z<\frac{0.2-0.1}{0.06}\right)-\mathrm{P}\left(Z<\frac{0.05-0.1}{0.06}\right)$$

$$=\mathrm{P}(Z<1.67)-\mathrm{P}(Z<-0.83)$$

= 0.95254 - 0.20327

= 0.74927

# CHAPTER 6: Using Sample Data to Make Estimations About Population Parameters

## 6.1 Introduction:

Statistical Inferences: (Estimation and Hypotheses Testing)

It is the procedure by which we reach a conclusion about a population on the basis of the information contained in a sample drawn from that population.

There are two main purposes of statistics;
- Descriptive Statistics: (Chapter 1 & 2): Organization & summarization of the data
- Statistical Inference: (Chapter 6 and 7): Answering research questions about some unknown population parameters.

**(1) Estimation:** (chapter 6)

Approximating (or estimating) the actual values of the unknown parameters:
- **Point Estimate:** A point estimate is single value used to estimate the corresponding population parameter.
- **Interval Estimate (or Confidence Interval):** An interval estimate consists of two numerical values defining a range of values that most likely includes the parameter being estimated with a specified degree of confidence.

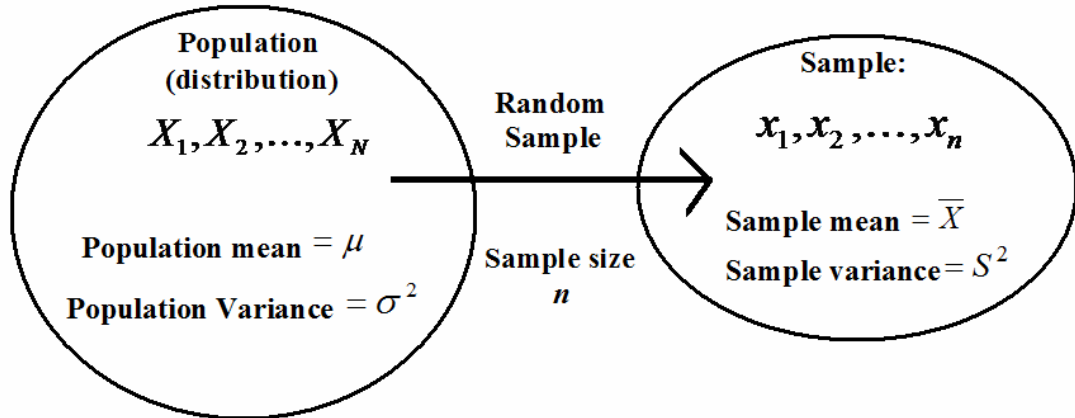**(2) Hypothesis Testing:** (chapter 7)

Answering research questions about the unknown parameters of the population (confirming or denying some conjectures or statements about the unknown parameters).

## 6.1: The Point Estimates of the Population Parameters:

| | Population Parameters | Point estimator |
|---|---|---|
| Mean | $\mu$ | $\bar{X}$ |
| Variance | $\sigma^2$ | $S^2$ |
| Standard Deviation | $\sigma$ | $S$ |
| Proportion | $P$ | $\hat{p}$ |
| The Difference between Two Means | $\mu_1 - \mu_2$ | $\bar{X_1} - \bar{X_2}$ |
| The Difference between Two Proportion | $P_1 - P_2$ | $\widehat{P_1} - \widehat{P_2}$ |

## 6.2 Confidence Interval for a Population Mean (μ) :

In this section we are interested in estimating the mean of a certain population $(\mu)$.



| Population: | Sample: |
|---|---|
| Population Size $= N$ | Sample Size $= n$ |
| Population Values: $X_1, X_2, \ldots, X_N$ | Sample values: $x_1, x_2, \ldots, x_n$ |
| Population Mean: $\mu = \dfrac{\sum\limits_{i=1}^{N} X_i}{N}$ | Sample Mean: $\overline{X} = \dfrac{\sum\limits_{i=1}^{n} x_i}{n}$ |
| Population Variance: $\sigma^2 = \dfrac{\sum\limits_{i=1}^{N}(X_i - \mu)^2}{N}$ | Sample Variance: $S^2 = \dfrac{\sum\limits_{i=1}^{n}(x_i - \overline{x})^2}{n-1}$ |

## (i) Point Estimation of μ:

A point estimate of the mean is a single number used to estimate (or approximate) the true value of $\mu$ .

– Draw a random sample of size $n$ from the population:

$$ -\ x_1, x_2, \ldots, x_n $$

– Compute the sample mean: $\overline{X} = \dfrac{1}{n}\sum\limits_{i=1}^{n} x_i$

## Result:

The sample mean $\overline{X} = \dfrac{1}{n}\sum\limits_{i=1}^{n} x_i$ is a "good" point estimator of the population mean $(\mu)$.

## (1-α) % confident level

- ### How to get α when confidence level (1-α) % known

### Example1 :

If we are 95% confident ,find α ?

$$\alpha = \frac{5}{100} = 0.05$$

### Example2 :

If we are 99% confident ,find α ?

$$\alpha = \frac{1}{100} = 0.01$$

### Example3 :

If we are 80% confident ,find α ?

$$\alpha = \frac{20}{100} = 0.20$$

### Example4 :

If we are 92% confident ,find α ?

$$\alpha = \frac{8}{100} = 0.08$$

---

## (ii) Confidence Interval (Interval Estimate) of μ:

An interval estimate of $\mu$ is an interval $(L,U)$ containing the true value of $\mu$ "with a probability of $1-\alpha$".

<span style="color:blue">(confidence level), degree of confidence</span>

* $1-\alpha$ = is called the confidence coefficient (level)
* L = lower limit of the confidence interval
* U = upper limit of the confidence interval

**Result:** (For the case when $\sigma$ is known)

(a) If $X_1, X_2 ..., X_n$ is a random sample of size $n$ from a normal distribution with mean $\mu$ and known variance $\sigma^2$, then:

A $(1-\alpha)100\%$ confidence interval for $\mu$ is:

$$\overline{X} \pm Z_{1-\frac{\alpha}{2}} \; \sigma_{\overline{X}}$$

$$\overline{X} \pm Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

$$\left( \overline{X} - Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \;,\; \overline{X} + Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right)$$

$$\overline{X} - Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \; < \mu < \; \overline{X} + Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

(b) If $X_1, X_2 ..., X_n$ is a random sample of size $n$ from a non-normal distribution with mean $\mu$ and known variance $\sigma^2$, and if the sample size $n$ is large $(n \geq 30)$, then:

An approximate $(1-\alpha)100\%$ confidence interval for $\mu$ is:

$$\overline{X} \pm Z_{1-\frac{\alpha}{2}} \; \sigma_{\overline{X}}$$

$$\overline{X} \pm Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

$$\left( \overline{X} - Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \;,\; \overline{X} + Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right)$$

$$\overline{X} - Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \; < \mu < \; \overline{X} + Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

Note that:

1. We are $(1-\alpha)100\%$ confident that the true value of $\mu$ belongs to the interval $(\overline{X} - Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} , \overline{X} + Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}})$ .

2. Upper limit of the confidence interval = $\overline{X} + Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$

3. Lower limit of the confidence interval = $\overline{X} - Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$

4. $Z_{1-\frac{\alpha}{2}}$ = Reliability Coefficient

5. $Z_{1-\frac{\alpha}{2}} \times \frac{\sigma}{\sqrt{n}}$ = margin of error = precision of the estimate

6. In general the interval estimate (confidence interval) may be expressed as follows:

$$\overline{X} \pm Z_{1-\frac{\alpha}{2}} \ \sigma_{\overline{X}}$$

estimator $\pm$ (reliability coefficient) $\times$ (standard Error)

estimator $\pm$ margin of error

## 6.3 The t Distribution:
## (Confidence Interval Using t)

We have already introduced and discussed the t distribution.

**Result:** (For the case when $\sigma$ is unknown + normal population) + n < 30
If $X_1, X_2 ..., X_n$ is a random sample of size $n$ from a normal distribution with mean $\mu$ and unknown variance $\sigma^2$ , then:
A $(1-\alpha)100\%$ confidence interval for $\mu$ is:

$$\overline{X} \pm t_{1-\frac{\alpha}{2}} \ \hat{\sigma}_{\overline{X}}$$

$$\overline{X} \pm t_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}}$$

$$\left( \overline{X} - t_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}} , \overline{X} + t_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}} \right)$$

where the degrees of freedom is:
$$df = v = n-1.$$
Note that:

1. We are $(1-\alpha)100\%$ confident that the true value of $\mu$ belongs

to the interval $\left( \overline{X} - t_{1-\frac{\alpha}{2}} \dfrac{S}{\sqrt{n}} , \overline{X} + t_{1-\frac{\alpha}{2}} \dfrac{S}{\sqrt{n}} \right)$.

2. $\hat{\sigma}_{\overline{X}} = \dfrac{S}{\sqrt{n}}$     (estimate of the standard error of $\overline{X}$ )

3. $t_{1-\frac{\alpha}{2}}$ = Reliability Coefficient

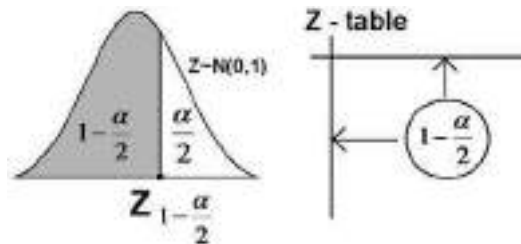4. In this case, we replace $\sigma$ by $S$ and Z by t.

5. In general the interval estimate (confidence interval) may be expressed as follows:

  Estimator $\pm$ (Reliability Coefficient) $\times$ (Estimate of the Standard Error)
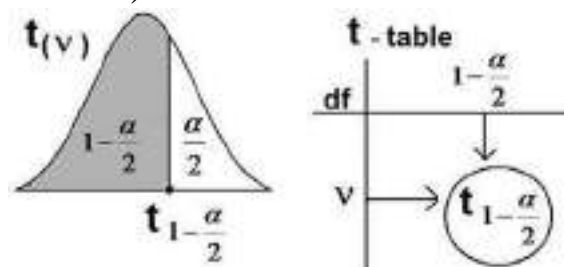$$\overline{X} \pm t_{1-\frac{\alpha}{2}} \; \hat{\sigma}_{\overline{X}}$$

**Notes: (Finding Reliability Coefficient)**

(1) We find the reliability coefficient $Z_{1-\frac{\alpha}{2}}$ from the Z-table as

follows:



(2) We find the reliability coefficient $t_{1-\frac{\alpha}{2}}$ from the t-table as

follows: $(df = v = n-1)$

**Example:**

Suppose that $Z \sim N(0,1)$. Find $Z_{1-\frac{\alpha}{2}}$ for the following cases:

(1) $\alpha = 0.1$      (2) $\alpha = 0.05$      (3) $\alpha = 0.01$

**Solution:**

(1) For $\alpha = 0.1$:
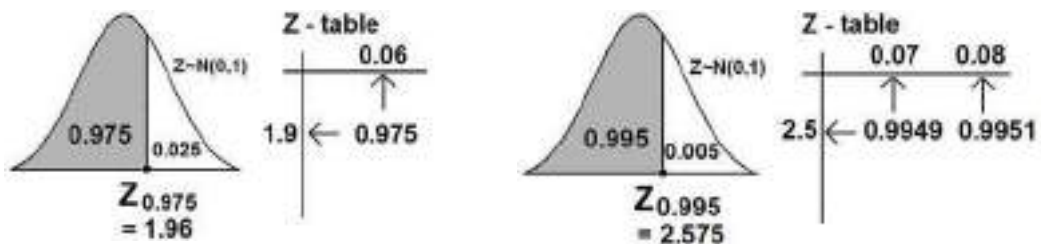
$$1 - \frac{\alpha}{2} = 1 - \frac{0.1}{2} = 0.95 \qquad \Rightarrow \qquad Z_{1-\frac{\alpha}{2}} = Z_{0.95} = 1.645$$

(2) For $\alpha = 0.05$:

$$1 - \frac{\alpha}{2} = 1 - \frac{0.05}{2} = 0.975 \qquad \Rightarrow \qquad Z_{1-\frac{\alpha}{2}} = Z_{0.975} = 1.96.$$

(3) For $\alpha = 0.01$:

$$1 - \frac{\alpha}{2} = 1 - \frac{0.01}{2} = 0.995 \qquad \Rightarrow \qquad Z_{1-\frac{\alpha}{2}} = Z_{0.995} = 2.575.$$
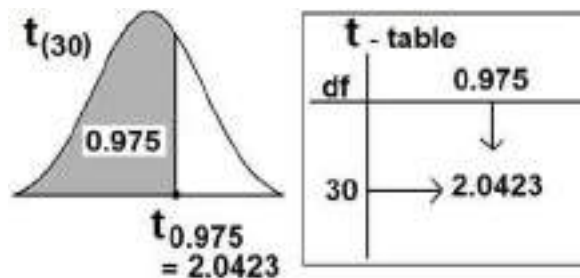


**Example:**

Suppose that $t \sim t(30)$. Find $t_{1-\frac{\alpha}{2}}$ for $\alpha = 0.05$.

**Solution:**

$df = \nu = 30$

$$1 - \frac{\alpha}{2} = 1 - \frac{0.05}{2} = 0.975 \qquad \Rightarrow \qquad t_{1-\frac{\alpha}{2}} = t_{0.975} = 2.0423$$

# The Confidence Interval (C.I) for the Population Mean $\mu$:

The $(1-\alpha)100\%$ Confidence Interval for the Population Mean $\mu$

1. Normal Distribution $+ \sigma^2$ is known.

2. Non-Normal Distribution $+ n \geq 30 + \sigma^2$ is known.

$$\bar{X} \pm Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

Normal Distribution $+ \sigma^2$ is unknown $+ n < 30$.

$$\bar{X} \pm t_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}}$$

**Example: (The case where $\sigma^2$ is known)**

Diabetic ketoacidosis is a potential fatal complication of diabetes mellitus throughout the world and is characterized in part by very high blood glucose levels. In a study on 123 patients living in Saudi Arabia of age 15 or more who were admitted for diabetic ketoacidosis, the mean blood glucose level was 26.2 mmol/l. Suppose that the blood glucose levels for such patients have a normal distribution with a standard deviation of 3.3 mmol/l.

(1) Find a point estimate for the mean blood glucose level of such diabetic ketoacidosis patients.

(2) Find a 90% confidence interval for the mean blood glucose level of such diabetic ketoacidosis patients.

**Solution:**

Variable = X = blood glucose level (Continuous quantitative variable).

Population = diabetic ketoacidosis patients in Saudi Arabia of age 15 or more.

Parameter of interest is: $\mu$ = the mean blood glucose level.

Distribution is normal with standard deviation $\sigma = 3.3$.

$\sigma^2$ is known ($\sigma^2 = 10.89$)

X ~ Normal($\mu$, 10.89)

$\mu$ = ?? (unknown- we need to estimate $\mu$)

Sample size: $n = 123$ (large)

Sample mean: $\overline{X} = 26.2$

(1) Point Estimation:

We need to find a point estimate for $\mu$.

$\overline{X} = 26.2$ is a point estimate for $\mu$.

$\mu \approx 26.2$

(2) Interval Estimation (Confidence Interval = C. I.):

We need to find 90% C. I. for $\mu$.

$90\% = (1 - \alpha)100\%$

$1 - \alpha = 0.9 \Leftrightarrow \alpha = 0.1 \Leftrightarrow \dfrac{\alpha}{2} = 0.05 \Leftrightarrow 1 - \dfrac{\alpha}{2} = 0.95$

The reliability coefficient is: $Z_{1-\frac{\alpha}{2}} = Z_{0.95} = 1.645$

90% confidence interval for $\mu$ is:

$$\left( \overline{X} - Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \; , \; \overline{X} + Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right)$$

$$\left( 26.2 - (1.645)\frac{3.3}{\sqrt{123}} \; , \; 26.2 + (1.645)\frac{3.3}{\sqrt{123}} \right)$$

$$(26.2 - 0.4894714 \; , \; 26.2 + 0.4894714)$$

$$(25.710529 \; , \; 26.689471)$$

We are 90% confident that the true value of the mean $\mu$ lies in the interval $(25.71 \, , \, 26.69)$, that is:

$$25.71 < \mu < 26.69$$

Note: for this example even if the distribution is not normal, we may use the same solution because the sample size n=123 is large.

**Example: (The case where $\sigma^2$ is unknown)**

A study was conducted to study the age characteristics of Saudi women having breast lump. A sample of **21** Saudi women gave a mean of 37 years with a standard deviation of 10 years. Assume that the ages of Saudi women having breast lumps are normally distributed.

(a) Find a point estimate for the mean age of Saudi women having breast lumps.

(b) Construct a 99% confidence interval for the mean age of Saudi women having breast lumps

**Solution:**

X = Variable = age of Saudi women having breast lumps (quantitative variable).

Population = All Saudi women having breast lumps.

Parameter of interest is: $\mu =$ the age mean of Saudi women having breast lumps.

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

$\mu = ??$ (unknown- we need to estimate $\mu$ )

$\sigma^2 = ??$ (unknown)

Sample size:      $n = 21$

Sample mean:    $\overline{X} = 37$

Sample standard deviation: $S = 10$

Degrees of freedom: df $= \nu =$ 21 - 1 =20

(a)  Point Estimation: We need to find a point estimate for $\mu$.

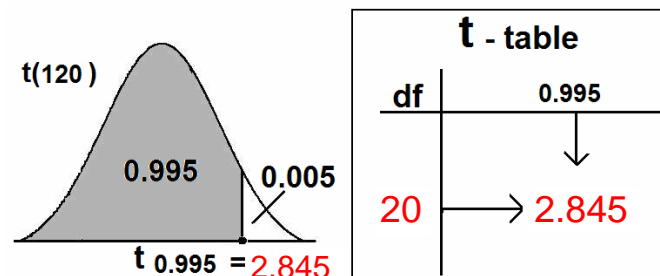$\overline{X} = 37$  is a "good" point estimate for $\mu$.

$\mu \approx 37 \; years$

(b) Interval Estimation (Confidence Interval = C. I.): We need to find 99% C. I. for $\mu$.

$99\% = (1-\alpha)100\%$

$1 - \alpha = 0.99 \Leftrightarrow \alpha = 0.01 \quad \Leftrightarrow \quad \dfrac{\alpha}{2} = 0.005 \quad \Leftrightarrow \quad 1 - \dfrac{\alpha}{2} = 0.995$

$\nu = df =$ 21-1=20

The reliability coefficient is: $t_{1-\frac{\alpha}{2}} = t_{0.995} =$ 2.845



99% confidence interval for $\mu$ is:

$$\overline{X} \pm t_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}}$$

$$37 \pm (2.845)\frac{10}{\sqrt{21}}$$

$$37 \pm 6.208$$

$$(37 - 6.208 , 37 + 6.208)$$
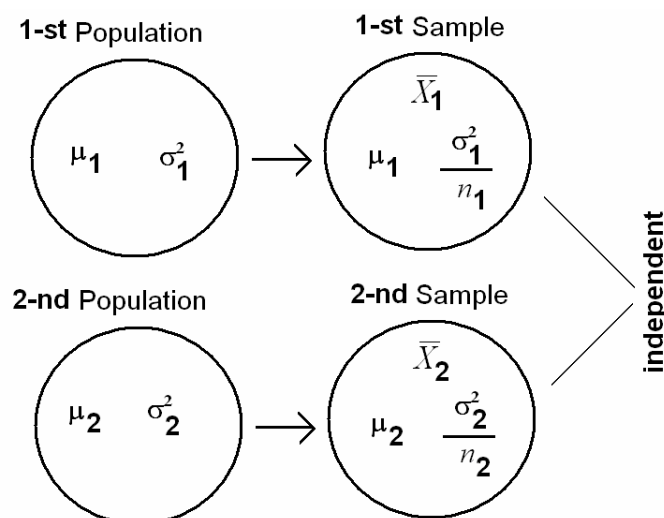
$$(30.792 , 43.208)$$

$$30.792 < \mu < 43.208$$

We are 99% confident that the true value of the mean μ lies in the interval $(30.792 , 43.208)$

112

## 6.4 Confidence Interval for the Difference between Two Population Means ($\mu_1 - \mu_2$):

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 $\mu_1$ and $\mu_2$, or equivalently, making inferences about the difference between the means ($\mu_1 - \mu_2$).
- We <u>independently</u> 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 $\overline{X}_1$ and $S_1^2$ be the sample mean and the sample variance of the 1-st sample.
- Let $\overline{X}_2$ and $S_2^2$ be the sample mean and the sample variance of the 2-nd sample.
- The sampling distribution of $\overline{X}_1 - \overline{X}_2$ is used to make inferences about $\mu_1 - \mu_2$.

**Recall:**

1. Mean of $\bar{X}_1 - \bar{X}_2$ is:

$$\mu_{\bar{X}_1 - \bar{X}_2} = \mu_1 - \mu_2$$

2. Variance of $\bar{X}_1 - \bar{X}_2$ is:

$$\sigma^2_{\bar{X}_1 - \bar{X}_2} = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

3. Standard error of $\bar{X}_1 - \bar{X}_2$ is:

$$\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

4. 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 $((\sigma_1^2/n_1) + (\sigma_2^2/n_2))$, that is:

- $\bar{X}_1 - \bar{X}_2 \sim N\left(\mu_1 - \mu_2 \;,\; \dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}\right)$

- $Z = \dfrac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}} \sim N(0,1)$

**Point Estimation of $\mu_1 - \mu_2$:**
**Result:**

$\bar{X}_1 - \bar{X}_2$ is a "good" point estimate for $\mu_1 - \mu_2$.

