Supervised Learning Algorithm

Supervised learning is a type of Machine learning in which the machine needs external supervision to learn. The supervised learning models are trained using the labeled dataset. Once the training and processing are done, the model is tested by providing a sample test data to check whether it predicts the correct output.

The goal of supervised learning is to map input data with the output data. Supervised learning is based on supervision, and it is the same as when a student learns things in the teacher's supervision. The example of supervised learning is **spam filtering**.

Supervised learning can be divided further into two categories of problem:
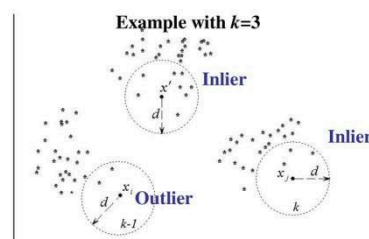
  o   Classification
  o   Regression

**Distance-basedmodels**

**Like Linear models, distance-based models are** based on the geometry of data**. As the name implies, distance-based models work on the concept of distance. In the context of Machine learning, the concept of distance is not based on merely the physical distance between two points.**



**Distance-based method**
[Knorr and Ng , CASCR 1997]

**Definition**: Data point *x* is an outlier if at most *k* points are within the distance *d* from *x*.

K-Nearest Neighbor(KNN) Algorithm for Machine Learning
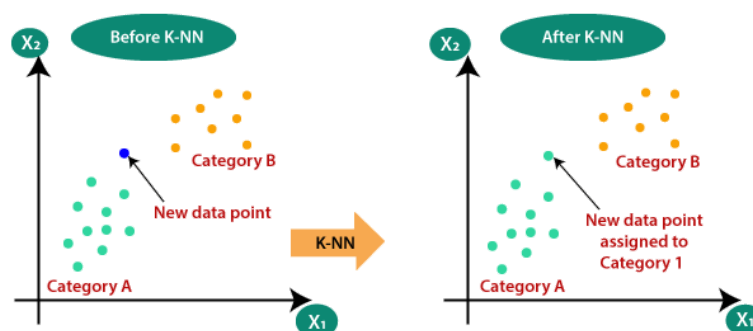
  o   K-Nearest Neighbour is one of the simplest Machine Learning algorithms based on Supervised Learning technique.
  o   K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories.
  o   K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well suite category by using K- NN algorithm.
  o   K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for the Classification problems.

- o K-NN is a **non-parametric algorithm**, which means it does not make any assumption on underlying data.
- o 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.
- o 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.
- o **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 x1, 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:
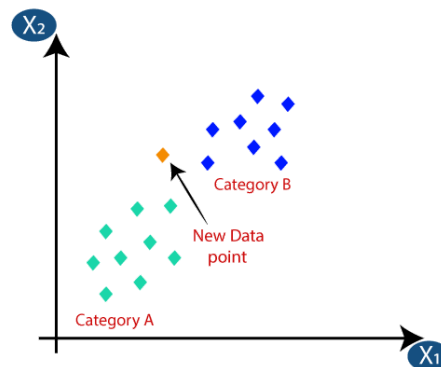


## How does K-NN work?

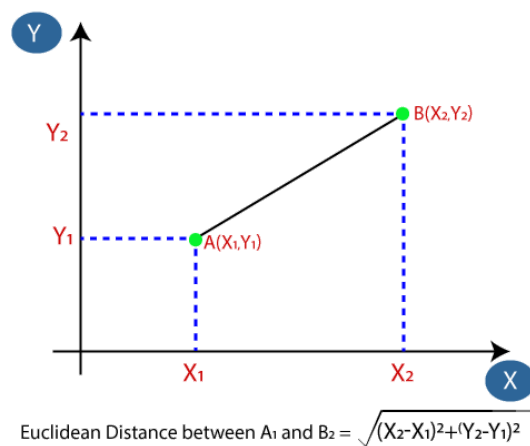The K-NN working can be explained on the basis of the below algorithm:

- o **Step-1:** Select the number K of the neighbors
- o **Step-2:** Calculate the Euclidean distance of **K number of neighbors**
- o **Step-3:** Take the K nearest neighbors as per the calculated Euclidean distance.

o **Step-4:** Among these k neighbors, count the number of the data points in each category.

o **Step-5:** Assign the new data points to that category for which the number of the neighbor is maximum.

o **Step-6:** Our model is ready.

Suppose we have a new data point and we need to put it in the required category. Consider the below image:



o Firstly, we will choose the number of neighbors, so we will choose the k=5.

o 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:



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

o By calculating the Euclidean distance we got the nearest neighbors, as three nearest neighbors in category A and two nearest neighbors in category B. Consider the below image:

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

<u>How to select the value of K in the K-NN Algorithm?</u>

Below are some points to remember while selecting the value of K in the K-NN algorithm:

- o 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.
- o 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.
- o Large values for K are good, but it may find some difficulties.

<u>Advantages of KNN Algorithm:</u>

- o It is simple to implement.
- o It is robust to the noisy training data
- o It can be more effective if the training data is large.
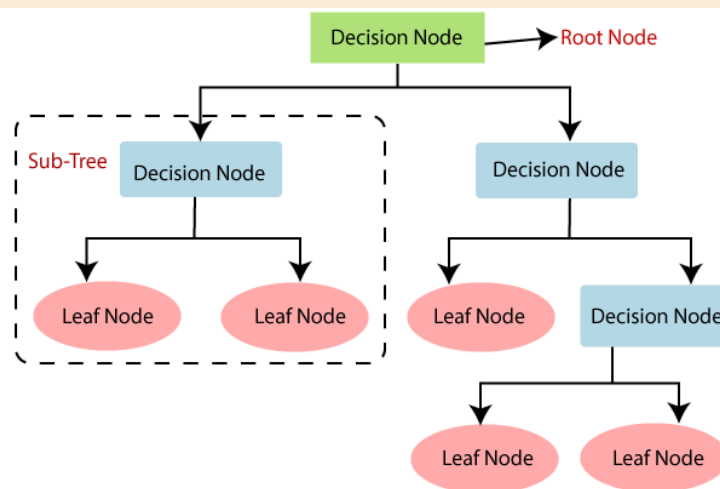
<u>Disadvantages of KNN Algorithm:</u>

