

# Descriptive Statistics

## Unit - I

### Statistics:

The science of collecting, organising, presenting, analyzing and interpreting data to assist in making more effective decisions.

statistical analysis - Used to Manipulate, summarise and Investigate data, so that Useful decision making information results.

### Why statistics Matters:

statistical procedures are used to advance our knowledge of ourselves and the World around us

groups of people, gathering scores, analyzing

statistics helps to strengthen your critical thinking skills and reasoning abilities

The sun will explode in less than 6 years (9/19/2002) - to evaluate the Validity of such assertions to see if they stand up to scientific scrutiny

statistics enable you to Understand and research results in the professional journals of your area of specialization.

Knowledgeable and competent professional in your area of study.

### STATISTICS AND OTHER STATISTICAL TERMS:

Research, a systematic inquiry in search of Knowledge, Involves use of statistics.

Set of scores, referred to as data.

Before the scores have undergone any type of statistical transformation or analysis, they are called raw scores.

A population consists of all elements - individuals, items or objects - whose characteristics are being studied. The population that is being studied is also called target population.

A portion of population selected for study is referred to as sample.

Usually populations are so large that a researcher cannot examine the entire group.

Therefore a sample is selected to represent the population in a research study.

The goal is to use the results obtained from the sample to help answer questions about the population.

A sample drawn in such a way that each element of the population has a chance of being selected is called random sample.

Samples are normally drawn from only the portion of the population that is accessible referred to as accessible population.

**Descriptive Statistics** Sum up and Condense a set of raw scores so that overall trends in the data become apparent. Percentages and averages are examples of descriptive statistics.

**Inferential Statistics** Involves predicting characteristics of a population based on data obtained from a sample.

# Scales of Measurements :

## Nominal Scales:

The least specific Measurement scale only (classify the observations into different categories)

eg: colours, types of trees, Religion

These Variables have no quantitative Values.

## Ordinal Scales:

In addition to classify the observations, this scale also (permits Ordering and ranking of these observations)

eg: first, second, third place winners

## Interval Scales:

With Interval scales, there is [equal distance between units] on the scale.

eg. Temperature in degrees is Measured on the Interval scale.

In the Interval scale there is an zero point.

## Ratio Scale:

In the ratio scale, the real absolute zero point is known.

X	Y	$X^2$	$XY$	$\Sigma(X-3)$	$Y+2$	$(Y+2)^2$
6	4	36	24	3	6	36
5	1	25	5	2	3	9
9	2	81	18	6	4	16
8	1	64	8	5	3	9
		<u>206</u>	<u>55</u>	<u>16</u>		<u>70</u>



## Organising data Using tables & graphs

Unorganised leadership Course of 35 Managers

54, 43, 48, 50, 52, 44, 46, 44, 51, 42, 50, 46, 51, 55,  
57, 48, 53, 51, 46, 52, 50, 45, 55, 49, 48, 50, 49,  
51, 48, 43 Frequency Vs cumulative frequency Vs relative  
frequency.

X	Tally	F
57	I	1
56	-	0
55	II	2
54	I	1
53	I	1
52	II	2
51	IIII	4
50	IIII	4
49	II	2
48	IIII	4
47	-	0
46	III	3
45	I	1
44	II	2
43	II	2
42	I	1
		<hr/> 30

### Relative Frequency Distribution:

The frequency refers to "how many times a particular score occurs", the relative frequency refers to the "portion of the time score occurs".

$$\text{Relative frequency} = \frac{f}{N}$$

$N \rightarrow$  No. of data in population

$$N = 35$$

X	f	Relative Frequency
10	4	$4/26 = 0.15 \Rightarrow 15\%$
9	2	$2/26 = 0.08 \Rightarrow 8\%$
8	6	$6/26 = 0.23 \Rightarrow 23\%$
7	4	$4/26 = 0.15 \Rightarrow 15\%$
6	5	$5/26 = 0.19 \Rightarrow 19\%$
5	3	$3/26 = 0.12 \Rightarrow 12\%$
4	0	$0/26 = 0 \Rightarrow 0\%$
3	2	$2/26 = 0.08 \Rightarrow 8\%$
	N = 26	1.00

### Cumulative Frequency Distribution:

It indicates the frequency of scores that fall at or below a particular score value. This type of table is useful if we want to know how many people scored below certain value on a test.

X	f	Cf	PR (%)
10	4	26	100
9	2	22	84.62
8	6	20	76.92
7	4	14	53.85
6	5	10	38.46
5	3	5	19.23
4	0	2	7.69
3	2	2	7.69

## Percentile Rank:

Percentile Rank gives a bit more information than frequency by indicating the percentage of score that fall at below a given score in a distribution.

$$\text{Percentile Rank} = \frac{cf}{N} \times 100$$

29	83	32	34
32	32	31	32
33	32	32	33
32	31	34	32
31	33	33	31

Complete the column with the raw score:

X	Tally	f	Relative frequency	cf	PR
29	I	1	$\frac{1}{20} = 0.05$	20	100
31	IIII	4	$\frac{4}{20} = 0.2$	19	95
32	IIII III	8	$\frac{8}{20} = 0.4$	15	75
33	IIII	5	$\frac{5}{20} = 0.25$	7	35
34	II	2	$\frac{2}{20} = 0.10$	2	10

$$N = 20$$

## Groups frequency distribution:

This combines scores into groups referred as class intervals thus condensing the data and making overall trends more apparent.

class interval	f
81 - 83	1
78 - 80	3
75 - 77	6

\* Interval size symbolised by  $I$  refers the no. of scores in the class interval.

For eg:  $I = 3$  for the class interval 21 to 23

\* The Range symbolised by  $R$  refers to the amount of spread in the distribution of the scores.

It is determined by subtracting the lower limit from the upper limit

\* Bowling scores of high school gym class

88	128	110	140	127	Upper limit highest no + 0.5
109	92	119	142	111	Lower limit lowest no - 0.5
126	85	124	138	92	
83	114	146	112	86	
98	132	128	95	120	
122	115	92	116	119	
79	118	81	112	115	

Step: 1 Range:

$$Up. L = 146, \quad L. L = 79$$

$$R = 146.5 - 78.5$$

$$R = 68$$

Step 2:

$$68 \div 2 = 34$$

$$68 \div 3 = 22$$

$$68 \div 5 = 13$$

Where  $i = 2, 3, 5$

Here  $i = 5$

class Interval:

Here  $i = 5$ , so we can take multiple of 5

class Interval

Tally

145 - 149

100 - 104

95 - 99

90 - 94

85 - 89

80 - 84

75 - 79

111

11

111

1

3

2

3

1

1  
1  
4  
3  
5  
6  
1  
2  
3  
2  
3  
1  
35



class intervals	Tally	Frequency $f$
145 - 149	I	1
140 - 144	II	2
135 - 139	I	1
130 - 134	I	1
125 - 129	IIII	4
120 - 124	II I	3
115 - 119	IIII	5
110 - 114	IIII I	6
105 - 109	I	1
100 - 104	-	0
95 - 99	II	2
90 - 94	II I	3
85 - 89	II	2
80 - 84	II I	3
75 - 79	I	1

$N = 35$

### Measures Of Central tendency:

Central tendency:

Central tendency is a single value that describes the most typical or representative score in an entire distribution.

eg: Mode, Median, Mean

The Mode:

The Mode (Mo) specifies the value with the highest frequency in a set of scores.

To determine the mode simply arrange scores in descending order (or create frequency distribution

If there are numerous scores

Once they are arranged so, it is easy to see at a glance which score occurred with the greatest frequency.

In the set of scores below, 73, 73, 72, 70, 68, 68, 68, 59, 59, 59, 55

$$Mo = 68$$

### Grouped Frequency Distribution:

The Mode would be the mid-point of the class interval with greatest frequency.

Class Interval	Mid point	f	Class Interval	Mid point	f
36 - 38	37	8	21 - 23	22	20
33 - 35	34	11	18 - 20	19	16
30 - 32	31	18	15 - 17	16	12
27 - 29	28	26	12 - 14	13	7
24 - 26	25	32	9 - 11	10	3

The Median:

$$Mode = 32$$

Median (Mdn): Which is the midpoint in a distribution.

To find the median for an odd number of scores  
⇒ Arrange the scores in descending order from highest to lowest

⇒ The location of the median will be the score that has an equal number of scores above and below as determined by:  $N + 1/2$

26, 25, 24, 20, 18, 17, 17, 15, 12

$$(9+1)/2 = 5$$

We are looking for 5th score. Thus the median is 18.

$$\text{Mdn} = 18$$

5 is not median but rather the location.

