

Project name : Adult census income prediction

In [1]:

```

1 # import libraries
2
3 import numpy as np
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 import imblearn
8 import warnings
9 warnings.filterwarnings('ignore')
10
11 from sklearn.preprocessing import LabelEncoder
12 from imblearn.over_sampling import RandomOverSampler
13 from sklearn.model_selection import train_test_split, GridSearchCV, RandomizedSearchCV
14 from sklearn.ensemble import RandomForestClassifier
15 from sklearn.tree import DecisionTreeClassifier
16 from sklearn.linear_model import LogisticRegression
17 from sklearn.tree import plot_tree
18 from sklearn.neighbors import KNeighborsClassifier
19 from sklearn.ensemble import AdaBoostClassifier
20 from sklearn.metrics import roc_curve, roc_auc_score, plot_confusion_matrix
21 from sklearn.preprocessing import MinMaxScaler, StandardScaler
22 from sklearn.metrics import confusion_matrix
23 from sklearn.metrics import confusion_matrix, classification_report, accuracy_score,
```

problem statement

The Goal is to predict whether a person has an income of more than 50K a year or not. This is basically a binary classification problem where a person is classified into the

>50K group or <=50K group.

- 1 Over the last two decades, humans have grown a lot of dependence on data and information in society and with this advent growth, technologies have evolved for their storage, analysis and processing on a huge scale. The fields of Data Mining and Machine Learning have not only exploited them for knowledge and discovery but also to explore certain hidden patterns and concepts which led to the prediction of future events, not easy to obtain.
- 2 The problem of income inequality has been of great concern in the recent years. Making the poor better off does not seem to be the sole criteria to be inquest for eradicating this issue. People of the United States believe that the advent of economic inequality is unacceptable and demands a fair share of wealth in the society. This model actually aims to conduct a comprehensive analysis to highlight the key factors that are necessary in improving an individual's income. Such an analysis helps to set focus on the important areas which can significantly improve the income levels of individuals.

In [2]:

```
1 !python --version
```

Python 3.10.0

Data Gathering

In [3]:

```
1 df = pd.read_csv('adult.csv')
2 df
```

Out[3]:

	age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	...
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	W
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	W
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	W
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	BI
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	BI
...
32556	27	Private	257302	Assoc-acdm	12	Married-civ-spouse	Tech-support	Wife	W
32557	40	Private	154374	HS-grad	9	Married-civ-spouse	Machine-op-insct	Husband	W
32558	58	Private	151910	HS-grad	9	Widowed	Adm-clerical	Unmarried	W
32559	22	Private	201490	HS-grad	9	Never-married	Adm-clerical	Own-child	W
32560	52	Self-emp-inc	287927	HS-grad	9	Married-civ-spouse	Exec-managerial	Wife	W

32561 rows × 15 columns

Exploratory data analysis

In [4]:

```
1 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32561 entries, 0 to 32560
Data columns (total 15 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   age         32561 non-null   int64  
 1   workclass   32561 non-null   object  
 2   fnlwgt     32561 non-null   int64  
 3   education   32561 non-null   object  
 4   education-num 32561 non-null   int64  
 5   marital-status 32561 non-null   object  
 6   occupation  32561 non-null   object  
 7   relationship 32561 non-null   object  
 8   race        32561 non-null   object  
 9   sex         32561 non-null   object  
 10  capital-gain 32561 non-null   int64  
 11  capital-loss 32561 non-null   int64  
 12  hours-per-week 32561 non-null   int64  
 13  country     32561 non-null   object  
 14  salary       32561 non-null   object  
dtypes: int64(6), object(9)
memory usage: 3.7+ MB
```

In [5]:

```
1 df.isna().sum()
```

Out[5]:

```
age          0
workclass    0
fnlwgt      0
education    0
education-num 0
marital-status 0
occupation   0
relationship  0
race         0
sex          0
capital-gain 0
capital-loss 0
hours-per-week 0
country      0
salary       0
dtype: int64
```

Here, we see that dataset didn't consist null values, so we didn't clean up our data according to missing values or null values

In [6]:

```
1 df.describe()
```

Out[6]:

	age	fnlwgt	education-num	capital-gain	capital-loss	hours-per-week
count	32561.000000	3.256100e+04	32561.000000	32561.000000	32561.000000	32561.000000
mean	38.581647	1.897784e+05	10.080679	1077.648844	87.303830	40.437456
std	13.640433	1.055500e+05	2.572720	7385.292085	402.960219	12.347429
min	17.000000	1.228500e+04	1.000000	0.000000	0.000000	1.000000
25%	28.000000	1.178270e+05	9.000000	0.000000	0.000000	40.000000
50%	37.000000	1.783560e+05	10.000000	0.000000	0.000000	40.000000
75%	48.000000	2.370510e+05	12.000000	0.000000	0.000000	45.000000
max	90.000000	1.484705e+06	16.000000	99999.000000	4356.000000	99.000000

In [7]:

```
1 df.dtypes
```

Out[7]:

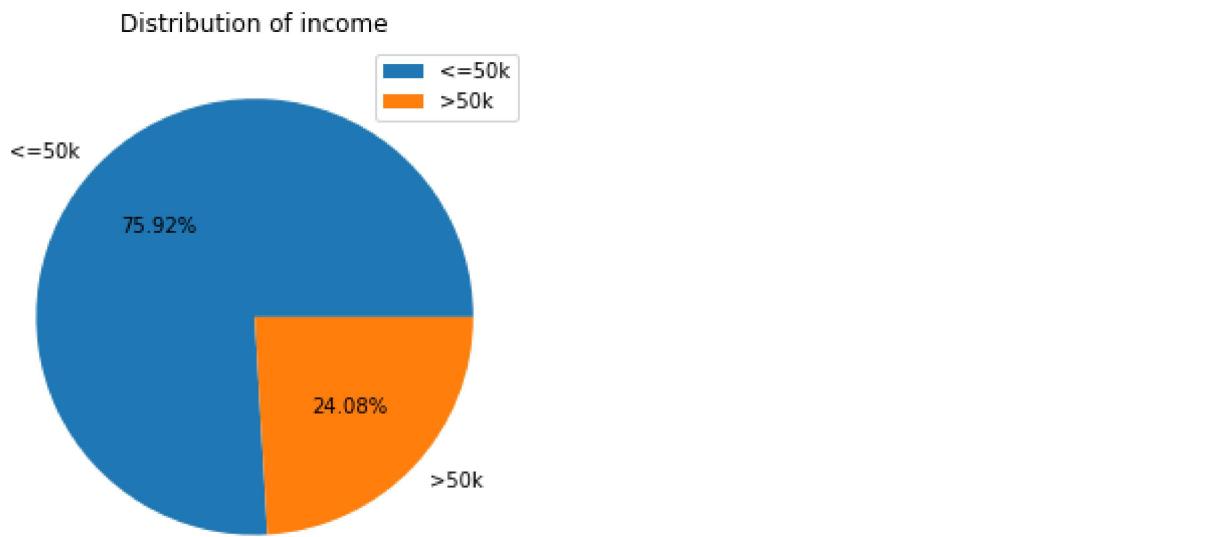
age	int64
workclass	object
fnlwgt	int64
education	object
education-num	int64
marital-status	object
occupation	object
relationship	object
race	object
sex	object
capital-gain	int64
capital-loss	int64
hours-per-week	int64
country	object
salary	object
dtype:	object

The target column in our dataset is binary, so we apply classification algorithm on dataset. but some of the input or independent variables are in the form of object data type, so we have to do feature engineering on our dataset so that We need to convert all variables to the same datatype for classification.. Also, We can observe that education-num and education are the same variables with different data types. We can drop any one of them.

i) To check the income proportion in dataset ($\leq 50k$ and $> 50k$)

In [8]:

```
1 # pie chart for income prediction
2
3 df1 = plt.subplots(figsize=(8,5))
4 plt.pie(df['salary'].value_counts(),labels=['<=50k','>50k'],autopct='%.2f%%')
5 plt.legend()
6 plt.title('Distribution of income')
7 plt.show()
```

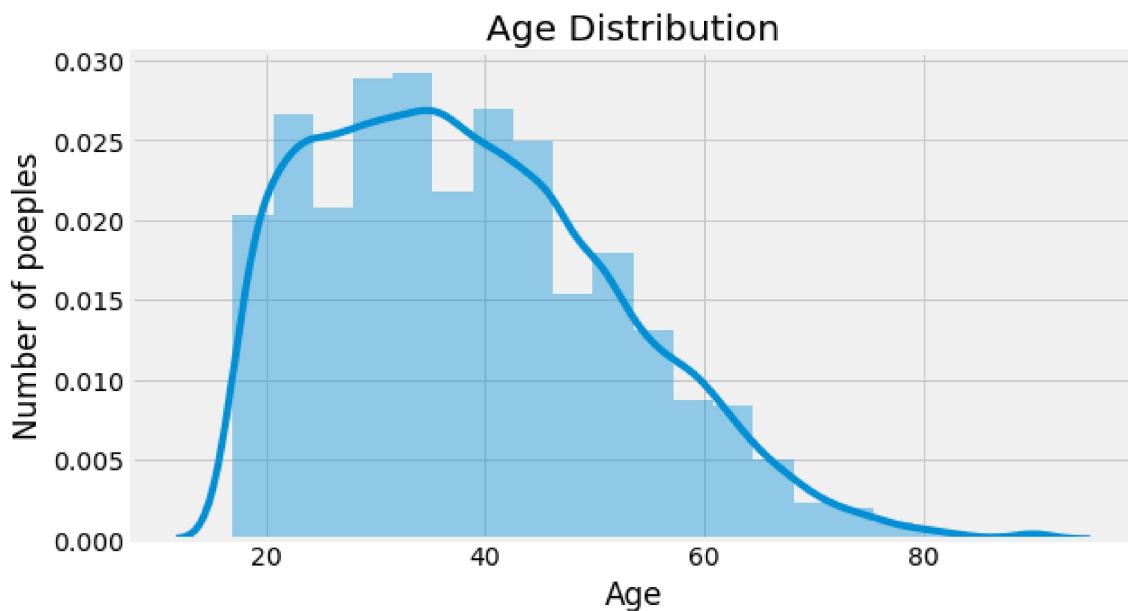


From above diagram we clearly conclude that the income is not distributed properly.

ii) To create a distribution plot for 'AGE'

In [9]:

```
1 age = df['age'].value_counts()
2 plt.figure(figsize=(10,5))
3 plt.style.use('fivethirtyeight')
4 sns.distplot(df['age'],bins=20)
5 plt.title('Age Distribution')
6 plt.xlabel('Age')
7 plt.ylabel('Number of poeple')
8 plt.show()
```



