

# Classification

# Agenda

## Key Takeaways-

- Classification, its types and algorithms
- k Nearest Neighbors (kNN)
- Types of distance measures
- Hyperparameter Tuning, GridSearchCV
- Cross validation and its types
- K-fold Cross Validation
- Evaluation measures for Classification

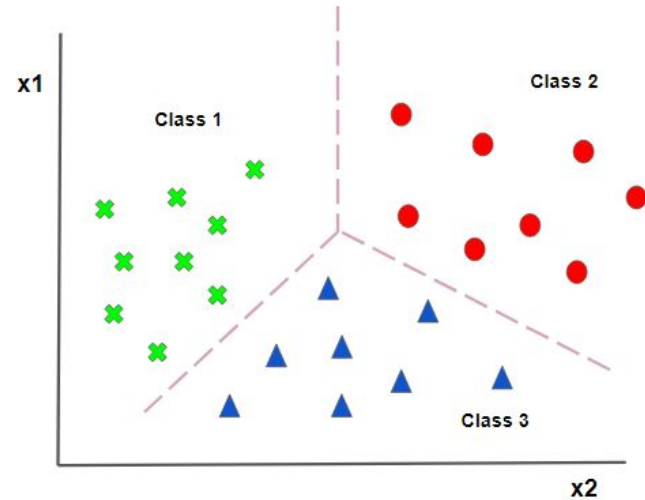
# Classification

Classification is a Supervised ML technique that helps in identifying/predicting the class/category for (new) instances.

If the target variable is categorical/qualitative in nature, then we use classification algorithms.

For example, Identifying

- Whether an email is spam or ham.
- Whether a transaction is fraudulent or not
- The handwritten digits 0-9
- Types of debit cards
- Whether a patient is cancerous or non-cancerous



# Types of Classification Tasks

Based upon the number of unique classes in the target variable, there are 3 types of classification tasks.

1. **Binary Classification** : The target variable has only two unique classes and the classification task is to identify/predict one of two classes for each instance.  
E.g. Identifying whether a transaction is fraudulent or not.
2. **Multi-class Classification** : The target variable has more than two unique classes and the classification task is to identify/predict one of more than two classes for each instance.  
E.g. Handwritten digit recognition
3. **Multi-label Classification** : Each instance of the labelled data belongs to one or more classes and the classification task is to identify/predict one or more classes for each instance.  
E.g. Predicting tags for Quora questions

# Types of Classification Algorithms

The most commonly used classification algorithms are -

1. k Nearest Neighbors (kNN)
2. Decision Tree
3. Random Forest
4. Logistic Regression
5. Support Vector Machine (SVM)

## k Nearest Neighbors (kNN)

k-NN algorithm assumes the similarity/distance between the new data point and available data points and put the new data point into a category having maximum votes in the neighborhood.

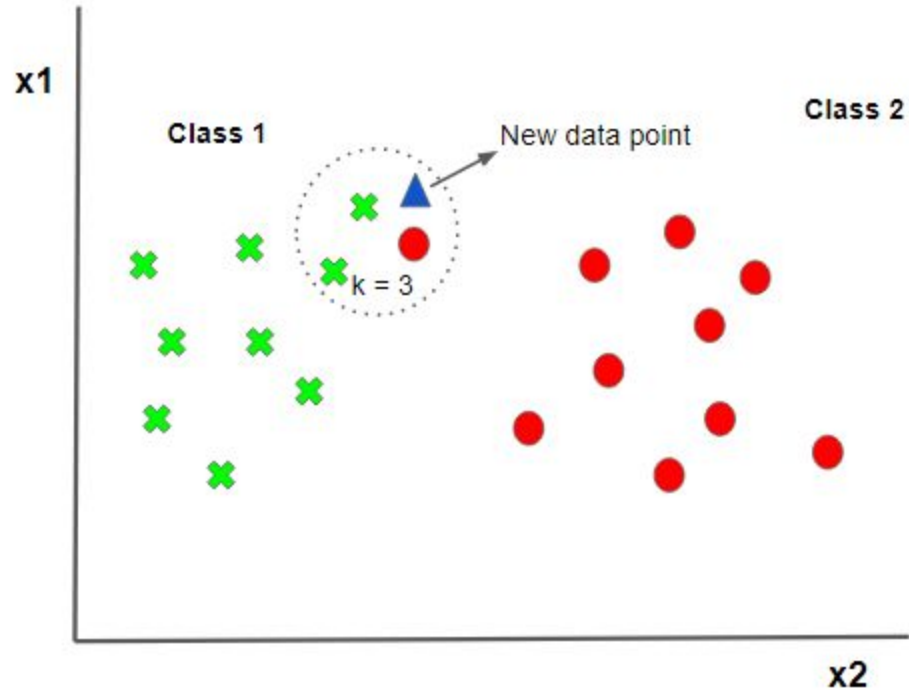
**Step 1.** Choose the  $k$  (number of nearest neighbors).

**Step 2.** Take the  $k$  nearest neighbors of new data point, based on distance (Euclidean)

**Step 3.** Among the  $k$  neighbors, count the number of data points in each category.

**Step 4.** Assign the new data point to the category with maximum counts.

## k Nearest Neighbors (kNN) [Contd.]



# Distance measures

The most commonly used distance measures in kNN are -

- Euclidean distance
- Manhattan distance
- Minkowski distance
- Cosine distance



# Euclidean distance

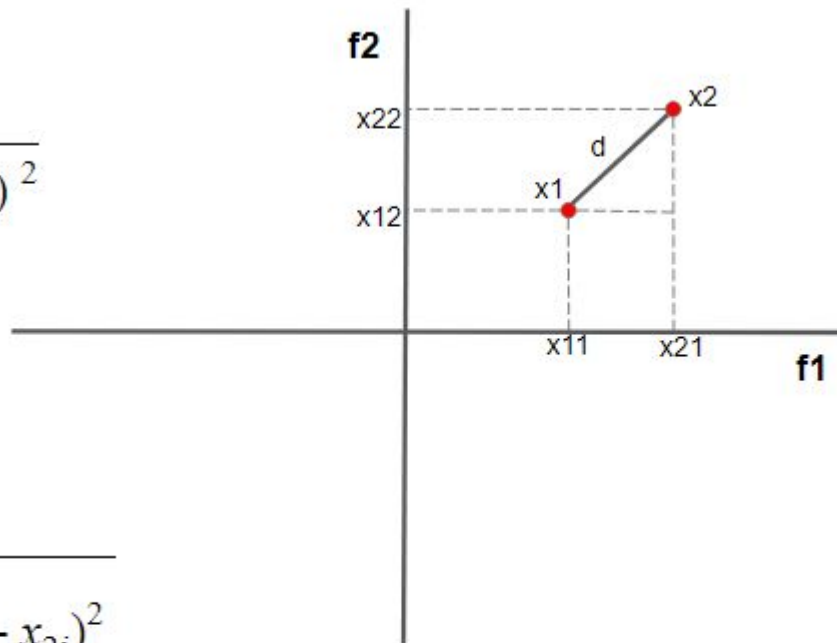
- Euclidean distance ( $d$ )

