# FA 590 Statistical Learning in Finance

## Assignment 3: Loan Default and Return Prediction using Nonlinear Models

**Federica Malamisura**

**2024 Fall**

**Abstract**

This assignment focuses on predicting loan defaults and returns using nonlinear models within a financial dataset from LendingClub, spanning 2008 to 2017. The dataset includes detailed information on loan and borrower characteristics. The analysis involves several key steps: Exploratory Data Analysis (EDA), data preprocessing, feature engineering, and model building. We employ advanced regression techniques, including Decision Trees, XGBoost, CatBoost, and LightGBM, to predict loan returns. Feature importance is analyzed using SHAP values to interpret model predictions. Models are evaluated using a time-series cross-validation approach to ensure robustness and reliability. This comprehensive analysis aims to enhance predictive accuracy and provide insights into factors influencing loan outcomes in the financial domain.

# Contents

# 1 Introduction

In the realm of finance, accurately predicting loan outcomes is crucial for risk management and strategic decision-making. This assignment leverages statistical learning techniques to analyze a comprehensive dataset from LendingClub, encompassing over one million personal loans originated between 2008 and 2017. The dataset provides rich insights into borrower profiles, loan characteristics, and subsequent loan performance.

The primary objective of this analysis is to develop predictive models for loan defaults and returns using nonlinear methodologies. By employing advanced machine learning algorithms such as Decision Trees, XGBoost, CatBoost, and LightGBM, we aim to enhance the precision of predictions. An integral part of the process involves Exploratory Data Analysis (EDA) and meticulous data preprocessing to handle missing values and engineer meaningful features.

Additionally, the use of SHAP (SHapley Additive exPlanations) for feature importance analysis allows us to interpret the contributions of different features towards model predictions. This assignment not only seeks to improve predictive accuracy but also to provide actionable insights into the factors influencing loan outcomes, thereby supporting better financial decision-making.

# 2 Importing Libraries and Loading Data

## 2.1 Import Libraries

```
import pandas as pd
import numpy as np
import matplotlib as mpl
import shap
from sklearn.model_selection import RandomizedSearchCV, TimeSeriesSplit,
    StratifiedKFold
from sklearn.metrics import mean_squared_error, r2_score, accuracy_score,
    f1_score, roc_auc_score, make_scorer
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
import xgboost as xgb
import catboost as cb
import lightgbm as lgb
```

Listing 1: Importing Libraries

## 2.2 Loading Data

```
# Load training data
loan_df = pd.read_csv('lc_loan.csv')

# Load test data
loan_test_df = pd.read_csv('lc_loan_test.csv')

# 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 Data Preprocessing

## 3.1 Sampling Strategy: TimeSeriesSplit

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 **TimeSeriesSplit** for the following reasons:

1. **Order Preservation**: TimeSeriesSplit maintains the chronological order of data, which is crucial for time series analysis.

2. **Incremental Training**: With each split, more recent data is used for testing, mimicking real-world scenarios where future data is unknown.

3. **Avoiding Data Leakage**: Ensures that future data is not used in training, maintaining the integrity of model evaluation.

4. **Efficient Resource Use**: Unlike expanding or rolling windows, TimeSeriesSplit focuses on efficient use of available data without unnecessary computational overhead.

This approach is particularly suited for financial data, where maintaining temporal order is key for accurate predictions.

```
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 3: Defining Expanding Window Splits

## 3.2 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 4: Handling Missing Values

## 3.3 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 5: Feature Engineering

# 4 Model Building

## 4.1 Regression: Predicting Loan Returns

I implemented and evaluated several linear regression models:

- **Decision Tree Regressor**: Captures nonlinear relationships and interactions without requiring feature scaling.

- **XGBoost Regressor**: An efficient and scalable gradient boosting framework.

- **CatBoost Regressor**: Handles categorical features automatically and reduces the need for extensive preprocessing.

