

generally used for attribute selection.

Dimensionality Reduction

- The number of input features, variables, or columns present in a given dataset is known as dimensionality, and the process to reduce this feature is called dimensionality reduction.
- It is a way of converting the higher dimension dataset into lower dimension dataset ensuring that it provides similar information.

Problems

Measures of Central Tendency

- Mean
- Median
- Mode
- Mid Range
- Range

$$\rightarrow \text{Mean} := (\bar{x})$$

$$\text{Mean} = \frac{\text{Sum of Observations}}{\text{No of Observations}}$$

→ Median :-

It means to identify the middle value from a set of observations.

i) When n is odd :-

$$\frac{(n+1)}{2} \text{ th position}$$

ii) When n is even -

$$\left(\frac{(n)}{2} + 1 \right) \text{ th}$$

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

→ Mode :-

Mode is the most frequently occurring value in the dataset.

a) Find the mean of the first 10 odd integers

Obs → 1, 3, 5, 7, 9, 11, 13, 15, 17, 19

$$\text{Mean} = \frac{\text{Sum}}{10} = \frac{100}{10} = 10$$

→ Mid Range :-

It is the sum of maximum and minimum ob / 2 in the data.

$$\text{i.e. Mid Range} = \frac{\text{max} + \text{min}}{2}$$

→ Range :-

It is the difference between the maximum and minimum value of observation in the data.

$$\text{i.e. Range} = \frac{\text{max} - \text{min}}{2}$$

→ Variance (σ^2)

$$\sigma^2 = \sum_{i=0}^n \frac{(x_i - \bar{x})^2}{n}$$

It is the measure of distance of observed values from the mean.

→ Standard Deviation (σ)

$$\sigma = \sqrt{\frac{\sum_{i=0}^n (x_i - \bar{x})^2}{n}}$$

a) What is the median of the following data set

32, 6, 21, 10, 8, 11, 12, 36, 17, 16, 15, 18, 40, 24, 21,
23, 24, 24, 29, 16, 32, 31, 10, 30, 35, 32, 18, 39, 12,

Find the mid range and range too

Step: Sort the data set (ascending)

After Sorting:-

6, 8, 10, 11, 12, 12, 15, 16, 16, 16, 17, 18, 18, 20, 21, 21,
23, 24, 24, 24, 29, 30, 31, 32, 32, 32, 35, 36, 39,
40

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

$$\frac{15^{\text{th}} + 16^{\text{th}}}{2} : \frac{21+21}{2} = 21$$

$$\text{Mid Range} = \frac{40-6}{2} : 23 \quad \left(\frac{\max - \min}{2} \right)$$

$$\text{Range} : \frac{30-6}{2} : 24 \quad \left(\frac{\max - \min}{2} \right)$$

a) Identify mode.

21, 19, 62, 21, 66, 28, 66, 48, 79, 59, 28, 62, 63, 63, 48, 66, 59, 66, 94, 79, 19, 94

Mode = 66 (repeated 4 times)

Modality - multimodal (66, 94) \rightarrow multiple values are repeating.

Q A dataset for analysis included only 1 attribute x .

$x: \{ 7, 12, 5, 8, 5, 9, 13, 12, 19, 7, 12, 12, 13, 3, 4, 5, 13, 8, 7, 6 \}$

a) Mean?

b) Median?

c) Standard Deviation

Sort the data

$x_r: \{ 3, 4, 5, 5, 5, 6, 7, 7, 7, 8, 8, 9, 12, 12, 12, 12, 13, 13, 13, 19 \}$

$$\text{Mean} = \frac{\text{Sum}}{\text{No of Obs}} = \frac{180}{20} = 9$$

b Median

$$\left(\frac{n}{2}\text{th}\right) + \left(\frac{n+1}{2}\text{th}\right)$$

$$\frac{10\text{th} + 11\text{th}}{2} : \frac{8+9}{2} = 8$$

c Standard Deviation:

$$\sigma^2 = S$$

Modality of Data:

- Data with only one mode is known as unimodal
- Data with two modes are known as bimodal
- Data with > 2 modes are known as multimodal

