# Real-Time Fraud Detection on Credit Cards Transactions

Prasanna Surianarayanan, Dipika Rajesh
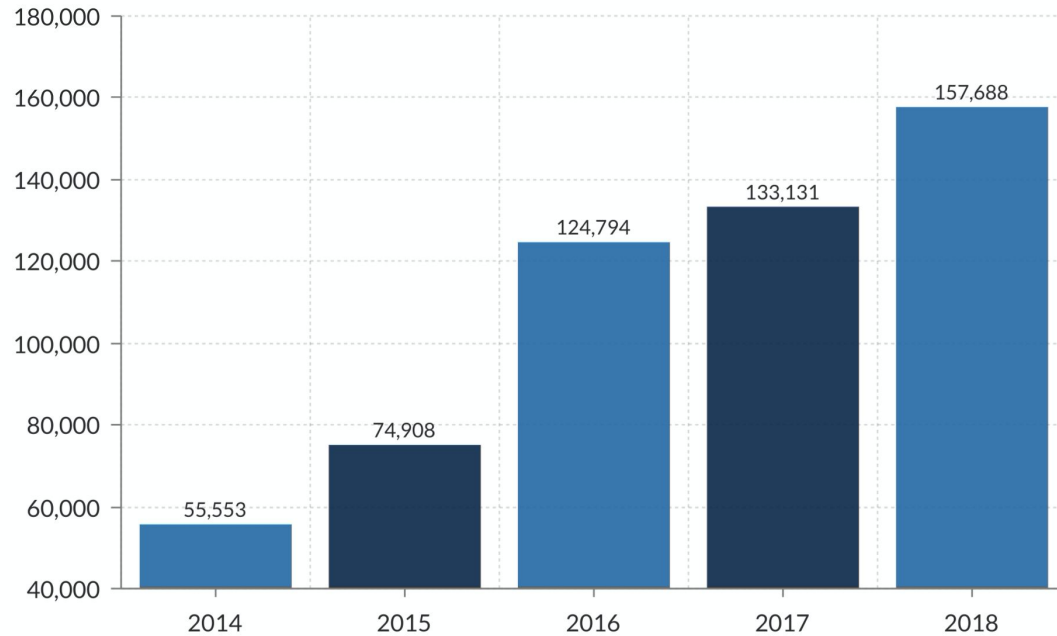
ps3703@nyu.edu, dr2898@nyu.edu

# Premise

- Fraud detection applies to a variety of sectors:
    - Banking and Finance
    - Insurance
    - Government agencies and law enforcement
- 2018 PwC Survey : fraud experienced by 49% of the 7,200 companies surveyed
- Fraud Detection : A billion-dollar business

# Need for Credit Card Fraud Detection

- Electronic commerce has rapidly increased over the years
- Increasing confidence of customers in electronic payments
- Customers increasingly vulnerable to new attacks
- Frauds need to be detected in (near) real time
- Scalable learning techniques needed to ingest and analyse massive amounts of streaming data.

