

① Discuss the privacy issues in data mining with examples.

Ans → Privacy Issues in data mining

1. minimal protection setup:

Most of the time data is protected by security measures like anti-viruses, usernames, passwords, etc. which doesn't protect the data in long term.

2. Access controls:

Access controls verifies the identity of the person ~~by~~ trying to access data. Single layer access control is not a very secured option.

3. Non-Verified Data updation:

Many times data is collected thoroughly and updated without the ~~verified~~ verification of source.

4. Security - Architect Evaluation:

To save money and time, a lot of Organizations skip the process of audit of the Security architect, which makes easier to hack the data collection.

5. Filtering and Validating External Sources:

Whenever an unauthorised device is able to connect to the security system, it gives an entry point for vulnerabilities.

Ex:- people bring office work home and access the official data via their personal devices, which can create a loophole.

② For a certain dataset, the values of the attribute Age is given as follows.

25, 30, 15, 16, 33, 35, 70, 52, 13, 25, 33, 25, 40,
36, 35, 19, 25, 16, 20, 13, 22, 45, 21, 35, 20.

(a) find the mean, median, ~~and~~ mode and mid-range of the data.

$$\text{mean} = \frac{1}{25} \sum_{i=1}^{25} x_i$$

$$\Rightarrow \bar{x} = \frac{1}{25} [2 \times 13 + 15 + 2 \times 16 + 19 + 2 \times 20 + \\ 21 + 22 + 4 \times 25 + 30 + 2 \times 33 + \\ 3 \times 35 + 36 + 40 + 45 + 52 + 70]$$

$$\Rightarrow \bar{x} = \frac{1}{25} \times 719$$

$$\Rightarrow \bar{x} = 28.76.$$

$\therefore n \rightarrow \text{odd}$

$$\text{median} = x_{\left[\frac{n+1}{2}\right]} = x_{\left[\frac{25+1}{2}\right]} = x_{13}$$

$\therefore \text{median} = 25 //$

age	no. of Students / people
13	2
15	1
16	2
19	1
20	2
21	1
22	1
25	4
30	1
33	2
35	3
36	1
40	1
45	1
52	1
70	1

$\rightarrow \therefore$ highest no. of people are of age 25.

$\therefore \text{Mode} = 25 //$

$$\text{Mid range} = \frac{\min + \max}{2}$$

$$\Rightarrow \text{midrange} = \frac{13 + 70}{2} = \frac{83}{2}$$

$$\therefore \text{midrange} = 41.5$$

(b) Calculate the Variance and Standard deviation.

$$n = 25$$

$$\Sigma x = 719$$

$$\bar{x} = 28.76$$

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

$$\Rightarrow \sigma^2 = \frac{1}{24} \left[(13 - 28.76)^2 \times 2 + (15 - 28.76)^2 + \right.$$

$$(16 - 28.76)^2 \times 2 + (19 - 28.76)^2 +$$

$$(20 - 28.76)^2 \times 2 + (21 - 28.76)^2 +$$

$$(22 - 28.76)^2 + (25 - 28.76)^2 \times 4 + (30 - 28.76)^2 +$$

$$(33 - 28.76)^2 \times 2 + (35 - 28.76)^2 \times 3 +$$

$$(36 - 28.76)^2 + (40 - 28.76)^2 + (45 - 28.76)^2 \\ \left. + (52 - 28.76)^2 + (70 - 28.76)^2 \right]$$

$$\Rightarrow \sigma^2 = \frac{1}{24} \left[2 \times 248 \cdot 3776 + 189 \cdot 3376 + 1628176 \times 2 \right. \\ + 95 \cdot 2576 + 2 \times 76 \cdot 7376 + 60 \cdot 2176 + \\ 45 \cdot 6976 + 4 \times 14 \cdot 1376 + 1 \cdot 5376 + \\ 2 \times 17 \cdot 9776 + 3 \times 38 \cdot 9376 + \\ 5 \cdot 2 \cdot 4176 + 126 \cdot 3376 + 263 \cdot 7376 + \\ \left. 540 \cdot 0976 + 1,700 \cdot 7376 \right]$$

$$\Rightarrow \sigma^2 = \frac{1}{24} \times 4260.56$$

$$\Rightarrow \sigma^2 = 177.5233 //$$

$$\text{Standard deviation } (\sigma) = \sqrt{\sigma^2}$$

$$\Rightarrow \sigma = \sqrt{177.5233}$$

$$\Rightarrow \sigma = 13.32$$

(c) Give the five-number summary of the data.

Solⁿ \Rightarrow 13, 13, 15, 16, 16, 19, 20, 20, 21, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 36, 40, 45, 52, 70.

