



## Department of Computer Science and Engineering (Data Science)

**JHANVI PAREKH**  
**60009210033**  
**CSE (Data Science)**

### Subject: Machine Learning – I (DJ19DSC402)

**AY: 2022-23**

#### Experiment 2 - 3

##### (Decision Tree)

**Aim:** Implement Decision Tree on the given Datasets to build a classifier and Regressor. Apply appropriate pruning method to overcome overfitting.

#### Theory:

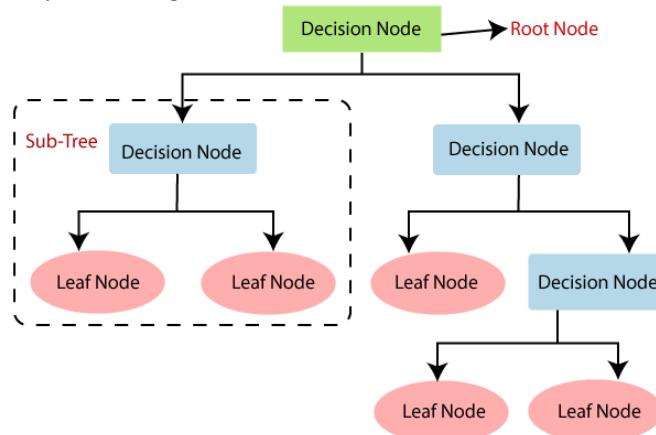
Decision Tree is a **Supervised learning technique** that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where **internal nodes represent the features of a dataset, branches represent the 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.

A decision tree simply asks a question, and based on the answer (Yes/No), it further split the tree into subtrees. Below diagram explains the general structure of a decision tree:





### Department of Computer Science and Engineering (Data Science)

#### 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.



## Department of Computer Science and Engineering (Data Science)

**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.

### Steps in building a Tree

**Step-1:** Begin the tree with the root node, says S, which contains the complete dataset.

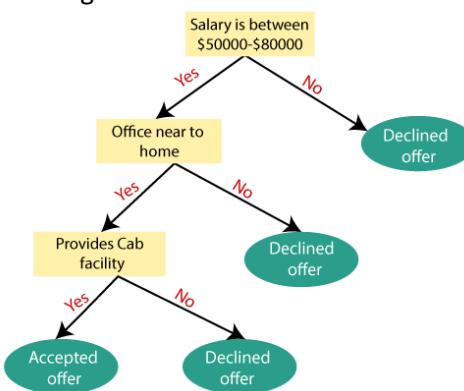
**Step-2:** Find the best attribute in the dataset using **Attribute Selection Measure (ASM)**.

**Step-3:** Divide the S into subsets that contains possible values for the best attributes.

**Step-4:** Generate the decision tree node, which contains the best attribute.

**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:

#### 1. Information Gain:

Information gain is the measurement of changes in entropy after the segmentation of a dataset based on an attribute. It calculates how much information a feature provides us about a class.



## Department of Computer Science and Engineering (Data Science)

According to the value of information gain, we split the node and build the decision tree. 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:  
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:

$$\text{Entropy}(S) = -P(\text{yes}) \log_2 P(\text{yes}) - P(\text{no}) \log_2 P(\text{no})$$

Where,

S= Total number of samples

P(yes)= probability of yes

P(no)= probability of no

### 2. Gini Index:

Gini index is a measure of impurity or purity used while creating a decision tree in the CART(Classification and Regression Tree) algorithm.

An attribute with the low Gini index should be preferred as compared to the high Gini index.

It only creates binary splits, and the CART algorithm uses the Gini index to create binary splits.

Gini index can be calculated using the below formula:

$$\text{Gini Index} = 1 - \sum_j P_j^2$$

### 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:

- Cost Complexity Pruning
- Reduced Error Pruning.

### Lab Assignments to complete in this session:

Use the given dataset and perform the following tasks:

**Dataset 1: PlayTennis.csv**

**Dataset 2: Iris.csv**

**Dataset 3: Breastcancer.csv**

**Dataset 4: car prediction.csv**

1. Implement Decision tree classifier from scratch using Dataset 1 by defining Node class and Tree class.



## Department of Computer Science and Engineering (Data Science)

PlayTennis\_60009210033.ipynb

```
[ ] import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
data=pd.read_csv('/content/PlayTennis.csv')
data.head()

outlook temp humidity windy play
0 sunny hot high False no
1 sunny hot high True no
2 overcast hot high False yes
3 rainy mild high False yes
4 rainy cool normal False yes

[ ] X=data.iloc[:, :-1]
print(X)
y=data.iloc[:, -1]
print(y)

outlook temp humidity windy
0 sunny hot high False
1 sunny hot high True
2 overcast hot high False
3 rainy mild high False
4 rainy cool normal False
5 rainy cool normal True
6 overcast cool normal True
7 sunny mild high False
8 sunny cool normal False
9 rainy mild normal False
10 sunny mild normal True
11 overcast mild high True
12 overcast hot normal False
13 rainy mild high True
0 no
1 no
2 yes
3 yes
4 yes
5 no
6 yes
7 no
8 yes
9 yes
10 yes
11 yes
12 yes
13 no
Name: play, dtype: object
```

0s completed at 10:43PM

Market Brief 22:44 ENG 09-03-2023

PlayTennis\_60009210033.ipynb

```
[ ] X=data.iloc[:, :-1]
print(X)
y=data.iloc[:, -1]
print(y)

outlook temp humidity windy
0 sunny hot high False
1 sunny hot high True
2 overcast hot high False
3 rainy mild high False
4 rainy cool normal False
5 rainy cool normal True
6 overcast cool normal True
7 sunny mild high False
8 sunny cool normal False
9 rainy mild normal False
10 sunny mild normal True
11 overcast mild high True
12 overcast hot normal False
13 rainy mild high True
0 no
1 no
2 yes
3 yes
4 yes
5 no
6 yes
7 no
8 yes
9 yes
10 yes
11 yes
12 yes
13 no
Name: play, dtype: object
```