# Credit Card Fraud Reports
# in the United States



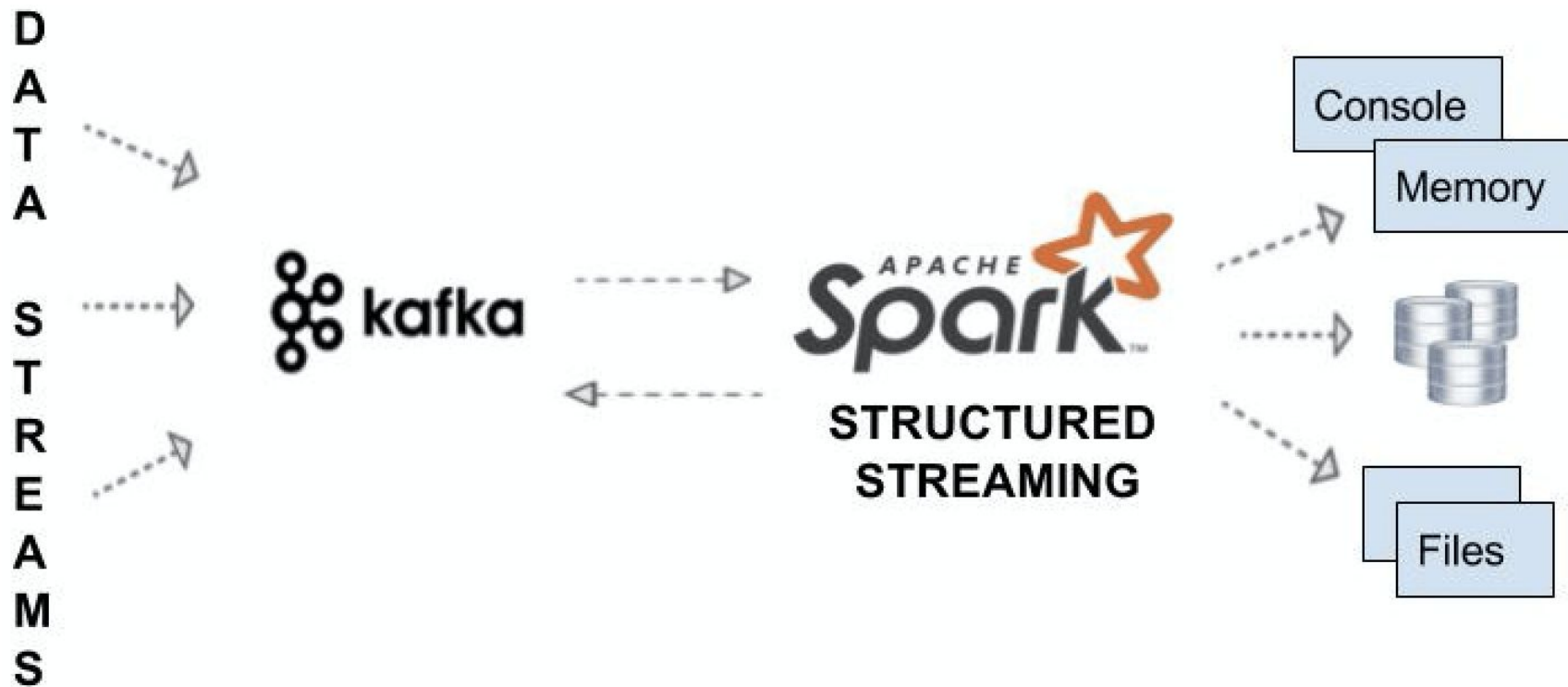| Year | Reports |
|------|---------|
| 2014 | 55,553 |
| 2015 | 74,908 |
| 2016 | 124,794 |
| 2017 | 133,131 |
| 2018 | 157,688 |

SHIFT
Credit Card Processing

# Leveraging Big Data for the Use Case

- Credit card transactions captured over time based on the cardholder's spending behaviour
- Events stored as transaction history for detection of fraud in the future
- Spark Structured Streaming and Apache Kafka work pretty well in capturing such incidents real time
- Ecosystem capable of handling vast data with good throughput and low latency
- SparkML used to develop the prediction model

# Technologies Used

- Apache Kafka for message ingestion
- Spark Structured Streaming to process datastreams from Kafka

- Spark ML to detect fraudulent transactions

# Proposed Architecture

# Ecosystem Overview

- **Apache Kafka** :
    - Transactions ingested to Kafka
    - Used for delivering streams of card transactions
    - Card transactions are published to a topic
    - Consumers read messages from subscribed topics

# Ecosystem Overview

- **Spark Structured Streaming** :
  - Reads stream of data from the Kafka topic
  - Creates dataframes in streaming mode
  - Treats arriving data as an unbounded input table
  - New items in streams treated like rows, appended to the input table
  - Queries applied to streaming dataframes just like static dataframes
  - Triggers used to control result updates

