

# Decision Tree Method-WORKSHOP

## PYTHON

# Introduction to Decision Tree:Recap

---

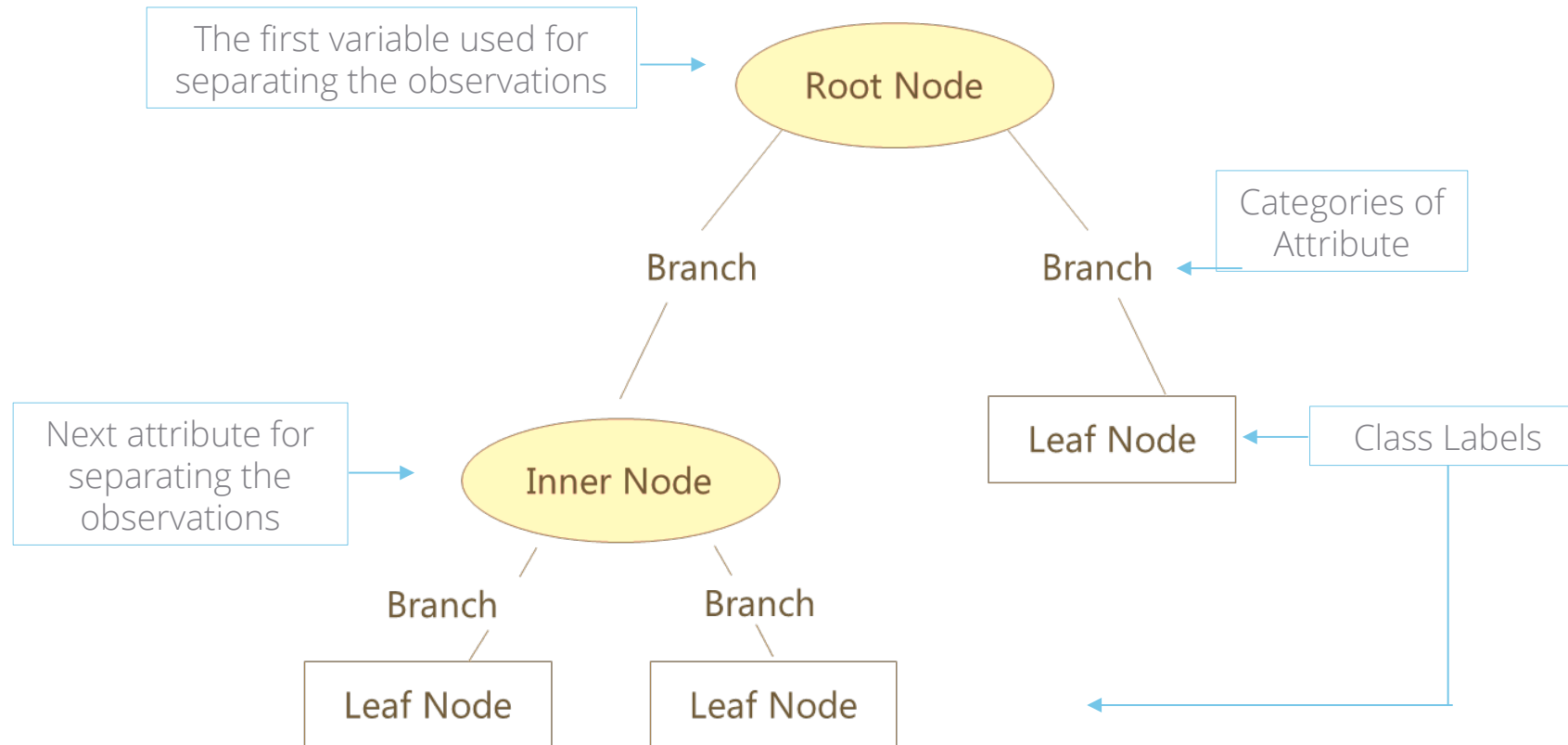
- One of the most robust predictive modeling techniques, **Decision Tree** uses data mining techniques for model building.
- Decision Tree **breaks down** a data set into smaller subsets and presents association **between target variable(dependent)** and independent variables as a tree structure.
- Final result is a **tree with Decision Nodes and Leaf Nodes**.
- A decision node has two or more branches and leaf node represents a classification or decision.

