

Descriptive Analysis / Statistics

- Descriptive statistics summarize & describe main features of a dataset.
- Provides simple summaries about the data and measures.
- Measures → mean
median
mode
variance
std deviation

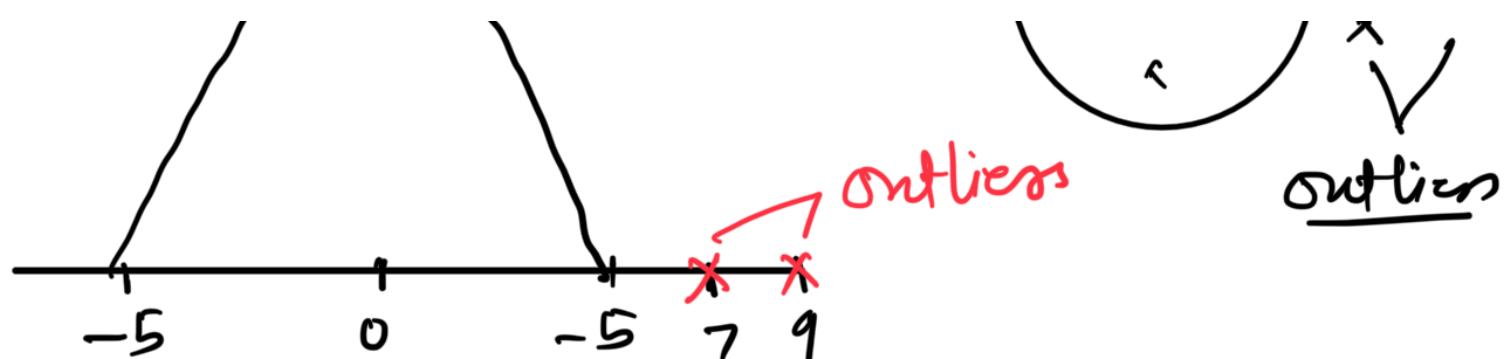
1. Mean : Average is the sum of all values divided by the number of values. It provides a measure of the central tendency of the data.

- Use it to normalize data
- helpful in summarizing the data

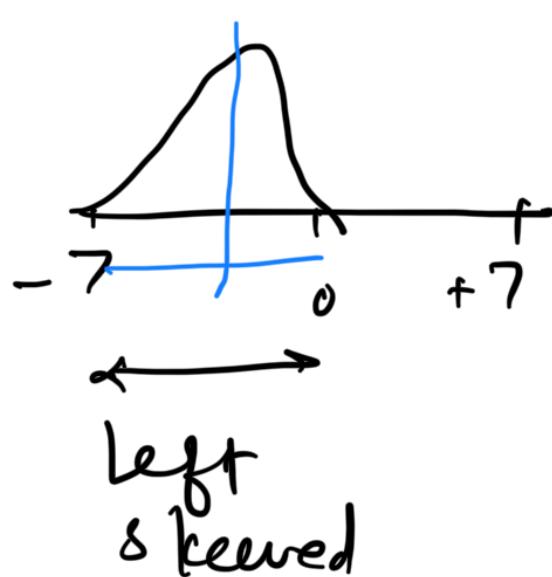
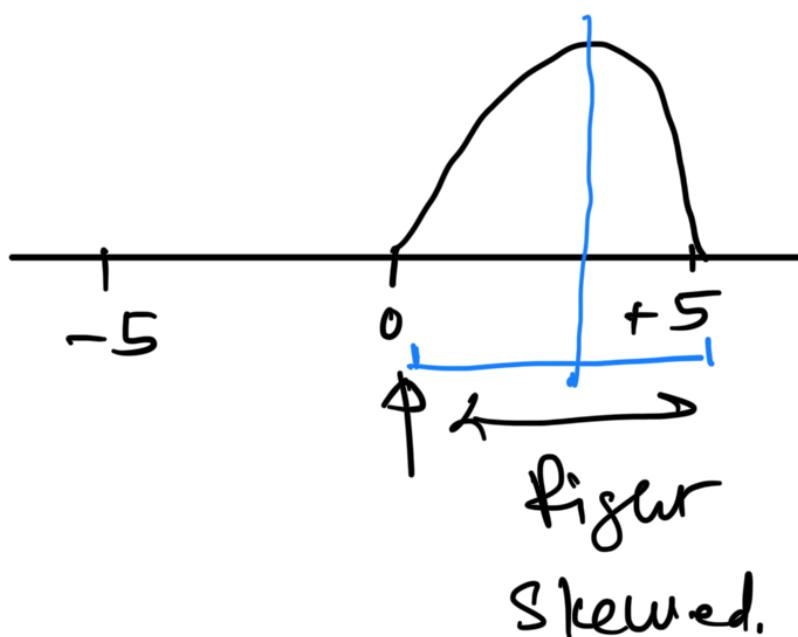
2. Median : The median is the

middle value of a sorted dataset





Correct \rightarrow Normal distributed
 $\{\pm 5\}$



- \rightarrow Used to handle skewed data
- \rightarrow Useful in descriptive statistics to understand data distribution.