# Ecosystem Overview

- **Spark ML :**
  - Used to create Models to detect fraudulent transactions
  - Historical transactions along with the class labels are used to train the model
  - Prediction is performed on the incoming data from the Kafka topic
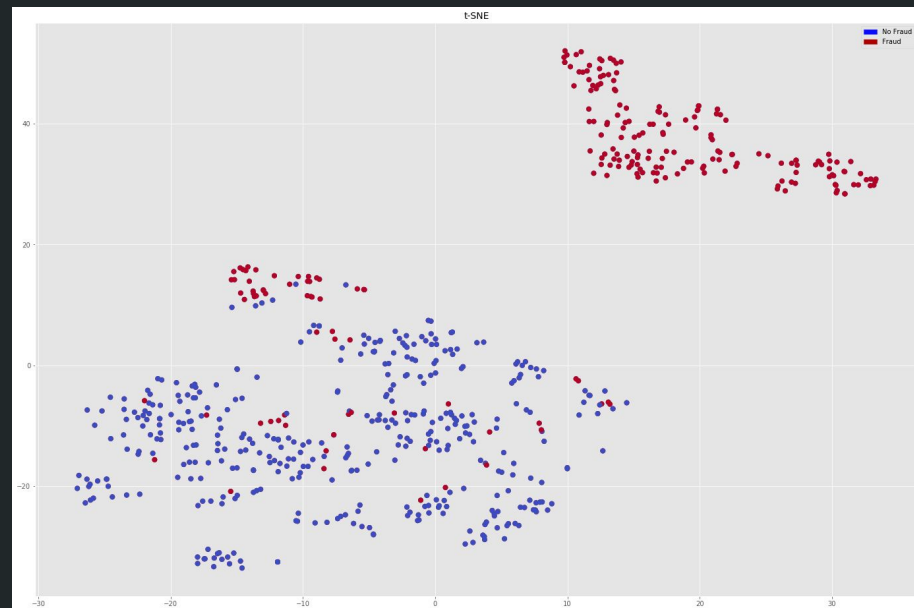  - Predicted results  written to an output sink

# Dataset Used

- Source of labelled data : Kaggle - https://www.kaggle.com/mlg-ulb/creditcardfraud
- Total features : 31
- 28 features transformed using PCA for confidentiality purposes
- 3 Non-anonymized features : Time, Amount, Class Label
- Dataset was selected due to non-synthetic, high-dimensional data.
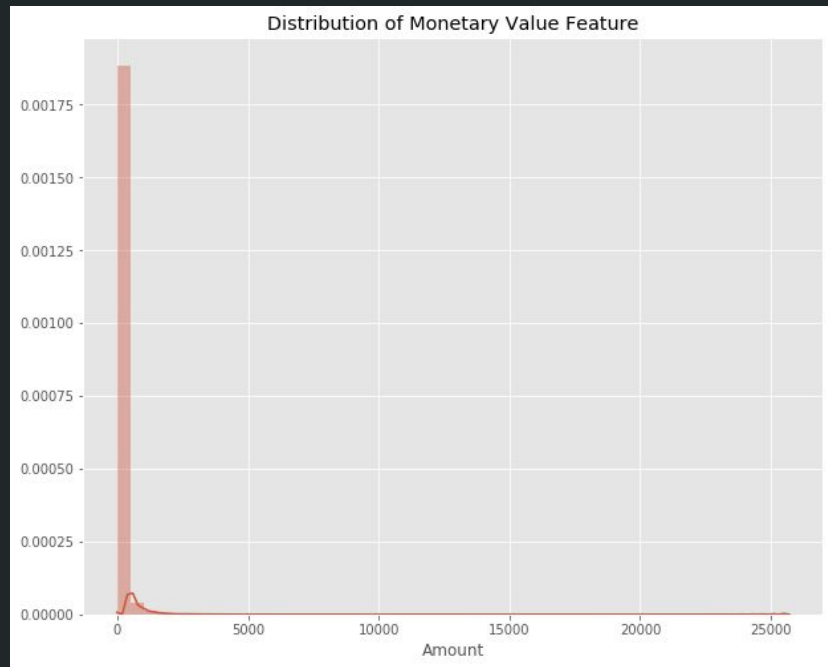
# Data Exploration!