**Interval Estimation (Confidence Interval) of $\mu_1 - \mu_2$:**
We will consider two cases.
**(i) First Case: $\sigma_1^2$ and $\sigma_2^2$ are known:**

If $\sigma_1^2$ and $\sigma_2^2$ are known, we use the following result to find an interval estimate for $\mu_1 - \mu_2$.
**Result:**
A $(1-\alpha)100\%$ confidence interval for $\mu_1 - \mu_2$ is:

$$(\bar{X}_1 - \bar{X}_2) \pm Z_{1-\frac{\alpha}{2}}\; \sigma_{\bar{X}_1 - \bar{X}_2}$$

$$(\overline{X}_1 - \overline{X}_2) \pm Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

$$\left( (\overline{X}_1 - \overline{X}_2) - Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad , \quad (\overline{X}_1 - \overline{X}_2) + Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

$$(\overline{X}_1 - \overline{X}_2) - Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2 < (\overline{X}_1 - \overline{X}_2) + Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Estimator $\pm$ (Reliability Coefficient) $\times$ (Standard Error)

## (ii) Second Case:
**Unknown equal Variances: ($\sigma_1^2 = \sigma_2^2 = \sigma^2$ is unknown):**

If $\sigma_1^2$ and $\sigma_2^2$ are equal but unknown ($\sigma_1^2 = \sigma_2^2 = \sigma^2$), then the pooled estimate of the common variance $\sigma^2$ is

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

where $S_1^2$ is the variance of the 1-st sample and $S_2^2$ is the variance of the 2-nd sample. The degrees of freedom of $S_p^2$ is

$$df = \nu = n_1 + n_2 - 2.$$

We use the following result to find an interval estimate for $\mu_1 - \mu_2$ when we have normal populations with unknown and equal variances.

**Result:**

A $(1-\alpha)100\%$ confidence interval for $\mu_1 - \mu_2$ is:

$$(\overline{X}_1 - \overline{X}_2) \pm t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_p^2}{n_1} + \frac{S_p^2}{n_2}}$$

$$\left( (\overline{X}_1 - \overline{X}_2) - t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_p^2}{n_1} + \frac{S_p^2}{n_2}} \quad , \quad (\overline{X}_1 - \overline{X}_2) + t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_p^2}{n_1} + \frac{S_p^2}{n_2}} \right)$$

where reliability coefficient $t_{1-\frac{\alpha}{2}}$ is the t-value with

df=$\nu$=$n_1$+$n_2$−2 degrees of freedom.

~~**Example:** (1st Case: $\sigma_1^2$ and $\sigma_2^2$ are known)~~

~~An experiment was conducted to compare time length~~

# The Confidence Interval ( C.I )for the Difference between two Population Means $\mu_1 - \mu_2$:

The $(1 - \alpha)100\%$

Confidence Interval for the Difference between two Population Means

$\mu_1 - \mu_2$

1. Normal Distribution + $\sigma_1^2$ and $\sigma_2^2$ are known

2. Non-Normal Distribution + $n_1, n_2 \geq 30$ + $\sigma_1^2$ and $\sigma_2^2$ are known.

$$(\bar{X}_1 - \bar{X}_2) \pm Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Normal Distribution + $n_1, n_2 < 30$ + $\sigma_1^2 = \sigma_2^2 = \sigma^2$ are unknown but equal

$$(\bar{X}_1 - \bar{X}_2) \pm t_{1-\frac{\alpha}{2}} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

**pooled variance:**

$$S_p = \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}} \quad , \text{d.f} = n_1 + n_2 - 2$$

An experiment was conducted to compare time length (duration time) of two types of surgeries (A) and (B). 75 surgeries of type (A) and 50 surgeries of type (B) were performed. The average time length for (A) was 42 minutes and the average for (B) was 36 minutes.

(1) Find a point estimate for $\mu_A - \mu_B$, where $\mu_A$ and $\mu_B$ are population means of the time length of surgeries of type (A) and (B), respectively.

(2) Find a 96% confidence interval for $\mu_A - \mu_B$. Assume that the population standard deviations are 8 and 6 for type (A) and (B), respectively.

**Solution:**

| Surgery | Type (A) | Type (B) |
|---|---|---|
| Sample Size | $n_A = 75$ | $n_B = 50$ |
| Sample Mean | $\overline{X}_A = 42$ | $\overline{X}_B = 36$ |
| Population Standard Deviation | $\sigma_A = 8$ | $\sigma_B = 6$ |

(1) A point estimate for $\mu_A - \mu_B$ is:
$$\overline{X}_A - \overline{X}_B = 42 - 36 = 6.$$

(2) Finding a 96% confidence interval for $\mu_A - \mu_B$:

$\alpha = ??$

$96\% = (1-\alpha)100\% \Leftrightarrow 0.96 = (1-\alpha) \Leftrightarrow \alpha = 0.04 \Leftrightarrow \alpha/2 = 0.02$

Reliability Coefficient: $Z_{1-\frac{\alpha}{2}} = Z_{0.98} = 2.055$

A 96% C.I. for $\mu_A - \mu_B$ is:

$$(\overline{X}_A - \overline{X}_B) \pm Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B}}$$

$$6 \pm Z_{0.98} \sqrt{\frac{8^2}{75} + \frac{6^2}{50}}$$

$$6 \pm (2.055) \sqrt{\frac{64}{75} + \frac{36}{50}}$$

$$6 \pm 2.578$$

$$3.422 < \mu_A - \mu_B < 8.58$$

We are 96% confident that $\mu_A - \mu_B \in (3.42, 8.58)$.
Note: Since the confidence interval does not include zero, we conclude that the two population means are not equal ($\mu_A - \mu_B \neq 0 \Leftrightarrow \mu_A \neq \mu_B$). Therefore, we may conclude that the mean time length is not the same for the two types of surgeries.

**Example:** ($2^{nd}$ Case: $\sigma_1^2 = \sigma_2^2$ unknown)

To compare the time length (duration time) of two types of surgeries (A) and (B), an experiment shows the following results based on two independent samples:

Type *A*:   140, 138, 143, 142, 144, 137
Type *B*:   135, 140, 136, 142, 138, 140

(1) Find a point estimate for $\mu_A - \mu_B$, where $\mu_A$ ($\mu_B$) is the mean time length of type *A* (*B*).
(2) Assuming normal populations with equal variances, find a 95% confidence interval for $\mu_A - \mu_B$.

**Solution:**

First we calculate the mean and the variances of the two samples, and we get:

| Surgery | Type (A) | Type (B) |
|---|---|---|
| Sample Size | $n_A = 6$ | $n_B = 6$ |
| Sample Mean | $\overline{X}_A = 140.67$ | $\overline{X}_B = 138.50$ |
| Sample Variance | $S_A^2 = 7.87$ | $S_B^2 = 7.10$ |

(1) A point estimate for $\mu_A - \mu_B$ is:
$$\overline{X}_A - \overline{X}_B = 140.67 - 138.50 = 2.17.$$

(2) Finding 95% Confidence interval for $\mu_A - \mu_B$:

95% = $(1-\alpha)100\% \Leftrightarrow 0.95 = (1-\alpha) \Leftrightarrow \alpha = 0.05 \Leftrightarrow \alpha/2 = 0.025$
.df $= \nu = n_A + n_B - 2 = 10$
Reliability Coefficient: $t_{1-\frac{\alpha}{2}} = t_{0.975} = 2.228$

The pooled estimate of the common variance is:

$$S_p^2 = \frac{(n_A - 1)S_A^2 + (n_B - 1)S_B^2}{n_A + n_B - 2}$$

$$= \frac{(6-1)(7.87) + (6-1)(7.1)}{6+6-2} = 7.485$$

A 95% C.I. for $\mu_A - \mu_B$ is:

$$(\bar{X}_A - \bar{X}_B) \pm t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_p^2}{n_A} + \frac{S_p^2}{n_B}}$$

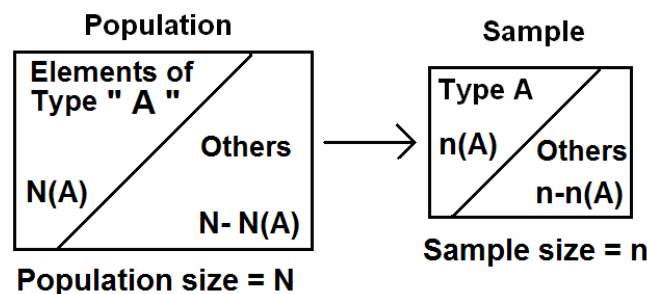$$2.17 \pm (2.228) \sqrt{\frac{7.485}{6} + \frac{7.485}{6}}$$

$$2.17 \pm 3.519$$

$$-1.35 < \mu_A - \mu_B < 5.69$$

We are 95% confident that $\mu_A - \mu_B \in (-1.35, 5.69)$.

Note: Since the confidence interval includes zero, we conclude that the two population means may be equal ($\mu_A - \mu_B = 0 \Leftrightarrow \mu_A = \mu_B$). Therefore, we may conclude that the mean time length is the same for both types of surgeries.

## 6.5 Confidence Interval for a Population Proportion (p):



**Recall:**

1. 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:

$$p = \frac{N(A)}{N} \qquad \text{(p is a parameter)}$$

2. For the sample:

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

$n =$ sample size

The sample proportion is:

$$\hat{p} = \frac{n(A)}{n}$$
($\hat{p}$ is a statistic)

3. The sampling distribution of the sample proportion ($\hat{p}$) is used to make inferences about the population proportion (p).

4. The mean of ($\hat{p}$) is: $\mu_{\hat{p}} = p$

5. The variance of ($\hat{p}$) is: $\sigma_{\hat{p}}^2 = \dfrac{p(1-p)}{n}$

6. The standard error (standard deviation) of ($\hat{p}$) is:

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}} .$$

7. For large sample size ($n \geq 30, np > 5, n(1-p) > 5$), the sample proportion ($\hat{p}$) has approximately a normal distribution with mean $\mu_{\hat{p}} = p$ and a variance $\sigma_{\hat{p}}^2 = p(1-p)/n$, that is:

$$\hat{p} \sim N\left( p, \frac{p(1-p)}{n} \right) \qquad \text{(approximately)}$$

$$Z = \frac{\hat{p} - p}{\sqrt{\dfrac{p(1-p)}{n}}} \sim N(0,1) \qquad \text{(approximately)}$$

**(i) Point Estimate for (p):**

**Result:**

A good point estimate for the population proportion (p) is the sample proportion ($\hat{p}$).

**(ii) Interval Estimation (Confidence Interval) for (p):**

**Result:**

For large sample size ($n \geq 30, np > 5, n(1-p) > 5$), an approximate $(1-\alpha)100\%$ confidence interval for (p) is:

$$\hat{p} \pm Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$\left( \hat{p} - Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \quad , \quad \hat{p} + Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right)$$

Estimator $\pm$ (Reliability Coefficient) $\times$ (Standard Error)

**Example:**

In a study on the obesity of Saudi women, a random sample of 950 Saudi women was taken. It was found that 611 of these women were obese (overweight by a certain percentage).
(1) Find a point estimate for the true proportion of Saudi women who are obese.
(2) Find a 95% confidence interval for the true proportion of Saudi women who are obese.

**Solution:**

Variable: whether or not a women is obese (qualitative variable)
Population: all Saudi women
Parameter:  p =the proportion of women who are obese.

Sample:
$n = 950$        (950 women in the sample)
$n(A) = 611$     (611 women in the sample who are obese)
The sample proportion (the proportion of women who are obese in the sample.) is:

$$\hat{p} = \frac{n(A)}{n} = \frac{611}{950} = 0.643$$

(1) A point estimate for p is:        $\hat{p} = 0.643$.
(2) We need to construct 95% C.I. for the proportion (p).

$$95\% = (1-\alpha)100\% \iff 0.95 = 1-\alpha \iff \alpha = 0.05 \iff \frac{\alpha}{2} = 0.025 \iff 1-\frac{\alpha}{2} = 0.975$$

The reliability coefficient:  $Z_{1-\frac{\alpha}{2}} = z_{0.975} = 1.96$.

