# High Level Machine Learning Classification Project Life Cycle

- 1. Domain Introduction
- 2. Problem statement
- 3. Data Source
- 4. Data Description
- 5. Identify the target variable
- 6. Read the data
- 7. Inspect the data
  - Check few samples
  - Check the data types
  - Check the initial summary
- 8. Data Manipulation
  - Check for missing values
  - Column string fomatting
  - Data fomatting
  - Imputation
- 9. Exploratory Data Analysis
  - univariate analysis
  - class ditribution in data
  - Varibles distribution according to class
  - Bucketing
  - Correlation Matrix
  - feature elimination / addition / transformation
- 10. Data preprocessing
  - Encoding categorical variable
  - Normalizing features
  - spliting train/val/test data
  - feature compression ()
- 11. Model Building
  - Baseline Model
  - Model Selection
  - Hyper parameter Selection

# 1.Domain Introduction

We have the customer data for a **telecom** company which offers many services like phone, internet, TV Streaming and Movie Streaming.

# 2.Problem Statement

"Find the Best model to predict behavior to retain customers. You can analyze all relevant customer data and develop focused customer retention programs."

# 3. Data Source

Available at : IBM watson analytics page (https://community.watsonanalytics.com/wp-content/uploads/2015/03/WA_Fn-UseC_-Telco-Customer-Churn.csv?cm_mc_uid=14714377267115403444551&cm_mc_sid_50200000=12578191540344455127&cm_mc_sid_52640000=36692891540344455130)

# 4. Data Description

This data set provides info to help you predict behavior to retain customers. You can analyze all relevant customer data and develop focused customer retention programs.

A telecommunications company is concerned about the number of customers leaving their landline business for cable competitors. They need to understand who is leaving. **Imagine that you're an analyst at this company and you have to find out who is leaving and why.**

> The data set includes information about:
>
> - Customers who left within the last month – the column is called Churn
> - Services that each customer has signed up for – phone, multiple lines, internet, online security, online backup, device protection, tech support, and streaming TV and movies
> - Customer account information – how long they've been a customer, contract, payment method, paperless billing, monthly charges, and total charges
> - Demographic info about customers – gender, age range, and if they have partners and dependents

# 5. Identify the target variable

**The Goal is to predict whether or not a particular customer is likely to retain services.** This is represented by the Churn column in dataset. Churn=Yes means customer leaves the company, whereas Churn=No implies customer is retained by the company.

# 6. Read the data

In [ ]: ▶|

```python
1  # Step1 : importing the librabries like Numpy, Pandas, matplotlib, seaborn
2  #          Reading the Data from the CSV File
3
4  import numpy as np
5  import pandas as pd
6  import matplotlib.pyplot as plt
7  %matplotlib inline
8  import seaborn as sns
9
10 import warnings
11
12
13 # to remove the warnings coming in the output
14 if __name__ == '__main__':
15     warnings.filterwarnings(action='ignore', category=UserWarning)
16     warnings.filterwarnings(action='ignore', category=DeprecationWarning)
17
18 # laoding The Data From the File
19 df = pd.read_csv('WA_Fn-UseC_-Telco-Customer-Churn.csv',index_col='customerID')
20
21 # the shapeof DataSet, number of rows and number of Columns
22 print(df.shape)
23 print(df.size)
24
25 print(" The Telecom Company Data Set contains %d rows and features %d"%(df.shape[0],df.shape[1]))
```

## 7. Inspect the data

https://www.kaggle.com/blastchar/telco-customer-churn# (https://www.kaggle.com/blastchar/telco-customer-churn)

In [ ]: ▶|

```python
1  # Inspecting the Data of the Customer , seeing the fields of the DataSet
2  print('-'*80,"\n The Records of Customer of Telecom \n",'-'*80)
3  print(df.head(3))
```

In [ ]: ▶|

```python
1  #Step2: Analysing the customer Data
2
3  print('-'*80,"\n Analysing the Data by using functions like info, describe, checking for missing values \n",'-'*
4  # the datatypes of the fields
5  print(df.info())
6
7  # checking the mean,min, max values for numeric fields
8  print(df.describe())
9
10 # checking if there are any null values
11 print(df.isnull().sum())
12
13 print('-'*80,'''\n We observe that
14      1. The attributes are mostly object type(Except -Tenure,monthly charges,Senior Citizen)
15      2. Numeric Fields are-Senior Citizen with binary value(0,1), Tenure is maximum 72 months and avg-32 months
16         monthly charges maximum is 118 rupees with average 64 Rupees
17      3. Senior Citizen should be Categorical field
18      4. And No Null Values ''')
```

In [ ]: ▶|

```python
1  # Analysing  for non numeric fields
2
3  # printing the information on Nonnumerical Fields
4  print('-'*80,"\n Printing the non numerical fields , their unique values , \n",'-'*80)
5  print(df.describe(include=object))
6
7  print('-'*80,'''\n From Non-numerical Values , We observe that
8      1. Gender is mostly balanced , mean both male and female are customers
9      2. Maximum Customers have taken phone Service, Internet Service via Fiber optic is
10        more preferred
11     3. 50% Customers have opted for month-to-month billing , Max opted for paperless
12        billing and mostly Electronic cheque
13     4. Total Charges should be Float field , instead of Object Type
14     ''')
15
```