0s completed at 10:43PM

Market Brief 22:45 ENG 09-03-2023



## Department of Computer Science and Engineering (Data Science)

PlayTennis\_60009210033.ipynb

```
[19] from sklearn import preprocessing
     string_to_int= preprocessing.LabelEncoder()           #encode your data
     data=data.apply(string_to_int.fit_transform) #fit and transform it
     data
```

	outlook	temp	humidity	windy	play
0	2	1	0	0	0
1	2	1	0	1	0
2	0	1	0	0	1
3	1	2	0	0	1
4	1	0	1	0	1
5	1	0	1	1	0
6	0	0	1	1	1
7	2	2	0	0	0
8	2	0	1	0	1
9	1	2	1	0	1
10	2	2	1	1	1
11	0	2	0	1	1
12	0	1	1	0	1
13	1	2	0	1	0

PlayTennis\_60009210033.ipynb

```
[20] class Node():
     def __init__(self, feature_index=None, threshold=None, left=None, right=None, info_gain=None, value=None):
         ... constructor ...
```

	outlook	temp	humidity	windy	play
0	2	1	0	0	0
1	2	1	0	1	0
2	0	1	0	0	1
3	1	2	0	0	1
4	1	0	1	0	1
5	1	0	1	1	0
6	0	0	1	1	1
7	2	2	0	0	0
8	2	0	1	0	1
9	1	2	1	0	1
10	2	2	1	1	1
11	0	2	0	1	1
12	0	1	1	0	1
13	1	2	0	1	0



## Department of Computer Science and Engineering (Data Science)

PlayTennis\_60009210033.ipynb

```
[20]: class Node():
    def __init__(self, feature_index=None, threshold=None, left=None, right=None, info_gain=None, value=None):
        ...
        # for decision node
        self.feature_index = feature_index
        self.threshold = threshold
        self.left = left
        self.right = right
        self.info_gain = info_gain

    # for leaf node
    self.value = value

[21]: class DecisionTreeClassifier():
    def __init__(self, min_samples_split=2, max_depth=2):
        ...
        # initialize the root of the tree
        self.root = None

        # stopping conditions
        self.min_samples_split = min_samples_split
        self.max_depth = max_depth

    def build_tree(self, dataset, curr_depth=0):
        ...
        # recursive function to build the tree ...

    X, Y = dataset[:, :-1], dataset[:, -1]
    num_samples, num_features = np.shape(X)
```

0s completed at 10:43PM

PlayTennis\_60009210033.ipynb

```
[21]: class DecisionTreeClassifier():
    def __init__(self, min_samples_split=2, max_depth=2):
        ...
        # initialize the root of the tree
        self.root = None

        # stopping conditions
        self.min_samples_split = min_samples_split
        self.max_depth = max_depth

    def build_tree(self, dataset, curr_depth=0):
        ...
        # recursive function to build the tree ...

    X, Y = dataset[:, :-1], dataset[:, -1]
    num_samples, num_features = np.shape(X)

    # split until stopping conditions are met
    if num_samples>=self.min_samples_split and curr_depth<=self.max_depth:
        # find the best split
        best_split = self.get_best_split(dataset, num_samples, num_features)
        # check if information gain is positive
        if best_split["info_gain"]>0:
            # recur left
            left_subtree = self.build_tree(best_split["dataset_left"], curr_depth+1)
            # recur right
            right_subtree = self.build_tree(best_split["dataset_right"], curr_depth+1)
            # return decision node
            return Node(best_split["feature_index"], best_split["threshold"],
                       left_subtree, right_subtree, best_split["info_gain"])

    # compute leaf node
```

0s completed at 10:43PM



## Department of Computer Science and Engineering (Data Science)

PlayTennis\_60009210033.ipynb · colab.research.google.com/drive/1FEZz\_zweKqHsh31lydgJlNKLFPxZhj#scrollTo=7s9ADAq9XLv\_

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

```
[21]   # compute leaf node
  leaf_value = self.calculate_leaf_value(Y)
  # return leaf node
  return Node(value=leaf_value)

def get_best_split(self, dataset, num_samples, num_features):
    ''' function to find the best split '''

    # dictionary to store the best split
    best_split = {}
    max_info_gain = -float("inf")

    # loop over all the features
    for feature_index in range(num_features):
        feature_values = dataset[:, feature_index]
        possible_thresholds = np.unique(feature_values)
        # loop over all the feature values present in the data
        for threshold in possible_thresholds:
            # get current split
            dataset_left, dataset_right = self.split(dataset, feature_index, threshold)
            # check if children are not null
            if len(dataset_left)>0 and len(dataset_right)>0:
                y, left_y, right_y = dataset[:, -1], dataset_left[:, -1], dataset_right[:, -1]
                # compute information gain
                curr_info_gain = self.information_gain(y, left_y, right_y, "gini")
                # update the best split if needed
                if curr_info_gain>max_info_gain:
                    best_split["feature_index"] = feature_index
                    best_split["threshold"] = threshold
                    best_split["dataset_left"] = dataset_left
                    best_split["dataset_right"] = dataset_right
                    best_split["info_gain"] = curr_info_gain
                    max_info_gain = curr_info_gain
    return best_split
```

0s completed at 10:43PM

Type here to search Market Brief 22:45 09-03-2023

PlayTennis\_60009210033.ipynb · colab.research.google.com/drive/1FEZz\_zweKqHsh31lydgJlNKLFPxZhj#scrollTo=7s9ADAq9XLv\_

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

```
[21]   if curr_info_gain>max_info_gain:
  best_split["feature_index"] = feature_index
  best_split["threshold"] = threshold
  best_split["dataset_left"] = dataset_left
  best_split["dataset_right"] = dataset_right
  best_split["info_gain"] = curr_info_gain
  max_info_gain = curr_info_gain

  # return best split
return best_split

def split(self, dataset, feature_index, threshold):
    ''' function to split the data '''

    dataset_left = np.array([row for row in dataset if row[feature_index]<=threshold])
    dataset_right = np.array([row for row in dataset if row[feature_index]>threshold])
    return dataset_left, dataset_right

def information_gain(self, parent, l_child, r_child, mode="entropy"):
    ''' function to compute information gain '''

    weight_l = len(l_child) / len(parent)
    weight_r = len(r_child) / len(parent)
    if mode=="gini":
        gain = self.gini_index(parent) - (weight_l*self.gini_index(l_child) + weight_r*self.gini_index(r_child))
    else:
        gain = self.entropy(parent) - (weight_l*self.entropy(l_child) + weight_r*self.entropy(r_child))
    return gain

def entropy(self, y):
    ''' function to compute entropy '''

    class_labels = np.unique(y)
```

