

Credit Card Default Prediction

Mohammed El-Bathy
Department of Computer Science and Engineering
Oakland University
Rochester, Michigan, USA
elbathy@oakland.edu

Abstract— This project is designed to predict credit card defaults by analyzing a variety of financial and demographic data attributes. The dataset undergoes multiple stages, including data preprocessing, feature engineering, and model training, which culminate in the evaluation and comparison of three machine learning models: Decision Tree, Logistic Regression, and Random Forest. Each model is tested through various evaluation metrics to determine its effectiveness, and the results are carefully documented to support model selection

I. INTRODUCTION

Credit card default prediction is a critical issue for financial institutions that aim to mitigate risk and minimize losses. By leveraging customer data such as payment history, credit limits, and demographic attributes, we can build predictive models to estimate the likelihood of default. This project explores the use of decision tree classifiers, focusing on preparing the data, handling class imbalance, and evaluating model performance.

II. DOMAIN KNOWLEDGE

In the domain of credit risk, **default** refers to the failure to meet the legal obligations of debt repayment. Credit card issuers face significant financial risks due to defaults, so it's crucial to develop predictive models to assess the likelihood of default before issuing credit..

The **features** used in this analysis cover various aspects of customer behavior and demographics, such as credit limit, payment history, bill amounts, and personal attributes (e.g., age, gender, education, and marital status). Understanding these variables is essential to modeling the risk of default

III. DATASET ANALYSIS & UNDERSTANDING

All the headings in the main body are numbered (automatically).

A. Data Characteristics

The dataset contains 30 variables, including both numerical and categorical data. Key variables include:

- **LIMIT_BAL**: The credit limit assigned to the customer.
- **AGE**: The age of the customer
- **PAY_0 to PAY_6**: Repayment status in different months.
- **BILL_AMT1 to BILL_AMT6**: Amount of bill statements from the past six months.
- **PAY_AMT1 to PAY_AMT6**: Amount paid in the past six months.
- **SEX, EDUCATION, MARRIAGE**: Categorical features indicating demographic details

The target variable is **default payment next month**, indicating whether the customer defaulted (1) or not (0)

B. Feature Analysis & Selection

- The key financial variables are credit limit (LIMIT_BAL), bill amounts (BILL_AMT1-6), and payment history (PAY_0-6). These features directly reflect the customer's ability to manage credit and meet payment obligations
- Demographic variables such as AGE, SEX, EDUCATION, and MARRIAGE also provide insights into the customer profile and their associated risk levels

IV. DATA PREPROCESSING FUNCTIONS

To ensure a smooth workflow, the project begins by importing the required libraries, which helps avoid any errors from missing dependencies. The `import_libraries` function accomplishes this initial step, setting the foundation for the project by confirming that all necessary packages are available. In addition we defined global variables.

```
import gradio as gr
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from io import BytesIO
from PIL import Image
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import StratifiedKFold, cross_val_score, GridSearchCV
from sklearn.metrics import confusion_matrix, classification_report, roc_auc_score
from sklearn.preprocessing import StandardScaler
from statsmodels.stats.outliers_influence import variance_inflation_factor
from statsmodels.tools.tools import add_constant
from imblearn.over_sampling import SMOTE
from sklearn import tree
```

```
# Global variables
data = None
encoded_data = None
X_resampled = None
y_resampled = None
dt_model = None
lr_model = None
rf_model = None
lr_results = None
rf_results = None
```

A. Loading Dataset.

The next step is loading the dataset through the `load_dataset` function. This function reads in the credit card client data, providing access to the information necessary for further processing and analysis.

```
# Function 2: Load the Dataset
def load_dataset():
    global data
    file_path = 'default_of_credit_card_clients.csv'
    data = pd.read_csv(file_path, header=1)
    return data
```

