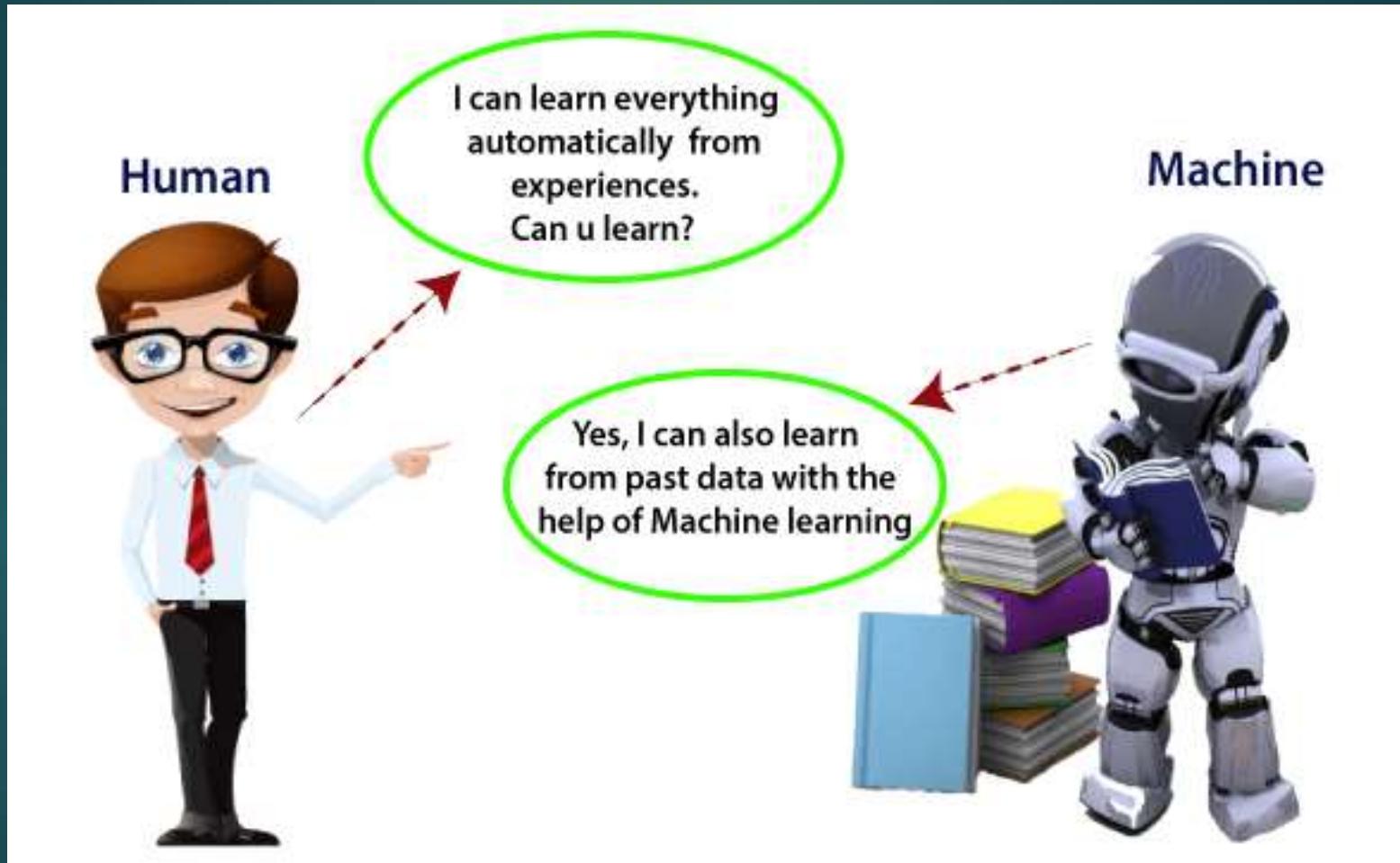


# Unit-3-5

# Introduction to Machine Learning

## CSE3007

# What is Machine Learning

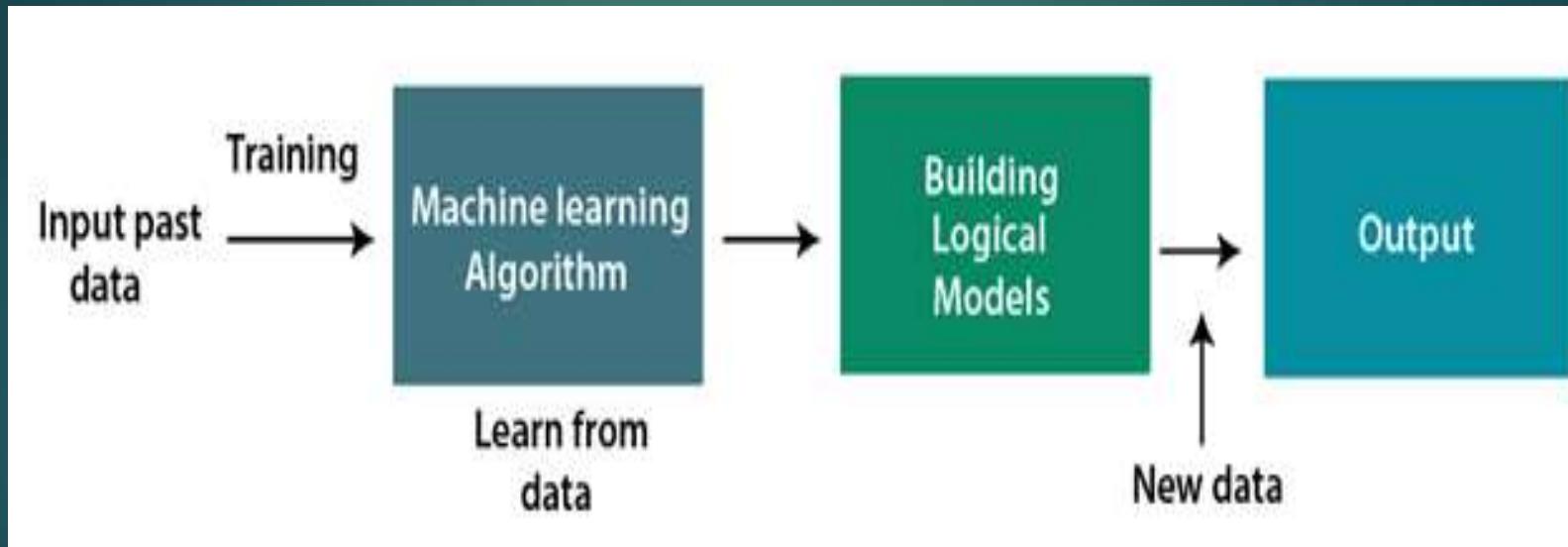


# Machine Learning

- ❑ Machine Learning is said as a subset of **artificial intelligence** that is mainly concerned with the development of algorithms which allow a computer to learn from the data and past experiences on their own.
- ❑ The term machine learning was first introduced by **Arthur Samuel** in **1959**. We can define it in a summarized way as:

*“Machine learning enables a machine to automatically learn from data, improve performance from experiences, and predict things without being explicitly programmed”*

# Block Diagram of ML



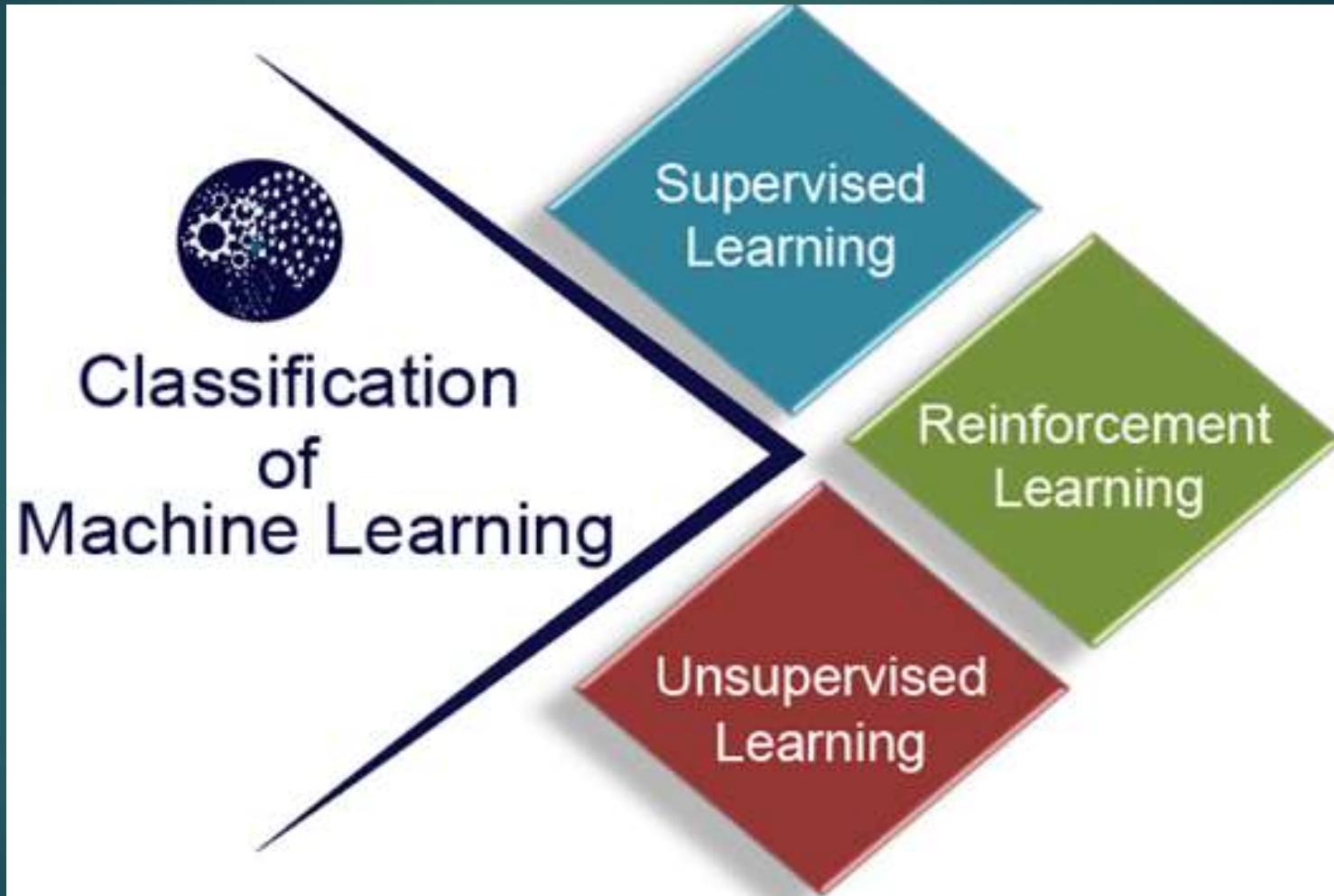
# Features of Machine Learning:

- ▶ Machine learning uses data to detect various patterns in a given dataset.
- ▶ It can learn from past data and improve automatically.
- ▶ It is a data-driven technology.
- ▶ Machine learning is much similar to data mining as it also deals with the huge amount of the data.

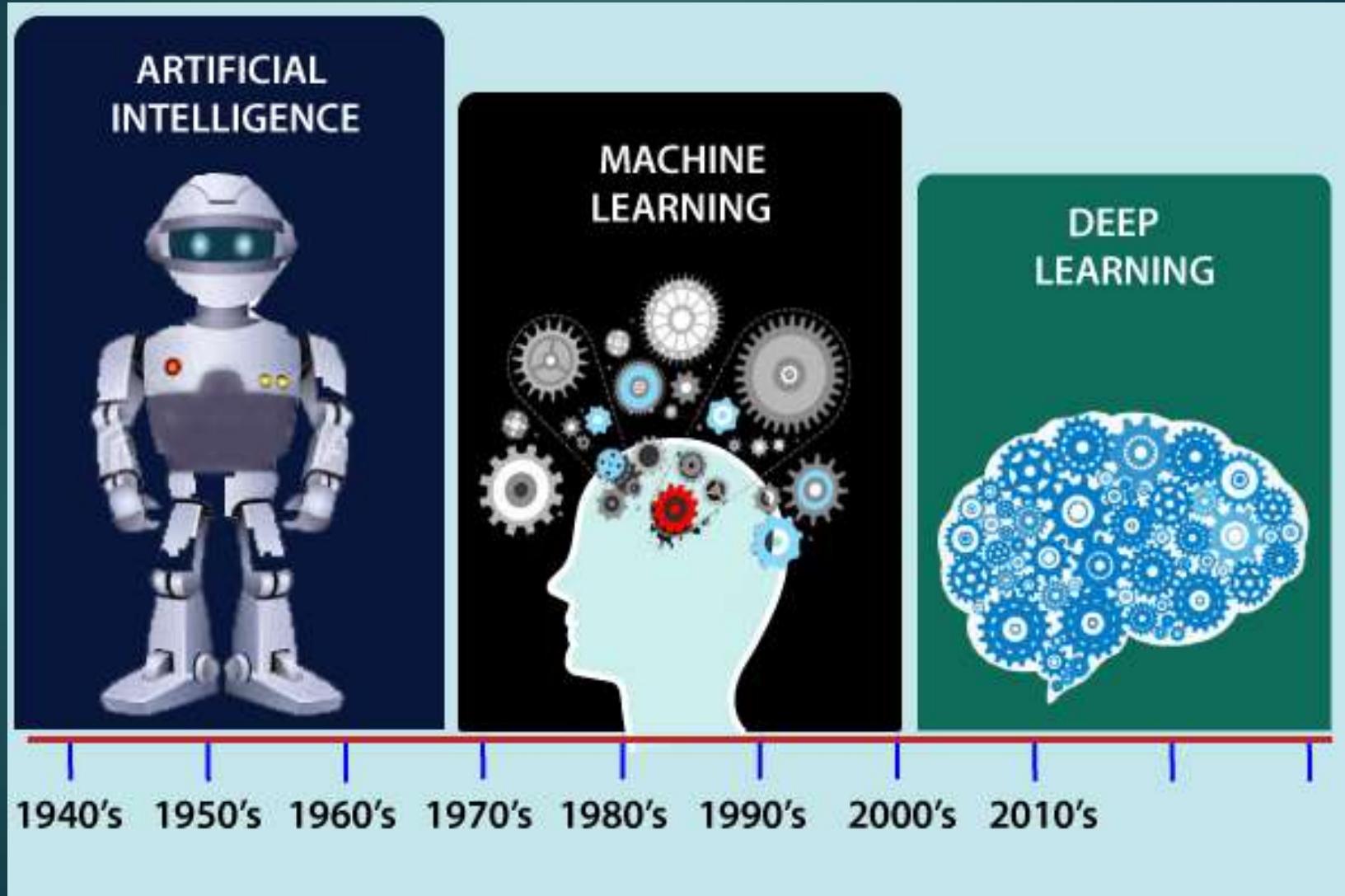
# Need for Machine Learning:

- Rapid increment in the production of data
- Solving complex problems, which are difficult for a human
- Decision making in various sector including finance
- Finding hidden patterns and extracting useful information from data.

# Classification of ML:



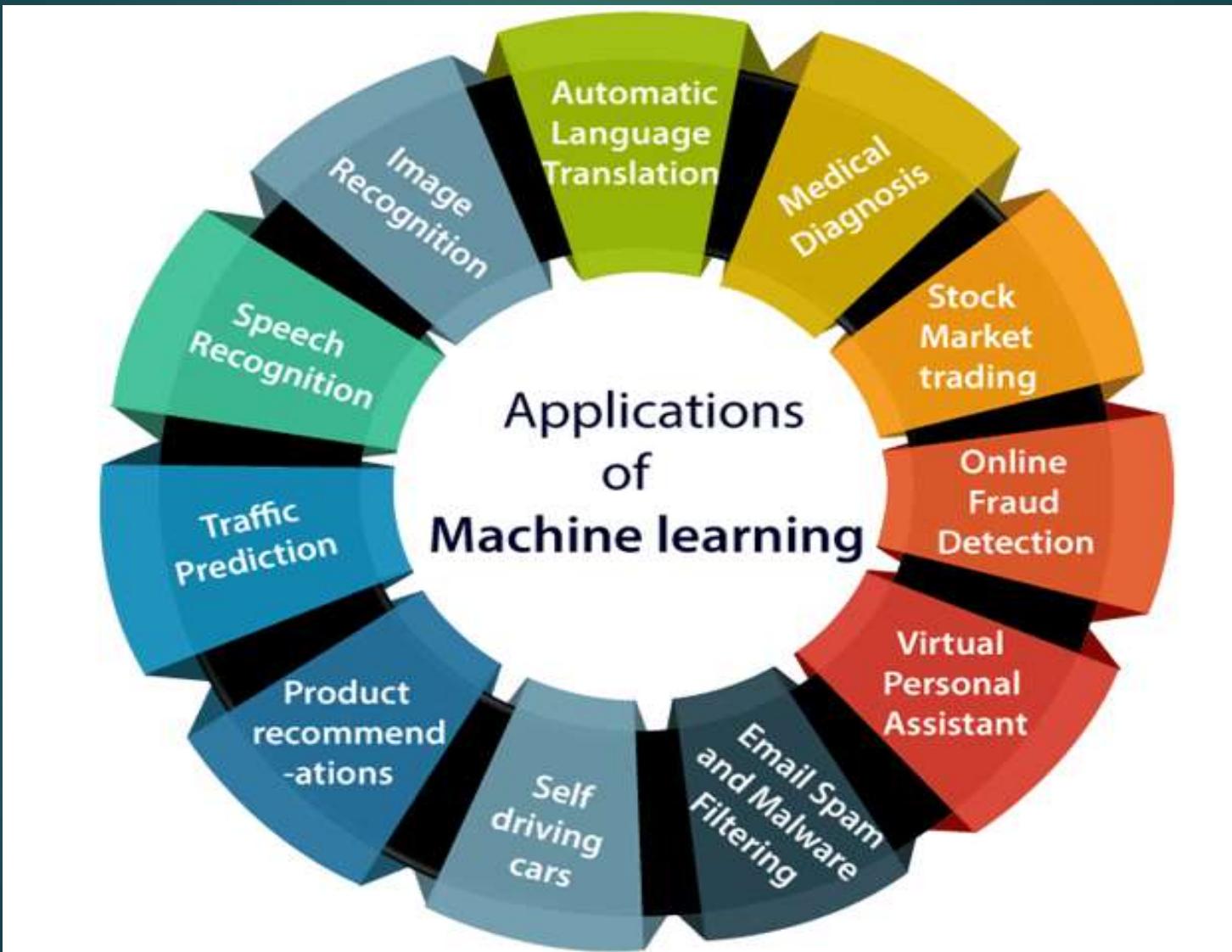
# History of ML:



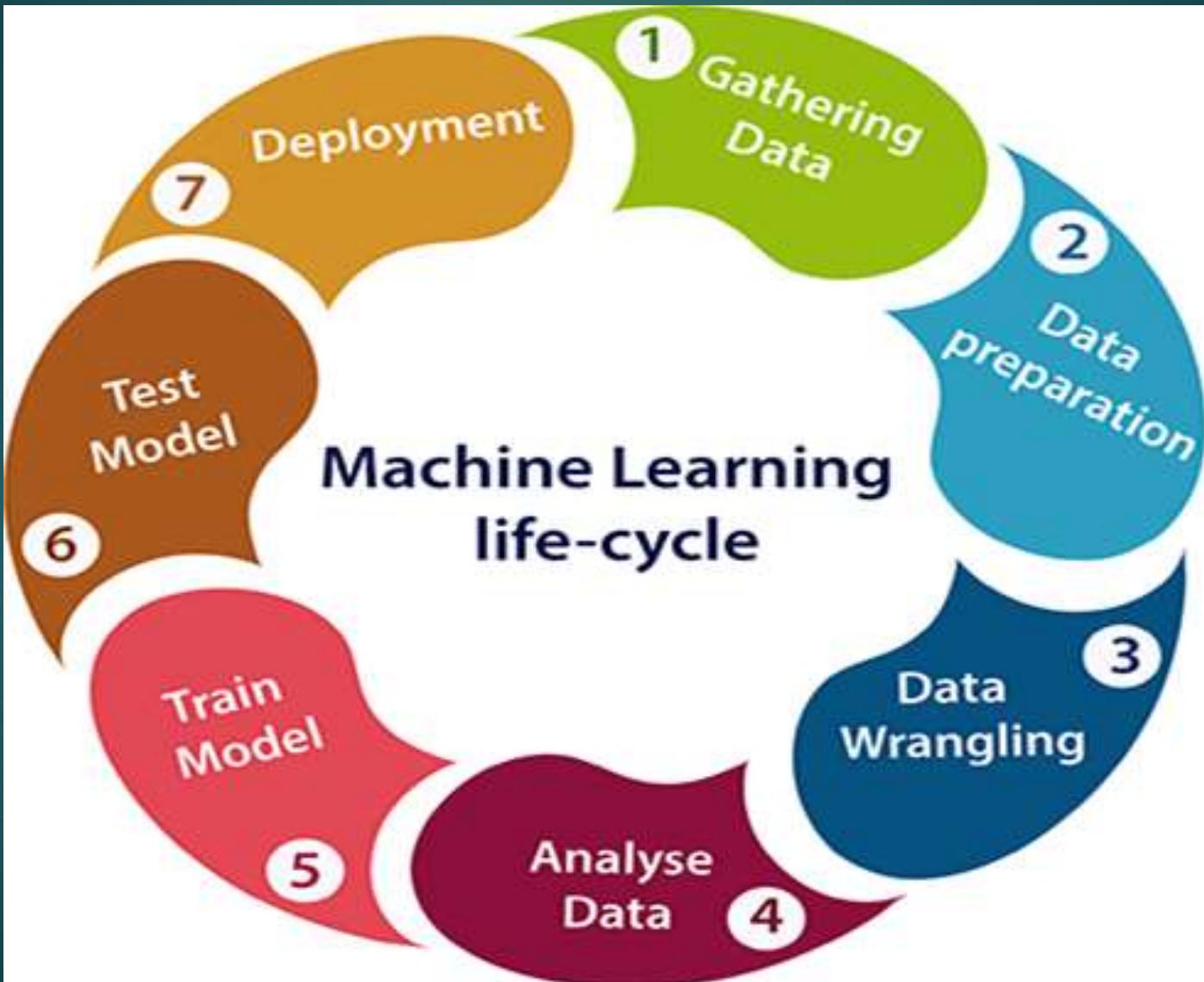
# History of ML:

- The early history of Machine Learning (Pre-1940)
- The era of stored program computers
- Computer machinery and intelligence
- Machine intelligence in Games
- The first "AI" winter
- Machine Learning from theory to reality
- Machine Learning at 21<sup>st</sup> century
- Machine Learning at present

# Applications of ML:



# ML Life cycle:



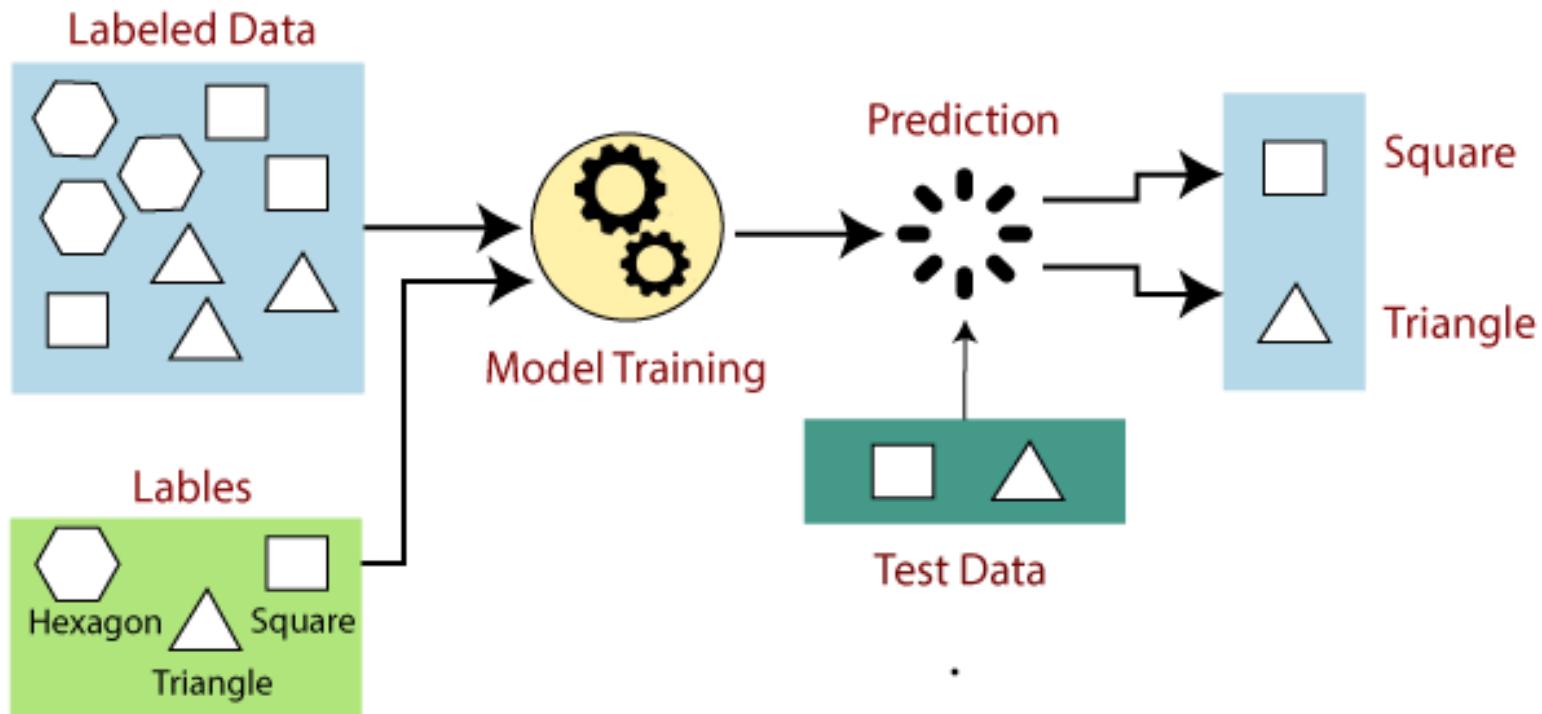
# Artificial intelligence VS Machine learning:

Artificial Intelligence	Machine learning
Artificial intelligence is a technology which enables a machine to simulate human behavior.	Machine learning is a subset of AI which allows a machine to automatically learn from past data without programming explicitly.
The goal of AI is to make a smart computer system like humans to solve complex problems.	The goal of ML is to allow machines to learn from data so that they can give accurate output.
In AI, we make intelligent systems to perform any task like a human.	In ML, we teach machines with data to perform a particular task and give an accurate result.
Machine learning and deep learning are the two main subsets of AI.	Deep learning is a main subset of machine learning.
AI has a very wide range of scope.	Machine learning has a limited scope.
AI is working to create an intelligent system which can perform various complex tasks.	Machine learning is working to create machines that can perform only those specific tasks for which they are trained.
AI system is concerned about maximizing the chances of success.	Machine learning is mainly concerned about accuracy and patterns.
The main applications of AI are <b>Siri, customer support using chatbots, Expert System, Online game playing, intelligent humanoid robot, etc.</b>	The main applications of machine learning are <b>Online recommender system, Google search algorithms, Facebook auto friend tagging suggestions, etc.</b>
On the basis of capabilities, AI can be divided into three types, which are, <b>Weak AI, General AI, and Strong AI.</b>	Machine learning can also be divided into mainly three types that are <b>Supervised learning, Unsupervised learning, and Reinforcement learning.</b>
It includes learning, reasoning, and self-correction.	It includes learning and self-correction when introduced with new data.
AI completely deals with Structured, semi-structured, and unstructured data.	Machine learning deals with Structured and semi-structured data.

# Supervised Machine Learning:

- Supervised learning is the types of machine learning in which machines are trained using well "labelled" training data, and on basis of that data, machines predict the output.
- The aim of a supervised learning algorithm is to **find a mapping function to map the input variable(x) with the output variable(y)**.
- Application:
  - **Risk Assessment**,
  - **Image classification**,
  - **Fraud Detection**,
  - **Spam filtering**, etc.

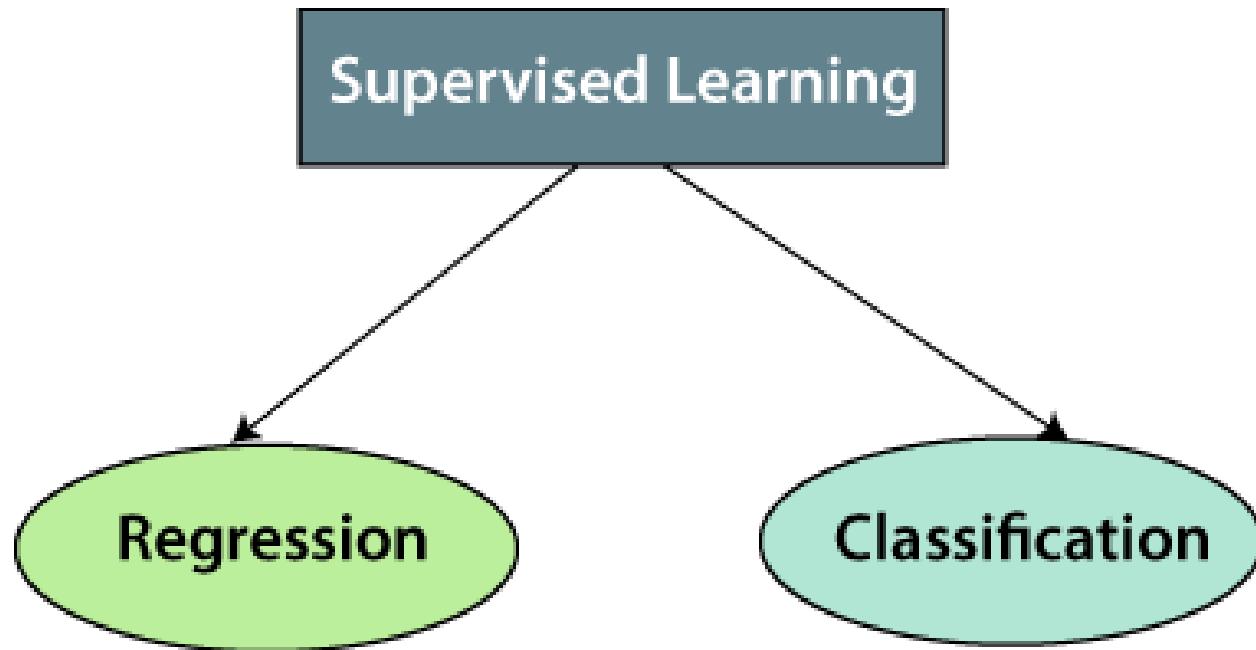
# How Supervised Learning Works:



# Steps Involved in Supervised Learning:

- First Determine the type of training dataset
- Collect/Gather the labelled training data.
- Split the training dataset into training **dataset, test dataset, and validation dataset.**
- Determine the input features of the training dataset, which should have enough knowledge so that the model can accurately predict the output.
- Determine the suitable algorithm for the model, such as support vector machine, decision tree, etc.
- Execute the algorithm on the training dataset. Sometimes we need validation sets as the control parameters, which are the subset of training datasets.
- Evaluate the accuracy of the model by providing the test set. If the model predicts the correct output, which means our model is accurate.