=length of shortest line from  $x_1$  to  $x_2$

$$d = \sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2}$$

- Euclidean distance is also represented as L2 norm.

$$d = \|x_1 - x_2\|_2 = \sqrt{\sum_{i=1}^d (x_{1i} - x_{2i})^2}$$



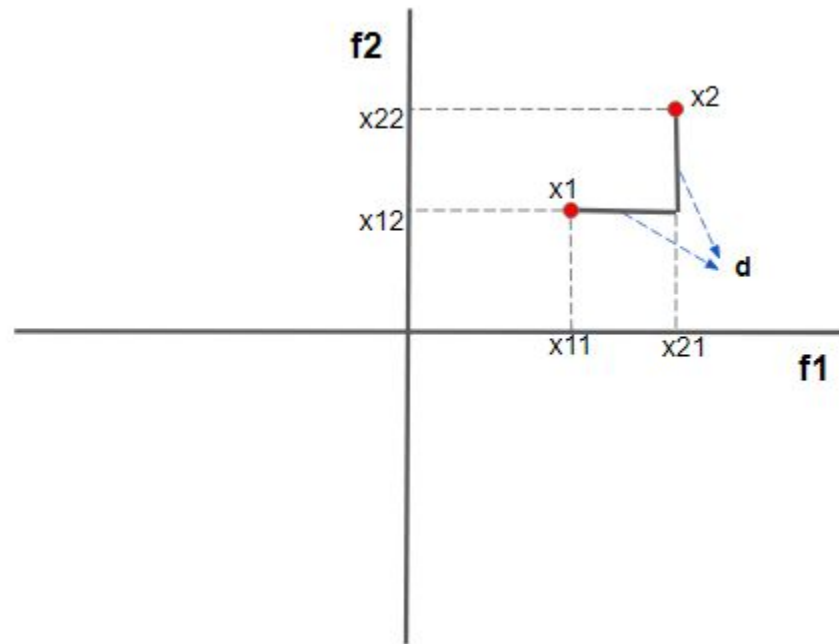
# Manhattan Distance

- Manhattan distance ( $d$ )  
= Sum of absolute differences

$$d = |x_{21} - x_{11}| + |x_{22} - x_{12}|$$

- Its name is from the city “Manhattan”  
Where roads are perpendicular to each other.
- Manhattan distance is also represented as  
L1 norm.

$$d = \|x_1 - x_2\|_1 = \sum_{i=1}^d |x_{1i} - x_{2i}|$$



# Minkowski Distance

- It is a generalized distance measure.
- Minkowski distance is also represented as Lp norm.

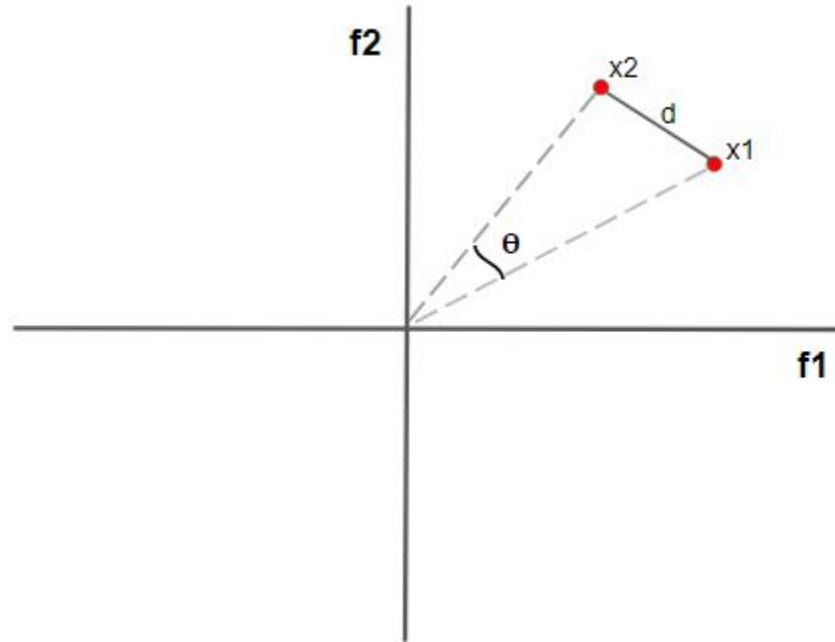
$$\left\| x_1 - x_2 \right\|_p = \left( \sum_{i=1}^d \left| x_{1i} - x_{2i} \right|^p \right)^{1/p}$$

- If  $p = 1$ , then Minkowski distance = Manhattan distance
- Similarly, if  $p = 2$  then Minkowski distance = Euclidean distance

# Cosine Distance

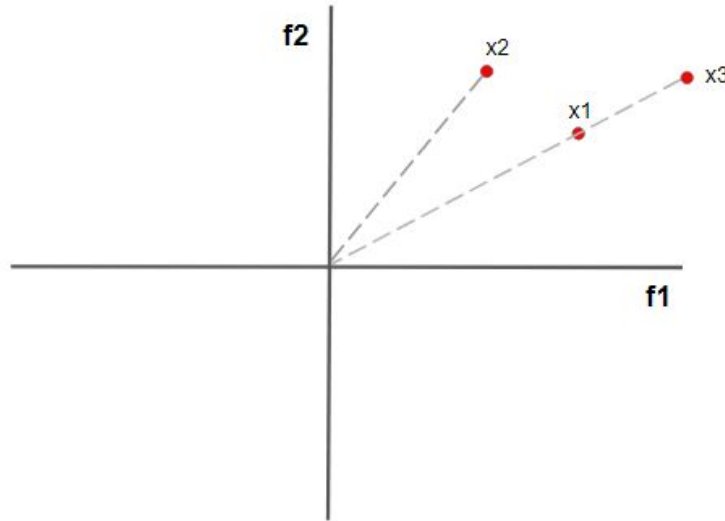
- Cosine distance =  $(1 - \text{Cosine similarity})$

Where Cosine similarity  $\equiv \cos \theta$  and  $\theta$  is the angle between two vectors.



# Quiz 1

Choose the correct statement(s) as per the following vectors.



- Cosine distance  $(x_1, x_3) <$  Euclidean distance  $(x_1, x_3)$
- Euclidean distance  $(x_2, x_3) >$  Euclidean distance  $(x_1, x_2)$
- Cosine distance  $(x_1, x_2) \neq 0$
- All of the above

## Quiz 2

Choose the correct statement(s).

