

Naive Bayes Classification

```
] # Splitting the dataset into the Training set and Test set
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(x, y, test_size = 0.20, random_state = 0)
```

```
]
# Feature Scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
```

```
▶ # Training the Naive Bayes model on the Training set
from sklearn.naive_bayes import GaussianNB
classifier = GaussianNB()
classifier.fit(X_train, y_train)

# Predicting the Test set results
y_pred = classifier.predict(X_test)

# Making the Confusion Matrix
from sklearn.metrics import confusion_matrix, accuracy_score
ac = accuracy_score(y_test, y_pred)
cm = confusion_matrix(y_test, y_pred)
```

```
☞ 0.7604166666666666
[[ 86   0   0  11   0   0]
 [  0 388  32  80   0   1]
 [  0  29  96   0   0   6]
 [ 71  38   0 289   0   1]
 [  0   0   0   0  31   2]
 [  0   0  21   0   7  59]]
```

```
[ ] classifier = BernoulliNB()
    classifier.fit(X_train, y_train)

BernoulliNB()

[ ] # Predicting the Test set results
    y_pred = classifier.predict(X_test)

    # Making the Confusion Matrix
    from sklearn.metrics import confusion_matrix, accuracy_score
    ac = accuracy_score(y_test, y_pred)
    cm = confusion_matrix(y_test, y_pred)

[ ] print(ac)
    print(cm)

0.6057692307692307
[[ 0  1  0 96  0  0]
 [ 0 304 55 120  0 22]
 [ 0 35 29  1  0 66]
 [ 0 35  0 364  0  0]
 [ 0  2  3  0  0 28]
 [ 0  4 24  0  0 59]]
```

Naive Bayes classifiers are a collection of classification algorithms based on Bayes' Theorem. It is not a single algorithm but a family of algorithms where all of them share a common principle, i.e. every pair of features being classified is independent of each other.

Bayes' Theorem

Bayes' Theorem finds the probability of an event occurring given the probability of another event that has already occurred. Bayes' theorem is stated mathematically as the following equation:

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

where A and B are events and $P(B) \neq 0$.

Basically, we are trying to find probability of event A, given the event B is true. Event B is also termed as evidence.

$P(A)$ is the priori of A (the prior probability, i.e. Probability of event before evidence is seen). The evidence is an attribute value of an unknown instance(here, it is event B).

$P(A|B)$ is a posteriori probability of B, i.e. probability of event after evidence is seen.

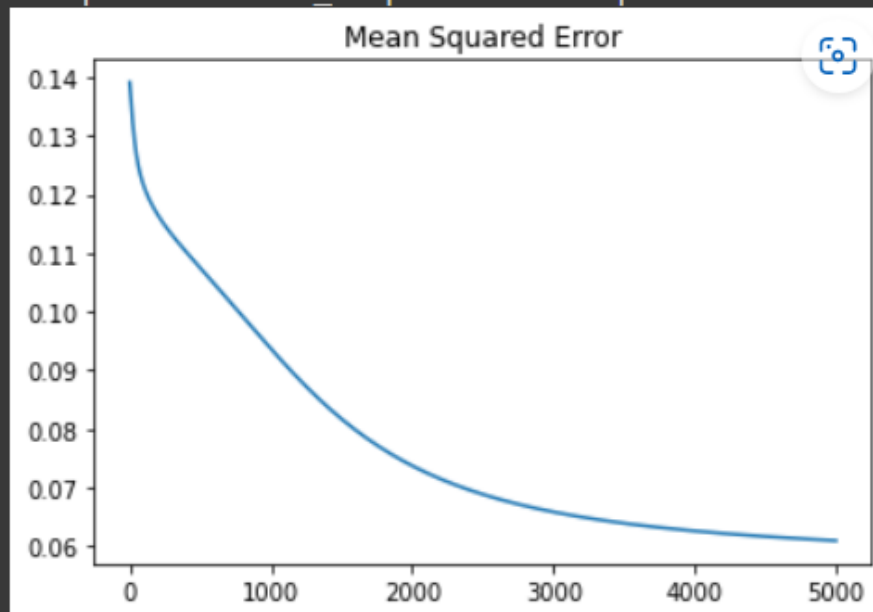
Now, with regards to our dataset, we can apply Bayes' theorem in following way:

$$P(y|X) = \{P(X|y) P(y)\}/\{P(X)\}$$

BPN -Back Propagation in ML

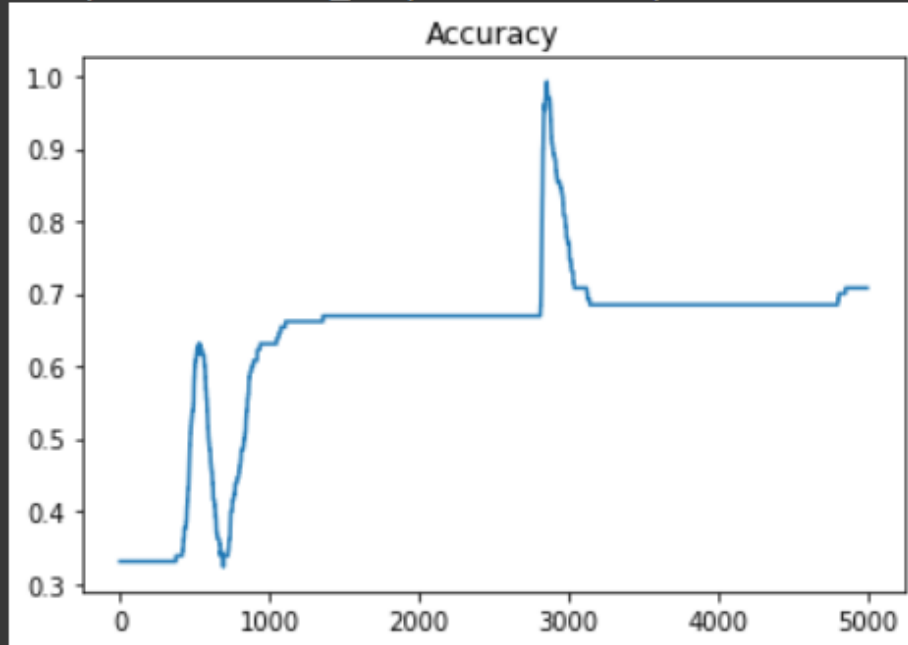
```
results.mse.plot(title="Mean Squared Error")
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f4896dee190>
```