Sample population contain more number of people in the age group of 20-40

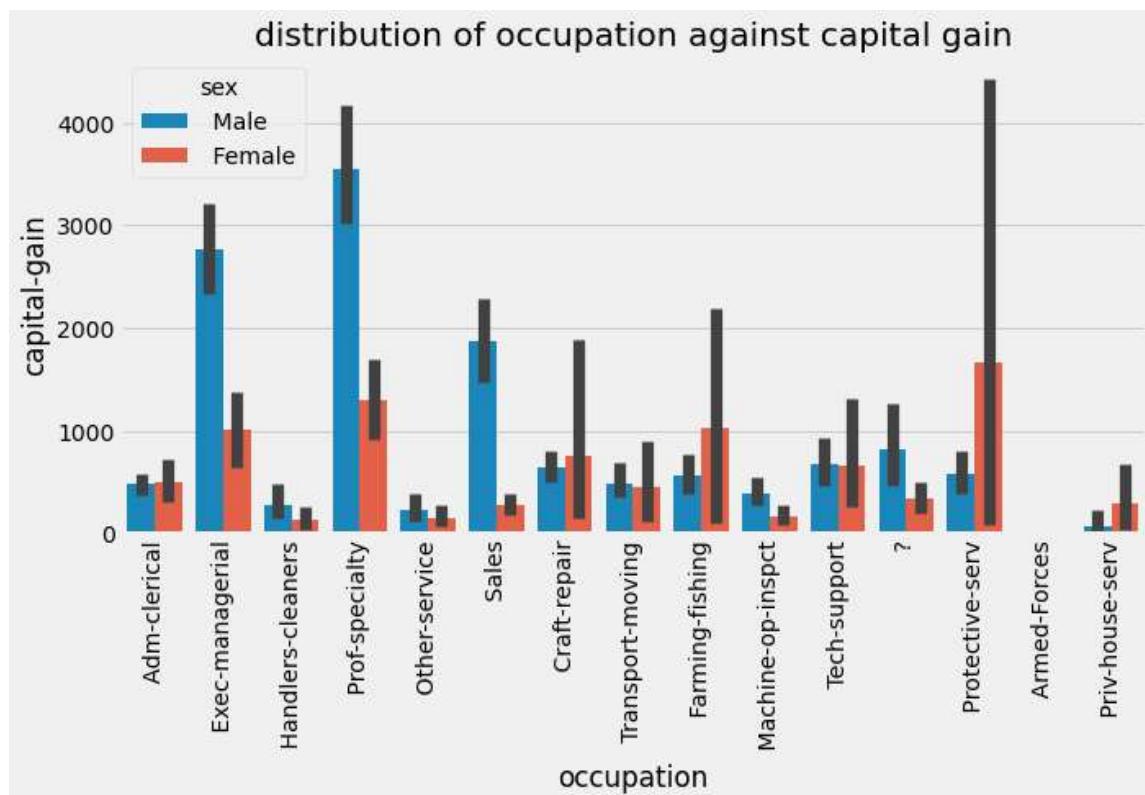
iii) To check the distribution of capital-gain and occupation as per gender

In [10]:

```
1 plt.figure(figsize=(10,5))
2 sns.barplot(x='occupation',y='capital-gain',data=df,hue='sex')
3 plt.xticks(rotation=90)
4 plt.title('distribution of occupation against capital gain')
```

Out[10]:

Text(0.5, 1.0, 'distribution of occupation against capital gain')



Capital gain is totally depend on occupation rather than income

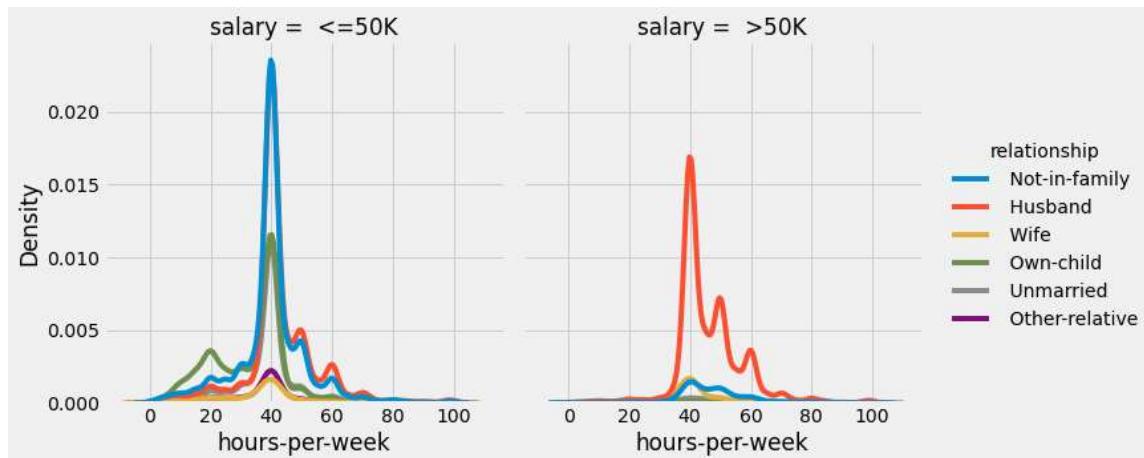
iv) To create a distribution plot of hours-per-week, relationship and salary

In [11]:

```
1 sns.displot(data=df, x=df["hours-per-week"], hue=df["relationship"], col=df["salary"])
```

Out[11]:

<seaborn.axisgrid.FacetGrid at 0x15c3dc18eb0>



peoples who are not in the family are worked for 40 hours and their salary is <=50k.

People who have a relationship of husband work for around 40 hours per week and receive a salary of >50 k.

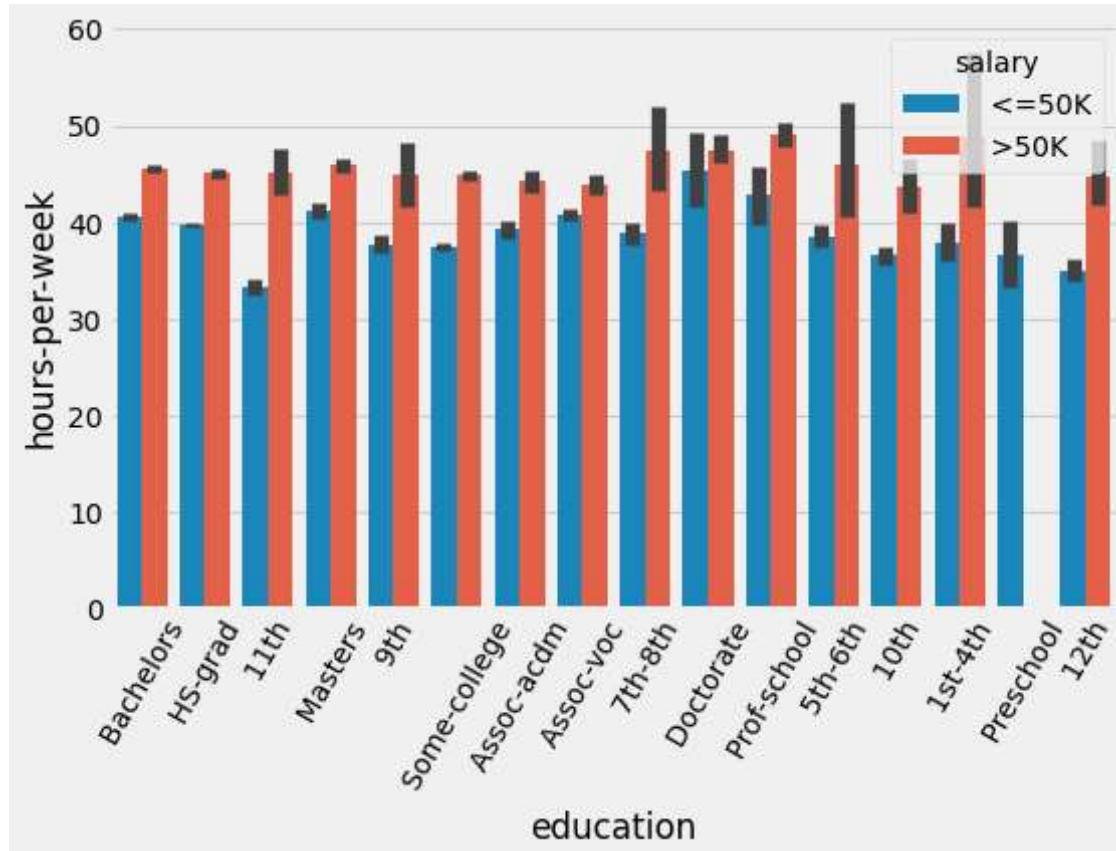
v) To create barplot of hours-per-week, education and salary

In [12]:

```
1 plt.figure(figsize=(8,5))
2 sns.barplot(x='education',y='hours-per-week',data=df,hue='salary')
3 plt.xticks(rotation=60)
```

Out[12]:

```
(array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15]),
 [Text(0, 0, ' Bachelors'),
  Text(1, 0, ' HS-grad'),
  Text(2, 0, ' 11th'),
  Text(3, 0, ' Masters'),
  Text(4, 0, ' 9th'),
  Text(5, 0, ' Some-college'),
  Text(6, 0, ' Assoc-acdm'),
  Text(7, 0, ' Assoc-voc'),
  Text(8, 0, ' 7th-8th'),
  Text(9, 0, ' Doctorate'),
  Text(10, 0, ' Prof-school'),
  Text(11, 0, ' 5th-6th'),
  Text(12, 0, ' 10th'),
  Text(13, 0, ' 1st-4th'),
  Text(14, 0, ' Preschool'),
  Text(15, 0, ' 12th')])
```



Here we clearly see that the people with more hours of work has more salary

Feature Engineering

In [13]:

```

1 for dataset in [df]:
2     dataset.loc[dataset['country'] != ' United-States', 'country'] = 0
3     dataset.loc[dataset['country'] == ' United-States', 'country'] = 1
4     dataset.loc[dataset['race'] != ' White', 'race'] = 0
5     dataset.loc[dataset['race'] == ' White', 'race'] = 1
6     dataset.loc[dataset['workclass'] != ' Private', 'workclass'] = 0
7     dataset.loc[dataset['workclass'] == ' Private', 'workclass'] = 1
8     dataset.loc[dataset['hours-per-week'] <= 40, 'hours-per-week'] = 0
9     dataset.loc[dataset['hours-per-week'] > 40, 'hours-per-week'] = 1
10    for col in df[df.columns]:
11        if df[col].dtypes == 'object':
12            le = LabelEncoder()
13            df[col] = le.fit_transform(df[col])
14    df = df.astype(int)
15    df=df.drop(["education"],axis=1)
16    df.head()

```

Out[13]:

	age	workclass	fnlwgt	education-num	marital-status	occupation	relationship	race	sex	capital-gain
0	39	0	77516	13	4	1	1	1	1	2174
1	50	0	83311	13	2	4	0	1	1	0
2	38	1	215646	9	0	6	1	1	1	0
3	53	1	234721	7	2	6	0	0	1	0
4	28	1	338409	13	2	10	5	0	0	0

train_test_split

In [14]:

```

1 x = df.drop('salary',axis=1)
2 y = df['salary']

```

In [15]:

```
1 df['salary'].value_counts()
```

Out[15]:

```

0    24720
1    7841
Name: salary, dtype: int64

```

As we see here, our data is imbalanced. Imbalanced data is a common problem in machine learning, which brings challenges to feature correlation, class separation and evaluation, and results in poor model performance. Classes that make up a large proportion of the data set are called majority

classes. Those that make up a smaller proportion are minority classes. so by using randomoversampler method we make our data balanced.