- **LightGBM Regressor**: Optimized for speed and performance, particularly with large datasets.

```
1
2 # Define target and features for regression
3 X = loan_df[feature_cols]
4 y_return = loan_df['return']
5
6 # Standardize features
7 X = X.select_dtypes(include=[np.number])  # Ensure only numeric data
8 scaler = StandardScaler()
9 X_scaled = scaler.fit_transform(X)
10
11 # Split data using TimeSeriesSplit
12 tscv = TimeSeriesSplit(n_splits=5)
13 best_regression_model = None
14 best_mse = float('inf')
15
16 for train_index, val_index in tscv.split(X_scaled):
17     X_train, X_val = X_scaled[train_index], X_scaled[val_index]
18     y_train, y_val = y_return.iloc[train_index], y_return.iloc[val_index]
19
20     # Define models
21     decision_tree = DecisionTreeRegressor(random_state=42)
22     xgboost_reg = xgb.XGBRegressor(random_state=42)
23     catboost_reg = cb.CatBoostRegressor(verbose=0, random_state=42)
24     lightgbm_reg = lgb.LGBMRegressor(random_state=42)
25
26     # Define reduced parameter grids
27     param_dist_tree = {'model__max_depth': [3, 5, 7], '
            model__min_samples_split': [2, 5]}
28     param_dist_xgb = {'model__n_estimators': [100], 'model__max_depth': [3,
            5], 'model__learning_rate': [0.01, 0.1]}
```

```
29    param_dist_cb = {'model__iterations': [100], 'model__depth': [3, 5], '
          model__learning_rate': [0.01, 0.1]}
30    param_dist_lgb = {'model__n_estimators': [100], 'model__max_depth': [3,
          5], 'model__learning_rate': [0.01, 0.1]}
31
32    # Create pipelines
33    pipeline_tree = Pipeline([('scaler', StandardScaler()), ('model',
          decision_tree)])
34    pipeline_xgb = Pipeline([('scaler', StandardScaler()), ('model',
          xgboost_reg)])
35    pipeline_cb = Pipeline([('scaler', StandardScaler()), ('model',
          catboost_reg)])
36    pipeline_lgb = Pipeline([('scaler', StandardScaler()), ('model',
          lightgbm_reg)])
37
38    # Use RandomizedSearchCV
39    randomized_tree = RandomizedSearchCV(pipeline_tree, param_distributions=
          param_dist_tree, cv=3, scoring='neg_mean_squared_error', n_iter=5,
          n_jobs=-1, random_state=42)
40    randomized_xgb = RandomizedSearchCV(pipeline_xgb, param_distributions=
          param_dist_xgb, cv=3, scoring='neg_mean_squared_error', n_iter=5,
          n_jobs=-1, random_state=42)
41    randomized_cb = RandomizedSearchCV(pipeline_cb, param_distributions=
          param_dist_cb, cv=3, scoring='neg_mean_squared_error', n_iter=5,
          n_jobs=-1, random_state=42)
42    randomized_lgb = RandomizedSearchCV(pipeline_lgb, param_distributions=
          param_dist_lgb, cv=3, scoring='neg_mean_squared_error', n_iter=5,
          n_jobs=-1, random_state=42)
43
44    # Fit models and collect results
45    models = {'DecisionTree': randomized_tree,
46              'XGBoost': randomized_xgb,
47              'CatBoost': randomized_cb,
48              'LightGBM': randomized_lgb
49              }
50
51    validation_results = {}
52    for name, random_search in models.items():
53        random_search.fit(X_train, y_train)
54        best_model = random_search.best_estimator_
55        y_pred = best_model.predict(X_val)
56        r2 = r2_score(y_val, y_pred)
57        mse = mean_squared_error(y_val, y_pred)
58        print(f'{name} on validation set:')
59        print(f' R-squared: {r2:.4f}')
60        print(f' MSE: {mse:.4f}')
61        # Store R-squared and MSE for model selection
62        validation_results[name] = {'model': best_model, 'r2': r2, 'mse':
              mse}
63
64    # Select the best model based on validation performance (lowest MSE)
65    for name, result in validation_results.items():
66    mse = result['mse']
67    if mse < best_mse:
```

```
68    best_mse = mse
69    best_regression_model = result['model']
70
71    print(f'Best regression model: {best_regression_model} with MSE: {
          best_mse:.4f}')
```

