# Telecom Churn Prediction Using PySpark

## Abstract

This report presents the development and evaluation of a churn prediction model using PySpark with the help of a 'telecom churn' dataset from IBM.. In the competitive telecom industry, customer retention is crucial for sustaining business growth and profitability. Telecom companies can leverage churn prediction models to identify customers at risk of leaving and take proactive measures to retain them. PySpark, with its scalability and speed, is a powerful tool for building such models. Its distributed computing capabilities enable processing of large datasets efficiently, while its machine learning library, MLlib, provides a wide range of algorithms for churn prediction.

## 1. Introduction

### 1.1 Problem Statement

Customer churn is a significant challenge for businesses in the telecommunications industry. Identifying customers at risk of churn is crucial for implementing proactive retention strategies and improving profitability. Traditional methods of churn prediction often rely on manual analysis and simple heuristics, which may not be effective in handling the large volumes of data generated in the telecom industry.

### 1.2 Motivation

The motivation behind this project is to develop a churn prediction model that can accurately identify customers at risk of churn. By doing so, businesses can implement targeted retention strategies, such as personalized offers or discounts, to retain these customers and improve overall profitability. This project also aims to demonstrate the capabilities of PySpark in handling large datasets and building machine learning models for churn prediction.

### 1.3 Objectives

The objectives of this project include:

- Developing a churn prediction model using PySpark.
- Evaluating the model's performance on the 'telecom churn' dataset.
- Providing insights and recommendations for businesses based on the model's predictions.

# 2. Background

## 2.1 Churn Prediction

Churn prediction is the task of identifying customers who are likely to stop using a product or service. It is an essential problem in customer relationship management, as retaining existing customers is often more cost-effective than acquiring new ones. Traditional approaches to churn prediction often rely on manual analysis and simple heuristics, which may not be effective in handling the large volumes of data generated in the telecom industry.

## 2.2 PySpark

PySpark is the Python API for Apache Spark, an open-source, distributed computing framework, and set of libraries for real-time, large-scale data processing. It is particularly advantageous for Python developers familiar with libraries like Pandas, as PySpark offers scalability for analyses and pipelines.

Apache Spark serves as a computational engine for processing huge datasets in parallel and batch systems. While Spark is primarily written in Scala, PySpark was developed to facilitate Spark's collaboration with Python.

## 2.3 Databricks

We used Databricks platform for this project. Databricks simplifies the process of working with PySpark by providing a user-friendly interface and a powerful set of tools. Here's how PySpark is used with Databricks:
1. **Workspace**: Databricks provides a workspace where users can create and manage notebooks.
2. **Cluster Management**: Databricks allows users to create and manage clusters for running PySpark jobs. A cluster is a set of virtual machines (VMs) that work together to process data. Databricks offers various types of clusters, such as Standard, High Concurrency, and GPU, to suit different needs.
3. **Integration with PySpark:** Databricks seamlessly integrates with PySpark, allowing users to write PySpark code in notebooks and run it on the Databricks cluster. PySpark libraries and functionalities are readily available in the Databricks environment, making it easy to leverage Spark's capabilities for data processing and analysis.
4. **Optimized Performance:** Databricks optimizes PySpark performance by automatically tuning the cluster based on the workload. It uses advanced techniques such as dynamic resource allocation and query optimization to ensure that jobs run efficiently and quickly.

## 2.4 PySpark and 'Telecom Churn' Dataset

The project utilises IBM's telecom churn dataset available on Kaggle. This data is provided by a fictional telco company that provided home phone and Internet services to 7043 customers in California**.** The dataset is a prime candidate for PySpark due to its complexity and the need for preprocessing. PySpark's ability to handle large datasets makes it well-suited for analyzing and processing it. Additionally, PySpark's integration with Python's data science libraries makes it easy to perform advanced analytics and machine learning tasks on the dataset.

PySpark's strengths lie in its ability to distribute data processing tasks across a cluster of machines, enabling parallel processing and efficient utilization of resources. This allows PySpark to handle large-scale datasets with ease, making it an ideal choice for churn prediction.

# 3. Methodology

## 3.1 Data Analysis

First, we started by understanding our data: determining whether our dataset is balanced or imbalanced, identifying our categorical and numerical columns, and analyzing how different columns affect our target variable.

## 3.2 Data Preprocessing

We built a pipeline to preprocess the data and train the model. The pipeline included steps for one-hot encoding and string indexing for categorical features, quant discretizing and vector assembling for model input.

