



University School of Automation and Robotics
GURU GOBIND SINGH INDRAPRASTHA UNIVERSITY
East Delhi Campus, Surajmal Vihar
Delhi - 110092



PROJECT REPORT

ARM-256

MACHINE LEARNING **LAB**

Title of Project: Telco Customer Churn Prediction
Using 5 different Classification Models

Student Name: Kriti Rastogi

Enrollment Number: 20119051623

Email ID:

kriti.919052023@ipu.ac.in

kritirastogi.usar@gmail.com

Contact Number: 9643339254

Google Drive Link:

<https://drive.google.com/drive/folders/1kXS9w9FFhY30ymWoAVNACcT-L8HBC5QI?usp=sharing>

TITLE: Telco Customer Churn Prediction Using 5 different Classification Models

ABSTRACT

This report seeks to estimate customer churn in the telecommunications sector by applying machine learning classification algorithms. The goal is to forecast whether or not a customer will churn based on a set of demographics, behavioural, and service usage features. The model doesn't consider any non-work-related factors, like customer feelings or life situations; strictly, the dataset's tabular data is what is used for model training. This project uses a dataset containing 7043 records and 21 attributes, which range from contract type and monthly bill to the internet service provided and the duration of the contract. The goal is to predict the target feature, Churn, using five popular classification models: Logistical Regression, Support Vector Machine (SVM), K-Nearest Neighbours (KNN), Decision Tree, and Random Forest. All these models are trained and evaluated using the accuracy metric and classification reports, also referred to as performance measures.

In order to improve the accuracy of the model, the model underwent hyperparameter tuning, which creates visual plots showing changes in accuracy with differing parameters. In this case, all the models were compared and their pros and cons discussed including the accuracy scores and the outcomes they produced. Moreover, the report provides an interactive segment where users can fill in all pertinent aspects, and based on the input given, the optimal model is chosen to churn prediction alongside displaying the predicted result with the accompanying model and accuracy. Even though the project showcases a good use-case of machine learning in churn prediction, currently, it seems suitable only for academic or experimental settings. This can be improved with more advanced or realistic data and could be implemented into actual strategies aimed at improving customer retention.

KEYWORDS

Customer Churn, Decision Tree, Random Forest, Support Vector Classifier, Hyperparameter Tuning

INTRODUCTION

Apart from offering quality service, companies nowadays have to retain customers particularly for competitive fields such as telecom. Customer churn, or the act of a customer choosing to leave the service, threatens a company's revenue and reputation. Companies can improve customer satisfaction and their reputation with the company by taking action beforehand and recognizing churn patterns early on. Just as humans recognize dissatisfaction through repetitive behaviour, patterns of customer interactions can be used to forecast whether a remote customer is likely to churn using historical customer data. While it may not be perfectly accurate all the time, it can be used for predicting which customers are high risk, allowing businesses to make

the best use of their resources during retention campaigns. Remote customers are grouped using structured customer details like billing information, service preferences, contract type, and even internet usage. With these details, machine learning models are built, accurately predicting relationships between the features and churn behaviour. These models provide users with major alerts in advance so adequate measures can be taken to tailor responses, alter service quality or provide better incentives to the at-risk customers. In this project, the implementation of five classification algorithms in building a system designed for churn prediction is studied. Regardless of the model's limitations with human emotion, its implementation encourages organizations to adopt a more data-based approach, thus improving customer retention efforts and enabling growth over time.

PROPOSED METHODOLOGY

The following is the methodology used in this project to predict the telecommunications customer churn:

1. The dataset has been taken from Kaggle (<https://www.kaggle.com/datasets/blastchar/telco-customer-churn/data>)
2. This dataset has details of 7043 customers having 21 different features which include their usage pattern, billing information, people who are dependent on them and other important information. The target variable is thus an indicator to know whether there is customer churn or not.
3. Initially, the missing values in the dataset are handled and categorical variables are converted to numerical using one hot encoding. The numerical features are also scaled for standardisation.
4. To understand the distribution of features and relationships between various attributes, Exploratory Data Analysis (EDA) was conducted along with many visualisations using libraries like Matplotlib and Seaborn.
5. 5 ML classification algorithms have been used- logistic regression, Support vector Machine, Decision Tree, Random Forest and K nearest neighbours. Each model was trained using the training data which was 70% of the entire dataset and tested using the remaining 30% of the test set.
6. To ensure high accuracy and visualise how the values of training and testing accuracies vary with different values of hyperparameters, Hyperparameter Tuning was applied using techniques like GridSearchCV and RandomisedSearchCV.
7. The model performance was evaluated using features like accuracy, f1-score, precision, recall and confusion matrix. Using these, any cases of underfitting or overfitting were handled.
8. A user input-based section was implemented at the end of the project, wherein based on the input values, the system checks for the best out of the 5 models, predicts and output along with mentioning the accuracy percentage of the same.

A.IMPORTING LIBRARIES

The first step taken to build the models was to import all the necessary modules into the jupyter notebook. Let us look at each of them one by one

- **Pandas:** Pandas is a Python library used for working with data sets. It has functions for analysing, cleaning, exploring, and manipulating data. The name "Pandas" has a reference to both "Panel Data", and "Python Data Analysis" and was created by Wes McKinney in 2008.
- **NumPy:** NumPy is a Python library used for working with arrays. It also has functions for working in domain of linear algebra, Fourier transform, and matrices. NumPy was created in 2005 by Travis Oliphant. It is an open-source project and you can use it freely. NumPy stands for Numerical Python.
- **Matplotlib:** Matplotlib is a low-level graph plotting library in python that serves as a visualization utility. Matplotlib was created by John D. Hunter. Matplotlib is open source and we can use it freely.
- **Seaborn:** Seaborn is a library that uses Matplotlib underneath to plot graphs. It will be used to visualize random distributions.
- **Scikit-learn (sklearn):** It is a widely used Python library for machine learning. It provides simple tools for tasks like classification, regression, clustering, and model evaluation. Built on top of NumPy and pandas, it includes features for data preprocessing, model training, and performance measurement, making it ideal for developing and testing machine learning models efficiently.

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RandomizedSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import classification_report, accuracy_score, confusion_matrix
```

B. LOADING DATASET AND EXPLORATORY DATA ANALYSIS(EDA)

Further, the dataset named as **“CustChurn.csv”** is loaded, having **7043 rows and 21 columns** and the necessary **Exploratory Data Analysis** is carried out to understand the characteristics of a dataset before applying formal modelling techniques.

```
df=pd.read_csv("CustChurn.csv")
df.head()
```

	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	...	DeviceProtection	TechSupport
0	7590-VHJEG	Female	0	Yes	No	1	No	No phone service	DSL	No	...	No	No
1	5575-GNVOE	Male	0	No	No	34	Yes	No	DSL	Yes	...	Yes	No
2	3668-QPYBK	Male	0	No	No	2	Yes	No	DSL	Yes	...	No	No
3	7795-CFOCW	Male	0	No	No	45	No	No phone service	DSL	Yes	...	Yes	Yes
4	9237-HQITU	Female	0	No	No	2	Yes	No	Fiber optic	No	...	No	No

5 rows x 21 columns

After loading the dataset, it is summarised and interpreted in various forms to understand the underlying trends, patterns and anomalies.

```
df.shape
[1] ✓ 0.0s
... (7043, 21)

