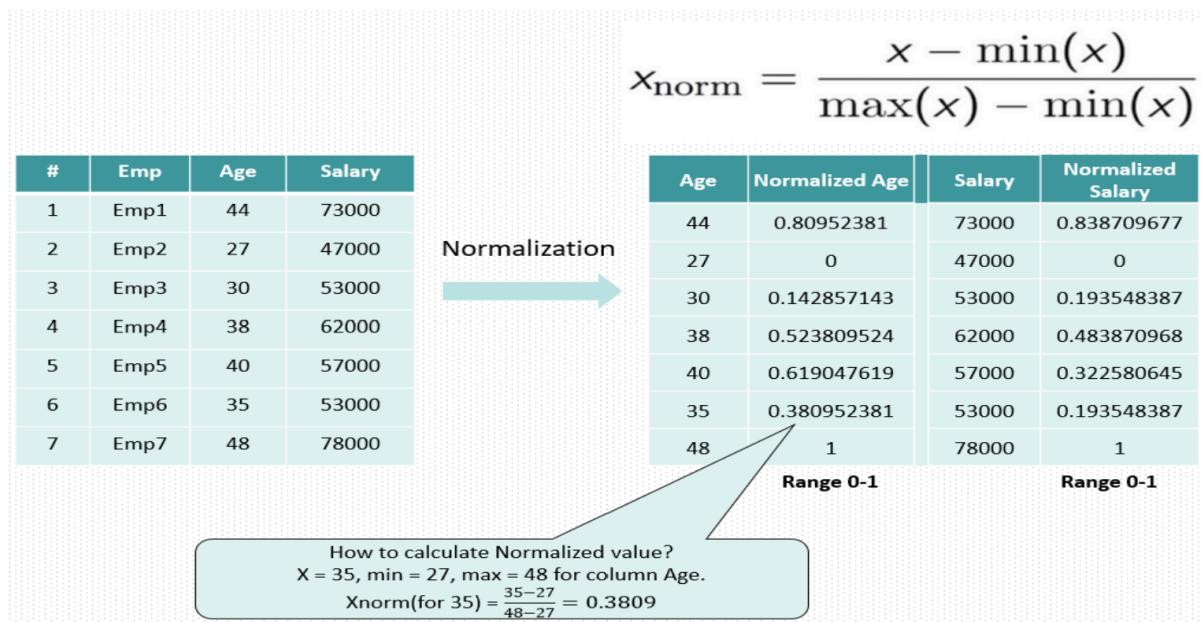


UNIT 4 DATA ANALYSIS

NORMALIZATION

Data Transformation: Apply transformations to make the data more suitable for analysis, such as normalisation or scaling.



MEAN

The most common measures the **mean** or average

1. **The mean or Average :** To calculate the average \bar{X} of the set of observation, add their value and divide by the number of observation

$$\text{Mean}, \bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

Let the mean of $x_1, x_2, x_3, \dots, x_n$ be A, then what is the mean of:

- 1. $(x_1+k), (x_2+k), (x_3+k), \dots, (x_n+k)$? Mean = A+K
- 2. $(x_1-k), (x_2-k), (x_3-k), \dots, (x_n-k)$? Mean = A-K
- 3. $kx_1, kx_2, kx_3, \dots, kx_n$? Mean = AK

PROBLEM 1

If the heights of 5 people are 142 cm, 150 cm, 149 cm, 156 cm and 153 cm.
Find the mean height.

$$\text{Mean height} = \frac{142 + 150 + 149 + 156 + 153}{5} = 150$$

$$\text{Mean, } \bar{x} = \frac{x_1 f_1 + x_2 f_2 + \dots + x_n f_n}{f_1 + f_2 + \dots + f_n}$$

PROBLEM 2

Find the mean of the following distribution

x_i	f_i	$x_i f_i$
4	5	20
6	10	60
9	10	90
10	7	70
15	8	120
	$\sum f_i = 40$	$\sum x_i f_i = 360$

$$\text{Mean} = \frac{\sum x_i f_i}{\sum f_i} = \frac{360}{40} = 9$$

PROBLEM 3

Find the average number of patients visiting the hospital in a day.

Number of patients	Number of days visiting hospital
0-10	2
10-20	6
20-30	9
30-40	7
40-50	4
50-60	2

Class mark (x_i)	frequency (f_i)	$x_i f_i$
5	2	10
15	6	90
25	9	225
35	7	245
45	4	180
55	2	110
Total	$\sum f_i = 30$	$\sum f_i x_i = 860$

$$\text{Mean} = \frac{\sum x_i f_i}{\sum f_i} = \frac{860}{30} = 28.97$$

MEDIAN

The median M is the midpoint of a distribution, the number such that half the observations are smaller, and the other half are larger.

STEPS

To Find the Median

1. Sort all the observation in order of size from smallest to largest
2. If the number of observations n is odd, the median M is the centre observation in the ordered list e.g. $M = n+1/2$ the obs.
3. If the number of observations n is even, the mean of two centre observation in the ordered list

Median (odd numbers)

$$\text{Median} = \left(\frac{n+1}{2} \right)^{\text{th}} \text{observation}$$

PROBLEM 4

56, 67, 54, 34, 78, 43, 23. What is the median?
Here, n (no. of observations) = 7

Arranging in ascending order, we get: 23, 34, 43, 54, 56, 67, 78.

$$\text{Median Position} = \frac{7+1}{2} = 4$$

$$\text{Median} = 54$$

Median (even)

$$\text{Median} = \frac{\frac{n}{2}^{\text{th}} \text{ obs.} + (\frac{n}{2} + 1)^{\text{th}} \text{ obs.}}{2}$$

PROBLEM 5

Let's consider the data: 50, 67, 24, 34, 78, 43.

What is the median?

Arranging in ascending order, we get: 24, 34, 43, 50, 67, 78.

$$\frac{6}{2} = 4$$

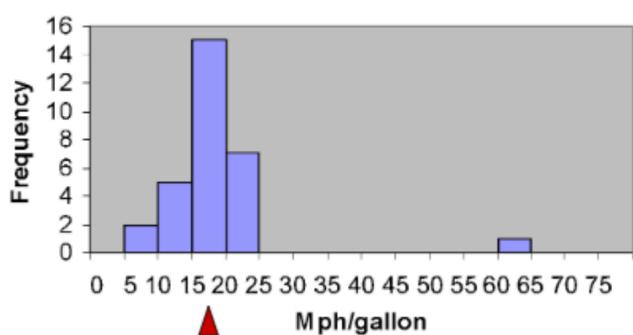
$$\text{Median} = \frac{43+50}{2} = 46.5$$

MODE

The Mode is the observation value with the highest frequency (observation with maximum frequency)

PROBLEM 6

Distribution of city fuel consumption



An average, the cars under study drive 18.9 miles per gallon, and 50% of cars under study drive at least 18 miles per gallon

PROBLEM 7

6, 8, 9, 3, 4, 6, 7, 6, 3 the value 6 appears the most number of times.
Thus, mode = 6.

Bimodal List

List A = {1, 2, 3, 3, 4, 4, 5, 6}

Mode [A] = {3, 4}

List A has 2 modes.
Therefore, it is a **bimodal list**.

Trimodal List

List B = {1, 2, 3, 3, 4, 4, 5, 5, 6}

Mode [B] = {3, 4, 5}

List B has 3 modes.
Therefore, it is a **trimodal list**.

$$\text{Mode} = l + \left[\frac{f_m - f_1}{2f_m - f_1 - f_2} \right] \times h$$

where, l= lower limit of modal class,
fm= frequency of modal class,
f1= frequency of class preceding modal class,
f2= frequency of class succeeding modal class,
h= class width

PROBLEM 8

Marks Obtained	0-20	20-40	40-60	60-80	80-100
Number of students	5	10	12	6	3

Find the mode of the given data:

The highest frequency = 12, so the modal class is 40-60.

l = lower limit of modal class = 40

f_m = frequency of modal class = 12

f_1 = frequency of class preceding modal class = 10

f_2 = frequency of class succeeding modal class = 6

h = class width = 20

$$\begin{aligned}\text{Mode} &= l + \left[\frac{f_m - f_1}{2f_m - f_1 - f_2} \right] \times h \\ &= 40 + \left[\frac{12 - 10}{2 \times 12 - 10 - 6} \right] \times 20 \\ &= 40 + \left[\frac{2}{8} \right] \times 20 \\ &= 45\end{aligned}$$

The relation between Mean, Median and Mode

PROBLEM 9

$$2 \text{ mean} + \text{Mode} = 3 \text{ Median}$$

We have a data whose mode == 65 and median == 61.6.
Find Mean

$$\begin{aligned}2 \text{Mean} + \text{Mode} &= 3 \text{ Median} \\ \therefore 2 \text{Mean} &= 3 \times 61.6 - 65 \\ \therefore 2 \text{Mean} &= 119.8 \\ \Rightarrow \text{Mean} &= \frac{119.8}{2} \\ \Rightarrow \text{Mean} &= 59.9\end{aligned}$$

Quartiles are statistical measures that divide a dataset into four equal parts, each representing 25% of the data. .

- ✓ **First Quartile (Q1):** Also known as the lower quartile, it represents the 25th percentile of the data.
- ✓ **Second Quartile (Q2):** This is the median of the dataset and represents the 50th percentile. It divides the dataset into two equal halves.
- ✓ **Third Quartile (Q3):** Also known as the upper quartile, it represents the 75th percentile of the data. This means that 75% of the data values fall below Q3.
- ✓ **Interquartile Range (IQR):** This is the range between the first and third quartiles (Q3 - Q1). The IQR measures the spread of the middle 50% of the data and is used to identify outliers.

n is the number of data points

Calculation of Quartiles

- ✓ $Q_1 = \left(\frac{1}{4}(n+1)\right)^{\text{th}} \text{ position}$
- ✓ $Q_2 = \left(\frac{1}{2}(n+1)\right)^{\text{th}} \text{ position}$
- ✓ $Q_3 = \left(\frac{3}{4}(n+1)\right)^{\text{th}} \text{ position}$

PROBLEM 10

Consider the dataset: 3, 7, 8, 5, 12, 7, 9, 10, 14, 21

Sort the Data: 3, 5, 7, 7, 8, 9, 10, 12, 14, 21

Find Q1: The position is $\frac{1}{4}(10+1)=2.75$ so Q_1 is between the 2nd and 3rd values, which is 5 and 7. Interpolating, $Q_1 = 5 + 0.75 \times (7 - 5) = 6.5$.

Find Q2: The position is $\frac{1}{2}(10+1)=5.5$, so Q_2 is between the 5th and 6th values, which is 8 and 9. Interpolating, $Q_2 = 8 + 0.5 \times (9 - 8) = 8.5$.

Find Q3: The position is $\frac{3}{4}(10+1)=8.25$ so Q_3 is between the 8th and 9th values, which is 12 and 14. Interpolating, $Q_3 = 12 + 0.25 \times (14 - 12) = 12.5$.

IQR: IQR=Q3–Q1=12.5–6.0 (**Interquartile Range**)

PROBLEM 11

Find Interquartile Range of 5 ,6, 1, 2, 15, 12, 27, 19, 18 (Odd numbers)

- ✓ Step 1: Put the numbers in order. 1, 2, 5, 6, 7, 9, 12, 15, 18, 19, 27.
- ✓ Step 2: Find the median. 1, 2, 5, 6, 7, **9**, 12, 15, 18, 19, 27.
- ✓ Step 3: Place parentheses around the numbers above and below the median. Not necessary statistically, but it makes Q1 and Q3 easier to spot. (1, 2, 5, 6, 7), **9**, (12, 15, 18, 19, 27).
- ✓ Step 4: Find Q1 and Q3 (Think of Q1 as a median in the lower half of the data and think of Q3 as a median for the upper half of data.) (1, 2, **5**, 6, 7), **9**, (12, 15, **18**, 19, 27). Q1 = 5 and Q3 = 18.
- ✓ Step 5: Subtract Q1 from Q3 to find the **interquartile range**. **18 – 5 = 13**.

