

# Machine Learning

What is Machine Learning?

Machine Learning is a branch of artificial intelligence that enables computers to learn and improve from experience without explicit programming.

It involves developing algorithms that analyze data, identify patterns and make decisions or predictions.

Categories :

1. Supervised Learning

2. Unsupervised Learning

3. Semi-Supervised Learning

4. Reinforcement Learning.

1. Supervised Learning: Uses labelled data to train models for classification

e.g.: Spam detection.

2. Unsupervised Learning: deals with unlabeled data to uncover hidden patterns, such as clustering customers by behaviour.

3. Semi-Supervised Learning: Combines small amount of labelled data and large amount of unlabeled datasets for improving learning.

4. Reinforcement Learning: Trains agents to make sequential decisions through rewards and penalties.

- playing games or controlling robots

## Steps in ML:

1. Collect and Preprocess data.
2. Split data into training and testing sets.
3. Train the model on the training set.
4. Evaluate and refine model.
5. Deploy the model for real-time predictions!

## 2 Applications Areas of Machine Learning:

- Spam Filtering: Identifies and blocks unwanted mails.
- Smart assistants: powering virtual assistants like Alexa, Siri and google Assistant to perform tasks using voice interaction.
- Social Media: Friend recommendations and Face Recognition.
- Navigation: Google Maps uses ML to analyze traffic and suggest optimal routes.
- Self Driving Cars: Helps vehicles to perceive their surroundings and make driving decisions.
- Recommendation System: Suggests movies/products or music based on user preferences.
- Cyber Security: Detects unusual activities and delivers traffic for early warnings.

### 3 Data Pre-processing :

- Data Pre-Processing is the essential step for preparing raw data for machine learning models.
- It improves the quality of data, making it suitable for analysis.

1. Data Cleaning : Handling Missing values, removing duplicates, and correcting errors.

2. Data Transformation : Converting data into a usable format, such as, scaling numerical values or encoding categorical variables.

3. Feature Selection : Choosing relevant features to reduce complexity and improve model performance.

4. Data Splitting : Dividing the dataset into training and testing sets.

## 4 Training and Choosing Predictive Models :

### 1. Training the model :

The model learns from a training dataset by identifying patterns and relationships between input features and output labels.

### 2. Choosing Predictive Models :

Different algorithms (e.g.: Decision Tree, SVM, Neural networks) are evaluated to find the best fit for the problem. The choice depends on the factors like datatype, size and complexity.

### 3. Optimization :

Adjust hyperparameters to improve the model's performance.

### 4. Testing :

Test the model on unseen data to ensure it can make accurate predictions in real-world scenarios.

## Unit - I

### 1. Machine Learning.

Machine learning is a branch of artificial intelligence where computers learn to perform tasks by analysing patterns in data instead of being explicitly programmed to follow specific rules.

How it works:

1. Data is fed into the machine
2. The machine finds patterns or relationships in the data.
3. It uses those patterns to make predictions or decisions without being told exactly what to do.

Example: To teach machine to recognize Cats, by providing labeled images of Cats and non-Cats. This machine learns patterns from these images and uses them to identify Cats in new pictures.

Categories:

1. Supervised
2. Unsupervised
3. Semi-supervised
4. Reinforcement

1. Supervised:

In Supervised learning, the machine is trained using labeled data. we provide input-output pairs, and machine learns the relationship between them.

• Example: predicting house prices based on features like size, location etc.

## 2. Unsupervised Learning:

Here, the data is not labeled. The machine tries to find hidden patterns or structures in the data without any guidance.

- Example: Grouping customers based on their purchase behaviour (clustering).

## 3. Semi-Supervised Learning:

A mix of both Supervised and Unsupervised learning. The machine is given a small amount of labeled data and a larger amount of unlabeled data.

- Example: A machine learning model trained to classify images with a small number of labeled photos and many unlabelled ones.

## 4. Self-Supervised Learning:

In this type, the machine generates its own labels from the data. It can predict one part of the data using another part of the same data.

- Example: Predicting the next word in a sentence during natural language processing.

## 5. Reinforcement Learning:

The machine learns by interacting with its environment and receiving feedback in the form of rewards or penalties.

- Example: A robot learning to navigate a maze by receiving rewards for reaching the destination.

## 2 Applications of Machine Learning:

Machine Learning has a wide range of applications across various industries.