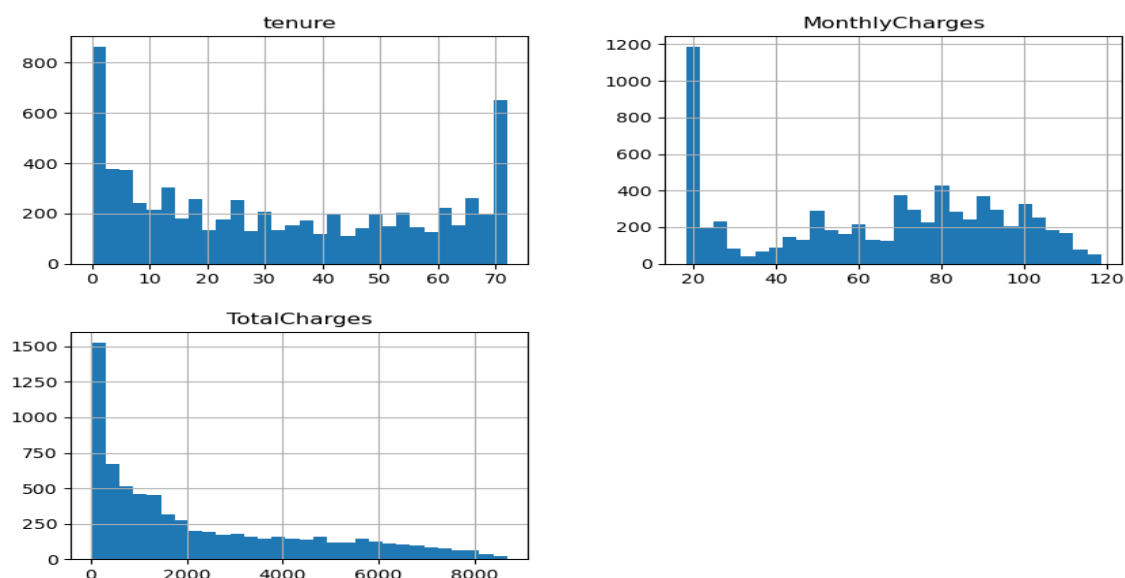
df.columns
[4] ✓ 0.0s
... Index(['customerID', 'gender', 'SeniorCitizen', 'Partner', 'Dependents',
         'tenure', 'PhoneService', 'MultipleLines', 'InternetService',
         'OnlineSecurity', 'OnlineBackup', 'DeviceProtection', 'TechSupport',
         'StreamingTV', 'StreamingMovies', 'Contract', 'PaperlessBilling',
         'PaymentMethod', 'MonthlyCharges', 'TotalCharges', 'Churn'],
        dtype='object')
```

The information as to if there are null values and what are the datatypes of the null values is also taken care of, which are further handled in the upcoming steps. Along with this, the basic mathematical calculations like calculating the mean, standard deviations, variances for various columns of the dataset are also carried out.

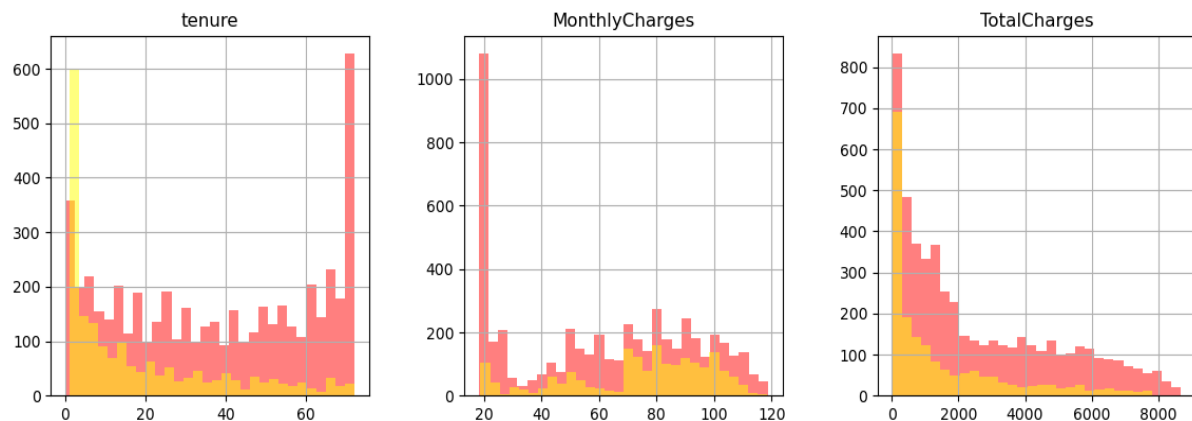
```
df.describe()
[6] ✓ 0.0s
...
   SeniorCitizen  tenure  MonthlyCharges
count  7043.000000  7043.000000  7043.000000
mean      0.162147    32.371149    64.761692
std       0.368612    24.559481    30.090047
min       0.000000     0.000000    18.250000
25%       0.000000     9.000000    35.500000
50%       0.000000    29.000000    70.350000
75%       0.000000    55.000000    89.850000
max       1.000000    72.000000   118.750000

df.info()
[5] ✓ 0.0s
...
<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
14  StreamingMovies       7043 non-null  object
15  Contract              7043 non-null  object
16  PaperlessBilling      7043 non-null  object
17  PaymentMethod         7043 non-null  object
18  MonthlyCharges        7043 non-null  float64
19  TotalCharges          7043 non-null  object
20  Churn                 7043 non-null  object
dtypes: float64(1), int64(2), object(18)
memory usage: 1.1+ MB
```

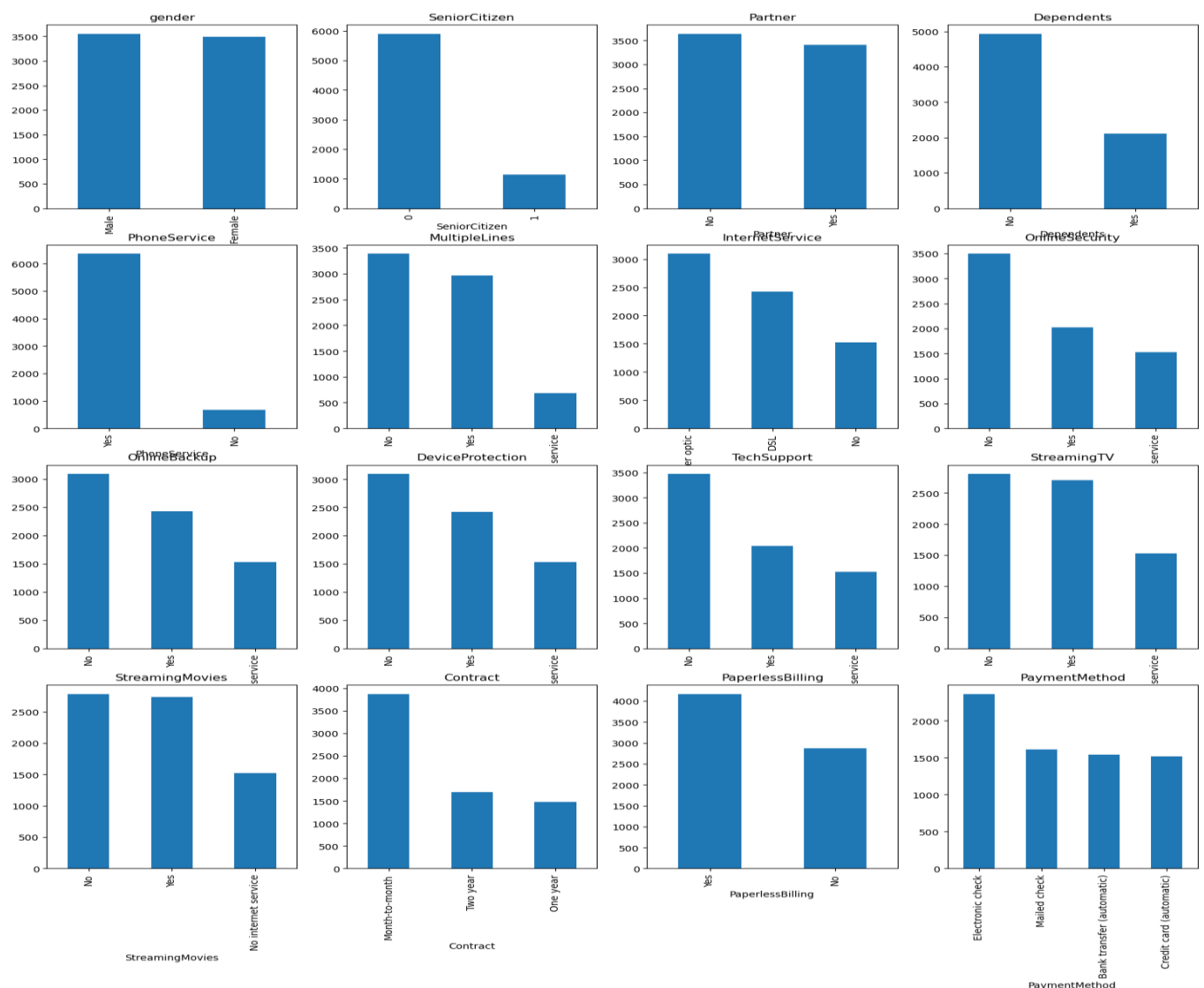
To further visualise the trends, various plots are drawn to show the number of customers involved for a particular tenure, who pay some particular monthly charges and a fixed amount as their annual charges to the telecommunications industry.



Based on the above trends, if the customer Churn taken place or not is also plotted using **double shaded histograms** to clearly understand the variation between the patterns.



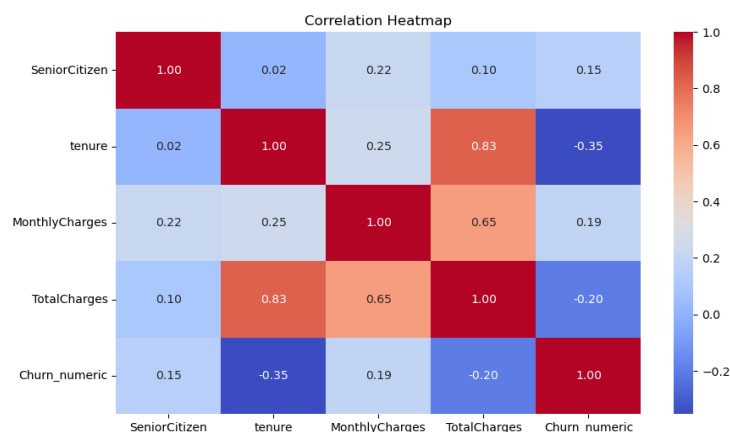
Furthermore, to yet clearly get to know about the exact count, for each column, the **bar graphs** are plot to show the category split among the columns themselves. This would further be helpful in pre processing of the data and also using the required models for training and testing of the dataset.



Pair plots are also drawn in 2 shades depicting if the churn took place or not based on 3 major columns- the monthly charges, the total charges and specifically the tenure in months. It gives a clear image of the ratio of churn Yes to the Churn No.



The **correlation heatmap** visualizes the relationships between different variables, with colours ranging from blue (negative correlation) to red (positive correlation). This helps identify patterns in customer behaviour, such as tenure affecting churn.



C. DATA PREPROCESSING

In this step, raw data is prepared for analysis by cleaning, transforming, and organizing it. It involves handling missing values, removing inconsistencies, encoding categorical variables, scaling features for uniformity, and filtering out noise. The goal is to ensure the data is accurate, structured, and suitable for meaningful insights or predictive modelling.

The null values are also looked for and are hence handled accordingly.