# Types of supervised Machine learning Algorithms:



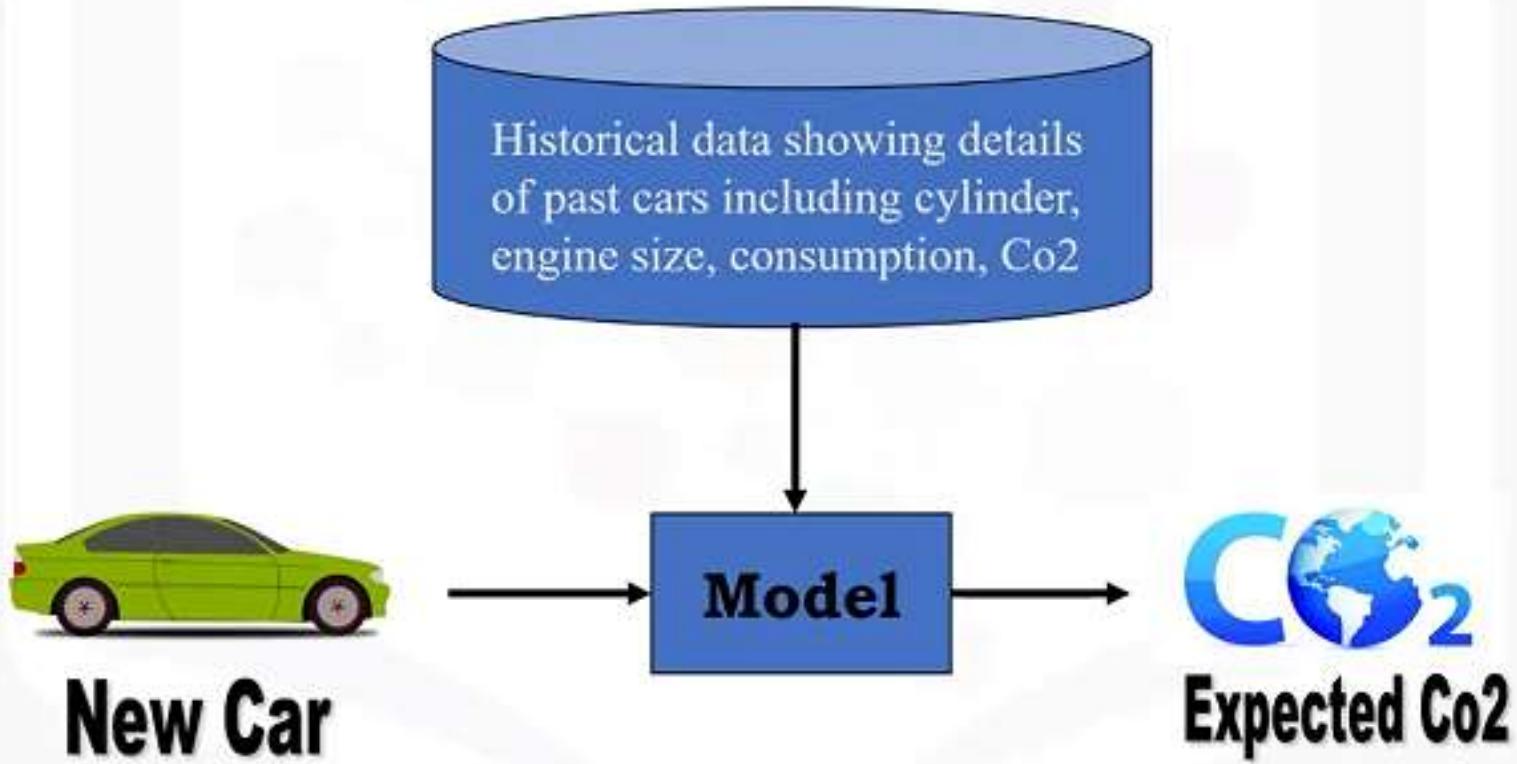
# Introduction to Regression

# What is Regression?

	X: Independent variable			Y: Dependent variable
	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

Regression is the process of predicting a continuous value

# What is a Regression model ?



# Types of Regression

- ❑ Regression algorithms are used if there is a relationship between the input variable and continuous output variable.
- ❑ Types:
  - Simple Linear Regression
  - Multiple Linear Regression

# Types of Regression models ?

- Simple Regression:

- Simple Linear Regression
- Simple Non-linear Regression

Predict `co2emission` vs `EngineSize` of all cars

- Multiple Regression:

- Multiple Linear Regression
- Multiple Non-linear Regression

Predict `co2emission` vs `EngineSize` and `Cylinders` of all cars

# Applications

- Sales forecasting
- Satisfaction analysis
- Price estimation
- Employment income

# Using Linear Regression to predict Continuous Values:

X: Independent variable

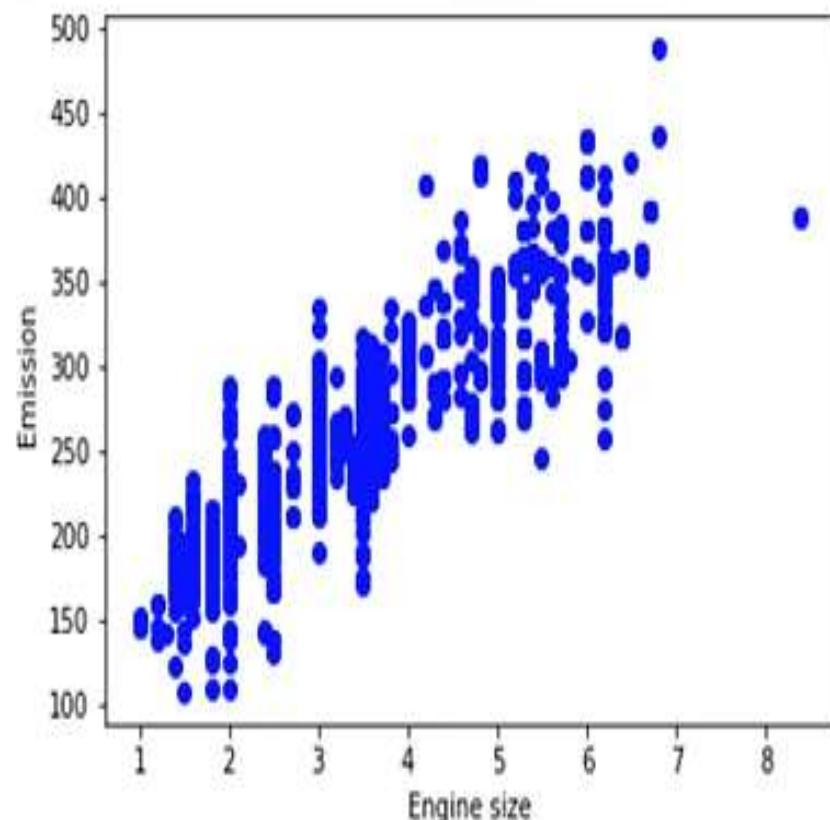
Y: Dependent variable

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

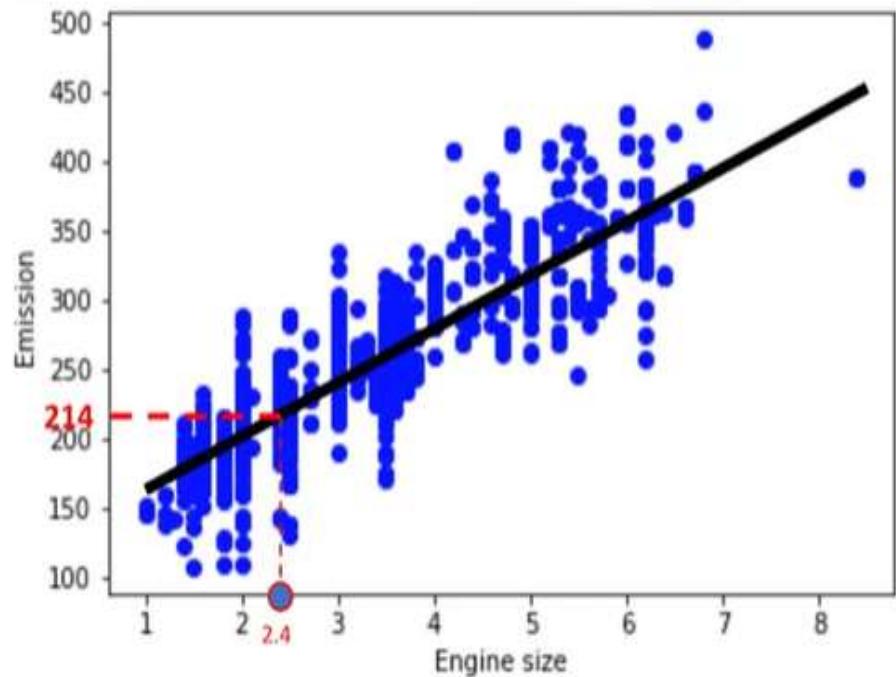
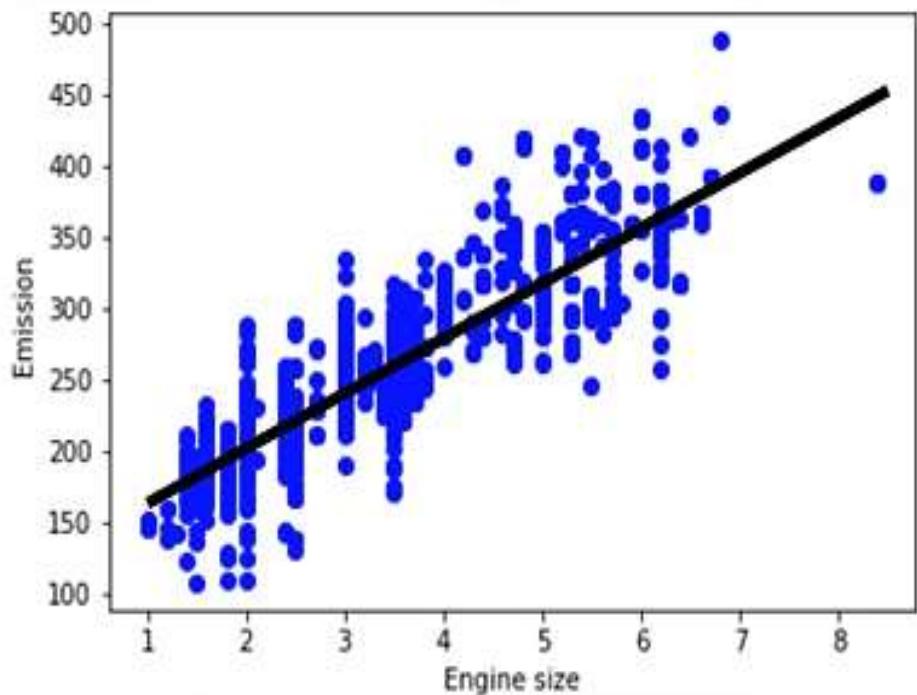
Continuous Values

# How does Linear Regression Work?

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?



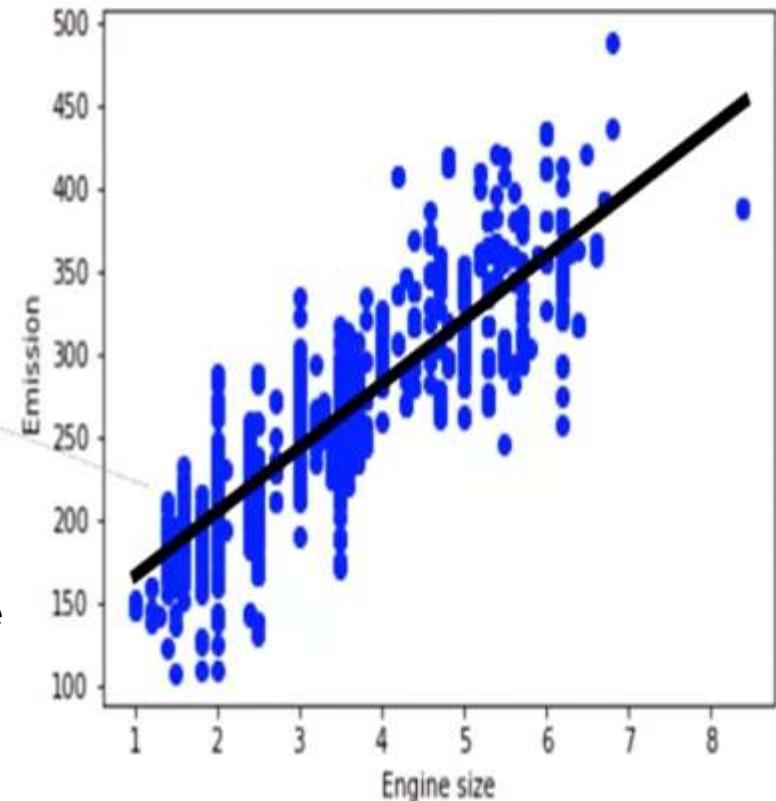
# How does Linear Regression Work?



# Linear Regression Model Representation:

$$\hat{y} = \theta_0 + \theta_1 x_1$$

- **y hat** is the dependent variable
- **x<sub>1</sub>** is the independent variable
- **Theta 0 and theta 1** are the parameters or coefficients of the line
- **Theta 1** is known as the slope or gradient
- **theta 0** is known as the intercept



x<sub>1</sub>

Activate Windows

# How to find the best fit:

$x_1 = 5.4$  independent variable

$y = 250$  actual Co2 emission of  $x_1$

$$\hat{y} = \theta_0 + \theta_1 x_1$$

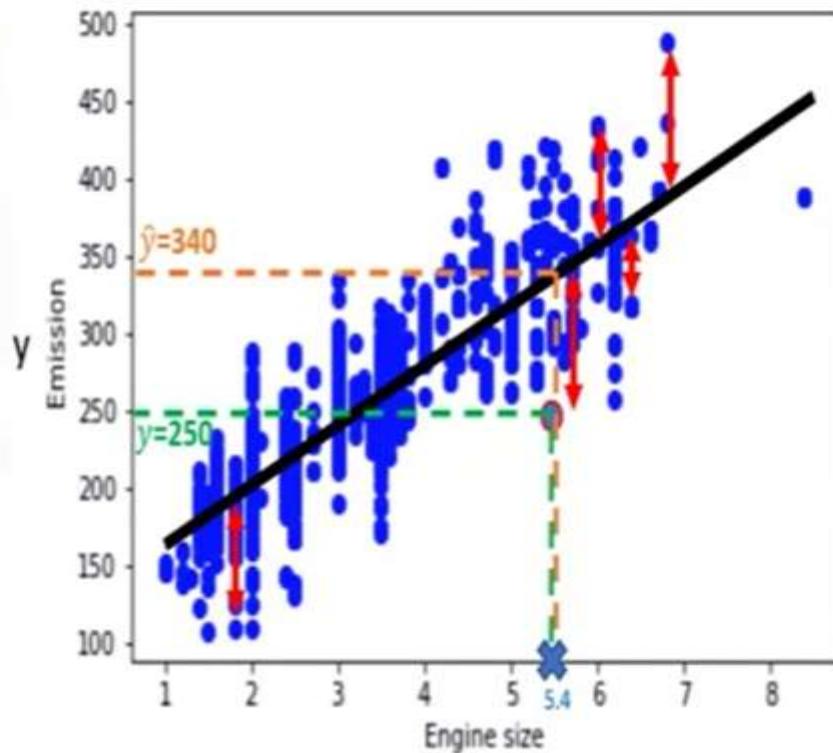
$\hat{y} = 340$  the predicted emission of  $x_1$

$$\text{Error} = y - \hat{y}$$

$$= 250 - 340$$

$$= -90$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



# Estimating parameters using mathematical approach:

$$\hat{y} = \theta_0 + \theta_1 x_1$$

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267

$$\theta_1 = \frac{\sum_{i=1}^s (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^s (x_i - \bar{x})^2}$$

$$\bar{x} = (2.0 + 2.4 + 1.5 + \dots) / 9 = 3.03$$

$$\bar{y} = (196 + 221 + 136 + \dots) / 9 = 226.22$$

$$\theta_1 = \frac{(2.0 - 3.03)(196 - 226.22) + (2.4 - 3.03)(221 - 226.22) + \dots}{(2.0 - 3.03)^2 + (2.4 - 3.03)^2 + \dots}$$

$$\theta_1 = 39$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\theta_0 = 226.22 - 39 * 3.03$$

$$\theta_0 = 125.74$$

$$\hat{y} = 125.74 + 39x_1$$

# Predictions with Simple Linear Regression:

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$$Co2Emission = \theta_0 + \theta_1 EngineSize$$

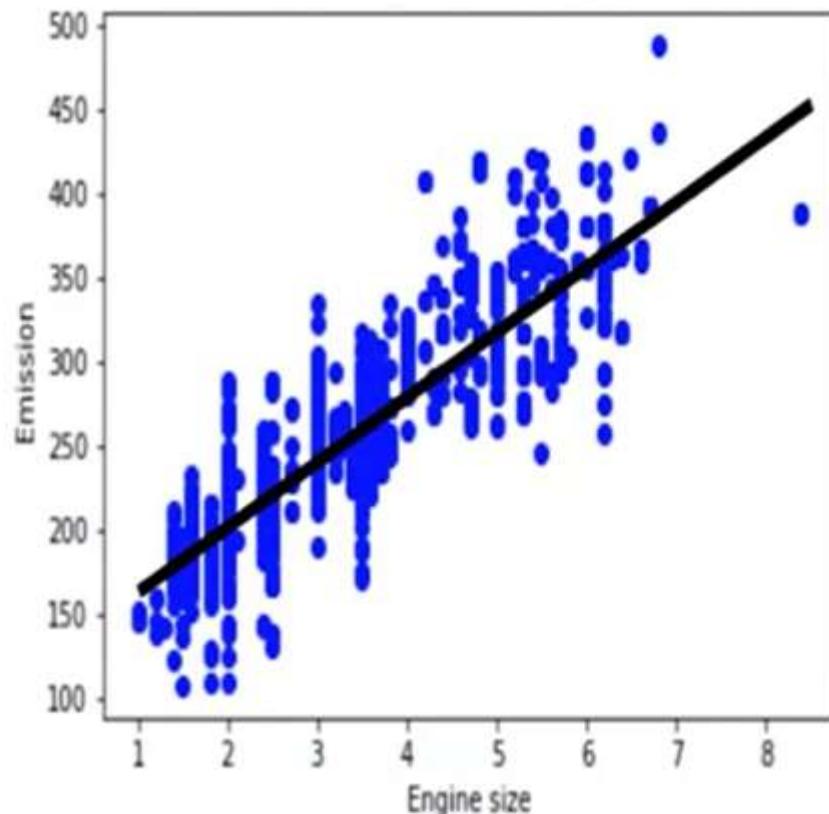
$$Co2Emission = 125 + 39 \cdot EngineSize$$

$$Co2Emission = 125 + 39 \times 2.4$$

$$Co2Emission = 218.6$$

# Pros of Simple Linear Regression:

- Very fast
- No parameter tuning
- Easy to understand, and highly interpretable



# Types of Linear Regression:

- Simple Linear Regression
  - Predict Co2emission vs EngineSize of all cars
    - Independent variable (x): EngineSize
    - Dependent variable (y): Co2emission
- • Multiple Linear Regression
  - Predict Co2emission vs EngineSize and Cylinders of all cars
    - Independent variable (x): EngineSize, Cylinders, etc.
    - Dependent variable (y): Co2emission

# Examples of Multiple Linear Regression:

- Independent variables effectiveness on prediction
  - Does revision time, test anxiety, lecture attendance and gender have any effect on the exam performance of students?
- • Predicting impacts of changes
  - How much does blood pressure go up (or down) for every unit increase (or decrease) in the BMI of a patient?

# Predicting continuous values with Multiple Linear Regression:

$$Co2\ Em = \theta_0 + \theta_1 \text{Engine size} + \theta_2 \text{Cylinders} + \dots$$

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\hat{y} = \theta^T X$$

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots] \quad X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \end{bmatrix}$$

X: Independent variable

Y: Dependent variable

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4		

Activate Windows

?

Go to Settings to activate Windows

# Using MSE to expose errors in the model:

$$\hat{y} = \theta^T X$$

$\hat{y}_i = 140$  the predicted emission of  $x_i$

$y_i = 196$  actual value of  $x_i$

$y_i - \hat{y}_i = 196 - 140 = 56$  residual error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267

# Estimating multiple linear regression parameters:

- How to estimate  $\theta$ ?
  - Ordinary Least Squares
    - Linear algebra operations
    - Takes a long time for large datasets (10K+ rows)
  - An optimization algorithm
    - Gradient Descent
    - Proper approach if you have a very large dataset

# Making predictions using multiple linear regression:

$$\hat{y} = \theta^T X$$

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

$$\theta^T = [125, 6.2, 14, \dots]$$

$$\hat{y} = 125 + 6.2x_1 + 14x_2 +$$

$$Co2Em = 125 + 6.2EngSize + 14 \text{ Cylinders} + \dots$$

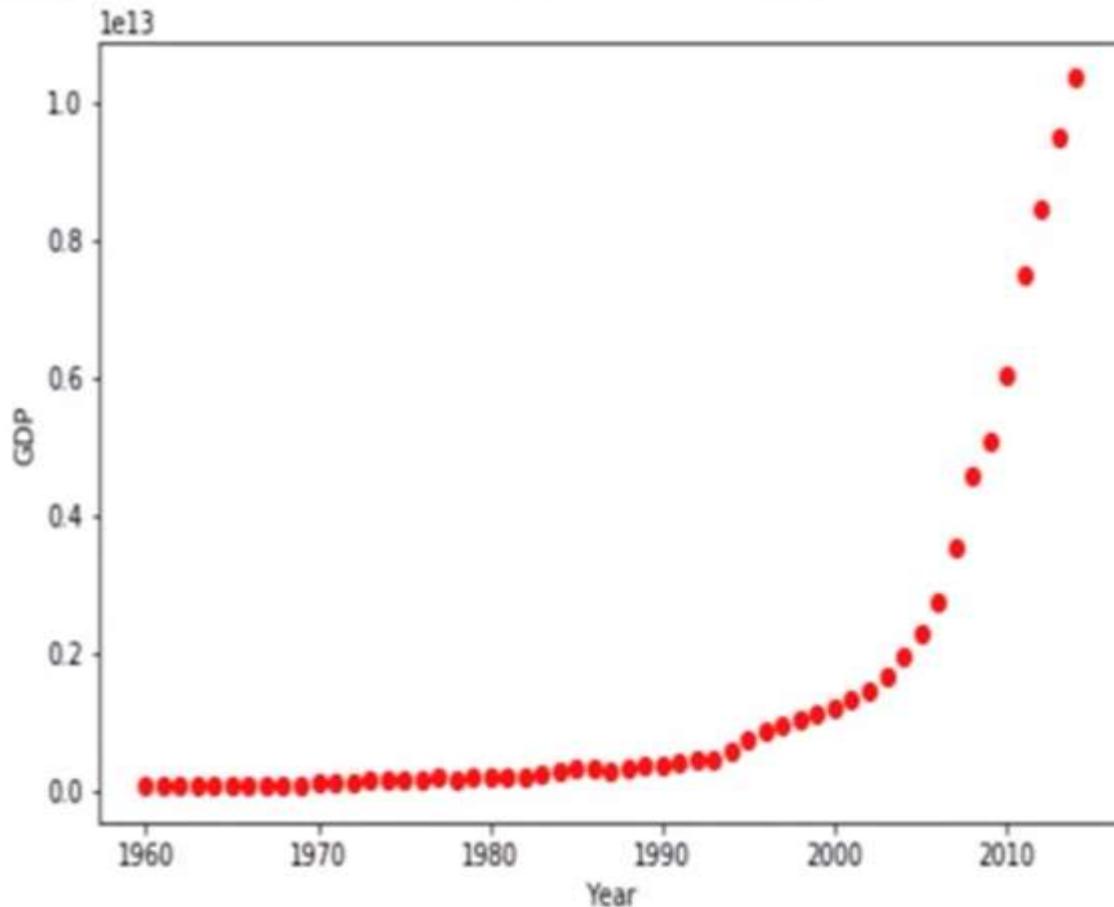
$$Co2Em = 125 + 6.2 \times 2.4 + 14 \times 4 + \dots$$

$$Co2Em = 214.1$$

# Non-Linear Regression:

- Should we use Linear Regression always?

Year	Value
0 1960	5.918412e+10
1 1961	4.955705e+10
2 1962	4.668518e+10
3 1963	5.009730e+10
4 1964	5.906225e+10
5 1965	6.970915e+10
6 1966	7.587943e+10
7 1967	7.205703e+10
8 1968	6.999350e+10
9 1969	7.871882e+10
...	.....

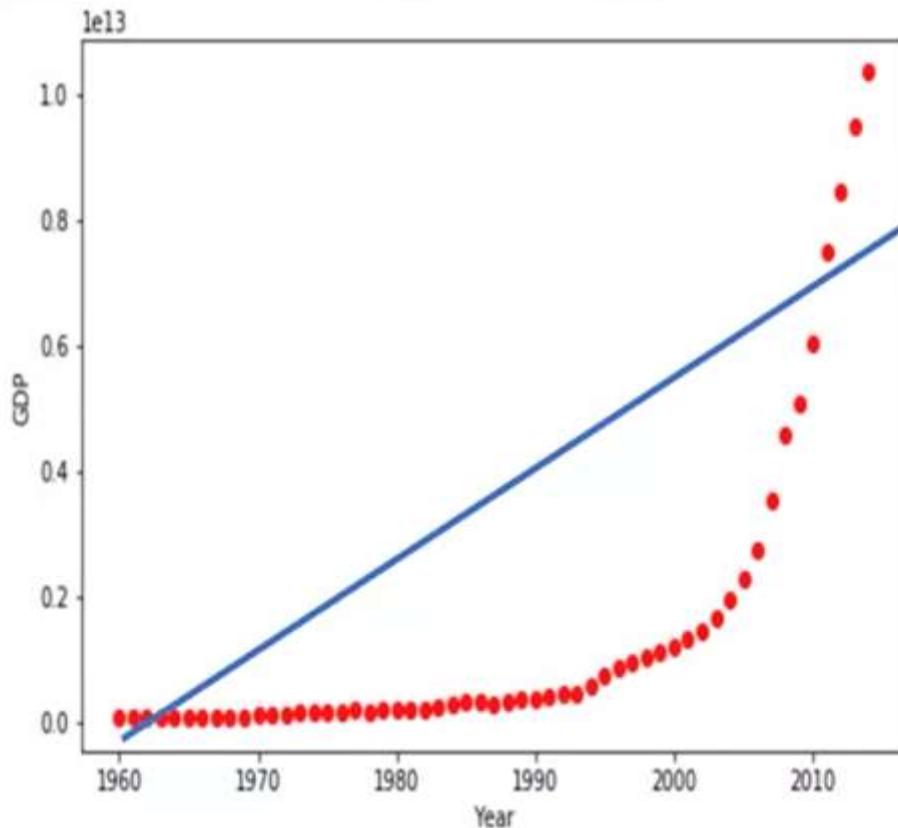


# Non-Linear Regression:

- Should we use Linear Regression always?

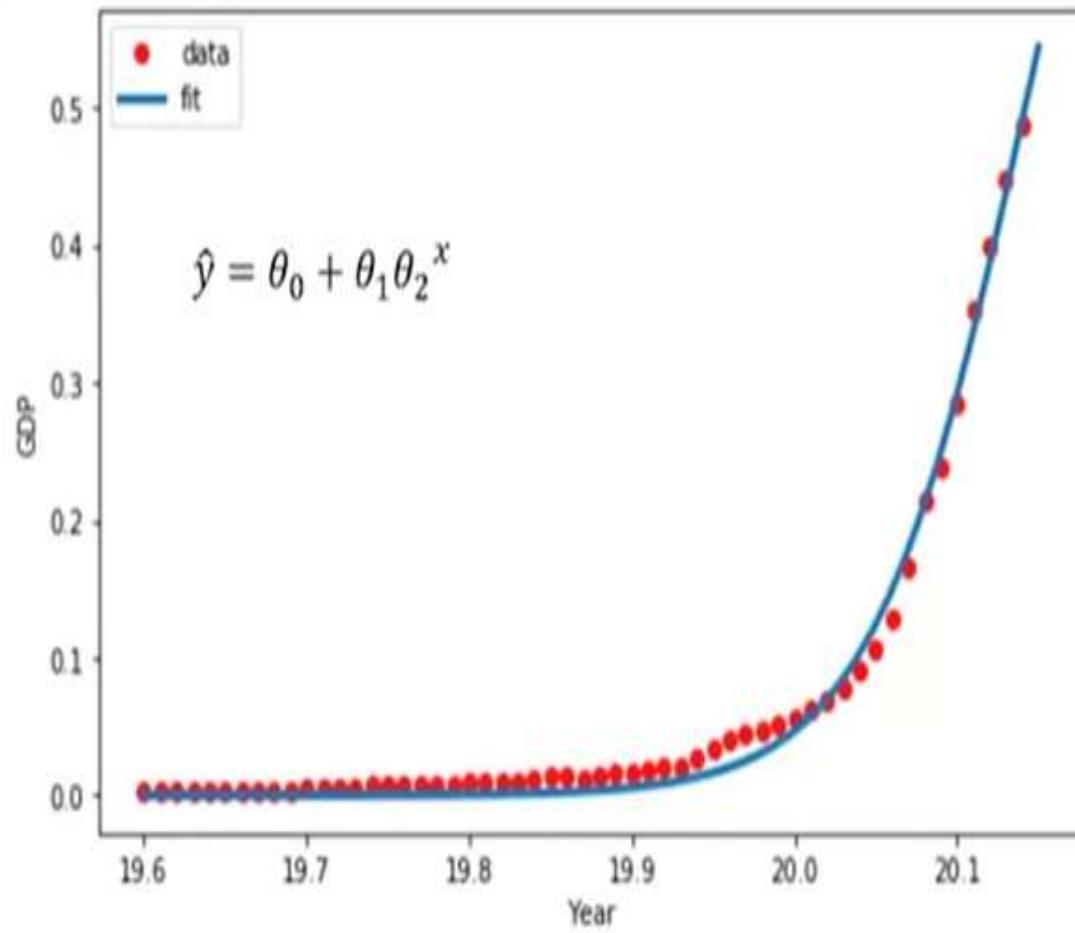
## Should we use linear regression?

Year	Value
0 1960	5.918412e+10
1 1961	4.955705e+10
2 1962	4.668518e+10
3 1963	5.009730e+10
4 1964	5.906225e+10
5 1965	6.970915e+10
6 1966	7.587943e+10
7 1967	7.205703e+10
8 1968	6.999350e+10
9 1969	7.871882e+10



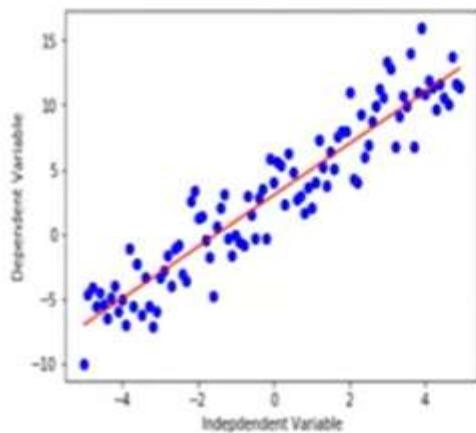
# Non-Linear Regression:

	Year	Value
0	1960	5.918412e+10
1	1961	4.955705e+10
2	1962	4.668518e+10
3	1963	5.009730e+10
4	1964	5.906225e+10
5	1965	6.970915e+10
6	1966	7.587943e+10
7	1967	7.205703e+10
8	1968	6.999350e+10
9	1969	7.871882e+10
...	.....	.....

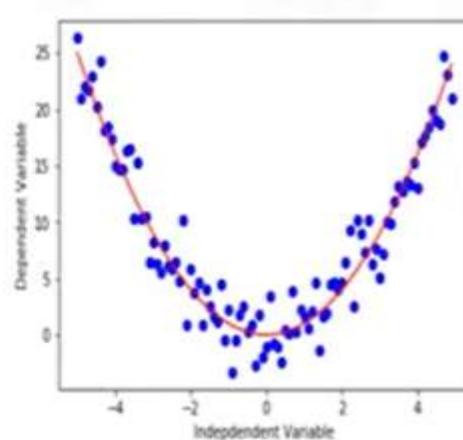


# Different types of Regression:

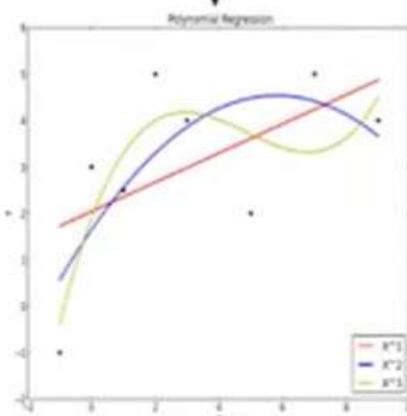
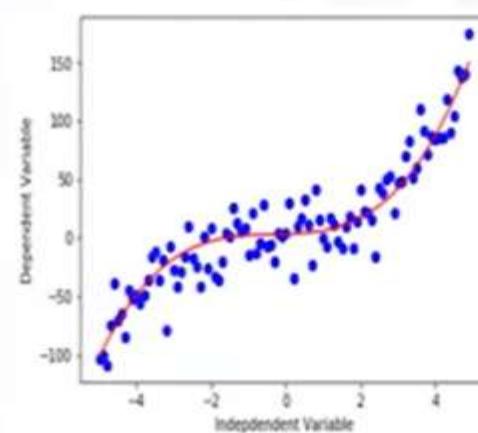
Linear Regression



Quadratic (Parabolic) Regression



Cubic Regression



# Linear Vs Non-Linear Regression:

- How can I know if a problem is linear or non-linear in an easy way?
  - Inspect visually
  - Based on accuracy
- How should I model my data, if it displays non-linear on a scatter plot?
  - Polynomial regression
  - Non-linear regression model
  - Transform your data

# What is Classification?

- A supervised learning approach
- Categorizing some unknown items into a discrete set of categories or “classes”
- The target attribute is a categorical variable

# How does Classification work?

Classification determines the class label for an unlabeled test case.

age	ed	employ	address	income	debtinc	creddebt	othdebt	default
41	3	17	12	176	9.3	11.359	5.009	1
27	1	10	6	31	17.3	1.362	4.001	0
40	1	15	14	55	5.5	0.856	2.169	0
41	1	15	14	120	2.9	2.659	0.821	0
24	2	2	0	28	17.3	1.787	3.057	1
41	2	5	5	25	10.2	0.393	2.157	0
39	1	20	9	67	30.6	3.834	16.668	0
43	1	12	11	38	3.6	0.129	1.239	0
24	1	3	4	19	24.4	1.358	3.278	1
36	1	0	13	25	19.7	2.778	2.147	0

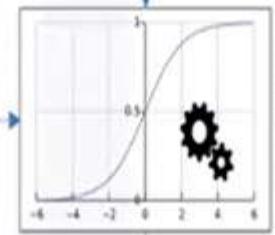
Categorical Variable

Loan default prediction

Modeling

age	ed	employ	address	income	debtinc	creddebt	othdebt	default
37	2	16	10	130	9.3	10.23	3.21	0

Prediction



Classifier

Predicted Labels

# Multi-class classification:

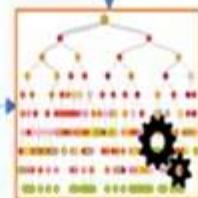
Age	Sex	BP	Cholesterol	Na	K	Drug
23	F	HIGH	HIGH	0.793	0.031	drugY
47	M	LOW	HIGH	0.739	0.056	drugC
47	M	LOW	HIGH	0.697	0.069	drugC
28	F	NORMAL	HIGH	0.564	0.072	drugX
61	F	LOW	HIGH	0.559	0.031	drugY
22	F	NORMAL	HIGH	0.677	0.079	drugX
49	F	NORMAL	HIGH	0.79	0.049	drugY
41	M	LOW	HIGH	0.767	0.069	drugC
60	M	NORMAL	HIGH	0.777	0.051	drugY
43	M	LOW	NORMAL	0.526	0.027	drugY

Categorical Variable

Modeling

Age	Sex	BP	Cholesterol	Na	K	Drug
36	F	LOW	HIGH	0.697	0.069	DrugX

Prediction

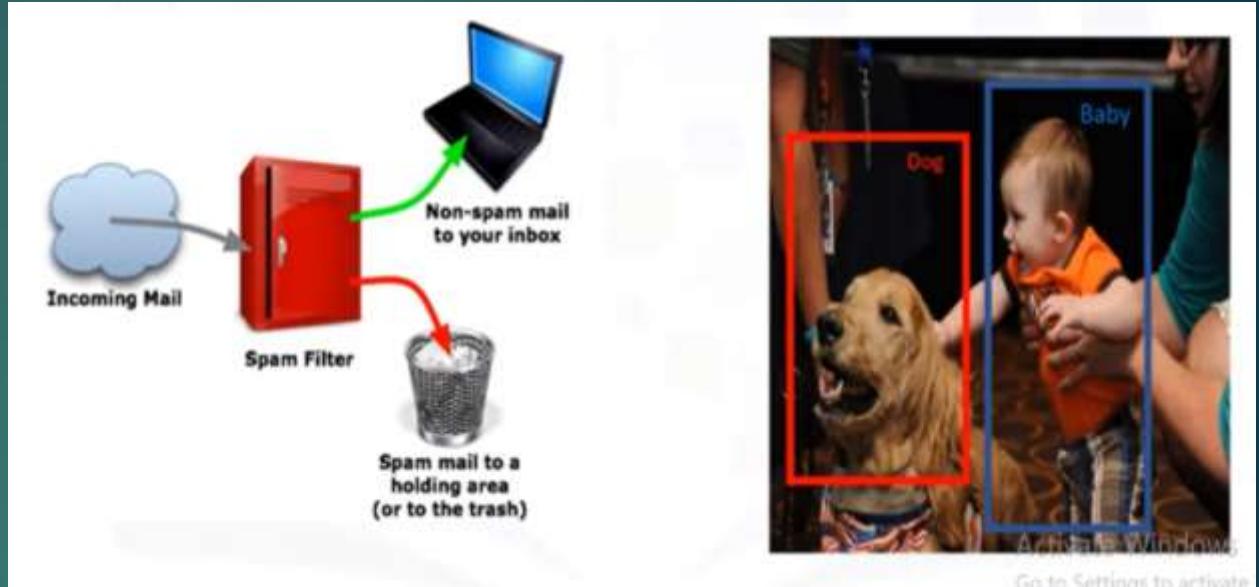


Classifier

Predicted Labels

# Classification Application

- Email filtering
- speech recognition
- Handwriting recognition
- Biometric identification
- Document classification etc.