Eg:- A: {1, 2, 3, 3, 4, 4, 5, 5}

$$\text{Mode}(A) = \{3, 4\}$$

A has two modes \therefore it is bimodal

a Calculate range, variance and standard deviation

$$-4, -2, 0, -2, 6, 4, 6, 0, -6, 4$$

Sorting

-6, -4, -2, -2, 0, 4, 4, 6, 6

c) Range. $\frac{\text{max} - \text{min}}{\#}$

$$= \frac{6 - (-6)}{10} = \underline{\underline{12}}$$

D) Variance = $\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$

Mean = $\frac{\text{Sum}}{\text{No of Obs}}$

$$= \frac{-6 - 4 - 2 - 2 + 4 + 4 + 6 + 6}{10}$$

$$\frac{6}{10} = 0.6$$

$$\sigma^2 = \frac{(-6 - 0.6)^2 + (-4 - 0.6)^2 + (-2 - 0.6)^2 + (-2 - 0.6)^2 + (0 - 0.6)^2 + (0 - 0.6)^2 + (4 - 0.6)^2 + (4 - 0.6)^2 + (6 - 0.6)^2 + (6 - 0.6)^2}{10}$$

$$= 17.82$$

c) $\sigma = \sqrt{\sigma^2} = \sqrt{17.82} = 4.22$

a) Calculate range, variance and standard deviation

-3 -3 -3 -3 0 3 3 3 3

$$\text{Range} = \frac{\text{max} - \text{min}}{2} = \frac{3 - (-3)}{2} = \frac{6}{2} = 3$$

$$\text{Variance} = \sum_{i=0}^n \left(\frac{x_i - \bar{x}}{n} \right)^2$$

$$\text{mean} = \frac{\text{sum of obs}}{10 \text{ obs}} = \frac{0}{9} = 0$$

$$\sigma^2 = \frac{(-3-0)^2 + (-3-0)^2 + (-3-0)^2 + (-3-0)^2 + (0-0)^2 + (3-0)^2 + (3-0)^2 + (3-0)^2 + (3-0)^2}{9}$$

$$= \frac{72}{9} = 8$$

$$\sigma = \sqrt{8}$$

$$= \sqrt{8}$$

=

Quartiles

Quartiles are the values from the dataset which divide the dataset into four equal parts where each part of the dataset contains an equal no. of observations.

aka. 1st quartile
 Lower Quartile (Q_1) = $\frac{(n+1)}{4}$ th term

Median (Q_2) = $\frac{2(n+1)}{4}$ th term

Upper Quartile (Q_3) = $\frac{3(n+1)}{4}$ th term
aka 3rd quartile

Inter-quartile range :- (IQR)

It is the difference between Q_3 and Q_1 .

$$IQR = Q_3 - Q_1$$

a) Find the quartiles of the following age:

23, 13, 37, 16, 26, 35, 26, 35

$$Q_1 = \left(\frac{n+1}{4}\right) \text{th term}$$

$$n=8$$

$$Q_1 = \frac{9}{4} = 2.25 \text{th term}$$

Step 1: After sorting

13 16 23 26 26 35 35 37

We can rewrite Q₁ as

$$\begin{aligned}
 Q_1 &= \text{2nd term} + 0.25 (\text{3rd term} - \text{2nd term}) \\
 &= 16 + 0.25 (23 - 16) \\
 &= 16 + 0.25 \times 7 \\
 &= 16 + \cancel{1.75} \\
 &= 17.75
 \end{aligned}$$

$$Q_2 = \left(\frac{2(n+1)}{4} \right) \text{th term}$$

$$\begin{aligned}
 &= \frac{2 \times 9}{4} \\
 &= \frac{18}{4} \cdot 4.5
 \end{aligned}$$

Rewriting

$$\begin{aligned}
 Q_2 &= \text{4th term} + 0.5 (\text{5th term} - \text{4th term}) \\
 &= 26 + 0.5 (26 - 26) \\
 &= 26
 \end{aligned}$$

$$Q_3 = \left(\frac{3(n+1)}{5} \right) \text{th term}$$

$$\frac{3 \times 9}{5} = \frac{27}{5} = 5.4$$

We can rewrite Q₁ as

$$\begin{aligned}Q_1 &= \text{2nd term} + 0.25(\text{3rd term} - \text{2nd term}) \\&= 16 + 0.25(23 - 16) \\&= 16 + 0.25 \times 7 \\&= 16 + \cancel{1.75} \\&= 17.75\end{aligned}$$

$$\begin{aligned}Q_2 &= \left(\frac{2(n+1)}{4}\right)\text{th term} \\&= \frac{2 \times 9}{4} \\&= \frac{18}{4} = 4.5\end{aligned}$$