# Dimensionality Reduction using t-SNE

**T-distributed Stochastic Neighbor Embedding (t-SNE)** nonlinear dimensionality reduction technique well-suited for embedding high-dimensional data for visualization in a low-dimensional space of two or three dimensions
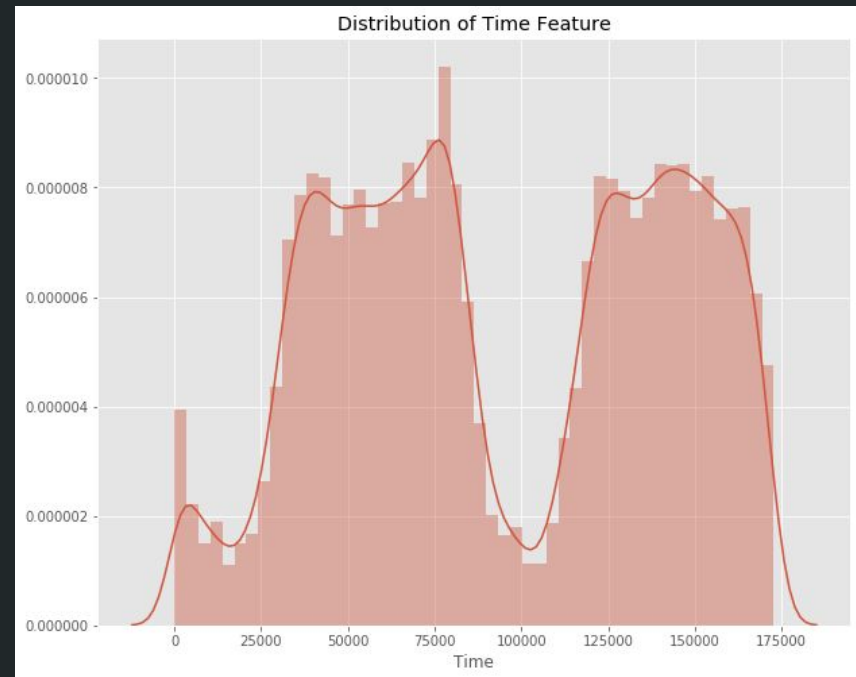
# Exploration of the Dataset

The vast majority of transactions are relatively small and only a tiny fraction of transactions comes even close to the maximum - needs to be scaled!
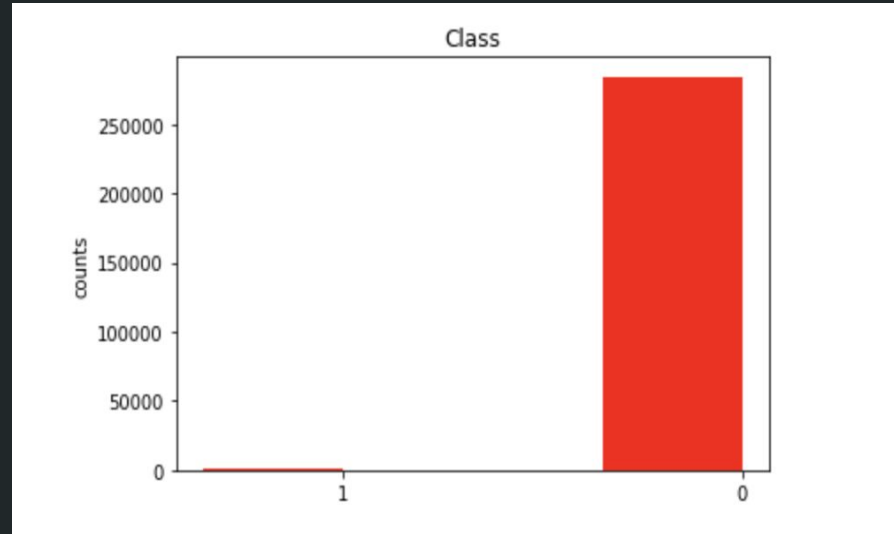


Distribution of Monetary Value Feature

- The time is recorded in the number of seconds since the first transaction in the data set.
- Can conclude that this data set includes all transactions recorded over the course of two days.



