

Customer Churn Analysis in the Telecom Industry

Akshay Danthi (0606752)

Uttam Gogineni (0605543)

MS in Business Analytics

Golden Gate University

INDEX

Executive Summary	3
Introduction	3
Data Procurement	4
Exploratory Data Analysis	5
Data Dictionary	5
Correlation Matrix	7
Univariate Analysis	8
Bivariate Analysis	11
Statistical Summary table	13
Observing churn distribution along other variables	13
Feature Engineering steps	14
Model Building	15
Scaling with StandardScaler()	15
Addressing Class Imbalance with ADASYN	15
Applying Logistic Regression	16
Applying KNN-Classifier	17
Applying Decision-Tree classifier	18
Applying Random Forest classifier	20
Applying Naïve Bias classifier	22
Applying XG-Boost classifier	23
Result Comparison	25
Conclusion	26
Key Insights from this paper	26
Business Implications and Recommendations	26
References	28

Executive Summary

In this competitive world with the growth of businesses and saturating the market, especially with multiple companies in the telecom industry fighting for market share, this causes complex challenges due to multiple service providers providing vibrant services. This becomes tough for companies to retain existing customers and stop customers to leave their providers. It costs lot more to acquire new customers it is needed to ensure that the retention level of customers is always high to justify the spending for customer acquisition costs.

Customer Churn in telecommunication is one of the most pressing matters in customer retention and generating revenue. It also affects the customer relationships with the service provider. Development of effective churn prediction infrastructure is important to improve customer acquisition costs and improve customer retainment. This study presents you a review of customer churn, behaviour, business reasons and introduces machine learning models to study customer churn behaviour and major reasons for customer churn.

Introduction

Telecom industry is the most important industry for modern communication, that has evolved since its inception from telegraphs to the current hand-held devices that provide ultra high-speed internet all over the country and the world. It has played an important role in shaping the communication infrastructure and had a profound impact on society and how we communicate with the world. Telecom industry is 1.85 trillion market in 2022 and is expected to surpass 2.65 trillion USD by 2030 with a registered CGAR of 4.85%. and in USA it is 427.43 billion USD and is expected to grow to 511 billion USD by 2028 at a CGAR of 3.67%.

The telecom industry started in 19th century with the invention of telegraph, which revolutionized the long-distance communication removing the need for unreliable letters. And telecom sector further boomed after telephone was invented by alexander graham bell enabling real-time voice communication. Now with the current technological improvements we can notice the changes that make our current telecom services that make day to day tasks with ease.

With the rise in competition from multinational companies the customer churn has become a critical challenge for the service providers. Customer churn refers to customers voluntarily cancelling their telecom service. As it influences market share and deteriorates the company's revenue and health of the telecom business. The major reasons for the customer churn include service quality, increased competition, customer dissatisfaction with service, change in customer preferences, lack of personalization, offerings from the rival providers.

Due to high churn rates the major loss is going to be the reduced revenue, long term effects include reduced brand reputation, customer loyalty and loss of market share. This causes the company to lose its market position and go into bankruptcy as their income dwindles.

To avoid it telecom operators, need to improve their service quality, provide personalized offers, better customer support, better pricing from customers and better offerings, use to machine learning algorithms to observe customer behaviour and provide better offerings based on the customer behaviour.

Data Procurement

The data sources of the churn data are difficult to procure as the telecom companies keep their data extremely confidential and do not let anyone have access to their databases and secure their databases and would go to extreme lengths to ensure their data does not fall into their competitor's hands. So, for this paper we will be working on data that has been collected from <https://www.kaggle.com/datasets/blastchar/telco-customer-churn>

A data set with 22 Features and 7043 records out of which many features are considered extremely important and help us understand the nuances of the customer behaviour. And their spending habits on their plans and requirements so the telecom companies can work on understanding customer's needs to improve their products.

The sample of the dataset:

```
data.head()
```

	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	...	TechSupport	StreamingTV
0	7590-VHVEG	Female	0	Yes	No	1	No	No phone service	DSL	No	...	No	No
1	5575-GNVDE	Male	0	No	No	34	Yes	No	DSL	Yes	...	No	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	No
4	9237-HQITU	Female	0	No	No	2	Yes	No	Fiber optic	No	...	No	No

5 rows × 22 columns

Exploratory Data Analysis

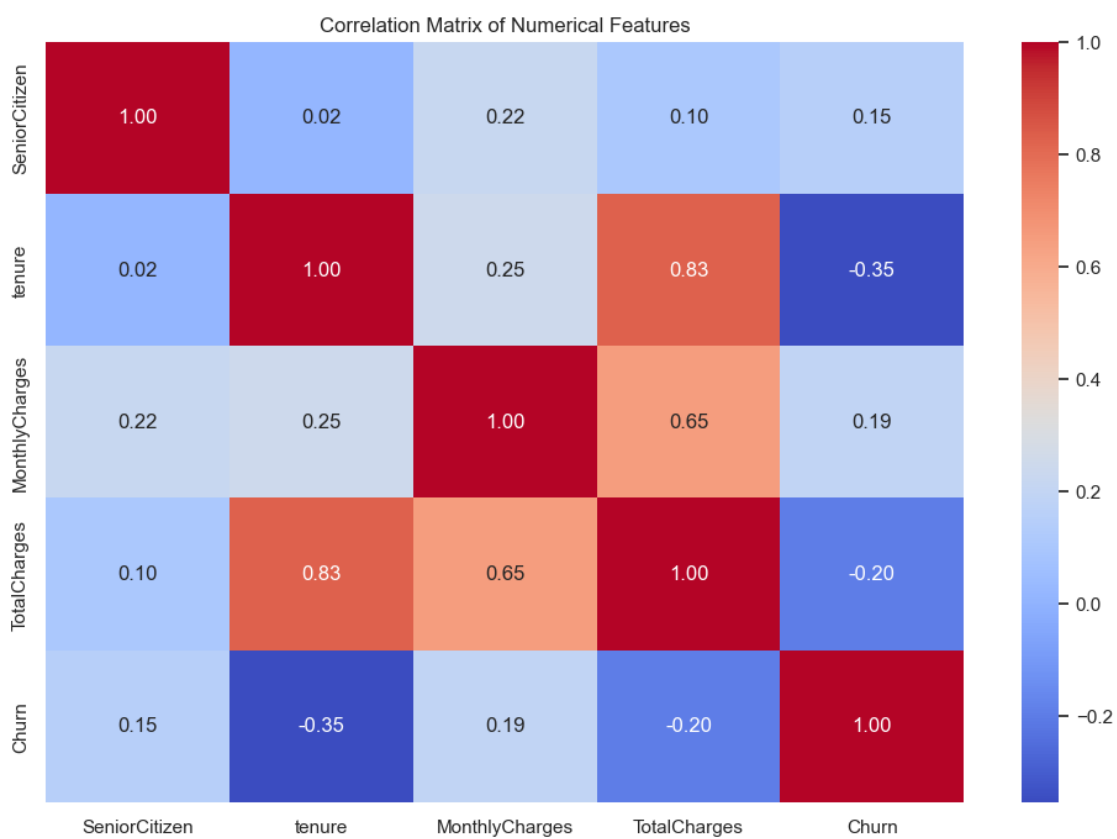
Data Dictionary of the data set:

Column Name	Data Type	Definition
customerID	Character	Unique id of the customer
gender	character	Sex of the customer
SeniorCitizen	integer	Flag to check if the customer is a senior citizen
Partner	Character	A yes or no flag to check if the customer has a partner.
Dependents	Character	A yes or no flag to check if the customer has dependants.
tenure	Integer	The customer's time in the company.
PhoneService	Character	A flag yes or no if the customer is opted for phone service
MultipleLines	Character	A yes or no flag to check if customer has multiple lines on his name.
InternetService	Character	A yes or no flag to check if customer has internet service from the company.
OnlineSecurity	Character	A yes or no flag to check if customer has opted for online security service.
OnlineBackup	Character	A yes or no flag to check if customer has opted for online backup.
DeviceProtection	Character	A yes or no flag to check if customer has opted for device protection service.
TechSupport	Character	A yes or no flag to check if customer has opted for tech support.
StreamingTV	Character	A yes or no flag to check if customer has opted for streaming tv service.
StreamingMovies	Character	A yes or no flag to check if customer has opted for streaming tv service.

