

Project Report
On
Telecom Customer Churn Prediction



Submitted
In partial fulfilment
For the award of the Degree of

PG-Diploma in Big Data Analytics

(C-DAC, ACTS (Pune))

Guided By:

Mr. Saurabh Pawar

Submitted By:

Abhay Singh Pundir (240340125001)

Prerna (240340125033)

Abhijeet Ashokrao Panchal (240340125002)

Prince Saini (240340125034)

Deepanker Sharma (240340125020)

Centre for Development of Advanced Computing

(C-DAC), ACTS (Pune- 411008)

Acknowledgement

This is to acknowledge our indebtedness to our Project Guide, **Mr. Saurabh Pawar**, C-DAC ACTS, Pune for his constant guidance and helpful suggestion for preparing this project **Telecom Customer Churn Prediction**. We express our deep gratitude towards him for inspiration, personal involvement, constructive criticism that he provided us along with technical guidance during the course of this project.

We take this opportunity to thank Head of the department **Mr. Gaur Sunder** for providing us such a great infrastructure and environment for our overall development.

We express sincere thanks to **Mrs. Namrata Ailawar** for their kind cooperation and extendible support towards the completion of our project.

It is our great pleasure in expressing sincere and deep gratitude towards **Mrs. Risha P R (Program Head)** and **Ms. Pratiksha Gacche (Course Coordinator, PG-DBDA)** for their valuable guidance and constant support throughout this work and help to pursue additional studies.

Also, our warm thanks to **C-DAC ACTS Pune**, which provided us this opportunity to carry out, this prestigious Project and enhance our learning in various technical fields.

Abhay Singh Pundir (240340125001)

Prerna (240340125033)

Abhijeet Ashokrao Panchal (240340125002)

Prince Saini (240340125034)

Deepanker Sharma (240340125020)

Abstract

Customer churn is a major problem and one of the most important concerns for large companies. Due to the direct effect on the revenues of the companies, especially in the telecom field, companies are seeking to develop means to predict potential customer to churn. Therefore, finding factors that increase customer churn is important to take necessary actions to reduce this churn. The Objective is to develop a churn prediction model which assists companies to predict customers who are most likely subject to churn. We are processing ETL operation on top of that data using Apache Spark and dumping required data into Mongo Db for data visualization through Tableau and build Machine Learning for estimating the customer churn rate. Prediction of customer churns help in planning process and strategic decision making in large companies.

Table of Contents

Abstract

1. Introduction

2. System Requirements

2.1 Software Requirements

2.2 Hardware Requirements

3. Functional Requirements

4. System Architecture

5. Methodology

6. Data Visualization and Representation

7. Conclusion And Future Scope

References

Chapter 1

Introduction

1.1 Introduction

Churn is the process of customers switching from one firm to another in given time. Retaining the existing customers is more profitable than fetching the new customers. The Companies concentrate to the extant customers to avert churn. A churn prediction model is needed to predict the churners.

The fast growth of marketplace in every business is giving rise to increased subscriber base. Accordingly, companies have recognized the significance of retaining the customers who is on hand. It has become necessary for service-providers to reduce the churn rate of customers since the inattention might negatively influence profitability of the company.

Churn prediction contributes to identify those users who are likely to switch a company over another. The dataset used for customer churn is of telecom industry. It is collected from www.github.com that contains customer and service information for telecom industry.

Datasets and features:

Data used in the project is structured in nature. It was collected from www.github.com. The Objective is to develop a churn prediction model which assists companies to predict customers who are most likely subject to churn. Logistic Regression, Decision Tree, Random Forest, Stacking Classifier and Voting Classifier models were used to predict customer churn rate.

Chapter 2

System Requirements

➤ **Hardware Requirements:**

- ☐ Platform – Windows 10
- ☐ RAM – 8 GB of RAM
- ☐ Peripheral Devices – Mouse, Keyboard, Monitor
- ☐ A network connection for data recovering over network.