Rewriting

$$\begin{aligned}Q_2 &= 4\text{th term} + 0.5(\text{5th term} - 4\text{th term}) \\&= 26 + 0.5(26 - 26) \\&= \underline{\underline{26}}\end{aligned}$$

$$Q_3 = \left(\frac{3(n+1)}{5}\right)\text{th term}$$

$$\frac{3 \times 9}{5} = \frac{27}{5} = 5.4$$

Rewriting

$$\begin{aligned}Q_3 &= 6\text{th term} + 0.75(\text{7th term} - \text{6th term}) \\&= 35 + 0.75(35 - 35) \\&= \underline{\underline{35}}\end{aligned}$$

a) Find Quartiles of the given data
10, 30, 5, 12, 20, 20, 25, 15, 18

After sorting

$$5, 10, 12, 15, 18, 20, 25, 30, 40$$

$$n=9$$

$$Q_1 = \left(\frac{n+1}{4}\right)\text{th term}$$

$$= \frac{10}{4} = 2.5$$

Rewriting

$$\begin{aligned}Q_1 &= 2\text{nd term} + 0.5(\text{3rd term} - \text{2nd term}) \\&= 10 + 0.5(12 - 10)\end{aligned}$$

$$= 10 + 0.5 \times 2$$

$$= \underline{\underline{11}}$$

$$Q_2 = \left(\frac{2(n+1)}{5}\right)\text{th term}$$

$$\frac{Q_1(10)}{4} = 5 = 18$$

$$Q_3 : \left[\frac{3(n+1)}{4} \right] \text{th term}$$

$$= \frac{30}{4} = 7.5$$

Rounding

$$Q_3 : 7^{\text{th}} \text{ term} + 0.5 (8^{\text{th}} \text{ term} - 7^{\text{th}} \text{ term})$$

$$= 25 + 0.5 (30 - 25)$$

$$= 25 + 0.25$$

$$= 25.25$$

HW Find 1st, and, 3rd quartile of
8, 5, 15, 20, 18, 30, 40, 25

Find IQR.

Solt.: 5, 8, 15, 18, 20, 25, 30, 40

Q₁: $\left[\frac{(n+1)}{4} \right] \text{th term}$

$$= \frac{9}{4}, \text{ or } 2.25^{\text{th}} \text{ term}$$

$$= \text{2nd term} + 0.25 (3^{\text{rd}} \text{ term} - \text{2nd term})$$

$$= 8 + 0.25 (15 - 8) = 9.75$$

$$Q_2 = [2(n+1)/4]^{th} \text{ term}$$

$$= \frac{18}{4} = 4.5^{\text{th}} \text{ item}$$

4^{th} term + as (5^{th} term - 4^{th} term)

$$= 18 + 0.5(20 - 18)$$

$$= 19$$

$$Q_3 = [3(n+1)/4]^{th} \text{ term}$$

$$= \frac{27}{4} = 6.75$$

6^{th} term + $0.75(7^{\text{th}}$ term - 6^{th} term)