In [ ]: ▶|

```python
print('-'*80)
print(" Identifying the Unique Fields for the Categorical Fields\n",'-'*80)

# printing the Unique values for the Categorical/Object  Fields
# As the Data contains Total charges numerical , instead of Float type , having max-val=10
# so that such columns are not considered.

def print_unique_values(df,max_val = 10):
    for col in df:
        if (len(df[col].unique()) < 10):
            print(df[col].name, ":" ,df[col].unique())

print_unique_values(df)
```

# 8. Data Manipulation

In [ ]: ▶

```python
#Defined A Common Function to remove spaces, special characters, brackets from column Names
# and Fields Values of a DataSet

def filter_df(df):

    import string
#     print(string.punctuation)

    def remove_punctuation(s):
        s = ''.join([i for i in s if i not in frozenset(string.punctuation) and i not in ' '])
        return s

    # To replace the spaces with '_', and removing brackets from Column NAmes
    df.columns = df.columns.str.strip().str.replace(' ', '_').str.replace('(', '').str.replace(')', '')

    # removing punctuation from the object type column values
    df_categorical = df.select_dtypes(include=object)
#     print(len(df_categorical))
    for col in df_categorical.columns:
        df[col] = df[col].apply(remove_punctuation)
    return df

# calling the function to Filter the Column Names and Column Values
df = filter_df(df)


print('-'*80,"\n The DataSet after removing the punctuations \n",'-'*80)
print(df.head())

```

In [ ]: ▶

```python
# Creating a Copy of the DataSet
df_org = df
print(df_org.head())
```

## Data Manipulation

**1. We Observe from the Data , that Total Charges Atrribute should be Float instead of object variable.**

**2. And Senior Citizen Atrribute should be Categorical , instead of Int**

```
In [ ]:    1  # We need to convert the Total Charges from object type to Numeric
           2  #df = df_org
           3
           4  # replacing the Space if found in Field with null value
           5  df['TotalCharges'] = df['TotalCharges'].replace(r'\s+', np.nan, regex=True)
           6  # converting to numeric
           7  #print(df['TotalCharges'].head())
           8  df['TotalCharges'] = pd.to_numeric(df['TotalCharges'])
           9
          10
          11  print('-'*80,"\n Total Charges converted to numeric Field \n",'-'*80)
          12  print(df.info())
```

```
In [ ]:    1  # making The SeniorCitizen Column , Categorical
           2
           3  df['SeniorCitizen'] = df['SeniorCitizen'].replace({1:'Yes',0:'No'})
           4
           5  print(" After Making Senior Citizen Field Categorical , checking for Nulls \n",'-'*80)
           6  print(df.isnull().sum())
           7  print(df.head())
```

```
In [ ]:    1  # We Observe the Conversions There are 11 values null in TotalCharges,
           2  #print(df[df['TotalCharges'].isnull()== True].head())
           3  print('-'*80,'''\n We Observe after converting TotalCharges and Senior Citizen columns
           4          1. Totalcharges Column has some null values
           5          2. Tenure for Such customers is 0 , hence we can fill null values by mean value so will not impact much \
           6  df['TotalCharges'] = df['TotalCharges'].fillna(df['TotalCharges'].mean())
           7  print(df.isnull().sum())
```

**We Observe that on converting the Field to Numeric , there are some null values for Total Charges created for customers who has not opted out of Service The Tenure for these customers is 0 ,but there may be additional charges so we can put the total charges as mean value**

**Imputation**

```python
# the Final Data Set
print(df.head())
```

# 9. Exploratory Data Analysis

```python
# IDentifying the Categorical and Numerical FEatures of the DataSet and creating two Datasets

df_categorical = df.select_dtypes(include=object)

column_categorical = df_categorical.columns
print('-'*80,"\n The Categorical Fields are \n",'-'*80)
print(column_categorical)
```

```python
# The Numerical Fields in the DataSet
df_numerical = df.select_dtypes(include=[np.float,np.int64])

print('-'*80,"\n The Numerical Fields are \n",'-'*80)
column_numerical = df_numerical.columns
print(column_numerical)
```

**Univariate Analysis**

In [ ]:

```python
# Function to display the distribution of each variable for
# Categorical Fields as well as Numerical fields , based on the object_mode value

def display_plot(df, col_to_exclude, object_mode = True):
    """
    This function plots the count or distribution of each column in the dataframe based on specified inputs
    @Args
      df: pandas dataframe
      col_to_exclude: specific column to exclude from the plot, used for excluded key
      object_mode: whether to plot on object data types or not (default: True)

    Return
      No object returned but visualized plot will return based on specified inputs
    """
    n = 0
    this = []

    if object_mode:
        nrows = 4
        ncols = 4
        width = 20
        height = 20

    else:
        nrows = 2
        ncols = 2
        width = 14
        height = 10


    for column in df.columns:
        if object_mode:
            if (df[column].dtypes == 'O') & (column != col_to_exclude):
                this.append(column)


        else:
            if (df[column].dtypes != 'O'):
                this.append(column)

    totcnt = len(this)
```

```
42          fig, ax = plt.subplots(nrows, ncols, sharex=False, sharey=False, figsize=(width, height))
43          for row in range(nrows):
44              for col in range(ncols):
45                  if (totcnt == 0):
46                      break
47                  if object_mode:
48                      g = sns.countplot(df[this[n]], ax=ax[row][col])
49                  else:
50                      g = sns.distplot(df[this[n]], ax = ax[row][col])
51                  totcnt -= 1
52
53
54
55
56                  ax[row,col].set_title("Column name: {}".format(this[n]))
57                  ax[row, col].set_xlabel("")
58                  ax[row, col].set_ylabel("")
59                  n += 1
60          plt.show();
61          return None
62
```