3. Mode — The mode is the most frequently occurring value in a dataset.

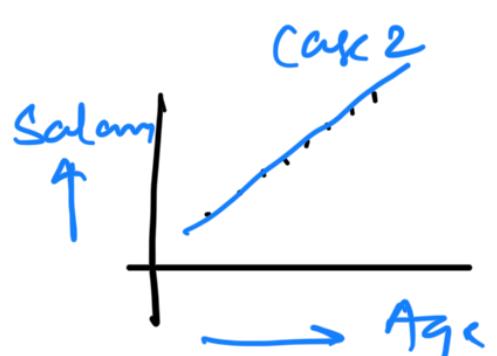
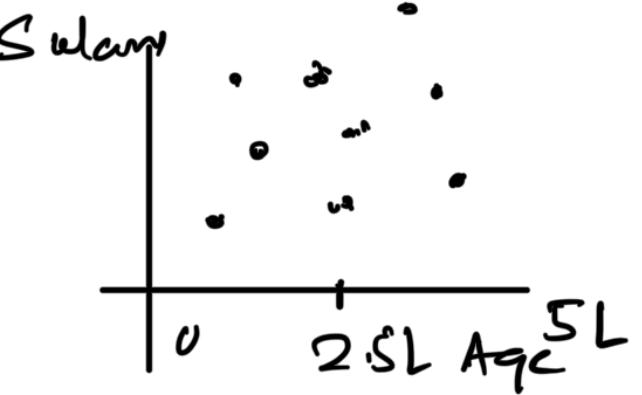
- Use to handle Categorical (text) data
- helpful to identify the most common value in the dataset

4. Variance -

Variance measures the dispersion of the data points from the mean

Case 1 - All profession

Case 2 - IT profession



→ Indicates the variability in data

→ Used to understand the spread of features

5. Standard Deviation - It is the

Square root of the variance, provide the measure of the average distance from the mean.

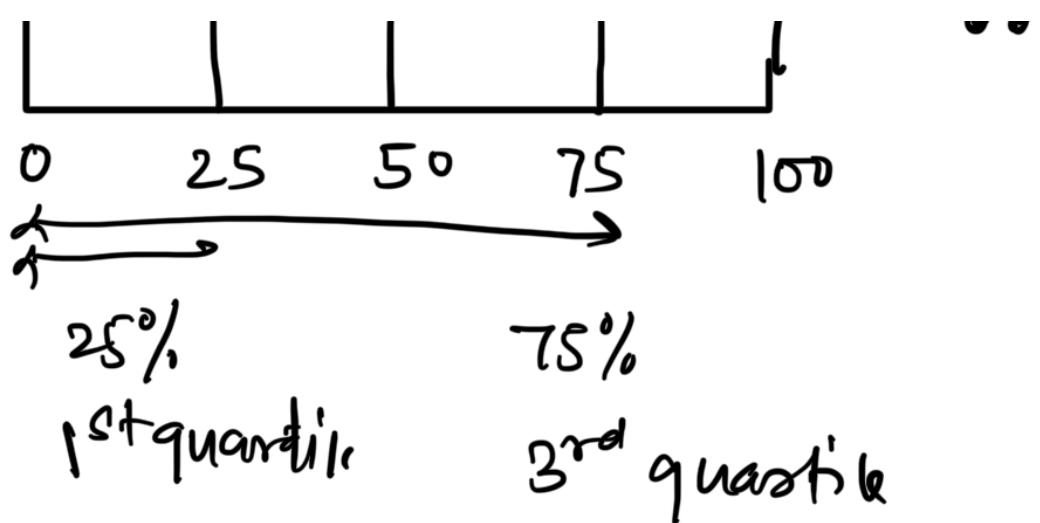
— Use it to scale the data

— help us to understand the spread of dataset.