B. Statistical Summary.

Once the data is loaded, a statistical summary is generated using `display_statistics`, which gives insight into each feature's characteristics, such as mean, median, and standard deviation. This summary helps in identifying potential issues, such as outliers, and provides a foundational understanding of the dataset's distributions.

```
# Function 3: Display statistical summary
def display_statistics():
    global data
    if data is not None:
        return data.describe().T
    else:
        return "Data not loaded."
```

C. Correlation Matrix

The `display_correlation` function then visualizes the correlation matrix, which reveals relationships between variables. This matrix assists in detecting multicollinearity among features, which is critical for certain models that can be sensitive to correlated inputs.

```
# Function 4: Display correlation matrix
def display_correlation():
    global data
    if data is not None:
        fig, ax = plt.subplots(figsize=(10, 6))
        correlation_matrix = data.corr()
        sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt='.2f', ax=ax)
        ax.set_title('Correlation Matrix of Features')

        buf = BytesIO()
        fig.savefig(buf, format='png')
        buf.seek(0)
        plt.close(fig)
        return Image.open(buf)
    else:
        return "Data not loaded."
```

D. Outlier Handling

Next, the `outlier_detection` function identifies and caps extreme values by setting limits at the 1st and 99th percentiles. This step is essential because outliers can skew the data and negatively impact model performance, particularly in distance-based models.

```
# Function 5: Outlier Detection and Handling
def outlier_detection():
    global data
    if data is not None:
        numeric_columns = ['LIMIT_BAL', 'AGE', 'BILL_AMT1', 'BILL_AMT2', 'BILL_AMT3',
                           'BILL_AMT4', 'BILL_AMT5', 'BILL_AMT6', 'PAY_AMT1',
                           'PAY_AMT2', 'PAY_AMT3', 'PAY_AMT4', 'PAY_AMT5', 'PAY_AMT6']

        for column in numeric_columns:
            upper_limit = data[column].quantile(0.99)
            lower_limit = data[column].quantile(0.01)
            data[column] = np.clip(data[column], lower_limit, upper_limit)
        return "Outliers capped at 1st and 99th percentiles."
    else:
        return "Data not loaded."
```

E. Encoding Categorical Variables

The `display_correlation` function then visualizes the correlation matrix, which reveals relationships between variables. This matrix assists in detecting multicollinearity

among features, which is critical for certain models that can be sensitive to correlated inputs.

```
# Function 6: Encoding Categorical Variables
def encoding_categorical():
    global data, encoded_data
    if data is not None:
        encoded_data = pd.get_dummies(data, columns=['SEX', 'EDUCATION', 'MARRIAGE'],
                                       drop_first=True)
        return encoded_data
    else:
        return "Data not loaded."
```

F. Feature Scaling

```
# Function 7: Feature Scaling
def feature_scaling():
    global encoded_data
    if encoded_data is not None:
        continuous_columns = ['LIMIT_BAL', 'AGE', 'PAY_AMT1', 'PAY_AMT2', 'PAY_AMT3',
                              'PAY_AMT4', 'PAY_AMT5', 'PAY_AMT6']

        scaler = StandardScaler()
        encoded_data[continuous_columns] = scaler.fit_transform(encoded_data[continuous_columns])
        return encoded_data
    else:
        return "Encoded data not available."
```

G. SMOTE

To address class imbalance, which is a common issue in binary classification, `apply_smote` applies Synthetic Minority Over-sampling Technique (SMOTE) to the dataset. This technique balances the data by oversampling the minority class, thereby improving the model's ability to learn accurate decision boundaries for both classes.

```
# Function 8: Apply SMOTE
def apply_smote():
    global encoded_data, X_resampled, y_resampled
    if encoded_data is not None:
        X = encoded_data.drop(columns=['default payment next month'])
        y = encoded_data['default payment next month']
        smote = SMOTE(random_state=42)
        X_resampled, y_resampled = smote.fit_resample(X, y)
        return f"SMOTE applied: Resampled dataset shape: {X_resampled.shape}"
    else:
        return "Encoded data not available."
```

