



# Random Forest Algorithm



# Random Forest

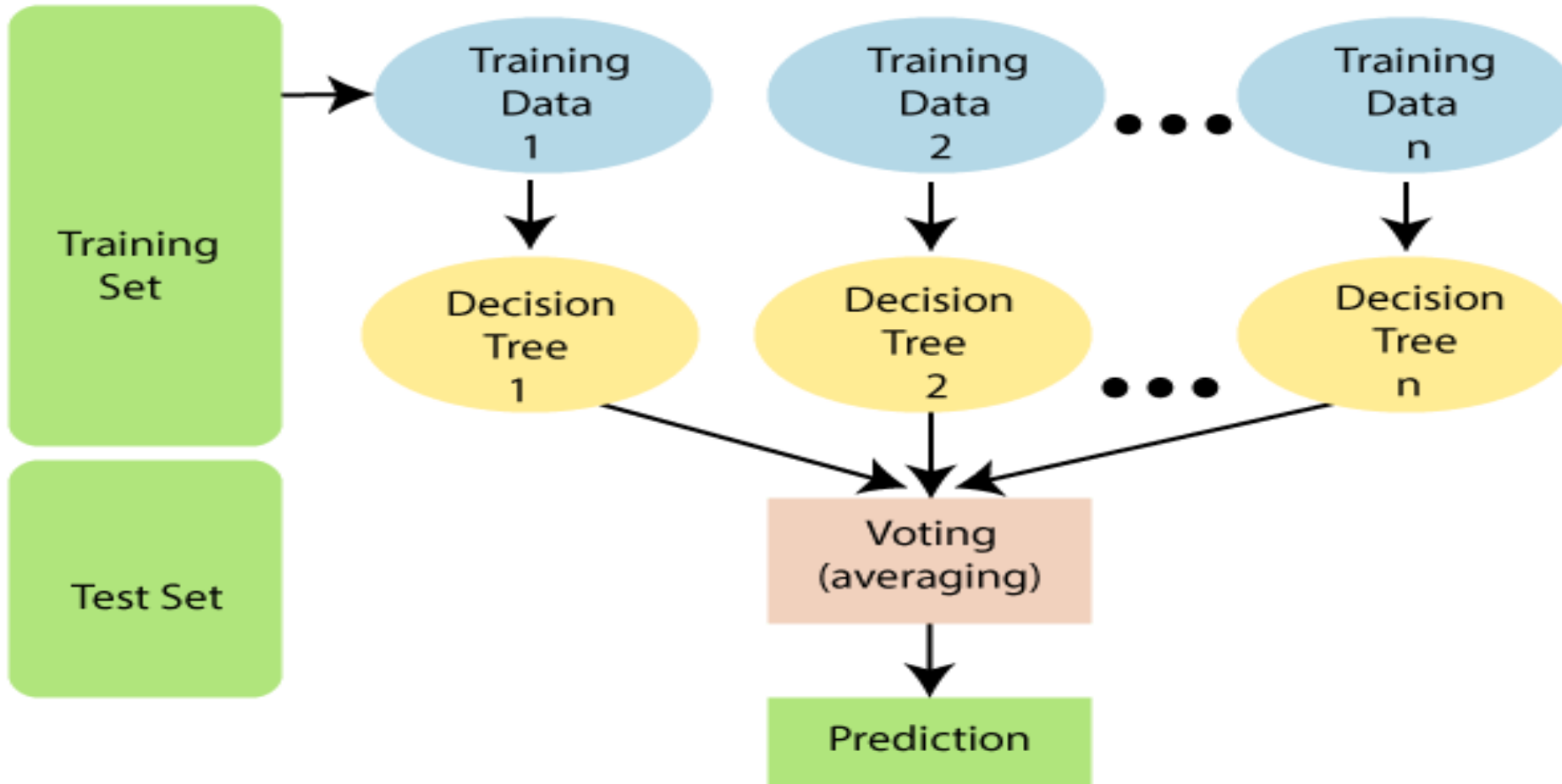
- Random Forest is a popular machine learning algorithm that belongs to the **supervised learning technique**.
- It can be used for both **Classification and Regression problems** in ML.
- It is based on the concept of **ensemble learning**, which is a process of *combining multiple classifiers to solve a complex problem and to improve the performance of the model*.