# Decision Tree – Basic Components

---

Component	Description	Alternate terms
Root node	Has no incoming edges and zero or more outgoing edges	Parent node
Internal nodes	Each has exactly one incoming edge and two or more outgoing edges	Decision nodes / Child nodes
Leaf node	Each has exactly one incoming edges and no outgoing edges	Terminal nodes
Branches	Categories of attributes	Edges

# Decision Tree – Basic Components



Class labels show observations belong to which class. The leaf node also shows Number of observations and Error rate (Actual classification vs classification given by the tree)



# Entropy

- **Entropy** measures the homogeneity of a sample. It is used as a parameter for checking the amount of uncertainty associated with a set of probabilities.
- **Entropy lies between 0 and 1**  
If the sample is completely homogeneous the entropy is 0 and if the sample is equally divided it has entropy of 1
- Entropy can be of two types, **for each category and at the variable level**
- **Entropy of a category** is calculated as:

$$-P1 * \log_2(P1) - P2 * \log_2(P2)$$

where,

P1 is the proportion of class 1

P2 is the proportion of class 2

# Entropy of a Category

Let us consider survey data from three cities depicting shopper's preferred brand

City	Brand A Voters	Brand B Voters	Number of Voters	% of votes for Brand A	% of votes for Brand B
Delhi	90	310	400	22.5%	77.5%
Chennai	10	90	100	10%	90%
Mumbai	100	100	200	50%	50%

Entropy for each city is calculated as:

Delhi:

Chennai:

Mumbai:

# Entropy at the Variable Level

---

- Entropy at the variable level can be derived by adding weighted averages of all category level entropy values
- Weights are the proportion of respondents in each category(here in each city)  
In the example under consideration,

Weights for the categories are

Delhi:

Chennai:

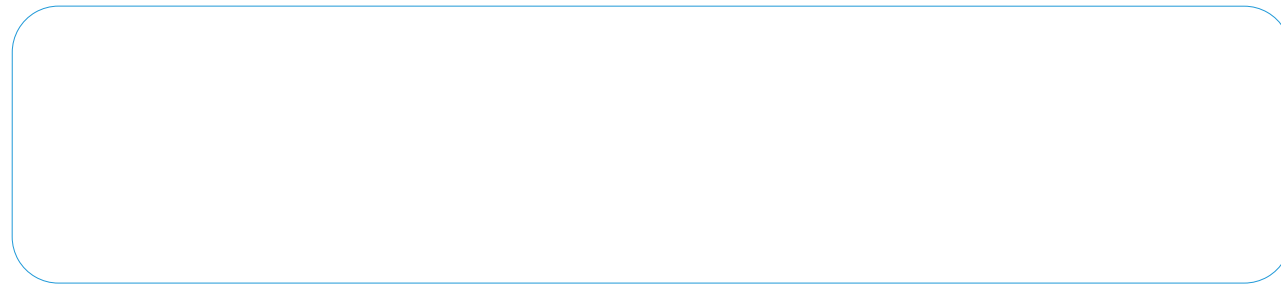
Mumbai:

Entropy at the variable level is

# Information Gain

---

- **Information Gain** is based on the decrease in entropy after a dataset is split on an attribute
- Constructing a decision tree is about finding attribute that returns the highest information gain



- Information gain can be interpreted as ability of reducing the uncertainty (Entropy) and hence increase predictability



# Information Gain

City	Brand A Voters	Brand B Voters	Number of Voters	% of votes for Brand A	% of votes for Brand B
Delhi	90	310	400	22.5%	77.5%
Chennai	10	90	100	10%	90%
Mumbai	100	100	200	50%	50%

Entropy for complete sample is calculated as follows:

$P1 = (\text{Total Brand A Voters} / \text{Total Voters})$

$P2 = (\text{Total Brand B Voters} / \text{Total Voters})$

Information Gain

Entropy at the variable level (Weighted average)



# Information Gain...

---

- Information Gain value is used to determine which attribute is the “best” – the attribute with most information gain is chosen
- Information gain for a variable is high when that variable has the low entropy at the variable level (Weighted average)
- Low entropy for a variable implies the classification based on that attribute is fairly homogenous , hence this attribute is selected as the first best attribute
- The same process is repeated till all attributes are used as split variables

# Case Study – Predicting Loan Defaulter

## Background

- The bank possesses demographic and transactional data of its loan customers. If the bank has a robust model to predict defaulters it can undertake better resource allocation.