In [ ]:

```
1  # Displaying the Plots for each variable for Categorical Fields
2  # Step3: Data Visualisation of the Customer Data
3
4  display_plot(df, 'customerid', object_mode = True)
5
6
7  print('-'*80, "\n The Main observations from the Categorical Data as follows \n",'-'*80)
8  print('''        1. Almost Equal Percentage of Male and Female Customers
9       2. Most of Customers have Phone Service
10      3. FiberOptics  is prefered InternetService way
11      4. Most of the Customers have taken Monthly Plan
12      5. Payment preferred by Customers is paperless and mostly pay through
13         Electronic Cheque
14      6. Very few Customers have taken online services, like Backup , Streaming ,
15         opportunity for company to sell that to customers
16      ''')
```

```
In [ ]:  ▶|    1  # Observing the Numerical Fields
              2  df.head()
              3
              4  display_plot(df, 'customerid', object_mode = False)
              5
              6  print('-'*80, "\n The Main observations from the Numerical Data as follows \n",'-'*80)
              7  print('''     1. Maximum Customers have taken month wise Tenure and followed by very long tenures,
              8          Hence Maximum customers pay monthly charges
              9
             10     ''')
```

## Feature Engineering

Based on the value of the services the subscribers subscribed to, there are **yes**, **no**, and **no phone / internet service**. These are somewhat related to primary products. Examples are illustrated through *panda crosstab* function below:

1. **Phone service (Primary) and Multiple lines (Secondary)**

   - If the subscribers have phone service, they may have multiple lines (yes or no).
   - But if the subscribers don't have phone service, the subscribers will never have multiple lines.

```
In [ ]:  ▶|    1  pd.crosstab(index = df["PhoneService"], columns = df["MultipleLines"])
```

2. **Internet Service (Primary) and other services, let's say streaming TV (secondary)**

   - If the subscribers have Internet services (either DSL or Fiber optic), the subscribers may opt to have other services related to Internet (i.e. onlineSecurity, OnlineBackup , DeviceProtection,TechSupport,StreamingTV, StreamingMovies,).
   - But if the subscribers don't have the Internet services, this secondary service will not be available for the subscribers.

```
In [ ]:  ▶|    1  # Show For StreamingTV Feature, which has yes, no , no Internet option
              2  pd.crosstab(index = df["InternetService"], columns = df["StreamingTV"])
```

So We can transform the Secondary Attributes(MultipleLines, StreamingTV...) to have only two values , Yes and No and Transform **No Phone / Internet service** to be the same **No**

In [ ]:

```python
1  # Function Defined to identify all the columns/Services that have more than 2 options
2  # as mentioned above and replace "No Phone Service" &"No Internet Service" to "No"
3
4  def convert_no_service (df):
5      col_to_transform = []
6      for col in df.columns:
7          if (df[col].dtype == 'O') & (col != 'customerid'):
8              if len(df[df[col].str.contains("No")][col].unique()) > 1:
9                  col_to_transform.append(col)
10
11     print("Total column(s) that will be transformed: {}".format(col_to_transform))
12     for col in col_to_transform:
13         df.loc[df[col].str.contains("No"), col] = 'No'
14
15     return df
```

In [ ]:

```python
1  df = convert_no_service(df)
```

In [ ]:

```python
1  # Let's see the data after transformation.
2  print("The Data Can be Visualised again after Transforming the Secondary Columns")
3  display_plot(df, 'customerid', object_mode = True)
```

**The Conclusions that can be made after transforming Columns are**

**1. Most customers have phone Service and Single line (opportunity 1 for multiline)**

**2. Customers having Internet Service , have very few Secondary services (StreamingTV, OnlineBackUp...) (Opportunity2 - for cross-sell)**