# What Is Random Forest?

- ***"Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset."***
- Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.
- **The greater number of trees in the forest leads to higher accuracy and prevents the problem of overfitting.**
- The below diagram explains the working of the Random Forest algorithm:

# Random Forest





# Why use Random Forest?

- It takes less training time as compared to other algorithms.
- It predicts output with high accuracy, even for the large dataset it runs efficiently.
- It can also maintain accuracy when a large proportion of data is missing.



# How does Random Forest algorithm work?

- Random Forest works in two-phase first is to create the random forest by combining  $N$  decision tree, and second is to make predictions for each tree created in the first phase.

The Working process can be explained in the below steps and diagram:

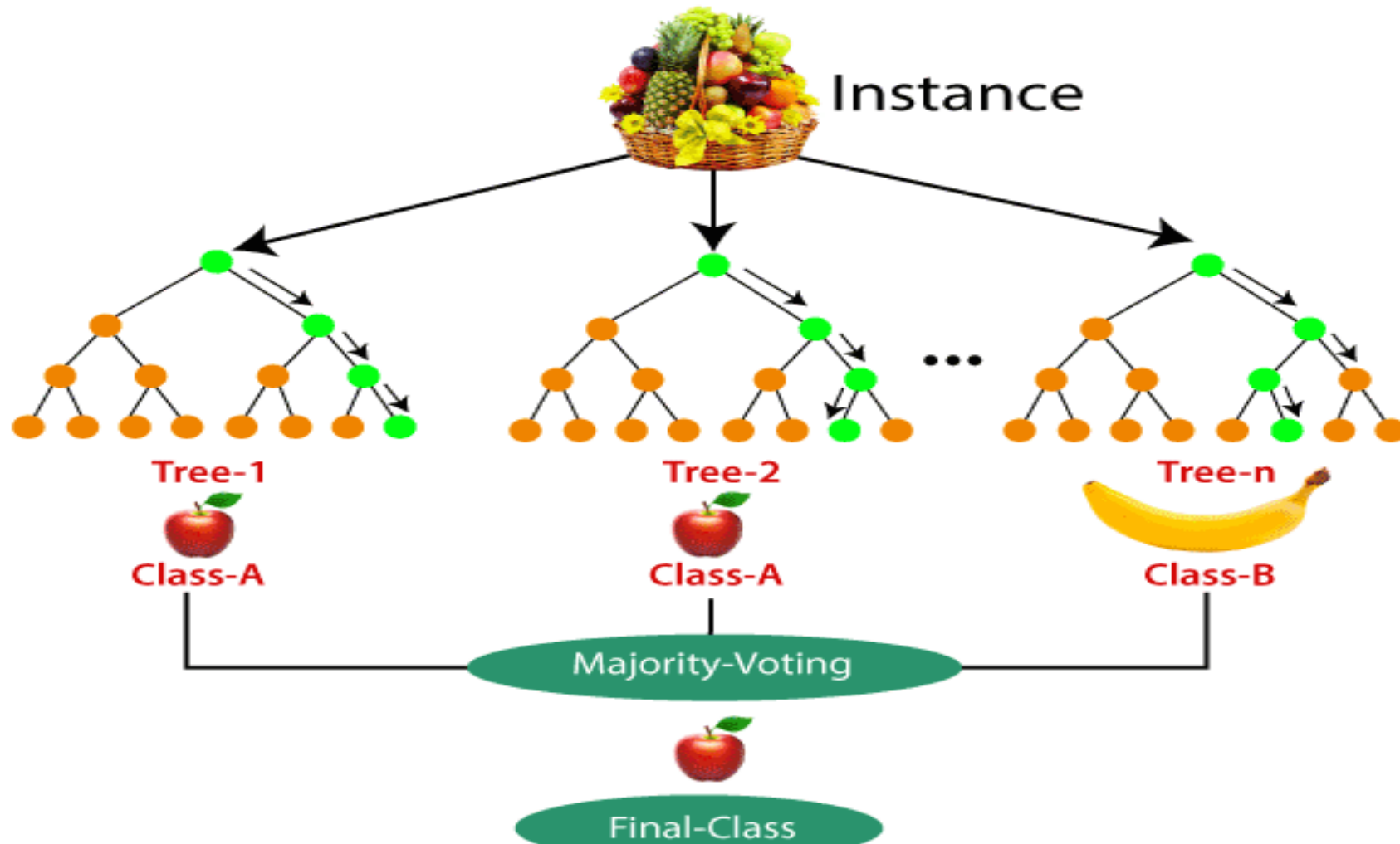
- **Step-1:** Select random  $K$  data points from the training set.
- **Step-2:** Build the decision trees associated with the selected data points (Subsets).
- **Step-3:** Choose the number  $N$  for decision trees that you want to build.
- **Step-4:** Repeat Step 1 & 2.
- **Step-5:** For new data points, find the predictions of each decision tree, and assign the new data points to the category that wins the majority votes.



