

A Fraud Detection Solution for Extremely Imbalanced Financial Data

Xiaoqing Miao

xmia665@aucklanduni.ac.nz

University of Auckland

Auckland, New Zealand

Abstract

In Iteration-3, our primary objective is minimizing expected cost under a business loss ratio of FP:FN = 1:25. We select the operating threshold from a validation-set cost-threshold curve and evaluate on the original-distribution test set. Using the PaySim dataset, we focus on TRANSFER/CASH_OUT transactions, apply log1p transformation + standardization, compare down-sampling LR (Iter-A) versus class-weighted LR (Iter-B), and adopt Iter-B as the current operating candidate.

On the test set, the model achieves ROC-AUC = 0.9756 (meeting the project-level AUC target of ≥ 0.95) and PR-AUC(AP) = 0.582. At the recommended operating point $\tau = 0.9488$, we obtain Precision = 27.42%, Recall = 65.95%, and FPR = 0.52%, delivering the lowest expected cost among candidates. This operating point does not yet satisfy the project-level stretch targets Precision $\geq 70\%$ / Recall $\geq 85\%$; thus, we position it as a cost-optimal baseline and outline next steps: targeted feature engineering and non-linear models (e.g., stronger tree-based learners and interaction terms), tighter threshold-alert quota coupling, and entity/time-grouped robustness checks. Overall, Iteration-3 secures strong ranking performance and a reproducible cost-recall trade-off, establishing a clear path to reach the project's stretch goals.

Keywords: fraud detection, class imbalance, CRISP-DM, logistic regression, PaySim dataset

ACM Reference Format:

Xiaoqing Miao. 2024. A Fraud Detection Solution for Extremely Imbalanced Financial Data. In *Proceedings of University of Auckland*. ACM, New York, NY, USA, 29 pages. <https://doi.org/XXXXXX.XXXXXXX>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org. *University of Auckland, 2025*

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/2024/XX

<https://doi.org/XXXXXX.XXXXXXX>

1 Business/Situation Understanding

1.1 Business/Situation Objectives

Business Context

The digital payment industry is facing severe fraud challenges. The PaySim dataset simulates one month of mobile payment transactions, containing a total of 6,362,620 records, of which only 0.13% (8,213 transactions) are fraudulent, reflecting the highly imbalanced nature of real-world scenarios.

Key Business Objectives

- **Reduce Financial Losses:** Reduce fraud-related losses by 75% through early detection.
- **Protect Customer Trust:** Prevent customers from becoming fraud victims and maintain customer satisfaction above 90%.
- **Optimize Operational Efficiency:** Reduce manual fraud investigation workload by 60%.

Business Success Criteria

- **Detection Performance:** Recall $\geq 85\%$, Precision $\geq 70\%$, AUC score ≥ 0.95 .
- **Operational Efficiency:** Transaction processing time < 100 milliseconds, system availability $\geq 99.9\%$.
- **Financial Impact:** ROI of 10:1 (for every \$1 invested, prevent \$10 in fraud losses).

Two-level Goals and Binding (Iteration-3 \leftrightarrow §8.4)

- **Iteration Goal (Iteration-3):** Min Cost @ FP:FN=1:25 (**Achieved**). Operating threshold $\tau_B=0.9488$ selected via validation-set cost-threshold curve under operational constraints (theoretical basis in [Elkan 2001](#)); all evaluation in Chapter 8 follows this τ_B .
- **Project-level Hard Targets:** AUC ≥ 0.95 (**Achieved, 0.9756**); Precision $\geq 70\%$ / Recall $\geq 85\%$ (**Not achieved**). This iteration serves as the cost-optimal baseline for subsequent improvements.

Alignment with Data Mining Objectives

These business objectives directly translate to:

- Building a binary classifier for fraud detection (fraud vs. legitimate).
- Identifying high-risk patterns in TRANSFER and CASH_OUT transactions.
- Developing real-time risk scoring for transaction approval decisions.

Expected Value

Through fraud prevention and control, annual savings of \$300,000–\$750,000 can be achieved, reducing investigation costs by \$75,000, and enhancing market competitive positioning as a secure payment platform.

1.2 Assessment of the Situation

1.2.1 Resource Inventory. As an individual student project, available resources include: PaySim open-source dataset (6.36 million transaction records), personal laptop (16GB RAM), Anaconda Python environment with Jupyter Notebook (with installed libraries including NumPy, Pandas, Matplotlib, Seaborn, scikit-learn, etc.), and 1 week of project time. Knowledge resources include course materials, Kaggle dataset documentation, online tutorials, and relevant academic papers. As a data mining course student, I possess foundational knowledge of classification algorithms and operational experience with Python data analysis and modeling tools.

1.2.2 Requirements and Assumptions. Requirements: Complete all phases of CRISP-DM, build a binary classification model capable of detecting fraudulent transactions with a performance target of $AUC \geq 0.95$, and submit a comprehensive project report covering data understanding, modeling process, and results interpretation.

Assumptions: PaySim synthetic data adequately reflects real-world fraud characteristics; personal laptop (16GB RAM) can handle **6,362,620** transaction records; two-week project cycle is sufficient for multiple iterations; Python and Jupyter Notebook with common data analysis libraries (such as Pandas, scikit-learn) have the functionality to meet project requirements.

1.2.3 Constraints.

- **Time Constraints:** Project must be completed within 2 weeks, with 15-20 hours available per week.
- **Technical Constraints:** Limited by personal computer's computational performance (16GB RAM), may require data sampling or batch processing.
- **Data Constraints:** Can only rely on publicly available PaySim synthetic dataset, lacking real business transaction data for further validation (PaySim simulator: López-Rojas et al. 2016).
- **Experience Constraints:** As a student, limited experience in actual business operations and risk control processes for financial fraud detection.

1.2.4 Risks and Mitigation Measures. Project risk assessment covers five main risk areas: technical, performance, time, software, and data. For each risk, specific response plans are developed to ensure the project can be completed successfully under adverse conditions, thereby maximizing the possibility of achieving project objectives.

Risk Assessment Table

Risk Category	Risk Description	Prob.	Impact	Mitigation Measures
Technical	Limited memory may cause bottlenecks processing 6.36M records	Med.	High	1. Stratified sampling 2. Batch processing 3. Use lab computers
Performance	Model may not reach $AUC \geq 0.95$ target	Med.	High	1. Multiple algorithms 2. Ensemble methods 3. Feature engineering
Time	2-week timeframe insufficient for complete iterations	Med.	Med.	1. Prioritize core phases 2. Optional optimization 3. Daily milestones
Software	Python environment or dependency conflicts	Low	High	1. Virtual environments 2. Regular saves 3. Google Colab backup
Data	Extreme imbalance (fraud only 0.13%) causes majority bias	High	Med.	1. SMOTE oversampling 2. Class weights 3. Precision-recall tuning

Table 1. Project Risk Assessment and Mitigation Strategies

See Table 1 for the project's risk assessment and mitigation strategies.