To find the Median for an even number of scores:

⇒ Arrange the scores in order from highest to lowest.

⇒ Divide the distribution in half and draw a line between the two scores that separate the distribution in to halves.

⇒ Add the two middle scores that surround the halfway point and divide by 2.

The resulting value will be the Median.

92, 91, 90, 90, 87, 82, 77, 75, 75, 70, 68, 60

Middle scores

$$\text{Mdn} = \frac{82 + 77}{2}$$

$$\text{Mdn} = 79.5$$

The formula method for determining the median used when working from grouped frequency distribution and cumulative frequency of scores.

$$\text{Mdn} = LL + \left( \frac{50\% \text{ of } N - C_f \text{ below}}{f_{wi}} \right) i$$

Where LL = Lower limit of class interval that contains Median.

$N$  = number of scores

$C_f$  below = cumulative frequency below the class interval that contains Median

$f_{wi}$  = frequency of scores in the interval that contains Median

$i$  = size of class interval

For example:

Class Interval	$f$	$C_f$
42 - 44	4	124 $\rightarrow N$
39 - 41	8	120
36 - 38	10	112
33 - 35	11	102
30 - 32	8	91
27 - 29	18	83
24 - 26	17	65
21 - 23	16	48
18 - 20	10	32
15 - 17	11	22
12 - 14	5	11
9 - 11	6	6

$$Mdn = 23.5 + \left[ \frac{62 - 48}{17} \right] 3$$

$$= 23.5 + 2.47$$

$$= 25.97$$

$$Mdn = 25.97$$



class intervals	f	cf
60 - 64	2	60
55 - 59	1	58
50 - 54	0	57
45 - 49	5	57
40 - 44	0	52
35 - 39	7	52
30 - 34	13	45
25 - 29	12	32
20 - 24	8	20
15 - 19	5	12
10 - 14	0	7
5 - 9	5	7
0 - 4	2	2

$$\text{Mdn} = LL + \left( \frac{50\% \cdot N - cf}{f_{wi}} \right) i$$

$$= 24.5 + \left( \frac{30 - 20}{12} \right) 5$$

$$= 24.5 + 4.16$$

$$\text{Mdn} = 28.67$$

The Mean:

The Mean is the sum total of all of the scores in a distribution divided by the total number of scores.

$$\text{For a population, } \mu = \frac{\sum X}{N}$$

$$\text{For a sample, } M = \frac{\sum x}{n}$$

Let us calculate the mean for the following set of scores from a population

78, 63, 42, 98, 87, 52, 72, 64, 75, 89

$$\mu = \frac{\sum x}{N} = \frac{720}{10} = 72$$

The Mean for the following set of scores from a sample involves the same set of calculations.

3, 8, 6, 9, 10, 17, 5, 8, 1

$$M = \frac{\sum x}{n} = \frac{67}{9} = 7.44$$

Mean for a simple Frequency Distribution:

To calculate the mean for scores that have been arranged into a simple frequency distribution the formula is modified as follows:

For a population,

$$\mu = \frac{\sum fx}{N}$$

For a sample

$$M = \frac{\sum fx}{n}$$

Where  $fx$  = frequency of score multiplied by the score itself.

Let us calculate the mean for the following scores from a sample arranged into simple frequency distribution table.

X	f	fX	X	f	fX
48	1	48	41	4	164
47	4	188	40	6	240
46	2	92	39	3	117
45	4	180	38	0	0
44	9	396	37	1	37
43	8	344	36	2	72
42	5	210	35	1	35
				n = 50	$\sum fX = 2123$

$$\therefore M = \frac{\sum fX}{n} = \frac{2123}{50}$$

$$M = 42.46$$

- 1) The annual incomes for the employees of a new hat company are listed below (in thousands)

65	33	52	47
72	73	68	65
58	49	18	59
31	62	61	74
42	27	54	52
80	78	66	22

Construct a grouped frequency distribution table Using the table, determine the median, mean and mode adding additional columns as needed.

Sol:

$$\text{Range} = U.L - L.L$$

$$= 80.5 - 17.5$$

$$R = 63$$

class interval:

$$i = 2, 3, 5$$

$$\frac{63}{2} = 31.5$$

$$\frac{63}{3} = 21$$

$$\frac{63}{5} = 12.6$$

$$i = 5$$

class interval	frequency		Midpoint	$\sum fx$
	$f$	$cf$	$x$	
80 - 84	1	24	82	82
75 - 79	1	23	77	77
70 - 74	3	22	72	216
65 - 69	4	19	67	268
60 - 64	2	15	62	124
55 - 59	2	13	57	114
50 - 54	3	11	52	156
45 - 49	2	8	47	94
40 - 44	1	6	42	42
35 - 39	0	5	37	0
30 - 34	2	5	32	64
25 - 29	1	3	27	27
20 - 24	1	2	22	22
15 - 19	1	1	17	17
	24			$\sum fx = 1303$

(i) Mode = 67

(ii) Median =  $L + \left[ \frac{50\% \cdot N - c.f \text{ below}}{f_w} \right] \times i$

$$= 54.5 + \left[ \frac{12 - 11}{2} \right] \times 5 = 54.5 + (0.5) \times 5$$

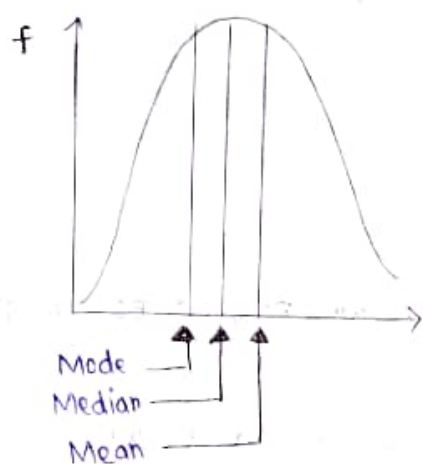
Median = 57

(iii) Mean =  $1303 \div 24 = 54.29$

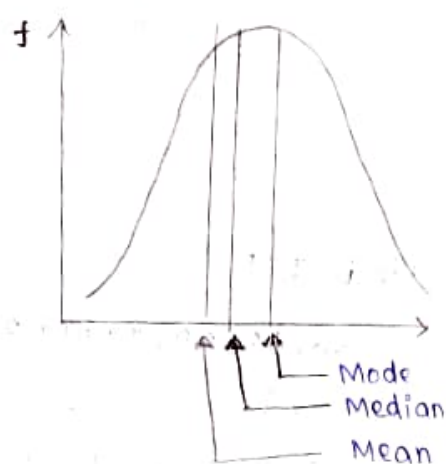


In skewed Distribution, the mode would again be at the peak; the mean would be located towards the tails in the direction of the skew (having been affected by either high or low extreme scores); and the median would be between the mode and the mean (so that half the scores lie above it and half below it).

### Positive skew

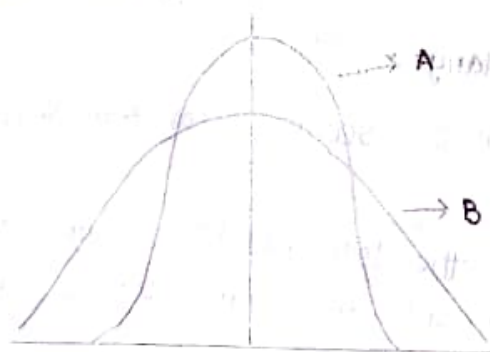


### Negative Skew



### Measures of Variability:

Measures of Variability provides information about how similar or different the scores are in relationship to other scores in the distribution.



Both of the distributions below have the same mean, however they are quite different in the variability of their scores.

Three common measures of variability are the Range, the Interquartile range, and the standard deviation.

## THE RANGE

Subtracting the lower limit of the lowest score from the upper limit of the highest score.

$$R = X_{UL - High} - X_{LL - Low}$$

For example:

For the following set of scores

3, 3, 5, 7, 7, 7, 8, 9

$$R = 9.5 - 2.5$$

$$R = 7$$

Find IQR

36, 42, 30, 7, 51, 29, 45, 35, 44, 53, 32, 50, 28, 43

33, 29.