```

df.dtypes
[7] ✓ 0.0s
...
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

df.isnull().sum()
[8] ✓ 0.0s
...
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

df.duplicated().sum()
[9] ✓ 0.0s
...
<bound method Series.sum of 0    False
1      False
2      False
3      False
4      False
...
7038   False
7039   False
7040   False
7041   False
7042   False
Length: 7043, dtype: bool>

df.duplicated().sum()
[10] ✓ 0.0s
...
0

```

To handle null values, **SimpleImputer** tool of the **sklearn** library is used. It is a univariate imputer for completing missing values with simple strategies. It replaces missing values using a descriptive statistic (e.g. mean, median, or most frequent) along each column, or using a constant value.

```

from sklearn.impute import SimpleImputer

# The imputer will replace missing values with the mean of the non-missing values for the respective columns
imputer = SimpleImputer(missing_values=np.nan, strategy='mean')

df1.TotalCharges = imputer.fit_transform(df1["TotalCharges"].values.reshape(-1, 1))
[28] ✓ 0.2s Python

from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
[29] ✓ 0.0s Python

```

The dataset which is now ready for applying the Classification Models is made ready and is described for the final time. It is encoded using **One Hot Encoding** - a method for converting categorical variables into a binary format. It creates new columns for each category where 1 means the category is present and 0 means it is not. The primary purpose of One Hot Encoding is to ensure that categorical data can be effectively used in machine learning models.

```

df1=pd.get_dummies(data=df,columns=['gender', 'Partner', 'Dependents',
'PhoneService', 'MultipleLines', 'InternetService', 'OnlineSecurity',
'OnlineBackup', 'DeviceProtection', 'TechSupport', 'StreamingTV',
'StreamingMovies', 'Contract', 'PaperlessBilling', 'PaymentMethod', 'Churn'], drop_first=True)
[23] ✓ 0.0s Python

df1.head()
[24] ✓ 0.0s Python
...

```

	customerID	SeniorCitizen	tenure	MonthlyCharges	TotalCharges	Churn_numeric	gender_Male	Partner_Yes	Dependents_Yes	PhoneService_Yes	...	StreamingTV_Y
0	7590-VHVEG	0	1	29.85	29.85	0	False	True	False	False	...	Fal
1	5575-GNVDE	0	34	56.95	1889.50	0	True	False	False	True	...	Fal
2	3668-QPYBK	0	2	53.85	108.15	1	True	False	False	True	...	Fal
3	7795-CFOCW	0	45	42.30	1840.75	0	True	False	False	False	...	Fal
4	9237-HQITU	0	2	70.70	151.65	1	False	False	False	True	...	Fal

5 rows x 33 columns

D. APPLYING THE 5 MACHINE LEARNING CLASSIFICATION MODELS

Before applying the models one by one, the entire dataset was standardised and split using the **Train_Test_Split** function. **70%** of the data was randomised and kept as the training data. Whereas the remaining **30%** was treated as the Testing data.

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()

[86] ✓ 0.0s Python

scaler.fit(df1.drop(['Churn_Yes'],axis = 1))
scaled_features = scaler.transform(df1.drop(['Churn_Yes'],axis = 1))

[87] ✓ 0.0s Python

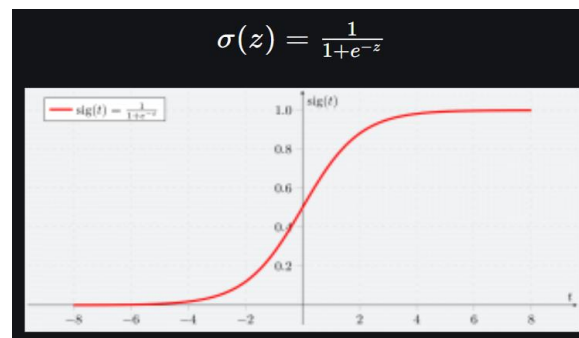
X = scaled_features
Y = df1['Churn_Yes']
X_train,X_test,Y_train,Y_test = train_test_split(X,Y,test_size = 0.3,random_state=44)

[88] ✓ 0.0s Python
```

With all the pre-requisites already understood, let us now discuss all the models one by one.

D.1. Logistic Regression

Logistic regression is a supervised machine learning algorithm used for classification tasks where the goal is to predict the probability that an instance belongs to a given class or not. Logistic regression is a statistical algorithm which analyse the relationship between two data factors. Logistic regression is used for binary classification where we use sigmoid function, that takes input as independent variables and produces a probability value between 0 and 1.



As per the syntax, Logistic Regression was applied on the train and test data, without tuning the hyperparameters at all. On forming the classification report, the train and test data accuracies were **80.65% and 80.03%** respectively, which showed that the dataset is neither underfit nor an overfit.

```

logmodel = logisticRegression()
logmodel.fit(X_train,Y_train)
predict_train=logmodel.predict(X_train)
predict_test=logmodel.predict(X_test)

acc_train=accuracy_score(Y_train, predict_train)
acc_test=accuracy_score(Y_test, predict_test)
print('Train data Accuracy Prediction:',round(acc_train*100,2),'%')
print('Test data Accuracy Prediction:',round(acc_test*100,2),'%')
print('\nConfusion Matrix:\n', confusion_matrix(Y_test, predict_test))
print('\nClassification Report:\n', classification_report(Y_test, predict_test))

```

```

139 | ✓ 0.0s
... | Train data Accuracy Prediction: 80.65 %
    | Test data Accuracy Prediction: 80.83 %
    |
    | Confusion Matrix:
    | [[1397 160]
    | [ 262 294]]
    |
    | Classification Report:
    |
    | precision    recall  f1-score   support
    |
    | False      0.84     0.90     0.87    1557
    | True       0.65     0.53     0.58     556
    |
    | accuracy    0.80     0.80    2113
    | macro avg   0.74     0.71     0.73    2113
    | weighted avg 0.79     0.80     0.79    2113

```

As per the classification report, the parameters are plotted on the major target categories- **Churn_No** and **Churn_Yes**. The features include-

Precision: measures the proportion of true positive predictions among all positive predictions made by the model. It is calculated as the ratio of true positives to the sum of true positives and false positives.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

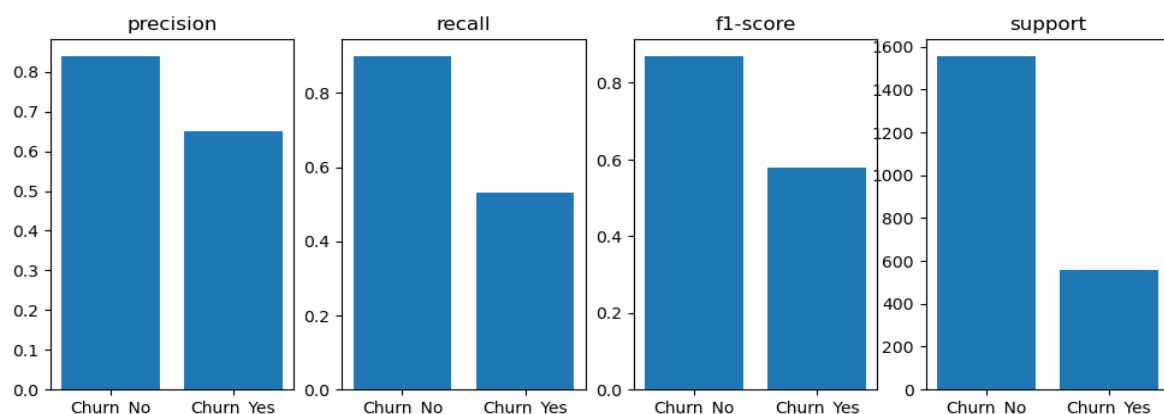
Recall (Sensitivity): measures the proportion of true positive predictions among all actual positive instances in the dataset. It is calculated as the ratio of true positives to the sum of true positives and false negatives.