# Example:

- Suppose there is a dataset that contains multiple fruit images.
- So, this dataset is given to the Random forest classifier.
- The dataset is divided into subsets and given to each decision tree.
- During the training phase, each decision tree produces a prediction result, and when a new data point occurs, then based on the majority of results, the Random Forest classifier predicts the final decision.
- Consider the below image:

# Example:RFA







# Applications of Random Forest

There are mainly four sectors where Random forest mostly used:

- **Banking:** Banking sector mostly uses this algorithm for the identification of loan risk.
- **Medicine:** With the help of this algorithm, disease trends and risks of the disease can be identified.
- **Land Use:** We can identify the areas of similar land use by this algorithm.
- **Marketing:** Marketing trends can be identified using this algorithm.



# Advantages & Disadvantages

## Advantages of Random Forest

- Random Forest is capable of performing both Classification and Regression tasks.
- It is capable of handling large datasets with high dimensionality.
- It enhances the accuracy of the model and prevents the overfitting issue.

## Disadvantages of Random Forest

- Although random forest can be used for both classification and regression tasks, it is not more suitable for Regression tasks.



# K-Nearest Neighbour(KNN) Algorithm

- K-Nearest Neighbor is one of the simplest Machine Learning algorithms based on Supervised Learning technique.
- K-NN algorithm stores all the available data and classifies a new data point based on the similarity.
- K-NN algorithm can be used for **Regression** as well as for **Classification** but mostly it is used for the Classification problems.
- K-NN is a **non-parametric algorithm**, which means it does not make any assumption on underlying data.
- It is also called a **lazy learner algorithm** because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset.
- KNN algorithm at the **training phase just stores the dataset** and when it gets new data, then it classifies that data into a category that is much similar to the new data.

- **Example:** Suppose, we have an image of a creature that looks similar to cat and dog, but we want to know either it is a cat or dog.
- So for this identification, we can use the KNN algorithm, as it works on a **similarity measure**.
- Our **KNN model will find the similar features of the new data set** to the cats and dogs images and based on the most similar features it will put it in either cat or dog category.

