

# Loan Approval Prediction Report

Paula Lozano Gonzalo

February 20, 2025

## 1 Introduction

This report outlines the process of developing and evaluating classification models for predicting loan approval using the Kaggle Playground Series - Season 4 - Episode 10 Loan Approval Prediction dataset. The goal was to train multiple classification models, perform hyperparameter tuning using cross-validation, and submit predictions to the Kaggle competition.

## 2 Data Understanding

The dataset contains the following features:

(I had to do some searches to understand what each thing was because I am not familiar with all the terms and there was no explanation of the features in kaggle like in the last assignment. That is why I am adding this into the report)

- **id**: Unique identifier for each record.
- **person\_age**: Age of the applicant.
- **person\_income**: Annual income of the applicant.
- **person\_home\_ownership**: Type of home ownership (e.g., 'own', 'rent', 'mortgage').
- **person\_emp\_length**: Employment length in years.
- **loan\_intent**: Purpose of the loan (e.g., 'education', 'medical', 'personal').
- **loan\_grade**: Loan risk grade (e.g., 'A', 'B', 'C').
- **loan\_amnt**: Loan amount requested.
- **loan\_int\_rate**: Loan interest rate (APR).
- **loan\_percent\_income**: Percentage of income allocated to loan repayment.
- **cb\_person\_default\_on\_file**: Binary indicator of default history (yes/no).
- **cb\_person\_cred\_hist\_length**: Length of credit history in years.
- **loan\_status**: Loan status (0 or 1).

### 3 Data Exploration and Visualization

When I plotted the histograms and the scatter plots I realized that since a lot of the data is polarized you can barely see the points. Because of this I used the histograms to visualize the distribution of features. (I will leave some examples bellow)