$$= 25 + 0.75(30 - 25)$$

$$= 28.75$$

$$\text{IQR} = Q_3 - Q_1$$

$$= 28.75 - 9.75$$

$$= \underline{\underline{19}}$$

Box Plot:-

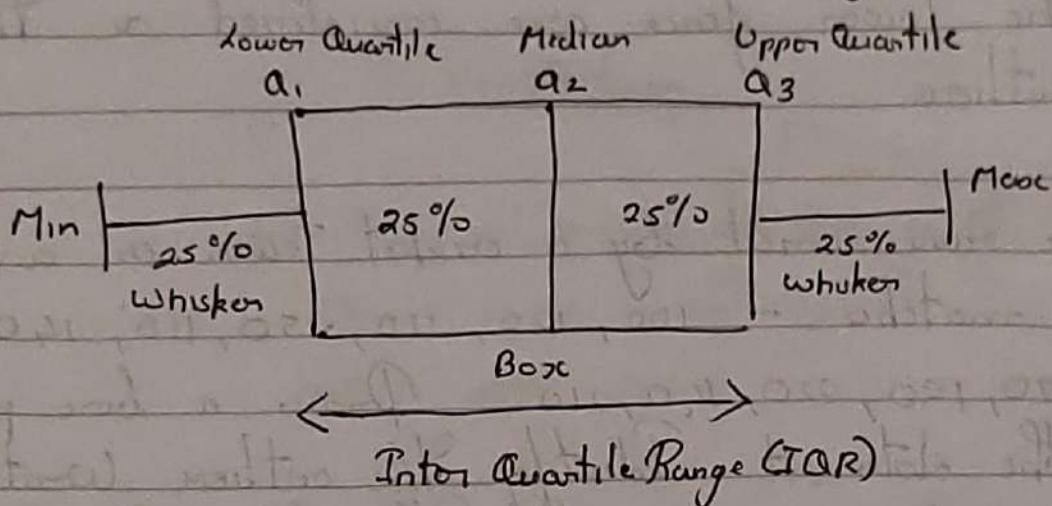
It is a graphical method to visualize data distributions for gaining insights and making informed decisions.

A box plot is also known as Whisker plot and is created to display the summary of the set of data values.

Elements of Box Plot:-

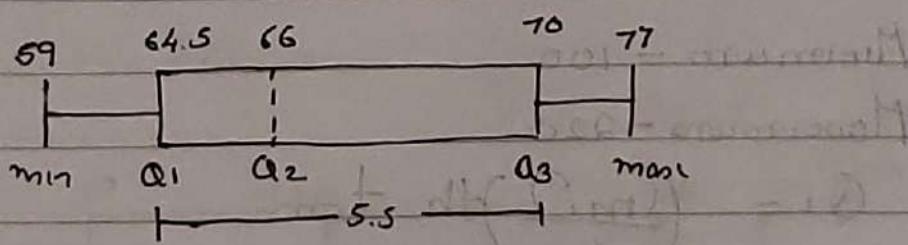
- i. Minimum :- It is the minimum value in the dataset
- ii. Q_1 :- 25% of data lies below the First (lower) Quartile.
- iii. Median (Q_2) :- Mid-point of the data set
- iv. Third Quartile (Q_3) :- 75% data lies below the 3rd quartile (upper quartile)
- v. Maximum :- Maximum value in the dataset

This is The 5-Number Summary Plot.



- a) Heights of 40 students having the 5-number summary as follows. Construct a box plot for the dataset

$$\text{min} = 59 \quad \text{max} = 77 \quad Q_1 = 64.5 \quad Q_2 = 66 \quad Q_3 = 70$$



Outlier :- (Finding out outliers)

$$\text{Upper Fence} = Q_3 + (1.5 * \text{IQR})$$

$$\text{Lower Fence} = Q_1 - (1.5 * \text{IQR})$$

The values above the upper fence and below the lower fence are considered as the outliers.

- Q The runs scored by a cricket team in a league of 13 matches - 100, 120, 110, 150, 110, 140, 130, 170, 120, 220, 140, 110. Draw a box plot for the data set. Identify the outliers. Write the 5 number summary of the data. Identify the outliers.