A 95% C.I. for the proportion (p) is:

$$\hat{p} \pm Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

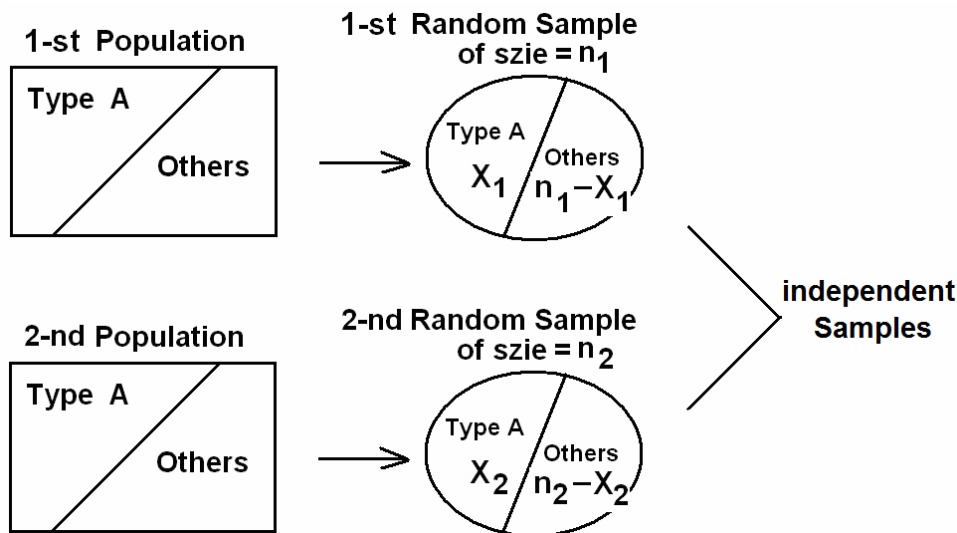$$0.643 \pm (1.96)\sqrt{\frac{(0.643)(1-0.643)}{950}}$$

$$0.643 \pm (1.96)(0.01554)$$

$$0.643 \pm 0.0305$$

$$(0.6127 \ , \ 0.6735)$$

We are 95% confident that the true value of the population proportion of obese women, p, lies in the interval $(0.61, 0.67)$, that is:

$$0.61 < p < 0.67$$

## 6.6 Confidence Interval for the Difference Between Two Population Proportions ( $p_1 - p_2$ ):



Suppose that we have two populations with:

- $p_1$ = population proportion of elements of type (A) in the 1-st population.
- $p_2$ = population 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.
- $\hat{p}_1 = \dfrac{X_1}{n_1}$ = the sample proportion of the 1-st sample
- $\hat{p}_2 = \dfrac{X_2}{n_2}$ = the sample proportion of the 2-nd sample
- The sampling distribution of $\hat{p}_1 - \hat{p}_2$ is used to make inferences about $p_1 - p_2$.

**Recall:**

1. Mean of $\hat{p}_1 - \hat{p}_2$ is: $\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2$

2. Variance of $\hat{p}_1 - \hat{p}_2$ is: $\sigma^2_{\hat{p}_1 - \hat{p}_2} = \dfrac{p_1 q_1}{n_1} + \dfrac{p_2 q_2}{n_2}$

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

4. 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 that $\hat{p}_1 - \hat{p}_2$ has approximately normal distribution with mean $\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2$ and variance $\sigma^2_{\hat{p}_1 - \hat{p}_2} = \dfrac{p_1 q_1}{n_1} + \dfrac{p_2 q_2}{n_2}$, that is:

$$\hat{p}_1 - \hat{p}_2 \sim N\left(p_1 - p_2, \dfrac{p_1 q_1}{n_1} + \dfrac{p_2 q_2}{n_2}\right) \quad \text{(Approximately)}$$

$$Z = \dfrac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\dfrac{p_1 q_1}{n_1} + \dfrac{p_2 q_2}{n_2}}} \sim N(0,1) \quad \text{(Approximately)}$$

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

## Point Estimation for $p_1 - p_2$:
**Result:**

A good point estimator for the difference between the two proportions, $p_1 - p_2$, is:

$$\hat{p}_1 - \hat{p}_2 = \frac{X_1}{n_1} - \frac{X_2}{n_2}$$

## **Interval Estimation (Confidence Interval) for _p₁− p₂_:**
**Result:**

For large $n_1$ and $n_2$, an approximate $(1-\alpha)100\%$ confidence interval for $p_1 - p_2$ is:

$$(\hat{p}_1 - \hat{p}_2) \pm Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}}$$

$$\left( (\hat{p}_1 - \hat{p}_2) - Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}} , \ (\hat{p}_1 - \hat{p}_2) + Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}} \right)$$

Estimator $\pm$ (Reliability Coefficient) $\times$ (Standard Error)

**Example:**

A researcher was interested in comparing the proportion of people having cancer disease in two cities (A) and (B). A random sample of 1500 people was taken from the first city (A), and another independent random sample of 2000 people was taken from the second city (B). It was found that 75 people in the first sample and 80 people in the second sample have cancer disease.

(1) Find a point estimate for the difference between the proportions of people having cancer disease in the two cities.

(2) Find a 90% confidence interval for the difference between the two proportions.

**Solution:**

$p_1$ = population proportion of people having cancer disease in the first city (A)

$p_2$ = population proportion of people having cancer disease in the second city (B)

$\hat{p}_1$ = sample proportion of the first sample

$\hat{p}_2$ = sample proportion of the second sample

$X_1$ = number of people with cancer in the first sample

$X_2$ = number of people with cancer in the second sample

For the first sample we have:

$n_1 = 1500$ , $X_1 = 75$

$$\hat{p}_1 = \frac{X_1}{n_1} = \frac{75}{1500} = 0.05 \quad , \qquad \hat{q}_1 = 1 - 0.05 = 0.95$$

For the second sample we have:

$$n_2 = 2000 \quad , \qquad X_2 = 80$$

$$\hat{p}_2 = \frac{X_2}{n_2} = \frac{80}{2000} = 0.04 \quad , \qquad \hat{q}_2 = 1 - 0.04 = 0.96$$

(1) Point Estimation for $p_1 - p_2$:

A good point estimate for the difference between the two proportions, $p_1 - p_2$, is:

$$\hat{p}_1 - \hat{p}_2 = 0.05 - 0.04$$
$$= 0.01$$

(2) Finding 90% Confidence Interval for $p_1 - p_2$:

90% = $(1-\alpha)100\% \Leftrightarrow 0.90 = (1-\alpha) \Leftrightarrow \alpha = 0.1 \Leftrightarrow \alpha/2 = 0.05$

The reliability coefficient: $Z_{1-\frac{\alpha}{2}} = z_{0.95} = 1.645$

A 90% confidence interval for $p_1 - p_2$ is:

$$(\hat{p}_1 - \hat{p}_2) \pm Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}}$$

$$(\hat{p}_1 - \hat{p}_2) \pm Z_{0.95} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}}$$

$$0.01 \pm 1.645 \sqrt{\frac{(0.05)(0.95)}{1500} + \frac{(0.04)(0.96)}{2000}}$$

$$0.01 \pm 0.01173$$

$$-0.0017 < p_1 - p_2 < 0.0217$$

We are 90% confident that $p_1 - p_2 \in (-0.0017, 0.0217)$.
Note: Since the confidence interval <mark>includes zero,</mark> we may conclude that the two population proportions are equal ($p_1 - p_2 = 0 \Leftrightarrow p_1 = p_2$). Therefore, we may conclude that the proportion of people having cancer is the same in both cities.

# CHAPTER 7: Using Sample Statistics To Test Hypotheses About Population Parameters:

In this chapter, we are interested in testing some hypotheses about the unknown population parameters.

## 7.1 Introduction:

Consider a population with some unknown parameter $\theta$. We are interested in testing (confirming or denying) some conjectures about $\theta$. For example, we might be interested in testing the conjecture that $\theta > \theta_o$, where $\theta_o$ is a given value.

- A hypothesis is a statement about one or more populations.
- A research hypothesis is the conjecture or supposition that motivates the research.
- A statistical hypothesis is a conjecture (or a statement) concerning the population which can be evaluated by appropriate statistical technique.
- For example, if $\theta$ is an unknown parameter of the population, we might be interested in testing the conjecture sating that $\theta \geq \theta_o$ against $\theta < \theta_o$ (for some specific value $\theta_o$).
- We usually test the null hypothesis ($H_o$) against the alternative (or the research) hypothesis ($H_1$ or $H_A$) by choosing one of the following situations:
  - (i)     $H_o: \theta = \theta_o$   against   $H_A: \theta \neq \theta_o$
  - (ii)    $H_o: \theta \geq \theta_o$   against   $H_A: \theta < \theta_o$
  - (iii)   $H_o: \theta \leq \theta_o$   against   $H_A: \theta > \theta_o$
- Equality sign must appear in the null hypothesis.
- $H_o$ is the null hypothesis and $H_A$ is the alternative hypothesis. ($H_o$ and $H_A$ are complement of each other)
- The null hypothesis ($H_o$) is also called "the hypothesis of no difference".
- The alternative hypothesis ($H_A$) is also called the research hypothesis.

- There are 4 possible situations in testing a statistical hypothesis:

|  |  | Condition of Null Hypothesis $H_o$ (Nature/reality) | |
|  |  | $H_o$ is true | $H_o$ is false |
| --- | --- | --- | --- |
| Possible Action (Decision) | Accepting $H_o$ | Correct Decision | Type II error $(\beta)$ |
|  | Rejecting $H_o$ | Type I error $(\alpha)$ | Correct Decision |

- There are two types of Errors:
  - Type I error = Rejecting $H_o$ when $H_o$ is true

    P(Type I error) = P(Rejecting Ho | Ho is true) = $\alpha$
  - Type II error = Accepting Ho when Ho is false

    P(Type II error) = P(Accepting Ho | Ho is false) = $\beta$

- The level of significance of the test is the probability of rejecting true $H_o$:

  $\alpha$ = P(Rejecting $H_o$ | $H_o$ is true) = P(Type I error)

- There are 2 types of alternative hypothesis:
  - One-sided alternative hypothesis:
    - $H_O: \theta \geq \theta_o$    against    $H_A: \theta < \theta_o$
    - $H_O: \theta \leq \theta_o$    against    $H_A: \theta > \theta_o$
  - Two-sided alternative hypothesis:
    - $H_O: \theta = \theta_o$    against    $H_A: \theta \neq \theta_o$

- We will use the terms "accepting" and "not rejecting" interchangeably. Also, we will use the terms "acceptance" and "nonrejection" interchangeably.
- We will use the terms "accept" and "fail to reject" interchangeably

**The Procedure of Testing $H_o$ (against $H_A$):**

The test procedure for rejecting $H_o$ (accepting $H_A$) or accepting $H_o$ (rejecting $H_A$) involves the following steps:

## 1- Determining Hypothesis

2. Determining a test statistic (T.S.)
We choose the appropriate test statistic based on the point estimator of the parameter.
The test statistic has the following form:

$$\text{Test statistic} = \frac{\textit{Estimate} - \textit{hypothesized parameter}}{\textit{Standard Error of the Estimate}}$$
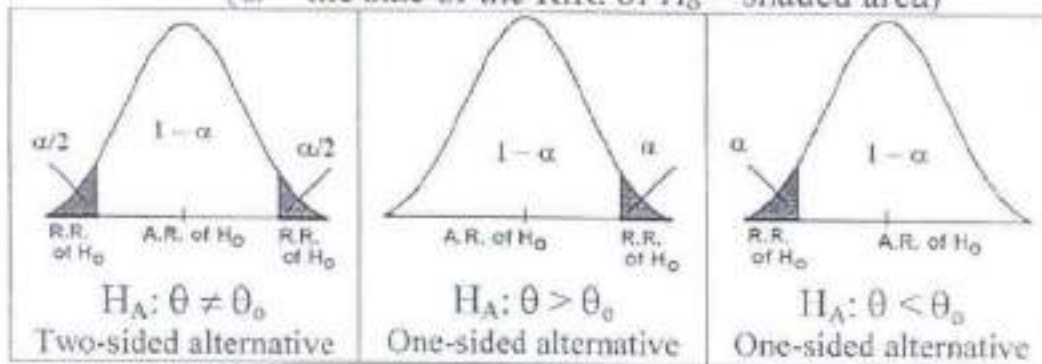
3. Determining the level of significance ($\alpha$):

$$\alpha = 0.01, 0.025, 0.05, 0.10$$

4. Determining the rejection region of $H_o$ (R.R.) and the acceptance region of $H_o$ (A.R.).
   The R.R. of $H_o$ depends on $H_A$ and $\alpha$:
   - $H_A$ determines the direction of the R.R. of $H_o$
   - $\alpha$ determines the size of the R.R. of $H_o$
     ($\alpha$ = the size of the R.R. of $H_o$ = shaded area)

| | | |
|---|---|---|
| $\alpha/2$  $1-\alpha$  $\alpha/2$<br>R.R. of $H_o$  A.R. of $H_o$  R.R. of $H_o$ | $1-\alpha$  $\alpha$<br>A.R. of $H_o$  R.R. of $H_o$ | $\alpha$  $1-\alpha$<br>R.R. of $H_o$  A.R. of $H_o$ |
| $H_A: \theta \neq \theta_o$ | $H_A: \theta > \theta_o$ | $H_A: \theta < \theta_o$ |
| Two-sided alternative | One-sided alternative | One-sided alternative |

5. Decision:
   We reject $H_o$ (and accept $H_A$) if the value of the test statistic (T.S.) belongs to the R.R. of $H_o$, and vice versa.

Notes:
1. The rejection region of $H_o$ (R.R.) is sometimes called "the critical region".
2. The values which separate the rejection region (R.R.) and the acceptance region (A.R.) are called "the critical values" or **Relibility Cofficient.**

### 7.2 Hypothesis Testing: A Single Population Mean ($\mu$):
Suppose that $X_1, X_2, \ldots, X_n$ is a random sample of size $n$ from a distribution (or population) with mean $\mu$ and variance $\sigma^2$.
We need to test some hypotheses (make some statistical inference) about the mean ($\mu$).

# Chapter 7 : Testing Hypothesis about population mean($\mu$):

| Hypothesis | $H_0: \mu = \mu_0$ $H_A: \mu \neq \mu_0$ | | $H_0: \mu \leq \mu_0$ $H_A: \mu > \mu_0$ | $H_0: \mu \geq \mu_0$ $H_A: \mu < \mu_0$ |
|---|---|---|---|---|
| First Case | $\sigma$ is known; Normal or [Non-normal Distribution($n \geq 30$)] | | | |
| Test Statistic (T.S.) | $Z = \dfrac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$ | | | |
| Rejection Region(R.R) & Acceptance Region(A.R) |  | |  |  |
| Reliability Coefficient | $-Z_{1-\alpha/2}$ or $Z_{1-\alpha/2}$ | | $Z_{1-\alpha}$ | $-Z_{1-\alpha}$ |
| Decision : Reject $H_0$ if the following condition satisfies | Reject $H_0$ (Accept $H_A$) at the significant level $\alpha$ if : | | | |
| | $Z > Z_{1-\alpha/2}$ Or $Z < -Z_{1-\alpha/2}$ | | $Z > Z_{1-\alpha}$ (one – Sided Test) | $Z < -Z_{1-\alpha}$ (one – Sided Test) |
| Second Case | $\sigma$ is unknown; Normal ,$n < 30$ (small) | | | |
| Test Statistic (T.S.) | $T = \dfrac{\bar{X} - \mu_0}{S/\sqrt{n}}, \quad df = n - 1$ | | | |
| Rejection Region(R.R) & Acceptance Region(A.R) |  | |  |  |
| Reliability Coefficient | $-t_{1-\alpha/2}$ or $t_{1-\alpha/2}$ | | $t_{1-\alpha}$ | $-t_{1-\alpha}$ |
| Decision : Reject $H_0$ if the Following condition satisfies | Reject $H_0$ (Accept $H_A$) at the significant level $\alpha$ if : | | | |
| | $T > t_{1-\alpha/2}$ Or $T < -t_{1-\alpha/2}$ (Two – Sided Test) | | $T > t_{1-\alpha}$ (one – Sided Test) | $T < -t_{1-\alpha}$ (one – Sided Test) |
| Special Case | $\sigma$ is unknown; Non-Normal ,$n \geq 30$ (Large) | | | |
| Test Statistic (T.S.) | $Z = \dfrac{\bar{X} - \mu_0}{S/\sqrt{n}}$ | | | |
| Rejection Region | Use the same R.R & A.R as in First Case(Z Case) | | | |

critical values

**Example: (first case: variance $\sigma^2$ is known)**

A random sample of 100 recorded deaths in the United States during the past year showed an average of 71.8 years. Assuming a population standard deviation of 8.9 year, does this seem to indicate that the mean life span today is greater than 70 years? Use a 0.05 level of significance.

**Solution:**

.$n$=100 (large),

$\sigma = 8.9$ ( $\sigma$ known)

$\bar{X}$=71.8, $\sigma$ =8.9 ( $\sigma$ is known )

$\mu$ =average (mean) life span

$\mu_o$= 70

$\alpha$=0.05


1) Hypotheses:

   $H_o$: $\mu \leq 70$  ($\mu_o$=70)

   $H_A$: $\mu > 70$   (research  hypothesis)

$H_o: \mu \le 70$    $(\mu_o = 70)$

$H_A: \mu > 70$    (research hypothesis)

Test statistics (T.S.) :

$$Z = \frac{\bar{X} - \mu_o}{\sigma/\sqrt{n}} = \frac{71.8 - 70}{8.9/\sqrt{100}} = 2.02$$

Level of significance:

$\alpha = 0.05$    1- $\alpha = 0.95$

Rejection Region of $H_o$ (R.R.): (critical region) is $Z_{1-\alpha} = Z_{0.95} = 1.645$

We should reject $H_o$ if: $Z(test) > Z_{1-\alpha}$
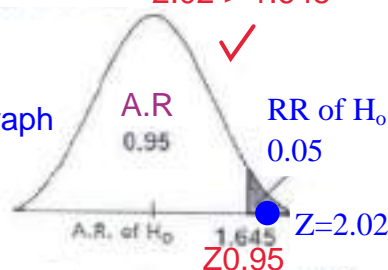
2.02 > 1.645     (condition satisfied)

✓

Decision : we reject $H_O$ accept $H_A$

Another solution :

From curve:

1) determine test value on graph

2) Z(test)=2.02 in R.R, then reject $H_o$

A.R 0.95

RR of $H_o$ 0.05

A.R. of $H_o$   1.645   Z=2.02

Z0.95

**Note: Using P- Value as a decision tool:**

P-value is the smallest value of $\alpha$ for which we can reject the null hypothesis $H_o$.

Calculating P-value:

* Calculating P-value depends on the alternative hypothesis $H_A$.

* Suppose that $Z_c = \frac{\bar{X} - \mu_o}{\sigma/\sqrt{n}}$ is the computed value of the test Statistic.

* The following table illustrates how to compute P-value, and how to use P-value for testing the null hypothesis:

$Z_c$ = Z (test)

| Alternative Hypothesis: | HA: $\mu \neq \mu_0$ | HA: $\mu > \mu_0$ | HA: $\mu < \mu_0$ |
|---|---|---|---|
| P - Value | $2 \times P(Z > \lvert Z_c \rvert)$ | $P(Z > Z_c)$ | $P(Z > -Z_c)$ |
| Significance Level = | $\alpha$ | | |
| Decision | Reject Ho if P-value $\leq \alpha$. | | |

## Example:
For the previous example, we have found that:

$$Z_c = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} = 2.02$$

The alternative hypothesis was HA: $\mu > 70$.

$P-Value = P(Z > Z_c)$
$= P(Z > 2.02) = 1 - P(Z < 2.02) = 1 - 0.9783 = 0.0217$

The level of significance was $\alpha = 0.05$.
Since P-value $\leq \alpha$, we reject $H_o$.

Decision :Reject $H_0$  If :

P-value < α
0.0217<0.05

✓

then
reject H0

## Example: (second case: variance is unknown)
The manager of a private clinic claims that the mean time of the patient-doctor visit in his clinic is 8 minutes. Test the hypothesis that $\mu = 8$ minutes against the alternative that $\mu \neq 8$ minutes if a random sample of 25 patient-doctor visits yielded a mean time of 7.8 minutes with a standard deviation of 0.5 minutes. It is assumed that the distribution of the time of this type of visits is normal. Use a 0.01 level of significance.

## Solution:
The distribution is normal.
$n = 25$ (small)
$\bar{X} = 7.8$

S=0.5 (sample standard deviation): $\sigma$ is unknown
$\mu$ = mean time of the visit , $\alpha = 0.01$

Hypotheses:

$H_o: \mu = 8 \quad (\mu_o = 8)$
$H_A: \mu \neq 8 \quad$ (research hypothesis)

Test statistics (T.S.):

$$T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} = \frac{7.8 - 8}{0.5/\sqrt{25}} = -2$$

$$df = v = n-1 = 25 - 1 = 24$$

Level of significance:

$\alpha = 0.01$ , $\alpha/2 = 0.005$ , $1 - \alpha/2 = 0.995$

Rejection Region of Ho (R.R.): (critical region)

$t_{1-\alpha/2} = t_{0.995} = 2.797$

We should reject Ho if:

$T < -t_{1-\alpha/2}$ or $T > + t_{1-\alpha/2}$

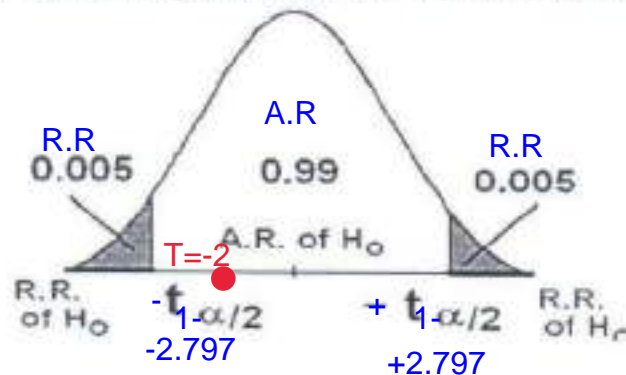-2 < -2.797   or   -2 > 2.797    (both conditions not satisfied)
✗ ✗                              Decision: Accept H0

Decision:

Since $T = -2 \in$ A.R., , we accept Ho: $\mu = 8$ at $\alpha = 0.01$ and reject
Ha: $\mu \neq 8$. Therefore, we conclude that the claim is correct.