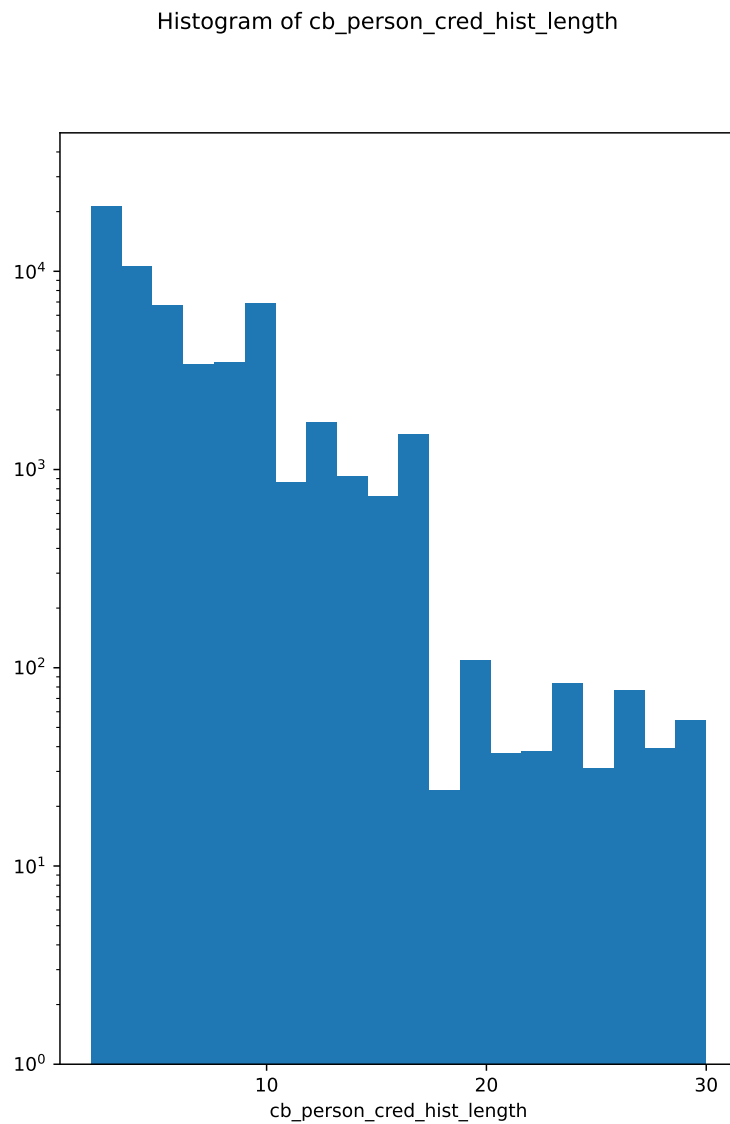


Figure 1: Credit History Lenght

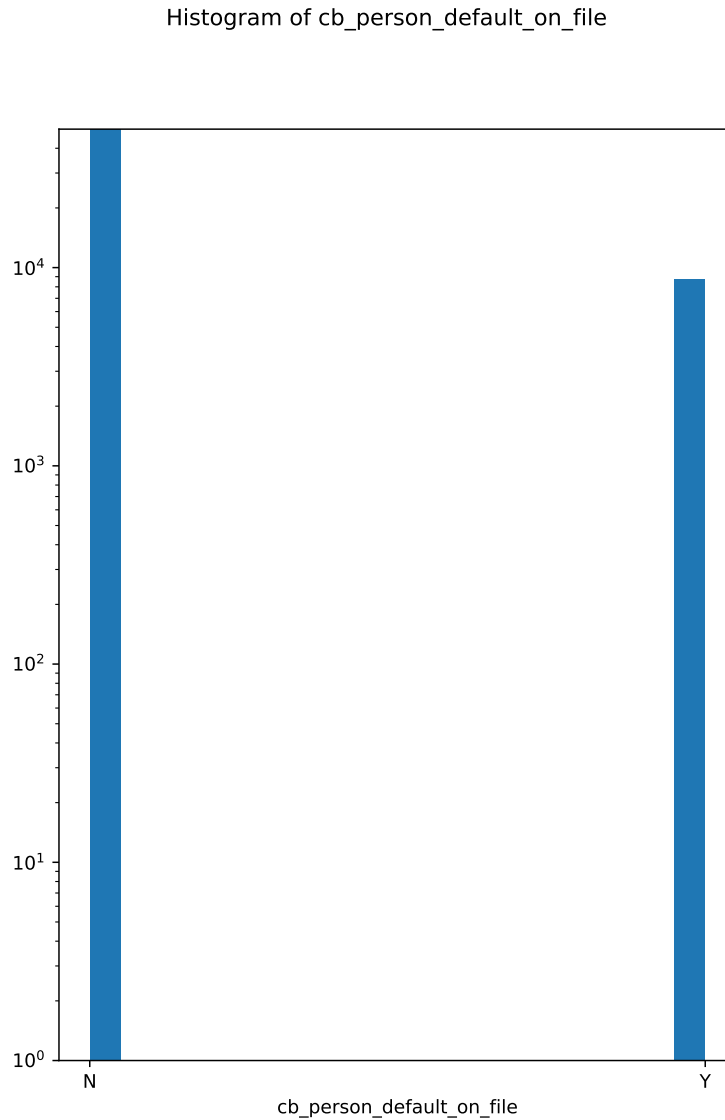


Figure 2: Does the person have default history?

## 4 Model Development and Evaluation

Three models were trained and evaluated:

1. **SGDClassifier**
2. **RandomForestClassifier**
3. **GradientBoostingClassifier**

### 4.1 SGDClassifier

- **Best Score (Cross-Validation):** 0.9144

- **Confusion Matrix:**

t/p	F	T
F	48509.0	1786.0
T	6518.0	1832.0

- **Precision:** 0.506

- **Recall:** 0.219

- **F1 Score:** 0.306

- **Analysis:** The model performed poorly, with only 50.6% of positive predictions being correct and 21.9% of actual positives identified.