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

<u>Decision Tree Classification Algorithm</u>

- o Decision Tree is a **Supervised learning technique** that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where **internal nodes represent the features of a dataset, branches represent the decision rules** and **each leaf node represents the outcome.**
- o In a Decision tree, there are two nodes, which are the **Decision Node** and **Leaf Node.** Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.

o   The decisions or the test are performed on the basis of features of the given dataset.

o   ***It is a graphical representation for getting all the possible solutions to a problem/decision based on given conditions.***

o   It is called a decision tree because, similar to a tree, it starts with the root node, which expands on further branches and constructs a tree-like structure.

o   In order to build a tree, we use the **CART algorithm,** which stands for **Classification and Regression Tree algorithm.**

o   A decision tree simply asks a question, and based on the answer (Yes/No), it further split the tree into subtrees.

o   Below diagram explains the general structure of a decision tree:

*Note: A decision tree can contain categorical data (YES/NO) as well as numeric data.*



## Why use Decision Trees?

There are various algorithms in Machine learning, so choosing the best algorithm for the given dataset and problem is the main point to remember while creating a machine learning model. Below are the two reasons for using the Decision tree:

o   Decision Trees usually mimic human thinking ability while making a decision, so it is easy to understand.

o   The logic behind the decision tree can be easily understood because it shows a tree-like structure.

## Decision Tree Terminologies

◈ **Root Node:** Root node is from where the decision tree starts. It represents the entire dataset, which further gets divided into two or more homogeneous sets.

◈ **Leaf Node:** Leaf nodes are the final output node, and the tree cannot be segregated further after getting a leaf node.

◈ **Splitting:** Splitting is the process of dividing the decision node/root node into sub-nodes according to the given conditions.

- **Branch/Sub Tree:** A tree formed by splitting the tree.

- **Pruning:** Pruning is the process of removing the unwanted branches from the tree.

- **Parent/Child node:** The root node of the tree is called the parent node, and other nodes are called the child nodes.
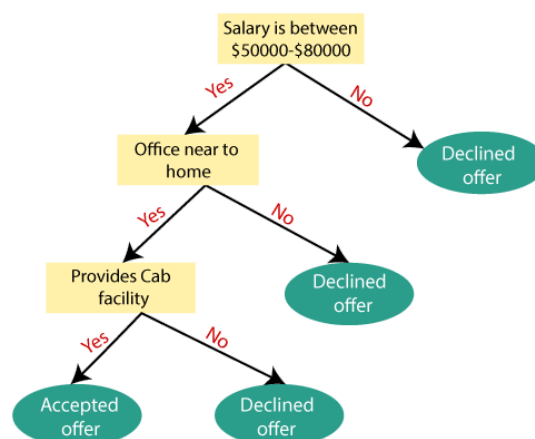
**How does the Decision Tree algorithm Work?**

In a decision tree, for predicting the class of the given dataset, the algorithm starts from the root node of the tree. This algorithm compares the values of root attribute with the record (real dataset) attribute and, based on the comparison, follows the branch and jumps to the next node.

For the next node, the algorithm again compares the attribute value with the other sub-nodes and move further. It continues the process until it reaches the leaf node of the tree. The complete process can be better understood using the below algorithm:

- o **Step-1:** Begin the tree with the root node, says S, which contains the complete dataset.
- o **Step-2:** Find the best attribute in the dataset using **Attribute Selection Measure (ASM).**
- o **Step-3:** Divide the S into subsets that contains possible values for the best attributes.
- o **Step-4:** Generate the decision tree node, which contains the best attribute.
- o **Step-5:** Recursively make new decision trees using the subsets of the dataset created in step -3. Continue this process until a stage is reached where you cannot further classify the nodes and called the final node as a leaf node.

**Example:** Suppose there is a candidate who has a job offer and wants to decide whether he should accept the offer or Not. So, to solve this problem, the decision tree starts with the root node (Salary attribute by ASM). The root node splits further into the next decision node (distance from the office) and one leaf node based on the corresponding labels. The next decision node further gets split into one decision node (Cab facility) and one leaf node. Finally, the decision node splits into two leaf nodes (Accepted offers and Declined offer). Consider the below diagram:

Attribute Selection Measures

While implementing a Decision tree, the main issue arises that how to select the best attribute for the root node and for sub-nodes. So, to solve such problems there is a technique which is called as **Attribute selection measure or ASM.** By this measurement, we can easily select the best attribute for the nodes of the tree. There are two popular techniques for ASM, which are:

- o **Information Gain**
- o **Gini Index**

## 1. Information Gain:

- o Information gain is the measurement of changes in entropy after the segmentation of a dataset based on an attribute.
- o It calculates how much information a feature provides us about a class.
- o According to the value of information gain, we split the node and build the decision tree.
- o A decision tree algorithm always tries to maximize the value of information gain, and a node/attribute having the highest information gain is split first. It can be calculated using the below formula:

1. Information Gain= Entropy(S)- [(Weighted Avg) *Entropy(each feature)
   **Entropy:** Entropy is a metric to measure the impurity in a given attribute. It specifies randomness in data. Entropy can be calculated as:

   Entropy(s)= -P(yes)log2 P(yes)- P(no) log2 P(no)
   **Where,**

   - o **S= Total number of samples**
   - o **P(yes)= probability of yes**
   - o **P(no)= probability of no**