- As the distance increases the similarity b/w two vectors(entities) also increases.
- As the distance decreases the similarity b/w two vectors(entities) also increases.
- Euclidean distance is similar to L2 norm.
- None of the above

# Model Parameter vs Hyperparameter

Parameters	Hyperparameters
It is a configuration <b>variable</b> that is <b>internal</b> to the <b>model</b>	It is a configuration that is <b>external</b> to the <b>model</b>
They are <b>estimated</b> or <b>learned</b> from <b>data</b> (often <b>not</b> set <b>manually</b> )	They are often <b>specified</b> by the <b>practitioner</b> (or often set <b>manually</b> )
They are <b>required</b> by the <b>model</b> when <b>making predictions</b>	They are often <b>tuned</b> for a given <b>predictive modeling problem</b> . This process is known as <b>Hyperparameter tuning</b> .
E.g. Regression <b>coefficients</b> , neural network <b>weights</b> , etc.	E.g. <b>test_size</b> in <b>train-test splitting</b> process, <b>k</b> (# of nearest neighbors) in <b>kNN</b> , <b>Kernel type</b> in <b>SVM</b> , <b>Maximum depth</b> of tree in <b>decision tree</b> , etc.

# Cross Validation

- Validation techniques are used to evaluate how well a model would generalize to new/unseen data.
- Cross-Validation is a validation technique used to estimate how well and accurately a predictive model will perform to an independent dataset (real time scenarios).
- Most commonly used cross-validation strategies are -
  - Validation set approach
  - Leave one out cross validation (LOOCV)
  - K-fold cross validation



# Cross Validation [Contd.]

## 1. Validation set approach

- This approach divides the dataset into two equal parts, 50% of the dataset is reserved for training purpose, whereas the remaining 50% is reserved for validation purpose.
- **Disadvantage** - This approach generally leads to a high bias model as there always remains a possibility of missing out on relevant and meaningful information due to considering only 50% of the data.

## 2. Leave one out cross validation (LOOCV)

- As the name suggests, here we reserve only one data point for validation purpose and use rest of the data for training purpose.
- In this approach, number of folds = number of data points in the dataset.
- **Disadvantages** - 1. Higher execution time as it is repeated for n times.  
2. This approach leads to a high variance model.

# K-fold Cross Validation

**Step 1.** Randomly **split** the entire **dataset** into **k folds/subsets**.

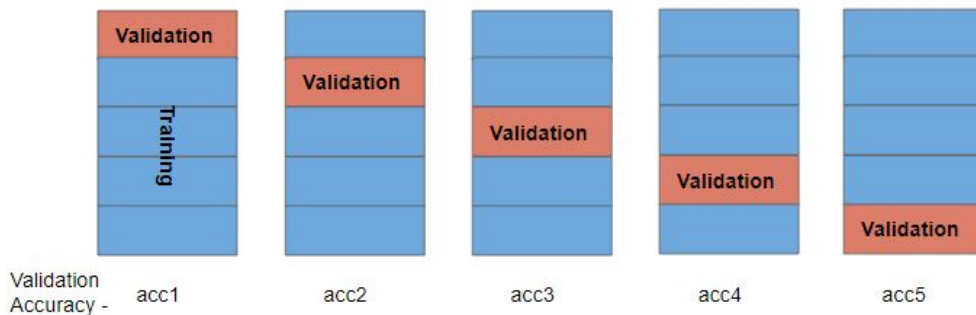
**Step 2.** In **each iteration**(or  $k$ th round), **train** the **model** using  $(k - 1)$  **folds** of the **dataset** and **validate/test** the **model** using the **kth fold**.

**Step 3.** **Calculate** the **accuracy** for this iteration.

**Step 4.** **Repeat** this **process** until each of the **k-folds** has served as the **validation/test** set.

**Step 5.** Take the **average** of all **k** such **accuracies** to get the **final validation** accuracy.

5-fold cross validation



**Final Accuracy** = **Average**(acc1,acc2,acc3,acc4,acc5)

## Quiz 3

In **k-fold cross validation**, how many **subsets** are used for **validation purpose** in each **iteration**?

- **(k - 1) subsets**
- **K subsets**
- **One fixed subset** in each **iteration**
- **One randomly selected subset** in each **iteration**

# GridSearchCV

- **Grid Search** is the **process** of performing **hyperparameter tuning** in order to **determine** the **optimal values** of the **hyperparameters** for a given **model**.
- **Manually** performing **hyperparameter tuning** -
  - is a **time consuming** process.
  - Also, it is **very hard** to keep **continuous track** of **hyperparameters** which we have tried and still have to try.
- **GridSearchCV** is a **built-in** function in **sklearn's model\_selection** package to **perform hyperparameter tuning** more **efficiently** and **effectively**.
- **GridSearchCV** tries all the **combinations** of the **values** passed and **evaluates** the **estimator(model)** for **each combination**.
- So **GridSearchCV** is used to find the **optimal hyperparameters** of a **model** which **results** in the most **accurate predictions**.

# kNN Time Complexity

- Test time complexity, for a dataset with  $n$  samples and  $d$  features =  $O(nd) + O(1) + O(1) \approx O(nd)$

Where,  $O(nd)$  for comparing  $n$  data points and each point is a  $d$  dimensional vector.

$O(1) + O(1)$  - to perform majority voting.

- So kNN consumes too much time in making predictions.

# Evaluation metrics for Classification

- The most commonly used and simplest evaluation measure for classification is Accuracy.
- Accuracy is the ratio of total number of correctly classified observations to the total observations (total predictions made).

Mathematically,

$$\text{Accuracy} = \frac{\text{Number of correct classifications}}{\text{Total number of predictions made}}$$

**Caveats** - Accuracy gives a false sense of evaluation while dealing with imbalanced data, where most of observation belongs to one class(majority class).

Therefore, we use various other types of evaluation metrics for classification tasks.

- Precision and Recall
- F1-score

# Confusion Matrix

- **Confusion matrix** provides a more **intuitive** way to **count/know** the **number** of **correct** and **incorrect classifications** for all the **classes**.
- **Confusion matrix** is a **NxN matrix**, where **N** is the **total number** of **classes**.

For a **binary classification** problem, the **confusion matrix** is as shown below.

		Predicted	
		Positives	Negatives
Actual	Positives	TP	FN
	Negatives	FP	TN

- **True Positive (TP)** - **Actual class** is **positive** and the **model** also **predicted** as **positive**.
- **False Positive (FP)** - **Actual class** is **negative** but the **model predicted** as **positive**.
- **False Negative (FN)** - **Actual class** is **positive** but the **model predicted** as **negative**.
- **True Negative (TN)** - **Actual class** is **negative** and the **model** also **predicted** as **negative**.

# Evaluation measures using Confusion Matrix

$$Accuracy = (TP + TN) / (TP + FP + FN + TN)$$

## Precision -

**Precision** tells how **precise**(**sure**) we are about the **predictions**. In other words, out of **total positive predictions** how many are **actually positive**.

$$Precision = TP / (TP + FP)$$

## Recall -

**Recall** is out of **total actual positives** how many our model **predicted** as **positive**.

$$Recall = TP / (TP + FN)$$

## F1-score -

**F1-score** combines both **precision** and **recall**. It is a **weighted average** (**harmonic mean** especially) of both **precision** and **recall**.

$$F1\ score = (2 * Precision * Recall) / (Precision + Recall)$$



## Quiz 4

What is the **precision** for **defaulter class** based on the below **confusion matrix**?