Distribution of Time Feature

# Class Distribution?

Most transactions are non-fraudulent. In fact, 99.83% of the transactions in this data set were not fraudulent while only 0.17% were fraudulent.
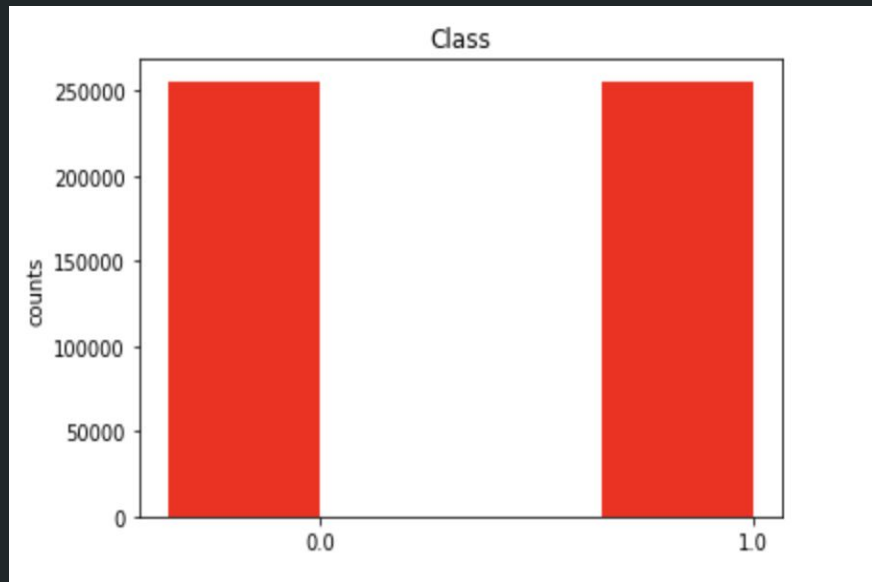
# Data Preprocessing!

- Imbalanced datasets can create biased machine learning models
- Oversampling methods (SMOTE) were applied on the dataset along with feature scaling.

# After applying SMOTE,

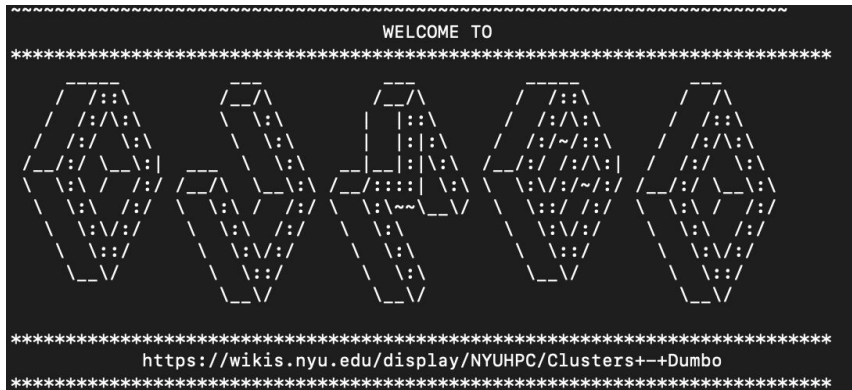the oversampling created balanced classes!

