

## ASSIGNMENT 4

**Name: Aakanksha Bhondve**

**Grno: 21810939**

**Rollno: 322004**

**Comp B1**

### **Aim –**

Build a Data model in Python using any classification model (Decision Tree or Naïve Bayes) and infer the result using accuracy score. Compare different classification models (not limited to NB and DT only) with respect to feature selection and accuracy. Infer the result - which model best suit for the dataset chosen.

### **Theory –**

**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.

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.

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

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

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.

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

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

**Naïve Bayes** algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.

It is mainly used in text classification that includes a high-dimensional training dataset.

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.

It is a probabilistic classifier, which means it predicts on the basis of the probability of an object.

Some popular examples of Naïve Bayes Algorithm are spam filtration, Sentimental analysis, and classifying articles.

## Code and Output –

### Classification using Decision Tree

```
In [1]: 1 # Importing the libraries
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 import pandas as pd

In [2]: 1 # Importing the dataset
        2 dataset = pd.read_csv('Social_Network_Ads.csv')
        3 X = dataset.iloc[:, [2, 3]].values
        4 y = dataset.iloc[:, 4].values
```

```
In [33]: 1 dataset
```

```
Out[33]:
```

	User ID	Gender	Age	EstimatedSalary	Purchased
0	15624510	Male	19	19000	0
1	15810944	Male	35	20000	0
2	15668575	Female	26	43000	0
3	15603246	Female	27	57000	0
4	15804002	Male	19	76000	0
...	...	...	...	...	...
395	15691863	Female	46	41000	1
396	15706071	Male	51	23000	1
397	15654296	Female	50	20000	1
398	15755018	Male	36	33000	0
399	15594041	Female	49	36000	1

400 rows × 5 columns

```
In [4]: 1 # Splitting the dataset into the Training set and Test set
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
```

```
In [5]: 1 # Feature Scaling
2 from sklearn.preprocessing import StandardScaler
3 sc = StandardScaler()
4 X_train = sc.fit_transform(X_train)
5 X_test = sc.transform(X_test)
```

```
In [6]: 1 # Fitting Decision Tree Classification to the Training set
2 from sklearn.tree import DecisionTreeClassifier
3 classifier = DecisionTreeClassifier(criterion = 'entropy', random_state = 0)
4 classifier.fit(X_train, y_train)
```

```
Out[6]: DecisionTreeClassifier(class_weight=None, criterion='entropy', max_depth=None,
                             max_features=None, max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=1, min_samples_split=2,
                             min_weight_fraction_leaf=0.0, presort=False,
                             random_state=0, splitter='best')
```

```
In [7]: 1 # Predicting the Test set results
2 y_pred = classifier.predict(X_test)
```

```
In [8]: 1 # Making the Confusion Matrix
2 from sklearn.metrics import confusion_matrix
3 cm = confusion_matrix(y_test, y_pred)
```

### Accuracy of desicion tree

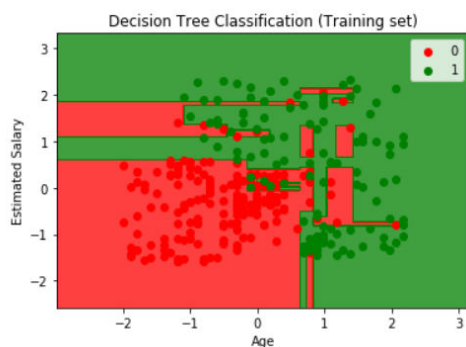
```
In [25]: 1 # Accuracy score of the classifier
2 accuracy_score = (cm[0,0]+cm[1,1])/100
3 print(accuracy_score)
```

0.91

```
In [9]: 1 # Visualising the Training set results
2 from matplotlib.colors import ListedColormap
3 X_set, y_set = X_train, y_train
4 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
5                       np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
6 plt.contourf(X1, X2, classifier.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
7              alpha = 0.75, cmap = ListedColormap(('red', 'green')))
8 plt.xlim(X1.min(), X1.max())
9 plt.ylim(X2.min(), X2.max())
10 for i, j in enumerate(np.unique(y_set)):
11     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
12                c = ListedColormap(('red', 'green'))(i), label = j)
13 plt.title('Decision Tree Classification (Training set)')
14 plt.xlabel('Age')
15 plt.ylabel('Estimated Salary')
16 plt.legend()
17 plt.show()
```

