**Skewness:-** Skewness is the deviation from Symmentry observed in data.

   * measure of asymmetry.



Positive Skewness
(right-modal)

$L < R$   $L = R$   $L > R$

Left    Middle value    Right

(No skewness)

Negative-Skewness
(left-modal)

Ex: data: $\{1, 2, 2, 3, 3, 3, 4, 4, 5\}$

frequency $\uparrow$

$L$ $R$ → Symmetrical

1 2 ③ 4 5 →
number

median = ③ (which divides data into 2 equal parts)

mean = median = mode
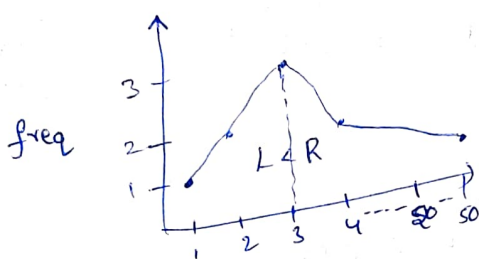
all the central tendency values are equal.

$\mu = \dfrac{27}{9} = 3.$

median = 3
mode = 3

Ex2: data = $\{1, 2, 2, 3, 3, 3, 4, 4, 50\}$

Positive skewed
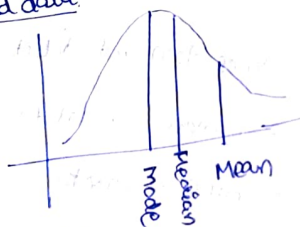


freq

L < R

$\mu = \dfrac{72}{9} = 8$, median = 3, mode = 3.

* mean is prone to outliers.

Ex3: data = $\{-8, 2, 2, 3, 3, 3, 4, 4, 5\}$

Left skewed data.

$\mu = \dfrac{18}{9} = 2$, median = 3, mode = 3

Positive (or) Right skewed data:



mode   median   Mean

mean > median > mode

Negative (or) left skewed data:



mean   median   mode

mean < median < mode

Outliers:- *Datapoint which is away from general pattern.
          *Odd one out.

mean $\begin{cases} 20 & 21 & 22 & 23 & 24 & 25 & 27 & \boxed{\mu = 23} \\ 20 & 21 & 22 & 23 & 24 & 25 & 70 & \boxed{\mu = 25} \\ & & & & & & \uparrow & \hookrightarrow \text{outlier} \end{cases}$

median $\begin{cases} 20 & 21 & 22 & (23) & 24 & 25 & 27 \\ 20 & 21 & 22 & (23) & 24 & 25 & 70 \end{cases}$
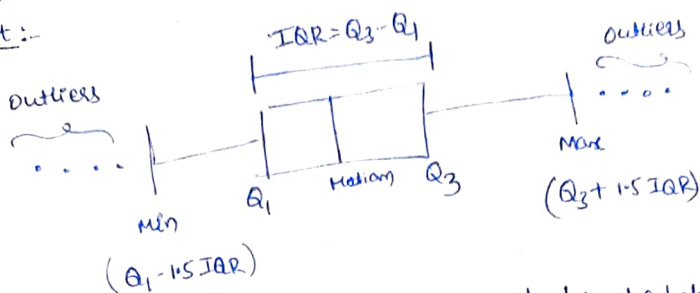
median = 23 for both data.

* It is suggestable to take median, when data have outliers.

outlier detection :- is the process of detecting data points which are not as per the general pattern in data.

Two ways: Boxplot & Z-score

Boxplot :-



IQR = $Q_3 - Q_1$

outliers

Max

$(Q_3 + 1.5 \, IQR)$

Median $Q_3$

outliers

Min $Q_1$

$(Q_1 - 1.5 \, IQR)$

* The outliers are the points that are present beyond & below the upper & lower whiskers.