# Classification Types

- ❑ Logistic Regression
- ❑ Decision Trees
- ❑ Support vector Machines
- ❑ Naive bayes,
- ❑ Linear discriminant analysis (LDA)
- ❑ k-nearest neighbor
- ❑ Neural networks

# Logistic Regression

Logistic regression is a classification algorithm for categorical variables.

	Independent variables										Dependent variable
	tenure	age	address	income	ed	employ	equip	callcard	wireless	churn	
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	Yes	
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	Yes	
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	No	
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	No	
4	7.0	35.0	14.0	80.0	2.0	15.0	0.0	1.0	0.0	?	

Continuous/Categorical variables

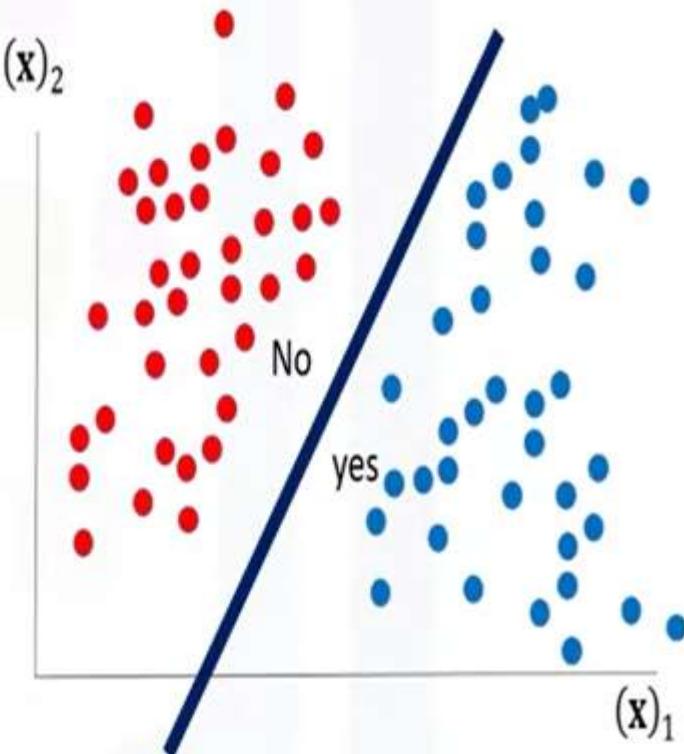
Categorical Variable

# Logistic Regression Applications

- Predicting the probability of a person having a heart attack
- Predicting the mortality in injured patients
- Predicting a customer's propensity to purchase a product or halt a subscription
- Predicting the probability of failure of a given process or product
- Predicting the likelihood of a homeowner defaulting on a mortgage

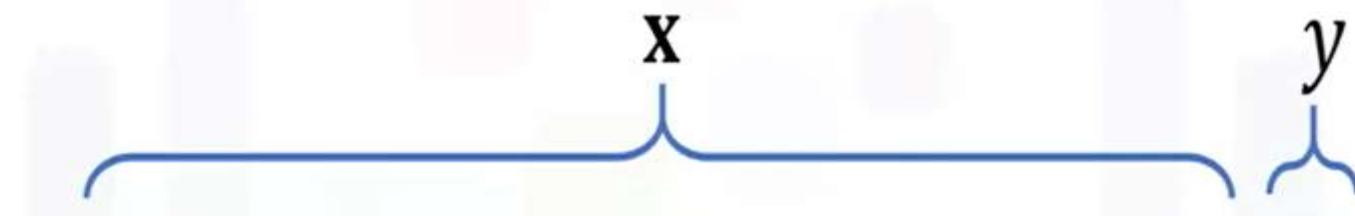
# When is Logistic Regression suitable?

- If your data is binary
  - 0/1, YES/NO, True/False
- If you need probabilistic results
- When you need a linear decision boundary
- If you need to understand the impact of a feature



$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 > 0$$

# Building a Model for customer churn:



The diagram shows a horizontal bracket above the feature names labeled 'X', and another bracket to the right of the target variable 'y'.

	tenure	age	address	income	ed	employ	equip	callcard	wireless	churn
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	1.0
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	1.0
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	0.0
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	0.0

$$X \in \mathbb{R}^{m \times n}$$
$$y \in \{0,1\}$$

$$\hat{y} = P(y=1|x)$$

$$P(y=0|x) = 1 - P(y=1|x)$$

Activate Windows

Go to Settings to activate Windows.

# Logistic Regression vs Linear Regression:

The diagram illustrates the inputs and output for a logistic regression model. On the left, a horizontal bracket labeled 'X' covers the columns of numerical values. On the right, a vertical bracket labeled 'y' covers the final column, which contains categorical labels.

	tenure	age	address	income	ed	employ	equip	callcard	wireless	churn
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	1
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	1
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	0
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	0
4	7.0	35.0	14.0	80.0	2.0	15.0	0.0	1.0	0.0	0

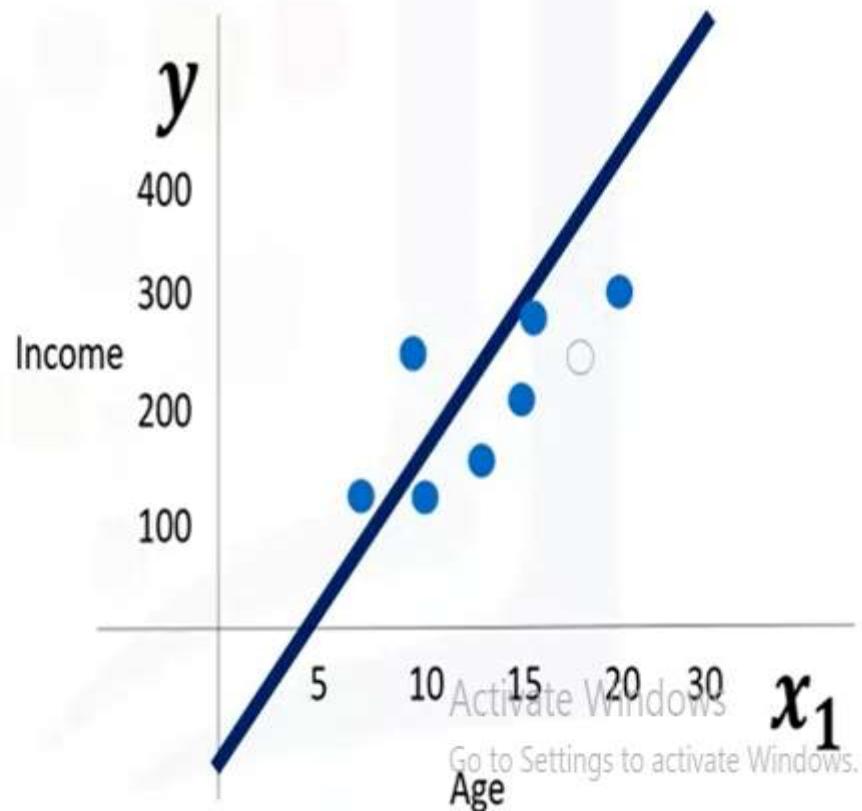
$$\hat{y} = P(y=1|x)$$

Activate Windows

Go to Settings to activate Windows

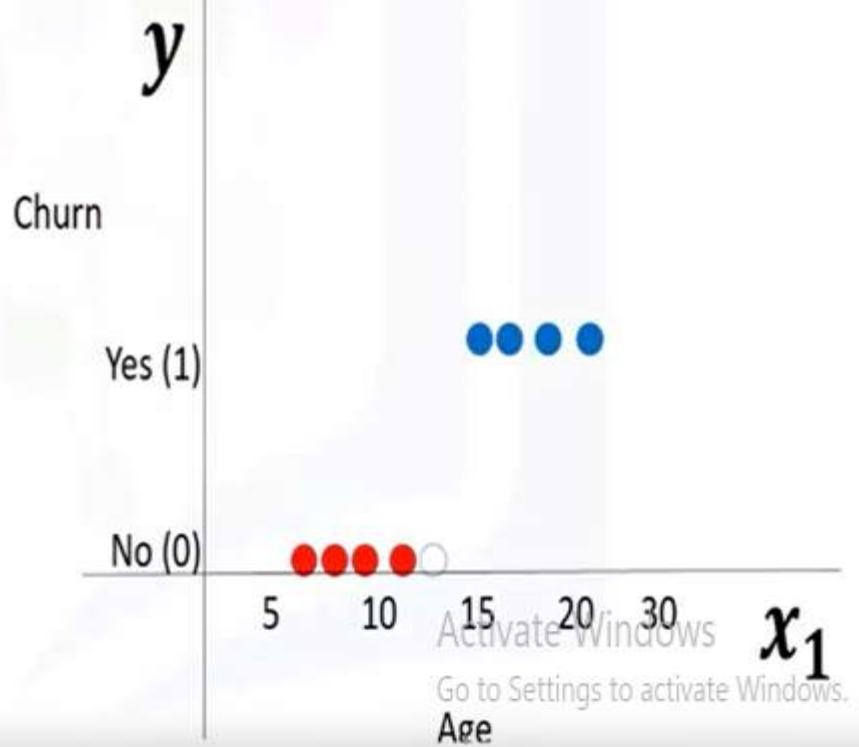
# Predicting a customer income:

	tenure	age	address	income	ed	employ	equip	callcard	wireless	churn
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	1
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	1
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	0
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	0
4	7.0	35.0	14.0	80.0	2.0	15.0	0.0	1.0	0.0	0



# Predicting churn using Linear Regression:

	tenure	age	address	income	ed	employ	equip	callcard	wireless	churn
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	1
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	1
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	0
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	0
4	7.0	35.0	14.0	80.0	2.0	15.0	0.0	1.0	0.0	0



# Predicting churn using Linear Regression:

$$\theta^T = [\theta_0, \theta_1]$$

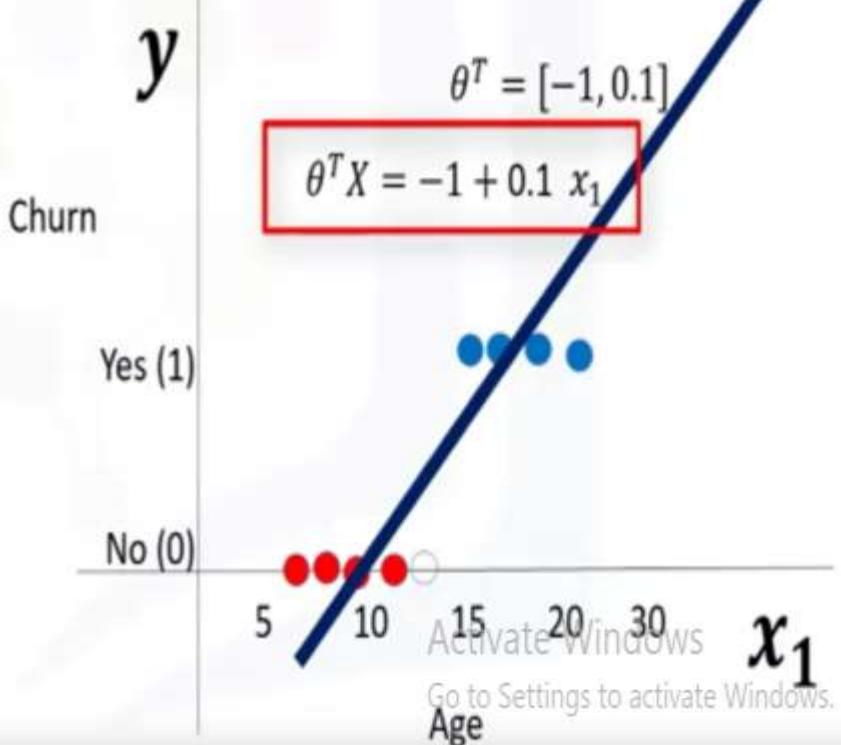
$$\theta_0 + \theta_1 x_1$$

$$a + b x_1$$

$$\theta^T X = \theta_0 + \theta_1 x_1$$

$$\theta^T X = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$$

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots] \quad X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \end{bmatrix}$$



# Linear Regression in classification problems:

$$\theta^T X = \theta_0 + \theta_1 x_1$$

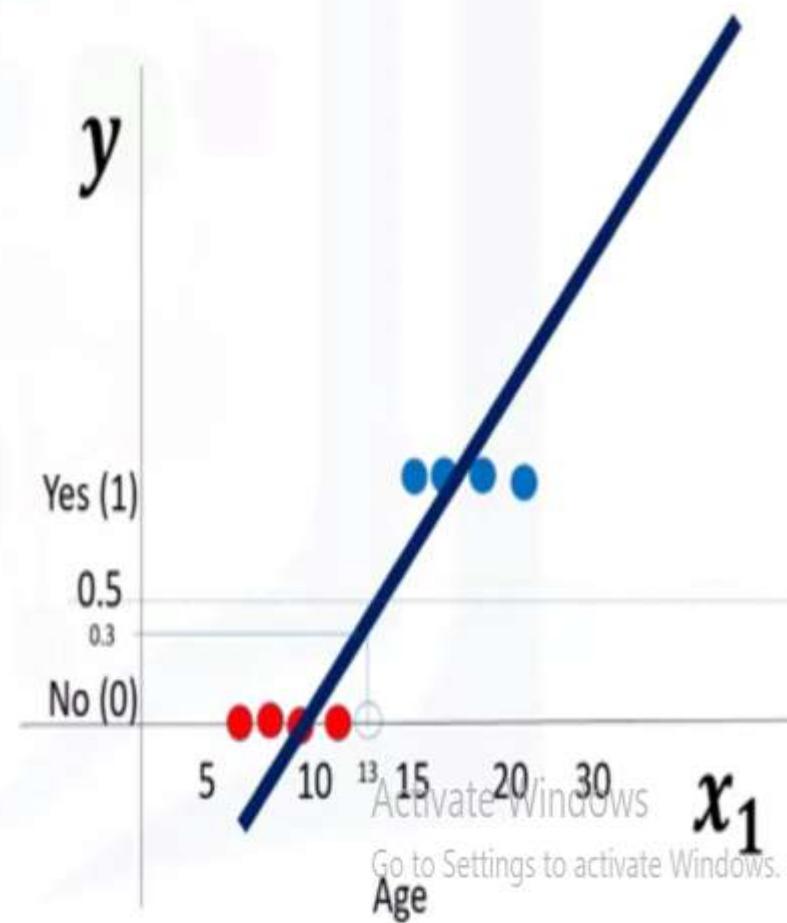
$$\theta^T X = -1 + 0.1 \cdot x_1$$

$$p_1 = [13] \rightarrow \theta^T X = -1 + 0.1 \cdot x_1 \\ = -1 + 0.1 \times 13 \\ = 0.3$$

$$\hat{y} = \begin{cases} 0 & \text{if } \theta^T X < 0.5 \\ 1 & \text{if } \theta^T X \geq 0.5 \end{cases}$$

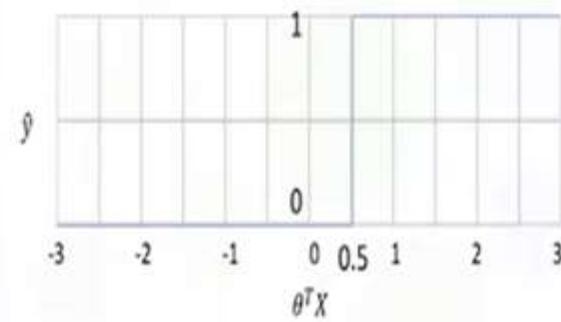
$$\theta^T X = 0.3$$

$\theta^T X < 0.5 \rightarrow \text{Class 0}$



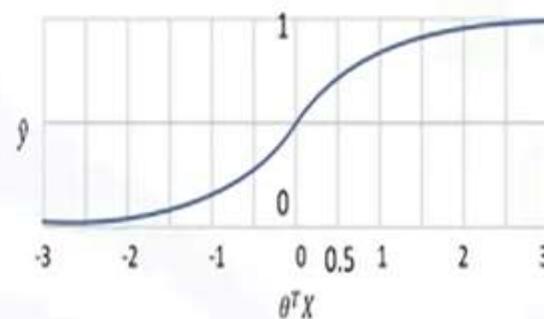
# The Problem with using Linear Regression:

$$\theta^T X = \theta_0 + \theta_1 x_1 + \dots$$



$$\hat{y} = \begin{cases} 0 & \text{if } \theta^T X < 0.5 \\ 1 & \text{if } \theta^T X \geq 0.5 \end{cases}$$

$$\sigma(\theta^T X) = \sigma(\theta_0 + \theta_1 x_1 + \dots)$$



$$\hat{y} = \sigma(\theta^T X)$$

P(y=1|x)

Activate Windows

# Sigmoid function in Logistic Regression:

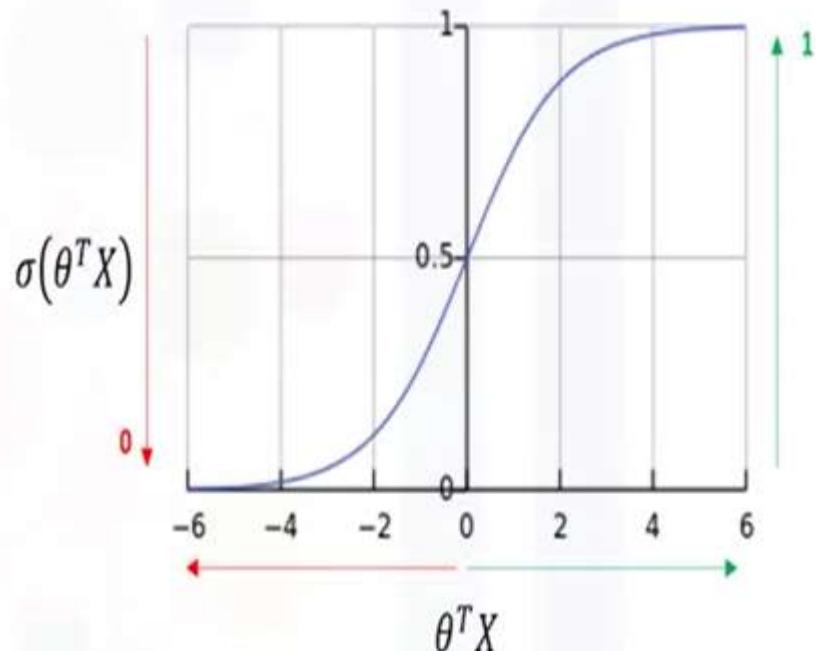
- Logistic Function

$$\sigma(\theta^T X) = \frac{1}{1 + e^{-\theta^T X}}$$

$$\sigma(\theta^T X) = 1$$

$$\sigma(\theta^T X) = 0$$

[0, 1]



$P(y=1|x)$



$P(y=1|x)$

# Clarification of the customer churn model:

What is the output of our model?

- $P(Y=1|X)$
  - $P(y=0|x) = 1 - P(y=1|x)$
- 
- $P(\text{Churn}=1|\text{income,age}) = 0.8$
  - $P(\text{Churn}=0|\text{income,age}) = 1 - 0.8 = 0.2$

$$\sigma(\theta^T X) \rightarrow P(y=1|x)$$

$$1 - \sigma(\theta^T X) \rightarrow P(y=0|x)$$

# The Initial training Process:

$$\sigma(\theta^T X) \rightarrow P(y=1|x)$$

1. Initialize  $\theta$ .  
 $\theta = [-1, 2]$
2. Calculate  $\hat{y} = \sigma(\theta^T X)$  for a customer.  
 $\hat{y} = \sigma([-1, 2] \times [2, 5]) = 0.7$
3. Compare the output of  $\hat{y}$  with actual output of customer,  $y$ , and record it as error.  
Error = 1 - 0.7 = 0.3
4. Calculate the error for all customers.  
 $Cost = J(\theta)$
5. Change the  $\theta$  to reduce the cost.  
 $\theta_{new}$
6. Go back to step 2.

# The actual training Process:

General Cost Function:

$$\sigma(\theta^T X) \longrightarrow P(y=1|x)$$

- Change the weight -> Reduce the cost

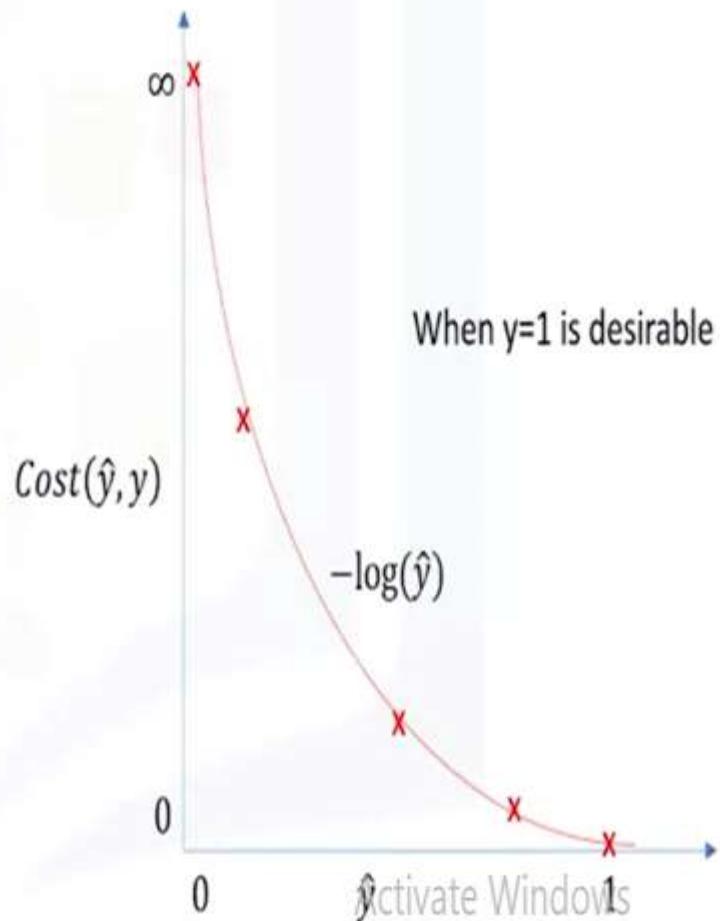
- Cost function

$$Cost(\hat{y}, y) = \frac{1}{2}(\sigma(\theta^T X) - y)^2$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(\hat{y}, y)$$

# Plotting the Cost Function of the Model:

- Model  $\hat{y}$
- Actual Value  $y=1$  or  $0$
- If  $Y=1$ , and  $\hat{y}=1 \rightarrow \text{cost} = 0$
- If  $Y=1$ , and  $\hat{y}=0 \rightarrow \text{cost} = \text{large}$



# Logistic Regression Cost Function:

- So, we will replace cost function with:

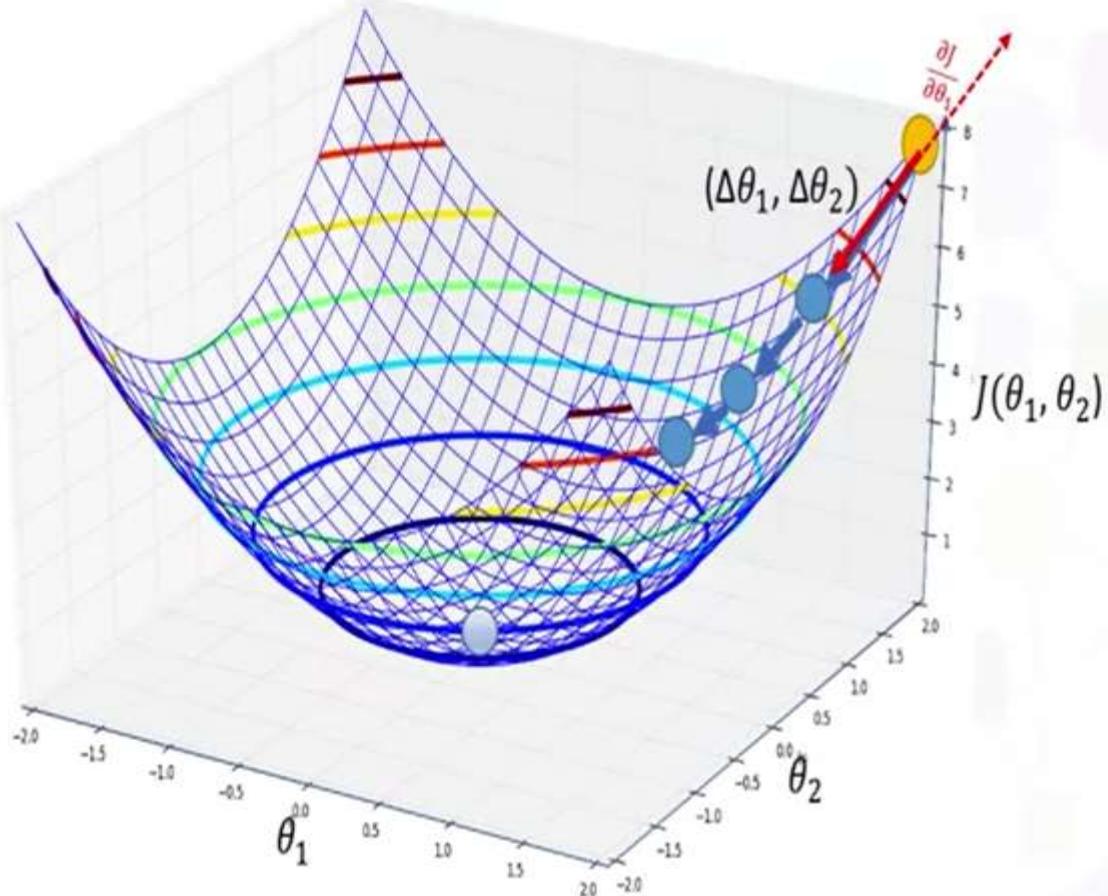
$$Cost(\hat{y}, y) = \frac{1}{2} (\sigma(\theta^T X) - y)^2$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(\hat{y}, y)$$

$$Cost(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)$$

# Using Gradient Descent to minimize the Cost:



$$\hat{y} = \sigma(\theta_1 x_1 + \theta_2 x_2)$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)$$

$$\frac{\partial J}{\partial \theta_1} = -\frac{1}{m} \sum_{i=1}^m (y^i - \hat{y}^i) x_1^i$$
$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \\ \frac{\partial J}{\partial \theta_3} \\ \vdots \\ \frac{\partial J}{\partial \theta_k} \end{bmatrix}$$

$$\text{New } \theta = \text{old } \theta - \eta \nabla J$$

Activate Windows

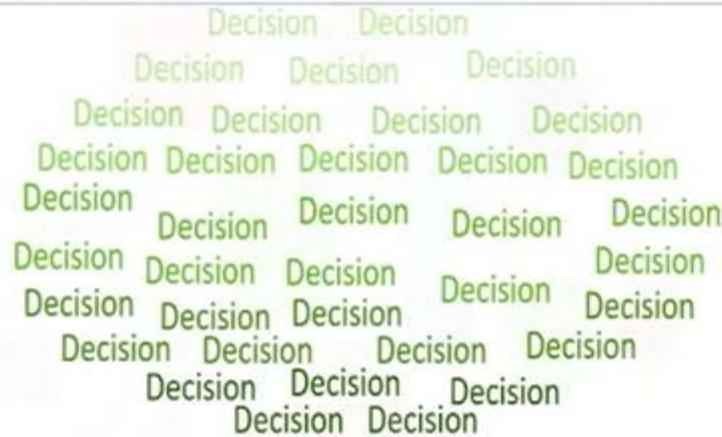
Go to Settings to activate Windows.

# Training Algorithm Recap:

1. initialize the parameters randomly.  
 $\theta^T = [\theta_0, \theta_1, \theta_2, \dots]$
2. Feed the cost function with training set, and calculate the error.  
$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)$$
3. Calculate the gradient of cost function.  
$$\nabla J = \left[ \frac{\partial J}{\partial \theta_1}, \frac{\partial J}{\partial \theta_2}, \frac{\partial J}{\partial \theta_3}, \dots, \frac{\partial J}{\partial \theta_k} \right]$$
4. Update weights with new values.  
$$\theta_{new} = \theta_{prev} - \eta \nabla J$$
5. Go to step 2 until cost is small enough.
6. Predict the new customer X.  
$$P(y=1|x) = \sigma(\theta^T X)$$

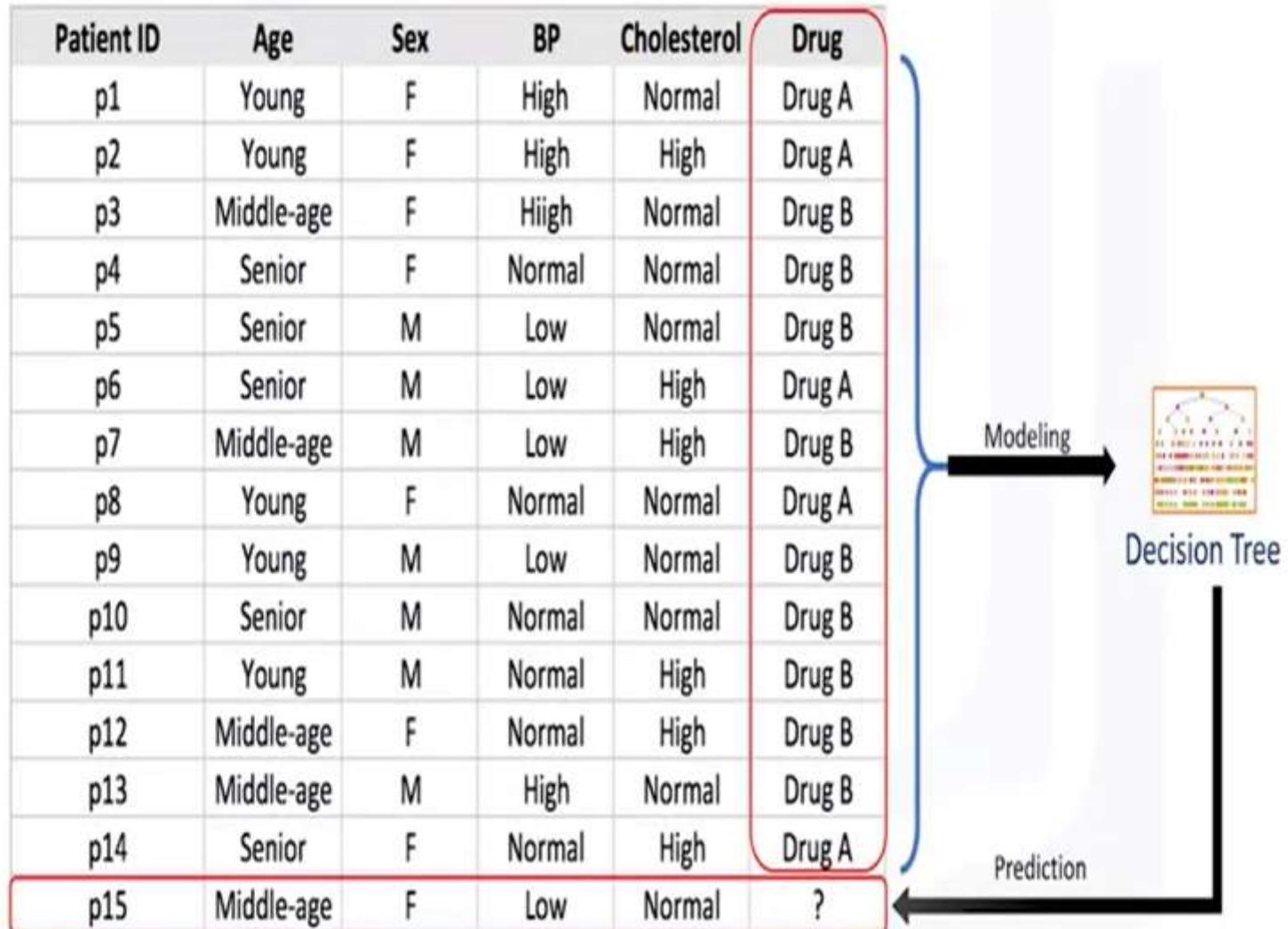
# Decision Tree:

## What is a decision tree?



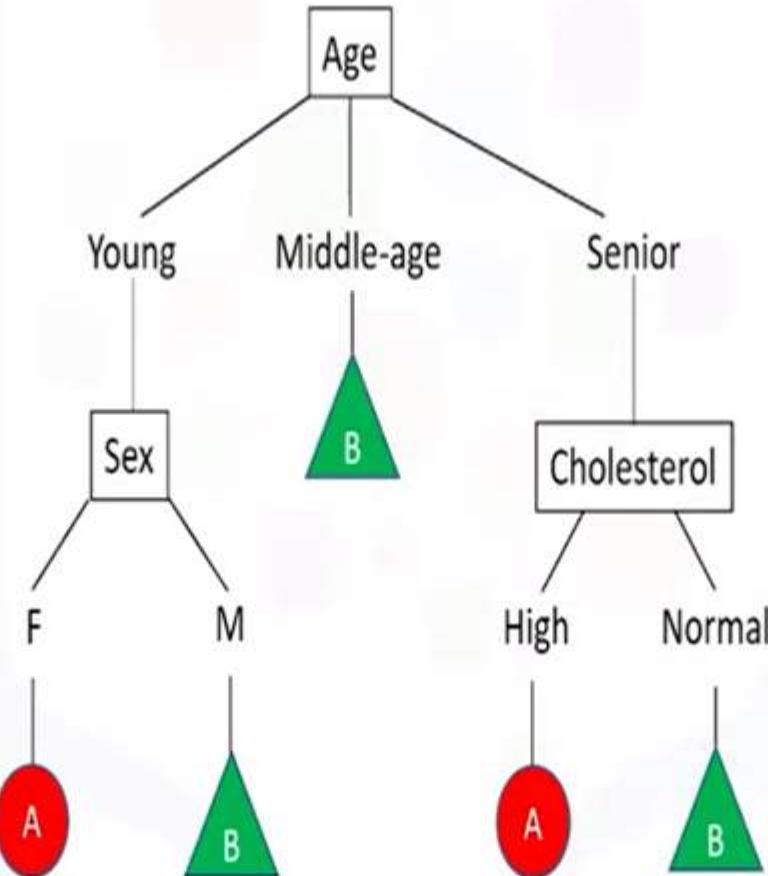
The basic intuition behind a decision tree is to map out all possible decision paths in the form of a tree.