7 28 29 29 | 30 32 33 35 | 36 42 43 44 | 45 50 51 53

$$Q_1 = \frac{29 + 30}{2}$$

$$Q_3 = \frac{44 + 45}{2}$$

$$Q_1 = 29.5$$

$$Q_3 = 44.5$$

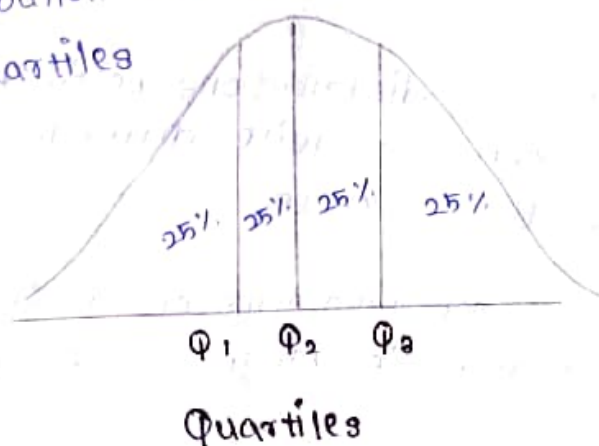
$$IQR = Q_3 - Q_1 = 44.5 - 29.5$$

$$IQR = 15$$

Interquartile Range:

The range of scores from the middle 50% of a distribution.

To determine the Interquartile range, We first divide the distribution into four equal parts, which produces three quartiles



$Q_1$  = the point at or below which 25 % of the scores lie.

$Q_2$  = the point at or below which 50 % of scores lie

$Q_3$  = the point at or below which 75 % of scores lie

## Standard Deviation:

Standard Deviation is the measure of the amount of Variation or dispersion of a set of values. A low standard deviation indicates that the values tend to be close to the mean of set, while the high standard deviation indicates that the values are spread out over wider range.

(i) Definitional formula - they are written the way that statistics are defined, these formulas usually involves more computations.

(ii) Computational formula - these are easier to use with a hand held calculator

### Definitional formula:

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$$

$\sigma$  = population standard deviation

$x$  = each value

$\mu$  = population mean

$N$  = number of values in population

1) 17, 24, 22, 26, 18, find S.D

$x$	$\mu$	$x - \mu$	$(x - \mu)^2$
17	21.4	-4.4	19.36
24	21.4	2.6	6.76
22	21.4	0.6	0.36
26	21.4	4.6	21.16
18	21.4	-3.4	11.56

$$\sum (x - \mu)^2 = 59.26$$

mean ( $\mu$ )

$$\frac{17 + 24 + 22 + 26 + 18}{5} = 21.4$$

$$\sigma = \sqrt{\frac{59.26}{5}}$$

$$\sigma = 3.44$$

2) 3, 3, 0, 1, 1, 2, 2, 2, 6, 0

X	$\mu$	$(x - \mu)$	$(x - \mu)^2$
3	2	1	1
3	2	1	1
0	2	-2	4
1	2	-1	1
1	2	-1	1
2	2	0	0
2	2	0	0
2	2	0	0
6	2	4	16
0	2	-2	4

$$\sum (x - \mu)^2 = 28$$

$$\sigma = \sqrt{\frac{28}{10}}$$

$$\sigma = 1.67$$

Sum of squares =  $\sum (x - \mu)^2$

$$\sigma^2 = \frac{99}{N} \text{ (Variance)}$$

iii) Computational formula:

$$\sigma = \sqrt{\frac{\sum x^2 - \frac{(\sum x)^2}{N}}{N}}$$

$$\frac{\sum (x - \mu)^2}{N}$$

$$\frac{\sum x^2 - \frac{(\sum x)^2}{N}}{N}$$



① 17, 24, 22, 26, 18

X	X <sup>2</sup>
17	289
24	576
22	484
26	676
18	324
$\sum X = 107$	$\sum X^2 = 2349$

$$\frac{(\sum X)^2}{N} = \frac{(107)^2}{5} = \frac{11449}{5} = 2,289.8$$

$$\sigma = \sqrt{\frac{\sum X^2 - \frac{(\sum X)^2}{N}}{N}}$$

$$= \sqrt{\frac{2349 - 2289.8}{5}}$$

$$= \sqrt{11.84}$$

$$\sigma = 3.44$$

## Unit - II

### Inferential Statistics - I.

#### Sampling Distribution:

A distribution of statistics obtained by selecting all the possible samples of specific size from a population.

#### Distribution of Sample Means:

The collection of sample means for all the possible random samples of a particular size ( $n$ ) that can be obtained from a population.

#### The central limit theorem:

The central limit theorem states that

(i) as the size of the sample ( $n$ ) increases, the shape of the sampling distribution of means approximates the shape of the normal distribution;

(ii) the mean of the sampling distribution of means will equal the population mean ( $\mu$ ); and

(iii) the standard deviation will equal  $\sigma/\sqrt{n}$ .

Shape: Even if the shape of the population distribution from which a sample is drawn is not normal, if the sample size is large (i.e.  $n = 60$  or more), the sampling distribution itself will be normal in shape.

Mean: The mean of the distribution of sample means is the mean of the population. Sample size does not affect the center of the distribution.

Standard Deviation: The standard deviation of the sampling distribution of means is called the standard error of the mean ( $\sigma_M$ ) or simply standard error. The standard error ( $\sigma_M$ )

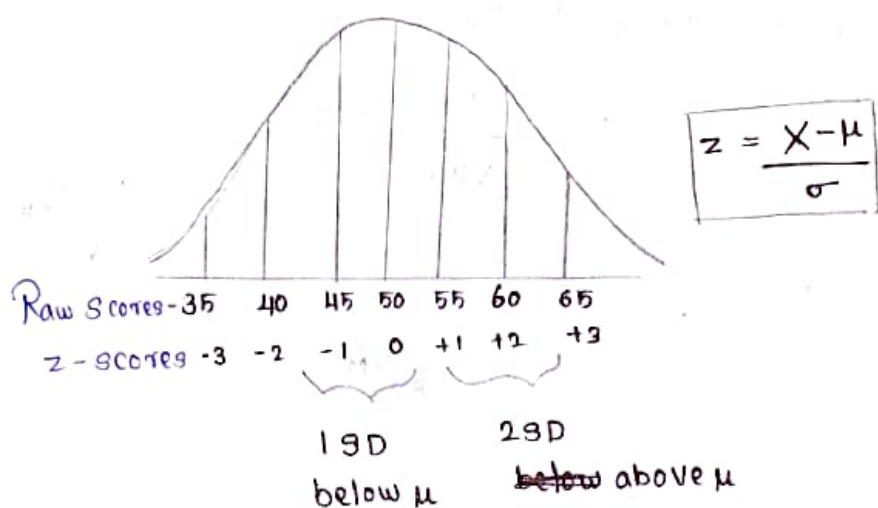
represents the amount that a sample mean ( $\bar{M}$ ) is expected to vary from the population mean ( $\mu$ ).

$$\sigma_M = \frac{\sigma}{\sqrt{n}}$$

## Probabilities, Proportions And Percentages of Sample Means:

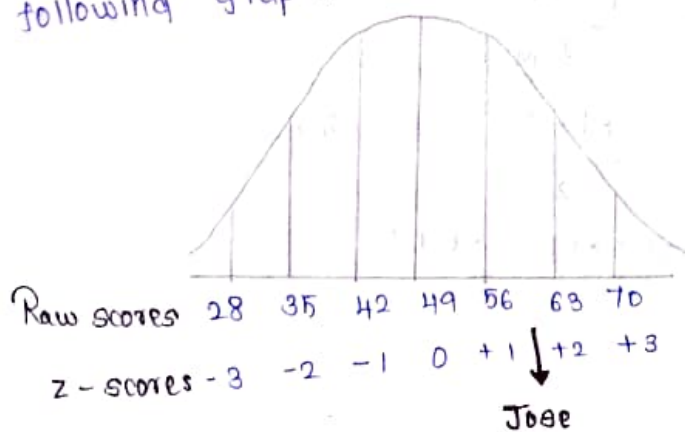
Z-score describes how far a particular raw score deviates from the mean in standard units.