6. Range and Interquartile Range (IQR)

2nd quartile 50%





- The range is the difference b/w max and min values.
- The IQR is the difference b/w the 75th and 25th quartile (percentile)
- helps to measure the spread of the middle 50% of the data

Use

- Identify outliers
- Understanding about distribution

$$\text{Variance } (\sigma^2) = \sigma^2 = \frac{\sum (x - \bar{x})^2}{n}$$

$$\text{std deviation } (\sigma) = \sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

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

$$\text{Coefficient of Variance} = \left(\frac{\sigma}{\bar{x}} \right) \times 100$$

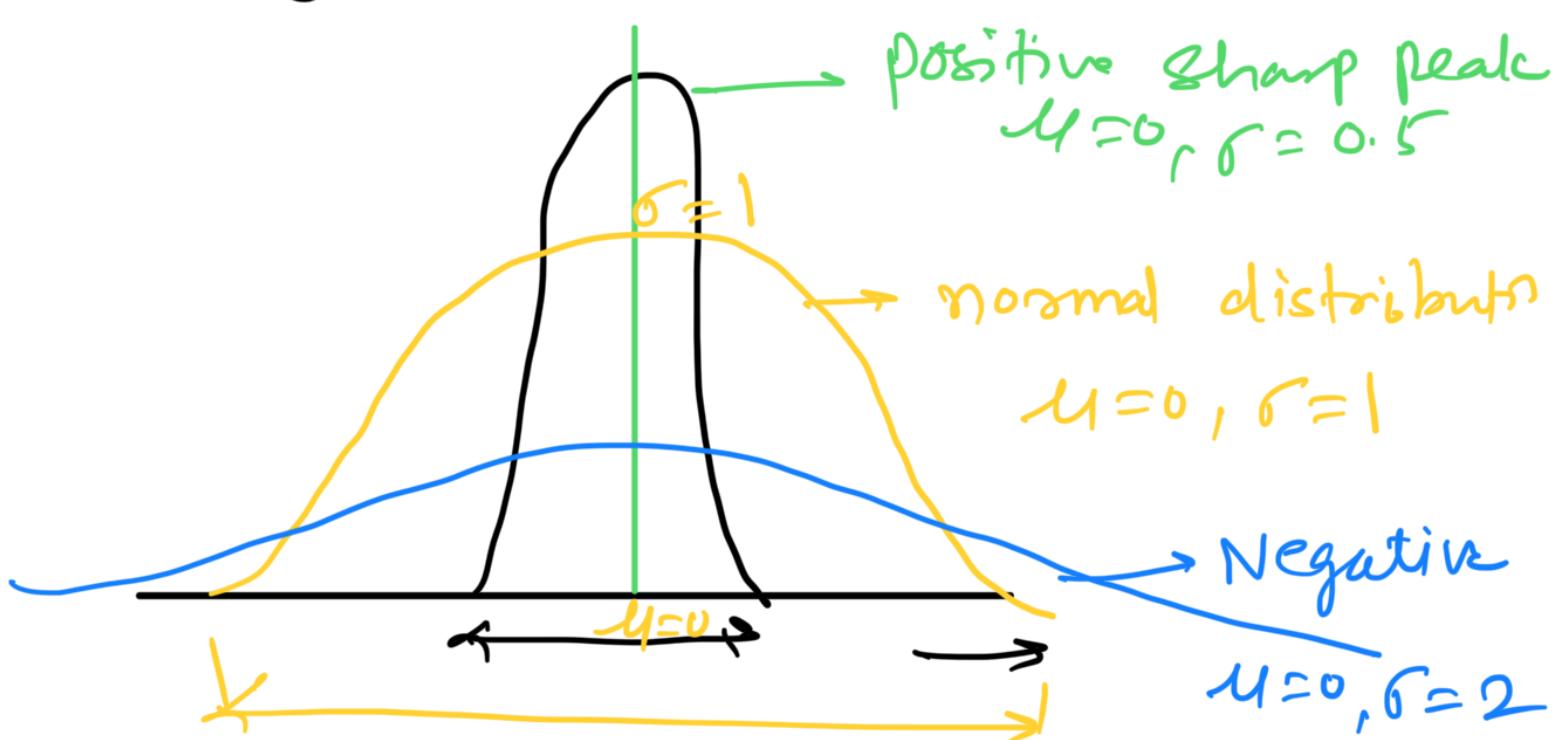
Measures of Dispersion

Range, σ , σ^2 , IQR, CV

Measures of Shape

Kurtosis - Measure the 'tailedness'

of the distribution



- Kurtosis talks @ the central peak and the dispersion of data set
- Types of kurtosis -
 1. Positive → Sharp peaks and

