# Praktikum 4.
# Anomaly Detection
## (Fraud Detection of Credit Card)
## Random Forest dan Logistic Regression

By: Rastri Prathivi, M.Kom.

```python
1 # import library
2 import warnings
3 warnings.filterwarnings('ignore')
4
5 import numpy as np
6 import pandas as pd
7 import seaborn as sns
8 import matplotlib.pyplot as plt
9 from imblearn.over_sampling import SMOTE
10 from sklearn.model_selection import train_test_split, GridSearchCV
11 from sklearn.preprocessing import StandardScaler
12 from sklearn.ensemble import RandomForestClassifier
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.ensemble import GradientBoostingClassifier
15 from sklearn.metrics import accuracy_score, roc_auc_score, plot_roc_curve, confusion_matrix, plot_confusion_matrix
16
```

UNTUK READ DATASET SILAKAN DISESUAIKAN SENDIRI CODINGNYA SEHINGGA DATA DAPAT DIBACA

```python
1 # Read the dataset
2 df = pd.read_csv("creditcard.csv")
3 df.head()
```

Out[2]:

| | Time | V1 | V2 | V3 | V4 | V5 | V6 | V7 | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | |
| 1 | 0.0 | 1.191857 | 0.266151 | 0.166480 | 0.448154 | 0.060018 | -0.082361 | -0.078803 | |
| 2 | 1.0 | -1.358354 | -1.340163 | 1.773209 | 0.379780 | -0.503198 | 1.800499 | 0.791461 | |
| 3 | 1.0 | -0.966272 | -0.185226 | 1.792993 | -0.863291 | -0.010309 | 1.247203 | 0.237609 | |
| 4 | 2.0 | -1.158233 | 0.877737 | 1.548718 | 0.403034 | -0.407193 | 0.095921 | 0.592941 | |

5 rows × 31 columns

```
1 # Printing quick information about the dataset
2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 284807 entries, 0 to 284806
Data columns (total 31 columns):
 #   Column  Non-Null Count   Dtype
---  ------  --------------   -----
 0   Time    284807 non-null  float64
 1   V1      284807 non-null  float64
 2   V2      284807 non-null  float64
 3   V3      284807 non-null  float64
 4   V4      284807 non-null  float64
 5   V5      284807 non-null  float64
 6   V6      284807 non-null  float64
 7   V7      284807 non-null  float64
 8   V8      284807 non-null  float64
 9   V9      284807 non-null  float64
 10  V10     284807 non-null  float64
```

```
1 # Checking missing values in each column
2 df.isnull().sum()
```

| Time | 0 |
|------|---|
| V1   | 0 |
| V2   | 0 |
| V3   | 0 |
| V4   | 0 |
| V5   | 0 |
| V6   | 0 |
| V7   | 0 |
| V8   | 0 |
| V9   | 0 |
| V10  | 0 |
| V11  | 0 |
| V12  | 0 |
| V13  | 0 |
| V14  | 0 |
| V15  | 0 |
| V16  | 0 |
| V17  | 0 |

```
1 # Identify duplicate values and mark all the duplicates as true
2 # https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.duplicated.html
3 df[df.duplicated(keep=False)]
```