H. Variance Inflation Factor

Subsequently, we calculate the Variance Inflation Factor assesses multicollinearity among the features by calculating the VIF for each feature. High VIF values indicate redundancy, which may degrade model performance, particularly in linear models. After any necessary modifications, `recalculate_vif` is used to verify that multicollinearity issues have been resolved, thereby ensuring model stability.

```
# Function 9: Calculate VIF
def calculate_vif():
    global encoded_data
    if encoded_data is not None:
        numeric_data = encoded_data.select_dtypes(include=[np.number])
        X_vif = add_constant(numeric_data)
        vif_data = pd.DataFrame()
        vif_data['feature'] = X_vif.columns
        vif_data['VIF'] = [variance_inflation_factor(X_vif.values, i) for i in range(X_vif.shape[1])]
        return vif_data
    else:
        return "Encoded data not available."
```

I. Recalculate Variance Inflation Factor

After any necessary modifications, `recalculate_vif` is used to verify that multicollinearity issues have been resolved, thereby ensuring model stability.

```
# Function 10: Recalculate VIF
def recalculate_vif():
    global encoded_data
    if encoded_data is not None:
        numeric_data = encoded_data.select_dtypes(include=[np.number])
        X_vif_updated = add_constant(numeric_data)
        vif_data_updated = pd.DataFrame()
        vif_data_updated['feature'] = X_vif_updated.columns
        vif_data_updated['VIF'] = [variance_inflation_factor(X_vif_updated.values, i) for i in range(X_vif_updated.shape[1])]
        return vif_data_updated
    else:
        return "Encoded data not available."
```

J. Visualize Distribution

To further explore the dataset, visualize distribution generates plots that display the distribution of key features. This visualization provides additional insights into data patterns, helping identify skewness or other anomalies that might necessitate transformation or further preprocessing.

```
# Function 13: Visualize Distribution
def visualize_distribution():
    global data, encoded_data
    if data is not None and encoded_data is not None:
        fig, axes = plt.subplots(2, 3, figsize=(10, 6))
        sns.histplot(data['LIMIT_BAL'], bins=30, kde=True, ax=axes[0, 0]).set_title('Distribution of Credit Limit')
        sns.histplot(data['AGE'], bins=30, kde=True, ax=axes[0, 1]).set_title('Distribution of Age')
        sns.countplot(x='SEX', data=encoded_data, ax=axes[0, 2]).set_title('Distribution of Gender')
        sns.countplot(x='EDUCATION', data=encoded_data, ax=axes[1, 0]).set_title('Distribution of Education')
        sns.countplot(x='MARRIAGE', data=encoded_data, ax=axes[1, 1]).set_title('Distribution of Marital Status')
        sns.countplot(x='default payment next month', data=encoded_data, ax=axes[1, 2]).set_title('Default vs Non-Default')
        plt.tight_layout()

        buf = BytesIO()
        fig.savefig(buf, format='png')
        buf.seek(0)
        plt.close(fig)
        return Image.open(buf)
    else:
        return "Data not loaded."
```

V. MODEL TRAINING AND EVALUATION FUNCTIONS

A. Decision Tree Model

1) The first machine learning model employed in this project is a Decision Tree. The initialize_decision_tree function initializes and trains a Decision Tree classifier on the resampled dataset, establishing a straightforward and interpretable baseline model.

```
# Function 14: Initialize and Train Decision Tree
def initialize_decision_tree():
    global dt_model, X_resampled, y_resampled
    if X_resampled is not None and y_resampled is not None:
        dt_model = DecisionTreeClassifier(random_state=42)
        dt_model.fit(X_resampled, y_resampled)
        return "Decision Tree Classifier initialized and trained."
    else:
        return "SMOTE resampled data not available."
```