1.3 Data Mining Objectives

The primary objective of this project is to **build an effective classification model to accurately predict the is-Fraud label in financial transactions**. Specifically, it aims to distinguish between normal transactions ($\text{isFraud} = 0$) and fraudulent ones ($\text{isFraud} = 1$) by analyzing various transaction features (e.g., transaction type, amount, account balance changes).

1.3.1 Model Evaluation Methods. Due to the significantly smaller number of fraud samples compared to normal samples, the dataset exhibits severe class imbalance. Therefore, the traditional **Accuracy metric is misleading** and will not be used as the primary evaluation criterion.

Instead, this study adopts the **ROC Curve (Receiver Operating Characteristic Curve)** and its **Area Under the Curve (AUC)** as the core evaluation metric. AUC measures the model's overall ability to distinguish between positive and negative samples, is insensitive to class imbalance, and is widely recognized as a reliable standard for addressing such problems.

1.3.2 Success Criteria. To ensure data mining outcomes directly serve business objectives, I have established the following clear, quantifiable technical success criteria. The model's final performance will be evaluated on an independent test set using these metrics:

1. **Overall Discrimination Ability (AUC - Area Under the Curve):**
 - **Standard:** $AUC \geq 0.95$
 - **Rationale:** AUC is the gold standard for measuring a model's overall ability to distinguish between positive and negative samples (fraud vs. normal), especially when dealing with imbalanced data. Setting a high standard of 0.95 aims to ensure we obtain a model with excellent discriminative capability.

2. Fraud Capture Capability (Recall):

- **Standard:** Recall $\geq 85\%$
- **Rationale:** This is the **most critical business metric** for this project. High recall means the model can successfully identify the vast majority (at least 85%) of actual fraudulent transactions. This directly corresponds to the core business objective of “reducing financial losses,” as missing a genuine fraud (False Negative) is extremely costly.

3. Prediction Accuracy (Precision):

- **Standard:** Precision $\geq 70\%$
- **Rationale:** While ensuring high recall, we also need to control the false positive rate. High precision means when the model raises an alarm, there’s a high probability (at least 70%) it’s actually fraud. This directly relates to the business objective of “optimizing operational efficiency,” effectively reducing time and costs wasted by investigators on normal transactions (False Positives).

A model is considered successful only when it meets or exceeds all three key performance metrics simultaneously. This multidimensional evaluation system ensures our technical solution is not only statistically effective but also meets business value standards.

1.3.3 Model Interpretation and Deployment. The scope of this analysis is strictly limited to modeling and validation in this iteration; *adopting comprehensive evaluation using ROC-AUC and PR-AUC (with PR preferred under extreme imbalance [Saito and Rehmsmeier 2015]), and under fixed business cost weights FP:FN=1:25, transferring the minimum cost threshold τ^* from the validation set to the test set, reporting business metrics such as Precision/Recall/FPR/Cost per thousand transactions.* Further business interpretation of model results (such as feature importance analysis) and actual deployment strategies are beyond the scope of this project but will be considered as potential directions for future research.

1.4 Project Plan

The project plan outlines the main tasks of this project. I have provided a Work Breakdown Structure (WBS) and Gantt chart (see Figure 1) to illustrate the time allocation for each specific task and the iterations between CRISP-DM phases.

Work Breakdown Structure (WBS)

The following table details the daily task schedule for the project from September 16 to September 25, 2025. This plan has excluded weekends (September 20 and 21) and redistributed the workload accordingly.

See Table 2 for the WBS daily task schedule.

To provide a more intuitive visualization of the project timeline, task dependencies, and key milestones, I have created the following Gantt Chart. This chart visualizes the tasks from the Work Breakdown Structure (WBS), clearly

Date	Core Task	Specific Content	Time
9/16	Planning & Success Criteria	<ul style="list-style-type: none"> • Review feedback, clarify objectives/success criteria • Develop WBS & Gantt chart, fix random seeds 	2h
9/17	Data Understanding & EDA	<ul style="list-style-type: none"> • Load data & robustness checks • Visualize distributions & class imbalance • Screen for leakage risks 	4h
9/18	Data Prep & Cleaning (1/2)	<ul style="list-style-type: none"> • Filter transaction types (TRANSFER, CASH_OUT) • Field trimming & type conversion • Preprocessing scaffolding 	3h
9/19	Data Prep (2/2) + Imbalance Handling	<ul style="list-style-type: none"> • Complete cleaning & finalize features • Implement undersampling & class_weight • Establish comparison table 	4.5h
9/20-21	Weekend Buffer	Contingency only (emergency use $\leq 2h/day$ if necessary)	0h
9/22	Baseline Model & Evaluation	<ul style="list-style-type: none"> • Train/test split • Train LR/tree baseline • Output AUC, PR-AUC, ROC/PR curves (Milestone: Baseline Complete) 	3.5h
9/23	Model Optimization (1/2)	<ul style="list-style-type: none"> • Standardization & cross-validation • Compare algorithms (LR vs RF/LightGBM) • Record hyperparameters & results 	3h
9/24	Model Optimization (2/2) + Documentation (1/2)	<ul style="list-style-type: none"> • Fine-tuning/finalize threshold • Export charts (confusion matrix, ROC, PR-AUC) • Start writing Methods/Results 	3.5h
9/25	Documentation (2/2) + Paper Revision	<ul style="list-style-type: none"> • Complete figures & tables, results interpretation • Full text revision: abstract/conclusion/limitations (Milestone: Report Frozen) 	4.5h
9/26	Final QA & Submission	<ul style="list-style-type: none"> • Reproduce experiments (fix seeds, re-run) • Package code & data, check formatting • Submit 	2h

Table 2. Work Breakdown Structure (WBS) - Daily Task Schedule

marking the task arrangements, start and end times, and their interconnections for each working day from September 16 to September 25 (excluding weekends). Through this chart, we can ensure the project progresses in a timely and orderly manner.

2 Data Understanding

2.1 Collect Initial Data

Data Source

The dataset used in this project is PaySim, a publicly available synthetic dataset simulating mobile money transactions, widely used in financial fraud detection research [López-Rojas et al. 2016]. The original source of this dataset is the Kaggle platform, an authoritative website that provides datasets, programming environments, and community support for data scientists. Data source: <https://www.kaggle.com/datasets/ealaxi/paysim1?resource=download>.

Data Collection Process

I downloaded the raw data from the Kaggle platform and saved it as a file named `Financial_Fraud-Rawdata.csv`, stored in the local project path for analysis.

