

* Measures of Dispersion

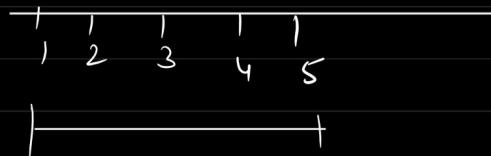
$$S_1 = 1, 2, \underline{3}, 4, 5$$

mean/median = 3

$$S_2 = 3, 3, 3, 3, 3$$

mean/median = 3

S_1



S_2



* How the data is spread?

→ Range

→ Percentage and percentile

→ Quartiles (Boxplot)

→ Variance

→ Standard deviation.

* Range — difference between maximum and minimum value.

$$\{1, 2, 3, 4, 5\}$$

$$\text{Range} = 5 - 1 = 4$$

$$\{1, 2, 3, 4, 1000\}$$

$$\text{Range} = 1000 - 1 = 999$$

* Outlier affects the range.

* Percentage and percentiles

$$1, 2, 3, 4, 5$$

What is the percentage of nos that are odd?

$$\frac{3}{5} \times 100 = 60\%$$

* Percentile

Defⁿ - A percentile is a value below which a certain percentage of observations lie.

$$\{1, 2, 3, 4, 4, 6, 7, 7, 8, 10\}$$

What is the percentile rank of 3?

$$\text{Percentile rank of a no} = \frac{\text{No. of values below that no}}{\text{Total nos (n)}} \times 100$$

$$= \frac{2}{10} \times 100 = 20^{\text{th}} \text{ percentile}$$

What value exists at 75th percentile?

$$\text{Value} = \frac{\text{Percentile}}{100} \times (n+1)$$

$$= \frac{75}{100} \times (10+1)$$

$$= \frac{3}{4} \times 11 = 8.25^{\text{th}} \text{ number}$$

\downarrow
8th number

(8.5th)
 \downarrow

75th percentile — 7.

avg of
8th & 9th no

→ 75th percentile is 8

It means 75% of the no. in the data is equals to or below 8.

* Quartile

→ Quartiles are values that divides a list of numbers into quarters.

* Put the no in order

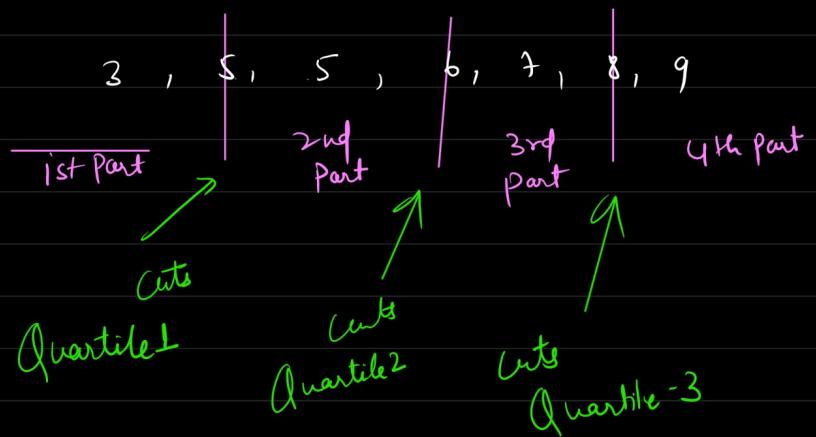
* then cut the number into 4 equal parts

* The quartiles are at the cut.

Ex. 6, 8, 5, 5, 7, 3, 9

order - 3, 5, 5, 6, 7, 8, 9

Cut the no into quarters.



$$Q_1 \rightarrow 5$$

$$Q_2 \rightarrow 6$$

$$Q_3 \rightarrow 8$$

Ex. - 1, 1, 1, 1, 2, 2, 2, 3, 3, 3, 4

total - 11 nos. \uparrow \uparrow \uparrow
 Q_1 Q_2 Q_3

if total no is odd

$$\checkmark Q_1 = \frac{n+1}{4}^{\text{th}} = \frac{11+1}{4}^{\text{th}} = \frac{12}{4}^{\text{th}} = 3^{\text{rd}} \text{ no.}$$

$$\checkmark Q_3 = \frac{3(n+1)}{4}^{\text{th}} = \frac{3(11+1)}{4}^{\text{th}} = \frac{3 \times 12}{4}^{\text{th}} = 9^{\text{th}} \text{ no.}$$

$$\checkmark Q_2 = \left(\frac{n+1}{2}\right)^{\text{th}} = \frac{11+1}{2}^{\text{th}} = 6^{\text{th}} \text{ no.}$$