2) To evaluate the Decision Tree, cross_validation_decision_tree performs cross-validation, calculating accuracy and ROC-AUC scores across folds. Cross-validation reduces the risk of overfitting and provides a reliable estimate of the model's generalization capabilities.

```
# Function 15: K-Fold Cross-Validation for Decision Tree
def cross_validation_decision_tree():
    global dt_model, X_resampled, y_resampled
    if dt_model is not None and X_resampled is not None:
        kf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
        accuracy_scores = cross_val_score(dt_model, X_resampled, y_resampled, cv=kf, scoring='accuracy')
        roc_auc_scores = cross_val_score(dt_model, X_resampled, y_resampled, cv=kf, scoring='roc_auc')
        return f"Mean Accuracy: {accuracy_scores.mean():.4f}, Std Dev: {accuracy_scores.std():.4f}\n" \
               f"Mean ROC-AUC: {roc_auc_scores.mean():.4f}, Std Dev: {roc_auc_scores.std():.4f}"
    else:
        return "Decision Tree or resampled data not available."
```

3) To enhance interpretability, visualize_decision_tree produces a graphical representation of the tree structure, displaying the decision-making paths and feature splits.

```
# Function 16: Visualize Decision Tree Structure
def visualize_decision_tree():
    global dt_model, X_resampled
    if dt_model is not None and X_resampled is not None:
        plt.figure(figsize=(20, 10))
        tree.plot_tree(dt_model, feature_names=X_resampled.columns, class_names=['No Default', 'Default'], filled=True)
        buf = BytesIO()
        plt.savefig(buf, format='png')
        buf.seek(0)
        plt.close()
        return Image.open(buf)
    else:
        return "Decision Tree model not trained."
```

4) The confusion matrix for the Decision Tree, generated by visualize_confusion_matrix_decision_tree, offers a detailed view of true versus predicted classifications.

This matrix highlights the model's accuracy on each class, helping to identify any biases in prediction.

```
# Function 17: Confusion Matrix Visualization for Decision Tree
def visualize_confusion_matrix_decision_tree():
    global dt_model, X_resampled, y_resampled
    if dt_model is not None and X_resampled is not None:
        plt.figure(figsize=(8, 4))
        sns.heatmap(confusion_matrix(y_resampled, dt_model.predict(X_resampled)), annot=True, fmt='d', cmap='Blues')
        buf = BytesIO()
        plt.savefig(buf, format='png')
        buf.seek(0)
        plt.close()
        return Image.open(buf)
    else:
        return "Decision Tree or resampled data not available."
```

5) Finally, classification_report_dt provides a classification report with precision, recall, and F1-score metrics for the Decision Tree. This report allows for a comprehensive evaluation of the model's performance across each class, especially in imbalanced datasets.

```
# Function 18: Classification Report for Decision Tree
def classification_report_dt():
    global dt_model, X_resampled, y_resampled
    if dt_model is not None and X_resampled is not None:
        y_pred = dt_model.predict(X_resampled)
        report = classification_report(y_resampled, y_pred)
        return report
    else:
        return "Decision Tree or resampled data not available."
```

B. Logistic Regression Model

1) Logistic Regression is the second model used in this analysis. The initialize_logistic_regression function initializes and trains the Logistic Regression model on the resampled data, providing a robust, interpretable model that serves as a benchmark.

```
# Function 19: Initialize and Train Logistic Regression
def initialize_logistic_regression():
    global lr_model, X_resampled, y_resampled
    if X_resampled is not None and y_resampled is not None:
        lr_model = LogisticRegression(max_iter=1000, random_state=42)
        lr_model.fit(X_resampled, y_resampled)
        return "Logistic Regression model initialized and trained."
    else:
        return "SMOTE resampled data not available."
```