PROBLEM 12

Find Interquartile Range of 4,4,10,11,15,7,14,12,6(Odd numbers)

- ✓ Step 1: Put the numbers in order. 4,4,6,7,10,11,12,14,15
- ✓ Step 2: Find the median. 4,4,6,7,10,11,12,14,15
- ✓ Step 3: Place parentheses around the numbers above and below the median. Not necessary statistically, but it makes Q1 and Q3 easier to spot. (4,4,6,7), 10, (11,12,14,15).
- ✓ Step 4: Find Q1 and Q3
- ✓ (4,4,6,7), 10, (11,12,14,15). Q1 = 5 and Q3 = 13.
- ✓ Step 5: Subtract Q1 from Q3 to find the interquartile range. $13 - 5 = 8$.

PROBLEM 13

Find the IQR for the following data set: 3, 5, 7, 8, 9, 11, 15, 16, 20, 21
(Even)

Step 1: Put the numbers in order: 3, 5, 7, 8, 9, 11, 15, 16, 20, 21.
Step 2: Make a mark in the centre of the data: 3, 5, 7, 8, 9, | 11, 15, 16, 20, 21.

Step 3: Place parentheses around the numbers above and below the mark you made in Step 2—it makes Q1 and Q3 easier to spot.
(3, 5, 7, 8, 9), | (11, 15, 16, 20, 21).

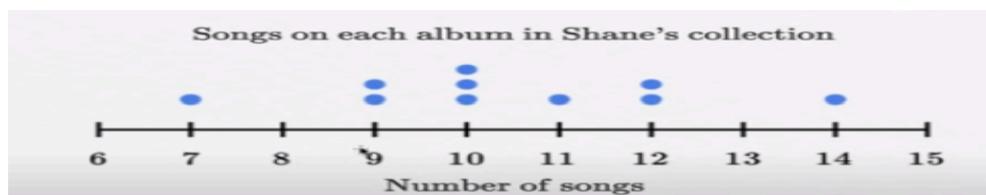
Step 4: Find Q1 and Q3

Q1 is the median (the middle) of the lower half of the data, and Q3 is the median (the middle) of the upper half of the data.

(3, 5, 7, 8, 9), | (11, 15, 16, 20, 21). Q1 = 7 and Q3 = 16.

Step 5: Subtract Q1 from Q3. $16 - 7 = 9$.

PROBLEM 14



✓ Step 1: Put the numbers in order. 7,9,9,10,10,10,10,11,12,12,12,14

✓ Step 2: Make a mark in the centre of the data:

7,9,9,10,10|10,11,12,12,14

Step 3: Place parentheses around the numbers above and below the mark you made in Step 2—it makes Q1 and Q3 easier to spot.

(7,9,9,10,10), | (10,11,12,12,14).

Step 4: Find Q1 and Q3

(7,9,**9**,10,10), | (10,11,**12**,12,14). Q1 = 9 and Q3 = 12.

Step 5: Subtract Q1 from Q3. $12 - 9 = 3$.

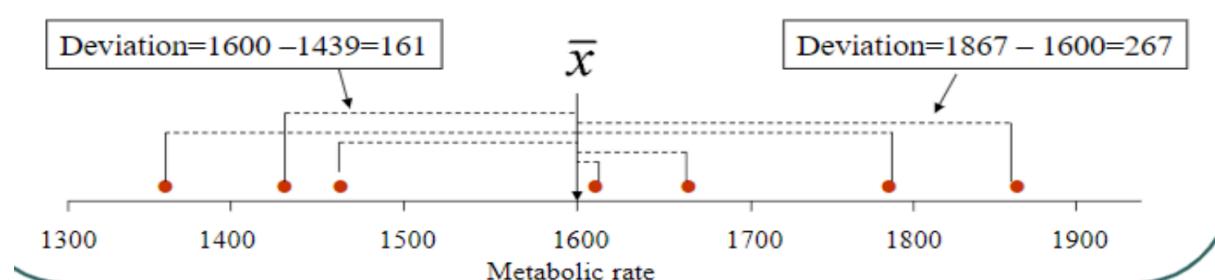
Standard Deviation

- ✓ If a distribution is symmetric
 - ✓ Use the average to measure the centre and the standard deviation to measure the spread
 - ✓ The SD measures how far the observations are from the average

PROBLEM 15

- ✓ E.g. A person's Metabolic rate = rate at which the body consumes energy Rates of 7 men in a study on dieting 1792, 1666, 1614, 1460, 1867, 1439, 1362.

$$\boxed{x = 1600 \text{ and } \text{sd} = 189.24}$$



Standard Deviation

$$s = \sqrt{\frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n-1}}$$

VARIANCE

The variance of an observed Variables is defined as the square of the standard deviation

$$\text{Variance} = s^2$$

Properties of Standard Deviation

- It measures the spread about the mean
- Only used in association with the mean, Good descriptive measure for Symmetric distributions
- If $s = 0$ all observations have the same values
- It is NOT a resistant measure, a few extreme observations may affect its values

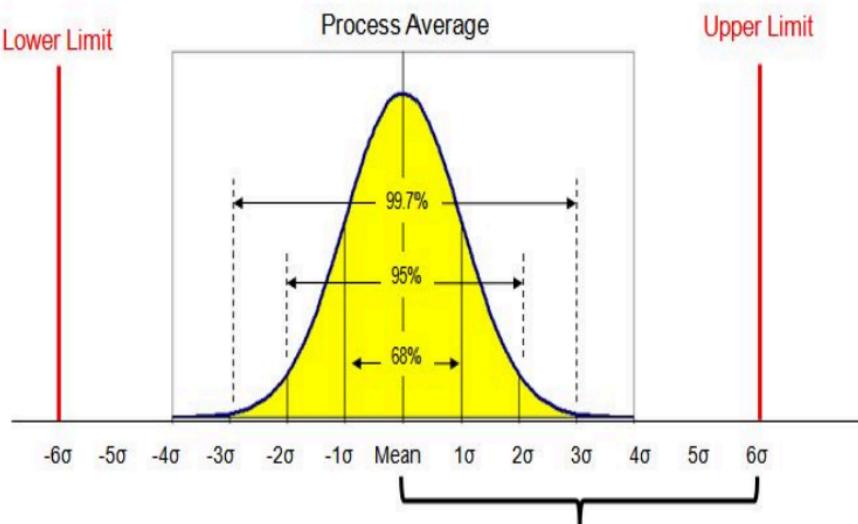
Interpreting the SD

PROBLEM 16

For many lists of observation- especially if their histogram is bell-shaped

1. Roughly 6* % of the observation in the list lie within 1 SD of the average
2. 95% of the observation lie within 2 SD of the average

Sigma	Defects	Confidence
1S	31.8%	68%
2S	6.7%	95%
3S	2.7%	99.7%
4S	0.6%	99.99%
5S	0.02%	
6S	0.00034%	



PROBLEM 17

Consider the following three data sets A, B and C and determine standard Deviation

$$A = \{9, 10, 11, 7, 13\}$$

$$B = \{10, 10, 10, 10, 10\}$$

$$C = \{1, 1, 10, 19, 19\}$$

Step 1: Mean

$$\text{Mean of Data set A} = (9+10+11+7+13)/5 = 10$$

$$\text{Mean of Data set B} = (10+10+10+10+10)/5 = 10$$

$$\text{Mean of Data set C} = (1+1+10+19+19)/5 = 10$$

$$\text{Standard Deviation Data set A} = \sqrt{[((9-10)^2 + (10-10)^2 + (11-10)^2 + (7-10)^2 + (13-10)^2)/5]} = 2$$

$$\text{Standard Deviation Data set B} = \sqrt{[((10-10)^2 + (10-10)^2 + (10-10)^2 + (10-10)^2 + (10-10)^2)/5]} = 0$$

$$\text{Standard Deviation Data set C} = \sqrt{[((1-10)^2 + (1-10)^2 + (10-10)^2 + (19-10)^2 + (19-10)^2)/5]} = 8.05$$

$$\text{Variance} = S^2$$

$$\text{Variance of A} = 2^2 = 4$$

$$\text{Variance of B} = 0^2 = 0$$

$$\text{Variance of C} = 8^2 = 64$$

$$\text{Coefficient of Variance} = CV = \text{standard deviation} / \text{Mean}$$

$$CV \text{ of A} = 2/10 = 0.2$$

$$CV \text{ of B} = 0/10 = 0$$

$$CV \text{ of C} = 8.05/10 = 0.805$$

PROBLEM 18

Consider the following three data sets A, B and C and determine standard Deviation

$$A = \{9, 10, 11, 7, 13\}$$

$$B = \{10, 10, 10, 10, 10\}$$

$$C = \{1, 1, 10, 19, 19\}$$

Step 1: Mean

$$\text{Mean of Data set A} = (9+10+11+7+13)/5 = 10$$

$$\text{Mean of Data set B} = (10+10+10+10+10)/5 = 10$$

$$\text{Mean of Data set C} = (1+1+10+19+19)/5 = 10$$

$$\text{Standard Deviation Data set A} = \sqrt{[((9-10)^2 + (10-10)^2 + (11-10)^2 + (7-10)^2 + (13-10)^2)/5]} = 2$$

$$\text{Standard Deviation Data set B} = \sqrt{[((10-10)^2 + (10-10)^2 + (10-10)^2 + (10-10)^2 + (10-10)^2)/5]} = 0$$

$$\text{Standard Deviation Data set C} = \sqrt{[((1-10)^2 + (1-10)^2 + (10-10)^2 + (19-10)^2 + (19-10)^2)/5]} = 8.05$$

$$\text{Variance} = S^2$$

$$\text{Variance of A} = 2^2 = 4$$

$$\text{Variance of B} = 0^2 = 0$$

$$\text{Variance of C} = 8^2 = 64$$

$$\text{Coefficient of Variance} = CV = \text{standard deviation / Mean}$$

$$CV \text{ of A} = 2/10 = 0.2$$

$$CV \text{ of B} = 0/10 = 0$$

$$CV \text{ of C} = 8.05/10 = 0.805$$

PROBLEM 19

Calculate the standard deviation from the following distribution of marks by using all the methods.

Marks	>No. of Students
>1-3	>40
>3-5	>30
>5-7	>20
>7-9	>10

$$\bar{X} = \frac{\sum fX}{\sum f} = \frac{400}{100} = 4$$

>Marks	> f	> X	> fX	> $(X - \bar{X})^2$	> $f(X - \bar{X})^2$
>1–3	>40	>2	>80	>4	>160
>3–5	>30	>4	>120	>0	>0
>5–7	>20	>6	>120	>4	>80
>7–9	>10	>8	>80	>16	>160
>Total	>100	>	>400	>	>400

$\sum fX = 400$

$$S = \sqrt{\frac{\sum f(X - \bar{X})^2}{\sum f}} = \sqrt{\frac{400}{100}} = \sqrt{4} = 2 \text{ marks}$$

$$\text{Variance} = S^2 = 4$$

$$\text{CV} = S/x = 2/4 = 0.5$$

Statistical Estimation

- ✓ **Parameter:** A numerical characteristic of a population, such as the mean (μ) or standard deviation (σ)
- ✓ **Estimator:** A statistic (a function of the sample data) used to estimate a population parameter.
- ✓ **Estimate:** The value obtained from an estimator. For instance, if the sample mean is 50, then 50 is the estimate of the population mean.
- ✓ **Point Estimate:** A single value given as an estimate of a population parameter. For example, using the sample mean to estimate the population mean.
- ✓ **Interval Estimate:** A range of values within which the population parameter is expected to lie, with a certain level of confidence. For example, a 95% confidence interval for the mean.

Factors Affecting Confidence Interval Estimates

- ✓ The factors that determine the width of a confidence interval are
 - ✓ The sample size, n
 - ✓ The variability in the population, usually SD estimated
 - ✓ The desired level of confidence

Confidence Intervals

Confidence intervals (CIs) offer a range of values where a parameter likely falls, quantifying uncertainty.