In [16]:

```
1 random_sampler = RandomOverSampler(random_state=30, sampling_strategy=0.95)
2 random_sampler.fit(x,y)
```

Out[16]:

RandomOverSampler(random_state=30, sampling_strategy=0.95)

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [17]:

```
1 x_new,y_new = random_sampler.fit_resample(x, y)
```

In [18]:

```
1 y_new.value_counts()
```

Out[18]:

```
0    24720
1    23484
Name: salary, dtype: int64
```

In [19]:

```
1 x_new_train, x_new_test, y_new_train,y_new_test = train_test_split(x_new,y_new,test_
2 x_new_train.shape, x_new_test.shape, y_new_train.shape, y_new_test.shape)
```

Out[19]:

```
((33742, 13), (14462, 13), (33742,), (14462,))
```

Model training

Random_forest

- 1 random forest doesn't make any assumption about the data or its distribution. Hence it generally requires minimal data transformations. Random forest algorithm makes use of random subsets of features and hence it can perform quite well with a high dimensional dataset (a dataset with a large number of features).
- 2 -----
- 3 While building a machine learning model, our aim is to generalize the model properly for giving predictions on unseen data.

4 The problem of overfitting takes place when we have a flexible model. A flexible model is having high variance because the learned parameters like the structure of the decision tree, etc will vary with the training data. On the contrary, an inflexible model is said to have a high bias as it makes assumptions about the training data and an inflexible model may not have the capacity to fit even the training data and in both situations, the model has high variance, and high bias implies the model is not able to generalize new and unseen data points properly.

5 So, we have to build a model carefully by keeping the bias-variance tradeoff in mind.

6 The main reason for the overfitting of the decision tree due to not put the limit on the maximum depth of the tree is because it has unlimited flexibility, which means it keeps growing unless, for every single observation, there is one leaf node present.

7 Moreover, instead of limiting the depth of the tree which results in reduced variance and an increase in bias, we can combine many decision trees that eventually convert into a forest, known as a single ensemble model (known as the random forest).

8 -----

9 The steps that are included while performing the random forest algorithm are as follows:

10

11 Step-1: Pick K random records from the dataset having a total of N records.

12

13 Step-2: Build and train a decision tree model on these K records.

14

15 Step-3: Choose the number of trees you want in your algorithm and repeat steps 1 and 2.

16

17 Step-4: In the case of a regression problem, for an unseen data point, each tree in the forest predicts a value for output. The final value can be calculated by taking the mean or average of all the values predicted by all the trees in the forest.

18

19 and, in the case of a classification problem, each tree in the forest predicts the class to which the new data point belongs. Finally, the new data point is assigned to the class that has the maximum votes among them i.e, wins the majority vote.

20 -----

21 n_estimators: We know that a random forest is nothing but a group of many decision trees, the n_estimator parameter controls the number of trees inside the classifier. We may think that using many trees to fit a model will help us to get a more generalized result, but this is not always the case. However, it will not cause any overfitting but can certainly increase the time complexity of the model. The default number of estimators is 100 in scikit-learn.

22

23 max_depth: It governs the maximum height upto which the trees inside the forest can grow. It is one of the most important hyperparameters when it comes to increasing the accuracy of the model, as we increase the depth of the tree the model accuracy increases upto a certain limit but then it will start to decrease gradually because of overfitting in the model. It is important to set its value appropriately to avoid overfitting. The default value is set to None, None specifies that the nodes inside the tree will continue to grow until all leaves become pure or all leaves contain less than

24

25 min_samples_split (another hyperparameter).

```
26 min_samples_split: It specifies the minimum amount of samples an internal node must hold in order to split into further nodes. If we have a very low value of min_samples_splits then, in this case, our tree will continue to grow and start overfitting. By increasing the value of min_samples_splits we can decrease the total number of splits thus limiting the number of parameters in the model and thus can aid in reducing the overfitting in the model. However, the value should not be kept very large that a number of parameters drop extremely causing the model to underfit. We generally keep
27
28 min_samples_split value between 2 and 6. However, the default value is set to 2.
29 min_samples_leaf: It specifies the minimum amount of samples that a node must hold after getting split. It also helps to reduce overfitting when we have ample amount of parameters. Less number of parameters can lead to overfitting also, we should keep in mind that increasing the value to a large number can lead to less number of parameters and in this case model can underfit also. The default value is set to 1.
30
31 max_features: Random forest takes random subsets of features and tries to find the best split. max_features helps to find the number of features to take into account in order to make the best split. It can take four values "auto", "sqrt", "log2" and None.
32 In case of auto: considers max_features = sqrt(n_features)
33 In case of sqrt: considers max_features = sqrt(n_features), it is same as auto
34 In case of log2: considers max_features = log2(n_features)
35 In case of None: considers max_features = n_features
36
37 max_leaf_nodes: It sets a limit on the splitting of the node and thus helps to reduce the depth of the tree, and effectively helps in reducing overfitting. If the value is set to None, the tree continues to grow infinitely.
38
39 max_samples: This hyperparameter helps to choose maximum number of samples from the training dataset to train each individual tree.
40 -----
41 Advantages:
42
43 -- Random Forest is unbiased as we train multiple decision trees and each tree is trained on a subset of the same training data.
44 -- It is very stable since if we introduce the new data points in the dataset, then it does not affect much as the new data point impacts one tree, and is pretty hard to impact all the trees.
45 -- Also, it works well when you have both categorical and numerical features in the problem statement.
46 -- It performs very well, with missing values in the dataset.
47 -----
48 Real life applications of Random Forests :
49
50 -- Fraud detection for bank accounts, credit card.
51
52 -- Detect and predict the drug sensitivity of a medicine.
53
54 -- Identify a patient's disease by analyzing their medical records.
55
56 -- Predict estimated loss or profit while purchasing a particular stock.
```

In [20]:

```
1 rf_model = RandomForestClassifier(random_state=15)
2 rf_model.fit(x_new_train,y_new_train)
```

Out[20]:

RandomForestClassifier(random_state=15)

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [21]:

```
1 y_pred = rf_model.predict(x_new_test)
```

Model evaluation

In [22]:

```
1 # Training data Evaluation
2 def model_eval(model,x,y):
3     y_pred = model.predict(x)
4
5     cm = confusion_matrix(y, y_pred)
6     print(f"Confusion Matrix = \n{cm}")
7     print("*"*50)
8
9     ac = accuracy_score(y, y_pred)
10    print(f"Accuracy Score = {ac}")
11    print("*"*50)
12
13    cr = classification_report(y, y_pred)
14    print(f"Classification Report = \n{cr}")
15
16    pr = precision_score(y,y_pred)
17    print(f'precision_score = {pr}')
18
19    rc = recall_score(y,y_pred)
20    print(f'recall_score = {rc}')
21
22    f1 = f1_score(y,y_pred)
23    print(f'f1_score = {f1}')
24
25    return "Model Eval Success"
```

```
1 >> Evaluation matrix :
2
3 Confusion Matrix :
4     - It is a N * N matrix used to evaluate performance of Classification model.
5     - It compares actual values and predicted values
6
7     |-----|
8     |           Actual Values           |
9     |-----|-----|
10    |           | Positive | Negative |
```

```

11 |Predicted | Positive | TP | FP |
12 |Values    | Negative | FN | TN |
13 |-----|-----|-----|-----|
14
15     ex. Positive, Negative
16     - True Positive (TP):
17         -Actual values is Positive and predicted as Positive.
18         -the pred(Positive) matches with actual (Positive) value.
19
20     - True Negative (TN):
21         -Actual values is Negative and predicted as Negative.
22         -the pred(Negative) matches with actual (Negative) value.
23
24     - False Positive (FP): Type-1 Error
25         -Actual values is Negative and predicted as Positive.
26         -the pred(Positive) not matches with actual (Negative) value.
27         -the pred value is falsely predicted.
28
29     - False Negative (FN): Type-2 Error
30         -Actual values is Positive and predicted as Negative.
31         -the pred(Negative) not matches with actual (Positive) value.
32         -the pred value is falsely predicted.
33
34 -----
35
36 We Define classification report using Precision,Recall,F1-Score and Support
37
38 1.Precision:
39
40     Precision = TP/TP+FP
41
42     We will focus on FP
43     It ranges from 0 to 1
44
45     Out of total predicted positive results(TP+FP) how many are actually positive
46     (TP)
47
48     When FP is important then Precision is the appropriate Metric
49     When money is on stake, Stock market, e-commerce, Finance
50
51 2.Recall:
52
53     Recall = TP/TP+FN
54
55     We will focus on FN
56     It ranges from 0 to 1
57
58     Out of total actual positive results(TP+FN) how many are predicted positive
59     (TP)
60
61     When FN is important then Recall is the appropriate Metric
62     When lives are on stake, Medical, Pharma
63
64 3.F1-Score:
65     f1-score = 2PR/P+R
66
67     It is harmonic mean of Precision and Recall
68     When Precision is low and Recall is high or vice-versa, its difficult to
69     compare
70     in such cases f1-score is a good metric to use.

```

```
68 | It strikes a balance between Precision and Recall.  
69 |  
70 | 4. Support: It is the total entries of each class in the actual dataset. It is  
simply the sum of rows for every class  
71 |  
72 | 5. Macro and weighted average of precision and recall:  
73 |  
74 |     -Precision Weighted avg= precision class0*(Support*Total)+Precision class1*  
(Support*Total)  
75 |     -Recall Weighted avg= Recall class0*(Support*Total)+Recall class1*  
(Support*Total)
```

In [23]:

```
1 model_eval(rf_model,x_new_train,y_new_train)
```

```
Confusion Matrix =  
[[17302    2]  
 [    0 16438]]  
*****  
Accuracy Score = 0.9999407266907712  
*****  
Classification Report =  
precision    recall    f1-score   support  
  
      0         1.00      1.00      1.00      17304  
      1         1.00      1.00      1.00      16438  
  
accuracy                  1.00      33742  
macro avg          1.00      1.00      1.00      33742  
weighted avg         1.00      1.00      1.00      33742
```

```
precision_score = 0.9998783454987834  
recall_score = 1.0  
f1_score = 0.9999391690492122
```

Out[23]:

```
'Model Eval Success'
```

In [24]:

```
1 model_eval(rf_model,x_new_test,y_new_test)
```

```
Confusion Matrix =
[[6550  866]
 [ 231 6815]]
*****
Accuracy Score = 0.9241460378924077
*****
Classification Report =
      precision    recall   f1-score   support
      0          0.97     0.88     0.92      7416
      1          0.89     0.97     0.93      7046
      accuracy           0.92      14462
      macro avg       0.93     0.93     0.92      14462
      weighted avg    0.93     0.92     0.92      14462

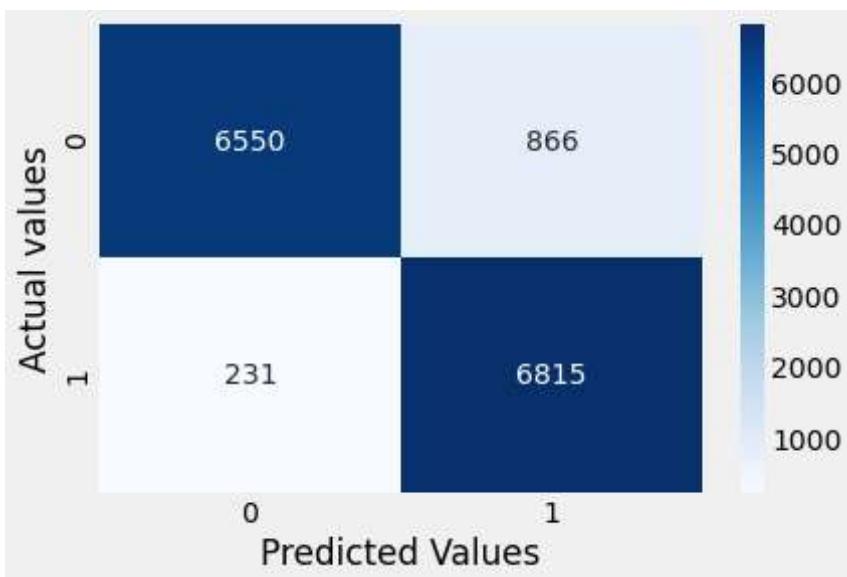
precision_score = 0.8872542637677385
recall_score = 0.9672154413851831
f1_score = 0.9255109662524613
```