The relationship between raw scores and z-scores is illustrated below:



Jose scored 60 on his history test. The class mean was 49, and the standard deviation was 7. What is Jose's equivalent z-score?  $z = \frac{60 - 49}{7} = +1.57$

Thus a raw score of 60 converts to a z-score of +1.57, meaning that Jose scored 1.57 standard deviations above the mean. This situation is illustrated in following graph:



## Probabilities, Proportions And Percentages of Sample Means:

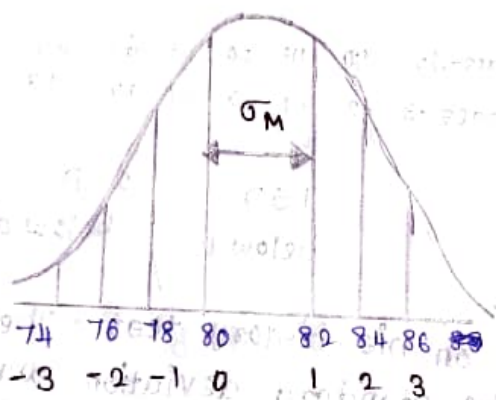
The Sampling distribution of means can be used to determine the probabilities, proportions and percentages associated with particular sample means.

$$z = \frac{M - \mu}{\sigma_M} \quad \text{and} \quad M = \mu + (z)(\sigma_M)$$

$\mu = 80$ ,  $\sigma = 14$  and  $n = 49$ . Using the appropriate formula we find the standard error ( $\sigma_M$ ) to be 2. Thus

$$\sigma_M = \frac{14}{\sqrt{49}} = 2$$

$$\sigma_M = \frac{\sigma}{\sqrt{n}}$$



- 1) Given normally shaped distribution with  $\mu = 80$  and  $\sigma_M = 2$ , A. What is the probability that an obtained sample mean will be below 81?

$$\begin{aligned} z &= \frac{M - \mu}{\sigma_M} \\ &= \frac{81 - 80}{2} = +.50 \end{aligned}$$

$$P(M < 81) = .6915$$

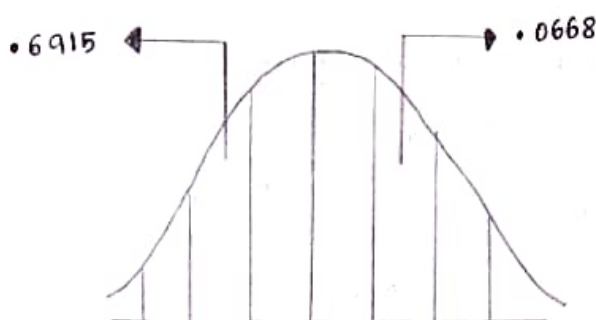


2) Greater than 83 :

$$Z = \frac{M - \mu}{\sigma_M}$$

$$= \frac{83 - 80}{2} = +1.50$$

$$P(M < 83) = \cancel{0.0668} \cdot 0.0668$$



3) Given a normally shaped population distribution with  $\mu = 95$  and  $\sigma = 5$ , a sample size of  $n = 25$  is drawn at random. The probability is  $0.05$  that mean of sample will above what value?

To find standard error

$$\sigma_M = \frac{\sigma}{\sqrt{n}} = 1$$

$$\sigma_M = \frac{\sigma}{\sqrt{n}}$$

To find  $z$ -score associated with closest proportion above  $0.0500$ , which is  $1.65$

Finally convert the  $z$ -score to a sample mean

$$M = \mu + (Z)(\sigma_M)$$

$$Z = \frac{M - \mu}{\sigma_M}$$

$$= 95 + (1.65)(1)$$

$$(Z)(\sigma_M) = M - \mu$$

$$M = 96.65$$

$$M = \mu + Z(\sigma_M)$$

To find  $z$ -score:

above (+)  $\rightarrow C$

above (-)  $\rightarrow B$

below (-)  $\rightarrow C$

below (+)  $\rightarrow B$

What range of sample mean to fall in 70% of the middle of the distribution?

$$\mu = 95, \sigma = 5, \sigma_M = 1$$

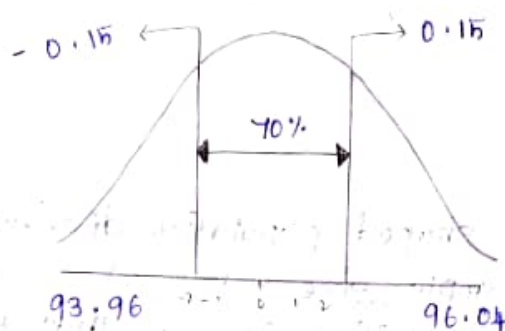
$$M = \mu + (z)(\sigma_M)$$

$$\Rightarrow M = 95 + (-1.04)(1)$$

$$M = 93.96$$

$$\Rightarrow M = 95 + (1.04)(1)$$

$$M = 96.04$$



$$100\% - 70\% = 30\%$$

$$= \frac{30}{100} = 0.3$$

middle  
tail

$$d = \frac{0.3}{2}$$

$$d = 0.15$$

- 1) Finding Probability or proportion of Given sample Means. Given a normal distribution with  $\mu = 100$  and  $\sigma = 12$ , a sample of  $n = 36$  is drawn at random.

A) What is the probability that sample mean will fall above 99?

B) What proportion of sample mean will have value less than 95?

$$z = \frac{M - \mu}{\sigma_M}$$

$$\sigma_M = \frac{\sigma}{\sqrt{n}} = \frac{12}{\sqrt{36}} = 2$$

$$M = z(\sigma_M) + \mu$$

$$99 = z(\sigma_M) + 100$$

$$-1 = 2z$$

$$z = -\frac{1}{2} = -0.5$$

## Hypothesis Testing:

=> Hypothesis testing is a procedure used in inferential statistics to estimate population parameters based on sample data.

=> Sample research question <sup>given</sup> => test of reading comprehension for 4th graders is normally distributed with mean  $\mu = 70$  and  $s.d = 10$

=> Random sample of  $N = 25$ , 4th graders are taught and then tested for reading comprehension

=> Sample  $M = 75$  is obtained  
mean

=> Thus the sample mean  $M$  differ enough from population mean ( $\mu$ ) to conclude that reading technique made a difference in level of comprehension.

### Testing steps:

=> Formulate the hypothesis

=> Indicate the  $\alpha$ -level and determine the critical values

=> calculate relevant statistics

=> Make decisions and report the result

### Formulating the hypothesis:

(i) Formulate the hypothesis, there will be to mutually exclusive hypothesis

(ii) Either the new reading technique does not have an effect on comprehension or it does.

(iii) The null hypothesis <sup>(H<sub>0</sub>)</sup> attributes any difference between obtained sample mean and population mean to chance.

\* chance: (due to chance or random sampling error)

\* equality:

## \* Ineffective treatment

(iv) The alternative hypothesis  $H_1$  called the research hypothesis is the opposite of the null.  $H_1$  describes true differences

(v) The effectiveness of treatment.

Step 1:

$$\Rightarrow H_0 = \mu = 70$$

$\Rightarrow H_0$  states that the new reading technique will not change the mean level of reading comprehension.

$\Rightarrow$  The population mean will be 70.

Step 2:

$$\Rightarrow H_1 = \mu \neq 70$$

$\Rightarrow H_1$  states that the new reading techniques does have an effect on comprehension

Step 2:

$\Rightarrow$  Indicate alpha level and determine the critical values

$\Rightarrow$  Most of an alpha level is set at 0.05 (or) 0.01 and occasionally at 0.001

Sample means that fall in this area have a 5% or less likelihood of occurring by chance alone if  $H_0$  is true.  $Z_{\text{critical values}} = \pm 1.96$

Step 3:

calculate the relevant statistics

$$Z_{\text{obt}} = \frac{M - \mu}{\sigma_M}$$

Step 4:

Make a decision and report the results  
Finally we must either we must reject the null hypothesis or accept the null hypothesis.



SRQ:

Given  $\mu = 70$ ,  $\sigma = 10$ ,  $n = 25$ , did the new reading technique have effect on Comprehension?

Step 1: Formulate hypothesis

$$H_0 : \mu = 70$$

$$H_1 : \mu \neq 70$$

Step 2 :

Indicate the alpha level and determine critical values

$$\alpha = 0.05$$

$$Z_{crit} = \pm 1.96$$

Step 3 :

Calculate Relevant Statistics

Given:  $n = 25$ ,  $\mu = 70$ ,  $\sigma = 10$ ,  $M = 75$

$$\sigma_M = \frac{\sigma}{\sqrt{n}} = \frac{10}{\sqrt{25}} = 2$$

$$Z_{obt} = \frac{M - \mu}{\sigma_M}$$

$$= \frac{75 - 70}{2}$$

$$= +2.50$$

Step 4: Make a decision and report the results.

The new reading program had a significant effect on reading comprehension.

Reject  $H_0$ ,  $Z = +2.50$  then  $p < 0.05$ .

Hypothesis Testing With z-Test Research Question:

A standardized productivity scale has a  $\mu = 25$  and a  $\sigma = 5$ . The CEO of company A wants to know if employee participation in company decisions has an effect on productivity. A sample of  $n = 75$  employees



participated in the decision making process was administrated the productivity scale with a result of  $M = 27$ . Does participation in company decisions have an effect on Productivity? Use  $\alpha = 0.01$  and a two tailed test.