## 2. Gini Index:

- o Gini index is a measure of impurity or purity used while creating a decision tree in the CART(Classification and Regression Tree) algorithm.
- o An attribute with the low Gini index should be preferred as compared to the high Gini index.
- o It only creates binary splits, and the CART algorithm uses the Gini index to create binary splits.
- o Gini index can be calculated using the below formula:

Gini Index= 1- $\sum_j P_j^2$

Pruning: Getting an Optimal Decision tree

*Pruning is a process of deleting the unnecessary nodes from a tree in order to get the optimal decision tree.*

A too-large tree increases the risk of overfitting, and a small tree may not capture all the important features of the dataset. Therefore, a technique that decreases the size of the learning tree without reducing accuracy is known as Pruning. There are mainly two types of tree **pruning** technology used:

- o **Cost Complexity Pruning**
- o **Reduced Error Pruning.**

## Advantages of the Decision Tree

- o It is simple to understand as it follows the same process which a human follow while making any decision in real-life.
- o It can be very useful for solving decision-related problems.
- o It helps to think about all the possible outcomes for a problem.
- o There is less requirement of data cleaning compared to other algorithms.

## Disadvantages of the Decision Tree

- o The decision tree contains lots of layers, which makes it complex.
- o It may have an overfitting issue, which can be resolved using the **Random Forest algorithm.**
- o For more class labels, the computational complexity of the decision tree may increase.

## Naïve Bayes Classifier Algorithm

- o Naïve Bayes algorithm is a supervised learning algorithm, which is based on **Bayes theorem** and used for solving classification problems.
- o It is mainly used in *text classification* that includes a high-dimensional training dataset.
- o Naïve Bayes Classifier is one of the simple and most effective Classification algorithms which helps in building the fast machine learning models that can make quick predictions.
- o **It is a probabilistic classifier, which means it predicts on the basis of the probability of an object**.
- o Some popular examples of Naïve Bayes Algorithm are **spam filtration, Sentimental analysis, and classifying articles**.

## Why is it called Naïve Bayes?

The Naïve Bayes algorithm is comprised of two words Naïve and Bayes, Which can be described as:

- o **Naïve**: It is called Naïve because it assumes that the occurrence of a certain feature is independent of the occurrence of other features. Such as if the fruit is identified on the

bases of color, shape, and taste, then red, spherical, and sweet fruit is recognized as an apple. Hence each feature individually contributes to identify that it is an apple without depending on each other.

o **Bayes**: It is called Bayes because it depends on the principle of Bayes' Theorem
.

Bayes' Theorem:

o Bayes' theorem is also known as **Bayes' Rule** or **Bayes' law**, which is used to determine the probability of a hypothesis with prior knowledge. It depends on the conditional probability.

o The formula for Bayes' theorem is given as:

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

| | Outlook | Play |
|---|---|---|
| 0 | Rainy | Yes |
| 1 | Sunny | Yes |
| 2 | Overcast | Yes |
| 3 | Overcast | Yes |
| 4 | Sunny | No |
| 5 | Rainy | Yes |
| 6 | Sunny | Yes |
| 7 | Overcast | Yes |
| 8 | Rainy | No |
| 9 | Sunny | No |

| 10 | Sunny | Yes |
|----|-------|-----|
| 11 | Rainy | No |

**Where,**

**P(A|B) is Posterior probability**: Probability of hypothesis A on the observed event B.

**P(B|A) is Likelihood probability**: Probability of the evidence given that the probability of a hypothesis is true.

**P(A) is Prior Probability**: Probability of hypothesis before observing the evidence.

**P(B) is Marginal Probability**: Probability of Evidence.

Working of Naïve Bayes' Classifier:

Working of Naïve Bayes' Classifier can be understood with the help of the below example:

Suppose we have a dataset of **weather conditions** and corresponding target variable **"Play"**. So using this dataset we need to decide that whether we should play or not on a particular day according to the weather conditions. So to solve this problem, we need to follow the below steps:

1. Convert the given dataset into frequency tables.
2. Generate Likelihood table by finding the probabilities of given features.
3. Now, use Bayes theorem to calculate the posterior probability.

**Problem**: If the weather is sunny, then the Player should play or not?

**Solution**: To solve this, first consider the below dataset:

**Frequency table for the Weather Conditions:**

| Weather | Yes | No |
|---------|-----|-----|
| Overcast | 5 | 0 |
| Rainy | 2 | 2 |
| Sunny | 3 | 2 |
| Total | 10 | 5 |

**Likelihood table weather condition:**

| Weather | No | Yes | |
|---------|-----|------|--------|
| Overcast | 0 | 5 | 5/14= 0.35 |
| Rainy | 2 | 2 | 4/14=0.29 |
| Sunny | 2 | 3 | 5/14=0.35 |
| All | 4/14=0.29 | 10/14=0.71 | |

**Applying Bayes'theorem:**

**P(Yes|Sunny)= P(Sunny|Yes)*P(Yes)/P(Sunny)**

P(Sunny|Yes)= 3/10= 0.3

P(Sunny)= 0.35

P(Yes)=0.71

So P(Yes|Sunny) = 0.3*0.71/0.35= **0.60**

**P(No|Sunny)= P(Sunny|No)*P(No)/P(Sunny)**

P(Sunny|NO)= 2/4=0.5

P(No)= 0.29

P(Sunny)= 0.35

So P(No|Sunny)= 0.5*0.29/0.35 = **0.41**

So as we can see from the above calculation that **P(Yes|Sunny)>P(No|Sunny)**

**Hence on a Sunny day, Player can play the game.**

Advantages of Naïve Bayes Classifier:

- o Naïve Bayes is one of the fast and easy ML algorithms to predict a class of datasets.
- o It can be used for Binary as well as Multi-class Classifications.
- o It performs well in Multi-class predictions as compared to the other Algorithms.
- o It is the most popular choice for **text classification problems**.

