

DTC Algorithm Implementation

1. Data Pre-Processing Step:

Below is the code for the pre-processing step:

```
1. # importing libraries
2. import numpy as nm
3. import matplotlib.pyplot as mtp
4. import pandas as pd
5.
6. #importing datasets
7. data_set= pd.read_csv('user_data.csv')
8.
9. #Extracting Independent and dependent Variable
10. x= data_set.iloc[:, [2,3]].values
11. y= data_set.iloc[:, 4].values
12.
13. # Splitting the dataset into training and test set.
14. from sklearn.model_selection import train_test_split
15. x_train, x_test, y_train, y_test= train_test_split(x, y, test_size= 0.25, random_state=0)
16.
17. #feature Scaling
18. from sklearn.preprocessing import StandardScaler
19. st_x= StandardScaler()
20. x_train= st_x.fit_transform(x_train)
21. x_test= st_x.transform(x_test)
```

In the above code, we have pre-processed the data. Where we have loaded the dataset, which is given as:

Index	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
5	15728773	Male	27	58000	0
6	15598044	Female	27	84000	0
7	15694829	Female	32	150000	1
8	15600575	Male	25	33000	0
9	15727311	Female	35	65000	0
10	15570769	Female	26	80000	0
11	15606274	Female	26	52000	0
12	15746139	Male	20	86000	0
13	15704987	Male	32	18000	0
14	15628972	Male	18	82000	0
15	15697686	Male	29	80000	0

2. Fitting a Decision-Tree algorithm to the Training set

- Now we will fit the model to the training set. For this, we will import the **DecisionTreeClassifier** class from **sklearn.tree** library.
- Below is the code for it:
 1. #Fitting Decision Tree classifier 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)

In the above code, we have created a classifier object, in which we have passed two main parameters;

- **"criterion='entropy'":** Criterion is used to measure the quality of split, which is calculated by information gain given by entropy.
- **random_state=0":** For generating the random states.

Below is the output for this:

```
Out[8]:
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')
```