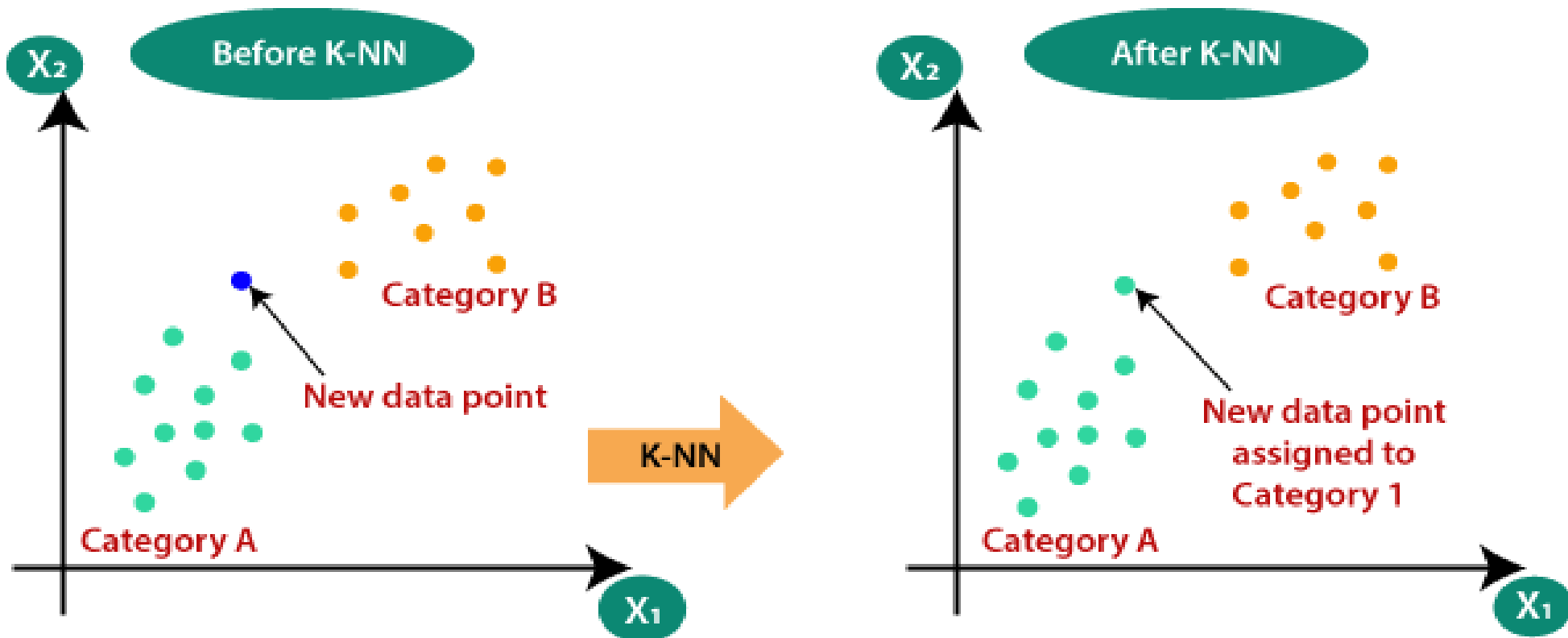




# Why do we need a K-NN Algorithm?

- Suppose there are two categories, i.e., Category A and Category B, and we have a new data point  $x_1$ , so this data point will lie in which of these categories.
- To solve this type of problem, we need a **K-NN algorithm**. With the help of K-NN, we can easily identify the category or class of a particular dataset. Consider the below diagram:

Example:





# How does K-NN work?

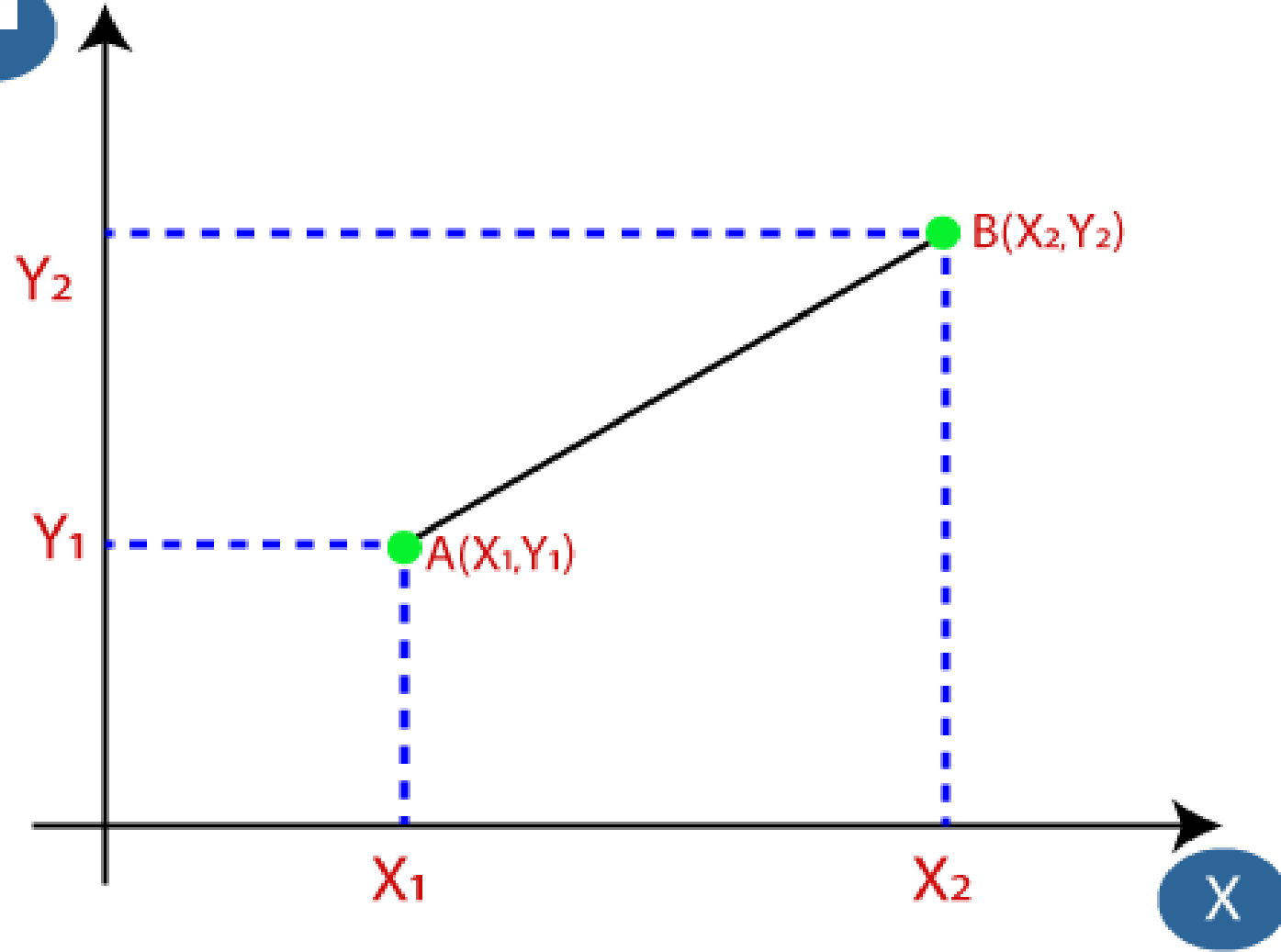
- **Step-1:** Select the number  $K$  of the neighbors
- **Step-2:** Calculate the **Euclidean distance** of  **$K$  number of neighbors**
- **Step-3:** Take the  $K$  nearest neighbors as per the calculated Euclidean distance.
- **Step-4:** Among these  $k$  neighbors, count the number of the data points in each category.
- **Step-5:** Assign the new data points to that category for which the number of the neighbor is maximum.
- **Step-6:** Our model is ready.



# How to calculate the Euclidean distance?

- Firstly, we will choose the number of neighbors, so we will choose the **k=5**.
- Next, we will calculate the **Euclidean distance** between the data points.
- The Euclidean distance is **the distance between two points**, which we have already studied in geometry. It can be calculated as:



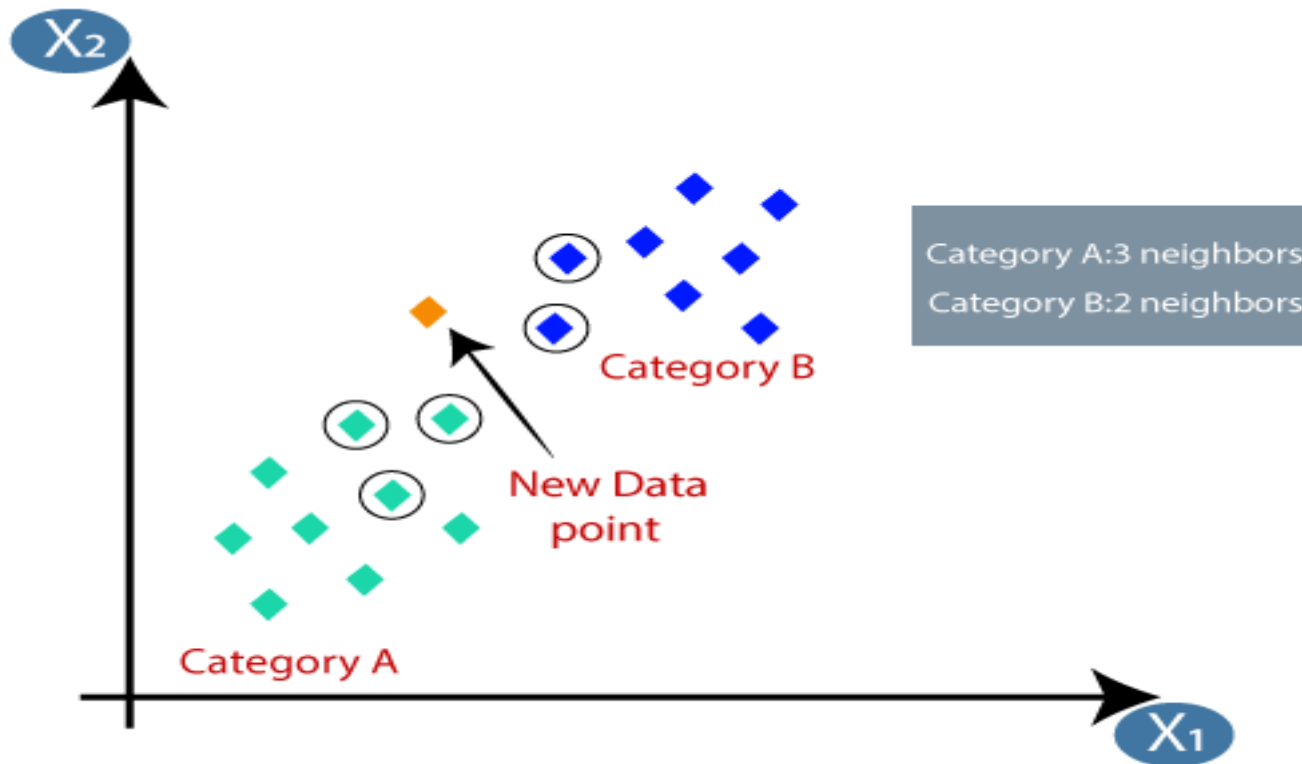


$$\text{Euclidean Distance between } A_1 \text{ and } B_2 = \sqrt{(X_2 - X_1)^2 + (Y_2 - Y_1)^2}$$



calculating the Euclidean distance we got the **nearest neighbor**, **two nearest neighbors in category A** and **two nearest neighbors in category B**. Consider the below image:

- As we can see the 3 nearest neighbors are from category A, hence this new data point must belong to category A.





# How to select the value of K in the K-NN algorithm?

- There is no particular way to determine the best value for "K", so we need to try some values to find the best out of them.
- The most preferred value for K is 5.
- A very low value for K such as  $K=1$  or  $K=2$ , **can be noisy and lead to the effects of outliers in the model.**
- Large values for K are good, but it may find some difficulties.



# Advantages & Disadvantages of KNN Algorithm



## Advantages of KNN Algorithm:

- It is simple to implement.
- It is robust to the noisy training data
- It can be more effective if the training data is large.

## Disadvantages of KNN Algorithm:

- Always needs to determine the value of K which may be complex some time.
- The **computation cost is high** because of calculating the distance between the data points for all the training samples.