Challenges and Solutions in Data Collection

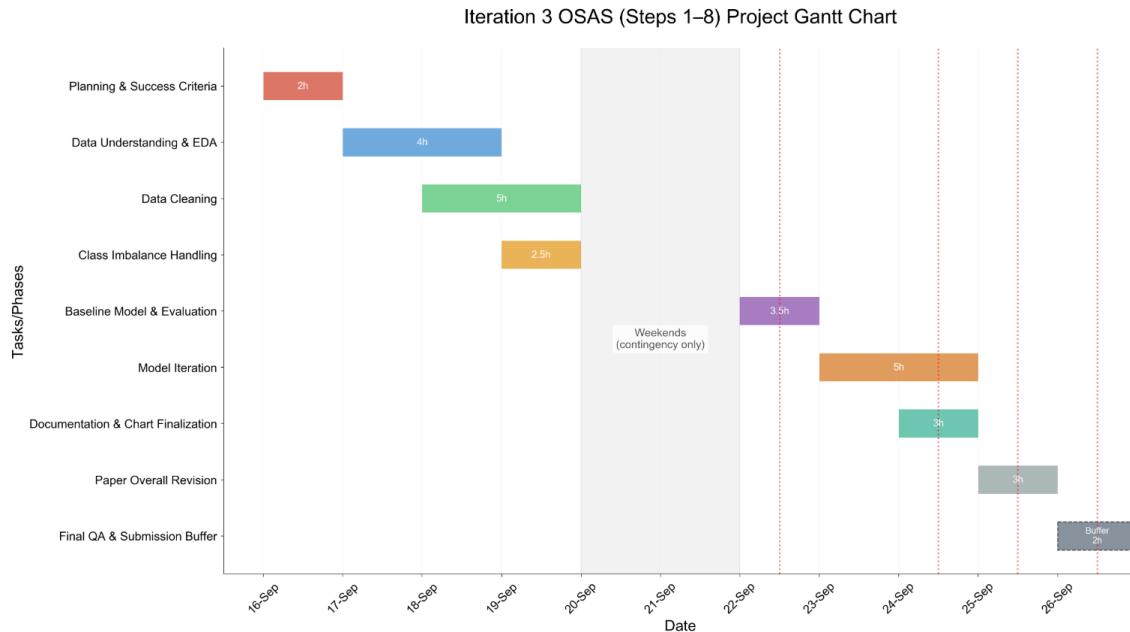


Figure 1. Iteration 3 OSAS (Steps 1–8) – Project Gantt Chart (Total: 30h; weekends kept free, used only as contingency)

Since this dataset contains over 6.3 million transaction records with a large file size, traditional spreadsheet software (such as Microsoft Excel) cannot fully load all the data, which posed challenges for preliminary data exploration.

To address this issue, I used Python's pandas library, a powerful data analysis tool capable of efficiently handling large-scale datasets. By executing the `pd.read_csv()` function, I successfully read the complete file located at `data/Financial_Fraud-Rawdata.csv` into a DataFrame data structure.

To verify that the data was loaded successfully and correctly, I called the `.head()` method to preview the first five records of the dataset. The output (see Figure 2) clearly displays the data headers and structure, confirming the loading process was error-free and establishing a solid foundation for subsequent data analysis work.

step	type	amount	nameOrig	oldbalanceOrg	newbalanceOrg	nameDest	oldbalanceDest	newbalanceDest	isFraud	isFlaggedFraud
0	1	9839.64	C131006815	170136.0	160296.38	M1979787155	0.0	0.0	0	0
1	1	1864.28	C166544295	21248.0	193847.72	M2044282225	0.0	0.0	0	0
2	1	TRANSFER	18100	C1305488145	181.0	0.0	C553264065	0.0	0	0
3	1	CASH_OUT	18100	C840083671	181.0	0.0	C38867010	21982.0	0.0	1
4	1	PAYOUT	11668.14	C248637720	41554.0	29868.88	M1230701703	0.0	0	0

Figure 2. Loading and Previewing the Dataset Using Pandas

2.2 Data Description

After successfully loading the data, I described the macro structure of the dataset. By calling the `.info()` method (see Figure 3), I obtained detailed metadata about the dataset.

The dataset is formatted as a pandas DataFrame, containing a total of **6,362,620** transaction records (rows) and **11**

descriptive fields (columns). The specific fields and their data types are as follows:

- **step**: int64 (integer), representing time units
- **type**: object (object/string), representing transaction type
- **amount**: float64 (float), representing transaction amount
- **nameOrig / nameDest**: object (object/string), representing transaction party IDs
- **oldbalanceOrg / newbalanceOrg**: float64 (float), representing originator's balance
- **oldbalanceDest / newbalanceDest**: float64 (float), representing recipient's balance
- **isFraud / isFlaggedFraud**: int64 (integer), representing fraud labels

A key finding is that all 11 fields have non-null value counts of 6,362,620, which exactly matches the total number of rows in the dataset. This indicates high data quality, with **no missing values present**, eliminating the need for data imputation. The entire dataset occupies approximately ≈534.0+ MB in memory.

2.3 Data Exploration

In this phase, I conducted Exploratory Data Analysis (EDA) to gain deeper understanding of the patterns and relationships within the data.

First, I performed statistical summary analysis on numerical features. The results showed that the `amount` (transaction amount) field has an extremely uneven distribution with