Another solution:
From curve:

t (test)= -2 (In A.R)

Decision :Accept H0



A.R
R.R          0.99         R.R
0.005                     0.005
                A.R. of Ho
            T=-2
R.R.      $-t_{1-\alpha/2}$    $+t_{1-\alpha/2}$   R.R.
of Ho                                             of Hr
          -2.797              +2.797

Note:

For the case of non-normal population with unknown variance,
and when the sample size is large (n ≥30), we may use the
following test statistic:

$$Z = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$$

That is, we replace the population standard deviation (σ) by the
sample standard deviation (S), and we conduct the test as
described for the first case.

## 7.3 Hypothesis Testing: The Difference Between Two Population Means: (Independent Populations)

Suppose that we have two (independent) 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 $\mu_1$ and $\mu_2$, 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 $\mu_1-\mu_2$.

We wish to test some hypotheses comparing the population means.

**Hypotheses:**

We choose one of the following situations:

(i)   $H_o: \mu_1 = \mu_2$ against $H_A: \mu_1 \neq \mu_2$
(ii)  $H_o: \mu_1 \geq \mu_2$ against $H_A: \mu_1 < \mu_2$
(iii) $H_o: \mu_1 \leq \mu_2$ against $H_A: \mu_1 > \mu_2$

or equivalently,

(i)   $H_o: \mu_1-\mu_2 = M_o$ against $H_A: \mu_1 - \mu_2 \neq M_o$
(ii)  $H_o: \mu_1-\mu_2 \geq M_o$ against $H_A: \mu_1 - \mu_2 < M_o$
(iii) $H_o: \mu_1-\mu_2 \leq M_o$ against $H_A: \mu_1 - \mu_2 > M_o$

**Test Statistic:**

**(1) First Case:**

For normal populations (or non-normal populations with large sample sizes), and if $\sigma_1^2$ and $\sigma_2^2$ are known, then the test statistic is:

$$Z = \frac{\bar{X}_1 - \bar{X}_2 - M_o}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}} \sim N(0,1)$$

### (2) Second Case:

For normal populations, and if $\sigma_1^2$ and $\sigma_2^2$ are unknown but equal ($\sigma_1^2 = \sigma_2^2 = \sigma^2$), then the test statistic is:

$$T = \frac{\bar{X}_1 - \bar{X}_2 - M_o}{\sqrt{\dfrac{S_p^2}{n_1} + \dfrac{S_p^2}{n_2}}} \sim t(n_1+n_2-2)$$

where the pooled estimate of $\sigma^2$ is

$$S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1+n_2-2}$$

and the degrees of freedom of $S_p^2$ is df$= v = n_1+n_2-2$.

## Testing Hypothesis about difference between two population means ($\mu_1 - \mu_2$) : (Independent population)

| Hypothesis | $H_0 : \mu_1 - \mu_2 = M_0$ <br> $H_A : \mu_1 - \mu_2 \neq M_0$ | $H_0 : \mu_1 - \mu_2 \leq M_0$ <br> $H_A : \mu_1 - \mu_2 > M_0$ | $H_0 : \mu_1 - \mu_2 \geq M_0$ <br> $H_A : \mu_1 - \mu_2 < M_0$ |
|---|---|---|---|
| | | $\sigma_1^2$ , $\sigma_2^2$ are known + Normal or Non-Normal with large samples | |
| Test Statistic (T.S.) | | $Z = \dfrac{\bar{X}_1 - \bar{X}_2 - M_0}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}} \sim N(0,1)$ | |
| Rejection Region(R.R) & Acceptance Region(A.R) |  |  |  |
| Reliability Coefficient | $-Z_{1-\alpha/2}$ or $Z_{1-\alpha/2}$ | $Z_{1-\alpha}$ | $-Z_{1-\alpha}$ |
| Decision : Reject $H_0$ if the following condition satisfies | Reject $H_0$ (Accept $H_A$) at the significant level $\alpha$ if : | | |
| | $Z > Z_{1-\alpha/2}$ <br> Or $Z < -Z_{1-\alpha/2}$ | $Z > Z_{1-\alpha}$ <br> (one – Sided Test) | $Z < -Z_{1-\alpha}$ <br> (one – Sided Test) |
| | | $\sigma_1^2$ , $\sigma_2^2$ are unknown but equal $(\sigma_1^2 = \sigma_2^2 = \sigma^2)$ + Normal | |
| Test Statistic (T.S.) | $T = \dfrac{\bar{X}_1 - \bar{X}_2 - M_0}{\sqrt{\dfrac{S_p^2}{n_1} + \dfrac{S_p^2}{n_2}}}$, | $S_p^2 = \dfrac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$ <br> df=$n_1$+$n_2$-2 | |
| Rejection Region(R.R) & Acceptance Region(A.R) |  |  |  |
| Reliability Coefficient | $-t_{1-\alpha/2}$ or $t_{1-\alpha/2}$ | $t_{1-\alpha}$ | $-t_{1-\alpha}$ |
| Decision : Reject $H_0$ if the Following condition satisfies | Reject $H_0$ (Accept $H_A$) at the significant level $\alpha$ if : | | |
| | $T > t_{1-\alpha/2}$ <br> Or $T < -t_{1-\alpha/2}$ <br> (Two –sided test) | $T > t_{1-\alpha}$ <br> (one – Sided Test) | $T < -t_{1-\alpha}$ <br> (one – Sided Test) |

**Example:** ($\sigma_1^2$, $\sigma_2^2$ are known)

Researchers wish to know if the data they have collected provide sufficient evidence to indicate the difference in mean serum uric acid levels between individuals with Down's syndrome and normal individuals. The data consist of serum uric acid on 12 individuals with Down's syndrome and 15 normal individuals. The sample means are

$\bar{X}_1 = 4.5$ mg/100ml
$\bar{X}_2 = 3.4$ mg/100ml

Assume the populations are normal with variances
$$\sigma_1^2 = 1$$
$$\sigma_2^2 = 1.5$$

. Use significance level $\alpha = 0.05$.

## Solution:

$\mu_1$ = mean serum uric acid levels for the individuals with Down's syndrome.

$\mu_2$ = mean serum uric acid levels for the normal individuals.

| | | |
|---|---|---|
| $n_1 = 12$ | $\bar{X}_1 = 4.5$ | $\sigma_1^2 = 1$ |
| $n_2 = 15$ | $\bar{X}_2 = 3.4$ | $\sigma_2^2 = 1.5$ |

### Hypotheses:

$H_o: \mu_1 = \mu_2$ against $H_A: \mu_1 \neq \mu_2$

or

$H_o: \mu_1 - \mu_2 = 0$ against $H_A: \mu_1 - \mu_2 \neq 0$

### Calculation:

1) $\alpha = 0.05$

2) $1 - \alpha/2 = 1 - 0.05/2 = 0.975$

3) $Z_{1-\alpha/2} = Z_{0.975} = 1.96$

### Test Statistic (T.S.):

$$Z = \frac{\bar{X}_1 - \bar{X}_2 - M_o}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = \frac{4.5 - 3.4 - 0}{\sqrt{\frac{1}{12} + \frac{1.5}{15}}} = 2.569$$

Decision : Reject $H_0$ If

$Z(test) > Z_{0.975}$ or $Z < -Z_{0.975}$

2.569 > 1.96

✓

first condition satisfied
no need to check the second one

Another solution:
from curve
Since
Z(test)=2.569 in R.R
Decision:
Reject $H_0$

A.R
0.95

R.R
0.025

R.R
0.025

R.R
of Ho

Decision : Reject $H_0$

Z=2.569
in (R.R)

$-Z_{0.975}$        $Z_{0.975}$

### Decision:

Since $Z = 2.569 \in$ R.R. we reject $H_o: \mu_1 = \mu_2$ and we accept (do not reject) $H_A: \mu_1 \neq \mu_2$ at $\alpha = 0.05$. Therefore, we conclude that the two population means are not equal.

### Notes:

1. We can easily show that a 95% confidence interval for $(\mu_1 - \mu_2)$ is (0.26, 1.94), that is:

$$0.26 < \mu_1 - \mu_2 < 1.94$$

| | | |
|---|---|---|
| King Saud University | 137 | Dr. Abdullah Al-Shiha |

Since this interval does not include 0, we say that 0 is not a candidate for the difference between the population means ($\mu_1 - \mu_2$), and we conclude that $\mu_1 - \mu_2 \neq 0$, i.e., $\mu_1 \neq \mu_2$. Thus we arrive **Another solution:** at the same conclusion by means of a confidence interval.
**by p-value:**

2. $P - Value = 2 \times P(Z > |Z_c|)$

$$= 2P(Z > 2.57) = 2[1 - P(Z < 2.57)] = 2(1 - 0.9949) = 0.0102$$

The level of significance was $\alpha = 0.05$.       **Reject H$_0$** If p-value < α

Since P-value $< \alpha$, we reject H$_0$.                     0.0102< 0.05

                                                                    ✓

                                                    **Decision: Reject Ho**

**Example:** ($\sigma_1^2 = \sigma_2^2 = \sigma^2$ is unknown)

An experiment was performed to compare the abrasive wear of two different materials used in making artificial teeth. 12 pieces of material 1 were tested by exposing each piece to a machine measuring wear. 10 pieces of material 2 were similarly tested. In each case, the depth of wear was observed. The samples of material 1 gave an average wear of 85 units with a sample standard deviation of 4, while the samples of materials 2 gave an average wear of 81 and a sample standard deviation of 5. Can we conclude at the 0.05 level of significance that the mean abrasive wear of material 1 is greater than that of material 2? Assume normal populations with equal variances.

**Solution:**

| Material 1 | material 2 |
|---|---|
| $n_1 = 12$ | $n_2 = 10$ |
| $\overline{X}_1 = 85$ | $X_2 = 81$ |
| $S_1 = 4$ | $S_2 = 5$ |

Hypotheses:
   H$_0$: $\mu_1 \leq \mu_2$
   H$_A$: $\mu_1 > \mu_2$
Or equivalently,
   H$_0$: $\mu_1 - \mu_2 \leq 0$
   H$_A$: $\mu_1 - \mu_2 > 0$
Calculation:
   $\alpha = 0.05$

0.95        0.05

A.R. of Ho    t $_{0.95}$    R.R. of Ho

= 1.725

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

$$= \frac{(12 - 1)4^2 + (10 - 1)5^2}{12 + 10 - 2} = 20.05$$

## Reliability Coefficient:

$$df = v = 12 + 10 - 2 = 20$$

$$\alpha = 0.05 \text{ ------- } 1 - \alpha = 0.95 \text{ ------- } t_{1-\alpha} = t_{0.95} = 1.725$$

## Test Statistic (T.S.):

$$T = \frac{\bar{X}_1 - \bar{X}_2 - M_0}{\sqrt{\frac{S_p^2}{n_1} + \frac{S_p^2}{n_2}}} = \frac{85 - 81 - 0}{\sqrt{\frac{20.05}{12} + \frac{20.05}{10}}} = 2.09$$

Reject $H_o$

If T(test)> $t_{1-\alpha}$

2.09 > 1.725

✓

Decision : Reject $H_o$

## Decision:

Since T= 2.09 ∈ R.R. (T= 2.09 > t $_{0.95}$ = 1.725), we reject H$_o$ and we accept **H$_A$: $\mu_1 - \mu_2 > 0$ (H$_A$: $\mu_1 > \mu_2$)** at $\alpha$ =0.05. Therefore, we conclude that the mean abrasive wear of material 1 is greater than that of material 2.