For example, a 95% confidence interval means that if you were to take 100 different samples and calculate a CI for each, approximately 95 of those intervals would contain the true population parameter

Point Estimate:

The point estimate is a single value from the sample used to estimate the population parameter. (Point estimate = mean)

$$\text{Point Estimation} = \bar{x} = \sum x/n$$

Margin of Error:

The margin of error is the range around the point estimate that accounts for variability and provides the interval

$$\text{Margin error} = z / \sqrt{n} \quad (\text{z-score or t-score})$$

Confidence intervals (CIs)

Confidence Interval for the Mean (known σ and z distribution)

$$\mu = \bar{x} \pm Z \frac{\sigma}{\sqrt{n}}$$

Confidence level	Critical (z) value to be used in confidence interval calculation
50%	0.67449
75%	1.15035
90%	1.64485
95%	1.95996
97%	2.17009
99%	2.57583
99.9%	3.29053

PROBLEM 20

1. A survey was taken of USA companies that do business with firms in India. One of the questions on the survey was: Approximately how many years has your company been trading with firms in India? A random sample of 44 responses to this question yielded a mean of 10.455 years. Suppose the population standard deviation for this question is 7.7 years. Using this information, construct a 90% confidence interval for the mean number of years that a company has been trading in India for the population of US companies trading with firms in India.

$\bar{x} = 10.455 \quad \sigma = 7.7 \quad n = 44 \quad z = 1.645$ (from table for 90% confidence)

$$10.455 - 1.645 \frac{7.7}{\sqrt{44}} \leq \mu \leq 10.45 + 1.645 \frac{7.7}{\sqrt{44}}$$
$$8.545 \leq \mu \leq 12.365$$

The 90% confidence interval for the mean number of years that a company has been trading with firms in India is approximately [8.545, 12.365] years

PROBLEM 21

2. A study is conducted in a company that employs 800 engineers. A random sample of 50 engineers reveals that the average sample age is 34.3 years. Historically, the population standard deviation of the age of the company's engineers is approximately 8 years. Construct a 98% confidence interval to estimate the average age of all the engineers in this company.

- ✓ Given data: $N = 800$ $n = 50$ $\bar{x} = 34.3$ $\sigma = 8$ $z = 2.33$ from the table for 98% confidence)

$$\bar{x} - Z \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}} \leq \mu \leq \bar{x} + Z \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$

$$34.3 - 2.33 \frac{8}{\sqrt{50}} \sqrt{\frac{800-50}{800-1}} \leq \mu \leq 34.3 + 2.33 \frac{8}{\sqrt{50}} \sqrt{\frac{800-50}{800-1}}$$

$$31.75 \leq \mu \leq 36.85$$

The 98% confidence interval to estimate the average age of all engineers in the company is approximately [31.67, 36.93] years

Estimating the Mean of a Normal Population: Sample Size is Small and Population is unknown

- ✓ When estimating the mean of a normal population with a small sample size and an unknown population standard deviation, you use the **t-distribution** instead of the normal distribution.
sample size (n) < 30

STEPS

- ✓ Steps to Construct the Confidence Interval
 - ✓ Determine the Sample Mean (\bar{x}) =
 - ✓ Determine the Sample Standard Deviation (s)
 - ✓ Find the t-Score based on degree of freedom and confidence level or $t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$
 - ✓ Construct the Confidence Interval (CI)

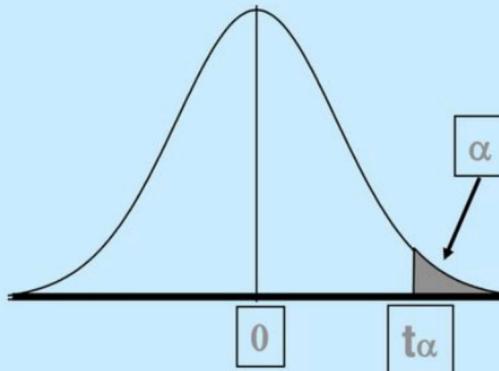
Table of Critical Values of t

Critical values of t for two-tailed tests

Significance level (α)

Degrees of freedom (df)	.2	.15	.1	.05	.025	.01	.005	.001
1	3.078	4.165	6.314	12.706	25.452	63.657	127.321	636.619
2	1.886	2.282	2.920	4.303	6.205	9.925	14.089	31.599
3	1.638	1.924	2.353	3.182	4.177	5.841	7.453	12.924
4	1.533	1.778	2.132	2.776	3.495	4.604	5.598	8.610
5	1.476	1.699	2.015	2.571	3.163	4.032	4.773	6.869
6	1.440	1.650	1.943	2.447	2.969	3.707	4.317	5.959
7	1.415	1.617	1.895	2.365	2.841	3.499	4.029	5.408
8	1.397	1.592	1.860	2.306	2.752	3.355	3.833	5.041
9	1.383	1.574	1.833	2.262	2.685	3.250	3.690	4.781
10	1.372	1.559	1.812	2.228	2.634	3.169	3.581	4.587
11	1.363	1.548	1.796	2.201	2.593	3.106	3.497	4.437
12	1.356	1.538	1.782	2.179	2.560	3.055	3.428	4.318
13	1.350	1.530	1.771	2.160	2.533	3.012	3.372	4.221
14	1.345	1.523	1.761	2.145	2.510	2.977	3.326	4.140
15	1.341	1.517	1.753	2.131	2.490	2.947	3.286	4.073
16	1.337	1.512	1.746	2.120	2.473	2.921	3.252	4.015
17	1.333	1.508	1.740	2.110	2.458	2.898	3.222	3.965
18	1.330	1.504	1.734	2.101	2.445	2.878	3.197	3.922
19	1.328	1.500	1.729	2.093	2.433	2.861	3.174	3.883
20	1.325	1.497	1.725	2.086	2.423	2.845	3.153	3.850
21	1.323	1.494	1.721	2.080	2.414	2.831	3.135	3.819
22	1.321	1.492	1.717	2.074	2.405	2.819	3.119	3.792
23	1.319	1.489	1.714	2.069	2.398	2.807	3.104	3.768
24	1.318	1.487	1.711	2.064	2.391	2.797	3.091	3.745
25	1.316	1.485	1.708	2.060	2.385	2.787	3.078	3.725
26	1.315	1.483	1.706	2.056	2.379	2.779	3.067	3.707
27	1.314	1.482	1.703	2.052	2.373	2.771	3.057	3.690
28	1.313	1.480	1.701	2.048	2.368	2.763	3.047	3.674
29	1.311	1.479	1.699	2.045	2.364	2.756	3.038	3.659
30	1.310	1.477	1.697	2.042	2.360	2.750	3.030	3.646
40	1.303	1.468	1.684	2.021	2.329	2.704	2.971	3.551
50	1.299	1.462	1.676	2.009	2.311	2.678	2.937	3.496
60	1.296	1.458	1.671	2.000	2.299	2.660	2.915	3.460
70	1.294	1.456	1.667	1.994	2.291	2.648	2.899	3.435
80	1.292	1.453	1.664	1.990	2.284	2.639	2.887	3.416
100	1.290	1.451	1.660	1.984	2.276	2.626	2.871	3.390
1000	1.282	1.441	1.646	1.962	2.245	2.581	2.813	3.300
Infinite	1.282	1.440	1.645	1.960	2.241	2.576	2.807	3.291

df	t _{0.100}	t _{0.050}	t _{0.025}	t _{0.010}	t _{0.005}
1	3.078	6.314	12.706	31.821	63.656
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
23	1.319	1.714	2.069	2.500	2.807
24	1.311	1.711	2.064	2.492	2.797
25	1.316	1.708	2.060	2.485	2.787
29	1.311	1.699	2.045	2.462	2.756
30	1.310	1.697	2.042	2.457	2.750
40	1.303	1.684	2.021	2.423	2.704
60	1.296	1.671	2.000	2.390	2.660
120	1.289	1.658	1.980	2.358	2.617
∞	1.282	1.645	1.960	2.327	2.576



With df = 24 and $\alpha = 0.05$,
 $t_\alpha = 1.711$.

PROBLEM 22

3. The owner of a large equipment rental company wants to make rather quick estimate of the average number of days a piece of ditch digging equipment is rented out per person per time. The company has records of all rentals, but the amount of time required to conduct an audit of all accounts would be prohibitive. The owner decides to take a random sample of rental invoices. Fourteen different rentals of ditch diggers are selected randomble from the files, yielding the following data. She uses these data construct a 99% confidence interval to estimate the average number of days that a ditch digger is rented and assumes that the number of datas per rental is normally distributed in the population. The data: 3, 1, 3, 2, 5, 1, 2, 1, 4, 2, 1, 3, 1, 1

Given information:

$$n = 14$$

Sample data: 3, 1, 3, 2, 5, 1, 2, 1, 4, 2, 1, 3, 1, 1

Confidence level = 99%

Calculate the Sample Mean

$$\bar{x} = \frac{\sum x_i}{n} = \frac{3 + 1 + 3 + 2 + 5 + 1 + 2 + 1 + 4 + 2 + 1 + 3 + 1 + 1}{14} = 1.79$$

Calculate the Sample Standard Deviation (s)

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} = \frac{(3-1.79)^2 + (1-1.79)^2 + (3-1.79)^2 + (2-1.79)^2 + (5-1.79)^2 + \dots + (1-1.79)^2}{14-1} = 0.52$$

$$\sigma = 0.72$$

$$Df = n-1 = 13$$

$$\alpha/2 = \frac{1-0.99}{2} = 0.005$$

$t = 3.012$ (from the table based on α and df)

t Table

cum. prob	$t_{.50}$	$t_{.75}$	$t_{.80}$	$t_{.85}$	$t_{.90}$	$t_{.95}$	$t_{.975}$	$t_{.99}$	$t_{.995}$	$t_{.999}$	$t_{.9995}$
one-tail	0.50	0.25	0.20	0.15	0.10	0.05	0.025	0.01	0.005	0.001	0.0005
two-tails	1.00	0.50	0.40	0.30	0.20	0.10	0.05	0.02	0.01	0.002	0.001
11	0.000	0.697	0.876	1.088	1.363	1.796	2.201	2.718	3.106	4.025	4.437
12	0.000	0.695	0.873	1.083	1.356	1.782	2.179	2.681	3.055	3.930	4.318
13	0.000	0.694	0.870	1.079	1.350	1.771	2.160	2.650	3.012	3.852	4.221
14	0.000	0.692	0.868	1.076	1.345	1.761	2.145	2.624	2.977	3.787	4.140
15	0.000	0.691	0.866	1.074	1.341	1.753	2.131	2.602	2.947	3.733	4.073

$$\bar{x} - t_{\alpha/2, n-1} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + t_{\alpha/2, n-1} \frac{\sigma}{\sqrt{n}}$$

$$1.79 - 3.012 \frac{0.72}{\sqrt{14}} \leq \mu \leq 1.79 + 3.012 \frac{0.72}{\sqrt{14}}$$

$$1.22 \leq \mu \leq 2.36$$

The 99% confidence interval for the average number of days that a ditch digger is rented is approximately [1.22, 2.36] days. This interval provides a range within which the true population mean is expected to lie with 99% confidence, assuming normal distribution of rental days.

Confidence Interval to Estimate the Population Proportion

To construct a confidence interval to estimate the population proportion, you use the sample proportion and the normal approximation if the sample size is large enough. $n > 30$

- ✓ Calculate the Sample Proportion ($\hat{P} = \frac{x}{n}$) where x is the number of successes and n is the sample size
- ✓ Determine the z-Score for the Desired Confidence Level based on confidential level
- ✓ Construct the Confidence Interval (CI)

$$\hat{P} - z\alpha_{/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq \mu \leq \hat{P} + z\alpha_{/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

PROBLEM 23

Assume you have a sample of 200 people where 80 of them favor a new policy. Construct a 95% confidence interval for the proportion of people who favor the policy.

1. Calculate the Sample Proportion ($\hat{P} = \frac{x}{n} = \frac{80}{200} = 0.40$)

2. Determine the z-Score for a 95% Confidence Level $z = 1.960$

3. Calculate the Interval $\hat{P} - z\alpha_{/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq \mu \leq \hat{P} + z\alpha_{/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

$$0.40 - 1.96 \sqrt{\frac{0.40(1-0.40)}{200}} \leq \mu \leq 0.40 + 1.96 \sqrt{\frac{0.40(1-0.40)}{200}}$$

$$0.332 \leq \mu \leq 0.468$$

The 95% confidence interval for the proportion of people who favor the new policy is approximately [0.332, 0.468]. This interval provides a range within which the true population proportion is expected to lie with 95% confidence.

PROBLEM 24

A clothing company produces men's jeans. The jeans are made and sold with either a regular cut or a boot cut. In an effort to estimate the proportion of their men's jeans market in Oklahoma city that prefers boot-cut jeans, the analyst takes a random sample of 423 jeans sales from the company two Oklahoma city retail outlets. Only 72 of the sales were for boot-cut jeans. Construct a 90% confidence interval to estimate the proportion of the population in Oklahoma city who prefer boot-cut jeans.

