

Machine Learning Internship **Assessment**

Customer Churn Prediction **By Varun Dhanashree**

Problem Statement

The telecom industry is highly competitive, with a constant influx of new players and evolving technologies. Maintaining customer loyalty and reducing churn is of paramount importance for telecom companies. Customer churn, or the rate at which customers discontinue their subscriptions, has a direct impact on revenue and profitability. The ability to predict and proactively address churn is essential for long-term sustainability and growth.

The primary objective of this report is to develop an accurate and effective predictive model to identify customers at risk of churning. By analyzing historical customer data, we aim to create a robust machine learning model that can foresee potential churn and help the telecom company take preventative actions.

Data Description

Dataset consists customer information for a customer churn prediction problem. It includes the following columns:

CustomerID: Unique identifier for each customer.

Name: Name of the customer.

Age: Age of the customer.

Gender: Gender of the customer (Male or Female).

Location: Location where the customer is based, with options including Houston, Los Angeles, Miami, Chicago, and New York.

Subscription_Length_Months: The number of months the customer has been subscribed.

Monthly_Bill: Monthly bill amount for the customer.

Total_Usage_GB: Total usage in gigabytes.

Churn: A binary indicator (1 or 0) representing whether the customer has churned (1) or not (0).

Data Pre-Processing and Exploratory Data Analysis (EDA)

The initial step involved exploring the dataset to understand its structure and characteristics.

- * The dataset contains information about 100,000 customers with 9 variables.
- * All variables have the correct data type, and there are no missing values or duplicate records.
- * Descriptive statistics were generated for each variable, revealing insights into customer demographics, subscription details, billing, usage, and churn behavior.
- * Gender and Location distributions were analysed, indicating the gender and location distribution of the customers.

Checking for Imbalance

We observe that Data is Balanced using a pie chart

Visualizing Distributions of data

Using Plots we see that All Features are well distributed

Outlier Detection Using Boxplots

The Data doesn't contain any Outliers

Feature Encoding

Categorical variables were encoded to numerical values to enable machine learning algorithms to process them effectively.

* One-Hot Encoding was applied to the 'Gender' and 'Location' variables.

*

Data Splitting

The dataset was divided into training and testing sets to enable model training and evaluation.

* Dataset is divided into 75:25 ratio.

Feature Scaling

Feature scaling was applied to ensure all variables were on the same scale, aiding model convergence.

* Standardized Scaling was applied to variables such as 'Age', 'Subscription_Length_Months', 'Monthly_Bill', and 'Total_Usage_GB'.

Feature Selection Using Random Forest Feature Importance

Identifying important features helps streamline the model and improve its interpretability.

* Random Forest Feature Importance was used to rank features based on their contribution to the target variable.

* The top features were 'Monthly_Bill', 'Total_Usage_GB', 'Age', and 'Subscription_Length_Months'.

<u>Feature</u>	<u>Importance</u>
Monthly_Bill	0.316383
Total_Usage_GB	0.290353
Age	0.194396
Subscription_Length_Months	0.142624
Gender_Male	0.016683
Location_Los Angeles	0.010595
Location_Houston	0.010007
Location_Miami	0.009792
Location_New York	0.009166

Model Building: Machine Learning Algorithms

Several machine learning algorithms were trained and evaluated using the dataset.

* Algorithms included Logistic Regression, Decision Tree, K-Nearest Neighbours, Gaussian Naive Bayes, AdaBoost, Gradient Boosting, Random Forest, XGBoost, and Support Vector Classifier (SVC).

* Training and test data performance metrics were calculated, revealing the strengths and weaknesses of each algorithm.

Model Building: Neural Network

An attempt was made to build a neural network model, but it did not yield satisfactory results.

Model Building: Ensembles

Ensemble models were evaluated, but no significant improvement was observed.

Model Building: Final Model Selection –

KNN Classifier and XGBoost

XGBoost Classifier was identified as the best-performing algorithm across various metrics and feature variations.

Hyperparameter Tuning

Hyperparameter tuning was explored to improve the model's performance, but no substantial gains were achieved.

We choose KNN Model as KNN works better (recall as a metric) on train compared to XGBoost

Cross-Validation

Cross-validation was performed to validate the model's performance and ensure it generalized well to new data.

(I) **Cross-Validation Scores (Accuracy)**: [0.49692857, 0.50057143, 0.49892857, 0.50478571, 0.505].

Mean Accuracy Score: 0.50

(II) **Cross-Validation Scores (Recall)**: [0.48990983, 0.49398798, 0.48869167, 0.50171772, 0.49427426].

Mean Recall Score: 0.49

Model Evaluation

(I) Train & Test Data Metrics

The final KNN model's performance was evaluated using various metrics on both the training and test datasets.

Train Metrics:

Accuracy: 0.69

Confusion Matrix:

[[26058 11585]

[11749 25608]]

Classification Report:

	precision	recall	f1-score	support
0	0.69	0.69	0.69	37643
1	0.69	0.69	0.69	37357
accuracy			0.69	75000
macro avg	0.69	0.69	0.69	75000
weighted avg	0.69	0.69	0.69	75000

Test Metrics:

Accuracy: 0.50

Confusion Matrix:

[[6271 6307]

[6290 6132]]

Classification Report:

	precision	recall	f1-score	support
0	0.50	0.50	0.50	12578
1	0.49	0.49	0.49	12422
accuracy			0.50	25000
macro avg	0.50	0.50	0.50	25000
weighted avg	0.50	0.50	0.50	25000

(II) ROC-AUC Curve

* Train ROC-AUC (area=0.75)

* Test ROC-AUC (area=0.49)

Saving Model

The final KNN model was saved as a pickle file for future use.

Model Deployment

Finally I've deployed model on the development environment IDE

Conclusion

The customer churn prediction project involved thorough exploratory data analysis, preprocessing, and the evaluation of various machine learning algorithms. The KNN Classifier was selected as the final model due to its superior performance across different metrics. While achieving optimal accuracy and recall is challenging, the insights gained from this project can guide the company's strategies for customer retention and business growth. Further analysis may involve gathering more data and exploring advanced techniques to improve model performance.