## 7.4 Paired Comparisons:

### Paired T-Test :

- In this section, we are interested in comparing the means of two related (non-independent/dependent) normal populations.
- In other words, we wish to make statistical inference for the difference between the means of two related normal populations.
- Paired t-Test concerns about testing the equality of the means of two related normal populations.

### Examples of related populations are:

1. Height of the father and height of his son.
2. Mark of the student in MATH and his mark in STAT.
3. Pulse rate of the patient before and after the medical treatment.
4. Hemoglobin level of the patient before and after the medical treatment.

---

### Test procedure:
#### Let

      X: Values of the first population

      Y: Values of the Second population

      D: Values of X – Values of Y

#### Means :

      $\mu_1$ = Mean of the first population

      $\mu_2$ = Mean of the Second population

      $\mu_D$ = Mean of X – Mean of Y    ($\mu_D = \mu_1 - \mu_2$)

# Confident Interval and Testing Hypothesis about difference between two population means ($\mu_D = \mu_1 - \mu_2$) : (Dependent/Related population)

| Calculate the following Quantities | • The difference (D-observation): $D_i = X_i - Y_i$ , $i=1,2,3,4,\ldots\ldots n$<br>• Sample mean of the D-Observations : $\bar{D} = \frac{\sum_{i=1}^{n} D_i}{n}$<br>• Sample Variance $S_D^2 = \frac{\sum_{i=1}^{n}(D_i - \bar{D})^2}{n-1}$<br>• Sample Standard Deviation $S_D = \sqrt{S_D^2}$ | | |
|---|---|---|---|
| | **Confident Interval for $\mu_D = \mu_1 - \mu_2$** | | |
| 100(1-α)% Confident interval for $\mu_D$ | $\bar{D} \pm t_{1-\frac{\alpha}{2}} \frac{S_D}{\sqrt{n}}$ , $df = n-1$ | | |
| | **Testing Hypothesis for $\mu_0 = \mu_1 - \mu_2$** | | |
| Hypothesis | $H_0:\mu_1 - \mu_2 = M_O$<br>$H_A:\mu_1 - \mu_2 \neq M_O$<br>Or<br>$H_0:\mu_D = M_O$ vs $H_A:\mu_D \neq M_O$ | $H_0:\mu_1 - \mu_2 \leq M_O$<br>$H_A:\mu_1 - \mu_2 > M_O$<br>Or<br>$H_0:\mu_D \leq M_O$ vs $H_A:\mu_D > M_O$ | $H_0:\mu_1 - \mu_2 \geq M_O$<br>$H_A:\mu_1 - \mu_2 < M_O$<br>Or<br>$H_0:\mu_D \geq M_O$ vs $H_A:\mu_D < M_O$ |
| Test Statistic (T.S.) | $T = \frac{\bar{D} - M_O}{S_D/\sqrt{n}}$ , $df = v = n-1$ | | |
| Rejection Region(R.R) & Acceptance Region(A.R) |  |  |  |
| Reliability Coefficient | $-t_{1-\alpha/2}$ or $t_{1-\alpha/2}$ | $t_{1-\alpha}$ | $-t_{1-\alpha}$ |
| Decision : Reject $H_0$ if the Following condition satisfies | Reject $H_0$ (Accept $H_A$) at the significant level α if : | | |
| | $T > t_{1-\alpha/2}$<br>Or $T < -t_{1-\alpha/2}$<br>(Two-sided test) | $T > t_{1-\alpha}$<br>(one-Sided Test) | $T < -t_{1-\alpha}$<br>(one-Sided Test) |

## Example:

Suppose that we are interested in studying the effectiveness of a
certain diet program on ten individual . Let the random variables X and Y given as
following table :

| Individual(i) | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Weight before (X.) | 86.6 | 80.2 | 91.5 | 80.6 | 82.3 | 81.9 | 88.4 | 85.3 | 83.1 | 82.1 |
| Weight After (Y.) | 79.7 | 85.9 | 81.7 | 82.5 | 77.9 | 85.8 | 81.3 | 74.7 | 68.3 | 69.7 |

Find :

1) A 95% Confident Interval for the difference between the mean of weights before the
   diet program ($\mu_1$) and the mean of weights after the diet program ($\mu_2$).
   [ $\mu_D = \mu_1 - \mu_2$]

2) Does the data provide sufficient evidence to allow us to conclude that the diet is
   good? Use $\alpha = 0.05$ and assume population is normal .

## Solution :

1-st population (X) = the weight of the individual before the diet program.
2-nd population (Y)= the weight of the same individual after the diet program.

We assume that the distributions of these random variables are normal with
means $\mu_1$ and $\mu_2$ , respectively.

These two variables are related (dependent/non-independent)because they are measured
on the same individual.

Calculate mean, standard deviation by calculator

$\overline{D}$ , $S_D$

| i | Xi | Yi | $D_i = Xi - Y_i$ |
|---|---|---|---|
| 1 | 86.6 | 79.7 | 6.9 |
| 2 | 80.2 | 85.9 | -5.7 |
| 3 | 91.5 | 81.7 | 9.8 |
| 4 | 80.6 | 82.5 | -1.9 |
| 5 | 82.3 | 77.9 | 4.4 |
| 6 | 81.9 | 85.8 | -3.9 |
| 7 | 88.4 | 81.3 | 7.1 |
| 8 | 85.3 | 74.7 | 10.6 |
| 9 | 83.1 | 68.3 | 14.8 |
| 10 | 82.1 | 69.7 | 12.4 |
| sum | $\Sigma x = 842$ | $\Sigma y = 787.5$ | $\Sigma_D = 54.5$ |

<u>First, we need to calculate :</u>(By calculator )

Sample Mean:
$$\bar{D} = \frac{\sum_{i=1}^{n} D_i}{n} = \frac{54.5}{10} = 5.45$$

Sample Variance :
$$S_D^2 = \frac{\sum_{i=1}^{n}(D_i - \bar{D})^2}{n-1} = \frac{(6.9-5.45)^2 + (-5.7 - 5.45)^2 + \ldots + (12.4 - 5.45)^2}{10-1} = 50.33$$

Sample Standard Deviation : $S_D = \sqrt{S_D^2} = \sqrt{50.33} = 7.09$

Reliability Coefficient : $t_{1-\alpha/2}$ :

$\alpha = 0.05$ ------ $1 - 0.05/2 = 1 - 0.025 = 0.975$  (df $= 10 - 1 = 9$)

$t_{1-\alpha/2} = t_{0.975} = 2.262$

Then 95% Confident Interval for $\mu_D = \mu_1 - \mu_2$

$$\bar{D} \pm t_{1-\frac{\alpha}{2}} \frac{S_D}{\sqrt{n}}$$

$$5.45 \pm 2.262 \frac{7.09}{\sqrt{10}}$$

$$5.45 \pm 5.0715$$

$$( 5.45 - 5.0715 , \ 5.45 + 5.0715 )$$

$$(0.38 , 10.52)$$

$$0.38 < \mu_D < 10.52$$

2)Does the data provide sufficient evidence to allow us to conclude that the diet is good? Use $\alpha = 0.05$ and assume population is normal .

Diet is good means --- weight after will be less than weight befor.

## Solution:

$\mu_1$= Mean of the first population
$\mu_2$= Mean of the second population
$\mu_D$ =Mean of X – Mean of Y    ($\mu_D = \mu_1 - \mu_2$ )

Hypothesis :

$$H_0: \mu_1 \leq \mu_2 \quad vs \quad H_A: \mu_1 > \mu_2$$

or    $$H_0: \mu_1 - \mu_2 \leq 0 \quad vs \quad H_A: \mu_1 - \mu_2 > 0$$

or    $$H_0: \mu_D \leq 0 \quad vs \quad H_A: \mu_D > 0$$

## Test Statistic:

$$\bar{D} = 5.45, \quad S_D = 7.09, n = 10$$

$$T = \frac{\overbrace{\bar{D}-M_O}}{\frac{S_D}{\sqrt{n}}} = \frac{\overbrace{5.45-0}}{\frac{7.09}{\sqrt{10}}} = 2.43$$

## Rejection Region(R.R):

$\alpha = 0.05$ ------ $1-\alpha = 0.95$ ----- $t_{1-\alpha} = t_{0.95} = 1.833$  (df =n-1 =9)

Reject $H_0$ if    $T > t_{1-\alpha}$

$$2.45 > 1.833 \quad \text{(condition satisfied )}$$

Then reject $H_0$ and accept $H_A$: $\mu_1 > \mu_2$

So, we have a good diet program .

## 7.5 Hypothesis Testing: A Single Population Proportion (p):

In this section, we are interested in testing some hypotheses about the population proportion (p).



**Recall:**

- $p$ = Population proportion of elements of Type $A$ in the population

$$p = \frac{no.\ of\ elements\ of\ type\ A\ in\ the\ population}{Total\ no.\ of\ elements\ in\ the\ population}$$

$$p = \frac{A}{N} \qquad (N = population\ size)$$

- $n$ = sample size
- $X$ = no. of elements of type $A$ in the sample of size $n$.
- $\hat{p}$ = Sample proportion elements of Type $A$ in the sample

$$\hat{p} = \frac{no.\ of\ elements\ of\ type\ A\ in\ the\ sample}{no.\ of\ elements\ in\ the\ sample}$$

$$\hat{p} = \frac{X}{n} \qquad (n = sample\ size = no.\ of\ elements\ in\ the\ sample)$$

- $\hat{p}$ is a "good" point estimate for $p$.
- For large $n$, $(n \geq 30,\ np > 5)$, we have

# Test Procedure:( $P_0$ is known number)

| Hypothesis | $H_0: P = P_0$ <br> $H_A: P \neq P_0$ | $H_0: P \leq P_0$ <br> $H_A: P > P_0$ | $H_c: P \geq P_0$ <br> $H_A: P < P_0$ |
|---|---|---|---|
| Test Statistic (T.S.) | $$Z = \dfrac{\hat{p} - P_0}{\sqrt{\dfrac{p_0 q_0}{n}}} \quad , \qquad q_0 = 1 - p_0$$ | | |
| Rejection Region(R.R) & Acceptance Region(A.R) |  |  |  |
| Reliability Coefficient | $-Z_{1-\alpha/2}$ or $Z_{1-\alpha/2}$ | $Z_{1-\alpha}$ | $-Z_{1-\alpha}$ |
| Decision : Reject $H_0$ if the following condition satisfies | Reject $H_0$ (Accept $H_A$) at the significant level $\alpha$ if : | | |
| | $Z > Z_{1-\alpha/2}$ <br> Or $Z < - Z_{1-\alpha/2}$ | $Z > Z_{1-\alpha}$ <br> (one – Sided Test) | $Z < - Z_{1-\alpha}$ <br> (one – Sided Test) |

### Example:

A researcher was interested in the proportion of females in the population of all patients visiting a certain clinic. The researcher claims that 70% of all patients in this population are females. Would you agree with this claim if a random survey shows that 24 out of 45 patients are females? Use a 0.10 level of significance.