(median)

if total no is even

$$Q_1 = \frac{n}{4}^{\text{th}} \text{ no}$$

$$Q_3 = \frac{3n}{4}^{\text{th}} \text{ no}$$

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

1, 1, 1, 1, 2, 2, 2, 3, 3, 3, 4

$$\longleftrightarrow Q_2 \longleftrightarrow$$

Median - the point that divides the data into two equal parts

* even no. of nos

$$1, \underline{\underline{2}}, 3, 3, 4, 4 \quad N = 6$$

$$Q_1 - \frac{6^{\text{th}}}{4} = \frac{3}{2} \Rightarrow 1.5^{\text{th}} \Rightarrow \frac{1+2}{2} = \frac{3}{2} = \underline{\underline{1.5}} \quad (Q_1)$$

Q_2 (median when the total nos are even)

$$Q_3 - \frac{3 \times 6^{\text{th}}}{4} = \frac{3 \times 6}{4} = \frac{18}{4} = 4.5^{\text{th}}$$

4th & 5th no

$$Q_3 = \frac{3+4}{2} = 3.5$$

$$\text{avg} \left(\frac{n}{2}^{\text{th}} \text{ no}, \frac{n}{2}^{\text{th}} + 1 \right)$$

$$\text{avg} \left(\frac{6}{2}^{\text{th}}, \frac{6}{2}^{\text{th}} + 1 \right)$$

$$\text{avg} \left(3^{\text{rd}}, 4^{\text{th}} \right)$$

$$= \frac{3+3}{2} \Rightarrow \underline{\underline{3}} \quad Q_2 = 3$$

Five point summary Q_0 (min) - L \longrightarrow 0% - 0th percentile

Q_1 - | \longrightarrow 25% - 25th percentile

transaction amount (median) Q_2 - 2 \longrightarrow 50% - 50th percentile

1000
2000
3000
-.
-.
-.
 $\left\{ \begin{array}{l} \min = 1000 \\ (Q_0) \\ Q_1 = 25^{\text{th}} \text{ percentile} \end{array} \right.$

Q_3 - 3 \longrightarrow 75% - 75th percentile

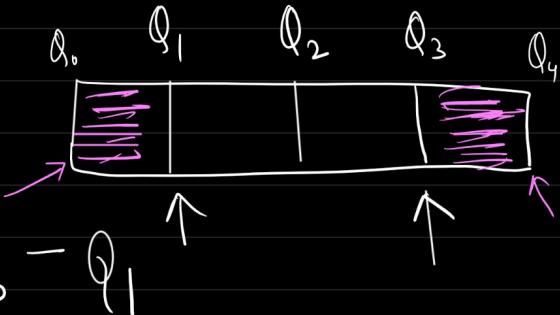
Q_4 (maximum) - 4 \longrightarrow 100% - 100% Percentile

$\hookrightarrow 5000 \rightarrow 25\% \text{ of transaction amount is equal to or below } 5000 \text{ in the data}$

$Q_2 = \underline{\underline{10000}} \rightarrow 50\% \text{ of transaction is equals to or below } 10000$

Q_3 - 75th percentile
 Q_4 - max no.

Inter Quartile range



2, 4, 4, 5, 6, 7, 8
↑ ↑
0 1

$$IQR = 7 - 4 = \underline{\underline{3}}$$

$$\{2, 3, 3, 3, 3, 3, 4, 4, 5, 5, 5, 5, 6, 6, 6, 6, 7, 8, 9, 9\}$$

$$Q_1 = 25^{\text{th}} \text{ percentile} = \frac{25}{100} \times 16 = \frac{1}{4} \times 16 = 4^{\text{th}} \text{ no}$$

$$Q_3 = 75^{\text{th}} \text{ percentile} = \frac{75}{100} \times 16 = \frac{3}{4} \times 16 = 12^{\text{th}} \text{ no.}$$