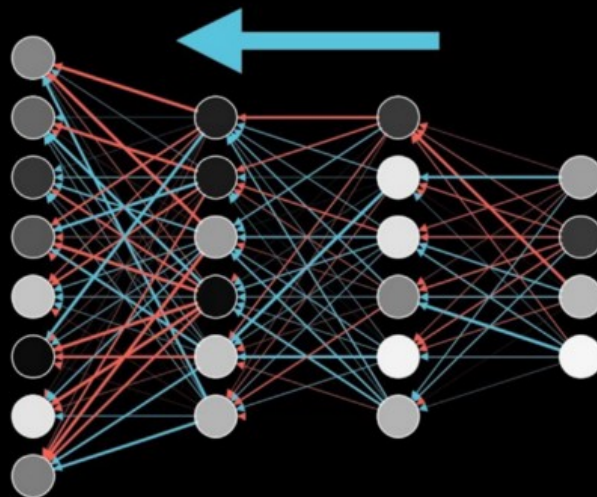
```
[ ] results.accuracy.plot(title="Accuracy")
```

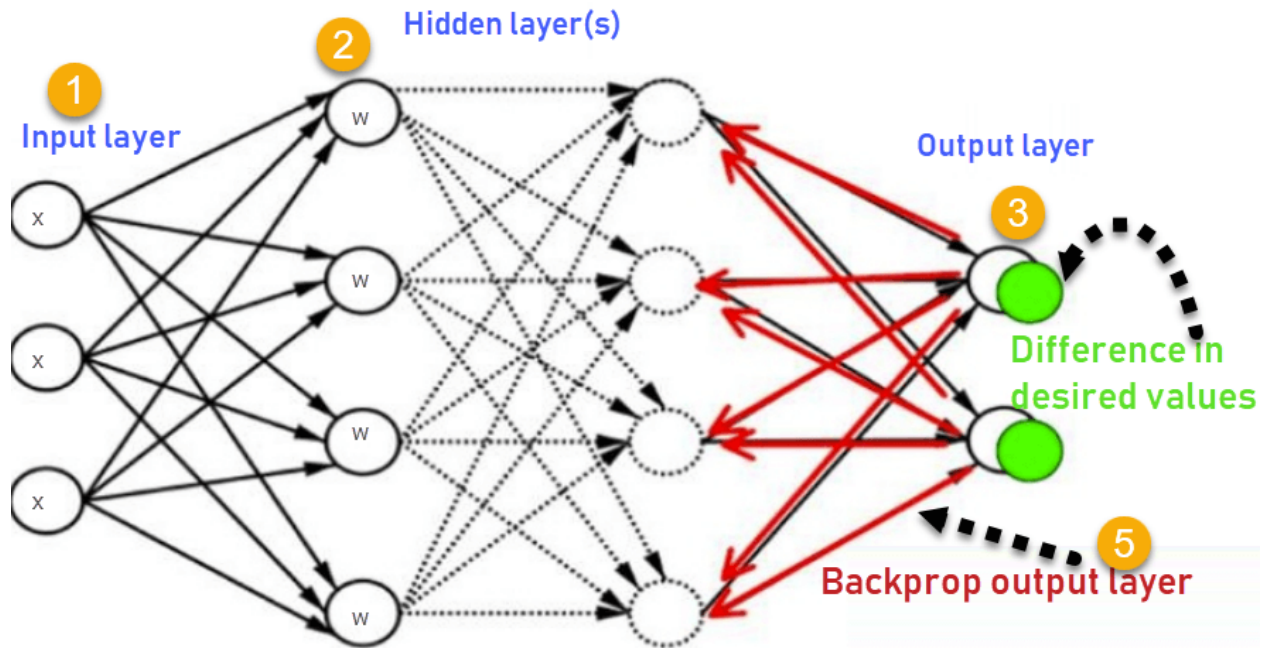
```
<matplotlib.axes._subplots.AxesSubplot at 0x7f4897461390>
```



Accuracy: 0.8

Backpropagation





Backpropagation is the essence of neural network training. It is the method of fine-tuning the weights of a neural network based on the error rate obtained in the previous epoch (i.e., iteration). Proper tuning of the weights allows you to reduce error rates and make the model reliable by increasing its generalization.

Backpropagation in neural network is a short form for “backward propagation of errors.” It is a standard method of training artificial neural networks. This method helps calculate the gradient of a loss function with respect to all the weights in the network.

Types of Backpropagation Networks

Two Types of Backpropagation Networks are:

- Static Back-propagation
- Recurrent Backpropagation

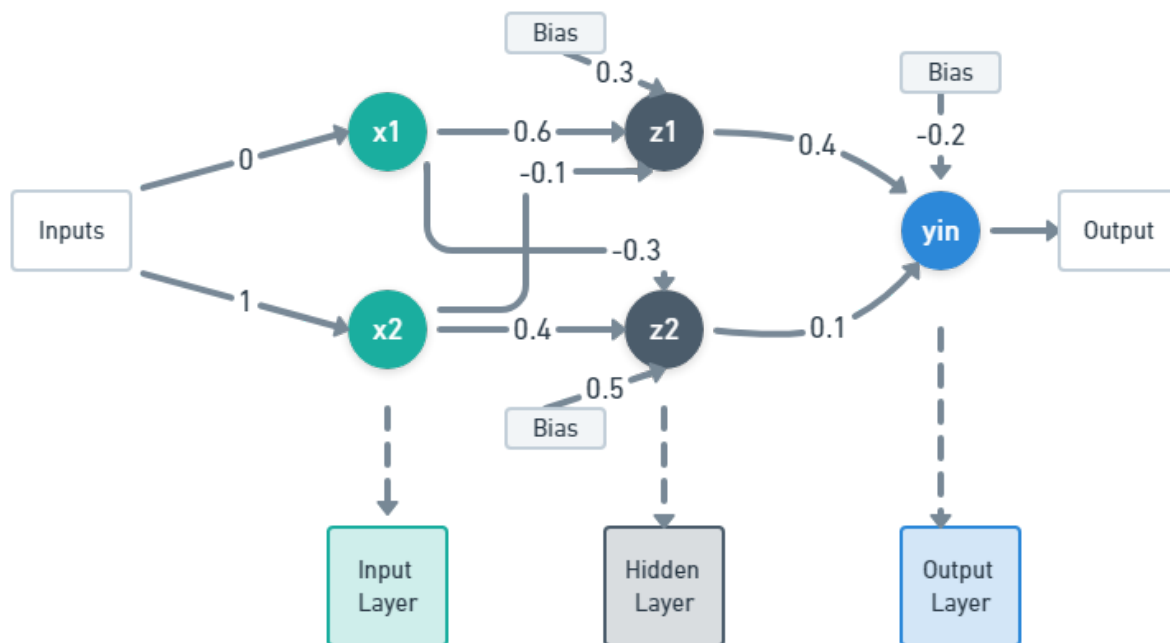
Static back-propagation:

It is one kind of backpropagation network which produces a mapping of a static input for static output. It is useful to solve static classification issues like optical character recognition.

Recurrent Backpropagation:

Recurrent Back propagation in data mining is fed forward until a fixed value is achieved. After that, the error is computed and propagated backward. The main difference between both of these methods is: that the mapping is rapid in static back-propagation while it is nonstatic in recurrent backpropagation.

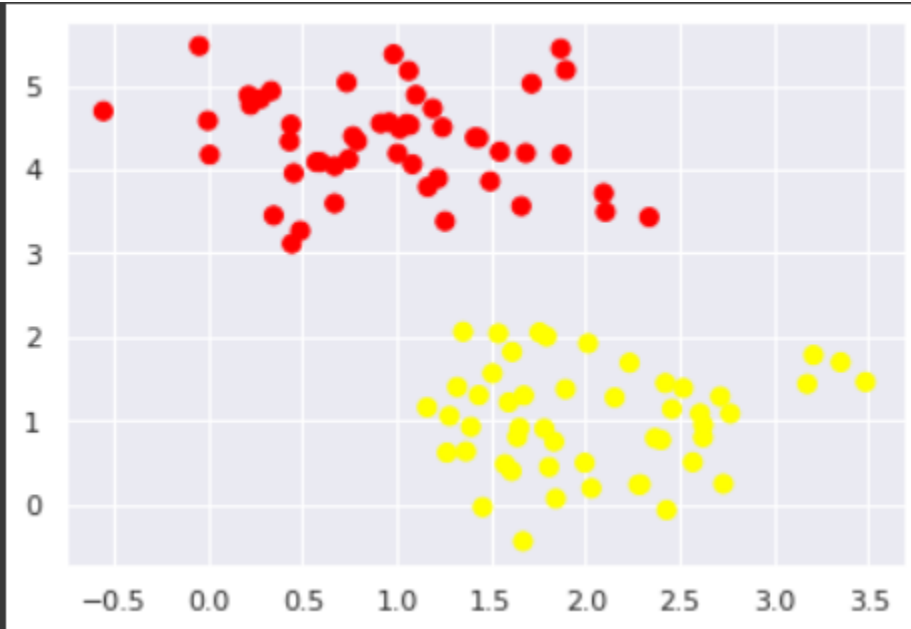
Multi layer feed forward neural network



```
Iteration: 500. Loss: 0.6242267489433289. Accuracy: 86.23999786376953
Iteration: 1000. Loss: 0.42309409379959106. Accuracy: 89.55999755859375
Iteration: 1500. Loss: 0.25730085372924805. Accuracy: 90.61000061035156
Iteration: 2000. Loss: 0.3684542179107666. Accuracy: 91.30000305175781
Iteration: 2500. Loss: 0.37141624093055725. Accuracy: 91.68000030517578
Iteration: 3000. Loss: 0.2659453749656677. Accuracy: 92.02999877929688
```

SVM

```
[ ] from sklearn.datasets import make_blobs
X,y= make_blobs(n_samples=100,centers=2,random_state=0,cluster_std=0.60)
plt.scatter(X[:,0],X[:,1],c=y,s=50,cmap="autumn");
print(X,y)
```