Disadvantages of Naïve Bayes Classifier:

- o Naive Bayes assumes that all features are independent or unrelated, so it cannot learn the relationship between features.

### Applications of Naïve Bayes Classifier:

- o It is used for **Credit Scoring**.
- o It is used in **medical data classification**.
- o It can be used in **real-time predictions** because Naïve Bayes Classifier is an eager learner.
- o It is used in Text classification such as **Spam filtering** and **Sentiment analysis**.

### Linear Regression

Linear regression is one of the most popular and simple machine learning algorithms that is used for predictive analysis. Here, **predictive analysis** defines prediction of something, and linear regression makes predictions for *continuous numbers* such as **salary, age, etc.**

It shows the linear relationship between the dependent and independent variables, and shows how the dependent variable(y) changes according to the independent variable (x).

It tries to best fit a line between the dependent and independent variables, and this best fit line is knowns as the regression line.

The equation for the regression line is:
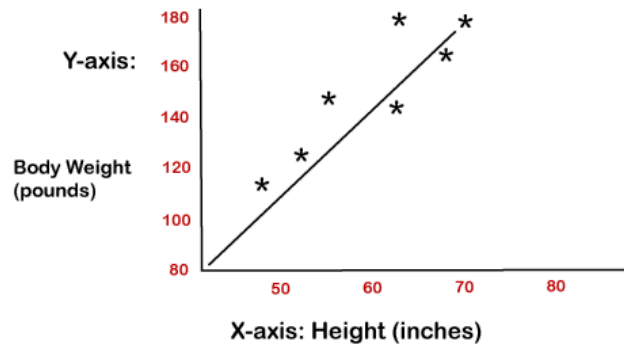
$y = a_0 + a*x + b$

Here, y= dependent variable

x= independent variable

$a_0$ = Intercept of line.

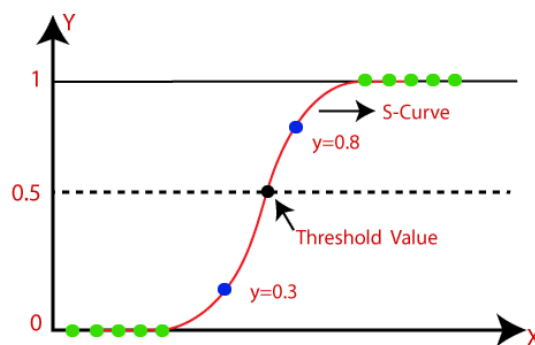Linear regression is further divided into two types:

- o **Simple Linear Regression:** In simple linear regression, a single independent variable is used to predict the value of the dependent variable.
- o **Multiple Linear Regression:** In multiple linear regression, more than one independent variables are used to predict the value of the dependent variable.

The below diagram shows the linear regression for prediction of weight according to height: Read more..

X-axis: Height (inches)

## Logistic Regression in Machine Learning

o Logistic regression is one of the most popular Machine Learning algorithms, which comes under the Supervised Learning technique. It is used for predicting the categorical dependent variable using a given set of independent variables.

o Logistic regression predicts the output of a categorical dependent variable. Therefore the outcome must be a categorical or discrete value. It can be either Yes or No, 0 or 1, true or False, etc. but instead of giving the exact value as 0 and 1, **it gives the probabilistic values which lie between 0 and 1**.

o Logistic Regression is much similar to the Linear Regression except that how they are used. Linear Regression is used for solving Regression problems, whereas **Logistic regression is used for solving the classification problems**.

o In Logistic regression, instead of fitting a regression line, we fit an "S" shaped logistic function, which predicts two maximum values (0 or 1).

o The curve from the logistic function indicates the likelihood of something such as whether the cells are cancerous or not, a mouse is obese or not based on its weight, etc.

o Logistic Regression is a significant machine learning algorithm because it has the ability to provide probabilities and classify new data using continuous and discrete datasets.

o Logistic Regression can be used to classify the observations using different types of data and can easily determine the most effective variables used for the classification. The below image is showing the logistic function:

Logistic Function (Sigmoid Function):

o The sigmoid function is a mathematical function used to map the predicted values to probabilities.

o It maps any real value into another value within a range of 0 and 1.

o The value of the logistic regression must be between 0 and 1, which cannot go beyond this limit, so it forms a curve like the "S" form. The S-form curve is called the Sigmoid function or the logistic function.

o In logistic regression, we use the concept of the threshold value, which defines the probability of either 0 or 1. Such as values above the threshold value tends to 1, and a value below the threshold values tends to 0.

Assumptions for Logistic Regression:

o The dependent variable must be categorical in nature.

o The independent variable should not have multi-collinearity.

Logistic Regression Equation:

The Logistic regression equation can be obtained from the Linear Regression equation. The mathematical steps to get Logistic Regression equations are given below:

o We know the equation of the straight line can be written as:

$$y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \cdots + b_n x_n$$

o In Logistic Regression y can be between 0 and 1 only, so for this let's divide the above equation by (1-y):

$$\frac{y}{1-y} ; \text{0 for y= 0, and infinity for y=1}$$

o But we need range between -[infinity] to +[infinity], then take logarithm of the equation it will become:

$$\log\left[\frac{y}{1-y}\right] = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \cdots + b_n x_n$$

The above equation is the final equation for Logistic Regression.