Sort \rightarrow 100, 110, 110, 110, 120, 120, 130, 140, 140, 150, 170, 220

Minimum - 100

Maximum - 220

$$Q_1 = \left(\frac{n+1}{4} \right) \text{th term}$$

$\therefore \frac{12+1}{4} \text{th term}$

$\therefore 3.25 \text{th term}$

$$\text{3rd term} + 0.25 (\text{4th term} - \text{3rd term})$$

$$110 + 0.25 (110 - 110)$$

110

$$Q_3 = \frac{3(n+1)}{4} \text{ th term}$$

$$\frac{39}{4} = 9.75 \text{ th term}$$

$$= 9 \text{ th term} + 0.75 (10^{\circ} \text{ th term} - 8 \text{ th term})$$

$$= 120 + 0.75 (140 - 110)$$

$$= 147.5$$

$$Q_1 = \frac{2(n+1)}{4}$$

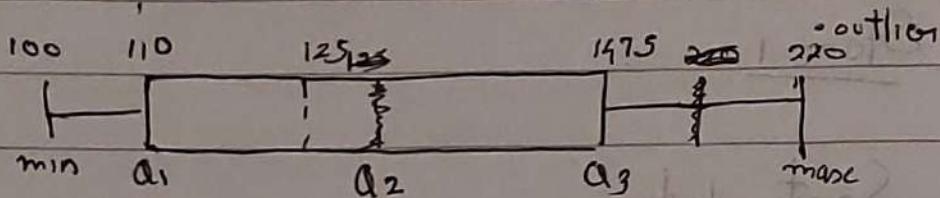
$$\frac{26}{4} = 6.5$$

$$= 6 \text{ th term} + 0.5 (7 \text{ th term} - 6 \text{ th term})$$

$$= 120 + 0.5 (130 - 120)$$

$$= 125$$

$$IQR = Q_3 - Q_1 = 37.5$$



$$\xleftarrow{-} 37.5 \xrightarrow{+}$$

$$\begin{aligned}\text{Upper Fence} &= Q_3 + (1.5 \times IQR) \\ &= 147.5 + (1.5 \times 37.5) = 203.75\end{aligned}$$

$$\text{Lower Fence} = Q_1 - (1.5 \times IQR) = 53.75$$

Oct 11 or 220

- Q An IT company has two stores that sell computers. The company record the number of sales each store made each month. In the past 12 months we have the following no. of sold computers.

Store 1 350, 460, 20, 160, 580, 250, ~~220~~ 210, 120, 200, 510, 290, 380

Store 2 520, 180, 260, 380, 80, 500, 630, 420, 210, 70, 440, 140.

Give the summary and draw bar plot.

Store 1

Sales data

20, 120, 160, 200, 210, 250, 280, 350, 380, 460, 510, 580

Minimum = 20

Maximum = 580

$$\text{Median} (Q_2) = \frac{(n+1)}{4} \text{th term}$$

$$= \frac{2 \times 13}{4} \text{th term}$$

$$= \frac{26}{4} \text{th term}$$

$$= 6.5$$

$$5\text{th term} + 0.5 (7\text{th term} - 5\text{th term}) \\ 250 + 0.5 (290 - 250)$$

$$250 + 0.5 \times 40$$

$$250 + 20$$

$$= 270$$

$$Q_1 = \frac{(n+1)}{4} \text{th term}$$

$$= \frac{13}{4} = 3.25$$

$$= 3\text{rd term} + 0.25 (4\text{th term} - 3\text{rd term})$$

$$= 160 + 0.25 (200 - 160)$$

$$= 160 + \frac{25 \times 10}{100}$$

$$= 170$$

$$Q_3 = \frac{3}{4}(n+1) \text{ th term}$$

$$\frac{3 \times 13}{4}$$

$$\frac{39}{4} \text{ th term}$$

$$= 9.75 \text{ th term}$$

$$9 \text{ th term} + 0.75(10 \text{ th term} - 9 \text{ th term})$$

$$380 + 0.75(460 - 380)$$

$$380 + 0.75 \times 80$$

$$380 + 60$$

$$440$$

$$\text{IQR. } Q_3 - Q_1$$

$$= 440 - 170$$

$$= 270$$

5-number summary

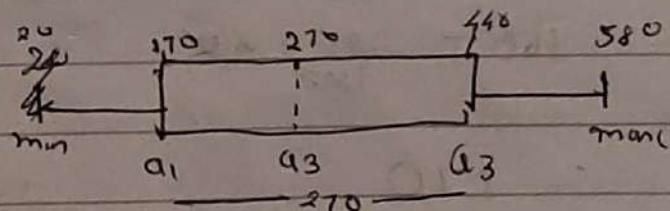
Minimum: 20

$Q_1: \cancel{100} 170$

Median: 270

$Q_3: \cancel{420} 440$

Maximum: 580



Stone 2

Sort the data

70, 80, 140, 180, 210, 260, ~~320~~, 380, 420, 440, 500,
520, 630

Minimum: 70 Maximum: 630

$$Q_1 = \frac{(n+1)}{4} \text{th term}$$

$\frac{13}{4}$ th term

3.25 th term

$$\begin{aligned} &= 3^{\text{rd}} \text{ term} + 0.25 (4^{\text{th}} \text{ term} - 3^{\text{rd}} \text{ term}) \\ &= 140 + 0.25 (180 - 140) \\ &= \underline{150} \end{aligned}$$

