

Problem Statement : Customer churn prediction refers to the process of identifying customers who are likely to stop using a product or service in the near future. It is a valuable predictive analytics technique used by businesses to forecast customer behavior and take proactive measures to retain customers.

Objective : objective of this project is to predict wather customer is about to churn or not.

Kaggle Dataset Link : <https://www.kaggle.com/datasets/blastchar/telco-customer-churn>  
(<https://www.kaggle.com/datasets/blastchar/telco-customer-churn>).

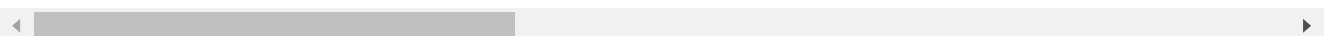
```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import pickle
5 from matplotlib import pyplot as plt
6 import scipy
7 from sklearn.model_selection import train_test_split, RandomizedSearchCV
8 from sklearn.preprocessing import LabelEncoder
9 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
10 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
11 from sklearn.feature_selection import SelectKBest
12 from collections import Counter
13 from imblearn.combine import SMOTEENN
14 plt.style.use('default')
15 import warnings
16 warnings.filterwarnings("ignore")
```

```
In [2]: 1 #import the dataset
2 df=pd.read_csv(r"Telco-Customer-Churn.csv")
3 df.head()
```

```
Out[2]:
```

	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetS
0	7590-VHVEG	Female	0	Yes	No	1	No	No phone service	
1	5575-GNVDE	Male	0	No	No	34	Yes	No	
2	3668-QPYBK	Male	0	No	No	2	Yes	No	
3	7795-CFOCW	Male	0	No	No	45	No	No phone service	
4	9237-HQITU	Female	0	No	No	2	Yes	No	Fibr

5 rows × 21 columns



```
In [3]: 1 #print concise summary of the dataset
        2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7043 entries, 0 to 7042
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   customerID            7043 non-null   object
1   gender                 7043 non-null   object
2   SeniorCitizen          7043 non-null   int64
3   Partner                7043 non-null   object
4   Dependents             7043 non-null   object
5   tenure                 7043 non-null   int64
6   PhoneService           7043 non-null   object
7   MultipleLines           7043 non-null   object
8   InternetService        7043 non-null   object
9   OnlineSecurity         7043 non-null   object
10  OnlineBackup           7043 non-null   object
11  DeviceProtection       7043 non-null   object
12  TechSupport            7043 non-null   object
13  StreamingTV            7043 non-null   object
..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..  ..
```

```
In [4]: 1 #check for missing values
        2 df.isnull().sum()
```

```
Out[4]: customerID            0
gender                        0
SeniorCitizen                 0
Partner                       0
Dependents                    0
tenure                        0
PhoneService                  0
MultipleLines                 0
InternetService               0
OnlineSecurity                0
OnlineBackup                  0
DeviceProtection              0
TechSupport                   0
StreamingTV                   0
StreamingMovies               0
Contract                      0
PaperlessBilling              0
PaymentMethod                 0
MonthlyCharges                0
TotalCharges                  0
Churn                         0
dtype: int64
```

```
In [5]: 1 #check for duplicate records
        2 df[df.duplicated()].shape[0]
```

```
Out[5]: 0
```

```
In [6]: 1 #check datatype
        2 df.dtypes
```

```
Out[6]: customerID      object
gender      object
SeniorCitizen  int64
Partner      object
Dependents   object
tenure       int64
PhoneService  object
MultipleLines object
InternetService object
OnlineSecurity object
OnlineBackup  object
DeviceProtection object
TechSupport   object
StreamingTV   object
StreamingMovies object
Contract      object
PaperlessBilling object
PaymentMethod object
MonthlyCharges float64
TotalCharges  object
Churn         object
dtype: object
```

```
In [7]: 1 #since customerId is not required for prediction so drop it
        2 df.drop('customerId',axis=1,inplace=True)
```

```
In [8]: 1 #since total charges is having numerical value but dtype is object to change it i
        2 df['TotalCharges']=pd.to_numeric(df['TotalCharges'],errors='coerce')
```

```
In [9]: 1 #print last 5 records of the dataset
        2 df.tail(5)
```

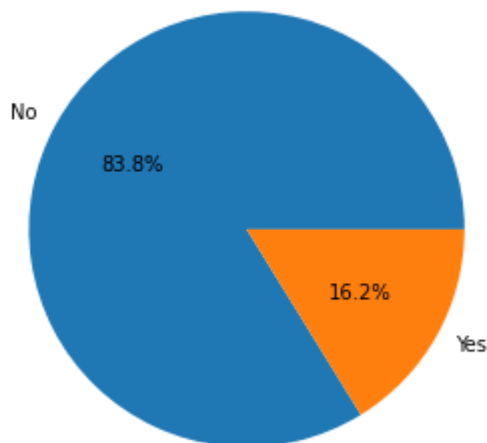
```
Out[9]:
```

	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	O
7038	Male	0	Yes	Yes	24	Yes	Yes	DSL	
7039	Female	0	Yes	Yes	72	Yes	Yes	Fiber optic	
7040	Female	0	Yes	Yes	11	No	No phone service	DSL	
7041	Male	1	Yes	No	4	Yes	Yes	Fiber optic	
7042	Male	0	No	No	66	Yes	No	Fiber optic	