## 3.3 Model Selection and Training

We split the dataset into training and test sets using a 80:20 ratio. This allowed us to train the model on a subset of the data and evaluate its performance on unseen data. We chose Logistic Regression as the model for churn prediction due to its simplicity and interpretability. We also considered the Random Forest model but ultimately decided on Logistic Regression for its ease of implementation and understanding.

## 3.4 Model Evaluation

Finally, after Training our models, we evaluated their performance using metrics such as accuracy, precision, recall, and F1 score. Since our dataset is imbalanced, precision and recall are more important metrics for evaluating our models than accuracy.

# 4. Data Preprocessing

## 4.1 Data Exploration

We conducted a thorough exploration of the 'telecom churn' dataset to understand its characteristics and identify potential patterns. This involved calculating summary statistics, visualizing data distributions, and identifying correlations between different variables. The exploration helped us gain insights into the dataset and identify features that are predictive of churn.

The data preprocessing phase was essential in preparing the 'telecom churn' dataset for model building. The following steps were performed to clean and process the data:

1. **Built a Pipeline:** A pipeline object, from pyspark.ml library, was created with stages set to include the preprocessing steps and to automate the data preprocessing, ensuring a streamlined workflow.
2. **Categorical Column Encoding:** All categorical columns in the dataset were identified and processed using the StringIndexer and OneHotEncoder from the pyspark.ml.feature module. This process converted categorical values into numerical format, making them suitable for model training.
3. **Processing the Target Column:** The target column 'Churn' was processed to convert it into binary values (0s and 1s) using StringIndexer. This step was crucial for the model to interpret the target variable correctly during training.
4. **Feature Assembly:** Features for model training were assembled using the VectorAssembler from pyspark.ml.feature. This step combined the encoded categorical columns and numerical columns into a single feature vector, which served as input for the model.
5. **Data Transformation:** Both the training and test datasets were transformed using the fitted pipeline. This applied the preprocessing steps to the data, preparing it for model evaluation.

# 5. Results

## 5.1 Model Performance

The best Logistic Regression model achieved an accuracy of 77.9% accuracy, 61.1% precision, and 49.4% recall.

Across various performance metrics, the Random Forest model consistently yielded similar results. It achieved an accuracy of 77.4%, precision of 60.3%, and recall of 48.1%.

## 5.2 Interpretation of Results

While both models exhibit similar accuracy, the emphasis on recall underscores the importance of minimizing false negatives (missed opportunities to retain customers) in customer retention efforts. The Logistic Regression model with the highest recall prioritizes identifying customers who are likely to churn, achieving a recall of 49.4% compared to the Random Forest model's 40.1%. Therefore, the Logistic Regression model is favored for its ability to capture potential churners effectively.

# 6. Conclusion and Future Work

## 6.1 Summary of Findings

In conclusion, we have successfully developed and evaluated a churn prediction model using PySpark and the 'telecom churn' dataset. The model shows promising results in predicting customer churn, which could be valuable for businesses in customer retention efforts. The model's accuracy, precision, and recall demonstrate its effectiveness in identifying customers at risk of churn.

## 6.2 Limitations

One limitation of our model is that it relies on historical data to predict future churn. Future work could involve incorporating real-time data and using more advanced machine learning techniques to improve the model's performance. Additionally, the model's performance may vary depending on the specific characteristics of the dataset and the business context.

## 6.3 Future Work

Future work could also involve exploring different features and modelling techniques to further improve the accuracy of the churn prediction model. Additionally, integrating the model into a real-time system for continuous monitoring and prediction of customer churn could be beneficial for businesses. Overall, the churn prediction model has the potential to be a valuable tool for businesses in the telecommunications industry in improving customer retention and profitability.

# 7. Appendix

## 7.1 Data Description

The dataset contains information about customers and their churn behavior. It includes the following important features:

- **gender**: Gender of the customer (Male/Female)
- **SeniorCitizen**: Whether the customer is a senior citizen or not (1, 0)
- **Partner**: Whether the customer has a partner or not (Yes, No)
- **Dependents**: Whether the customer has dependents or not (Yes, No)
- **tenure**: Number of months the customer has stayed with the company
- **PhoneService**: Whether the customer has a phone service or not (Yes, No)
- **InternetService**: Customer's internet service provider (DSL, Fiber optic, No)
- **Contract**: The contract term of the customer (Month-to-month, One year, Two year)
- **PaperlessBilling**: Whether the customer has paperless billing or not (Yes, No)
- **PaymentMethod**: The customer's payment method (Electronic check, Mailed check, Bank transfer (automatic), Credit card (automatic))
- **MonthlyCharges**: The amount charged to the customer monthly
- **TotalCharges**: The total amount charged to the customer
- **Churn**: Whether the customer churned or not (Yes, No)