$$Q_2 = \frac{2(n+1)}{4} \text{th term}$$

$$\frac{2 \times 13}{4} = 6.5 \text{ th term}$$

$$\begin{aligned} &\cdot 6^{\text{th}} \text{ term} + 0.5 (7^{\text{th}} \text{ term} - 6^{\text{th}} \text{ term}) \\ &\cdot 260 + 0.5 (380 - 260) \\ &\cdot 260 + 0.5 (60) \quad \cdot \underline{320} \\ &= \end{aligned}$$

$$Q_3: \frac{3(15+24)}{4} \text{ ft from } \frac{3+15}{4} \text{ ft from}$$

9.75 ft from

$$9\text{th term} + 0.75 \left(\frac{10}{4}\text{th term} - 9\text{th term} \right)$$

$$420 + 0.75 \left(\frac{50}{4} - 42 \right)$$

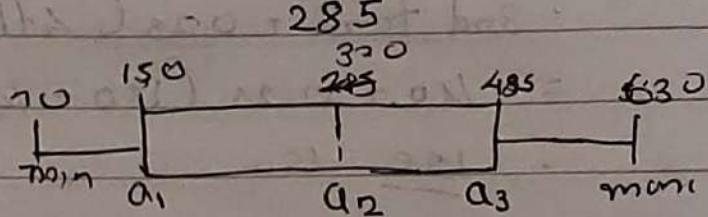
$$420 + 0.75 \times 20$$

$$420 + 15$$

$$\underline{\underline{435}} \quad 485$$

IQR. $Q_3 - Q_1$

$$= 485 - 150$$



s-number summary

Minimum: 70

Maximum: 630

$Q_1: 150$

$Q_2: 225$

$Q_3: 485$

2 Consider the data set: 180, 156, 9, 176, 163, 1827, 166, 171
 Draw box plot and identify outliers

After sorting:

9, 156, 163, 166, 171, 176, 180, 1827

a. $\frac{(n+1)}{4}$ th term

$$\therefore \left(\frac{9}{4}\right)\text{th term} = 2.25\text{th term}$$

$$= 2\text{nd term} + 0.25(3\text{rd term} - 2\text{nd term})$$

$$= 156 + 0.25(163 - 156)$$

$$= 156 + 0.25 \times 7$$

$$= 156.75$$

=

a₂. $\frac{2(n+1)}{4}$ th term

$$\therefore \frac{9}{2}\text{th term}$$

$$= 4.5\text{th term}$$

$$= 4\text{th term} + 0.5(5\text{th term} - 4\text{th term})$$

$$= 166 + 0.5(171 - 166)$$

$$= 166 + 0.5 \times 5 \quad : \quad 168.5$$

=

$$Q_3 = \frac{3(C_{n+1})}{5}$$

$$\cdot \frac{3 \times 9}{5} : \frac{27}{5} \text{ th term}$$

6.75 th term

6^{th} term + 0.75 (7^{th} term - 6^{th} term)

