



CZ4041 Machine Learning

IEEE – CIS Fraud Detection Kaggle Project

Agenda

1. Introduction

2. Exploratory Data Analysis

3. Data Pre-processing

4. Modelling and Results

Team Members



Jason



Kenneth



Yong Wei



Huan Zhang

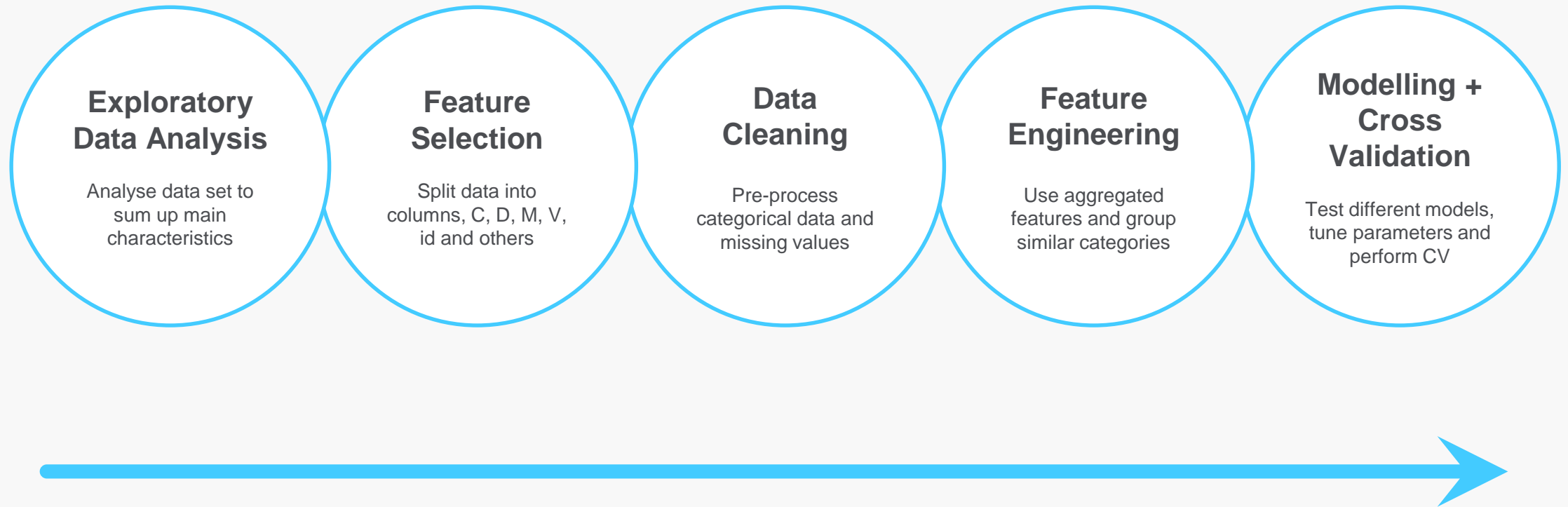
Problem Statement

The aim of the challenge is to benchmark machine learning models on a challenging large-scale dataset to **predict if a transaction made is fraudulent**.

The dataset of credit card transactions is provided by the Vesta Corporation, described as the world's leading payment service company.

The dataset includes identity and transaction CSV files for both test and train.

Overview





Exploratory Data Analysis

Exploratory Data Analysis

- Train dataset: 590,540 x 434
- Fraud transactions: 20663
- Target variable 'isFraud'