The best hyperparameters found for the SGD classifier were:

```

1 {
2   'features__categorical__categorical-features-only__do_numerical': False,
3   'features__categorical__categorical-features-only__do_predictors': True,
4   'features__categorical__encode-category-bits__categories': 'auto',
5   'features__categorical__encode-category-bits__handle_unknown': 'ignore',
6   'features__categorical__missing-data__strategy': 'most_frequent',
7   'features__numerical__missing-data__strategy': 'median',
8   'features__numerical__numerical-features-only__do_numerical': True,
9   'features__numerical__numerical-features-only__do_predictors': True,
10  'features__numerical__polynomial-features__degree': 2,
11  'model__alpha': 0.0001,
12  'model__average': False,
13  'model__class_weight': None,
14  'model__early_stopping': False,
15  'model__epsilon': 0.1,
16  'model__eta0': 0.0,
17  'model__fit_intercept': True,
18  'model__l1_ratio': 0.15,
19  'model__learning_rate': 'optimal',
20  'model__loss': 'log_loss',
21  'model__max_iter': 1000,
22  'model__n_iter_no_change': 5,
23  'model__n_jobs': -1,
24  'model__penalty': None,
25  'model__power_t': 0.5,
26  'model__random_state': None,
27  'model__shuffle': True,
28  'model__tol': 0.001,
29  'model__validation_fraction': 0.1,
30  'model__verbose': 0,
31  'model__warm_start': False
32 }
```

## 4.2 RandomForestClassifier

- **Best Score (Cross-Validation):** 0.9499

- **Confusion Matrix:**

t/p	F	T
F	49754.0	541.0

T 2355.0 5995.0

- **Precision:** 0.917
- **Recall:** 0.718
- **F1 Score:** 0.805
- **Kaggle Score:** 0.93810
- **Analysis:** The model performed significantly better, with 91.7% of positive predictions being correct and 71.8% of actual positives identified.

The best hyperparameters found for the RandomForest classifier were:

```
1 {
2   'features__categorical__categorical-features-only__do_numerical': False,
3   'features__categorical__categorical-features-only__do_predictors': True,
4   'features__categorical__encode-category-bits__categories': 'auto',
5   'features__categorical__encode-category-bits__handle_unknown': 'ignore',
6   'features__categorical__missing-data__strategy': 'most_frequent',
7   'features__numerical__missing-data__strategy': 'median',
8   'features__numerical__numerical-features-only__do_numerical': True,
9   'features__numerical__numerical-features-only__do_predictors': True,
10  'features__numerical__polynomial-features__degree': 2,
11  'model__bootstrap': True,
12  'model__ccp_alpha': 0.0,
13  'model__class_weight': None,
14  'model__criterion': 'gini',
15  'model__max_depth': None,
16  'model__max_features': None,
17  'model__max_leaf_nodes': None,
18  'model__max_samples': None,
19  'model__min_impurity_decrease': 0.0,
20  'model__min_samples_leaf': 1,
21  'model__min_samples_split': 2,
22  'model__min_weight_fraction_leaf': 0.0,
23  'model__monotonic_cst': None,
24  'model__n_estimators': 500,
25  'model__n_jobs': -1,
26  'model__oob_score': False,
27  'model__random_state': None,
28  'model__verbose': 0,
29  'model__warm_start': False
30 }
```

