

# **FA 590 Statistical Learning in Finance**

## **Assignment 2: Loan Default and Return Prediction Using Linear Models**

**Federica Malamisura**

**2024 Fall**

## **Abstract**

In this assignment, we analyze a dataset containing over one million personal loans originated between 2008 and 2017 on LendingClub. The dataset includes detailed loan characteristics, borrower characteristics, and loan outcomes. The tasks include Exploratory Data Analysis (EDA), Data Preprocessing, Building Regression and Classification Models using Linear Models, Feature Importance Analysis using SHAP, and Making Predictions on the Test Set. We utilize linear regression and logistic regression models to predict loan returns and early defaults, respectively, and assess feature importance to understand the factors influencing these loan outcomes. The models are evaluated using appropriate metrics, and predictions are made on the test set for submission.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Importing Libraries and Loading Data</b>	<b>1</b>
2.1	Import Libraries . . . . .	1
2.2	Loading Data . . . . .	1
<b>3</b>	<b>Exploratory Data Analysis</b>	<b>2</b>
3.1	Understanding the Data Structure . . . . .	2
3.2	Distribution of Target Variables . . . . .	2
3.3	Correlation Analysis . . . . .	2
3.4	Additional Insights . . . . .	3
<b>4</b>	<b>Data Preprocessing</b>	<b>3</b>
4.1	Sampling Strategy: Expanding-Window Approach . . . . .	3
4.2	Handling Missing Values . . . . .	4
4.3	Feature Engineering . . . . .	4
<b>5</b>	<b>Model Building</b>	<b>5</b>
5.1	Regression: Predicting Loan Returns . . . . .	5
5.1.1	Advantages of Using MSE for This Dataset . . . . .	6
5.1.2	4. Specific Justifications for This Loan Return Prediction Task . . . . .	7
5.2	Classification: Predicting Early Defaults . . . . .	8
<b>6</b>	<b>Feature Importance</b>	<b>9</b>
6.1	Feature Importance Using SHAP Values . . . . .	9
<b>7</b>	<b>Predicting on the Test Set</b>	<b>10</b>
7.1	Preparing the Test Data . . . . .	10
7.2	Making Predictions . . . . .	11
7.3	Preparing the Submission File . . . . .	12
<b>8</b>	<b>Conclusion</b>	<b>12</b>

# 1 Introduction

In this assignment, we focus on analyzing and predicting loan defaults and returns using data from LendingClub. Accurately predicting these outcomes is crucial for financial institutions to manage risk effectively and make informed lending decisions. The dataset comprises over one million personal loans issued between 2008 and 2017, providing a rich source of information for analysis.

We perform Exploratory Data Analysis to understand the data structure and identify key relationships. Data preprocessing steps include handling missing values, feature engineering, and encoding categorical variables. We employ linear regression models to predict loan returns and logistic regression models for predicting early defaults. Feature importance is assessed using SHAP values, allowing us to interpret the models and identify the most influential features. Finally, we make predictions on the test set and prepare the results for submission.

## 2 Importing Libraries and Loading Data

### 2.1 Import Libraries

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split, GridSearchCV,
    TimeSeriesSplit
6 from sklearn.linear_model import LinearRegression, Lasso, Ridge, ElasticNet,
    LogisticRegression
7 from sklearn.metrics import mean_squared_error, r2_score, accuracy_score,
    f1_score, roc_auc_score, classification_report, confusion_matrix
8 import shap
9 from sklearn.preprocessing import StandardScaler
```

Listing 1: Importing Libraries

### 2.2 Loading Data

```
1 # Load training data
2 loan_df = pd.read_csv('lc_loan.csv')
3
4 # Load test data
5 loan_test_df = pd.read_csv('lc_loan_test.csv')
6
7 # Load submission template
8 submission_df = pd.read_csv('submission_test.csv')
```

Listing 2: Loading the Dataset

We start by loading the required libraries for data manipulation, visualization, modeling, and interpretability. The dataset is loaded into pandas DataFrames for both the training and test sets. We also load the submission template which will be used later to prepare our predictions for submission.

## 3 Exploratory Data Analysis

### 3.1 Understanding the Data Structure

```
1 # Display summary statistics for key variables
2 key_vars = ['loan_amnt', 'int_rate', 'fico_range_high', 'annual_inc']
3 loan_df[key_vars].describe()
```

Listing 3: Exploring the Dataset

We explore the dataset to understand its structure, the types of variables it contains, and the presence of missing values.