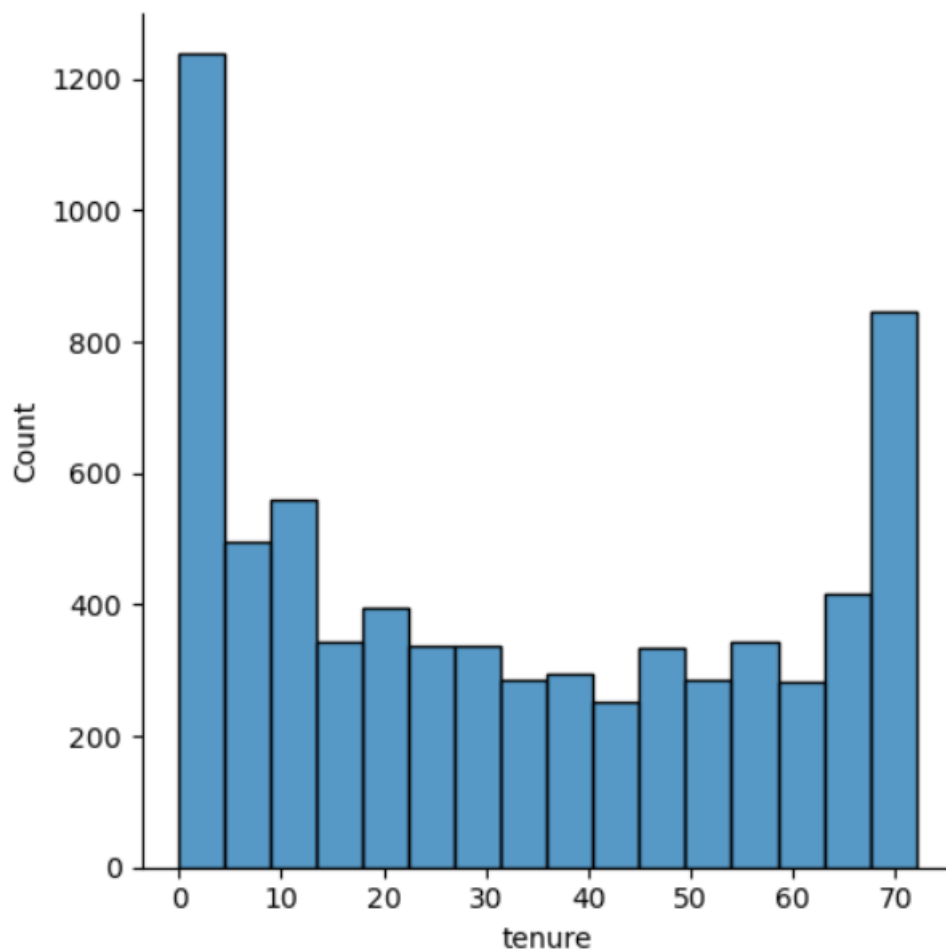
Contract	Character	A variable to check the customer's monthly or yearly contract.
PaperlessBilling	Character	A yes or no flag to check if customer has opted for paperless billing.
PaymentMethod	Character	Payment method for the bill payment
MonthlyCharges	Float	Average charge for the month
TotalCharges	Float	Total charges for month for monthly contract, total charges for year for yearly contract
CustFinancialStatus (<i>Derived column</i>)	Character	A variable that designates the user financial status based on the customer spending
Churn	Integer	A 0,1 flag to check if the customer has churned or not.

In this dataset, nulls were found only in TotalCharges and they were imputed with 0.

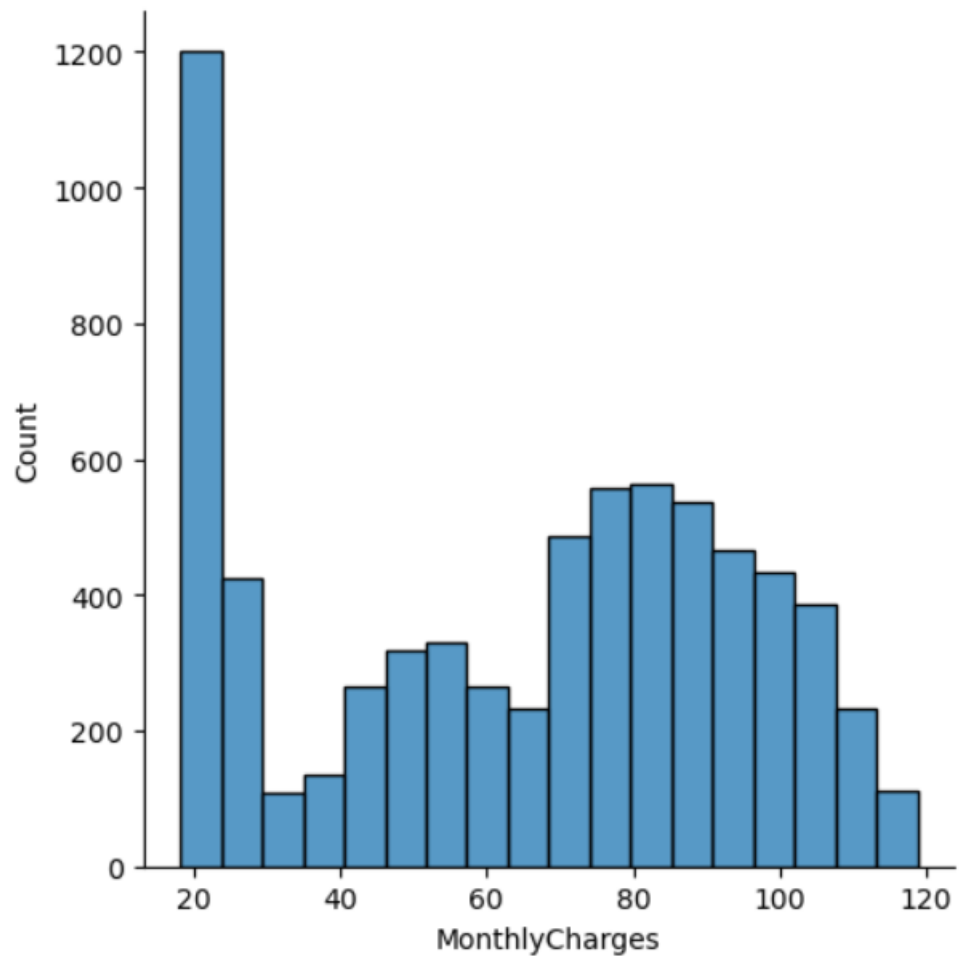
Correlation Matrix between numerical features:



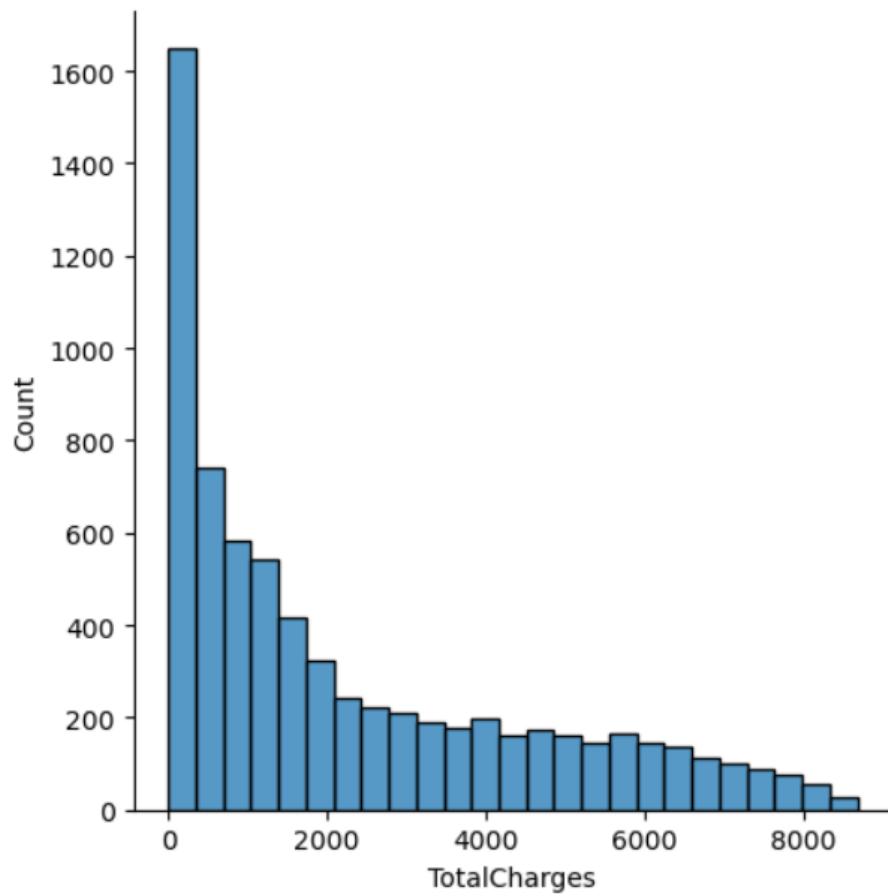
Univariate Analysis



From the above plot we can notice that the plot is right skewed and does not follow a proper bell curve to be classified as normal distribution. And from the plot we notice that most observations are at 0-10 and 60-70 months of tenure.

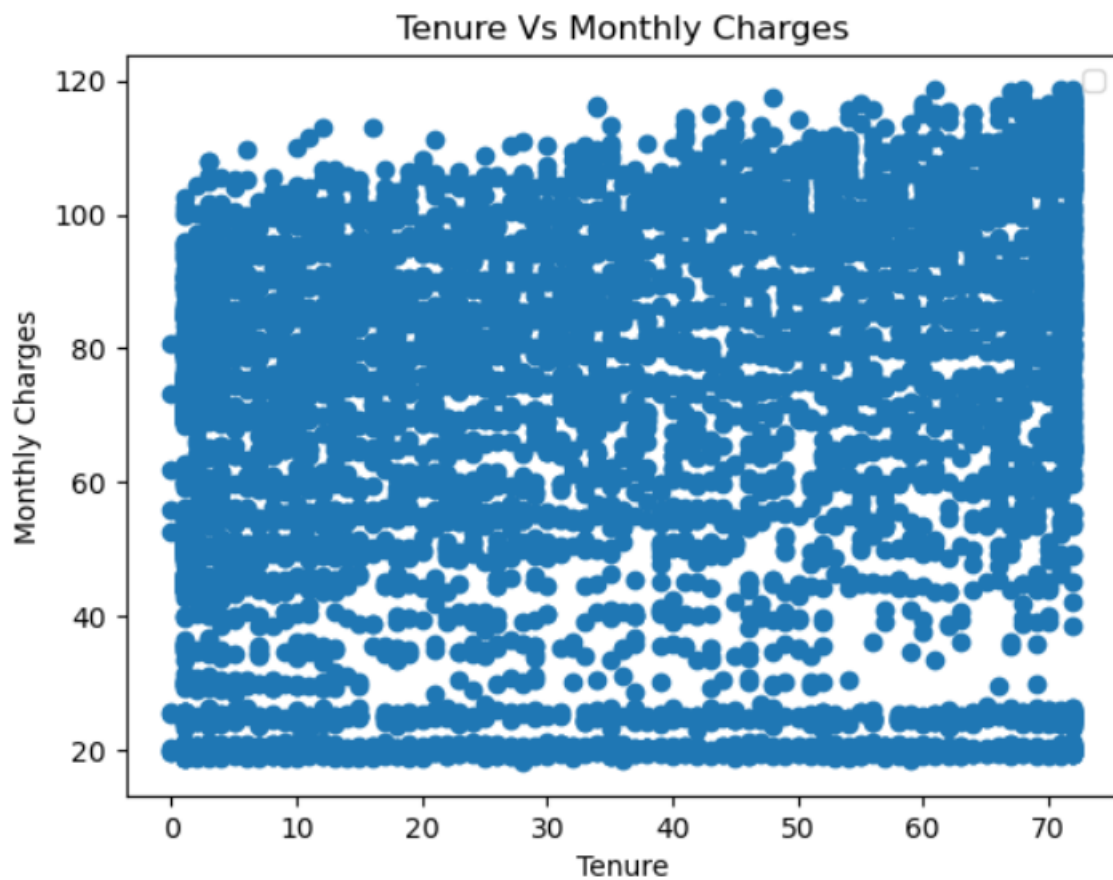


From the above plot we can notice that the plot is right skewed and does not follow a proper bell curve to be classified as normal distribution. And most observations are between the monthly charges are between 20 to 40 \$ per month.

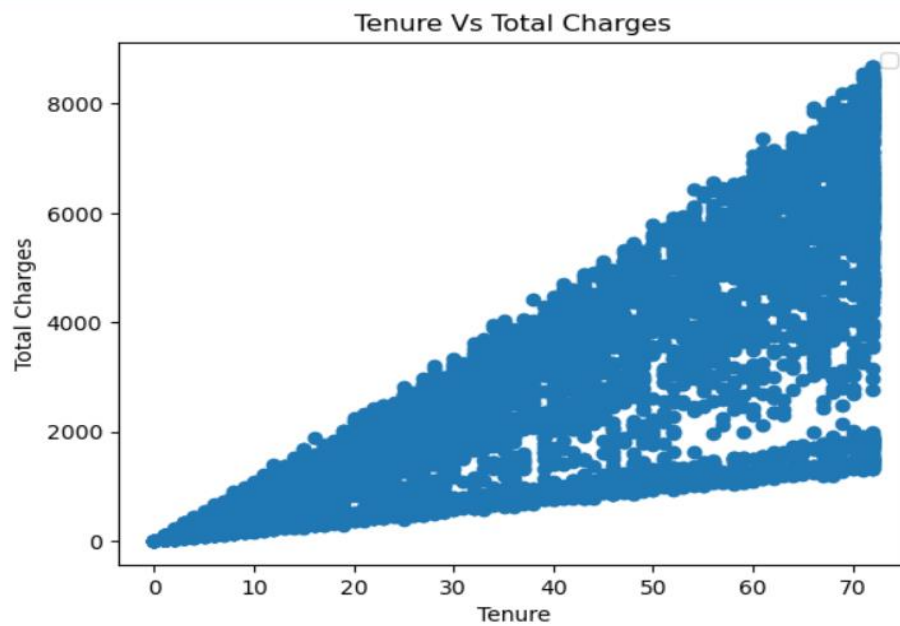


From the above plot we can notice that the plot is right skewed and does not follow a proper bell curve to be classified as normal distribution. And from the plot we can see that most people charges are between 0 to 2000.