Step 1:

$$H_0: \mu = 25$$

$H_0$  states that the employee participation in company decisions will not change the Productivity.

$$H_1: \mu \neq 25$$

It states that employee participation in company decision will have an effect on productivity

Step 2:

Indicate alpha level and determine the critical values.

$$\alpha = 0.01$$

$$Z_{\text{crit}} = \pm 2.58$$

Step 3:

$$\text{Given } n = 75, \mu = 25, \sigma = 5, M = 27$$

Calculate the relevant statistics

$$\begin{aligned} Z_{\text{obt}} &= \frac{M - \mu}{\sigma_M} & \sigma_M &= \frac{\sigma}{\sqrt{n}} \\ & & &= \frac{5}{\sqrt{75}} \\ & & &= 0.58 \\ &= \frac{27 - 25}{0.58} \\ &= \frac{2}{0.58} \end{aligned}$$

$$Z = 3.45$$

Step 4:

Make a decision and report the results. The ~~mean~~ mean steading technique. The employee participation has a significant effect on productivity.

Reject  $H_0$ .  $Z = +3.45$ ,  $P < 0.01$

Effect Size for a z-test:

Effect size indicates magnitude of a treatment effect:

General formula,

$$d = \frac{|M - \mu|}{\sigma}$$

For our research Problem,

$$d = \frac{|75 - 70|}{10}$$

$$d = 0.5$$

$\Rightarrow d = 0.20$  to  $0.49$  - Small effect

$\Rightarrow d = 0.50$  to  $0.79$  - Moderate effect

$\Rightarrow d = 0.80$  and above - Large effect

One Sample t-Test:

The One-Sample t test is a test of hypothesis about a population mean ( $\mu$ ) when the population standard deviation ( $\sigma$ ) is not known.

This test is used when researchers want to know (1) if a sample is representative of a population and/or (2) if a particular treatment or condition has a significant effect.  $df = n - 1$

- 1) The population mean of a standardized test of critical thinking is  $\mu = 53$ . A group of faculty members at a small community college underwent a - 10 - week training program to learn techniques designed to help students develop their critical thinking skills. After the training, the new techniques were implemented in the classrooms. The mean critical thinking score for a sample of  $n = 87$  students exposed to the new techniques was  $M = 55$  with  $SS = 6013$ . Do the results suggest that the training program had a significant effect? Use a two-tailed test and  $\alpha = 0.5$ .

Step 1: Formulate hypothesis

$$H_0: \mu = 53$$

$$H_1: \mu \neq 53$$

The null hypothesis

Step 2: Indicate the alpha level and determine critical values

$$\alpha = 0.05$$

$$df = 86$$

$$t_{crit} = \pm 2.000$$

Step 3:

calculate relevant statistics

$$S = \sqrt{\frac{SS}{n-1}}$$

$$= \sqrt{\frac{6013}{87-1}}$$

$$S = 8.36$$

Before we can determine our obtained  $t$ -value, we first need to calculate standard error using estimated standard deviation above:

$$S_M = \frac{S}{\sqrt{n}} = \frac{8.36}{\sqrt{87}} = 0.90$$

Finally, we can calculate the  $t$ -statistic:

$$t_{obt} = \frac{M - \mu}{S_M}$$

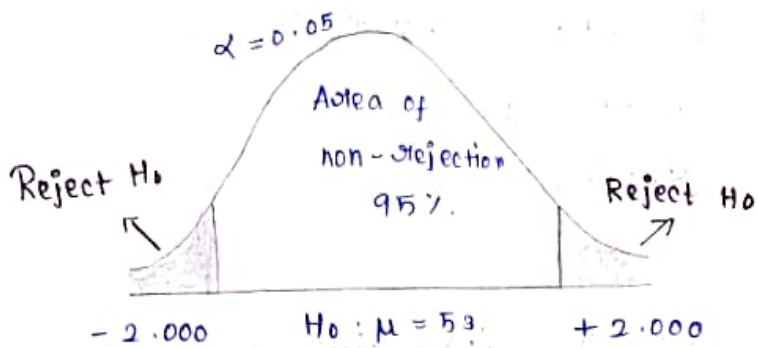
$$= \frac{55 - 53}{0.90}$$

$$= +2.22$$

Step 4:

Make a decision and report the results

The  $t$ -distribution for this example is below:



The  $t_{crit}$  values are  $\pm 2.000$ . To reject  $H_0$ , we need a sample mean with an obtained  $t$ -value beyond 2.000 in either direction. Our obtained  $t$ -value was +2.22 which falls in the direction of the rejection. Thus we reject the null-hypothesis.

Test	Hypothesis	$t_{crit}$
Two-tailed	$H_0: \mu = \text{value specified in pblm}$ $H_1: \mu \neq \text{value specified in pblm}$	$t_{crit} = \pm \text{value in chart}$
one-tailed, left	$H_0: \mu \geq \text{value spec in pblm}$ $H_1: \mu < \text{value spec in pblm}$	$t_{crit} = - \text{value in chart}$
one-tailed, right	$H_0: \mu \leq \text{value spec}$ $H_1: \mu > \text{value spec}$	$t_{crit} = + \text{value in chart}$



Formulate the hypothesis

$$H_0: \mu \leq 9$$

$$H_1: \mu > 9$$

Step 2: Indicating alpha level and determine the critical values

$$\alpha = 0.05$$

$$df = n - 1 = 25 - 1$$

$$df = 24$$

$$t_{crit} = +1.711$$

Step 3:

$$S_M = \frac{s}{\sqrt{n}} = \frac{6}{\sqrt{25}}$$

$$= \frac{6}{5} = 1.2$$

$$t_{obt} = \frac{11 - 9}{1.2}$$

$$= 1.66666$$

$$t_{obt} = +1.67$$

Step 4:

The study did not show that extroverts have significantly more friends than national average.

Fail to reject  $H_0$ ,  $t(24) = +1.67$ ,  $P > 0.05$

Research Question:

The mean scores



Step 1:

Formulate hypothesis

$$H_0: \mu \geq 83$$

$$H_1: \mu < 83$$

Lower scores reflect  
better performance  
in golf

Step 2: Indicating alpha level and determine the critical values.

$$\alpha = 0.05$$

$$df = n - 1$$

$$df = 24$$

$$t_{crit} = -1.711$$

Step 3:

$$s = \sqrt{\frac{99}{n-1}} = \sqrt{\frac{788}{24}}$$

$$s = 5.73$$

$$s_M = \frac{s}{\sqrt{n}} = 1.15$$

$$t_{obt} = \frac{81 - 83}{1.15}$$

$$t_{obt} = -1.74$$



Step 4:

The flag has significantly improved golf scores.

Reject  $H_0$ ,  $t_{crit} > t_{obt}$ ,  $P > 0.05$

$$t(24) = -1.74$$

Two - Sample t - test:

Standard Error of Difference:

$$S_{M_1 - M_2} = \sqrt{\left( \frac{SS_1 + SS_2}{n_1 + n_2 - 2} \right) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

The standard deviation of sampling distribution of difference between means is called the standard error of difference between means.

Formula for T-test, Independent sample Designs:

$$t_{\text{obt}} = \frac{M_1 - M_2}{S_{M_1 - M_2}}$$

Degrees of freedom:

$$df = n_1 + n_2 - 2$$

A political

Step 1: Formulate hypothesis

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

Step 2: Indicate the alpha level and determine critical values

$$\alpha = 0.05$$

$$df = n_1 + n_2 - 2$$

$$= 11 + 10 - 2$$

$$= 19$$

$$t_{crit} = \pm 2.093$$

Step 3: Calculate relevant statistics

$X_1$	$X_1^2$	$X_2$	$X_2^2$
59	3481	36	1296
48	2308	42	1764
45	2025	50	2500
39	1521	37	1369
52	2704	51	2601
56	3136	32	1024
50	2500	47	2209
41	1681	38	1444
46	2116	40	1600
45	2025	31	961
48	2304		
$\Sigma X_1 = 529$	$\Sigma X_1^2 = 25,797$	$\Sigma X_2 = 404$	$\Sigma X_2^2 = 16,768$

$$n_1 = 11$$

$$n_2 = 10$$

$$M_1 = 48.09$$

$$M_2 = 40.40$$

Calculate the sum of squares for each sample

$$SS = \Sigma X^2 - \frac{(\Sigma X)^2}{n}$$

$$SS_1 = 25,797 - \frac{(529)^2}{11}$$

$$= 356.91$$

$$SS_2 = 16,768 - \frac{(404)^2}{10}$$

$$= 446.40$$

calculate the standard error of difference.  
calculate the  $t$ -statistic.

$$S_{M_1 - M_2} = \sqrt{\left( \frac{SS_1 + SS_2}{n_1 + n_2 - 2} \right) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

$$S_{M_1 - M_2} = \sqrt{\left( \frac{356.91 + 446.4}{11 + 10 - 2} \right) \left( \frac{1}{11} + \frac{1}{10} \right)}$$

$$= \sqrt{(42.28)(.19)}$$

$$= +2.83$$

$$t_{obt} = \frac{(M_1 - M_2)}{S_{M_1 - M_2}} = \frac{(48.09 - 40.40)}{2.83}$$

$$t_{obt} = +2.72$$

Step 4:

Make decision and report the results

Students who read the newspaper scored significantly higher on knowledge of current events than students who watched the evening news. Reject  $H_0$ .

$$t(19) = +2.72, P < 0.05$$

Research Ques:

An Industrial.