```
1 # Check data types and missing values
2 loan_df.info()
3 print('\nMissing values in training data:')
4 print(loan_df.isnull().sum())
```

Listing 4: Data Information and Missing Values

### 3.2 Distribution of Target Variables

```
1 # Distribution of 'early_default'
2 sns.countplot(x='early_default', data=loan_df)
3 plt.title('Distribution of Early Default')
4 plt.show()
```

Listing 5: Distribution of Early Default

```
1 # Distribution of 'return'
2 sns.histplot(loan_df['return'], kde=True)
3 plt.title('Distribution of Loan Returns')
4 plt.show()
```

Listing 6: Distribution of Loan Returns

We examine the distribution of the target variables `early_default` and `return` to gain insights into their characteristics.

### 3.3 Correlation Analysis

```
1 # Correlation matrix of numerical features
2 corr = loan_df.select_dtypes(include=['float64', 'int64']).corr()
3
4 # Heatmap of the correlation matrix
5 plt.figure(figsize=(12, 10))
6 sns.heatmap(corr, annot=True, fmt=".2f", cmap='coolwarm')
7 plt.title('Correlation Matrix')
8 plt.show()
```

Listing 7: Correlation Matrix

We generate a correlation matrix for numerical features to identify potential relationships between variables.

## 3.4 Additional Insights

```
1 # Relationship between FICO score and interest rate
2 plt.figure(figsize=(8,6))
3 sns.scatterplot(x='fico_range_high', y='int_rate', data=loan_df, alpha=0.3)
4 plt.title('FICO Score vs Interest Rate')
5 plt.show()
```

Listing 8: FICO Score vs. Interest Rate

An exploratory analysis of the relationship between the borrower's FICO score and the interest rate of the loan is conducted.

## 4 Data Preprocessing

### 4.1 Sampling Strategy: Expanding-Window Approach

In time series data, it is essential to consider the temporal order of observations when splitting data into training and testing sets. I chose to implement an **expanding window approach** for the following reasons:

1. **Data Utilization:** The expanding window approach utilizes all available historical data up to the current point for training. This means that as time progresses, the model is trained on an increasing amount of data, potentially capturing more comprehensive patterns.
2. **Real-World Scenario Representation:** In financial contexts, especially in credit risk modeling, it is realistic to assume that models are updated as more data becomes available. The expanding window mirrors this practice.
3. **Model Stability:** Using more data for training can lead to more stable and robust models, reducing variance and improving generalization.
4. **Avoiding Data Leakage:** Since future data is not used in training past predictions, the expanding window prevents data leakage, maintaining the integrity of model evaluation.

In contrast, a rolling window approach uses a fixed-size window that moves forward in time, potentially discarding valuable historical data. While rolling windows are useful for non-stationary data where recent observations are more relevant, the expanding window is more appropriate when historical data remains informative over time.

```
1 # Sort data by origination date to ensure temporal ordering
2 loan_df['year'] = pd.to_datetime(loan_df['year'], format='%Y')
3 loan_df = loan_df.sort_values('year')
4
5 # Define number of splits for expanding window
6 n_splits = 5
7 tscv = TimeSeriesSplit(n_splits=n_splits)
8 print(f'Number of splits: {n_splits}')
```

Listing 9: Defining Expanding Window Splits

## 4.2 Handling Missing Values

I handled missing values to prevent biases and inaccuracies in our models. For numerical features, we use the median value for imputation, and for categorical features, we use the mode.

```
1 # Summarize missing values
2 missing_values = loan_df.isnull().sum()
3 missing_values = missing_values[missing_values > 0].sort_values(ascending=
    False)
4 print('Missing values in training data:')
5 print(missing_values)
6
7 # Impute missing values
8 for column in missing_values.index:
9     if loan_df[column].dtype in ['float64', 'int64']:
10         median = loan_df[column].median()
11         loan_df[column].fillna(median, inplace=True)
12     else:
13         mode = loan_df[column].mode()[0]
14         loan_df[column].fillna(mode, inplace=True)
15
16 # Verify no missing values remain
17 print('\nMissing values after imputation:')
18 print(loan_df.isnull().sum().sum())
```

Listing 10: Handling Missing Values

## 4.3 Feature Engineering

I performed feature engineering to enhance the predictive power of our models. This includes creating new features and transforming existing ones.