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

F1-score: is the harmonic mean of precision and recall. It provides a single metric that combines both precision and recall, balancing the trade-off between the two. It is calculated as:

$$\text{F1-score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

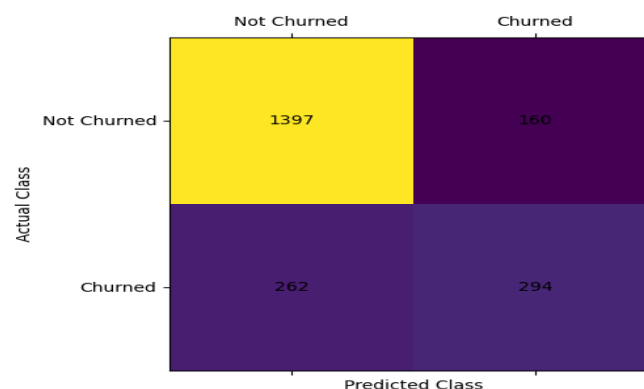
Support: refers to the number of actual occurrences of the class in the dataset. It is the number of instances in each class.



Along with the bar plots, confusion matrix is also drawn. It evaluates a classification model predicting customer churn. It consists of four sections:

- **True Negatives (1397):** Correctly predicted as "Not Churned."
- **False Positives (160):** Incorrectly predicted as "Churned" when they actually did not.
- **False Negatives (262):** Incorrectly predicted as "Not Churned" when they actually did.
- **True Positives (294):** Correctly predicted as "Churned."

This matrix helps assess performance, revealing its accuracy, misclassification rates, and areas for improvement of the logistic regression model used.



To Optimize the Performance of Logistic Regression hyperparameters can be tuned. **GridSearchCV** is a method provided by the scikit-learn library in Python, which is used to systematically work through multiple combinations of parameter tunes, cross-validating as it goes to determine which tune gives the best performance.

The parametric values used are:

- **L1 Regularization:** it is regularization technique which is used to reduce the values of the coefficients of the equation. This regularization technique is used as feature elimination method.
- **L2 Regularization:** it is a regularization technique which is used to reduce the values of the coefficients not exactly but nearly to zero.
- **Elastic Net Regularization:** Elastic Net regularization is a combination of both L1 and L2 regularization techniques. It adds both the L1 and L2 penalty terms to the loss function during training.

Solvers algorithms are also used to optimize the model's performance by finding the best-fit parameters that minimize the loss function.

The **max_iter** parameter specifies the maximum number of iterations that the solver is allowed to take in order to converge to a solution, it being 1000 in our case.

The output train and test accuracies received after varying many hyperparameters is noted and hence a line graph is also plotted to clearly depict the knowledge.

```

... Fitting 5 folds for each of 56 candidates, totalling 280 fits
C = 0.001: Train Accuracy = 79.27%, Test Accuracy = 78.51%
C = 0.01: Train Accuracy = 80.51%, Test Accuracy = 80.50%
C = 0.1: Train Accuracy = 80.57%, Test Accuracy = 80.22%
C = 1: Train Accuracy = 80.65%, Test Accuracy = 80.03%
C = 10: Train Accuracy = 80.65%, Test Accuracy = 80.03%
C = 100: Train Accuracy = 80.59%, Test Accuracy = 80.03%
C = 1000: Train Accuracy = 80.57%, Test Accuracy = 80.03%

Best Model Accuracy after tuning:
Train data Accuracy Prediction: 80.59 %
Test data Accuracy Prediction: 80.03 %

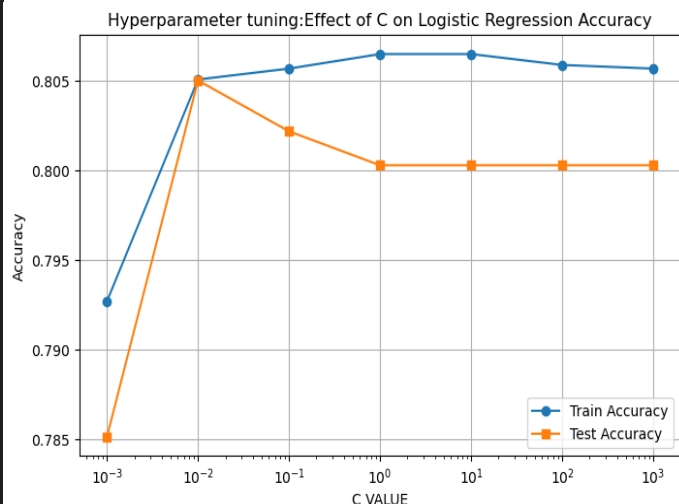
Confusion Matrix:
[[1400 157]
 [ 265 291]]

Classification Report:
              precision    recall  f1-score   support

     False      0.84      0.90      0.87     1557
     True       0.65      0.52      0.58      556

 accuracy      0.80      0.80      0.80     2113
 macro avg     0.75      0.71      0.72     2113
 weighted avg  0.79      0.80      0.79     2113

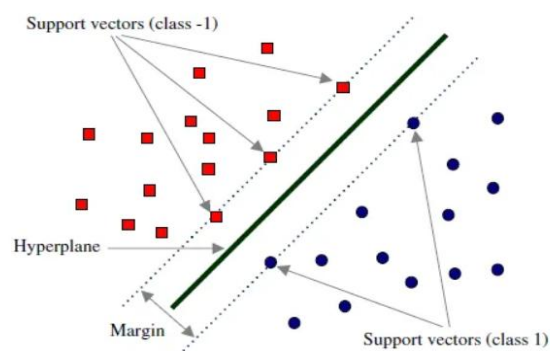
```



Hence, the best model accuracy was similar to the one calculated without tuning the hyperparameter.

D.2. Support Vector Classification

In machine learning, SVC stands for **Support Vector Classifier**, which is a type of supervised learning model used for classification tasks. Specifically, SVC is a variant of **Support Vector Machine (SVM)** that is tailored for classification. SVC is widely used in various domains, including text classification, image classification, bioinformatics, and more. It is particularly effective for binary classification tasks and can be extended to handle multi-class classification using strategies such as one-vs-one or one-vs-rest.



The data was loaded and extracted in the same way as was done for the logistic regression model and the classification report is also generated.

```

svc = SVC()
svc.fit(X_train, Y_train)

y_pred_svc = svc.predict(X_test)
print(classification_report(Y_test, y_pred_svc))
s=svc.score(X_train,Y_train)
print("svc score= ",s)
acc=accuracy_score(Y_test, y_pred_svc)
print("accuracy score=",acc)

```

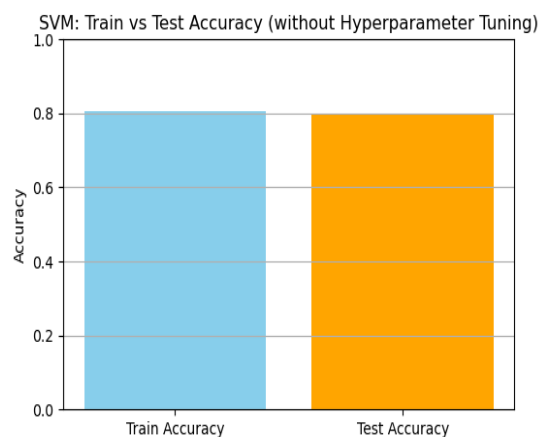
[36] ✓ 1.4s

	precision	recall	f1-score	support
False	0.83	0.92	0.87	1557
True	0.67	0.48	0.56	556
accuracy			0.80	2113
macro avg	0.75	0.70	0.71	2113
weighted avg	0.79	0.80	0.79	2113

svc score= 0.8170385395537525
accuracy score= 0.8012304779933743

It is noted that the train and test accuracies without tuning the hyperparameters at all was **81.70% and 80.12%**, which is slightly more than that of the Logistic regression model. This is because it finds an optimal decision boundary by maximizing the margin between classes, making it more effective for complex and high-dimensional data.