## Exploratory Data Analysis :

In [10]:

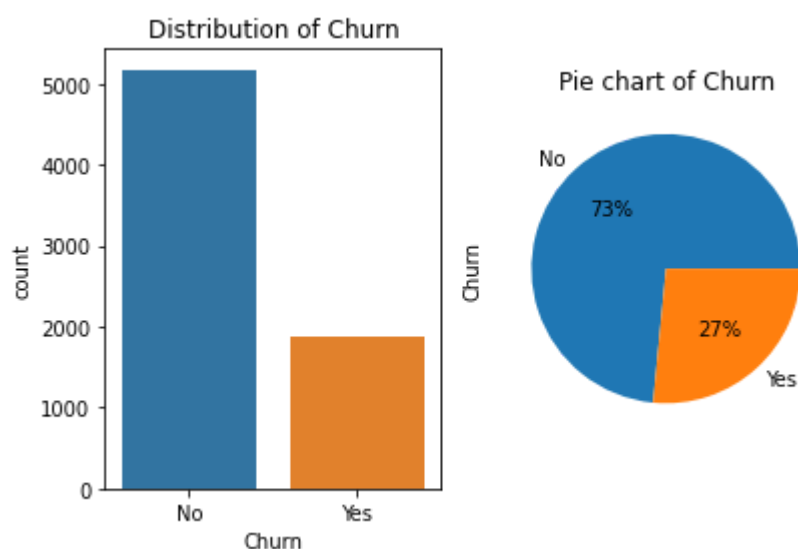
```
1 #pie chart to count senior citizen
2 plt.figure(figsize=(10,5))
3 plt.pie(df["SeniorCitizen"].value_counts(),autopct="%.1f%%",labels=["No", "Yes"])
4 plt.show()
```



as we can see 83.8 % of the customers are senior citizen and only 16.2% are adult customer.

In [11]:

```
1 #check the distribution of churn class
2 plt.subplot(121)
3 sns.countplot(data=df,x="Churn")
4 plt.title("Distribution of Churn")
5 plt.subplot(122)
6 df['Churn'].value_counts().plot(kind='pie',autopct="%1.f%%",labels=['No', 'Yes'])
7 plt.title('Pie chart of Churn')
8 plt.tight_layout()
9 plt.show()
```



In [12]:

```
1 #perentage of each class sample distribution
2 print("Customer Churn : {}".format(np.round((len(df[df["Churn"]=="Yes"])/len(df))
3 print("Customer Not Churn : {}".format(np.round((len(df[df["Churn"]=="No"])/len(df))
```