```
1 # Create debt-to-income ratio
2 loan_df['debt_to_income'] = loan_df['loan_amnt'] / loan_df['annual_inc']
3 loan_test_df['debt_to_income'] = loan_test_df['loan_amnt'] / loan_test_df['
    annual_inc']
4
5 # Log transformation for skewed features
6 loan_df['log_loan_amount'] = np.log1p(loan_df['loan_amnt'])
7 loan_test_df['log_loan_amount'] = np.log1p(loan_test_df['loan_amnt'])
8
9 # Encode categorical variables
10 categorical_cols = loan_df.select_dtypes(include=['object']).columns
11
12 for col in categorical_cols:
13     if col in loan_test_df.columns:
14         loan_df[col] = loan_df[col].astype('category').cat.codes
15         loan_test_df[col] = loan_test_df[col].astype('category').cat.codes
16     else:
17         loan_df[col] = loan_df[col].astype('category').cat.codes
18
19 # Define feature list
20 feature_cols = [col for col in loan_df.columns if col not in ['id', '
    early_default', 'return', 'loan_date', 'loan_status', 'year']]
```

```

21
22 # Update feature_cols to include only columns present in both training and
    test data
23 train_features = set(loan_df.columns)
24 test_features = set(loan_test_df.columns)
25 common_features = train_features.intersection(test_features)
26 feature_cols = [col for col in feature_cols if col in common_features]

```

Listing 11: Feature Engineering

## 5 Model Building

### 5.1 Regression: Predicting Loan Returns

I implemented and evaluated several linear regression models:

- **Linear Regression:** Serves as a baseline model.
- **Ridge Regression:** Addresses multicollinearity by introducing L2 regularization.
- **Lasso Regression:** Performs feature selection through L1 regularization.
- **ElasticNet:** Combines L1 and L2 regularization.

```

1 # Define target and features for regression
2 X = loan_df[feature_cols]
3 y_return = loan_df['return']
4
5 # Standardize features
6 X = X.select_dtypes(include=[np.number]) # Ensure only numeric data
7 scaler = StandardScaler()
8 X_scaled = scaler.fit_transform(X)
9
10 # Split data using expanding-window approach
11 for train_index, val_index in tscv.split(X_scaled):
12     X_train, X_val = X_scaled[train_index], X_scaled[val_index]
13     y_train, y_val = y_return.iloc[train_index], y_return.iloc[val_index]
14
15 # Hyperparameter tuning with GridSearchCV
16 ridge = Ridge()
17 lasso = Lasso()
18 elastic = ElasticNet()
19
20 # Define parameter grid
21 param_grid = {'alpha': [0.01, 0.1, 1, 10, 100]}
22
23 # Ridge Regression
24 grid_ridge = GridSearchCV(ridge, param_grid, cv=5, scoring='
    neg_mean_squared_error', n_jobs=-1)
25 grid_ridge.fit(X_train, y_train)
26 best_ridge = grid_ridge.best_estimator_

```