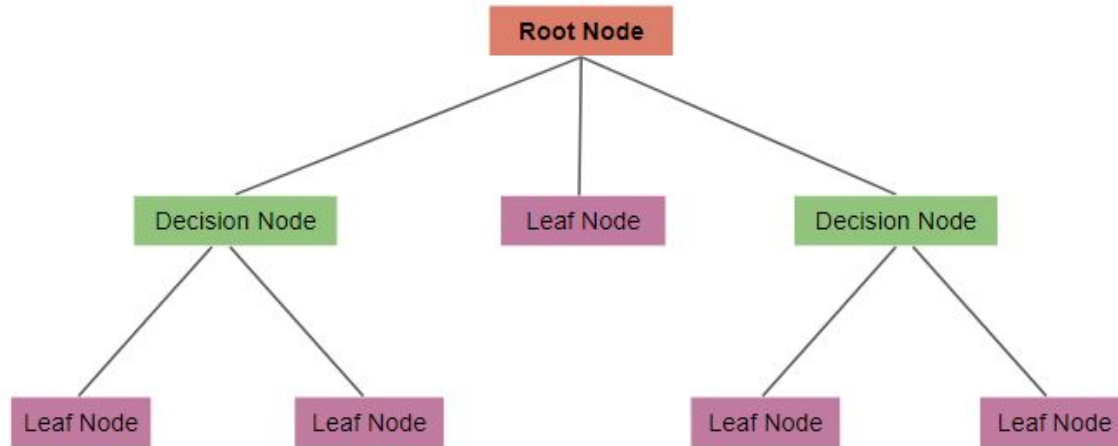
Here, the **rows** (**horizontal records**) indicate the **actual values** and the **columns** (**vertical records**) indicate the **predicted values**.

	Defaulter	Non-defaulter
Defaulter	55	8
Non-defaulter	12	25

- $55 / (55+8)$
- $55 / (55 + 12)$
- $12 / (55 + 12)$
- $12 / (12 + 25)$

# Decision Tree

- A tree-like structure having decision nodes and leaf nodes.
  - Decision nodes - To make decisions based on certain conditions
  - Leaf nodes - contain the outcome(class labels)



- To make a prediction for an instance, the path from the root to the leaf is followed as per the conditions.

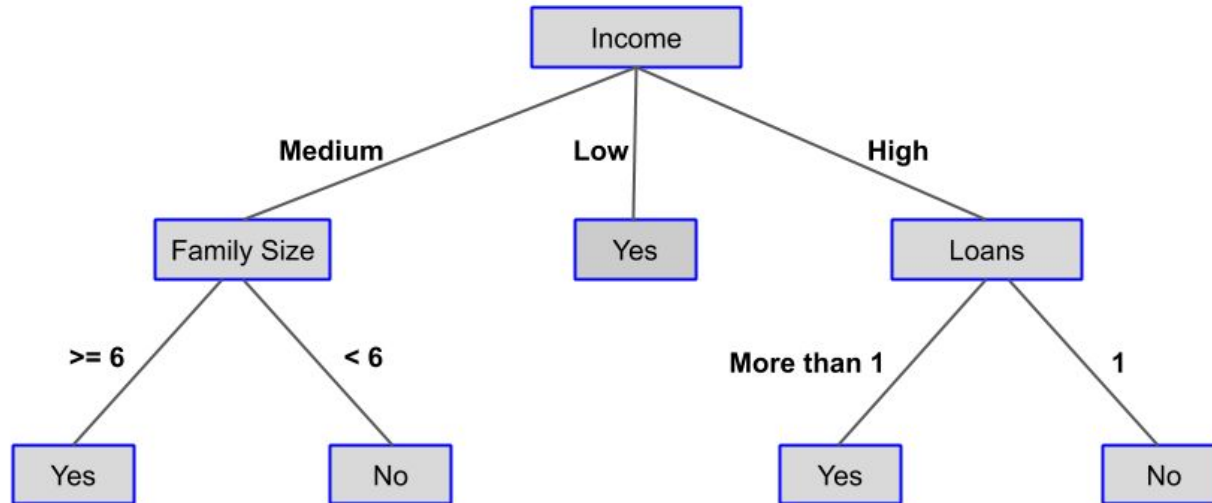
# Decision Tree [Contd.]

Consider a **case** where a **bank** has to **decide** whether a **customer** will **enroll** for **credit card** or **not** based on the **features** like **Income**, **Family Size** and **Loans**.

Income	Family Size	Loans	Credit Card (Yes/No)
Medium	$\geq 6$	1	Yes
Medium	$\geq 6$	More than 1	Yes
Low	$\geq 6$	1	Yes
High	$\geq 6$	1	No
High	$< 6$	More than 1	Yes
Low	$< 6$	More than 1	Yes
Medium	$< 6$	1	No
High	$< 6$	1	No
Medium	$< 6$	More than 1	No
Low	$\geq 6$	More than 1	Yes
Low	$< 6$	1	Yes
High	$\geq 6$	More than 1	Yes
Medium	$< 6$	More than 1	No
High	$\geq 6$	More than 1	Yes

## Decision Tree [Contd.]

So if we try to **predict** the **credit card enrollment** based on the **features** like **Income**, **Family Size** and **Loans** using **decision tree** then the **tree** may look like - .



## Decision Tree [Contd.]

The previous **decision tree** can be further **simplified** as -

```
if(Income == "Low"):
    credit_card = "Yes"
elif(Income == "Medium"):
    if(Family Size == ">=6"):
        Credit_card = "Yes"
    elif(Humidity == "<6"):
        Credit_card = "No"
elif(Income == "High"):
    if(Loans == "More than 1"):
        Credit_card = "Yes"
    elif(Loans == "1"):
        Credit_card = "No"
```

So the **decision tree** can also be **referred** as a **nested if-else classifier**.

## Quiz 5

Which of the following nodes in Decision Tree contains the conditions (decisions).

- Root Node
- Internal Nodes
- Leaf Nodes
- All the nodes

# How does the Decision Tree algorithm work?

**Step 1.** Select the **best attribute** using **Attribute Selection Measures(ASM)**, for e.g. **Information Gain** (mostly used).

**Step 2.** Make that **attribute** a **decision node** and **split** the **dataset** into **smaller subsets**.

**Step 3.** Start **tree** building by **repeating** this **process recursively** for each **child** until one of the **condition** is satisfied -

- There are **no** more **remaining attributes**.
- All the **tuples** belong to the **same attribute value**.
- There are **no** more **instances**.

# Attribute Selection Measures (ASM)

- To build a decision tree for a dataset having  $d$  features, the major challenge is to decide which feature to use at root node and at other decision nodes. To fix this, Attribute Selection Measures (ASM) are used.
- Attribute Selection Measures (ASM) compare different attributes/predictors and rank them for the purpose of tree/model building.
- The most commonly used attribute selection measures are -
  - Information Gain
  - Entropy
  - Gini Index



# Information Gain

- Information Gain measures how much information about the classes(categories) can be gained using a feature.
- The feature with highest information gain is taken at the root node (or decision node).