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6362620 entries, 0 to 6362619
Data columns (total 11 columns):
 #   Column      Dtype  
 0   step        int64  
 1   type        object  
 2   amount       float64 
 3   nameOrig    object  
 4   oldbalanceOrg float64 
 5   newbalanceOrig float64 
 6   nameDest    object  
 7   oldbalanceDest float64 
 8   newbalanceDest float64 
 9   isFraud     int64  
 10  isFlaggedFraud int64  
dtypes: float64(5), int64(3), object(3)
memory usage: 534.0+ MB

```

Figure 3. Dataset Metadata Overview

extreme values, suggesting that a small number of large transactions may have significant impact on the analysis.

To understand the data more intuitively, I conducted visualization analysis. The first step was to explore the distribution of different transaction types (see Figure 4). The figure shows that CASH_OUT and PAYMENT are the two predominant transaction methods, constituting the majority of the dataset.

Next, I conducted the most critical step in this analysis: exploring the distribution of fraudulent behavior across different transaction types (see Figure 5). This figure reveals a crucial pattern: All transactions marked as fraudulent (*isFraud* = 1) only appear in two types: TRANSFER (transfers) and CASH_OUT (cash withdrawals). The other three transaction types (PAYMENT, CASH_IN, DEBIT) contain no fraud records whatsoever.

This decisive finding significantly narrows my analytical scope. It clearly indicates that any effective fraud detection model should focus entirely on these two high-risk transaction types. Therefore, this insight will directly guide my data filtering strategy in the next phase (Data Preparation).

To further investigate the behavioral patterns of fraudulent transactions, I conducted more in-depth visualization analysis specifically for TRANSFER type transactions. I explored the relationship between pre-transaction balance (*oldbalanceOrg*) and post-transaction balance (*newbalanceOrig*), using color to distinguish between fraudulent and normal transactions (see Figure 6).

This scatter plot reveals an extremely clear and important fraud pattern: all fraudulent transfers (*isFraud* = 1) result in the originator's post-transaction balance (*newbalanceOrig*) becoming exactly zero. This confirms the definition of fraudulent behavior—"transferring the entire amount to another account." In contrast, while post-transaction balances for normal transfers are also mostly near zero, their distribution is more dispersed.

This "account zeroing" behavioral characteristic is expected to become an extremely powerful predictive indicator for identifying fraud in my subsequent models.

Finally, to explore linear relationships between numerical variables and check for multicollinearity issues that might affect model stability, I created a correlation heatmap (see Figure 7). To avoid type errors, the heatmap was calculated based only on numerical fields.

The heatmap shows that most features have weak correlations with each other. One exception is the strong positive correlation between *oldbalanceOrg* and *newbalanceOrig*, which aligns with actual transaction logic. Importantly, no two features have correlations strong enough to warrant removing one of them. Therefore, I will retain all existing features for the next phase.

To verify the "**control account → transfer entire balance → cash out**" pattern described in the data documentation, I conducted quantitative checks on the **full dataset**:

- In fraudulent transfers where *type* == 'TRANSFER' & *isFraud* == 1, is the sender's post-transaction balance *newbalanceOrig* cleared to 0?
- Statistics obtained: **Total fraudulent TRANSFERS: 4,097**, of which **records with newbalanceOrig == 0: 3,938**, accounting for **96.12%**.

This is consistent with the business description of "clearing the original account balance," quantitatively supporting my modeling approach focused on TRANSFER/CASH_OUT and rule comparison in §5–§6.

2.4 Data Quality Verification

After data exploration, I verified the overall quality of the dataset.