Outliers are extreme values

$$IQR = 6 - 3 = 3$$

Lower fence

$$\text{Upper fence} \quad \left(Q_3 \right) - \left(Q_1 \right)$$

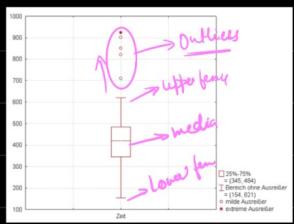
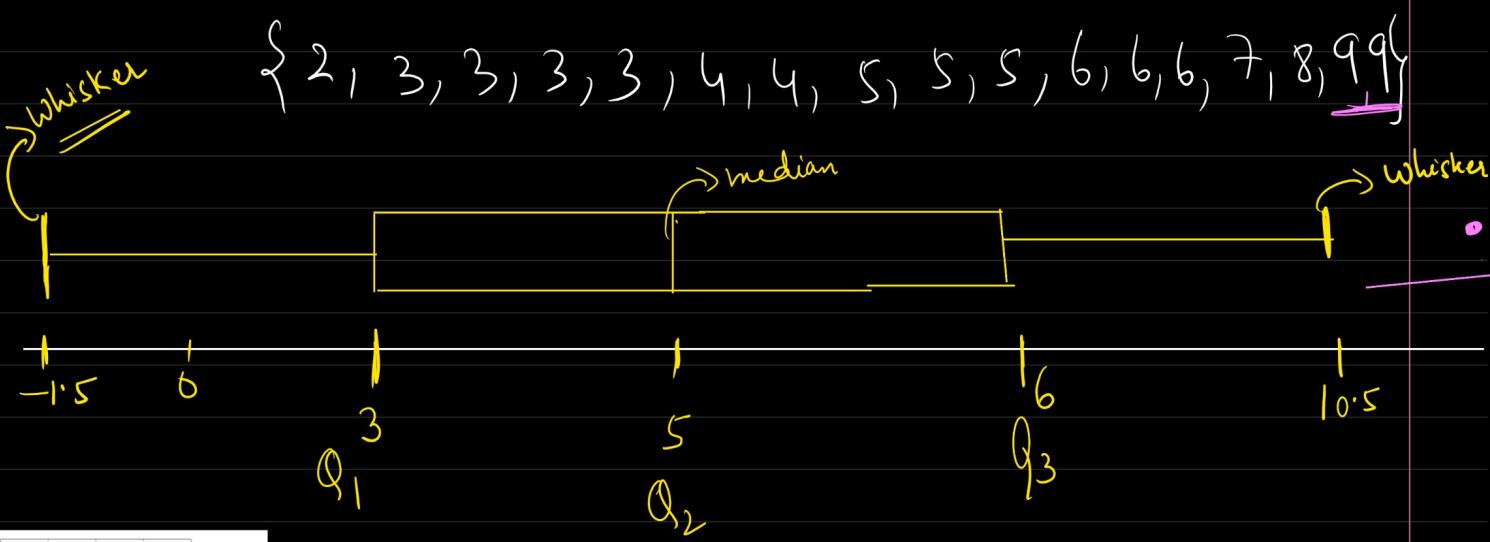
$$\text{Lower fence} = Q_1 - 1.5 \times IQR$$

$$\text{Upper fence} = Q_3 + 1.5 \times IQR$$

$$L \cdot F = 3 - 1.5 \times 3 = -1.5$$

Box - whisker plot

$$U.F = 6 + 1.5 \times 3 = 4.5 + 6 = 10.5$$



Measures of spread / dispersion

① Variance

② Standard deviation

* Mean-deviation

$$\begin{array}{c}
 \text{0} \\
 | \quad \text{1 unit} \quad | \quad \text{1 unit} \quad | \quad \text{2 units} \\
 | \quad \text{1,} \quad \text{2,} \quad \text{3,} \quad \text{4,} \quad \text{5} \\
 \hline
 | \quad \text{2 units} \\
 \uparrow \\
 \text{mean}
 \end{array}
 \Rightarrow \frac{2+1+0+1+2}{5} = \frac{6}{5} = 1.2$$

→ On an avg each of the data is 1.2 units away from mean value

* Variance — The average of the squared differences from the mean.

Population Variance

$$\sigma^2 = \frac{N}{\sum_{i=1}^N (x_i - \mu)^2}$$

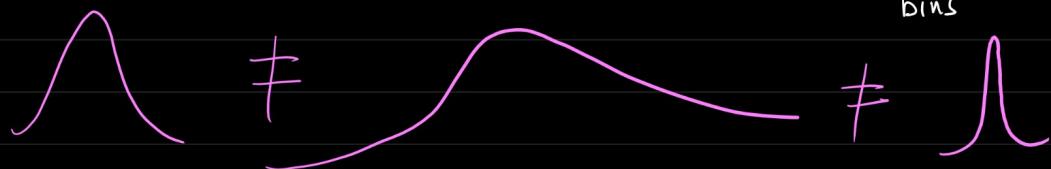
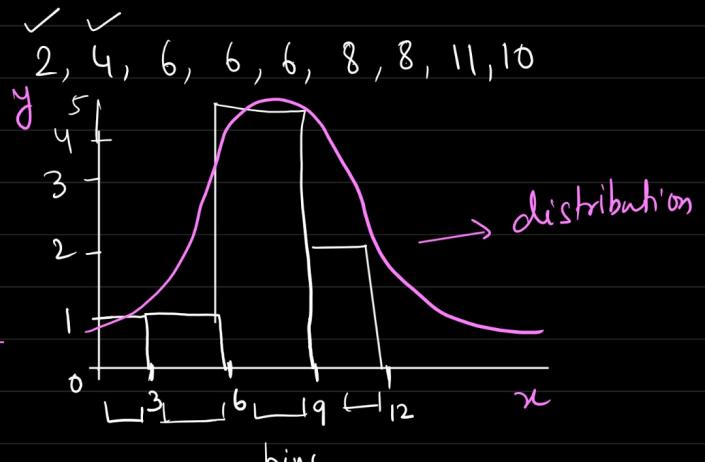
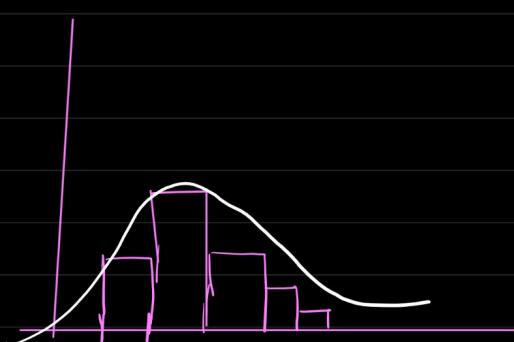
population mean