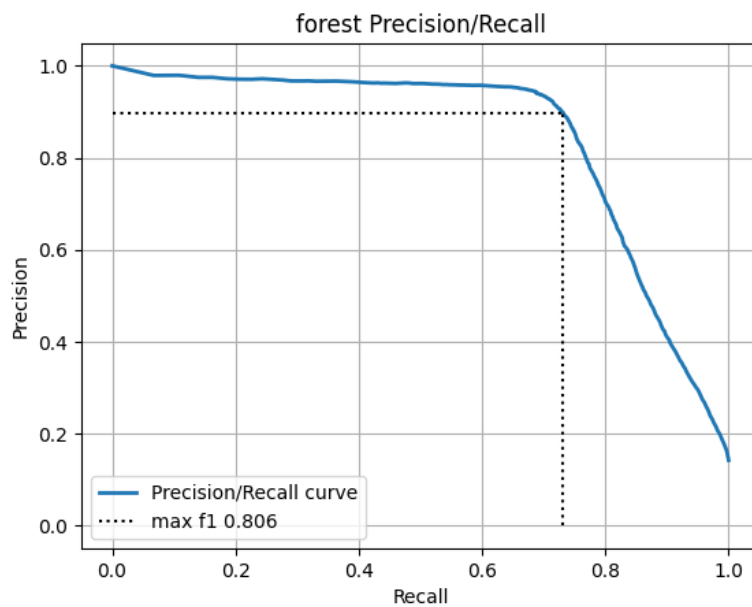


Figure 3: PR

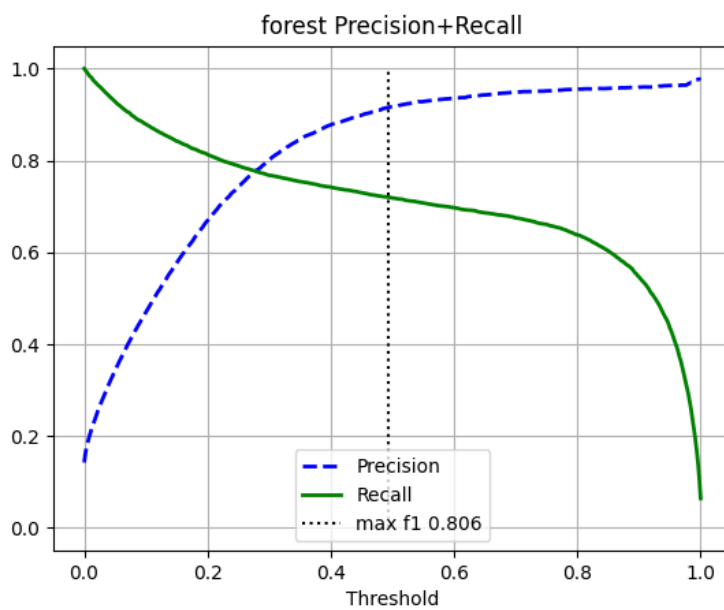


Figure 4: Precision/Recall

I have to say that I don't fully understand these graphs because I am not really familiar with them but I still wanted to plot them and look at them to see how good/bad they were.

### 4.3 GradientBoostingClassifier

- **Best Score (Cross-Validation):** 0.9490

- **Confusion Matrix:**

t/p	F	T
F	49775.0	520.0
T	2359.0	5991.0

- **Precision:** 0.920

- **Recall:** 0.717

- **F1 Score:** 0.806

- **Kaggle Score:** 0.94164

- **Analysis:** The model performed slightly better than RandomForestClassifier, with 92.0% of positive predictions being correct and 71.7% of actual positives identified.

The best hyperparameters found for the GradientBoosting classifier were:

```

1 {
2   'features__categorical__categorical-features-only__do_numerical': False,
3   'features__categorical__categorical-features-only__do_predictors': True,
4   'features__categorical__encode-category-bits__categories': 'auto',
5   'features__categorical__encode-category-bits__handle_unknown': 'ignore',
6   'features__categorical__missing-data__strategy': 'most_frequent',
7   'features__numerical__missing-data__strategy': 'median',
8   'features__numerical__numerical-features-only__do_numerical': True,
9   'features__numerical__numerical-features-only__do_predictors': True,
10  'features__numerical__polynomial-features__degree': 2,
11  'model__ccp_alpha': 0.0,
12  'model__criterion': 'squared_error',
13  'model__init': None,
14  'model__learning_rate': 0.1,
15  'model__loss': 'exponential',
16  'model__max_depth': None,
17  'model__max_features': None,
18  'model__max_leaf_nodes': None,
19  'model__min_impurity_decrease': 0.0,
20  'model__min_samples_leaf': 1,
21  'model__min_samples_split': 2,
22  'model__min_weight_fraction_leaf': 0.0,
23  'model__n_estimators': 200,
24  'model__n_iter_no_change': None,
25  'model__random_state': None,
26  'model__subsample': 0.9,
27  'model__tol': 0.0001,
28  'model__validation_fraction': 0.1,
29  'model__verbose': 0,
30  'model__warm_start': False
31 }