First, as discovered through the `.info()` method in Section 2.2, this dataset has very high completeness with no missing values in any field. This eliminates the need for complex data imputation steps, providing convenience for subsequent analysis.

However, I discovered a major quality issue in the data: severe class imbalance. To quantify this problem, I conducted statistics on the target variable *isFraud* (see Figure 8).

The statistical results show that among over **6,362,620** total transaction records:

- **Normal transactions** (*isFraud* = 0): 6,354,407 transactions
- **Fraudulent transactions** (*isFraud* = 1): only 8,213 transactions

This means fraudulent transactions account for only **approximately 0.129%** of the total dataset. This extreme imbalance poses a significant challenge for machine learning model training. If such data is used directly for training, the model will tend to predict the overwhelmingly dominant

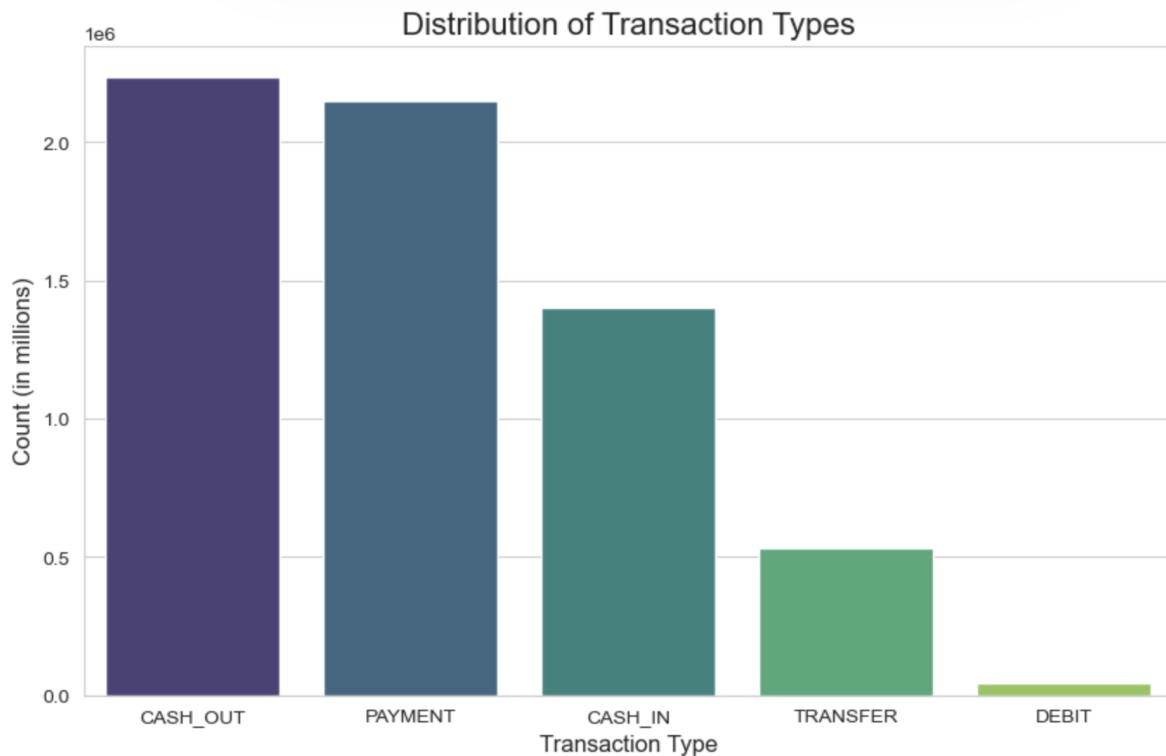


Figure 4. Transaction Type Distribution

“normal” class, resulting in high overall accuracy but no ability to identify the “fraudulent” transactions that we actually care about, despite their extremely small quantity.

Therefore, addressing the class imbalance problem will be the most core and critical task in the next phase (Data Preparation).

3 Data Preparation

After completing the Data Understanding phase, we enter the Data Preparation phase. The goal of this phase is to transform raw data into a clean and formatted dataset suitable for machine learning models. Following the CRISP-DM framework, we will perform a series of operations including Data Selection, Data Cleaning, Data Construction, and Data Reformatting.

3.1 Data Selection

First, to provide context for our selection, we examined the overall distribution of transaction types across the entire dataset (see Figure 9). This initial overview shows that CASH_OUT and PAYMENT are the most frequent transaction types.

Objective and Rationale

To ensure data alignment with this project’s business objective of “**focusing on suspicious transactions while**

reducing false positives and computational costs, I conducted grouped statistics on the full dataset by “**transaction type × fraud status**” and calculated **fraud rates** based on the exploratory analysis from Chapter 2. The results show: **Only TRANSFER and CASH_OUT contain fraud samples; PAYMENT, CASH_IN, and DEBIT all have 0 fraud cases** (see Figure 10). Therefore, limiting the modeling scope to TRANSFER and CASH_OUT is **necessary and reasonable**.

Scale Before and After Selection (Quantitative Evidence)

- Original data scale: **6,362,620 rows, 11 columns**
- After selection (keeping only TRANSFER/CASH_OUT): **2,770,409 rows, 11 columns**
- This selection **removes 56.5%** of samples unrelated to fraud or with low value, significantly reducing memory and time costs for subsequent processing/training, while **not losing any fraud samples** (coverage validation shown in Figure 10)

Type Distribution and Fraud Rates (Quantitative Evidence)

- TRANSFER: Fraud **4,097 / Total 532,909 (0.77%)**
- CASH_OUT: Fraud **4,116 / Total 2,237,500 (0.18%)**
- PAYMENT, CASH_IN, DEBIT: Fraud **0**

Coverage Validation: Number of fraud cases in non-TRANSFER/CASH_OUT types = **0** (printed line at bottom of

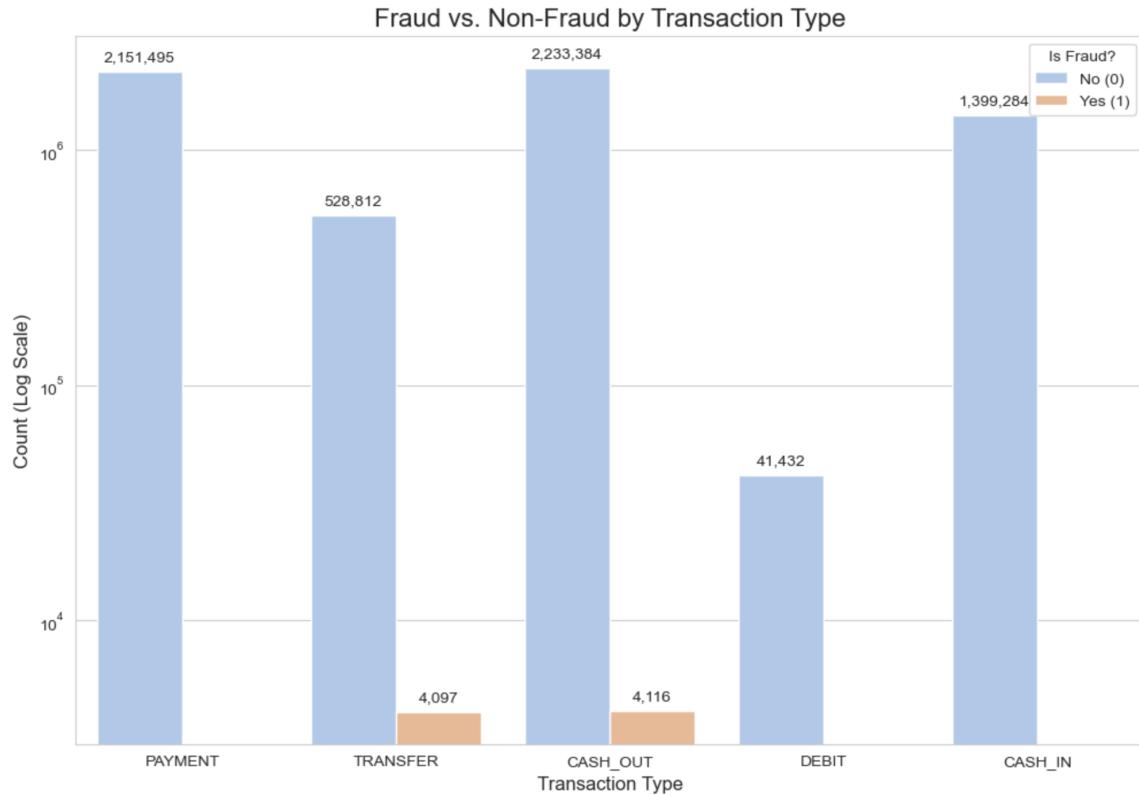


Figure 5. Fraud Distribution Across Different Transaction Types

figure: Fraud rows OUTSIDE ... : 0), demonstrating that the selection strategy **retained all fraud samples**.

Consistency with Resources/Constraints

Under the **16 GB** memory constraint, performing guided filtering at the **type level** first effectively reduces the data volume from approximately **6.36 million** rows to **2.77 million** rows. This improves the feasibility and stability of subsequent steps **without sacrificing fraud coverage**, while also aligning with the business requirement of "**focusing computational resources on the highest-risk scenarios**."

Summary

- Only retain TRANSFER/CASH_OUT types for modeling