Z-score :- Z score tells how many standard deviations (sd) away a data point is from the mean standard score.
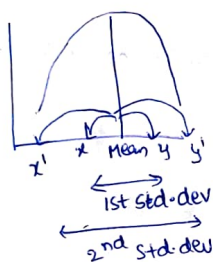
$$Z\text{-Score} = \frac{X - \mu}{\sigma}$$

$X$ ← (obs)
$\mu$ → (mean)
$\sigma$ → (std dev)

Ex:
$\mu = 30$
$X = 45$
$\sigma = 5$

$Z = \frac{45-30}{5} = 3$

The point $x$ is 3 std devs away from mean



$x'$ * Mean $y$ $y'$
← Ist std-dev →
← 2nd std dev →

$x = \mu - SD$
$y = \mu + SD$

$x' = \mu - 2SD$
$y' = \mu + 2SD$

| 3 std dev → 99% of data |
| 2 std dev → 95% of data |
| 1 std dev → 65 - 60% of data will be covered |

Ex 2:

| 20 | 22 | 25 | 26 | 28 | 59 | → outlier

$\mu = 30$
$\sigma = 13$

| $x$ | $x_i - \mu$ | $(x_i - \mu)^2$ | z score |
|---|---|---|---|
| 20 | -10 | 100 | -0.76 |
| 22 | -8 | 64 | -0.61 |
| 25 | -5 | 25 | -0.38 |
| 26 | -4 | 16 | -0.31 → values falls in 1st SD |
| 28 | -2 | 4 | -0.15 |
| 59 | 29 | 841 | 2.23 → falls in 3rd SD |

more compared to other values.

**Covariance** : measure of linear association blw two features:

numeric

proportionality

$$Cov(x,y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

$(x_i - \bar{x})$ : diff of $i^{th}$ obs w.r.t mean $(\bar{x})$

Ex: hours: $\{1, 2, 3, 4, 5\}$

marks: $\{20, 40, 60, 80, 100\}$

$\bar{x} = 3 \qquad \bar{y} = 60$

$$Cov(x,y) \Rightarrow \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

| $x$ | $y$ | $x_i - \bar{x}$ | $y_i - \bar{y}$ | $(x_i - \bar{x})(y_i - \bar{y})$ |
|---|---|---|---|---|
| 1 | 20 | -2 | -40 | 80 |
| 2 | 40 | -1 | -20 | 20 |
| 3 | 60 | 0 | 0 | 0 |
| 4 | 80 | 1 | 20 | 20 |
| 5 | 100 | 2 | 40 | 80 |
| | | | | 200 |

$$Cov(x,y) = \frac{200}{5}$$

$$\boxed{Cov(x,y) = 40}$$

+ve Cov → the two features have +ve relationship.

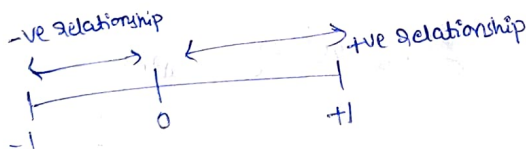-ve Cov → the " " will " -ve relationship.

**Correlation:** Standardized measure of strength and direction of the linear relationship blw two variables.

Range $\Rightarrow \{-1 \text{ to } +1\}$

$Cor(x,y) = +1$, very high +ve corr

$Cor(x,y) = -1$, very high -ve corr

$$Correlation = \frac{Cov(x,y)}{\sigma x * \sigma y}$$

-ve relationship $\qquad$ +ve relationship



Pick above example:-

| $(x_i - \bar{x})^2$ | $(y_i - \bar{y})^2$ |
|---|---|
| 4 | 1600 |
| 1 | 400 |
| 0 | 0 |
| 1 | 400 |
| 4 | 1600 |
| 10 | 4000 |

$\sigma x = \frac{10}{5} = \sqrt{2}$

$\sigma y = \frac{4000}{5} = \sqrt{800}$

$correlation = \frac{40}{\sqrt{800} \sqrt{2}} = \frac{40}{\sqrt{800 \times \sqrt{2} \times \sqrt{2}}}$

$= 0.025 \frac{40}{40}$

$= 1$

∴ features are highly correlated.