1. Health Care
2. Finance
3. Retail and E-commerce
4. Autonomous Vehicles
5. Natural language Processing
6. Image and video Analysis
7. Robotics
8. Cybersecurity
9. Gaming
10. Marketing and Advertising

### 1. Health Care:

- disease prediction and diagnosis
- personalized treatment recommendations
- detecting tumors in X-rays

### 2. Finance:

- Fraud detection
- Stock Market Predictions
- Credit Scoring and Risk Management

### 3. Retail and E-commerce:

- Product Recommendations (Amazon, Netflix)
- Customer behaviour Analysis

### 4. Autonomous Vehicles:

- Self-driving Cars
- Traffic pattern Analysis

## 5. Natural Language Processing:

- Speech Recognition (Siri, Alexa)

- Language translation (Google Translate)

- Chatbots and Virtual assistants

## 6. Image and Video Analysis:

- Facial Recognition (Unlocking Phone)

- Object detection in Images and videos

## 7. Robotics:

- Industrial Automation and robot control

- Reinforcement learning in robot movement and task performance.

## 8. Cyber Security:

- Detecting and Preventing Cyberattacks

- Spam filtering and malware detection

## 9. Gaming:

- AI opponents in games (Chess)

- Game world simulation and player interaction.

## 10. Marketing and Advert Advertising:

- Target Advertising.

- Ad Performance Optimizations

- Predictive analytics for sales and conversions.

### 3 Data Pre-Processing

Data Preprocessing in machine learning is the step where raw data is cleaned up and prepared to be used by models.

#### Tasks:

- Handling Missing Data (filling in or removing incomplete entries)
- Normalizing or Scaling data (Putting it on the same scale)
- Encoding Categorical Data (turning text labels into state numbers)
- Removing noise or irrelevant data

### 4 Choosing Predictive Models

Choosing a predictive model in machine learning depends on the type of problem and the data.

1. Regression Model (like Linear Regression) are used when predicting continuous values (eg. house prices)
2. Classification Models (like Decision Trees or Logistic Regression) are used when predicting categories (eg: whether an email is spam or not)
3. Clustering models (like K-means) are used to group similar data points without predefined labels.

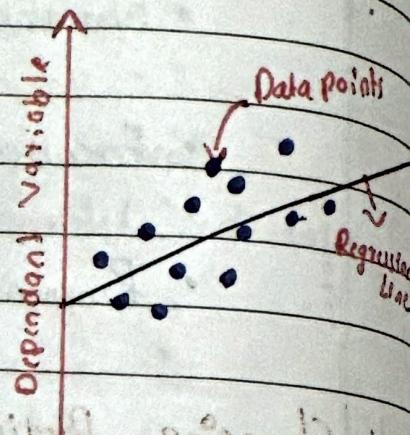
## Unit-II

### 1 Linear Regression

Linear Regression is a machine learning algorithm used to predict a continuous value based on the relationship between input features and the output. It finds the best-fitting straight line (regression line) through the data points by minimizing the differences between the predicted values and the actual values.

Eg:

- weight of the car and miles per gallon
- Area of the house and price of the house.



### Linear Model

- They make the prediction using a linear function of the input features.
- For Linear Regression, the general prediction formula is

$$b_0 = (\sum y_i)(\sum x_i^2) - (\sum x_i)(\sum x_i y_i)$$

$$\sum x_i^2 - (\sum x_i)^2$$

$$b_1 = \frac{(\sum x_i y_i) - (\sum x_i)(\sum y_i)}{\sum x_i^2 - (\sum x_i)^2}$$

$$\therefore b_1 = \frac{(\sum x_i y_i) - (\sum x_i)(\sum y_i)}{\sum x_i^2 - (\sum x_i)^2}$$