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

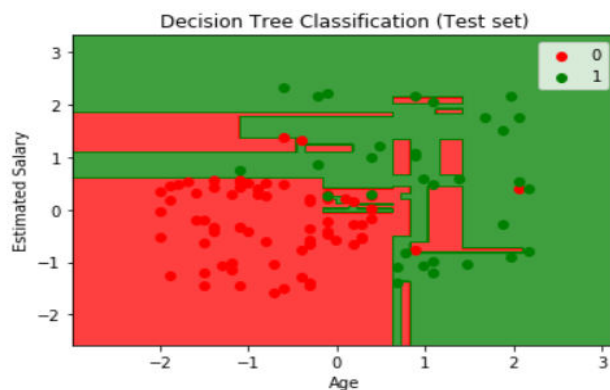
'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.



```
In [10]: 1 # Visualising the Test set results
2 from matplotlib.colors import ListedColormap
3 X_set, y_set = X_test, y_test
4 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
5                       np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
6 plt.contourf(X1, X2, classifier.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
7              alpha = 0.75, cmap = ListedColormap(('red', 'green')))
8 plt.xlim(X1.min(), X1.max())
9 plt.ylim(X2.min(), X2.max())
10 for i, j in enumerate(np.unique(y_set)):
11     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
12                c = ListedColormap(('red', 'green'))(i), label = j)
13 plt.title('Decision Tree Classification (Test set)')
14 plt.xlabel('Age')
15 plt.ylabel('Estimated Salary')
16 plt.legend()
17 plt.show()
```

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.



## Classification using Naive Bayes

```
In [26]: 1 # Fitting Naive Bayes to the Training set
2 from sklearn.naive_bayes import GaussianNB
3 classifier = GaussianNB()
4 classifier.fit(X_train, y_train)
```

Out[26]: GaussianNB(priors=None, var\_smoothing=1e-09)

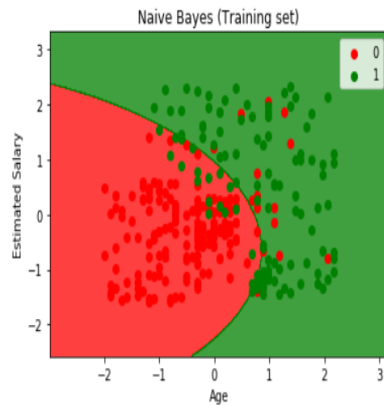
```
In [27]: 1 # Predicting the Test set results
2 y_pred = classifier.predict(X_test)
```

```
In [28]: 1 # Making the Confusion Matrix
2 from sklearn.metrics import confusion_matrix
3 cm1 = confusion_matrix(y_test, y_pred)
```

```
In [30]: 1 # Visualising the Training set results
2 from matplotlib.colors import ListedColormap
3 X_set, y_set = X_train, y_train
4 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
5                       np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
6 plt.contourf(X1, X2, classifier.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
7              alpha = 0.75, cmap = ListedColormap(('red', 'green')))
8 plt.xlim(X1.min(), X1.max())
9 plt.ylim(X2.min(), X2.max())
10 for i, j in enumerate(np.unique(y_set)):
11     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
12                c = ListedColormap(('red', 'green'))(i), label = j)
13 plt.title('Naive Bayes (Training set)')
14 plt.xlabel('Age')
15 plt.ylabel('Estimated Salary')
16 plt.legend()
17 plt.show()
18
```

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

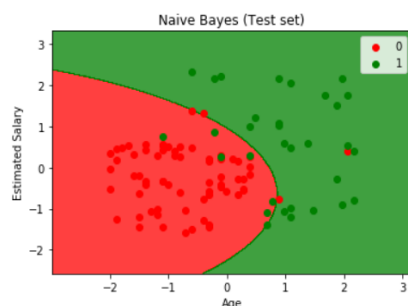
'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.



```
In [31]: 1 # Visualising the Test set results
2 from matplotlib.colors import ListedColormap
3 X_set, y_set = X_test, y_test
4 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
5                      np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
6 plt.contourf(X1, X2, classifier.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
7             alpha = 0.75, cmap = ListedColormap(('red', 'green')))
8 plt.xlim(X1.min(), X1.max())
9 plt.ylim(X2.min(), X2.max())
10 for i, j in enumerate(np.unique(y_set)):
11     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
12               c = ListedColormap(('red', 'green'))(i), label = j)
13 plt.title('Naive Bayes (Test set)')
14 plt.xlabel('Age')
15 plt.ylabel('Estimated Salary')
16 plt.legend()
17 plt.show()
```

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.



### Accuracy of Naive Bayes

```
In [32]: 1 # Accuracy score of the classifier
2 accuracy_score1= (cm1[0,0]+cm1[1,1])/100
3 print(accuracy_score1)
```

0.9

*The difference in the accuracy score by 1% is seen, thus Decision Tree is more accurate*