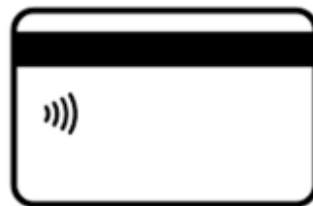# Credit Card Fraud Detection

Group 4: Jiayi Yang, Neal Xu, Suenkei Chan, Yao Zhang, Zhuofan Dong

# Agenda

Business Problem

Definitions of Variables & Data Exploration

Methodology & Analytics Approach

Results & Recommendation

# Business Problem

➢ **Background:** Fraudulent activity—**problem** for financial institutions that issue **credit cards**.

➢ **Problem: Insufficient detection** ➡ **loss of customers' confidence** in the security mechanisms

➡ **churn** of customers.

**Excessive alerts** ➡ **lack of confidence in** the detection algorithms as unreliable

➢ **Value of Solution: increase** customers' **stickiness**

➢ **Purpose and Objective:** develop **anomaly detection mechanisms** （ (Extreme Imbalanced Dataset): make the right fraud detection and reduce false fraud alerts.

# Exploratory Analysis

Default Correlation with other Variables

Feature Selection

| | |
|---|---|
| Class | 1.000000 |
| V17 | 0.326481 |
| V14 | 0.302544 |
| V12 | 0.260593 |
| V10 | 0.216883 |
| V16 | 0.196539 |
| V3 | 0.192961 |
| V7 | 0.187257 |
| V11 | 0.154876 |
| V4 | 0.133447 |
| V18 | 0.111485 |
| V1 | 0.101347 |
| V9 | 0.097733 |
| V5 | 0.094974 |
| V2 | 0.091289 |
| V6 | 0.043643 |
| V21 | 0.040413 |
| V19 | 0.034783 |
| V20 | 0.020090 |
| V8 | 0.019875 |
| V27 | 0.017580 |
| Time | 0.017082 |
| V28 | 0.009536 |
| V24 | 0.007221 |
| Amount | 0.005632 |
| V13 | 0.004570 |
| V26 | 0.004455 |
| V15 | 0.004223 |
| V25 | 0.003308 |
| V23 | 0.002685 |
| V22 | 0.000805 |



Feature Importances

Generally, the features revealed from Random Forest are more convincing as it takes linear and nonlinear relationship into consideration.

THE UNIVERSITY OF CHICAGO

# Models Performance

**No Resampling**

**Over-Sampling (SMOTE)**

**Weighted Learning**

Nonlinear:
  XGboosting
  Decision Tree
  Random Forest
  Neural Network
Linear:
  Logistic Regression

```
df['Class'].value_counts()

0    284315
1       492
Name: Class, dtype: int64
```

0.16%

CONTROLLING CLASS WEIGHTS FOR IMBALANCED DATASETS

All Leads to Overfitting:

```
Classification report:
              precision    recall  f1-score

           0       1.00      1.00      1.00
           1       1.00      1.00      1.00

    accuracy                           1.00
   macro avg       1.00      1.00      1.00
weighted avg       1.00      1.00      1.00
```

THE UNIVERSITY OF CHICAGO

# Model(for loop)

Advantages:

- Balanced data for both categories
- Train and test the models with different data in each loop without replacement
- Average performance score tells the truth of the model performance

Disadvantage:

- Size of training set may not be enough to fully train the models
- Hard to use the models with large computing resources such as deep learning

Undersampling

↓

Sample Without Replacement

↓

Train/Test the model

↓

Record the Performance

↓

Take the Average(LLN)

THE UNIVERSITY OF
CHICAGO

# Model(for loop)

Label 0: 250          Label1: 492

|  | accuracy | precision | recall | f1 |
|---|---|---|---|---|
| **Logistic Regression** | 0.930269 | 0.968517 | 0.925163 | 0.946048 |
| **Random Forest** | 0.931480 | 0.969500 | 0.926187 | 0.947122 |
| **XGBOOST** | 0.929507 | 0.968245 | 0.924734 | 0.945681 |
| **Decision Tree** | 0.932735 | 0.970212 | 0.926731 | 0.947818 |
| **Neural Network** | 0.931166 | 0.968773 | 0.926345 | 0.946941 |

Label 0: 492          Label1: 492

|  | accuracy | precision | recall | f1 |
|---|---|---|---|---|
| **Logistic Regression** | 0.937601 | 0.962064 | 0.911859 | 0.936071 |
| **Random Forest** | 0.937095 | 0.972221 | 0.899037 | 0.933930 |
| **XGBOOST** | 0.924764 | 0.938818 | 0.907956 | 0.922817 |
| **Decision Tree** | 0.917804 | 0.965918 | 0.867990 | 0.913875 |
| **Neural Network** | 0.919493 | 0.960431 | 0.875718 | 0.915931 |

|  | accuracy | precision | recall | f1 |
|---|---|---|---|---|
| **Logistic Regression** | 0.954509 | 0.978263 | 0.882575 | 0.927693 |
| **Random Forest** | 0.954308 | 0.977512 | 0.882115 | 0.927119 |
| **XGBOOST** | 0.953125 | 0.979982 | 0.876492 | 0.925142 |
| **Decision Tree** | 0.954665 | 0.977125 | 0.883342 | 0.927575 |
| **Neural Network** | 0.930762 | 0.968548 | 0.925877 | 0.946482 |

Label 0: 1000          Label1: 492

# Model in Semi-Supervised Learning

```
0     279580
1       5227
Name: cluster_labels, dtype: int64
```

```
Average cross-validation score for XG
Standard deviation for XGboosting : (
Accuracy score for XGboosting : 1.0
Classification report:
              precision    recall  f

           0       1.00      1.00
           1       1.00      1.00

    accuracy
   macro avg       1.00      1.00
weighted avg       1.00      1.00

Confusion matrix for XGboosting
[[ 1515      0]
 [    0 83928]]
```

```
Accuracy score for Logistic Regression : 0.9922404409957516
Classification report:
              precision    recall  f1-score   support

           0       0.70      1.00      0.82      1515
           1       1.00      0.99      1.00     83928

    accuracy                           0.99     85443
   macro avg       0.85      0.99      0.91     85443
weighted avg       0.99      0.99      0.99     85443

Confusion matrix for Logistic Regression
[[ 1510      5]
 [  658 83270]]
```

THE UNIVERSITY OF
CHICAGO

# Result Summary

- Useful in situation where fraud has a high cost
  <u>(We have a high recall around 0.98)</u>
- Approaches and models in supervised learning
- Approaches in semi-supervised learning
- Problematic

- Limited size of positive case; deeply imbalance dataset
- The dataset is continent-based
- The label is inaccurate

# Future Improvement

- Different Dataset

- Advanced Models to detect the hidden patterns (Deep Learning)

- Expected Loss

- Novelty Detection