Out[24]:

'Model Eval Success'

In [25]:

```
1 cm = confusion_matrix(y_new_test, y_pred)
2 plt.figure(figsize=(6,4))
3 sns.heatmap(cm, annot=True, fmt='d', cmap="Blues")
4 plt.xlabel("Predicted Values")
5 plt.ylabel("Actual values")
6 plt.show()
```



In random forest algorithm, it consist low bias and low variance, so it is generalised model. it doesn't need of hyperparameter tuning

In [26]:

```
1 model_eval_test = pd.DataFrame({'recall':[recall_score(y_new_test,y_pred)],'precision':[precision_score(y_new_test,y_pred)],'f1_score':[f1_score(y_new_test,y_pred)],'accuracy':[accuracy_score(y_new_test,y_pred)]})
2 model_eval_test
```

Out[26]:

	recall	precision	f1_score	accuracy
0	0.967215	0.887254	0.925511	0.924146

Feature_selection

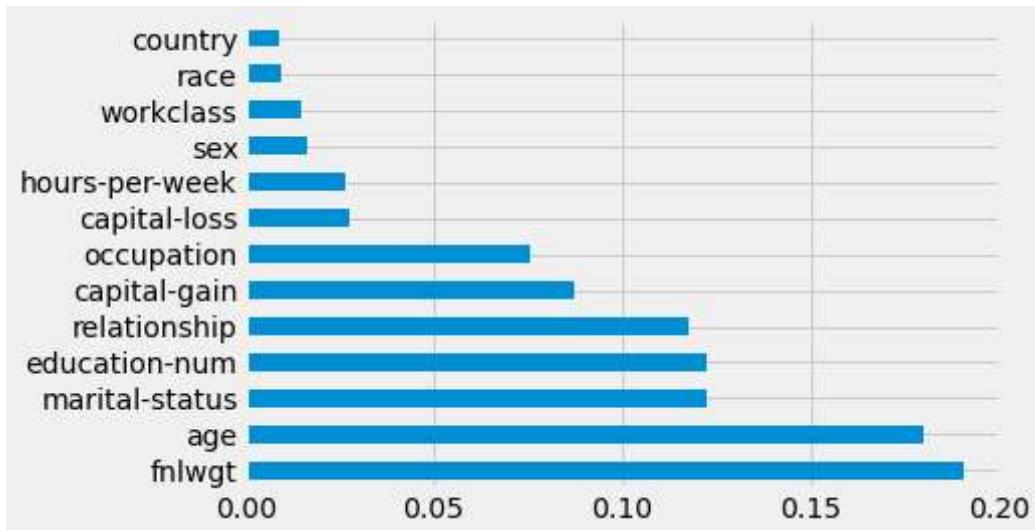
Feature Selection is the method of reducing the input variable to your model by using only relevant data and getting rid of noise in data. It is the process of automatically choosing relevant features for your machine learning model based on the type of problem you are trying to solve.

In [27]:

```
1 s1 = pd.Series(rf_model.feature_importances_, index = x_new.columns)
2 s1.sort_values(ascending=False).plot(kind="barh")
```

Out[27]:

<AxesSubplot:>



As we see the diagram, every input variable has their own impact on dataset, so we didn't ignore or neglect any of the variable. we will work by including all the variables.

Decision_tree

- A decision tree is a type of supervised learning algorithm that is commonly used in machine learning to model and predict outcomes based on input data. It is a tree-like structure where each internal node represents a decision or test on a specific feature or attribute, each branch represents the outcome of that decision, and each leaf node represents the final decision or prediction.
- Example of a decision tree

4

5 Suppose we want to build a decision tree to predict whether a person is likely to buy a new car based on their demographic and behavior data. The decision tree starts with the root node, which represents the entire dataset. The root node splits the dataset based on the “income” attribute. If the person’s income is less than or equal to \$50,000, the decision tree follows the left branch, and if the income is greater than \$50,000, the decision tree follows the right branch.

6

7 The left branch leads to a node that represents the “age” attribute. If the person’s age is less than or equal to 30, the decision tree follows the left branch, and if the age is greater than 30, the decision tree follows the right branch. The right branch leads to a leaf node that predicts that the person is unlikely to buy a new car.

8

9 The left branch leads to another node that represents the “education” attribute. If the person’s education level is less than or equal to high school, the decision tree follows the left branch, and if the education level is greater than high school, the decision tree follows the right branch. The left branch leads to a leaf node that predicts that the person is unlikely to buy a new car. The right branch leads to another node that represents the “credit score” attribute. If the person’s credit score is less than or equal to 650, the decision tree follows the left branch, and if the credit score is greater than 650, the decision tree follows the right branch. The left branch leads to a leaf node that predicts that the person is unlikely to buy a new car. The right branch leads to a leaf node that predicts that the person is likely to buy a new car.

10

11 In summary, a decision tree is a graphical representation of all the possible outcomes of a decision based on the input data. It is a powerful tool for modeling and predicting outcomes in a wide range of domains, including business, finance, healthcare, and more.

12

13 -- Below are some assumptions that we made while using the decision tree:

14

15 i) At the beginning, we consider the whole training set as the root.

16 ii) Feature values are preferred to be categorical. If the values are continuous then they are discretized prior to building the model.

17 iii) On the basis of attribute values, records are distributed recursively.

18 iv) We use statistical methods for ordering attributes as root or the internal node.

19

20 In the Decision Tree, the major challenge is the identification of the attribute for the root node at each level. This process is known as attribute selection. We have two popular attribute selection measures:

21

22 i) Information Gain

23 ii) Gini Index

24

25 -- 1. Information Gain :

26

27 When we use a node in a decision tree to partition the training instances into smaller subsets the entropy changes. Information gain is a measure of this change in entropy. Definition: Suppose S is a set of instances, A is an attribute, S_v is the subset of S with $A = v$, and $\text{Values}(A)$ is the set of all possible values of A, then $\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|\text{left } | S_v|}{|\text{right } | S|} \text{Entropy}(S_v)$. Entropy is the measure of uncertainty of a random variable, it characterizes the impurity of an arbitrary collection of examples. The higher the entropy more the information content. Definition: Suppose S is a set of instances, A is an attribute, S_v is the subset of S with $A = v$, and $\text{Values}(A)$ is the set of all possible values of A, then

28

29 -- $\text{Information_gain} = E(S) - \sum(\text{weighted_average} * \text{entropy of each feature})$

```
30
31 -- Building Decision Tree using Information Gain The essentials:
32
33 i) Start with all training instances associated with the root node
34 ii) Use info gain to choose which attribute to label each node with
35 iii) Note: No root-to-leaf path should contain the same discrete attribute twice
36 iv) Recursively construct each subtree on the subset of training instances that
37 would be classified down that path in the tree.
38 v) If all positive or all negative training instances remain, the label that node
39 "yes" or "no" accordingly
40 vi) If no attributes remain, label with a majority vote of training instances left
41 at that node
42 vii) If no instances remain, label with a majority vote of the parent's training
43 instances.
44
45 -- Gini Index :
46
47 i) Gini Index is a metric to measure how often a randomly chosen element would be
48 incorrectly identified.
49 ii) It means an attribute with a lower Gini index should be preferred.
50 iii) Sklearn supports "Gini" criteria for Gini Index and by default, it takes
51 "gini" value.
52 iv) The Formula for the calculation of the Gini Index is given below.
53 v) The Gini Index is a measure of the inequality or impurity of a distribution,
54 commonly used in decision trees and other machine learning algorithms. It ranges
55 from 0 to 1, where 0 represents perfect equality (all values are the same) and 1
56 represents perfect inequality (all values are different).
57
58 -- Some additional features and characteristics of the Gini Index are:
59
60 i) It is calculated by summing the squared probabilities of each outcome in a
61 distribution and subtracting the result from 1.
62 ii) A lower Gini Index indicates a more homogeneous or pure distribution, while a
63 higher Gini Index indicates a more heterogeneous or impure distribution.
64 iii) In decision trees, the Gini Index is used to evaluate the quality of a split
65 by measuring the difference between the impurity of the parent node and the
66 weighted impurity of the child nodes.
67 Compared to other impurity measures like entropy, the Gini Index is faster to
68 compute and more sensitive to changes in class probabilities.
69 iv) One disadvantage of the Gini Index is that it tends to favor splits that
70 create equally sized child nodes, even if they are not optimal for classification
71 accuracy.
72 v) In practice, the choice between using the Gini Index or other impurity measures
73 depends on the specific problem and dataset, and often requires experimentation
74 and tuning
```

In [28]:

```
1 tree = DecisionTreeClassifier(criterion='entropy', max_depth= 6, min_samples_leaf= 9, n
2 tree.fit(x_new_train,y_new_train)
```

Out[28]:

```
DecisionTreeClassifier(criterion='entropy', max_depth=6, min_samples_leaf=
9,
min_samples_split=7)
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [29]:

```
1 y_pred = tree.predict(x_new_test)
```

In [30]:

```
1 print(model_eval(tree,x_new_train,y_new_train))
```

```
Confusion Matrix =
[[13617 3687]
 [ 2716 13722]]
*****
Accuracy Score = 0.8102365005038231
*****
Classification Report =
      precision    recall   f1-score   support
          0        0.83     0.79     0.81     17304
          1        0.79     0.83     0.81     16438
          accuracy           0.81     33742
          macro avg       0.81     0.81     0.81     33742
          weighted avg     0.81     0.81     0.81     33742
precision_score = 0.7882129932793382
recall_score = 0.8347730867502129
f1_score = 0.8108251839158567
Model Eval Success
```

In [31]:

```
1 print(model_eval(tree,x_new_test,y_new_test))
```

Confusion Matrix =
[[5831 1585]
[1144 5902]]

Accuracy Score = 0.8112985755773752

Classification Report =

	precision	recall	f1-score	support
0	0.84	0.79	0.81	7416
1	0.79	0.84	0.81	7046
accuracy			0.81	14462
macro avg	0.81	0.81	0.81	14462
weighted avg	0.81	0.81	0.81	14462

precision_score = 0.7882997195138239

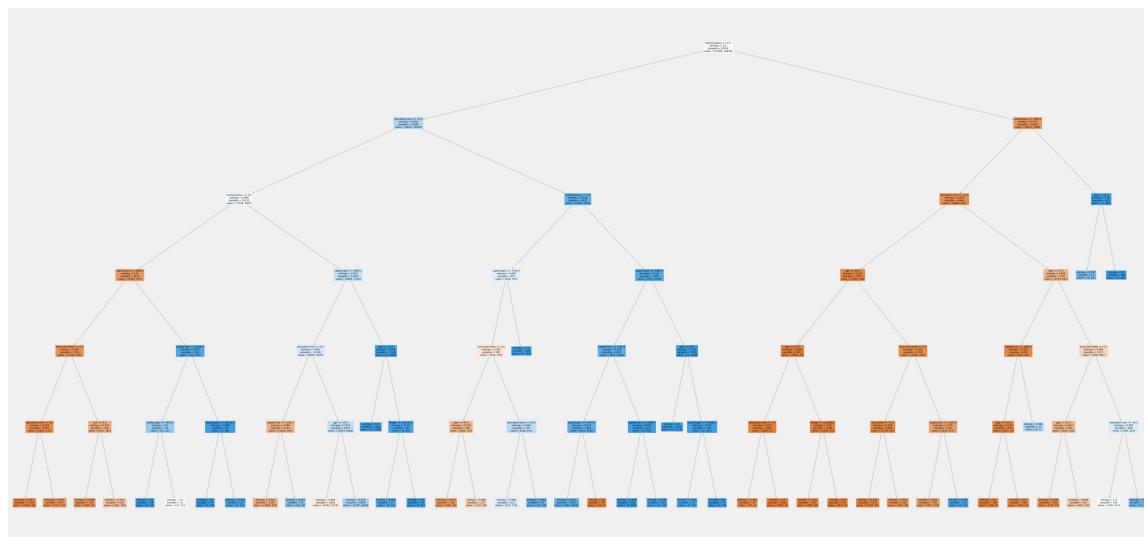
recall_score = 0.8376383763837638

f1_score = 0.8122204637721049

Model Eval Success

In [32]:

```
1 plt.figure(figsize=(100,50))
2 decision_tree_ = plot_tree(tree,feature_names=x_new.columns,filled=True)
3 plt.savefig('decision_tree.png')
```



In [33]:

```
1 tree = DecisionTreeClassifier()
2 hype = {'criterion':['entropy','gini index'],
3         'min_samples_split':np.arange(2,10),
4         'min_samples_leaf':np.arange(1,10),
5         'max_depth':[4,5,6,7]}
6 rscv = RandomizedSearchCV(tree,hype,cv=7)
7 rscv.fit(x_new_train,y_new_train)
8 rscv.best_params_
```

Out[33]:

```
{'min_samples_split': 8,
'min_samples_leaf': 1,
'max_depth': 7,
'criterion': 'entropy'}
```

In [34]:

```
1 tree = DecisionTreeClassifier(criterion='entropy',max_depth= 7,min_samples_leaf= 7,r
2 tree.fit(x_new_train,y_new_train)
```

Out[34]:

```
DecisionTreeClassifier(criterion='entropy', max_depth=7, min_samples_leaf=
7,
                     min_samples_split=8)
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [35]:

```
1 y_pred = tree.predict(x_new_test)
```

In [36]:

```
1 print(model_eval(tree,x_new_train,y_new_train))
```

```
Confusion Matrix =  
[[13969  3335]  
 [ 2977 13461]]  
*****  
Accuracy Score = 0.812933436073736  
*****  
Classification Report =  
precision    recall   f1-score   support  
  
      0       0.82      0.81      0.82      17304  
      1       0.80      0.82      0.81      16438  
  
accuracy          0.81      0.81      0.81      33742  
macro avg       0.81      0.81      0.81      33742  
weighted avg     0.81      0.81      0.81      33742  
  
precision_score = 0.8014408192426768  
recall_score = 0.8188952427302592  
f1_score = 0.8100740205813323  
Model Eval Success
```

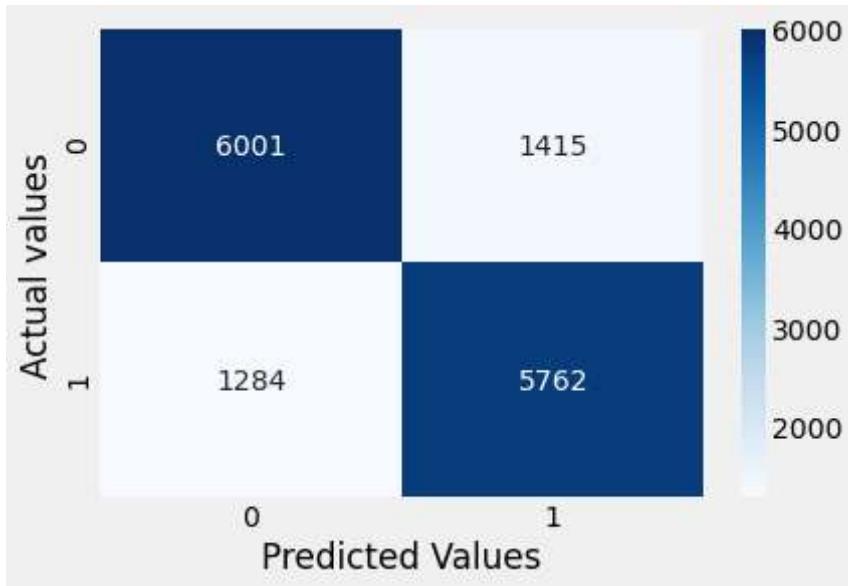
In [37]:

```
1 print(model_eval(tree,x_new_test,y_new_test))
```

```
Confusion Matrix =  
[[6001 1415]  
 [1284 5762]]  
*****  
Accuracy Score = 0.8133729774581663  
*****  
Classification Report =  
precision    recall   f1-score   support  
  
      0       0.82      0.81      0.82      7416  
      1       0.80      0.82      0.81      7046  
  
accuracy          0.81      0.81      0.81      14462  
macro avg       0.81      0.81      0.81      14462  
weighted avg     0.81      0.81      0.81      14462  
  
precision_score = 0.8028424132645953  
recall_score = 0.8177689469202384  
f1_score = 0.8102369401673346  
Model Eval Success
```

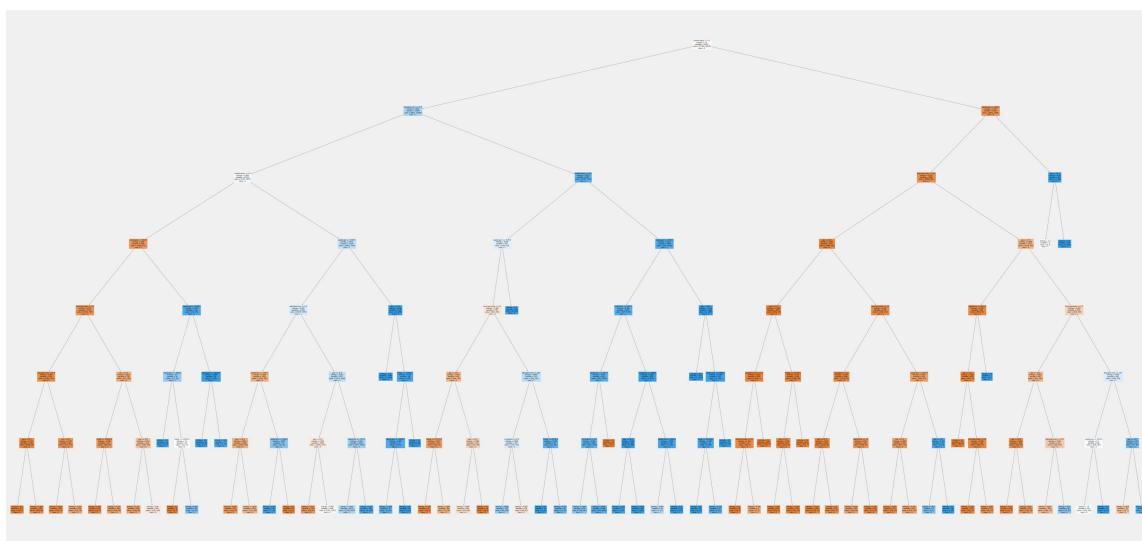
In [38]:

```
1 cm = confusion_matrix(y_new_test, y_pred)
2 plt.figure(figsize=(6,4))
3 sns.heatmap(cm, annot=True, fmt='d', cmap="Blues")
4 plt.xlabel("Predicted Values")
5 plt.ylabel("Actual values")
6 plt.show()
```



In [39]:

```
1 plt.figure(figsize=(100,50))
2 plot_tree(tree, feature_names=x_new.columns, filled=True, class_names=['0', '1'])
3 plt.savefig('dec_tree.png')
```



In [40]:

```

1 model_eval_test.loc[1] = [recall_score(y_new_test,y_pred),precision_score(y_new_test
2 model_eval_test

```

Out[40]:

	recall	precision	f1_score	accuracy
0	0.967215	0.887254	0.925511	0.924146
1	0.817769	0.802842	0.810237	0.813373

Logistic_regression

```

1 Logistic Regression:
2
3 >> It is a Supervised ML algorithm which predicts categorical values based on
probabilities of target class/event.
4 >> It is a probabilistic algorithm(0 to 1)
5 >> It is a Classification algorithm.
6 >> It is a parametric algorithm
7 -----
8 >> Logistic regression is a simple and more efficient method for binary and linear
classification problems.
9 >> It is a classification model which is very easy to realize and achieves very
good performance with linearly separable
10 classes.
11 >> It is an extensively employed algorithm for classification in industry.
12 >> Logistic regression is famous because it can convert the values of logits (log-
odds), which can range from  $-\infty$  to  $+\infty$  to
13 a range between 0 and 1.
14 >> As logistic functions output the probability of occurrence of an event it can
be applied to many real-life scenarios.
15 >> It is for this reason that the logistic regression model is very popular.
Another reason why logistic fairs in comparison
16 to linear regression is that it is able to handle the categorical variables.
17 -----
18
19 Assumptions of logistic regression :
20
21 1. Target var shoulf be categorical in nature (binary, ordinal , multiclass)
22 2. Constructs Linear Boundaries
23 3. Requires large datasets
24
25 >> Some others:
26 4. No Multicolinearity
27 5. Independence
28
29 -----
30
31 Logistic Function (Sigmoid Function):
32
33 >>The sigmoid function is a special form of the logistic function and is usually
denoted by  $\sigma(x)$  or  $\text{sig}(x)$ .
34 It is given by:  $\sigma(x) = 1/(1+\exp(-x))$ 
35
36 >> A few properties include:

```

```
37 -Domain: (-∞, +∞)
38 -Range: (0, +1)
39 - $\sigma(0) = 0.5$ 
40 -The function is monotonically increasing.
41 -The function is continuous everywhere.
42 -The function is differentiable everywhere in its domain.
43 -Numerically, it is enough to compute this function's value over a small range of
  numbers,
44   e.g., [-10, +10].
45   For values less than -10, the function's value is almost zero.
46   For values greater than 10, the function's values are almost one.
47
48 >> The sigmoid function is a mathematical function used to map the predicted
  values to probabilities.
49 >> It maps any real value into another value within a range of 0 and 1.
50 >> The value of the logistic regression must be between 0 and 1, which cannot go
  beyond this limit, so it forms a curve
51   like the "S" form. The S-form curve is called the Sigmoid function or the
  logistic function.
52 >> In logistic regression we use the concept of the threshold value which defines
  the probability of either 0 or 1.
53   Such as values above the threshold value tends to 1, and a value below the
  threshold values tends to 0.
54 >> With the help of this sigmoid function, we can successfully predict the output
  in terms of a probability
```

In [41]:

```
1 log_reg = LogisticRegression()
2 log_reg.fit(x_new_train,y_new_train)
```

Out[41]:

LogisticRegression()

In a Jupyter environment, please rerun this cell to show the HTML representation or
trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page
with nbviewer.org.

