

\\## Building Decision Tree using sklearn

Data set : users.csv

✓ Step 1: Loading the data

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

↗ Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5 #reading data from csv file
6 user=pd.read_csv('/content/drive/MyDrive/Subjects/ Python_AI_ML_DL_NLP/3 Machine Learning/2. Supervised Machine Learning/2_Classific
7 user.info()
```

↗

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 99 entries, 0 to 98
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  -
0   User ID                99 non-null    int64
1   Gender                 99 non-null    object
2   Age                   99 non-null    int64
3   EstimatedSalary       99 non-null    int64
4   Purchased              99 non-null    int64
dtypes: int64(4), object(1)
memory usage: 4.0+ KB
```

```
1 user.head()
```

↗

	User ID	Gender	Age	EstimatedSalary	Purchased
0	15624510	Male	12	19000	0
1	15810944	Male	11	20000	0
2	15668575	Female	1	43000	0
3	15603246	Female	27	57000	0
4	15804002	Male	10	76000	0

```
1 user['Gender'].unique()
```

↗ array(['Male', 'Female'], dtype=object)

✓ Setting the predictors and target:

```
1 #selecting predictor attributes
2 X=user.columns.drop("Purchased")
3 #selecting target attribute
4 y=user['Purchased']
```

✓ Step 2: Feature engineering - Encoding

```
1 # Encoding all the predictor variables to convert the categorical values to numerical values.
2 user_encoded = pd.get_dummies(user[X])
3 print("Total number of predictors after encoding = ", len(user_encoded.columns))
4 # Printing the list of columns after encoding to understand the encoding process
5 user_encoded.columns
6 user_encoded
```

↩ Total number of predictors after encoding = 5

	User ID	Age	EstimatedSalary	Gender_Female	Gender_Male
0	15624510	12	19000	0	1
1	15810944	11	20000	0	1
2	15668575	1	43000	1	0
3	15603246	27	57000	1	0
4	15804002	10	76000	0	1
...
94	15786993	45	83000	1	0
95	15709441	35	44000	1	0
96	15710257	4	25000	1	0
97	15582492	5	123000	0	1
98	15575694	35	73000	0	1

99 rows x 5 columns

✓ Step 3: Splitting the dataset into train and test data

```

1 # Import the required function
2 from sklearn.model_selection import train_test_split
3 # Splitting data into train and test datasets
4 X_train, X_test, y_train, y_test = train_test_split(user_encoded, y, test_size=0.15, random_state=0)
5 # Printing the shape of the resulting datasets
6 print("Shape of X_train and y_train are:", X_train.shape, "and", y_train.shape, " respectively")
7 print("Shape of X_test and y_test are:", X_test.shape, "and", y_test.shape, " respectively")
8

```

↩ Shape of X_train and y_train are: (84, 5) and (84,) respectively
Shape of X_test and y_test are: (15, 5) and (15,) respectively

✓ Step 4: Building the model

```

1 # from scipy.sparse import random
2 # Importing required class
3 from sklearn.tree import DecisionTreeClassifier
4 # create object
5 model = DecisionTreeClassifier(criterion='entropy', random_state=1)
6 model.fit(X_train, y_train)

```

↩

```

DecisionTreeClassifier
DecisionTreeClassifier(criterion='entropy', random_state=1)