Step 1: Formulate hypothesis

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

Step 2: Indicate alpha level and determine critical values

$$\alpha = 0.05$$

$$df = n_1 + n_2 - 2$$

$$= 7 + 6 - 2$$

$$df = 11$$

$$\pm t_{crit} = 2.201$$

Step 3:

$X_1$	$X_1^2$	$X_2$	$X_2^2$
8	64	7	49
10	100	6	36
9	81	8	64
7	49	10	100
8	64	7	49
11	121	9	81
12	144		
$\Sigma X_1 = 65$	$\Sigma X_1^2 = 623$	$\Sigma X_2 = 47$	$\Sigma X_2^2 = 379$

$$n_1 = 7$$

$$n_2 = 6$$

$$M_1 = \frac{65}{7} = 9.285$$

$$M_2 = 7.8333$$

calculate sum of squares

$$SS_1 = 623 - \frac{(65)^2}{7}$$

$$SS_1 = 19.428$$

$$SS_2 = 379 - \frac{(47)^2}{6}$$

$$SS_2 = 10.8333$$



$$S_{M_1 - M_2} = \sqrt{\left( \frac{SS_1 + SS_2}{n_1 + n_2 - 2} \right) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

$$S_{M_1 - M_2} = \sqrt{\left( \frac{19.428 + 10.833}{7 + 6 - 2} \right) \left( \frac{1}{7} + \frac{1}{6} \right)}$$

$$= \sqrt{2.751 \times 0.309}$$

$$= 0.9219$$

$$t_{obt} = \frac{(M_1 - M_2)}{S_{M_1 - M_2}} = \frac{9.285 - 7.8333}{0.9219}$$

$$t_{obt} = 1.5746$$

step 4:

Make decision and report the results

Listening to smooth jazz music did not significantly affect widget productivity. Failed to reject  $H_0$ ,  $t(11) = +1.59$ ,  $p < 0.05$

## Two Sample t-Test Related Sample Designs:

Standard Deviation of the difference:

$$S_D = \sqrt{\frac{\sum D^2 - \frac{(\sum D)^2}{n}}{n-1}}$$

$$M_D = \frac{\sum D}{n}$$

$M_D \rightarrow$  Mean difference

Standard error:

$$S_{M_D} = \frac{S_D}{\sqrt{n}}$$

$$t_{obt} = \frac{M_D}{S_{M_D}}$$

Sample Research Question:

A dietician

step 1:

Formulate the hypothesis

$$H_0 : M_D = 0$$

$$H_1 : M_D \neq 0$$

Step 2:

Indicate the alpha level and determine the critical values.

$$\alpha = 0.05$$

$$df = n - 1 = 7$$

$$t_{crit} = \pm 2.365$$

Step 3:

calculate relevant statistics

subject	Weight before	Weight after	D'	D <sup>2</sup>
1	156	142	-14	196
2	192	173	-19	361
3	138	140	2	4
4	167	151	-16	256
5	110	109	-1	1
6	159	151	-8	64
7	171	154	-17	289
8	129	133	4	16
			$\Sigma D = -69$	$\Sigma D^2 = 1187$

$$MD = \frac{\Sigma D}{n} = \frac{-69}{8} = -8.63$$

$$S_{MD} = \frac{SD}{\sqrt{n}} = \frac{9.20}{\sqrt{8}} = 3.25$$

$$s_p = \sqrt{\frac{\Sigma D^2 - \frac{(\Sigma D)^2}{n}}{n-1}} = \sqrt{\frac{1187 - \frac{(-69)^2}{8}}{8-1}} = 9$$

$$t_{obt} = \frac{MD}{S_{MD}} = \frac{-8.63}{3.25} = -2.66$$

Step 4:

Make a decision and report the results  
subjects on the "no nighttime snack" diet showed  
a significant weight loss. Reject  $H_0$ ,  
 $t(7) = -2.66$ ,  $P < 0.05$ .

Research Ques:

A psychology Professor

Step 1:

Formulate hypothesis

$$H_0: \mu_D \geq 0$$

$$H_1: \mu_D \neq 0$$

Step 2:

Indicate alpha level and determine the critical values

$$\alpha = 0.05$$

$$df = n - 1 = 18 - 1 = 17$$

$$t_{crit} = -1.740$$

Step 3:

$$SD = \sqrt{\frac{SSD}{n-1}} = \sqrt{\frac{286}{17}} = 4.1016$$

$$S_{MD} = \frac{SD}{\sqrt{n}} = \frac{4.1016}{\sqrt{18}} = 0.9667$$

$$t_{obs} = \frac{MD}{S_{MD}} = \frac{-3}{0.9667} = -3.10334$$

## One - Way Analysis of Variance

ANOVA is used to abbreviate analysis of Variance.  
 Variance in ANOVA:

⇒ Total Variance is separated into two kinds:  
 Within - treatment Variance and between - treatments Variance.

⇒ Within - treatment Variance refers to be the Variability within a particular sample

⇒ Between treatments Variance refers to the variability between the treatment groups.

⇒ Between - treatments Variance = Individual differences + Experimental error + Treatment

⇒ Within - treatment Variance = Individual differences + Experimental error.

Sum of Squares:

$$SS_{wi} = \sum \left[ \sum X_{tj}^2 - \frac{(\sum x_t)^2}{n_t} \right]$$

$$SS_{bet} = \sum \left[ \frac{(\sum x_t)^2}{n_t} \right] - \frac{(\sum X_{tot})^2}{N}$$

$$SS_{total} = \sum X_{tot}^2 - \frac{(\sum X_{tot})^2}{N}$$

Degrees of freedom:

$$df_{wi} = N - K$$

$$df_{bet} = K - 1$$

$$df_{total} = N - 1$$

$$df_{wi} + df_{bet} = df_{total}$$



Mean Square:

$$MS_{\text{bet}} = \frac{SS_{\text{bet}}}{df_{\text{bet}}}$$

$$MS_{\text{wi}} = \frac{SS_{\text{wi}}}{df_{\text{wi}}}$$

F - statistic:

$$F_{\text{obt}} = \frac{MS_{\text{bet}}}{MS_{\text{wi}}}$$

Research Question:

A counselling psychologist

Step 1: Formulate hypothesis

$$H_0: \mu_1 = \mu_2 = \mu_3$$

$H_1$ : Some  $\mu$ 's are not equal.