In [42]:

```
1 y_pred = log_reg.predict(x_new_test)
```

In [43]:

```
1 print(model_eval(log_reg,x_new_train,y_new_train))
```

```
Confusion Matrix =  
[[16108 1196]  
 [11268 5170]]  
*****
```

```
Accuracy Score = 0.6306087368857803  
*****
```

```
Classification Report =  
precision recall f1-score support  
  
 0       0.59    0.93    0.72    17304  
 1       0.81    0.31    0.45    16438  
  
accuracy          0.63    33742  
macro avg       0.70    0.62    0.59    33742  
weighted avg     0.70    0.63    0.59    33742
```

```
precision_score = 0.8121269242852655
```

```
recall_score = 0.31451514782820295
```

```
f1_score = 0.4534292229433432
```

```
Model Eval Success
```

In [44]:

```
1 print(model_eval(log_reg,x_new_test,y_new_test))
```

```
Confusion Matrix =  
[[6890 526]  
 [4855 2191]]  
*****
```

```
Accuracy Score = 0.6279214493154474  
*****
```

```
Classification Report =  
precision recall f1-score support  
  
 0       0.59    0.93    0.72    7416  
 1       0.81    0.31    0.45    7046  
  
accuracy          0.63    14462  
macro avg       0.70    0.62    0.58    14462  
weighted avg     0.69    0.63    0.59    14462
```

```
precision_score = 0.8064041221935959
```

```
recall_score = 0.3109565711041726
```

```
f1_score = 0.44883744750588955
```

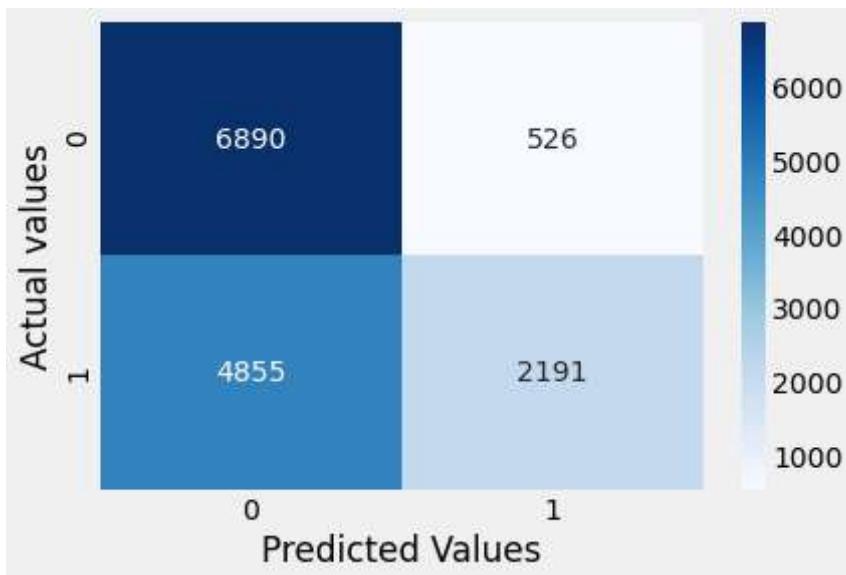
```
Model Eval Success
```

In [45]:

```

1 cm = confusion_matrix(y_new_test, y_pred)
2 plt.figure(figsize=(6,4))
3 sns.heatmap(cm, annot=True, fmt='d', cmap="Blues")
4 plt.xlabel("Predicted Values")
5 plt.ylabel("Actual values")
6 plt.show()

```



In [46]:

```

1 y_pred_train_prob = log_reg.predict_proba(x_new_train)
2 y_pred_train_prob

```

Out[46]:

```
array([[0.5443506 , 0.4556494 ],
       [0.51773696, 0.48226304],
       [0.54235258, 0.45764742],
       ...,
       [0.52969392, 0.47030608],
       [0.5674072 , 0.4325928 ],
       [0.59543042, 0.40456958]])
```

In [47]:

```

1 y_pred_train_prob_class1 = y_pred_train_prob[:,1]
2 y_pred_train_prob_class1

```

Out[47]:

```
array([0.4556494 , 0.48226304, 0.45764742, ..., 0.47030608, 0.4325928 ,
       0.40456958])
```

AUC-ROC curve

```

1 ROC :
2     >>The Receiver Operator Characteristic (ROC) curve is an evaluation metric for
3     binary classification problems.
4     >>It is a probability curve that plots the TPR against FPR at various
      threshold values and essentially separates

```

```

4     the ‘signal’ from the ‘noise.’
5     >>In other words, it shows the performance of a classification model at all
6     classification thresholds.
7 -----
8     Importance of ROC curve :
9
10    >> A useful tool when predicting the probability of a binary outcome is the
11     Receiver Operating Characteristic curve,
12     or ROC curve.
13    >> The curves of different models can be compared directly in general or for
14     different thresholds.
15    >> The area under the curve (AUC) can be used as a summary of the model skill.
16    >> The shape of the curve contains a lot of information, including what we might
     care about most for a problem,
     the expected false positive rate, and the false negative rate
17    >> An operator may plot the ROC curve for the final model and choose a threshold
     that gives a desirable balance
18     between the false positives and false negatives.

```

In [48]:

```

1 fpr,tpr,threshold = roc_curve(y_new_train, y_pred_train_prob_class1)
2 np.around(threshold,3)

```

Out[48]:

```
array([2.    , 1.    , 1.    , ..., 0.132, 0.124, 0.086])
```

In [49]:

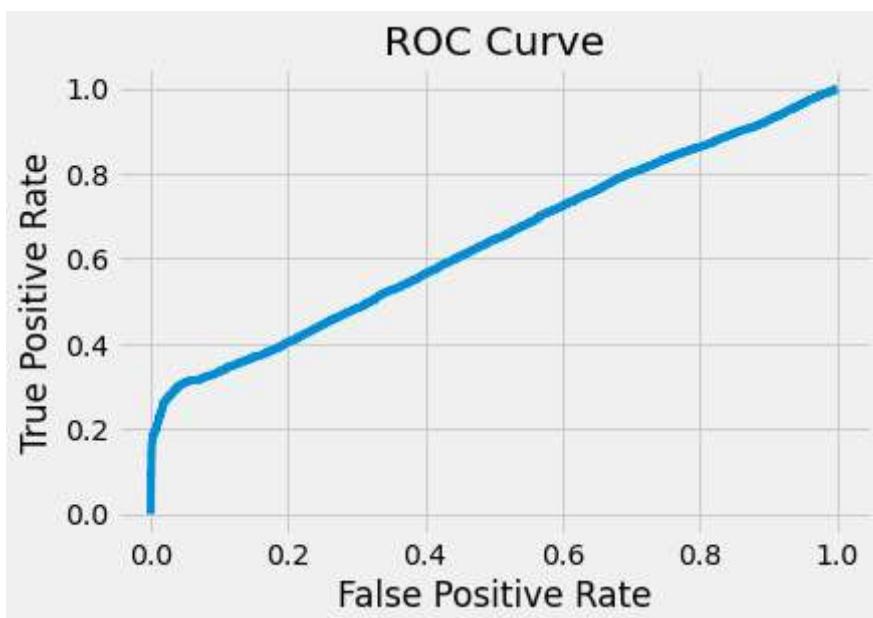
```

1 plt.plot(fpr,tpr)
2 plt.xlabel("False Positive Rate")
3 plt.ylabel("True Positive Rate")
4 plt.title("ROC Curve")

```

Out[49]:

```
Text(0.5, 1.0, 'ROC Curve')
```



```

1 >>The Receiver Operator Characteristic (ROC) curve is an evaluation metric for
     binary classification problems.

```

```

2 >>It is a probability curve that plots the TPR against FPR at various threshold
3 values and essentially separates the ‘signal’
4 from the ‘noise.’ In other words, it shows the performance of a classification
5 model at all classification thresholds.
6 >>The Area Under the Curve (AUC) is the measure of the ability of a binary
7 classifier to distinguish between classes and is
8 used as a summary of the ROC curve.
9 >>The higher the AUC, the better the model’s performance at distinguishing between
10 the positive and negative classes.
11 >>When AUC = 1, the classifier can correctly distinguish between all the Positive
12 and the Negative class points.
13 If, however, the AUC had been 0, then the classifier would predict all Negatives
14 as Positives and all Positives as Negatives.
15 >>When  $0.5 < \text{AUC} < 1$ , there is a high chance that the classifier will be able to
16 distinguish the positive class values from
17 the negative ones. This is so because the classifier is able to detect more
18 numbers of True positives and True negatives
19 than False negatives and False positives
20 >>When  $\text{AUC} = 0.5$ , then the classifier is not able to distinguish between Positive
21 and Negative class points.
22 Meaning that the classifier either predicts a random class or a constant class
23 for all the data points.
24
25 So, the higher the AUC value for a classifier, the better its ability to
26 distinguish between positive and negative classes.

```

In [50]:

```
1 roc_auc_score(y_new, log_reg.predict_proba(x_new)[:,1])
```

Out[50]:

0.6353182858025213

In [51]:

```
1 roc_auc_score(y_new, log_reg.decision_function(x_new))
```

Out[51]:

0.6353182858025213

Here the value for roc_auc_curve is not much better, is far from 1. so regarding to assumptions of roc_auc the classifier is not able to distinguish between Positive and Negative class points. Meaning that the classifier either predicts a random class or a constant class for all the data points.

In [52]:

```
1 model_eval_test.loc[2] = [recall_score(y_new_test,y_pred),precision_score(y_new_test
2 model_eval_test
```

Out[52]:

	recall	precision	f1_score	accuracy
0	0.967215	0.887254	0.925511	0.924146
1	0.817769	0.802842	0.810237	0.813373
2	0.310957	0.806404	0.448837	0.627921

K-Nearest Neighbors

```
1 The k-nearest neighbors algorithm, also known as KNN or k-NN, is a non-parametric,
2 supervised learning classifier, which uses proximity to make classifications or
3 predictions about the grouping of an individual data point.
4
5 The K-Nearest Neighbor algorithm is a supervised learning algorithm that can be
6 used for both classification and regression tasks. The algorithm works by finding
7 the K nearest neighbors to a given data point, and then using those neighbors to
8 predict the class or value of the data point.
9
10 It uses data in which there is a target column present i.e., labelled data to model
11 a function to produce an output for the unseen data. It uses the euclidean
12 distance formula to compute the distance between the data points for
13 classification or prediction.
14
15 The main objective of this algorithm is that similar data points must be close to
16 each other so it uses the distance to calculate the similar points that are close
17 to each other.
18
19 >>> The following operations have happened during each iteration of the
20 algorithm. For each of the unseen or test data point, the kNN classifier must:
21
22 Step-1: Calculate the distances of test point to all points in the training set
23 and store them
24
25 Step-2: Sort the calculated distances in increasing order
26
27 Step-3: Store the K nearest points from our training dataset
28
29 Step-4: Calculate the proportions of each class
30
31 Step-5: Assign the class with the highest proportion
32
33 >>> Advantages:
34
35 1. No Training Period: It does not learn anything during the training period since
36 it does not find any discriminative function with the help of the training data.
37 In simple words, actually, there is no training period for the KNN algorithm. It
38 stores the training dataset and learns from it only when we use the algorithm for
39 making the real-time predictions on the test dataset. As a result, the KNN
40 algorithm is much faster than other algorithms which require training. For
41 Example, SupportVector Machines(SVMs), Linear Regression, etc. Moreover, since the
42 KNN algorithm does not require any training before making predictions as a result
43 new data can be added seamlessly without impacting the accuracy of the algorithm.
44
45 2. Easy to implement and understand: To implement the KNN algorithm, we need only
46 two
47 parameters i.e. the value of K and the distance metric(e.g. Euclidean or
48 Manhattan, etc.).
49 Since both the parameters are easily interpretable therefore they are easy to
50 understand.
51
52 Applications of KNN >>
53
54 1. Text mining
55 2. Agriculture
56 3. Finance
```

```
34 4. Medical
35 5. Facial recognition
36 6. Recommendation systems (Amazon, Hulu, Netflix, etc)
37
38 The various real-life applications of the KNN Algorithm includes:
39
40 1. KNN allows the calculation of the credit rating. By collecting the financial
   characteristics vs. comparing people having similar financial features to a
   database we can calculate the same. Moreover, the very nature of a credit rating
   where people who have similar financial details would be given similar credit
   ratings also plays an important role. Hence the existing database can then be used
   to predict a new customer's credit rating, without having to perform all the
   calculations.
41
42 2. In political science: KNN can also be used to predict whether a potential voter
   "will vote" or "will not vote", or to "vote Democrat" or "vote Republican" in an
   election. Apart from the above-mentioned use cases, KNN algorithms are also used
   for handwriting detection (like OCR), Image recognition, and video recognition.
43
```