In [ ]: ▶

```python
#1 We have observed above each Columns Behavior for the customers, now we can see the impact
#2 of these columns on Target Column Churn
3
4
#5 Function to compare each Column with Churn
def compare_Features(cols_cmp,Target = 'Churn'):
    for col in cols_cmp:
        print(col)
        pd.crosstab(index = df[col],columns=df[Target],margins=True).plot(kind='bar',figsize=(5,3))
        for lbl in df[col].unique():
            Fea_cnt = len(df[(df[col] == lbl)&(df['Churn']== 'Yes')])
            print(" Percent of ({} {}) Left Company {}".format(lbl ,col, (Fea_cnt/Churn_cnt)*100))


#15 compare_Features(df,df.Churn)
cols = column_categorical
Churn_cnt = len(df[(df['Churn'] == 'Yes')])
cols_cmp = cols[0:5]

#20 Comparing the First 5 Columns /FEatures with Churn
compare_Features(cols_cmp)
print('-'*80)
print(''' The Conclusion that can be made from Comparing with Churn are:
    1. Churning of Customers was independent of Gender (i.e no impact)
    2. More percentage Young Customers, Non Partnered , having no Dependents and have
       PhoneService left Company
    3. Company could make customers partners, and see if there is any issue with Phone Service ''')
```

```python
1  # Comparing the Next 5 Columns /FEatures with Churn
2  cols_cmp = cols[5:10]
3  compare_Features(cols_cmp)
4  print('-'*80)
5  print(''' The Conclusion that can be made from Comparing with Churn are:
6          1. Churning of Customers was independent of whether they had Multiple Lines
7          2. More percentage Customers taking FiberOptic Internet Service, not having Online
8              Services ...left Company
9          3. Company could check if some issue with FiberOptic Service, and try selling online
10             services ''')
```

```python
1  # Comparing the Next 5 Columns /FEatures with Churn
2  cols_cmp = cols[10:16]
3
4  compare_Features(cols_cmp)
5  print('-'*80)
6  print(''' The Conclusion that can be made from Comparing with Churn are:
7          1. Percentage of Customers having no TechSupport left Company
8          2. Having Streaming Tv, Movies or not had same impact on leaving Company
9          3. Maximum percentage of Customers with MonthtoMonth Contract Left
10         3. Company could try making customers with longer contracts to retain them  ''')
```

```python
1  # To check the relation between Tenure and Partner Atrributes
2  plt.figure(figsize=(8,4))
3  sns.countplot(x=df['tenure'],hue=df.Partner);
```

**We Observer , Most of the Customers that are Partners Stay Longer with The Company. So Being a Partner is a Plus-Point For the Company as they will Stay Longer with Them.**

**Let's Check for Outliers in Monthly Charges And Total Charges Using Box Plots**

```python
1  df.boxplot('MonthlyCharges');
2  print(" No Outliers in Monthly Charges")
```

In [ ]:
```python
# For Total Charges Attribute , we observe there are outliers

# df.boxplot('TotalCharges')
# print(df['TotalCharges'].describe())
Tot_mean = df['TotalCharges'].mean()
Tot_std  = df['TotalCharges'].std()
std_low =  Tot_mean - 3 * Tot_std
std_upp =  Tot_mean + 3 * Tot_std
print("The number of TotalCharges below (mean -2 std)" + str(std_low) + " are " , len(df[df['TotalCharges'] < st
print("The number of TotalCharges above (mean +2 std)  " + str(std_upp) + " are ", len(df[df['TotalCharges'] > s
print("We Observe 99percent confidence of Total Charges values lie between %0.2f and %0.2f " % (std_low, std_upp
# and %0.3f,std_upp))
```

In [ ]:
```python
#df.drop(df[df['TotalCharges'] > std_upp],axis=0, inplace=True)
df = df_org
Out = df[df['TotalCharges'] > std_upp]
print(df.shape)
Tot_new =  df[(df['TotalCharges'] > std_upp)].index
print(Tot_new)
#print(df.drop(Out, inplace=True))
df = df.drop(Tot_new,axis=0)
print(df.shape)
```

**Monthly Charges don't have any Outliers so we don't have to Get into Extracting Information from Outliers.**

In [ ]:
```python
## correlation matrix

# Let's Check the Correaltion Matrix in Seaborn
sns.heatmap(df.corr(),xticklabels=df.corr().columns.values,yticklabels=df.corr().columns.values,annot=True);
```

**Here We observe Tenure and Total Charges Atrributes are correlated 50% and also Monthly charges and Total Charges are also correlated 31% with each other. But It is not above 0.90 , so we can retain columns**

**we can assume from our domain expertise that , Total Charges ~ Monthly Charges * Tenure + Additional Charges(Tax).**

## Bucketing

In [ ]:
```python
# Bucketting the Tenure Feature into slabs
df['tenure'].describe()
```

In [ ]:
```python
# Changing the Tenure to categorical column

print(" Min and Max Values of tenure  = %0.2f & %0.2f  "%(min(df['tenure']), max(df['tenure'])))

def tenure_lab(telcom) :

    if telcom["tenure"] <= 12 :
        return "Tenure_0-12"
    elif (telcom["tenure"] > 12) & (telcom["tenure"] <= 24 ):
        return "Tenure_12-24"
    elif (telcom["tenure"] > 24) & (telcom["tenure"] <= 48) :
        return "Tenure_24-48"
    elif (telcom["tenure"] > 48) & (telcom["tenure"] <= 60) :
        return "Tenure_48-60"
    elif telcom["tenure"] > 60 :
        return "Tenure_gt_60"


df["tenure_group"] = df.apply(lambda x:tenure_lab(x),axis = 1)
```

In [ ]:
```python
print(" Telecom Data with Tenure in Buckets\n",'-'*80)
print(df.head())
```