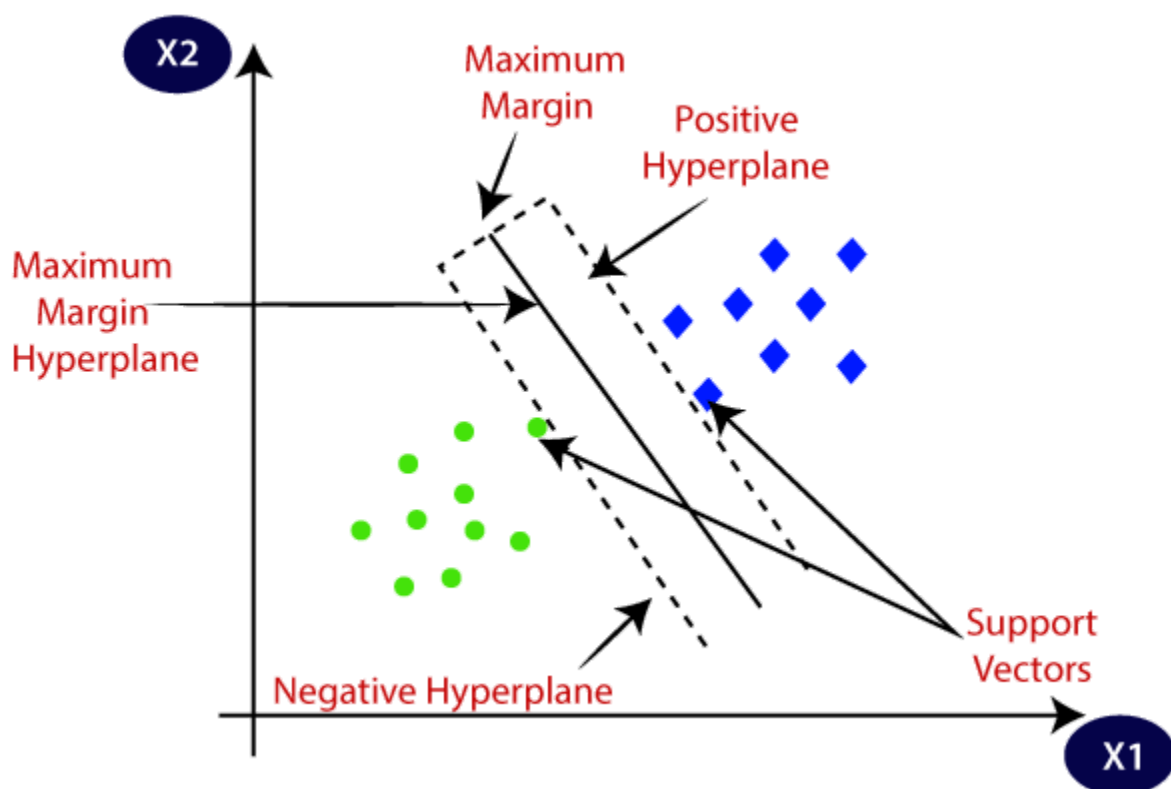
```
[ ] [[11  0]
     [ 0  9]]
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	11
1	1.00	1.00	1.00	9
accuracy			1.00	20
macro avg	1.00	1.00	1.00	20
weighted avg	1.00	1.00	1.00	20

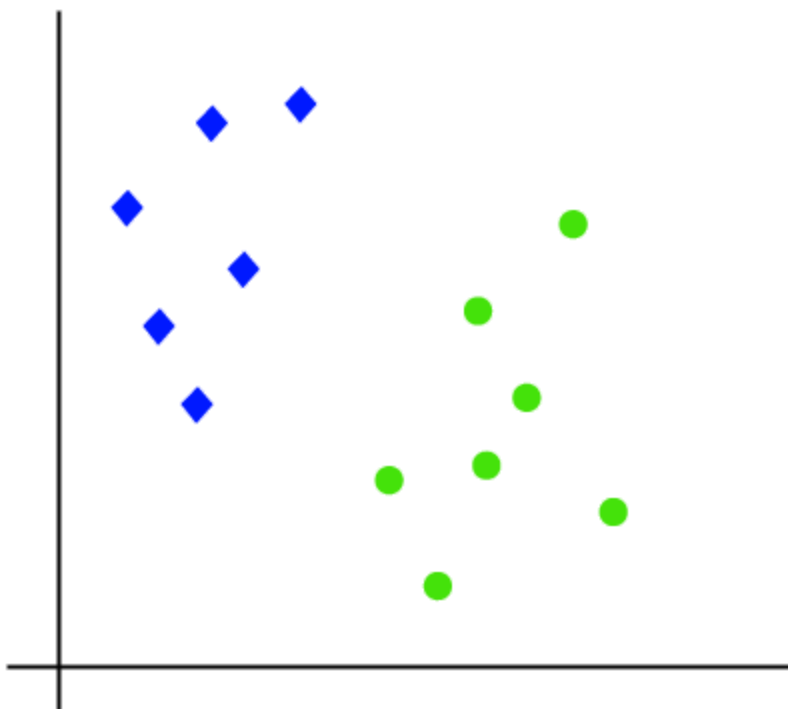
HYPER PARAMETERS

Hyperparameters are parameters whose values control the learning process and determine the values of model parameters that a learning algorithm ends up learning. The prefix 'hyper_' suggests that they are 'top-level' parameters that control the learning process and the model parameters that result from it.

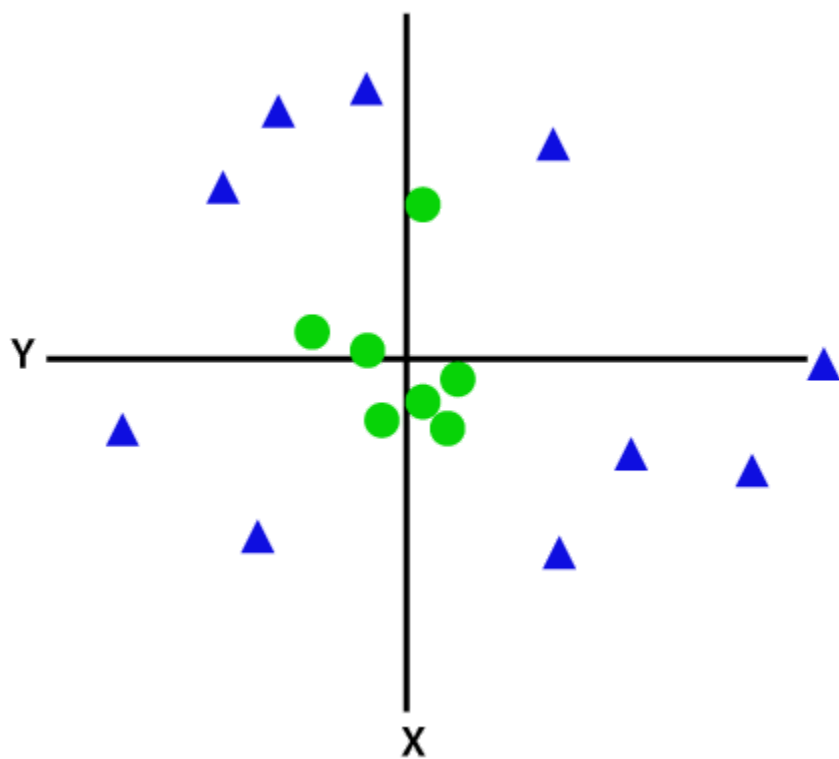
As a machine learning engineer designing a model, you choose and set hyperparameter values that your learning algorithm will use before the training of the model even begins. In this light, hyperparameters are said to be external to the model because the model cannot change its values during learning/training.



Linera svm



Non linear svm



Logistic Regression

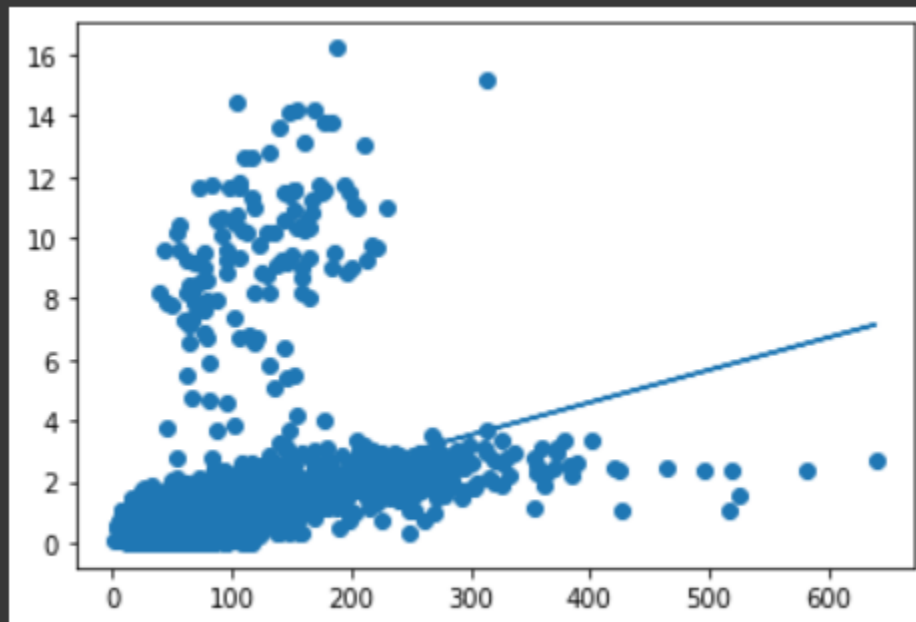
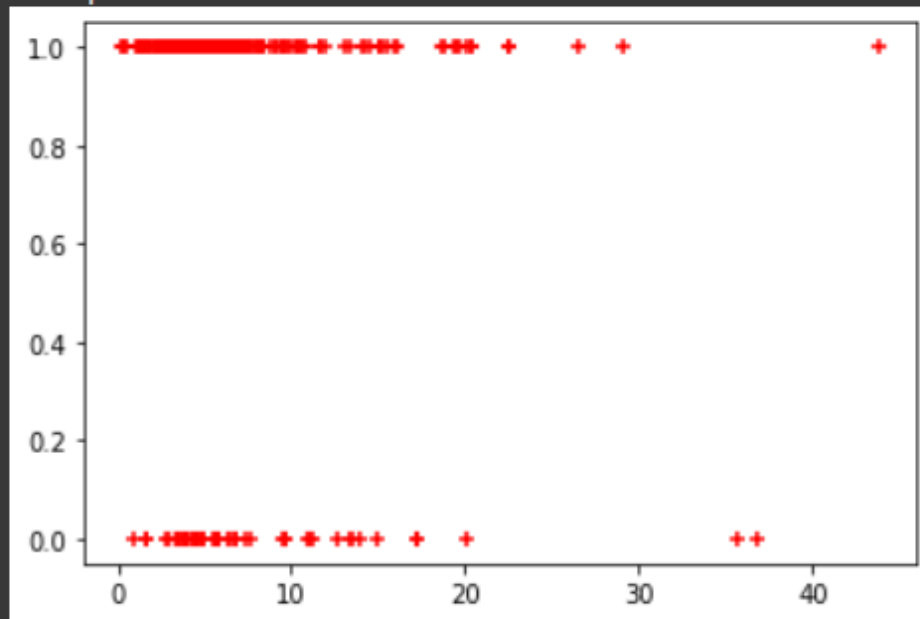
This type of statistical model (also known as *logit model*) is often used for classification and predictive analytics. Logistic regression estimates the probability of an event occurring, such as voted or didn't vote, based on a given dataset of independent variables. Since the outcome is a probability, the dependent variable is bounded between 0 and 1. In logistic regression, a logit transformation is applied