Type of Logistic Regression:

On the basis of the categories, Logistic Regression can be classified into three types:
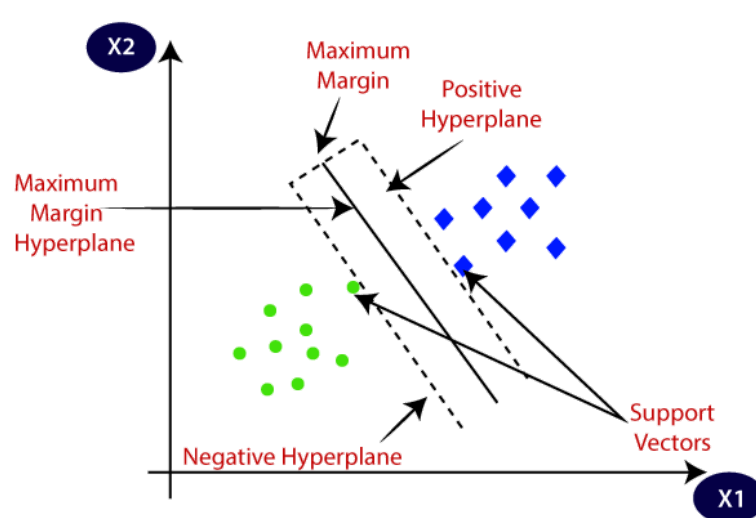
- o **Binomial:** In binomial Logistic regression, there can be only two possible types of the dependent variables, such as 0 or 1, Pass or Fail, etc.
- o **Multinomial:** In multinomial Logistic regression, there can be 3 or more possible unordered types of the dependent variable, such as "cat", "dogs", or "sheep"
- o **Ordinal:** In ordinal Logistic regression, there can be 3 or more possible ordered types of dependent variables, such as "low", "Medium", or "High".
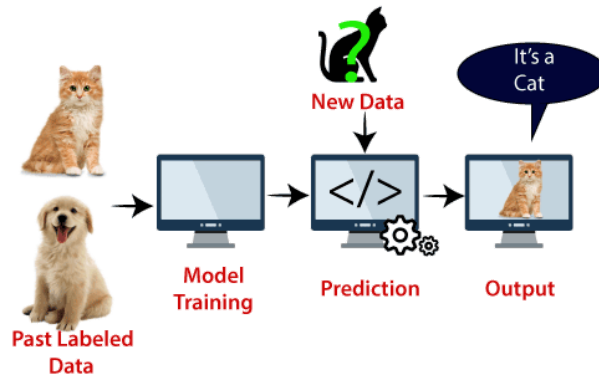
## Support Vector Machine Algorithm

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. However, primarily, it is used for Classification problems in Machine Learning.

The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane.

SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called as support vectors, and hence algorithm is termed as Support Vector Machine. Consider the below diagram in which there are two different categories that are classified using a decision boundary or hyperplane:



**Example:** SVM can be understood with the example that we have used in the KNN classifier. Suppose we see a strange cat that also has some features of dogs, so if we want a model that can accurately identify whether it is a cat or dog, so such a model can be created by using the SVM algorithm. We will first train our model with lots of images of cats and dogs so that it can learn about different features of cats and dogs, and then we test it with this strange creature. So as support vector creates a decision boundary between these two data (cat and dog) and choose extreme cases (support vectors), it will see the extreme case of cat and dog. On the basis of the support vectors, it will classify it as a cat. Consider the below diagram:

SVM algorithm can be used for **Face detection, image classification, text categorization,** etc.

## Types of SVM

**SVM can be of two types:**

o **Linear SVM:** Linear SVM is used for linearly separable data, which means if a dataset can be classified into two classes by using a single straight line, then such data is termed as linearly separable data, and classifier is used called as Linear SVM classifier.

o **Non-linear SVM:** Non-Linear SVM is used for non-linearly separated data, which means if a dataset cannot be classified by using a straight line, then such data is termed as non-linear data and classifier used is called as Non-linear SVM classifier.

## Hyperplane and Support Vectors in the SVM algorithm:

**Hyperplane:** There can be multiple lines/decision boundaries to segregate the classes in n-dimensional space, but we need to find out the best decision boundary that helps to classify the data points. This best boundary is known as the hyperplane of SVM.

The dimensions of the hyperplane depend on the features present in the dataset, which means if there are 2 features (as shown in image), then hyperplane will be a straight line. And if there are 3 features, then hyperplane will be a 2-dimension plane.

We always create a hyperplane that has a maximum margin, which means the maximum distance between the data points.
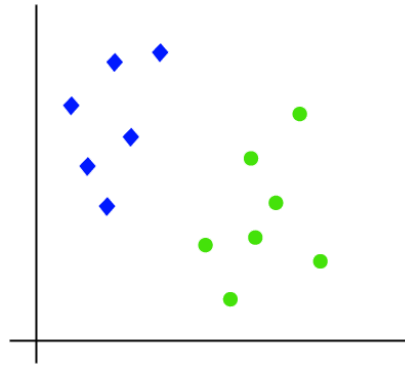
**Support Vectors:**

The data points or vectors that are the closest to the hyperplane and which affect the position of the hyperplane are termed as Support Vector. Since these vectors support the hyperplane, hence called a Support vector.
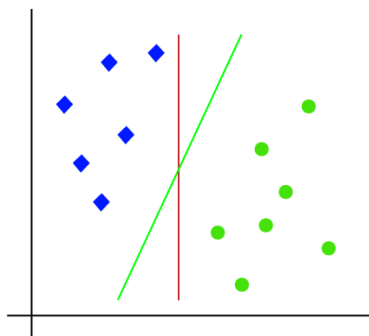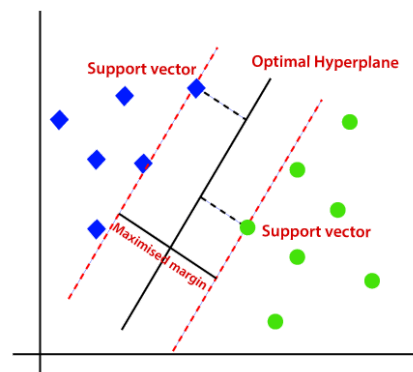
## How does SVM works?