```

27
28 # Lasso Regression
29 grid_lasso = GridSearchCV(lasso, param_grid, cv=5, scoring='
    neg_mean_squared_error', n_jobs=-1)
30 grid_lasso.fit(X_train, y_train)
31 best_lasso = grid_lasso.best_estimator_
32
33 # ElasticNet
34 param_grid_elastic = {'alpha': [0.01, 0.1, 1, 10], 'l1_ratio': [0.2, 0.5,
    0.8]}
35 grid_elastic = GridSearchCV(elastic, param_grid_elastic, cv=5, scoring='
    neg_mean_squared_error', n_jobs=-1)
36 grid_elastic.fit(X_train, y_train)
37 best_elastic = grid_elastic.best_estimator_
38
39 # Evaluate models
40 models = {'Ridge': best_ridge, 'Lasso': best_lasso, 'ElasticNet':
    best_elastic}
41 validation_results = {} # To store results for each model
42
43 for name, model in models.items():
44     y_pred = model.predict(X_val)
45     r2 = r2_score(y_val, y_pred)
46     mse = mean_squared_error(y_val, y_pred)
47     print(f'{name} on validation set:')
48     print(f'    R-squared: {r2:.4f}')
49     print(f'    MSE: {mse:.4f}')
50     # Store MSE for model selection
51     validation_results[name] = {'model': model, 'mse': mse}
52
53 # Select the best model based on validation performance (lowest MSE)
54 best_model = None
55 best_mse = float('inf')
56
57 for name, result in validation_results.items():
58     mse = result['mse']
59     if mse < best_mse:
60         best_mse = mse
61         best_model = result['model']
62
63 print(f'Best model: {best_model} with MSE: {best_mse:.4f}')

```

Listing 12: Regression Modeling

### 5.1.1 Advantages of Using MSE for This Dataset

#### A. Absolute Measure of Prediction Error:

- **Direct Interpretation:** MSE provides an absolute measure of the average squared prediction error, directly indicating how far the model's predictions deviate from the actual values.
- **Same Units (Squared):** Although in squared units, MSE is interpretable in the context of the target variable's units, making it useful for understanding the magnitude of errors.

## B. Sensitivity to Large Errors:

- **Penalizes Large Deviations:** By squaring the errors, MSE places a higher penalty on larger errors, which is crucial in financial contexts where large prediction errors can have significant consequences.
- **Risk Management:** Prioritizing models with lower MSE can help in selecting models that minimize the risk of large losses.

## C. Model Comparison Across Different Models:

- **Consistency:** MSE allows for consistent comparison across different models, even if the models use different transformations or have different complexities.
- **Hyperparameter Tuning:** Using MSE aligns the evaluation metric with the loss function optimized during model training (e.g., in Ridge, Lasso, and ElasticNet regression), facilitating more effective hyperparameter tuning.

## D. Focus on Prediction Accuracy:

- **Practical Relevance:** In financial applications, the practical goal is often to minimize prediction errors rather than maximizing the explained variance.
- **Actionable Insights:** Lower MSE directly translates to more accurate predictions, which is critical for decision-making processes like loan pricing, risk assessment, and investment strategies.

### 5.1.2 4. Specific Justifications for This Loan Return Prediction Task

#### A. Financial Impact of Prediction Errors:

- **Loan Returns:** Predicting loan returns accurately is vital for assessing profitability and risk. Large errors in prediction can lead to incorrect investment decisions.
- **Emphasis on Error Magnitude:** Using MSE ensures that the model selection process emphasizes minimizing large prediction errors that could have significant financial repercussions.

#### B. Alignment with Business Objectives:

- **Risk-Averse Strategies:** Financial institutions often adopt risk-averse strategies, preferring models that minimize the potential for large losses.
- **Model Evaluation Criteria:** Selecting models based on MSE aligns with the goal of minimizing financial risk, as it prioritizes models with smaller average squared errors.

#### C. Suitability for Imbalanced and Noisy Data:

- **Handling Outliers:** Financial datasets may contain outliers due to extreme market movements or anomalies.
- **Robustness:** While MSE is sensitive to outliers, in this context, we want to be aware of models' performance concerning extreme values. This sensitivity encourages the selection of models that perform better even in the presence of outliers.

## 5.2 Classification: Predicting Early Defaults

I used logistic regression to predict early defaults. Due to computational considerations, we sample 10% of the dataset. By performing hyperparameter tuning and threshold optimization to maximize the F1 score.

```
1 # Sample 10% of the dataset to speed up computation
2 loan_df_sampled = loan_df.sample(frac=0.1, random_state=42)
3
4 # Define target and features for default prediction using the sampled data
5 X_default = loan_df_sampled[feature_cols]
6 X_default = X_default.select_dtypes(include=[np.number])
7 y_default = loan_df_sampled['early_default']
8
9 # Standardize features using the same scaler
10 X_default_scaled = scaler.transform(X_default)
11
12 # Define the model with a reduced parameter grid
13 logreg = LogisticRegression()
14 param_grid = {
15     'C': [0.1, 1, 10],          # Reduced number of C values
16     'penalty': ['l1', 'l2'],    # Penalty types
17     'solver': ['liblinear']     # 'liblinear' solver supports 'l1' and 'l2'
18 }
19
20 # Define custom scorer to maximize F1 score
21 def custom_scorer(estimator, X, y):
22     y_pred_prob = estimator.predict_proba(X)[:, 1]
23     thresholds = np.arange(0.1, 0.9, 0.1)
24     best_f1 = -np.inf
25     for thresh in thresholds:
26         y_pred = (y_pred_prob >= thresh).astype(int)
27         f1 = f1_score(y, y_pred)
28         if f1 > best_f1:
29             best_f1 = f1
30     return best_f1
31
32 fl_scorer = make_scorer(custom_scorer, needs_proba=False, greater_is_better=
    True)
33
34 # Perform GridSearchCV with fewer folds (e.g., cv=3)
35 grid_search = GridSearchCV(logreg, param_grid, cv=3, scoring=fl_scorer, n_jobs
    =-1, verbose=1)
36 grid_search.fit(X_default_scaled, y_default)
37
38 # Get the best model
39 best_logreg = grid_search.best_estimator_
40
41 # Predict probabilities on the entire sampled dataset
42 y_pred_prob = best_logreg.predict_proba(X_default_scaled)[:, 1]
43
44 # Tune threshold to maximize F1 score on the sampled data
45 thresholds = np.arange(0.1, 0.9, 0.1)
46 best_f1 = -np.inf
```