Given $n = 423$ $x = 72$ $\hat{p} = \frac{x}{n} = \frac{72}{423} = 0.17$ $z = 1.645$ for 90% confidence from the table

$$\hat{P} - z\alpha/2 \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq \mu \leq \hat{P} + z\alpha/2 \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$0.17 - 1.642 \sqrt{\frac{0.17(1-0.17)}{423}} \leq \mu \leq 0.17 + 1.642 \sqrt{\frac{0.17(1-0.17)}{423}} \quad \mathbf{0.14 \leq \mu \leq 0.20}$$

The 90% confidence interval to estimate the proportion of the population in Oklahoma City who prefer boot-cut jeans is approximately [0.14, 0.2]. This interval provides a range within which the true proportion of the population is expected to lie with 90% confidence.

Hypothesis Testing

- Hypothesis is considered as intelligent guess or prediction, that gives directional to the researcher to answer the research question.
- Hypothesis are defined as formal statement of the tentative or expected prediction or explanation of the relationship between two or more variables in a specified population.
- A hypothesis is formal tentative statement of the expected relationship between two or more variable under study.
- A hypothesis helps to translate the research problem and objective into a clear explaining or prediction of the expected results or outcome of the study.

Hypothesis testing refers to

1. Making an assumption about a population parameters
2. Collecting sample data
3. Calculating a sample statistic
4. Using the sample statistic to evaluate the hypothesis.

Null Hypothesis (H_0): State the hypothesized value of the parameter before sampling. The assumption we wish to test or the assumption we are trying to reject.

Population mean $\mu = 20$ There is no difference between coke and diet coke

Alternative Hypothesis (H_A): All possible alternatives other than the null hypothesis

e.g $\mu > 20$ and $\mu < 20$

There is a difference between coke and diet coke.

Selecting and Interpreting Significance Level

- ✓ Selecting and interpreting the significance level (α) is crucial in hypothesis testing.
- ✓ The significance level, denoted as α , is the probability of rejecting the null hypothesis when it is actually true. It represents the threshold for determining whether the observed results are statistically significant.
- ✓ Common Choices for α
 - ✓ 0.10: Indicates a 10% risk of a Type I error
 - ✓ 0.05: Indicates a 5% risk of a Type I error (most commonly used).
 - ✓ 0.01: Indicates a 1% risk of a Type I error.
- ✓ Interpreting the Significance Level (Comparison with α)
 - ✓ If $p \leq \alpha$, reject the null hypothesis
 - ✓ If $p > \alpha$, do not reject the null hypothesis
- ✓ Implications of α Choices
 - ✓ **Higher α :** More likely to reject the null hypothesis, increasing the risk of Type I errors.
 - ✓ **Lower α :** Less likely to reject the null hypothesis, reducing the risk of Type I errors but increasing the risk of Type II errors.

Parametric Tests

1. Parametric Test
2. Non Parametric tests

Parametric test Parametric tests that make specific assumptions about the parameters (the mean and variance) of the population distribution from which the data is drawn.

Applications

- ✓ Used for continuous variables ($y = f(x)$)
- ✓ Used when data are measured on appropriate interval or ratio scale of measurement
- ✓ Data should follow normal distribution

Parametric tests that make specific assumptions about the parameters (the mean and variance) of the population distribution from which the data is drawn.

Normality: Many parametric tests assume that the data follows a normal distribution (bell-shaped curve)

Homogeneity of Variance: Some parametric tests assume that different groups have similar variances (σ^2)

- ✓ **t-Test:** Used to compare the means of two groups to determine if they are statistically significantly different from each other.
 - ✓ Independent t-test (comparing two separate groups)
 - ✓ Paired t-test (for comparing two related groups)
 - ✓ One-Sample t-Test (To compare the mean of a single group to a known value or population mean)
- ✓ **ANOVA (Analysis of Variance):** Used to compare the means of three or more groups to see if there is a significant difference among them.
 - ✓ One-way ANOVA (for one factor)
 - ✓ Two-way ANOVA (for two factors)
- ✓ **Pearson Correlation Coefficient:** Measures the strength and direction of the linear relationship between two continuous variables.

- ✓ **Regression Analysis:** Used to understand the relationship between a dependent variable and one or more independent variables.
 - ✓ Linear regression
 - ✓ Multiple and logistic regressions

- ✓ **Z-Test:** Used to compare the sample mean to the population mean when the sample size is large and the population variance is known.

Parametric Tests-t test

- ✓ Developed by Prof. W S Gosseti
- ✓ A t-test compares the difference between two mean of different groups to determine whether the difference is statically significant.

One Sample t-test

Assumptions:

- ✓ Population is normally distributed
- ✓ Sample is drawn from the population, and it should be random
- ✓ Known the population mean

Condition

- ✓ Population standard deviation is not known
- ✓ Size of the sample is small (<30)

Parametric Tests-t test-one sample

Steps to Perform a One-Sample t-Test:

1. State the Hypotheses

- **Null Hypothesis (H_0):** $\mu = \mu_0$ (the sample mean is equal to the hypothesized population mean)
- **Alternative Hypothesis (H_1):** $\mu \neq \mu_0$ (the sample mean is not equal to the hypothesized population mean)

2. Calculate the t-Statistic: $t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$

3. Determine the Degrees of Freedom (df = n-1)

4. Find the p-Value: Compare the t-statistic to the t-distribution with the appropriate degrees of freedom to find the p-value

5. Make a Decision

- **Reject H_0** if the p-value is less than the significance level (α), commonly set at 0.05.
- **Fail to Reject H_0** if the p-value is greater than α .

6. Interpret the Results

If you reject the null hypothesis, it means there is a significant difference between the sample mean and the hypothesized population mean. If you fail to reject it, there is not enough evidence to suggest a significant difference.

Parametric Tests-t test

PROBLEM 25

1. The following Data represents haemoglobin values in gm/dl for 10 patient
10.5 , 9, 6.5, 8, 11, 7, 7.5, 8.5, 9.5, 12 Is the mean value for patients significantly differ from the mean values of general population (12 gm/dl). Evaluate the role of change

Mean

$$\bar{x} = \frac{10.5 + 9 + 6.5 + 8 + 11 + 7 + 7.5 + 8.5 + 9.5 + 12}{10}$$

Standard

$$\sigma = \sqrt{\frac{(10.5-8.95)^2 + (9-8.95)^2 + (6.5-8.95)^2 + (8-8.95)^2 + (11-8.95)^2 + (7-8.95)^2 + (7.5-8.95)^2 + (8.5-8.95)^2 + (9.5-8.95)^2 + (12-8.95)^2}{10-1}}$$
$$= 1.802$$

Step 1: State the Hypotheses

Null Hypothesis (H_0): The mean hemoglobin value of the patient sample is equal to the mean hemoglobin value of the general population.

$$H_0: \mu = 12 \text{ gm/dl}$$

Alternative Hypothesis (H_1): The mean hemoglobin value of the patient sample is different from the mean hemoglobin value of the general population.

$$H_1: \mu \neq 12 \text{ gm/dl}$$

Step 2: Calculate the t-Statistic

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{8.95 - 12}{\frac{1.80201}{\sqrt{10}}} = -5.352$$

Step 3 Determine the Degrees of Freedom (df = n-1)

$$df = 10 - 1 = 9$$

Degrees of freedom (df)	.2	.15	.1	.05	.025	.01	.005	.001
1	1.378	1.963	3.078	6.314	12.708	31.821	63.857	318.309
2	1.061	1.386	1.886	2.920	4.203	6.965	9.925	22.327
3	0.978	1.250	1.638	2.353	3.182	4.541	5.841	10.215
4	0.941	1.190	1.533	2.132	2.776	3.747	4.604	7.773
5	0.920	1.156	1.476	2.015	2.571	3.365	4.032	5.893
6	0.906	1.134	1.440	1.943	2.447	3.143	3.707	5.208
7	0.896	1.119	1.415	1.895	2.365	2.998	3.499	4.785
8	0.889	1.108	1.387	1.860	2.306	2.896	3.355	4.501
9	0.883	1.100	1.363	1.833	2.281	2.865	3.312	4.297
10	0.879	1.093	1.371	1.812	2.260	2.764	3.298	4.144
11	0.876	1.088	1.363	1.796	2.201	2.718	3.206	4.025
12	0.873	1.083	1.356	1.782	2.179	2.681	3.055	3.930
13	0.870	1.079	1.350	1.771	2.160	2.650	3.012	3.852
14	0.868	1.076	1.345	1.761	2.145	2.624	2.977	3.787
15	0.866	1.074	1.341	1.753	2.131	2.602	2.947	3.733
16	0.865	1.071	1.337	1.746	2.120	2.583	2.921	3.686
17	0.863	1.069	1.333	1.740	2.110	2.567	2.898	3.646
18	0.862	1.067	1.330	1.734	2.101	2.552	2.878	3.610
19	0.861	1.066	1.328	1.728	2.093	2.537	2.858	3.580

Step 4: Then compare with tabulated value for 9 df and 5% level of significance it is = -1.8331. The calculated values > tabulated value

5. Make a Decision

Reject H₀ and concluded that there is a statistically significant difference between the mean of sample ad population mean and this difference is unlikely due to chance

6. Interpret the Results

The significant difference in mean hemoglobin values suggests that the patient group has a lower average hemoglobin level compared to the general population.

Degrees of freedom (df)	.2	.15	.1	.05	.025	.01	.005	.001
1	1.378	1.963	3.078	6.314	12.708	31.821	63.857	318.309
2	1.061	1.386	1.886	2.920	4.203	6.965	9.925	22.327
3	0.978	1.250	1.638	2.353	3.182	4.541	5.841	10.215
4	0.941	1.190	1.533	2.132	2.776	3.747	4.604	7.773
5	0.920	1.156	1.476	2.015	2.571	3.365	4.032	5.893
6	0.906	1.134	1.440	1.943	2.447	3.143	3.707	5.208
7	0.896	1.119	1.415	1.895	2.365	2.998	3.499	4.785
8	0.889	1.108	1.387	1.860	2.306	2.896	3.355	4.501
9	0.883	1.100	1.363	1.833	2.262	2.821	3.250	4.297
10	0.879	1.093	1.372	1.812	2.228	2.764	3.169	4.144
11	0.876	1.088	1.363	1.796	2.201	2.718	3.106	4.025
12	0.873	1.083	1.356	1.782	2.179	2.681	3.055	3.930
13	0.870	1.079	1.350	1.771	2.160	2.650	3.012	3.852
14	0.868	1.076	1.345	1.761	2.145	2.624	2.977	3.787
15	0.866	1.074	1.341	1.753	2.131	2.602	2.947	3.733
16	0.865	1.071	1.337	1.746	2.120	2.583	2.921	3.686
17	0.863	1.069	1.333	1.740	2.110	2.567	2.898	3.646
18	0.862	1.067	1.330	1.734	2.101	2.552	2.878	3.610
19	0.861	1.066	1.328	1.728	2.093	2.537	2.858	3.580

Paired t-test

Used when measurements are taken from the same subject before and after some manipulation or treatment.

e.g. To determine the significance of a difference in blood pressure before and after administration of an experimental pressure substance.