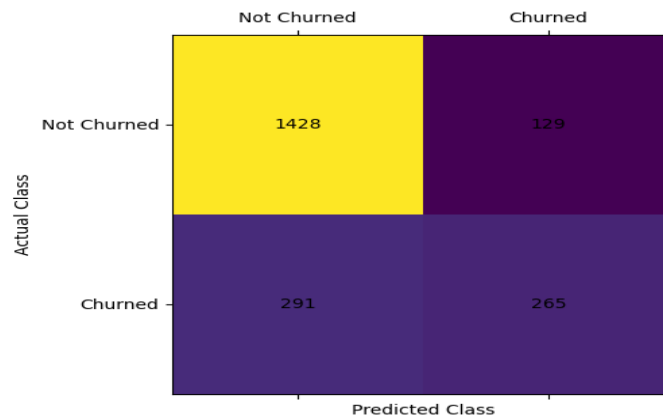
A simple bar graph is plotted to visualise the train and test accuracies for both categories of the classification target variable.



Along with the bar plot, confusion matrix is also drawn which evaluates the SVM classifier predicting customer churn. It shows:

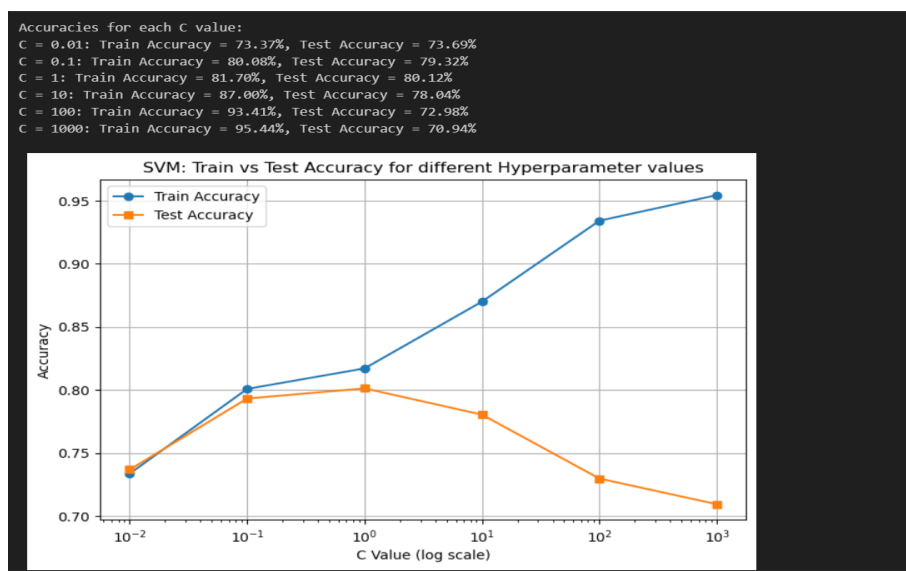
- **True Negatives (1397):** Correctly classified as "Not Churned."
- **False Positives (160):** Incorrectly predicted as "Churned."
- **False Negatives (262):** Incorrectly predicted as "Not Churned."
- **True Positives (294):** Correctly classified as "Churned."

It helps assess the model's accuracy, precision, and recall, showing how well it differentiates between churned and non-churned customers.



SVM also has some hyper-parameters (like what C or gamma values to use). **GridSearchCV** takes a dictionary that describes the parameters that could be tried on a model to train it. The grid of parameters is defined as a dictionary, where the keys are the parameters and the values are the settings to be tested.

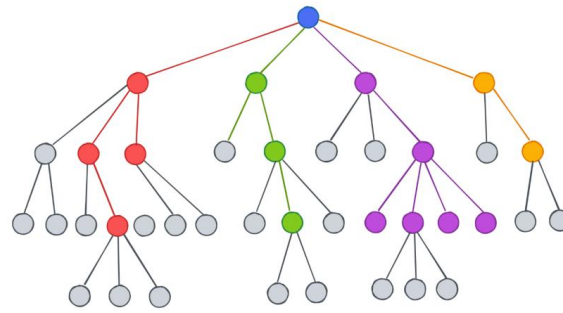
The output train and test accuracies received after varying many C Values is noted and hence a line graph is also plotted to clearly depict the knowledge gained.



The trend between the train and test accuracies after tuning the hyperparameters can be seen above.

D.3. Decision Tree

A decision tree is one of the most powerful tools of supervised learning algorithms used for both classification and regression tasks. It builds a flowchart-like tree structure where each internal node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (terminal node) holds a class label. It is constructed by recursively splitting the training data into subsets based on the values of the attributes until a stopping criterion is met, such as the maximum depth of the tree or the minimum number of samples required to split a node.



The **Impurity and Entropy** is calculated so as to find a root node as well as generate the entire tree, using the **ID3** algorithm.

$$Entropy = - \sum_{i=1}^n p_i * \log(p_i)$$

The data was loaded and extracted in the same way as was done for logistic regression, and SVM models and the classification report is also generated.

```
dtc = DecisionTreeClassifier()

dtc.fit(X_train, Y_train)
y_pred_dtc = dtc.predict(X_test)
print(classification_report(Y_test, y_pred_dtc))

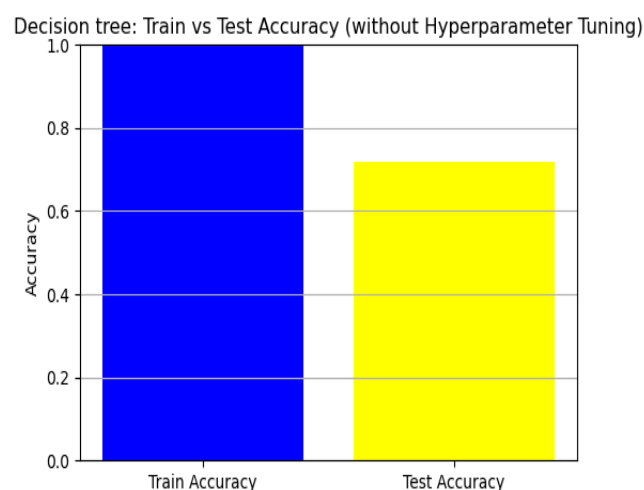
D=dtc.score(X_train,Y_train)
print("DTC score= ",D)
A=accuracy_score(Y_test, y_pred_dtc)
print("accuracy score= ",A)
```

[40] ✓ 0.0s Python

	precision	recall	f1-score	support
False	0.81	0.81	0.81	1557
True	0.47	0.48	0.47	556
accuracy			0.72	2113
macro avg	0.64	0.64	0.64	2113
weighted avg	0.72	0.72	0.72	2113

DTC score= 0.9987829614604462
accuracy score= 0.7188831045906294

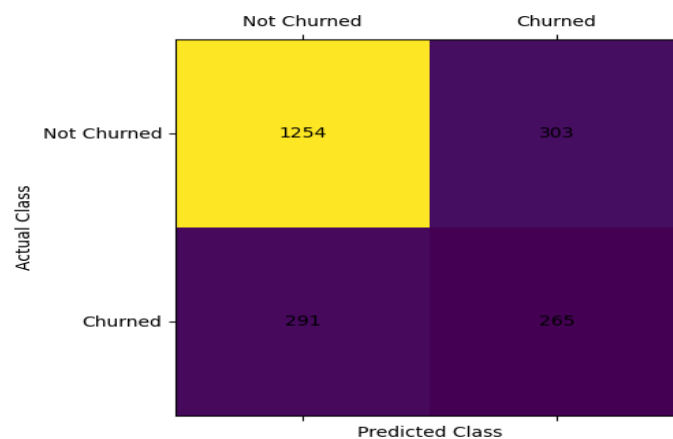
Here, the train data accuracy without tuning any hyperparameters was **99.87%**, while that of testing data was only **71.88%**. This clearly depicted a case of **overfitting** of the model. Still the Bar plot of these accuracies is obtained.



A confusion matrix was also plot to evaluate the **Decision Tree classifier** predicting customer churn. It consists of:

- **True Negatives (1397)**: Correctly classified as "Not Churned."
- **False Positives (160)**: Incorrectly predicted as "Churned."
- **False Negatives (262)**: Incorrectly predicted as "Not Churned."
- **True Positives (294)**: Correctly classified as "Churned."

Decision trees tend to capture complex patterns well but can be prone to overfitting. This matrix highlights where the model performs well and where it misclassifies churned customers.