|        | Time     | V1        | V2       | V3        | V4        | V5        | V6        | V7 |
|--------|----------|-----------|----------|-----------|-----------|-----------|-----------|----|
| 32     | 26.0     | -0.529912 | 0.873892 | 1.347247  | 0.145457  | 0.414209  | 0.100223  | 0.7 |
| 33     | 26.0     | -0.529912 | 0.873892 | 1.347247  | 0.145457  | 0.414209  | 0.100223  | 0.7 |
| 34     | 26.0     | -0.535388 | 0.865268 | 1.351076  | 0.147575  | 0.433680  | 0.086983  | 0.6 |
| 35     | 26.0     | -0.535388 | 0.865268 | 1.351076  | 0.147575  | 0.433680  | 0.086983  | 0.6 |
| 112    | 74.0     | 1.038370  | 0.127486 | 0.184456  | 1.109950  | 0.441699  | 0.945283  | -0. |
| ...    | ...      | ...       | ...      | ...       | ...       | ...       | ...       | ... |
| 283485 | 171627.0 | -1.457978 | 1.378203 | 0.811515  | -0.603760 | -0.711883 | -0.471672 | -0. |
| 284190 | 172233.0 | -2.667936 | 3.160505 | -3.355984 | 1.007845  | -0.377397 | -0.109730 | -0. |
| 284191 | 172233.0 | -2.667936 | 3.160505 | -3.355984 | 1.007845  | -0.377397 | -0.109730 | -0. |
| 284192 | 172233.0 | -2.691642 | 3.123168 | -3.339407 | 1.017018  | -0.293095 | -0.167054 | -0. |
| 284193 | 172233.0 | -2.691642 | 3.123168 | -3.339407 | 1.017018  | -0.293095 | -0.167054 | -0. |

1854 rows × 31 columns

```
[ ]    1 # https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.drop_duplicates.html
       2 # drop data duplicated
       3 df = df.drop_duplicates(keep='first')
```

```
[ ]    1 # Check the distribution of the credit card fraud cases
       2 class_proportion = df['Class'].value_counts()
       3 class_proportion
```
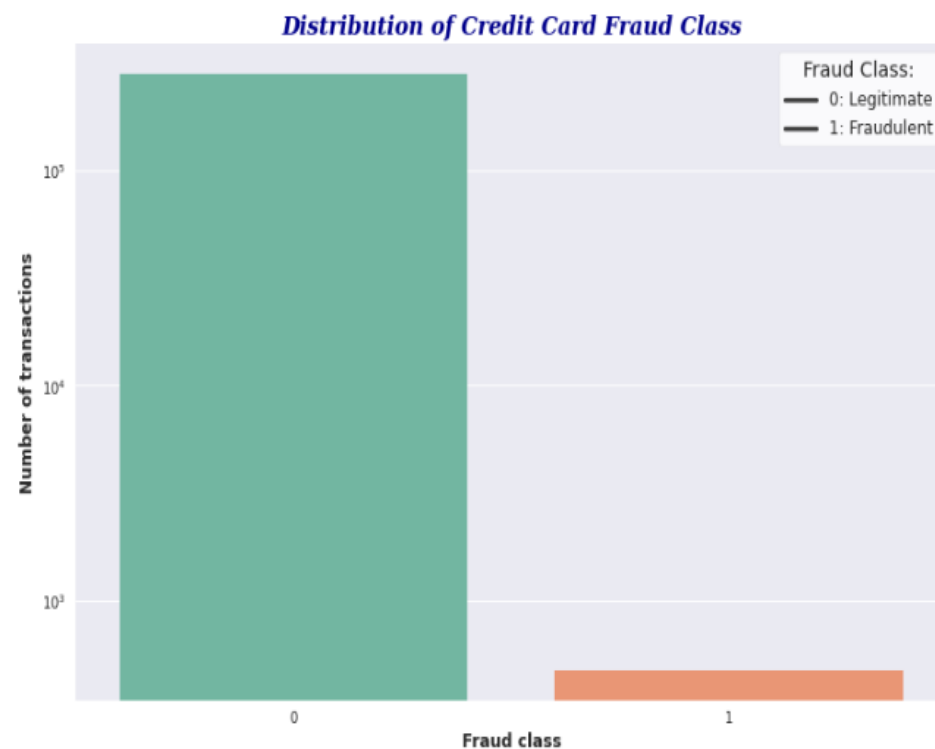
Out[7]:
```
        0    283253
        1       473
        Name: Class, dtype: int64
```

```python
1 # Plotting a barchart to see the the distribution of the credit card fraud cases
2 plt.style.use('seaborn')
3 font1 = {'family': 'serif',
4          'fontstyle': 'italic',
5          'fontsize': 16,
6          'fontweight': 'bold',
7          'color': 'DarkBlue'}
8 font2 = {'weight': 'bold', 'size': 12}
9 font3 = {'weight': 'normal', 'size': 12}
10
11 fig, ax = plt.subplots(figsize=(12, 8))
12 sns.barplot(class_proportion.index, class_proportion.values, palette='Set2')
13 ax.set_title('Distribution of Credit Card Fraud Class', fontdict=font1)
14 ax.set_xlabel('Fraud class', fontdict=font2)
15 ax.set_xticklabels(ax.get_xticklabels(), rotation=0)
16 ax.set_ylabel('Number of transactions', fontdict=font2)
17 ax.set_yscale('log')
18 handles, labels = ax.get_legend_handles_labels()
19 ax.legend(handles, labels=['0: Legitimate', '1: Fraudulent'], prop= font3,
20          title ='Fraud Class:', title_fontsize=14,
21          frameon=True, facecolor='white')
22 plt.show()
```