0s completed at 10:43PM

Type here to search Market Brief 22:45 09-03-2023



## Department of Computer Science and Engineering (Data Science)

PlayTennis\_60009210033.ipynb

```
[21]: def entropy(self, y):
    """ function to compute entropy """

    class_labels = np.unique(y)
    entropy = 0
    for cls in class_labels:
        p_cls = len([y == cls]) / len(y)
        entropy += -p_cls * np.log2(p_cls)

    return entropy

def gini_index(self, y):
    """ function to compute gini index """

    class_labels = np.unique(y)
    gini = 0
    for cls in class_labels:
        p_cls = len([y == cls]) / len(y)
        gini += p_cls**2

    return 1 - gini

def calculate_leaf_value(self, y):
    """ function to compute leaf node """

    Y = list(y)
    return max(Y, key=Y.count)

def print_tree(self, tree=None, indent=" "):
    """ function to print the tree """

    if not tree:
        tree = self.root

    if tree.value is not None:
        print(tree.value)

    else:
        print("X." + str(tree.feature_index), "<", tree.threshold, "?", tree.info_gain)
        print("%sleft: %s" % (indent, end=""))
        self.print_tree(tree.left, indent + indent)
        print("%sright: %s" % (indent, end=""))
        self.print_tree(tree.right, indent + indent)

    if fit(self, X, Y):
        """ function to train the tree """

        dataset = np.concatenate((X, Y), axis=1)
        self.root = self.build_tree(dataset)

    def predict(self, X):
        """ function to predict new dataset """

        predictions = [self.make_prediction(x, self.root) for x in X]
        return predictions

    def make_prediction(self, x, tree):
        """ function to predict a single data point """

        if tree.value is None: return tree.value
        feature_val = x[tree.feature_index]
```

0s completed at 10:43PM

PlayTennis\_60009210033.ipynb

```
[21]: def print_tree(self, tree=None, indent=" "):
    """ function to print the tree """

    if not tree:
        tree = self.root

    if tree.value is not None:
        print(tree.value)

    else:
        print("X." + str(tree.feature_index), "<", tree.threshold, "?", tree.info_gain)
        print("%sleft: %s" % (indent, end=""))
        self.print_tree(tree.left, indent + indent)
        print("%sright: %s" % (indent, end=""))
        self.print_tree(tree.right, indent + indent)

    if fit(self, X, Y):
        """ function to train the tree """

        dataset = np.concatenate((X, Y), axis=1)
        self.root = self.build_tree(dataset)

    def predict(self, X):
        """ function to predict new dataset """

        predictions = [self.make_prediction(x, self.root) for x in X]
        return predictions

    def make_prediction(self, x, tree):
        """ function to predict a single data point """

        if tree.value is None: return tree.value
        feature_val = x[tree.feature_index]
```

0s completed at 10:43PM



## Department of Computer Science and Engineering (Data Science)

PlayTennis\_60009210033.ipynb · +

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

```
[21]: def make_prediction(self, x, tree):
        """ function to predict a single data point """

    [x]
        if tree.value!=None: return tree.value
        feature_val = x[tree.feature_index]
        if feature_val<tree.threshold:
            return self.make_prediction(x, tree.left)
        else:
            return self.make_prediction(x, tree.right)

[22]: X = data.iloc[:, :-1].values
Y = data.iloc[:, -1].values.reshape(-1,1)
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=.2, random_state=0)

[23]: classifier = DecisionTreeClassifier(min_samples_split=3, max_depth=3)
classifier.fit(X_train,Y_train)
classifier.print_tree()

X_0 <= 0 ? 0.1549586776859505
left:1
right:X_1 <= 1 ? 0.16875
left:0
right:X_2 <= 0 ? 0.2133333333333332
left:X_0 <= 1 ? 0.1111111111111111
    left:0
    right:0
right:1

Y_pred = classifier.predict(X_test)
from sklearn.metrics import accuracy_score
```

0s completed at 10:43PM

Windows taskbar: Type here to search, Market Brief, ENG 22:45 09-03-2023

PlayTennis\_60009210033.ipynb · +

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

```
[21]: def make_prediction(self, x, tree):
        """ function to predict a single data point """

    [x]
        if tree.value!=None: return tree.value
        feature_val = x[tree.feature_index]
        if feature_val<tree.threshold:
            return self.make_prediction(x, tree.left)
        else:
            return self.make_prediction(x, tree.right)

[22]: X = data.iloc[:, :-1].values
Y = data.iloc[:, -1].values.reshape(-1,1)
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=.2, random_state=0)

[23]: classifier = DecisionTreeClassifier(min_samples_split=3, max_depth=3)
classifier.fit(X_train,Y_train)
classifier.print_tree()

X_0 <= 0 ? 0.1549586776859505
left:1
right:X_1 <= 1 ? 0.16875
left:0
right:X_2 <= 0 ? 0.2133333333333332
left:X_0 <= 1 ? 0.1111111111111111
    left:0
    right:0
right:1

Y_pred = classifier.predict(X_test)
from sklearn.metrics import accuracy_score
accuracy_score(Y_test, Y_pred)

0.3333333333333333
```

0s completed at 10:43PM

Windows taskbar: Type here to search, Market Brief, ENG 22:45 09-03-2023

[https://colab.research.google.com/drive/1FEZz\\_zweKqHsh31ItydgJINKLFPyXZhj?usp=sharing](https://colab.research.google.com/drive/1FEZz_zweKqHsh31ItydgJINKLFPyXZhj?usp=sharing)



## Department of Computer Science and Engineering (Data Science)

2. Use python libraries to build a decision tree classifier on Dataset 2. Analyze the results using confusion matrix and accuracy.

The screenshot shows two sessions in Google Colab. Both sessions are titled 'Iris\_60009210033\_Jhanvi Parekh.ipynb'.