Listing 6: Regression Modeling

### 4.1.1 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.

### 4.1.2 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.

## 4.2 Classification: Predicting Early Defaults

I used non-linear models as well to predict early defaults. By performing hyperparameter tuning and threshold optimization to maximize the F1 score.

```
1  # Prepare data for classification
2
3  # Define target and features for classification
4  X_default = loan_df[feature_cols]
5  X_default = X_default.select_dtypes(include=[np.number])  # Ensure numerical
       columns
6  y_default = loan_df['early_default']
7
8  # Handle class imbalance if necessary
9  class_balance = y_default.value_counts(normalize=True)
10 print('Class distribution:', class_balance)
11
12 # Standardize features
13 scaler_default = StandardScaler()
14 X_default_scaled = scaler_default.fit_transform(X_default)
15
16 # Use StratifiedKFold for cross-validation
```

```python
17  skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
18
19  # Define models
20  decision_tree_clf = DecisionTreeClassifier(random_state=42)
21  xgboost_clf = xgb.XGBClassifier(random_state=42, use_label_encoder=False,
        eval_metric='logloss')
22  catboost_clf = cb.CatBoostClassifier(verbose=0, random_state=42)
23  lightgbm_clf = lgb.LGBMClassifier(random_state=42)
24
25  # Define parameter grids
26  param_dist_tree_clf = {'model__max_depth': [3, 5, 7], '
        model__min_samples_split': [2, 5]}
27  param_dist_xgb_clf = {'model__n_estimators': [100], 'model__max_depth': [3,
        5], 'model__learning_rate': [0.01, 0.1]}
28  param_dist_cb_clf = {'model__iterations': [100], 'model__depth': [3, 5], '
        model__learning_rate': [0.01, 0.1]}
29  param_dist_lgb_clf = {'model__n_estimators': [100], 'model__max_depth': [3,
        5], 'model__learning_rate': [0.01, 0.1]}
30
31  # Create pipelines
32  pipeline_tree_clf = Pipeline([('scaler', StandardScaler()), ('model',
        decision_tree_clf)])
33  pipeline_xgb_clf = Pipeline([('scaler', StandardScaler()), ('model',
        xgboost_clf)])
34  pipeline_cb_clf = Pipeline([('scaler', StandardScaler()), ('model',
        catboost_clf)])
35  pipeline_lgb_clf = Pipeline([('scaler', StandardScaler()), ('model',
        lightgbm_clf)])
36
37  # Define RandomizedSearchCV
38  models_clf = {
39      'DecisionTree': (pipeline_tree_clf, param_dist_tree_clf),
40      'XGBoost': (pipeline_xgb_clf, param_dist_xgb_clf),
41      'CatBoost': (pipeline_cb_clf, param_dist_cb_clf),
42      'LightGBM': (pipeline_lgb_clf, param_dist_lgb_clf)
43  }
44
45  # Use f1_score as scoring metric
46  f1_scorer = make_scorer(f1_score, average='binary')
47
48  validation_results_clf = {}
49
50  for name, (pipeline, param_dist) in models_clf.items():
51      random_search = RandomizedSearchCV(pipeline, param_distributions=
            param_dist, cv=skf, scoring=f1_scorer, n_iter=5, n_jobs=-1,
            random_state=42)
52      random_search.fit(X_default_scaled, y_default)
53      best_model = random_search.best_estimator_
54      y_pred = best_model.predict(X_default_scaled)
55      f1 = f1_score(y_default, y_pred)
56      print(f'{name}:')
57      print(f'  F1 Score (training data): {f1:.4f}')
58      validation_results_clf[name] = {'model': best_model, 'f1': f1}
59
```

```python
60 # Select best model based on highest F1 score
61 best_model_name_clf = None
62 best_f1_clf = 0
63
64 for name, result in validation_results_clf.items():
65     f1 = result['f1']
66     if f1 > best_f1_clf:
67         best_f1_clf = f1
68         best_model_name_clf = name
69         best_classifier = result['model']
70
71 print(f'Best classification model: {best_model_name_clf} with F1 Score: {
        best_f1_clf:.4f}')
72
73 # Predict probabilities on training data
74 y_pred_prob = best_classifier.predict_proba(X_default_scaled)[:, 1]
75
76 # Tune threshold to maximize F1 score
77 thresholds = np.arange(0.1, 0.9, 0.1)
78 best_threshold = 0.5
79 best_f1_score = 0
80
81 for thresh in thresholds:
82     y_pred_thresh = (y_pred_prob >= thresh).astype(int)
83     f1 = f1_score(y_default, y_pred_thresh)
84     if f1 > best_f1_score:
85         best_f1_score = f1
86         best_threshold = thresh
87
88 print(f'Final threshold set to: {best_threshold:.2f} with F1 score: {
        best_f1_score:.4f}')
89
90 # Define function to apply threshold
91 def apply_threshold(probs, threshold):
92     return (probs >= threshold).astype(int)
93
94 # Evaluate on the training set
95 y_train_pred = apply_threshold(y_pred_prob, best_threshold)
96 train_accuracy = accuracy_score(y_default, y_train_pred)
97 train_auc = roc_auc_score(y_default, y_pred_prob)
98
99 print('\nIn-sample Performance for Early Default Prediction:')
100 print(f'Accuracy: {train_accuracy:.4f}')
101 print(f'F1 Score: {best_f1_score:.4f}')
102 print(f'AUC Score: {train_auc:.4f}')
```