Assumptions:

1. Populations are distributed normally
2. Samples are drawn independently and at Random

Condition

1. Samples are related with each other
2. Sizes of the sample are small and equal

PROBLEM 26

5. Compare t calculated and theoretical values
6. Conclusion

Blood pressure of 8 patients before and after treatment

BP before	BP after
180	140
200	145
230	150
240	155
170	120
190	130
200	140
165	130

Parametric Tests-t test
(One Sample two occasion)

1. State the Hypotheses

- ✓ **Null Hypothesis (H0):** The mean difference in blood pressure before and after treatment is zero, implying no effect of the treatment. $H_0: \mu_d = 0$
- ✓ **Alternative Hypothesis (H1):** The mean difference in blood pressure before and after treatment is not zero, implying an effect of the treatment. $H_1: \mu_d \neq 0$

2. Calculate the Differences

BP before	BP after	Dif (x)
180	140	40
200	145	55
230	150	80
240	155	85
170	120	50
190	130	60
200	140	60
165	130	35
		$\sum x = 465$

3. Calculate the Mean and Standard Deviation of the Differences

Mean of difference

$$\bar{d} = \frac{\sum d_i}{n} = \frac{40 + 55 + 80 + 85 + 50 + 60 + 60 + 35}{8} = \frac{465}{8} = 58.125$$

Standard deviation of the differences

$$(40 - 58.125)^2 = 331.64$$

$$(55 - 58.125)^2 = 9.78$$

$$(80 - 58.125)^2 = 480.64$$

$$(85 - 58.125)^2 = 727.64$$

$$(50 - 58.125)^2 = 65.14$$

$$(60 - 58.125)^2 = 3.52$$

$$(60 - 58.125)^2 = 3.52$$

$$(35 - 58.125)^2 = 532.14$$

$$331.64 + 9.78 + 480.64 + 727.64 + 65.14 + 3.52 + 3.52 + 532.14 = 2152.44$$

$$s_d = \sqrt{\frac{\sum (d_i - \bar{d})^2}{n-1}} = \sqrt{\frac{2152.44}{7}} \approx \sqrt{307.49} \approx 17.54$$

4. Calculate the t-Statistic

$$t = \frac{\bar{d}}{s_d/\sqrt{n}} = \frac{58.125}{17.54/\sqrt{8}} = \frac{58.125}{6.20} \approx 9.37$$

Degrees of freedom (df)	Significance level (α)								
	.10	.05	.025	.01	.005	.001	.0001	.00001	.000001
2	3.078	4.207	9.279	12.708	24.452	63.657	377.731	630.219	
3	1.886	2.302	3.030	4.813	8.255	18.559	34.009	51.199	
4	1.645	2.132	2.776	4.171	7.478	15.892	26.129	37.434	
5	1.532	1.979	2.322	2.776	3.848	6.864	9.890	14.689	
6	1.476	1.886	2.050	2.571	3.650	4.830	6.775	9.600	
7	1.423	1.812	2.015	2.447	3.499	4.604	6.531	8.439	
8	1.383	1.753	1.984	2.358	3.355	4.461	6.351	8.185	
9	1.351	1.703	1.937	2.282	3.250	4.358	6.173	7.927	
10	1.325	1.660	1.892	2.228	3.169	4.236	5.991	7.657	
11	1.303	1.626	1.854	2.178	3.078	4.127	5.813	7.387	
12	1.287	1.593	1.821	2.135	3.000	4.015	5.640	7.120	
13	1.274	1.563	1.791	2.093	2.929	3.930	5.470	6.857	
14	1.263	1.536	1.765	2.058	2.863	3.849	5.303	6.690	
15	1.254	1.512	1.743	2.028	2.807	3.768	5.140	6.527	
16	1.247	1.491	1.724	2.000	2.754	3.690	5.075	6.367	
17	1.241	1.472	1.707	1.975	2.707	3.618	4.912	6.213	
18	1.236	1.455	1.691	1.952	2.665	3.547	4.846	6.059	
19	1.232	1.439	1.677	1.931	2.628	3.478	4.777	5.906	
20	1.229	1.424	1.664	1.911	2.592	3.412	4.710	5.757	
21	1.226	1.411	1.652	1.892	2.558	3.347	4.645	5.610	
22	1.224	1.4	1.641	1.874	2.525	3.284	4.580	5.465	
23	1.222	1.389	1.630	1.857	2.494	3.223	4.515	5.322	
24	1.220	1.378	1.620	1.841	2.464	3.163	4.450	5.182	
25	1.219	1.368	1.610	1.826	2.436	3.105	4.386	5.045	
26	1.218	1.358	1.601	1.811	2.409	3.049	4.323	4.909	
27	1.217	1.349	1.592	1.797	2.384	3.094	4.261	4.775	
28	1.216	1.340	1.584	1.784	2.360	3.039	4.199	4.643	
29	1.215	1.331	1.576	1.772	2.338	3.085	4.138	4.513	
30	1.215	1.323	1.568	1.760	2.317	3.032	4.077	4.384	

5. Compare t calculated and theoretical values

Tabulated (Dof = 7), with level of significance 0.05, two tails = 2.36

6. Conclusion

't' calculated value is higher than table values hence reject the null hypothesis. This indicates that there is a significant difference in blood pressure before and after treatment, suggesting that the treatment had a substantial effect.

Hypothesis Testing for ρ (Repeated Measures t-test)

PROBLEM 27

Two researcher is conducting a study to evaluate the frequency of sheep collisions during a driving test. Over the course of a day driving test, the number of sheep hit by vehicles is recorded. The data collected from multiple driving tests is as follows:

Subject	Score on test A	Score on test B
1	28	25
2	26	27
3	33	28
4	30	31
5	32	29
6	30	30
7	31	32
8	18	21
9	22	25
10	24	20

Parametric Tests-t test (One Sample two occasion)

1. State the Hypotheses

Null Hypothesis (H_0): There is no significant difference in the mean number of sheep collisions between Test A and Test B. $H_0: \mu_A = \mu_B$

Alternative Hypothesis (H_1): There is a significant difference in the mean number of sheep collisions between Test A and Test B. $H_1: \mu_A \neq \mu_B$

2. Calculate the Differences

Subject	Score on test A	Score on test B	Difference D
1	28	25	3
2	26	27	-1
3	33	28	5
4	30	31	-1
5	32	29	3
6	30	30	0
7	31	32	-1
8	18	21	-3
9	22	25	-3
10	24	20	4
			$\sum D = 6$

3. Calculate the Mean and Standard Deviation of the Differences

Mean of difference

$$\bar{d} = \frac{\sum d_i}{n} = \frac{3 + (-1) + 5 + (-1) + 3 + 0 + (-1) + (-3) + (-3) + 4}{10} = \frac{7}{10} = 0.7$$

Standard deviation of the differences

$$\begin{aligned}
 (3 - 0.7)^2 &= 5.29 & 5.29 + 2.89 + 18.49 + 2.89 + 5.29 + 0.49 + 2.89 + 13.69 + 13.69 + 11.29 &= 75.75 \\
 (-1 - 0.7)^2 &= 2.89 \\
 (5 - 0.7)^2 &= 18.49 \\
 (-1 - 0.7)^2 &= 2.89 \\
 (3 - 0.7)^2 &= 5.29 \\
 (0 - 0.7)^2 &= 0.49 \\
 (-1 - 0.7)^2 &= 2.89 \\
 (-3 - 0.7)^2 &= 13.69 \\
 (-3 - 0.7)^2 &= 13.69 \\
 (4 - 0.7)^2 &= 11.29
 \end{aligned}$$

$$s_d = \sqrt{\frac{75.75}{9}} \approx \sqrt{8.42} \approx 2.90$$

4. Calculate the t-Statistic

$$t = \frac{\bar{d}}{s_d/\sqrt{n}} = \frac{0.7}{2.90/\sqrt{10}} = \frac{0.7}{0.92} \approx 0.76$$

Degrees of freedom (df)	Significance level (α)							
	.2	.15	.1	.05	.025	.01	.005	.001
1	3.078	4.168	6.314	12.708	25.482	63.657	177.727	836.619
2	1.888	2.082	2.924	4.293	4.205	9.025	14.088	31.999
3	1.680	1.802	2.351	3.182	3.078	5.842	7.483	12.824
4	1.533	1.645	2.133	2.776	2.499	4.604	5.896	10.820
5	1.476	1.699	2.016	2.571	3.393	4.032	4.773	6.889
6	1.440	1.650	1.943	2.447	3.069	3.707	4.317	5.558
7	1.415	1.617	1.895	2.348	3.441	3.498	4.029	5.408
8	1.397	1.587	1.865	2.295	3.251	3.552	3.921	4.741
9	1.383	1.574	1.833	2.240	3.095	3.250	3.690	4.791
10	1.372	1.559	1.812	2.238	3.034	3.168	3.581	4.587
11	1.363	1.548	1.798	2.201	2.993	3.106	3.487	4.437
12	1.357	1.539	1.785	2.179	2.950	3.066	3.409	4.318
13	1.350	1.530	1.770	2.158	2.930	3.032	3.372	4.227
14	1.345	1.523	1.761	2.145	2.910	2.977	3.328	4.160

5. Compare t calculated and theoretical values

Tabulated (Dof = 9), with level of significance 0.05, two tails = 2.262

6. Conclusion

't' calculated value is Lower than table values hence accept the null hypothesis.

In this paired t-test, you would conclude that there is no significant difference in the number of sheep collisions between Test A and Test B

Two sample t-test

- ✓ Used when the two independent random samples come from the normal populations having unknown or same variance
- ✓ We test the null hypothesis that the two population means are same i.e $\mu_1 = \mu_2$

Assumptions

1. Population are distributed normally
2. Samples are drawn independently and at random

Conditions

1. Standard deviation in the populations are same and not known
2. Size of the sample is small

Two sample t-test-Procedure

1. State the Hypotheses

Null Hypothesis (H0): There is no significant difference between the means of the two populations. $H_0: \mu_1 = \mu_2$

Alternative Hypothesis (H1): There is a significant difference between the means of the two populations. $H_1: \mu_1 \neq \mu_2$

2. Collect and Summarize Data

Sample Means: \bar{x}_1 for group 1 and \bar{x}_2 for group 2

Sample variance s_1^2 for group 1 and s_2^2 for group 2

Sample size n_1 for group 1 and n_2 for group 2

3. Calculate the Test Statistic

Equal Variances Assumed (Calculate the pooled variance)

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

Calculate the test statistic:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

4. Calculate Degrees of Freedom **Equal Variances Assumed**

$$df = n_1 + n_2 - 2$$

5. Find t value in table

6. Conclusions

PROBLEM 28

Two sample t-test

The following data represent weight in kg for 10 males and 12 females

Males 80, 75, 95, 55, 60, 70, 75, 72, 80, 65

Females 60, 70, 50, 85, 45, 60, 80, 65, 70, 62, 77, 82 Show that hypothesis

Step 1: Hypotheses:

Null Hypothesis (H_0): There is no significant difference between the mean weight of males and females. $H_0: \mu_{\text{males}} = \mu_{\text{females}}$

Alternative Hypothesis (H_1): There is a significant difference between the mean weight of males and females. $H_1: \mu_{\text{males}} \neq \mu_{\text{females}}$

Step 3: Calculate Sample Statistics:

Males:

$$\bar{x}_1 = \frac{80 + 75 + 95 + 55 + 60 + 70 + 75 + 72 + 80 + 65}{10} = \frac{630}{10} = 63 \quad \text{Deviations} = [(80 - 63)^2, (75 - 63)^2, (95 - 63)^2, (55 - 63)^2, (60 - 63)^2, (70 - 63)^2,$$

Deviations = [289, 144, 1024, 64, 9, 49, 144, 81, 289, 4]

$$s_1^2 = \frac{289 + 144 + 1024 + 64 + 9 + 49 + 144 + 81 + 289 + 4}{10 - 1} = \frac{2097}{9} \approx 233$$

Sample Size (n_1) = 10

For Females

$$\bar{x}_2 = \frac{60 + 70 + 50 + 85 + 45 + 60 + 80 + 65 + 70 + 62 + 77 + 82}{12} = \frac{734}{12} \approx 61.17$$

$$\text{Deviations} = [(60 - 61.17)^2, (70 - 61.17)^2, (50 - 61.17)^2, (85 - 61.17)^2, (45 - 61.17)^2, (60 - 61.17)^2, (80 - 61.17)^2, (65 - 61.17)^2, (70 - 61.17)^2, (62 - 61.17)^2, (77 - 61.17)^2, (82 - 61.17)^2]$$

$$s_2^2 = \frac{1.37 + 77.87 + 124.63 + 573.39 + 260.59 + 1.37 + 353.79 + 14.82 + 77.87 + (45 - 61.17)^2 + (60 - 61.17)^2 + (80 - 61.17)^2 + (65 - 61.17)^2 + (70 - 61.17)^2 + (62 - 61.17)^2 + (77 - 61.17)^2 + (82 - 61.17)^2}{12 - 1}$$

$$S_2^2 = 192.23$$

Sample size = 12

4. Calculate t values

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} = \frac{(10 - 1) \cdot 233 + (12 - 1) \cdot 192.23}{10 + 12 - 2} = 210.56$$
$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} = \frac{63 - 61.17}{\sqrt{210.56 \left(\frac{1}{10} + \frac{1}{12} \right)}} = 0.29$$

4. Calculate 5.Degrees of Freedom $df = n_1 + n_2 - 2 = 20$

5. Find the p-value

Use a t-distribution table or calculator to find the p-value for $t=0.29$ $df=20$.

