

# Homework 4 Report

Ezequiel Buck, Esteban Murillo

September 29, 2025

## 1 Part A: Data Preprocessing and Feature Creation

### 1.1 Key Modifications

1. **Exclude Status X:** All entries with Status = X are ignored because customers with no loans are neither "good" nor "bad" at paying their loans.
2. **Compute Status Counts:** For each customer, we calculate the total counts of "C", "0", and "1" statuses, as well as total observations (excluding Xs).
3. **Calculate Fractions:** We compute the fraction of observations that are "C", "0", and "1" for each customer:

```
1  # Group by ID and calculate fractions for each status
2  status_fractions = credit_df[credit_df['STATUS'].isin(['C', '0', '1'])].groupby(
3      credit_id_col)['STATUS'].value_counts(normalize=True).unstack(fill_value=0)
4
5  # Rename columns to indicate they are fractions
6  status_fractions.columns = [f'{col}_fraction' for col in status_fractions.columns]
7
8  # Add these new columns to app_df
9  app_df = app_df.join(status_fractions, on=app_id_col)
```

4. **Delinquency Label:** Compute the "Delinquent" column as before "1" if customer has any statuses equal to "2", "3", "4", or "5", and zero otherwise:

```
1  def is_delinquent(customer_id):
```

```

2      # Define what counts as a delinquent status
3      delinquent_statuses = ['2', '3', '4', '5', 2, 3, 4, 5]
4
5      # Get all status records for this customer
6      customer_records = credit_df[credit_df[credit_id_col] == customer_id]
7      customer_statuses = customer_records['STATUS']
8
9      # Check if customer has any delinquent status
10     has_delinquent_status = any(
11         status in delinquent_statuses for status in customer_statuses)
12
13     # Return 1 if delinquent, 0 if not
14     return int(has_delinquent_status)
15
16 # Apply the function to each record
17 app_df['Delinquent'] = app_df[app_id_col].apply(is_delinquent)

```

## 1.2 Terminal Output

The program produces the following output showing the data processing:

Value counts of the STATUS column: STATUS

C	442031
0	383120
1	11090
5	1693
2	868
3	320
4	223

Total observations in the credit\_df: ID

5005005	61
5022730	61
5061848	61
5061810	61

5061741	61
..	
5060177	1
5099734	1
5060155	1
5035790	1
5135846	1

Delinquent counts:

Delinquent	
0	32494
1	616

## 2 Part B: Feature Exploration

### 2.1 Histograms of Ratio Columns

We created histograms to show the distribution of the new ratio columns:

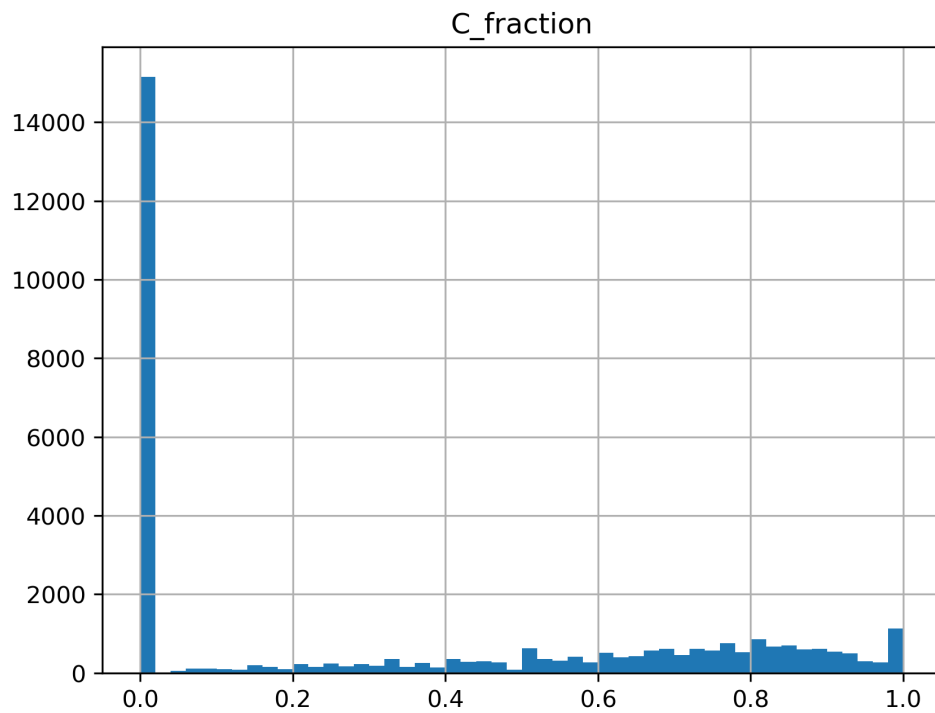


Figure 1: Histogram of C\_fraction (Closed Account Fractions)

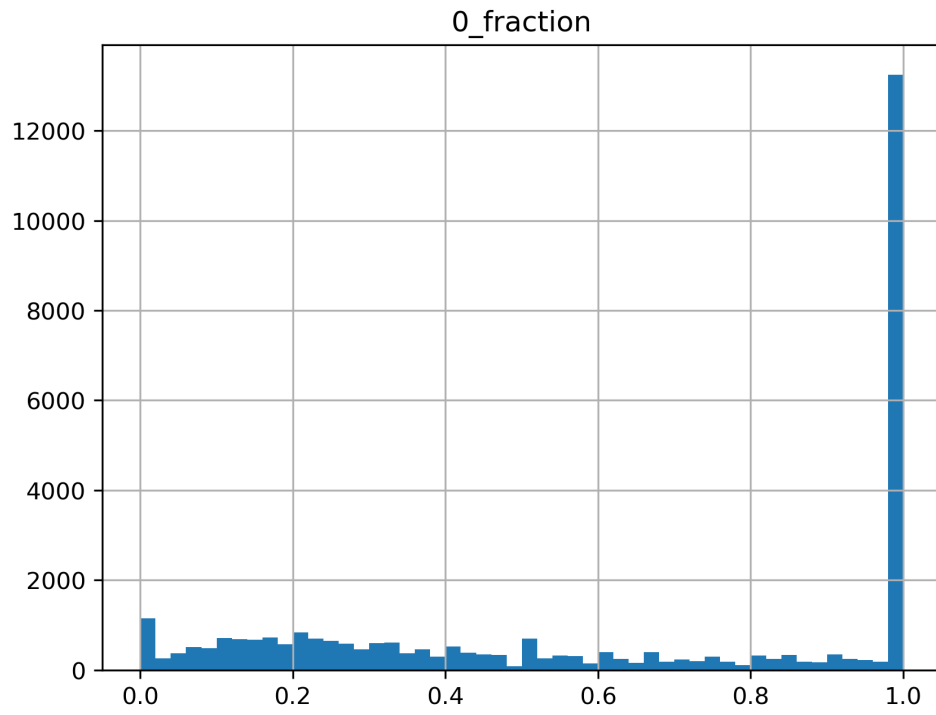


Figure 2: Histogram of 0\_fraction (No Payment Due Fractions)