Qn1) Glucose level and Age of Six Subjects are given in table below. Predict the glucose levels of a person of age 55. Using Linear Regression:

| Subject | Age | GLlevel |
|---------|-----|---------|
| 1       | 43  | 99      |
| 2       | 21  | 65      |
| 3       | 25  | 79      |
| 4       | 49  | 75      |
| 5       | 57  | 87      |
| 6       | 59  | 81      |

| X   | Y   | $x^2$ | XY    |
|-----|-----|-------|-------|
| 43  | 99  | 1849  | 4257  |
| 21  | 65  | 441   | 1365  |
| 25  | 79  | 625   | 1975  |
| 42  | 75  | 1764  | 3150  |
| 57  | 87  | 3249  | 4959  |
| 59  | 81  | 3481  | 4779  |
| 247 | 486 | 11409 | 20485 |

$$b_0 = (\sum Y)(\sum x^2) - (\sum x)(\sum xy)$$

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

$$= \frac{486 \times 11409 - 247 \times 20485}{6(11409) - (247)^2}$$

$$= \frac{5,544,774 - 5,059,795}{68454 - 61009}$$

$$= \frac{484,979}{7445}$$

$$b_0 = 65.142$$

Date:

$$b_1 = \frac{1}{n} (\sum xy) - (\bar{x})(\bar{y})$$

$$\frac{1}{n} (\sum x^2) - (\bar{x})^2$$

$$= \frac{1}{6} (20485) - (347)(486)$$

$$= \frac{1}{6} (11409) - (347)^2$$

$$\begin{array}{r} 2868 \\ - 7445 \\ \hline 0223 \end{array}$$

$$b_1 = 0.3852$$

$$y = b_0 + b_1 x$$

$$= 65.142 + 0.28852 (58)$$

$$y =$$

| X  | $\bar{x}$ | $y$ | $\bar{y}$ | X  |
|----|-----------|-----|-----------|----|
| 10 | 65.142    | 58  | 28        | 84 |
| 12 | 67.027    | 62  | 30        | 16 |
| 14 | 68.912    | 65  | 32        | 26 |
| 16 | 70.797    | 71  | 34        | 46 |
| 18 | 72.682    | 75  | 36        | 56 |
| 20 | 74.567    | 78  | 38        | 76 |
| 22 | 76.452    | 82  | 40        | 96 |

$$(\bar{x}\bar{y}) - (\bar{xy})$$

$$(\bar{x}\bar{y}) - (\bar{xy})$$

$$23.106 \times 28.6 - 40.811 \times 38.8$$

$$(\bar{x}\bar{y}) - (\bar{xy})$$

### 3. Advantage of Linear Regression:

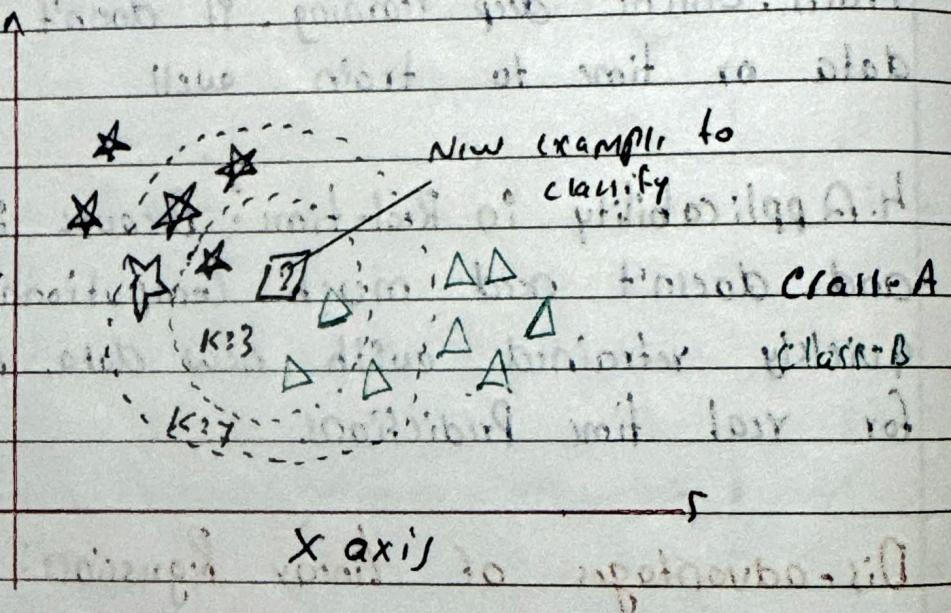
1. Easy to Implement : It is easy to implement and maintain because it doesn't need much effort or complex engineering.
2. Scalability : Since it's not computationally expensive, it can be used for situations that need scaling like big data applications.
3. Interpretability : It is easy to understand and quick to train unlike deep learning. It doesn't need as much data or time to train well.
4. Applicability in Real-time : Because it's easy to train and doesn't need much computational power, it can be quickly retrained with new data, making it useful for real time Predictions.