Customer Churn : 26.54%

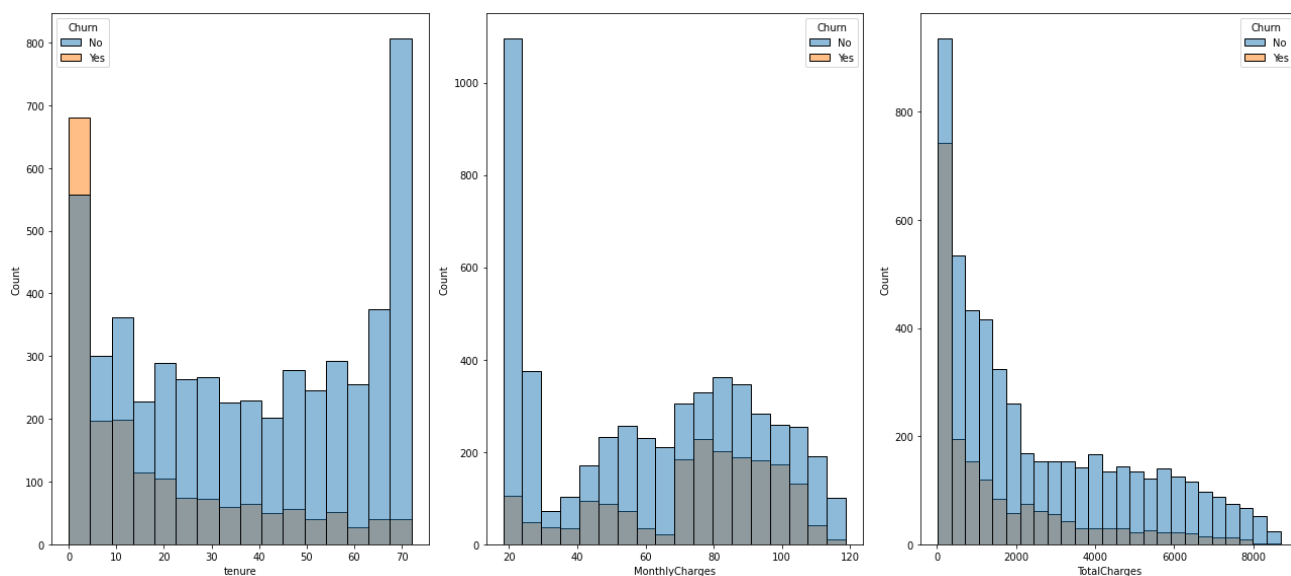
Customer Not Churn : 73.46%

Since our dataset is highly imbalance we need to balance before fitting it into model

```
In [13]: 1 #how much loss we are having because of customer churn
2 churn_customers=df[df["Churn"]=="Yes"]
3 loss=churn_customers["TotalCharges"].sum()
4 total_revenue=df["TotalCharges"].sum()
5 print("We have lost around {}$ due to customer churn".format(loss))
6 print("We have lost around {} percentage of revenue due to customer churn".form
```

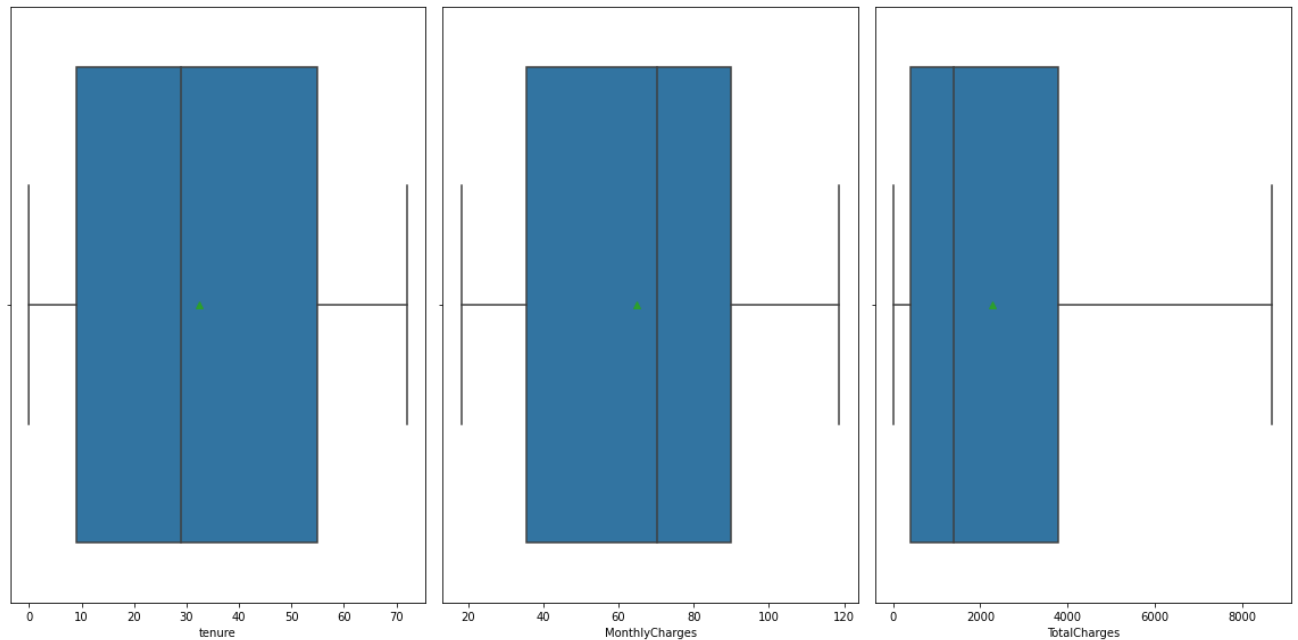
We have lost around 2862926.9\$ due to customer churn  
 We have lost around 17.83 percentage of revenue due to customer churn

```
In [14]: 1 #plot numerical features with histogram
2 fig,axs=plt.subplots(nrows=1,ncols=3,figsize=(18,8))
3 axes=axes.flatten()
4 num_columns=['tenure', 'MonthlyCharges', 'TotalCharges']
5 for i,col in enumerate(num_columns):
6     if(col!='SeniorCitizen'):
7         sns.histplot(x=col,data=df,hue='Churn',ax=axes[i])
8 fig.tight_layout()
9 plt.show()
```



In [15]:

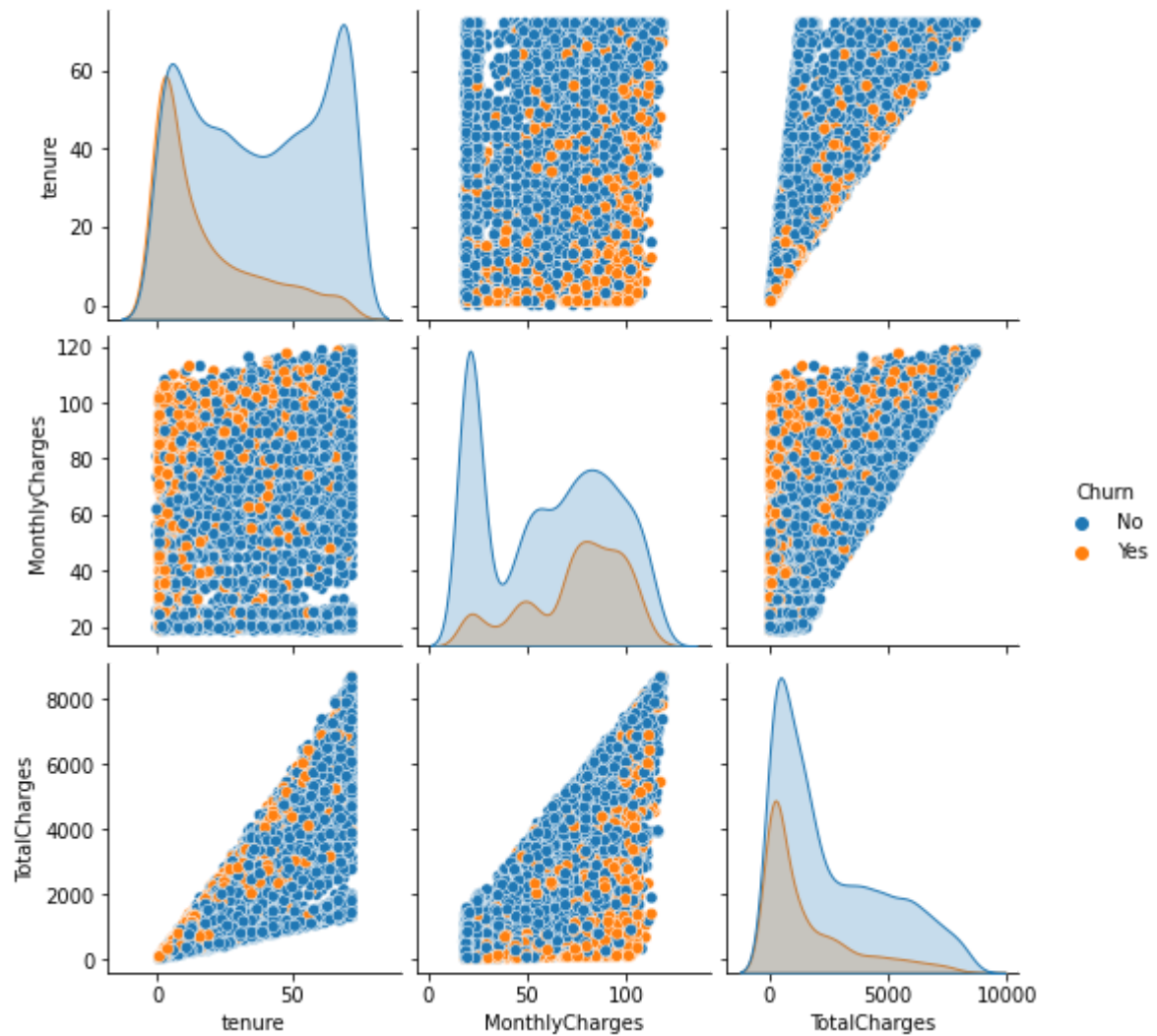
```
1 #plot numerical features with boxplot
2 fig,axs=plt.subplots(nrows=1,ncols=3,figsize=(16,8))
3 axes=axes.flatten()
4 num_columns=['tenure', 'MonthlyCharges', 'TotalCharges']
5 for i,col in enumerate(num_columns):
6     if(col!='SeniorCitizen'):
7         sns.boxplot(x=col,data=df,showmeans=True,ax=axes[i])
8 fig.tight_layout()
9 plt.show()
```



after plotting histogram and boxplot we found that there is no outlier present in numeric dataset so we don't need to do any kind of outlier treatment.

In [16]:

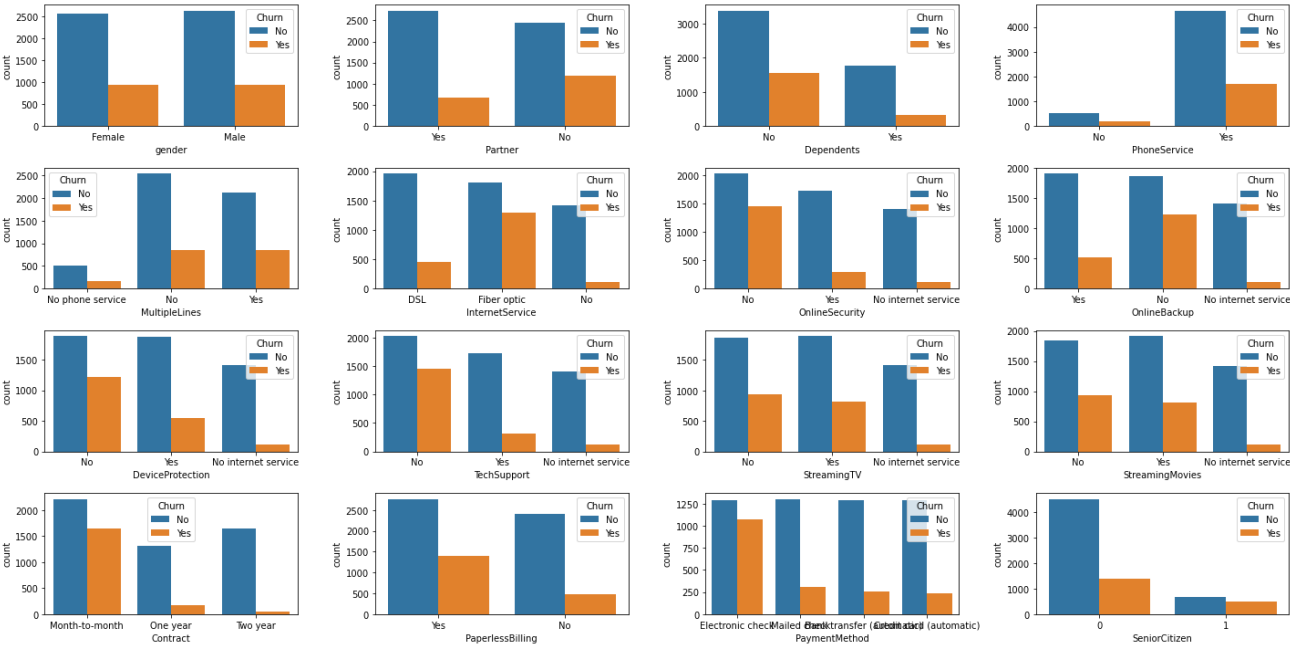
```
1 sns.pairplot(df.drop(columns="SeniorCitizen"),hue="Churn",kind="scatter")  
2 plt.show()
```



## Univariate Analysis

In [17]:

```
1 #plot categorical features :
2 cat_features=list(df.select_dtypes(include='object').columns)
3 cat_features.remove('Churn')
4 cat_features.append('SeniorCitizen')
5
6 fig,axs=plt.subplots(nrows=4,ncols=4,figsize=(20,10))
7 axes=axes.flatten()
8 for i,col in enumerate(cat_features):
9     sns.countplot(x=col,hue="Churn",data=df,ax=axes[i])
10 #adjust spacing between subplots
11 fig.tight_layout()
12 plt.show()
```



# Data Cleaning

In [18]:

```
1 df.head(5)
```

Out[18]:

	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	Online
0	Female	0	Yes	No	1	No	No phone service	DSL	
1	Male	0	No	No	34	Yes	No	DSL	
2	Male	0	No	No	2	Yes	No	DSL	
3	Male	0	No	No	45	No	No phone service	DSL	
4	Female	0	No	No	2	Yes	No	Fiber optic	



```
In [19]: 1 #check for null values
         2 df.isnull().sum()
```

```
Out[19]: gender                0
SeniorCitizen                0
Partner                      0
Dependents                   0
tenure                       0
PhoneService                 0
MultipleLines                0
InternetService              0
OnlineSecurity               0
OnlineBackup                 0
DeviceProtection             0
TechSupport                  0
StreamingTV                  0
StreamingMovies              0
Contract                     0
PaperlessBilling             0
PaymentMethod                0
MonthlyCharges               0
TotalCharges                 11
Churn                        0
dtype: int64
```

```
In [20]: 1 df["TotalCharges"].fillna(df["TotalCharges"].mean(),inplace=True)
```

```
In [21]: 1 df.isnull().sum().sum()
```

```
Out[21]: 0
```

```
In [22]: 1 #encoding categorical values into numeric using label encoder
         2 encoder=LabelEncoder()
         3 for feature in df.select_dtypes(include='object').columns:
         4     df[feature]=encoder.fit_transform(df[feature])
```

```
In [23]: 1 df.head()
```