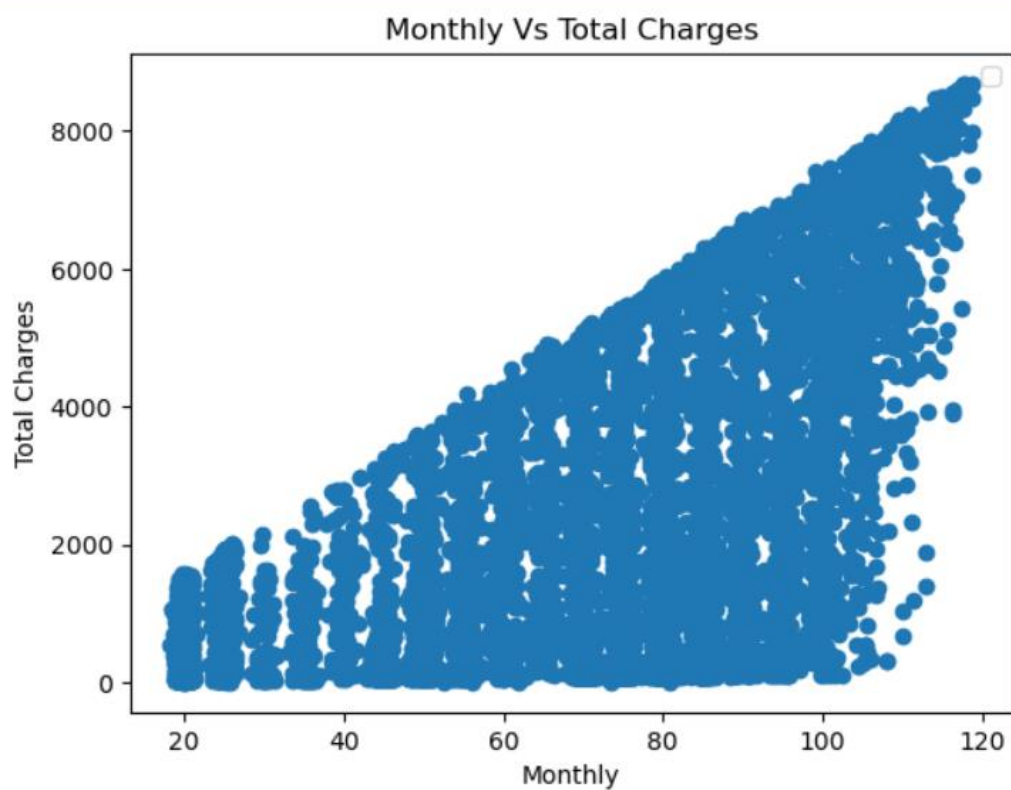
Bivariate analysis:



From the above plot with tenure against monthly charges we can see that there is no correlation between the points and no there seems to be no relation ship between both the variables.



From the above figure we can notice that there is a positive correlation between total charges and tenure. We can see that with increase in tenure there is a positive increase in the total charges. We can also see this in the correlation matrix later.

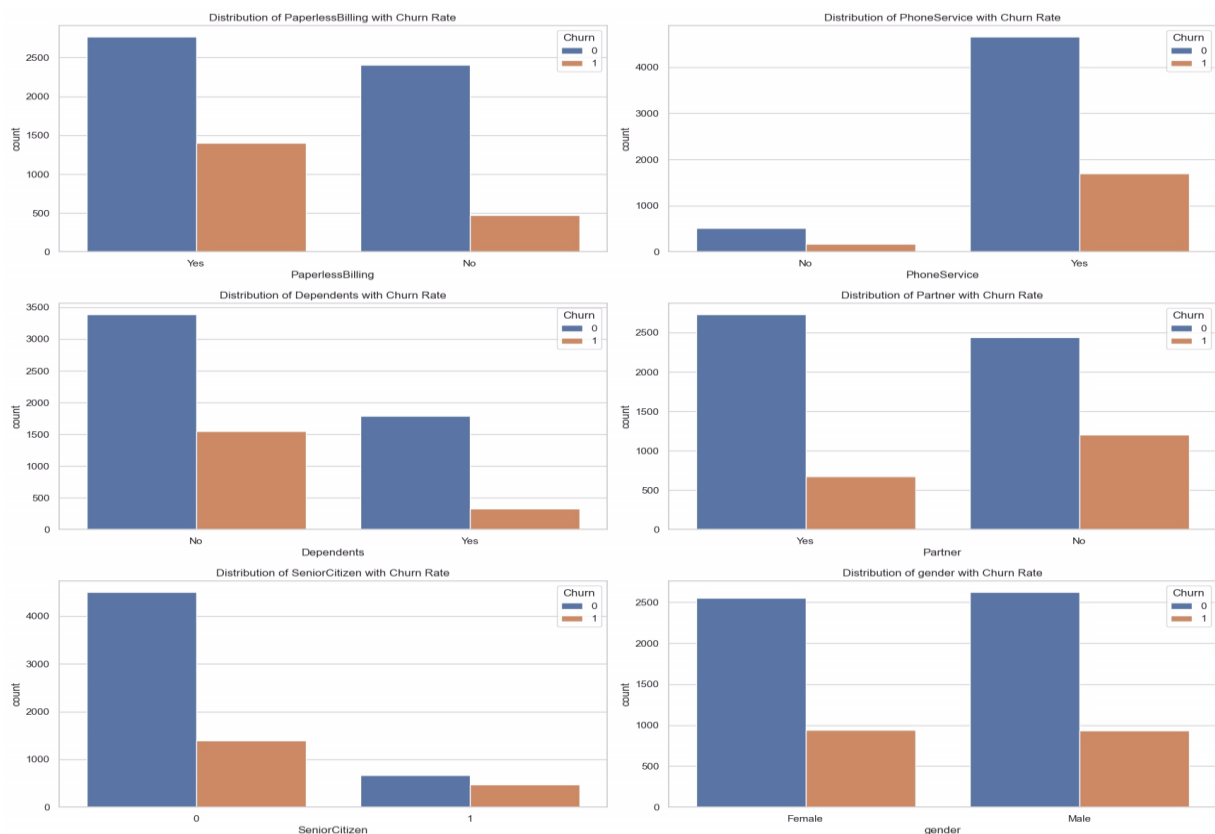


From the above figure we can notice that there is a positive correlation between total charges and Monthly charges. We can see that with increase in Monthly charges there is a positive increase in the total charges. We can also see this in the correlation matrix later.

	tenure	MonthlyCharges	TotalCharges
count	7043.000000	7043.000000	7043.000000
mean	32.371149	64.761692	2279.734304
std	24.559481	30.090047	2266.794470
min	0.000000	18.250000	0.000000
25%	9.000000	35.500000	398.550000
50%	29.000000	70.350000	1394.550000
75%	55.000000	89.850000	3786.600000
max	72.000000	118.750000	8684.800000

From the above summary of numerical columns, we can see that there are no null values so going down we can discard the notion of creating a strategy for treating the null values.

Observing Churn along PaperlessBilling (Y/N), PhoneService (Y/N), Dependents (Y/N), Partner (Y/N), SeniorCitizen (Y/N), gender (M/F):



Feature Engineering:

CustFinancialStatus:

A new variable called CustFinancialStatus was created. It classified customers into 3 categories ("Lower Income", "Middle Class", "Upper Class") based on the services they use and their payment method.

The rules prioritize the `PaymentMethod` column as follows:

1. If `PaymentMethod` is "Electronic check", then classify as "Middle Class".
2. If `PaymentMethod` is "Mailed check", then classify as "Lower Income".
3. If `PaymentMethod` is "Bank transfer (automatic)" or "Credit card (automatic)", then classify as "Upper Class".
4. If `Contract` is "Two year", classify as "Upper Class" unless already classified by `PaymentMethod`.
5. If a customer has multiple premium services and is not already classified as "Upper Class" by `PaymentMethod`, classify as "Middle Class".
6. If a customer has minimal services, classify as "Lower Income".