on the odds—that is, the probability of success divided by the probability of failure. This is also commonly known as the log odds, or the natural logarithm of odds, and this logistic function is represented by the following formulas:

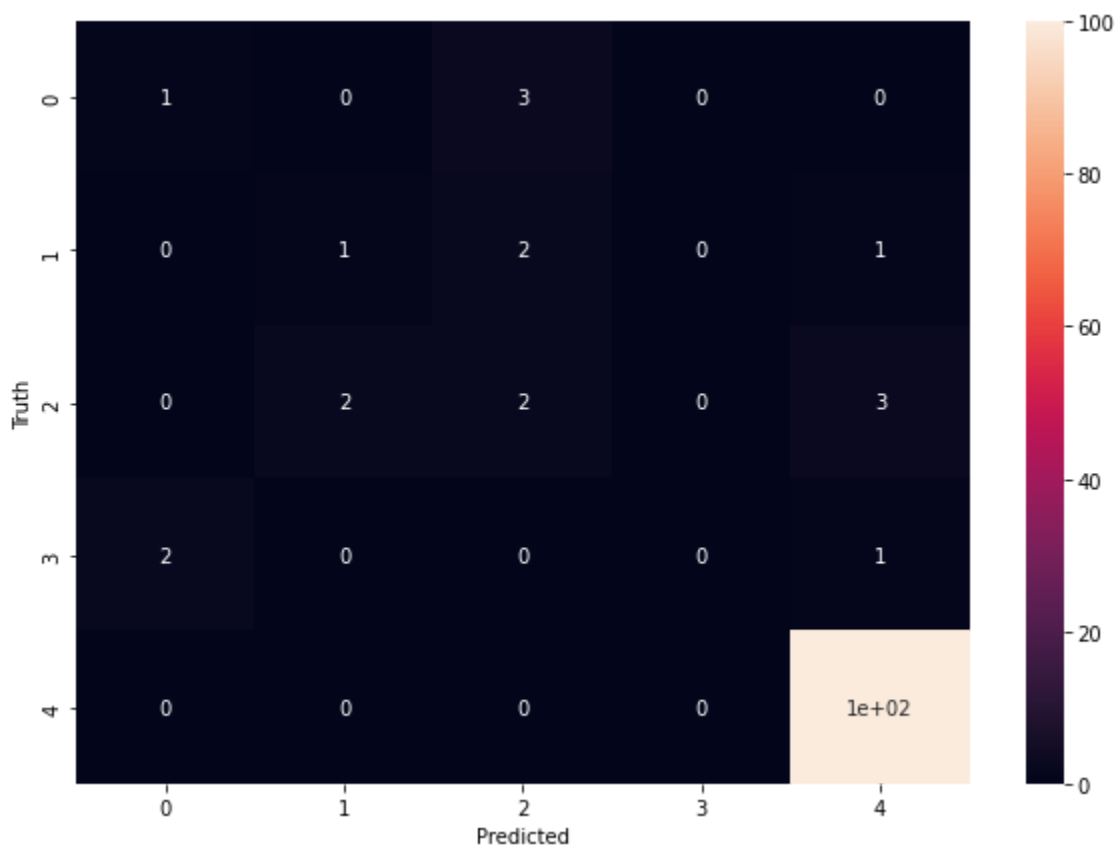
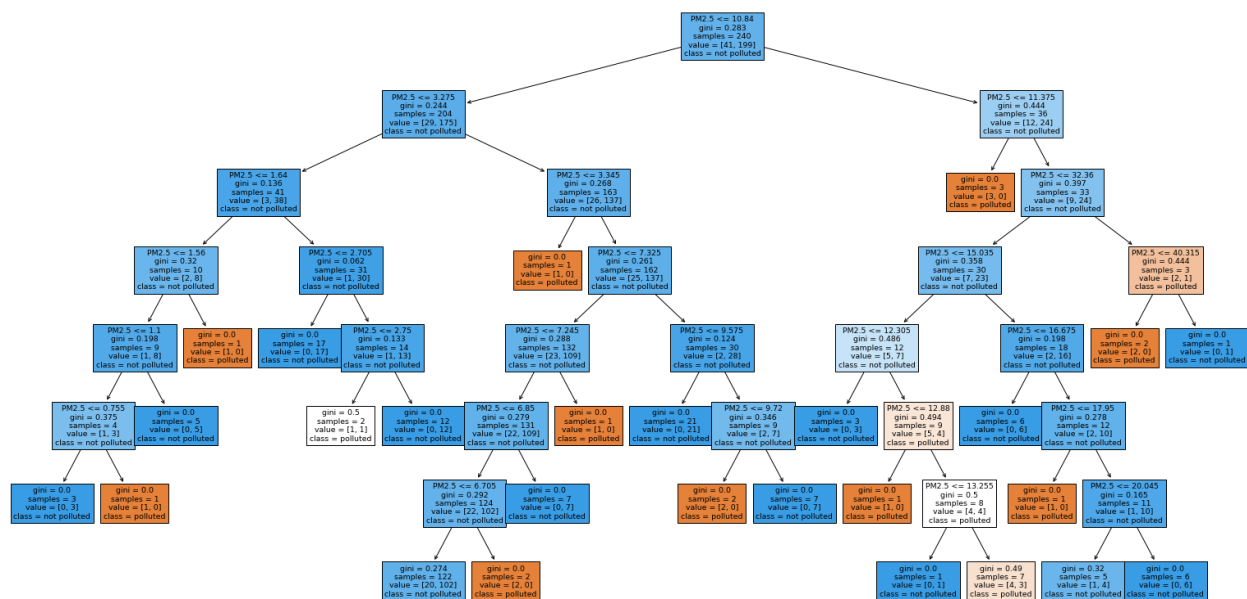
$$\text{Logit}(\pi) = 1/(1 + \exp(-\pi))$$

$$\ln(\pi/(1-\pi)) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

In this logistic regression equation, $\text{logit}(\pi)$ is the dependent or response variable and x is the independent variable. The beta parameter, or coefficient, in this model is commonly estimated via maximum likelihood estimation (MLE). This method tests different values of beta through multiple iterations to optimize for the best fit of log odds. All of these iterations produce the log likelihood function, and logistic regression seeks to maximize this function to find the best parameter estimate. Once the optimal coefficient (or coefficients if there is more than one independent variable) is found, the conditional probabilities for each observation can be calculated, logged, and summed together to yield a predicted probability. For binary classification, a probability less than .5 will predict 0 while a probability greater than 0 will predict 1. After the model has been computed, it's best practice to evaluate the how well the model predicts the dependent variable, which is called goodness of fit. The Hosmer–Lemeshow test is a popular method to assess model fit.