Minimum = 13.

$$\text{Quartile 1} = \left(\frac{N+1}{4} \right)^{\text{th}} \text{ term}$$

$$= (6.5)^{\text{th}} \text{ term}$$

i.e. between 19 & 20 [6th - 7th]

$$= \frac{19+20}{2} = 19.5.$$

Median = 25.

$$\text{Quartile 3} = \left(\frac{N+1}{4} \right) 3^{\text{th}} \text{ term}$$

$$= 19.5^{\text{th}} \text{ term}$$

i.e. between 35 & 35 [19th - 20th]

$$= \frac{35+35}{2} = 35 //$$

Maximum = 70.

d) Create a box plot of the data.

Solⁿ

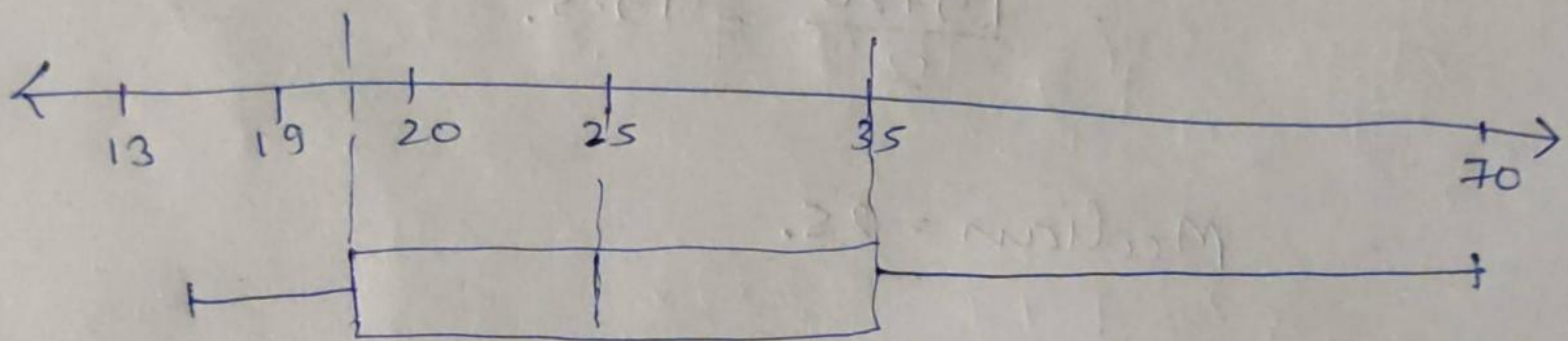
Minimum = 13.