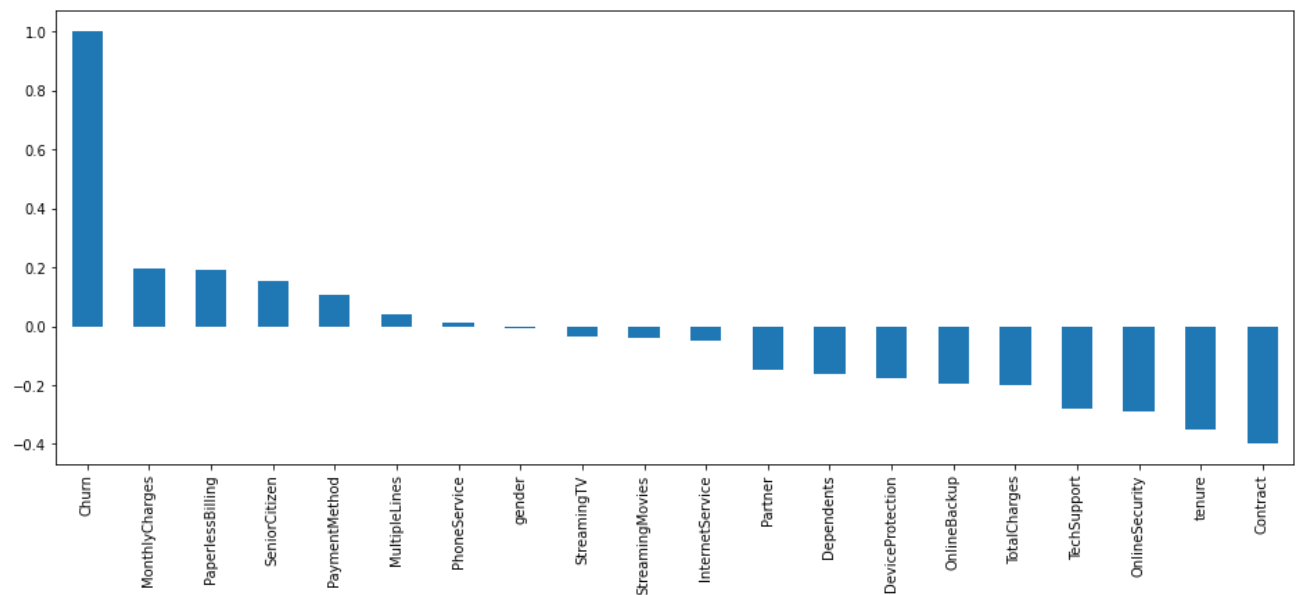
```
Out[23]:
```

	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	Online
0	0	0	1	0	1	0	1	0	
1	1	0	0	0	34	1	0	0	
2	1	0	0	0	2	1	0	0	
3	1	0	0	0	45	0	1	0	
4	0	0	0	0	2	1	0	1	

```
In [24]: 1 df.dtypes
```

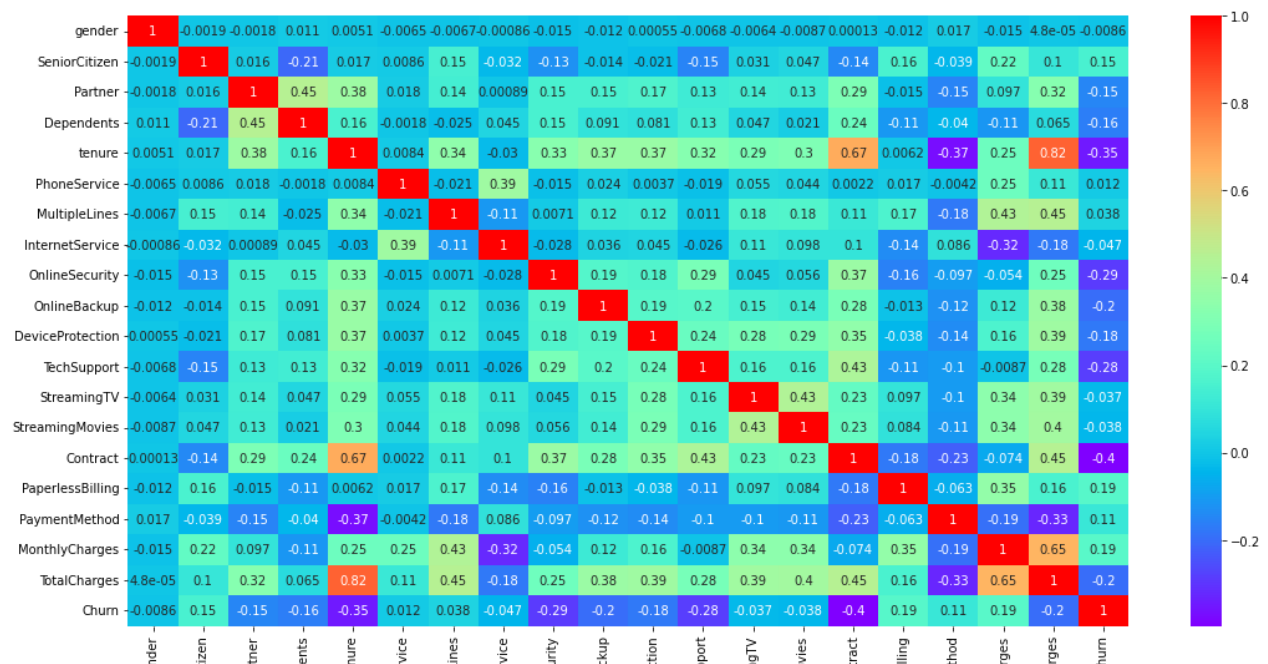
```
Out[24]: gender                int32
SeniorCitizen                int64
Partner                      int32
Dependents                   int32
tenure                       int64
PhoneService                 int32
MultipleLines                int32
InternetService              int32
OnlineSecurity               int32
OnlineBackup                 int32
DeviceProtection             int32
TechSupport                  int32
StreamingTV                  int32
StreamingMovies              int32
Contract                     int32
PaperlessBilling              int32
PaymentMethod                int32
MonthlyCharges                float64
TotalCharges                  float64
Churn                        int32
dtype: object
```

```
In [25]: 1 #get correlation of churn with other variables
2 plt.figure(figsize=(16,6))
3 df.corr()["Churn"].sort_values(ascending=False).plot(kind="bar")
4 plt.show()
```