<matplotlib.collections.PathCollection at 0x7fa8447f5c10>





Decision Tree

A tree has many analogies in real life, and turns out that it has influenced a wide area of machine learning, covering both classification and regression. In decision analysis, a decision tree can be used to visually and explicitly represent decisions and decision making. As the name goes, it uses a tree-like model of decisions. Though a commonly used tool in data mining for deriving a strategy to reach a particular goal, its also widely used in machine learning, which will be the main focus of this article.

How can an algorithm be represented as a tree?

For this let's consider a very basic example that uses titanic data set for predicting whether a passenger will survive or not. Below model uses 3 features/attributes/columns from the data set, namely sex, age and sibsp (number of spouses or children along).

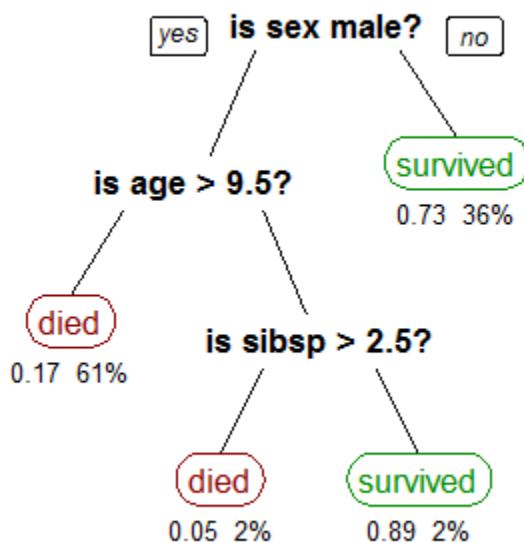
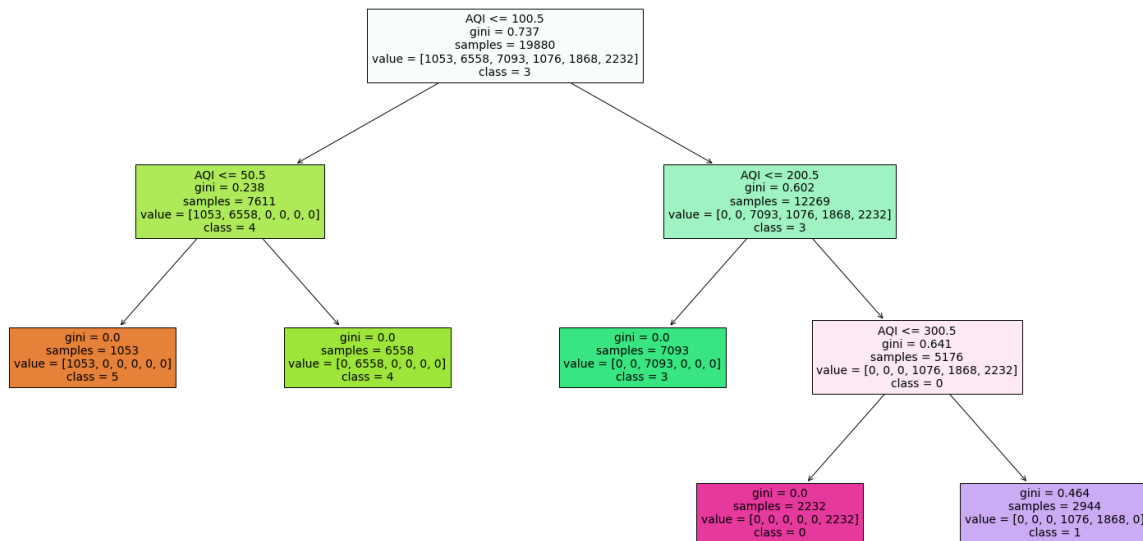


Image taken from wikipedia

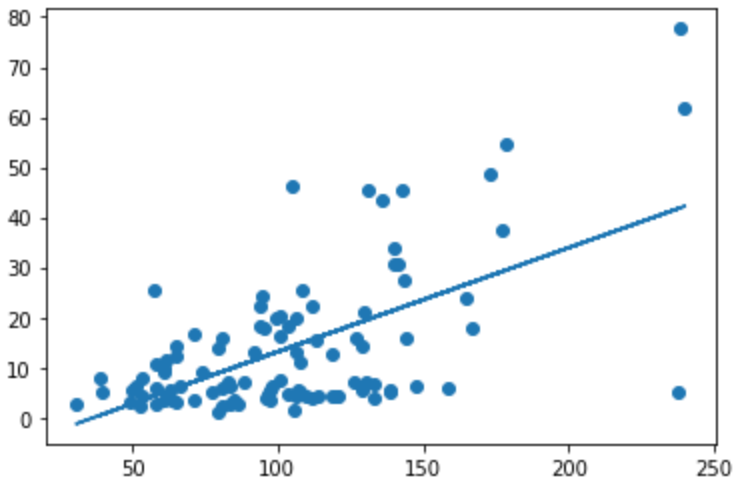
A decision tree is drawn upside down with its root at the top. In the image on the left, the bold text in black represents a condition/internal node, based on which the tree splits into branches/ edges. The end of the branch that doesn't split anymore is the decision/leaf, in this case, whether the passenger died or survived, represented as red and green text respectively.



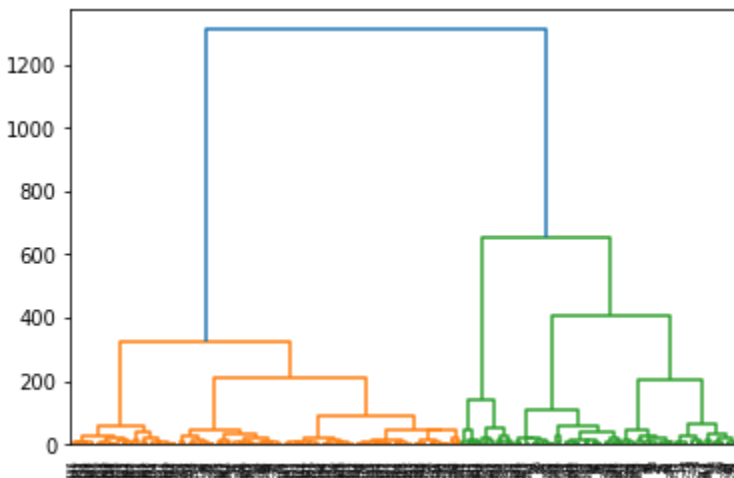