In [53]:

```
1 x_new_train, x_new_test, y_new_train,y_new_test = train_test_split(x_new,y_new,test_
2 x_new_train.shape, x_new_test.shape, y_new_train.shape, y_new_test.shape
```

Out[53]:

```
((33742, 13), (14462, 13), (33742,), (14462,))
```

In [54]:

```

1 normal_scaler = MinMaxScaler() # x_norm = x - x_min / x_max-x_min
2 normal_scaler.fit(x_new) ### x_min && x_max, Normalization
3 normalized_x = normal_scaler.transform(x_new) ### Apply the transformed values and r
4 x_norm = pd.DataFrame(normalized_x,columns=x_new.columns)
5 x_norm

```

Out[54]:

	age	workclass	fnlwgt	education-num	marital-status	occupation	relationship	race	≤
0	0.301370	0.0	0.044302	0.800000	0.666667	0.071429	0.2	1.0	
1	0.452055	0.0	0.048238	0.800000	0.333333	0.285714	0.0	1.0	
2	0.287671	1.0	0.138113	0.533333	0.000000	0.428571	0.2	1.0	
3	0.493151	1.0	0.151068	0.400000	0.333333	0.428571	0.0	0.0	
4	0.150685	1.0	0.221488	0.800000	0.333333	0.714286	1.0	0.0	
...
48199	0.410959	1.0	0.099157	0.866667	0.333333	0.285714	0.0	1.0	
48200	0.438356	0.0	0.077157	0.800000	0.333333	0.285714	0.0	1.0	
48201	0.383562	1.0	0.075655	0.533333	0.333333	0.857143	0.0	1.0	
48202	0.232877	1.0	0.098791	0.800000	0.333333	0.285714	0.0	1.0	
48203	0.602740	0.0	0.020818	0.800000	0.666667	0.000000	0.2	1.0	

48204 rows × 13 columns

In [55]:

```

1 x_norm_train, x_norm_test, y_new_train,y_new_test = train_test_split(x_norm,y_new,test_size=0.2,random_state=42)

```

In [56]:

```

1 knn_clf_model = KNeighborsClassifier(n_neighbors=5, p=2)
2 knn_clf_model.fit(x_new_train,y_new_train) ## data will be loaded ### LAZY Learner

```

Out[56]:

KNeighborsClassifier()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [57]:

```

1 # Training Data model Evaluation
2 model_eval(knn_clf_model,x_new_train,y_new_train)

```

Confusion Matrix =
[[13054 4250]
 [1956 14482]]

Accuracy Score = 0.8160749214628653

Classification Report =

	precision	recall	f1-score	support
0	0.87	0.75	0.81	17304
1	0.77	0.88	0.82	16438
accuracy			0.82	33742
macro avg	0.82	0.82	0.82	33742
weighted avg	0.82	0.82	0.82	33742

precision_score = 0.7731155242366005
recall_score = 0.8810074218274729
f1_score = 0.8235427921524027

Out[57]:

'Model Eval Success'

In [58]:

```

1 # Testing Data model Evaluation
2 model_eval(knn_clf_model,x_new_test,y_new_test)

```

Confusion Matrix =
[[4859 2557]
 [1575 5471]]

Accuracy Score = 0.7142857142857143

Classification Report =

	precision	recall	f1-score	support
0	0.76	0.66	0.70	7416
1	0.68	0.78	0.73	7046
accuracy			0.71	14462
macro avg	0.72	0.72	0.71	14462
weighted avg	0.72	0.71	0.71	14462

precision_score = 0.6814897857498754
recall_score = 0.7764689185353392
f1_score = 0.7258856308876209

Out[58]:

'Model Eval Success'

Hyperparameter tuning :

```
1 The two best strategies for Hyperparameter tuning are: GridSearchCV.  
RandomizedSearchCV.  
2  
3 >>> GridSearchCV :  
4  
5 In GridSearchCV approach, the machine learning model is evaluated for a range of  
hyperparameter values. This approach is called GridSearchCV, because it searches  
for the best set of hyperparameters from a grid of hyperparameters values.  
6  
7 For example, if we want to set two hyperparameters C and Alpha of the Logistic  
Regression Classifier model, with different sets of values. The grid search  
technique will construct many versions of the model with all possible combinations  
of hyperparameters and will return the best one.  
8  
9 >>> RandomizedSearchCV :  
10  
11 RandomizedSearchCV solves the drawbacks of GridSearchCV, as it goes through only a  
fixed number of hyperparameter settings. It moves within the grid in a random  
fashion to find the best set of hyperparameters. This approach reduces unnecessary  
computation.
```

In [59]:

```
1 knn_model = KNeighborsClassifier(n_neighbors=5, p = 2)  
2  
3 hyp = {  
4     "n_neighbors":np.arange(3,20),  
5     "p": [1,2]  
6 }  
7  
8 gscv_knn = GridSearchCV(knn_model,hyp,cv= 5)  
9 gscv_knn.fit(x_new_train,y_new_train)  
10 gscv_knn.best_estimator_
```

Out[59]:

KNeighborsClassifier(n_neighbors=3, p=1)

In a Jupyter environment, please rerun this cell to show the HTML representation or
trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page
with nbviewer.org.

In [60]:

```
1 knn_hyp = gscv_knn.best_estimator_  
2 knn_hyp
```

Out[60]:

KNeighborsClassifier(n_neighbors=3, p=1)

In a Jupyter environment, please rerun this cell to show the HTML representation or
trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page
with nbviewer.org.

In [61]:

```
1 y_pred = knn_hyp.predict(x_new_test)
```

In [62]:

```
1 # Training Data evalaution
2 model_eval(knn_hyp,x_new_train,y_new_train)
```

```
Confusion Matrix =
[[14148  3156]
 [ 701 15737]]
*****
```

```
Accuracy Score = 0.8856914231521545
*****
Classification Report =
```

	precision	recall	f1-score	support
0	0.95	0.82	0.88	17304
1	0.83	0.96	0.89	16438
accuracy			0.89	33742
macro avg	0.89	0.89	0.89	33742
weighted avg	0.89	0.89	0.89	33742

```
precision_score = 0.8329540041285132
recall_score = 0.9573549093563694
f1_score = 0.8908324134612663
```

Out[62]:

```
'Model Eval Success'
```

In [63]:

```

1 # Testing data evaluation
2 model_eval(knn_hyp,x_new_test,y_new_test)

```

```

Confusion Matrix =
[[5072 2344]
 [1048 5998]]
*****
Accuracy Score = 0.7654542940118932
*****
Classification Report =
      precision    recall   f1-score   support
          0       0.83     0.68     0.75     7416
          1       0.72     0.85     0.78     7046
          accuracy           0.77
          macro avg       0.77     0.77     0.76     14462
          weighted avg     0.78     0.77     0.76     14462

precision_score = 0.719012227283625
recall_score = 0.8512631280158955
f1_score = 0.7795684949311151

```

Out[63]:

'Model Eval Success'

In [64]:

```

1 cm = confusion_matrix(y_new_test, y_pred)
2 plt.figure(figsize=(6,4))
3 sns.heatmap(cm, annot=True, fmt='d', cmap="Blues")
4 plt.xlabel("Predicted Values")
5 plt.ylabel("Actual values")
6 plt.show()

```



```

1 four different distance functions, which are
2 Euclidean distance,
3 cosine similarity measure,
4 Minkowsky,
5 correlation and

```

6 Chi square are used in the k-NN classifier respectively

To finding the different values of k and p

In [65]:

```
1 distanc_param=1
2 train_acc_list = []
3 test_acc_list = []
4 k_range = np.arange(3,20)
5 for k in k_range:
6     model = KNeighborsClassifier(n_neighbors=k,p=distanc_param)
7     model.fit(x_new_train,y_new_train)
8     # Training Acc
9     train_acc = model.score(x_new_train,y_new_train)
10    train_acc_list.append(train_acc)
11    # Testing Acc
12    test_acc_list.append(model.score(x_new_test,y_new_test))
13 print(train_acc_list)
14 print(test_acc_list)
```

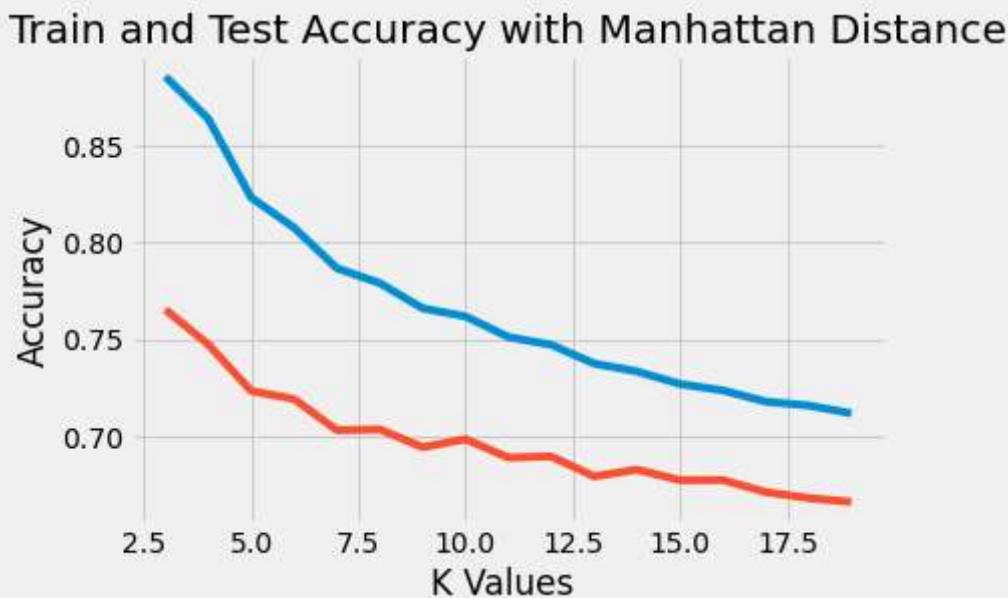
```
[0.8856914231521545, 0.8638195720467073, 0.8232173552249422, 0.80768774820
69824, 0.7867346333945824, 0.7790587398494458, 0.7660482484737123, 0.76189
91168276925, 0.7512299211664987, 0.747228972793551, 0.737508150080019, 0.7
334775650524569, 0.7269871376918974, 0.7236974690296959, 0.717918321379882
7, 0.7159030288661016, 0.7118428071839251]
[0.7654542940118932, 0.7474069976490112, 0.7234130825611949, 0.71926427879
96127, 0.7031530908588024, 0.7035679712349606, 0.6943714562301203, 0.69852
02599917024, 0.6889780113400636, 0.6896003319043009, 0.6792974692297055,
0.6828239524270502, 0.6773613608076338, 0.6774305075369935, 0.671207301894
6204, 0.6680956990734338, 0.6660904439220025]
```