```

47 best_thresh = 0.5
48
49 for thresh in thresholds:
50     y_pred = (y_pred_prob >= thresh).astype(int)
51     f1 = f1_score(y_default, y_pred)
52     if f1 > best_f1:
53         best_f1 = f1
54         best_thresh = thresh
55
56 print(f'Final threshold set to: {best_thresh:.2f} with F1 score: {best_f1:.4f}
    ')
57
58 # Define a function to apply the threshold
59 def apply_threshold(probs, threshold):
60     return (probs >= threshold).astype(int)
61
62 # Evaluate on the sampled training set
63 y_train_pred = apply_threshold(y_pred_prob, best_thresh)
64 train_accuracy = accuracy_score(y_default, y_train_pred)
65 train_f1 = f1_score(y_default, y_train_pred)
66 train_auc = roc_auc_score(y_default, y_pred_prob)
67
68 print('\nIn-sample Performance for Early Default Prediction:')
69 print(f'Accuracy: {train_accuracy:.4f}')
70 print(f'F1 Score: {train_f1:.4f}')
71 print(f'AUC Score: {train_auc:.4f}')

```

Listing 13: Classification Modeling

## 6 Feature Importance

### 6.1 Feature Importance Using SHAP Values

We use SHAP (SHapley Additive exPlanations) values to interpret our models and identify the most significant features influencing the predictions.

```

1 # Convert scaled arrays back to DataFrames with column names
2 X_scaled_df = pd.DataFrame(X_scaled, columns=X.columns)
3 X_default_scaled_df = pd.DataFrame(X_default_scaled, columns=X_default.columns
    )
4
5 # SHAP for Regression Model
6 explainer_reg = shap.Explainer(best_model, X_scaled_df)
7 shap_values_reg = explainer_reg(X_scaled_df)
8
9 # Plot feature importance for regression
10 shap.summary_plot(shap_values_reg, features=X_scaled_df, plot_type='bar',
    max_display=10)
11
12 # Get top 10 features for regression
13 shap_importance_reg = pd.DataFrame({
14     'feature': X_scaled_df.columns,

```

```

15     'importance': np.abs(shap_values_reg.values).mean(axis=0)
16 })
17 shap_importance_reg = shap_importance_reg.sort_values(by='importance',
18     ascending=False).head(10)
19 print('Top 10 Important Features for Regression:')
20 print(shap_importance_reg)
21
22 # SHAP for Classification Model
23 explainer_clf = shap.Explainer(best_logreg, X_default_scaled_df)
24 shap_values_clf = explainer_clf(X_default_scaled_df)
25
26 # Plot feature importance for classification
27 shap.summary_plot(shap_values_clf, features=X_default_scaled_df, plot_type='
28     bar', max_display=10)
29
30 # Get top 10 features for classification
31 shap_importance_clf = pd.DataFrame({
32     'feature': X_default_scaled_df.columns,
33     'importance': np.abs(shap_values_clf.values).mean(axis=0)
34 })
35 shap_importance_clf = shap_importance_clf.sort_values(by='importance',
36     ascending=False).head(10)
37 print('Top 10 Important Features for Classification:')
38 print(shap_importance_clf)

```

Listing 14: Feature Importance with SHAP

## 7 Predicting on the Test Set

### 7.1 Preparing the Test Data

I ensured that the test data undergoes the same preprocessing and feature engineering steps as the training data.