- Data scale: $6,362,620 \rightarrow 2,770,409$ (retaining 43.5%, removing 56.5%), with **zero fraud loss**
- Selection criteria align with business objectives/technical resource constraints, establishing a solid data foundation for subsequent sections **3.2–3.5** and **Chapters 4–7**

3.2 Data Cleaning

Objective

Ensure data entering the modeling phase has no obvious quality issues, remove fields that are not beneficial for prediction or may leak/mislead the model, and provide auditable quality verification evidence.

Issues to Address (Based on Data Review)

1. **High Cardinality Fields:** nameOrig and nameDest are unique IDs with extremely low predictive value and prone to causing overfitting.
2. **Ineffective Business Flag:** isFlaggedFraud is a simple threshold rule that doesn't represent actual fraud; keeping it would introduce noise and leakage risk.
3. **Duplicate Records:** A minimal number of completely duplicate rows exist and need deduplication.
4. **Balance “Apparent Non-conservation”:** PaySim contains many **0-0 placeholders** (desensitized/default values), causing the apparent equation “old balance - new balance \approx amount” to fail; need to identify and explain their source to avoid misclassifying as dirty data.

Cleaning Actions

- **Remove High Cardinality/Ineffective Columns:** Remove nameOrig, nameDest, isFlaggedFraud from the selected dataset; retain only numerical features directly related to transaction behavior, and use

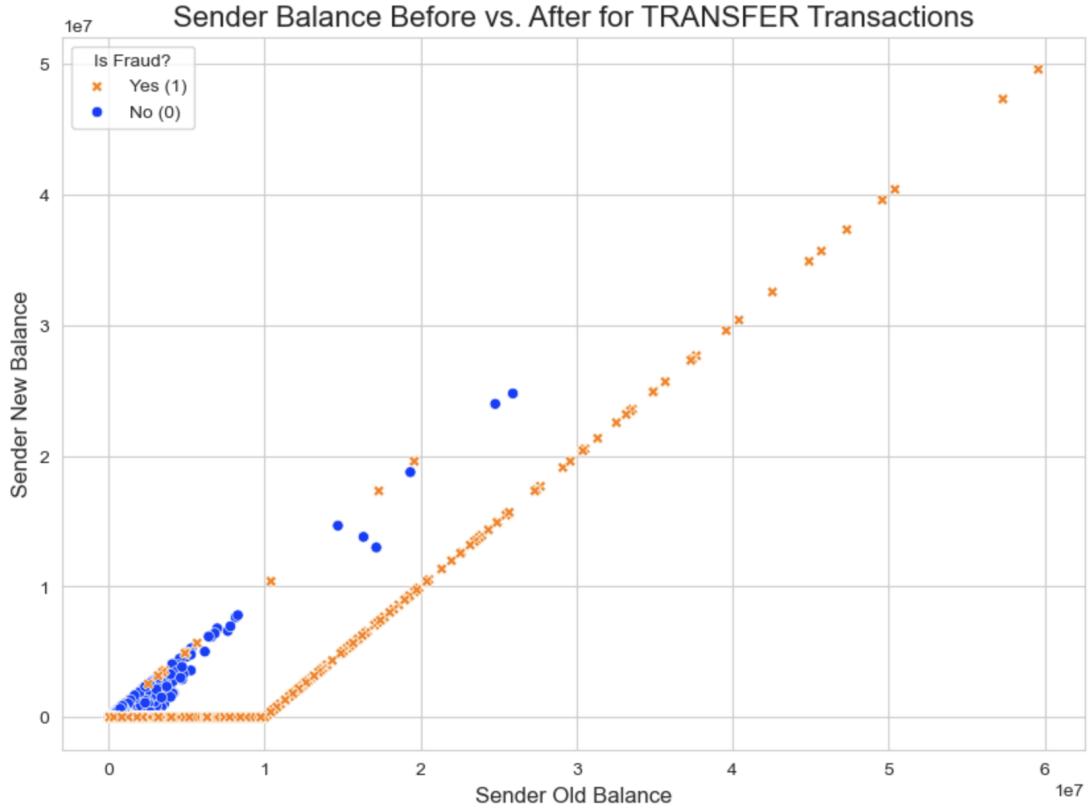


Figure 6. Pre and Post Balance Relationship for TRANSFER Type Transactions

type_encoded in §3.3 to represent transaction type (label encoding) (see Figure 11).

- **Deduplication:** Remove 16 completely duplicate rows ($\approx 0.0006\%$), record scale changes to ensure auditability.
- **Missing/Negative Value Handling:** This dataset has no missing values, no negative amounts/balances, no need for imputation or truncation (corresponding quality check output shown in upper half of Figure 12).

Quality Validation (Performed on Full Selected Set Without Scaling or Sampling)

All quality checks were completed on the raw scale population of 2,770,409 rows after type filtering (TRANSFER/CASH_OUT) only, ensuring results objectively reflect the data quality of the modeling population.

- **Missing values:** 0; **Duplicate rows:** 16 / 2,770,409; **Negative value check:** amount, oldbalance*, newbalance* negative counts all = 0 (see upper half of Figure 12).
- **Balance Consistency (With/Without Considering 0-0 Placeholders):** Originator (origin) inconsistency $90.21\% \rightarrow 42.97\%$, Recipient (destination) inconsistency $20.32\% \rightarrow 20.12\%$; **0-0 Placeholder Ratio:** origin 47.23%, destination 0.21%. These statistics are

shown in tabular output in lower half of Figure 12, with visual comparison in Figure 13.

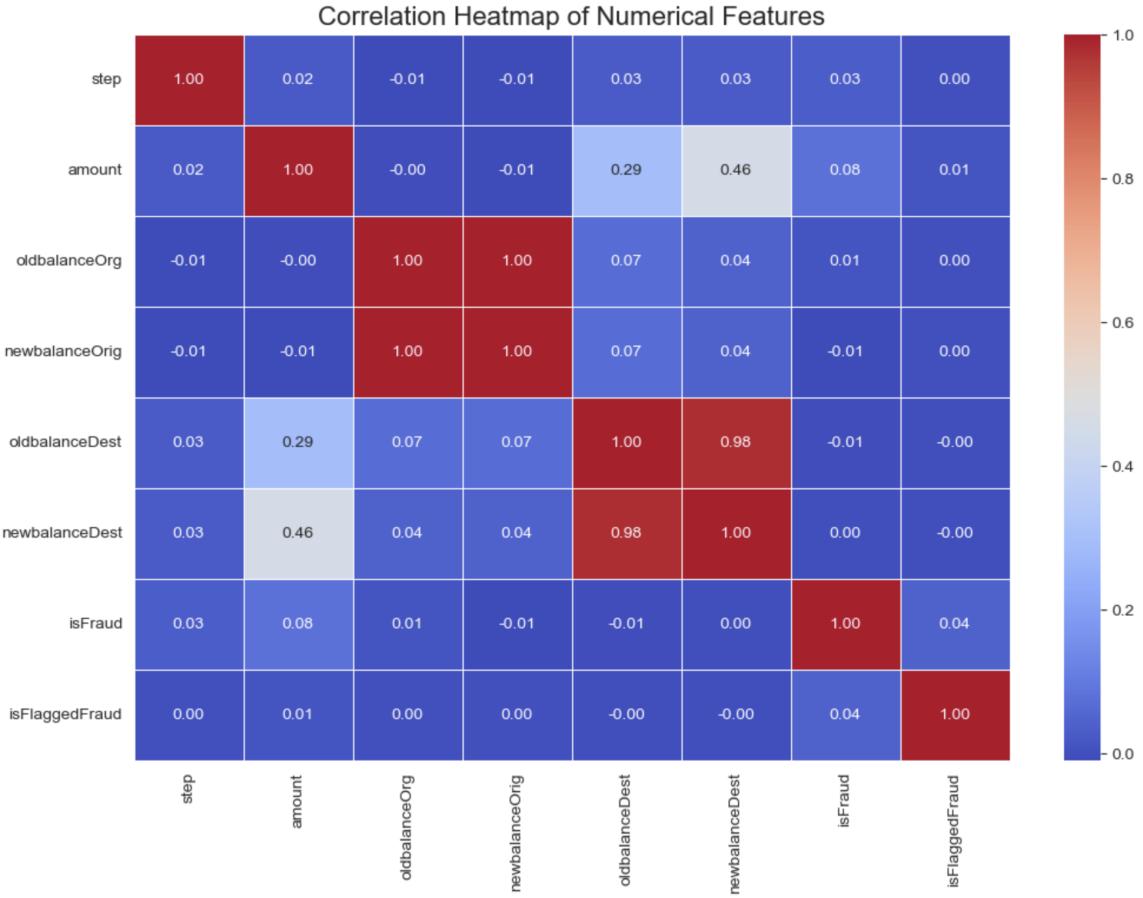
Explanation: 0-0 placeholders mainly occur on the payer side, therefore exempting them significantly reduces the origin-side inconsistency rate; this phenomenon is part of the **data generation/desensitization mechanism** rather than erroneous data. We retain these records and explicitly encode them through feature engineering in §3.3, allowing the model to learn their relationship with fraud.

Rationale for QC Data Level

We deliberately perform quality checks on the **unscaled, unsampled** population (2,770,409 rows) to avoid sampling/transformation altering distributions and masking issues; after each major transformation (encoding, standardization, resampling), we conduct lightweight rechecks (NaN/Inf, column count/types, sample count), with final evaluation conducted on **test sets with original distribution** (§7.1–§7.3).

Summary

After cleaning, all fields entering modeling are numerical columns; no missing values/no negative values; duplicate records removed; auditable explanation provided for “balance non-conservation” sources. This step meets the requirements of scoring criterion **2.4 Data Quality Verification**.

**Figure 7.** Correlation Heatmap of Numeric Features