Step 2: Indicate alpha level and determine critical values.

$$N = 18, K = 3$$

$$df_{\text{wi}} = N - K = 18 - 3 = 15$$

$$df_{\text{bet}} = K - 1 = 3 - 1 = 2$$

$$df_{\text{total}} = N - 1 = 17$$

$$t_{\text{crit}} = 3.68$$

step 9:

$X_1$	$X_1^2$	$X_2$	$X_2^2$	$X_3$	$X_3^2$
2	4	3	9	8	64
4	16	11	16	7	49
5	25	2	4	10	100
3	9	4	16	6	36
4	16	3	9	8	64
6	36	2	4	9	81

$$\sum X_1 = 24 \quad \sum X_1^2 = 106 \quad \sum X_2 = 18 \quad \sum X_2^2 = 58 \quad \sum X_3 = 48 \quad \sum X_3^2 = 394$$

$$SS_{wi} = \sum \left[ \frac{\sum X_t^2}{n_t} - \frac{(\sum X_t)^2}{n_t} \right]$$

$$= \left[ 106 - \frac{(24)^2}{6} \right] + \left[ 58 - \frac{(18)^2}{6} \right] + \left[ 394 - \frac{(48)^2}{6} \right]$$

$$= [10 + 4 + 10]$$

$$SS_{wi} = 24$$

$$SS_{bet} = \sum \left[ \frac{(\sum X_t)^2}{n_t} \right] - \frac{(\sum X_{tot})^2}{N} \Rightarrow \frac{(90)^2}{18}$$

$$= \left[ \frac{(24)^2}{6} + \frac{(18)^2}{6} + \frac{(48)^2}{6} \right] - 450$$

$$= [96 + 54 + 384] - 450$$

$$SS_{bet} = 84$$

$$SS_{\text{total}} = \sum X_{\text{tot}}^2 - \frac{(\sum X_{\text{tot}})^2}{N}$$

$$= 558 - 450$$

$$SS_{\text{tot}} = 108$$

$$MS_{\text{bet}} = \frac{SS_{\text{bet}}}{df_{\text{bet}}} = \frac{84}{2} = 42$$

$$MS_{\text{wi}} = \frac{SS_{\text{wi}}}{df_{\text{wi}}} = \frac{24}{15} = 1.6$$

$$F = \frac{MS_{\text{bet}}}{MS_{\text{wi}}} = 2.625$$

$$F_{\text{obt}} = 2.625$$

Summary Data:

Questions:

Three different techniques meditation

Step 1: Formulate the hypothesis

$$H_0: \mu_1 = \mu_2 = \mu_3$$

$H_1$ : Some  $\mu$ 's are not equal

Step 2: Indicate alpha level and determine critical values

$$N = 15, K = 3$$

$$df_{wi} = N - K = 15 - 3 = 12$$

$$df_{bet} = K - 1 = 3 - 1 = 2$$

$$df_{tot} = N - 1 = 15 - 1 = 14$$

$$t_{crit} = 3.88$$

Step 3:

$X_1$	$X_1^2$	$X_2$	$X_2^2$	$X_3$	$X_3^2$
10	100	6	36	5	25
12	144	8	64	9	81
9	81	3	9	12	144
15	225	0	0	8	64
13	169	2	4	4	16
$\Sigma X_1 = 59$	$\Sigma X_1^2 = 719$	$\Sigma X_2 = 19$	$\Sigma X_2^2 = 113$	$\Sigma X_3 = 38$	$\Sigma X_3^2 = 330$

$$SS_{wi} = \sum \left[ \sum x_t^2 - \frac{(\sum x_t)^2}{n_t} \right]$$

$$= \left[ 719 - \frac{(59)^2}{5} \right] + ~~22.8~~ \left[ 113 - \frac{(19)^2}{5} \right]$$

$$+ ~~22.8~~ \left[ 330 - \frac{(38)^2}{5} \right]$$

~~SS bet~~

$$= 22.8 + 40.8 + 41.2$$

$$SS_{wi} = 104.8$$

$$SS_{bet} = \sum \left[ \frac{(\sum x_t)^2}{n_t} \right] - \frac{(\sum x_{tot})^2}{N}$$

$$= \frac{(59)^2}{5} + \frac{(19)^2}{5} + \frac{(38)^2}{5} - \frac{(116)^2}{15}$$

$$= 696.2 + 72.2 + 288.8 - 897.066$$

$$SS_{bet} = 160.134$$

$$SS_{tot} = \sum x_{tot}^2 - \frac{(\sum x_{tot})^2}{N}$$

$$= 1162 - \frac{(116)^2}{15}$$

$$= 1162 - 897.066$$

$$SS_{tot} = 264.934$$



$$MS_{bet} = \frac{SS_{bet}}{df_{bet}} = \frac{160.134}{2} = 80.067$$

$$MS_{wi} = \frac{SS_{wi}}{df_{wi}} = \frac{104.8}{12} = 8.733$$

$$F = \frac{MS_{bet}}{MS_{wi}} = \frac{80.067}{8.733} = 9.1683$$

$$F_{obt} = 9.1683$$

Step 4:

Summary Data

Source	SS	df	MS	F	P
Medication	104.8	2	80.067	9.168	<0.05
Exercise	160.134	12	8.733		
Diet	264.934	14			

Step 4: Conclusion:

There is significant difference in mean reduction of blood pressure among three techniques.

Reject  $H_0$ .  $t(15) = 9.1683$

## Chi-Square test

### THE CHI-SQUARE STATISTIC:

The chi-square statistic measures the amount of discrepancy between observed frequencies and the frequencies that would be expected due to chance, or random sampling error. The formula for chi-square is

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

Where :  $f_o$  = observed frequencies

$f_e$  = expected frequencies

- ⇒ Observed frequencies are actual frequencies from our sample that fall into each category
- ⇒ Expected frequencies are the frequency values that would be expected if the null hypothesis is true.
- ⇒ If we find no differences between our observed and expected frequencies, then our obtained  $\chi^2$  value will equal 0.
- ⇒ But if our observed frequencies differ from those that would be expected, then  $\chi^2$  will be greater than 0.

### (i) GOODNESS OF FIT FOR KNOWN PROPORTIONS:

Sample Research Question:

A new professor at a midsize college wanted to see if her grade distribution after her first year of teaching was comparable to the overall grade distribution, which has the following percentages : A - 10%, B - 22%, C - 40%, D - 21% and F - 7%. The distribution of the new Professor's grade for 323 students at the end of

her first year as follows: 38 recieved A's, 78 recieved B's, 139 recieved C's, 55 recieved D's and 13 recieved F's. Does the new professor's grade distribution fit the overall college's distribution? Test Using  $\alpha = 0.05$ .

Sol:

Step 1: Formulate the hypothesis

$\left\{ \begin{array}{l} \text{chi-square is} \\ \text{Ho-null hypothesis} \\ \text{always fits} \\ \text{H}_1 - \text{not fit} \end{array} \right\}$

H<sub>0</sub>: The distribution of grades for the new professor fits overall grade distribution of the college.

H<sub>1</sub>: The distribution of grades for the new professor does not fit the overall grade distribution of the college.

Step 2: Indicate alpha level and determine the critical Values.

$$\alpha = 0.05$$

$$df = k - 1$$

$$= 5 - 1 = 4$$

$$df = 4$$

$$\chi^2_{crit} = 9.488$$

Step 3: calculate Relevant statistics

$f_e = (\text{Known proportion}) (n)$

$$A = 0.10 \times 323 = 32.30$$

$$B = 0.22 \times 323 = 71.06$$

$$C = 0.40 \times 323 = 129.20$$

$$D = 0.21 \times 323 = 67.83$$

$$F = 0.07 \times 323 = 22.61$$

$$\chi^2_{obt} = \sum \frac{(f_o - f_e)^2}{f_e}$$

$$= \left( \frac{(38 - 32.3)^2}{32.3} \right) + \frac{(78 - 71.06)^2}{71.06} + \frac{(139 - 129.2)^2}{129.2} + \frac{(55 - 67.83)^2}{67.83} + \left( \frac{(13 - 22.61)^2}{22.61} \right)$$

$$= 1.01 + 0.68 + 0.74 + 2.43 + 4.08$$

$$= 8.94$$

Step 4 :

(ii) TEST OF INDEPENDENCE :

\* The other type of problem, where chi-square has its greatest usefulness is in testing for independence of variables.

\* This test examines the frequencies for two variables at different levels in order to determine whether one variable is independent of the other or if the variables are related.

Sample Research Question:

A professor at a veterinary school of medicine is curious about whether or not a relationship exists between gender and type of pets owned in childhood. She randomly asks a sample of  $n = 260$  students (132 female and 128 male) about their pet ownership, the frequency of which is recorded in the table below. Is there a relationship between gender and type of pet ownership? Test at  $\alpha = 0.05$ .

Gender	Type of pet					Row total
	Dogs	Cats	Birds	Reptiles	Rodents	
Female	58	36	22	4	12	132
Male	62	22	14	10	20	128
Column total	120	58	36	14	32	$n = 260$



Expected frequency :-

$$f_e = \frac{(f_c)(f_r)}{n}$$

$f_c$  = column total

$f_r$  = row total

$n$  = sample size

$$f_e = \frac{(f_c)(f_r)}{n} = \frac{(120)(132)}{260} = 60.92$$

Type of Pets

Gender	Dogs	Cats	Birds	Reptiles	Rodents	Row
Female	58 (60.92)	36 (29.45)	22 (18.28)	4 (7.11)	12 (16.24)	132
Male	62 (59.08)	22 (28.55)	14 (17.72)	10 (6.89)	20 (15.75)	128
column	120	58	36	14	32	$n = 260$

Step 1: Formulate hypothesis

$H_0$  : Gender and pet ownership are independent and Unrelated.