Listing 7: Classification Modeling

# 5 Predicting on the Test Set

## 5.1 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(ascending=False)
4  print('Missing values in test data:')
5  print(missing_values_test)
6
7  # Impute missing values in test data
8  for column in missing_values_test.index:
9      if loan_test_df[column].dtype in ['float64', 'int64']:
10         if column in loan_df.columns:
11             median = loan_df[column].median()
12         else:
13             median = loan_test_df[column].median()
14         loan_test_df[column].fillna(median, inplace=True)
15     else:
16         if column in loan_df.columns:
17             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]],
               axis=0)
34         combined_categories = combined_categories.astype('category').cat.
               categories
35         loan_df[col] = pd.Categorical(loan_df[col], categories=
               combined_categories).codes
36         loan_test_df[col] = pd.Categorical(loan_test_df[col], categories=
               combined_categories).codes
37     else:
38         loan_test_df[col] = 0  # Or an appropriate default value
39
40 # Prepare test data
```

```
41 X_test = loan_test_df[feature_cols]
42 X_test = X_test.select_dtypes(include=[np.number])
43 X_test = X_test.reindex(columns=feature_cols)
44
45 # Apply the same scaling
46 X_test_scaled = scaler.transform(X_test)
```
Listing 8: Preparing Test Data

## 5.2 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 9: Making Predictions

## 5.3 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 10: Preparing Submission File

# 6 Generative AI Implementation

In this assignement, Generative AI, in particular o1-mini from OpenAI and Claude Sonnet 3.5 from Anthropics were used to get a more precise view of the parameters to be used for training the various Gradient Boost models, in particular the strengths and weaknesses of XGBoost, CatBoost, and LightGBM, in order to use for testing and predictive purposes the most effective model based on the available dataset

# 7 Conclusion

In this assignment, I performed a comprehensive analysis of LendingClub loan data to predict loan returns and early defaults using non-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.

Machine Learning Gradient-Boost models were employed for predicting loan returns, with hyperparameter tuning to select the best model. They were used as well for predicting early defaults, with careful threshold tuning to optimize the F1 score.

The goal of the research is to highlight how linear and nonlinear models, even of different types from each other, can give different results depending on the specific dataset. This makes it clear that there is no absolute best model, but that careful and meticulous research is needed on which model is actually the best fit for the specific dataset and the variables that we want to predict