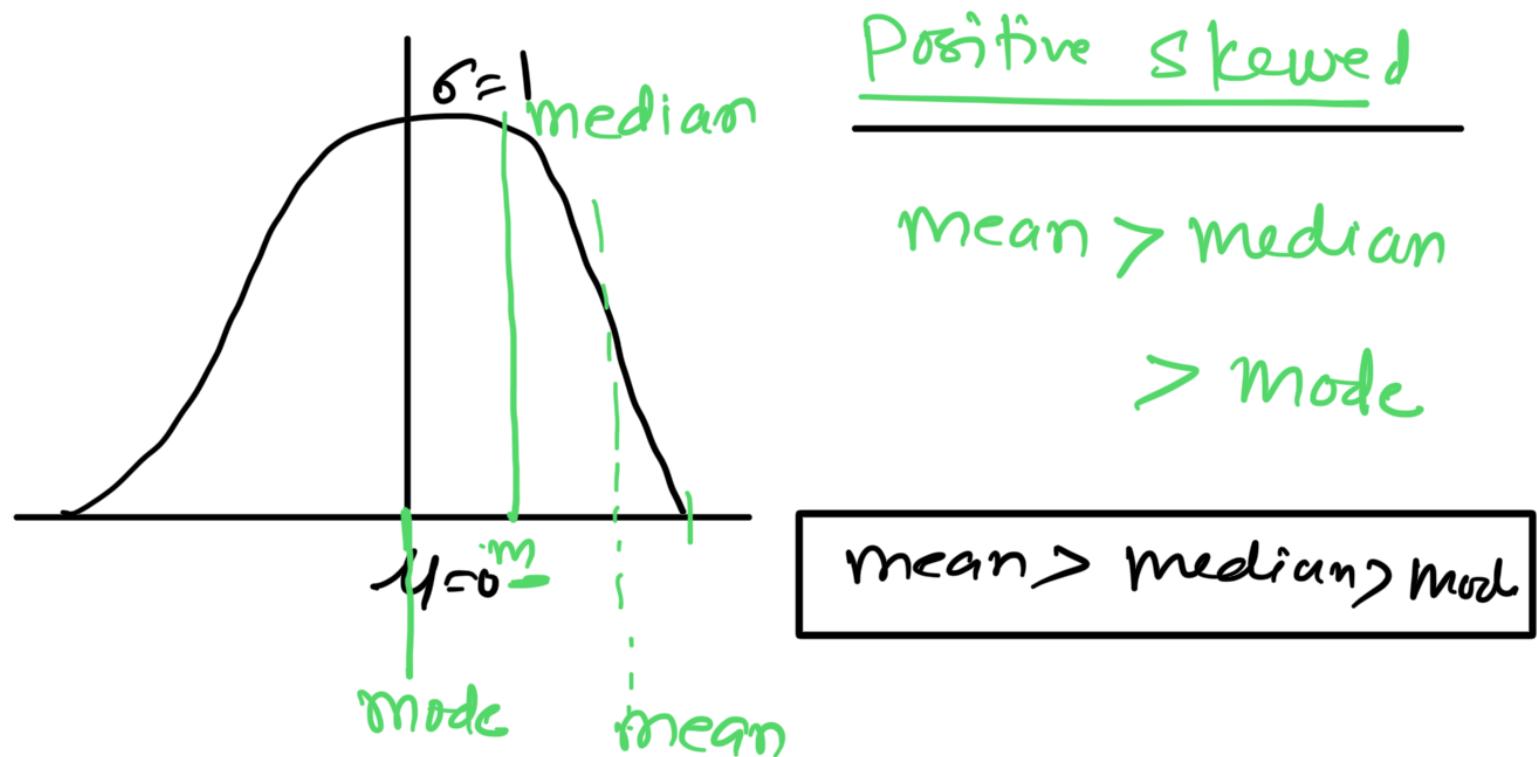
Lighter dispersion

2. Negative \rightarrow wide peaks and thicker dispersion

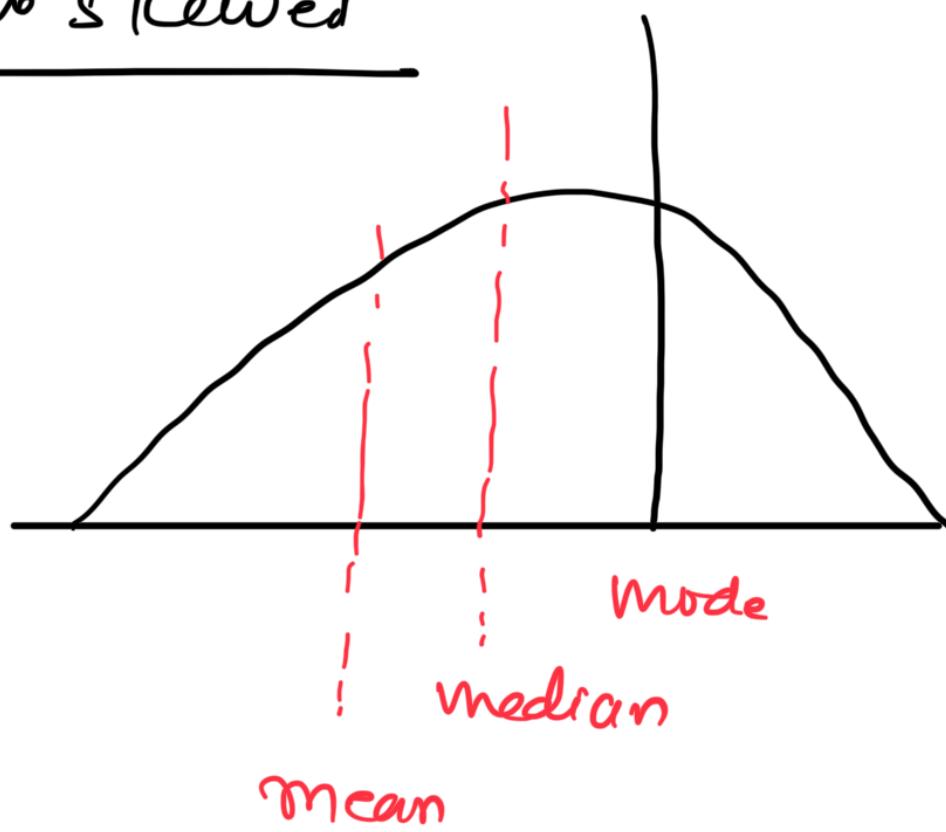
3. Normal distributions
-

Skewness 

Positive skewed
Negative skewed



Negative skewed



mean < median < mode

Summary

Kurtosis - Measures the "tailedness"