Information gain further depends on Entropy.

**Entropy** : It measures randomness in the data. Higher the entropy, more difficult to draw conclusions from the information.

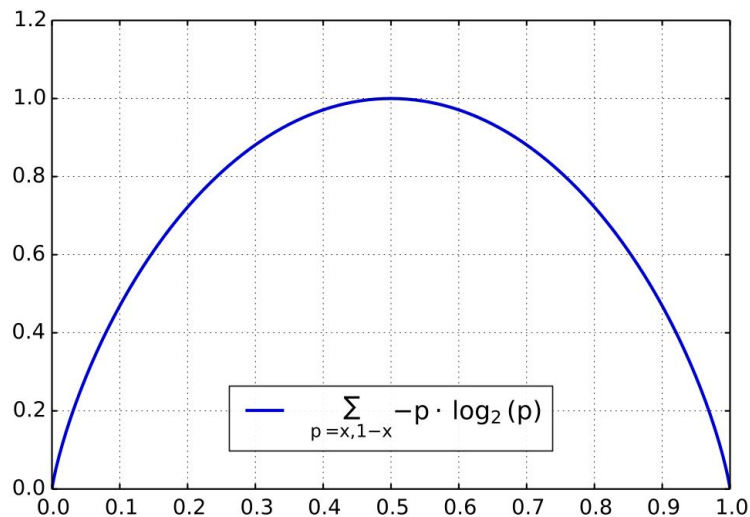
- Entropy for a dataset(S) with C classes is computed as -

$$Entropy(S) = - \sum_i^C p_i \log_2(p_i)$$

Where,  $p_i$  is the probability of class  $i$ .

# Properties of Entropy

- If **all** the **class labels** are **equally probable** (or **equally divided**) then **entropy** is **maximum**.  
E.g. In a **binary classification**, if (**y+** = 50% and **y-** = 50%) then  
 $H(y) = -0.5 \log(0.5) - 0.5 \log(0.5) = 1$
- If only **one class** is fully **dominating** (**most probable**) then **entropy** is **minimal**.  
E.g. In a **binary classification**, if (**y+** = 99% and **y-** = 1%) then  
 $H(y) = -0.99 \log(0.99) - 0.01 \log(0.01) = 0.08$



# Attribute Selection Measures [Contd.]

- **Information Gain** is computed as -

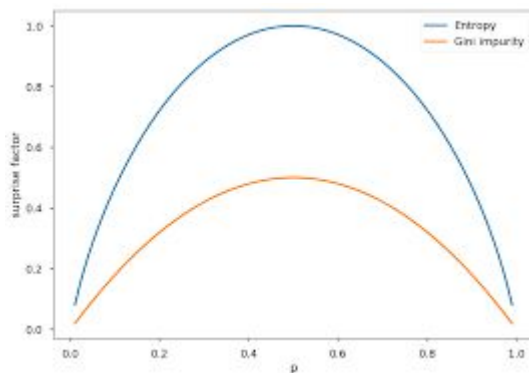
$$\text{Information Gain} = \text{Entropy}(S) - [\{\text{weighted average}\} * \text{Entropy}(\text{each feature})]$$

## Gini Index

- It **measures** how often a **randomly picked instance** would be **incorrectly classified**.
- It is a **measure** of **impurity** and an **attribute** with **low Gini index** is preferred **first**.

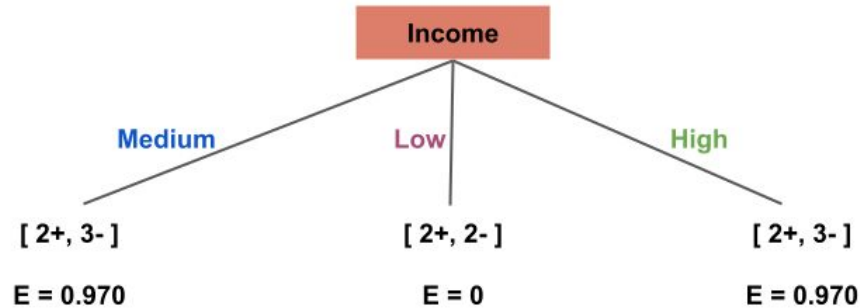
Mathematically,

$$\text{Gini Index} = 1 - \sum_i^C (p_i)^2$$



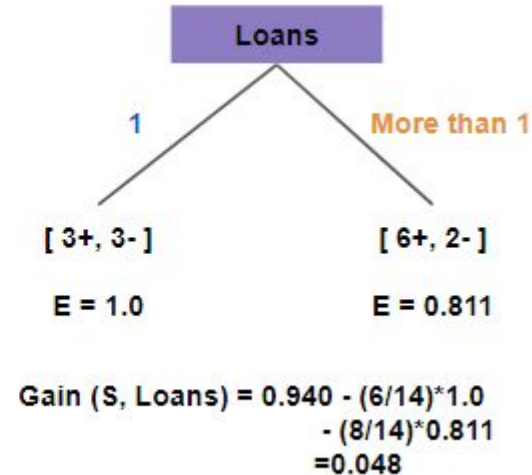
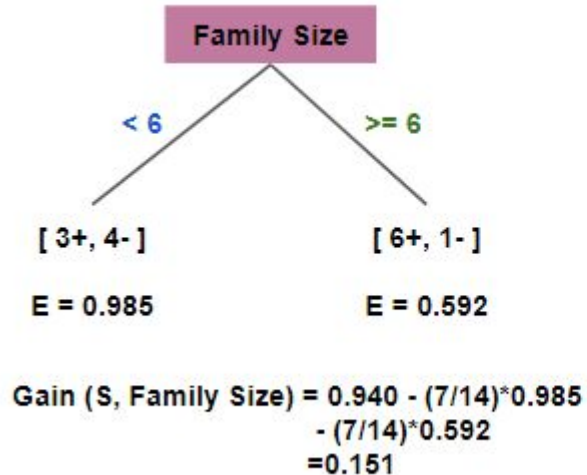
# Attribute at the Root Node

- Attribute with highest Information Gain is selected at the root node.



$$\text{Gain (S, Income)} = 0.940 - (5/14)*0.970 - (4/14)*0 - (5/14)*0.970 = 0.247$$

## Attribute at the Root Node [Contd.]



- So the attribute with highest information gain is Income.

## Quiz 6

Which of the following age group subset has entropy = 0?