Model Building:

Step 1 – Scaling the data with StandardScaler():

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
data[num_cols] = scaler.fit_transform(data[num_cols])
cat_columns = list(data.select_dtypes(include=['object']).columns)
cat_columns.remove('customerID')
data = pd.get_dummies(data, columns=cat_columns, drop_first=True)
data.drop("customerID",axis=1,inplace=True)
X = data.drop("Churn",axis = 1)
y = data["Churn"]
```

In this code snippet, data preprocessing steps are implemented for a machine learning model. Firstly, StandardScaler from Scikit-learn is used to standardize numerical columns (num_cols) in the dataset, ensuring that each feature contributes equally to the model. Then, categorical columns are identified, and dummy variables are created using pd.get_dummies, which helps in converting categorical data into a numerical format suitable for machine learning algorithms. The 'customerID' column is removed as it's likely a unique identifier and not useful for prediction. Finally, the dataset is split into features (X) and the target variable (y), which is 'Churn' in this case, preparing the data for model training.

Step 2 - Addressing Class Imbalance with ADASYN:

```
from imblearn.over_sampling import ADASYN
X_resampled, y_resampled = ADASYN().fit_resample(X, y)
X_resampled
```

In this segment, the code addresses the issue of class imbalance in the dataset using the ADASYN (Adaptive Synthetic Sampling) method from the imblearn library. ADASYN generates synthetic samples for the minority class, making the class distribution more balanced. The fit_resample method is applied to the feature set X and target variable y, resulting in a new, resampled dataset (X_resampled, y_resampled) where the class distribution is more uniform. This preprocessing step is crucial for improving the performance of machine learning models on imbalanced datasets.

Step 3 - Splitting dataset into Training and Testing dataset:

```
X_train, X_test, y_train, y_test = train_test_split(X_resampled, y_resampled, test_size=0.2,
random_state=42)
```

This code line is utilizing the `train_test_split` function from Scikit-learn to divide the resampled dataset (`X_resampled`, `y_resampled`) into training and testing sets. It allocates 20% of the data to the test set (`test_size=0.2`) and the remaining 80% to the training set. The `random_state` parameter set to 42 ensures reproducibility of the split. This process is crucial for evaluating the machine learning model's performance on unseen data, maintaining the integrity and generalizability of the model.

Step 4 – (i) a) Applying Logistic Regression:

```
from sklearn.linear_model import LogisticRegression
# Initialize the possible hyperparameters for the LogisticRegression Model
param_grid_lr = {
    'C': [0.1, 1, 10],
    'solver': ['saga']
}

# Initialize grid search to find best possible hyperparameters for the model
grid_search_lr = GridSearchCV(
    estimator=LogisticRegression(max_iter=100, solver='saga'), # Reduced iterations for grid search
    param_grid=param_grid_lr,
    cv=3,
    verbose=1,
    n_jobs=-1
)

# Fit the model with the training data
grid_search_lr.fit(X_train, y_train)
```

b) Explanation:

In this segment, a Logistic Regression model is instantiated with specific configuration parameters, notably the number of iterations and the optimization solver. Logistic Regression, a fundamental classification technique in machine learning, is adept at binary classification tasks. It models the probability of a binary response based on one or more predictor variables. The model is then trained, and its predictive performance is assessed using standard evaluation metrics, which are crucial in gauging its effectiveness in accurately classifying the data into relevant categories.

c) Analysis of Performance metrics of the Logistic Regression model:

```
Best Parameters: {'C': 10, 'solver': 'saga'}
Best Score: 0.7845614624901888
```

	precision	recall	f1-score	support
0	0.80	0.75	0.77	1036
1	0.76	0.81	0.78	1024

accuracy			0.78	2060
macro avg	0.78	0.78	0.78	2060
weighted avg	0.78	0.78	0.78	2060

```
[[774 262]
 [194 830]]
```

F1-Score:

Class 0: The F1-score, a harmonic mean of precision and recall, was 77% for non-churn predictions, indicating a balanced performance.

Class 1: For churn predictions, the F1-score was marginally higher at 78%, suggesting a slightly better balance in predicting churn cases.

Accuracy: The overall accuracy of the model was 77%, denoting a robust performance across both classes.

Confusion Matrix Analysis:

True Negatives (TN): 768 customers were correctly identified as not at risk of churning.

False Positives (FP): 268 customers were mistakenly predicted as likely to churn.

False Negatives (FN): 200 customers who were churning weren't flagged by the model.

True Positives (TP): 824 churning customers were accurately identified.

Step 4 – (ii) a) Applying KNN-Classifier:

```
from sklearn.neighbors import KNeighborsClassifier
```

```
# Initialize possible hyperparameters
```

```
param_grid_knn = {
    'n_neighbors': [3, 5, 7, 9, 11, 13, 15],
    'weights': ['uniform', 'distance'],
    'metric': ['euclidean', 'manhattan', 'minkowski']
}
```

```
# Initialize grid search to get best hyperparameters
```

```
grid_search_knn = GridSearchCV(
    estimator=KNeighborsClassifier(),
    param_grid=param_grid_knn,
    cv=5,
    verbose=1,
    n_jobs=-1 # Use all processors
)
```

```
# Fit the model with the training data
```

```
grid_search_knn.fit(X_train_np, y_train_np)
```

b) Explanation:

In this section, the K-Nearest Neighbors (KNN) algorithm, a simple yet effective machine learning technique, is employed for classification. The code begins by ensuring that the data,

originally in Pandas DataFrame or Series format, is converted into NumPy arrays. This conversion is crucial for optimizing computational efficiency, as KNN involves intensive distance calculations that are more efficiently executed with NumPy's optimized structures. The KNN model, set to consider the 5 nearest neighbors, is trained on this transformed data.

c) Analysis of Performance metrics of the KNN-Classifer model:

Fitting 5 folds for each of 42 candidates, totalling 210 fits
 Best Parameters: {'metric': 'manhattan', 'n_neighbors': 3, 'weights': 'distance'}

Best Score: 0.791237598811608

	precision	recall	f1-score	support
0	0.89	0.68	0.77	1036
1	0.74	0.92	0.82	1024
accuracy			0.80	2060
macro avg	0.82	0.80	0.80	2060
weighted avg	0.82	0.80	0.80	2060

```
[[708 328]
 [ 86 938]]
```

F1-Score:

Class 0: The F1-score is 75%, reflecting a balance between precision and recall for non-churn predictions.

Class 1: An F1-score of 81% shows a strong balance between precision and recall in churn predictions.

Confusion Matrix Analysis:

True Negative (TN): 663 instances where the model correctly identified non-churning customers.

False Positive (FP): 373 instances where the model incorrectly labeled non-churning customers as churning. False Negative

(FN): 65 instances where the model failed to identify churning customers.

True Positive (TP): 959 instances where the model correctly identified churning customers.

Step 4 – (iii) a) Applying Decision Tree Classifier:

```
from sklearn.tree import DecisionTreeClassifier
```

```
# Initialize the Hyperparameters for the Decision Tree
```

```
param_grid_dt = {
    'max_depth': [None, 10, 20, 30, 40, 50],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4],
    'criterion': ['gini', 'entropy']
}
```

```
# Initialize the grid-search to find best hyperparameters
```

```
grid_search_dt = GridSearchCV(
    estimator=DecisionTreeClassifier(), # Initialize DecisionTreeClassifier()
```

```

param_grid=param_grid_dt,
cv=5,
verbose=1,
n_jobs=-1 # Use all processors
)

# Fit the training data
grid_search_dt.fit(X_train, y_train)

```

b) Explanation:

The Decision Tree Classifier operates by splitting the dataset into smaller subsets based on feature value conditions. It starts at the tree's root and splits the data on the feature that results in the largest information gain, repeating this process at each node. The tree's branches represent these decisions, leading to leaf nodes that signify the outcome or class prediction. This model is favored for its interpretability, as the path from the root to a leaf represents a decision-making process, making it relatively straightforward to understand and visualize.

c) Analysis of Performance metrics of the Decision Tree classifier model:

```

Fitting 5 folds for each of 108 candidates, totalling 540 fits
Best Parameters: {'criterion': 'entropy', 'max_depth': 20, 'min_samples_lea
f': 1, 'min_samples_split': 5}
Best Score: 0.7655062013310461

```

	precision	recall	f1-score	support
0	0.76	0.74	0.75	1036
1	0.74	0.76	0.75	1024
accuracy			0.75	2060
macro avg	0.75	0.75	0.75	2060
weighted avg	0.75	0.75	0.75	2060

```

[[764 272]
 [242 782]]

```

Confusion Matrix Analysis:

F1-Score:

Class 0: The F1-score for non-churn predictions was 76%, reflecting a balance between precision and recall.

Class 1: A similar F1-score of 76% for churn predictions indicates an effective balance in this category as well.

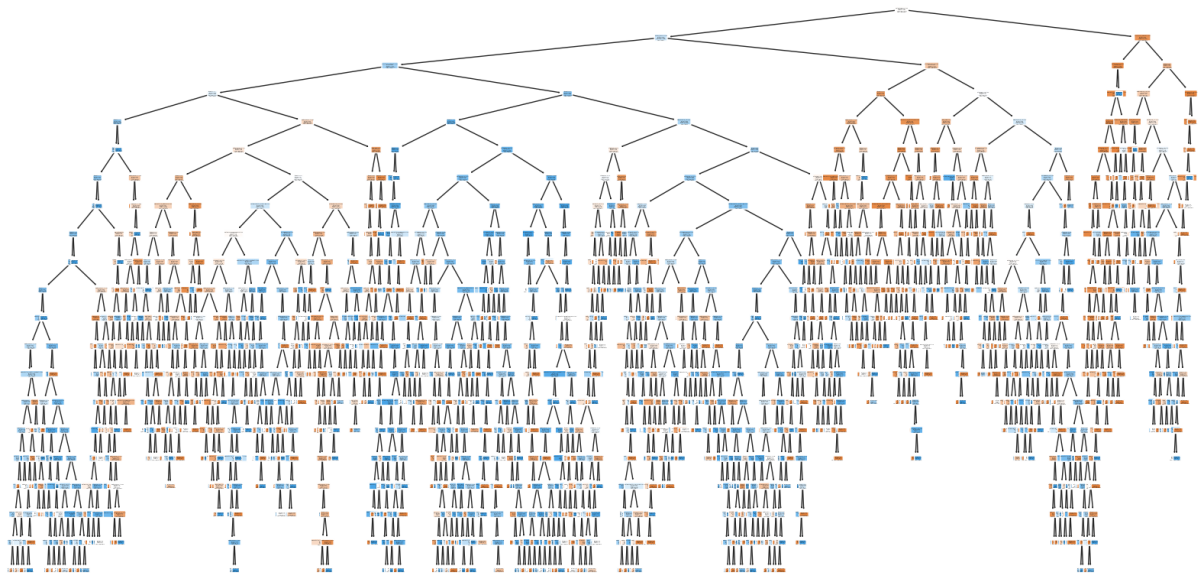
True Negatives (TN): 777 customers were correctly identified as not churning.

False Positives (FP): 259 customers were incorrectly labelled as at risk of churning.

False Negatives (FN): 235 customers who were churning were not detected by the model.

True Positives (TP): 789 churning customers were accurately identified.

d) Decision Tree visualization:



Step 4 – (iv) a) Applying Random Forest Classifier:

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, confusion_matrix

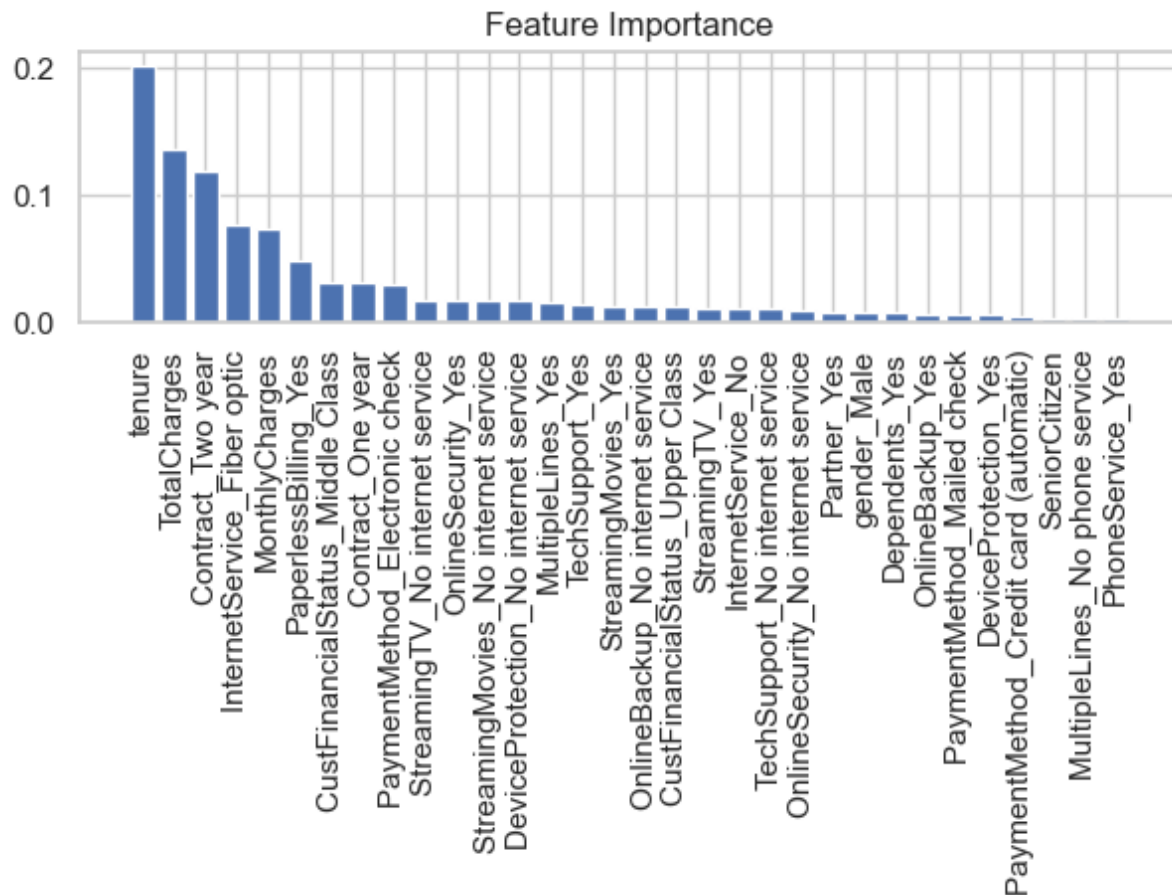
# Create hyperparameters and train the Random Forest Classifier
param_grid_rf = {
    'n_estimators': [100, 200],
    'max_features': ['sqrt', 'log2'],
    'max_depth': [4, 6, 8],
    'min_samples_split': [2, 10]
}

# Initialize Grid search CV for the Random Forest model
grid_search_rf = GridSearchCV(
    estimator=RandomForestClassifier(),
    param_grid=param_grid_rf,
    cv=3,
    verbose=1,
    n_jobs=-1
)
grid_search_rf.fit(X_train, y_train)
```

b) Explanation:

The Random Forest model works by creating a 'forest' of decision trees during training, each constructed using a random subset of features and training data points. When making predictions, each tree in the forest votes, and the most common outcome (for classification tasks) is taken as the final prediction. This ensemble approach makes the Random Forest model robust against overfitting and effective in handling classification tasks.

c) Feature Importance of the Random Forest model:



d) Analysis of Performance metrics of the Random Forest model:

Fitting 3 folds for each of 24 candidates, totalling 72 fits

Best Parameters: {'max_depth': 8, 'max_features': 'sqrt', 'min_samples_split': 2, 'n_estimators': 200}

Best Score: 0.772666431755033

	precision	recall	f1-score	support
0	0.80	0.70	0.74	1036
1	0.73	0.82	0.77	1024
accuracy			0.76	2060
macro avg	0.76	0.76	0.76	2060
weighted avg	0.76	0.76	0.76	2060

```
[[722 314]
 [183 841]]
```

Confusion Matrix Analysis:

F1-Score:

Class 0: An F1-score of 82% for non-churn predictions indicates a good balance between precision and recall.