## 10. Data preprocessing

## Encoding categorical variable

```
In [ ]:    1  # Categorical Features and count of unique values
           2  #print(df.nunique())
           3  bin_cols = df.nunique()[df.nunique() < 3].keys().tolist()
           4
           5  multi_cols = df.nunique()[(df.nunique() < 7) & (df.nunique() > 2)].keys().tolist()
           6  num_cols = df.nunique()[(df.nunique() > 7) ].keys().tolist()
           7  print("The Binary Categorical Columns are : \n",'-'*80)
           8  print("Binary Value Categorical Columns \n",'-'*80)
           9  print(bin_cols)
          10  print("Multiple Value Categorical Columns \n",'-'*80)
          11  print(multi_cols)
          12  print("Numerical Columns \n",'-'*80)
          13  print(num_cols)
```

```
In [ ]:    1  #Label encoding for all the Categorical Values
           2  from sklearn.preprocessing import LabelEncoder
           3  from sklearn.preprocessing import MinMaxScaler
           4
           5  #Encoding the binary columns
           6
           7
           8  le = LabelEncoder()
           9  for i in bin_cols :
          10      df[i] = le.fit_transform(df[i])
          11
          12  #Duplicating columns for multi value columns
          13  df = pd.get_dummies(data = df,columns = multi_cols )
```

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

```
In [ ]:    1  list(df.columns)
```

## Normalizing features

```
In [ ]:  ▶  1  telcom = df
             2
             3  #Scaling Numerical columns as there is large difference in values , to make them comparable
             4  '''
             5  ransforms features by scaling each feature to a given range.
             6
             7  This estimator scales and translates each feature individually such that it is in the given range on the trainin
             8  '''
             9
            10  std = MinMaxScaler()
            11
            12
            13  scaled = std.fit_transform(telcom[num_cols])
            14  scaled = pd.DataFrame(scaled,columns=num_cols)
```

```
In [ ]:  ▶  1  print(scaled.shape)
             2  scaled.head(2)
```

```
In [ ]:  ▶  1  #dropping original values and  merging new scaled values for numerical columns
             2  df_telcom_og = telcom.copy()
             3  telcom = telcom.drop(columns = num_cols,axis = 1)
             4
             5  telcom.reset_index(drop=False, inplace=True)
             6
             7  telcom = pd.concat([telcom, scaled], axis=1)
             8
             9  telcom.set_index('customerID', inplace=True)
            10
            11  print("The DataSet with scaled numerical values \n",'-'*80)
            12
            13  print(telcom.shape)
            14  print(telcom.head())
```

## Spliting The Data Set into Train/Val/Test data

In [ ]: ▶|    1  telcom.info()

In [ ]: ▶|
```
1   # importing the metric for performance calculation
2
3   from sklearn.model_selection import train_test_split
4   from sklearn.linear_model import LogisticRegression
5   from sklearn.metrics import confusion_matrix,accuracy_score,classification_report
6   from sklearn.metrics import roc_auc_score,roc_curve,scorer
7   from sklearn.metrics import f1_score
8   import statsmodels.api as sm
9   from sklearn.metrics import precision_score,recall_score
10  #splitting train and test data
11
12  # telcom = df
13  target_col = telcom["Churn"]
14  telcom = telcom.drop('Churn',axis=1)
15  #target_col = target_col.values.reshape(-1,1)
16
17  X_train,X_test,y_train,y_test = train_test_split(telcom,target_col,test_size = 0.25 ,random_state = 111)
18
19  print(" The Training and Test Data Set size after Splitting \n",'-'*80)
20  print(X_train.shape)
21  print(X_test.shape)
22  print(y_train.shape)
23  print(y_test.shape)
```

# 11. Model Building

In [ ]:
```python
# importing the Machine Learning Models
from sklearn.dummy import DummyClassifier

# Machine learning
from sklearn import tree , linear_model
from sklearn.svm import LinearSVC
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.linear_model import LinearRegression, LogisticRegression, Ridge, Lasso, SGDClassifier
from sklearn.tree import DecisionTreeClassifier
from xgboost.sklearn import XGBClassifier
```

In [ ]:
```python
# validation
from sklearn import datasets, model_selection, metrics , preprocessing
```

In [ ]:
```python
# Grid and Random Search
import scipy.stats as st
from scipy.stats import randint as sp_randint
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import RandomizedSearchCV
```

In [ ]:
```python
# Metrics
from sklearn.metrics import precision_recall_fscore_support, roc_curve, auc
```

In [ ]:
```python
#utilities
import time
import io, os, sys, types, time, datetime, math, random
```

In [ ]:

```python
# calculate the fpr and tpr for all thresholds of the classification
def plot_roc_curve(y_test, preds):
    fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
    roc_auc = metrics.auc(fpr, tpr)
    plt.title('Receiver Operating Characteristic')
    plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
    plt.legend(loc = 'lower right')
    plt.plot([0, 1], [0, 1],'r--')
    plt.xlim([-0.01, 1.01])
    plt.ylim([-0.01, 1.01])
    plt.ylabel('True Positive Rate')
    plt.xlabel('False Positive Rate')
    plt.show()



# Function that runs the requested algorithm and returns the accuracy metrics
def fit_ml_algo(algo, X_train, y_train, X_test, cv):
    # One Pass
    model = algo.fit(X_train, y_train)
    test_pred = model.predict(X_test)
    if (isinstance(algo, (LogisticRegression,
                          KNeighborsClassifier,
                          GaussianNB,
                          DecisionTreeClassifier,
                          RandomForestClassifier,
                          GradientBoostingClassifier))):
        probs = model.predict_proba(X_test)[:,1]
    else:
        probs = "Not Available"
    acc = round(model.score(X_test, y_test) * 100, 2)
    # CV
    train_pred = model_selection.cross_val_predict(algo,
                                                   X_train,
                                                   y_train,
                                                   cv=cv,
                                                   n_jobs = -1)
    acc_cv = round(metrics.accuracy_score(y_train, train_pred) * 100, 2)
    return train_pred, test_pred, acc, acc_cv, probs

# Utility function to report best scores
```

```
42  def report(results, n_top=5):
43      for i in range(1, n_top + 1):
44          candidates = np.flatnonzero(results['rank_test_score'] == i)
45          for candidate in candidates:
46              print("Model with rank: {0}".format(i))
47              print("Mean validation score: {0:.3f} (std: {1:.3f})".format(
48                      results['mean_test_score'][candidate],
49                      results['std_test_score'][candidate]))
50              print("Parameters: {0}".format(results['params'][candidate]))
51              print("")
52
```