$$176 + 0.75 (180 - 176)$$

$$\begin{array}{r} 176 + 0.75 \times 5 \\ \hline 179 \end{array}$$

$$\text{IQR} = Q_3 - Q_1$$

$$= 179 - 157.75$$

$$= 21.25$$

~~157.75 + 21.25~~

Outlier

Upper Fence = $Q_3 + (1.5 \times \text{IQR})$

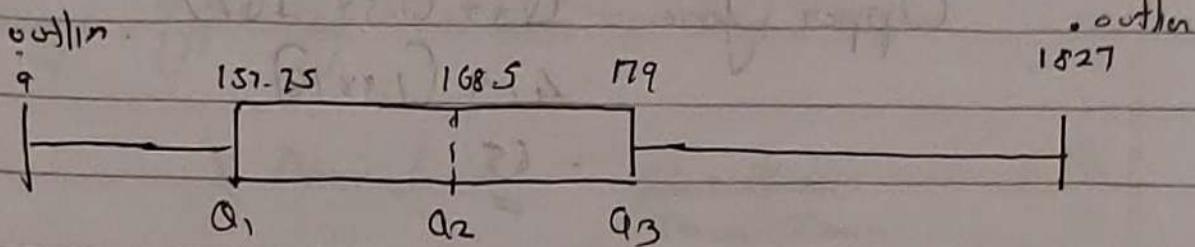
$$= 179 + (1.5 \times 21.25)$$

$$= 179 + 31.875$$

$$= 210.875$$

$$\text{Lower Fence} = Q_1 - (1.5 \times IQR)$$

$$= 157.75 - 81.875 \\ = 75.875$$



Option: - 1827, 9

- 3 22, 24, 25, 28, 29, 31, 33, 37, 41, 53, 64

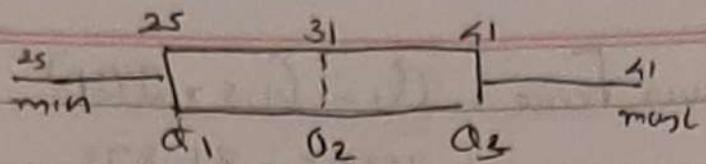
min - 22

$$Q_1 = ((n+1)/4)^{\text{th}} \text{ term} \\ = 12/4 = 3 \rightarrow 25$$

$$Q_2 = ((n+1)/4)^{\text{th}} \text{ term} \\ = 24/4 = 6 \rightarrow 31$$

$$Q_3 = ((3(n+1)/4)^{\text{th}} \text{ term} \\ = 36/4 = 9 \rightarrow 41$$

$$\text{IQR} = Q_3 - Q_1 = 41 - 25 = 16$$



Upper fence = $Q_3 + (1.5 \times IQR)$
 $= 41 + (1.5 \times 16)$
 $= 65$

Lower fence = $Q_1 - (1.5 \times IQR)$
 $= 25 - 24$

No outlier

Normalization:-

Preprocessing technique

Scaling technique in machine learning, applied during data preparation to change the values of numeric column in the dataset to use a common scale

Methods of Normalization

- Min-Max Normalization
- Z-Score Normalization
- Decimal Scaling

Min-Max Normalization

→ Min-max scaling transforms features to a specified range, typically between 0 and 1.

$$X_{\text{Normalized}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

- a) Use the min-max method to normalize the data
 200, 300, 400, 600, 1000

$$\min = 200$$

$$\max = 1000$$

$$x \quad \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

$$200 \quad \frac{200 - 200}{1000 - 200} = 0$$

$$300 \quad \frac{300 - 200}{1000 - 200} = 0.125$$

$$400 \quad \frac{400 - 200}{1000 - 200} = 0.25$$

$$600 \quad \frac{600 - 200}{1000 - 200}$$

$$1000 = \frac{1000 - 200}{1000 - 200} \quad \textcircled{1}$$

2 Z-Score Normalization

- Also known as standardization.
- Transforms the data so that the mean is 0 and standard deviation is 1.
- Process adjust the data values based on how far they deviate from the mean measured in units of standard deviation.

$$Z = \frac{x - \mu}{\sigma} \quad \sigma = \text{S.D.} \quad \mu = \text{mean}$$

a Use z-Score to normalize

200, 300, 400, 600, 1000

$$\mu = \frac{200 + 300 + 400 + 600 + 1000}{5}$$

$$\cdot \frac{2500}{5} = 500$$

$$\text{Variance} = \frac{(200 - 500)^2 + (300 - 500)^2 + (400 - 500)^2 + (600 - 500)^2 + (1000 - 500)^2}{5}$$

$$\begin{array}{r}
 \text{S.D.} \\
 \left. \begin{array}{r}
 90000 + 40000 + 10000 + 1000 + \\
 \hline
 250000
 \end{array} \right\} 5 \\
 = 209.79 \quad (282.8)
 \end{array}$$

Decimal Scaling Normalization

- In this technique, we move the decimal point of value of the attribute
- The movement of the decimal points totally depends on the maximum value among all values in the attribute.