2) Evaluation of the Logistic Regression model is conducted through evaluate_logistic_regression, which generates a classification report and calculates the ROC-AUC score. These metrics offer a thorough view of the model's ability to discriminate between classes and provide a basis for comparing it to other models.

```
# Function 20: Logistic Regression Evaluation
def evaluate_logistic_regression():
    global lr_model, X_resampled, y_resampled
    if lr_model is not None and X_resampled is not None:
        y_pred = lr_model.predict(X_resampled)
        report = classification_report(y_resampled, y_pred)
        auc = roc_auc_score(y_resampled, lr_model.predict_proba(X_resampled)[:, 1])
        return f"Classification Report:\n{report}\nROC-AUC Score: {auc:.4f}"
    else:
        return "Logistic Regression model or resampled data not available."
```

C. Random Forest Model

1) The third and final model in this project is a Random Forest, which is both flexible and powerful for handling complex data patterns. The initialize_random_forest function initializes and tunes the Random Forest model using GridSearchCV, optimizing hyperparameters based on ROC-

AUC scores. This tuning process enhances the model's performance by selecting the best combination of parameters.

```
# Function 21: Initialize and Train Random Forest with Hyperparameter Tuning
def initialize_random_forest():
    global rf_model, X_resampled, y_resampled, rf_results
    if X_resampled is not None and y_resampled is not None:
        rf_model = RandomForestClassifier(random_state=42)
        param_grid = {
            'n_estimators': [50, 100, 200],
            'max_depth': [None, 10, 20, 30],
            'min_samples_split': [2, 5, 10]
        }
        grid_search = GridSearchCV(rf_model, param_grid, cv=5, scoring='roc_auc', n_jobs=-1)
        grid_search.fit(X_resampled, y_resampled)

        rf_model = grid_search.best_estimator_
        rf_results = {
            "Best Parameters": grid_search.best_params_,
            "Best ROC-AUC": grid_search.best_score_
        }
        return ("Random Forest model initialized and trained with tuning. Best ROC-AUC: (rf_results['Best ROC-AUC']):.4f")
    else:
        return "SMOTE resampled data not available."
```

2) The evaluation function evaluate_random_forest then generates a classification report and ROC-AUC score, which reveal the model's performance in distinguishing between classes. This assessment is critical for comparing the Random Forest's predictive power against the other models.

```
# Function 22: Random Forest Evaluation
def evaluate_random_forest():
    global rf_model, X_resampled, y_resampled
    if rf_model is not None and X_resampled is not None:
        y_pred = rf_model.predict(X_resampled)
        report = classification_report(y_resampled, y_pred)
        auc = roc_auc_score(y_resampled, rf_model.predict_proba(X_resampled)[:, 1])
        return f"Classification Report:\n{report}\nROC-AUC Score: {auc:.4f}"
    else:
        return "Random Forest model or resampled data not available."
```