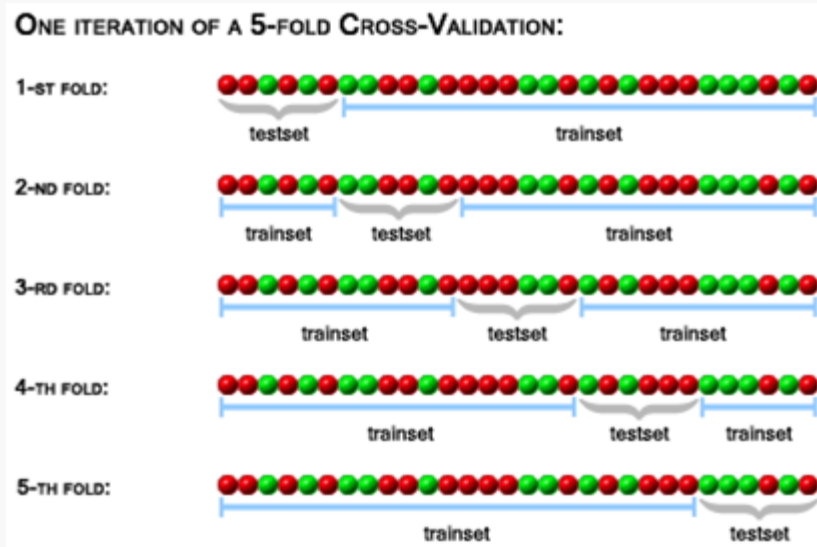
Explanation of Variables

FEATURE	DESCRIPTION
<i>TransactionDT</i>	Timedelta from a given reference datetime (not an actual timestamp)
<i>TransactionAmt</i>	Transaction amount in USD
<i>ProductCD</i>	product code, the product for each transaction
<i>card1-card6</i>	Card used for payment
<i>addr</i>	address
<i>dist</i>	distance
<i>P_ and (R_) email domain</i>	purchaser and recipient email domain
<i>C1-C14</i>	counting, address and other things, actual meaning is masked
<i>D1-D15</i>	timedelta, such as days between previous transactions. etc.
<i>M1-M9</i>	match, such as names on card and address, etc
<i>V1-V339</i>	Vesta engineered rich features, including ranking, counter and other entity relations