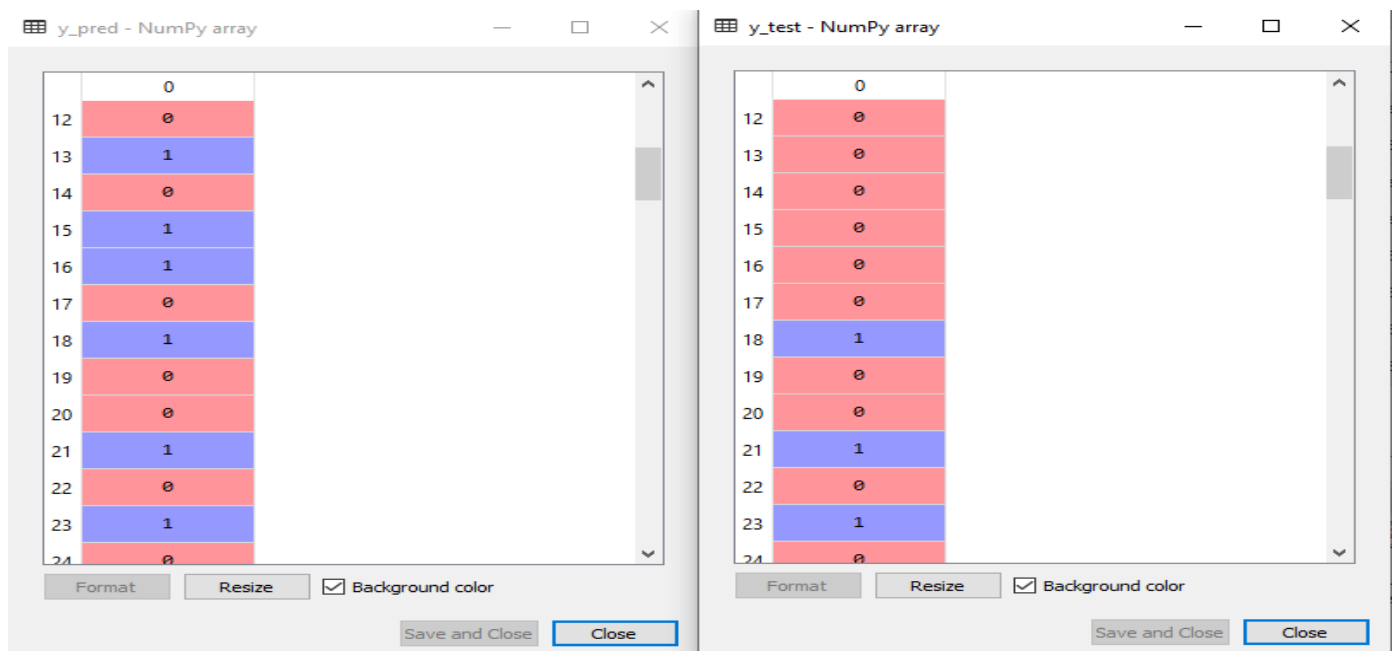
3. Predicting the test result

Now we will predict the test set result. We will create a new prediction vector **y_pred**. Below is the code for it:

1. #Predicting the test set result
2. `y_pred= classifier.predict(x_test)`

Output:

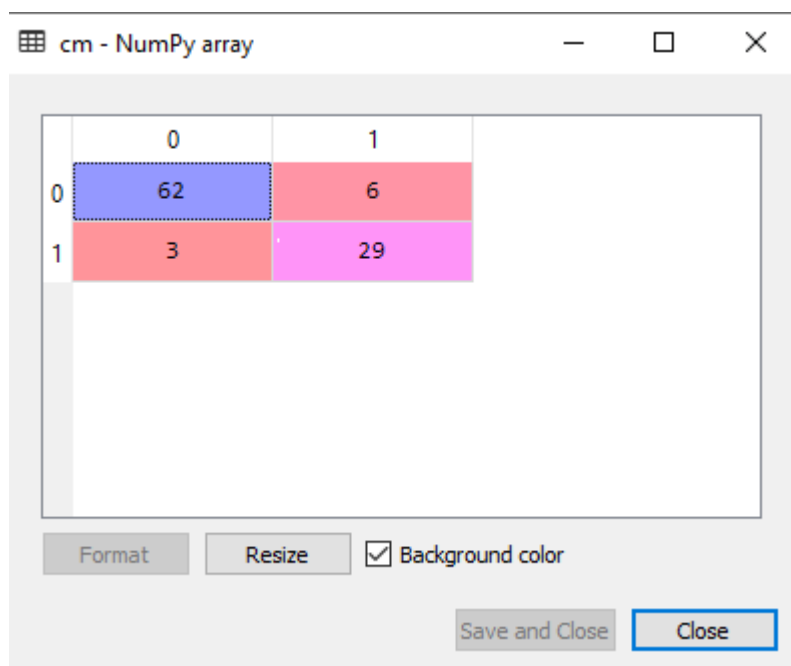
- In the below output image, the predicted output and real test output are given.
- We can clearly see that there are some values in the prediction vector, which are different from the real vector values. These are prediction errors.



4. Test accuracy of the result (Creation of Confusion matrix)

- In the above output, we have seen that there were some incorrect predictions, so if we want to know the number of correct and incorrect predictions, we need to use the confusion matrix.
- Below is the code for it:
 1. #Creating the Confusion matrix
 2. from sklearn.metrics **import** confusion_matrix
 3. cm= confusion_matrix(y_test, y_pred)

Output:



In the above output image, we can see the confusion matrix, which has **6+3= 9 incorrect predictions** and **62+29=91 correct predictions**.

Therefore, we can say that compared to other classification models, the Decision Tree classifier made a good prediction.