- Above 60
- 30 to 60
- Below 30
- None of these

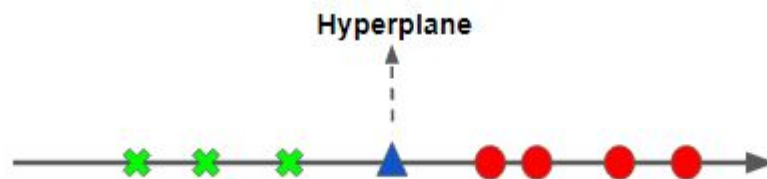
Age Group	Income	Deaulter (Yes/No)
Below 30	High	No
Below 30	High	No
30 to 60	High	Yes
Above 60	High	Yes
Above 60	Low	No
30 to 60	Low	Yes
Below 30	Low	Yes
Above 60	High	No
30 to 60	Low	Yes
Above 60	High	Yes
30 to 60	Low	Yes
Below 30	High	No

# Decision Tree Time & Space Complexity

- Training time complexity
  - for a dataset with  $n$  samples and  $d$  features =  $O(n (\log n) d)$ 
    - Where,  $n (\log n)$  corresponds to Sorting
    - $d$  - to evaluate  $d$  features every time
- Run (test) complexity
  - Time complexity :  $O(\text{depth}) = O(k)$ , where  $k$  is the maximum depth of the tree.
    - In order to predict the class for an instance we need to make  $k$  decisions.
  - Space complexity :  $O(\text{nodes}) = \# \text{ of internal nodes} + \# \text{ of leaf nodes}$
- So decision tree can handle large data with small(significant) number of features.

# Support Vector Machine (SVM)

- **Support Vector Machine (SVM)** tries to **find out** an **optimal hyperplane** (in a **d-dimensional** space) that can easily **separate** the **data points** into **classes**.
- **Hyperplane** is a **decision boundary** that **classifies** the **data points**. Its **dimension** varies as per the **number of features** varies. E.g. For a **dataset** having **two features**, the **hyperplane** will be a **line** whereas for a **dataset** having **3 features**, the **hyperplane** will be a **plane** in **3d space**.



1d Hyperplane (Point)

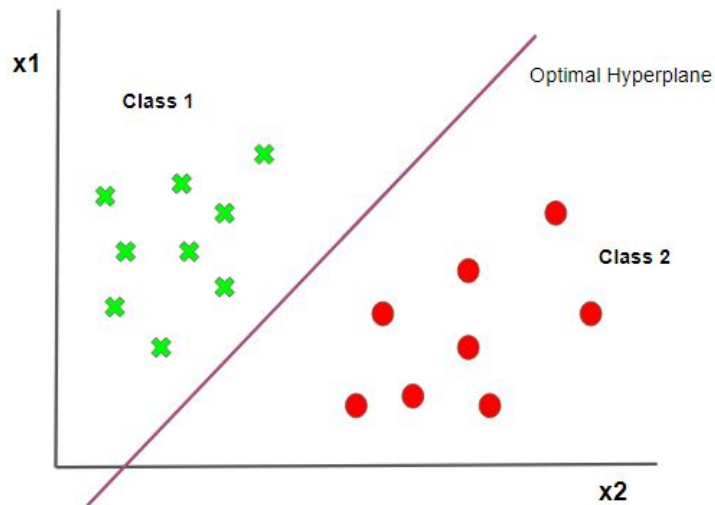
Class 1, for any datapoint to its left

Class 2, for any datapoint to its right



# Hyperplane [Contd.]

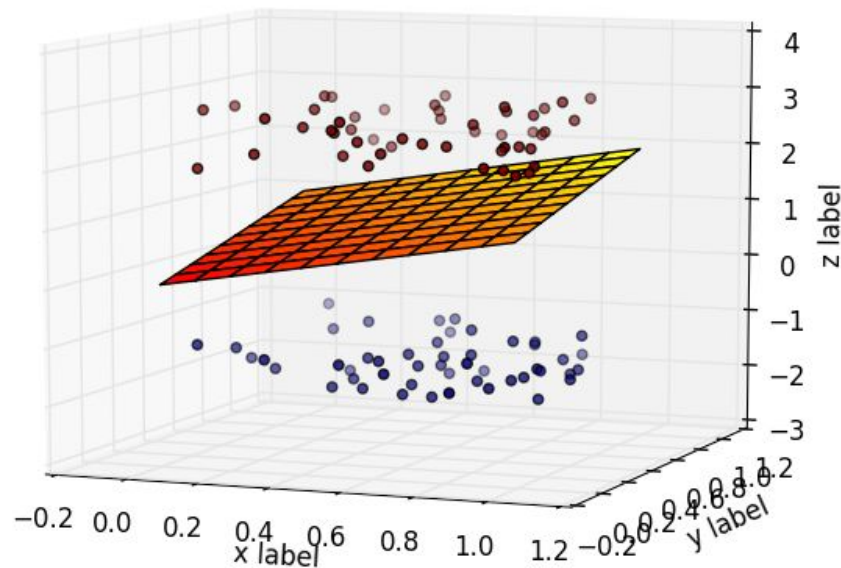
2d Hyperplane (Line)



Class 1, for any datapoint to its left

Class 2, for any datapoint to its right

3d Hyperplane (Plane)



Class 1, for any datapoint above the hyperplane.

Class 2, for any datapoint below the hyperplane.

# Quiz 9

Choose the correct statement(s) about Figure 1 and 2.



Figure 1

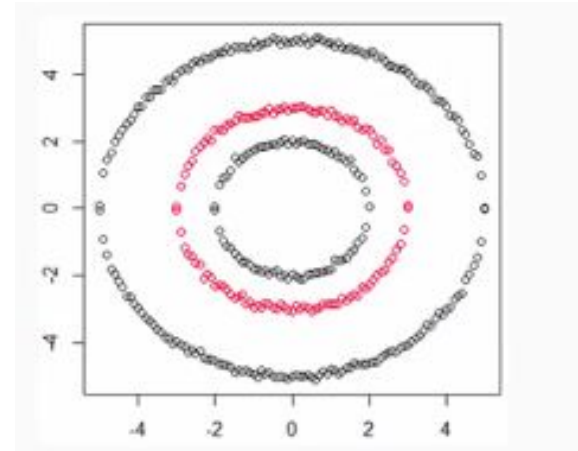
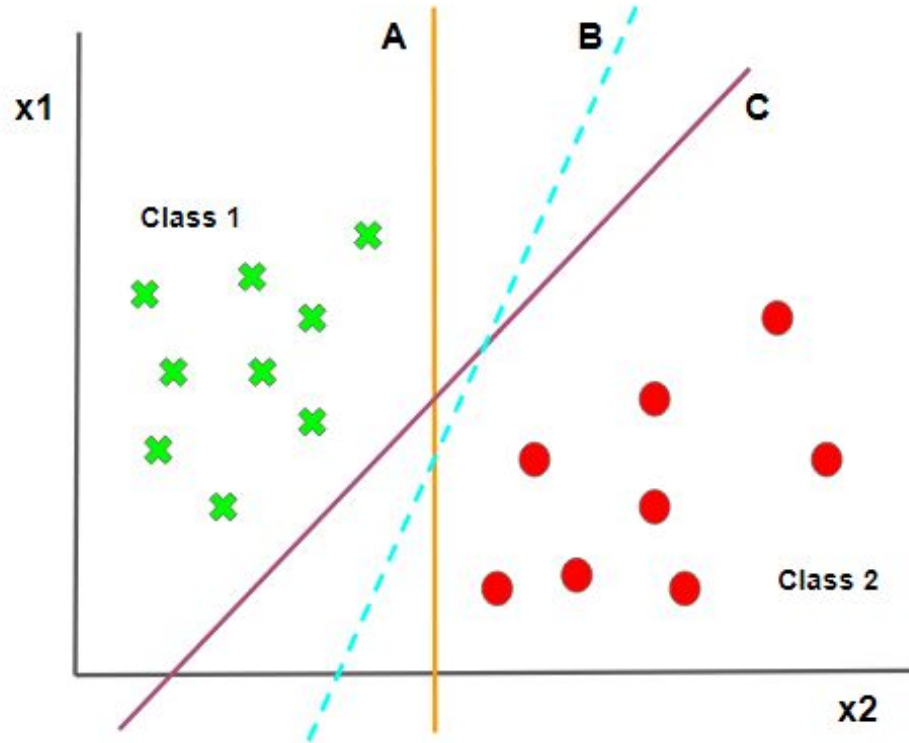