➤ **Software Requirements:**

- Python 3
- Apache Spark
- MongoDB
- Tableau
- OS – Window

Chapter 3

Functional Requirements

(1) Python 3:

- Python is a general purpose and high-level programming language.
- It is use for developing desktop GUI applications, websites and web applications.
- Python allows to focus on core functionality of the application by taking care of common programming tasks.

(2) Apache Spark:

- Apache Spark is an open-source cluster computing system that provides high-level API in Java, Scala, Python and R.
- Apache Spark is one of the fastest-growing big data projects in the history of the Apache Software Foundation. With its memory-oriented architecture, flexible processing libraries, and ease-of-use, Spark has emerged as a leading distributed computing framework for real time analytics.
- Spark is used for many types of data processing – it comes packaged with support for machine learning, interactive queries (SQL), statistical queries with R, graph processing, ETL, and streaming.
- For loading and storing data, Spark integrates with a number of storage MongoDB, and more.

(3) MongoDB:

- MongoDB, the most popular NoSQL database, is an open-source document-oriented database.
- MongoDB allows a highly flexible and scalable document structure.
- MongoDB has built in solution for partitioning and sharing your database.
- MongoDB provides a variety of storage engines, allowing you to choose one most suited to your application.
- A real-life scenario for this kind of data manipulation is storing and querying real-

time, intraday market data in MongoDB.

(4) Tableau:

- **Data Visualization:** - Tableau is a data visualization tool, and provides complex computation, data blending, and dashboarding for creating beautiful data visualizations.
- **Quickly Create Interactive Visualization:** - Users can create a very interactive visual by using drag n drop functionalities of Tableau.
- **Comfortable in Implementation:** - Many types of visualization options are available in Tableau, which enhances the user experience. Tableau is very easy to learn in comparison to Python. Who don't have any idea about coding, they also can quickly learn Tableau.
- **Tableau can Handle Large Amounts of Data:** - Tableau can easily handle millions of rows of data. A large amount of data can create different types of visualization without disturbing the performance of the dashboards. As well as, there is an option in Tableau where the user can make 'live' to connect different data sources like SQL, etc.
- It helps create interactive graphs and charts in the form of dashboards and worksheets to gain business insights.
- All of this is made possible with gestures as simple as drag and drop.

Data Cleaning Process:



Data cleansing or data cleaning is the process of detecting and correcting (or removing) corrupt or inaccurate records from a record set, table, or database and refers to identifying incomplete, incorrect, inaccurate or irrelevant parts of the data and then replacing, modifying, or deleting the dirty or coarse data. Data cleansing may be performed interactively with data wrangling tools, or as batch processing through scripting. After cleansing, a data set should be consistent with other similar data sets in the system. The inconsistencies detected or removed may have been originally caused by user entry errors, by corruption in transmission or storage, or by different data dictionary definitions of similar entities in different stores. Data cleaning differs from data