5. Visualizing the training set result:

- To visualize the training set result we will plot a graph for the decision tree classifier.
- The classifier will predict yes or No for the users who have either Purchased or Not purchased the SUV car as we did in [Logistic Regression](#).
- Below is the code for it:

1. #Visulaizing the trianing set result
2. from matplotlib.colors **import** ListedColormap
3. x_set, y_set = x_train, y_train
4. x1, x2 = nm.meshgrid(nm.arange(start = x_set[:, 0].min() - 1, stop = x_set[:, 0].max() + 1, step = 0.01),
5. nm.arange(start = x_set[:, 1].min() - 1, stop = x_set[:, 1].max() + 1, step = 0.01))
6. mtp.contourf(x1, x2, classifier.predict(nm.array([x1.ravel(), x2.ravel()]).T).reshape(x1.shape),
7. alpha = 0.75, cmap = ListedColormap(('purple','green')))
8. mtp.xlim(x1.min(), x1.max())
9. mtp.ylim(x2.min(), x2.max())
10. for i, j in enumerate(nm.unique(y_set)):
11. mtp.scatter(x_set[y_set == j, 0], x_set[y_set == j, 1],
12. c = ListedColormap(('purple', 'green'))(i), label = j)
13. mtp.title('Decision Tree Algorithm (Training set)')
14. mtp.xlabel('Age')
15. mtp.ylabel('Estimated Salary')
16. mtp.legend()
17. mtp.show()

Output:



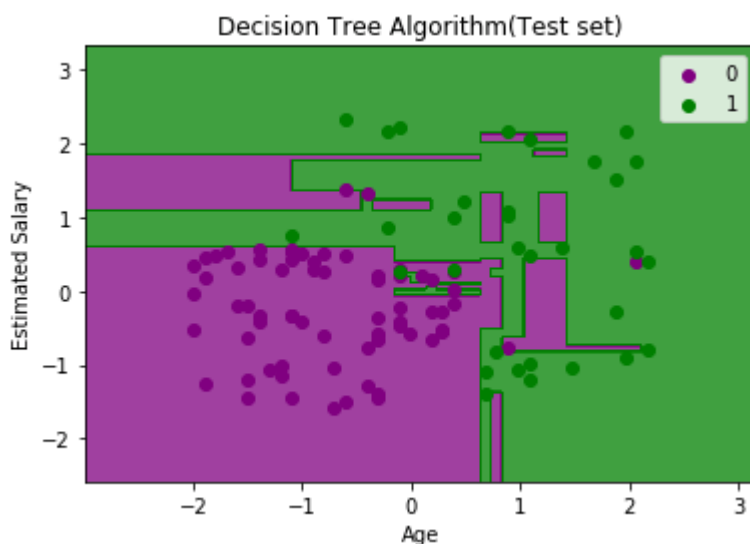
- The above output is completely different from the rest classification models.
- It has both vertical and horizontal lines that are splitting the dataset according to the age and estimated salary variable.
- As we can see, the tree is trying to capture each dataset, which is the case of overfitting.

6. Visualizing the test set result:

Visualization of test set result will be similar to the visualization of the training set except that the training set will be replaced with the test set.

1. #Visulaizing the test set result
2. from matplotlib.colors **import** ListedColormap
3. x_set, y_set = x_test, y_test
4. x1, x2 = nm.meshgrid(nm.arange(start = x_set[:, 0].min() - 1, stop = x_set[:, 0].max() + 1, step = 0.01),
nm.arange(start = x_set[:, 1].min() - 1, stop = x_set[:, 1].max() + 1, step = 0.01))
5. mtp.contourf(x1, x2, classifier.predict(nm.array([x1.ravel(), x2.ravel()]).T).reshape(x1.shape),
alpha = 0.75, cmap = ListedColormap(('purple','green')))
6. mtp.xlim(x1.min(), x1.max())
7. mtp.ylim(x2.min(), x2.max())
8. for i, j in enumerate(nm.unique(y_set)):
9. mtp.scatter(x_set[y_set == j, 0], x_set[y_set == j, 1],
c = ListedColormap(('purple', 'green'))(i), label = j)
10. mtp.title('Decision Tree Algorithm(Test set)')
11. mtp.xlabel('Age')
12. mtp.ylabel('Estimated Salary')
13. mtp.legend()
14. mtp.show()

Output:



- As we can see in the above image that there are some green data points within the purple region and vice versa.
- So, these are the incorrect predictions which we have discussed in the confusion matrix.