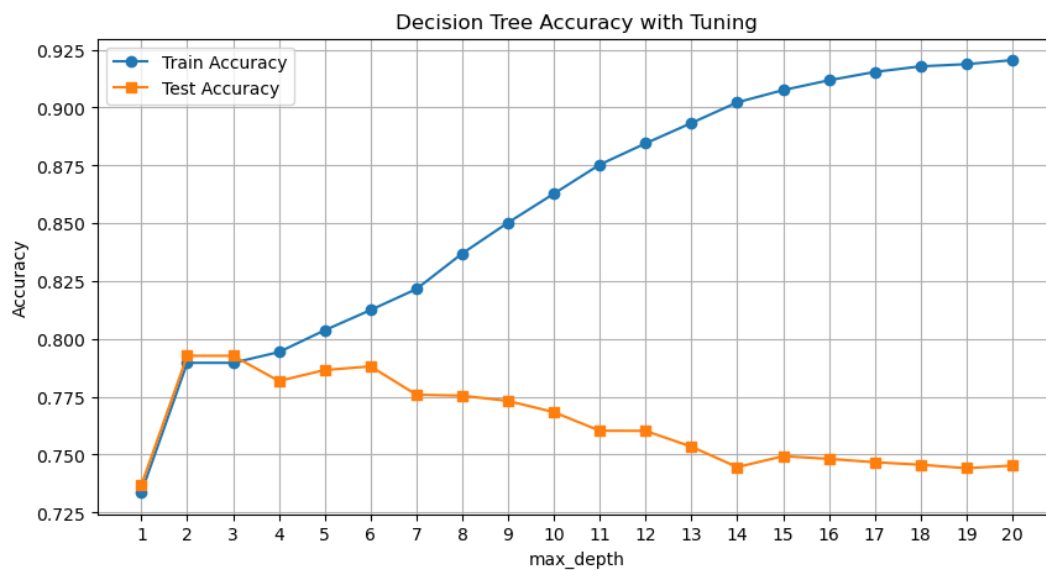


Now, to avoid the overfitting occurred, the hyperparameters were tuned here as well. They control the behaviour and the structure of the model during the training phase.

*The **max_depth** hyperparameter* controls the maximum depth to which the decision tree is allowed to grow. It is very important to tune *max_depth* carefully to find the right balance between model complexity and generalization performance.

The output train and test accuracies received after varying many Max_Depth hyperparameter Values is noted and hence a line graph is also plotted to clearly depict the knowledge gained

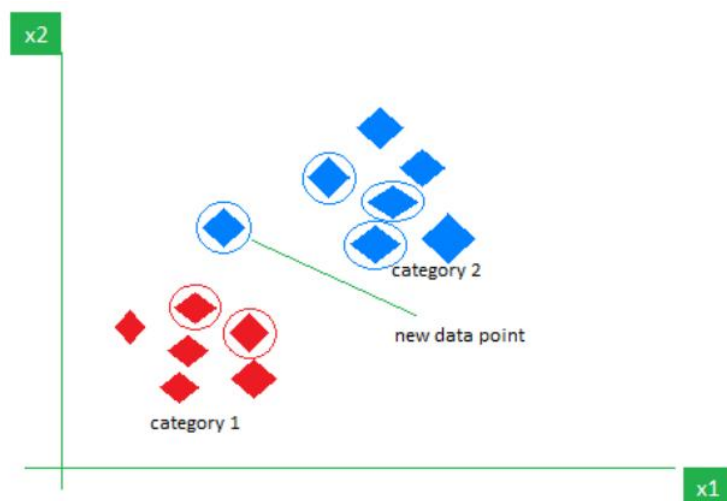
```
... max_depth = 1: Train Accuracy = 0.734, Test Accuracy = 0.737
max_depth = 2: Train Accuracy = 0.790, Test Accuracy = 0.793
max_depth = 3: Train Accuracy = 0.790, Test Accuracy = 0.793
max_depth = 4: Train Accuracy = 0.794, Test Accuracy = 0.782
max_depth = 5: Train Accuracy = 0.804, Test Accuracy = 0.787
max_depth = 6: Train Accuracy = 0.812, Test Accuracy = 0.788
max_depth = 7: Train Accuracy = 0.822, Test Accuracy = 0.776
max_depth = 8: Train Accuracy = 0.837, Test Accuracy = 0.775
max_depth = 9: Train Accuracy = 0.850, Test Accuracy = 0.773
max_depth = 10: Train Accuracy = 0.863, Test Accuracy = 0.768
max_depth = 11: Train Accuracy = 0.875, Test Accuracy = 0.760
max_depth = 12: Train Accuracy = 0.884, Test Accuracy = 0.760
max_depth = 13: Train Accuracy = 0.893, Test Accuracy = 0.753
max_depth = 14: Train Accuracy = 0.902, Test Accuracy = 0.745
max_depth = 15: Train Accuracy = 0.908, Test Accuracy = 0.749
max_depth = 16: Train Accuracy = 0.912, Test Accuracy = 0.748
max_depth = 17: Train Accuracy = 0.915, Test Accuracy = 0.747
max_depth = 18: Train Accuracy = 0.918, Test Accuracy = 0.746
max_depth = 19: Train Accuracy = 0.919, Test Accuracy = 0.744
max_depth = 20: Train Accuracy = 0.920, Test Accuracy = 0.745
```



The trend between the train and test accuracies after tuning the hyperparameters can be seen above and the model which was originally an overfit is now bettered up.

D.4. K Nearest Neighbours (KNN)

The K-Nearest Neighbours (KNN) algorithm is a supervised machine learning method employed to tackle classification and regression problems. Evelyn Fix and Joseph Hodges developed this algorithm in 1951, which was subsequently expanded by Thomas Cover.



The data was loaded and extracted in the same way as was done for the previous models and the classification report is also generated, keeping the count of nearest **datapoints to be 30**.

```

knn = KNeighborsClassifier(n_neighbors = 30)
knn.fit(X_train,Y_train)
pred_knn = knn.predict(X_test)

[44] ✓ 0.2s

print(classification_report(Y_test,pred_knn))
k=knn.score(X_train,Y_train)
print("knn score=",k)
a=accuracy_score(Y_test, pred_knn)
print("accuracy score=",a)

[45] ✓ 0.0s

...      precision    recall  f1-score   support

      False    0.84    0.88    0.86    1557
       True    0.62    0.55    0.58    556

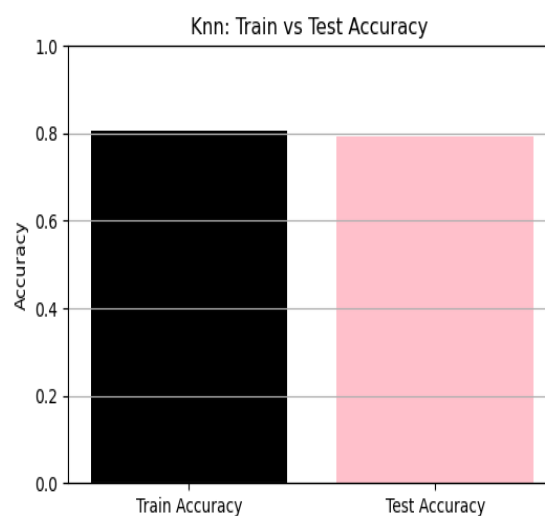
 accuracy          0.79    2113
  macro avg       0.73    0.71    0.72    2113
  weighted avg     0.79    0.79    0.79    2113

knn score= 0.804868154158215
accuracy score= 0.792238523426408

```

Here, the train and test accuracies are **80.48% and 79.22%** respectively, showing no major signs of overfitting or underfitting.

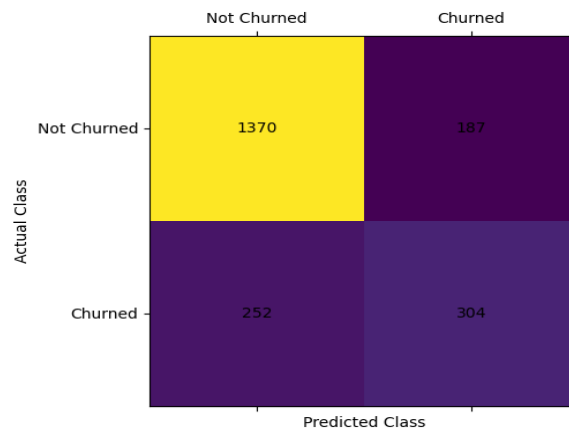
The bar plot is also obtained to visualise the same.