## Objective

- To predict whether the customer applying for the loan will be a defaulter

## Available Information

- Sample size is 700
- Age group, Years at current address, Years at current employer, Debt to Income Ratio, Credit Card Debts, Other Debts are the independent variables
- **Defaulter** (=1 if defaulter, 0 otherwise) is the dependent variable





## Independent Variables

### Dependent Variable



**DATA SCIENCE**  
INSTITUTE

# Classification Tree in Python

---

# Importing the Libraries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier,
DecisionTreeRegressor, plot_tree

from sklearn.metrics import confusion_matrix, precision_score,
recall_score, accuracy_score, roc_curve, roc_auc_score
```

- ❑ **sklearn.tree** module includes Decision Tree – based models for classification and regression



# Classification Tree in Python

```
# Importing and Readyng the Data for Modeling
```

```
bankloan = pd.read_csv("BANK LOAN.csv")
```

```
bankloan1 = bankloan.drop(['SN'], axis = 1)
```

```
bankloan1['AGE'] = bankloan1['AGE'].astype('cat')
```

```
bankloan2 = pd.get_dummies(bankloan1)  
bankloan2.head()
```

**drop()** is used to remove unwanted variables.

**pd.get\_dummies()** converts categorical variables into dummy variables. Since AGE is a categorical variable, it is converted into three dummy variables: AGE\_1, AGE\_2, and AGE\_3.

```
# Output
```

	EMPLOY	ADDRESS	DEBTINC	CREDDEBT	OTHDEBT	DEFAULTER	AGE_1	AGE_2	AGE_3
0	17	12	9.3	11.36	5.01	1	0	0	1
1	10	6	17.3	1.36	4.00	0	1	0	0
2	15	14	5.5	0.86	2.17	0	0	1	0
3	15	14	2.9	2.66	0.82	0	0	0	1
4	2	0	17.3	1.79	3.06	1	1	0	0



# Classification Tree Using Information Gain

# Creating Data Partitions

```
X = bankloan2.loc[:,bankloan2.columns != 'DEFAULTER']  
y = bankloan2.loc[:, 'DEFAULTER']  
  
X_train, X_test, y_train, y_test = train_test_split(X, y,  
                                                    test_size = 0.30,  
                                                    random_state = 999)
```

- ❑ **train\_test\_split()** from `sklearn.model_selection` is used to split dataset into random train and test sets.
- ❑ **test\_size** represents the proportion of dataset to be included in the test set.
- ❑ **random\_state** sets the seed for the random number generator.

# Classification Tree Using Information Gain

# Classification Tree Using Information Gain

```
dtcl = DecisionTreeClassifier(criterion='entropy',  
                             min_samples_split= int(len(X_train)*.10))  
dtcl.fit(X_train, y_train)
```

- ❑ **DecisionTreeClassifier()** from sklearn.tree fits a classification tree.
- ❑ **criterion= 'entropy'** specifies the function to measure the split. Default is 'gini' for Gini impurity. 'entropy' stands for information gain.
- ❑ **min\_samples\_split=** minimum number of samples required to split an internal node. This number is set to be 10% of the sample size.
- ❑ The output displays model specifications.