For a two-tailed test, you compare the p-value to your significance level (α typically 0.05)

6. Make a Decision

If the p-value is large (e.g., greater than 0.05), fail to reject the null hypothesis.
There is no significant difference in average weights between males and females.

Perform the above calculations to finalize the results. The t-test result indicates whether there is a significant difference in average weights between the two groups.

Z-test

- ✓ A Z-test is a statistical test used to determine whether there is a significant difference between sample and population means or between the means of two independent samples. It is based on the Z-distribution, which is a special case of the normal distribution with a mean of 0 and a standard deviation of 1. The Z-test is appropriate when the sample size is large (typically $n>30$) or when the population variance is known

Assumption

- Population is normally distributed
- The sample is drawn at random

Condition

- ✓ Population standard deviation σ is known
- ✓ Size of the sample is large ($n>30$)

One-Sample Z-Test

$$Z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$$

Two-Sample Z-Test

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

Z-test-steps

1. State the Hypotheses

Null Hypothesis (H_0): There is no effect or difference.

Alternative Hypothesis (H_1): There is a significant effect or difference.

2. Calculate the Z-Statistic:

3. Determine the p-value:

4. Make a Decision

5. Report the Results

Hypothesis Testing Z- test

PROBLEM 29

1. One of Real Estate advertises that the mean selling time of a residential home is 40 days or less after it is listed in their company. A sample of 50 recently sold homes shows a sample mean selling time of 45 days and a standard deviation of 20 days. Using a 0.02 level of significance, test the validity of the company's claim.

Given: $n = 50$, mean = 45 days, standard deviation = 20 days, significant level = 0.02

Step 1: Hypotheses

- a) **Null Hypothesis (H_0):** The mean selling time is 40 days or less. $H_0: \mu \leq 40$
- b) **Alternative Hypothesis (H_1):** The mean selling time is greater than 40 days.

$$H_1: \mu > 40$$

This is a one-tailed test with $\alpha=0.02$ then critical value is $Z_{0.02} = 2.05$

Step 2: Calculate the Test Statistic

$$Z = \frac{\bar{x} - \mu_0}{S_{\bar{x}}} = \frac{45 - 40}{\frac{20}{\sqrt{50}}} = 1.77$$

d) Since $z= 1.77$ is less than 2.05, we do not reject the null hypothesis. A sample selling time 45 days is not significantly greater than population mean of 40 days at 0.02 level of confidence.

The sample does not provide sufficient evidence to suggest that the mean selling time is greater than 40 days.

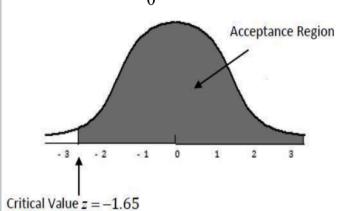
PROBLEM 30

2. The height of adults in a certain town is found to have a mean of 166.17 cm with standard deviation of 5.89 cm. A random sample of 144 adults in the slum section of the town is discovered to have a mean height of 164.65 cm. Does this height indicate that the residents of the slum area are significantly shorter in height at 0.05 level of significance?

- a) Set the null and alternative hypothesis

$$H_0: \mu \geq 166.17 \quad H_a: \mu < 166.17 \text{ cm}$$

- b) Population standard deviation σ is 5.89 and sample size is large at $n = 144$ use the Z-static. This is also a one tailed test with $\alpha=0.05$. the critical point is $Z_{0.05} = 1.65$ We shall reject H_0 if $Z < -1.65$



- c) Computation

$$Z = \frac{\bar{x} - \mu_0}{S_{\bar{x}}} = \frac{164.65 - 166.17}{\frac{5.89}{\sqrt{144}}} = -3.93$$

- d) Since $z= -3.93$ is less than -1.65, we reject the null hypothesis. This means the height of adults in the slum section of the town significantly lower than 166.7 cm at 0.05 level of significance.

Pearson's 'r' Correlation

Pearson's correlation coefficient, often denoted by r, measures the strength and direction of the linear relationship between two continuous variables

1. Formula

The formula for Pearson's correlation coefficient r is:

$$r = \frac{n \sum(XY) - (\sum X)(\sum Y)}{\sqrt{[n \sum X^2 - (\sum X)^2][n \sum Y^2 - (\sum Y)^2]}}$$

2. Interpretation

Type of correlation	Correlation coefficient
Perfect positive correlation	$r = +1$
Partial positive correlation	$0 < r < +1$
No correlation	$r = 0$
Partial negative correlation	$0 > r > -1$
Perfect negative correlation	$r = -1$

3. Significance Testing

To determine if the correlation is statistically significant, you can use a t-test:

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$

The t-value can then be compared to the critical value from the t-distribution table with $n-2$ degrees of freedom to determine significance.

PROBLEM 31

Pearson's 'r' Correlation-Problem

We have collected data on the length of hands (in centimeters) and height (in centimeters) for 5 individuals. The data is as follows:

Calculate and Interpret the Correlation coefficient of the two variable		
Person	Hand	Height
A	17	150
B	15	154
C	19	169
D	17	172
E	21	175
Total	89	820

	x	y	xy	x^2	y^2
A	17	150	2550	289	22500
B	15	154	2310	225	23716
C	19	169	3211	361	28561
D	17	172	2924	289	29584
E	21	175	3675	441	30625
Total	89	820	14670	1605	134986

$$r = \frac{n \sum xy - (\sum x)(\sum y)}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}. \quad r = 0.72$$

The formula gives us a correlation coefficient of **0.72**, (Critical value is 832) which is a high, positive correlation. Meaning that in this data set, as height increases, so does hand height.

PROBLEM 32

Pearson's 'r' Correlation-Problem

Calculate and Interpret the Correlation coefficient of the two variable		
Person	Weight	Blood Pressure
A	150	125
B	169	130
C	175	160
D	180	169
E	200	150

	x	y	xy	x ²	y ²
A	150	125	18750	22500	15625
B	169	130	21970	28561	16900
C	175	160	28000	30625	25600
D	180	169	30420	32400	28561
E	200	150	30000	40000	22500
Total	874	734	12914	15408	109186

$$r = \frac{n \sum xy - (\sum x)(\sum y)}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}. \quad r = 0.61$$

Weight and blood pressure have a moderate, positive correlation.

PROBLEM 33

Pearson's 'r' Correlation-Problem

1. Suppose you computed $r = 0.801$ using $n = 10$ data. Explain significant of the data

From the Table for $n = 10$ the critical values ± 0.685 then

r is not significant between $+ 0.685$ and $- 0.685$ but it is significant > 0.685

2. Suppose you computed $r = -0.624$ using $n = 14$ data. Explain significant of the data

From the Table for $n = 14$ the critical values ± 0.592 then

r is not significant between $+ 0.592$ and $- 0.592$ but -0.624 it is significant

3. Suppose you computed $r = 0.776$ using $n = 6$ data. Explain significant of the data

From the Table for $n = 6$ the critical values ± 0.832 then

r is not significant between $+ 0.832$ and $- 0.832$ but it is significant > 0.811

Non-Parametric Tests

Chi-Square test

1. **Chi-square test (χ^2)** The Chi-Square test is used to assess whether observed frequencies in categorical data differ significantly from expected frequencies. It is commonly applied in two main types:

Chi-Square Test of Independence: Evaluates if two categorical variables are independent or associated.

Chi-Square Test of Goodness of Fit: Determines if a sample data matches the expected distribution.

Chi-Square Test of Independence: Procedure

1. Set Up the Hypotheses:

Null Hypothesis (H_0): The variables are independent (no association).

Alternative Hypothesis (H_1): The variables are dependent (there is an association).

2. Create a Contingency Table: Tabulate the observed frequencies for each combination of the variables.

3. Calculate the Expected Frequencies: $E_{ij} = \frac{(\text{row total} \times \text{column total})}{\text{grand total}}$

4. Compute the Chi-Square Statistic: (O_{ij} observe frequency)

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

5. Determine the Degrees of Freedom (df) = (number of rows-1) (number of column-1)

6. Compare to the Critical Value:

If $\chi^2_{\text{calculated}} > \chi^2_{\text{critical}}$, reject the null hypothesis.

PROBLEM 34

Chi-Square test

Suppose we want to test if there is an association between gender and voting preference in an election. We collect the following

	Vote for A	Vote for B	Total
Male	30	20	50
Female	25	25	50
Total	55	45	100

Step 1: Hypothesis

H0 :Gender and Preference are independent

H1 :Gender and Preference are dependent

Create a Contingency

$$E_{\text{Male, A}} = \frac{50 \times 55}{100} = 27.5$$

$$E_{\text{Male, B}} = \frac{50 \times 45}{100} = 22.5$$

$$E_{\text{Female, A}} = \frac{50 \times 55}{100} = 27.5$$

$$h \quad E_{\text{Female, B}} = \frac{50 \times 45}{100} = 22.5$$

Table Chi-Square Calculation

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

$$\chi^2 = \frac{(30 - 27.5)^2}{27.5} + \frac{(20 - 22.5)^2}{22.5} + \frac{(25 - 27.5)^2}{27.5} + \frac{(25 - 22.5)^2}{22.5}$$

$$= 1.01$$

Degree of freedom = $(2-1)(2-1) = 1$

Critical value $\alpha = 0.05$ and df = 1 is 3.841

Decision: Since $\chi^2 = 1.01 < 3.841$ we fail to reject the null

suggesting no significant association between gender and voting preference.

PROBLEM 35

Chi-Square test

Example 1 A die is thrown with following results. Is the die unbiased?

Number of turned up	1	2	3	4	5	6
Frequency	16	20	25	14	29	28

Solution: Let us take the hypothesis that the die is unbiased. If that is so, the probability of obtaining any one of the six number is $1/6$ and as such the expected frequency of any one number coming upward is $E_f = 132 * 1/6 = 22$. Now we can write the observed frequencies along with expected frequencies and work out the values of λ^2 as follows

PROBLEM 36

Chi-Square test

No. of turned up	Observed Frequency, O_f	Expected Frequency, E_f	$(O_f - E_f)$	$(O_f - E_f)^2$	$(O_f - E_f)^2/E_f$
1	16	22	-6	36	36/22
2	20	22	-2	4	4/22
3	25	22	3	9	9/22
4	14	22	-8	64	64/22
5	29	22	7	49	49/22
6	28	22	6	36	36/22
Σ	132				9

$$\sum [(O_f - E_f)^2 / E] = 9$$

Hence the calculated value of $\lambda^2 = 9$ then Degree of freedom in the given problem is $(n-1) = (6-1) = 5$

The Table value of λ^2 for 5 DoF at 5% level of significance is **11.071**. Comparing computed values and table values of λ^2 , we find that computed values is less than the table values and as such could have arisen due to fluctuations of sampling. The result, thus, support the hypothesis and it can be concluded that the die is unbiased.