## Feature Engineering and Data Modeling

```python
1 # Check the proportion of the fraud cases and identify the imbalance
2 df['Class'].value_counts(normalize=True)
```
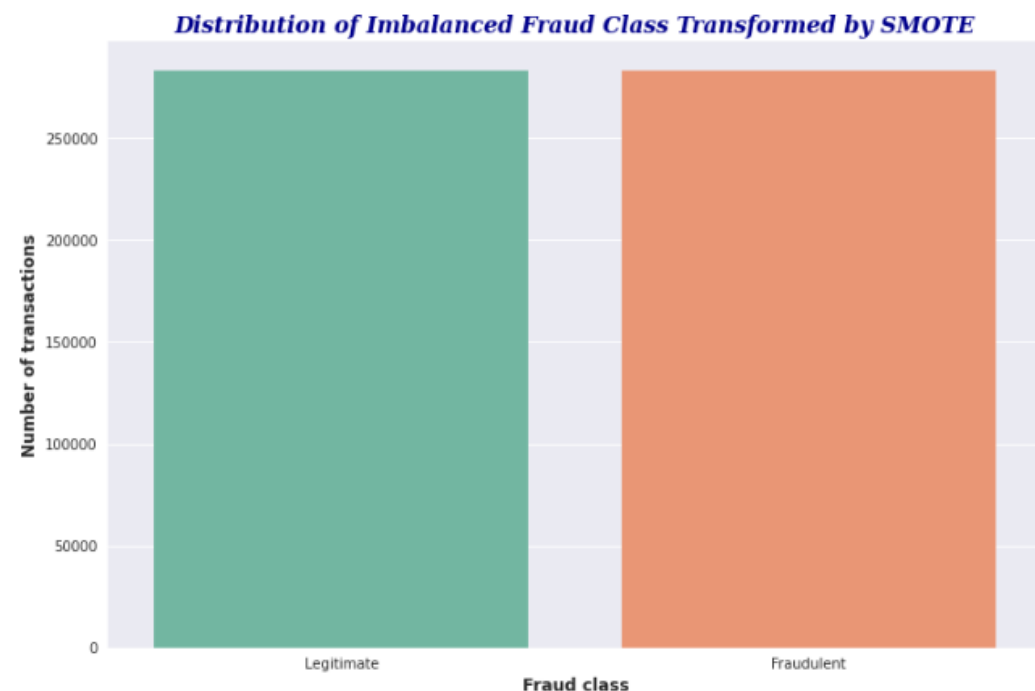
```python
1 # Arrange the dataset into features matrix and target vector
2 # Drop the 'Time' variable as it does not that much help our analysis
3 X = df.drop(columns=(['Time', 'Class']))
4 y = df['Class']
```

```python
1 # Make a SMOTE instance, then fit and apply it in one step
2 # to create an oversampled version of our dataset.
3
4 sm = SMOTE(sampling_strategy='auto', random_state=3, k_neighbors=5)
5 X_oversampled , y_oversampled = sm.fit_resample(X, y)
```

```python
1 # Summarize the fraud class distribution of the new SMOTE-transformed dataset
2 unique_original, counts_original = np.unique(y, return_counts=True)
3 unique_oversampled, counts_oversampled = np.unique(y_oversampled, return_counts=True)
4
5 print('Original fraud class distribution:', dict(zip(unique_original, counts_original)))
6 print('New transformed fraud class distribution:',dict(zip(unique_oversampled, counts_oversampled)))
```

```python
# Visualize the SMOTE-transformed target variable
plt.style.use('seaborn')
font1 = {'family': 'serif',
        'fontstyle': 'italic',
        'fontsize': 16,
        'fontweight': 'bold',
        'color': 'DarkBlue'}
font2 = {'weight': 'bold', 'size': 12}

fig, ax = plt.subplots(figsize=(12, 8))
sns.countplot(y_oversampled, palette='Set2', ax=ax)
ax.set_title('Distribution of Imbalanced Fraud Class Transformed by SMOTE', fontdict=font1)
ax.set_xlabel('Fraud class', fontdict=font2)
ax.set_xticklabels(['Legitimate', 'Fraudulent'])
ax.set_ylabel('Number of transactions', fontdict=font2)
plt.show()
```