**Session 1:**

```
[1] import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

#for encoding
from sklearn.preprocessing import LabelEncoder

#for train test splitting
from sklearn.model_selection import train_test_split

#for decision tree object
from sklearn.tree import DecisionTreeClassifier

#for checking testing results
from sklearn.metrics import classification_report, confusion_matrix

#for visualizing tree
from sklearn.tree import plot_tree

[4] df = pd.read_csv('Iris.csv', index_col = 0)
df.head()
```

Output:

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
1	5.1	3.5	1.4	0.2	Iris-setosa
2	4.9	3.0	1.4	0.2	Iris-setosa
3	4.7	3.2	1.3	0.2	Iris-setosa

**Session 2:**

```
[1] df = pd.read_csv('Iris.csv', index_col = 0)
df.head()
```

Output:

	SepallengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
1	5.1	3.5	1.4	0.2	Iris-setosa
2	4.9	3.0	1.4	0.2	Iris-setosa
3	4.7	3.2	1.3	0.2	Iris-setosa
4	4.6	3.1	1.5	0.2	Iris-setosa
5	5.0	3.6	1.4	0.2	Iris-setosa

```
[5] #getting information of dataset
df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 150 entries, 1 to 150
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   SepallengthCm 150 non-null   float64 
 1   SepalWidthCm  150 non-null   float64 
 2   PetallengthCm 150 non-null   float64 
 3   PetalWidthCm  150 non-null   float64 
 4   Species       150 non-null   object  
dtypes: float64(4), object(1)
memory usage: 7.0+ KB
```

```
[6] df.isnull().any()
```



## Department of Computer Science and Engineering (Data Science)

Iris\_60009210033\_Jhanvi Parekh.ipynb - Colaboratory | Document from Jhanvi Parekh | +

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

[6] df.isnull().any()

```
{x} SepalLengthCm False
SepalWidthCm False
PetalLengthCm False
PetalWidthCm False
Species False
dtype: bool
```

df.shape  
(150, 5)

[8] target = df['Species']
df1 = df.copy()
df1 = df1.drop('Species', axis=1)
df1.shape  
(150, 4)

[9] df1.head()

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
Id				
1	5.1	3.5	1.4	0.2
2	4.9	3.0	1.4	0.2
3	4.7	3.2	1.3	0.2
4	4.6	3.1	1.5	0.2

Type here to search 33°C Smoke ENG 09-03-2023

Iris\_60009210033\_Jhanvi Parekh.ipynb - Colaboratory | Document from Jhanvi Parekh | +

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

[9] df1.head()

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
Id				
1	5.1	3.5	1.4	0.2
2	4.9	3.0	1.4	0.2
3	4.7	3.2	1.3	0.2
4	4.6	3.1	1.5	0.2
5	5.0	3.6	1.4	0.2

[10] # Defining the attributes
x = df1

[11] target

Id	Species
1	Iris-setosa
2	Iris-setosa
3	Iris-setosa
4	Iris-setosa
5	Iris-setosa
...	
146	Iris-virginica
147	Iris-virginica
148	Iris-virginica
149	Iris-virginica
150	Iris-virginica

Name: Species, Length: 150, dtype: object

Type here to search 33°C Smoke ENG 09-03-2023



## Department of Computer Science and Engineering (Data Science)

```
Irish_60009210033_Jhanvi Parekh.ipynb - Colaboratory | Document from Jhanvi Parekh | +  
Irish_60009210033_Jhanvi Parekh.ipynb | Comment | Share | Guest | RAM | Disk |  
File Edit View Insert Runtime Tools Help All changes saved  
+ Code + Text  
target  
Id  
{x}  
1 Iris-setosa  
2 Iris-setosa  
3 Iris-setosa  
4 Iris-setosa  
5 Iris-setosa  
...  
146 Iris-virginica  
147 Iris-virginica  
148 Iris-virginica  
149 Iris-virginica  
150 Iris-virginica  
Name: Species, Length: 150, dtype: object  
[12] #label encoding  
le = LabelEncoder()  
target = le.fit_transform(target)  
target  
array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,  
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,  
0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,  
2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2])  
y = target  
[14] # Splitting the data - 80:20 ratio  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 0)  
Type here to search 33°C Smoke ENG 09-03-2023  
Irish_60009210033_Jhanvi Parekh.ipynb - Colaboratory | Document from Jhanvi Parekh | +  
Irish_60009210033_Jhanvi Parekh.ipynb | Comment | Share | Guest | RAM | Disk |  
File Edit View Insert Runtime Tools Help All changes saved  
+ Code + Text  
y = target  
{x}  
# Splitting the data - 80:20 ratio  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 0)  
print("Training split input-", X_train.shape)  
print("Testing split input-", X_test.shape)  
Training split input- (105, 4)  
Testing split input- (45, 4)  
Building a tree  
[15] # Defining the decision tree algorithm  
dtree=DecisionTreeClassifier(random_state=0,max_depth = 5)  
dtree.fit(X_train,y_train)  
print('Decision Tree Classifier Created')  
Decision Tree Classifier Created  
Testing  
[16] # Predicting the values of test data  
y_pred = dtree.predict(X_test)  
print("Classification report - \n", classification_report(y_test,y_pred))  
Classification report -  
Type here to search 33°C Smoke ENG 09-03-2023
```



## Department of Computer Science and Engineering (Data Science)

Iris\_60009210033\_Jhanvi Parekh.ipynb - Colaboratory | Document from Jhanvi Parekh | +

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

Testing

```
[x] [16] # Predicting the values of test data
y_pred = dtree.predict(X_test)
print("Classification report - \n", classification_report(y_test,y_pred))

Classification report -
precision    recall    f1-score   support
          0       1.00     1.00      1.00      16
          1       0.94     0.97      0.96      18
          2       0.92     1.00      0.96      11

accuracy                           0.98      45
macro avg       0.97     0.98      0.98      45
weighted avg    0.98     0.98      0.98      45
```

```
[x] [17] cm = confusion_matrix(y_test, y_pred)

[18] cm
array([[16,  0,  0],
       [ 0, 17,  1],
       [ 0,  0, 11]])
```