In [26]:

```
1 plt.figure(figsize=(18,9))
2 sns.heatmap(df.corr(),annot=True,cmap="rainbow")
3 plt.show()
```



since we are using ensemble methods for model building so there is no need of feature scaling as its prediction is based on creating multiple decision tree

In [27]:

```
1 #seperating independent variables and target variable
2 x=df.drop("Churn",axis=1)
3 y=df["Churn"]
```

In [28]:

```
1 x.shape
```

Out[28]: (7043, 19)

## Feature Selection

selecting only 10 features which has higher correlation with churn

In [29]:

```
1 select_feature=SelectKBest(k=10) #no of features to be select
2 select_feature.fit(x,y)
```

Out[29]:

```
▼ SelectKBest
SelectKBest()
```

In [30]:

```
1 #Top 10 high correlated features
2 select_feature.get_feature_names_out()
```

Out[30]:

```
array(['Dependents', 'tenure', 'OnlineSecurity', 'OnlineBackup',
       'DeviceProtection', 'TechSupport', 'Contract', 'PaperlessBilling',
       'MonthlyCharges', 'TotalCharges'], dtype=object)
```

In [31]:

```
1 x=x[select_feature.get_feature_names_out()]
```

```
In [32]: 1 x.shape
```

```
Out[32]: (7043, 10)
```

according to the feature selection we have selected 10 top features out of 19 features

split data into training and validation set in 80:20 ratio

```
In [33]: 1 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=42)
```

```
In [34]: 1 x_train.shape,y_train.shape,x_test.shape,y_test.shape
```

```
Out[34]: ((5634, 10), (5634,), (1409, 10), (1409,))
```

```
In [35]: 1 #its imbalance dataset
2 y.value_counts()
```

```
Out[35]: 0    5174
1    1869
Name: Churn, dtype: int64
```

```
In [36]: 1 def evaluate_model_performance(model,test_data):
2     prediction=model.predict(test_data)
3     #print("Training Accuracy : ",model.score(x_train,y_train))
4     print("Validation Accuracy : {:.2f} %".format(accuracy_score(y_test,prediction)))
5     print("Precision Score : {:.2f} %".format(precision_score(y_test,prediction)))
6     print("Recall Score : {:.2f} %".format(recall_score(y_test,prediction)))
7     print("F1 Score : {:.2f} %".format(f1_score(y_test,prediction)))
8     print(classification_report(y_test,prediction))
```

```
In [37]: 1 #Random Forest Model without balancing dataset and without hyper paramter tuning
2 rand_forest=RandomForestClassifier()
3 rand_forest.fit(x_train,y_train)
```

```
Out[37]: ▼ RandomForestClassifier
RandomForestClassifier()
```

```
In [38]: 1 #measure the performance of random forest model
2 evaluate_model_performance(rand_forest,x_test)
```

Validation Accuracy : 0.79 %

Precision Score : 0.64 %

Recall Score : 0.47 %

F1 Score : 0.54 %

	precision	recall	f1-score	support
0	0.83	0.90	0.86	1036
1	0.64	0.47	0.54	373
accuracy			0.79	1409
macro avg	0.73	0.69	0.70	1409
weighted avg	0.78	0.79	0.78	1409

```
In [39]: 1 #GradientBoostingClassifier without balancing dataset and without hyper paramter
2 gbc_model=GradientBoostingClassifier( )
3 gbc_model.fit(x_train,y_train)
```

```
Out[39]: ▾ GradientBoostingClassifier
GradientBoostingClassifier()
```

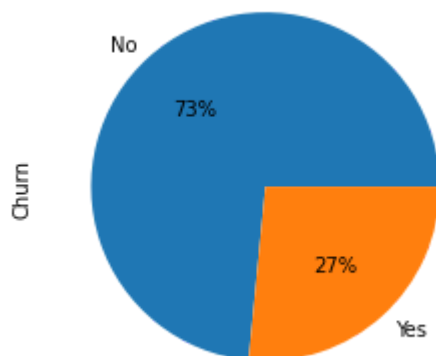
```
In [40]: 1 #measure the performance of GradientBoostingClassifier
2 evaluate_model_performance(gbc_model,x_test)
```