## Random Forest Classifier(RFC)

```
[ ]    1 # Separate the transformed features matrix and target vector into random train and test subsets
       2 X_train, X_test, y_train, y_test = train_test_split(X_oversampled, y_oversampled, random_state=3)
```

```
[ ]    1 # Instantiate and fit the model
       2 rfc = RandomForestClassifier(n_estimators=150)
       3 rfc.fit(X_train, y_train)
```

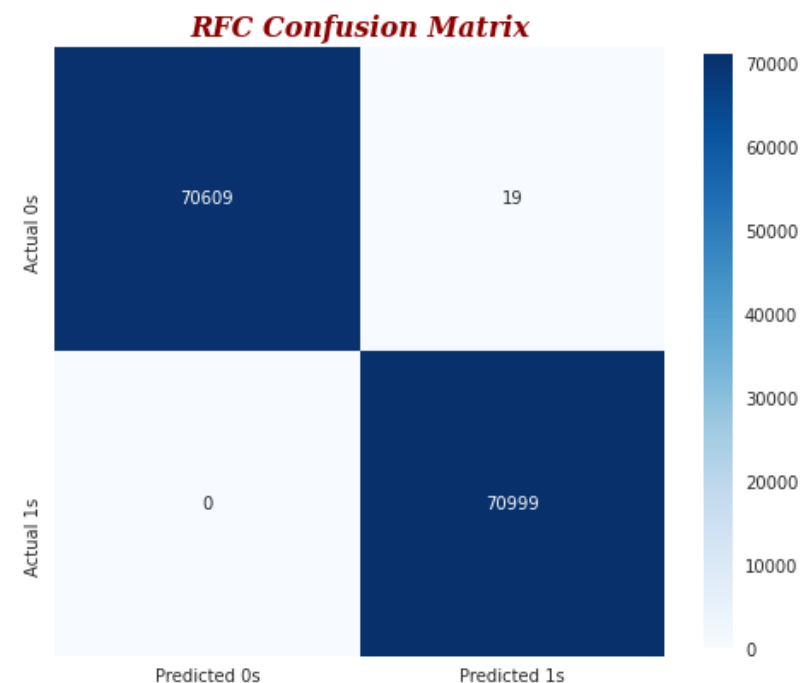RandomForestClassifier(n_estimators=150)

```
[ ]    1 # Model Evalution -classification accuracy
       2 training_rfc_accuracy = rfc.score(X_train, y_train)
       3 testing_rfc_accuracy = rfc.score(X_test, y_test)
       4
       5 print("Training RFC Accuracy:", training_rfc_accuracy)
       6 print("Testing RFC Accuracy:", testing_rfc_accuracy )
```

```
Training RFC Accuracy: 1.0
Testing RFC Accuracy: 0.9998658447894823
```

```python
# Plotting the confusion matrix
fig, ax = plt.subplots(figsize=(8, 8))
font1 = {'family': 'serif',
         'fontstyle': 'italic',
         'fontsize': 16,
         'fontweight': 'bold',
         'color': 'DarkRed'}
font2 = {'weight': 'bold', 'size': 12}

sns.heatmap(confusion_matrix(y_test, rfc.predict(X_test)),
            cmap='Blues',
            square=True,
            annot=True,
            fmt='d',
            cbar_kws={'shrink': 0.8},
            xticklabels=['Predicted 0s', 'Predicted 1s'],
            yticklabels=['Actual 0s', 'Actual 1s'])
ax.set_title('RFC Confusion Matrix', fontdict=font1)
plt.show()
```

```python
# Model evaluation - Sensitivity, Specificity and Precision

TN, FP, FN, TP = confusion_matrix(y_test, rfc.predict(X_test)).flatten()
print("True Negatives:", TN)
print("False Positives:", FP)
print("False Negatives:", FN)
print("True Positives:", TP)

sensitivity = TP/(TP + FN)
specificity = TN/(TN + FP)
precision = TP/(TP + FP)
print("\nSensitivity:", sensitivity)
print("Specificity:", specificity)
print("Precision:", precision)
```

```
True Negatives: 70609
False Positives: 19
False Negatives: 0
True Positives: 70999

Sensitivity: 1.0
Specificity: 0.999730984878184
Precision: 0.9997324621926835
```

```python
1 # Check the predicted probabilities for every observation in the test data subset
2 # Note that the default classification threshold is 0.5
3
4 testing_probabilities= rfc.predict_proba(X_test)
5 testing_probabilities
```

```python
1 # Convert the testing probabilities into a dataframe
2 testing_probabilities_df = pd.DataFrame(testing_probabilities, columns=['1 - p(X_test)', 'p(X_test)'])
3 testing_probabilities_df.head()
```

```python
1 # Get predictions
2 rfc.predict(X_test)
```

