# ID2221 Project Report: Distributed Fraud Detection using Spark ML

Emma Lind and Mattias Stahre

October 2020

## 1 Problem description

Fraud detection is a set of activities undertaken to prevent for instance, money from being obtained through false pretenses. Fraud detection can be applied in banking, fraud may include forging checks or using stolen credit cards. The focus of this project is credit card fraud detection. The goal with this project is to classify transactions using machine learning models, the important point is that we will analyse big data and need to distribute the processing on different nodes. To handle this problem we will use Spark ML. In other words, we will classify credit card transactions on a distributed system using machine learning and spark processing. This is a typical classification problem, hence we will test two different classification machine learning algorithms. More precisely, Decision Tree and Random Forest and compare their accuracy.

## 2 Dataset

The data set used in this project is from Kaggle [1]. The data is labeled and in a csv format with each row correspond to a transaction and each column correspond to a feature of the transaction, for instance, time. Number of features used were 28.

## 3 Method

Our method for solving this problem was the following;

- Firstly, upload the data to hdfs.

- Then create a program `models.scala` that evaluates how many features to use and then train the models and upload them to hdfs.

- Further, create another program `main.scala` that reads the models from hdfs and then detects frauds and prints those results to the user. It also saves the results/predictions to hdfs.

- Lastly, create a program `evalResults.py` that evaluates the results by comparing the predictions with the labeled data using a couple of metrics such as false positive, false negative and error rate.

The main development tools were Scala, the Apache Spark API and hdfs for data storage [2].

# 4    Results

When evaluating the two models a number of techniques were used. Error rate in percentage, false positive and false negative. Additionally we wanted to print the detected frauds from the data set.

When developing the code we started by simply using all the parameters in the data set to train our model. This gave poor results so we implemented our own algorithm to test how many of the parameters to use. The result from this was better but it turned out that the models we used were not deterministic which means that the suggested number of features could vary. We decided to use an arbitrary model from the training program.

As for the results of the models the following was found.

For the Decision Tree Classifier the error rate was low at 0.0747%. The false positive was also low at 0.0152%. False negative were much higher at 34% which seem to indicate that the model misses out on many transactions that were actually frauds but does not often label a non-fraudulent transaction as fraud.

As for Random Forest Classifier the error rate was low at 0.0666%. The false positives were also low at 0.0175%. Looking at the false negatives though the percentage was even higher than for the Decision Tree Classifier at 28%. This again means that the model tends to not get it wrong it terms of when frauds are not fraudulent but it misses out, in this model, on almost half of the transactions that were actually frauds.

This comparison shows that even though the error rate in terms of the total number of transactions is low for both of the models the number of missed fraudulent transactions is quite high. The reason for this is probably because of how often non-frauds occur in the data set compared to frauds. Because most of the transactions in the data set are non-fraudulent and both models are quite good at not giving false positives, evaluating the models against all transactions is not a good way to evaluate them.

Having trained a model our program for analysis takes a csv-file, which contains transactions, and returns a list of the suspected frauds. The identifier for each transaction corresponds to the row in the input data.

The figure below show that often the models predict fraud to be true (indicated by integer 1) although the label is actually false (indicated by integer 0).



Figure 1: Figure showing part of the false positives for Decision Tree Classifier.

# 5 How to run the code

First start the namenode and the datanode:
```
hdfs --daemon start namenode
hdfs --daemon start datanode
```

Upload the data to hdfs:
```
hdfs dfs -mkdir /kth
hdfs dfs -put creditcard.csv /kth
```

Check so it is uploaded:
```
hdfs dfs -ls /kth
```

Then run the first program `models.scala` to determine how many features to use and to train models: cd/trainModels
```
sbt run
```

Then run the main program `main.scala` to test the models: cd/id2221
```
sbt run
```

Lastly run the python program `evalResults.py` to evaluate the results:
```
python3 evalResults.py
```

# References

[1] Dileep, "Credit card fraud detection," Available at https://www.kaggle.com/dileep070/anomaly-detection.

[2] Available at https://spark.apache.org/.