

Pre-Processing

1 Arranged Data : 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70

a) Mean & Median

$$\therefore \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$= \frac{13 + 15 + \dots + 46 + 52 + 70}{27}$$

$$= 809 / 27 \quad \boxed{\bar{x} = 29.96}$$

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

$$= \frac{27+1}{2}$$

$$= \frac{28}{2} = 14^{\text{th}} \text{ value}$$

$$\boxed{\text{Median} = 25}$$

b) Mode (Also comment on modality)

13, 15, 19, 21, 30, 36, 40, 45, 46, 52, 70 occur 1 time

16, 20, 22, 33 occur 2 times

25, 35 occur 4 times. Hence the dataset is bimodal

FOR EDUCATIONAL USE with 25 & 35 as modes

c) MidRange

$$= \frac{\text{Max} - \text{Min}}{2}$$

$$= \frac{70 - 13}{2}$$

$$= 41.5$$

d) Find Q_1 & Q_3

$$Q_1 \text{ (First Quartile)} = 25\% \text{ of data}$$

$$= \frac{25}{100} \times 27$$

$$= 6.75$$

$$\approx 7^{\text{th}} \text{ term}$$

$$Q_1 = 20$$

$$Q_3 \text{ (Third Quartile)} = 75\% \text{ of data}$$

$$= \frac{75}{100} \times 27$$

$$= 20.25$$

$$\approx 20^{\text{th}} \text{ term}$$

$$Q_3 = 35$$

e) Five Number Summary

$$\text{Min} = 13$$

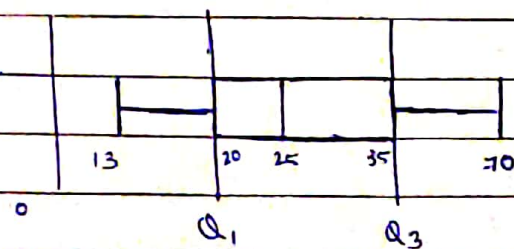
$$Q_1 = 20$$

$$\text{Median} = 25$$

$$Q_3 = 35$$

$$\text{Max} = 70$$

4) Boxplot



9) How is quantile quantile graph different from quantile plot

A quantile plot is a graphical method used to show the approximate percentage of values below or equal to the independent variable in an univariate distribution. Thus it displays quantile information for all the data, where the values measured for the independent variables are plotted against their corresponding quantile.

A quantile quant plot graphs the quantile on univariate distribution against the corresponding quantiles of another univariate distribution. Both axes display range of values measured for their respective corresponding distribution and points are plotted that correspond to the quantile values of the two distributions.

Q2

Age	Frequency
1-5	200
6-15	450
16-20	300
21-50	1500
51-80	700
81-110	44

Compute an approximate median value for the data.

From the table $N = 3194$

Median = $\frac{N}{2} = 1597^{\text{th}}$ value \therefore Median = 20.5 - 50.5

$$\text{Median} = L_1 + \left(\frac{N/2 - (\sum f_{\text{area}})}{f_{\text{area}}(\text{median})} \right) \cdot \text{width}$$

$$L_1 = 20.5$$

$$N = 3194$$

$$\sum f_{\text{area}} = 200 + 456 + 300$$

$$f_{\text{area median}} = 1500$$

$$\text{width} = 30$$

$$\begin{aligned} \text{Median} &= 20.5 + \left[\frac{1597 - 950}{1500} \right] \times 30 \\ &= \underline{\underline{33.44 \text{ years}}} \end{aligned}$$

3	age	23	23	27	27	39	41	47	49	50
	fat %	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2	31.2

	age	52	54	54	56	57	58	58	60	61
	% fat	34.6	42.5	28.8	33.4	30.2	34.1	32.9	41.2	35.7

a) Calculate mean, median, SD of age & fat

$$\text{Mean for age} = \frac{836}{18} = 46.44$$

$$\text{Median} = \frac{50 + 52}{2} = 51$$

$$SD = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}} = \sqrt{\frac{2972.2}{18}} = 12.85$$

For fat

$$\begin{aligned}\text{Mean} &= \frac{518}{18} \\ &= 28.78\end{aligned}$$

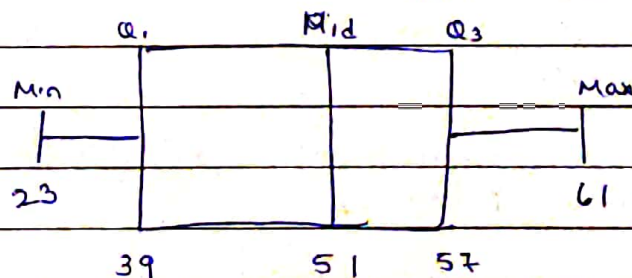
$$\begin{aligned}\text{Median} &= \frac{30.2 + 31.2}{2} \\ &= 30.7\end{aligned}$$

$$\begin{aligned}\text{SD} &= \sqrt{\frac{1456.70}{18}} \\ &= 8.99\end{aligned}$$

b) Boxplot :