**Linear SVM:**

The working of the SVM algorithm can be understood by using an example. Suppose we have a dataset that has two tags (green and blue), and the dataset has two features x1 and x2. We want a classifier that can classify the pair(x1, x2) of coordinates in either green or blue. Consider the below image:

So as it is 2-d space so by just using a straight line, we can easily separate these two classes. But there can be multiple lines that can separate these classes. Consider the below image:
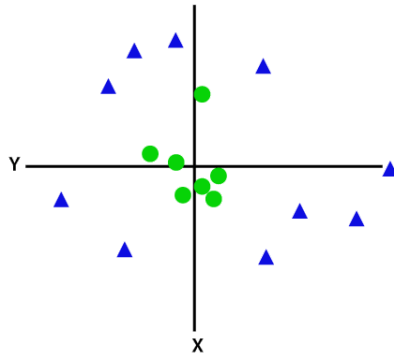
Hence, the SVM algorithm helps to find the best line or decision boundary; this best boundary or region is called as a **hyperplane**. SVM algorithm finds the closest point of the lines from both the classes. These points are called support vectors. The distance between the vectors and the hyperplane is called as **margin**. And the goal of SVM is to maximize this margin. The **hyperplane** with maximum margin is called the **optimal hyperplane**.
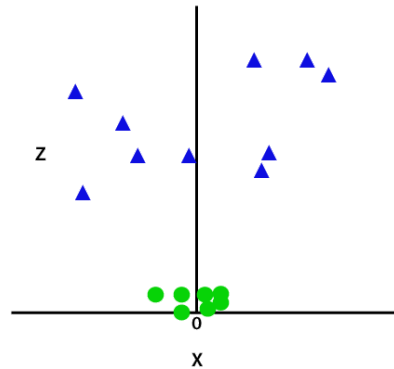
**Non-Linear SVM:**

If data is linearly arranged, then we can separate it by using a straight line, but for non-linear data, we cannot draw a single straight line. Consider the below image:
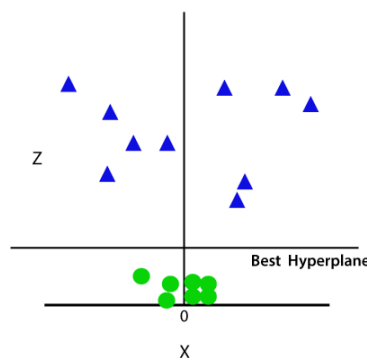
So to separate these data points, we need to add one more dimension. For linear data, we have used two dimensions x and y, so for non-linear data, we will add a third dimension z. It can be calculated as:
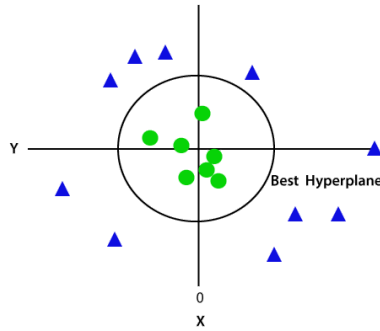
$$z = x^2 + y^2$$

By adding the third dimension, the sample space will become as below image:



So now, SVM will divide the datasets into classes in the following way. Consider the below image:



Since we are in 3-d Space, hence it is looking like a plane parallel to the x-axis. If we convert it in 2d space with z=1, then it will become as:

Hence we get a circumference of radius 1 in case of non-linear data.

**Generalized Linear Models**

**Prerequisite:**

- Linear Regression
- Logistic Regression

The following article discusses the **Generalized linear models** (GLMs) which explains how Linear regression and Logistic regression are a member of a much broader class of models. GLMs can be used to construct the models for regression and classification problems by using the type of distribution which best describes the data or labels given for training the model. Below given are some types of datasets and the corresponding distributions which would help us in constructing the model for a particular type of data (The term data specified here refers to the **output data** or the **labels** of the dataset).

1. Binary classification data – Bernoulli distribution
2. Real valued data – Gaussian distribution
3. Count-data – Poisson distribution

To understand GLMs we will begin by defining exponential families. Exponential families are a class of distributions whose probability density function(PDF) can be molded into the following form:

**Proof** – **Bernoulli distribution** is a member of the exponential family.

Therefore, on comparing **Eq1** and **Eq2** :

*Note: As mentioned above the value of **phi** (which is the same as the activation or **sigmoid function** for Logistic regression) is not a coincidence. And it will be proved later in the article how Logistic regression model can be derived from the Bernoulli distribution.*

 proof – **Gaussian distribution** is a member of the exponential family.

Therefore, on comparing **Eq1** and **Eq3**:

**Constructing GLMs:**

To construct GLMs for a particular type of data or more generally for linear or logistic classification problems the following three assumptions or design choices are to be considered:

The first assumption is that if **x** is the input data parameterized by theta the resulting output or y will be a member of the exponential family. This means if we are provided with some labeled data our goal is to find the right parameters theta which fits the given model as closely as possible. The third assumption is the least justified and can be considered as a design choice.

**Linear Regression Model:**

To show that Linear Regression is a special case of the GLMs. It is considered that the output labels are **continuous values** and are therefore a **Gaussian distribution**. So, we have

The first equation above corresponds to the first assumption that the output labels (or **target**

**variables**) should be the member of an exponential family, Second equation corresponds to the assumption that the **hypothesis** is equal the **expected value** or **mean** of the distribution and lastly, the third equation corresponds to the assumption that natural parameter and the input parameters follow a linear relationship.

