## Report: Task 10.2HD

- Thus, the goal of this work is to enhance the model by evaluating and utilizing a variety of machine learning techniques on an unbalanced credit card dataset. These include of methods like features extraction, majority class undersampling, and minority class oversampling. Preliminary boosting methods including LightGBM, CatBoost, and XGBoost are included for further model improvements. The main evaluation metric in this instance was the same Area Under the Receiver Operating Characteristic Curve, or AUC-ROC, which is frequently used in classification issues, especially in imbalanced data sets like the one in this instance. Other measures were also calculated for information and classification pattern analysis.

- ✓ Since the first step in this approach was to undertake exploratory data analysis, or EDA, I started by attempting to comprehend the data as it is with the ultimate goal of identifying different factors that might affect the model's performance. I used undersampling and oversampling techniques because of a serious class imbalance issue, in which the number of fraudulent transactions is far lower than the total amount of data. In order to compare with other models, the baseline model was initially assessed using the Random Forest classifier.

- ✓ Although we might have thought about utilizing Principal Component Analysis (PCA), this was immediately rejected as the Therefore, there wasn't much to gain from dimensionality reduction in the data set. Rather, my emphasis was on employing SMOTE (Synthetic Minority Over-sampling Technique) to equalize the data set and identify instances of fraud related to the minority class.

- ✓ I created a synthetic data set with 1128 samples, 10 attributes, a trained Random Forest classifier, and a goal variable: binary classes for the experiment. The precision of the model was not as significant as the correctness of the probabilistic forecasts. At first, I saw that even if the model has a high classification accuracy, there are certain problems with non-varying probability estimates. Thus, I used the sigmoid and isotonic calibration procedures in an effort to improve the probability calibration.

- ✓ To be more specific, I divided the provided dataset into test, calibration, and training sets in order to carry out the problem-solving procedure. I then followed the training set's instructions to train the Random Forest model, using the calibration set to modify the probability estimates. After training the models, I calibrated the models using two calibrations techniques utilizing classifiers that are pre-fitted in scikit-learn's CalibratedClassifierCV class. The decision was made to proceed with probability calibration instead of repeating the models that produced these probabilities because probability calibration was the primary focus.

**The Output screenshot of 10.2HD ipynb Task**

## 1.2 Importing and fixing the csv file

Based off 10.1D, let's normalise the Amount column and drop the Time column before continuing.

```
[3]: # Load the dataset
     df = pd.read_csv('creditcard.csv')

     # Rescale the 'Amount' column using StandardScaler and create a new column 'normAmount'
     df['normAmount'] = StandardScaler().fit_transform(df['Amount'].values.reshape(-1, 1))

     # Drop the 'Time' and original 'Amount' columns as they are not required for further analysis
     df = df.drop(['Time', 'Amount'], axis=1)

     # Display the first few rows to ensure the changes are applied correctly
     df.head()
```

| | V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 | V9 | V10 | ... | V21 | V22 | V23 | V24 | V25 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 | 0.363787 | 0.090794 | ... | -0.018307 | 0.277838 | -0.110474 | 0.066928 | 0.128539 |
| 1 | 1.191857 | 0.266151 | 0.166480 | 0.448154 | 0.060018 | -0.082361 | -0.078803 | 0.085102 | -0.255425 | -0.166974 | ... | -0.225775 | -0.638672 | 0.101288 | -0.339846 | 0.167170 |
| 2 | -1.358354 | -1.340163 | 1.773209 | 0.379780 | -0.503198 | 1.800499 | 0.791461 | 0.247676 | -1.514654 | 0.207643 | ... | 0.247998 | 0.771679 | 0.909412 | -0.689281 | -0.327642 |
| 3 | -0.966272 | -0.185226 | 1.792993 | -0.863291 | -0.010309 | 1.247203 | 0.237609 | 0.377436 | -1.387024 | -0.054952 | ... | -0.108300 | 0.005274 | -0.190321 | -1.175575 | 0.647376 |
| 4 | -1.158233 | 0.877737 | 1.548718 | 0.403034 | -0.407193 | 0.095921 | 0.592941 | -0.270533 | 0.817739 | 0.753074 | ... | -0.009431 | 0.798278 | -0.137458 | 0.141267 | -0.206010 |

5 rows × 30 columns

---

Jupyter 10.2HD sample Last Checkpoint: 1 hour ago

File   Edit   View   Run   Kernel   Settings   Help                                    Not Trusted

🖫  +  ✂  📋  📋  ▶  ■  C  ▶▶  Code  ∨                          JupyterLab 📄  ⊙  Python 3 (ipykernel) (

**Requirements:** Split the dataset as suggested (e.g. at ratio 0.3) with random_state set. Statistical info about the split datasets expected, such as shape, sample count.

## 1.3 Data

Let's split the dataset into training and testing sets. Also check their shapes.

```
[5]: import pandas as pd
     from sklearn.model_selection import train_test_split

     # Assuming df is your DataFrame already loaded with data

     # Get data points and labels
     X = df.iloc[:, df.columns != 'Class']
     y = df.iloc[:, df.columns == 'Class'].values.ravel()  # Flatten y if it's a 2D array

     # Split into training and testing sets
     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.1, random_state=0)

     # Reset indices to avoid issues with non-existent indices
     X_train = X_train.reset_index(drop=True)
     y_train = pd.Series(y_train).reset_index(drop=True)  # Convert to Series for consistency
     X_test = X_test.reset_index(drop=True)
     y_test = pd.Series(y_test).reset_index(drop=True)  # Convert to Series for consistency

     # View shapes
     print(f'X_train\'s shape : {X_train.shape}')
     print(f'y_train\'s shape : {y_train.shape}')
     print(f'X_test\'s shape  : {X_test.shape}')
     print(f'y_test\'s shape  : {y_test.shape}')

     # Statistical info
     print("\nStatistical info of training set:")
     print(X_train.describe())
     print("\nStatistical info of testing set:")
     print(X_test.describe())
```

```
X_train's shape : (199364, 29)
y_train's shape : (199364,)
X_test's shape  : (85443, 29)
y_test's shape  : (85443,)
```

Jupyter 10.2HD sample Last Checkpoint: 1 hour ago

File Edit View Run Kernel Settings Help

JupyterLab   Python 3 (ipykernel)

Not Trusted

|       |            |            |            |            |
|-------|------------|------------|------------|------------|
| std   | 1.963554   | 1.657379   | 1.516716   | 1.417138   |
| min   | -46.855047 | -63.344608 | -33.680984 | -5.560118  |
| 25%   | -0.921539  | -0.601213  | -0.892838  | -0.848835  |
| 50%   | 0.019705   | 0.063784   | 0.177888   | -0.017852  |
| 75%   | 1.316707   | 0.882437   | 1.025529   | 0.745566   |
| max   | 2.451888   | 22.957729  | 9.382558   | 16.715537  |

|       | V5            | V6            | V7            | V8 \          |
|-------|---------------|---------------|---------------|---------------|
| count | 199364.000000 | 199364.000000 | 199364.000000 | 199364.000000 |
| mean  | -0.001494     | -0.000218     | -0.000870     | -0.001080     |
| std   | 1.368744      | 1.328673      | 1.226018      | 1.212338      |
| min   | -42.147898    | -23.496714    | -43.557242    | -73.216738    |
| 25%   | -0.602874     | -0.769177     | -0.554220     | -0.209986     |
| 50%   | -0.055832     | -0.274307     | 0.039228      | 0.021881      |
| 75%   | 0.609149      | 0.397928      | 0.569638      | 0.327023      |
| max   | 34.099309     | 23.917837     | 44.054461     | 20.007208     |

|       | V9            | V10           | ...  | V20           | V21 \         |
|-------|---------------|---------------|------|---------------|---------------|
| count | 199364.000000 | 199364.000000 | ...  | 199364.000000 | 199364.000000 |
| mean  | 0.000212      | 0.001357      | ...  | 0.000430      | -0.000014     |
| std   | 1.102021      | 1.092801      | ...  | 0.770257      | 0.743450      |
| min   | -13.434066    | -24.588262    | ...  | -23.646898    | -34.830382    |
| 25%   | -0.644753     | -0.535403     | ...  | -0.211662     | -0.229272     |
| 50%   | -0.049633     | -0.092060     | ...  | -0.062889     | -0.029045     |
| 75%   | 0.507098      | 0.458129      | ...  | 0.132834      | 0.187095      |
| max   | 15.504095     | 23.745136     | ...  | 39.420904     | 27.202839     |

|       | V22           | V23           | V24           | V25 \         |
|-------|---------------|---------------|---------------|---------------|
| count | 199364.000000 | 199364.000000 | 199364.000000 | 199364.000000 |
| mean  | -0.000022     | -0.000258     | 0.000362      | 0.000395      |
| std   | 0.727625      | 0.629145      | 0.605208      | 0.521175      |
| min   | -10.933144    | -44.807735    | -2.822684     | -10.295397    |
| 25%   | -0.544345     | -0.162021     | -0.354179     | -0.316088     |
| 50%   | 0.006744      | -0.010015     | 0.040074      | 0.018014      |
| 75%   | 0.531017      | 0.147583      | 0.438953      | 0.350802      |
| max   | 10.503090     | 22.528412     | 4.022866      | 7.519589      |

|       | V26           | V27           | V28           | normAmount    |
|-------|---------------|---------------|---------------|---------------|
| count | 199364.000000 | 199364.000000 | 199364.000000 | 199364.000000 |
| mean  | -0.000094     | -0.000027     | 0.000015      | 0.001271      |
| std   | 0.481842      | 0.401042      | 0.324848      | 0.983948      |
| min   | -2.534330     | -22.565679    | -11.710896    | -0.353229     |
| 25%   | -0.327127     | -0.070864     | -0.052907     | -0.330640     |
| 50%   | -0.052287     | 0.001066      | 0.011110      | -0.265271     |
| 75%   | 0.241092      | 0.090491      | 0.077989      | -0.043058     |
| max   | 3.463246      | 12.152481     | 22.620072     | 78.235272     |

[8 rows x 29 columns]

Statistical info of testing set:

|       | V1           | V2           | V3           | V4           | V5 \         |
|-------|--------------|--------------|--------------|--------------|--------------|
| count | 85443.000000 | 85443.000000 | 85443.000000 | 85443.000000 | 85443.000000 |
| mean  | -0.000734    | 0.006277     | 0.003574     | -0.001682    | 0.003486     |
| std   | 1.947325     | 1.637050     | 1.515182     | 1.412908     | 1.606722     |
| min   | -56.407510   | -72.715728   | -48.325580   | -5.683171    | -113.743307  |
| 25%   | -0.916858    | -0.591858    | -0.883828    | -0.848202    | -0.688288    |
| 50%   | 0.013238     | 0.070185     | 0.185047     | -0.024189    | -0.051627    |
| 75%   | 1.313257     | 0.806615     | 1.031155     | 0.737784     | 0.618067     |
| max   | 2.454330     | 15.876023    | 4.079168     | 16.875344    | 34.801666    |

Jupyter 10.2HD sample Last Checkpoint: 1 hour ago

File  Edit  View  Run  Kernel  Settings  Help

Not Trusted

JupyterLab ⬚ ⊙ Python 3 (ipykernel)