# How to build a Decision Tree:



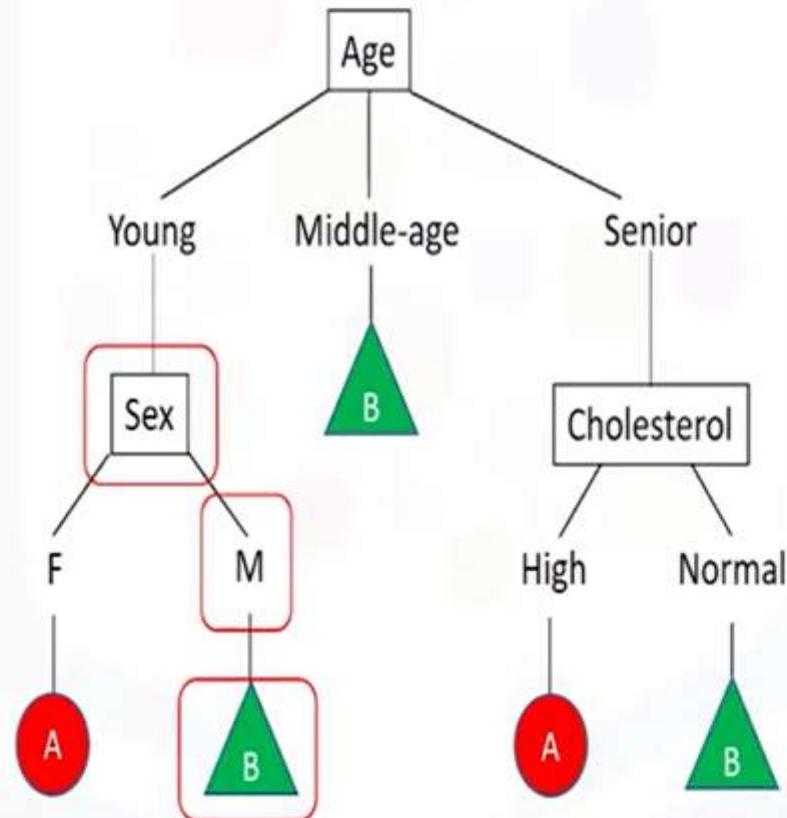
# Building the decision tree with training set:

▲ Drug B  
● Drug A



# Building the decision tree with training set:

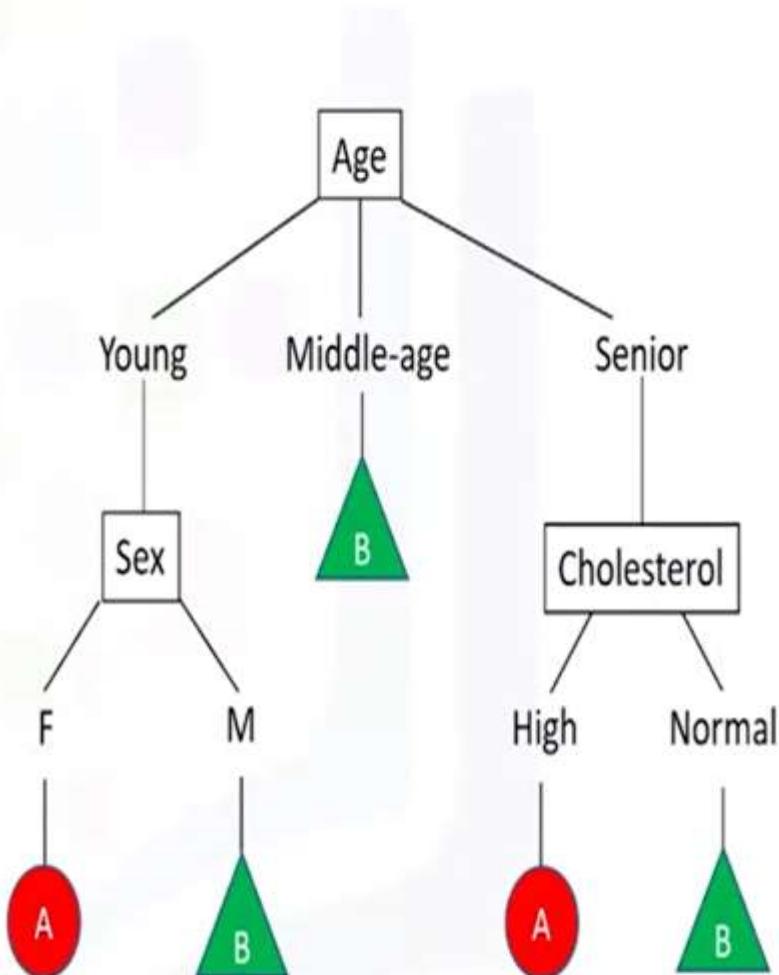
Drug B  
Drug A



- Each **internal node** corresponds to a test
- Each **branch** corresponds to a result of the test
- Each **leaf node** assigns a classification

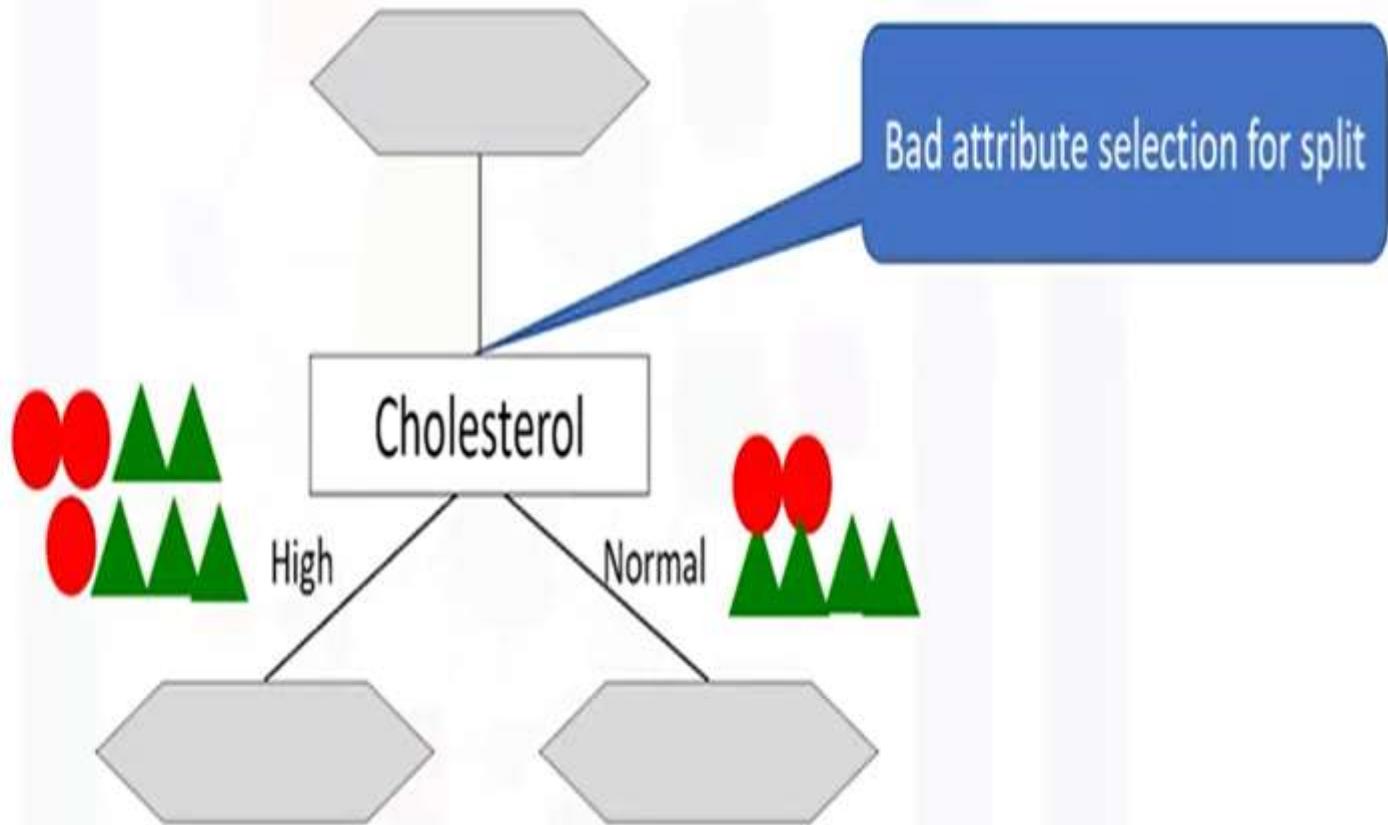
# Decision Tree Learning Algorithm:

1. Choose an attribute from your dataset.
2. Calculate the significance of attribute in splitting of data.
3. Split data based on the value of the best attribute.
4. Go to step 1.



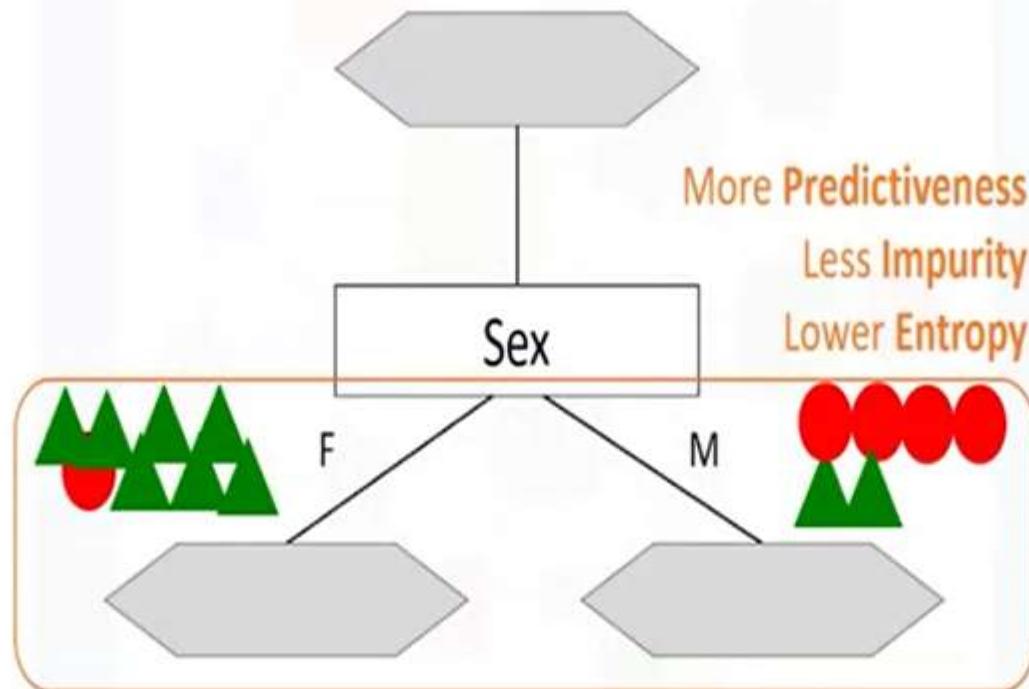
# Which attribute is the best?

▲ Drug B  
● Drug A



# Which attribute is the best?

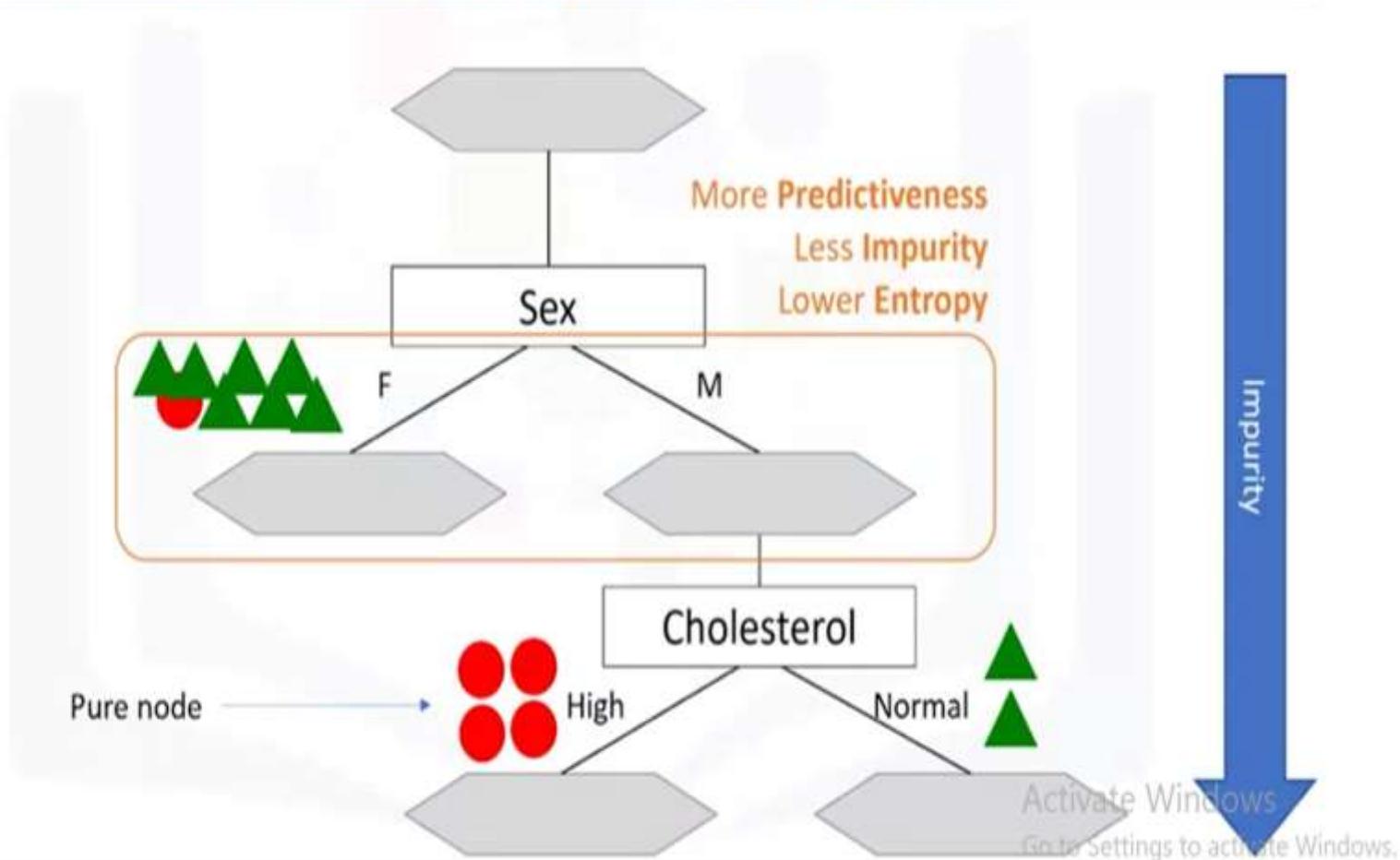
▲ Drug B  
● Drug A



# Which attribute is the best?

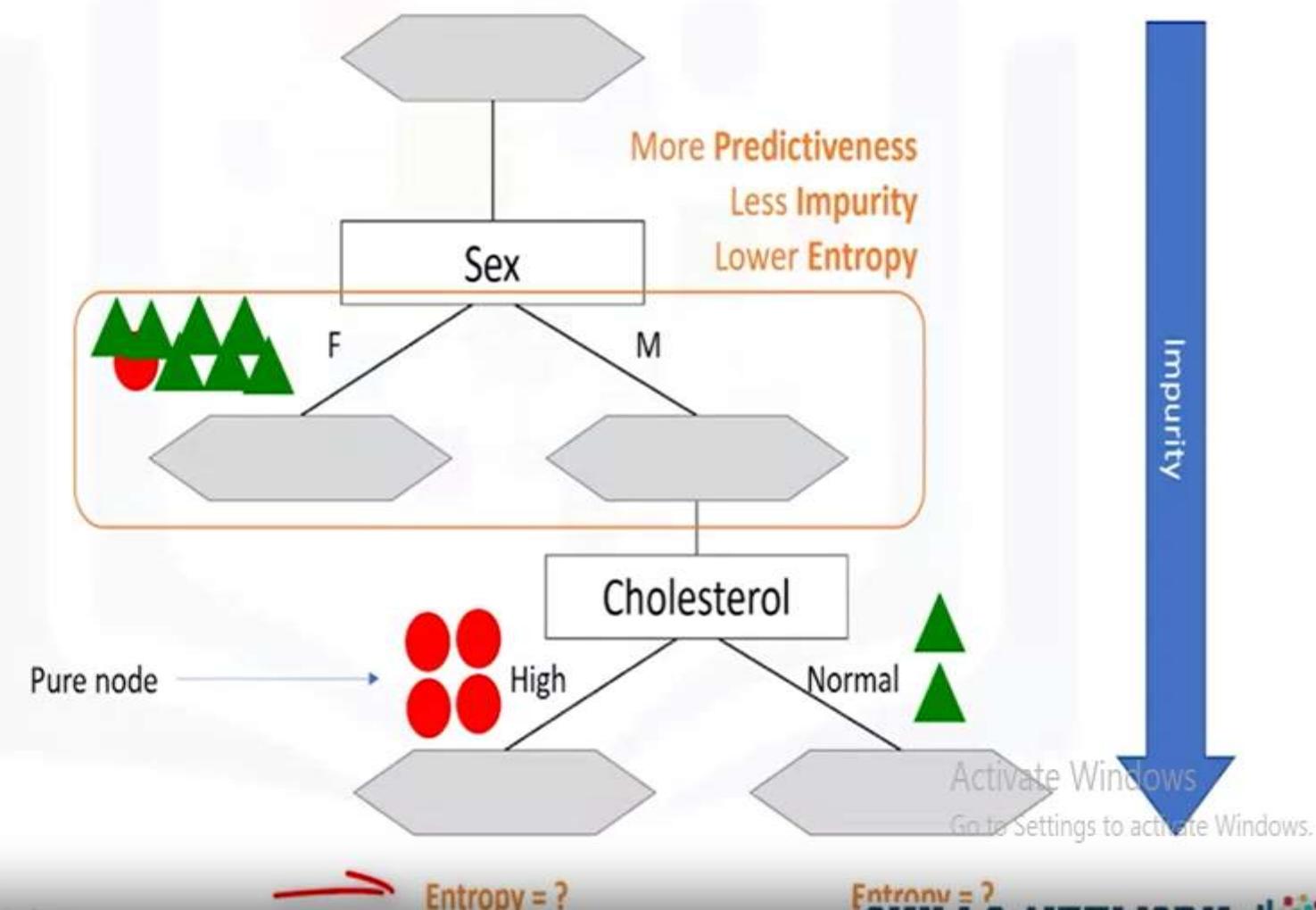
Which attribute is the best ?

▲ Drug B  
● Drug A



# Which attribute is the best?

Drug B  
Drug A

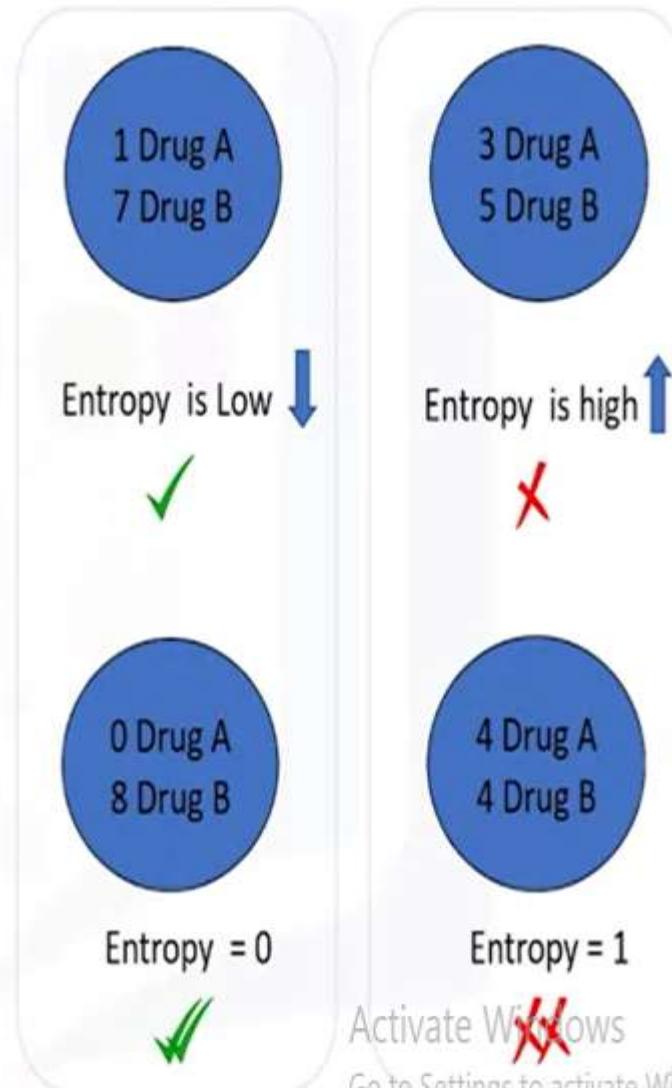


# What is Entropy ?

- Measure of randomness or uncertainty

$$\text{Entropy} = - p(A)\log(p(A)) - p(B)\log(p(B))$$

The lower the Entropy, the less uniform the distribution, the purer the node.



# Which attribute is best one to use?

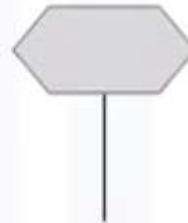
Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	High	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A

S: [9 B, 5 A]

$$E = - p(B) \log(p(B)) - p(A) \log(p(A))$$

$$E = - (9/14) \log(9/14) - (5/14) \log(5/14)$$

$$E = 0.940$$

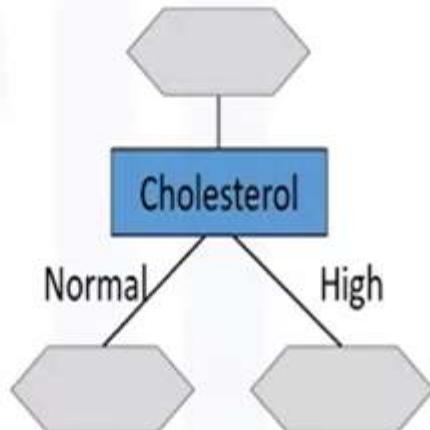


# Is cholesterol the best attribute?

Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	Hiigh	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A

S: [9 B, 5 A]

E = 0.940



S: [6 B, 2 A]

E = 0.811

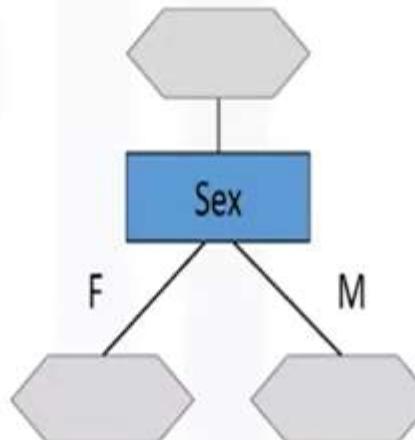
S: [3 B, 3 A]

E = 1.00

# What about ‘Sex’ attribute?

Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	High	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A

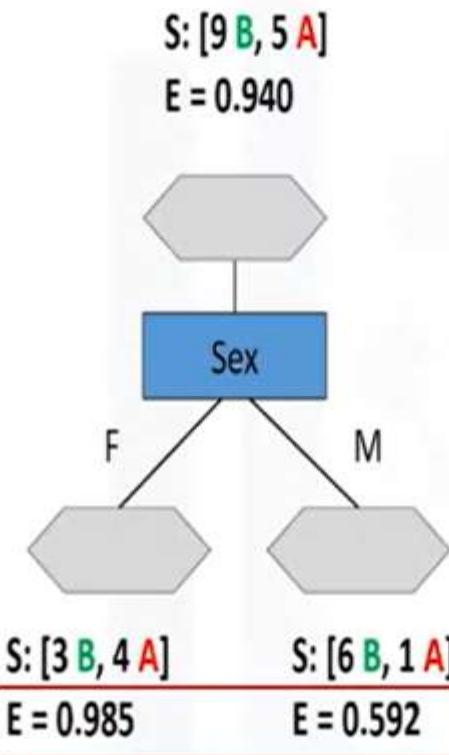
S: [9 B, 5 A]  
E = 0.940



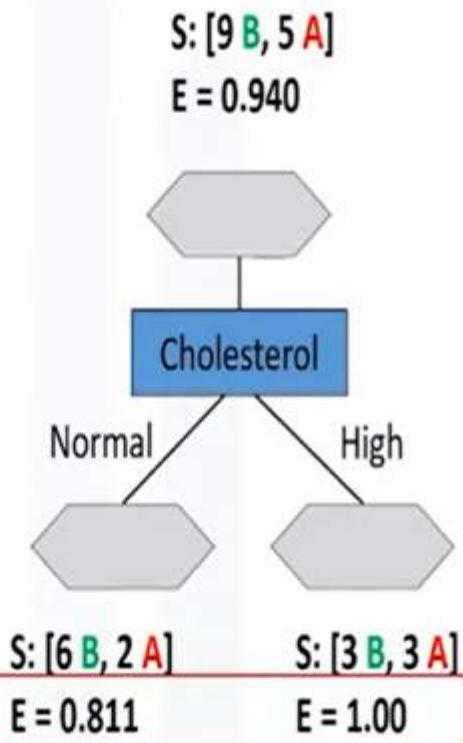
S: [3 B, 4 A]  
E = 0.985

S: [6 B, 1 A]  
E = 0.592

# Which attribute is the best?



Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	High	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A



?

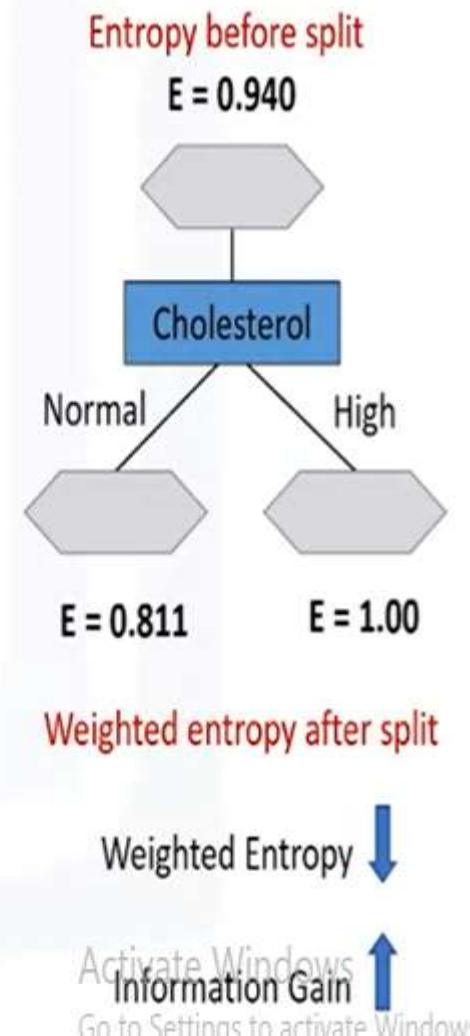
The tree with the higher Information Gain after splitting.

Activate Windows

# What is Information Gain?

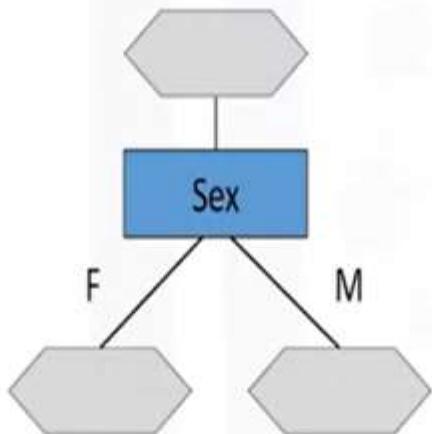
Information gain is the information that can increase the level of certainty after splitting.

$$\text{Information Gain} = (\text{Entropy before split}) - (\text{weighted entropy after split})$$



# Which attribute is the best?

S: [9 B, 5 A]  
E = 0.940



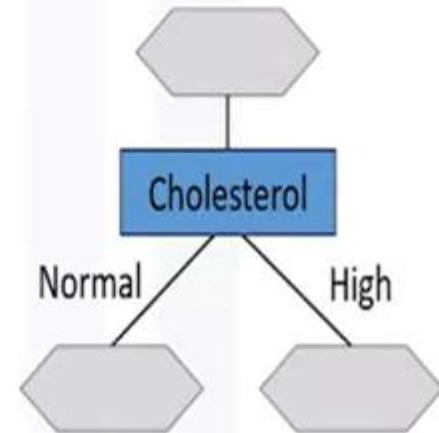
S: [3 B, 4 A]  
E = 0.985

S: [6 B, 1 A]  
E = 0.592

Gain (s, Sex)  
 $= 0.940 - [(7/14)0.985 + (7/14)0.592]$   
 $= 0.151$

Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	High	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A

S: [9 B, 5 A]  
E = 0.940



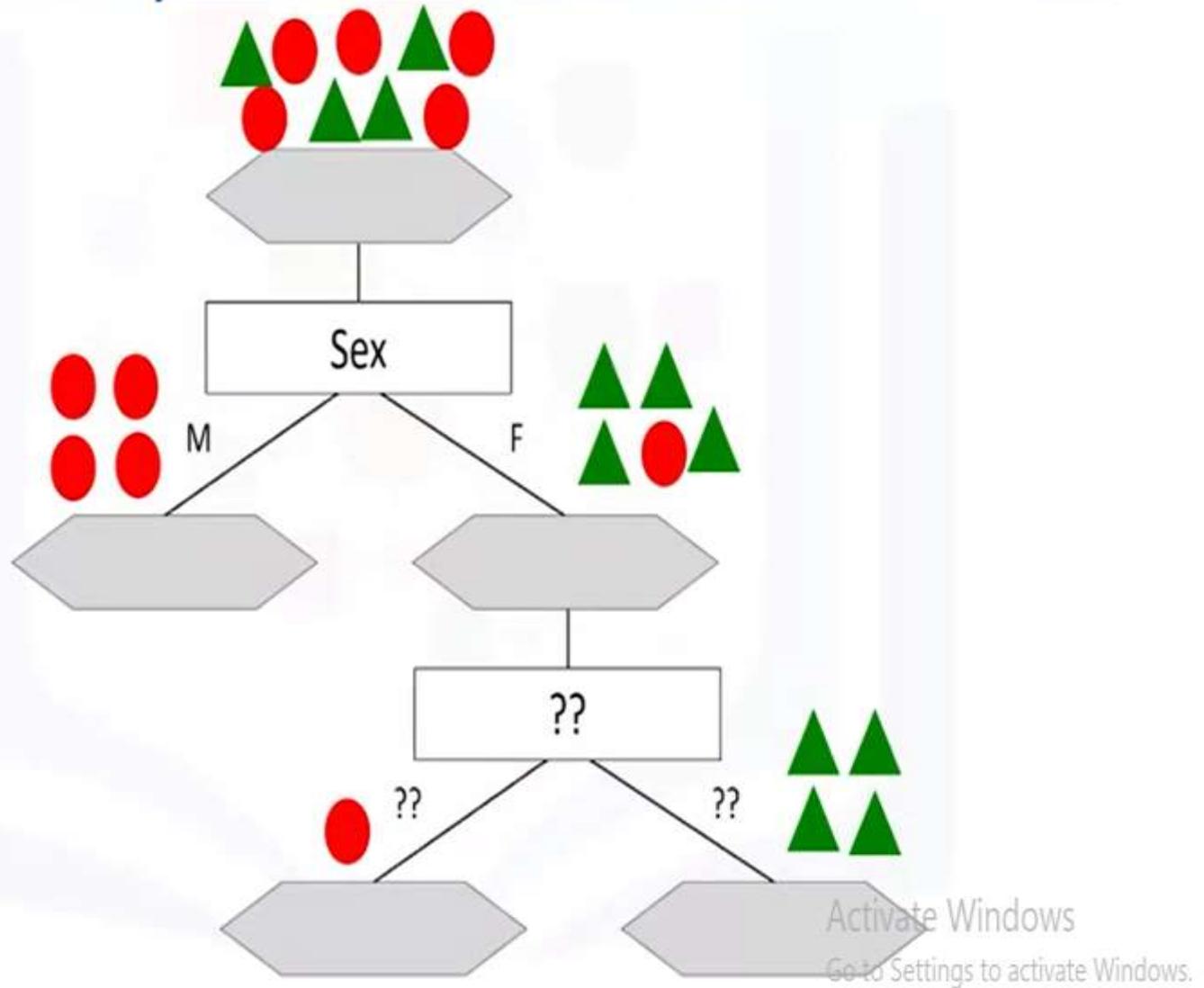
S: [6 B, 2 A]  
E = 0.811

S: [3 B, 3 A]  
E = 1.00

?