**Logistic Regression Model:**

To show that Logistic Regression is a special case of the GLMs. It is considered that the output labels are **Binary valued** and are therefore a **Bernoulli distribution**. So, we have From the third assumption, it is proven that:

The function that maps the natural parameter to the canonical parameter is known as the **canonical response function** (here, the log-partition function) and the inverse of it is known as the **canonical link function**.
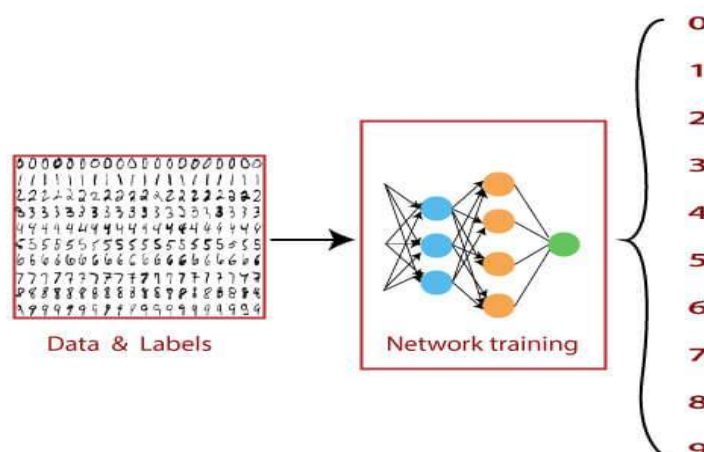
Therefore by using the three assumptions mentioned before it can be proved that the **Logistic** and **Linear Regression** belongs to a much larger family of models known as GLMs.

## MNIST Dataset

The MNIST (**Modified National Institute of Standards and Technology**) database is a large database of handwritten numbers or digits that are used for training various image processing systems. The dataset also widely used for training and testing in the field of **machine learning**. The set of images in the MNIST database are a combination of two of NIST's databases: Special Database 1 and Special Database 3.

The MNIST dataset has **60,000** training images and **10,000** testing images.

The **MNIST** dataset can be online, and it is essentially a database of various handwritten digits. The MNIST dataset has a large amount of data and is commonly used to demonstrate the real power of deep neural networks. Our brain and eyes work together to recognize any numbered image. Our mind is a potent tool, and it's capable of categorizing any image quickly. There are so many shapes of a number, and our mind can easily recognize these shapes and determine what number is it, but the same task is not simple for a computer to complete. There is only one way to do this, which is the use of deep neural network which allows us to train a computer to classify the handwritten digits effectively.
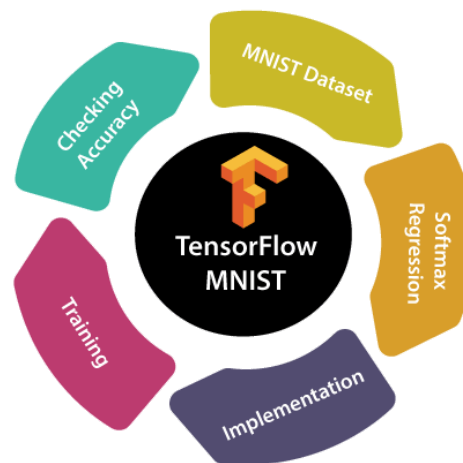


So, we have only dealt with data which contains simple data points on a **Cartesian coordinate system**. From starting till now, we have distributed with binary class datasets. And when we use multiclass datasets, we will use the **Softmax activation** function is quite useful for classifying binary datasets. And it was quite effective in arranging values between

0 and 1. The sigmoid function is not effective for multicausal datasets, and for this purpose, we use the softmax activation function, which is capable of dealing with it.

The MNIST dataset is a multilevel dataset consisting of 10 classes in which we can classify numbers from 0 to 9. The major difference between the datasets that we have used before and the MNIST dataset is the method in which MNIST data is inputted in a neural network.

In the perceptual model and linear regression model, each of the data points was defined by a simple x and y coordinate. This means that the input layer needs two nodes to input single data points.

In the MNIST dataset, a single data point comes in the form of an image. These images included in the **MNIST** dataset are typical of **28*28** pixels such as 28 pixels crossing the horizontal axis and 28 pixels crossing the vertical axis. This means that a single image from the MNIST database has a total of 784 pixels that must be analyzed. The input layer of our neural network has 784 nodes to explain one of these images.
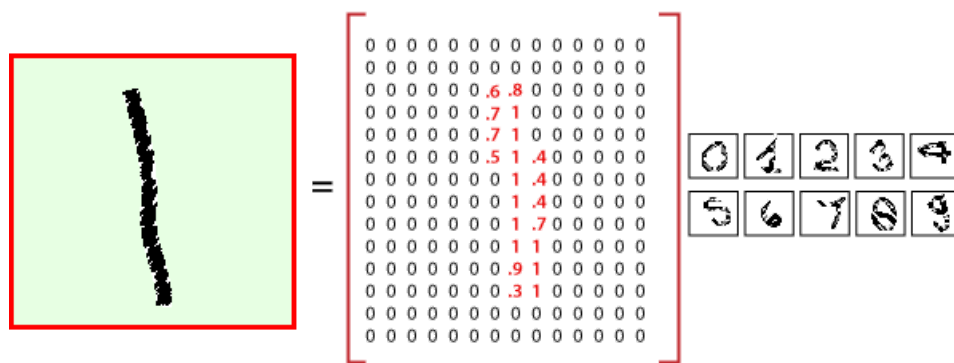


Here, we will see how to create a function that is a model for recognizing handwritten digits by looking at each pixel in the image. Then using TensorFlow to train the model to predict the image by making it look at thousands of examples which are already labeled. We will then check the model's accuracy with a test dataset.

**MNIST dataset in TensorFlow, containing information of handwritten digits spitted into three parts:**

- o **Training Data (mnist.train) -55000 datapoints**
- o **Validation Data (mnist.validate) -5000 datapoint**
- o **Test Data (mnist.test) -10000 datapoints**

Now before we start, it is important to note that every data point has two parts: an image (x) and a corresponding label (y) describing the actual image and each image is a 28x28 array, i.e., 784 numbers. The label of the image is a number between 0 and 9 corresponding to the TensorFlow MNIST image. To download and use MNIST dataset, use the following **commands**:

1. from tensorflow.examples.tutorials.mnist **import** input_data
2. mnist = input_data.read_data_sets("MNIST_data/", one_hot=True)

## Ranking

Ranking is a regression machine learning technique.