Code ∨

|        |           |            |            |            |             |   |
|--------|-----------|------------|------------|------------|-------------|---|
| count  | 85443.000000 | 85443.000000 | 85443.000000 | 85443.000000 | 85443.000000 |   |
| mean   | -0.000734 | 0.006277   | 0.003574   | -0.001682  | 0.003486    |   |
| std    | 1.947325  | 1.637050   | 1.515182   | 1.412908   | 1.406722    |   |
| min    | -56.407518 | -72.715728 | -48.325589 | -5.683171  | -113.743387 |   |
| 25%    | -0.916858 | -0.591858  | -0.883828  | -0.848282  | -0.688280   |   |
| 50%    | 0.013238  | 0.070185   | 0.185047   | -0.024189  | -0.051627   |   |
| 75%    | 1.313257  | 0.806615   | 1.031155   | 0.737784   | 0.618867    |   |
| max    | 2.454038  | 15.876923  | 4.079168   | 16.875344  | 34.801666   |   |

|        | V6         | V7         | V8         | V9         | V10        | \ |
|--------|-----------|------------|------------|------------|-------------|---|
| count  | 85443.000000 | 85443.000000 | 85443.000000 | 85443.000000 | 85443.000000 |   |
| mean   | 0.000489  | 0.002030   | 0.004620   | -0.000495  | -0.003167   |   |
| std    | 1.300636  | 1.262562   | 1.151291   | 1.090691   | 1.079574    |   |
| min    | -26.160506 | -28.215112 | -50.943369 | -9.481456  | -20.949192  |   |
| 25%    | -0.766664 | -0.553479  | -0.207216  | -0.638926  | -0.535400   |   |
| 50%    | -0.273666 | 0.042343   | 0.023782   | -0.053821  | -0.094949   |   |
| 75%    | 0.399864  | 0.572423   | 0.328337   | 0.507388   | 0.443126    |   |
| max    | 73.301628 | 120.589494 | 18.748872  | 9.272376   | 15.331742   |   |

|        | ...   | V20        | V21        | V22        | V23        | \ |
|--------|-------|------------|------------|------------|------------|---|
| count  | ...   | 85443.000000 | 85443.000000 | 85443.000000 | 85443.000000 |   |
| mean   | ...   | -0.001004  | 0.000033   | 0.000052   | 0.000602   |   |
| std    | ...   | 0.772484   | 0.713266   | 0.721198   | 0.613304   |   |
| min    | ...   | -54.497720 | -22.665685 | -9.409423  | -12.828095 |   |
| 25%    | ...   | -0.211881  | -0.226184  | -0.537784  | -0.161490  |   |
| 50%    | ...   | -0.061529  | -0.030687  | 0.000971   | -0.011789  |   |
| 75%    | ...   | 0.133600   | 0.184846   | 0.523600   | 0.147023   |   |
| max    | ...   | 38.117289  | 22.579714  | 7.220158   | 20.803344  |   |

|        | V24        | V25        | V26        | V27        | V28        | \ |
|--------|-----------|------------|------------|------------|------------|---|
| count  | 85443.000000 | 85443.000000 | 85443.000000 | 85443.000000 | 85443.000000 |   |
| mean   | -0.000845 | -0.000922  | 0.000220   | 0.000062   | -0.000036  |   |
| std    | 0.606464  | 0.521520   | 0.483126   | 0.409616   | 0.341987   |   |
| min    | -2.836827 | -8.696627  | -2.604551  | -9.793568  | -15.430084 |   |
| 25%    | -0.355671 | -0.319716  | -0.326068  | -0.070797  | -0.053129  |   |
| 50%    | 0.040976  | 0.013508   | -0.051695  | 0.001984   | 0.011561   |   |
| 75%    | 0.441093  | 0.350617   | 0.240657   | 0.092224   | 0.070900   |   |
| max    | 4.584548  | 5.826159   | 3.517346   | 31.612198  | 33.847808  |   |

|        | normAmount |
|--------|-----------|
| count  | 85443.000000 |
| mean   | -0.002966 |
| std    | 1.036492  |
| min    | -0.353129 |
| 25%    | -0.331288 |
| 50%    | -0.265271 |
| 75%    | -0.047356 |
| max    | 102.362243 |

[8 rows x 29 columns]

4

Jupyter 10.2HD sample Last Checkpoint: 1 hour ago

File Edit View Run Kernel Settings Help

Not Trusted

JupyterLab ⬚ ⚙ Python 3 (ipykernel)

## 1.4 Functions

Here we define a few functions that we will use later in the project.

```python
import numpy as np
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
from sklearn.calibration import CalibratedClassifierCV
from sklearn.metrics import roc_curve, roc_auc_score, brier_score_loss
from sklearn.model_selection import GridSearchCV


def apply_PCA(X, plot=True):
    ''' Fit PCA on data, then plot explained variance if needed
        Parameter
        ----------
            X : DataFrame that will be scaled and transformed according to PCA
            plot : Boolean. If True, displays the explained variance by principal components of PCA

        Output
        -------
            Array of principal components resulting from PCA
            Vector of explained variance
    '''

    scaler = StandardScaler()
    X_scaled = scaler.fit_transform(X)

    # By default keeps all components
    pca = PCA()
    pca_comp = pca.fit_transform(X_scaled)
    explained_var = pca.explained_variance_ratio_

    if plot:
        fig, ax = plt.subplots(1, 1, figsize=(10, 5))
        ax.plot(range(0, len(explained_var)), explained_var.cumsum(), marker='o')
        ax.set(xlabel='Principal components', ylabel='Inertia')
        plt.axvline(x=11, color='r', linestyle='--')
        plt.title('Amount of inertia explained by principal components', fontweight='bold')
        plt.show()

    return pca_comp, explained_var


def plot_pca_3D(pca_components, labels):
    ''' Plot values of principal components of PCA
        Parameters
        ----------
            pca_components : Array resulting from fitting PCA to data
            labels : ground truth labels from data
    '''

    not_fraud = pca_components[labels['Class'] == 0]
    fraud = pca_components[labels['Class'] == 1]

    plot3D = plt.figure().add_subplot(projection='3d')

    # First, second and third principal components in each x, y and z variable
    plot3D.scatter(fraud[:, 2], fraud[:, 0], fraud[:, 1], label='Fraudulent transactions', color='r', alpha=1)
    plot3D.scatter(not_fraud[:, 2], not_fraud[:, 0], not_fraud[:, 1], label='Not fraudulent', color='g', alpha=0.2)
```

File   Edit   View   Run   Kernel   Settings   Help

🔲  +  ✂  🗋  🗋  ▶  ■  C  ⏩  Code        ∨                                    JupyterLab ⍐   ⊚   Python 3 (i

```python
def plot_learning_auc(cv, model_name, ax):
    ''' Plot mean AUC learning curve based on cross validation result table (Gradient Boosting methods)
        Parameters
        -----------
            cv : Gradient Boosting cross validation output DataFrame
            model_name : String for plot title
            ax : Location of plot in the figure (ex: ax[0])

        Output
        -----------
            ax : Return ax to fill figure
    '''

    ax.plot(range(cv.shape[0]), cv['train-auc-mean'], label='Train AUC', color='b')
    ax.plot(range(cv.shape[0]), cv['test-auc-mean'], label='Valid AUC', color='g')
    ax.set_xlabel('Iterations')
    ax.set_ylabel('Mean AUC')
    ax.set_title(f'Learning curve - {model_name}', fontweight='bold')
    ax.legend()
    return ax

def plot_roc_curve_manual(y, y_pred, model_name, ax, title):
    ''' Plot ROC Curve based on model predictions, with focus on AUC metric
        Parameters
        -----------
            y : Array of ground truth values to predict
            y_pred : Array of predictions (i.e output of model)
            model_name : String for plot title
            ax : Location of plot in the figure (ex. ax[0])

        Output
        -----------
            ax : Return ax to fill figure
    '''

    fpr, tpr, _ = roc_curve(y, y_pred)
    ax.plot(fpr, tpr, label=f'AUC {model_name} : {roc_auc_score(y, y_pred):.5f}')
    ax.plot([0, 1], [0, 1], ls='--', color='black')
    ax.set_xlabel('False Positive Rate')
    ax.set_ylabel('True Positive Rate')
    ax.set_title(title, fontweight='bold')
    ax.legend(prop={'size': 11})
    return ax

def calibrate_predictions(classifier, X_train, y_train, X_test, y_test, sample_weight_train, sample_weight_test, cv, method):
    ''' Calibrate predictions in order to get more intuition about the probability of an item being an anomaly
        Parameters
        -----------
            classifier : Classifier model to calibrate
            X_train, y_train, X_test, y_test : Datasets for training and test (or validation)
            sample_weight_train, sample_weight_test : Weights of target value in the original dataset
            cv : Cross-validation set
            method : Method used to calibrate predictions : ['Isotonic', 'Sigmoid']

        Output
        -----------
            clf_calib : Fitted classifier
            preds_calib : Calibrated predictions with respect to the method (output is base on predict_proba() function)
            clf_brier_score : Brier Score of calibration
```

Jupyter  10.2HD sample Last Checkpoint 1 hour ago

File  Edit  View  Run  Kernel  Settings  Help

Not Trusted

+  ✂  □  □  ▶  ■  C  ▶▶  Code  ∨

JupyterLab  ◎  Python 3 (ipykernel)

```python
    # Predictions
    preds_calib = clf_calib.predict_proba(X_test)[:, 1]

    # Brier Score
    clf_brier_score = brier_score_loss(y_test, preds_calib, sample_weight=sample_weight_test.ravel())

    return clf_calib, preds_calib, clf_brier_score

def plot_predictions(y_pred, calibration_type, brier_score, ax):
    ''' Plot calibrated predictions
        Parameter
        ----------
            y_pred : Array of predictions (i.e output of model)
            calibration_type : Method used for calibration
            brier_score : Brier score (output of calibration)
            ax : Location of plot in the figure (ex. ax[0])

        Output
        ----------
            ax : Return ax to fill figure
    '''

    ax.plot(range(len(y_pred)), np.sort(y_pred), label=f'{calibration_type} ({brier_score:.4f})')
    ax.legend()
    ax.set_title('Model predictions by calibration type (Brier score)', fontweight='bold')
    ax.set_ylabel('$P(Y=1)$')
    ax.set_xlabel("Instances (ordered)")
    return ax

def clf_plot_predict(clf, params, model_name, title, x_train, y_train, x_test, y_test):
    ''' Train and plot predictions using the classifier and parameters provided
        Parameters
        ----------
            clf : Classifier to train
            params : Dictionary of parameters for GridSearchCV
            model_name : String for plot title
            title : String for figure title
            x_train, y_train : Training data
            x_test, y_test : Test data
    '''

    searchCV = GridSearchCV(
        estimator=clf,
        param_grid=params,
        scoring='roc_auc',
        cv=5,
        verbose=False,
        n_jobs=-1  # Utilise all CPU cores to reduce waiting times
    )

    searchCV.fit(x_train, y_train)
    print(f'Best AUC score (training) : {searchCV.best_score_}')
    print(f'Best params : {searchCV.best_params_}')

    # Predictions
    preds_cv = searchCV.best_estimator_.predict_proba(x_test)
    fig, ax = plt.subplots(1, 1, figsize=(7, 5))
    plot_roc_curve_manual(y_test, preds_cv[:, 1], model_name, ax, title)
    plt.show()
```

```python
[9]: import pandas as pd

     # Assuming the datasets are already loaded as X_train, X_test, y_train, y_test
     # Example:
     # X_train = pd.read_csv('X_train.csv')
     # X_test = pd.read_csv('X_test.csv')
     # y_train = pd.read_csv('y_train.csv')
     # y_test = pd.read_csv('y_test.csv')

     # Checking for Missing Values in the datasets
     print(f'Missing values in X_train : {X_train.isnull().sum().sum()}')
     print(f'Missing values in y_train : {y_train.isnull().sum().sum()}')
     print(f'Missing values in X_test  : {X_test.isnull().sum().sum()}')
     print(f'Missing values in y_test  : {y_test.isnull().sum().sum()}')

     Missing values in X_train : 0
     Missing values in y_train : 0
     Missing values in X_test  : 0
     Missing values in y_test  : 0
```

**Duplicated values**

Some values are duplicated in the dataset, and we will focus on them in the *Feature Engineering* part.

[11]:
```python
import pandas as pd

# Assuming the datasets are already loaded as X_train, X_test, y_train, y_test
# Example:
# X_train = pd.read_csv('X_train.csv')
# X_test = pd.read_csv('X_test.csv')
# y_train = pd.read_csv('y_train.csv')
# y_test = pd.read_csv('y_test.csv')

# Checking for Duplicate Values in the datasets
print(f'Duplicated values in X_train : {X_train.duplicated().sum()}')
print(f'Duplicated values in y_train : {y_train.duplicated().sum()}')
print(f'Duplicated values in X_test  : {X_test.duplicated().sum()}')
print(f'Duplicated values in y_test  : {y_test.duplicated().sum()}')
```

```
Duplicated values in X_train : 5161
Duplicated values in y_train : 199362
Duplicated values in X_test  : 1375
Duplicated values in y_test  : 85441
```

**Imbalanced classes**

We can see that the classes are severely imbalanced in the overall dataset, with fraudulent transactions accounting for only 492 samples in comparison to regular transactions at 284315 samples. This means fraudulent transactions only account for 0.173% out of the total transactions.