Class 1: The F1-score for churn predictions was slightly higher at 84%, indicating a very effective balance for this class.

Accuracy: The overall accuracy of the model was 83%, signifying its strong performance across both churn and non-churn predictions.

Confusion Matrix Analysis:

True Negatives (TN): 807 customers were accurately identified as not at risk of churning.

False Positives (FP): 229 customers were mistakenly labelled as likely to churn.

False Negatives (FN): 118 customers who were churning were not flagged by the model.

True Positives (TP): 906 churning customers were correctly identified.

Step 4 – (v) a) Applying Naïve Bias Classifier:

```
from sklearn.naive_bayes import GaussianNB

# Parameter grid for var_smoothing
param_grid_nb = {
    'var_smoothing': np.logspace(0, -9, num=100)
}

# Setting up the GridSearchCV for Naive Bayes
grid_search_nb = GridSearchCV(estimator=GaussianNB(), param_grid=param_grid_nb,
                               cv=5, verbose=1, scoring='accuracy')

# Fitting GridSearchCV to the training data
grid_search_nb.fit(X_train, y_train)
```

b) Explanation:

The Naive Bayes model is a probabilistic classifier that applies Bayes' theorem with the assumption of independence between features. It calculates the probability of each class based on the input features and classifies a new instance by choosing the class with the highest posterior probability.

c) Analysis of Performance metrics of the Random Forest model:

Fitting 5 folds for each of 100 candidates, totalling 500 fits

Best Parameters: {'var_smoothing': 0.02310129700083159}

Best Score: 0.6860057415365389

	precision	recall	f1-score	support
0	0.78	0.51	0.62	1036
1	0.63	0.86	0.73	1024
accuracy			0.68	2060
macro avg	0.71	0.68	0.67	2060
weighted avg	0.71	0.68	0.67	2060

```
[[528 508]
 [147 877]]
```

Confusion Matrix Analysis:

F1-Score:

Class 0: The F1-score for non-churn predictions is 62%, reflecting a balance between precision and recall.

Class 1: For churn predictions, the F1-score is higher at 72%, showing a better balance in this category.

Accuracy: The overall accuracy of the model stands at 68%, which is moderate and suggests room for improvement.

Confusion Matrix Insights:

True Negatives (TN): 529 customers were correctly identified as non-churning.

False Positives (FP): 507 customers were incorrectly predicted as churning.

False Negatives (FN): 154 customers who were churning were not identified by the model.

True Positives (TP): 870 churning customers were correctly identified.

Step 4 – (vi) a) Applying XG-Boost Classifier:

```
import xgboost as xgb

# Create hyperparameters and train the XGBoost model
parameters = {
    'n_estimators': [100, 200, 300],
    'learning_rate': [0.01, 0.05, 0.1],
    'max_depth': [3, 4, 5],
    # Add more parameters here
}
grid_search = GridSearchCV(xgb.XGBClassifier(use_label_encoder=False,
                                             eval_metric='logloss'),
                           parameters, cv=3, scoring='accuracy')

grid_search.fit(X_train, y_train)

# Find best parameters
best_parameters = grid_search.best_params_
print("Best Parameters:", best_parameters)

# Create and train the XGBoost model with optimized hyperparameters
xgb_classifier = xgb.XGBClassifier(
    use_label_encoder=False,
    eval_metric='logloss',
    learning_rate=0.1, # Optimized learning rate
    max_depth=5,      # Optimized max depth
    n_estimators=300   # Optimized number of estimators
)

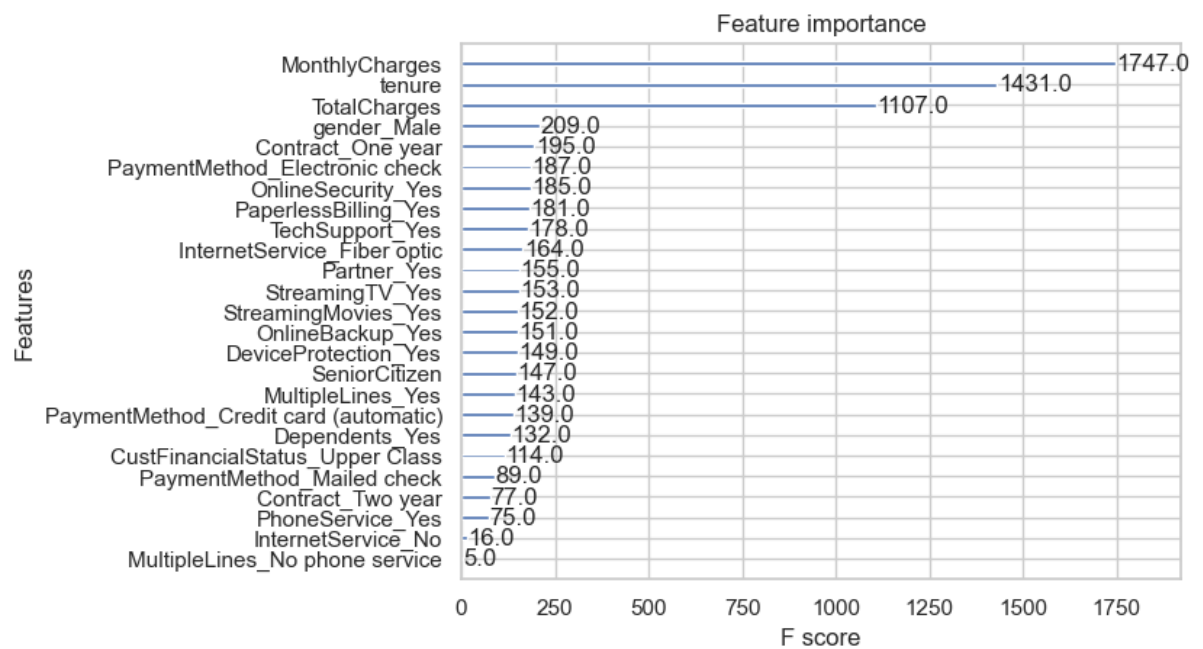
# Fit the model to the training data
xgb_classifier.fit(X_train, y_train)
```

b) Explanation:

The XG-Boost (eXtreme Gradient Boosting) model is an advanced implementation of gradient boosting algorithms, known for its efficiency and performance. It builds a series of decision trees sequentially, where each new tree attempts to correct the errors made by the

previous ones. XG-Boost incorporates techniques like regularization and tree pruning to prevent overfitting, and it efficiently handles large datasets and complex classification tasks, often delivering superior predictive performance.

c) Feature Importance of the XG-Boost model:



d) Analysis of Performance metrics of the Random Forest model:

Best Parameters: {'learning_rate': 0.1, 'max_depth': 5, 'n_estimators': 300 }

	precision	recall	f1-score	support
0	0.85	0.78	0.81	1036
1	0.79	0.86	0.82	1024
accuracy			0.82	2060
macro avg	0.82	0.82	0.82	2060
weighted avg	0.82	0.82	0.82	2060

[[806 230]
[146 878]]

Confusion Matrix Analysis:

Accuracy: The overall accuracy of the model is 81%, indicating a high level of correctness across both classes.