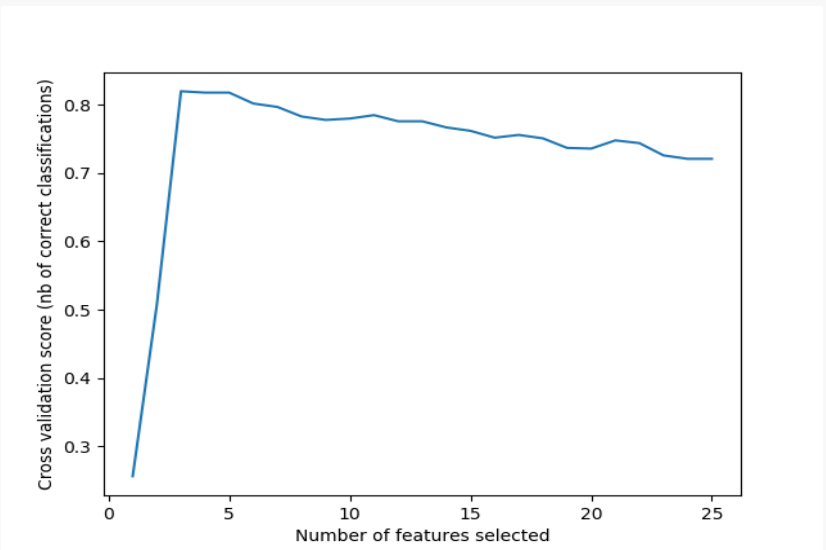


Data Pre-processing

Feature Selection method



Cross validation



Recursive feature elimination

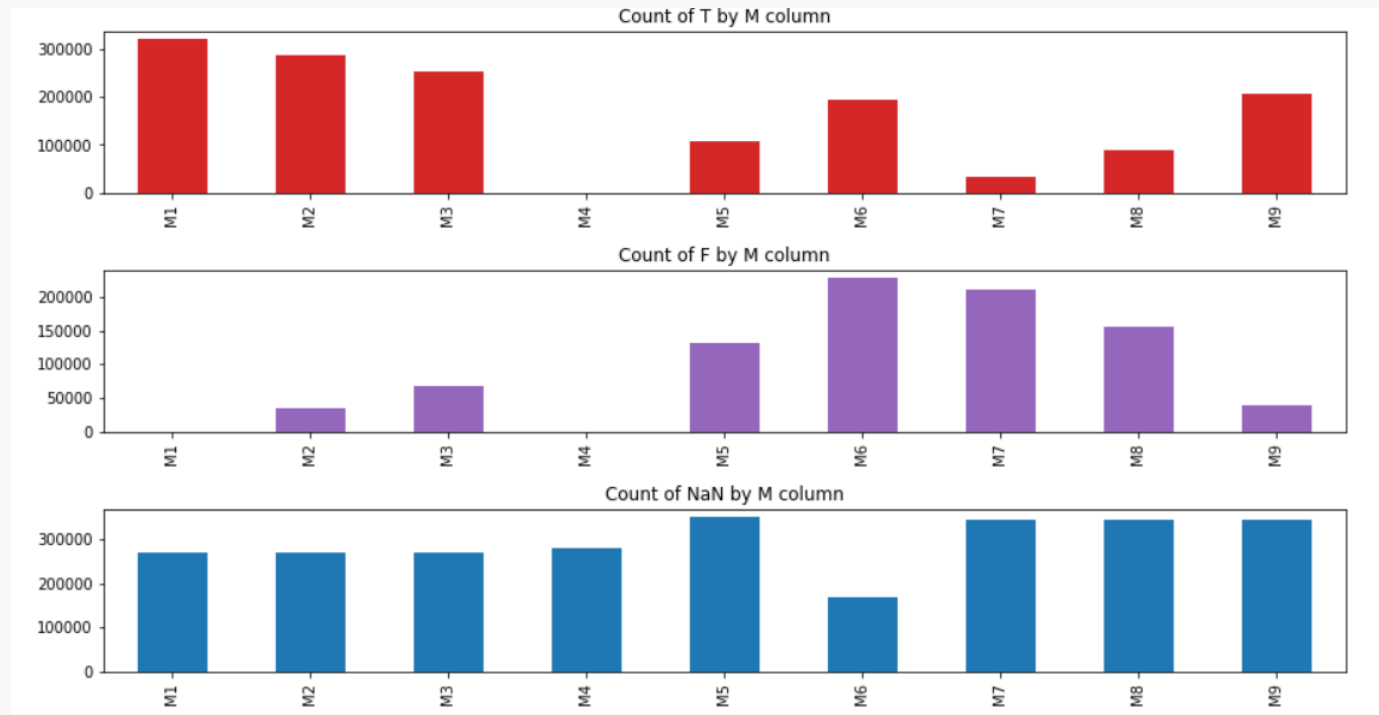
Feature Selection method

Correlation Analysis

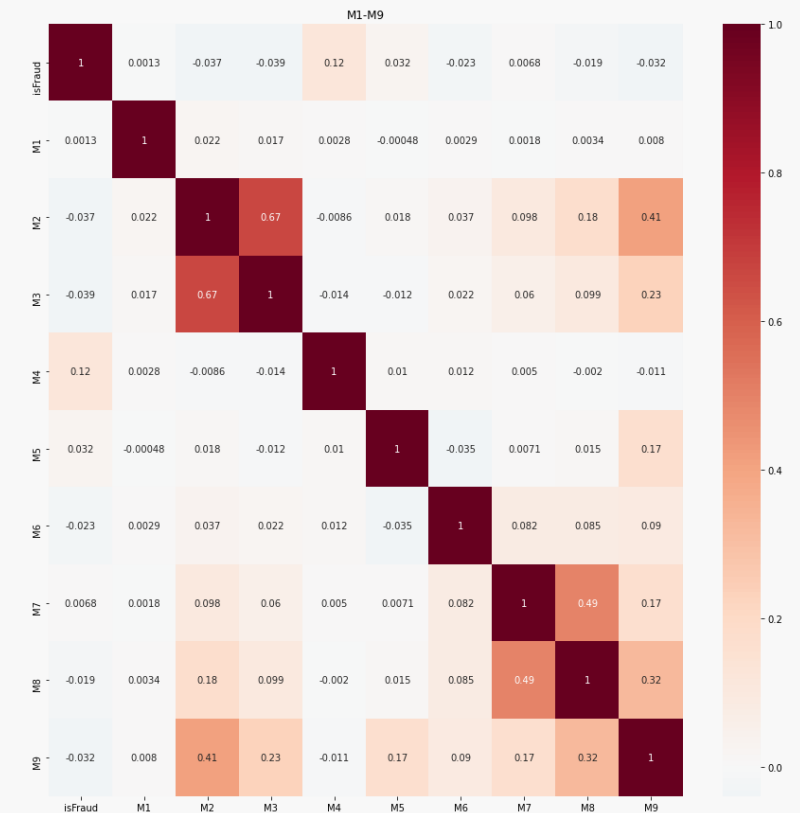
V Feature Selection

- V features have too many columns
- Relationship between features are vague
- Use RFE with cross validation to siphon out V features
- Preprocessing:
 - remove all the columns with high Nan values
- V features of RFE ranking of 1 are selected