## Baseline model with DummyClassifier

In [ ]:

```python
1  # Instantiating a Dummy Classifier Model ,that will be the baseline for
2  # comparisions for all other models
3
4  clf = DummyClassifier(strategy='most_frequent',random_state=0)
5  clf.fit(X_train, y_train)
6  accuracy = clf.score(X_test, y_test)
7
```

```
In [ ]: ▶|    1  preds = clf.predict(X_test)
              2
              3
              4  # dummyistic Regression
              5  start_time = time.time()
              6  train_pred_dummy, test_pred_dummy, acc_dummy, acc_cv_dummy, probs_dummy = fit_ml_algo(DummyClassifier(strategy='
              7                                                                          X_train,
              8                                                                          y_train,
              9                                                                          X_test,
             10                                                                          10)
             11  dummy_time = (time.time() - start_time)
             12  print("Accuracy for Dummy Classfier Model: %s" % acc_dummy)
             13  print("Accuracy CV 10-Fold: %s" % acc_cv_dummy)
             14  print("Running Time: %s" % datetime.timedelta(seconds=dummy_time))
             15
             16  print(" Classification report for Training Data\n",'-'*80 )
             17  print (metrics.classification_report(y_train, train_pred_dummy))
             18
             19  print(" Classification report for Test Data\n",'-'*80 )
             20  print (metrics.classification_report(y_test, test_pred_dummy))
             21
```

## Select Candidate Algorithms

**1. KNN**

**2. Logistic Regression**

**3. Naive Bayes**

**4. Decision Tree**

**5. Random Forest**

## 6. Gradient Boosted Trees

In [ ]: ▶|

```python
# Model 1 - k-Nearest Neighbors

start_time = time.time()
train_pred_knn, test_pred_knn, acc_knn, acc_cv_knn, probs_knn = fit_ml_algo(KNeighborsClassifier(n_neighbors = 3,
                                                                                                 n_jobs = -1),
                                                                            X_train,
                                                                            y_train,
                                                                            X_test,
                                                                            10)

knn_time = (time.time() - start_time)
print("K Nearest Neighbour Model \n",'-'*80)
print("Accuracy for KNN Model: %s" % acc_knn)
print("Accuracy CV 10-Fold: %s" % acc_cv_knn)
print("Running Time: %s" % datetime.timedelta(seconds=knn_time))

print (metrics.classification_report(y_train, train_pred_knn))

print (metrics.classification_report(y_test, test_pred_knn))
print("ROC_AUC Curve for KNN \n",'-'*80)
plot_roc_curve(y_test, probs_knn)
```

In [ ]:  ▶|

```python
1   # Specify parameters and distributions to sample from
2   param_dist = {'penalty': ['l2', 'l1'],
3                                'class_weight': [None, 'balanced'],
4                                'C': np.logspace(-20, 20, 10000),
5                                'intercept_scaling': np.logspace(-20, 20, 10000)}
6
7
8
9   # Run Randomized Search
10  n_iter_search = 10
11  lrc = LogisticRegression()
12  random_search = RandomizedSearchCV(lrc,
13                                      n_jobs=-1,
14                                      param_distributions=param_dist,
15                                      n_iter=n_iter_search)
16
17  start = time.time()
18  random_search.fit(X_train, y_train)
19  print("RandomizedSearchCV took %.2f seconds for %d candidates"
20        " parameter settings." % ((time.time() - start), n_iter_search))
21  print(report(random_search.cv_results_))
22  print('-'*80)
23  print(random_search.best_estimator_)
24  print('-'*80)
25  print(random_search.best_params_)
26
```

In [ ]: ▶|

```python
 1  # Modeling the Data with various Models , one by one , since we want to see individual results
 2  # for each model in detail , hence not defined a function for same
 3
 4  # Model 1 - Logistic Regression
 5  start_time = time.time()
 6  train_pred_log, test_pred_log, acc_log, acc_cv_log, probs_log = fit_ml_algo(LogisticRegression(n_jobs = 1,penalt
 7                                                                 X_train,
 8                                                                 y_train,
 9                                                                 X_test,
10                                                                 10)
11  log_time = (time.time() - start_time)
12  print(" Logistic REgression Model \n",'-'*80)
13  print("Accuracy for Logistic Regression : %s" % acc_log)
14  print("Accuracy for CV 10-Fold: %s" % acc_cv_log)
15  print("Running Time: %s" % datetime.timedelta(seconds=log_time))
16
17  print (metrics.classification_report(y_train, train_pred_log))
18
19  print (metrics.classification_report(y_test, test_pred_log))
20  print(" ROC-AUC Curve for Logistic Regression Model\n",'-'*80)
21  plot_roc_curve(y_test, probs_log)
```