Chapter 4

System Architecture

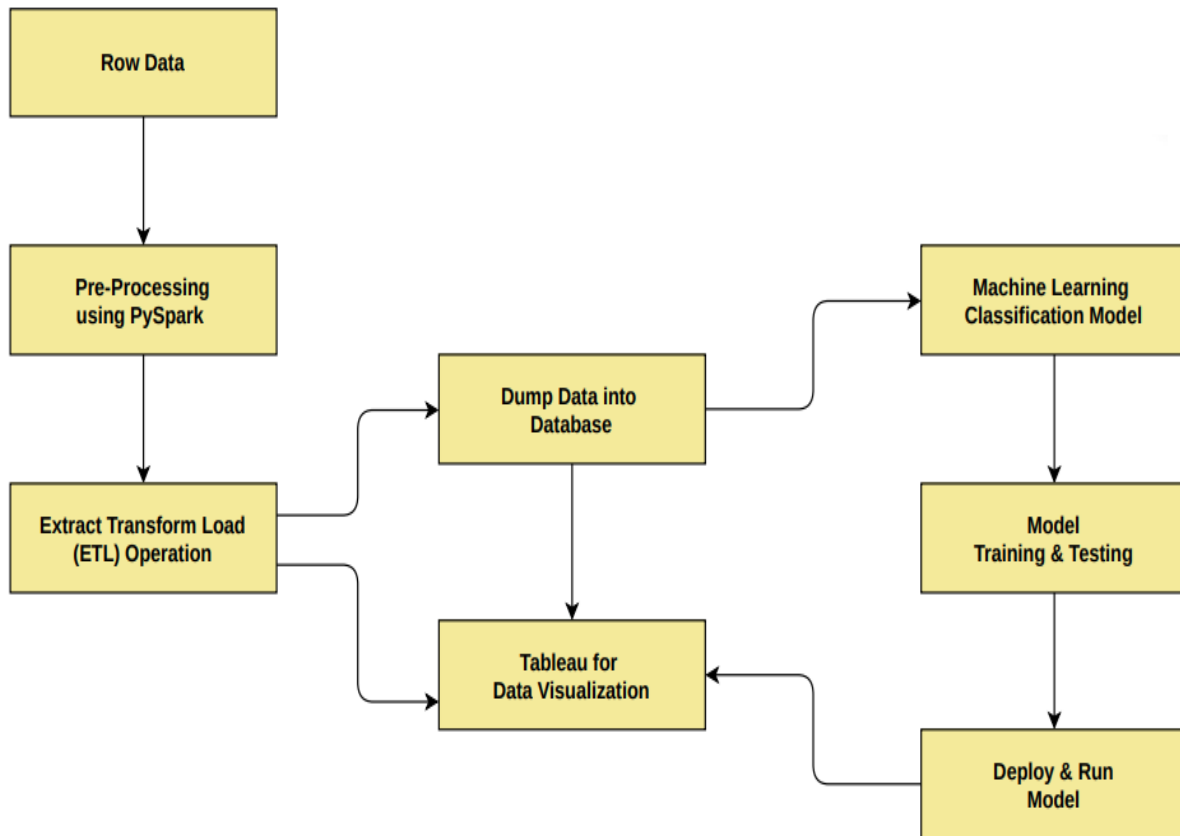


Fig: System Architecture of Telecom Customer Churn Prediction

Chapter 5

Methodology

In this project we have applied various different types of Classification.

- Logistic Regression
- Decision Tree Classifier
- Random Forest Classifier
- Stacking Classifier
- Voting Classifier
- Boosting
- Principal component analysis (PCA)
- SMOTE

During the implementation we analyze the accuracy of all the algorithms.

Machine Learning Algorithms

Logistic Regression

Logistic regression is used for solving the classification problems. In Logistic regression, instead of fitting a regression line, we fit an "S" shaped logistic function, which predicts two maximum values (0 or 1).

Decision Tree Classifier

Decision Tree is a Supervised learning technique that can be used for both classification and

Regression problems, but mostly it is preferred for solving Classification problems. It is a tree structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome. In a Decision tree, there are two nodes, which are the Decision Node and Leaf Node. Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.

Random Forest Classifier

It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model. As the name suggests, "Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset."

Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.

Stacking Classifier

In stacking, a new model is trained from the predictions from various models Predicted columns act as features with response variable as the original one Usually we stack weaker models at lower level and stronger models at upper level But, there is no such rule as to which model to be used at lower level and higher level.

Voting Classifier

The voting classifier aggregates the predicted class or predicted probability on basis of hard voting or soft voting. So if we feed a variety of base models to the voting classifier it makes sure to resolve the error by any model.

Hard Voting: Voting is calculated on the predicted output class.

Soft Voting: Voting is calculated on the predicted probability of the output class.