[13]:
```python
import pandas as pd

# Assuming y_train and y_test contain the target labels, where 1 represents the anomaly/fraud and 0 represents normal transactions

# Checking class distribution in y_train
class_distribution_train = y_train.value_counts(normalize=True)
class_distribution_test = y_test.value_counts(normalize=True)

# Printing the imbalance in percentages
print(f'Class distribution in y_train:\n{class_distribution_train * 100}\n')
print(f'Class distribution in y_test:\n{class_distribution_test * 100}\n')

# Checking total counts
fraud_train = y_train.value_counts()[1]
not_fraud_train = y_train.value_counts()[0]
fraud_test = y_test.value_counts()[1]
not_fraud_test = y_test.value_counts()[0]

print(f'Total samples in y_train: {len(y_train)}')
print(f'Fraudulent transactions in y_train: {fraud_train} ({(fraud_train / len(y_train)) * 100:.2f}%)')
print(f'Regular transactions in y_train: {not_fraud_train} ({(not_fraud_train / len(y_train)) * 100:.2f}%)\n')

print(f'Total samples in y_test: {len(y_test)}')
print(f'Fraudulent transactions in y_test: {fraud_test} ({(fraud_test / len(y_test)) * 100:.2f}%)')
print(f'Regular transactions in y_test: {not_fraud_test} ({(not_fraud_test / len(y_test)) * 100:.2f}%)\n')
```

```
Class distribution in y_train:
0    99.82695
1     0.17305
Name: proportion, dtype: float64

Class distribution in y_test:
0    99.827955
1     0.172045
Name: proportion, dtype: float64

Total samples in y_train: 199364
Fraudulent transactions in y_train: 345 (0.17%)
Regular transactions in y_train: 199019 (99.83%)

Total samples in y_test: 85443
Fraudulent transactions in y_test: 147 (0.17%)
Regular transactions in y_test: 85296 (99.83%)
```

The training and test sets both contain 29 features, and these features are the same.

```
[15]: # Check if features in the training and test sets are the same
      train_features = X_train.columns
      test_features = X_test.columns

      if train_features.equals(test_features):
          print("The training and test sets contain the same features.")
          print(f"Number of features in train and test sets: {len(train_features)}")
      else:
          print("The training and test sets have different features.")
          print(f"Features in training set: {list(train_features)}")
          print(f"Features in test set: {list(test_features)}")
```

```
The training and test sets contain the same features.
Number of features in train and test sets: 29
```

Jupyter 10.2HD sample Last Checkpoint: 1 hour ago

File   Edit   View   Run   Kernel   Settings   Help

+  ✂  ▢  ▭  ▶  ■  C  ⏭  Code  ∨          JupyterLab ☐  ⊕  Python 3 (ipyke

### Features distributions by classes

We take a quick look a the features distribution by target class. By comparing boxplots, we would be able to spot features which have specific distribution if the item is fraudulent or not. In n general, we can see that many of the normal samples have significant outlier values.

```
[17]: import pandas as pd
      import matplotlib.pyplot as plt

      # Assuming X_train and y_train are already defined and contain your training data
      # y_train should be a Series with 1 for fraudulent and 0 for normal transactions

      # Ensure y_train is treated as a Series and get its values
      if isinstance(y_train, pd.Series):
          y_train_values = y_train.values
      else:
          y_train_values = y_train['Class'].values  # Assuming y_train is a DataFrame

      # Separate fraudulent and benign transactions
      fraudulent = X_train[y_train_values == 1]
      benign = X_train[y_train_values == 0]

      # Loop through all features to create boxplots
      for feature in X_train.columns:
          # Check if the feature is numeric
          if pd.api.types.is_numeric_dtype(X_train[feature]):
              fig, ax = plt.subplots(1, 1, figsize=(8, 6))
              plt.boxplot([fraudulent[feature], benign[feature]])
              plt.title(f'Distribution of {feature} by Transaction Class')
              plt.xticks([1, 2], ['Fraudulent Transactions', 'Normal Transactions'])
              plt.ylabel(feature)  # Add y-axis label for clarity
              plt.grid()  # Optional: Add grid for better readability
              plt.show()
```

Distribution of V9 by Transaction Class



Distribution of V10 by Transaction Class

Distribution of normAmount by Transaction Class

Jupyter 10.2HD sample Last Checkpoint: 1 hour ago

File   Edit   View   Run   Kernel   Settings   Help

JupyterLab   Python 3 (ipyk

## 2.2 Feature correlations

It can be interesting to get an intuition about the correlations between features, and the correlations between each feature and the target. We can see that the most divers correlations seem to be related to normAccount and every other feature.

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Assuming X_train is already defined and contains your training data

# Calculate the correlation matrix
corr = X_train.corr()

# Create a heatmap to visualize the correlations
plt.figure(figsize=(16, 12))
sns.heatmap(corr, annot=True, fmt=".2f", cmap='coolwarm', square=True, cbar_kws={"shrink": .8})
plt.title('Feature Correlation Heatmap')
plt.show()
```

JupyterLab   Python 3 (ipykernel)



Feature Correlation Heatmap

**Feature correlation with target**

Roughly half of the features have a somewhat significant correlation to the targets, but none are highly correlated.

```python
[21]: # Ensure y_train is a DataFrame with the correct column name
if isinstance(y_train, pd.Series):
    y_train = y_train.to_frame(name='Class')  # Convert Series to DataFrame with 'Class' column

# Concatenate X_train and y_train
df_combined = pd.concat([X_train, y_train], axis=1)

# Check if 'Class' column exists in the concatenated DataFrame
if 'Class' not in df_combined.columns:
    print("Error: 'Class' column not found in concatenated DataFrame")

# Calculate absolute correlations with the target variable 'Class'
corr_target = np.abs(df_combined.corr()['Class']).sort_values(ascending=False)[1:]

# Plotting the bar plot for feature correlations
fig, ax = plt.subplots(1, 1, figsize=(12, 5))
corr_target.plot(kind='bar', ax=ax)
plt.title('Features Correlation with Target', Fontweight='bold')
plt.ylabel('Correlation with Target')
plt.xlabel('Feature')
plt.show()
```

**Features Correlation with Target**

```python
# Importing necessary libraries
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import roc_auc_score, roc_curve, auc
import matplotlib.pyplot as plt

# Manually specifying random forest model parameters
rf_model = RandomForestClassifier(n_estimators=50, max_depth=5, class_weight='balanced', random_state=0)

# Fitting the model to training data
rf_model.fit(X_train, y_train)

# Get predicted probabilities for test set
y_test_pred_proba_rf = rf_model.predict_proba(X_test)[:, 1]

# Calculate ROC-AUC score
roc_auc_rf = roc_auc_score(y_test, y_test_pred_proba_rf)
print("ROC AUC score for Random Forest:", roc_auc_rf)

# Plot ROC curve
fpr, tpr, _ = roc_curve(y_test, y_test_pred_proba_rf)
roc_auc = auc(fpr, tpr)

# Plot ROC Curve
plt.figure(figsize=(8, 6))  # Increase figure size for better visibility
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')  # Diagonal line for random chance
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate', fontsize=12)
plt.ylabel('True Positive Rate', fontsize=12)
plt.title('ROC Curve for Random Forest', fontsize=14)
plt.legend(loc="lower right")
plt.grid(True)  # Adding a grid to make the plot clearer
plt.tight_layout()  # Adjust layout for better visibility of all elements
plt.show()
```

ROC Curve for Random Forest

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import GridSearchCV, train_test_split
from sklearn.linear_model import SGDClassifier
from sklearn.metrics import roc_auc_score, roc_curve, auc
import time

# Example dataset (replace with your actual dataset)
# X, y = pd.read_csv('your_dataset.csv'), pd.read_csv('your_labels.csv')

# Splitting dataset into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Define hyperparameter grid for SVM (using SGDClassifier)
params = {
    'loss': ['modified_huber'],
    'alpha': [0.00001, 0.0001, 0.001, 0.01, 0.1],
    'epsilon': [0.000001, 0.00001, 0.0001, 0.001, 0.01, 0.1, 1],
    'learning_rate': ['invscaling'],
    'eta0': [0.00001, 0.0001, 0.001],
    'class_weight': ['balanced']
}

# Initialize GridSearchCV
grid_search = GridSearchCV(SGDClassifier(random_state=0, fit_intercept=False),
                           param_grid=params, cv=3, scoring='roc_auc', n_jobs=-1, verbose=1)

# Start timer
start_time = time.time()

# Fit GridSearchCV
print("Starting grid search...")
grid_search.fit(X_train, y_train)
print("Grid search completed.")

# Output best parameters from grid search
best_params = grid_search.best_params_
print("Grid search best params")
print(best_params)

# Using best params to build SVM model
best_svm = SGDClassifier(**best_params, random_state=0, fit_intercept=False)
best_svm.fit(X_train, y_train)

# Calculate AUC score on training data
y_train_pred_proba = best_svm.decision_function(X_train)
best_auc_score = roc_auc_score(y_train, y_train_pred_proba)

# Display best AUC score and parameters
print("Using best grid search params to build SVM model.")
print(f"Best AUC score (training) : {best_auc_score}")
print(f"Best params : {best_params}")
```

```
print(f"CPU times: total: {cpu_time:.1f} s")

# Plotting ROC Curve
y_test_pred_proba = best_svm.decision_function(X_test)
fpr, tpr, _ = roc_curve(y_test, y_test_pred_proba)
roc_auc = auc(fpr, tpr)

plt.figure(figsize=(8, 6))  # Increase figure size
plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'SVM (AUC = {roc_auc:.2f})')
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')  # Diagonal line for random chance
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate', fontsize=12)
plt.ylabel('True Positive Rate', fontsize=12)
plt.title('ROC Curve for SVM', fontsize=14)
plt.legend(loc="lower right")
plt.grid(True)  # Adding a grid to make the plot clearer
plt.tight_layout()  # Adjust layout for better visibility of all elements
plt.show()
```

```
Starting grid search...
Fitting 3 folds for each of 105 candidates, totalling 315 fits
Grid search completed.
Grid search best params
{'alpha': 0.1, 'class_weight': 'balanced', 'epsilon': 1e-06, 'eta0': 1e-05, 'learning_rate': 'invscaling', 'loss': 'modified_huber'}
Using best grid search params to build SVM model.
Best AUC score (training) : 0.979084578645343
Best params : {'alpha': 0.1, 'class_weight': 'balanced', 'epsilon': 1e-06, 'eta0': 1e-05, 'learning_rate': 'invscaling', 'loss': 'modified_hube
CPU times: total: 43.9 s
```



ROC Curve for SVM

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import GridSearchCV, train_test_split
from xgboost import XGBClassifier
from sklearn.metrics import roc_auc_score, roc_curve, auc
import time

# Example dataset (replace with your actual dataset)
# X, y = pd.read_csv('your_dataset.csv'), pd.read_csv('your_labels.csv')

# Splitting dataset into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Define hyperparameter grid for XGBoost
params = {
    'max_depth': [4],
    'learning_rate': [0.00001, 0.0001, 0.001],
    'gamma': [0.1],
    'subsample': [0.2, 0.3, 0.4, 0.5],
    'objective': ['binary:logistic'],
    'eval_metric': ['auc'],
    'scale_pos_weight': [577],  # Value is sum(negative class)/sum(positive class)
    'n_estimators': [80]
}

# Initialize GridSearchCV
grid_search = GridSearchCV(XGBClassifier(booster='gbtree'), param_grid=params, cv=5, scoring='roc_auc', n_jobs=-1, verbose=1)

# Start timer
start_time = time.time()

# Fit GridSearchCV
print("Starting grid search...")
grid_search.fit(X_train, y_train)
print("Grid search completed.")

# Output best parameters from grid search
best_params = grid_search.best_params_
print("Grid search best params")
print(best_params)

# Using best params to build XGBoost model
xgb_model = XGBClassifier(
    max_depth=best_params['max_depth'],
    learning_rate=best_params['learning_rate'],  # Use the value directly
    gamma=best_params['gamma'],  # Use the value directly
    subsample=best_params['subsample'],  # Use the value directly
    objective='binary:logistic',
    eval_metric='auc',
    scale_pos_weight=577,
    n_estimators=80
)

# Fit the XGBoost model
```

```
print("Best AUC score (training) for XGBoost:", roc_auc_xgb_train)

# Plotting ROC Curve
y_test_pred_proba = xgb_model.predict_proba(X_test)[:, 1]
fpr, tpr, _ = roc_curve(y_test, y_test_pred_proba)
roc_auc_xgb = auc(fpr, tpr)

plt.figure(figsize=(8, 6))  # Increase figure size for better visibility
plt.plot(fpr, tpr, color='darkorange', lw=2, label='XGBoost (AUC = %0.2f)' % roc_auc_xgb)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')  # Diagonal line for random chance
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate', fontsize=12)
plt.ylabel('True Positive Rate', fontsize=12)
plt.title('ROC Curve for XGBoost', fontsize=14)
plt.legend(loc="lower right")
plt.grid(True)  # Adding a grid for clarity
plt.tight_layout()  # Adjust layout for better visibility of all elements
plt.show()
```