### 4. Dis-advantages of Linear Regression:

- It assumes a linear relationship between the dependent and independent variables, which is rarely seen in real-world data.
- It assumes ~~the~~ a straight line relationship between variables which isn't the case.
- It is prone to noise and overfitting.
- It's not good choice for data sets with fewer observations than attributes as it can cause overfitting.

## KNN - K-Nearest Neighbour

- KNN is a lazy learning non-parametric algorithm uses data with several classes to predict the classification of the new sample point.
- This called a lazy learning algorithm or lazy learner because it doesn't perform any training when you supply the training data. Instead it just stores the data during the training time and doesn't perform any calculations.



formula:  $\sqrt{(x-a)^2 + (y-b)^2}$

|            | a          | b     | c |
|------------|------------|-------|---|
| Brightness | Saturation | Class |   |

|                |    |    |     |
|----------------|----|----|-----|
| d <sub>1</sub> | 40 | 30 | Red |
|----------------|----|----|-----|

|  |    |    |      |
|--|----|----|------|
|  | 50 | 50 | Blue |
|--|----|----|------|

|  |    |    |      |
|--|----|----|------|
|  | 60 | 90 | Blue |
|--|----|----|------|

|  |    |    |     |
|--|----|----|-----|
|  | 25 | 70 | Red |
|--|----|----|-----|

|  |    |    |      |
|--|----|----|------|
|  | 70 | 70 | Blue |
|--|----|----|------|

|  |    |    |     |
|--|----|----|-----|
|  | 60 | 10 | Red |
|--|----|----|-----|

|  |    |    |      |
|--|----|----|------|
|  | 25 | 80 | Blue |
|--|----|----|------|

|  |    |    |     |
|--|----|----|-----|
|  | 10 | 90 | Red |
|--|----|----|-----|

|  |    |    |      |
|--|----|----|------|
|  | 10 | 90 | Blue |
|--|----|----|------|

|  |    |    |     |
|--|----|----|-----|
|  | 10 | 90 | Red |
|--|----|----|-----|

|  |    |    |      |
|--|----|----|------|
|  | 10 | 90 | Blue |
|--|----|----|------|

|  |    |    |     |
|--|----|----|-----|
|  | 10 | 90 | Red |
|--|----|----|-----|

|  |    |    |      |
|--|----|----|------|
|  | 10 | 90 | Blue |
|--|----|----|------|

|  |    |    |     |
|--|----|----|-----|
|  | 10 | 90 | Red |
|--|----|----|-----|

|  |    |    |      |
|--|----|----|------|
|  | 10 | 90 | Blue |
|--|----|----|------|

formula  $\sqrt{(x-a)^2 + (y-b)^2}$

$$d_1 = \sqrt{(20-40)^2 + (35-30)^2} \quad d_2 = \sqrt{(20-50)^2 + (35-50)^2}$$

$$d_1 = \sqrt{(20-40)^2 + (35-30)^2} = \sqrt{400 + 25} = \sqrt{425} = 25$$

$$d_1 = \sqrt{400 + 25} = 25$$

$$d_1 = \sqrt{625} = 25$$

$$d_2 = \sqrt{(20-60)^2 + (35-90)^2} \quad d_4 = \sqrt{(20-10)^2 + (35-25)^2}$$

$$= \sqrt{(40)^2 + (55)^2} \quad = \sqrt{100 + 100}$$

~~$d_2 = \sqrt{400 + 3025} = 55$~~

$$= \sqrt{1600 + 3025} = 55$$

$$d_2 = \sqrt{(20-70)^2 + (35-70)^2} = \sqrt{2500 + 1225} = 50$$

$$= 61.8$$

$$d_6 = \sqrt{(20-60)^2 + (35-10)^2} = 47.17\text{mm}$$

$$d_7 = \sqrt{(20-25)^2 + (35-80)^2} = 45\text{mm}$$

| Brightness | Saturation | Class | Distance | Rank  |
|------------|------------|-------|----------|-------|
| 40         | 80         | Red   | 25       | 10    |
| 50         | 50         | Blue  | 33.54    | 25    |
| 60         | 90         | Blue  | 68.01    | 33.54 |
| 10         | 25         | Red   | 10       | 45    |
| 70         | 70         | Blue  | 61.03    | 47.17 |
| 60         | 10         | Red   | 47.17    | 61.03 |
| 25         | 80         | Blue  | 45       | 68.01 |

The new entry is classified as Red.