Principal component analysis (PCA)

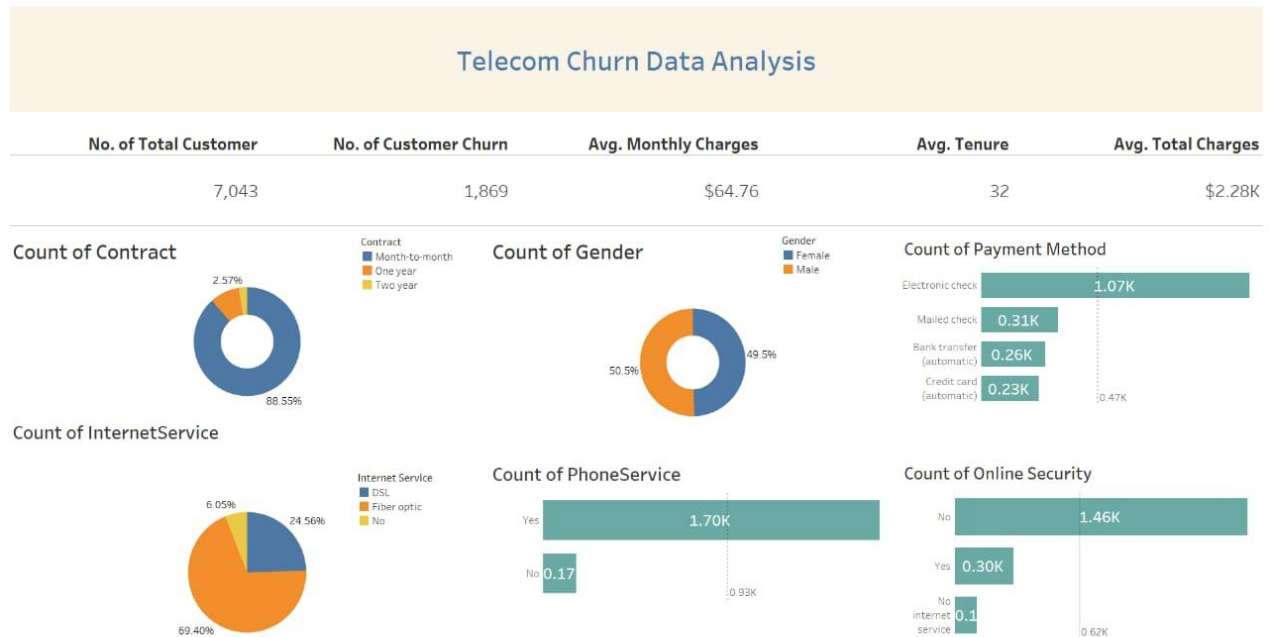
Principal component analysis (PCA) is the process of computing the principal components and using them to perform a change of basis on the data, sometimes using only the first few principal components and ignoring the rest. The importance of each component decreases when going to 1 to n, it means the 1 PC has the most importance, and n PC will have the least importance. In this technique we selected 16 PCA and performed various Machine Learning algorithms.

- 1.Decision Tree
- 2.Random forest
- 3.SVC (linear, rbf, poly)