Out[20]:

|   | 1 - p(X_test) | p(X_test) |
|---|---------------|-----------|
| 0 | 0.0 | 1.0 |
| 1 | 0.0 | 1.0 |
| 2 | 0.0 | 1.0 |
| 3 | 1.0 | 0.0 |
| 4 | 0.0 | 1.0 |

```python
1 # Model evaluation -AUC
2 # Calculate AUC for both training and testing subsets
3 # Only probabilities being in the positive class is needed for the calculation, that is the second column
4 training_rfc_AUC = roc_auc_score(y_train, rfc.predict_proba(X_train)[:, 1])
5 testing_rfc_AUC = roc_auc_score(y_test, rfc.predict_proba(X_test)[:, 1])
6
7 print("Training RFC AUC:", training_rfc_AUC)
8 print("Testing RFC AUC:", testing_rfc_AUC)
```

```
Training RFC AUC: 1.0
Testing RFC AUC: 0.9999915933086132
```

## Logist Regression (LGR) Model

```python
1 # Separate the transformed features matrix and target vector into random train and test subsets
2 X_train, X_test, y_train, y_test = train_test_split(X_oversampled, y_oversampled, random_state=3)
```

```python
1 # define dictionary of hyperparameters
2 params = {'penalty': ['l1', 'l2'],
3          'C': [0.0001, 0.001, 0.01, 10, 50, 100],
4          'class_weight': [None, 'balanced']}
```

```python
1 # Instantiate Logistic Regression model. N.B: the default solver doesn't support l1 regularization
2 # Instantiate Grid Search to find the best hyperparameters and fit the model
3 lgr = LogisticRegression(solver='liblinear')
4 gs = GridSearchCV(lgr, params, cv = 5)
5 gs.fit(X_train, y_train)
```
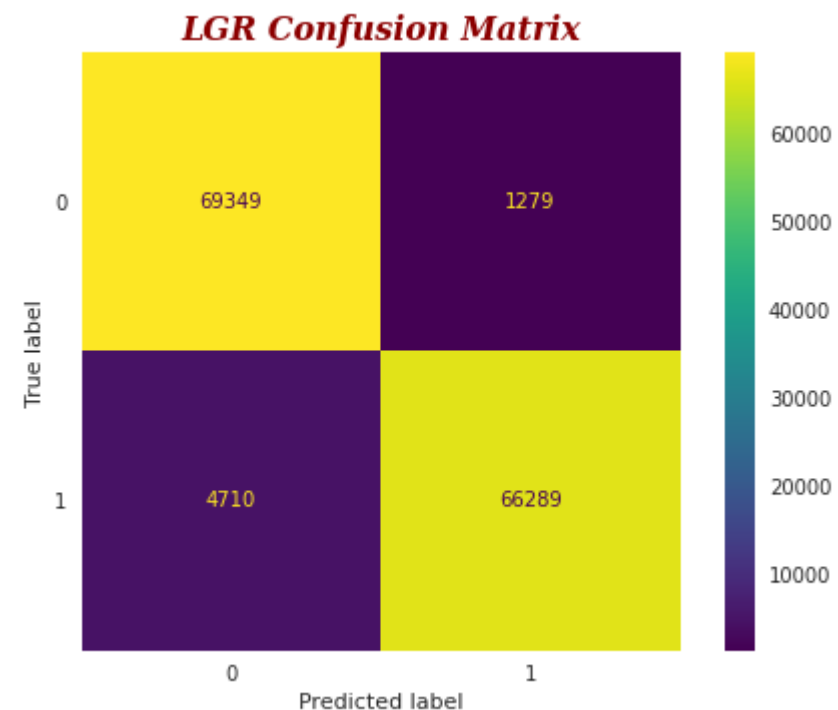
```python
1 # Model evaluation - accuracy
2 training_lgr_accuracy = gs.score(X_train, y_train)
3 testing_lgr_accuracy = gs.score(X_test, y_test)
4
5 print("Training LGR Accuracy:", training_lgr_accuracy)