Steps involved in KNN (Algorithm)

- Calculate the distance between test data and each row of training data with the help of any of any the method.

The most common method to calculate distance is euclidean.

- Now, based on the distance value, sort them in ascending order.

- Now, it will choose the top  $K$  rows from the sorted array.

- Now, it will assign a class to the test point based on most frequent class of these rows.

## 5 Naive Bayes Classifier (probability based)

$$\text{Posterior Prob. } P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$P(A|B)$  = posterior probability (likelihood of  $A$ )

$P(B|A)$  = Likelihood involving a number of other

$P(A)$  = Class Prior

$P(B)$  = prediction prior (with very small  $\epsilon$ )

① Collect the data

② Convert it into frequency table

③ Calculate prior / class predictions

④ Apply it on Bayes to fix the predicted class

Distance  
Red  
Red

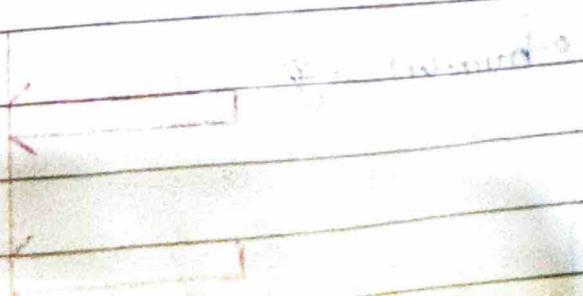
Blue  
Blue

Red

15

20

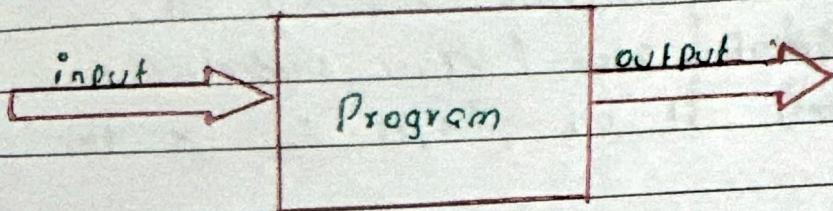
the lines starting from around point midum at 15  
and ending at 20 and not meeting (at border)  
at the bottom



## 6 Traditional Programming vs Machine Learning.

### Traditional Programming.

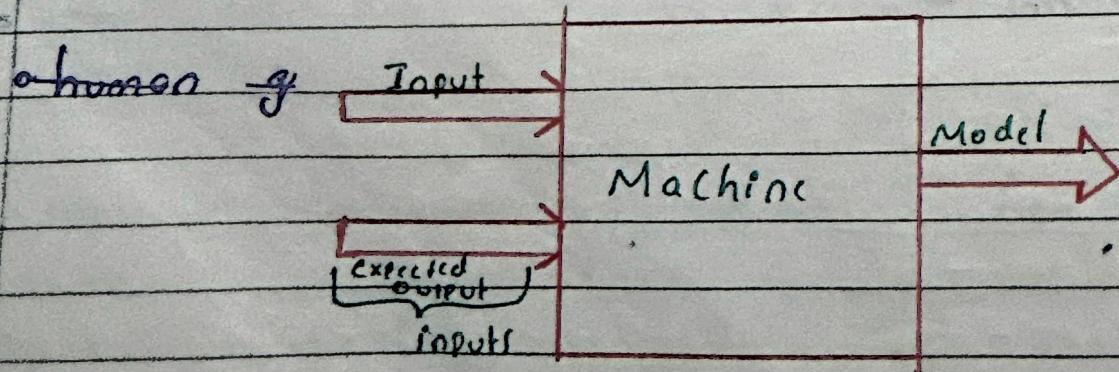
- In traditional programming, humans write explicit rules to solve a problem.
- humans give the computer inputs, and define the exact steps.



- The programs written once will not be changed until and unless they are updated manually by somebody.

### Machine Learning.

- In machine learning, humans don't explicitly write rules. Instead, the computer learns the rules on its own by analyzing data.