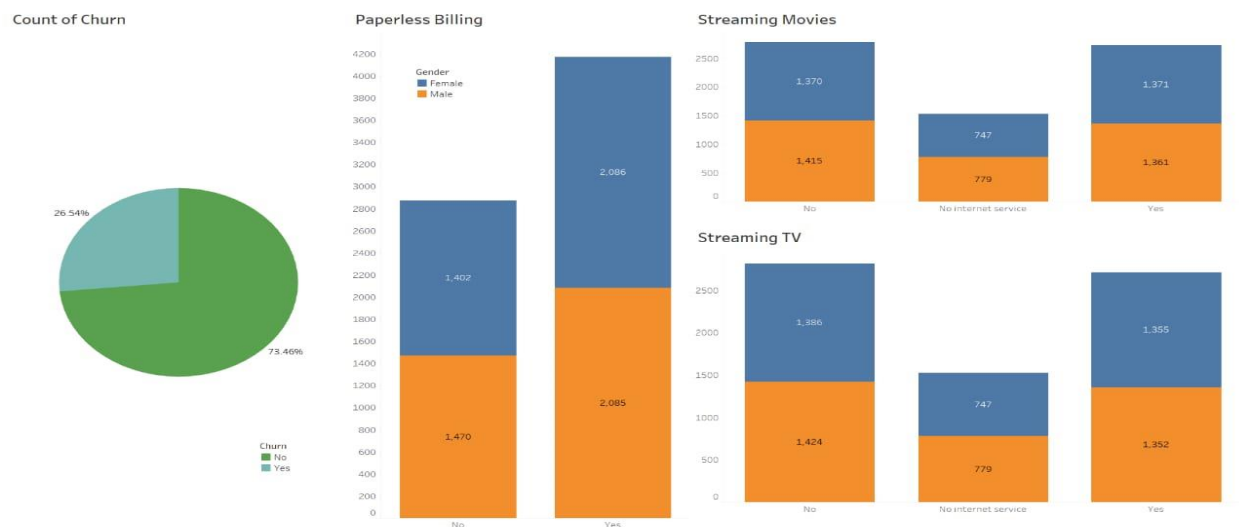
Chapter 6

Data Visualization and Representation

- Dashboard based on Key Performance Indicators (KPI) effects on contract, payment method, internet service, phone service, online security



- Dashboard based on churn effect on paperless billing, streaming Movies, Streaming TV Gender wise



- Dashboard based on churns percentage of features partner, senior citizen, gender, Dependent



➤ Output Screenshots

- Spark Session Creation

```
File Edit View Run Kernel Tabs Settings Help
ModelBuilding.ipynb Data_Cleaning_EDA.ipynb
[3]: # Create a Spark session
spark = SparkSession.\
    builder.\
    appName("TelecomCustomerChurn_Cleaning_EDA").\
    config("spark.jars.packages", "org.mongodb.spark:mongo-spark-connector_2.12:3.0.0").\
    getOrCreate()

[4]: spark

[4]: SparkSession - in-memory

SparkContext
Spark UI
Version v3.4.2
Master local[*]
AppName TelecomCustomerChurn_Cleaning_EDA
```

- Spark UI

Spark Jobs (?)

User: 7869s
Total Uptime: 2.0 min
Scheduling Mode: FIFO
Active Jobs: 1
Completed Jobs: 21

Event Timeline

Active Jobs (1)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
21	rdd at MongoSpark.scala:154 rdd at MongoSpark.scala:154 (kill)	2024/08/16 01:00:04	2 s	0/1	0/5

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Completed Jobs (21)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Executors

[Show Additional Metrics](#)

Summary

	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Excluded
Active(1)	0	62.6 KiB / 434.4 MiB	0.0 B	8	0	0	88	88	4.1 min (1 s)	136.7 MiB	67.4 MiB	67.4 MiB	0
Dead(0)	0	0.0 B / 0.0 B	0.0 B	0	0	0	0	0	0.0 ms (0.0 ms)	0.0 B	0.0 B	0.0 B	0
Total(1)	0	62.6 KiB / 434.4 MiB	0.0 B	8	0	0	88	88	4.1 min (1 s)	136.7 MiB	67.4 MiB	67.4 MiB	0