```
Starting grid search...
Fitting 5 folds for each of 12 candidates, totalling 60 fits
Grid search completed.
Grid search best params
{'eval_metric': 'auc', 'gamma': 0.1, 'learning_rate': 0.0001, 'max_depth': 4, 'n_estimators': 80, 'objective': 'binary:logistic', 'scale_pos_weigh
7, 'subsample': 0.4}
Best AUC score (training) for XGBoost: 0.9952162620087504
```



ROC Curve for XGBoost

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import GridSearchCV, train_test_split
from lightgbm import LGBMClassifier
from sklearn.metrics import roc_auc_score, roc_curve, auc
import time

# Example dataset (replace with your actual dataset)
# X, y = pd.read_csv('your_dataset.csv'), pd.read_csv('your_labels.csv')

# Splitting dataset into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Define hyperparameter grid for LightGBM
params = {
    'learning_rate': [0.0001],
    'max_depth': [9],
    'reg_alpha': [0.001],
    'reg_lambda': [0.001],
    'subsample': [0.0001],
    'objective': ['binary'],
    'scale_pos_weight': [577],
    'n_estimators': [60]
}

# Initialize GridSearchCV for LightGBM
grid_search = GridSearchCV(LGBMClassifier(random_state=0), param_grid=params, cv=5, scoring='roc_auc', n_jobs=-1, verbose=1)

# Start timer
start_time = time.time()

# Fit GridSearchCV
print("Starting grid search...")
grid_search.fit(X_train, y_train)
print("Grid search completed.")

# Output best parameters from grid search
best_params = grid_search.best_params_
print("Grid search best params")
print(best_params)

# Using best params to build LightGBM model
lgbm_model = LGBMClassifier(
    learning_rate=best_params['learning_rate'],  # Use the value directly
    max_depth=best_params['max_depth'],  # Use the value directly
    reg_alpha=best_params['reg_alpha'],  # Use the value directly
    reg_lambda=best_params['reg_lambda'],  # Use the value directly
    subsample=best_params['subsample'],  # Use the value directly
    objective='binary',
    scale_pos_weight=577,
    n_estimators=best_params['n_estimators']  # Use the value directly
)

# Fit the LightGBM model
```

```
plt.title( ROC Curve for LightGBM , fontsize=14)
plt.legend(loc="lower right")
plt.grid(True)  # Adding a grid for clarity
plt.tight_layout()  # Adjust layout for better visibility of all elements
plt.show()
```

```
Starting grid search...
Fitting 5 folds for each of 1 candidates, totalling 5 fits
[LightGBM] [Info] Number of positive: 394, number of negative: 227451
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.015376 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 7395
[LightGBM] [Info] Number of data points in the train set: 227845, number of used features: 29
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.001729 -> initscore=-6.358339
[LightGBM] [Info] Start training from score -6.358339
Grid search completed.
Grid search best params
{'learning_rate': 0.0001, 'max_depth': 9, 'n_estimators': 60, 'objective': 'binary', 'reg_alpha': 0.001, 'reg_lambda': 0.001, 'scale_pos_weight':
'subsample': 0.0001}
[LightGBM] [Info] Number of positive: 394, number of negative: 227451
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.016902 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 7395
[LightGBM] [Info] Number of data points in the train set: 227845, number of used features: 29
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.001729 -> initscore=-6.358339
[LightGBM] [Info] Start training from score -6.358339
Best AUC score (training) for LightGBM: 0.9996051138096415
```



ROC Curve for LightGBM

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import GridSearchCV, train_test_split
from catboost import CatBoostClassifier
from sklearn.metrics import roc_auc_score, roc_curve, auc
import time

# Example dataset (replace with your actual dataset)
# X, y = pd.read_csv('your_dataset.csv'), pd.read_csv('your_labels.csv')

# Splitting dataset into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Define hyperparameter grid for CatBoost
params = {
    'learning_rate': [0.001],
    'max_depth': [3],
    'reg_lambda': [0.001],
    'bootstrap_type': ['Bernoulli'],
    'subsample': [0.1],
    'eval_metric': ['AUC'],
    'scale_pos_weight': [577],
    'logging_level': ['Silent'],
    'n_estimators': [90]
}

# Initialize GridSearchCV for CatBoost
grid_search = GridSearchCV(CatBoostClassifier(random_state=0), param_grid=params, cv=5, scoring='roc_auc', n_jobs=-1, verbose=1)

# Start timer
start_time = time.time()

# Fit GridSearchCV
print("Starting grid search...")
grid_search.fit(X_train, y_train)
print("Grid search completed.")

# Output best parameters from grid search
best_params = grid_search.best_params_
print("Grid search best params")
print(best_params)

# Best score from grid search
print("Best AUC score from grid search:", grid_search.best_score_)

# Using best params to build CatBoost model
cat_model = CatBoostClassifier(
    learning_rate=best_params['learning_rate'],  # Use the value directly
    max_depth=best_params['max_depth'],  # Use the value directly
    reg_lambda=best_params['reg_lambda'],  # Use the value directly
    bootstrap_type=best_params['bootstrap_type'],  # Use the value directly
    subsample=best_params['subsample'],  # Use the value directly
    eval_metric='AUC',  # Directly assign since it's a string
    scale_pos_weight=577,
```

```
# Plotting ROC Curve
fpr, tpr, _ = roc_curve(y_test, y_test_pred_proba)
roc_auc_cat_test = auc(fpr, tpr)

plt.figure(figsize=(8, 6))  # Increase figure size for better visibility
plt.plot(fpr, tpr, color='darkorange', lw=2, label='CatBoost (AUC = %0.2f)' % roc_auc_cat_test)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')  # Diagonal line for random chance
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate', fontsize=12)
plt.ylabel('True Positive Rate', fontsize=12)
plt.title('ROC Curve for CatBoost', fontsize=14)
plt.legend(loc="lower right")
plt.grid(True)  # Adding a grid for clarity
plt.tight_layout()  # Adjust layout for better visibility of all elements
plt.show()
```
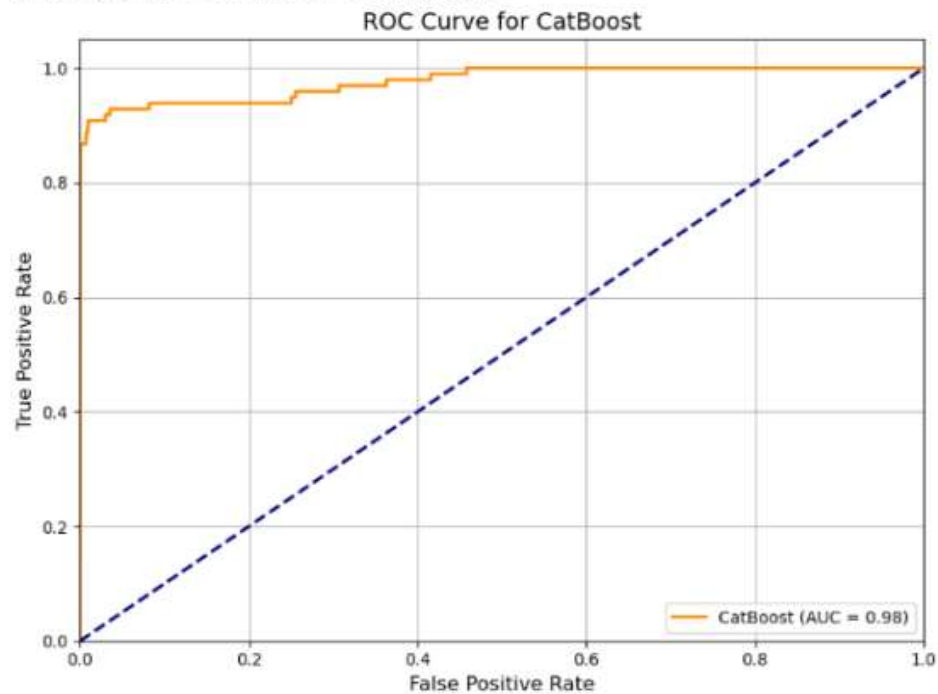
```
Starting grid search...
Fitting 5 folds for each of 1 candidates, totalling 5 fits
Grid search completed.
Grid search best params
{'bootstrap_type': 'Bernoulli', 'eval_metric': 'AUC', 'learning_rate': 0.001, 'logging_level': 'Silent', 'max_depth': 3, 'n_estimators': 90, 'reg_
a': 0.001, 'scale_pos_weight': 577, 'subsample': 0.1}
Best AUC score from grid search: 0.9784795899452634
Best AUC score (testing) for CatBoost: 0.9771116979431052
```



ROC Curve for CatBoost

```python
# Import necessary libraries
import numpy as np
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler

# Assuming X_train is your training dataset

# Standardize the dataset
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X_train)

# Fit PCA to the scaled data
pca = PCA()
pca.fit(X_scaled)

# Variance explained by each principal component
explained_variance = pca.explained_variance_ratio_

# Plot the explained variance
plt.figure(figsize=(10, 6))
plt.plot(range(1, len(explained_variance) + 1), explained_variance, marker='o', linestyle='--')
plt.axvline(x=11, color='r', linestyle='--', label='11 Components (45% Variance)')
plt.axvline(x=25, color='orange', linestyle='--', label='25 Components (90% Variance)')
plt.title('Explained Variance by Principal Components')
plt.xlabel('Principal Component')
plt.ylabel('Explained Variance Ratio')
plt.xticks(range(1, len(explained_variance) + 1))
plt.legend()
plt.grid()
plt.show()