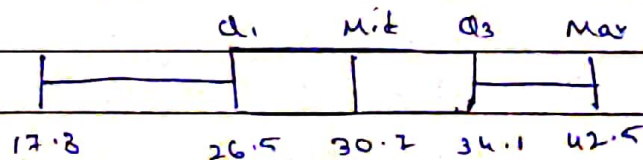
Max : 61

Age : min = 23 $Q_1 = 39$ Med = 51 $Q_3 = 57$



Max = 42.5

Fat : min = 17.8 $Q_1 = 26.5$ median = 30.7 $Q_3 = 34.1$



$$4 \quad (22, 1, 42, 10) \quad \& \quad (20, 0, 36, 8)$$

\therefore Euclidean Distance

$$\begin{aligned}
 &= \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + (x_{i3} - x_{j3})^2 + (x_{i4} - x_{j4})^2} \\
 &= \sqrt{(22 - 20)^2 + (1 - 0)^2 + (42 - 36)^2 + (10 - 8)^2} \\
 &= \sqrt{45} \\
 &= 6.708
 \end{aligned}$$

\therefore Minowski distance

$$\begin{aligned}
 &= \sqrt[n]{|x_{i1} - x_{j1}|^n + |x_{i2} - x_{j2}|^n + \dots + |x_{in} - x_{jn}|^n} \\
 &\quad n = 3 \quad \text{given} \\
 &= \sqrt[3]{|22 - 20|^3 + |1 - 0|^3 + |42 - 36|^3 + |10 - 8|^3} \\
 &= \sqrt[3]{8 + 1 + 216 + 8} \\
 &= \sqrt[3]{233} \\
 &= \underline{\underline{6.153}}
 \end{aligned}$$

\therefore Manhattan Distance

$$\begin{aligned}
 &= |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots \\
 &= |22 - 20| + |1 - 0| + |42 - 36| + |10 - 8| \\
 &= 2 + 1 + 6 + 2 \\
 &= 11
 \end{aligned}$$

∴ Supremum Distance

$$\begin{aligned}
 &= \lim_{n \rightarrow \infty} \left(\sum_{j=1}^p |x_{ij} - x_{jt}|^2 \right)^{1/n} \\
 &= \max_f |x_{ij} - x_{jt}| \\
 &= \max (2, 1, 6, 2) \\
 &= \underline{6}
 \end{aligned}$$

5 Suppose we have 2-D data set

x_i	A_1	A_2
x_1	1.5	1.7
x_2	2	1.9
x_3	1.6	1.8
x_4	1.2	1.5
x_5	1.5	1.0

Consider the data as 2-D data points. Given a new data point, $x_c = 1.4, 1.6$ as a query, rank the database points based on similarity with the query using Euclidean distance, Manhattan distance, supremum distance and cosine similarity.

→ cosine similarity = $\frac{x^t \cdot y}{\|x\| \|y\|}$ x^t - transpose of x

$\|x\|$ = Euclidean norm $\|y\|$ = Euclidean norm

From points (1.4, 1.6) we get

	Euclidean	Manhattan	Supremum	cosine similarity
x_1	0.1414	0.2	0.1	0.99999
x_2	0.6708	0.9	0.6	0.99575
x_3	0.2828	0.4	0.2	0.99997
x_4	0.2236	0.3	0.2	0.99903
x_5	0.6083	0.7	0.6	0.96536

\therefore Ranks.

Euclidean: x_1, x_4, x_3, x_5, x_2

Manhattan: x_1, x_4, x_3, x_5, x_2

Supremum: x_1, x_4, x_3, x_5, x_2

Cosine: x_1, x_3, x_4, x_2, x_5

\rightarrow Normalise the data set

\therefore	A_1	A_2
x_1	0.661682	0.74984
x_2	0.72500	0.68275
x_3	0.66436	0.74741
x_4	0.62470	0.78087
x_5	0.83250	0.55470

\therefore Recompute Euclidean

Eud Dist

	x_1	0.00415
\therefore Rank	x_2	0.09217
x_1, x_3, x_4, x_2, x_5	x_3	0.00781
	x_4	0.04409
	x_5	0.26320

FOR EDUCATIONAL USE