Validation Accuracy : 0.81 %  
Precision Score : 0.67 %  
Recall Score : 0.53 %  
F1 Score : 0.59 %

	precision	recall	f1-score	support
0	0.84	0.90	0.87	1036
1	0.67	0.53	0.59	373
accuracy			0.81	1409
macro avg	0.75	0.72	0.73	1409
weighted avg	0.80	0.81	0.80	1409

as we can see our model is not performing up to the mark because of imbalance nature of dataset so we will balance it to reduce TN,FN and increase TP,FP

```
In [41]: 1 plt.figure(figsize=(8,4))
2 y.value_counts().plot(kind="pie",autopct="%1.f%%",labels=['No', 'Yes'])
3 plt.show()
```



we have 2 classes class 0 and class 1. class 0 - majority class class 1 -minority class

```
In [42]: 1 smote=SMOTEENN()
2 x_st,y_st=smote.fit_resample(x,y)
```

```
In [43]: 1 y_st.value_counts().plot(kind="bar")
2 plt.title("target class distribution after under sampling")
3 plt.show()
```



```
In [44]: 1 y_st.value_counts()
```

```
Out[44]: 1    3101
0    2666
Name: Churn, dtype: int64
```

since we have performed SMOTEENN (combination of Smote + ENN) sampling method and we can see our dataset is nearly balanced

```
In [45]: 1 #now split training and validation set using balanced dataset
2 x_train,x_test,y_train,y_test=train_test_split(x_st,y_st,test_size=0.2,random_state=42)
```

```
In [46]: 1 x_train.shape,y_train.shape,x_test.shape,y_test.shape
```

```
Out[46]: ((4613, 10), (4613,)), (1154, 10), (1154,))
```

Building Model with Balanced Dataset and performance hyper parameter tuning using RandomSearchCV

```
In [47]: 1 param_grid={'n_estimators':[40,80,120,160,200],
2             'max_depth':[2,4,6,8,10],
3             'criterion':['gini'],
4             'random_state':[27,42,43]}
5         }
6 random_search_cv=RandomizedSearchCV( estimator=RandomForestClassifier(), param_grid=param_grid, cv=5)
7 random_search_cv.fit(x_train,y_train)
```

Fitting 5 folds for each of 12 candidates, totalling 60 fits

```
Out[47]: RandomizedSearchCV
estimator: RandomForestClassifier
RandomForestClassifier
```

```
In [48]: 1 random_search_cv.best_params_
```

```
Out[48]: {'random_state': 27, 'n_estimators': 160, 'max_depth': 10, 'criterion': 'gini'}
```

```
In [49]: 1 #Get final model with best param from RandomizedSearchCV
2 rf_final_model=random_search_cv.best_estimator_
```

```
In [50]: 1 #evaluate Random Forest Classifier
2 evaluate_model_performance(rf_final_model,x_test)
```

Validation Accuracy : 0.97 %

Precision Score : 0.95 %

Recall Score : 0.98 %

F1 Score : 0.97 %

	precision	recall	f1-score	support
0	0.98	0.95	0.96	541
1	0.95	0.98	0.97	613
accuracy			0.97	1154
macro avg	0.97	0.97	0.97	1154
weighted avg	0.97	0.97	0.97	1154

```
In [51]: 1 param_grid2 = {'n_estimators':[100, 150, 200, 250, 300],
2                  'criterion': ['friedman_mse', 'squared_error', 'mse', 'mae'],
3                  'max_depth': [2,4,6,8],
4                  'learning_rate': [0.001, 0.01, 0.1, 0.2],
5                  'loss': ['deviance', 'exponential']
6                  }
```

```
In [52]: 1 random_search_cv2=RandomizedSearchCV(estimator=GradientBoostingClassifier(random_
2 random_search_cv2.fit(x_train,y_train)
```

Fitting 5 folds for each of 12 candidates, totalling 60 fits

```
Out[52]: RandomizedSearchCV
estimator: GradientBoostingClassifier
GradientBoostingClassifier
```

```
In [53]: 1 random_search_cv2.best_params_
```

```
Out[53]: {'n_estimators': 250,
'max_depth': 6,
'loss': 'deviance',
'learning_rate': 0.2,
'criterion': 'mse'}
```

```
In [54]: 1 gb_final_model=random_search_cv2.best_estimator_
```

In [55]:

```
1 #evaluate final GradientBoostingClassifier Performance
2 evaluate_model_performance.gb_final_model,x_test)
```

Validation Accuracy : 0.97 %

Precision Score : 0.97 %

Recall Score : 0.98 %

F1 Score : 0.97 %

	precision	recall	f1-score	support
0	0.98	0.96	0.97	541
1	0.97	0.98	0.97	613
accuracy			0.97	1154
macro avg	0.97	0.97	0.97	1154
weighted avg	0.97	0.97	0.97	1154

Save Final Model Integration with application

In [56]:

```
1 file=open("trained_model.pkl","wb")
2 pickle.dump.gb_final_model,file)
3 file.close()
```

Conclusion : after balancing the dataset using smootenn and hyper paramter tuning model performance has increase and the highest f1 score we are getting is 97%.