Gain (s, Cholesterol)  
 $= 0.940 - [(8/14)0.811 + (6/14)1.0]$   
 $= 0.048$

# Correct way to build a decision tree:



# Support Vector Machine(SVM):

## Classification with SVM

ID	Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class
1000025	5	1	1	1	2	1	3	1	1	benign
1002945	5	4	4	5	7	10	3	2	1	benign
1015425	3	1	1	1	2	2	3	1	1	malignant
1016277	6	8	8	1	3	4	3	7	1	benign
1017023	4	1	1	3	2	1	3	1	1	benign
1017122	8	10	10	8	7	10		7	1	malignant
1018099	1	1	1	1	2	10	3	1	1	benign
1018561	2	1	2	H	2	1	3	1	1	benign
1033078	2	1	1	1	2	1	1	1	5	benign
1033078	4	2	1	1	2	1	2	1	1	benign

ID	Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class
1000015	6	1	1	1	7	1	3	1	1	Benign

Modeling

Prediction

Accuracy = 89% Windows

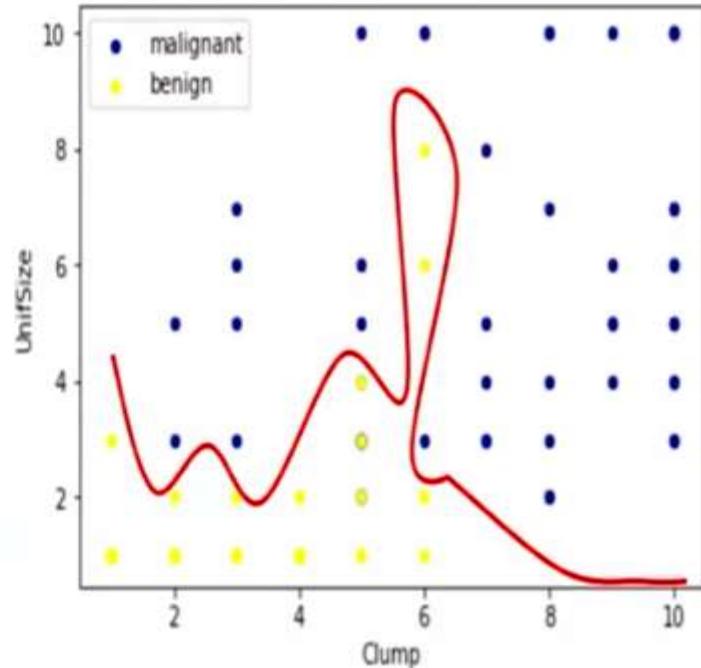


# What is SVM?

SVM is a supervised algorithm that classifies cases by finding a separator.

1. Mapping data to a **high-dimensional** feature space
2. Finding a **separator**

Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class
5	1	1	1	2	1	3	1	1	benign
5	4	4	5	7	10	3	2	1	benign
3	1	1	1	2	2	3	1	1	malignant
6	8	8	1	3	4	3	7	1	benign
4	1	1	3	2	1	3	1	1	benign
8	10	10	8	7	10		7	1	malignant
1	1	1	1	2	10	3	1	1	benign
2	1	2	H	2	1	3	1	1	benign
2	1	1	1	2	1	1	1	5	benign
4	2	1	1	2	1	2	1	1	benign

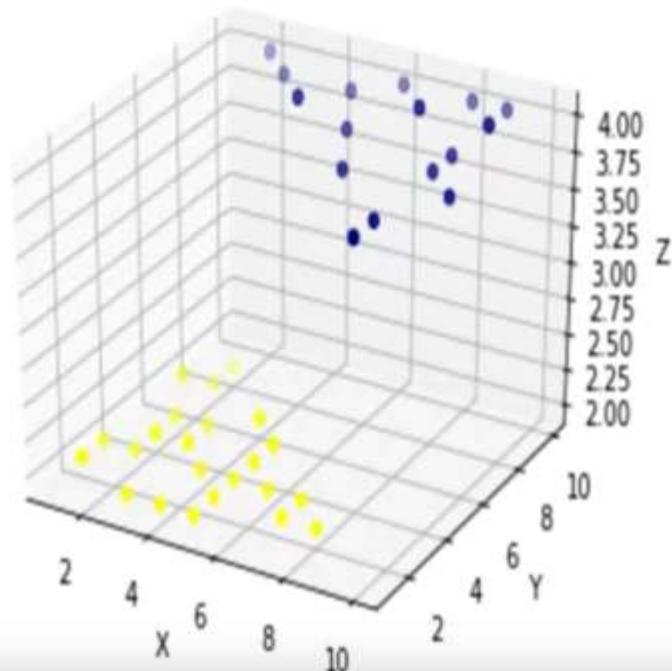


# What is SVM?

SVM is a supervised algorithm that classifies cases by finding a separator.

1. Mapping data to a **high-dimensional** feature space
2. Finding a **separator**

Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class
5	1	1	1	2	1	3	1	1	benign
5	4	4	5	7	10	3	2	1	benign
3	1	1	1	2	2	3	1	1	malignant
6	8	8	1	3	4	3	7	1	benign
4	1	1	3	2	1	3	1	1	benign
8	10	10	8	7	10		7	1	malignant
1	1	1	1	2	10	3	1	1	benign
2	1	2	H	2	1	3	1	1	benign
2	1	1	1	2	1	1	1	5	benign
4	2	1	1	2	1	2	1	1	benign

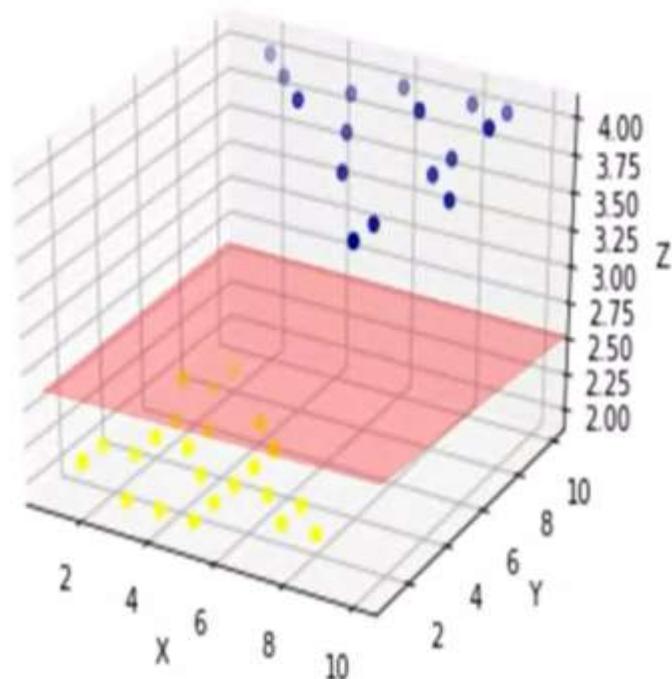


# What is SVM?

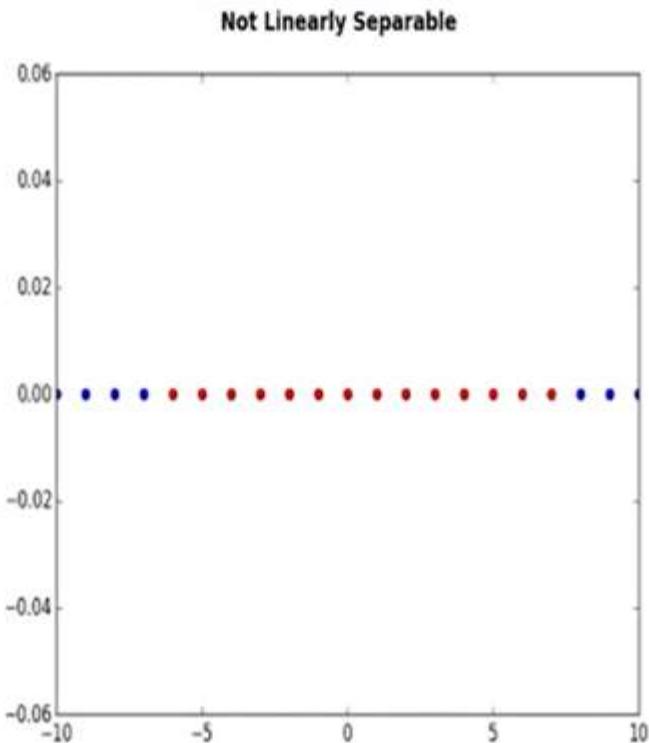
SVM is a supervised algorithm that classifies cases by finding a separator.

1. Mapping data to a **high-dimensional** feature space
2. Finding a **separator**

Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class
5	1	1	1	2	1	3	1	1	benign
5	4	4	5	7	10	3	2	1	benign
3	1	1	1	2	2	3	1	1	malignant
6	8	8	1	3	4	3	7	1	benign
4	1	1	3	2	1	3	1	1	benign
8	10	10	8	7	10		7	1	malignant
1	1	1	1	2	10	3	1	1	benign
2	1	2	H	2	1	3	1	1	benign
2	1	1	1	2	1	1	1	5	benign
4	2	1	1	2	1	2	1	1	benign

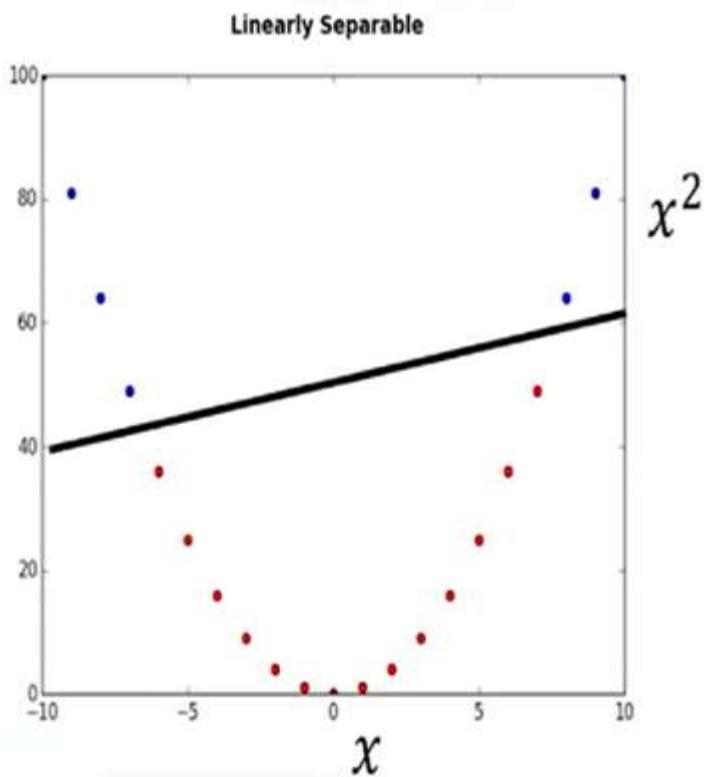
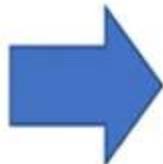


# Data Transformation:



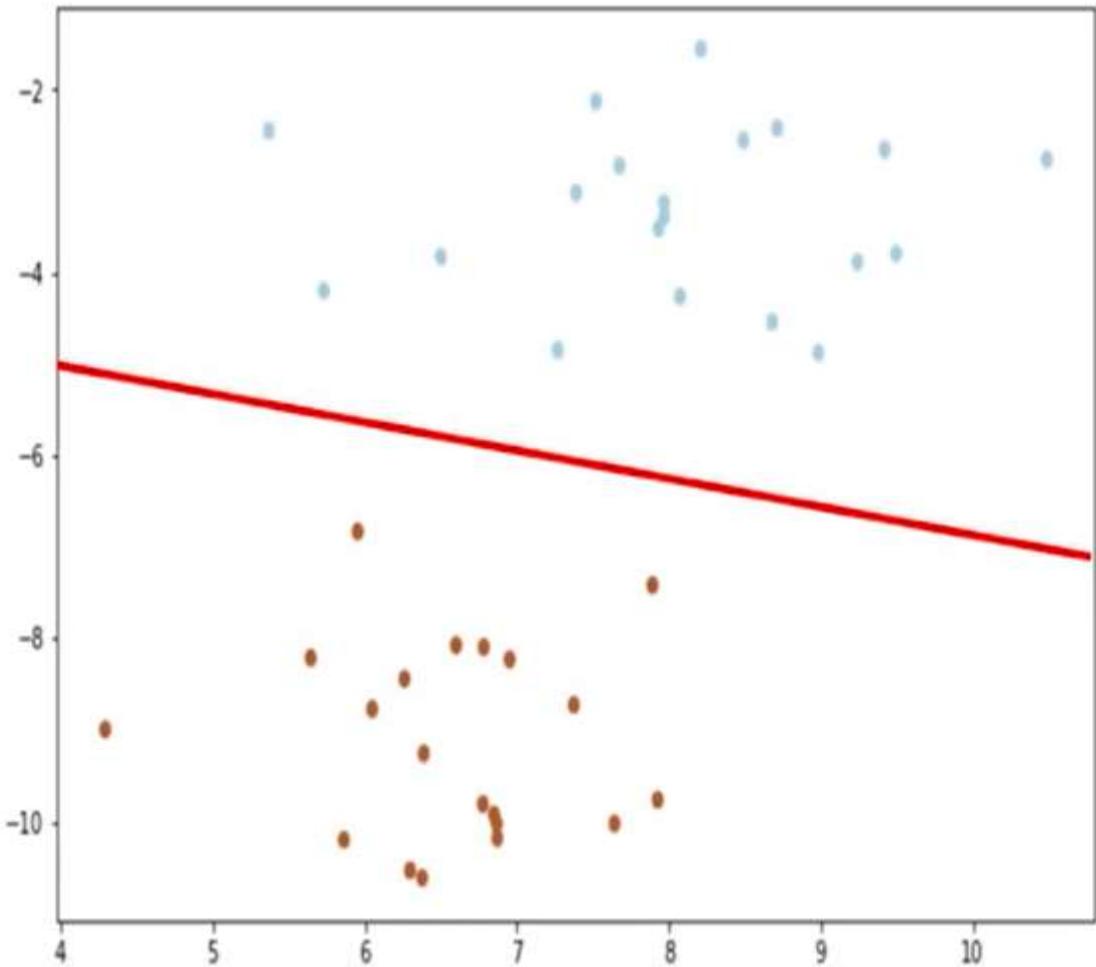
Kernelling:

- Linear
- Polynomial
- RBF
- Sigmoid

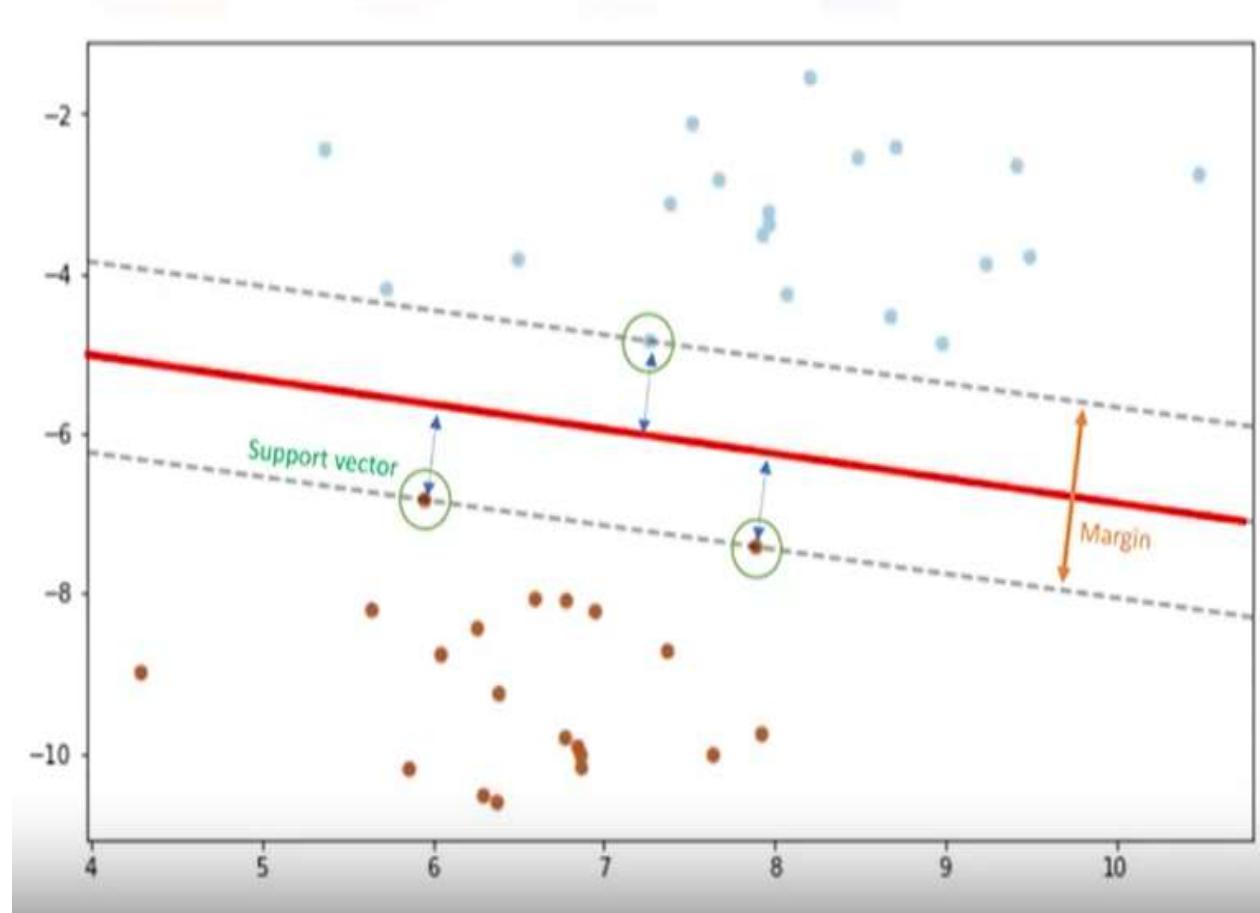


$$\phi(x) = [x, x^2]$$

# Using SVM to find the hyperplane:

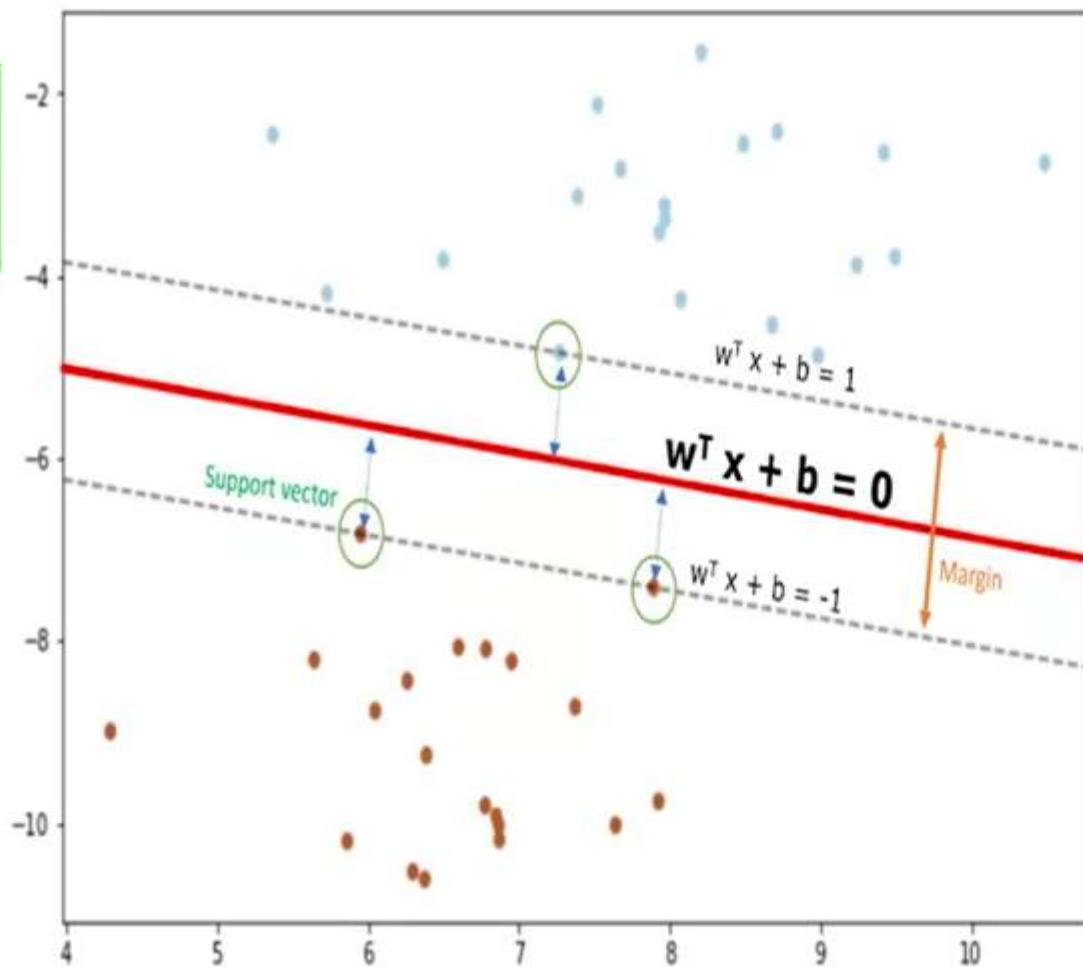


# Using SVM to find the hyperplane:



# Using SVM to find the hyperplane:

Find  $w$  and  $b$  such that  
 $\Phi(w) = \frac{1}{2} w^T w$  is minimized;  
and for all  $\{(x_i, y_i)\}$ :  $y_i(w^T x_i + b) \geq 1$



# Pros and Cons of SVM:

- Advantages:
  - Accurate in high-dimensional spaces
  - Memory efficient
- Disadvantages:
  - Prone to over-fitting
  - No probability estimation
  - Small datasets

# SVM Applications:

- Image recognition
- Text category assignment
- Detecting spam
- Sentiment analysis
- Gene Expression Classification
- Regression, outlier detection and clustering

# Unsupervised Machine Learning:

*“Unsupervised learning is a type of machine learning in which models are trained using unlabeled dataset and are allowed to act on that data without any supervision”*

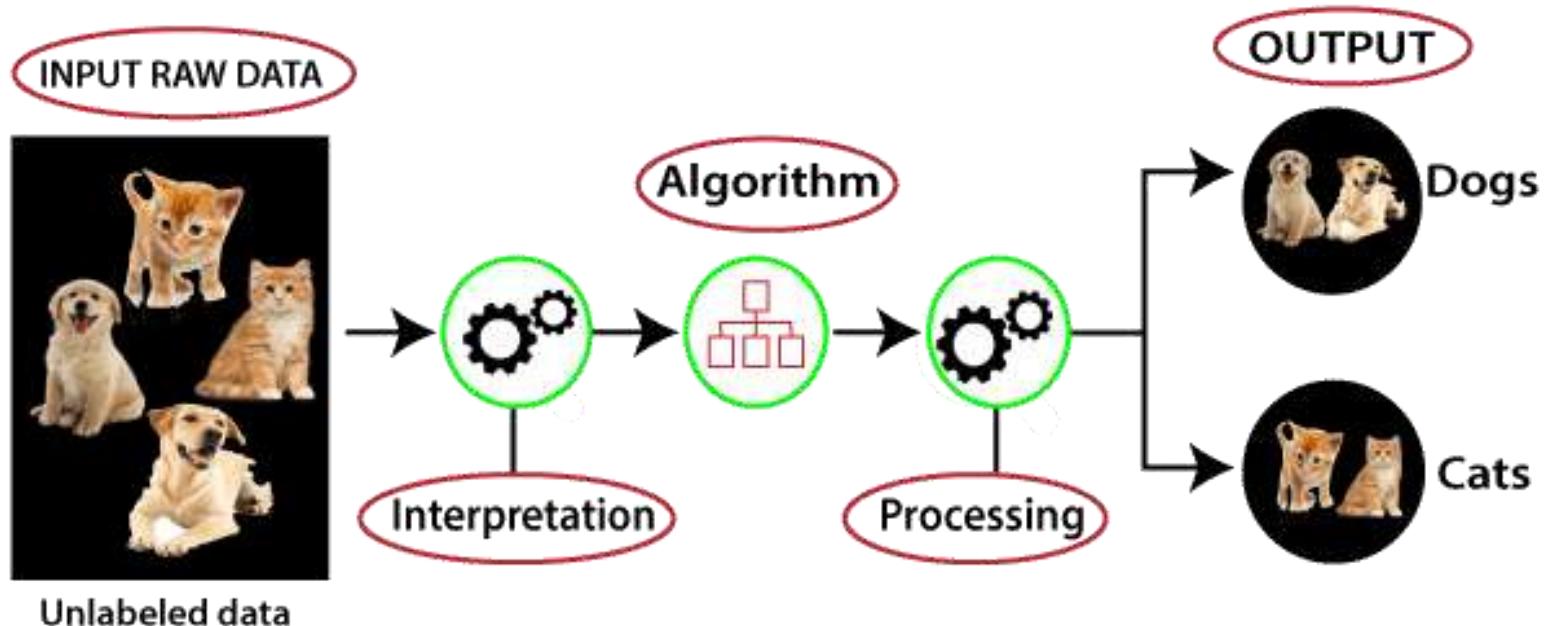
## □ Goal:

- Unsupervised learning is to **find the underlying structure of dataset, group that data according to similarities, and represent that dataset in a compressed format.**

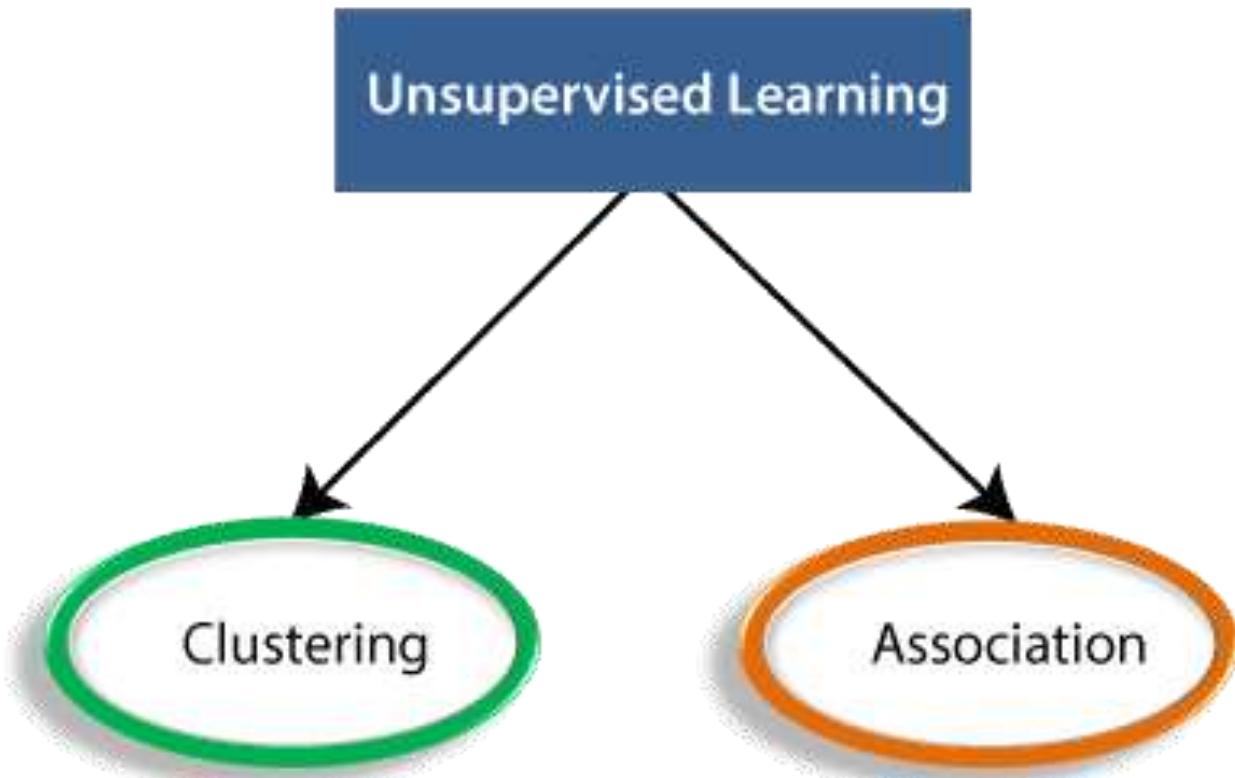
## □ Example: Clustering the dog and cat image dataset into the groups according to similarities between images.



# Working of Unsupervised Learning:



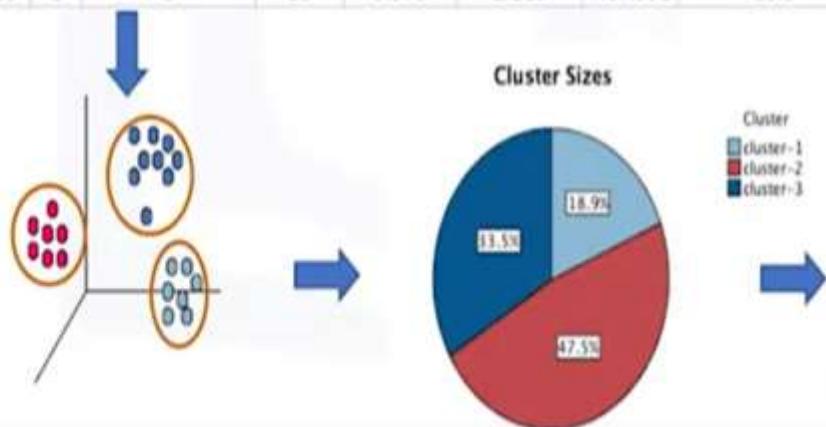
# Types of Unsupervised Learning Algorithm:



# Introduction to Clustering:

## Clustering for segmentation

Customer Id	Age	Edu	Years Employed	Income	Card Debt	Other Debt	Address	DebtIncomeRatio	Defaulted
1	41	2	6	19	0.124	1.073	NBA001	6.3	0
2	47	1	26	100	4.582	8.218	NBA021	12.8	0
3	33	2	10	57	6.111	5.802	NBA013	20.9	1
4	29	2	4	19	0.681	0.516	NBA009	6.3	0
5	47	1	31	253	9.308	8.908	NBA008	7.2	0
6	40	1	23	81	0.998	7.831	NBA016	10.9	1
7	38	2	4	56	0.442	0.454	NBA013	1.6	0
8	42	3	0	64	0.279	3.945	NBA009	6.6	0
9	26	1	5	18	0.575	2.215	NBA006	15.5	1

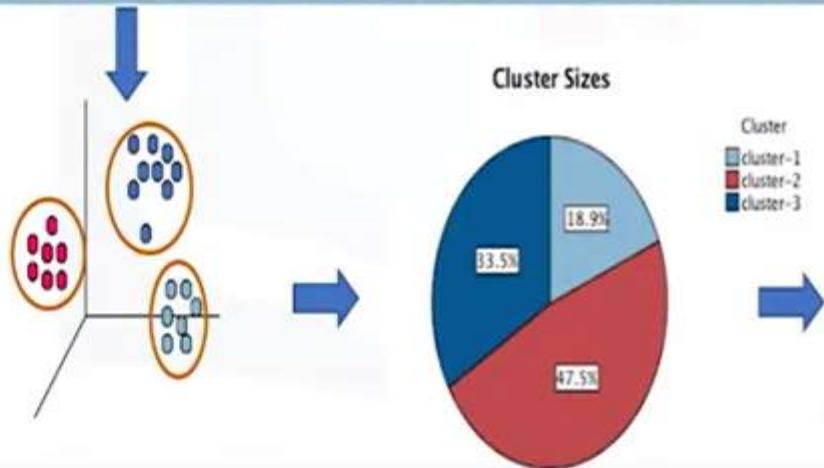


Activate Windows  
Go to Settings to activate Windows.