M Feature Selection

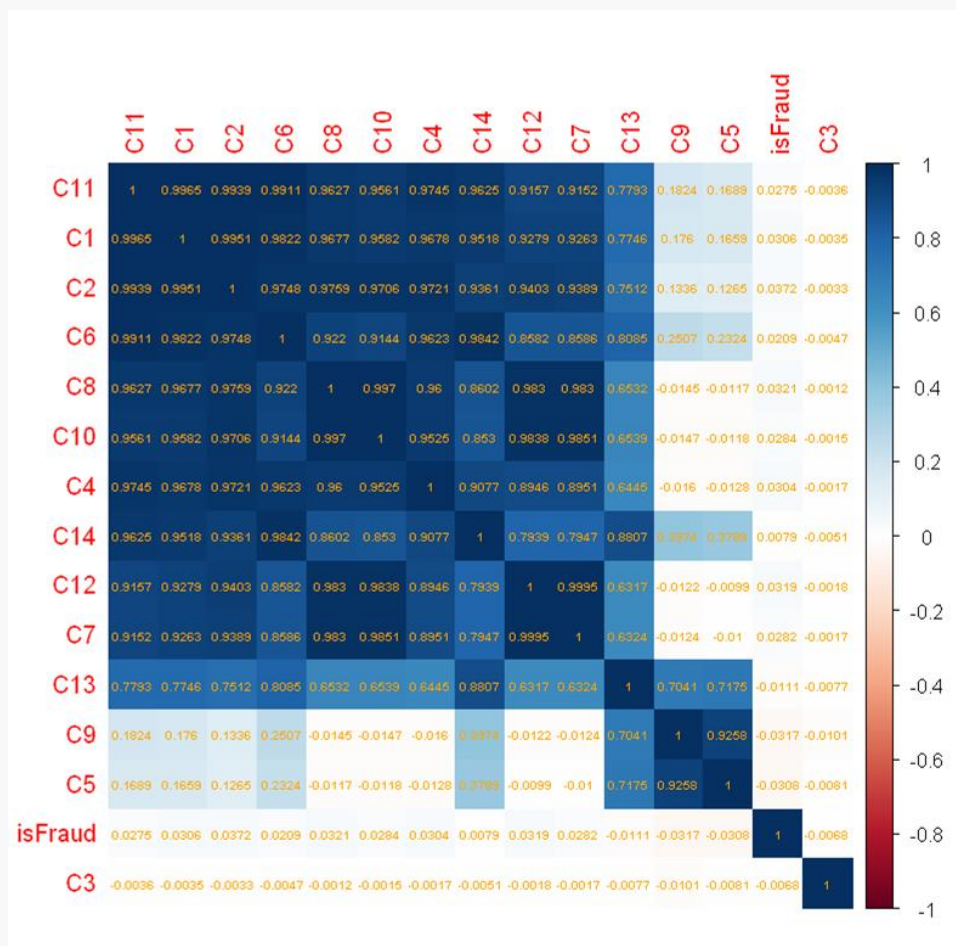


M feature bar chart



Correlation Heatmap for