# Cumulative explained variance
cumulative_variance = np.cumsum(explained_variance)

# Plot cumulative explained variance
plt.figure(figsize=(10, 6))
plt.plot(range(1, len(cumulative_variance) + 1), cumulative_variance, marker='o', linestyle='-')
plt.axhline(y=0.90, color='orange', linestyle='--', label='90% Variance')
plt.title('Cumulative Explained Variance by Principal Components')
plt.xlabel('Principal Component')
plt.ylabel('Cumulative Explained Variance')
plt.xticks(range(1, len(cumulative_variance) + 1))
plt.legend()
plt.grid()
plt.show()
```
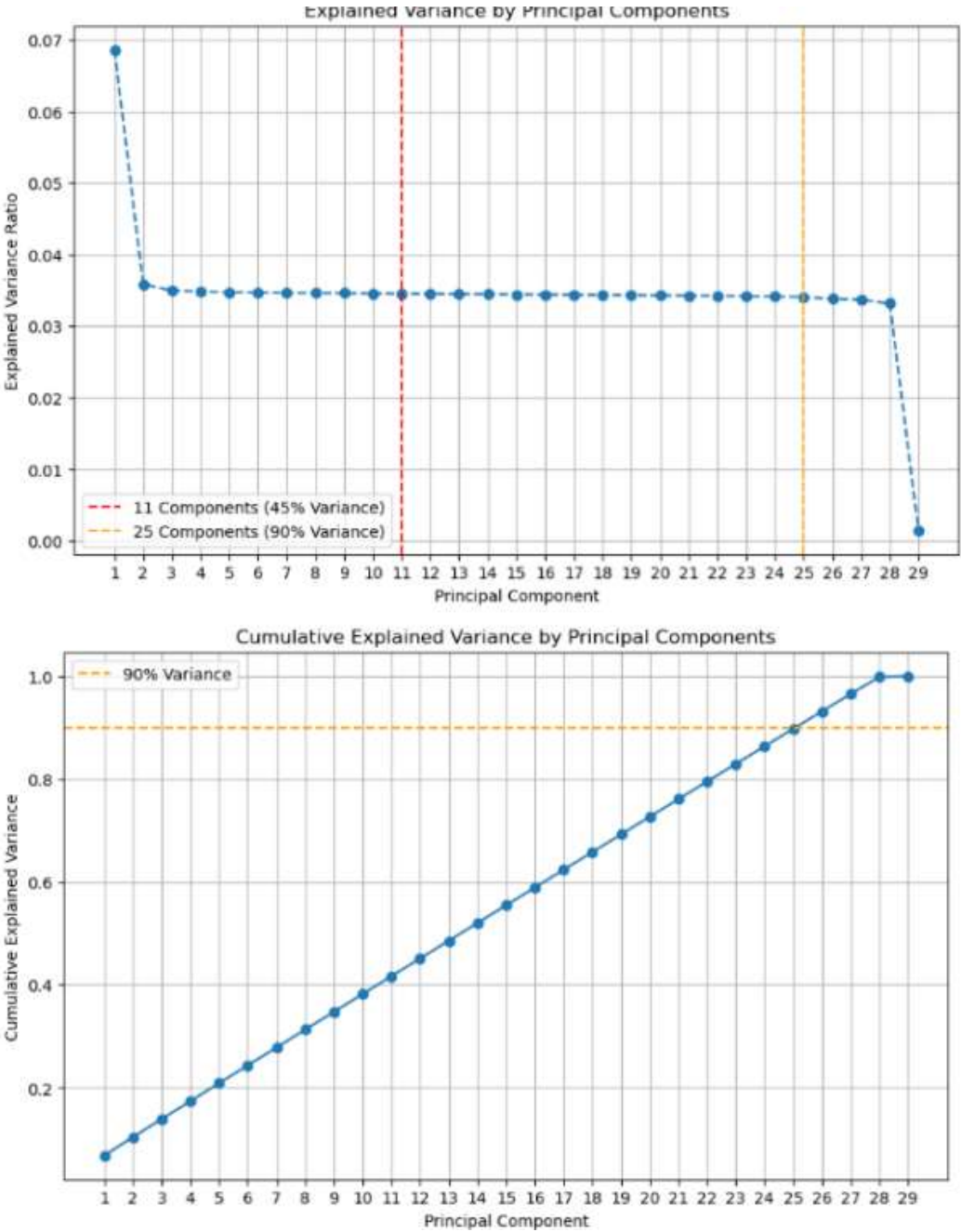
Explained Variance by Principal Components



Cumulative Explained Variance by Principal Components

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
from sklearn.decomposition import PCA

# Assuming X_train and y_train are already defined

def apply_PCA(X, n_components=3, plot=True):
    """Fit PCA on the data and optionally plot the explained variance."""
    pca = PCA(n_components=n_components)
    pca_components = pca.fit_transform(X)

    if plot:
        # Plot the explained variance
        plt.figure(figsize=(10, 6))
        plt.plot(np.cumsum(pca.explained_variance_ratio_), marker='o', linestyle='--')
        plt.axhline(y=0.90, color='r', linestyle='-')
        plt.axvline(x=11, color='g', linestyle='--')  # Example line for 11 components
        plt.title('Cumulative Explained Variance by Principal Components')
        plt.xlabel('Number of Principal Components')
        plt.ylabel('Cumulative Explained Variance')
        plt.grid()
        plt.show()

    return pca_components, pca

def plot_pca_3D(components, labels):
    """Plot the PCA results in a 3D scatter plot."""
    fig = plt.figure(figsize=(10, 8))
    ax = fig.add_subplot(111, projection='3d')

    # Map labels to colors: fraud = red, non-fraud = green
    colors = np.where(labels == 1, 'red', 'green')  # Assuming '1' indicates fraud

    # Plot the points
    scatter = ax.scatter(components[:, 0], components[:, 1], components[:, 2],
                         c=colors, marker='o', alpha=0.6)

    ax.set_title('3D PCA Plot of Transactions')
    ax.set_xlabel('1st Principal Component')
    ax.set_ylabel('2nd Principal Component')
    ax.set_zlabel('3rd Principal Component')

    # Legend
    red_patch = plt.Line2D([0], [0], marker='o', color='w', markerfacecolor='red', markersize=10, label='Fraudulent')
    green_patch = plt.Line2D([0], [0], marker='o', color='w', markerfacecolor='green', markersize=10, label='Non-Fraudulent')
    ax.legend(handles=[red_patch, green_patch])

    plt.show()

# Fit PCA to X_train data
pca_components, pca = apply_PCA(X_train, plot=True)
```
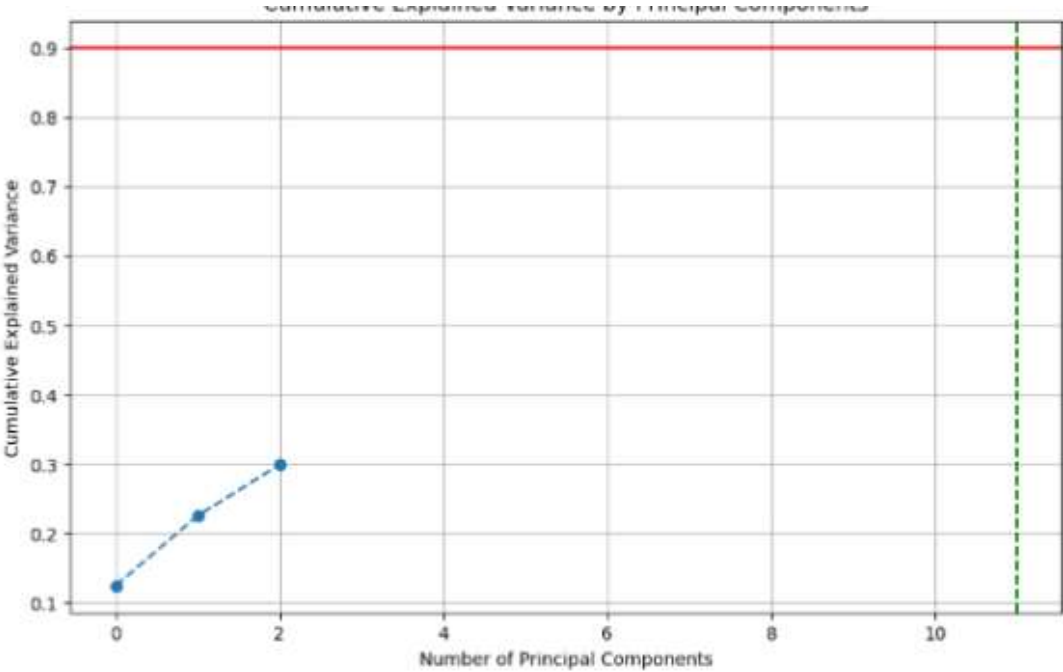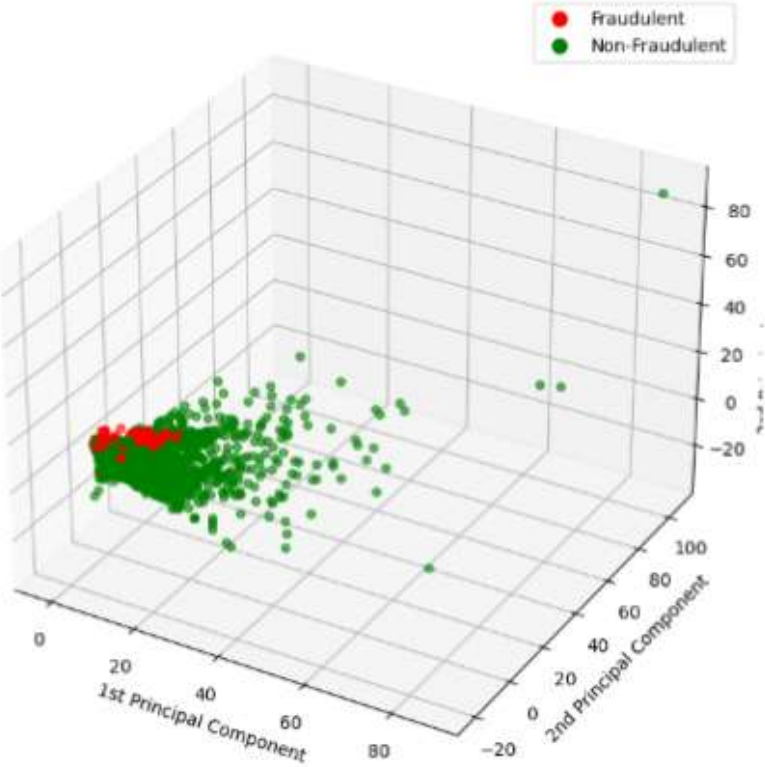
Cumulative Explained Variance by Principal Components



3D PCA Plot of Transactions

```
import numpy as np
import pandas as pd

# Assuming X_train and y_train are already defined
# Assuming pca_components were already calculated from previous PCA

# Find outlier indices based on the third principal component (z-axis)
# We will take the two biggest outliers for the z-axis
indexes_outliers = np.argsort(pca_components[:, 2])[-2:][::-1]

# Remove the outliers from X_train
X_train_removed = X_train.drop(indexes_outliers)

# Remove the outliers from y_train using NumPy indexing
y_train_removed = np.delete(y_train, indexes_outliers)

# Fit PCA on the cleaned data
pca_components_cleaned, _ = apply_PCA(X_train_removed, plot=False)

# Plot the PCA 3D result without the outliers
plot_pca_3D(pca_components_cleaned, y_train_removed)
```

3D PCA Plot of Transactions

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

# Apply PCA to the test data
pca_components_test, _ = apply_PCA(X_test, plot=False)

# Identify outliers in the test set
indexes_outliers_test = np.argsort(pca_components_test[:, 1])[-2:][::-1]

# Create a mask to remove outliers from X_test
mask = ~X_test.index.isin(indexes_outliers_test)
X_test_removed = X_test[mask]

# Apply PCA on the cleaned test data
X_test_values, _ = apply_PCA(X_test_removed, plot=False)

# Separate cleaned fraud and not fraud classes from training data
fraud_cleaned = pca_components_cleaned[y_train_removed == 1]
not_fraud_cleaned = pca_components_cleaned[y_train_removed == 0]

# Plotting the 2D PCA results
fig, ax = plt.subplots(1, 3, figsize=(18, 5))

# Training set plot (Not Fraud)
ax[0].scatter(not_fraud_cleaned[:, 0], not_fraud_cleaned[:, 1], label='Not Fraud', color='g', alpha=0.2)
ax[0].set_title('First Two Components of PCA - Train Set (Not Fraud)', fontweight='bold')
ax[0].set_xlabel('First Principal Component')
ax[0].set_ylabel('Second Principal Component')
ax[0].legend(loc='upper right')

# Training set plot (Fraud)
ax[1].scatter(not_fraud_cleaned[:, 0], not_fraud_cleaned[:, 1], label='Not Fraud', color='g', alpha=0.1)
ax[1].scatter(fraud_cleaned[:, 0], fraud_cleaned[:, 1], label='Fraud', color='r', alpha=0.9)
ax[1].set_title('First Two Components of PCA - Train Set (Fraud Highlighted)', fontweight='bold')
ax[1].set_xlabel('First Principal Component')
ax[1].set_ylabel('Second Principal Component')
ax[1].legend(loc='upper right')

# Testing set plot
ax[2].scatter(X_test_values[:, 0], X_test_values[:, 1], label='Test Values', color='blue', alpha=0.2)
ax[2].set_title('First Two Components of PCA - Test Set', fontweight='bold')
ax[2].set_xlabel('First Principal Component')
ax[2].set_ylabel('Second Principal Component')
ax[2].legend(loc='upper right')

plt.tight_layout()
plt.show()
```
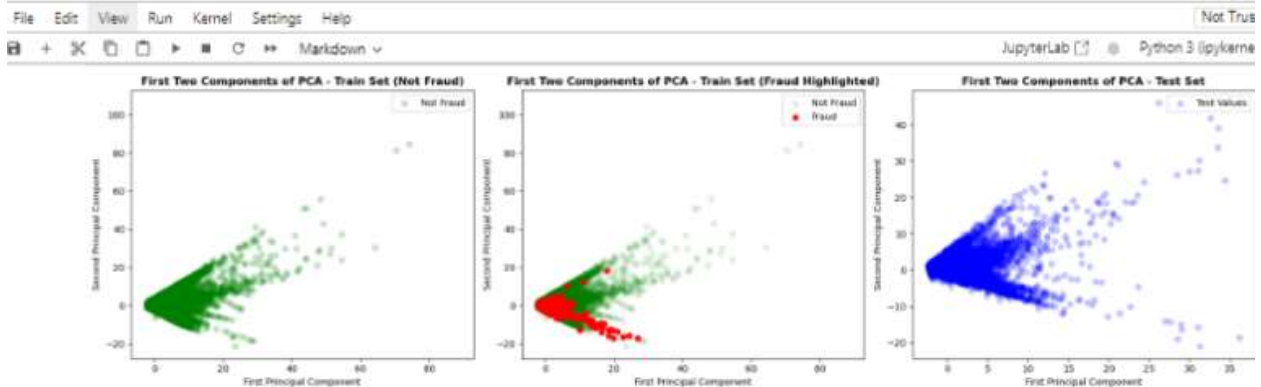
File   Edit   View   Run   Kernel   Settings   Help                                                                                  Not Trus

＋   ✂   ▢   ▢   ▶   ■   C   ▸▸   Markdown ⌄                                                       JupyterLab ⬀   ●   Python 3 (ipykerne

```python
from sklearn.decomposition import PCA
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from xgboost import XGBClassifier
from lightgbm import LGBMClassifier
from catboost import CatBoostClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import roc_auc_score, roc_curve, auc
import matplotlib.pyplot as plt
import time