- These Model are nothing but programs generated by algorithms without being explicitly programmed by humans
- These models continuously evolve into better models automatically as the amount of quality data increases

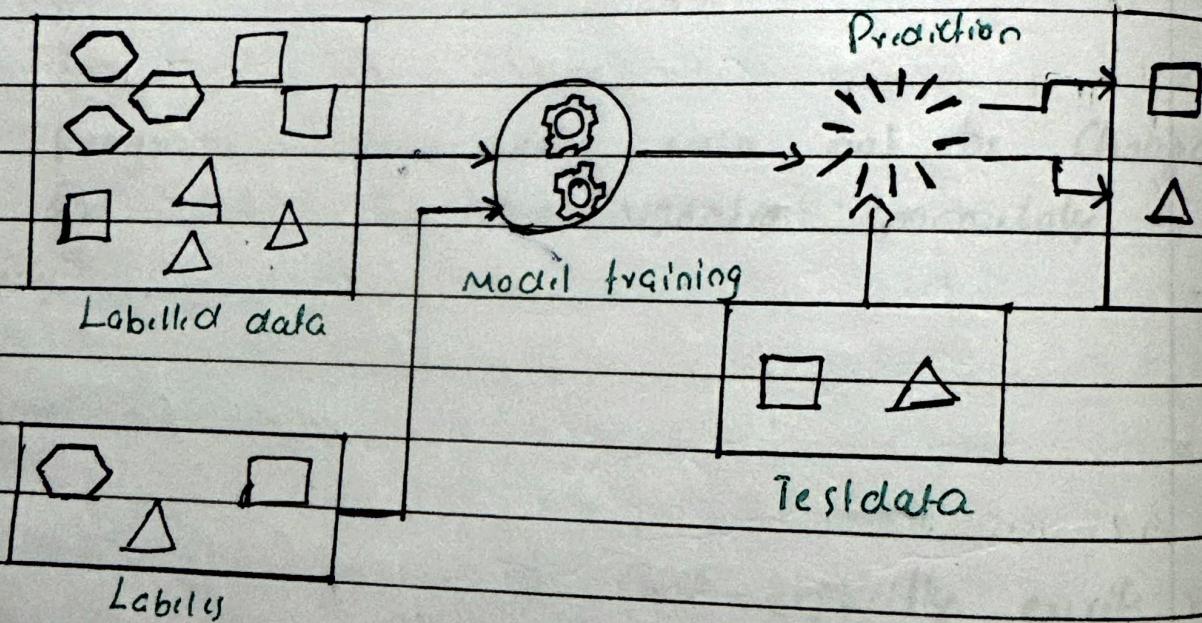
## Unit - II

Supervised Learning is a machine learning approach where models are trained using labelled data.

Each input has corresponding output, allowing the algorithm to learn the relationship between them.

It uses this learned relationship to make predictions on new, unseen data.

Eg:- Identifying if an email is spam or not

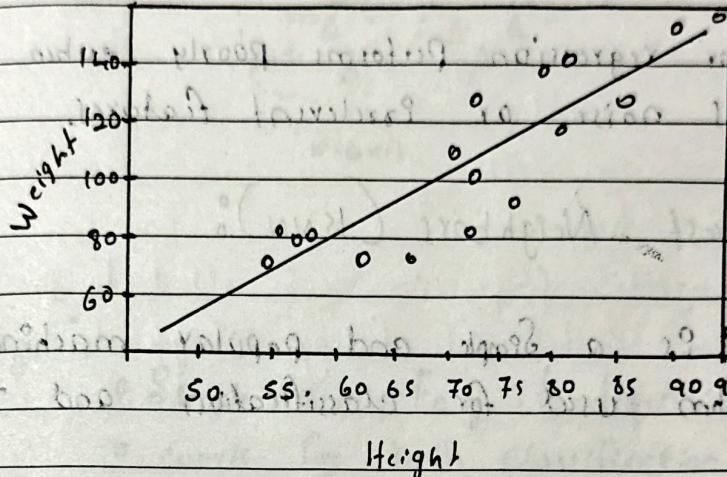


## 2 Linear Regression:

Linear Regression is a Supervised Learning technique used to predict continuous numerical value by finding the relationship between inputs and outputs.

$$\text{Formula: } Y = mX + C.$$

where  $m$  is slope and  $C$  is the intercept



e.g. weight of a car based on its height.

- weight of a car based on its height.

- Area of the house and Price of the house.