Sample Variance

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

$\boxed{n-1}$

sample mean



$$\text{data} = \{1, 2, 3, 3, 4, 4\}$$

How to Calculate Variance

- Calculate mean
- for each no in data, Subtract the mean and the no
- Square of difference
- Calculate the avg of square of difference

x	\bar{x}	$x - \bar{x}$	$(x - \bar{x})^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.68
3	2.83	0.17	0.03
3	2.83	0.17	0.03
4	2.83	1.17	1.37
4	2.83	1.17	1.37
<u>2.83</u>			<u>6.82</u>

$$s^2 = \frac{6.82}{n-1} = \frac{6.82}{5} = \underline{\underline{1.37}}$$

Variance \uparrow Spread \uparrow



* Standard deviation

Standard deviation is a measure of how spread out numbers are.

↳ Square root of Variance
 $s = \sqrt{\text{Var}}$
 $s = \sqrt{1.37} = 1.17$

Std dev of population $\Rightarrow \sigma = \sqrt{\text{Var}_p}$

Std dev of sample $\Rightarrow s = \sqrt{\text{Var}_s}$

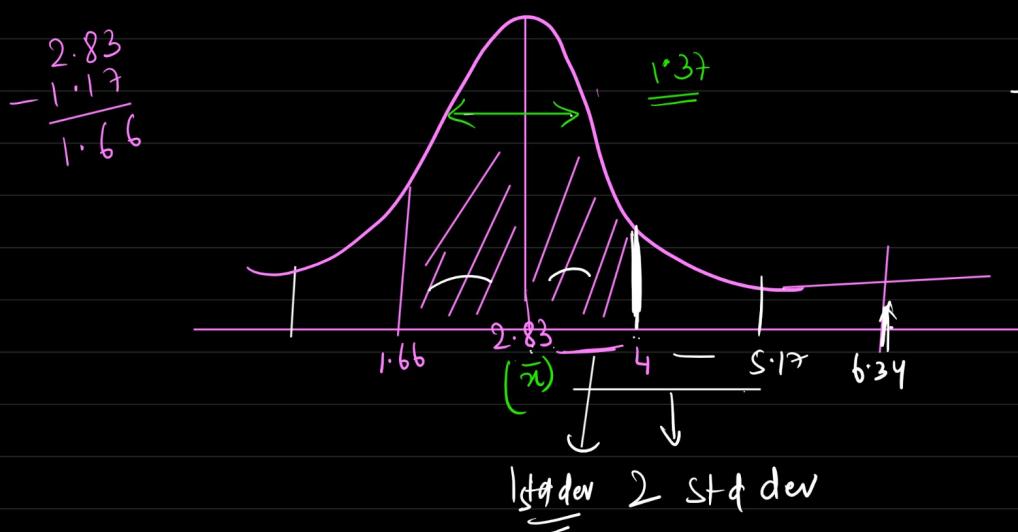
✓ 1, 2, 3, 3, 4, 4

$\uparrow \uparrow \uparrow$

$$\text{Var} = 1.37 \quad \checkmark$$

$$\text{Std dev} = \sqrt{1.37}$$

→ Standard way of knowing where your data lies



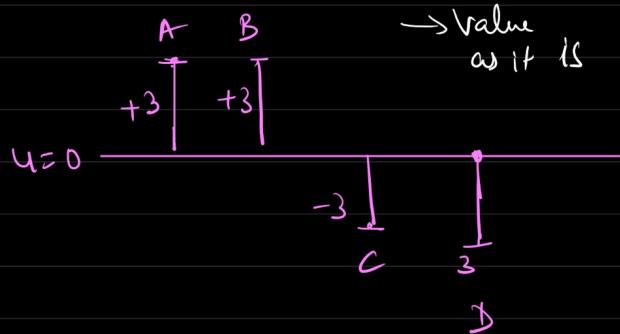
$$\begin{array}{r}
 2.83 \\
 +1.17 \\
 \hline
 4.00 \\
 +1.17 \\
 \hline
 5.17 \\
 +1.17 \\
 \hline
 6.34
 \end{array}$$

68 - 1 std
95% - 2 std
99.7% - 3 std.