Normalized Value of Attribute = $\frac{v}{10^j}$

Step 1 find maximum absolute value.

Step 2 Find the value of j (if max value is 100 Then j is i.e. the no. of digits)

Step 3 $\frac{v}{10^j}$

Q Normalize using decimal scaling

200, 300, 400, 600, 1000

$$\frac{v}{10^j}$$

$\max = 1000 \therefore j = 4$
(max absolute value)

$x \rightarrow x\text{-normalized } (v)$

$$200 \quad \frac{200}{10^4} = \frac{200}{10000} = 0.02$$

$$300 \quad \frac{300}{10^4} = \frac{300}{10000} = 0.03$$

$$400 \quad \frac{400}{10^4} = \frac{400}{10000} = 0.04$$

$$600 \quad \frac{600}{10^4} = \frac{600}{10000} = 0.06$$

$$1000 \quad \frac{1000}{10^4} = \frac{1000}{10000} = 0.1$$

Data Similarity and Dissimilarity

Data Similarity is The measure of how alike two data objects are.

- Data dissimilarity is a measure of how different the data objects are.
 - The similarity measure is usually expressed as a numerical value
 - The value is high when the objects are more alike.
- * Proximity Refers to Similarity or Dissimilarity
- Methods for Measuring Similarity / Dissimilarity,
- Euclidean Distance :-
 - Manhattan Distance
 - Minkowski Distance
 - Supreme Distance

Euclidean Distance:

Let the two points i and j ,
The distance b/w the points is given by

$$d(i, j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{in} - x_{jn})^2}$$

- a) Compute the Euclidean Distance b/w the two objects $(22, 1, 42, 14)$ and $(20, 0, 36, 8)$

$$\begin{aligned}
 d_{C_1, j} &= \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{in} - x_{jn})^2} \\
 &= \sqrt{(2-20)^2 + (1-0)^2 + (42-36)^2 + (10-8)^2} \\
 &= \sqrt{2^2 + 1^2 + 6^2 + 2^2} \\
 &= \sqrt{45} \\
 &= 6.71
 \end{aligned}$$

Manhattan Distance

$$d_{C_1, j} = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{in} - x_{jn}|$$

a) $(22, 1, 42, 10)$ and $(20, 0, 36, 8)$

$$\begin{aligned}
 d_{C_1, j} &= |22-20| + |1-0| + |42-36| + |10-8| \\
 &= |2+1+6+2| \\
 &= \underline{\underline{11}}
 \end{aligned}$$

Minkowski Distance

$$d_{C_1, j} = \left(|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{in} - x_{jn}|^p \right)^{\frac{1}{p}}$$

a) Compute The Minkowski Distance

$(22, 1, 42, 10)$ and $(20, 0, 36, 8)$ $p=3$

$$d_{(i,j)} = \left(|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{in} - x_{jn}|^p \right)^{1/p}$$

$$= \sqrt[3]{(22-20)^3 + (1-0)^3 + (42-36)^3 + (10-8)^3}$$

$$= \sqrt[3]{(8+1+6^3+2^3)}$$

$$= \sqrt[3]{233}$$

$$= 6.15$$

Supreme Distance :-

$$d_{(i,j)} = \max(|x_{i1} - x_{j1}|, |x_{i2} - x_{j2}|, \dots, |x_{in} - x_{jn}|)$$

a) Find The supreme distance between 2 objects,
 $(22, 1, 42, 10)$ and $(20, 0, 36, 8)$

$$d_{(i,j)} = \max(|x_{i1} - x_{j1}|, |x_{i2} - x_{j2}|, \dots, |x_{in} - x_{jn}|)$$

$$= \max(22-20, 1-0, 42-36, 10-8)$$

$$= \max(2, 1, 6, 2)$$

$$= \underline{\underline{6}}$$

Module 2

Database:-

- Stores data in tables
- Deals with operational or transactional data (current data).
- Can store MBs to GBs of data.
- Used for OLTP (Online Transaction Processing)