Figure 2

- Figure 1 and 2 both are linearly separable.
- Figure 1 is linearly separable but 2 is not.
- Figure 2 is linearly separable but 1 is not.
- Figure 1 and 2 both are non-linearly separable.

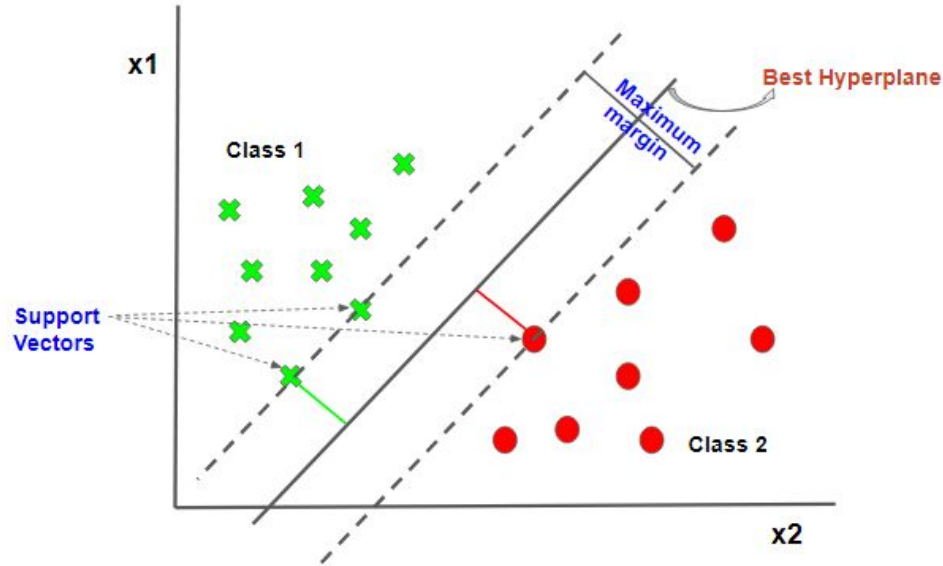
# How to find an optimal Hyperplane

Multiple hyperplanes like A, B and C are possible then how to find the optimal hyperplane?



# How to find an optimal Hyperplane

In finding the **best hyperplane**, the **distance** between the **closest point** of **each class** and the **hyperplane** plays an **important** role. This **distance** is known as **margin** and **SVM** tries to find a **hyperplane** with the **maximum margin**.

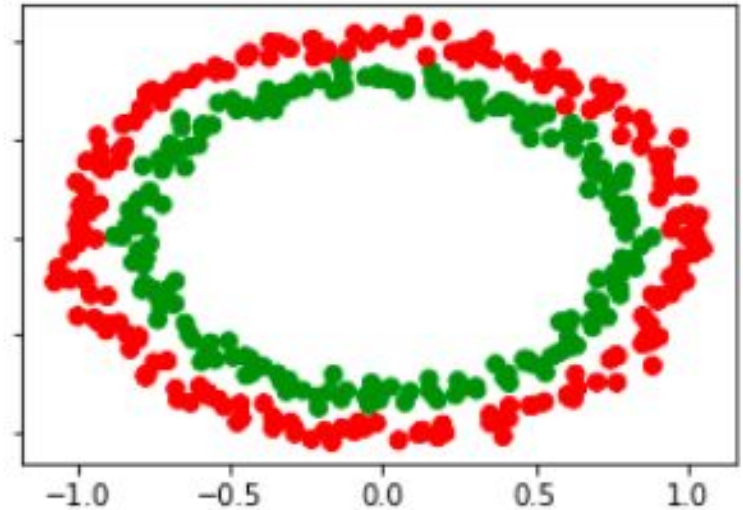
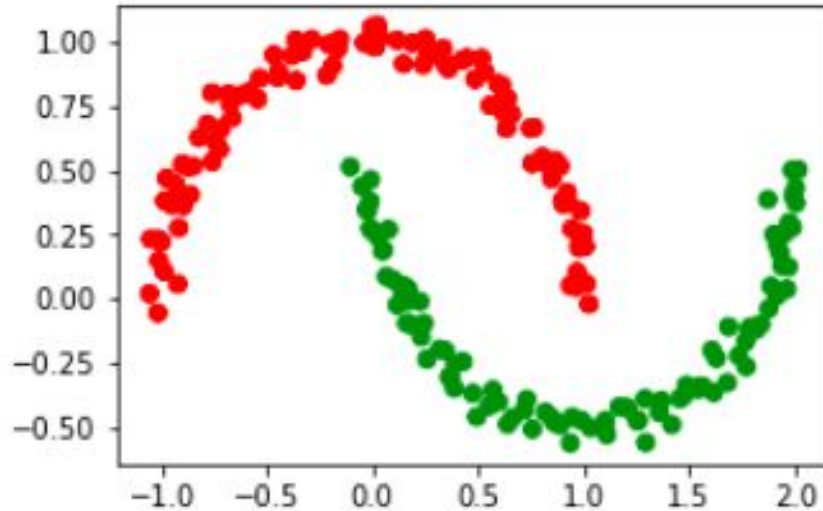


**Support vectors** - The **data points**(**vectors**) that are **closest** to the **hyperplane** and **dominate** in deciding the **position** and **orientation** of the **hyperplane**.

# Non-linearly Separable data

- The datasets which can not be easily separated into classes using a single line.

Non-linearly seperable data



# Types of SVM

SVM can be further classified in two categories.