of the distribution

Skewness - Measures the asymmetry

of the distribution

1) Positive skew - Mean > Median

2) Negative skew - Mean < Median

3) Zero skew - Mean = Median

Covariance & Correlation

Covariance: $\text{Cov}(x, y)$

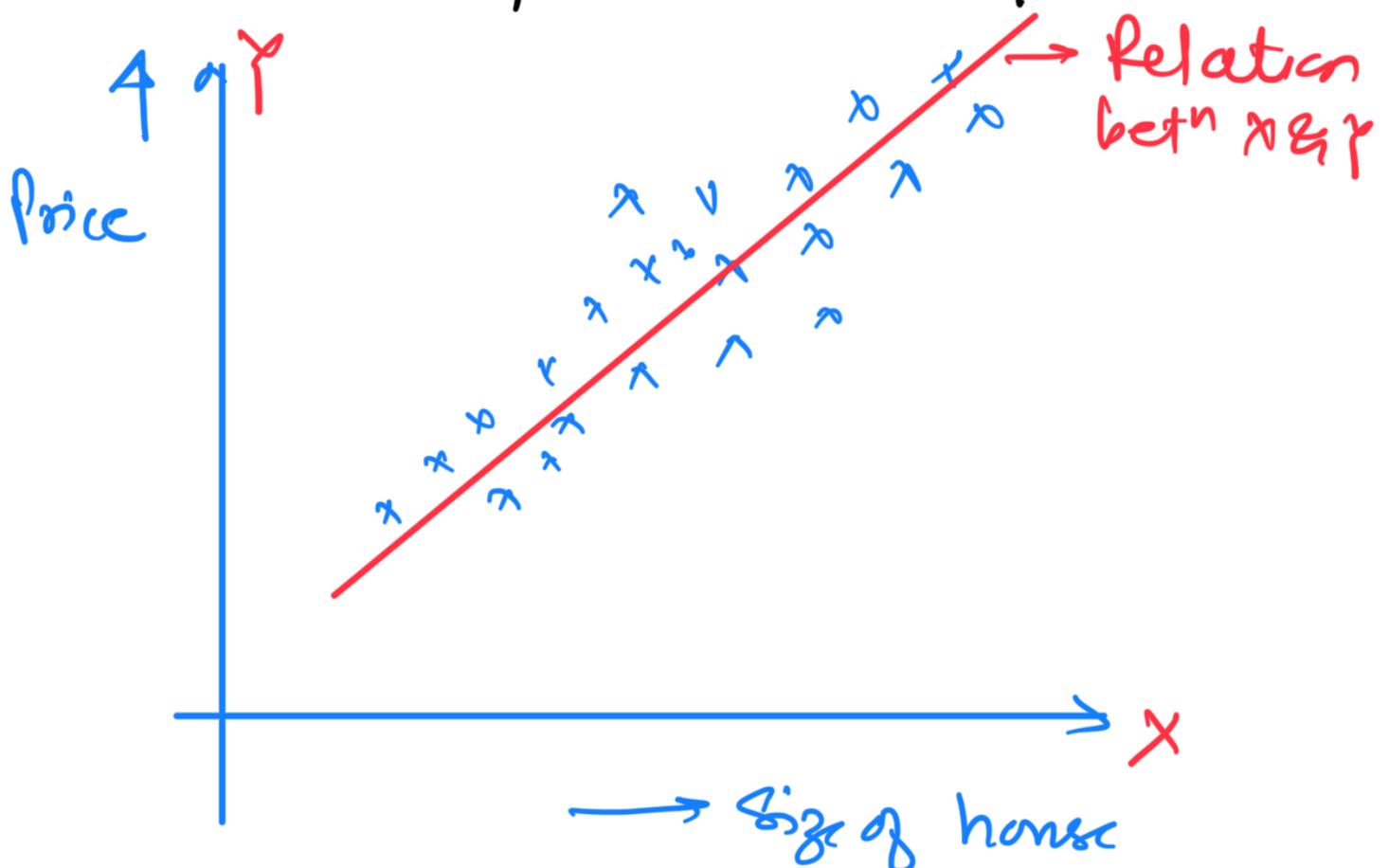
$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