* Variance →

$$\text{Var}_p = \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

why square?



$X - M$

$$\frac{+3 + 3 + (-3) + (-3)}{4} = 0$$

+ve and -ve
are negating
each other

absolute value

$$\rightarrow \frac{|3| + |3| + |-3| + |-3|}{4} = \frac{12}{4} = 3$$

mean deviation



$$\rightarrow \frac{|+8| + |+1| + |-2| + |-1|}{4} = \frac{12}{4} = 3$$



$$\sqrt{\frac{3^2 + 3^2 + (-3)^2 + (-3)^2}{4}} = \sqrt{\frac{36}{4}} = 3$$

$$\sqrt{\frac{8^2 + 1^2 + (-2)^2 + (-1)^2}{4}} = \sqrt{\frac{64 + 1 + 4 + 1}{4}} = \sqrt{\frac{70}{4}} = 4.184$$

* Variance_{sample} = $\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$ → why?

Bessel correction

We use $n-1$ rather than n is because sample variance will be unbiased

Estimator



We are estimating variance of population using variance of sample



✓ $(x - \bar{x})^2$

Population $\frac{(x - M)^2}{n}$

$(x - M)^2 > (x - \bar{x})^2$ = var.

$n \rightsquigarrow n-1$

$\textcircled{2} = \textcircled{1} - \frac{8}{4} = 2$

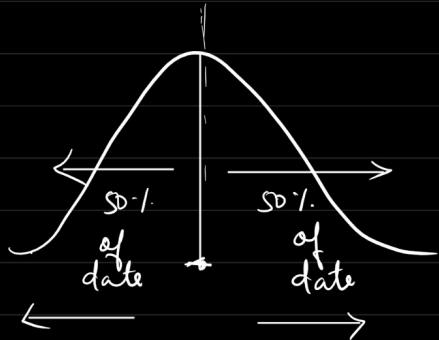
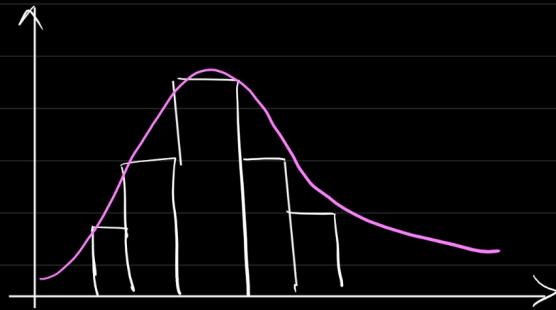
In most of the cases

$\Rightarrow S^2 = \frac{\sum (x - \bar{x})^2}{n-1} \rightarrow \text{smaller}$

You are reducing numerator

Measures of Symmetry

→ If anything is exactly towards the left and right.

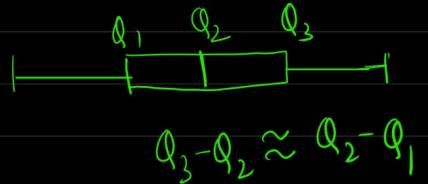
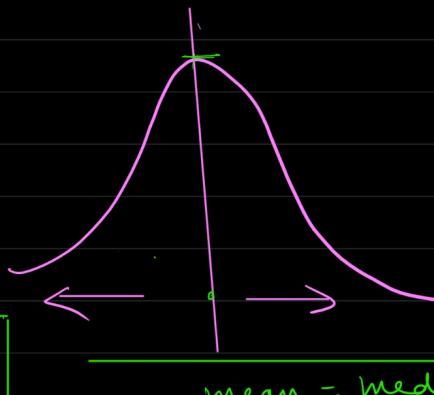