### Solution:

p = Proportion of female in the population.

.$n = 45$ (large)

X= no. of female in the sample = 24

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

$$\hat{p} = \frac{X}{n} = \frac{24}{45} = 0.5333$$

$$p_o = \frac{70}{100} = 0.7$$

$$\alpha = 0.10$$

Hypotheses:

$$H_o: p = 0.7 \quad (p_o = 0.7)$$
$$H_A: p \neq 0.7$$

Level of significance:

$$\alpha = 0.10 \qquad 1 - \alpha/2 = 1 - 0.10/2 = 0.95$$

Test Statistic (T.S.):

$$Z = \frac{\hat{p} - p_o}{\sqrt{\frac{p_o(1-p_o)}{n}}}$$

$$= \frac{0.5333 - 0.70}{\sqrt{\frac{(0.7)(0.3)}{45}}} = -2.44$$

Rejection Region of $H_o$ (R.R.):
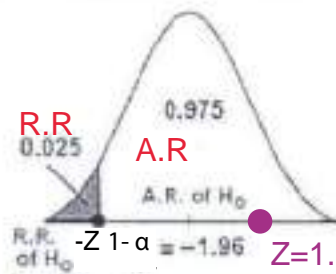
Critical values:

$$Z_{1-\alpha/2} = Z_{0.95} = 1.645$$

We reject $H_0$ If:

$$Z(test) > Z_{0.95} \quad or \quad Z(test) < -Z_{0.95}$$
$$-2.44 > 1.645 \qquad or \qquad -2.44 < -1.645$$
✗ ✓

Since one of the conditions is valid then,

Decision : Reject $H_0$

From curve

Since Z(test)= -2.44 in R.R

Decision :

Reject $H_0$

R.R 0.05   A.R 0.9   R.R 0.05

Z=-2.44

A.R. of $H_a$

R.R. of $H_o$   -Z 1-α/2   Z 1-α/2   R.R. of $H_a$

-1.645   1.645

- Z 1- α/2 = - Z 0.95 = - 1.645

Decision:

Since Z= -2.44 ∈Rejection Region of $H_o$ (R.R), we reject

$H_o:p=0.7$ and accept $H_A:p \neq 0.7$ at $\alpha=0.1$. Therefore, we do not agree with the claim stating that 70% of the patients in this population are females.

## Example:

In a study on the fear of dental care in a certain city, a survey showed that 60 out of 200 adults said that they would hesitate to take a dental appointment due to fear. Test whether the proportion of adults in this city who hesitate to take dental appointment is less than 0.25. Use a level of significance of 0.025.

## Solution:

p = Proportion of adults in the city who hesitate to take a dental appointment.

$n = 200$ (large)

X= no. of adults who hesitate in the sample = 60

$\hat{p}$ = proportion of adults who hesitate in the sample

$$\hat{p} = \frac{X}{n} = \frac{60}{200} = 0.3$$

$p_0 = 0.25$

$\alpha = 0.025$

Hypotheses:

$H_o: p \geq 0.25$    ($p_o=0.25$)

$H_A: p < 0.25$    (research hypothesis)

Level of significance:

$\alpha = 0.025$

Test Statistic (T.S.):

$$Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0.3 - 0.25}{\sqrt{\frac{(0.25)(0.75)}{200}}} = 1.633$$

Rejection Region of $H_o$ (R.R.):

Critical value:    $Z_{1-\alpha} = Z_{0.975} = 1.96$

Critical Region:

We reject $H_o$ if:

$Z < -Z_{1-\alpha}$

$1.633 < -1.96$

$\times$ $\times$

Accept $H_0$ (condition not satisfy)

Another solution:
From curve :

Since Z(test)=1.633  in A.R

Decision:

Accept H₀



R.R
0.025

0.975

A.R

A.R. of H₀

R.R
of H₀    -Z 1- α =−1.96    Z=1.633 (in A.R)

**Decision:**

Since $Z=1.633 \in$ Acceptance Region of $H_o$ (A.R.), we accept (do not reject) $H_o: p \geq 0.25$ and we reject $H_A: p < 0.25$ at $\alpha = 0.025$. Therefore, we do not agree with claim stating that the proportion of adults in this city who hesitate to take dental appointment is less than 0.25.

## 7.6 Hypothesis Testing: The Difference Between Two Population Proportions $(p_1 - p_2)$:

In this section, we are interested in testing some hypotheses about the difference between two population proportions $(p_1 - p_2)$.



Suppose that we have two populations:

- $p_1$ = population proportion of the 1-st population.
- $p_2$ = population proportion of 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.
- $\hat{p}_1 = \dfrac{X_1}{n_1}$ = the sample proportion of the 1-st sample
- $\hat{p}_2 = \dfrac{X_2}{n_2}$ = the sample proportion of the 2-nd sample
- The sampling distribution of $\hat{p}_1 - \hat{p}_2$ is used to make inferences about $p_1 - p_2$.
- For large $n_1$ and $n_2$, we have

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\dfrac{p_1 q_1}{n_1} + \dfrac{p_2 q_2}{n_2}}} \sim N(0,1) \quad \text{(Approximately)}$$

- $q = 1 - p$

## Hypotheses:

We choose one of the following situations:

    (i)   $H_o$: $p_1 = p_2$   against   $H_A$: $p_1 \neq p_2$

    (ii)  $H_o$: $p_1 \geq p_2$   against   $H_A$: $p_1 < p_2$

    (iii) $H_o$: $p_1 \leq p_2$   against   $H_A$: $p_1 > p_2$

or equivalently,

    (i)   $H_o$: $p_1 - p_2 = 0$   against   $H_A$: $p_1 - p_2 \neq 0$

    (ii)  $H_o$: $p_1 - p_2 \geq 0$   against   $H_A$: $p_1 - p_2 < 0$

    (iii) $H_o$: $p_1 - p_2 \leq 0$   against   $H_A$: $p_1 - p_2 > 0$

Note, under the assumption of the equality of the two population proportions ($H_o$: $p_1 = p_2 = p$), the pooled estimate of the common proportion $p$ is:

$$\overline{p} = \frac{X_1 + X_2}{n_1 + n_2} \qquad\qquad (\overline{q} = 1 - \overline{p})$$

The test statistic (T.S.) is

$$Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\dfrac{\bar{p}\bar{q}}{n_1} + \dfrac{\bar{p}\bar{q}}{n_2}}} \sim N(0,1)$$

## Test Procedure:

| Hypothesis | $H_0: P_1 - P_2 = 0$ <br> $H_A: P_1 - P_2 \neq 0$ | $H_0: P_1 - P_2 \leq 0$ <br> $H_A: P_1 - P_2 > 0$ | $H_0: P_1 - P_2 \geq 0$ <br> $H_A: P_1 - P_2 < 0$ |
|---|---|---|---|
| Test Statistic (T.S.) | $Z = \dfrac{\hat{p}_1 - \hat{p}_2}{\sqrt{\dfrac{\bar{p}\bar{q}}{n_1} + \dfrac{\bar{p}\bar{q}}{n_2}}}$ , Pooled proportion: $\bar{p} = \dfrac{x_1 + x_2}{n_1 + n_2}$ <br> where $\bar{q} = 1 - \bar{p}$ | | |
| Rejection Region(R.R) & Acceptance Region(A.R) |  |  |  |
| Reliability Coefficient | $-Z_{1-\alpha/2}$ or $Z_{1-\alpha/2}$ | $Z_{1-\alpha}$ | $-Z_{1-\alpha}$ |
| Decision : Reject $H_0$ if the following condition satisfies | Reject $H_0$ (Accept $H_A$) at the significant level $\alpha$ if : | | |
| | $Z > Z_{1-\alpha/2}$ <br> Or $Z < -Z_{1-\alpha/2}$ | $Z > Z_{1-\alpha}$ <br> (one – Sided Test) | $Z < -Z_{1-\alpha}$ <br> (one – Sided Test) |

## Example:

In a study about the obesity (overweight), a researcher was interested in comparing the proportion of obesity between males and females. The researcher has obtained a random sample of 150 males and another independent random sample of 200 females. The following results were obtained from this study.

| | n | Number of obese people(X) |
|---|---|---|
| Males | 150 | 21 |
| Females | 200 | 48 |

Can we conclude from these data that there is a difference between the proportion of obese males and proportion of obese females?
Use $\alpha = 0.05$ and assume that the two population proportions are equal.

## Solution

$p_1$ = population proportion of obese males
$p_2$ = population proportion of obese females
$\hat{p}_1$ = sample proportion of obese males
$\hat{p}_2$ = sample proportion of obese females

| Males | Females |
|---|---|
| $n_1 = 150$ | $n_2 = 200$ |
| $X_1 = 21$ | $X_2 = 48$ |

$$\hat{p}_1 = \frac{X_1}{n_1} = \frac{21}{150} = 0.14 \qquad \hat{p}_2 = \frac{X_2}{n_2} = \frac{48}{200} = 0.24$$

The pooled estimate of the common proportion $p$ is:

$$\bar{p} = \frac{X_1 + X_2}{n_1 + n_2} = \frac{21 + 48}{150 + 200} = 0.197$$

Hypotheses:

$$H_o: p_1 = p_2$$
$$H_A: p_1 \neq p_2$$

or

$$H_o: p_1 - p_2 = 0$$
$$H_A: p_1 - p_2 \neq 0$$

Level of significance: $\alpha = 0.05$    1- α/2=1 -0.05/2=0.975

Test Statistic (T.S.):

$$Z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{\bar{p}(1-\bar{p})}{n_1} + \frac{\bar{p}(1-\bar{p})}{n_2}}} = \frac{(0.14 - 0.24)}{\sqrt{\frac{0.197 \times 0.803}{150} + \frac{0.197 \times 0.803}{200}}} = -2.328$$

Rejection Region (R.R.) of $H_o$:
Critical values:

$$Z_{1-\alpha/2} = Z_{0.975} = 1.96$$

Critical region:
Reject $H_o$ if: $Z_{test} < -1.96$ or $Z_{test} > 1.96$

-2.328<-1.96

✓

Decision: Reject $H_0$   (Since one of the conditions satisfied)

R.R
0.025

A.R
0.95

R.R
0.025

Another solution:
From curve
Since Z(test) in R.R

Z=-2.328
(in R.R)

A.R. of H$_0$

Decision:
Reject H$_0$

R.R.
of H$_0$

-Z$_{1-\alpha/2}$

Z$_{1-\alpha/2}$

R.R.
of H$_0$

Reject H$_0$

-1.96

1.96

Decision:

Since $Z = -2.328 \in R.R.$, we reject H$_o$: $p_1 = p_2$ and accept H$_A$: $p_1 \neq p_2$ at $\alpha=0.05$. Therefore, we conclude that there is a difference between the proportion of obese males and the proportion of obese females. Additionally, since, $\hat{p}_1 = 0.14 < \hat{p}_2 = 0.24$, we may conclude that the proportion of obesity for females is larger than that for males.

# TABLE OF CONTENTS

# TABLE OF CONTENTS

# TABLE OF CONTENTS