# Introduction to Clustering:

Customer ID	Age	Edu	Years Employed	Income	Card Debt	Other Debt	Address	DebtIncomeRatio	Defaulted
1	41	2	6	19	0.124	1.073	NBA001	6.3	0
2	47	1	26	100	4.582	8.218	NBA021	12.8	0
3	33	2	10	57	6.111	5.802	NBA013	20.9	1
4	29	2	4	19	0.681	0.516	NBA009	6.3	0
5	47	1	31	253	9.308	8.908	NBA008	7.2	0
6	40	1	23	81	0.998	7.831	NBA016	10.9	1
7	38	2	4	56	0.442	0.454	NBA013	1.6	0
8	42	3	0	64	0.279	3.945	NBA009	6.6	0
9	26	1	5	18	0.575	2.215	NBA006	15.5	1

Customer ID	Segment
1	YOUNG AND LOW INCOME
2	AFFLUENT AND MIDDLE AGED
3	AFFLUENT AND MIDDLE AGED
4	YOUNG AND LOW INCOME
5	AFFLUENT AND MIDDLE AGED
6	AFFLUENT AND MIDDLE AGED
7	YOUNG AND LOW INCOME
8	YOUNG AND LOW INCOME
9	AFFLUENT AND MIDDLE AGED

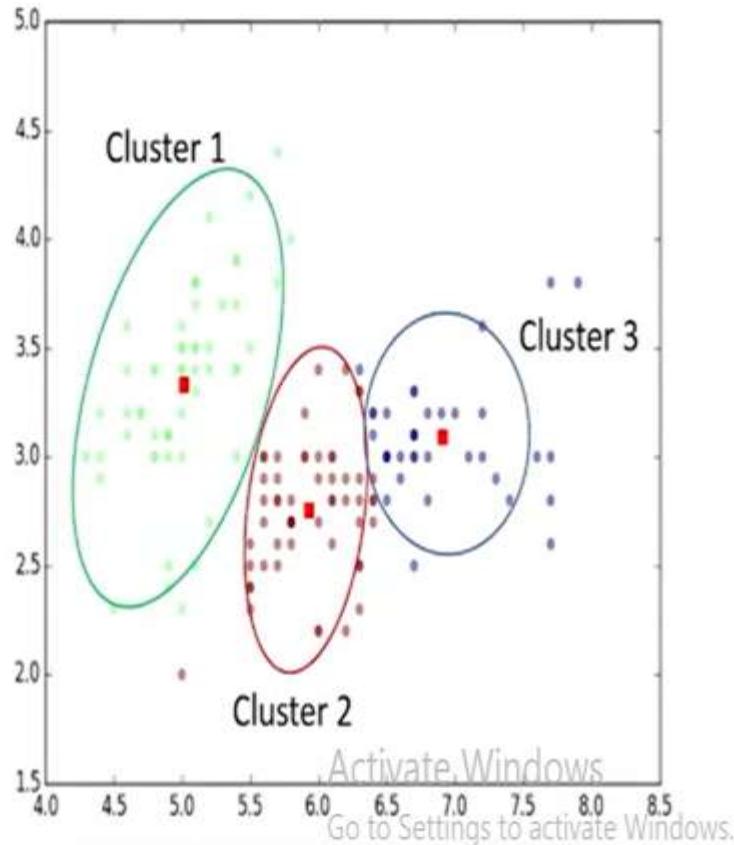


Go to Settings to activate Windows.

# Clustering:

## What is a cluster?

A group of objects that are **similar to other objects** in the cluster, and **dissimilar to data points** in other clusters.



# Clustering Vs. Classification:

Labeled dataset

Customer Id	Age	Edu	Years Employed	Income	Card Debt	Other Debt	Address	DebtIncomeRatio	Defaulted
1	41	2	6	19	0.124	1.073	NBA001	6.3	0
2	47	1	26	100	4.582	8.218	NBA021	12.8	0
3	33	2	10	57	6.111	5.802	NBA013	20.9	1
4	29	2	4	19	0.681	0.516	NBA009	6.3	0
5	47	1	31	253	9.308	8.908	NBA008	7.2	0
6	40	1	23	81	0.998	7.831	NBA016	10.9	1
7	38	2	4	56	0.442	0.454	NBA013	1.6	0
8	42	3	0	64	0.279	3.945	NBA009	6.6	0
9	26	1	5	18	0.575	2.215	NBA006	15.5	1
10	47	3	23	115	0.653	3.947	NBA011	4	0

Modeling

Decision Tree

Prediction

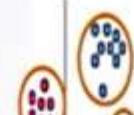


Unlabeled dataset

Customer Id	Age	Edu	Years Employed	Income	Card Debt	Other Debt	Address	DebtIncomeRatio	Defaulted
1	41	2	6	19	0.124	1.073	NBA001	6.3	0
2	47	1	26	100	4.582	8.218	NBA021	12.8	0
3	33	2	10	57	6.111	5.802	NBA013	20.9	1
4	29	2	4	19	0.681	0.516	NBA009	6.3	0
5	47	1	31	253	9.308	8.908	NBA008	7.2	0
6	40	1	23	81	0.998	7.831	NBA016	10.9	1
7	38	2	4	56	0.442	0.454	NBA013	1.6	0
8	42	3	0	64	0.279	3.945	NBA009	6.6	0
9	26	1	5	18	0.575	2.215	NBA006	15.5	1

Modeling

Segmentation



K-Means

Activate Windows

Go to Settings to activate Windows

# Clustering Applications:

- RETAIL/MARKETING:
  - Identifying buying patterns of customers
  - Recommending new books or movies to new customers
- BANKING:
  - Fraud detection in credit card use
  - Identifying clusters of customers (e.g., loyal)
- INSURANCE:
  - Fraud detection in claims analysis
  - Insurance risk of customers

# Clustering Applications:

- PUBLICATION:

- Auto-categorizing news based on their content
- Recommending similar news articles

- MEDICINE:

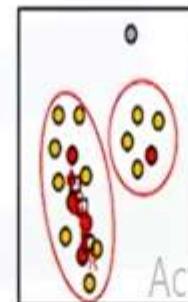
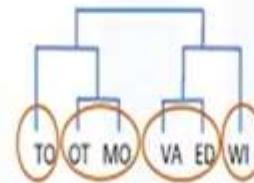
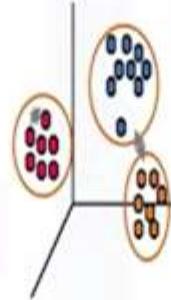
- Characterizing patient behavior

- BIOLOGY:

- Clustering genetic markers to identify family ties

# Clustering Algorithms?

- Partitioned-based Clustering
  - Relatively efficient
  - E.g. k-Means, k-Median, Fuzzy c-Means
- Hierarchical Clustering
  - Produces trees of clusters
  - E.g. Agglomerative, Divisive
- Density-based Clustering
  - Produces arbitrary shaped clusters
  - E.g. DBSCAN

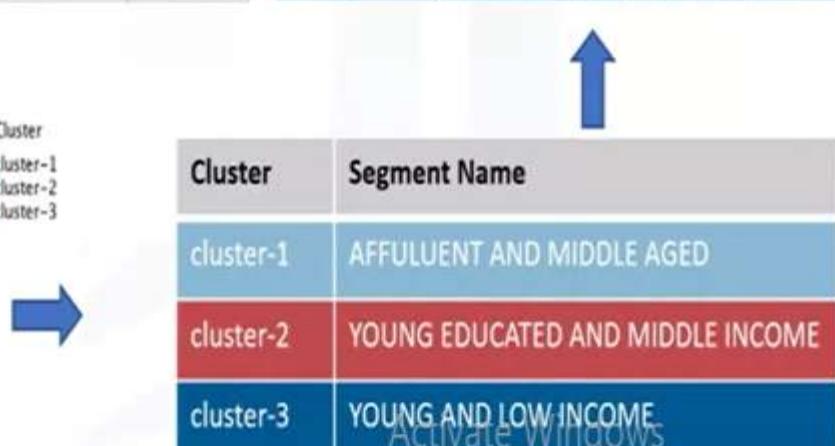
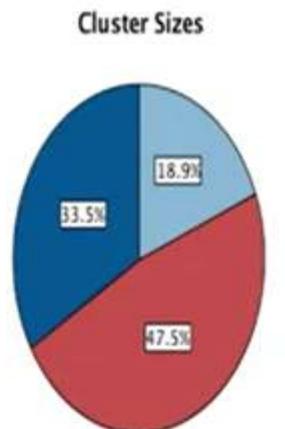
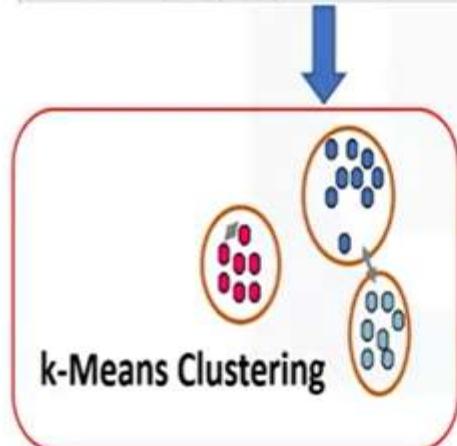


Activate Windows

# Introduction to k-Means Clustering:

Customer ID	Age	Edu	Years Employed	Income	Card Debt	Other Debt	Address	DebtIncomeRatio	Defaulted
1	41	2	6	19	0.124	1.073	NBA001	6.3	0
2	47	1	26	100	4.582	8.218	NBA021	12.8	0
3	33	2	10	57	6.111	5.802	NBA013	20.9	1
4	29	2	4	19	0.681	0.516	NBA009	6.3	0
5	47	1	31	253	9.308	8.908	NBA008	7.2	0
6	40	1	23	81	0.998	7.831	NBA016	10.9	1
7	38	2	4	56	0.442	0.454	NBA013	1.6	0
8	42	3	0	64	0.279	3.945	NBA009	6.6	0
9	26	1	5	18	0.575	2.215	NBA006	15.5	1

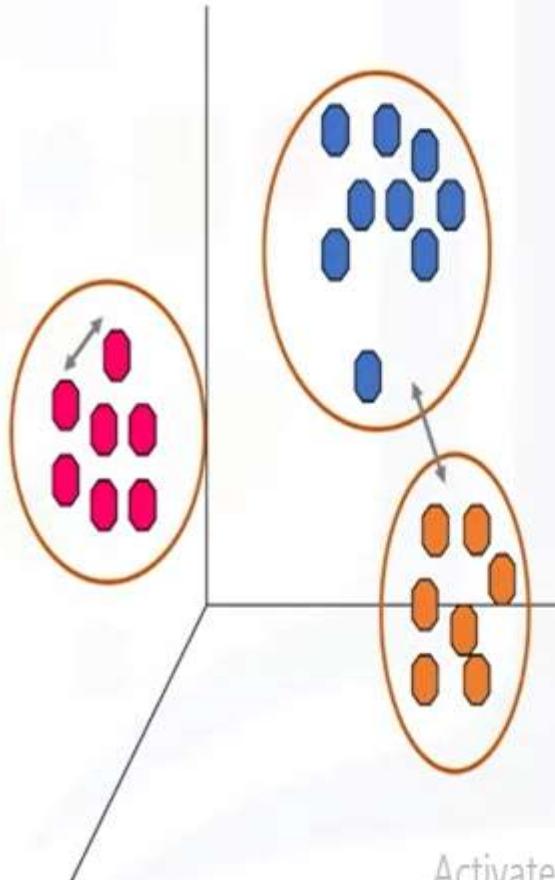
Customer ID	Segment
1	YOUNG AND LOW INCOME
2	AFFLUENT AND MIDDLE AGED
3	AFFLUENT AND MIDDLE AGED
4	YOUNG AND LOW INCOME
5	AFFLUENT AND MIDDLE AGED
6	AFFLUENT AND MIDDLE AGED
7	YOUNG AND LOW INCOME
8	YOUNG AND LOW INCOME
9	AFFLUENT AND MIDDLE AGED



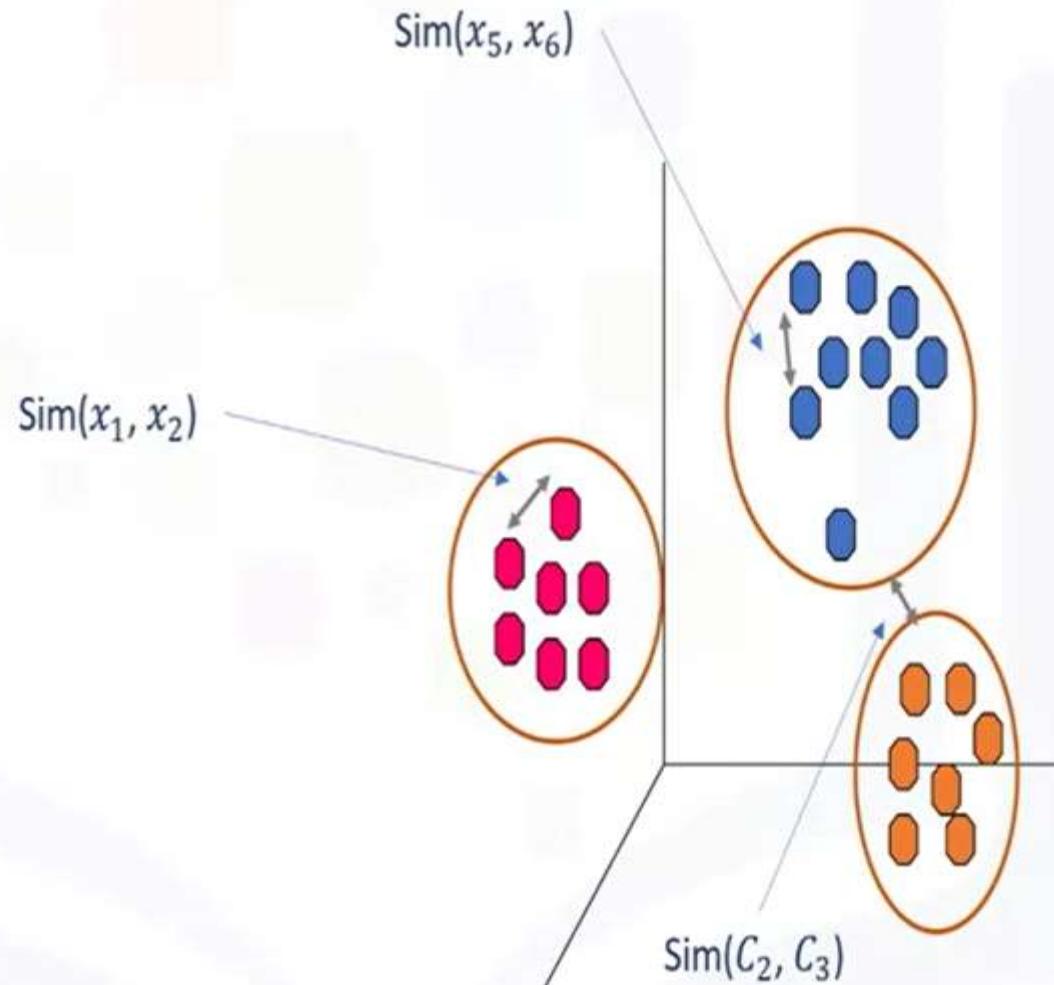
Go to Settings to activate Windows.

# K-Means Algorithm:

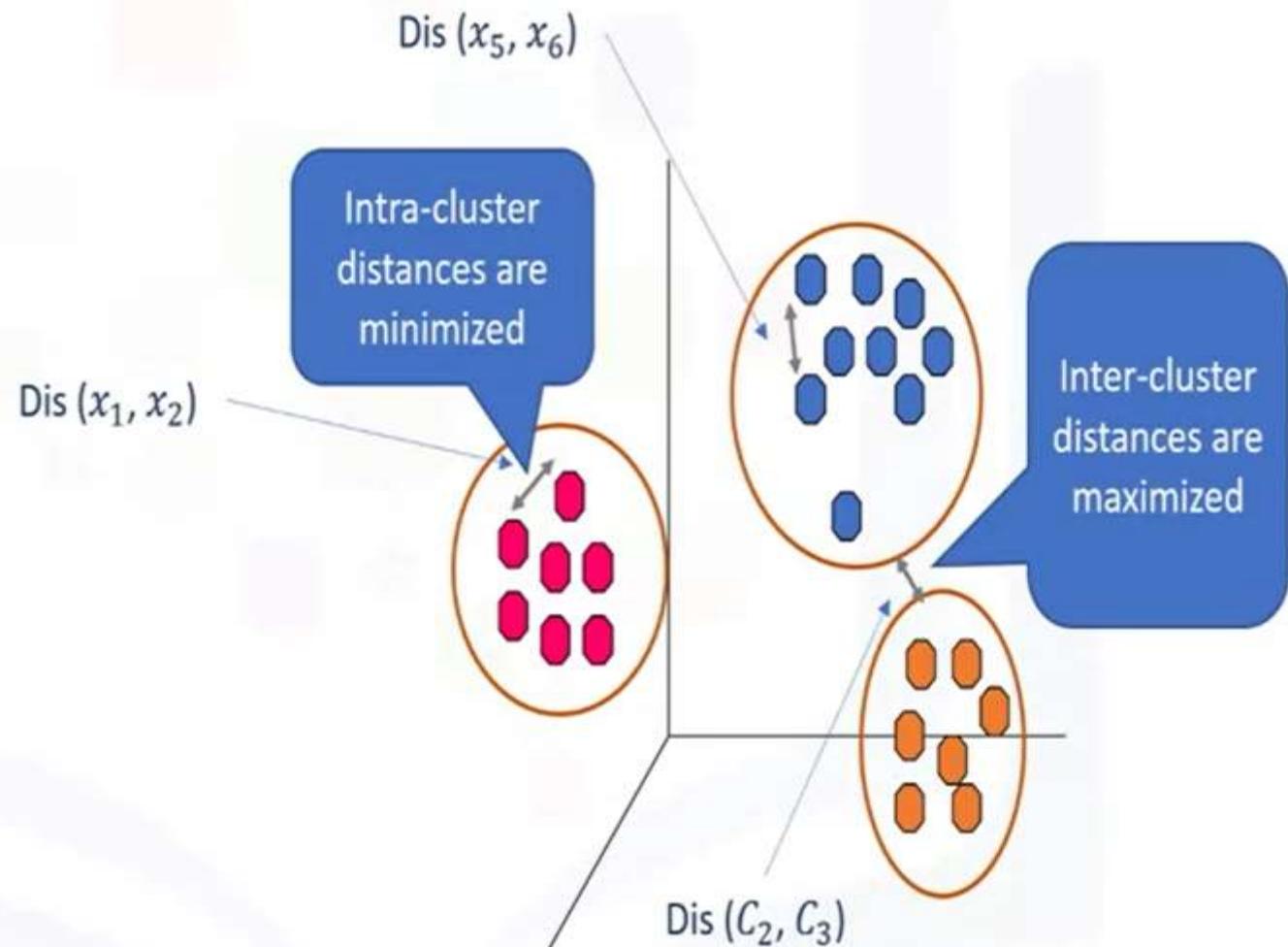
- Partitioning Clustering
- K-means divides the data into non-overlapping subsets (clusters) without any cluster-internal structure
- Examples within a cluster are very similar
- Examples across different clusters are very different



# Determine the similarity or dissimilarity:



# Determine the similarity or dissimilarity:



# 1-dimensional Similarity/distance:



Customer 1

Age

54

Customer 2

Age

50

$$\text{Dis}(x_1, x_2) = \sqrt{\sum_{i=0}^n (x_{1i} - x_{2i})^2}$$

$$\text{Dis}(x_1, x_2) = \sqrt{(54 - 50)^2} = 4$$

## 2-dimensional Similarity/distance:

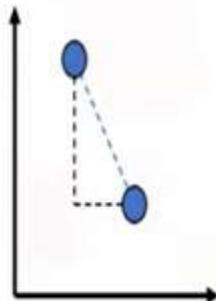


Customer 1

Age	Income
54	190

Customer 2

Age	Income
50	200



$$\text{Dis}(x_1, x_2) = \sqrt{\sum_{i=0}^n (x_{1i} - x_{2i})^2}$$

$$= \sqrt{(54 - 50)^2 + (190 - 200)^2} = 10.77$$

Activate Windows

# Multi-dimensional Similarity/distance:



Customer 1

Age	Income	education
54	190	3



Customer 2

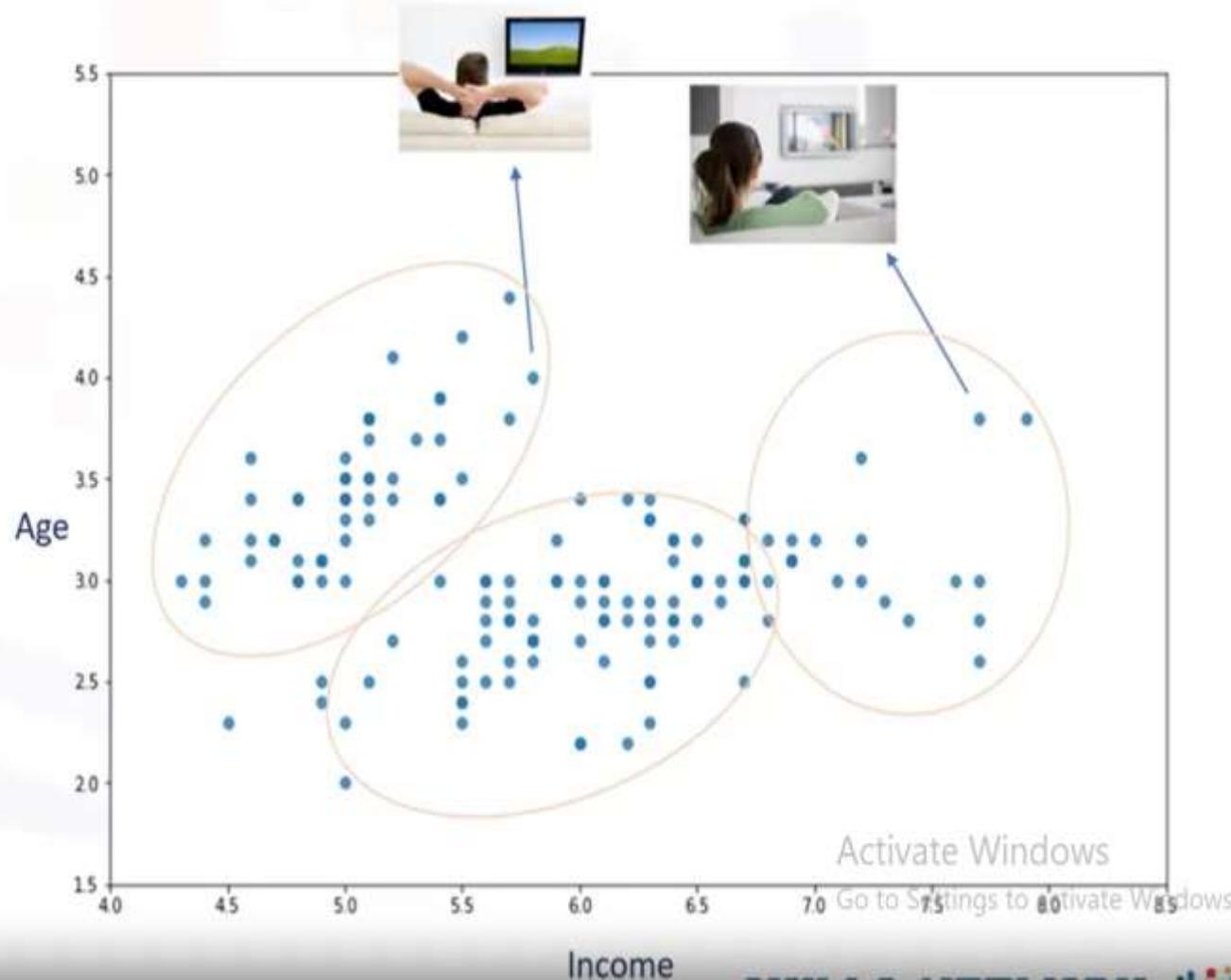
Age	Income	education
50	200	8

$$\text{Dis}(x_1, x_2) = \sqrt{\sum_{i=0}^n (x_{1i} - x_{2i})^2}$$

$$= \sqrt{(54 - 50)^2 + (190 - 200)^2 + (3 - 8)^2} = 11.87$$

# How does k-Means clustering work:

Customer ID	Age	Income
1	3	4
2	2	6
3	3.5	2
...	...	..



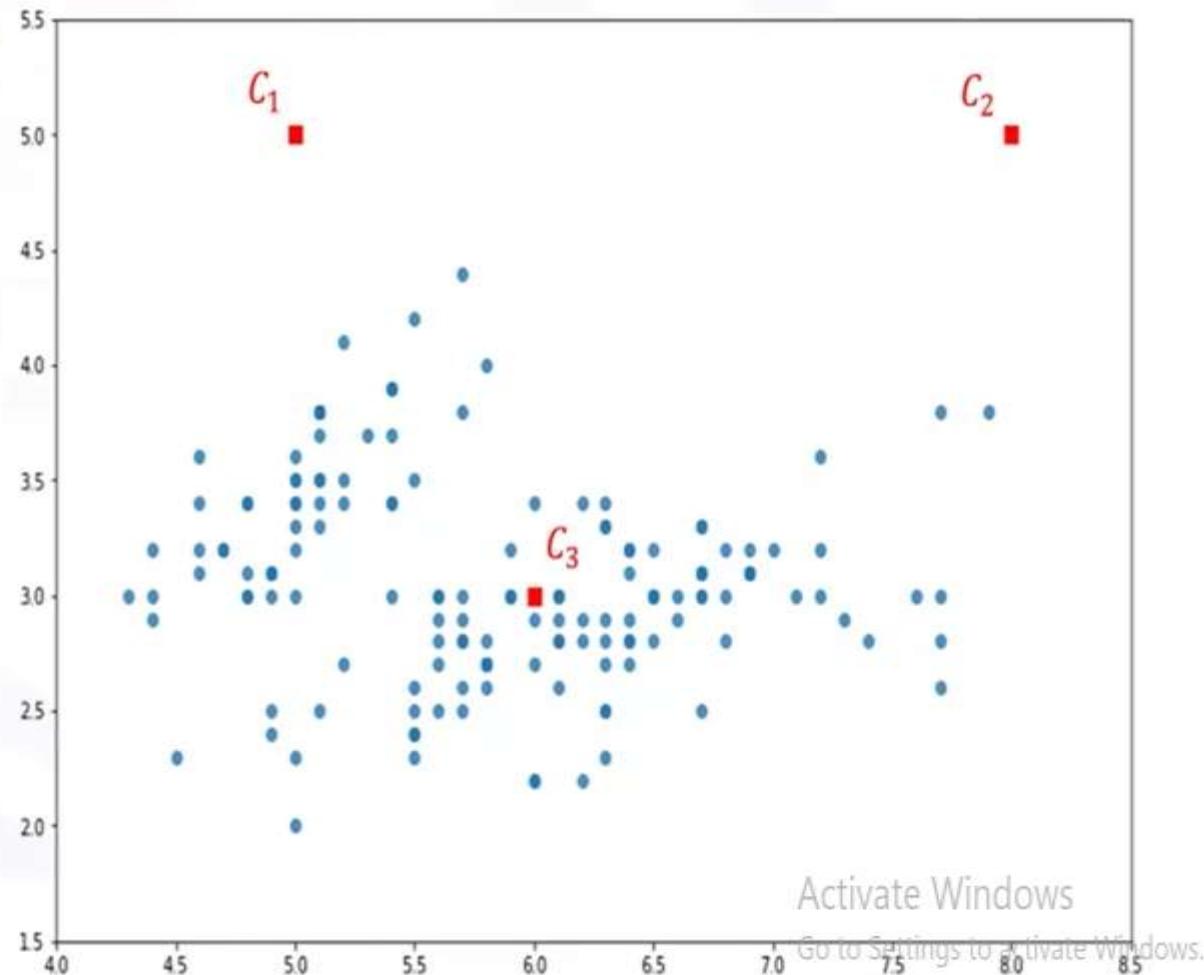
# k-Means clustering – initialize k:

1) Initialize k=3  
centroids randomly

$$C_1 = [8., 5.]$$

$$C_2 = [5., 5.]$$

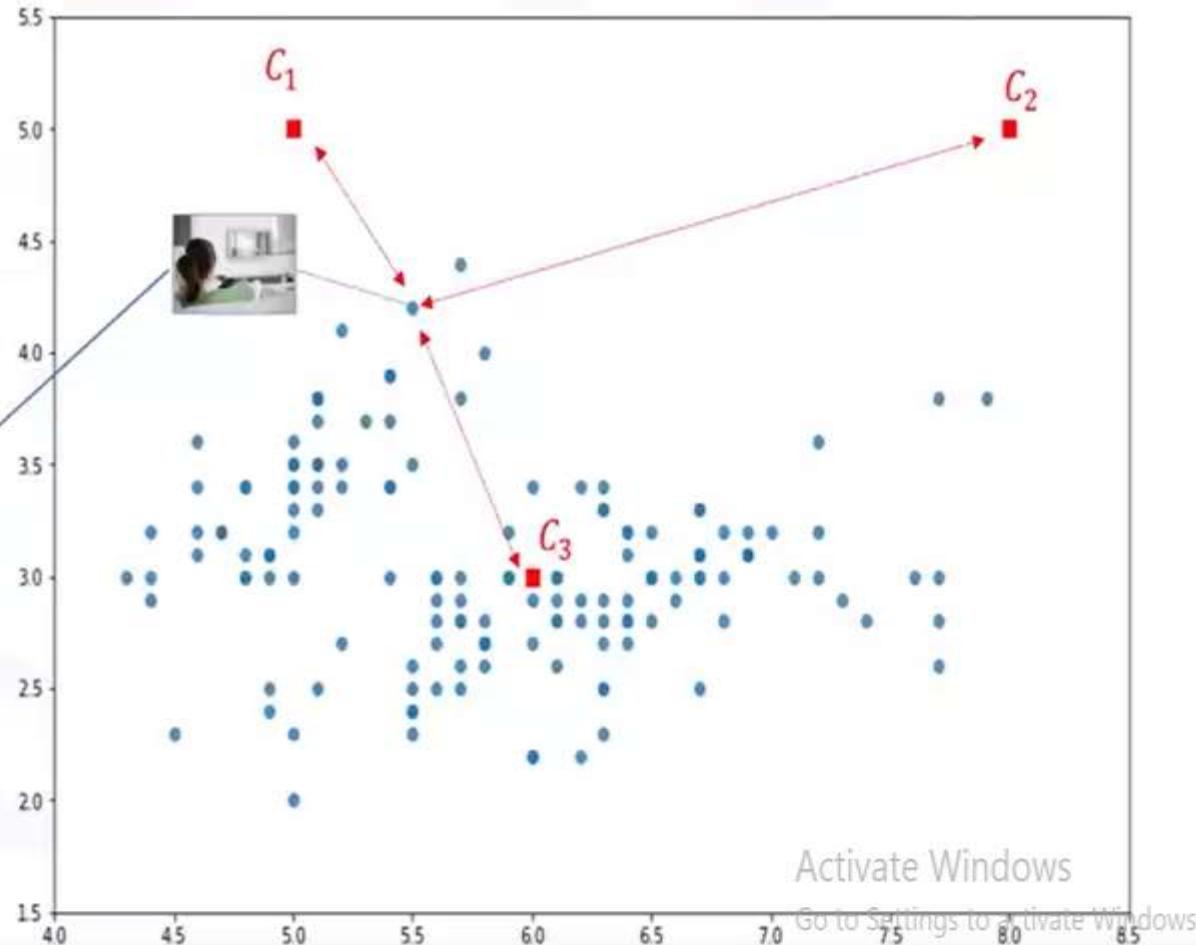
$$C_3 = [6., 3.]$$



# k-Means clustering – calculate the distance:

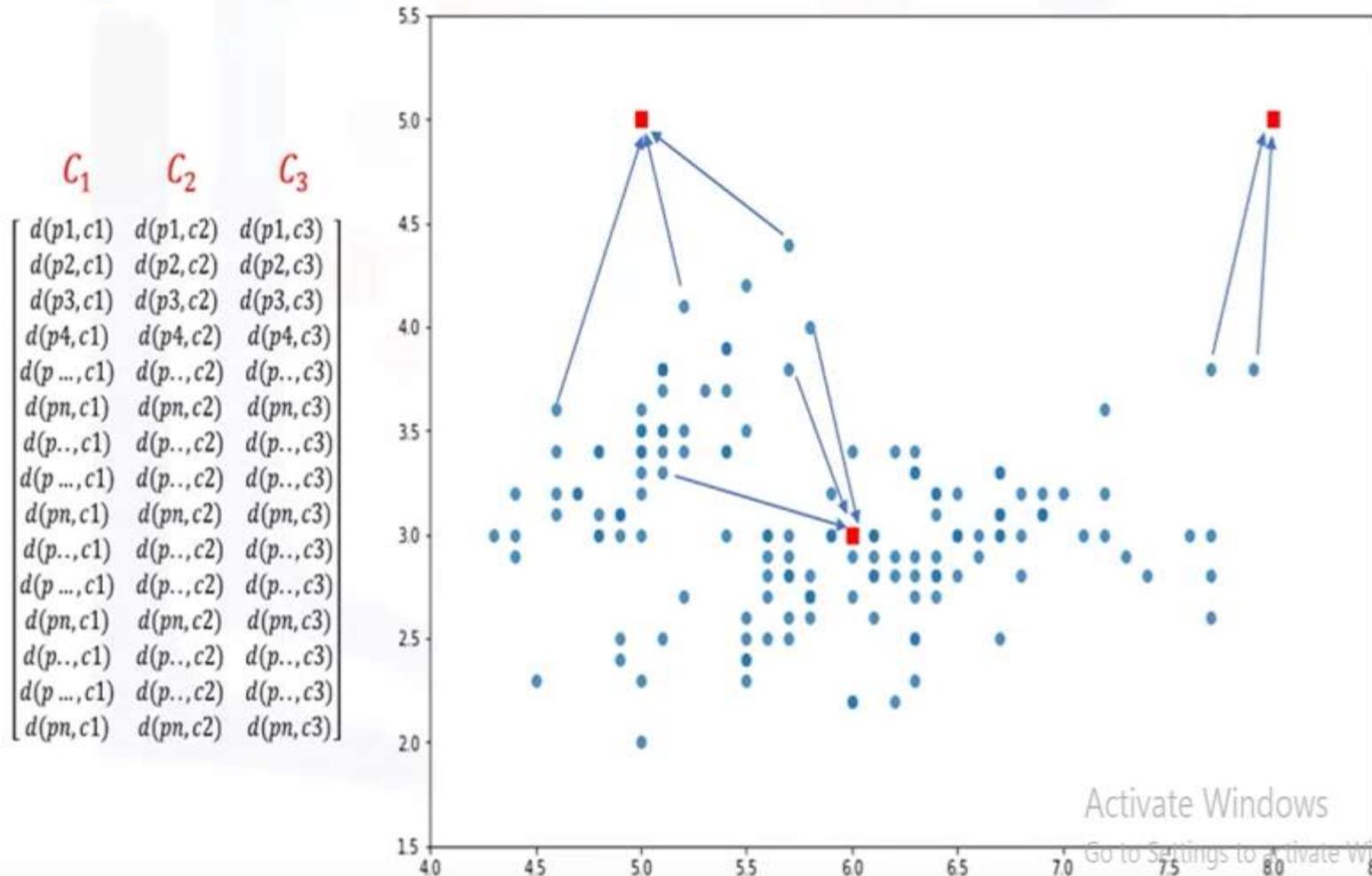
## 2) Distance calculation

$C_1$	$C_2$	$C_3$
$d(p1, c1)$	$d(p1, c2)$	$d(p1, c3)$
$d(p2, c1)$	$d(p2, c2)$	$d(p2, c3)$
$d(p3, c1)$	$d(p3, c2)$	$d(p3, c3)$
$d(p4, c1)$	$d(p4, c2)$	$d(p4, c3)$
$d(p\ldots, c1)$	$d(p\ldots, c2)$	$d(p\ldots, c3)$
$d(pn, c1)$	$d(pn, c2)$	$d(pn, c3)$
$d(p\ldots, c1)$	$d(p\ldots, c2)$	$d(p\ldots, c3)$
$d(p\ldots, c1)$	$d(p\ldots, c2)$	$d(p\ldots, c3)$
$d(pn, c1)$	$d(pn, c2)$	$d(pn, c3)$
$d(p\ldots, c1)$	$d(p\ldots, c2)$	$d(p\ldots, c3)$
$d(p\ldots, c1)$	$d(p\ldots, c2)$	$d(p\ldots, c3)$
$d(pn, c1)$	$d(pn, c2)$	$d(pn, c3)$
$d(p\ldots, c1)$	$d(p\ldots, c2)$	$d(p\ldots, c3)$
$d(p\ldots, c1)$	$d(p\ldots, c2)$	$d(p\ldots, c3)$
$d(pn, c1)$	$d(pn, c2)$	$d(pn, c3)$
$d(p\ldots, c1)$	$d(p\ldots, c2)$	$d(p\ldots, c3)$
$d(p\ldots, c1)$	$d(p\ldots, c2)$	$d(p\ldots, c3)$
$d(pn, c1)$	$d(pn, c2)$	$d(pn, c3)$