* Skewness → measure of data symmetry

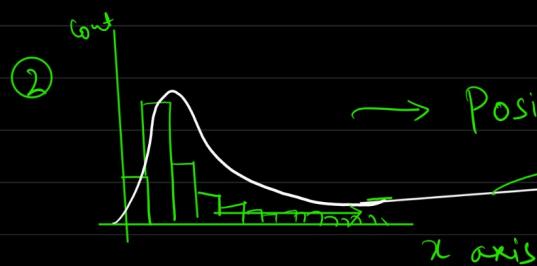
① No Skewness

$$\text{Skewness} = 0$$

$$\boxed{\text{Skewness} = \frac{\bar{x} - \mu}{n \sigma^3}}$$



1, 2, 3, 4, 5

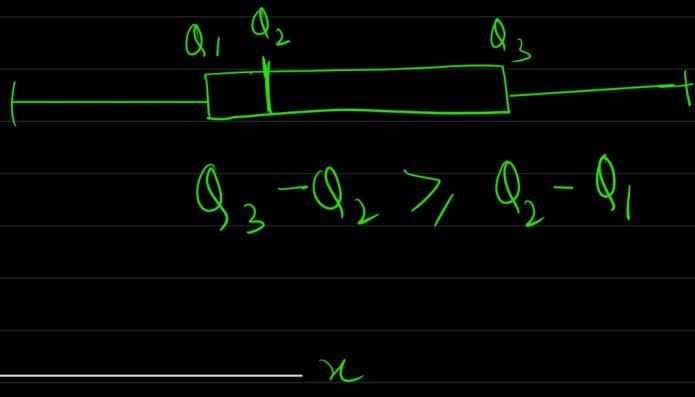
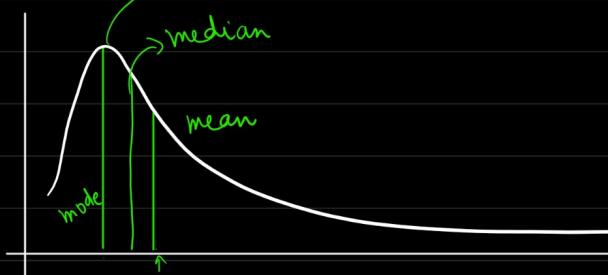


→ Positive Skewed

tail is on right side of distribution

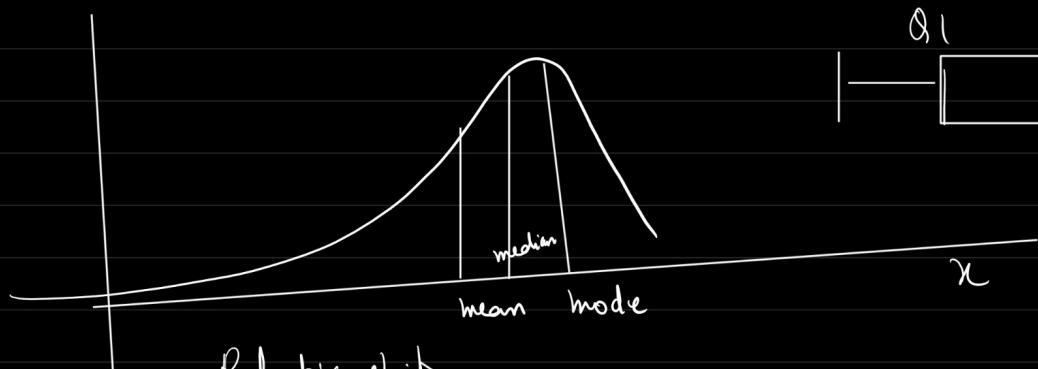


$$Q_3 - Q_2 > Q_2 - Q_1$$



mean > median > mode

③ Left skewed distribution | neg skewed distribution



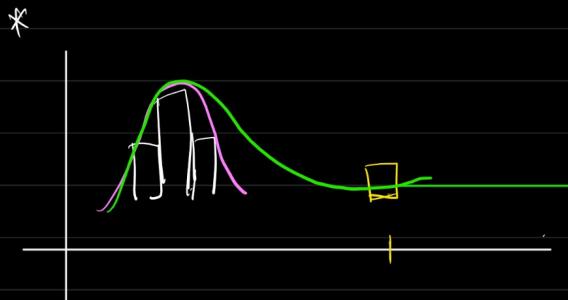
$$Q_2 - Q_1 > Q_3 - Q_2$$

Relationship

mode > median > mean

Transformations

- Log transform
- Box-cox transformation
- Exponential tran
- reciprocal "
- Outlier treatment



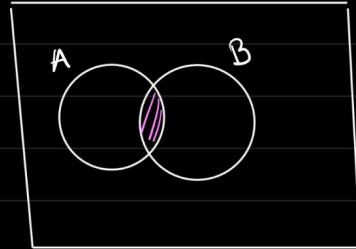
Set

$$A = \{1, 2, 4, 5, 7\}$$

$$B = \{2, 4, 5\}$$

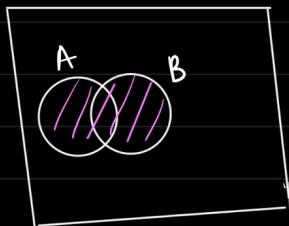
① Intersection (common elements)

$$A \cap B = \{2, 4, 5\}$$



② Union (all distinct elements from both the sample)

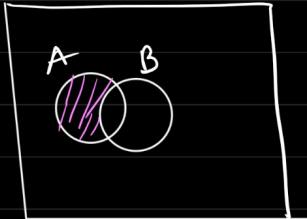
$$A \cup B = \{1, 2, 4, 5, 7\}$$



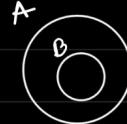
③ Difference (items which are present only in first set)

$$A - B = \{1, 7\}$$

$$A - B$$



④ Subset (all the element of B is present in A, then we say B is subset of A)



