# CISC5950_Final Project Report

Haoran Xue, Sheng Yun

May 7th, 2023

# Introduction

In this project, we would be leveraging the power of Apache Spark for data processing and machine learning tasks, focusing on four distinct parts, including a basic text cleaning process using the toxic comment text classification dataset, applying Logistic Regression to identify the most significant risk factors associated with heart disease and predict overall risk levels, analyzing the Census Dataset from the UCI Machine Learning Repository, as well as employing Random Forest Classifier and Decision Tree Classifier using Apache Spark ML/MLlib to train the dataset. The main objectives of the project were to provide an informative explanation of the data cleaning process, the pipeline, and the explanation and evaluation of model performances. The code utilized various libraries, such as PySpark, Pandas, and Sys, and functionalities to execute the data pre-processing, feature engineering, and model training tasks.

# P1: Toxic Comment Classification

In this part, we were going to explore basic text processing using the toxic comment text classification dataset. Although the machine learning and text processing techniques employed were not sophisticated, our primary objective was to convert the comment text column into a sparse vector representation that could be utilized by a classification algorithm within the Spark ML library. According to the instructions, we only need to make the sample code work, so the main problem was just to convert the code to the spark code that looked like ours. The biggest challenge here was the resource setup because the sample code allocated the amount we could not give on our setup clusters(e2-standard), which only have 2GB and an 8-core CPU. This gave us a hard time testing out the settings(larger resources, the scheduler would not accept, but smaller ones would make the server idling). Then we gave up on setting up our own and chose to use the default settings(no configs at all!), the program run perfectly after making that changes.

## Data Description and Data Cleaning

The Toxic Comment Text Classification dataset is a large dataset originally released by Jigsaw and Google. The data is formatted as a 'csv' file, with each row representing a unique comment. The columns in the dataset include id, comment_text, and other binary labels (0 or 1) indicating whether the comment falls into the respective toxicity category.

The data cleaning process began with the conversion of input 'csv' files into Spark DataFrames using the 'to_spark_df()' function. This function read a 'csv' file using the Pandas library, replaced missing values with empty strings, and then converted the Pandas DataFrame into a Spark DataFrame. This process ensured that the data was cleaned and suitable for further analysis within the PySpark environment. The cleaned DataFrames were subsequently used as training and testing datasets.

## Model Explanation and Evaluation

Their designed pipeline includes several steps of to do this task, including tokenization, feature extraction, and model training. First, the Tokenizer was applied to the 'comment_text' column to split the text into a list of words. Then the HashingTF (hashing term frequency) technique was employed to transform the list of words into numerical features. Following this, the IDF (inverse document frequency) was calculated to weigh the importance of words in the dataset. Finally, the product of TF and IDF was computed to generate the TF-IDF features, which represented the significance of words in each comment. The chosen model for this task was Logistic Regression, which was implemented using the PySpark 'LogisticRegression' class. The model was trained separately for each label column in the 'out_cols' list, excluding the 'id' and 'comment_text' columns. The regularization parameter (REG) was set to 0.1 to prevent overfitting. For each label, the model was fit on the training dataset with TF-IDF features and then used to make predictions on the test dataset. The probability of a comment being classified as toxic was extracted and added to the 'test_res' DataFrame. This process was repeated for all labels in the 'out_cols' list. By incorporating a well-structured data-cleaning process and a comprehensive pipeline, the program successfully prepared the data, extracted relevant features, and trained the model for accurate classification. The use of Logistic Regression ensured that

Figure 1: Toxic Comment with LR

the model was both interpretative and efficient in identifying toxic comments. We got the first 20 rows with the comment ids and the predicted probabilities for each class as output, shown in Figure 1.

## P2: Heart Disease Prediction using Logistic Regression

In this part, we converted a Python script to PySpark to predict the likelihood of heart disease within a 10-year period using logistic regression. The Python script used Logistic Regression to identify the most significant risk factors associated with heart disease and predict overall risk levels. In this question, we would be using the Framingham Heart dataset.

### Data Description and Data Cleaning

The Framingham Heart Dataset is originated from the Framingham Heart Study, with an initial cohort of 5209 subjects. The Framingham Heart dataset comprises various demographic, clinical, and laboratory attributes collected from the study participants. The data typically includes various demographic information about patients, such as age, sex, BMI, the measurement information of patients such as heart rate, total cholesterol level, systolic blood pressure, diastolic blood pressure, and fasting blood pressure etc., and the outcome of the presence or absence of coronary heart disease. This dataset enables researchers to investigate the relationship between various risk factors.

The data cleaning process began with loading the dataset using the 'csv' format and inferring the schema as usual. The 'education' column was dropped because it is not relevant to the prediction task. Additionally, we renamed the 'male' column to 'Sex_male' for clarity. Missing values were removed using the 'dropna()' function, and the dataset was split into training and test sets with an 8-2 ratio. The 'cigsPerDay', 'totChol', and 'glucose' columns were then cast to the 'FloatType', and any remaining missing values were dropped.