$Q_1 = 19.5$

Q_2 (Median) = 25

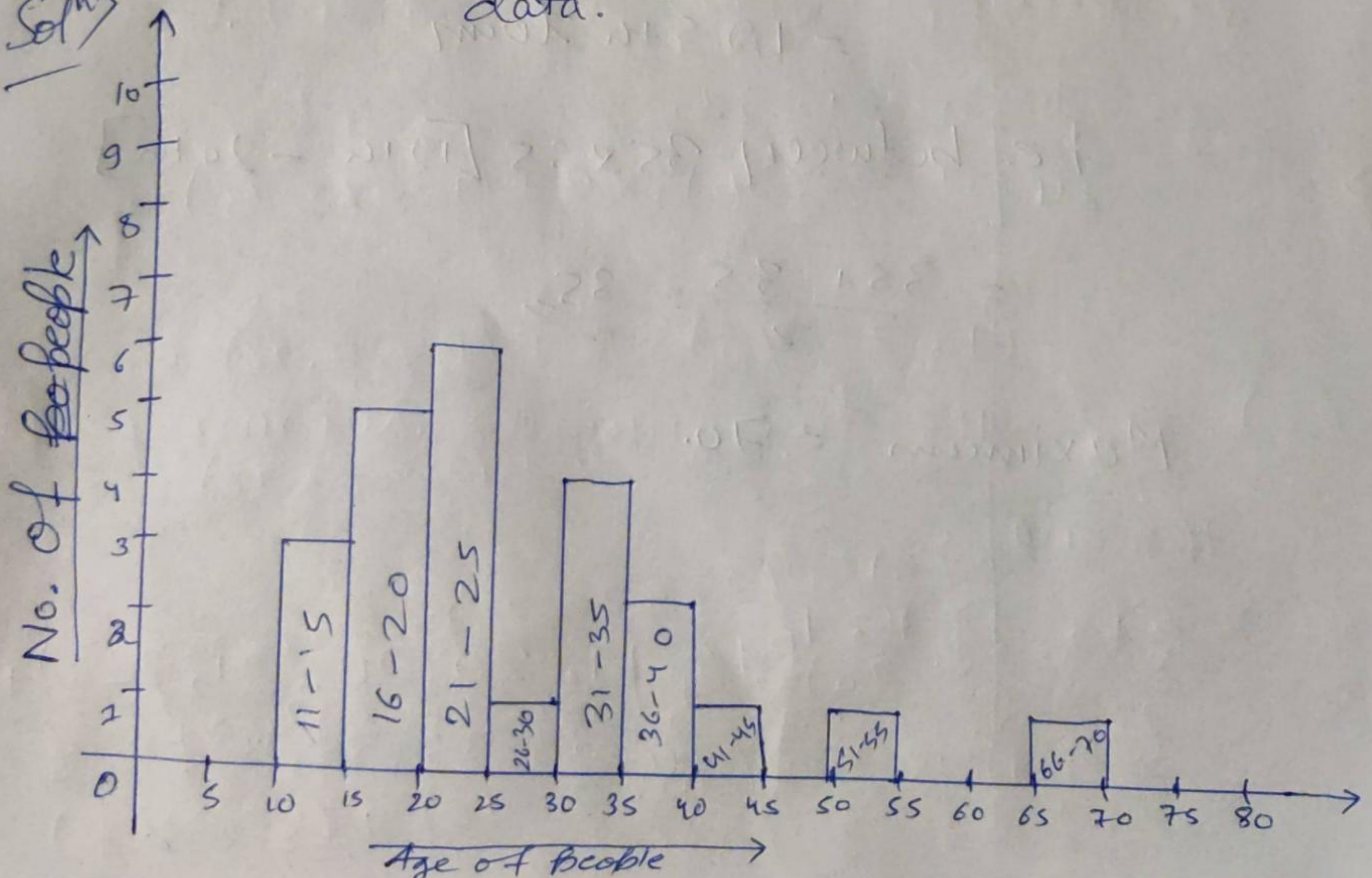
$Q_3 = 35$

Maximum = 70.



e) Draw a histogram to represent the data.

Solⁿ



(3) Create a dissimilarity matrix for the given data:-

Item	Colour	Price	Size
1	Blue	500	Small
2	Green	100	Large
3	Red	300	Small
4	Green	600	Medium

$$[d = \frac{P-M}{P}]$$

Solⁿ

$P \rightarrow$ no. of attributes = 03

$$d(2,1) = \frac{3-0}{3} = 1$$

$$d(3,1) = \frac{3-1}{3} = 0.67$$

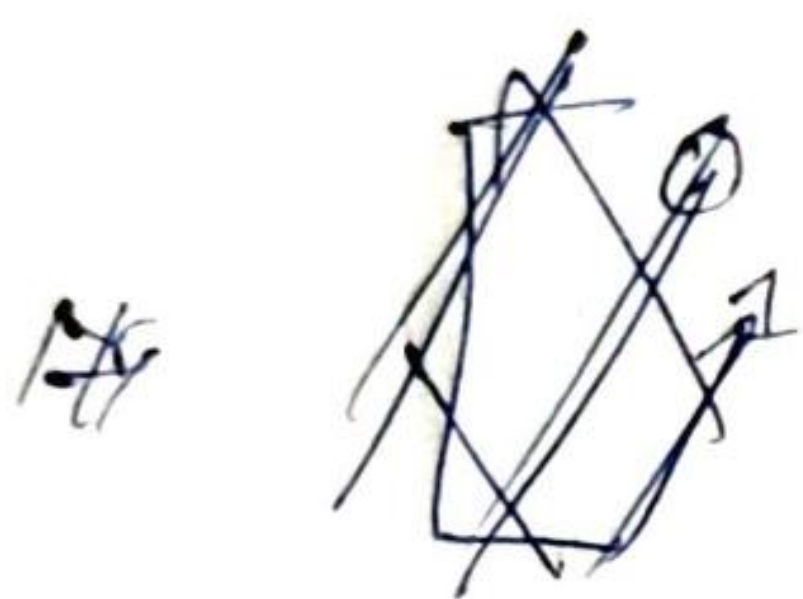
$$d(3,2) = \frac{3-0}{3} = 1$$

$$d(4,1) = \frac{3-0}{3} = 1$$

$$d(4,2) = \frac{3-1}{3} = 0.67$$

$$d(4,3) = \frac{3-0}{3} = 1.$$

$$\begin{bmatrix} d(1,1) \\ d(2,1) & d(2,2) \\ d(3,1) & d(3,2) & d(3,3) \\ d(4,1) & d(4,2) & d(4,3) & d(4,4) \end{bmatrix}$$



$$= \begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 0.67 & 1 & 0 & \\ 1 & 0.67 & 1 & 0 \end{bmatrix}$$

Ag