Correlation: ($\rho(x, y)$)

$$\rho(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{s_x s_y}}$$

Regression Analysis

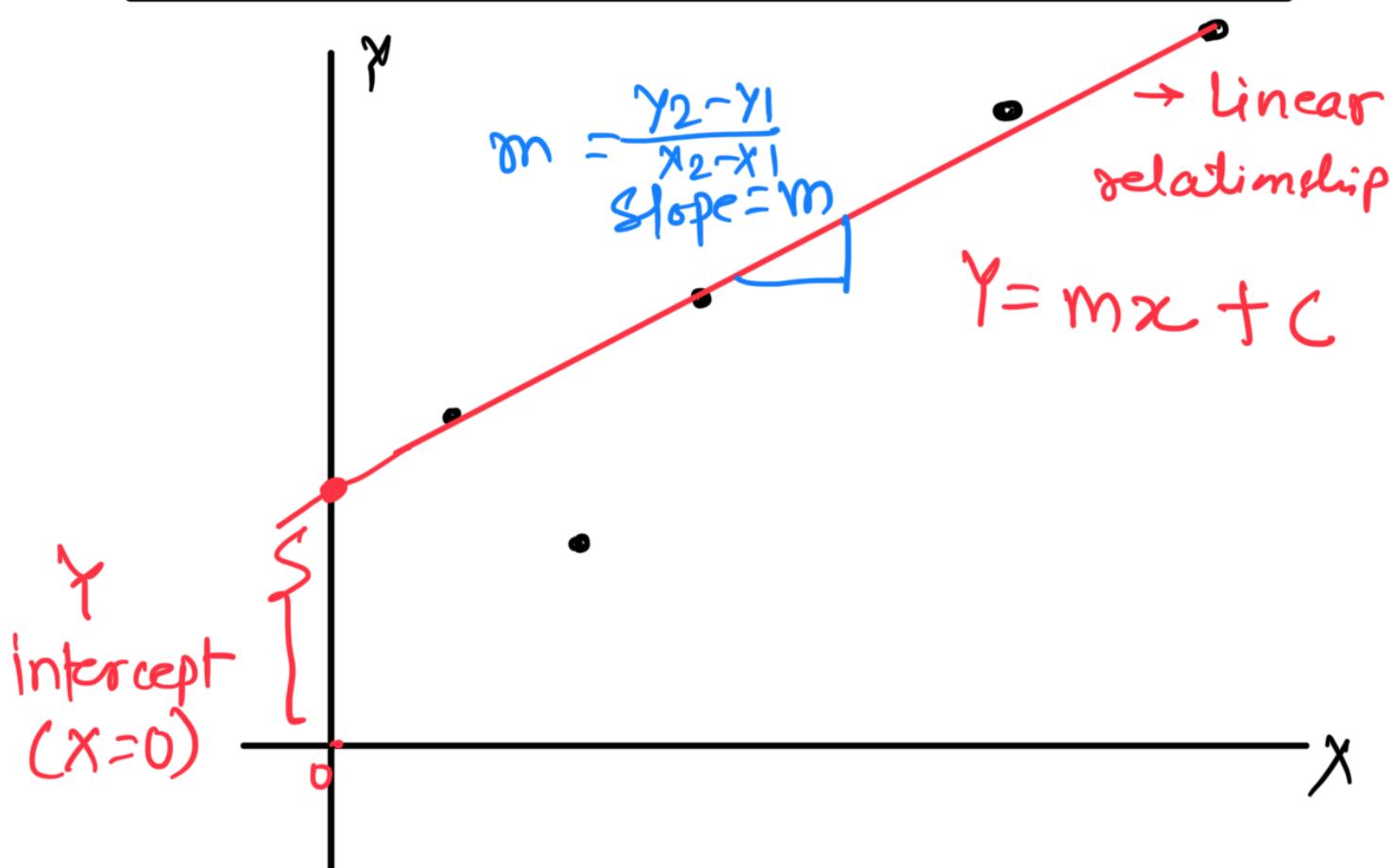
Example : Regression Techniques



UK

To define relationship we can use
Regression Technique

X	1	2	3	4	5
Y	5	3	6	8	10



Get the value of y

$$Y = mX + C$$

↓ ↓

Slope Value of x
from table

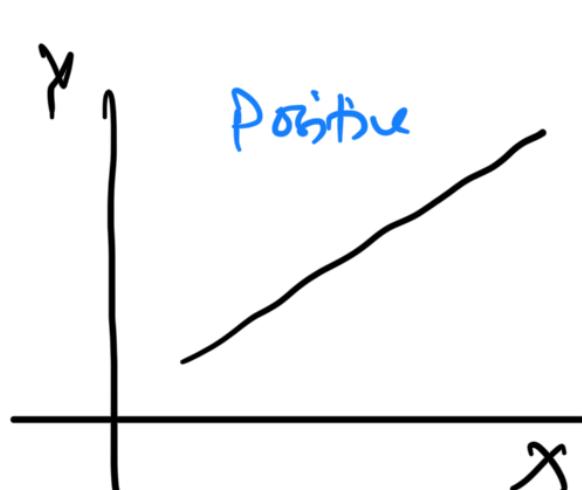
Constant term

$$x = 10$$

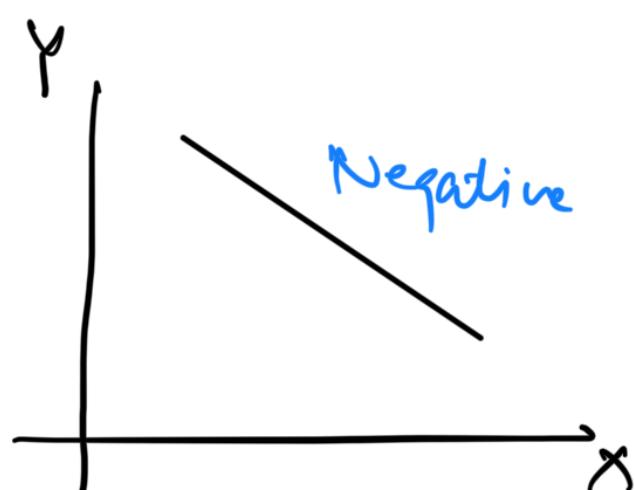
$$x = 20$$

Linear Regression Model

↓
Identify the relationship between features/columns of the dataset