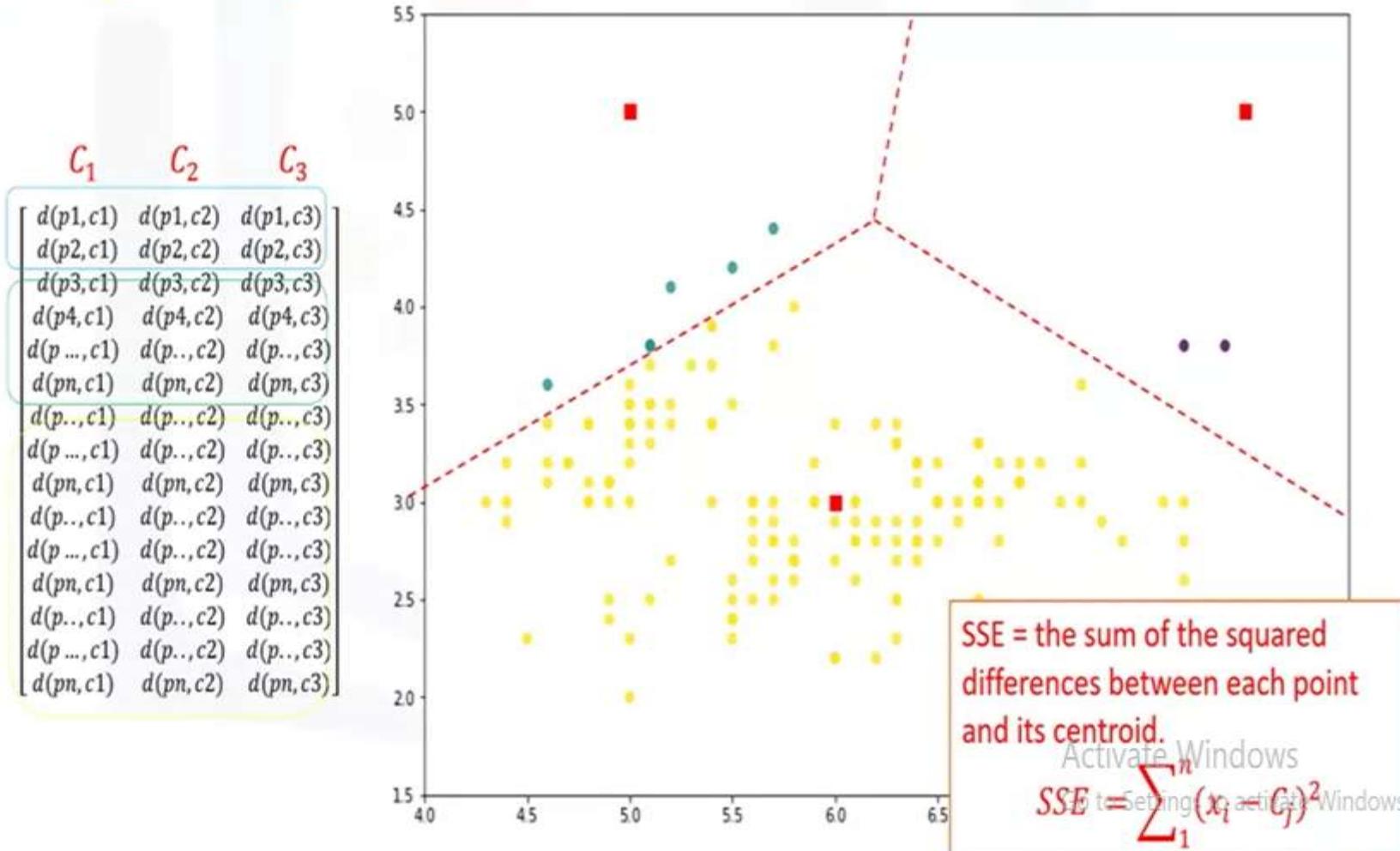
# k-Means clustering – assign to centroid:

## 3) Assign each point to the closest centroid



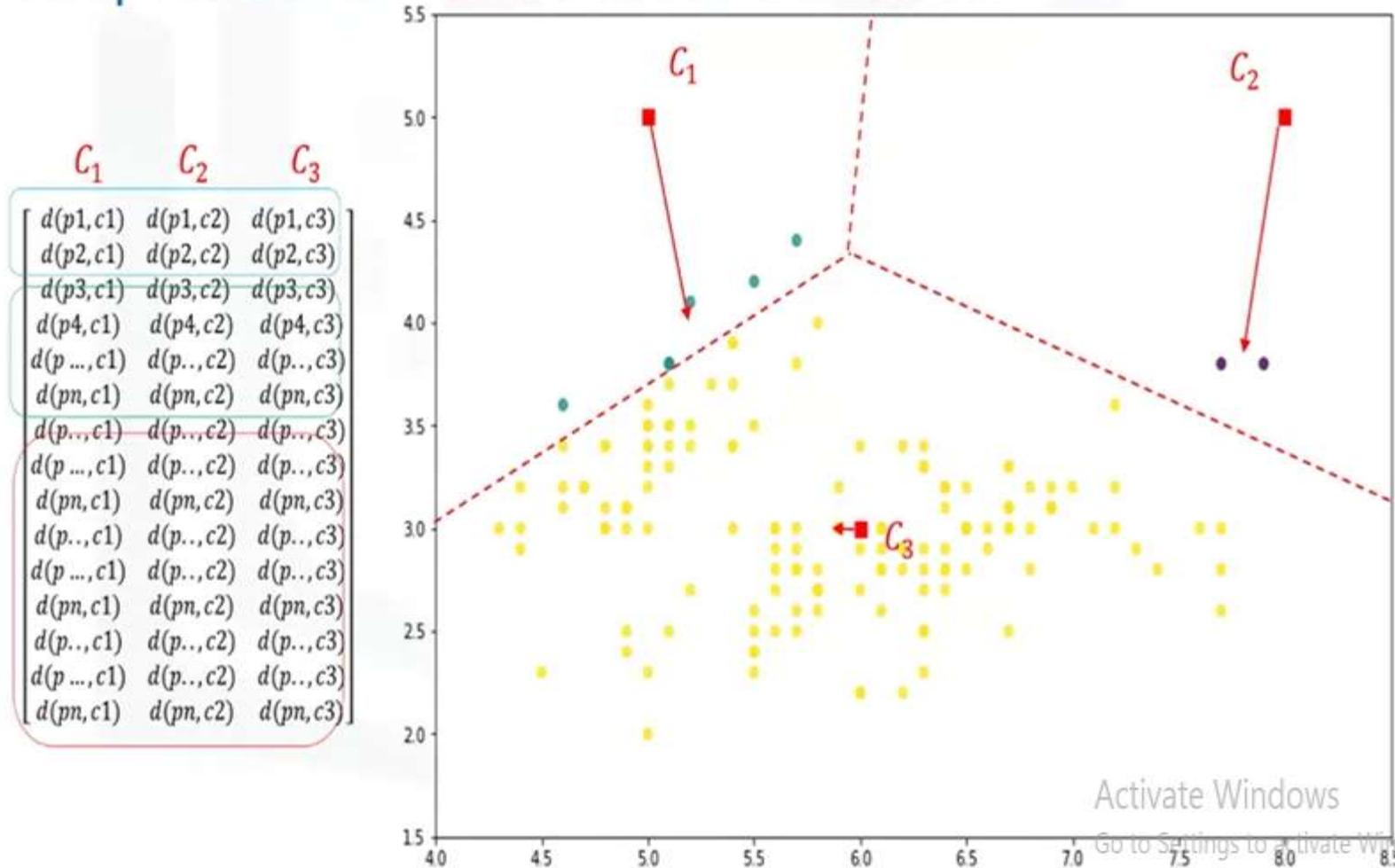
## k-Means clustering – assign to centroid:

3) Assign each point to the closest centroid



# k-Means clustering – compute new centroids:

4) Compute the new centroids for each cluster.

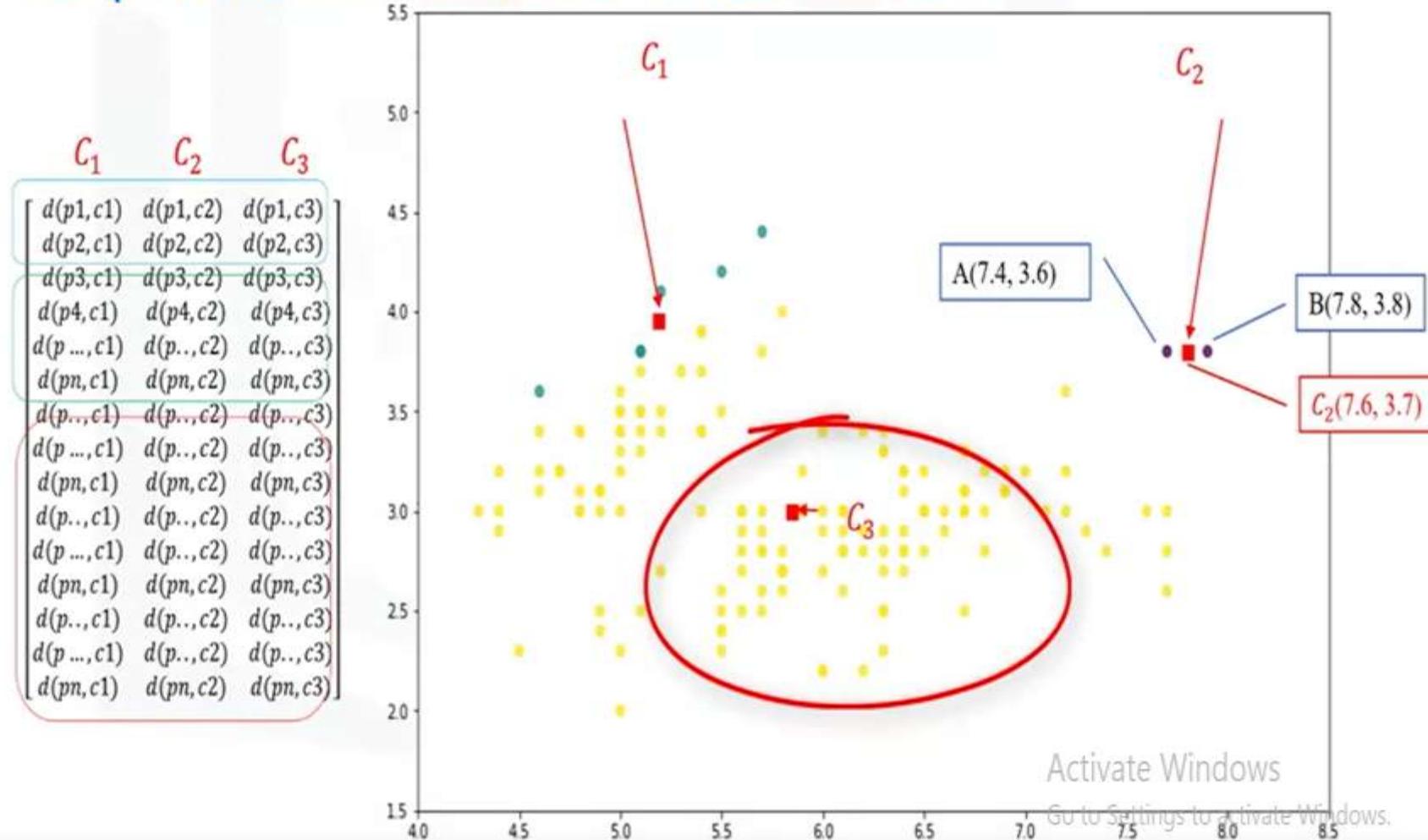


Activate Windows

Go to Settings to activate Windows.

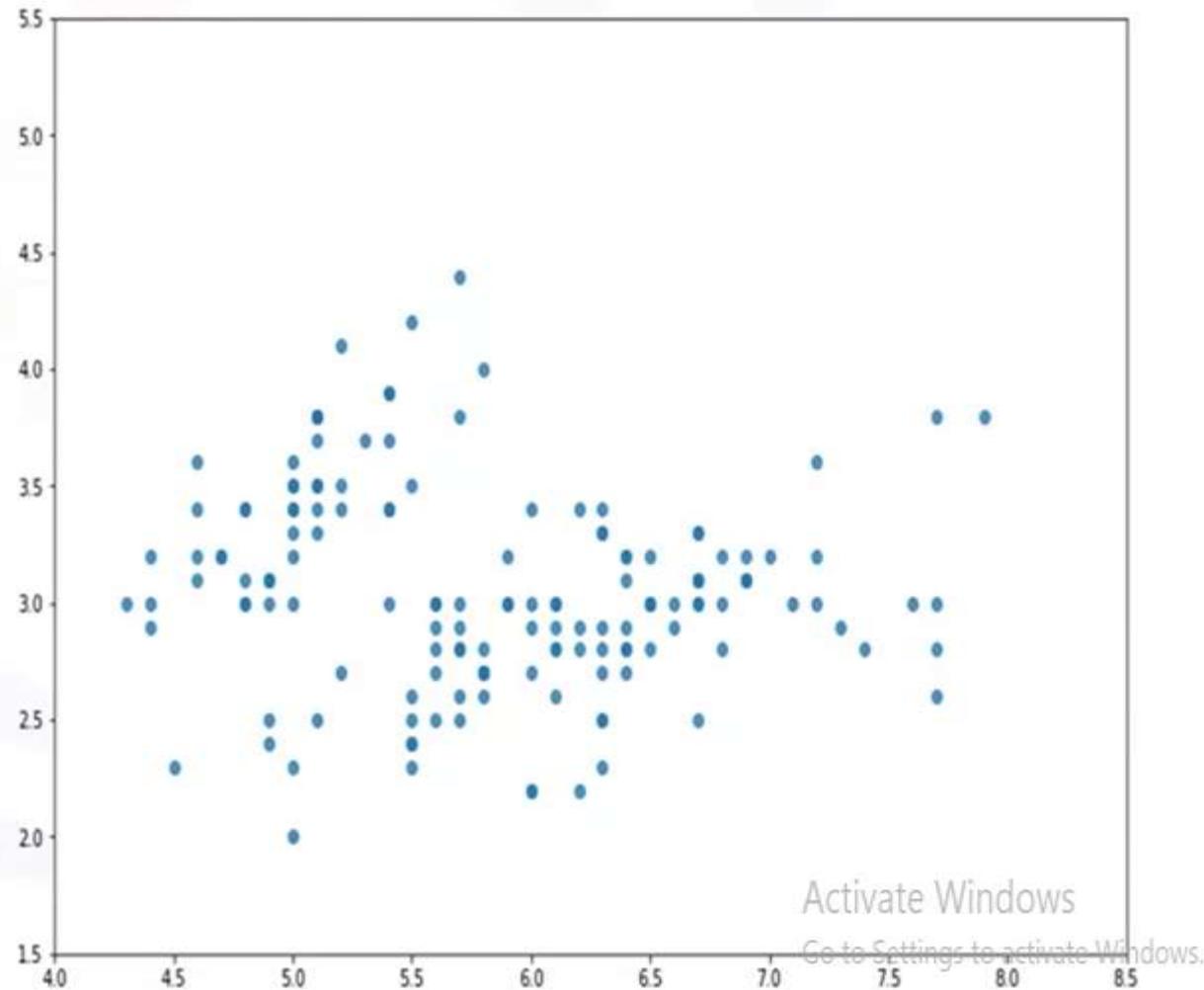
# k-Means clustering – compute new centroids:

4) Compute the new centroids for each cluster.



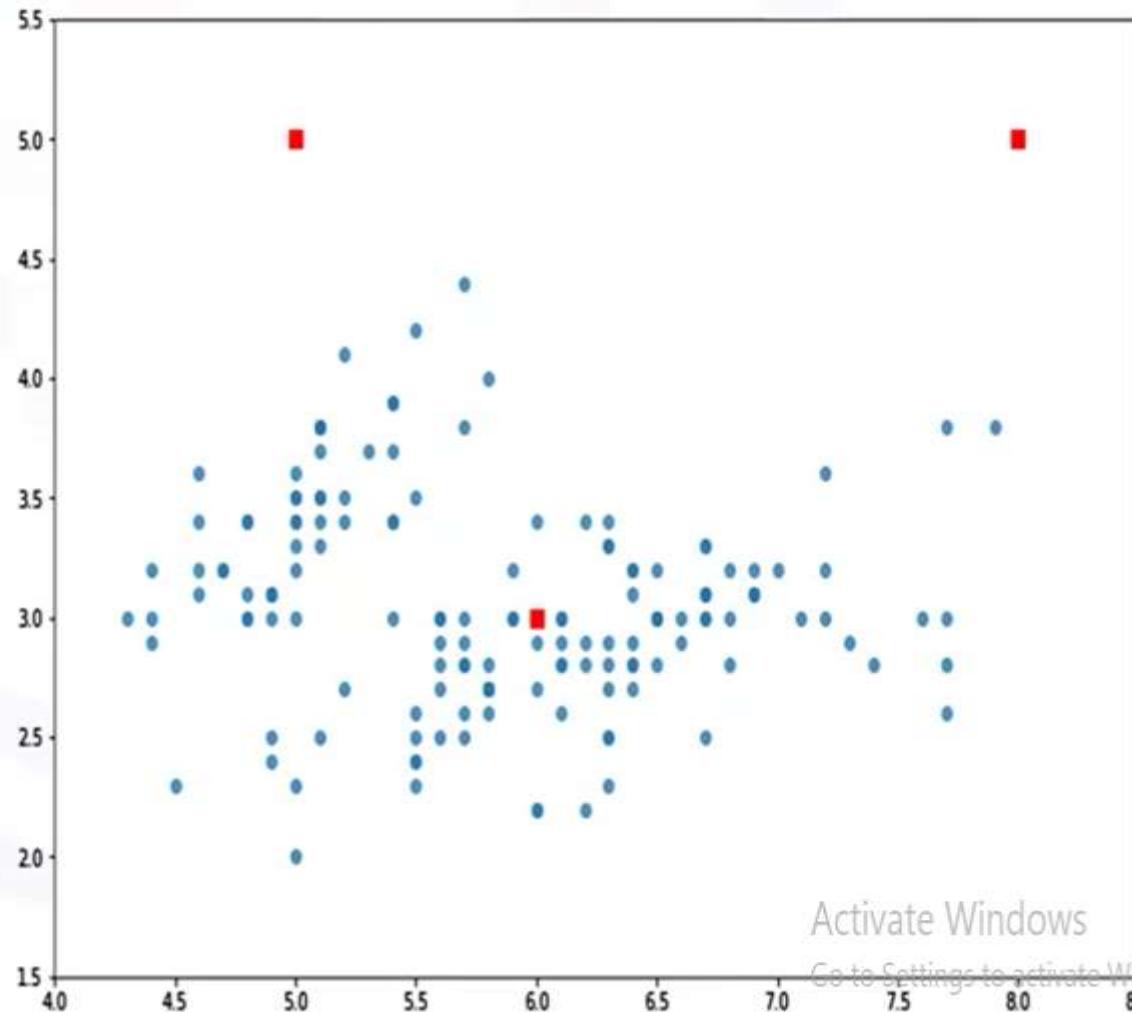
# k-Means clustering – repeat:

5) Repeat until there  
are no more changes.



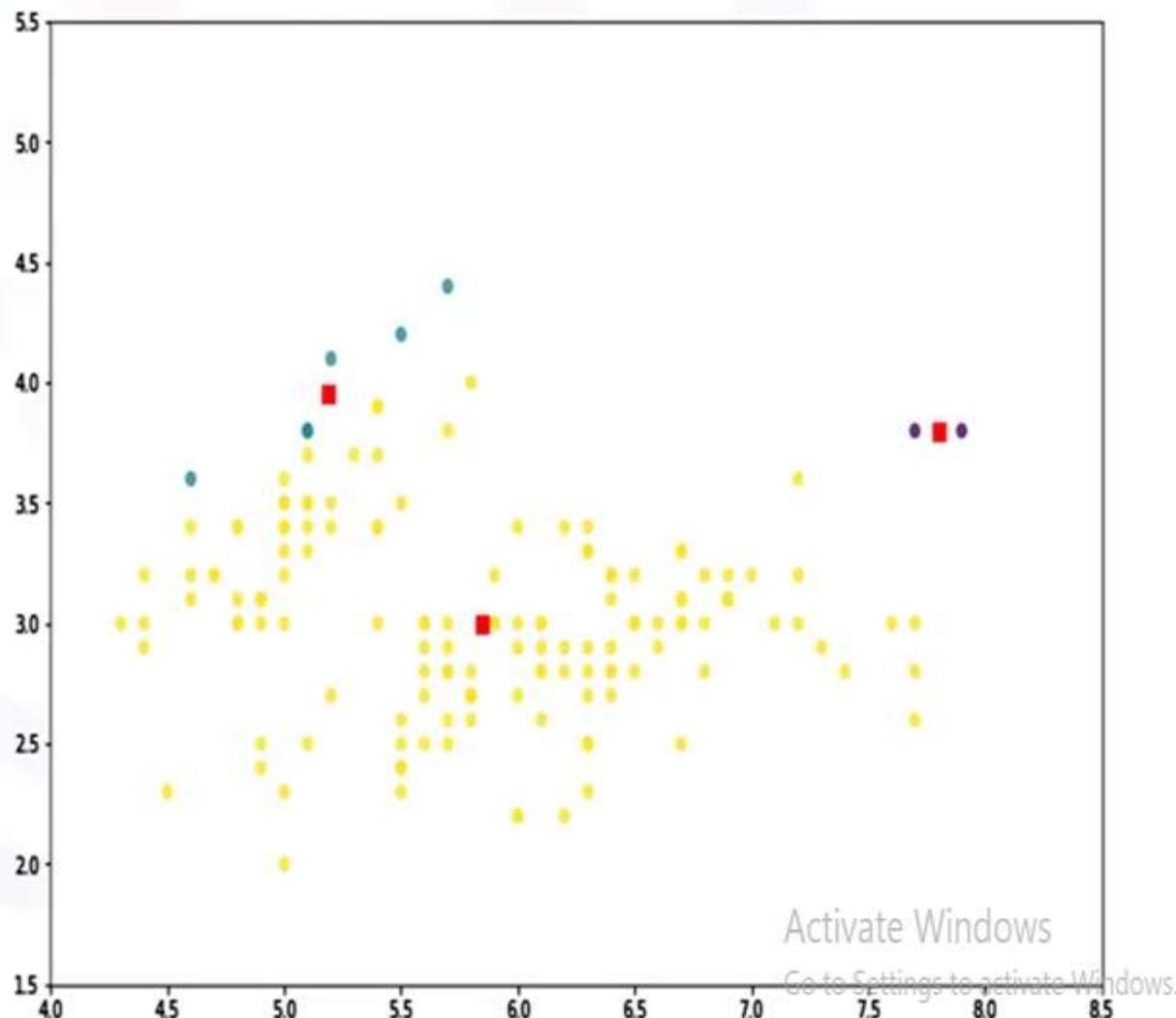
# k-Means clustering – repeat:

5) Repeat until there  
are no more changes.



# k-Means clustering – repeat:

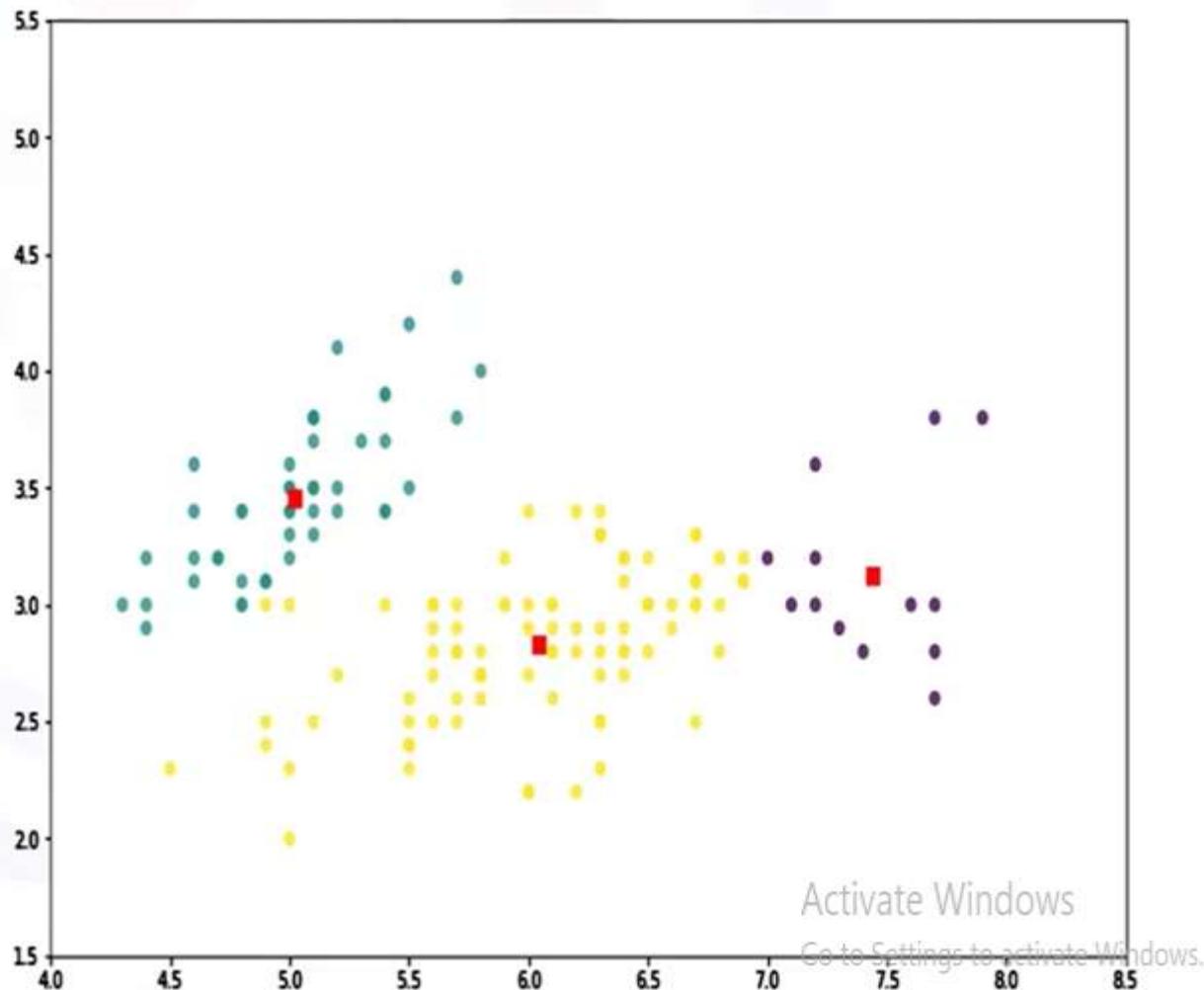
5) Repeat until there  
are no more changes.



Activate Windows  
Go to Settings to activate Windows.

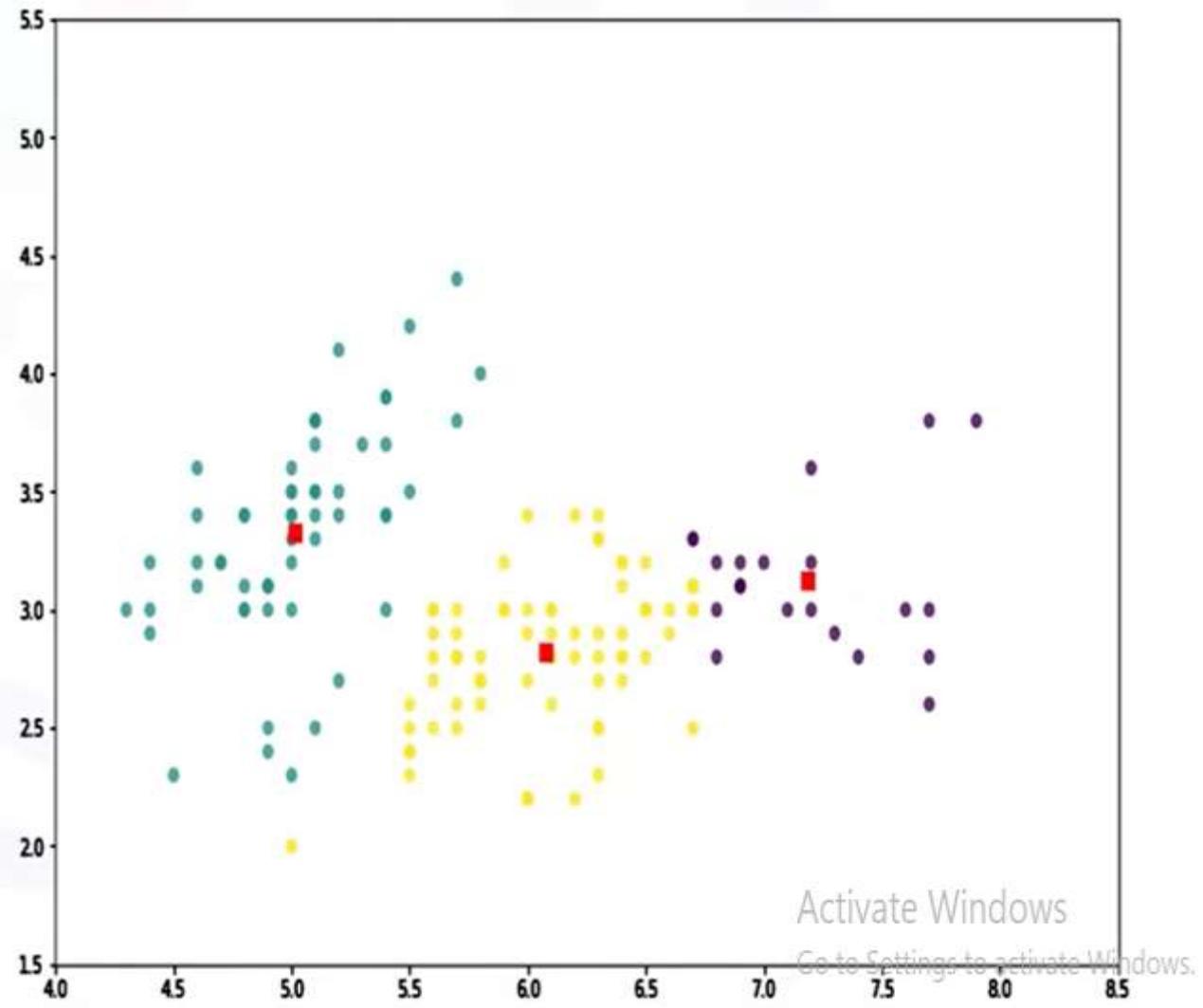
# k-Means clustering – repeat:

5) Repeat until there  
are no more changes.