In [ ]:

```python
# Model 3 - Gaussian Naive Bayes
start_time = time.time()
train_pred_gaussian, test_pred_gaussian, acc_gaussian, acc_cv_gaussian, probs_gau = fit_ml_algo(GaussianNB(),
                                                                                                X_train,
                                                                                                y_train,
                                                                                                X_test,
                                                                                                10)
gaussian_time = (time.time() - start_time)
print(" Naives Bayes Gausian Model \n",'-'*80)
print("Accuracy: %s" % acc_gaussian)
print("Accuracy CV 10-Fold: %s" % acc_cv_gaussian)
print("Running Time: %s" % datetime.timedelta(seconds=gaussian_time))

print (metrics.classification_report(y_train, train_pred_gaussian))

print (metrics.classification_report(y_test, test_pred_gaussian))
print("ROC_AUC Curve for NB Model\n",'-'*80)
plot_roc_curve(y_test, probs_gau)
```

In [ ]:

```python
# Model 4 - Decision Tree Classifier
start_time = time.time()
train_pred_dt, test_pred_dt, acc_dt, acc_cv_dt, probs_dt = fit_ml_algo(DecisionTreeClassifier(),
                                                                       X_train,
                                                                       y_train,
                                                                       X_test,
                                                                       10)
dt_time = (time.time() - start_time)
print(" Decision Tree Model \n",'-'*80)
print("Accuracy: %s" % acc_dt)
print("Accuracy CV 10-Fold: %s" % acc_cv_dt)
print("Running Time: %s" % datetime.timedelta(seconds=dt_time))

print (metrics.classification_report(y_train, train_pred_dt))

print (metrics.classification_report(y_test, test_pred_dt))
print("ROC_AUC curve for Decision Tree Model \n",'-'*80)
plot_roc_curve(y_test, probs_dt)
```

In [ ]:

```python
# Model 5 - Random Forest Classifier - Random Search for Hyperparameters

# Utility function to report best scores
# def report(results, n_top=5):
#     for i in range(1, n_top + 1):
#         candidates = np.flatnonzero(results['rank_test_score'] == i)
#         for candidate in candidates:
#             print("Model with rank: {0}".format(i))
#             print("Mean validation score: {0:.3f} (std: {1:.3f})".format(
#                 results['mean_test_score'][candidate],
#                 results['std_test_score'][candidate]))
#             print("Parameters: {0}".format(results['params'][candidate]))
#             print("")



# Specify parameters and distributions to sample from
param_dist = {"max_depth": [10, None],
              "max_features": sp_randint(1, 11),
              "min_samples_split": sp_randint(2, 20),
              "min_samples_leaf": sp_randint(1, 11),
              "bootstrap": [True, False],
              "criterion": ["gini", "entropy"]}


# Run Randomized Search
n_iter_search = 10
rfc = RandomForestClassifier(n_estimators=10)
random_search = RandomizedSearchCV(rfc,
                                   n_jobs = -1,
                                   param_distributions=param_dist,
                                   n_iter=n_iter_search)

start = time.time()
random_search.fit(X_train, y_train)
print("RandomizedSearchCV took %.2f seconds for %d candidates"
      " parameter settings." % ((time.time() - start), n_iter_search))
print(report(random_search.cv_results_))
print('-'*80)
print(random_search.best_estimator_)
print('-'*80)
```

```
42  print(random_search.best_params_)
```

In [ ]:

```python
1   #Final Model 5  Random Forest Classifier
2   start_time = time.time()
3   rfc = RandomForestClassifier(n_estimators=10,
4                                bootstrap=True,
5                                max_depth=10,
6                                min_samples_leaf=9,
7                                min_samples_split=7,
8                                criterion='entropy',
9                                max_features=2)
10  train_pred_rf, test_pred_rf, acc_rf, acc_cv_rf, probs_rf = fit_ml_algo(rfc,X_train,y_train,X_test,10)
11  rf_time = (time.time() - start_time)
12  print(" Random Forest Model  \n",'-'*80)
13  print("Accuracy for RF: %s" % acc_rf)
14  print("Accuracy CV 10-Fold: %s" % acc_cv_rf)
15  print("Running Time: %s" % datetime.timedelta(seconds=rf_time))
16
17  print (metrics.classification_report(y_train, train_pred_rf))
18
19  print (metrics.classification_report(y_test, test_pred_rf))
20  print(" ROC-AUC Curve for RF\n",'-'*80)
21  plot_roc_curve(y_test, probs_rf)
```