# Choosing and Fitting the Machine Learning model

# Model Choice - Random Forest Classifier

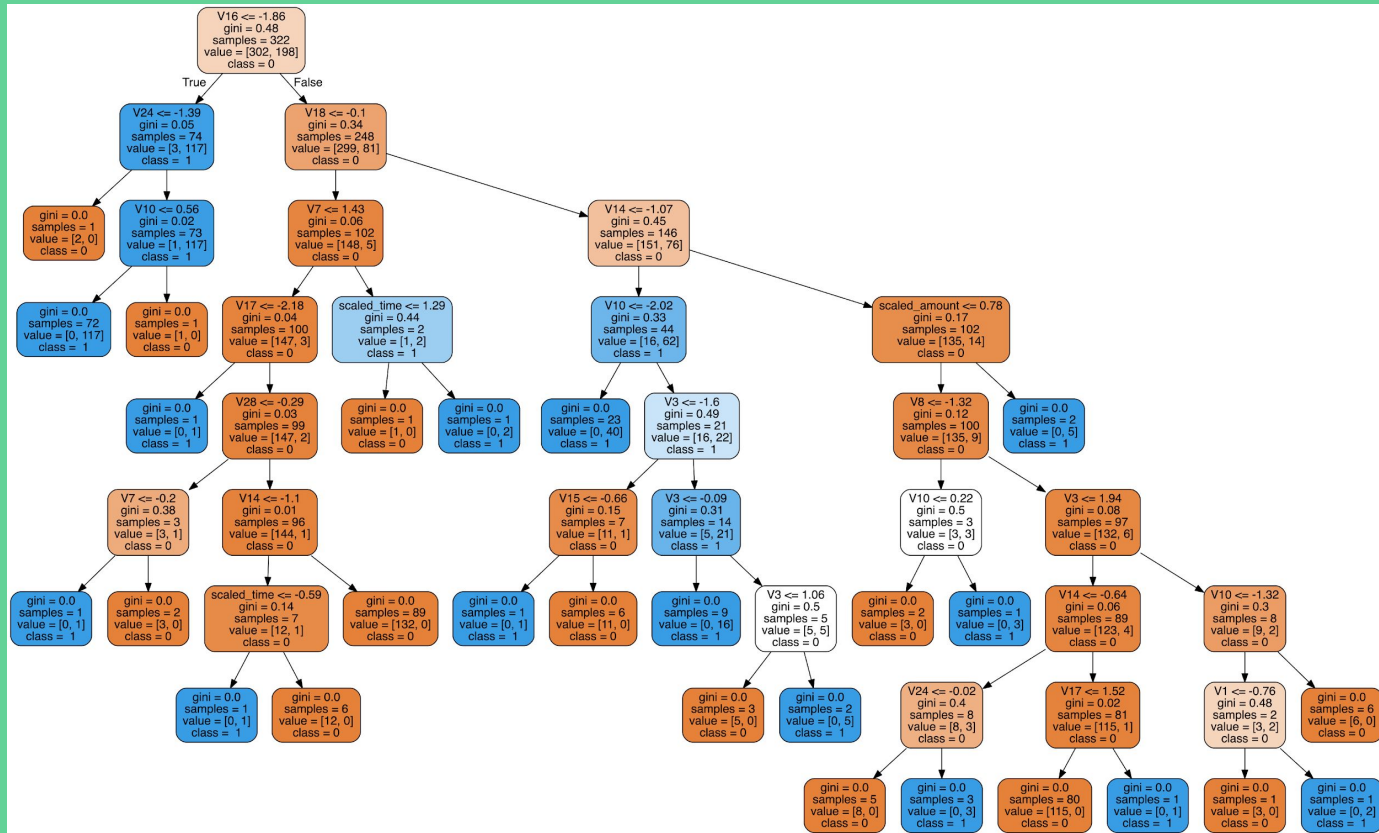**Advantages of Random Forest Classifiers:**

- Generalizes patterns well
- Robust to different input types
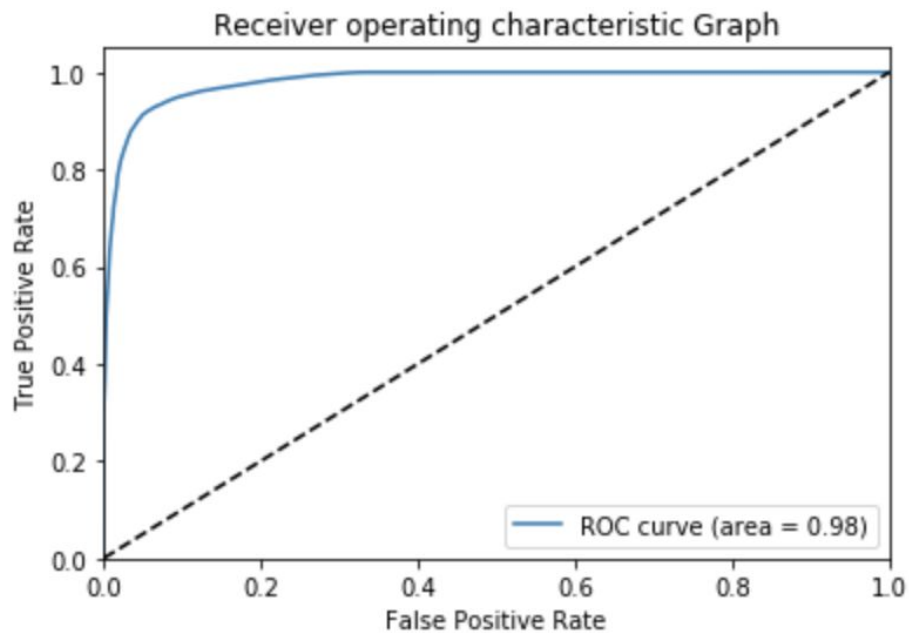- Robust to outliers
- Highly scalable