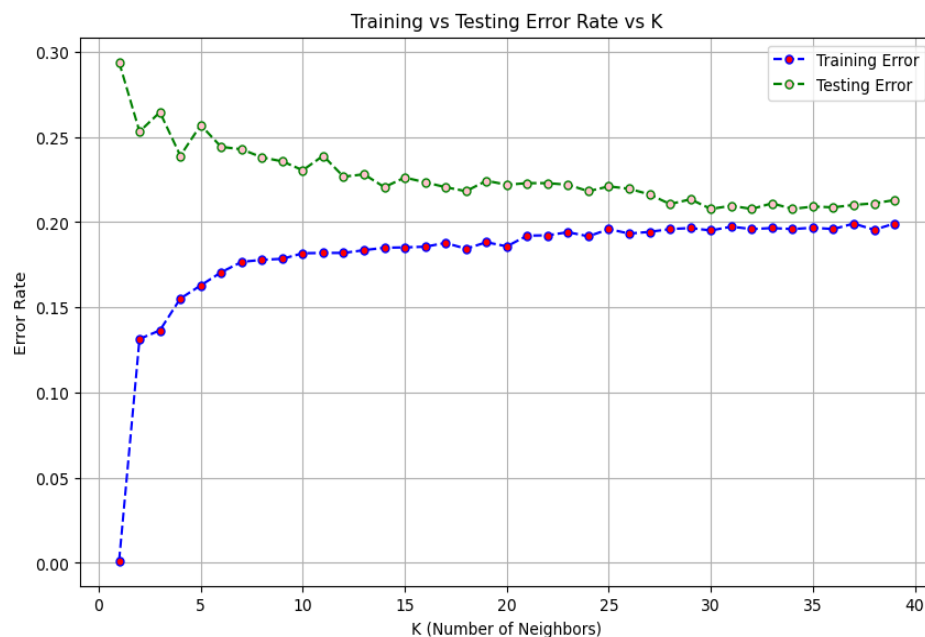
The confusion matrix evaluates the **K-Nearest Neighbours (KNN) classifier** predicting customer churn. It consists of:

- **True Negatives (1397):** Correctly classified as "Not Churned."
- **False Positives (160):** Incorrectly predicted as "Churned."
- **False Negatives (262):** Incorrectly predicted as "Not Churned."
- **True Positives (294):** Correctly classified as "Churned."

KNN makes predictions based on the majority class among its nearest neighbours. This matrix helps analyse its performance, including accuracy and misclassification rates.

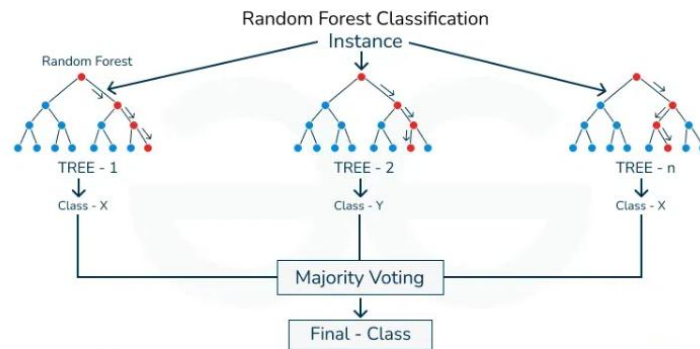


To see further changes in the accuracies and the error rates, the hyperparameters were tuned. In this case, the value of k acts as hyperparameter. It should neither be too large such that it gives a wrong classification, nor be too small also. Here, varying the k value between 1 to 40 and plotting the graph for train and test errors which occurred helped identify the trend and know the accuracies in a better way.



D.5 Random Forest

Random Forest Classification predicts categorical outcomes based on the input data. It uses multiple decision trees and outputs the label that has the maximum votes among all the individual tree predictions. Random Forest can handle large datasets and high-dimensional data. By combining predictions from many decision trees, it reduces the risk of overfitting compared to a single decision tree. It is robust to noisy data and works well with categorical data.



The data was loaded and extracted in the same way as done for the previous models and the classification report is also generated.

```

rf = RandomForestClassifier(n_estimators=100, random_state=42)
rf.fit(X_train, Y_train)

pred_rf = rf.predict(X_test)

[49] ✓ 0.5s

print(classification_report(Y_test, pred_rf))

train_score = rf.score(X_train, Y_train)
print("Random Forest Train Score:", train_score)

test_accuracy = accuracy_score(Y_test, pred_rf)
print("Random Forest Test Accuracy:", test_accuracy)

[50] ✓ 0.0s

...

```

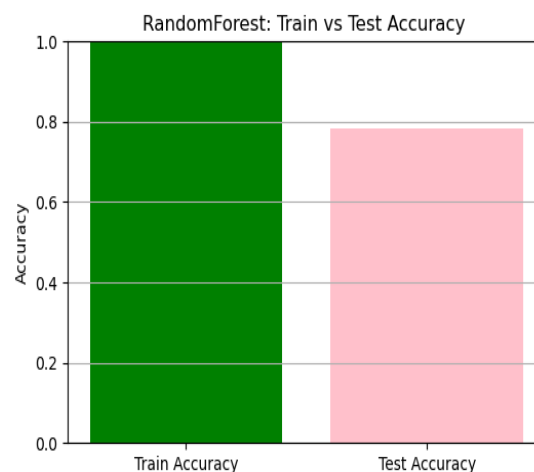
	precision	recall	f1-score	support
False	0.82	0.90	0.86	1557
True	0.62	0.46	0.53	556
accuracy			0.78	2113
macro avg	0.72	0.68	0.69	2113
weighted avg	0.77	0.78	0.77	2113

```

Random Forest Train Score: 0.9987829614604462
Random Forest Test Accuracy: 0.7841930903928065

```

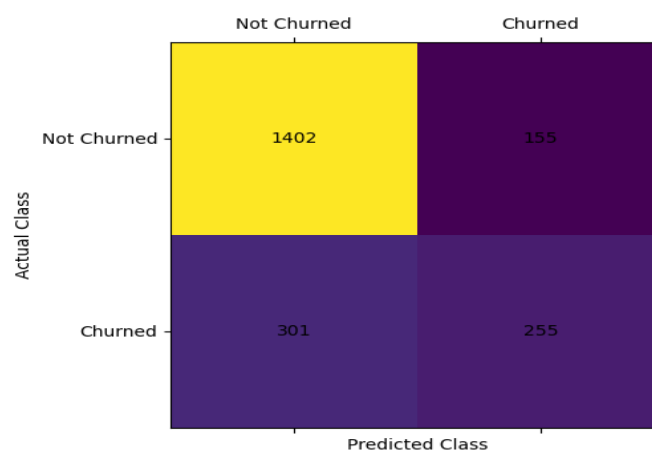
Here, the train data accuracy without tuning any hyperparameters was **99.87%**, while that of testing data was only **78.4%**, similar to the trend obtained in the decision tree model. This clearly depicted a case of **overfitting** of the model. Still the Bar plot of these accuracies is obtained.



The confusion matrix evaluates the **Random Forest classifier** predicting customer churn. It consists of:

- **True Negatives (1397):** Correctly classified as "Not Churned."
- **False Positives (160):** Incorrectly predicted as "Churned."
- **False Negatives (262):** Incorrectly predicted as "Not Churned."
- **True Positives (294):** Correctly classified as "Churned."

Random Forest enhances predictive accuracy by aggregating multiple decision trees, reducing overfitting while maintaining interpretability. This matrix shows its strengths and misclassification patterns.



While Random Forest is already a robust model fine-tuning its hyperparameters such as the number of trees, maximum depth and feature selection can improve its prediction. **RandomisedSearchCV** is used which randomly selects combinations and evaluates the model often leading to faster results especially when there are many hyperparameters.

n_estimators: Defines the number of trees in the forest. More trees typically improve model performance but increase computational cost.

max_features: Limits the number of features to consider when splitting a node. This helps control overfitting.

max_depth: Controls the maximum depth of each tree. A shallow tree may underfit while a deep tree may overfit.

max_leaf_nodes: Limits the number of leaf nodes in the tree hence controlling its size and complexity.