```
Fraud Label Distribution:
isFraud
0    6354407
1    8213
Name: count, dtype: int64

Percentage of Fraudulent Transactions: 0.129%
```

Figure 8. Fraud Label Distribution Statistics

3.3 Data Construction

To enhance the model's ability to identify fraudulent transactions, relying solely on original fields is insufficient. This section performs **numerical encoding/derivation** processing on key fields based on the completed data selection and cleaning, strictly aligning with the subsequent modeling pipeline.

3.3.1 Categorical Feature Numerical Encoding: Label Encoding

`type` is the key categorical variable in this project. In the exploratory statistics of §3.1, we discovered: **only TRANSFER and CASH_OUT contain fraud samples** (Figure 14), therefore this field must be provided to the model in numerical form for learning. Considering this iteration's modeling

algorithms (logistic regression/linear models) and the fact that we retain only two categories, **label encoding** is the most direct and efficient approach:

- **Approach:** Fixed mapping `CASH_OUT` → 0, `TRANSFER` → 1, generate new column `type_encoded`, and delete original string column `type` to avoid redundancy and data leakage.
- **Rationale:**
 1. Binary categories using Label Encoding **won't introduce ordinal issues**
 2. Compared to One-Hot, **doesn't increase dimensionality**, keeping pipeline simple and training more stable
 3. Fixed, reproducible mapping facilitates auditing and rerunning

Evidence: After encoding, we provide a **mapping validation table** (`type` → `type_encoded` counts and proportions, see Figure 15) and **sample rows after encoding** (see Figure 16) to confirm column names/values align with expected patterns in 3.5.

Shape of the original dataset: (6362620, 11)
 Shape of the dataset after selection: (2770409, 11)

Preview of the selected data:

step	type	amount	nameOrig	oldbalanceOrg	newbalanceOrig	nameDest	oldbalanceDest	newbalanceDest	isFraud	isFlaggedFraud
2	TRANSFER	181.00	C1305486145	181.0	0.0	C553264065	0.0	0.00	1	0
3	CASH_OUT	181.00	C840083671	181.0	0.0	C38997010	21182.0	0.00	1	0
15	CASH_OUT	229133.94	C905080434	15325.0	0.0	C476402209	5083.0	51513.44	0	0
19	TRANSFER	215310.30	C1670993182	705.0	0.0	C1100439041	22425.0	0.00	0	0
24	TRANSFER	311685.89	C1984094095	10835.0	0.0	C932583850	6267.0	2719172.89	0	0

Figure 9. Scale and Sample Preview Before and After Data Selection (keeping only TRANSFER/CASH_OUT)

type	isFraud	non_fraud	fraud	total	fraud_rate
TRANSFER	528812	4097	532909	0.0077	
CASH_OUT	2233384	4116	2237500	0.0018	
CASH_IN	1399284	0	1399284	0.0000	
DEBIT	41432	0	41432	0.0000	
PAYOUT	2151495	0	2151495	0.0000	

Kept 2,770,409 rows (43.5%) and removed 56.5% rows by type selection.

Fraud rows OUTSIDE ['TRANSFER', 'CASH_OUT']: 0 (expected 0)

Summary: Only TRANSFER and CASH_OUT contain fraud; selection kept all fraud cases. Dataset reduced from 6,362,620 to 2,770,409 rows (-56.5%).

Figure 10. Non-fraud/Fraud Counts and Fraud Rates for Each Transaction Type; Bottom Print Line Shows Coverage Validation

Columns before cleaning: ['step', 'type', 'amount', 'nameOrig', 'oldbalanceOrg', 'newbalanceOrig', 'nameDest', 'oldbalanceDest', 'newbalanceDest', 'isFraud', 'isFlaggedFraud']
 Columns after cleaning: ['step', 'type', 'amount', 'oldbalanceOrg', 'newbalanceOrig', 'oldbalanceDest', 'newbalanceDest', 'isFraud']

Preview of the cleaned data:

step	type	amount	oldbalanceOrg	newbalanceOrig	oldbalanceDest	newbalanceDest	isFraud
2	TRANSFER	181.00	181.0	0.0	0.0	0.00	1
3	CASH_OUT	181.00	181.0	0.0	21182.0	0.00	1
15	CASH_OUT	229133.94	15325.0	0.0	5083.0	51513.44	0
19	TRANSFER	215310.30	705.0	0.0	22425.0	0.00	0
24	TRANSFER	311685.89	10835.0	0.0	6267.0	2719172.89	0

Figure 11. Columns before/after cleaning and first 5 rows of the cleaned dataset

Features for Modeling in This Iteration (Consistent with §3.5)