```

1 # Handle missing values in test data
2 missing_values_test = loan_test_df.isnull().sum()
3 missing_values_test = missing_values_test[missing_values_test > 0].sort_values
4     (ascending=False)
5 print('Missing values in test data:')
6 print(missing_values_test)
7
8 # Impute missing values in test data
9 for column in missing_values_test.index:
10     if loan_test_df[column].dtype in ['float64', 'int64']:
11         if column in loan_df.columns:
12             median = loan_df[column].median()
13         else:
14             median = loan_test_df[column].median()
15         loan_test_df[column].fillna(median, inplace=True)
16     else:
17         if column in loan_df.columns:
18             mode = loan_df[column].mode()[0]

```

```

18         else:
19             mode = loan_test_df[column].mode()[0]
20             loan_test_df[column].fillna(mode, inplace=True)
21
22 # Verify no missing values remain in test data
23 print('\nMissing values in test data after imputation:')
24 print(loan_test_df.isnull().sum().sum())
25
26 # Feature engineering on test data
27 loan_test_df['debt_to_income'] = loan_test_df['loan_amnt'] / loan_test_df['
    annual_inc']
28 loan_test_df['log_loan_amount'] = np.log1p(loan_test_df['loan_amnt'])
29
30 # Encode categorical variables
31 for col in categorical_cols:
32     if col in loan_test_df.columns:
33         combined_categories = pd.concat([loan_df[col], loan_test_df[col]],
34             axis=0)
35         combined_categories = combined_categories.astype('category').cat.
36             categories
37         loan_df[col] = pd.Categorical(loan_df[col], categories=
38             combined_categories).codes
39         loan_test_df[col] = pd.Categorical(loan_test_df[col], categories=
40             combined_categories).codes
41     else:
42         loan_test_df[col] = 0 # Or an appropriate default value
43
44 # Prepare test data
45 X_test = loan_test_df[feature_cols]
46 X_test = X_test.select_dtypes(include=[np.number])
47 X_test = X_test.reindex(columns=feature_cols)
48
49 # Apply the same scaling
50 X_test_scaled = scaler.transform(X_test)

```

Listing 15: Preparing Test Data

## 7.2 Making Predictions

I used the trained models to make predictions on the test data.

```

1 # Predict loan returns using the regression model
2 test_returns = best_model.predict(X_test_scaled)
3
4 # Predict early defaults using the classification model
5 test_default_prob = best_logreg.predict_proba(X_test_scaled)[: , 1]
6 test_defaults = apply_threshold(test_default_prob, best_thresh)
7 test_defaults = test_defaults.astype(int)

```

Listing 16: Making Predictions

## 7.3 Preparing the Submission File

I prepared the submission file by combining the predictions with the test set IDs.

```
1 # Prepare Submission DataFrame
2 if 'id' in loan_test_df.columns:
3     submission_df = pd.DataFrame({'id': loan_test_df['id']})
4 else:
5     submission_df = pd.DataFrame({'id': loan_test_df.index})
6
7 # Add predictions to submission dataframe
8 submission_df['early_default'] = test_defaults
9 submission_df['return'] = test_returns
10
11 # Reorder columns to match submission requirements
12 submission_df = submission_df[['id', 'early_default', 'return']]
13
14 # Save predictions to CSV
15 submission_df.to_csv('submission_test.csv', index=False)
16 print('Predictions saved to submission_test.csv')
```

Listing 17: Preparing Submission File

## 8 Conclusion

In this assignment, I performed a comprehensive analysis of LendingClub loan data to predict loan returns and early defaults using linear models. Through exploratory data analysis, we gained insights into the data structure and relationships between variables. We addressed missing values and performed feature engineering to enhance our models.

Linear regression models were employed for predicting loan returns, with hyperparameter tuning to select the best model. Logistic regression was used for predicting early defaults, with careful threshold tuning to optimize the F1 score. Feature importance analysis using SHAP values provided interpretability and insights into the significant factors influencing loan outcomes.

While linear models offer simplicity and interpretability, they may not capture complex non-linear relationships inherent in financial data. Future work could involve exploring non-linear models, ensemble methods, and advanced feature engineering to improve predictive performance. Overall, this assignment highlights the challenges and considerations in building predictive models in finance using linear approaches.