# PCA with reduced components to limit the load
pca = PCA(n_components=5, whiten=True, random_state=0).fit(X_train_removed)
X_train_pca = pca.transform(X_train_removed)
X_test_pca = pca.transform(X_test)

# Function to evaluate models with reduced parameter space
def evaluate_model(model, params, X_train, y_train, X_test, y_test):
    grid_search = GridSearchCV(model, param_grid=params, cv=3, scoring='roc_auc', n_jobs=-1, verbose=1)

    start_time = time.time()
    grid_search.fit(X_train, y_train)
    end_time = time.time()

    # Best model and AUC score on test set
    best_model = grid_search.best_estimator_
    y_test_probs = best_model.predict_proba(X_test)[:, 1]
    test_auc = roc_auc_score(y_test, y_test_probs)

    print(f"\nBest parameters: {grid_search.best_params_}")
    print(f"Best AUC score (training): {grid_search.best_score_}")
    print(f"Test AUC score: {test_auc}")
    print(f"Total time: {end_time - start_time:.2f} seconds")

    # ROC Curve
    fpr, tpr, _ = roc_curve(y_test, y_test_probs)
    roc_auc = auc(fpr, tpr)

    plt.figure()
    plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC curve (AUC = {roc_auc:.2f})')
    plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
    plt.xlim([0.0, 1.0])
    plt.ylim([0.0, 1.05])
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
    plt.title(f'ROC Curve - Test Set ({model.__class__.__name__})')
    plt.legend(loc="lower right")
    plt.show()

    return test_auc

# Optimized model parameters for faster computation
models_params = {
    'RandomForestClassifier': {
        'model': RandomForestClassifier(random_state=0),
        'params': {'n_estimators': [50], 'max_depth': [7], 'class_weight': ['balanced']}
    },
    'SVM': {
        'model': SVC(probability=True, random_state=0),
        'params': {'C': [1], 'kernel': ['linear']}
    },
    'XGBoost': {
        'model': XGBClassifier(random_state=0),
        'params': {'n_estimators': [50], 'max_depth': [7], 'learning_rate': [0.1]}
    },
    'LightGBM': {
        'model': LGBMClassifier(random_state=0),
        'params': {'n_estimators': [50], 'max_depth': [7], 'learning_rate': [0.1]}
```
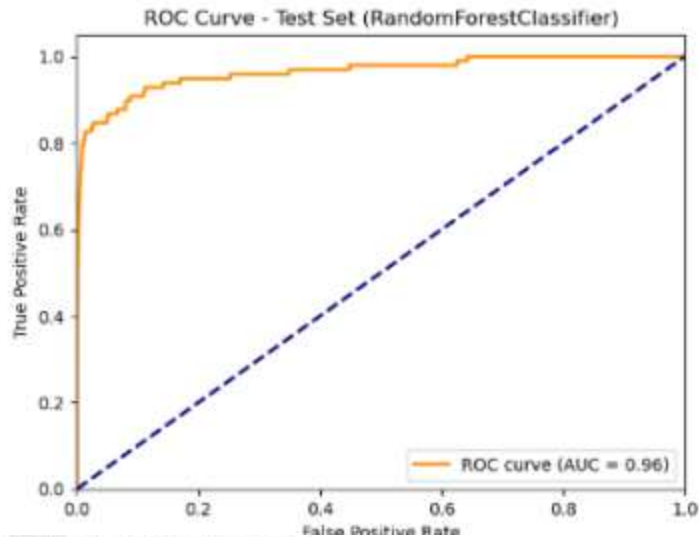
le Edit View Run Kernel Settings Help

+ ✕ ▢ ▢ ▶ ■ C ⋙ Code ⌄ JupyterLab

Evaluating RandomForestClassifier
Fitting 3 folds for each of 1 candidates, totalling 3 fits

Best parameters: {'class_weight': 'balanced', 'max_depth': 7, 'n_estimators': 50}
Best AUC score (training): 0.7975946330561571
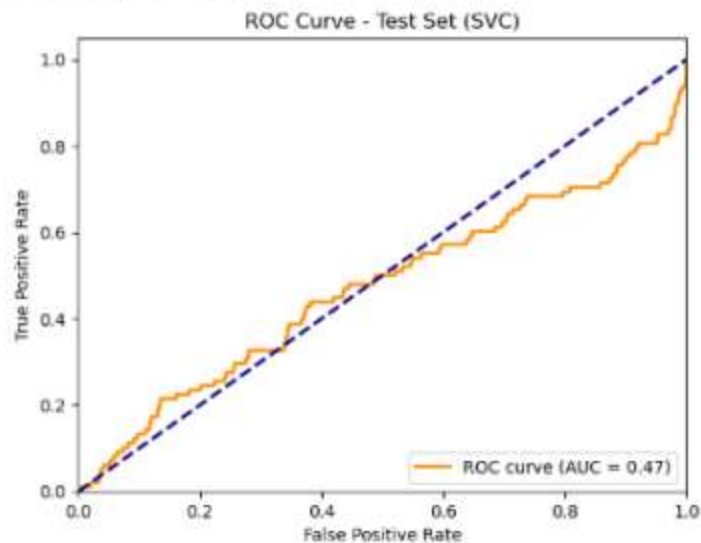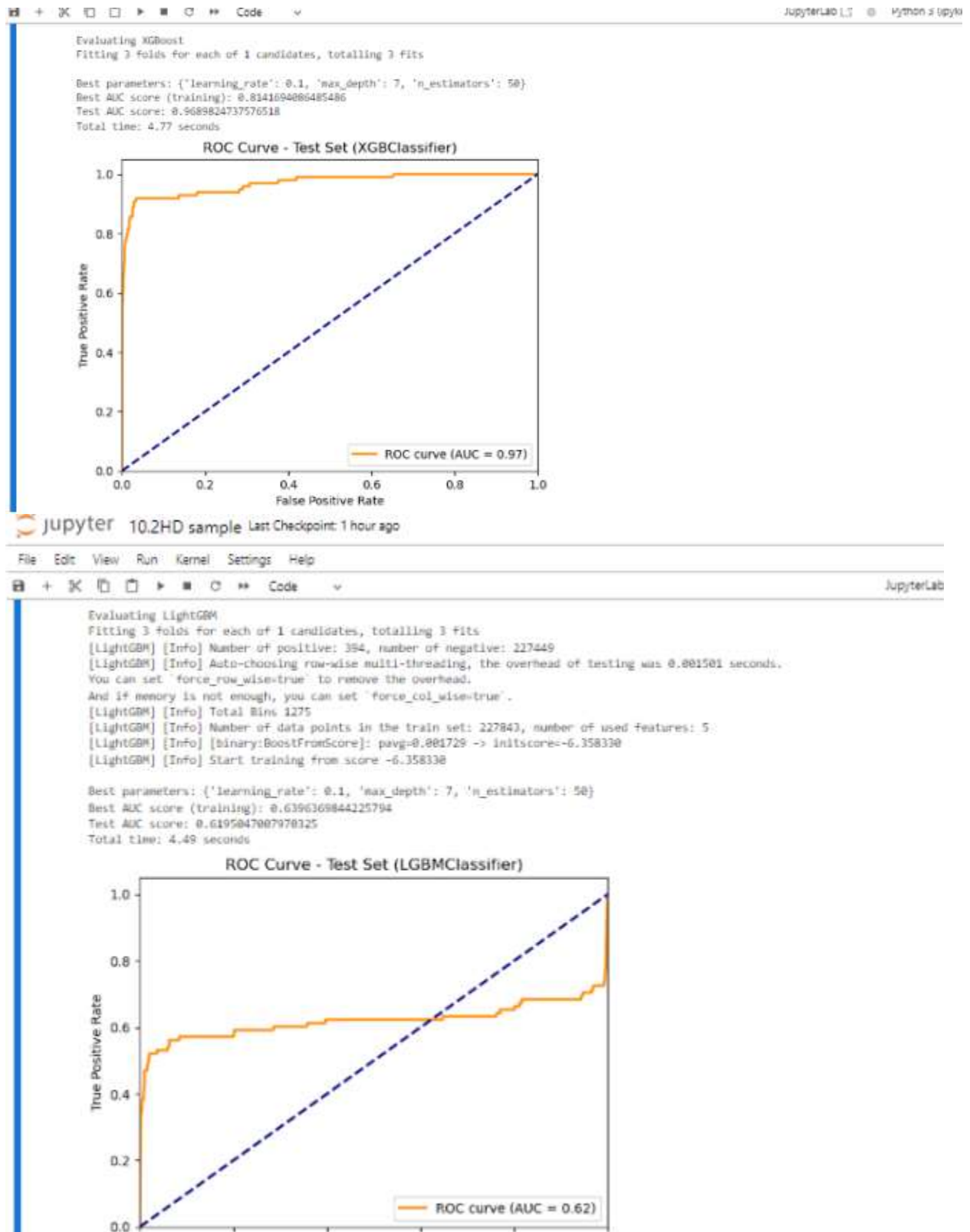Test AUC score: 0.9644805938695118
Total time: 47.67 seconds

### ROC Curve - Test Set (RandomForestClassifier)



le Edit View Run Kernel Settings Help

+ ✕ ▢ ▢ ▶ ■ C ⋙ Code ⌄

Evaluating SVM
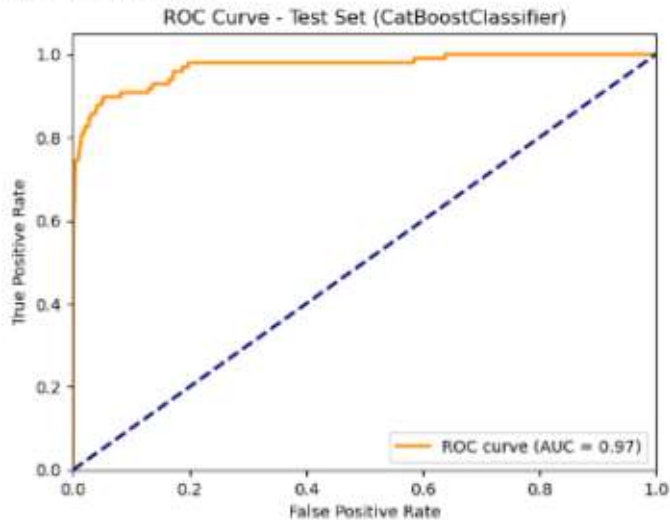Fitting 3 folds for each of 1 candidates, totalling 3 fits

Best parameters: {'C': 1, 'kernel': 'linear'}
Best AUC score (training): 0.6922595033802231
Test AUC score: 0.47208129431626335
Total time: 393.21 seconds

### ROC Curve - Test Set (SVC)

▣ + ✂ ▢ ▢ ▶ ■ ⟳ ⤚ Code ⌄                                    JupyterLab |☐  ⊚  Python 3 (ipyk

Evaluating XGBoost
Fitting 3 folds for each of 1 candidates, totalling 3 fits

Best parameters: {'learning_rate': 0.1, 'max_depth': 7, 'n_estimators': 50}
Best AUC score (training): 0.8141694086485486
Test AUC score: 0.9689824737576518
Total time: 4.77 seconds

**ROC Curve - Test Set (XGBClassifier)**



🌀 Jupyter  10.2HD sample  Last Checkpoint: 1 hour ago

File   Edit   View   Run   Kernel   Settings   Help

▣ + ✂ ▢ ▢ ▶ ■ ⟳ ⤚ Code ⌄                                    JupyterLab

Evaluating LightGBM
Fitting 3 folds for each of 1 candidates, totalling 3 fits
[LightGBM] [Info] Number of positive: 394, number of negative: 227449
[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing was 0.001501 seconds.
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 1275
[LightGBM] [Info] Number of data points in the train set: 227843, number of used features: 5
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.001729 -> initscore=-6.358330
[LightGBM] [Info] Start training from score -6.358330

Best parameters: {'learning_rate': 0.1, 'max_depth': 7, 'n_estimators': 50}
Best AUC score (training): 0.6396369844225794
Test AUC score: 0.6195047007970325
Total time: 4.49 seconds

**ROC Curve - Test Set (LGBMClassifier)**

```
Evaluating CatBoost
Fitting 3 folds for each of 1 candidates, totalling 3 fits

Best parameters: {'depth': 7, 'iterations': 50, 'learning_rate': 0.1}
Best AUC score (training): 0.833825918809065
Test AUC score: 0.9705405234688135
Total time: 7.12 seconds
```

### ROC Curve - Test Set (CatBoostClassifier)