6 print("Testing LGR Accuracy:", testing_lgr_accuracy)
```

```
Training LGR Accuracy: 0.958647050101323
Testing LGR Accuracy: 0.9577128654846887
```

```python
1 # Plotting the confusion matrix
2 font1 = {'family': 'serif',
3          'fontstyle': 'italic',
4          'fontsize': 16,
5          'fontweight': 'bold',
6          'color': 'DarkRed'}
7
8 plot_confusion_matrix(gs, X_test, y_test, values_format='d')
9 plt.title('LGR Confusion Matrix', fontdict=font1)
10 plt.grid(False)
11 plt.show()
```

```python
1 # Model evaluation - Sensitivity, Specificity and Precision
2 TN, FP, FN, TP = confusion_matrix(y_test, gs.predict(X_test)).flatten()
3 print("True Negatives:", TN)
4 print("False Positives:", FP)
5 print("False Negatives:", FN)
6 print("True Positives:", TP)
7
8 sensitivity = TP/(TP + FN)
9 specificity = TN/(TN + FP)
10 precision = TP/(TP + FP)
11 print("\nSensitivity:", sensitivity)
12 print("Specificity:", specificity)
13 print("Precision:", precision)
```

```
True Negatives: 69349
False Positives: 1279
False Negatives: 4710
True Positives: 66289

Sensitivity: 0.933661037940011
Specificity: 0.981891034717094
Precision: 0.981070921146046
```

```python
1  # Model evaluation -AUC
2  # Calculate AUC for both training and testing subsets
3  # Only probabilities being in the positive class is needed for the calculation, that is the second column
4  training_lgr_AUC = roc_auc_score(y_train, gs.predict_proba(X_train)[:, 1])
5  testing_lgr_AUC = roc_auc_score(y_test, gs.predict_proba(X_test)[:, 1])
6
7  print("Training LGR AUC:", training_lgr_AUC)
8  print("Testing LGR AUC:", testing_lgr_AUC)
```

Training LGR AUC: 0.991880379650 2905
Testing LGR AUC: 0.9915963855793377