M

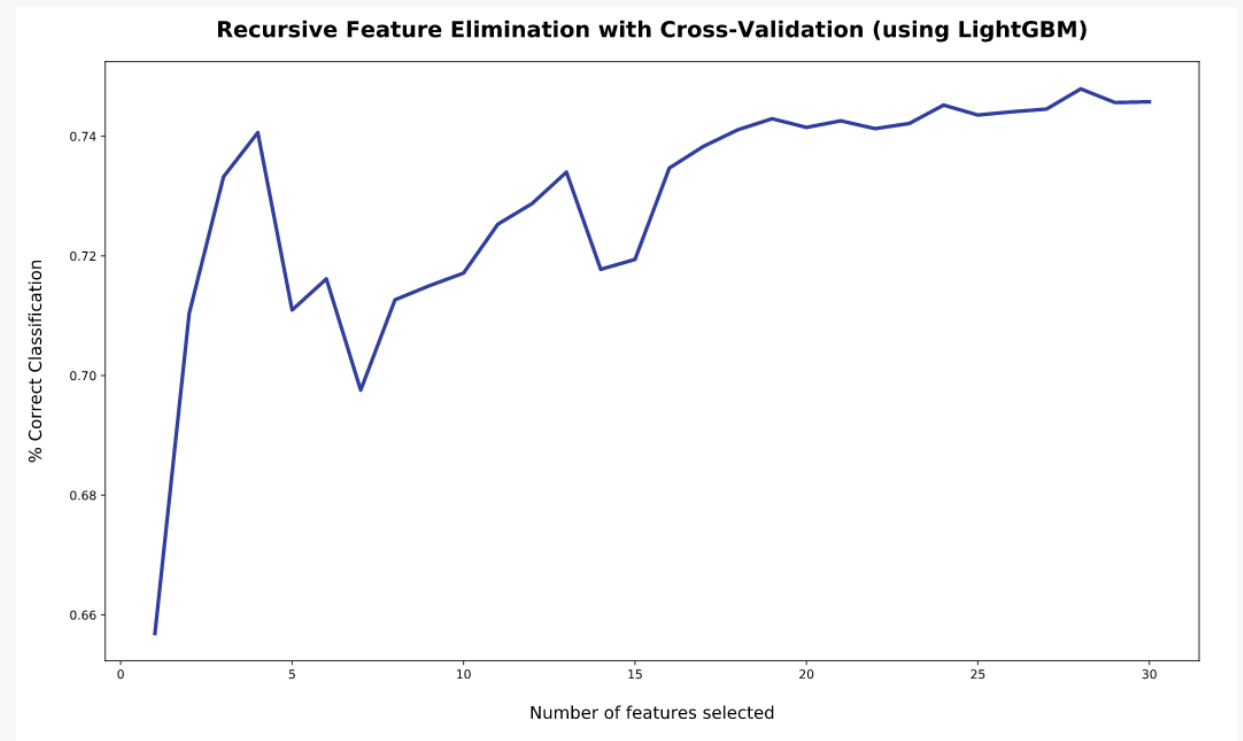
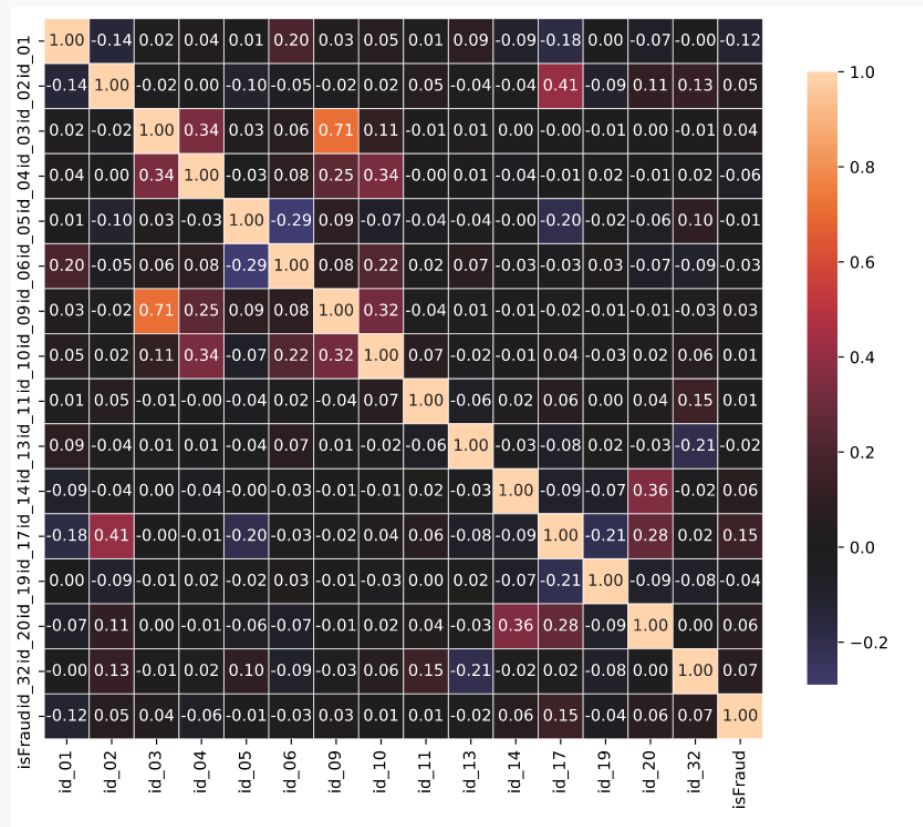
C Feature Selection