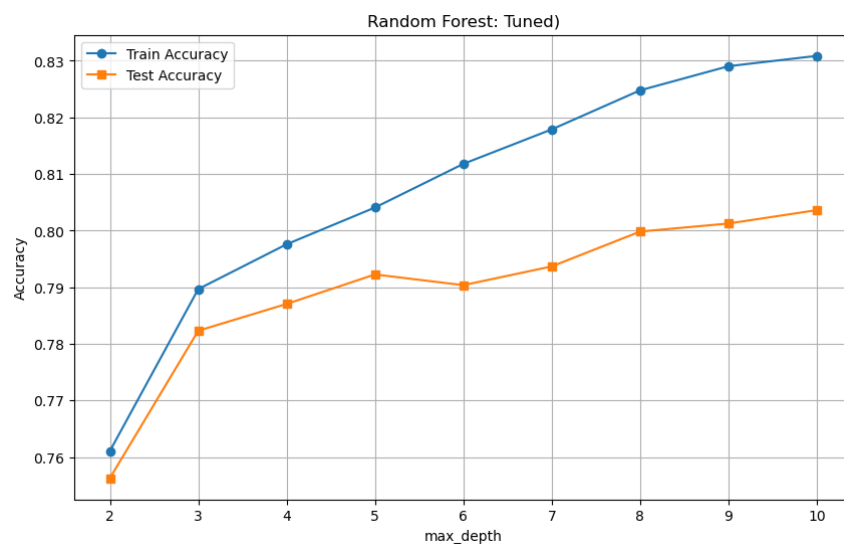
max_sample: Apart from the features, we have a large set of training datasets. max_sample determines how much of the dataset is given to each individual tree.

min_sample_split: Specifies the minimum number of samples required to split an internal node.

```
Fitting 3 folds for each of 20 candidates, totalling 60 fits
Best hyperparameters: {'n_estimators': 100, 'min_samples_split': 10, 'min_samples_leaf': 10, 'max_features': 'sqrt', 'max_depth': 10, 'bootstrap': True}

Final Train Accuracy: 0.8308
Final Test Accuracy: 0.8036
max_depth = 2: Train Accuracy = 0.7611, Test Accuracy = 0.7563
max_depth = 3: Train Accuracy = 0.7897, Test Accuracy = 0.7823
max_depth = 4: Train Accuracy = 0.7976, Test Accuracy = 0.7870
max_depth = 5: Train Accuracy = 0.8041, Test Accuracy = 0.7922
max_depth = 6: Train Accuracy = 0.8118, Test Accuracy = 0.7903
max_depth = 7: Train Accuracy = 0.8178, Test Accuracy = 0.7937
max_depth = 8: Train Accuracy = 0.8247, Test Accuracy = 0.7998
max_depth = 9: Train Accuracy = 0.8290, Test Accuracy = 0.8012
max_depth = 10: Train Accuracy = 0.8308, Test Accuracy = 0.8036
```

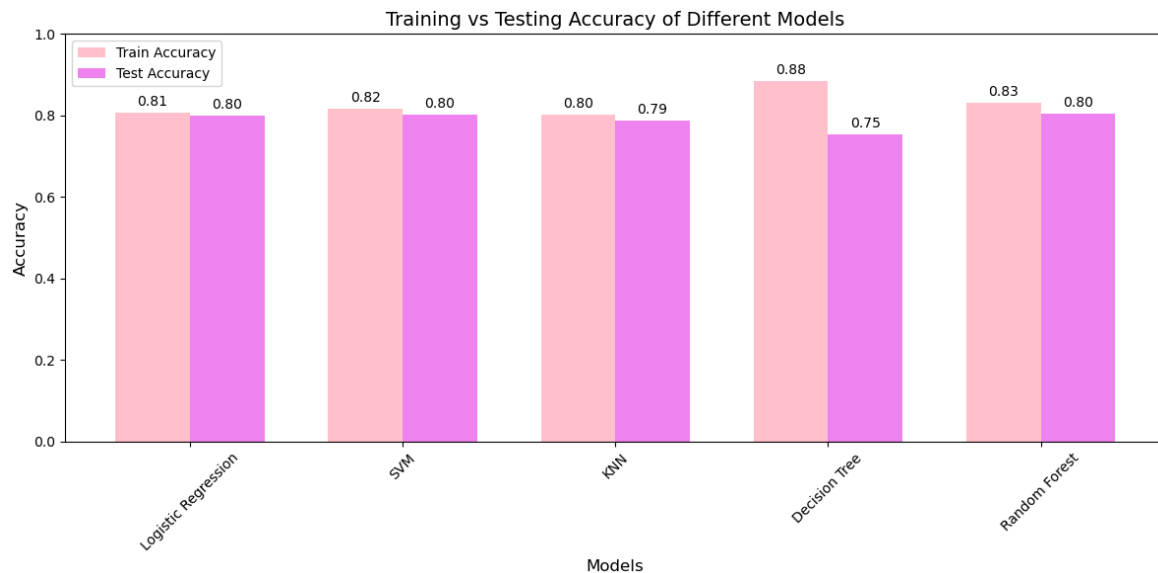
Thus, tuning the hyperparameters improved the accuracy and avoided the model to overfit. The points are plot and similar to other models, the trends can be understood.



The final accuracies were thus **under 85%** respectively.

RESULT AND DISCUSSION

After having applied all the models, and to have a comparative study for all the models clearly, bar plots are made using double bar graph to show train and test accuracies for each of the 5 models.



Among all models, Decision tree classifier achieved an accuracy of 88% for the training data, though its testing accuracy was just 75%, showing that it is slightly an overfit model. But on the testing data, the highest accuracy attained was 80%, shown by models like logistic regression, SVM and Random Forest.

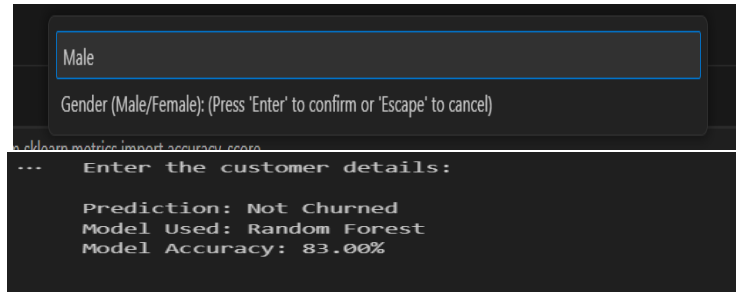
The other models except Decision Tree show a very small gap between the training and testing data, reflecting more consistency and reliability. Overall, the performance differences are very small, but SVM and Random Forest models provide the best trade-off between accuracy and overfitting.

CONCLUSION

ML Classification techniques were effectively applied on the real-world dataset so as to predict the customer churn in the telecom sector. After preprocessing and exploring the data, five classification models—**Logistic Regression, SVM, KNN, Decision Tree, and Random Forest**—were trained and evaluated. SVM and Random Forest models showed the most balanced performance. Hyperparameter Tuning for all the models also helped improve the efficiency and remove any chances of underfitting or overfitting.

CLASSIFIER	ACCURACY ON TRAINING DATA	ACCURACY ON TESTING DATA
Logistic Regression	0.81	0.80
Support vector machine	0.82	0.80
K Nearest Neighbours	0.80	0.79
Decision tree	0.88	0.75
Random Forest	0.83	0.80

Additionally, a very basic user input enabled section also helps the end user to input various features as asked. Based on the inputs received, the system identifies the best suited model and predicts the final output as to whether there will be a churn or not. Along with the result prediction, the name of the model and accuracy for the prediction obtained is also shown as output so as to give better clarity to the user.



```
Male
Gender (Male/Female): (Press 'Enter' to confirm or 'Escape' to cancel)

... Enter the customer details:
Prediction: Not Churned
Model Used: Random Forest
Model Accuracy: 83.00%
```

This work demonstrates that with the right preprocessing and model selection, machine learning can be a powerful tool for customer retention strategy in the telecom industry.

FUTURE WORK

The first and foremost thing is to improve all the models and try to achieve as much accuracy as possible. There is a possibility that one model may be overpower the other in the future. The project, after completion can be deployed and be used to identify the customer churn not only in the telecommunications industry, but also in many other sections based on the features and the consumer behaviour. It also has potential to reduce the churn and make the industries more profitable in the long term.

REFERENCES

- <https://www.kaggle.com/datasets/blastchar/telco-customer-churn/data>
- https://www.researchgate.net/publication/387816970_Predicting_Customer_Churn_in_Telecommunications_with_Machine_Learning_Models
- https://www.researchgate.net/publication/386372608_Telecommunication_Customer_Churn_with_Responsibile_AI_A_Predictive_Model_Debugging_and_Business_Decision_Making
- <https://www.youtube.com/watch?v=9iD8DMF6odw>
- <https://www.youtube.com/watch?v=n40hS9tQmcY&t=357s>
- <https://www.geeksforgeeks.org/decision-tree/>
- https://www.w3schools.com/python/pandas/pandas_intro.asp
- <https://www.geeksforgeeks.org/random-forest-hyperparameter-tuning-in-python/>
- <https://www.ibm.com/think/topics/random-forest>
- <https://scikitlearn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>
- <https://medium.com/data-science/k-nearest-neighbor-classifier-explained-a-visual-guide-with-code-examples-for-beginners-a3d85cad00e1>