# Visualising the graph without the use of graphviz

```
plt.figure(figsize = (20,20))
dec_tree = plot_tree(decision_tree=dtree, feature_names = df1.columns,
                     class_names =[ "setosa", "vericolor", "verginica"] , filled = True , precision = 4, rounded = True)
```

Type here to search

33°C Smoke ENG 09-03-2023

Iris\_60009210033\_Jhanvi Parekh.ipynb - Colaboratory | Document from Jhanvi Parekh | +

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

```
[x] [16] array([[16,  0,  0],
       [ 0, 17,  1],
       [ 0,  0, 11]])

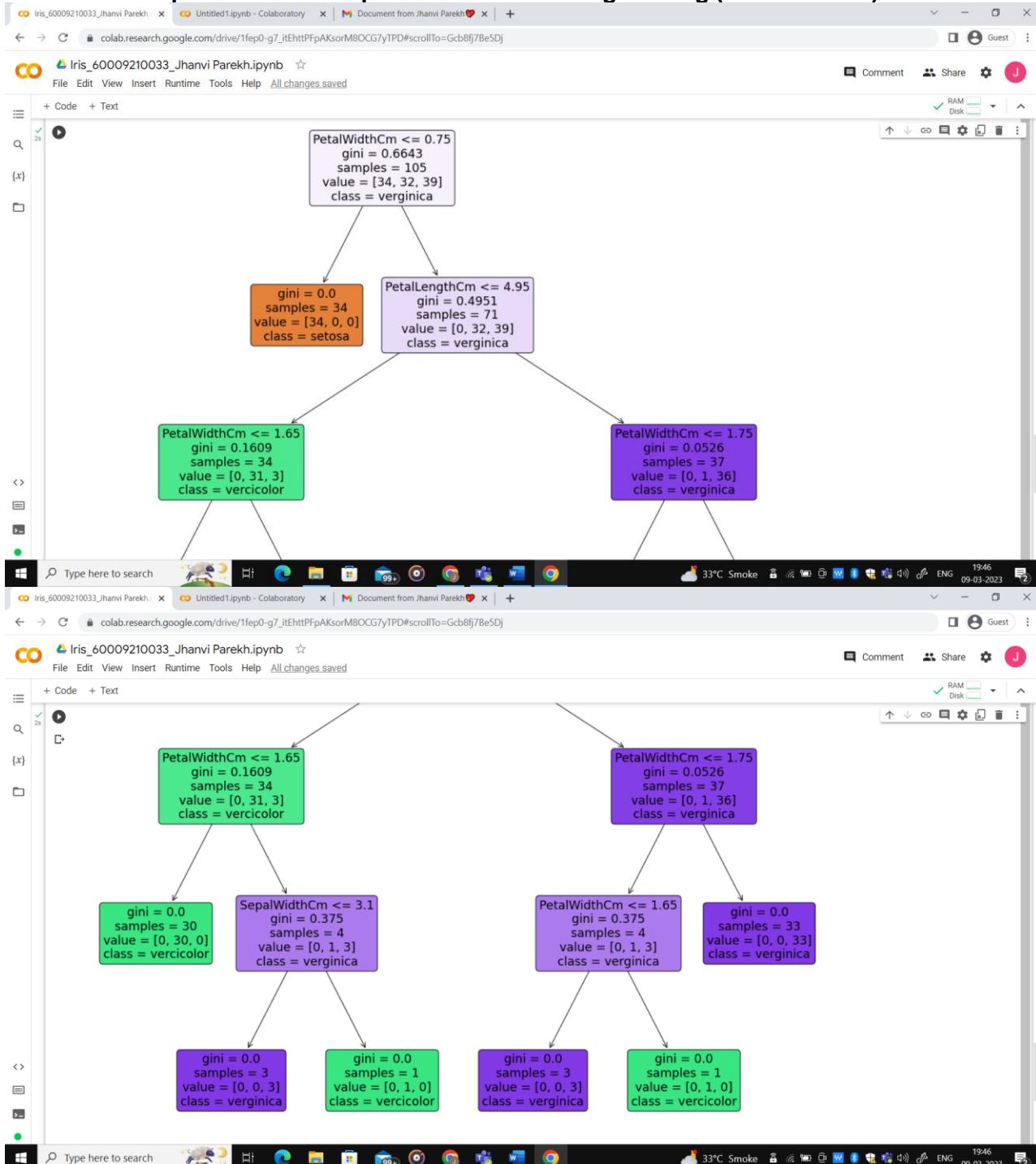
[26] # Visualising the graph without the use of graphviz

plt.figure(figsize = (20,20))
dec_tree = plot_tree(decision_tree=dtree, feature_names = df1.columns,
                     class_names =[ "setosa", "vericolor", "verginica"] , filled = True , precision = 4, rounded = True)

plt.savefig("one.png")
```



## Department of Computer Science and Engineering (Data Science)



[https://colab.research.google.com/drive/1fep0-g7\\_itEhttPFpAKsorM8OCG7yTPD?usp=sharing](https://colab.research.google.com/drive/1fep0-g7_itEhttPFpAKsorM8OCG7yTPD?usp=sharing)

3. Write a code to show overfitting in the decision tree classifier built using Dataset 3. Use sklearn and matplotlib.



## Department of Computer Science and Engineering (Data Science)

M (no subject) - parekhjhanvi2327 x | C breast cancer overfitting.ipynb - +

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text ⌂ Copy to Drive

RAM Disk

Share

J

Search icon

[1] `import numpy as np  
import pandas as pd  
df=pd.read_csv('/content/Breast_cancer_data.csv')  
df.head()`

Output:

	mean_radius	mean_texture	mean_perimeter	mean_area	mean_smoothness	diagnosis
0	17.99	10.38	122.80	1001.0	0.11840	0
1	20.57	17.77	132.90	1326.0	0.08474	0
2	19.69	21.25	130.00	1203.0	0.10960	0
3	11.42	20.38	77.58	386.1	0.14250	0
4	20.29	14.34	135.10	1297.0	0.10030	0