```
from sklearn.tree import plot_tree
plt.figure(figsize = (16,10))
plot_tree(dtc1, filled = True, feature_names = list(X.columns))
plt.show();
```

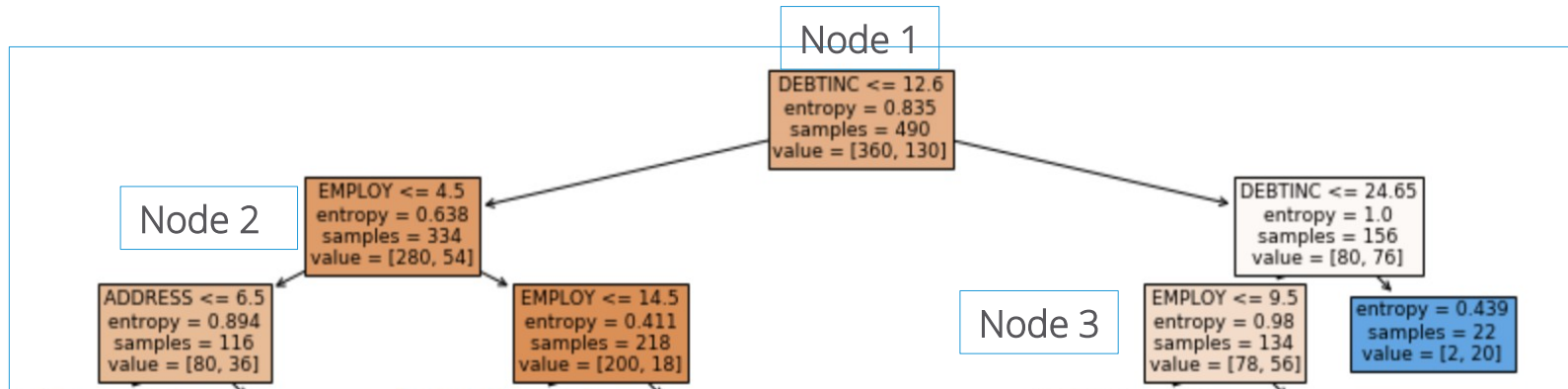
□ **plot\_tree** is used to plot the decision

**feature\_names** is used to mention feature names.

- 
- ```

graph TD
    Root["CREDDEBT <= 3.2  
entropy = 0.607  
samples = 121  
value = [103, 18]"]
    Root --> L1["ADDRESS <= 1.5  
entropy = 0.57  
samples = 119  
value = [103, 16]"]
    Root --> R1["CREDDEBT <= 0.235  
entropy = 0.949  
samples = 57  
value = [21, 36]"]
    
    L1 --> L2["ADDRESS <= 5.5  
entropy = 0.469  
samples = 100  
value = [90, 10]"]
    L1 --> L3["entropy = 0.0  
samples = 2  
value = [0, 2]"]
    
    L2 --> L4["entropy = 0.0  
samples = 29  
value = [29, 0]"]
    L2 --> L5["CREDDEBT <= 0.395  
entropy = 0.586  
samples = 71  
value = [61, 10]"]
    
    L5 --> L6["entropy = 0.991  
samples = 9  
value = [5, 4]"]
    L5 --> L7["OTHDEBT <= 0.77  
entropy = 0.459  
samples = 62  
value = [56, 6]"]
    
    L7 --> L8["entropy = 0.918  
samples = 3  
value = [1, 2]"]
    L7 --> L9["DEBTINC <= 6.05  
entropy = 0.358  
samples = 59  
value = [55, 4]"]
    
    L9 --> L10["entropy = 0.0  
samples = 18  
value = [18, 0]"]
    L9 --> L11["entropy = 0.461  
samples = 41  
value = [37, 4]"]
    
    R1 --> R2["entropy = 0.894  
samples = 29  
value = [20, 9]"]
    R1 --> R3["OTHDEBT <= 4.435  
entropy = 0.905  
samples = 53  
value = [17, 36]"]
    
    R3 --> R4["entropy = 0.962  
samples = 44  
value = [17, 27]"]
    R3 --> R5["entropy = 0.0  
samples = 9  
value = [0, 9]"]
  
```

# Classification Tree Interpretation



## Interpretation :

- Due to a large number of continuous predictors, a tree with several nodes and branches is generated.
- Tree starts with all 490 observations (Train set). 360 are non-defaulters (0) and the remaining 130 are defaulters (1).
- DEBTINC is the first split variable, left branch is  $\leq 12.6$  and right branch is  $> 12.6$ . 334/490 have  $\text{DEBTINC} \leq 12.6$ .
- EMPLOY is the second split on left branch, which further divides 334 obs. into 280 non-defaulters (0) and the remaining 54 as defaulters (1).
- The algorithm progresses till no further variable split is left.



# Classification Tree in Python – Prediction

# Generating Predictions for the model

```
y_pred = dtcl.predict(X_test)
y_pred_probs = dtcl.predict_proba(X_test)

cutoff = 0.3
pred_test = np.where(y_pred_probs[:,1] > cutoff, 1, 0)
pred_test
```

# Output