```

```

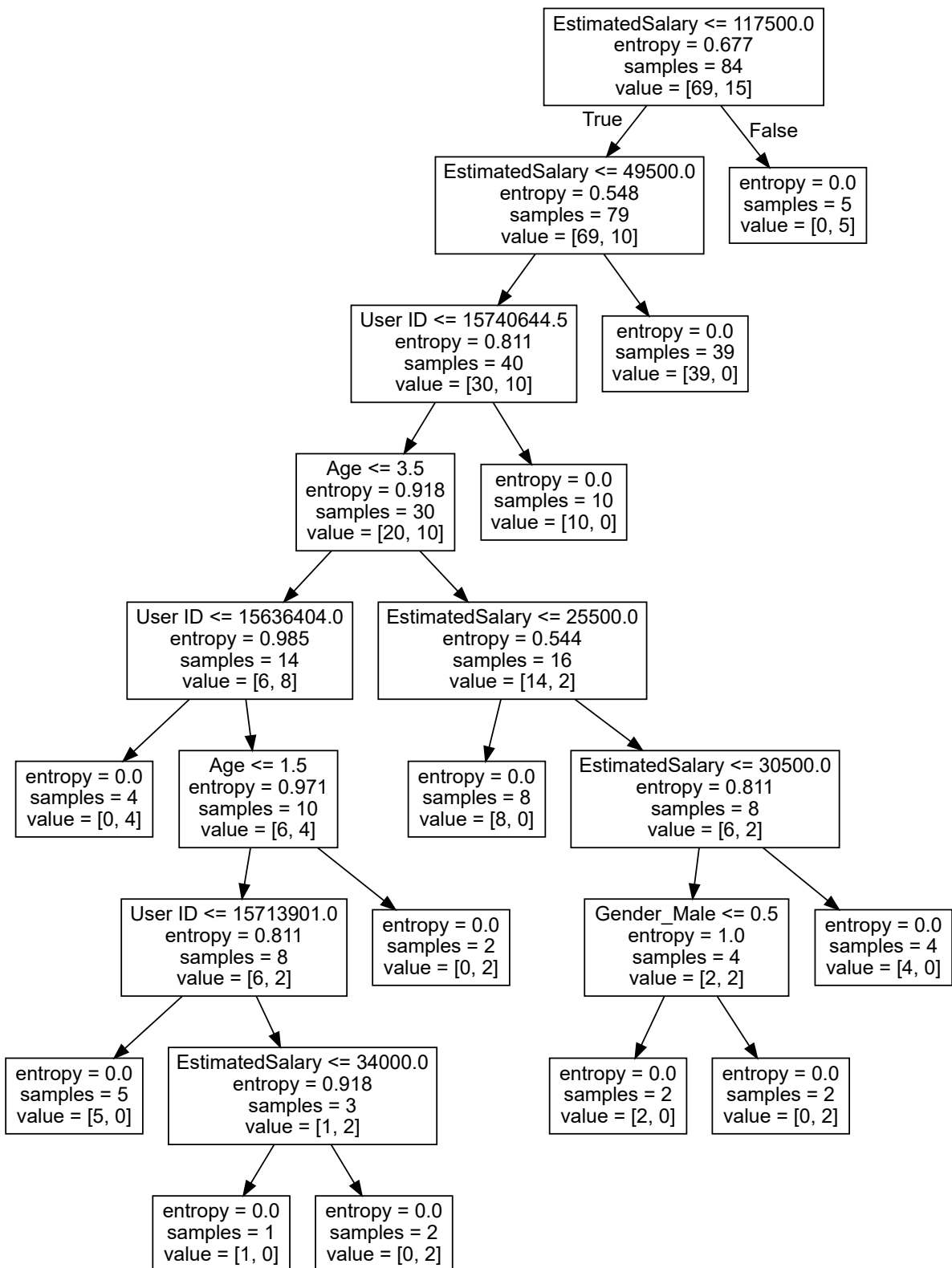
1 # Predicting target values using the model built on training data
2 train_predictions = model.predict(X_train)
3 test_predictions = model.predict(X_test)

```

```

1 # Importing the required libraries (Ensure that they are already installed.)
2 from sklearn.tree import export_graphviz
3 import graphviz
4 # Generating the tree
5 data = export_graphviz(model, out_file=None,
6                        feature_names=user_encoded.columns)
7 graph = graphviz.Source(data)
8 graph

```



✓ Step 5: Evaluate model performance on train and test sets

```
1 # Getting the accuracy on train data
2 train_accuracy = model.score(X_train,y_train)
3 print("Accuracy of the model on train data = ",train_accuracy)
4 # Getting the accuracy on test data
5 test_accuracy = model.score(X_test,y_test)
6 print("Accuracy of the model on test data = ",test_accuracy)
```

➡ Accuracy of the model on train data = 1.0
Accuracy of the model on test data = 0.8

You can observe that the training accuracy is 100% and the test accuracy is approximately 80%.

This could mean that the model is overfit to the training data and is not a good approximation of the input to output mapping.