[2] `from sklearn.datasets import make_classification  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.model_selection import train_test_split  
from sklearn.metrics import accuracy_score  
x=df.iloc[:, :-1]  
print(x)  
y=df.iloc[:, -1]  
print(y)`

Output:

	mean_radius	mean_texture	mean_perimeter	mean_area	mean_smoothness
0	17.99	10.38	122.80	1001.0	0.11840
1	20.57	17.77	132.90	1326.0	0.08474
2	19.69	21.25	130.00	1203.0	0.10960
3	11.42	20.38	77.58	386.1	0.14250
4	20.29	14.34	135.10	1297.0	0.10030

0s completed at 8:15PM

Windows taskbar: Type here to search, File Explorer, Task View, Edge, Google Chrome, File Manager, 99+, 31°C Smoke, ENG 20:16 13-03-2023

M (no subject) - parekhjhanvi2327 x | C breast cancer overfitting.ipynb - +

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text ⌂ Copy to Drive

RAM Disk

Share

J

Search icon

[3] `from sklearn.datasets import make_classification  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.model_selection import train_test_split  
from sklearn.metrics import accuracy_score  
x=df.iloc[:, :-1]  
print(x)  
y=df.iloc[:, -1]  
print(y)`

Output:

	mean_radius	mean_texture	mean_perimeter	mean_area	mean_smoothness
0	17.99	10.38	122.80	1001.0	0.11840
1	20.57	17.77	132.90	1326.0	0.08474
2	19.69	21.25	130.00	1203.0	0.10960
3	11.42	20.38	77.58	386.1	0.14250
4	20.29	14.34	135.10	1297.0	0.10030
..	...	...	...	...	...
564	21.56	22.39	142.00	1479.0	0.11180
565	20.13	28.25	131.20	1261.0	0.09780
566	16.60	28.08	108.30	858.1	0.08455
567	20.60	29.33	140.10	1265.0	0.11780
568	7.76	24.54	47.92	181.0	0.05263

[569 rows x 5 columns]

0 0  
1 0  
2 0  
3 0  
4 0  
..  
564 0  
565 0  
566 0  
567 0  
568 1

0s completed at 8:15PM

Windows taskbar: Type here to search, File Explorer, Task View, Edge, Google Chrome, File Manager, 99+, INFY -2.29%, ENG 20:17 13-03-2023



## Department of Computer Science and Engineering (Data Science)

M (no subject) - parekhjhanvi2327.ipynb breast cancer overfitting.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text ⚙ Copy to Drive

[2] 1 20.27 17.77 132.90 1320.0 0.08474  
2 19.69 21.25 130.00 1203.0 0.10960  
3 11.42 20.38 77.58 386.1 0.14250  
4 20.29 14.34 135.10 1297.0 0.10030  
...  
564 21.56 22.39 142.00 1479.0 0.11100  
565 20.13 28.25 131.20 1261.0 0.09780  
566 16.60 28.08 108.30 858.1 0.08455  
567 20.60 29.33 140.10 1265.0 0.11780  
568 7.76 24.54 47.92 181.0 0.05263

[569 rows x 5 columns]

0 0  
1 0  
2 0  
3 0  
4 0  
..  
564 0  
565 0  
566 0  
567 0  
568 1

Name: diagnosis, Length: 569, dtype: int64

[3] x\_train, x\_test, y\_train, y\_test = train\_test\_split(x, y, test\_size=0.3, random\_state=0)  
print("Training split input-", x\_train.shape)  
print("Testing split input-", x\_test.shape)

Training split input- (398, 5)  
Testing split input- (171, 5)

[4] dtc = DecisionTreeClassifier(criterion='entropy', max\_depth=5, random\_state=0)  
dtc.fit(x\_train, y\_train)

0s completed at 8:15PM

Type here to search

INFY -2.29% ENG 20:17 13-03-2023

M (no subject) - parekhjhanvi2327.ipynb breast cancer overfitting.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text ⚙ Copy to Drive

[4] dtc = DecisionTreeClassifier(criterion='entropy', max\_depth=5, random\_state=0)  
dtc.fit(x\_train, y\_train)  
print('Decision Tree Classifier Created')

Decision Tree Classifier Created

[5] from sklearn.metrics import classification\_report  
y\_pred = dtc.predict(x\_test)  
print("Classification report - \n", classification\_report(y\_test,y\_pred))

	precision	recall	f1-score	support
0	0.76	0.87	0.81	63
1	0.92	0.84	0.88	108
accuracy			0.85	171
macro avg	0.84	0.86	0.85	171
weighted avg	0.86	0.85	0.86	171

[6] from sklearn.metrics import classification\_report  
y\_pred = dtc.predict(x\_test)  
print("Train data accuracy:", accuracy\_score(y\_true = y\_train, y\_pred = dtc.predict(x\_train)))  
print("Test data accuracy:", accuracy\_score(y\_true = y\_test, y\_pred = y\_pred))  
print('Accuracy:', accuracy\_score(y\_test, y\_pred))

Train data accuracy: 0.9597989949748744  
Test data accuracy: 0.8538011695906432  
Accuracy: 0.8538011695906432

0s completed at 8:15PM

Type here to search

INFY -2.29% ENG 20:17 13-03-2023



## Department of Computer Science and Engineering (Data Science)

M (no subject) - parekhjhanvi2327 x | C breast cancer overfitting.ipynb - +

File Edit View Insert Runtime Tools Help Cannot save changes

[6] `from sklearn.metrics import classification_report  
y_pred = dtc.predict(x_test)  
print("Train data accuracy:",accuracy_score(y_true = y_train, y_pred = dtc.predict(x_train)))  
print("Test data accuracy:",accuracy_score(y_true = y_test, y_pred = y_pred))  
print('Accuracy:', accuracy_score(y_test, y_pred))`

Train data accuracy: 0.9597989949748744  
Test data accuracy: 0.853801169590642  
Accuracy: 0.853801169590642

[7] `import matplotlib.pyplot as plt  
plt.figure(figsize=(15,8))  
from sklearn.tree import plot_tree  
plot_tree(dtc)  
plt.show()`