- About Ranking
- Ranking Methods
- Ranking Algorithms
    - XGBoost

## About Ranking

Ranking is a machine learning technique to rank items.

Ranking is useful for many applications in information retrieval such as e-commerce, social networks, recommendation systems, and so on. For example, a user searches for an article or an item to buy online. To build a recommendation system, it becomes important that similar articles or items of relevance appear to the user such that the user clicks or purchases the item. A simple regression model can predict the probability of a user to click an article or buy an item. However, it is more practical to use ranking technique and be able to order or rank the articles or items to maximize the chances of getting a click or purchase. The prioritization of the articles or the items influence the decision of the users.

The ranking technique directly ranks items by training a model to predict the ranking of one item over another item. In the training model, it is possible to have items, ranking one over the other by having a "score" for each item. Higher ranked items have higher scores and lower ranked items have lower scores. Using these scores, a model is built to predict which item ranks higher than the other.

## Ranking Methods

Oracle Machine Learning supports pairwise and listwise ranking methods through XGBoost.

For a training data set, in a number of sets, each set consists of objects and labels representing their ranking. A ranking function is constructed by minimizing a certain loss function on the training data. Using test data, the ranking function is applied to get a ranked list of objects.

Ranking is enabled for XGBoost using the regression function. OML4SQL supports pairwise and listwise ranking methods through XGBoost.

Pairwise ranking: This approach regards a pair of objects as the learning instance. The pairs and lists are defined by supplying the same case_id value. Given a pair of objects, this approach gives an optimal ordering for that pair. Pairwise losses are defined by the order of the two objects. In OML4SQL, the algorithm uses LambdaMART to perform pairwise ranking with the goal of minimizing the average number of inversions in ranking.

Listwise ranking: This approach takes multiple lists of ranked objects as learning instance. The items in a list must have the same case_id. The algorithm uses LambdaMART to perform list-wise ranking.

See Also:

- "Ranking Measures and Loss Functions in Learning to Rank" a research paper presentation at https://www.researchgate.net/
- *Oracle Database PL/SQL Packages and Types Reference* for a listing and explanation of the available model settings for XGBoost.

Note:

The term hyperparameter is also interchangeably used for model setting.

**Related Topics**

- XGBoost
- DBMS_DATA_MINING — Algorithm Settings: XGBoost

**Ranking Algorithms**

Ranking falls under the Regression function.

OML4SQL supports XGBoost algorithm for ranking.

**Related Topics**

- XGBoost

## Structured outputs

**Structured prediction** or **structured (output) learning** is an umbrella term for supervised machine learning techniques that involves predicting structured objects, rather than scalar discrete or real values.[1]

Similar to commonly used supervised learning techniques, structured prediction models are typically trained by means of observed data in which the true prediction value is used to adjust model parameters. Due to the complexity of the model and the interrelations of predicted variables the process of prediction using a trained model and of training itself is often computationally infeasible and approximate inference and learning methods are used.

## Applications

For example, the problem of translating a [natural language](#) sentence into a syntactic representation such as a [parse tree](#) can be seen as a structured prediction problem[2] in which the structured output domain is the set of all possible parse trees. Structured prediction is also used in a wide variety of application domains including [bioinformatics](#), [natural language processing](#), [speech recognition](#), and [computer vision](#).

### Example: sequence tagging

Sequence tagging is a class of problems prevalent in [natural language processing](#), where input data are often sequences (e.g. sentences of text). The sequence tagging problem appears in several guises, e.g. [part-of-speech tagging](#) and [named entity recognition](#). In POS tagging, for example, each word in a sequence must receive a "tag" (class label) that expresses its "type" of word:

| This | [DT](#) |
|------|------|
| is | [VBZ](#) |
| a | [DT](#) |
| tagged | [JJ](#) |
| sentence | [NN](#) |

The main challenge of this problem is to resolve [ambiguity](#): the word "sentence" can also be [a verb](#) in English, and so can "tagged".

While this problem can be solved by simply performing [classification](#) of individual tokens, that approach does not take into account the empirical fact that tags do not occur independently; instead, each tag displays a strong [conditional dependence](#) on the tag of the previous word. This fact can be exploited in a sequence model such as a [hidden Markov model](#) or [conditional random field](#)[2] that predicts the entire tag sequence for a sentence, rather than just individual tags, by means of the [Viterbi algorithm](#).

## Techniques

Probabilistic [graphical models](#) form a large class of structured prediction models. In particular, [Bayesian networks](#) and [random fields](#) are popular. Other algorithms and models for structured prediction include [inductive logic programming](#), [case-based reasoning](#), [structured SVMs](#), [Markov logic networks](#), [Probabilistic Soft Logic](#), and [constrained conditional models](#). Main techniques:

- [Conditional random field](#)
- [Structured support vector machine](#)
- [Structured k-Nearest Neighbours](#)
- [Recurrent neural network](#), in particular [Elman network](#)

### Structured perceptron

One of the easiest ways to understand algorithms for general structured prediction is the structured perceptron of Collins.[3] This algorithm combines the perceptron algorithm for learning linear classifiers with an inference algorithm (classically the Viterbi algorithm when used on sequence data) and can be described abstractly as follows. First define a "joint feature function" $\Phi(\mathbf{x}, \mathbf{y})$ that maps a training sample $\mathbf{x}$ and a candidate prediction $\mathbf{y}$ to a vector of length $n$ ($\mathbf{x}$ and $\mathbf{y}$ may have any structure; $n$ is problem-dependent, but must be fixed for each model). Let GEN be a function that generates candidate predictions.