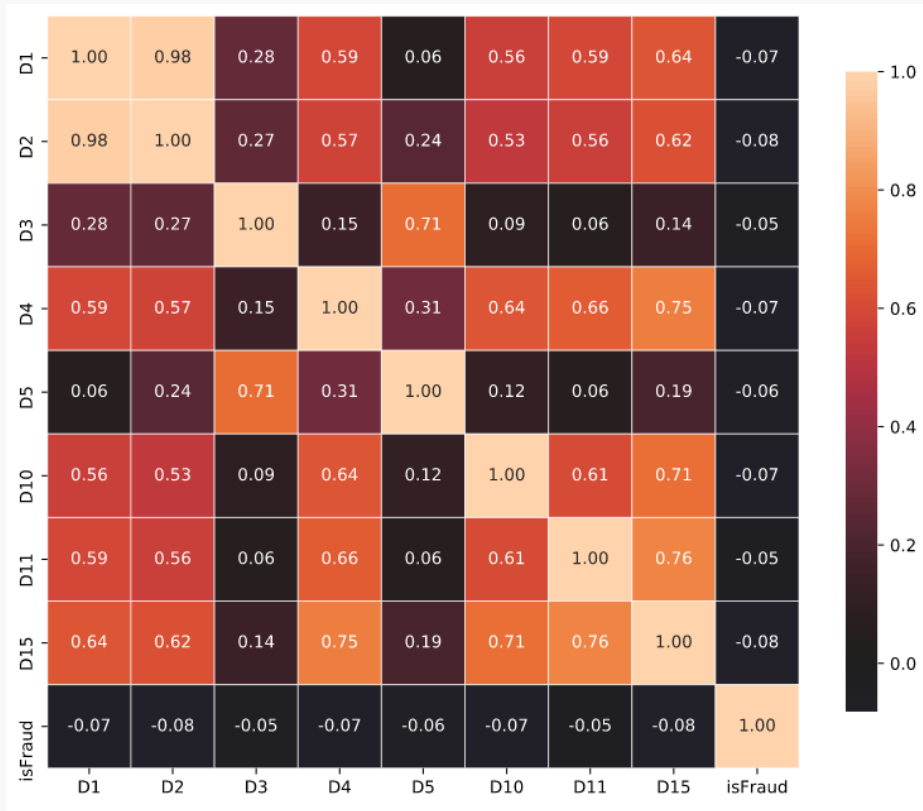
Heatmap for C features

- Just like V, C features are vague
- We can use correlation heatmap to determine relationship of each features

ID Feature Selection



D Feature Selection



Explanation for D features

- Just like V, C features are vague
- We can use correlation heatmap to determine relationship of each features

Pre-processing Features

Categorical Feature Representation

- Treat as categorical features
- Label Encoding
- Mixture of both

Normalisation of Numerical Data

Mean Normalisation
of numerical columns

Handling Missing Values

- Let LightGBM handle
- Replace with value outside feature range (-999)
- Mean Imputation
- Median Imputation
- Mode Imputation