```

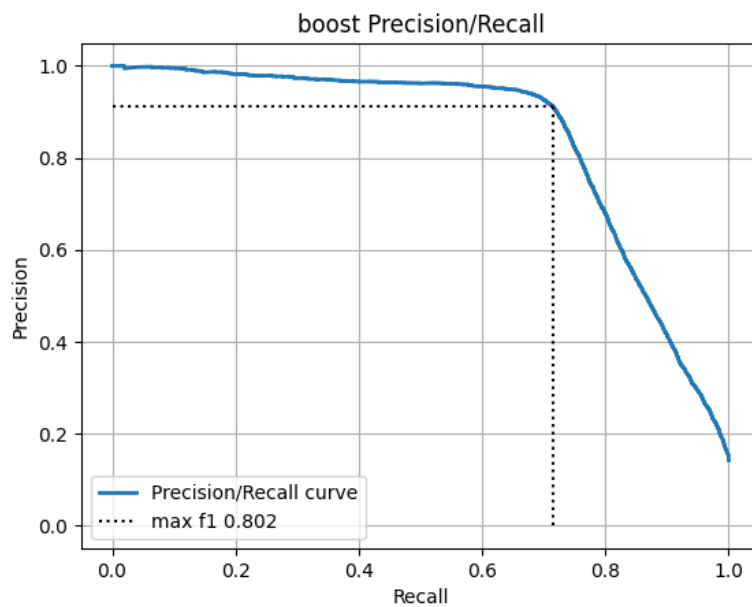


Figure 5: PR

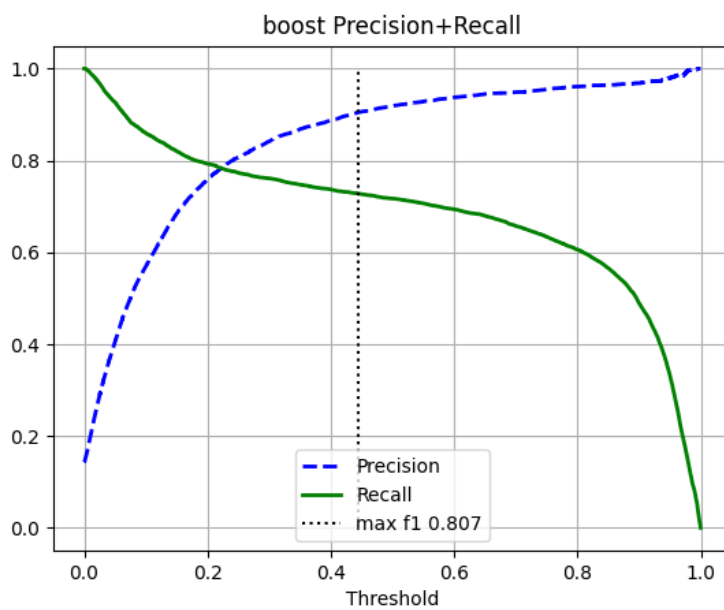


Figure 6: Precision/Recall

## 5 Model Analysis

- **SGDClassifier**: Performed poorly, probably because to the complexity of the data.



```
terminate kill signals hide
Status: Elapsed: IO/R: IO/W: Parent:
Dead 02:44:06 113 MiB 355 MiB bash

Memory: 8.3% ..... 1.33 GiB
C python3 ./pipeline.py random-search --model-type
M boost --train-file data/train.csv --model-file
D models/GradientBoostingClassifier.joblib --search
per-core reverse-tree-cpu-lazy >
```

Figure 7: Fun fact. Time boost took to run.

- **RandomForestClassifier**: Performed well, achieving a high Kaggle score and demonstrating good precision and recall.
- **GradientBoostingClassifier**: Slightly outperformed RandomForestClassifier, indicating that boosting techniques may be more effective for this dataset.

Submission and Description		Private Score	Public Score	Selected
	<b>predictionsGradientBoostingProba.csv</b> Complete (after deadline) · 21s ago · GradientBoostingClassifier	<b>0.94164</b>	<b>0.94214</b>	<input type="checkbox"/>
	<b>predictionsRandomForestProba.csv</b> Complete (after deadline) · 13h ago · RandomForestClassifier	<b>0.93810</b>	<b>0.93987</b>	<input type="checkbox"/>

Figure 8: Kaggle scores.