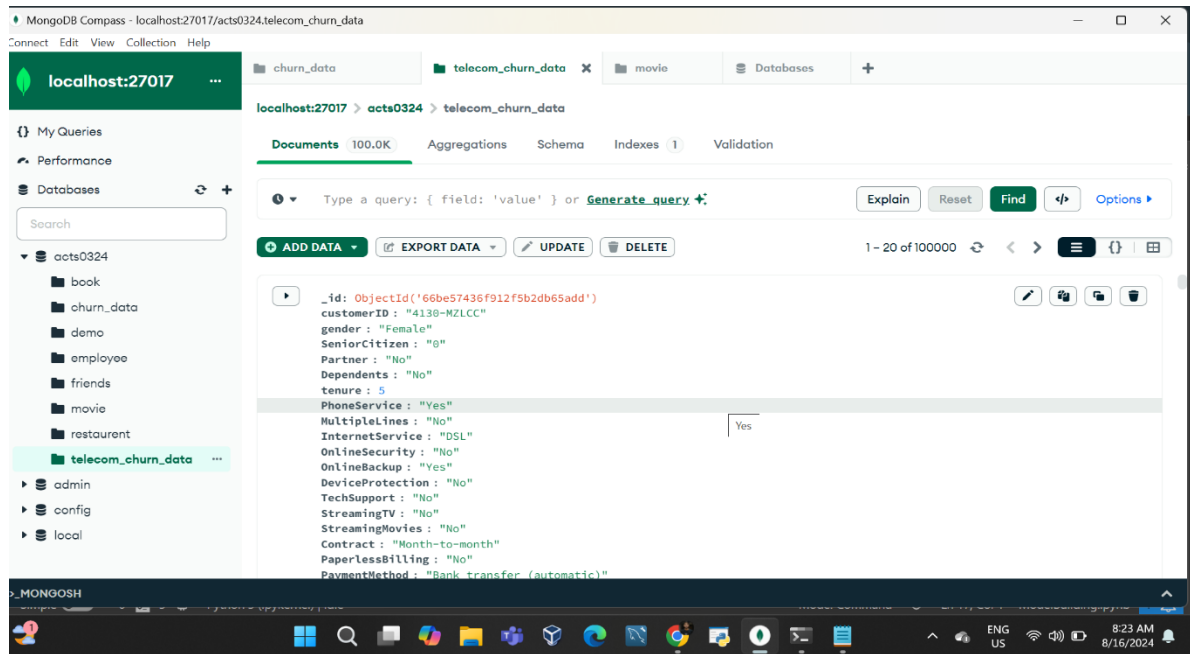
Executors

Show 20 entries

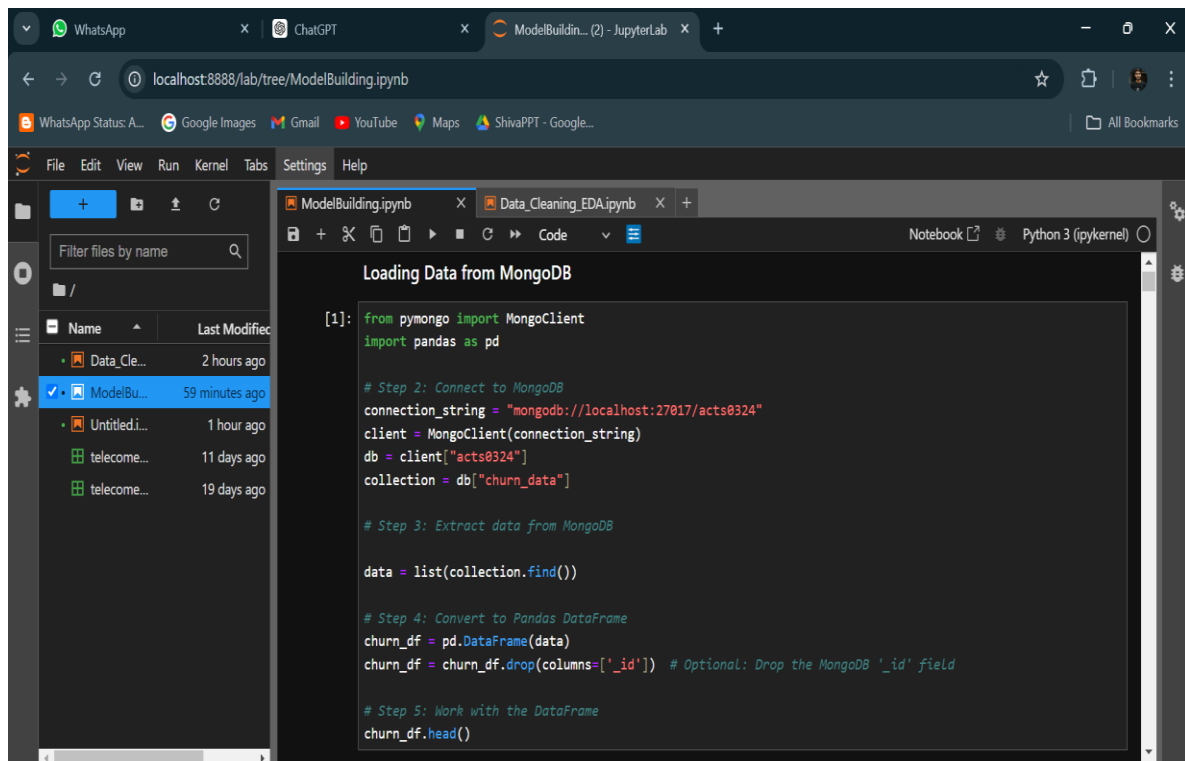
Search:

Executor ID	Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Th
driver	localhost:64899	Active	0	62.6 KiB /	0.0 B	8	0	0	88	88	4.1 min (1	136.7 MiB	67.4 MiB	67.4 MiB	Thre

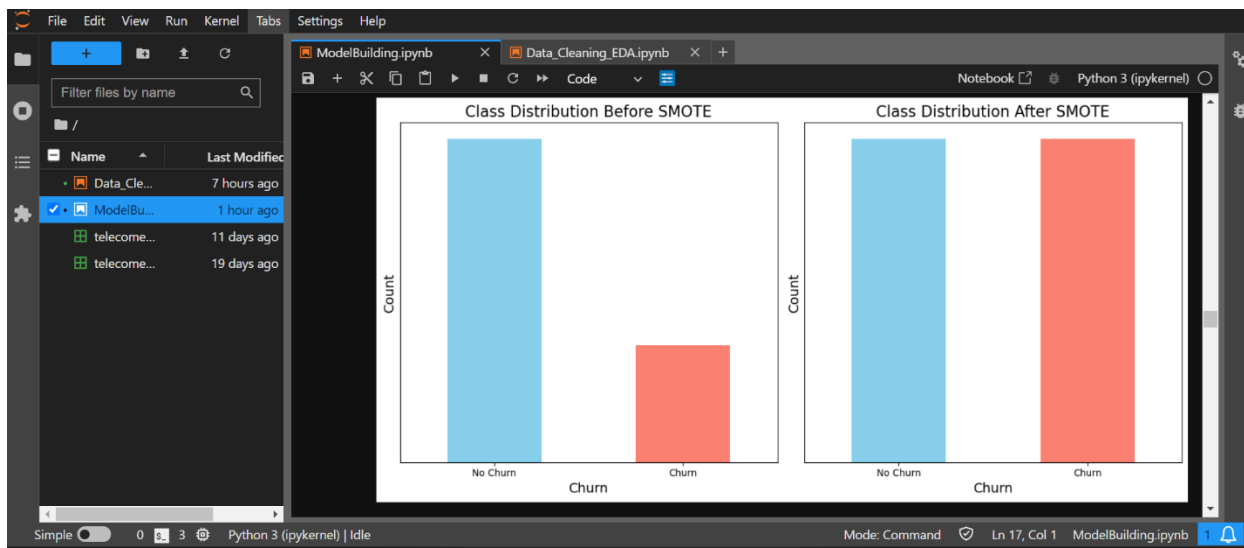
• Loaded Data in MongoDB



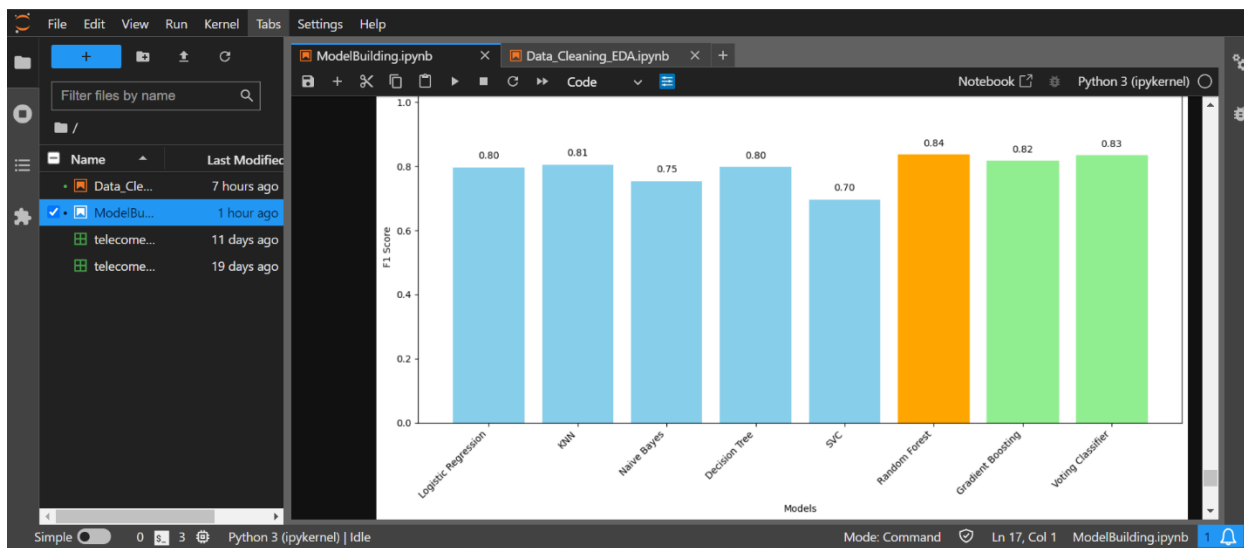
• Extracting Data from MongoDB



• Manage Imbalance Data (SMOTE)



• Models Performance



Conclusion And Future Scope

Telecom customer churn is a central issue for telecom companies, since it decreases profits. Furthermore, preventing customer churn is crucial. As the global telecom industry is becoming more dissolving and companies are increasingly struggling to retain customers. Currently, most companies invest heavily in marketing to attract new customers. However, keeping existing customers is cheaper than acquiring new customers. Thus, it is becoming more critical and a significant concern for telecommunication companies to prevent customer churn. We used various classification models to predict telecom churn using customer churn data of telecom industry.

Moreover, the results of this project will give us the ability to predict customer behavior and loss accurately and to optimize their strategies to improve customer retention rates. Meanwhile, the findings will help companies reduce costs and optimize their budgets. Furthermore, for telecom companies, it will be possible to improve customer targeting through the results of this project and to increase the profits of telecom companies.

References

- <https://spark.apache.org/docs/latest/api/python/index.html>
- <https://raw.githubusercontent.com/IBM/telco-customer-churn-onicp4d/master/data/Telco-Customer-Churn.csv>
- https://imbalanced-learn.org/stable/references/generated/imblearn.over_sampling.RandomOverSampler.html
- <https://www.investopedia.com/terms/c/churnrate.asp>
- <https://www.analyticsvidhya.com/blog/2021/06/understanding-random-forest>
- <https://towardsdatascience.com/using-principal-component-analysis-pca-for-machine-learning-b6e803f5bf1e>