Advantage: It is easy to understand and implement.

- Linear regression is easy to understand and implement.
- The model clearly shows the relationship between inputs and outputs.
- It requires minimal computational resources, making it efficient for real-time applications.
- Can handle large datasets.
- Useful in fields like finance, healthcare and marketing.
- The model can be trained quickly with new data.

## DPS-Advantages:

- It assumes straight line relationships between dependent and independent variables which is rare in real world.
- It is prone to noise and overfitting.
- Linear regression performs poorly when the data contains noise or irrelevant features.

## 3. K-Nearest Neighbors (KNN):

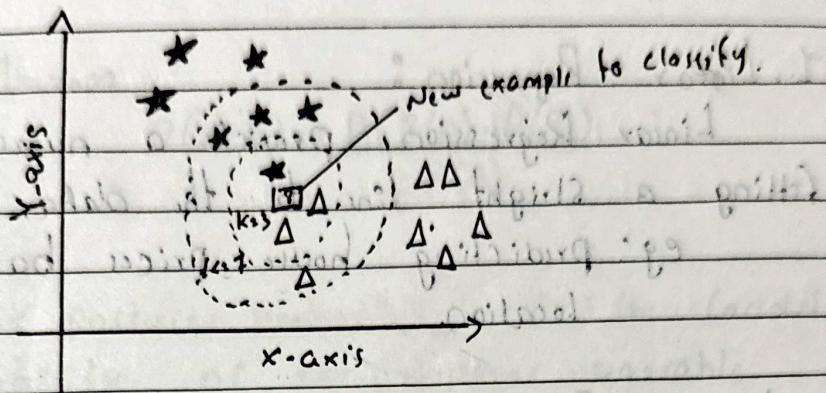
- KNN is a simple and popular machine learning algorithm used for classification and regression.
- It predicts results based on the K-nearest data points from the training datasets.
- It's called a lazy learning or Lazy learner because it doesn't perform any training even we supply the data, instead it just stores the data during time.

### How it works?

1. Distance Calculation: Compute the distance between the new data and training point.
2. Sort the Neighbors: Sort the distances in ascending order.
3. Select the K-neighbors: Pick the top K-nearest neighbors.

#### 4. Mals. Predictions:

- o For classification: Assign the most common value
- o For regression: Take the average value.



- Advantage:
- o Simple and Easy to understand
  - o works for both classification and regression
  - o Handles multi-class problems efficiently
  - o No need for training model
- Ds-advantage:
- o Computationally expensive
  - o struggles with large datasets
  - o Requires proper datasets for good accuracy

#### Example:

Predict Persons T-Shirt size based on their height and weight by calculating distance

Distance =  $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

## 4 Linear Models in Machine Learning:

Linear models are a type of machine learning algorithm that makes predictions based on a linear relationship between input features and the target variable.

### 1. Linear Regression:

Linear Regression predicts a numerical value by fitting a straight line to the data.  
e.g. Predicting house prices based on size and location.

### 2. Logistic Regression:

Logistic Regression is used to predict binary outcomes instead of fitting straight lines.  
e.g. predicting whether tumour is benign or malignant.

#### Linear Regression

- Used to predict the continuous dependent variables
- Used for solving regression problems
- Predict the value of continuous variables
- Least square estimation method or maximum likelihood estimation is used
- Output must be continuous

#### Logistic Regression

- Used to predict the categorical dependent variable
- Used for solving classification problems
- Predict the value of categorical variables
- It is not required to have linear relationships between dependant and independent variables.

## 5 Naive Bayes Classifier:

Naive Bayes assumes that all features in the datasets are conditionally independent given each class. It classifies data using Bayes theorem.

### 1. Bayes Theorem.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

- $P(A|B)$  or Posterior Probability: is the conditional probability of the response variable.
- $P(B)$ : or Class Prior: is the prior probability of the response variable.
- $P(B|A)$  or Likelihood: is basically the reverse of the posterior probability.
- $P(A)$  or Marginal Probability: is the evidence or the probability of training data.

How it works:

- Step 1: Convert raw data into frequency tables.
- Step 2: Calculate prior, likelihood, and posterior probability.
- Step 3: Classify based on highest posterior probability.

Example:

Predicting if someone play tennis based on weather conditions like "Rainy", "Hot", "High humidity" and "No wind".