Using RandomForestClassifier from SparkML library, trained on NYU DUMBO HPC with 500,000 data instances!

```python
algo = RandomForestClassifier(featuresCol='features', labelCol='Class')
start_time = time()
model = algo.fit(trainingData)
end_time = time()
predictions = model.transform(testData)
```

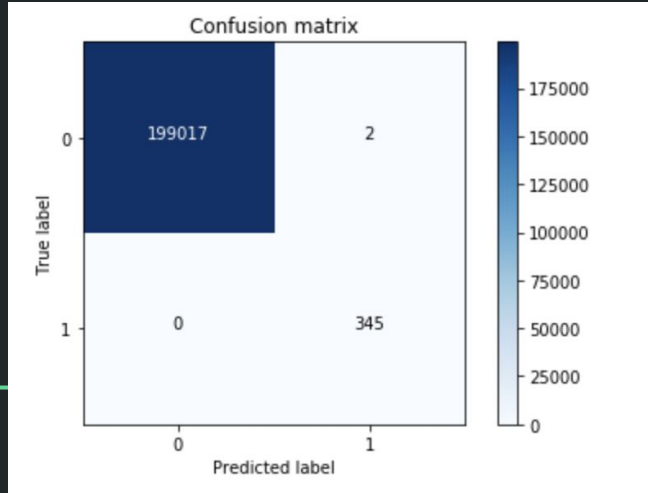# Visualizing the Random Forest model using Eureka Trees library

# Results! - Area Under the Receiver-Operating Characteristic (AUROC) metric
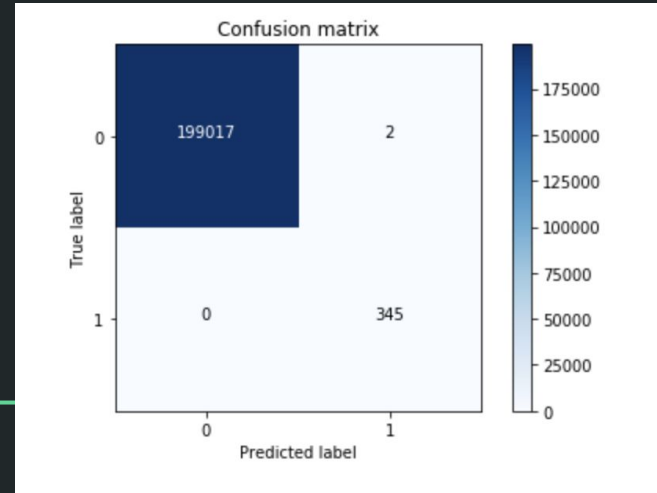


AUC — ROC Curve

# Evaluation of the random forest model - confusion matrix



Recall in the train set = 1.0

Recall in the test set = 0.8095

Never use accuracy score as an evaluation metric on imbalanced datasets - will be high and misleading.
Never test the model on the sampled (over or under) data.

# Future Work

- Adding ElasticSearch and Kibana to the pipeline for visualization and trend analysis
- Expanding model to include Time-Series Analysis of Customer Credit Card Transactions
- Online Machine Learning with incoming real-time data

# Thank you!