Augment OS and Browser data

- Group OS Type
- Group Browser Type
- Group OS Version & Browser Type

Extract transaction day and time

- Extract Time
- Extract Day
- Extract Day and Time

Feature Aggregation

Use aggregated features

Selected Methods

Categorical Feature Representation

- Treat as categorical features
- Label Encoding
- Mixture of both

Normalisation of Numerical Data

Mean Normalisation
of numerical columns

Handling Missing Values

- Let LightGBM handle
- Replace with value outside feature range (-999)
- Mean Imputation
- Median Imputation
- Mode Imputation

Augment OS and Browser data

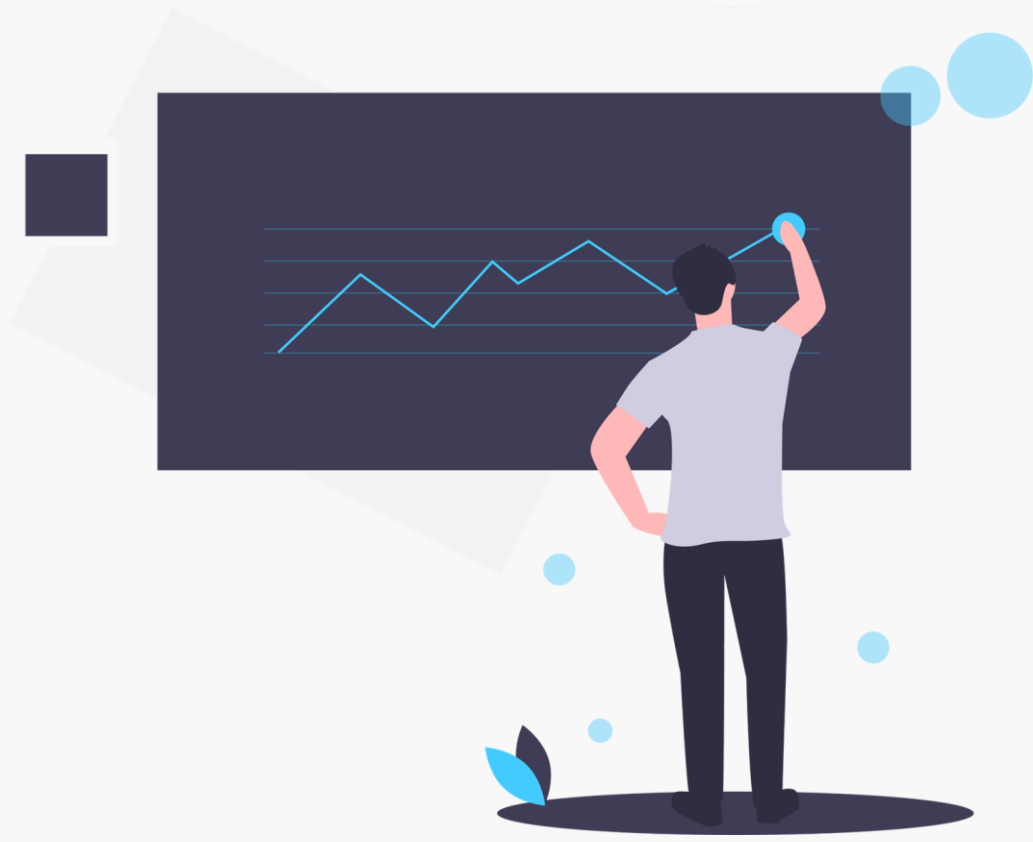
- Group OS Type
- Group Browser Type
- Group OS Version & Browser Type

Extract transaction day and time

- Extract Time
- Extract Day
- Extract Day and Time

Feature Aggregation

Use aggregated features



Modelling and Results

Models Used

Type / Model	LightGBM	XGBoost	Random Forest
<i>Validation Score (AUC)</i>	● 0.971931	● 0.975871	● 0.999903
<i>Kaggle Public Score</i>	● 0.939755	● 0.937507	● 0.898593
<i>Kaggle Private Score</i>	● 0.911989	● 0.908053	● 0.868237

Ensemble Learning

TODO: diagram

Kaggle Score and Rank

CATEGORY	PUBLIC LB SCORE	PRIVATE LB SCORE
<i>Score</i>	0.943850	0.917212
<i>Rank</i>	2744 / 6381	2451 / 6381
<i>Percentile</i>	~43	~38.41

Conclusion

Conclusion