```
Test AUC Scores for All Models:
RandomForestClassifier: 0.9645
SVM: 0.4726
XGBoost: 0.9690
LightGBM: 0.6195
```

```python
import pandas as pd
from imblearn.over_sampling import SMOTE
from imblearn.under_sampling import RandomUnderSampler
from imblearn.pipeline import Pipeline
from sklearn.model_selection import train_test_split


# SMOTE to oversample the minority class to 10% of the majority
over = SMOTE(sampling_strategy=0.1, random_state=0)

# RandomUnderSampler to undersample the majority class to 30%
under = RandomUnderSampler(sampling_strategy=0.3, random_state=0)

# Combine SMOTE and RandomUnderSampler in a pipeline
steps = [('oversample', over), ('undersample', under)]
pipeline = Pipeline(steps=steps)

# Resample the training dataset
X_train_smote, y_train_smote = pipeline.fit_resample(X_train, y_train)

# Convert y_train_smote (NumPy array) to Pandas Series to use value_counts()
y_train_smote_series = pd.Series(y_train_smote)

# Print new class distribution
print(f'New distribution of target after SMOTE and undersampling:\n{y_train_smote_series.value_counts()}')

# Split the resampled dataset into training and validation sets
X_train_smote, X_valid_nonused, y_train_smote, y_valid_nonused = train_test_split(
    X_train_smote, y_train_smote, test_size=0.1, random_state=0, stratify=y_train_smote)

# No scaling needed in this step for random forest (if needed, uncomment and apply scaling)
# scaler = StandardScaler()
# X_train_scaled = scaler.fit_transform(X_train_smote)
# X_test_scaled = scaler.transform(X_test)
```

```
New distribution of target after SMOTE and undersampling:
0    75816
1    22745
Name: count, dtype: int64
```

```python
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from xgboost import XGBClassifier
from lightgbm import LGBMClassifier
from catboost import CatBoostClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import roc_auc_score, roc_curve, auc
import matplotlib.pyplot as plt

# Function to evaluate a model and print results
def evaluate_model(model, params, X_train, y_train, X_test, y_test, model_name):
    # Perform grid search with fewer parameters for faster computation
    grid_search = GridSearchCV(model, param_grid=params, cv=2, scoring='roc_auc', n_jobs=-1)
    grid_search.fit(X_train, y_train)

    print(f"\n{model_name} - Best params: {grid_search.best_params_}")
    print(f"{model_name} - Best AUC score (training set): {grid_search.best_score_:.4f}")

    # Train best model
    best_model = grid_search.best_estimator_
    best_model.fit(X_train, y_train)

    # Predict probabilities for the test set
    y_test_proba = best_model.predict_proba(X_test)[:, 1]

    # Calculate AUC score for the test set
    auc_score = roc_auc_score(y_test, y_test_proba)
    print(f"{model_name} - AUC score (test set): {auc_score:.4f}")

    # Plot ROC curve
    plot_roc_curve(y_test, y_test_proba, f"ROC Curve for {model_name}")

    return auc_score

# Function to plot ROC curve
def plot_roc_curve(y_test, y_test_proba, title):
    fpr, tpr, _ = roc_curve(y_test, y_test_proba)
    roc_auc = auc(fpr, tpr)

    plt.figure()
    plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC curve (AUC = {roc_auc:.2f})')
    plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
    plt.xlim([0.0, 1.0])
    plt.ylim([0.0, 1.05])
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
    plt.title(title)
    plt.legend(loc="lower right")
    plt.show()

# Define the parameters for each model with reduced ranges
params_random_forest = {
    'n_estimators': [50, 70],   # Reduced to fewer values
    'max_depth': [7, 8],        # Narrowed down the search range
    'class_weight': ['balanced']
}

params_svm = {
    'C': [1],                   # Fixed for faster evaluation
    'kernel': ['linear'],       # Using linear kernel for speed
}
```
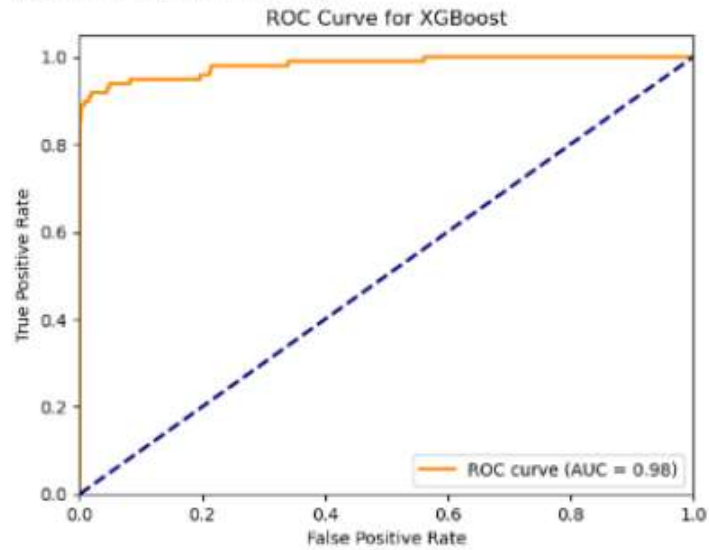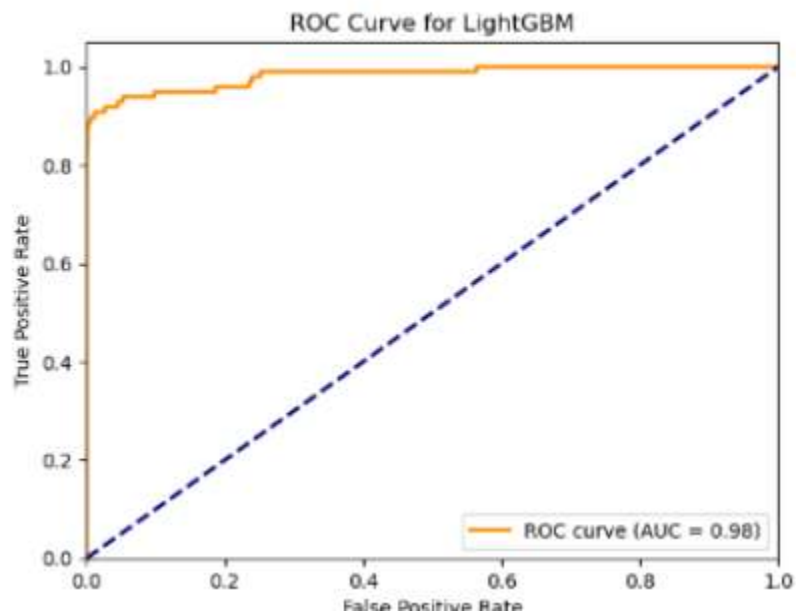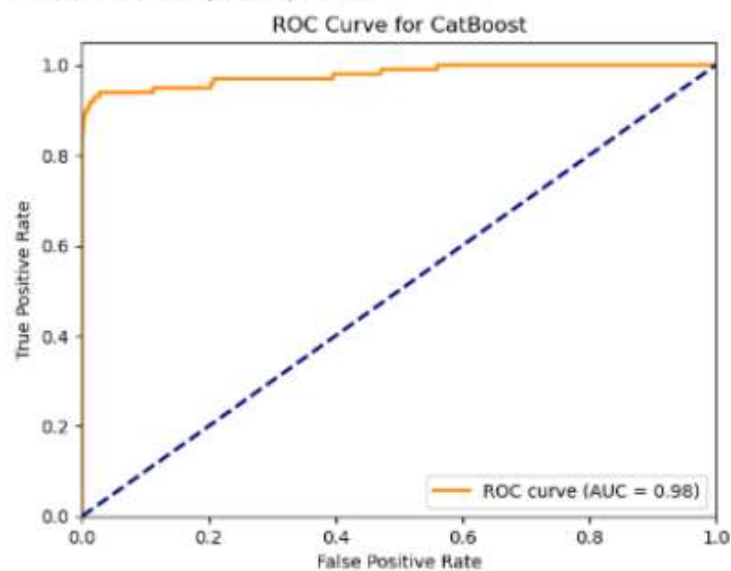
33

File    Edit    View    Run    Kernel    Settings    Help

🖫    +    ✂    🗐    🗂    ▶    ■    C    ⏩    Code        ⌄

Random Forest - Best params: {'class_weight': 'balanced', 'max_depth': 8, 'n_estimators': 70}
Random Forest - Best AUC score (training set): 0.9976
Random Forest - AUC score (test set): 0.9830

### ROC Curve for Random Forest



SVM - Best params: {'C': 1, 'kernel': 'linear'}
SVM - Best AUC score (training set): 0.9869
SVM - AUC score (test set): 0.9721

### ROC Curve for SVM

XGBoost - Best params: {'learning_rate': 0.1, 'max_depth': 3, 'n_estimators': 50}
XGBoost - Best AUC score (training set): 0.9942
XGBoost - AUC score (test set): 0.9821



ROC Curve for XGBoost

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
LightGBM - AUC score (test set): 0.9823



ROC Curve for LightGBM

CatBoost - Best params: {'depth': 3, 'iterations': 50, 'learning_rate': 0.1}
CatBoost - Best AUC score (training set): 0.9924
CatBoost - AUC score (test set): 0.9791



ROC Curve for CatBoost

Final AUC Scores for All Models:
Random Forest: 0.9830
SVM: 0.9721
XGBoost: 0.9821
LightGBM: 0.9823

```python
import numpy as np
import pandas as pd
import time
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from xgboost import XGBClassifier
from lightgbm import LGBMClassifier
from catboost import CatBoostClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import roc_auc_score, roc_curve
import matplotlib.pyplot as plt

# Start timing
start_time = time.time()

# Simulated data - replace these with your actual training and test datasets
X_train_smote = pd.DataFrame(np.random.rand(100, 10), columns=[f'feature_{i}' for i in range(10)])  # Example training features
y_train_smote = pd.Series(np.random.randint(0, 2, size=100))  # Example binary target
X_test = pd.DataFrame(np.random.rand(50, 10), columns=[f'feature_{i}' for i in range(10)])  # Example test features
y_test = pd.Series(np.random.randint(0, 2, size=50))  # Example binary target

# Simulated code for anomaly ratio in duplicated training samples
counts = y_train_smote[X_train_smote.duplicated(keep=False)].value_counts().values
anomaly_ratio = 100 * (counts[1] / (counts[0] + counts[1])) if len(counts) > 1 else 0.0
print(f'Anomaly ratio in duplicated training samples: {anomaly_ratio:.2f}%')

# Feature engineering
column_names = X_train_smote.columns.tolist()
X_train_smote_eng = X_train_smote.copy()
X_test_eng = X_test.copy()
for col in column_names:
    X_train_smote_eng[f'{col}_sq'] = X_train_smote_eng[col] ** 2
    X_test_eng[f'{col}_sq'] = X_test_eng[col] ** 2

# Duplicate feature for duplicate samples
X_train_smote_eng['dup'] = X_train_smote_eng.duplicated(keep=False).astype(int)
X_test_eng['dup'] = X_test_eng.duplicated(keep=False).astype(int)

# Number of features post engineering
num_features = X_train_smote_eng.shape[1]
print(f'Number of features post engineering: {num_features}')

# Function to evaluate a model
def evaluate_model(model, params, model_name):
    grid_search = GridSearchCV(model, param_grid=params, cv=5, scoring='roc_auc', n_jobs=-1)
    grid_search.fit(X_train_smote_eng, y_train_smote)

    print(f"{model_name} - Grid search best params: {grid_search.best_params_}")
    best_auc_train = grid_search.best_score_
    print(f"{model_name} - Best AUC score (training): {best_auc_train:.4f}")
```
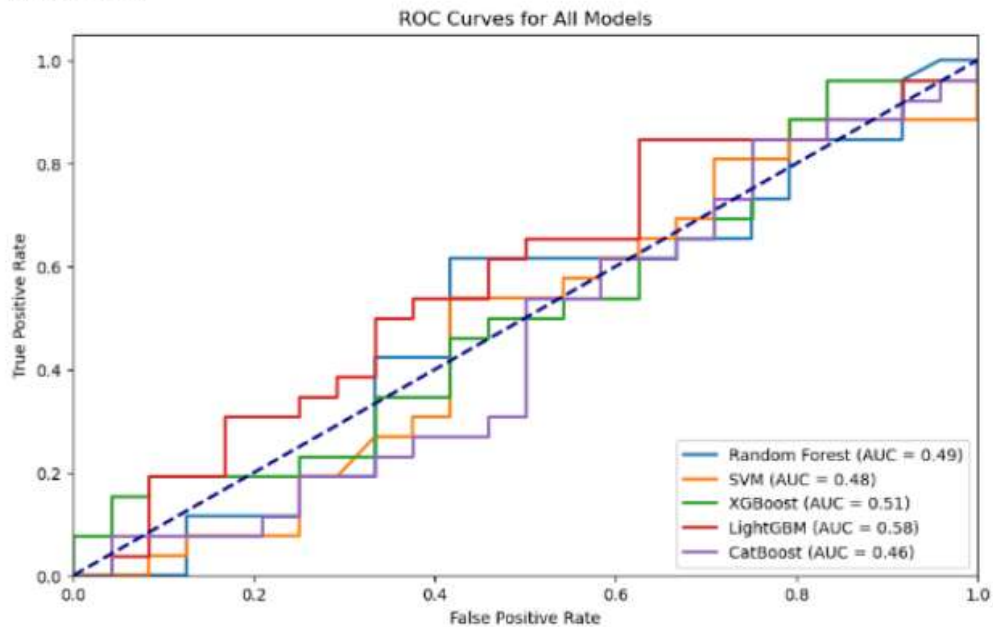
```
Test AUC Scores for All Models:
Random Forest: 0.4928
SVM: 0.4776
XGBoost: 0.5096
LightGBM: 0.5769
CatBoost: 0.4551
```



ROC Curves for All Models