In [ ]: ▶|

```python
# Model 6 -  Gradient Boosting Trees
start_time = time.time()
train_pred_gbt, test_pred_gbt, acc_gbt, acc_cv_gbt, probs_gbt = fit_ml_algo(GradientBoostingClassifier(),
                                                                            X_train,
                                                                            y_train,
                                                                            X_test,
                                                                            10)
gbt_time = (time.time() - start_time)
print(" Gradient Boost Tree\n",'-'*80)
print("Accuracy: %s" % acc_gbt)
print("Accuracy CV 10-Fold: %s" % acc_cv_gbt)
print("Running Time: %s" % datetime.timedelta(seconds=gbt_time))

print (metrics.classification_report(y_train, train_pred_gbt))

print (metrics.classification_report(y_test, test_pred_gbt))
print(" ROC-AUC Curve for GBT \n",'-'*80)
plot_roc_curve(y_test, probs_gbt)
```

In [ ]: ▶|

```python
def xgb_f1(y, t):
    #
    # Function to evaluate the prediction based on F1 score, this will be used as evaluation metric when trainin
    # Args:
    #   y: label
    #   t: predicted
    #
    # Return:
    #   f1: F1 score of the actual and predicted
    #
    t = t.get_label()
    y_bin = [1. if y_cont > 0.5 else 0. for y_cont in y]   # change the prob to class output
    return 'f1', f1_score(t, y_bin)

best_xgb = XGBClassifier(objective = 'binary:logistic',
                         colsample_bylevel = 0.7,
                         colsample_bytree = 0.8,
                         gamma = 1,
                         learning_rate = 0.15,
                         max_delta_step = 3,
                         max_depth = 4,
                         min_child_weight = 1,
                         n_estimators = 50,
                         reg_lambda = 10,
                         scale_pos_weight = 1.5,
                         subsample = 0.9,
                         silent = True,
                         n_jobs = 4
                        )

xgbst = best_xgb.fit(X_train, y_train, eval_metric = xgb_f1, eval_set = [(X_train, y_train), (X_test, y_test)],
        early_stopping_rounds = 20)
```

```
In [ ]:  ▶|   1  train_pred_xgbst, test_pred_xgbst, acc_xgbst, acc_cv_xgbst, probs_xgbst = fit_ml_algo(xgbst,
         2                                                       X_train,
         3                                                       y_train,
         4                                                       X_test,
         5                                                       10)
```

```
In [ ]:  ▶|   1  import xgboost as xgb
         2  xgb.plot_importance(best_xgb, max_num_features = 15)
         3  plt.show();
```

## Compare all models

```
In [ ]:  ▶|   1  models = pd.DataFrame({
         2      'Model': ['KNN', 'Logistic Regression',
         3                'Random Forest', 'Naive Bayes',
         4                'Decision Tree',
         5                'Gradient Boosting Trees'],
         6      'Score': [
         7          acc_knn,
         8          acc_log,
         9          acc_rf,
        10          acc_gaussian,
        11          acc_dt,
        12          acc_gbt,
        13
        14      ]})
        15
        16  print(" WE observe from the Table , GBT, Logistic Regression,  Random Forest have good accuracy\n",'-'*80)
        17  models.sort_values(by='Score', ascending=False)
```

In [ ]:

```python
# printing a ROC-AUC Curve for the All models in same plot to compare
models = [
    'KNN',
    'Logistic Regression',
    'Random Forest',
    'Naive Bayes',
    'Decision Tree',
    'Gradient Boosting Trees',

]
probs = [
    probs_knn,
    probs_log,
    probs_rf,
    probs_gau,
    probs_dt,
    probs_gbt
]
colors = [
    'blue',
    'green',
    'red',
    'cyan',
    'magenta',
    'yellow',
    'black',
]
```

In [ ]:

```python
def plot_roc_curves(y_test, prob, model):
    fpr, tpr, threshold = metrics.roc_curve(y_test, prob)
    roc_auc = metrics.auc(fpr, tpr)
    plt.plot(fpr, tpr, 'b', label = model + ' AUC = %0.2f' % roc_auc, color=colors[i])
    plt.legend(loc = 'lower right')

for i, model in list(enumerate(models)):
    plot_roc_curves(y_test, probs[i], models[i])

plt.show()
```

## Interpretation

# 1. The Conclusions that can be made from the graph and Probability scores from the test dataset

**1. The DummyClassifier , used For Baseline , gave accuracy of 72%**

**2. We Observe that based on ROC_AUC, the Algoritms that the Model is trained the best accuracy is given by Gradient Boost , Logistic Regression, Random Forest and Naives Bayes Algoritms**

**3. The Probability score metric also reconfirms the conclusion.**

**Hence , The Best Model Trained for This test DataSet is Gradient Boost with score 80%**

# 2. The Observations Made From DataSet after Analysis

**1. Customer Churning is irrepective of Gender**

**2. Maximum Customers have Phone Service with single Line , and churn out more in them**

**3. Maximum Customers having internet, prefer Fiber Optic, But Churn High too in this category**

**4. Maximum Customers availing month-monthBilling and high churn as compared to long tenure customers**

# 3.Customer Retension Programs Suggested

**1. Analyse issues faced by customer in FiberOptics Internet offering and resolve them at earliest, to prevent Churning in this segment**

**2. Convert Customers from Month to Month billing to long tenure to Churn Less**

**3. Analyse if Phone Service having any issues and resolve at the earliest, to prevent Churning**

**4. Less Customers having Multiline and Online Services - Company can focus and sell models**

**5. Company to offer Partnership to Customers ,they would stay longer duration.**