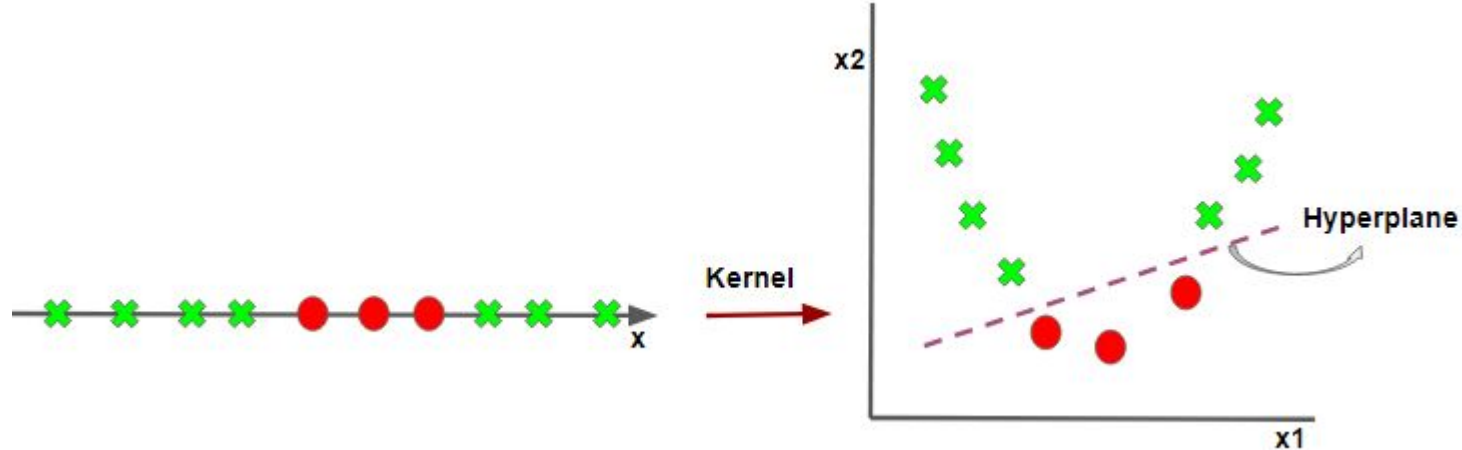
- **Linear SVM** - Best suited to linearly separable data. If a dataset can be easily separated in two classes using a single line then it is known as linearly separable data.
- **Non-Linear SVM** - As the name suggests, it is best suited to non-linearly separable data. In this, the dataset is mapped to a higher dimensional space using kernel trick to get an optimal hyperplane.

Most commonly used kernels are polynomial and radial.

# Kernel Trick

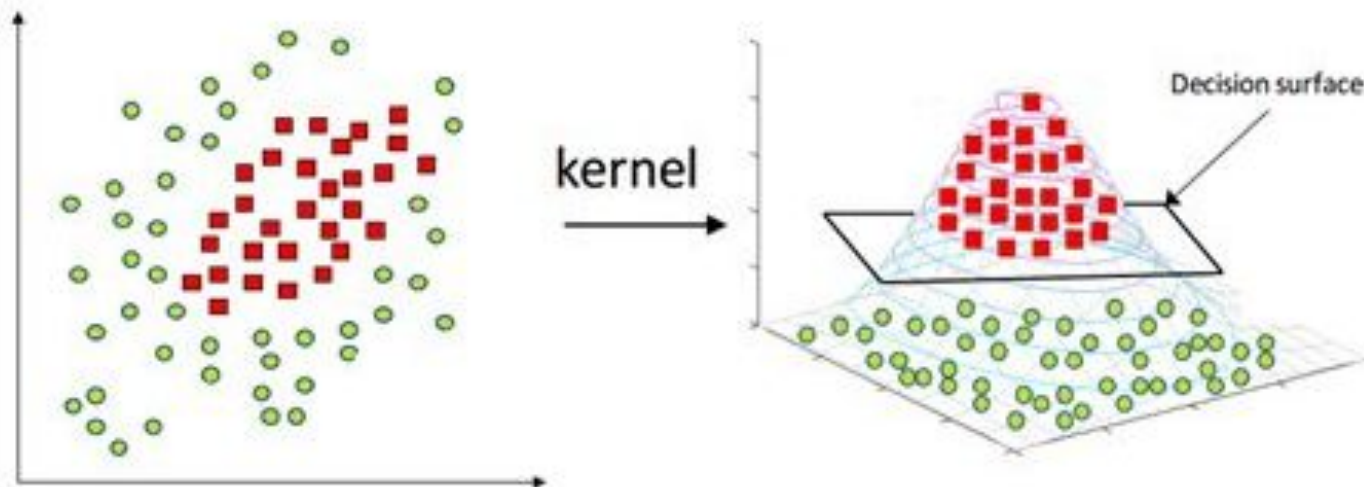
- The **non-linearly separable dataset** often becomes **separable** after **mapping** the **original feature space** to a **higher dimensional space**. This entire **mapping process** is known as **kernelization**.
- So **Kernel methods** are used to **transform** the **data** into a **higher dimensional space** such that the **data** becomes **separable** in the **newer space**.

## Example 1 -



# Kernel Trick [Contd.]

Example 2 -





# Types of Kernel

The most commonly used kernels in SVM are -

1. Linear Kernel
2. Polynomial Kernel
3. Radial Basis Kernel (RBF)

## Polynomial Kernel

In polynomial kernel, the dot product by increasing the power of the kernel is computed.

$$K(X_1, X_2) = (1 + X_1^T X_2)^b$$

Here, **b** = degree of polynomial kernel

## Polynomial Kernel [Contd.]

Let's say there are two points  $x_a$ ,  $x_b$  in a original feature space (2-dimensional) as

$$x_a = (a_1, a_2) \quad x_b = (b_1, b_2)$$

Now if we map them to a higher dimensional space using a polynomial kernel of degree 2

$$\begin{aligned} K(x_a, x_b) &= (1 + x_a^T x_b)^2 \\ &= (1 + a_1 b_1 + a_2 b_2)^2 \\ &= (1 + a_1^2 b_1^2 + a_2^2 b_2^2 + 2 a_1 b_1 + 2 a_2 b_2 + 2 a_1 b_1 a_2 b_2) \\ &\Rightarrow [1, a_1^2, a_2^2, a_1, a_2, a_1 * a_2] \rightarrow x'_a \\ &\quad [1, b_1^2, b_2^2, b_1, b_2, b_1 * b_2] \rightarrow x'_b \\ &= (x'_a)^T (x'_b) \end{aligned}$$

This new feature space is 6 dimensional.

# RBF (Radial Basis Function) Kernel

- This is the most popular and general purpose kernel.

Mathematically,

$$K_{RBF}(X_1, X_2) = e^{\left(\frac{-\|X_1 - X_2\|^2}{2\sigma^2}\right)}$$

Where  $\|X_1 - X_2\|$  is the Euclidean distance between  $X_1$  and  $X_2$ .

- This kernel is also known as Gaussian Kernel.

# Quiz 10

Choose the correct statement(s).

- In Kernelization process, the data is mapped to a lower dimensional space.
- SVM without any kernel methods is analogous to linear SVM.
- As the Euclidean distance increases in RBF kernel, its value decreases.
- All of the above