VI. GARDIO USER INTERFACE

```
# Gardio Interface Setup
with gr.Blocks() as credit_default_interface:
    with gr.Tab("Import Libraries"):
        gr.Interface(fn=import_libraries, inputs=[], outputs="text").render()
    with gr.Tab("Load Dataset"):
        gr.Interface(fn=load_dataset, inputs=[], outputs="dataframe").render()
    with gr.Tab("Display Statistics"):
        gr.Interface(fn=display_statistics, inputs=[], outputs="dataframe").render()
    with gr.Tab("Display Correlation"):
        gr.Interface(fn=display_correlation, inputs=[], outputs=gr.Image(type="pil")).render()
    with gr.Tab("Outlier Detection"):
        gr.Interface(fn=outlier_detection, inputs=[], outputs="text").render()
    with gr.Tab("Encoding Categorical Variables"):
        gr.Interface(fn=encoding_categorical, inputs=[], outputs="dataframe").render()
    with gr.Tab("Feature Scaling"):
        gr.Interface(fn=feature_scaling, inputs=[], outputs="dataframe").render()
    with gr.Tab("Apply SMOTE"):
        gr.Interface(fn=apply_smote, inputs=[], outputs="text").render()
    with gr.Tab("Calculate VIF"):
        gr.Interface(fn=calculate_vif, inputs=[], outputs="dataframe").render()
    with gr.Tab("Recalculate VIF"):
        gr.Interface(fn=recalculate_vif, inputs=[], outputs="dataframe").render()
    with gr.Tab("Visualize Distribution"):
        gr.Interface(fn=visualize_distribution, inputs=[], outputs=gr.Image(type="pil")).render()
    with gr.Tab("Initialize and Train Decision Tree"):
        gr.Interface(fn=initialize_decision_tree, inputs=[], outputs="text").render()
    with gr.Tab("Decision Tree Cross-Validation"):
        gr.Interface(fn=cross_validation_decision_tree, inputs=[], outputs="text").render()
    with gr.Tab("Visualize Decision Tree"):
        gr.Interface(fn=visualize_decision_tree, inputs=[], outputs=gr.Image(type="pil")).render()
    with gr.Tab("Decision Tree Classification Report"):
        gr.Interface(fn=classification_report_dt, inputs=[], outputs="text").render()
    with gr.Tab("Initialize and Train Logistic Regression"):
        gr.Interface(fn=initialize_logistic_regression, inputs=[], outputs="text").render()
    with gr.Tab("Evaluate Logistic Regression"):
        gr.Interface(fn=evaluate_logistic_regression, inputs=[], outputs="text").render()
    with gr.Tab("Initialize and Train Random Forest"):
        gr.Interface(fn=initialize_random_forest, inputs=[], outputs="text").render()
    with gr.Tab("Evaluate Random Forest"):
        gr.Interface(fn=evaluate_random_forest, inputs=[], outputs="text").render()

# Launch Interface
credit_default_interface.launch(share=True)
```

VII. CONCLUSION

This project successfully built a **Credit Card Default Prediction** model using decision trees, providing valuable insights into the factors that contribute to credit default risk. The preprocessing steps, including outlier handling, one-hot encoding, feature scaling, and SMOTE, ensured that the model was well-prepared to handle the data complexities. While the model performed moderately well, future improvements could include hyperparameter tuning and trying other models like **Random Forest** or **Gradient Boosting** for potentially better accuracy.

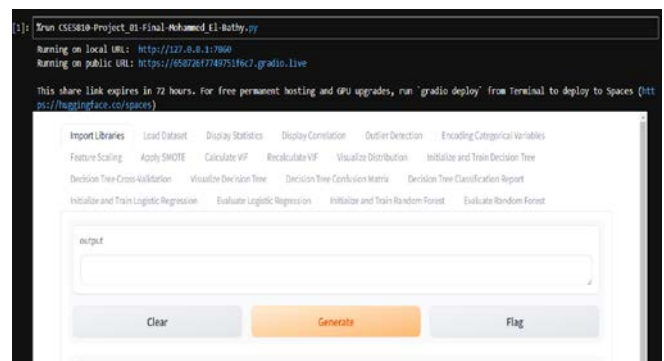
By analyzing the **confusion matrix** and the **classification report**, we can identify areas where the model's performance could be improved, particularly in reducing false positives and false negatives.

This solution is highly interpretable, making it a useful tool for financial institutions to assess and mitigate credit risk.

APPENDIX A

The Jupyter Notebook, the Python file and this report with the data file are uploaded are also available at the URL <https://github.com/MohammedEl-Bath/CreditCardDefaultPrediction>

STEPS TO RUN NOTEBOOK AND GRADIO APP



REFERENCES

- [1] Breiman, L. (2001). Random forests. Machine learning, 45(1), 5-32.
- [2] Chawla, N. V., et al. (2002). SMOTE: Synthetic Minority Over-sampling Technique. Journal of Artificial Intelligence Research..
- [3] Handling Imbalanced Datasets. Towards Data Science. Retrieved from <https://towardsdatascience.com/handling-imbalanced-datasets-in-machine-learning-7a0e84220f28>.