```
array([0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0,
       1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1,
       0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 1, 1, 0,
       1, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1,
       0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1,
       0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0,
       0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0,
       1, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0,
       0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 1, 0,
       0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1])
```



# Classification Tree in Python – ROC Curve

# Area Under ROC Curve

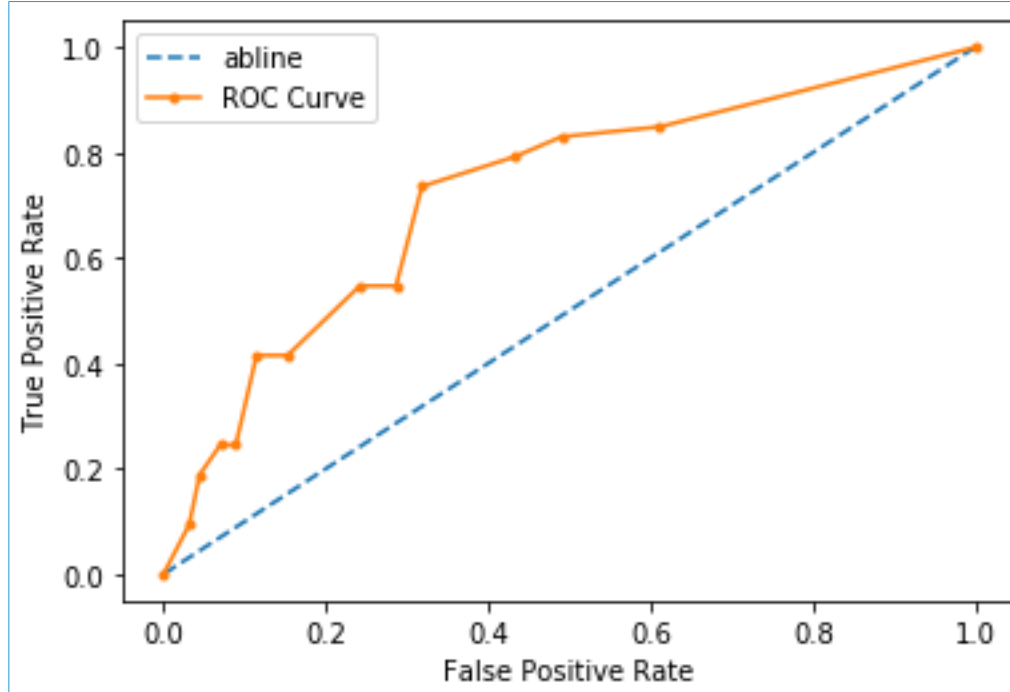
```
DTfpr, DTtpr, thresholds = roc_curve(y_test, y_pred_probs[:,1])

abline_probs = [0 for _ in range(len(y_test))]
abline_auc = roc_auc_score(y_test, abline_probs)
abline_fpr, abline_tpr, _ = roc_curve(y_test, abline_probs)

plt.plot(abline_fpr, abline_tpr, linestyle='--', label='abline')
plt.plot(DTfpr, DTtpr, marker='.', label='ROC Curve')
plt.xlabel('False Positive Rate');plt.ylabel('True Positive Rate')
plt.legend(); plt.show()
```

# Classification Tree in Python – ROC Curve

```
# Output
```



```
# Plotting The Tree
```

# Classification Tree in Python – Confusion Matrix

```
# Confusion Matrix
```

```
confusion_matrix(y_test, pred_test, labels=[0, 1])
```

```
array([[107,  50],  
       [ 14,  39]], dtype=int64)
```

```
accuracy_score(y_test, pred_test)
```

```
0.6952380952380952
```

```
precision_score(y_test, pred_test)
```

```
0.43820224719101125
```

```
recall_score(y_test, pred_test)
```

```
0.7358490566037735
```

- ❑ **accuracy\_score()** = number of correct predictions out of total predictions
- ❑ **precision\_score()** = true positives / (true positives + false positives)
- ❑ **recall\_score()** also known as 'Sensitivity' = true positives / (true positives + false negatives)

```
# Area Under ROC Curve
```

```
auc = roc_auc_score(y_test, y_pred_probs[:,1])
```

```
print('AUC: %.3f' % auc)
```

```
AUC: 0.720
```

THANK YOU!!