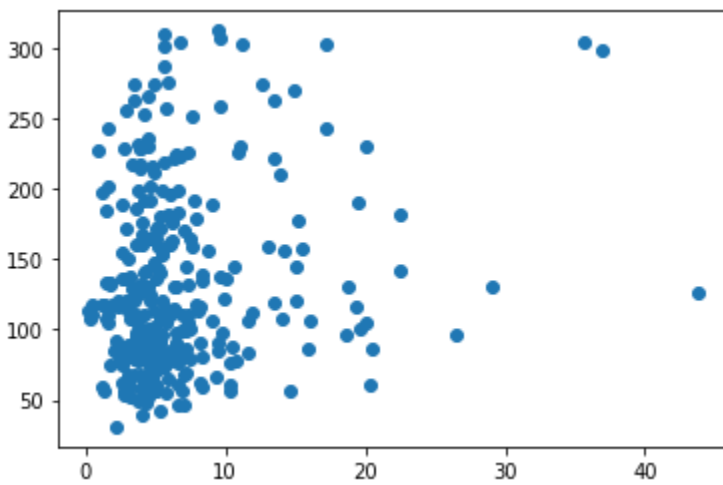
Linear Regression

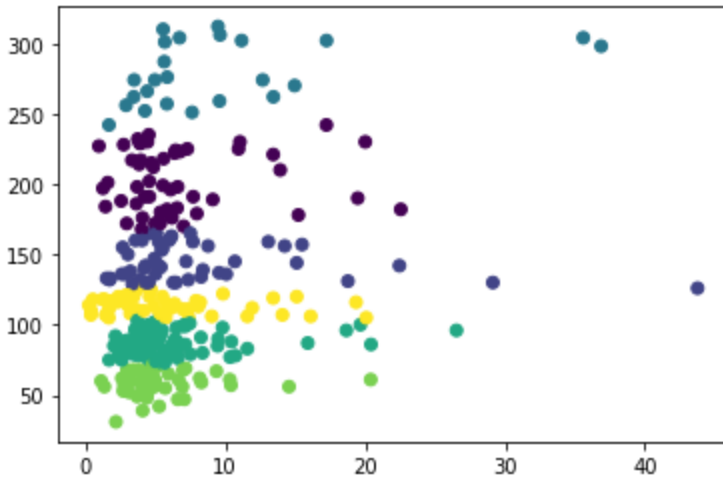
Linear regression analysis is used to predict the value of a variable based on the value of another variable. The variable you want to predict is called the dependent variable. The variable you are using to predict the other variable's value is called the independent variable.

This form of analysis estimates the coefficients of the linear equation, involving one or more independent variables that best predict the value of the dependent variable. Linear regression fits a straight line or surface that minimizes the discrepancies between predicted and actual output values. There are simple linear regression calculators that use a “least squares” method to discover the best-fit line for a set of paired data. You then estimate the value of X (dependent variable) from Y (independent variable).



Hierarchical clustering algorithm





HC VS KMEANS

K-Means is that it needs us to pre-enter the number of clusters (K) but Hierarchical clustering has no such requirements. The algorithm on itself deduces the optimum number of cluster and displays it form of dendrogram.

Performance of K-Means on spherical data is better than that of HCA Hierarchical clustering is a purely agglomerative approach and goes on to build one giant cluster.

K-Means algorithm in all its iterations has same number of clusters.

K-Means need circular data, while Hierarchical clustering has no such requirement.