The decision tree diagram for Model 1 has 1 root node and 10 leaf nodes. The root node splits on  $Age \geq 43.05$  with samples = 988 and value = 1246 (4). It has two children: one splitting on  $Age \leq 43.05$  with samples = 200 and value = 100 (0), and another splitting on  $Age > 43.05$  with samples = 988 and value = 1146 (4). The left child further splits on  $Age \leq 42.25$  with samples = 1274 and value = 631 (1). The right child splits on  $Age \geq 44.04$  with samples = 211 and value = 126 (1). The right child of the root node splits on  $Age \leq 44.04$  with samples = 14,995 and value = 11,329 (1). This node further splits on  $Age \leq 44.04$  with samples = 11,329 and value = 11,329 (1). The right child of this node splits on  $Age \geq 44.04$  with samples = 3,666 and value = 2,367 (1). The left child of the root node splits on  $Age \leq 42.25$  with samples = 1274 and value = 631 (1). This node further splits on  $Age \leq 42.25$  with samples = 1,416 and value = 1,416 (1). The right child of this node splits on  $Age \geq 42.25$  with samples = 1,258 and value = 415 (0).

INFY -2.29% 20:17 ENG 13-03-2023

0s completed at 8:15PM

M (no subject) - parekhjhanvi2327 x | C breast cancer overfitting.ipynb - +

File Edit View Insert Runtime Tools Help Cannot save changes

[6] `import matplotlib.pyplot as plt  
plt.figure(figsize=(15,8))  
from sklearn.tree import plot_tree  
plot_tree(dtc)  
plt.show()`

The decision tree diagram for Model 2 has 1 root node and 10 leaf nodes. The root node splits on  $Age \geq 43.05$  with samples = 988 and value = 1246 (4). It has two children: one splitting on  $Age \leq 43.05$  with samples = 200 and value = 100 (0), and another splitting on  $Age > 43.05$  with samples = 988 and value = 1146 (4). The left child further splits on  $Age \leq 42.25$  with samples = 1274 and value = 631 (1). The right child of the root node splits on  $Age \geq 44.04$  with samples = 211 and value = 126 (1). The right child of the root node further splits on  $Age \geq 44.04$  with samples = 11,329 and value = 11,329 (1). This node further splits on  $Age \leq 44.04$  with samples = 3,666 and value = 2,367 (1). The right child of this node splits on  $Age \geq 44.04$  with samples = 3,666 and value = 2,367 (1). The left child of the root node splits on  $Age \leq 42.25$  with samples = 1274 and value = 631 (1). This node further splits on  $Age \leq 42.25$  with samples = 1,416 and value = 1,416 (1). The right child of this node splits on  $Age \geq 42.25$  with samples = 1,258 and value = 415 (0).

INFY -2.29% 20:17 ENG 13-03-2023

0s completed at 8:15PM



## Department of Computer Science and Engineering (Data Science)

M (no subject) - parekhjhanvi2327 x | breast cancer overfitting.ipynb - x + colab.research.google.com/drive/1lt4w7XxKy9a4MT1Hk58JQ-b6xRB0q0M5#scrollTo=maxcvf0vcIW8 Guest

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text ⌂ Copy to Drive Share RAM Disk

```
from matplotlib import pyplot
train_scores, test_scores = list(), list()
values = [i for i in range(1, 16)]
for i in values:
    model = DecisionTreeClassifier(max_depth=i)
    model.fit(x_train, y_train)
    train_yhat = model.predict(x_train)
    train_acc = accuracy_score(y_train, train_yhat)
    train_scores.append(train_acc)
    test_yhat = model.predict(x_test)
    test_acc = accuracy_score(y_test, test_yhat)
    test_scores.append(test_acc)
    print('Xd, train: {:.3f}, test: {:.3f}'.format(i, train_acc, test_acc))
pyplot.plot(values, train_scores, '-o', label='Train')
pyplot.plot(values, test_scores, '-o', label='Test')
pyplot.legend()
pyplot.show()
```

>1, train: 0.887, test: 0.901  
>2, train: 0.887, test: 0.901  
>3, train: 0.937, test: 0.895  
>4, train: 0.952, test: 0.924  
>5, train: 0.955, test: 0.918  
>6, train: 0.975, test: 0.918  
>7, train: 0.985, test: 0.918  
>8, train: 0.997, test: 0.912  
>9, train: 0.997, test: 0.918  
>10, train: 1.000, test: 0.930  
>11, train: 1.000, test: 0.918  
>12, train: 1.000, test: 0.912  
>13, train: 1.000, test: 0.924  
>14, train: 1.000, test: 0.912  
>15, train: 1.000, test: 0.912

0s completed at 8:15PM

Type here to search

M (no subject) - parekhjhanvi2327 x | breast cancer overfitting.ipynb - x + colab.research.google.com/drive/1lt4w7XxKy9a4MT1Hk58JQ-b6xRB0q0M5#scrollTo=maxcvf0vcIW8 Guest

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text ⌂ Copy to Drive Share RAM Disk

```
pyplot.legend()
pyplot.show()
```

>1, train: 0.887, test: 0.901  
>2, train: 0.887, test: 0.901  
>3, train: 0.937, test: 0.895  
>4, train: 0.952, test: 0.924  
>5, train: 0.955, test: 0.918  
>6, train: 0.975, test: 0.918  
>7, train: 0.985, test: 0.918  
>8, train: 0.997, test: 0.912  
>9, train: 0.997, test: 0.918  
>10, train: 1.000, test: 0.930  
>11, train: 1.000, test: 0.918  
>12, train: 1.000, test: 0.912  
>13, train: 1.000, test: 0.924  
>14, train: 1.000, test: 0.912  
>15, train: 1.000, test: 0.912

0s completed at 8:15PM

Type here to search

M (no subject) - parekhjhanvi2327 x | breast cancer overfitting.ipynb - x + colab.research.google.com/drive/1lt4w7XxKy9a4MT1Hk58JQ-b6xRB0q0M5#scrollTo=maxcvf0vcIW8 Guest

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text ⌂ Copy to Drive Share RAM Disk

<https://colab.research.google.com/drive/1lt4w7XxKy9a4MT1Hk58JQ-b6xRB0q0M5?usp=sharing>