Chi-Square test

2. One sample sign test: It is based on the direction or the plus or minus signs of observation in a sample and not on their numerical magnitudes

PROBLEM 37

Problem: The compressive strength of insulating blocks used in the construction of new houses is tested by a civil engineer. The engineer needs to be certain at the 5% level of significance that the median compressive strength is at least 1000 psi. Twenty randomly selected blocks give the following results:

Observation	Compressive Strength						
1	1128.7	6	718.4	11	1167.1	16	1153.6
2	679.1	7	787.4	12	1387.5	17	1423.3
3	1317.2	8	1562.3	13	679.9	18	1122.6
4	1001.3	9	1356.9	14	1323.2	19	1644.3
5	1107.6	10	1153.2	15	788.4	20	737.4

Test (at the 5% level of significance) the null hypothesis that the median compressive strength of the insulting blocks is 1000 psi against the alternative that it is greater

Solution The hypotheses are $H_0 : \theta = 1000$ and $H_1 : \theta > 1000$

Comp. Strength	Sign						
1128.7	+	718.4	-	1167.1	+	1153.6	+
679.1	-	787.4	-	1387.5	+	1423.3	+
1317.2	+	1562.3	+	679.9	-	1122.6	+
1001.3	+	1356.9	+	1323.2	+	1644.3	+
1107.6	+	1153.2	+	788.4	-	737.4	-

We have 14 plus signs and the required probability value is calculated directly from the binomial formula as (for binomial distribution $p = 1/2$)

$$\begin{aligned}
 P(X = r) &= \binom{n}{r} q^{n-r} p^r = \binom{n}{r} (1-p)^{n-r} p^r \\
 P(X \geq 14) &= \sum_{r=14}^{20} \binom{20}{r} \left(\frac{1}{2}\right)^{20-r} \left(\frac{1}{2}\right)^r \\
 &= \frac{20.19.18.17.16.15}{1.2.3.4.5.6} \left(\frac{1}{2}\right)^{20} + \frac{20.19.18.17.16}{1.2.3.4.5} \left(\frac{1}{2}\right)^{20} + \frac{20.19.18.17}{1.2.3.4} \left(\frac{1}{2}\right)^{20} + \frac{20.19.18}{1.2.3} \left(\frac{1}{2}\right)^{20} + \frac{20.19}{1.2} \left(\frac{1}{2}\right)^{20} + \frac{20}{1} \left(\frac{1}{2}\right)^{20} + \left(\frac{1}{20}\right)^{20} \\
 &= \left(\frac{1}{2}\right)^{20} (38760 + 15504 + 4845 + 1140 + 190 + 20 + 1) \\
 &= 0.05766
 \end{aligned}$$

Since we are performing a one-tailed test, we must compare the calculated value with the value 0.05. Since $0.05 < 0.05766$ we conclude that we cannot reject the null hypothesis and that on the basis of the available evidence, we cannot conclude that the median compressive strength of the insulating blocks is greater than 1000 psi.

PROBLEM 38

Chi-Square test

2. A certain type of solid rocket fuel is manufactured by bonding an igniter with a propellant. In order that the fuel burns smoothly and does not suffer either "flame-out" or become unstable it is essential that the material bonding the two components of the fuel has a shear strength of 2000 psi. The results arising from tests performed on 20 randomly selected samples of fuel are as follows:

Observation	Shear Strength						
1	2128.7	6	1718.4	11	2167.1	16	2153.6
2	1679.1	7	1787.4	12	2387.5	17	2423.3
3	2317.2	8	2562.3	13	1679.9	18	2122.6
4	2001.3	9	2356.9	14	2323.2	19	2644.3
5	2107.6	10	2153.2	15	1788.4	20	1737.4

Using the 5% level of significance, test the null hypothesis that the median shear strength is 2000 psi.

The hypotheses are $H_0 : \theta = 2000$ $H_1 : \theta \neq 2000$. We determine the signs associated with each observation as shown below and perform a two-tailed test.

Shear Strength	Sign						
2128.7	+	1718.4	-	2167.1	+	2153.6	+
1679.1	-	1787.4	-	2387.5	+	2423.3	+
2317.2	+	2562.3	+	1679.9	-	2122.6	+
2001.3	+	2356.9	+	2323.2	+	2644.3	+
2107.6	+	2153.2	+	1788.4	-	1737.4	-

We have 14 plus signs and the required probability value is calculated directly from the binomial formula:

$$\begin{aligned}
 P(X \geq 14) &= \sum_{r=14}^{20} \binom{20}{r} \left(\frac{1}{2}\right)^{20-r} \left(\frac{1}{2}\right)^r \\
 &= \frac{20.19.18.17.16.15}{1.2.3.4.5.6} \left(\frac{1}{2}\right)^{20} + \frac{20.19.18.17.16}{1.2.3.4.5} \left(\frac{1}{2}\right)^{20} + \frac{20.19.18.17}{1.2.3.4} \left(\frac{1}{2}\right)^{20} + \frac{20.19.18}{1.2.3} \left(\frac{1}{2}\right)^{20} + \frac{20.19}{1.2} \left(\frac{1}{2}\right)^{20} + \frac{20}{1} \left(\frac{1}{2}\right)^{20} + \left(\frac{1}{2}\right)^{20} \\
 &= \left(\frac{1}{2}\right)^{20} (38760 + 15504 + 4845 + 1140 + 190 + 20 + 1) = 0.05766
 \end{aligned}$$

Since we are performing a two-tailed test, we must compare the calculated value with 0.025. Since $0.025 < 0.05766$ we cannot reject the null hypothesis on the basis of the evidence and conclude that the median shear strength is not significantly different from 2000 psi.

PROBLEM 39

Chi-Square test

3. A certain type of solid rocket fuel is manufactured by binding an igniter with a propellant. In order that the fuel burns smoothly and does not suffer either "flame-out" or become unstable it is essential that the material bonding the two components of the fuel has a shear strength of 2000 psi. The results arising from tests performed on 10 randomly selected samples of fuel are as follows

Observation	Shear Strength	Observation	Shear Strength
1	2128.7	6	1718.4
2	1679.1	7	1787.4
3	2317.2	8	2562.3
4	2001.3	9	2356.9
5	2107.6	10	2153.2

Using the 5% level of significance, test the null hypothesis that the median shear strength is 2000 psi.

Answer The hypotheses are $H_0 : \theta = 2000$ $H_1 : \theta \neq 2000$

Shear Strength	Sign	Shear Strength	Sign
2128.7	+	1718.4	-
1679.1	-	1787.4	-
2317.2	+	2562.3	+
2001.3	+	2356.9	+
2107.6	+	2153.2	+

We have 7 plus signs and the required probability value is calculated directly from the binomial formula as

$$\begin{aligned} P(X \geq 7) &= \sum_{r=7}^{10} \binom{10}{r} \left(\frac{1}{2}\right)^{10-r} \left(\frac{1}{2}\right)^r = \frac{10.9.8}{1.2.3} \left(\frac{1}{2}\right)^{10} + \frac{10.9}{1.2} \left(\frac{1}{2}\right)^{10} + \frac{10}{1} \left(\frac{1}{2}\right)^{10} + \left(\frac{1}{2}\right)^{10} \\ &= \frac{10.9.8}{1.2.3} \left(\frac{1}{2}\right)^{10} + \frac{10.9}{1.2} \left(\frac{1}{2}\right)^{10} + \frac{10}{1} \left(\frac{1}{2}\right)^{10} + \left(\frac{1}{2}\right)^{10} = \left(\frac{1}{2}\right)^{10} (120 + 45 + 10 + 1) \simeq 0.172 \end{aligned}$$

Since we are performing a two-tailed test, we must compare the calculate value with the value 0.025. Since $0.025 < 0.172$ we cannot reject the null hypothesis on the basis of the available evidence and we cannot conclude that the median shear strength is different to 2000 psi.

Chi-Square test

The sign test for paired data

PROBLEM 40

The sign test for paired data

Automotive development engineers are testing the properties of two anti-lock braking systems in order to determine whether they exhibit any significant difference in the stopping distance achieved by different cars. The systems are fitted to 10 cars and a test is run ensuring that each system is used on each car under conditions which are as uniform as possible. The stopping distances (in yards) obtained are given in the table below

Car	Anti-lock Braking System	
	1	2
1	27.7	26.3
2	32.1	31.0
3	29.6	28.1
4	29.2	28.1
5	27.8	27.9
6	26.9	25.8
7	29.7	28.2
8	28.9	27.6
9	27.3	26.5
10	29.9	28.3

Use a sign test to test the null hypothesis that the mean difference between the hardnesses produced by the two methods is zero against the alternative that it is not zero. Use the 1% level of significance.

Solution We are testing to find any differences in the median stopping distance figures for each braking system. The null and alternative hypotheses are:

$$H_0 : \theta_1 = \theta_2 \text{ or } H_0 : \theta \text{ differences} = 0$$

$$H_1 : \theta_1 \neq \theta_2 \text{ or } H_1 : \theta \text{ differences} \neq 0$$

We perform a two-tailed test. The signed differences shown by the two systems are shown in the table below:

Car	Anti-lock Braking System		Sign
	1	2	
1	27.7	26.3	+
2	32.1	31.0	+
3	29.6	28.1	+
4	29.2	28.1	+
5	27.8	27.9	-
6	26.9	25.8	+
7	29.7	28.2	+
8	28.9	27.6	+
9	27.3	26.5	+
10	29.9	28.3	+

We have 9 plus signs and the required probability value is calculated directly from the binomial formula as

$$P(X \geq 9) = \sum_{r=9}^{10} \binom{10}{r} \left(\frac{1}{2}\right)^{10-r} \left(\frac{1}{2}\right)^r = \frac{10}{1} \left(\frac{1}{2}\right)^{10} + \left(\frac{1}{2}\right)^{10} = 11 \times \left(\frac{1}{2}\right)^{10} \simeq 0.011$$

Since we are performing a two-tailed test, we must compare the calculated value with the value 0.025. Since $0.011 < 0.025$ we reject the null hypothesis on the basis of the available evidence and conclude that the differences in the median stopping distances recorded is significant at the 5% level

Mann-Whitney test

3. Mann-Whitney Test (U test): It is a nonparametric counterpart of the t test used to compare the means of two independent populations.

The following assumptions underlie the use of the Mann-Whitney test

1. The samples are independent
2. The level of data is at least ordinal

The two-tailed hypothesis being tested

H_0 : The two populations are identical

H_a : The two populations are not identical

The problem can be solved in two categories

Small sample case: When both n_1 and $n_2 \leq 10$

Larger sample case: When both n_1 and $n_2 > 10$

Small sample case: When both n_1 and $n_2 \leq 10$

$$U_1 = n_1 n_2 + \frac{n_1(n_1 + 1)}{2} - W \quad U_2 = n_1 n_2 + \frac{n_2(n_2 + 1)}{2} - W$$

PROBLEM 41

Mann-Whitney test

Is there a difference between health service workers and police service workers in the amount of compensation employers pay them per hour during Corona effect time. (randomly selected).

Health Service Worker (Rupees)	Police Service workers (Rupees)
20.10	26.19
19.80	23.88
22.36	25.50
18.75	2164
21.90	24.85
22.96	25.30
20.75	24.12
	23.45

Hypothesis

H_0 : The health service population is identical to the police service population on employee compensation

H_a : The health service population is not identical to the police service population on employee compensation

Total Employee Compensation	Rank	Group
18.75	1	H
19.80	2	H
20.10	3	H
20.75	4	H
21.64	5	P
21.90	6	H
22.36	7	H
22.96	8	H
23.45	9	P
23.88	10	P
24.12	11	P
24.85	12	P
25.30	13	P
25.50	14	P
26.19	15	P

$$W_1 = 1+2+3+4+6+7+8 = 31$$

$$W_2 = 5+9+10+11+12+13+14+15 = 89$$

$$U_1 = n_1 n_2 + \frac{n_1(n_1 + 1)}{2} - W_1 \quad U_1 = 7 * 8 + \frac{7(7+1)}{2} - 31 = 53$$

$$U_2 = n_1 n_2 + \frac{n_2(n_2 + 1)}{2} - W_2 \quad U_2 = 7 * 8 + \frac{8(8+1)}{2} - 89 = 3$$

Because of U_2 is the smallest values of U , we use $U = 3$ as the test statistic. Because it is the smallest size let $n_1 = 7$ and $n_2 = 8$

From table P values of 0.0011 x2 = 0.0022 (because two tailed)

The statistical conclusion is that the populations are not identical

PROBLEM 42

Larger sample case: When both n_1 and $n_2 > 10$

$$U = n_1 n_2 + \frac{n_1(n_1+1)}{2} - W_1 \quad \mu_U = \frac{n_1 n_2}{2} \quad \sigma_U = \sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}$$

Do construction workers who purchase lunch from street vendors spend less per meal than construction workers who go to restaurants for lunch?