$H_1$  : There is a relationship between gender and pet ownership.

Step 2: Indicate alpha level and determine the critical values.

$$df = (R-1)(C-1)$$

$$\alpha = 0.05$$

$R$  = no. of Rows

$C$  = no. of Columns

$$df = (2-1)(5-1) = 4$$

$$\chi^2_{crit} = 9.488$$



Step 3: Calculate Relevant Statistics

$$\chi^2_{\text{obt}} = \sum \frac{(f_o - f_e)^2}{f_e}$$

$$= \frac{(58 - 60.92)^2}{60.92} + \frac{(36 - 29.45)^2}{29.45} + \frac{(22 - 18.28)^2}{18.28} +$$

$$+ \frac{(20 - 15.75)^2}{15.75}$$

$$= 0.14 + 1.46 + 0.76 + 1.36 + 1.11 + 0.14 +$$

$$1.50 + 0.78 + 1.40 + 1.15$$

$$= 9.80$$

Step 4: Make a decision and Report the results

There is a relationship between gender and type of pet owned. Reject  $H_0$ .  $\chi^2 (4 \text{ df}, n = 260)$

$$9.80, p < 0.05$$

A Researcher is Interested

## Penmanship Quality

Handedness	Low	Medium	High	Row total
Left-handed	8 (7.31)	29 (28.38)	6 (7.31)	43
Right-handed	9 (9.69)	37 (37.62)	11 (9.69)	57
column total	17	66	17	

Formulate hypothesis

Step 1:

$H_0$ : Handedness and penmanship quality are Unrelated

$H_1$ : There is a relationship between handedness and penmanship quality.

Step 2: Indicate alpha level and determine critical values

$$df = (R-1)(C-1)$$

$$= (2-1)(3-1)$$

$$df = 2$$

$$\alpha = 0.05$$

$$\chi^2_{crit} = 5.991$$

Step 3: calculate Relevant Statistics

$$\chi^2_{obt} = \sum \frac{(f_o - f_e)^2}{f_e}$$

$$= \frac{(8 - 7.31)^2}{7.31} + \frac{(29 - 28.38)^2}{28.38} + \frac{(6 - 7.31)^2}{7.31} +$$

$$\frac{(9 - 9.69)^2}{9.69} + \frac{(37 - 37.62)^2}{37.62} + \frac{(11 - 9.69)^2}{9.69}$$

$$= 0.065 + 0.01 + 0.23 + 0.04 + 0.01 + 0.17$$

$$= 0.52$$

## Non-Parametric Test

### Mann Whitney U Test

Two Sample T-test:

#### Procedure:

- 1) Rank data (1 being the lowest)
- 2) Find the sum of the ranks for each group
- 3) Find U for both groups

$$U_1 = N_1 N_2 + \frac{N_1 (N_1 + 1)}{2} - T_1 \quad (\text{or})$$

$$U_2 = N_1 N_2 + \frac{N_2 (N_2 + 1)}{2} - T_2$$

- 4) Using a table for Mann Whitney U, find the critical value
- 5) Compare the critical value with the computed value (smaller value between  $U_1$  and  $U_2$ )

### Mann Whitney U Test

FLUSHES	3	4	2	6	2	5
KLEEROUT	9	7	5	10	6	8

FLUSHM	3	4	2	6	2	5	Total
RANK	3	4	1.5	7.5	1.5	5.5	23
KLEEROUT	9	7	5	10	6	8	
RANK	11	9	5.5	12	7.5	10	55

PROCEDURE:

Find  $U$  for both groups

$$U_1 = (6)(6) + \frac{6(7)}{2} - 23$$

$$U_2 = (6)(6) + \frac{6(7)}{2} - 55$$

Thus,

$$U_1 = 34 \quad \text{and} \quad U_2 = 2$$

Hence the computed value  $U = 2$

Find critical value

$$n_1 = 6, \quad n_2 = 6$$

The critical value is 5

\*  $\Rightarrow$  critical value: 5

computed value: 2

critical value - opposite to other test

Since the computed value is less than the critical value we reject the null hypothesis.

Therefore, there is a significant difference in the rated effectiveness of two laxatives.

### KRUSKAL - WALLIS TEST PROCEDURE:

- 1) Rank all the scores (1 being the lowest)
  - 2) Find  $T_c$  (the total of all the ranks of each group)
  - 3) solve for  $H = \frac{12}{N(N+1)} \sum \frac{T_c^2}{n_c} - 3(N+1)$
  - 4) Using the Appropriate Table, find the critical value
  - 5) Compare the critical value with the computed value
- Note: Reject null if computed value is greater than or equal to the critical value.

a	Rank	b	Rank	c	Rank	d	Rank
68		78		94		54	
63		69		82		51	
58		58		73			
51		57		67			
41		53		66			
				61			



## Unit - IV

### Simple Linear Regression

The two main objectives:

- (i) Establish if there is a relationship between two variables. More specifically establish, if there is a statistically significant relationship between the two.

Ex: Income and Spending, Wage & gender, student height and exam score.

- (ii) Forecasting new observations, we can use what we know about the relationship to forecast. Unobserved values, there are two variables

\* Dependent Variable - This is the variable whose values we want to explain our forecast, denoted by  $y$ .

\* Independent Variable - This is the variable that explains the other one, we denote it by  $x$

Linear Equation:

$$y = mx + c$$

Simple Linear Regression:

$$y = \beta_0 + \beta_1 x$$

$\beta_0 \rightarrow$  Intercept

$\beta_1 \rightarrow$  slope

$$y = \beta_0 + \beta_1 x + \epsilon \text{ (epsilon)}$$

$$\beta_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

$\bar{x} \rightarrow$  Mean of  $x$

$\bar{y} \rightarrow$  Mean of  $y$

$x$	$y$	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(x - \bar{x})(y - \bar{y})$
1	2	-2	-2	4	4
2	4	-1	0	1	0
3	5	0	1	0	0
4	4	1	0	1	0
5	5	2	1	4	2
15	20			$\sum (x - \bar{x})^2$	$\sum (x - \bar{x})(y - \bar{y})$
$\bar{x} = 3$	$\bar{y} = 4$			10	6

$$B_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

$$B_1 = \frac{6}{10}$$

$$B_1 = 0.6$$

$$y^{\wedge} = 2.2 + 0.6(x)$$

$$= 2.2 + 0.6(3)$$

$$y = B_0 + B_1 x$$

$$4 = B_0 + (0.6)(3)$$

$$B_0 = 4 - 0.6(3)$$

$$= 4 - 1.8$$

$$B_0 = 2.2$$

$$y = 2.2 + 0.6x$$

Standard Error:

$$\text{Standard error} = \sqrt{\frac{\sum (y^{\wedge} - y)^2}{(n-2)}}$$

$y^1$	$y^1 - y$	$(y^1 - y)^2$
2.8	0.8	0.64
3.4	-0.6	0.36
4	-1	0.16
4.6	0.6	0.36
5.2	0.2	0.04
		$\Sigma (y^1 - y)^2 = 2.4$

$$\text{standard error} = \sqrt{\frac{\Sigma (y^1 - y)^2}{n - 2}}$$

$$= \sqrt{\frac{2.4}{3}}$$

$$= 0.89$$

It should be less than 1

- 2) The Values of  $x$  and corresponding values of  $y$  are shown below.

$x$       0      1      2      3      4

$y$       2      3      5      4      6

- (i) Find the least square regression line  $y = ax + b$   
(ii) Estimate the value of  $y$ , when  $x = 10$ .

$x$	$y$	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(x - \bar{x})(y - \bar{y})$
0	2	-2	-2	4	4
1	3	-1	-1	1	1
2	5	0	1	0	0
3	4	1	0	1	0
4	6	2	2	4	4

$\Sigma x: 10$      $\Sigma y: 20$

$\Sigma (x - \bar{x})^2$

10

$\Sigma (x - \bar{x})(y - \bar{y})$

9

$\bar{x}: 2$      $\bar{y}: 4$

$$B_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

$$B_1 = \frac{9}{10}$$

$$B_1 = 0.9$$

$$y = B_0 + B_1 x$$

$$4 = B_0 + (0.9) \times 2$$

$$4 = B_0 + 1.8$$

$$B_0 = 2.2$$

$$y = 2.2 + 0.9x$$

$$y = 2.2 + 0.9(10)$$

$$y = 11.2$$

(ii)