K-Means uses median or mean to compute centroid for representing cluster while HCA has various linkage method that may or may not employ the centroid.

With introduction of mini batches K-Means can work with very large datasets but HCA lacks in this regard.

Hierarchical methods are suited for cases which require arrangement of the clusters into a natural hierarchy. In K-means all clusters are on same level i.e. similar WCSS or cohesiveness.

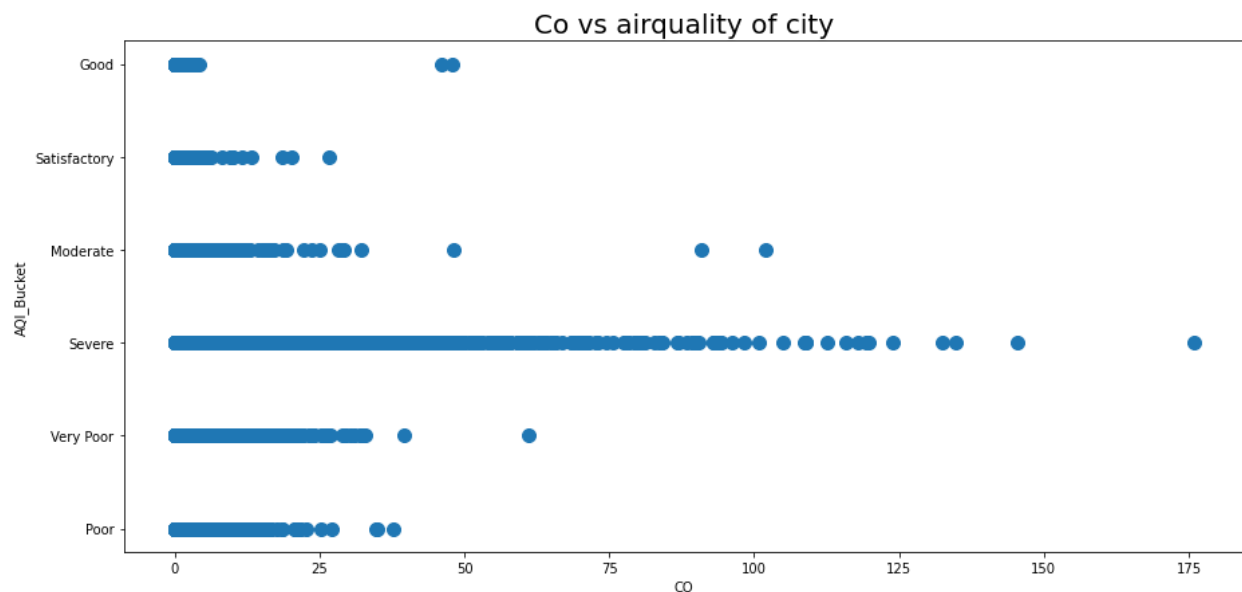
HCA can produce reproducible results while older versions of K-Means can't

K-Means simply divides data into mutually exclusive subsets while HCA arranges it into a tree format.

K-Means

K-means clustering is one of the simplest and popular unsupervised machine learning algorithms.

Typically, unsupervised algorithms make inferences from datasets using only input vectors without referring to known, or labelled, outcomes.



A cluster refers to a collection of data points aggregated together because of certain similarities.

You'll define a target number k , which refers to the number of centroids you need in the dataset. A centroid is the imaginary or real location representing the center of the cluster.

Every data point is allocated to each of the clusters through reducing the in-cluster sum of squares.

In other words, the K-means algorithm identifies k number of centroids, and then allocates every data point to the nearest cluster, while keeping the centroids as small as possible.