```
CPU times: total: 1855.75 s
Wall time: 1.21 s
```

## 3.3 Train test split

In the rest of the document, we will use the SMOTE techniques. As per the case study, a different training set for the probability calibration step was used. I suspect this is
that a different set can be fitted onto CalibratedClassifierCV and avoid using the same sets for training and calibration.

Since using a separate validation set with GridSearchCV isn't necessary, I will just use the smaller training dataset.

- A smaller training dataset : (X_train_cut, y_train_cut)
- (unused) A validation set (X_valid, y_valid), which will be used for hyperparameters tunning

I make sure to use the `stratify` parameter to split the data, in order to keep the inital proportion of each classes in the splitted datasets.

```python
[56]: from sklearn.model_selection import train_test_split

# Perform a train-test split with stratification
X_train_cut, X_valid_unused, y_train_cut, y_valid_unused = train_test_split(
    X_train_smote,
    y_train_smote,
    test_size=0.1,  # 10% validation set
    random_state=41,
    stratify=y_train_smote.values  # Keep class proportions the same
)

# Print sizes of the split datasets
print(f"Training set size: {X_train_cut.shape}, Validation set size (unused): {X_valid_unused.shape}")
```

```
Training set size: (90, 10), Validation set size (unused): (10, 10)
```

```python
import numpy as np

# Calculate the ratio of non-anomalous to anomalous instances
# Non-anomaly ratio
num = y_train_smote.value_counts()[0] / len(y_train_smote)
print(f"Non-anomaly ratio: {num:.2%}")

# Anomaly ratio
denom = y_train_smote.value_counts()[1] / len(y_train_smote)
print(f"Anomaly ratio: {denom:.2%}")

# Ratio of non-anomaly instances to anomaly instances
res = num / denom
print(f"Ratio of non-anomaly to anomaly: {res:.2f}")

# Create array of weights: ratio for anomalous instances, 1 for non-anomalous
sample_weight = np.array([res if i == 1 else 1 for i in y_train_smote.values.ravel()])

# Normalize the weights
train_res_array = sample_weight / len(sample_weight)
print(f"Sample weights (first 10): {train_res_array[:10]}")
```

```python
from sklearn.model_selection import train_test_split

# Create calibration split out of SMOTE data
X_train_calib, X_valid_calib, y_train_calib, y_valid_calib, sw_train, sw_valid = \
    train_test_split(X_train_smote, y_train_smote, train_res_array, test_size=0.1, random_state=42, stratify=y_train_smote.values)

# Verify the shapes of the splits
print(f"Training data shape (calibration split): {X_train_calib.shape}")
print(f"Validation data shape (calibration split): {X_valid_calib.shape}")
print(f"Training weights shape: {sw_train.shape}")
print(f"Validation weights shape: {sw_valid.shape}")
```
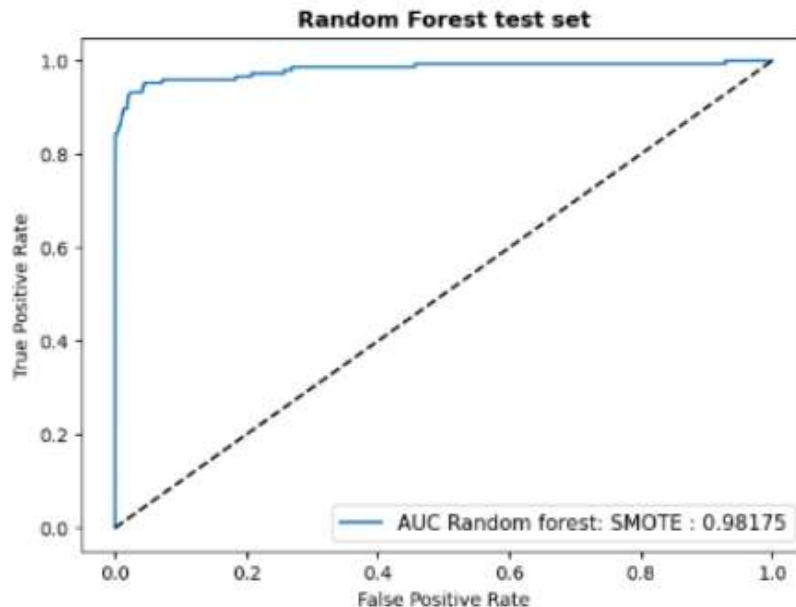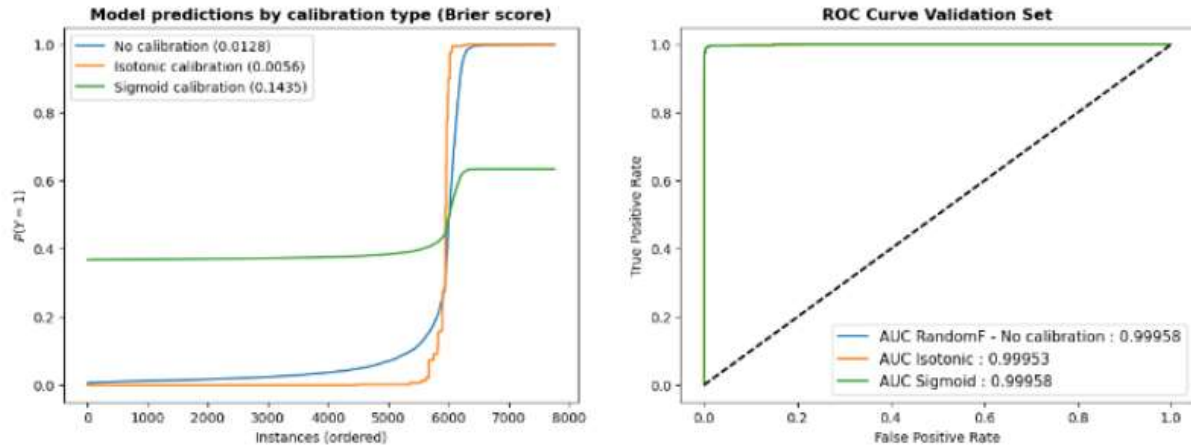
```
Training data shape (calibration split): (90, 10)
Validation data shape (calibration split): (10, 10)
Training weights shape: (90,)
Validation weights shape: (10,)

Best params : {'class_weight': 'balanced', 'max_depth': 9, 'n_estimators': 70}
CPU times: total: 11.7 s
Wall time: 33.2 s
```



Random Forest test set

AUC Random forest: SMOTE : 0.98175

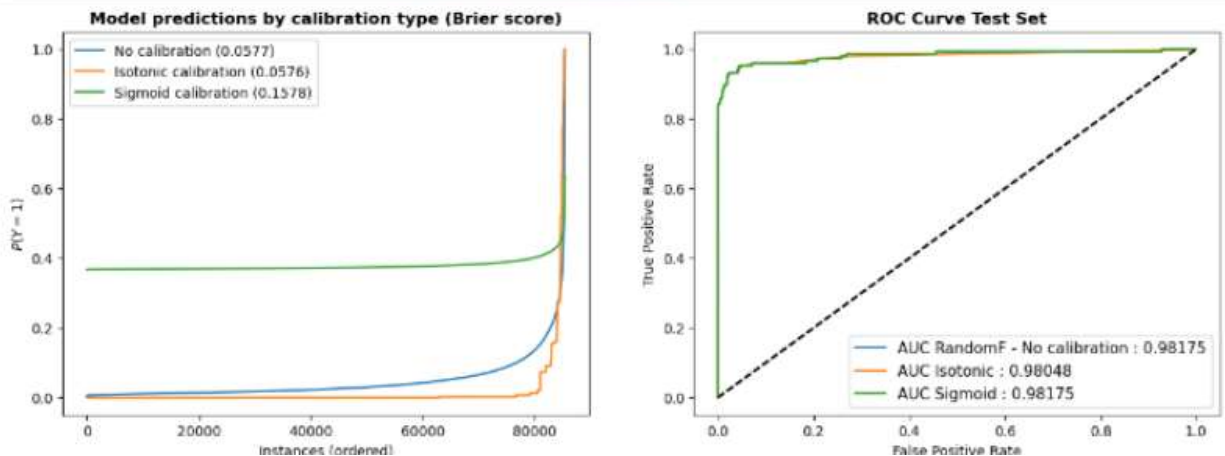Model predictions by calibration type (Brier score) — ROC Curve Validation Set

```
# Store model
loaded_model = searchCV.best_estimator_

# No calibration
preds_clf = loaded_model.predict_proba(X_test)[:, 1]
clf_brier_score = brier_score_loss(y_test, preds_clf, sample_weight=test_res_array.ravel())

# Isotonic calibration
clf_calib_iso, preds_calib_iso, brier_score_iso = calibrate_predictions(loaded_model,
                                                 X_train_calib, y_train_calib,
                                                 X_test, y_test,
                                                 sw_train.ravel(), test_res_array.ravel(),
                                                 'prefit', 'isotonic')

# Sigmoid calibration
clf_calib_sig, preds_calib_sig, brier_score_sig = calibrate_predictions(loaded_model,
                                                 X_train_calib, y_train_calib,
                                                 X_test, y_test,
                                                 sw_train.ravel(), test_res_array.ravel(),
                                                 'prefit', 'sigmoid')

# Plots
fig, ax = plt.subplots(1, 2, figsize=(15, 5))
plot_predictions(preds_clf, 'No calibration', clf_brier_score, ax[0])
plot_predictions(preds_calib_iso, 'Isotonic calibration', brier_score_iso, ax[0])
plot_predictions(preds_calib_sig, 'Sigmoid calibration', brier_score_sig, ax[0])
plot_roc_curve_manual(y_test, preds_clf, 'RandomF - No calibration', ax[1], 'ROC Curve Test Set')
plot_roc_curve_manual(y_test, preds_calib_iso, 'Isotonic', ax[1], 'ROC Curve Test Set')
plot_roc_curve_manual(y_test, preds_calib_sig, 'Sigmoid', ax[1], 'ROC Curve Test Set');
```



Model predictions by calibration type (Brier score) — ROC Curve Test Set

✓ For this project, I investigated several machine learning strategies and acceleration techniques using the highly unbalanced credit card fraud dataset. To address the class imbalance issue indicated in the flow chart, I created an unbalanced dataset by undersampling and oversampling following Random Forest tuning. Regarding the fundamental challenges in the models' operation resulting from the dataset's imbalance, the models achieved excellent outcomes by employing AUC-ROC as their principal evaluation metric.

✓ By using boosting techniques like LightGBM, Catboost, and XGBoost, the model's capability for fraud case categorisation was greatly enhanced. Of them, the LightGBM model was the most notable since it allowed for the achievement of 99.960% of the AUC score. This illustrates how effectively gradient boosting models can train on complex patterns and data with unequal class distributions. Furthermore, approaches for probability calibration such as sigmoid and isotonic calibration were used to produce more accurate probabilities, ensuring that the models had both credible probabilities and strong classification performance.

✓ The results showed that when boosting methods and resampling approaches were used to unbalanced datasets, they performed effectively. But when I attempted to use the same process in this particular instance, the effects of using feature space with PCA were incredibly insignificant. This data set's optimal classification model was determined by utilizing a variety of intricate and advanced modeling techniques in conjunction with appropriate feature engineering and EDA. Ultimately, LightGBM emerged as the most successful model after achieving the best AUC score and demonstrating strong performance in this task.