$B \rightarrow A$ — True

(B is subset of A)

$\times A - B$ — False

⑤ Superset (A is containing all elements of B, A is Superset of B)

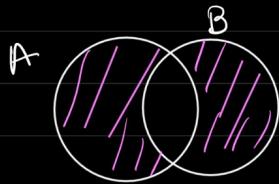
$A \rightarrow B \Rightarrow$ True

$B \not\rightarrow A \Rightarrow$ False

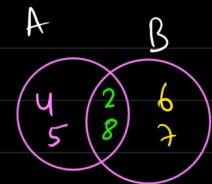
⑥ Symmetric difference (Opposite of intersection)

The elements that are distinct in both sets excluding intersection

$$A \Delta B = \{1, 7\}$$



$$\begin{aligned} A &= \{4, 5, \underline{2}, 8\} \\ B &= \{\underline{2}, 6, 7\} \end{aligned}$$



$$A \Delta B = \{4, 5, 6, 7\}$$

* Covariance and correlation

	(X)	(Y)
	transaction amount	transaction count
Y↑ X↑	-	-
Y↑ X↓	-	-
Y↓ X↑	-	-
Y↓ X↓	-	-

Understanding the relationship



* Example (direct)

Predict price of house based on area of house

Area of house (sqft)	Price	AH ↑ P ↑
1100	80	AH ↓ P ↓
1200	85	
→	=	
=	=	

Area is sqrt

* Example (indirect)



① Covariance

$$\checkmark \text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

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

Variance was spread of data \Rightarrow relationship of a feature with itself

Cov mean, you are trying to understand the relationship of a feature with respect to other features.

\rightarrow Covariance \rightarrow relationship b/w two variable.

\rightarrow

$$\begin{array}{ll} X \uparrow Y \uparrow & \text{or} \\ X \downarrow Y \downarrow & X \uparrow Y \downarrow \\ (+ve) \text{ Cov} & X \downarrow Y \uparrow \\ (-ve) \text{ Cov} & \end{array}$$

$$\begin{array}{cc} X & Y \\ 2 & 3 \\ 3 & 5 \\ 6 & 6 \\ 1 & 8 \end{array}$$

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$(2-3)(3-5.5) + (3-3)(5-5.5) + (6-3)(6-5.5) + (1-3)(8-5.5)$$

$$\Rightarrow \frac{(-1)(-2.5) + 0 + 3 \times 0.5 + (-2) \times 2.5}{3} = \frac{-1 + 0 + 1.5 - 5}{3} = \frac{-4.5}{3} = -1.5$$

$$\Rightarrow \frac{2.5 + 0 + 1.5 - 5}{3} = \frac{-1}{3} = -0.33$$

\rightarrow The two features X and Y are negatively related.

$$\begin{array}{cc}
 X & Y \\
 2 & 3 \\
 4 & 5 \\
 6 & 7
 \end{array}
 \quad \text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1} \\
 = \frac{(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)}{3-1} \\
 \bar{x} = 4 \quad \bar{y} = 5 \\
 \Rightarrow \frac{4+0+4}{2} = 4 = \text{(+ve)}$$

$\rightarrow x$ & y are having a positive relation.

Advantage

\rightarrow Relationship b/w X, Y

+ve or -ve

* Disadvantage

$$\textcircled{1} \quad \begin{array}{cc} \overline{X} & \overline{Y} \\ A & B \end{array}$$

$$\text{Cov}(x, y) = 50$$

$$\text{Cov}(A, B) = 100$$

\downarrow
 \rightarrow No comparison of strength of relationship in Covariance

\rightarrow No any standardized scale to interpret the strength.

② Covariance has dimension

$$\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$x \rightarrow$ height \rightarrow ft

$y \rightarrow$ wt \rightarrow kg.

$\text{Cov.} \Rightarrow \text{ft. kg}$

(x)

(y)

(z)

transaction Amount (Rs)	height (ft)	weight (kg)
-------------------------------	----------------	----------------

$$\text{Cov}(\text{trAmt}, \text{ht}) \Rightarrow \text{Rs} \cdot \text{ft} \Rightarrow 450 \text{ Re.ft}$$

$$\text{Cov}(\text{height}, \text{wt}) \Rightarrow \text{ft} \cdot \text{kg} \Rightarrow 600 \text{ ft.kg}$$