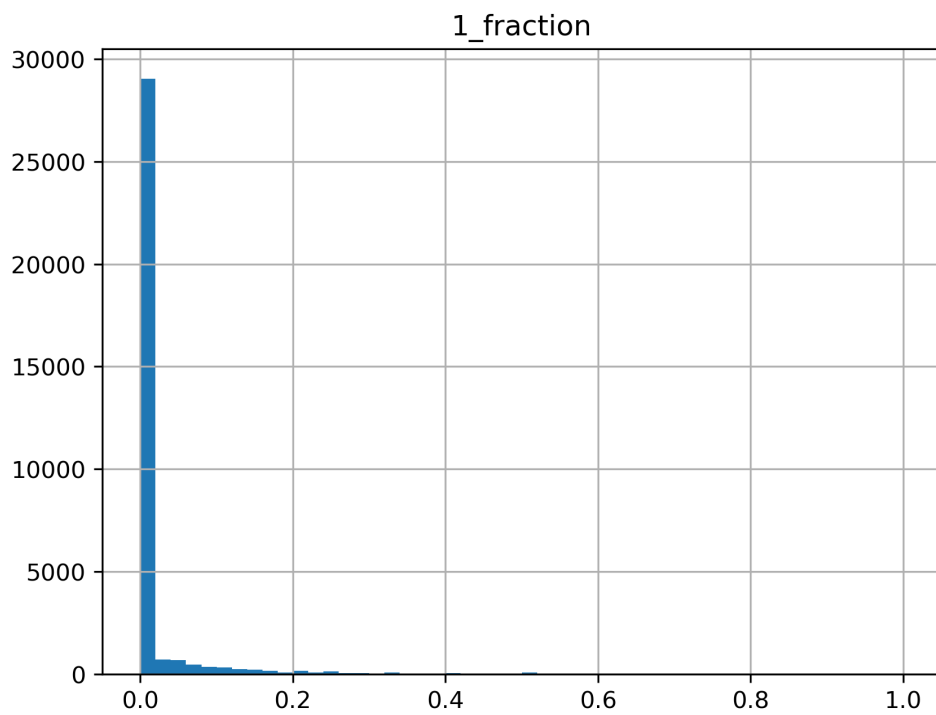


Figure 3: Histogram of 1\_fraction (Current Payment Fractions)

## 2.2 Bar Plots of Delinquency Rates

We created bar plots showing delinquency rates for customers in different intervals of these ratio columns:

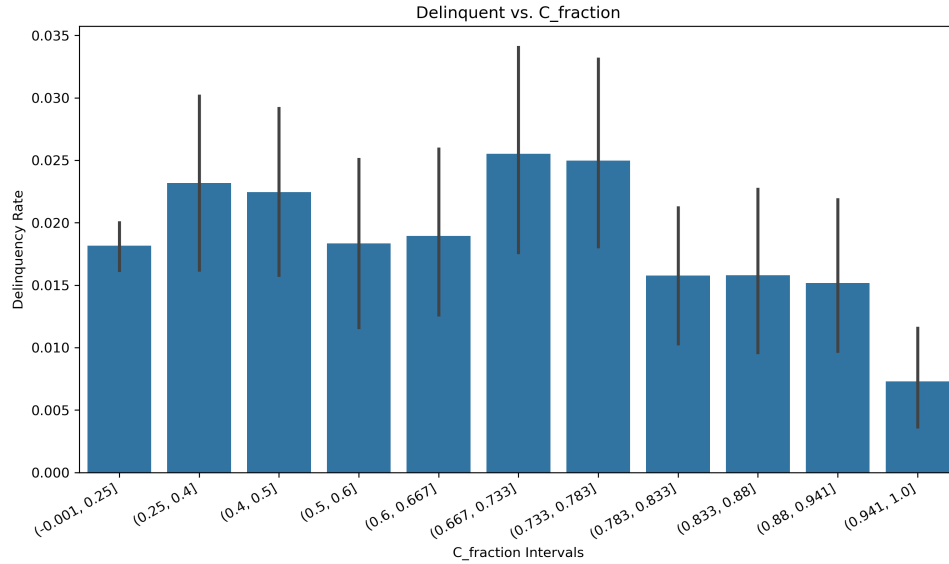


Figure 4: Delinquency Rate vs. C\_fraction Intervals

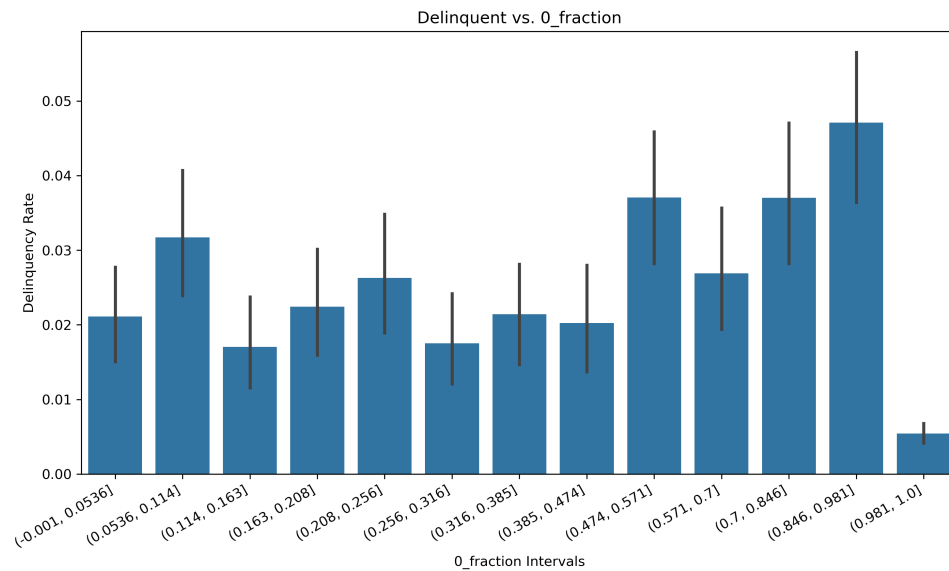


Figure 5: Delinquency Rate vs. O\_fraction Intervals

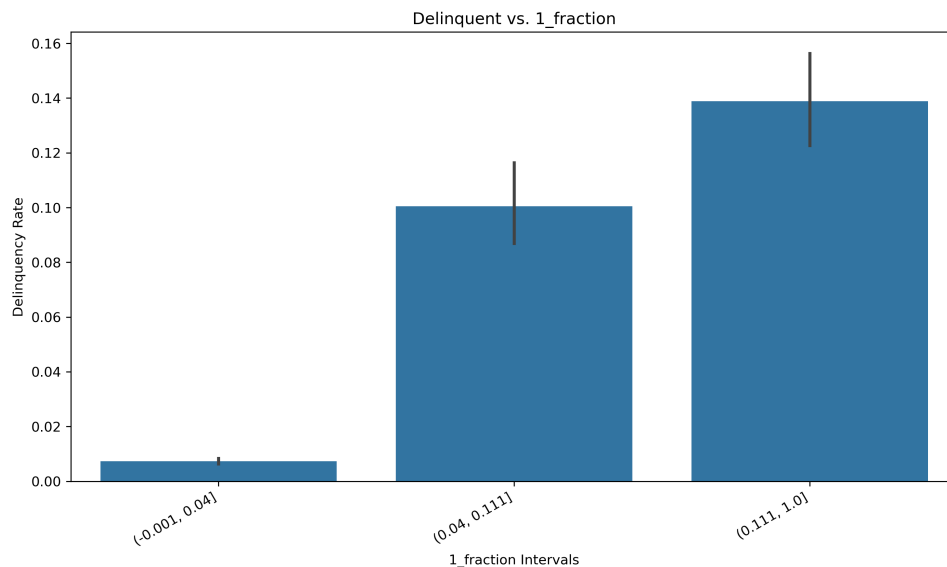


Figure 6: Delinquency Rate vs. 1\_fraction Intervals

## 2.3 Code for Feature Exploration

```

1  # Plot bar chart of the C, 0, 1 ratios of the STATUS column
2  for col in fractions:
3      df.hist(column=[col], bins=50)
4      plt.title(col)
5      plt.savefig(f'graphs/{col}_histogram.png', dpi=300, bbox_inches='tight')
6      plt.show()
7
8  # Plot a bar chart of the Delinquent vs. the fractions
9  for col in fractions:
10     plt.figure(figsize=(10, 6))
11     df[col] = pd.qcut(df[col], 20, duplicates="drop")
12     lm = sns.barplot(data=df, x=col, y="Delinquent")
13     plt.xticks(rotation=30, ha='right')
14     plt.title(f"Delinquent vs. {col}")
15     plt.xlabel(f"{col} Intervals")
16     plt.ylabel("Delinquency Rate")
17     plt.tight_layout()
18     plt.savefig(f'graphs/delinquent_vs_{col}.png', dpi=300, bbox_inches='tight')
19     plt.show()