Wenders ($n_1 = 14$)	2.75	3.29	4.63	3.61	3.10	4.29	2.25	2.97	4.01	3.68	3.15	2.97	4.05	3.60		
Restaurant ($n_2 = 16$)	4.10	4.75	3.95	3.50	4.25	4.98	5.75	4.10	2.70	3.65	5.11	4.80	6.25	3.89	4.80	5.50

Hypothesis

H_0 = The populations of construction worker spending for lunch at vendors and restaurants are the same

H_1 = The populations of construction worker spending at vendors is shifted to the left of the population of

Value	Rank	Group	Value	Rank	Group
2.25	1	V	4.01	16	V
2.70	2	R	4.05	17	V
2.75	3	V	4.10	18.5	R
2.97	4.5	V	4.10	18.5	R
2.97	4.5	v	4.25	20	R
3.10	6	V	4.29	21	V
3.15	7	V	4.53	22	V
3.29	8	V	4.75	23	R
3.50	9	R	4.80	24.5	R
3.60	10	V	4.80	24.5	R
3.61	11	V	4.98	26	R
3.65	12	R	5.11	27	R
3.68	13	V	5.50	28	R
3.89	14	R	5.75	29	R
3.95	15	R	6.25	30	R

$$W_1 = 1+3+4.5+4.5+6+7+8+10+11+13+16+17+21+22 = 144$$

Solving for U, μ_U and σ_U using formula

$$U = n_1 n_2 + \frac{n_1(n_1+1)}{2} - W_1 = 185$$

$$\mu_U = \frac{n_1 n_2}{2} = 112$$

$$\sigma_U = \sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}} = 24.1$$

$$Z = \frac{U - \mu_U}{\sigma_U} = 3.03$$

The p-value associated with Z = 3.03 is 0.0012 < 0.0022. the null hypothesis is rejected

Kruskal Wallies Test

4. Kruskal Wallies Test: This test determines whether all of the groups come from the same or equal populations or whether at least one group comes from a different population.

The process of computing a Kruskal-Wallis K Statistic begins with ranking the data **in** the groups together, as though they were from one group.

$$K = \frac{12}{n(n+1)} \left(\sum_{i=1}^c \frac{T_i^2}{n_i} \right) - 3(n+1)$$

Where

C = number of groups

n = total number of items

T_i = Total of ranks in a group

n_i = number of items in a group

K = X² with Dof = c-1

PROBLEM 43

Kruskal Wallies Test

Problem: Agribusiness researcher are interested in determining the conditions under which Christmas trees grow fastest. A random sample of equivalent size seedling is divided into four groups. Use the Kruskal-Wallis test to determine whether there is a significant difference in the growth of trees in the groups use

Group 1 (Native)	Group 2 (+ water)	Group 3 (+ Fertilizer)	Group 4 (+ water and Fertilizer)
8	1	11	18
5	12	14	20
7	11	10	16
11	9	16	15
9	13	17	14
6	12	12	22

Hypothesis

$$H_0: \text{Group 1} = \text{Group 2} = \text{Group 3} = \text{Group 4}$$

H_a : At least one group is different

$\text{DoF} = C - 1 = 4 - 1 = 3$ and $\alpha = 0.01$ (the critical values of chi-square is $\chi^2_{0.01, 3} = 11.3449$ if $K > 11.3449$ the decision is to reject the null hypothesis)

Assigning Raking																								
5	6	7	8	9	9	10	10	11	11	11	12	12	12	13	14	14	15	15	16	16	17	18	20	22
1	2	3	4	5.5	5.5	7.5	7.5	10	10	10	13	13	13	15	16.5	16.5	18	19.5	19.5	21	22	23	24	
Group 1 (Native)						Group 2 (+ water)						Group 3 (+ Fertilizer)						Group 4 (+ water and Fertilizer)						
4						7.5						10						22						
1						13						16.5						23						
3						10						7.5						19.5						
10						5.5						19.5						18						
5.5						15						21						16.5						
2						13						13						24						
$T_1 = 25.5 (n_1=6)$						$T_2 = 64.0 (n_2=6)$						$T_3 = 87.5 (n_3=6)$						$T_4 = 123.0 (n_4=6)$						

$$\sum_{i=1}^c \frac{T_i^2}{n_i} = \frac{25.5^2}{6} + \frac{64^2}{6} + \frac{87.5^2}{6} + \frac{123^2}{6} = 4588.6$$

$$K = \frac{12}{n(n+1)} \left(\sum_{i=1}^c \frac{T_i^2}{n_i} \right) - 3(n+1)$$

$$K = \frac{12}{24(24+1)} (4588.6) - 3(24+1) = 16.77$$

The observed K value is 16.77 and the critical value $\chi^2_{0.01, 3} = 11.3449$. because the observed value is greater than the table value, the null hypothesis is rejected. There is a significant difference in the way the trees grow.

PROBLEM 44

Linear Regression

Example:

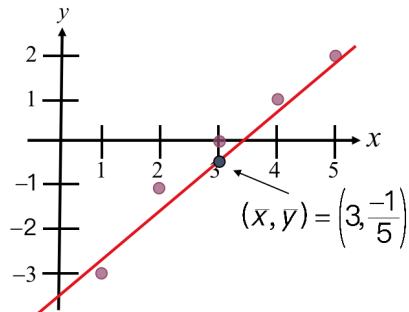
Find the equation of the regression line.

$$b = \bar{y} - mx = \frac{-1}{5} - (1.2) \frac{15}{5} = -3.8$$

x	y	xy	x^2	y^2
1	-3	-3	1	9
2	-1	-2	4	1
3	0	0	9	0
4	1	4	16	1
5	2	10	25	4
$\sum x = 15$	$\sum y = -1$	$\sum xy = 9$	$\sum x^2 = 55$	$\sum y^2 = 15$

$$m = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum x^2 - (\sum x)^2} = \frac{5(9) - (15)(-1)}{5(55) - (15)^2} = \frac{60}{50} = 1.2$$

The equation of the regression line is
 $\hat{y} = 1.2x - 3.8$.



PROBLEM 45

Linear Regression

The following data represents the number of hours 12 different students watched television during the weekend and the scores of each student who took a test the following Monday.

- Find the equation of the regression line.
- Use the equation to find the expected test score for a student who watches 9 hours of TV.

Hours, x	0	1	2	3	3	5	5	5	6	7	7	10
Test score, y	96	85	82	74	95	68	76	84	58	65	75	50
xy	0	85	164	222	285	340	380	420	348	455	525	500
x^2	0	1	4	9	9	25	25	25	36	49	49	100
y^2	9216	7225	6724	5476	9025	4624	5776	7056	3364	4225	5625	2500

$$\sum x = 54$$

$$\sum y = 908$$

$$\sum xy = 3724$$

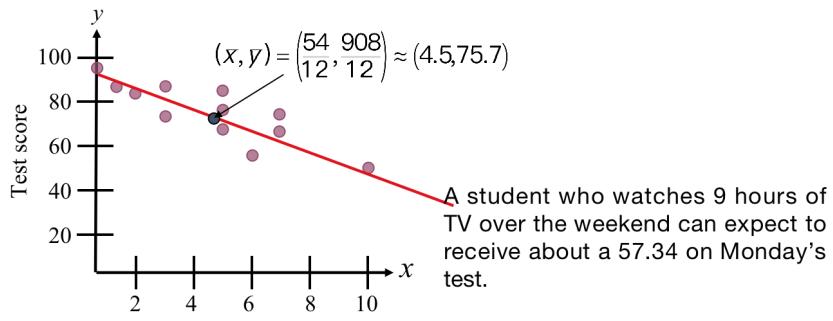
$$\sum x^2 = 332$$

$$\sum y^2 = 70836$$

$$m = \frac{n\sum xy - (\sum x)(\sum y)}{n\sum x^2 - (\sum x)^2} = \frac{12(3724) - (54)(908)}{12(332) - (54)^2} \approx -4.067$$

$$\begin{aligned} b &= \bar{y} - mx \\ &= \frac{908}{12} - (-4.067)\frac{54}{12} \\ &\approx 93.97 \end{aligned}$$

$$\hat{y} = -4.07x + 93.97$$



Using the equation $\hat{y} = -4.07x + 93.97$, we can predict the test score for a student who watches 9 hours of TV.

$$\hat{y} = -4.07(9) + 93.97 = -4.07(9) + 93.97 = 57.34$$

PROBLEM 46

Multiple Regression

Problem: A Distributor of frozen desert pies want to evaluate factors thought to influence demand

Week	Pie Sales	Price (\$)	Advertising (\$100s)
1	350	5.50	3.3
2	460	7.50	3.3
3	350	8.00	3.0
4	430	8.00	4.5
5	350	6.80	3.0
6	380	7.50	4.0
7	430	4.50	3.0
8	470	6.40	3.7
9	450	7.00	3.5
10	490	5.00	4.0
11	340	7.20	3.5
12	300	7.90	3.2
13	440	5.90	4.0
14	450	5.00	3.5
15	300	7.00	2.7

Depend Variable: Pie Sales (units per week)
 Independent Variables: Price
 Advertising
 Data are collected for 15 weeks

Solution

$$\text{Sales} = b_0 + b_1 (\text{price}) + b_2 (\text{Advertising})$$

	y	x ₁	x ₂	x ₁ x ₂	x ₁ ²	x ₂ ²
1	350	5.5	3.3	1925	1155	18.15
2	460	7.5	3.3	3450	1518	24.75
3	350	8	3	2800	1050	24
4	430	8	4.5	3440	1935	36
5	350	6.8	3	2380	1050	20.4
6	380	7.5	4	2850	1520	30
7	430	4.5	3	1935	1290	13.5
8	470	6.4	3.7	3008	1739	23.68
9	450	7	3.5	3150	1575	24.5
10	490	5	4	2450	1960	20
11	340	7.2	3.5	2448	1190	25.2
12	300	7.9	3.2	2370	960	25.28
13	440	5.9	4	2596	1760	23.6
14	450	5	3.5	2250	1575	17.5
15	300	7	2.7	2100	810	18.9
	Σ	99.2	52.2	39152	21087	345.4
					6	675.26
						185

$$b_1 = \frac{\sum x_1 y \sum x_2^2 - \sum x_1 x_2 \sum x_2 y}{\sum x_1^2 \sum x_2^2 - (\sum x_1 x_2)^2}$$

$$b_2 = \frac{\sum x_2 y \sum x_1^2 - \sum x_1 x_2 \sum x_1 y}{\sum x_1^2 \sum x_2^2 - (\sum x_1 x_2)^2}$$

$$(\sum x_1 x_2)^2 = 26814169$$

$$b_1 = -24.97509$$

$$b_2 = 74.131$$

$$b_0 = \frac{\sum y}{n} - \frac{\sum x_1}{n} - \frac{\sum x_2}{n}$$

$$B_0 = 306.526$$

$$\text{Sales} = 306.526 - 24.975(\text{Price}) + 74.131(\text{Advertising})$$

	y	x ₁	x ₂	x ₁ ,y	x ₂ ,y	x ₁ ,x ₂	x ₁ ²	x ₂ ²
1	350	5.5	3.3	1925	1155	18.15	30.25	10.89
2	460	7.5	3.3	3450	1518	24.75	56.25	10.89
3	350	8	3	2800	1050	24	64	9
4	430	8	4.5	3440	1935	36	64	20.25
5	350	6.8	3	2380	1050	20.4	46.24	9
6	380	7.5	4	2850	1520	30	56.25	16
7	430	4.5	3	1935	1290	13.5	20.25	9
8	470	6.4	3.7	3008	1739	23.68	40.96	13.69
9	450	7	3.5	3150	1575	24.5	49	12.25
10	490	5	4	2450	1960	20	25	16
11	340	7.2	3.5	2448	1190	25.2	51.84	12.25
12	300	7.9	3.2	2370	960	25.28	62.41	10.24
13	440	5.9	4	2596	1760	23.6	34.81	16
14	450	5	3.5	2250	1575	17.5	25	12.25
15	300	7	2.7	2100	810	18.9	49	7.29
	Σ	99.2	52.2	39152	21087	345.4	675.26	185

$$b_1 = \frac{\sum x_1 y \sum x_2^2 - \sum x_1 x_2 \sum x_2 y}{\sum x_1^2 \sum x_2^2 - (\sum x_1 x_2)^2}$$

$$b_2 = \frac{\sum x_2 y \sum x_1^2 - \sum x_1 x_2 \sum x_1 y}{\sum x_1^2 \sum x_2^2 - (\sum x_1 x_2)^2}$$

$$(\sum x_1 x_2)^2 = 26814169$$

$$b_1 = -24.97509$$

$$b_2 = 74.131$$

$$b_0 = \frac{\sum y}{n} - \frac{\sum x_1}{n} - \frac{\sum x_2}{n}$$

$$B_0 = 306.526$$

$$\text{Sales} = 306.526 - 24.975(\text{Price}) + 74.131(\text{Advertising})$$