Positive linear
relationship



Negative linear
relationship

$X \rightarrow \text{increase}$
$Y \rightarrow \text{increase}$
<hr/>
$X \rightarrow \text{decrease}$
$Y \rightarrow \text{decrease}$

$X \rightarrow \text{increase}$
$Y \rightarrow \text{decrease}$
<hr/>
$X \rightarrow \text{decrease}$
$Y \rightarrow \text{decrease}$

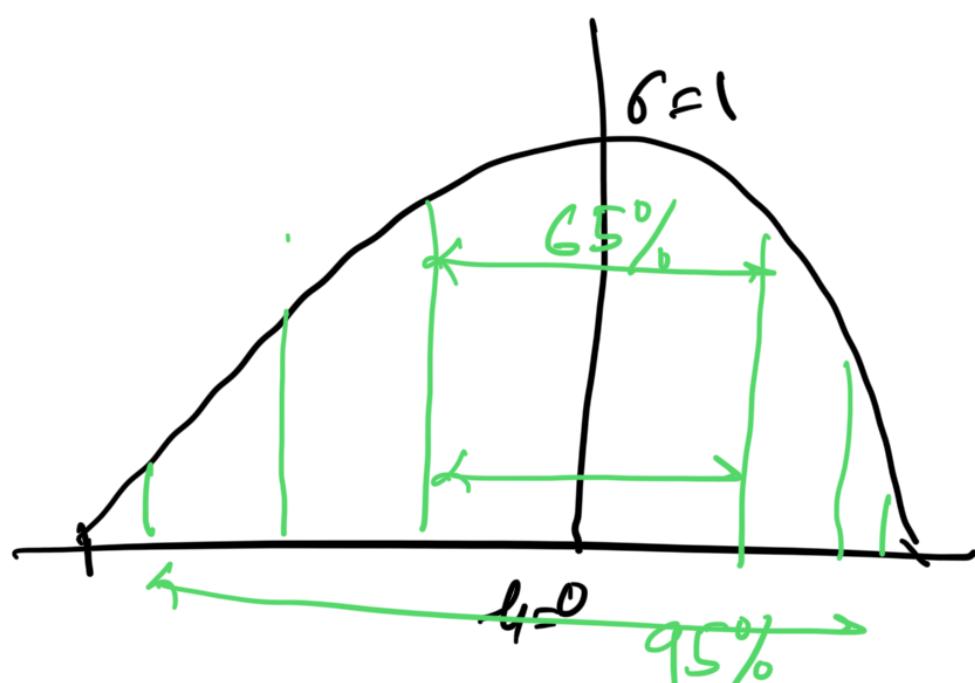
Probability & Distribution

Probability distribution - Understanding

Probability distribution is crucial for many machine learning algorithms.