In [66]:

```

1 plt.plot(k_range,train_acc_list)
2 plt.plot(k_range,test_acc_list)
3 plt.xlabel("K Values ")
4 plt.ylabel("Accuracy")
5 plt.title("Train and Test Accuracy with Manhattan Distance")
6 plt.savefig("Accuracy_MD.png")

```



In [67]:

```

1 print(train_acc_list[13])
2 test_acc_list[13]

```

0.7236974690296959

Out[67]:

0.6774305075369935

In [68]:

```

1 model_eval_test.loc[3] = [recall_score(y_new_test,y_pred),precision_score(y_new_test
2 model_eval_test

```

Out[68]:

	recall	precision	f1_score	accuracy
0	0.967215	0.887254	0.925511	0.924146
1	0.817769	0.802842	0.810237	0.813373
2	0.310957	0.806404	0.448837	0.627921
3	0.851263	0.719012	0.779568	0.765454

Adaboost

```
1 AdaBoost is the first designed boosting algorithm with a particular loss function.
```

```
2
3 -----
4
5 boosting also creates an ensemble of classifiers by resampling the data, which are
then combined by majority voting. However, in boosting, resampling is
strategically geared to provide the most informative training data for each
consecutive classifier. In essence, each iteration of boosting creates three weak
classifiers:
6
7 the first classifier C1 is trained with a random subset of the available training
data. The training data subset for the second classifier C2 is chosen as the most
informative subset, given C1 .
8
9 Specifically, C2 is trained on a training data only half of which is correctly
classified by C1 , and the other half is misclassified.
10
11 The third classifier C3 is trained with instances on which C1 and C2 disagree. The
three classifiers are combined through a three-way majority vote.
12
13 -----
14
15 Boosting has been a prevalent technique for tackling binary classification
problems. These algorithms improve the prediction power by converting a number of
weak learners to strong learners.
16
17 -----
18
19 AdaBoost is an ensemble learning method (also known as “meta-learning”) which was
initially created to increase the efficiency of binary classifiers. AdaBoost uses
an iterative approach to learn from the mistakes of weak classifiers, and turn
them into strong ones.
20
21 1. Assisgn The sample weights to all datapoints
22     Sample weights = 1/ Total no of samples
23
24 2. Construct Decision Stump
25
26 3. Calculate TOTAL ERROR
27     TE = no of misclassified samples / total no of samples
28
29 4. Calculate the performance
30     performacnce = 1/2 * log((1-TE)/TE)
31
32     importance, Amount of Say, Siginificance, influence etc
33
34 5. Calculate New Sample Weights
35     NSW = old Sample Weight * e^(+/-performance)
36
37     NSW correct = old Sample Weight * e ^(-performace)
38     NSW Incorrect = old Sample Weight * e ^(+performace)
39
40 6. TOTAL weight = (NSW correct * No correct Samples) + (NSW Incorrect * No
incorrect Samples)
41
42 7. Normalization of NSW
43     Norm_NSW = NSW / Total weight
44
45     Norm_NSW_correct = NSW correct / Total weight
46     Norm_NSW_Incorrect = NSW Incorrect / Total weight
47
```

```
48 8. Create a buckets using normalized weights. [0 to 1]
49
50 9. create a random array between 0 to 1, create new dataset by picking values from
buckets
51
52 10. Repeat the process 1 to 9 until same result previous stump or we reach to
n_estimator = max_limit
53   n_estimator = 50
54
55 we go through this procedure, and if classification problem is arise then by using
accuracy_score, confusion_matrix, classification report is used for evaluation and
for regression mean_squared_error, mean_absolute_error, r2_score is used for
evaluation.
```

In [69]:

```
1 ada_model = AdaBoostClassifier(base_estimator=None,
2                               n_estimators=50,
3                               learning_rate=1.0,
4                               random_state=5)
5 ada_model.fit(x_new_train,y_new_train)
```

Out[69]:

AdaBoostClassifier(random_state=5)

**In a Jupyter environment, please rerun this cell to show the HTML representation or
trust the notebook.**

**On GitHub, the HTML representation is unable to render, please try loading this page
with nbviewer.org.**

In [70]:

```
1 y_pred = ada_model.predict(x_new_test)
```

In [71]:

```

1 # training data Evaluation
2
3 model_eval(ada_model,x_new_train,y_new_train)

```

```

Confusion Matrix =
[[14081 3223]
 [ 2692 13746]]
*****
Accuracy Score = 0.8246991879556635
*****
Classification Report =
      precision    recall   f1-score   support
0          0.84     0.81     0.83     17304
1          0.81     0.84     0.82     16438
               accuracy           0.82     33742
macro avg       0.82     0.82     0.82     33742
weighted avg    0.83     0.82     0.82     33742

precision_score = 0.8100654134009075
recall_score = 0.8362331183842316
f1_score = 0.8229412997276019

```

Out[71]:

'Model Eval Success'

In [72]:

```

1 # testing data Evaluation
2
3 model_eval(ada_model,x_new_test,y_new_test)

```

```

Confusion Matrix =
[[6043 1373]
 [1125 5921]]
*****
Accuracy Score = 0.8272714700594662
*****
Classification Report =
      precision    recall   f1-score   support
0          0.84     0.81     0.83     7416
1          0.81     0.84     0.83     7046
               accuracy           0.83     14462
macro avg       0.83     0.83     0.83     14462
weighted avg    0.83     0.83     0.83     14462

precision_score = 0.8117630929531121
recall_score = 0.8403349418109566
f1_score = 0.8258019525801954

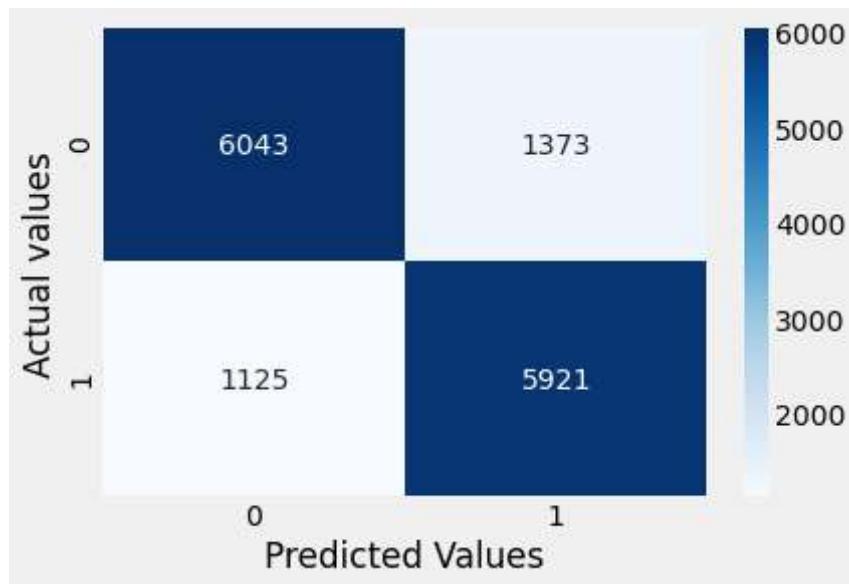
```

Out[72]:

'Model Eval Success'

In [73]:

```
1 cm = confusion_matrix(y_new_test, y_pred)
2 plt.figure(figsize=(6,4))
3 sns.heatmap(cm, annot=True, fmt='d', cmap="Blues")
4 plt.xlabel("Predicted Values")
5 plt.ylabel("Actual values")
6 plt.show()
```



Adaboost with different max_depth Decision tree

In [74]:

```

1 train_acc_list = []
2 test_acc_list = []
3 for i in range(1,10):
4     base_model = DecisionTreeClassifier()
5     base_model = DecisionTreeClassifier(max_depth=i,random_state=7)
6     model = AdaBoostClassifier(base_estimator=base_model,
7                               n_estimators=50,
8                               learning_rate=1.0,
9                               random_state=10)
10
11     model.fit(x_new_train,y_new_train)
12
13     train_acc = model.score(x_new_train,y_new_train)
14     train_acc_list.append(train_acc)
15
16     test_acc = model.score(x_new_test,y_new_test)
17     test_acc_list.append(test_acc)
18
19     print(f"Iteration >> {i}, Training Acc = {train_acc}, Testing Acc = {test_acc}")

```

Iteration >> 1, Training Acc = 0.8246991879556635, Testing Acc = 0.8272714
 700594662
 Iteration >> 2, Training Acc = 0.8399027917728646, Testing Acc = 0.8417922
 832250034
 Iteration >> 3, Training Acc = 0.8508683539802027, Testing Acc = 0.8467708
 47738902
 Iteration >> 4, Training Acc = 0.862545195898287, Testing Acc = 0.84843036
 92435348
 Iteration >> 5, Training Acc = 0.883498310710687, Testing Acc = 0.85562162
 90969437
 Iteration >> 6, Training Acc = 0.9184102898464821, Testing Acc = 0.8715945
 235790347
 Iteration >> 7, Training Acc = 0.9494102305731729, Testing Acc = 0.8857696
 030977735
 Iteration >> 8, Training Acc = 0.9836109299982218, Testing Acc = 0.9058221
 546120868
 Iteration >> 9, Training Acc = 0.9999110900361567, Testing Acc = 0.9131517
 079242152

Rather than being a model in itself, AdaBoost can be applied on top of any classifier to learn from its shortcomings and propose a more accurate model. It is usually called the “best out-of-the-box classifier” for this reason

In [75]:

```
1 model_eval_test.loc[4] = [recall_score(y_new_test,y_pred),precision_score(y_new_test
2 model_eval_test
```

Out[75]:

	recall	precision	f1_score	accuracy
0	0.967215	0.887254	0.925511	0.924146
1	0.817769	0.802842	0.810237	0.813373
2	0.310957	0.806404	0.448837	0.627921
3	0.851263	0.719012	0.779568	0.765454
4	0.840335	0.811763	0.825802	0.827271

In [76]:

```
1 model_eval_test.index = ['Random_forest','Decision_tree','Logistic_regression','KNN']
```

In [77]:

```
1 model_eval_test
```

Out[77]:

	recall	precision	f1_score	accuracy
Random_forest	0.967215	0.887254	0.925511	0.924146
Decision_tree	0.817769	0.802842	0.810237	0.813373
Logistic_regression	0.310957	0.806404	0.448837	0.627921
KNN	0.851263	0.719012	0.779568	0.765454
Adaboost	0.840335	0.811763	0.825802	0.827271

Our dataset is balanced, so we conclude our final assumption by considering accuracies of algorithms.

Here, random forest algorithm gives better accuracy rather than remaining algorithms. so we select random forest algorithm as a final model with 92.41% accuracy.