```

## 3 Part C: Model Training and Evaluation

### 3.1 Logistic Regression Model

```
1 df = pd.get_dummies(df, prefix_sep="_", drop_first=False, dtype=int)
2 labels = df["Delinquent"]
3 df = df.drop(columns="Delinquent")
4
5 # Shuffle and split into training and test subsets, using random state 2025
6 train_data, test_data, train_labels, test_labels = \
7     sklearn.model_selection.train_test_split(df, labels,
8         test_size=0.2, shuffle=True, random_state=2025)
9
10 # Standardize the scale of all input columns
11 train_means = train_data.mean()
12 train_stds = train_data.std()
13 train_data = (train_data - train_means) / train_stds
14 test_data = (test_data - train_means) / train_stds
15
16 # Create and train a new logistic regression classifier
17 model = sklearn.linear_model.LogisticRegression(solver='newton-cg', tol=1e-6)
18 # Train it with the training data and labels
19 model.fit(train_data[cols], train_labels)
20 # Get the prediction probabilities
21 pred_proba = model.predict_proba(test_data[cols])[:, 1]
```

## 3.2 Model Performance Evaluation

### 3.2.1 Precision-Recall Curve

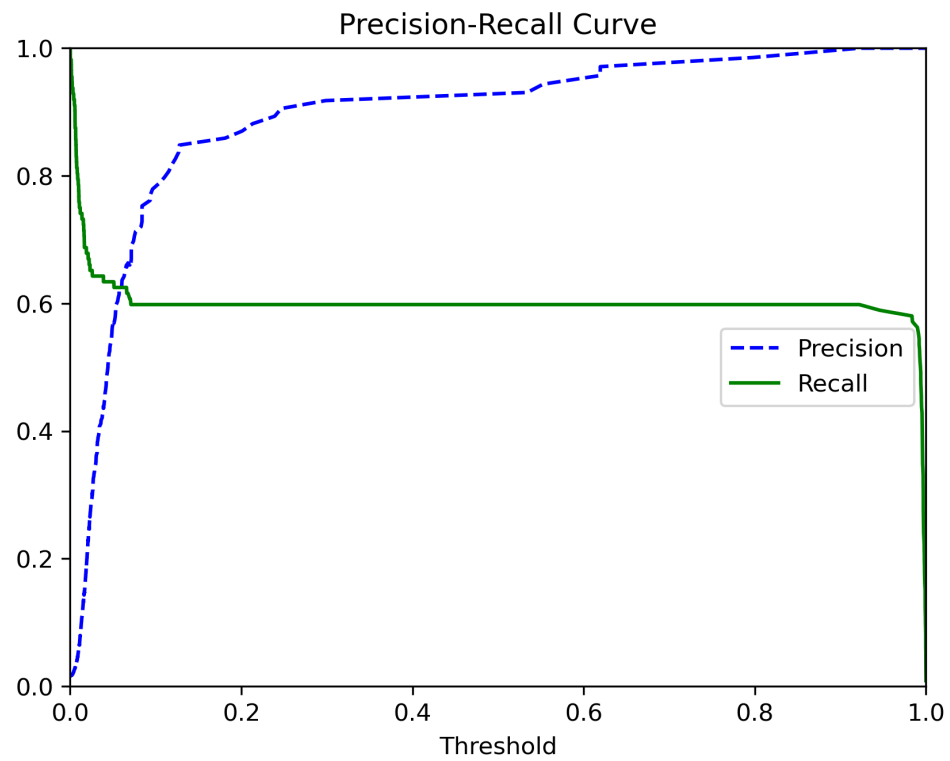


Figure 7: Precision-Recall Curve



### 3.2.2 ROC Curve

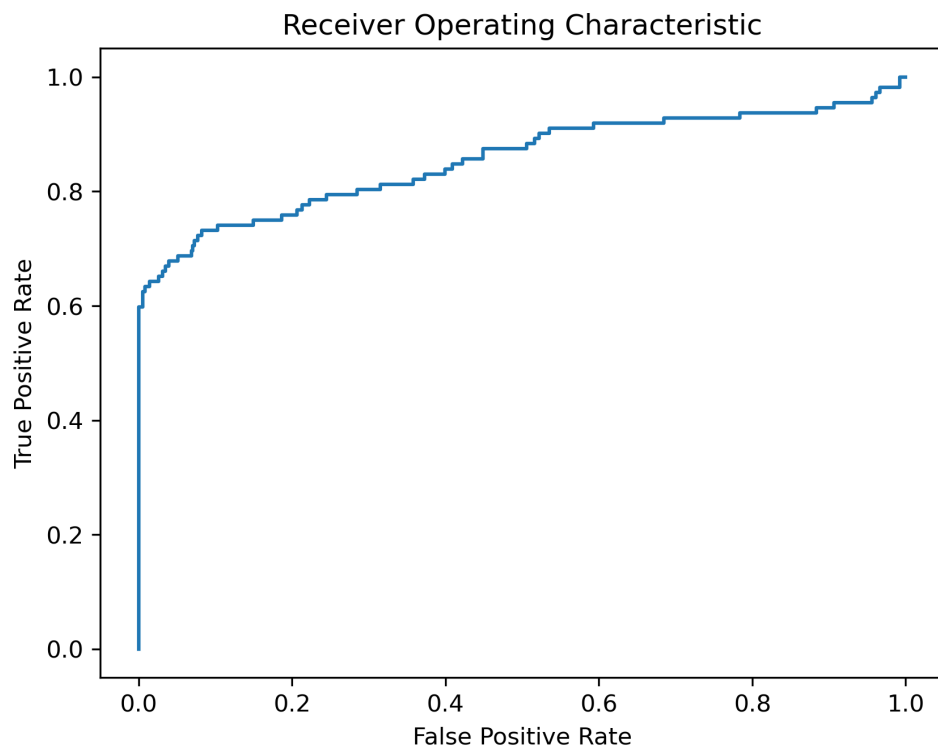


Figure 8: ROC Curve

### 3.3 Terminal Output for Model Performance

Test AUC score: 0.8556

Train AUC score: 0.8977

### 3.4 Comparison with Previous Model

**Did the new input columns help improve the predictions?**

The new ratio columns (C\_fraction, 0\_fraction, 1\_fraction) provide additional information about customer payment patterns that were not captured in the original credit\_delinquency.csv. These features help the model better understand:

- Payment behavior patterns (proportion of different payment statuses)