F1-Score:

Class 0: The F1-score is 80%, which is a balance between precision and recall for predictions of non-churn.

Class 1: Similarly, the F1-score for churn predictions is 82% which shows that the model is able to predict the Churning customers efficiently.

Confusion Matrix:

The model predicted 790 true negatives (correct non-churn predictions) and 880 true positives (correct churn predictions).

However, there were 246 false positives (non-churn customers incorrectly labeled as churn) and 144 false negatives (churn customers incorrectly labeled as non-churn).

Results:

Comparing the Accuracy and F1-scores of every model:

```
# Parsing the accuracies
accuracies = []
for report in classification_reports:
    lines = report.split('\n')
    for line in lines:
        if 'accuracy' in line:
            accuracy = float(line.strip().split()[1])
            accuracies.append(accuracy)
            break

# For adding model names to visualization in X-axis
model_names = ['Logistic Regression', 'KNN Classifier', 'Decision Tree',
               'Random Forest', 'Naive Bias', 'XGBoost']

# Plot for accuracies across all models
plt.figure(figsize=(10, 6))
plt.bar(model_names, accuracies, color='green')
plt.xlabel('Models')
plt.ylabel('Accuracies')
plt.title('Accuracies for Class 1 Across Models after Hyperparameter
tuning')
plt.ylim(0, 1)
plt.show()
```

Conclusion:

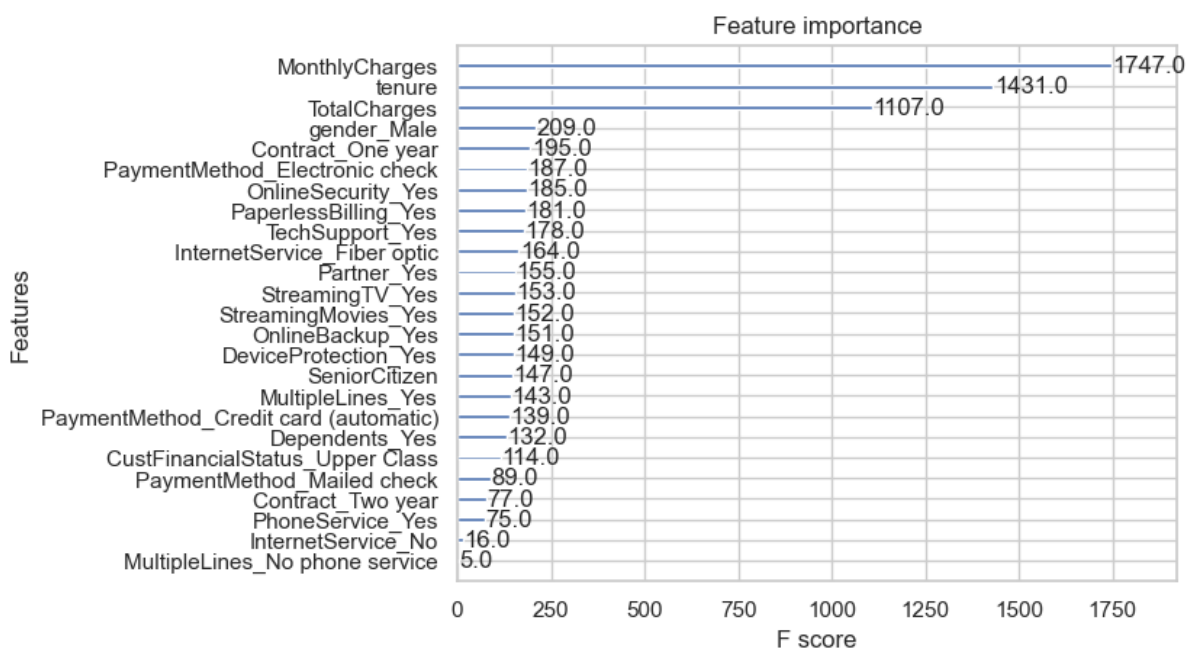
This study rigorously evaluated multiple machine learning models to predict customer churn for a telecommunications company. Among the models tested, including Logistic Regression, KNN-Classifier, Decision Tree, Random Forest, Naïve Bias, and XG-Boost, the XG-Boost model emerged as the most effective. It demonstrated a commendable balance between precision and recall, with an F1-Score of 82% for predicting churning customers (Class 1) and 81% for non-churning customers (Class 0).

Key Insights from our research:

Precision and Recall: The XG-Boost model achieved 78% precision for churning customers, indicating its effectiveness in correctly identifying those at risk of churn. For non-churning customers, the precision was even higher at 85%. In terms of recall, the model successfully identified 86% of the churning customers and 76% of the non-churning customers.

Overall Accuracy: With an overall accuracy of 81%, the model shows a high level of correctness across both classes.

Most Influential Variables: Analysis of feature importance revealed that TotalCharges, tenure, and MonthlyCharges are crucial factors influencing customer churn. Their impact is significantly higher compared to other variables in the dataset.



Business Implications and Recommendations:

Targeted Customer Retention Strategies: Given the high influence of TotalCharges, tenure, and MonthlyCharges on churn likelihood, targeted retention strategies can be developed. For instance, tailored communication and special offers can be directed at customers with high TotalCharges and long tenure, as they might be more prone to churn.

Customized Pricing Plans: Understanding the significant impact of MonthlyCharges on churn, the company might consider revising its pricing plans. This could involve offering more competitive or personalized pricing options to retain customers.

Proactive Engagement: The high recall rate for churning customers suggests that the model can effectively identify at-risk customers. The company can use this information for proactive engagement, such as reaching out with customer satisfaction surveys or special offers before these customers decide to leave.

Continuous Model Monitoring and Improvement: While the XG-Boost model shows promise, continuous monitoring and refinement are crucial. Regularly updating the model with new data will ensure its ongoing effectiveness.

Further Research and Development: The study also opens avenues for further research. Exploring additional variables or employing more advanced machine learning techniques could yield even better predictive performance.

In conclusion, the application of the XG-Boost model in predicting customer churn presents a valuable opportunity for the telecommunications company to enhance its customer retention strategies. By leveraging the insights gained from the model, the company can take a more data-driven approach to address churn, ultimately leading to improved customer satisfaction and loyalty.

References

- Bingquan Huang, Mohand Tahar Kechadi, Brian Buckley (2012) Customer churn prediction in telecommunications,
<https://www.sciencedirect.com/science/article/pii/S0957417411011353>
- Meng Hui (2012) The Research on Development Trend of Telecommunications Industry,
<https://www.sciencedirect.com/science/article/pii/S221266781200247X>
- S. A. Qureshi, A. S. Rehman, A. M. Qamar, A. Kamal and A. Rehman, "Telecommunication subscribers' churn prediction model using machine learning," Eighth International Conference on Digital Information Management (ICDIM 2013), Islamabad, Pakistan, 2013, pp. 131-136, doi: 10.1109/ICDIM.2013.6693977.
- Wei CP, Chiu IT. Turning telecommunications call details to churn prediction: a data mining approach. *Expert Syst Appl.* 2002;23(2):103–12.
- Umayaparvathi V, Iyakutti K. A survey on customer churn prediction in telecom industry: datasets, methods and metric. *Int Res J Eng Technol.* 2016;3(4):1065–70.
- Kiss C, Bichler M. Identification of influencers—measuring influence in customer networks. *Decis Support Syst.* 2008;46(1):233–53. <https://doi.org/10.1016/j.dss.2008.06.007>.
- Chen T, Guestrin C. Xgboost. A scalable tree boosting system. *CoRR.* 2016. [arXiv:1603.02754](https://arxiv.org/abs/1603.02754)