① Normal Distribution (Gaussian distribution)

- bell-shaped curve, symmetric around the mean
- Used in many algorithms
 - e.g. - Linear Regression
 - Logistic Regression
- Used in statistical test



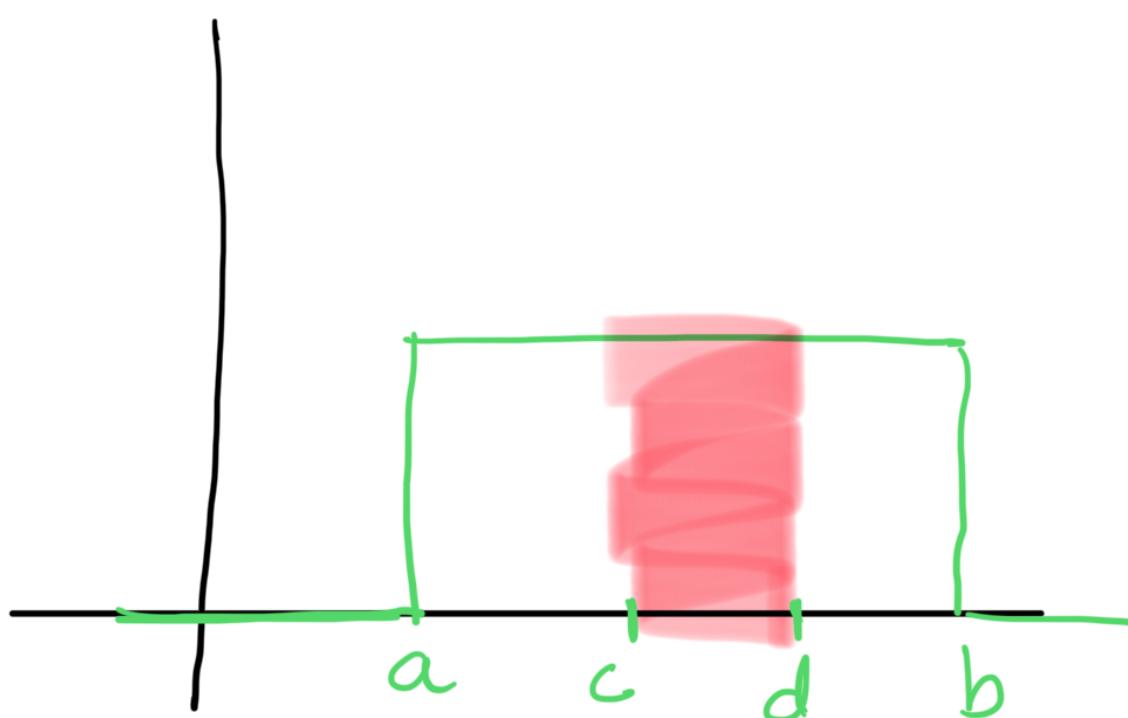
$f(x) = \text{pdf} = \text{probability distribution function}$

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

2. Uniform distribution

- has constant probability.
- Used for random Sampling
- Useful in simulation



$$f(x) = \frac{1}{b-a}$$

$$\mu = \frac{a+b}{2}$$

$$P(c \leq x \leq d) = \frac{d-c}{b-a}$$

$$\sigma = \sqrt{\frac{(b-a)^2}{12}}$$

③ Binomial Distribution

It describes the number of success in a fixed number of trials.

- used in binary classification
- helpful in modeling binary outcomes

$$f(x) = \binom{n}{x} p^x q^{(n-x)}$$

n = no. of trials p = prob of success

x = # of success $q = (1-p)$
= prob of failure

$$\binom{n}{x} = \frac{n!}{x!(n-x)!}$$

Hypothesis Testing - is used to make

decisions about a population based on sample data. It helps in determining if there is enough evidence to reject a null hypothesis.

① t-test : compares the mean,

of two groups to see if they are significantly different from each other.

→ feature selection

→ Comparing 2 models

② Chi - Square Test: Checks the

independence of two categorical variable

- feature selection

- Understanding association in Categorical data

③ Pearson Correlation Coefficient:

measures the linear correlation between two variables.

- quantifying the degree of linear relationship

- Identify highly correlated features.