→ We can not compare two different dimension

$$X \quad Y \quad Z$$

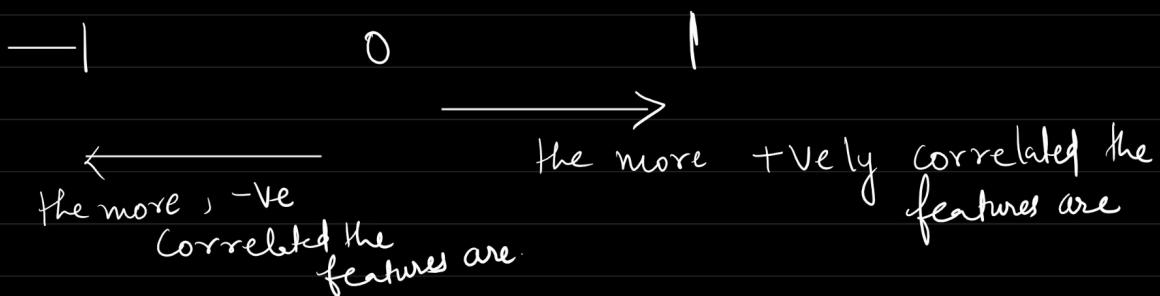
$$\text{Cov}(x,y) \quad \text{Cov}(y,z)$$

Not comparable → different dimensions.

- * Soln
 - -1 to 1
 - dimensionless quantity

② Pearson Correlation Coefficient [-1 to 1]

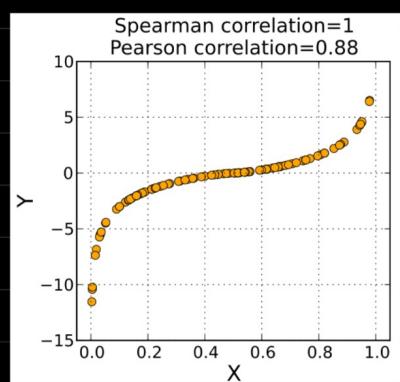
$$f_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y} = [-1 \text{ to } 1]$$



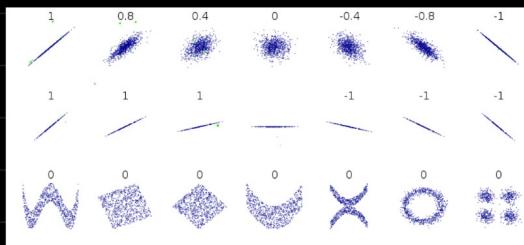
$$f_{x,y} = 0.4 \quad , \quad f_{A_1 B} = 0.8$$

→ feature $A_1 B$ is highly correlated as compared to x, y .

→ Pearson Correlation coefficient always measures the linear relationship



Pearson correlation
↓
 $x \uparrow y \uparrow$ Linear relationship



What to do for Non linear relationship

Spearman Rank Correlation

$$\gamma_s = \frac{\text{Cov}(R(x), R(y))}{\sigma_{R(x)} * \sigma_{R(y)}}$$

$R(x)$ — Rank of x
 $R(y)$ — Rank of y

x	y	$R(x)$	$R(y)$
5	6	3	1
7	4	2	2
8	3	1	3
1	1	5	5
2	2	4	4

5 4 3 2 1 ↑ 1st

$x \rightarrow 5, 7, 8, 1, 2 \rightarrow$ Sort the value $\rightarrow 1, 2, 5, 7, 8$

Sort the no \rightarrow Highest no will be rank 1.

Dataset \rightarrow 1000 feature

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6 \dots x_{1000}$ Price

y
↑
off

$$\left\{ \begin{array}{l} x_1 - y \\ x_2 - y \\ x_3 - y \\ x_4 - y \end{array} \right. \quad \text{Corr} \approx 0$$

* Random Variable

↳ A set of Possible values from a random experiment.

* Tossing a coin → Experiment is random,

↳ H, T
 ↓
 Outcomes will be random.

Quantify these random values.

$$X = \begin{cases} 0 & \text{--- tail (T)} \\ 1 & \text{--- Head (H)} \end{cases}$$

$$X = \underline{\{0, 1\}}$$

↑
Random variable
 ↓

If can
 take any
 value from
 the set of values

* We have an experiment
 (tossing a coin)

* quantify each event

$$\{0, 1\}$$

* This value is Random variable.

$$\boxed{x + 5 = 10 \\ x = 10 - 5 = x = 5}$$

In Algebra, a variable value is fixed

$$X = \{0, 1, 2, 3\}$$

X could be 0, 1, 2, or 3 randomly

$$X = \{1, 2, 3, 4, 5, 6\} \xrightarrow{\text{sample space}}$$

$$P(X=1) = 1/6$$

$$(\text{Tossing a coin}) P(X=1) = 1/2 \longrightarrow \left\{ \begin{matrix} \text{tail} \rightarrow \text{Head} \\ \{0, 1\} \end{matrix} \right.$$



$$\underline{P(X \leq 4)}$$