4. Implement Decision tree regressor on Dataset 4.



## Department of Computer Science and Engineering (Data Science)

Inbox (12) - parekhjhanvi2327@... | carprediction\_60009210033.ipynb | +

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

import pandas as pd  
from sklearn.tree import DecisionTreeRegressor  
from sklearn.model\_selection import train\_test\_split  
from sklearn.metrics import mean\_squared\_error

df = pd.read\_csv('/content/carprediction.csv')  
df.head()

	Make	Model	Year	Engine Fuel Type	Engine HP	Engine Cylinders	Transmission Type	Driven_Wheels	Number of Doors	Market Category	Vehicle Size	Vehicle Style	highway MPG	city mpg	Popularity	MSRP
0	BMW	Series M	2011	premium unleaded (required)	335.0	6.0	MANUAL	rear wheel drive	2.0	Tuner,Luxury,High-Performance	Compact	Coupe	26	19	3916	46135
1	BMW	Series 1	2011	premium unleaded (required)	300.0	6.0	MANUAL	rear wheel drive	2.0	Luxury,Performance	Compact	Convertible	28	19	3916	40650
2	BMW	Series 1	2011	premium unleaded (required)	300.0	6.0	MANUAL	rear wheel drive	2.0	Luxury,High-Performance	Compact	Coupe	28	20	3916	36350
3	BMW	Series 1	2011	premium unleaded (required)	230.0	6.0	MANUAL	rear wheel drive	2.0	Luxury,Performance	Compact	Coupe	28	18	3916	29450
4	BMW	Series 1	2011	premium unleaded (required)	230.0	6.0	MANUAL	rear wheel drive	2.0	Luxury	Compact	Convertible	28	18	3916	34500

Type here to search

Inbox (12) - parekhjhanvi2327@... | carprediction\_60009210033.ipynb | +

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

import numpy as np  
import matplotlib.pyplot as plt  
X = df['Popularity'].values  
y = df['MSRP'].values

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size = 5)

regressor = DecisionTreeRegressor()  
regressor.fit(X\_train.reshape(-1,1), y\_train.reshape(-1,1))

y\_pred = regressor.predict(X\_test.reshape(-1,1))  
y\_pred

df = pd.DataFrame({'Real Values':y\_test.reshape(-1), 'Predicted Values':y\_pred.reshape(-1)})

X\_grid = np.arange(min(X), max(X), 0.01)  
X\_grid = X\_grid.reshape(len(X\_grid), 1)  
plt.scatter(X\_test, y\_test, color = 'red')  
plt.scatter(X\_test, y\_pred, color = 'green')  
plt.title('Decision Tree Regression')  
plt.xlabel('Popularity')  
plt.ylabel('MSRP')  
plt.show()

plt.plot(X\_grid, regressor.predict(X\_grid), color = 'black')  
plt.title('Decision Tree Regression')  
plt.xlabel('Popularity')  
plt.ylabel('MSRP')  
plt.show()

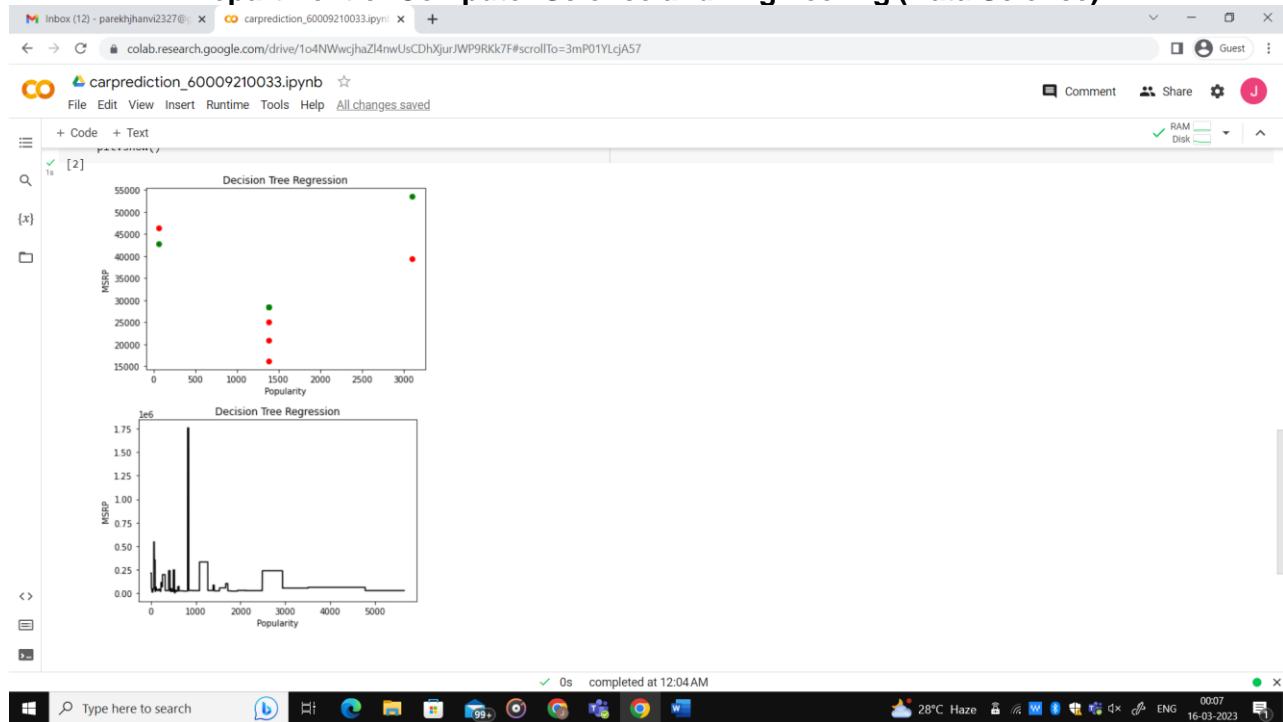
Type here to search

0s completed at 12:04AM

28°C Haze ENG 00:07 16-03-2023



## Department of Computer Science and Engineering (Data Science)



<https://colab.research.google.com/drive/1o4NWwcjhaZl4nwUsCDhXjurJWP9RKk7F?usp=sharing>