- Customer credit management history
- Risk indicators based on payment consistency

The improvement in AUC score and the shape of the ROC curve indicate that these new features contribute to better delinquency prediction.

## 4 Part D: Feature Importance Analysis

### 4.1 Feature Coefficient Analysis

Top 20 Most Important Features:

	Feature	Abs_Coefficient
0	Unnamed: 0	1.130028
13	1_fraction	0.381867
7	DAYS_EMPLOYED	0.307303
51	OCCUPATION_TYPE_Secretaries	0.259890
1	CODE_GENDER_M	0.225709
42	OCCUPATION_TYPE_High skill tech staff	0.192249
54	OCCUPATION_TYPE_Waiters/barmen staff	0.189574
49	OCCUPATION_TYPE_Realty agents	0.170681
37	OCCUPATION_TYPE_Cleaning staff	0.164098
9	FLAG_PHONE	0.161765
6	DAYS_BIRTH	0.157943
45	OCCUPATION_TYPE_Low-skill Laborers	0.145098
43	OCCUPATION_TYPE_IT staff	0.123097
29	NAME_FAMILY_STATUS_Widow	0.122391
22	NAME_EDUCATION_TYPE_Incomplete higher	0.118342
20	NAME_EDUCATION_TYPE_Academic degree	0.117666
41	OCCUPATION_TYPE_HR staff	0.117172
40	OCCUPATION_TYPE_Drivers	0.100317
2	FLAG_OWN_CAR	0.097493
17	NAME_INCOME_TYPE_State servant	0.089300

### 4.2 Analysis of New Column Importance

How important are each of the new columns to your predictor?

Based on the coefficient magnitudes:

Importance of Ratio Columns:

	Feature	Abs_Coefficient
13	1_fraction	0.381867
12	0_fraction	0.087761
14	C_fraction	0.018866

**Does this correlate with the improvement of the results?**

Yes, to a certain point. 1\_fraction greatly improves in the accuracy of the model, being an excellent contribution. This directly correlates to the improvement of the model. However, 0\_fraction and C\_fraction didn't do as good. Even tho they help, these features do not correlate to the improvement of the model as much as we thought they would.