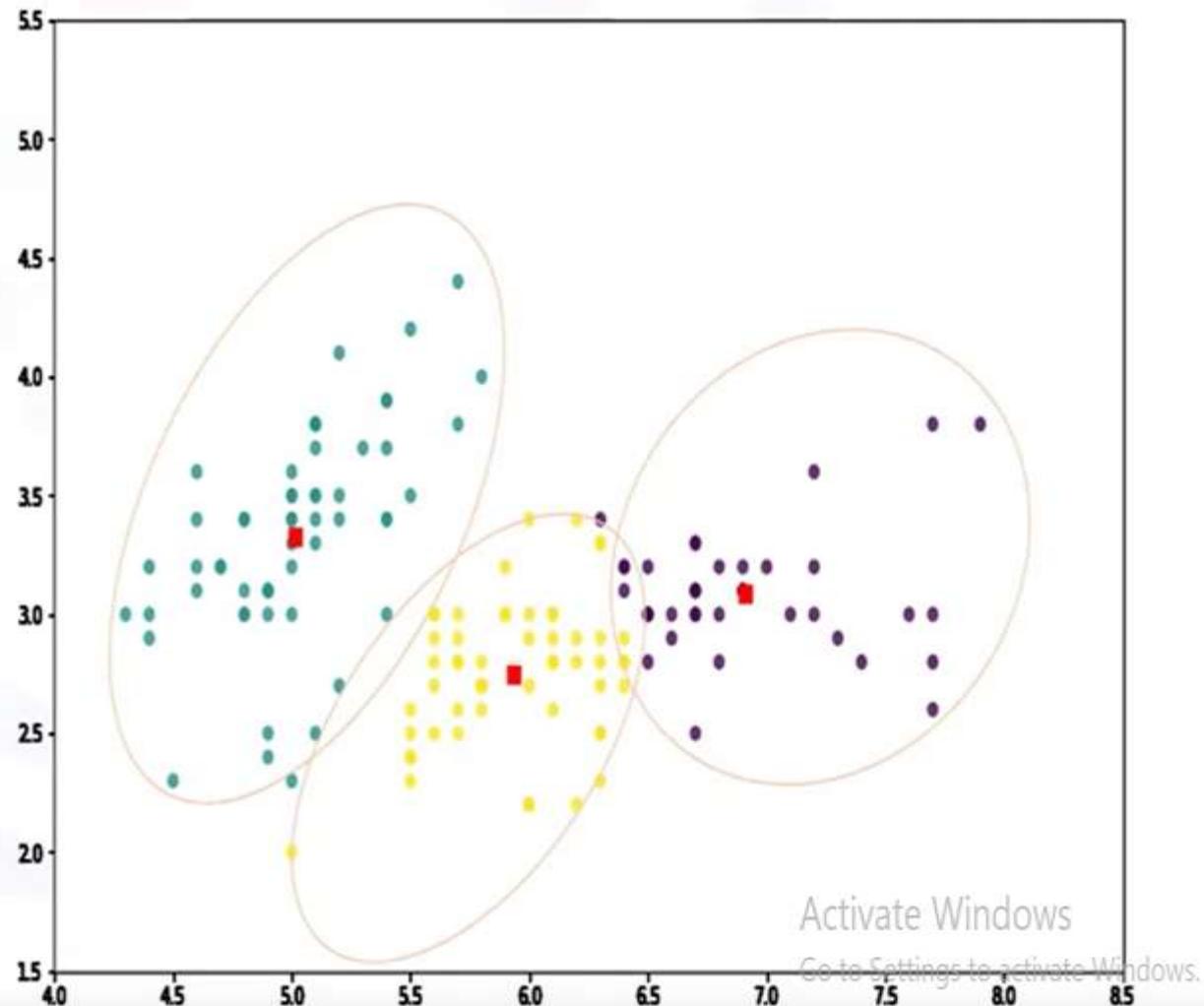
# k-Means clustering – repeat:

5) Repeat until there  
are no more changes.



# k-Means clustering – repeat:

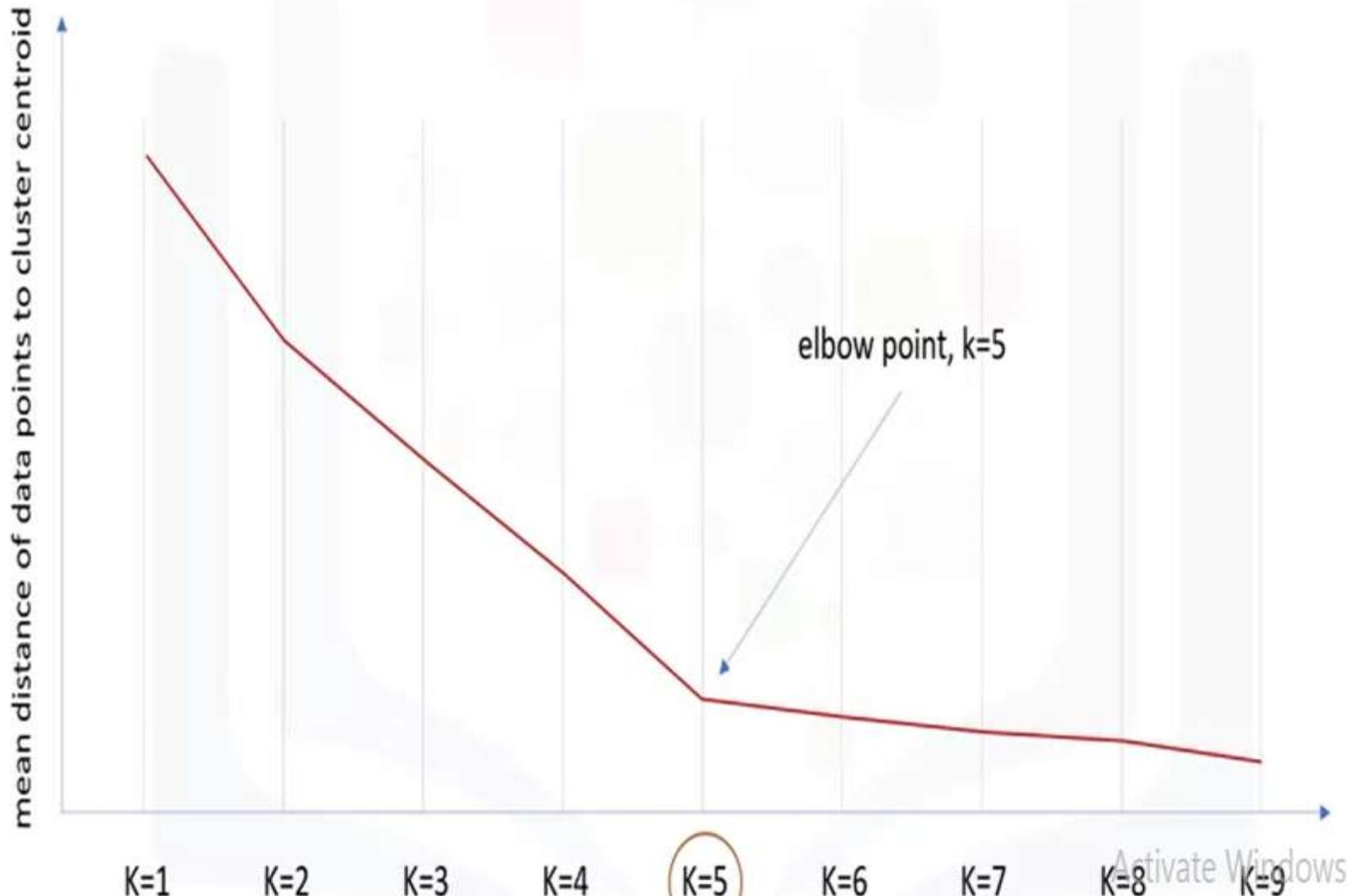
5) Repeat until there  
are no more changes.



## k-Means clustering algorithm:

1. Randomly placing  $k$  centroids, one for each cluster.
2. Calculate the distance of each point from each centroid.
3. Assign each data point (object) to its closest centroid, creating a cluster.
4. Recalculate the position of the  $k$  centroids.
5. Repeat the steps 2-4, until the centroids no longer move.

# Choosing k:



## K-Means recap:

- Med and Large sized databases (*Relatively efficient*)
- Produces sphere-like clusters
- Needs number of clusters (k)



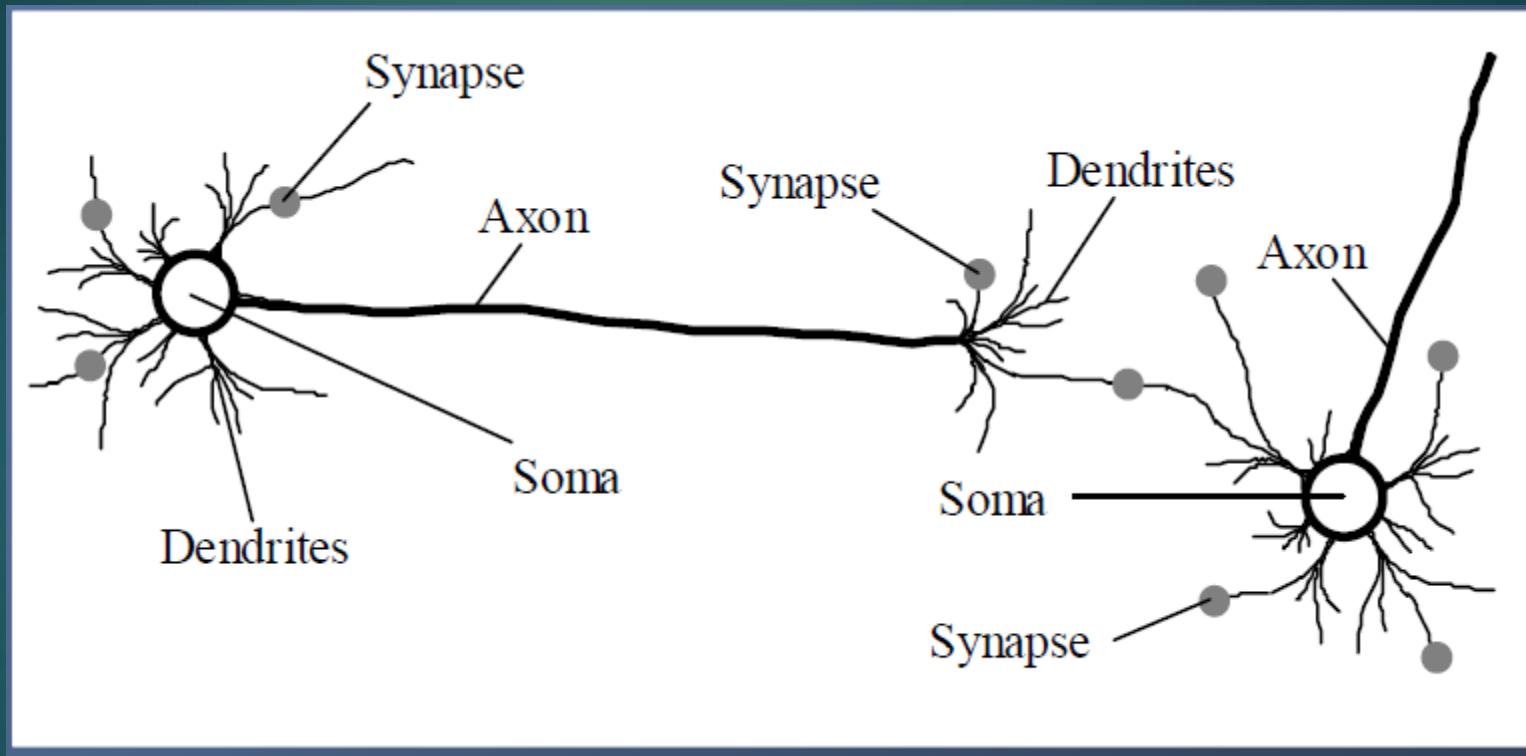
Neural  
network

# Introduction, or how the brain works

- Machine learning involves adaptive mechanisms that enable computers to learn from experience, learn by example and learn by analogy.
- Learning capabilities can improve the performance of an intelligent system over time. The most popular approaches to machine learning are **artificial neural networks** and **optimization algorithms**.

- A **neural network** can be defined as a model of reasoning based on the human brain.
- The brain consists of a densely interconnected set of nerve cells called *neurons*. These neurons are also known as basic information-processing units.
- The human brain incorporates nearly 10 billion neurons and 60 trillion connections, *synapses*, between them.
- By using multiple neurons simultaneously, the brain can perform its functions much faster than the fastest computers in existence today.

# Biological neural network



- Each neuron has a very simple structure, but an army of such elements constitutes a tremendous processing power.
- A neuron consists of a cell body, **soma**, a number of fibers called **dendrites**, and a single long fiber called the **axon**.
- **Dendrites** are thin and widely branching fibers, reaching out in different directions to make connections to a larger number of cells within the cluster.
- Input connection are made from the **axons** of other cells to the **dendrites** or directly to the body of the cell. These are known as axondentritic and axonsomatic **synapses**.

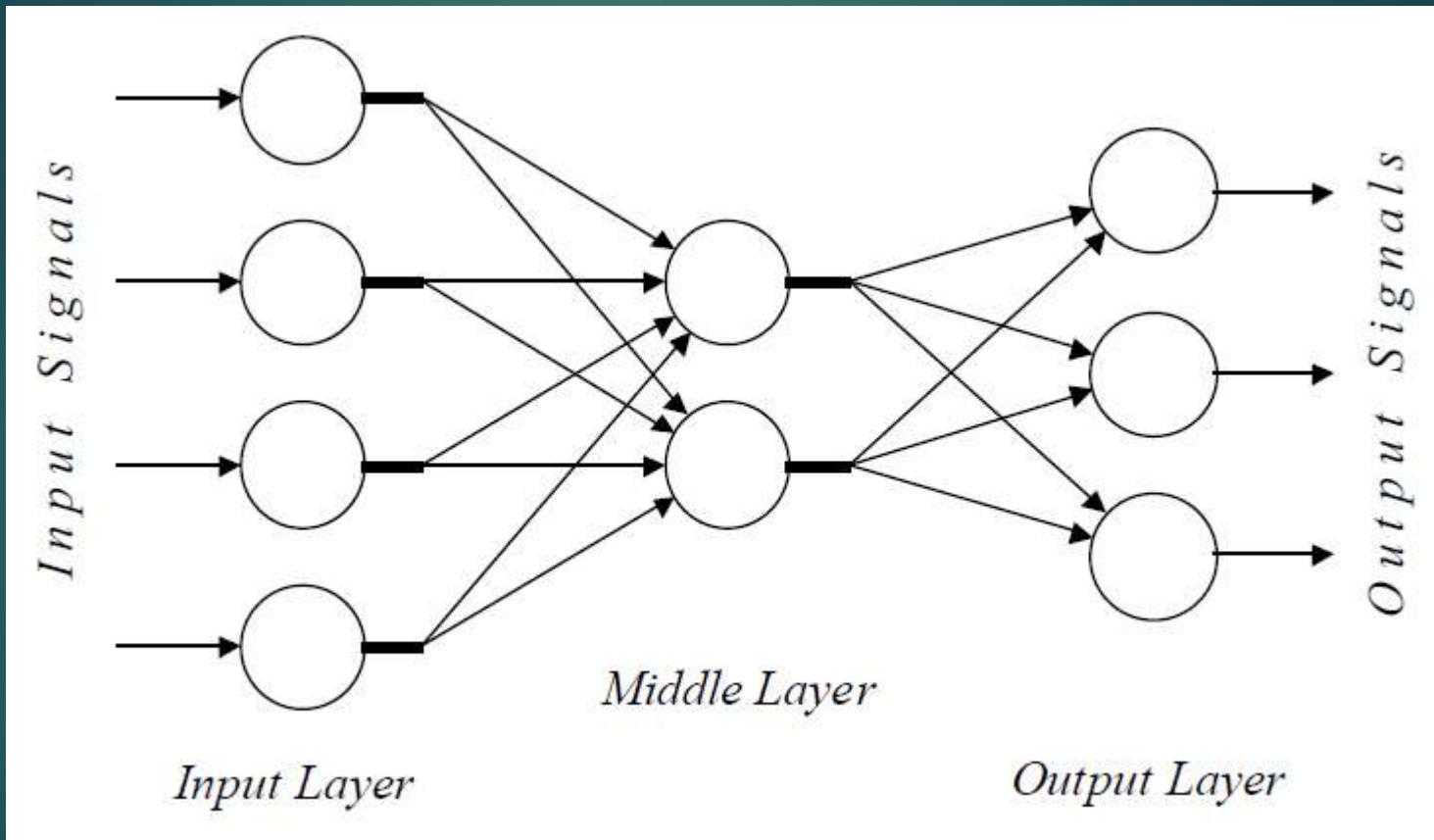
- There is only one axon per neuron.
- *Axon* is a single and long fiber, which transports the output signal of the cell as electrical impulses (action potential) along its length. The end of the axon may divide in many branches, which are then connected to other cells. The branches have the function to fan out the signal to many other inputs.

- ▶ Our brain can be considered as a highly complex, non-linear and parallel information-processing system.
- ▶ Information is stored and processed in a neural network simultaneously throughout the whole network, rather than at specific locations. In other words, in neural networks, both data and its processing are **global** rather than local.
- ▶ Learning is a fundamental and essential characteristic of biological neural networks. The ease with which they can learn led to attempts to emulate a biological neural network in a computer.

# Artificial neural network

- ▶ An artificial neural network consists of a number of very simple processors, also called **neurons**, which are analogous to the biological neurons in the brain.
- ▶ The neurons are connected by weighted links passing signals from one neuron to another.
- ▶ The output signal is transmitted through the neuron's outgoing connection. The outgoing connection splits into a number of branches that transmit the same signal. The outgoing branches terminate at the incoming connections of other neurons in the network.

# Architecture of a typical artificial neural network

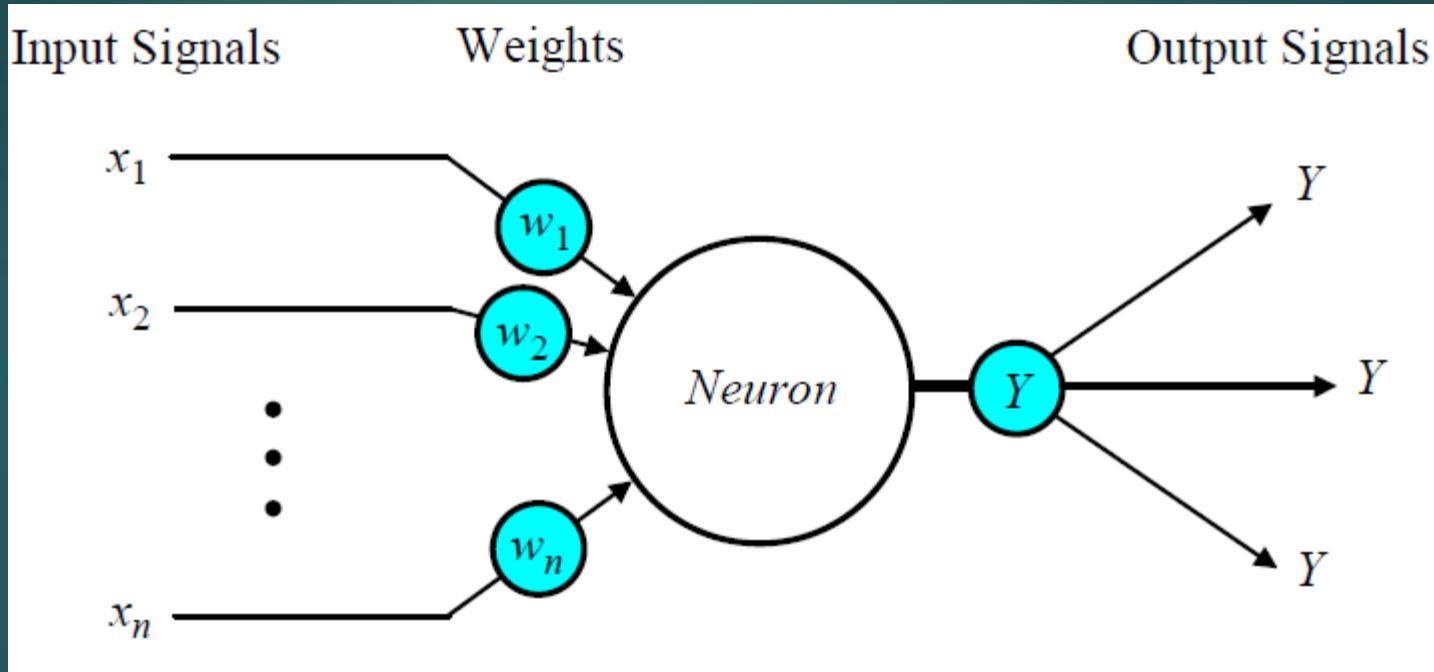


# Analogy between biological and artificial neural networks

<i>Biological Neural Network</i>	<i>Artificial Neural Network</i>
Soma	Neuron
Dendrite	Input
Axon	Output
Synapse	Weight

# The neuron as a simple computing element

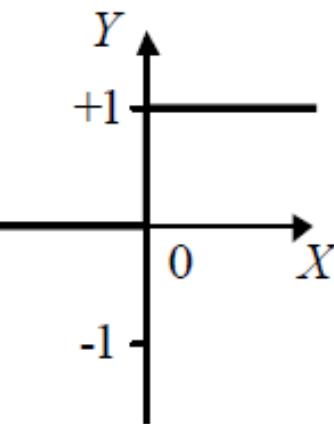
Diagram of a neuron



- ▶ The neuron computes the weighted sum of the input signals and compares the result with a **threshold value**,  $\theta$  (theta).
- ▶ If the net input is less than the threshold, the neuron output is  $-1$ . But if the net input is greater than or equal to the threshold, the neuron becomes activated and its output attains a value  $+1$ .
- ▶ The neuron uses the following transfer or **activation function**:  
$$X = \sum_{i=1}^n x_i w_i \quad Y = \begin{cases} +1, & \text{if } X \geq \theta \\ -1, & \text{if } X < \theta \end{cases}$$
- ▶ This type of activation function is called a **sign function**.

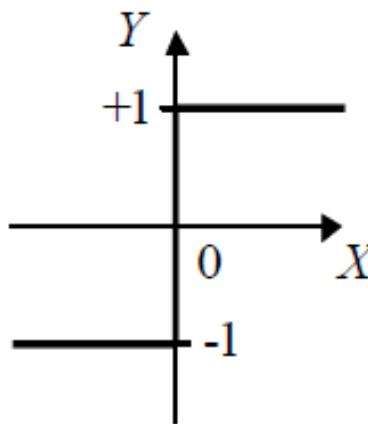
# Activation functions of a neuron

*Step function*



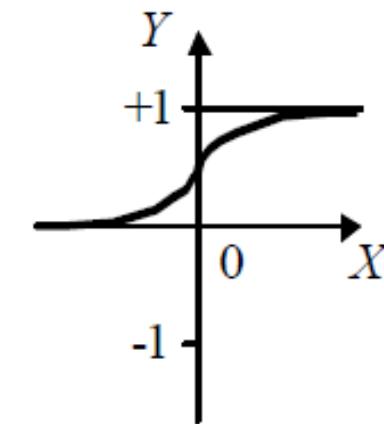
$$Y^{step} = \begin{cases} 1, & \text{if } X \geq 0 \\ 0, & \text{if } X < 0 \end{cases}$$

*Sign function*



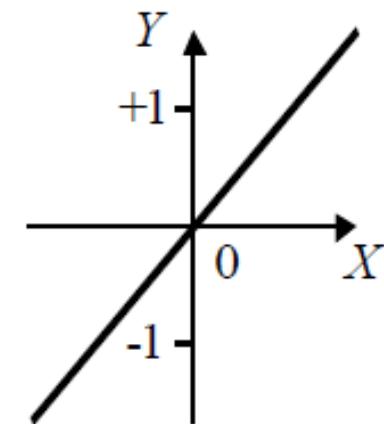
$$Y^{sign} = \begin{cases} +1, & \text{if } X \geq 0 \\ -1, & \text{if } X < 0 \end{cases}$$

*Sigmoid function*



$$Y^{sigmoid} = \frac{1}{1+e^{-X}}$$

*Linear function*

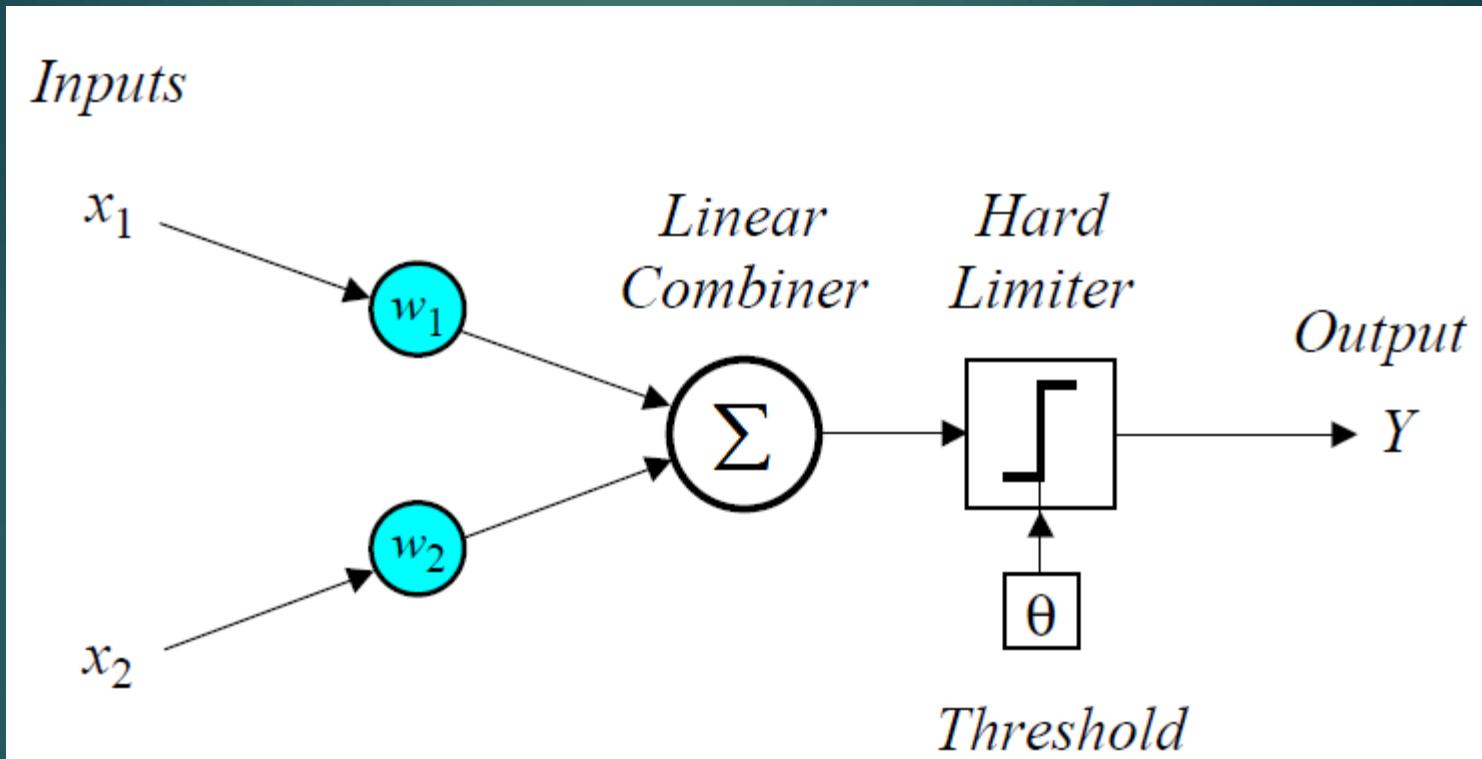


$$Y^{linear} = X$$

# Can a single neuron learn a task? -The Perceptron

- ▶ In 1958, **Frank Rosenblatt** introduced a training algorithm that provided the first procedure for training a simple ANN: a **perceptron**.
- ▶ The perceptron is the simplest form of a neural network. It consists of a single neuron with *adjustable* synaptic weights and a *hard limiter*.
- ▶ The model consists of a linear combiner followed by a hard limiter.
- ▶ The weighted sum of the inputs is applied to the hard limiter, which produces an output equal to +1 if its input is positive and -1 if it is negative.

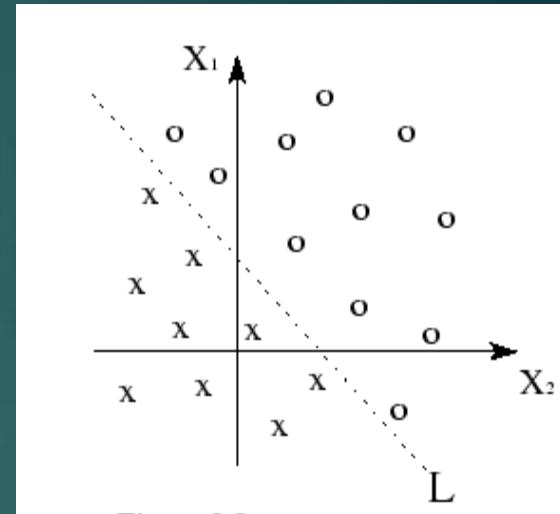
# Single-layer two-input perceptron



- In the case of an elementary perceptron, the  $n$  dimensional space is divided by a *hyperplane* into two decision regions. The hyperplane is defined by the ***linearly separable function***:

# Linear Separability

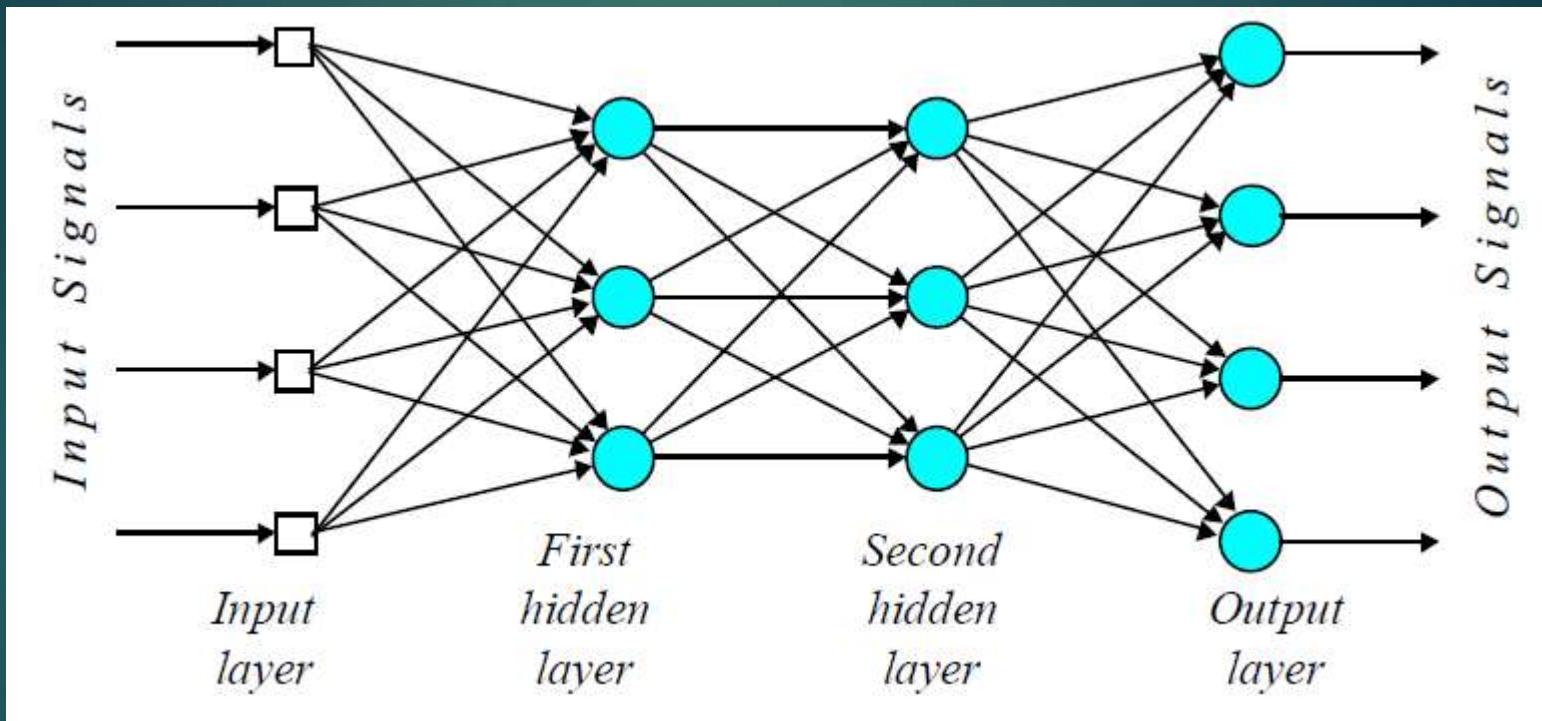
- Consider two-input patterns being classified into two classes.
- Each point with either symbol of X or O represents a pattern with a set of values.
- Each pattern is classified into one of two classes. Notice that these classes can be separated with a single line. They are known as *linearly separable* patterns.
- ***Linear separability*** refers to the fact that classes of patterns with n-dimensional vector can be separated with a single *decision surface*.



# Multilayer neural networks

- ▶ A multilayer perceptron is a feedforward neural network with one or more hidden layers.
- ▶ The network consists of an **input layer** of source neurons, at least one middle or **hidden layer** of computational neurons, and an **output layer** of computational neurons.
- ▶ The input signals are propagated in a forward direction on a layer-by-layer basis.

# Multilayer perceptron with two hidden layers



# Opportunities and Challenges

- ▶ The field of machine learning is quite new, unexplored and rapidly expanding due to new formalizations of the learning problems. For example, many of the *algorithms rely on a lot of parameters*. Thus an innovative way of tuning them can lead to significant rise in throughput.
- ▶ There are many more scopes of advancements which involve the learning of natural way of processing things like humans or animals do. However it raises concerns like *privacy issues* also since it involves manipulation of private data generated by users, which can be a huge security risk.
- ▶ In spite of these insecurities, ML is bounded to be more popular since it allows to solve problems by learning them in place of traditional ways.