Figure 2: Hear Dsease with LR

## Model Explanation and Evaluation

We used several steps of pipeline, including feature engineering, model training, and evaluation. First, the 'VectorAssembler' was used to combine the selected feature columns into a single 'features' column. The features were then standardized using the 'StandardScaler' to ensure that they had equal importance during model training. Next, a logistic regression model was defined with 'scaled_features' as input and 'TenYearCHD' as the label column. The pipeline was then built using the 'Pipeline' class, with the assembler, scaler, and logistic regression model as stages.

The chosen model for this task was logistic regression, so we implemented using the PySpark 'LogisticRegression' class. Logistic regression is a popular choice for binary classification tasks as it provides an interpretable and efficient method for predicting probabilities. The model is trained on the preprocessed and standardized features in the training dataset, and predictions are made on the test dataset using the 'transform()' function. The performance of the logistic regression model was evaluated using the area under the receiver operating characteristic (ROC) curve, which is a widely accepted metric for binary classification tasks. The 'BinaryClassificationEvaluator' class was used to calculate the area under the ROC curve, which provided a measure of the model's ability to discriminate between positive and negative cases. The resulting value was then printed to provide insight into the model's performance. As the result shown in Figure 2, we got the AUC score with 0.752, which meant that the model was moderate successful in correctly predicting instances of both classes.

## P3: Logistic Regression classifier on Census Income Data

In this part, we aimed at training a binary classifier to predict an individual's income level which had two possible values, ¿50K and ¡50K, using the Census dataset from the UCI Machine Learning Repository, based on various demographic and economic features.

### Dataset Description and Data Cleaning

The dataset on Census Dataset consists of anonymous information such as age, occupation, education, and working class, with 48,842 instances and 14 attributes, the dataset offers a diverse mix of categorical, numerical, and missing values. We would utilize Spark ML/MLlib's logistic regression capabilities to train our model and evaluate its performance, aiming to achieve a high accuracy rate.

Our data cleaning process started with loading the train and test datasets using the 'csv' format and inferring the schema. Columns are renamed to more descriptive names, and the 'income' column is processed such that values were converted to binary indicators (1 for "¿50K" and 0 otherwise). We have found there was a leading space for "¿50" in both train and test set, and a "." in the end of the "¿50K" for test set. The train and test datasets were then combined for pre-processing purposes.

Figure 3: Census with LR

Categorical columns were identified and indexed using the 'StringIndexer' class, with the '_index' suffix added to the output column names. OneHotEncoder was used to encode the indexed categorical columns, appending an '_encoded' suffix to the output column names.

## Model Explanation and Evaluation

In order to achieve our prediction goal, we implemented several processing steps in a pipeline that as the same as we did in part2. First, we employed the 'VectorAssembler' to merge the encoded categorical columns with the numerical columns, creating a unified 'features' column. Next, we utilized the 'StandardScaler' to standardize the features, ensuring that each feature carried equal significance during the model training phase. We then established a logistic regression model with 'scaled_features' as input and 'income' as the label column. Finally, we constructed the pipeline using the 'Pipeline' class, incorporating the indexers, encoders, assembler, scaler, and logistic regression model as stages in the process.

What we need to do for this task was logistic regression, which was also implemented using the PySpark 'LogisticRegression' class. The output of the logistic regression model was measured using the area under the receiver operating characteristic curve (ROC AUC) that could measure the ability of the model to distinguish between individuals with high and low income based on their characteristics. In this case, the code had produced an AUC score of 0.904, which indicated that the Logistic Regression model had good performance in predicting the income of the individuals.

# P4: Random Forest and Decision Tree on Census Income Data

In this part, we would revisit the third part's problem but instead employ Random Forest Classifier and Decision Tree Classifier using Apache Spark ML/MLlib. This approach would enable us to compare the performance of different machine learning algorithms on the same dataset, further enhancing our understanding of their effectiveness and suitability for specific tasks.

## Mmodel Explanation and Evaluation

The script utilized another two different classification models: Random Forest, and Decision Tree classifiers. Random forest is an ensemble method that combines multiple decision trees and averages their predictions, making it more robust and less prone to overfitting. Decision tree classifier is simple yet powerful model that recursively split the feature space based on feature importance. Both of the two models were trained on the pre-processed and standardized features in the training dataset,

Figure 4: Census with RF&DT

and predictions were made on the test dataset using the 'transform()' function. The performance of the random forest, and decision tree models was evaluated using the area under the receiver operating characteristic (ROC) curve, we got the AUC of Random Forest was 0.89, and the AUC of Decision Tree was 0.59. Compared with the AUC score of Logistic Regression we got from P3(0.904), we could say that Logistic Regression performed best among these three models in predicting the personal income, and the Decision Tree made the worth performance. This could be because Decision Trees were prone to overfitting, and could be affected by the randomness introduced during the tree construction process. However, for the results we got for now, we could not conclude that Logistic Regression was the most suitable model for the dataset, we still need to consider other metrics such as precision, recall and F1 scores.

## Conclusion

In this project, we were expected to utilize Apache Spark for data processing jobs and applied various machine learning algorithms to different datasets, and then showed their performance with Area Under ROC curve score. The project was divided into four distinct parts, each focused on a specific aspect of data processing and machine learning method. The project showed us the power and versatility of Apache Spark for handling data processing jobs.