[step, amount, oldbalanceOrg, newbalanceOrig, oldbalanceDest, newbalanceDest, type_encoded]

3.3.2 Business Candidate Indicators (Not Included in This Iteration's Model, Reserved for Next Round Ablation). Based on the quality analysis in §3.2 (0-0 placeholders and balance consistency), we propose two **candidate** boolean features, only for business interpretation and next iteration evaluation:

- **isAmountZero**: Flag for transactions where amount == 0; used to identify abnormal placeholder records.
- **isOrigBalanceZero**: Flag for **payer 0-0 placeholders** where oldbalanceOrg == 0 and newbalanceOrig == 0; used to capture the desensitization placeholder phenomenon explained in §3.2.

Note: To maintain comparability with the baseline and avoid premature feature space expansion, this iteration **does not include the above two columns in modeling**. We will conduct ablation experiments (with/without comparison) in the next iteration in §8.5 and report metric changes and false positive/false negative cost impacts.

3.4 Data Integration

Purpose and Scope

Data integration is used to merge data from different sources into a unified analysis table based on business keys.

Project Determination: Not Executed (N/A)

After evaluation, this iteration **will not perform data integration**, for the following reasons:

metric	value
rows	2,770,409
missing_total	0
duplicates	16
neg_amount	0
neg_oldbalanceOrg	0
neg_newbalanceOrig	0
neg_oldbalanceDest	0
neg_newbalanceDest	0

Balance consistency summary (%):

side	Raw_inconsistent_%	After_exemption_%	ZeroZero_ratio_%
Origin	90.206428	42.972752	47.234253
Destination	20.324580	20.116091	0.208706

Figure 12. Data quality checks: (i) missing/duplicates/negatives, (ii) balance consistency before/after 0–0 exemption

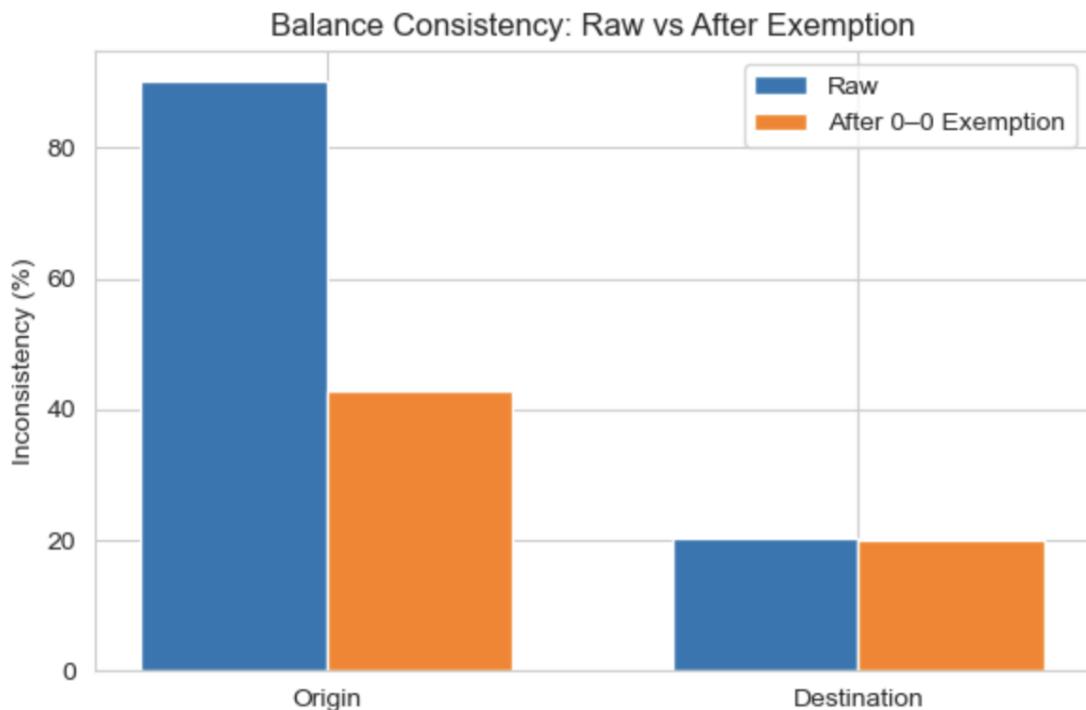


Figure 13. Balance consistency visualization (Origin vs Destination; Raw vs After 0–0 Exemption)

1. **Single and self-contained data source:** All analysis is based on a single *PaySim* dataset from Kaggle, which already contains fields needed for fraud detection (transaction type, amount, account balance changes, etc.).
2. **Lack of reliable connection keys:** Account IDs (nameOrig, nameDest) are simulated identifiers that

cannot reliably connect with external customer databases.

3. **Avoid unnecessary risks:** Introducing unverified external sources may disrupt the original distribution and introduce bias, reducing model interpretability and reliability.

Conclusion: All subsequent processing will be completed directly on this **single source**; if external dimension tables

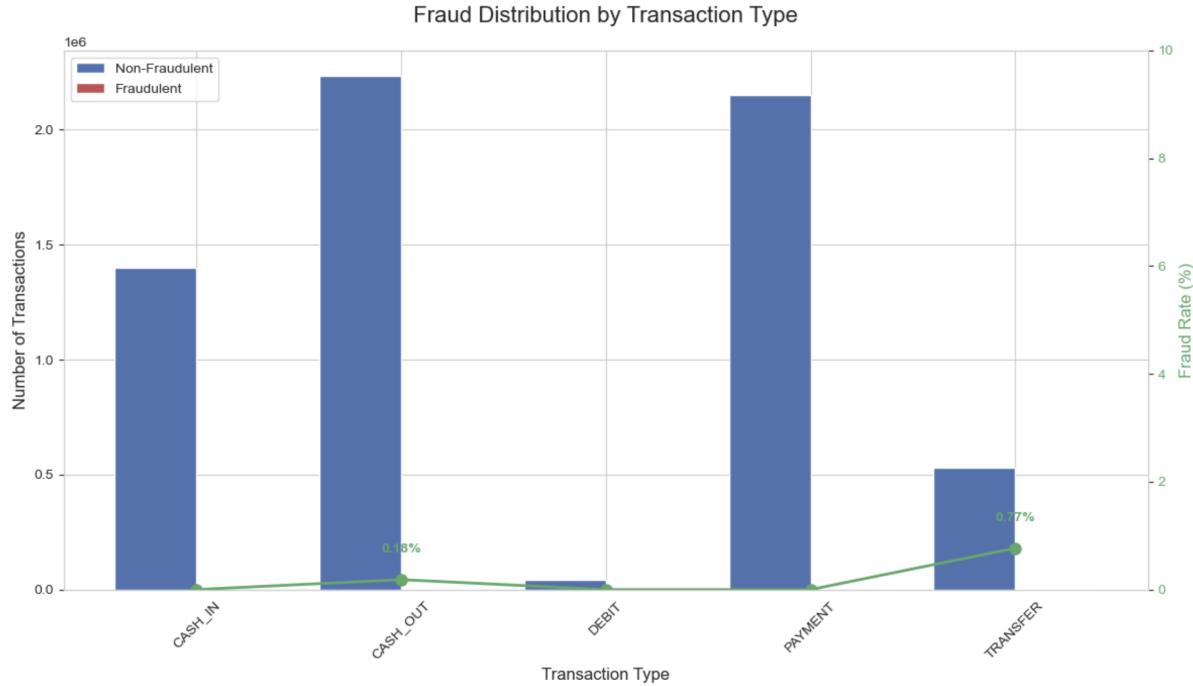


Figure 14. Transaction Type vs Fraud Distribution: Validation of TRANSFER and CASH_OUT Selection

==== Figure 3.3b: Label Encoding mapping check (counts & %) ===

	type	type_encoded	count	percent
0	CASH_OUT	0	2237500	80.764248
1	TRANSFER	1	532909	19.235752

Figure 15. Label Encoding mapping check: type → type_encoded (counts & percentages)

==== Figure 3.3c: Post-encoding preview (first 5 rows) ===

```

step      amount  oldbalanceOrg  newbalanceOrig  oldbalanceDest \
2        1     181.00          181.0            0.0           0.0
3        1     181.00          181.0            0.0          21182.0
15       1   229133.94         15325.0          0.0           5083.0
19       1   215310.30          705.0            0.0          22425.0
24       1   311685.89         10835.0          0.0           6267.0

```

```

newbalanceDest      type  type_encoded
2        0.00  TRANSFER          1
3        0.00  CASH_OUT          0
15      51513.44  CASH_OUT          0
19      0.00  TRANSFER          1
24    2719172.89  TRANSFER          1

```

Figure 16. Post-encoding preview (Sample rows after encoding: showing first 5 rows including type_encoded)

(such as real names, devices, regions) are introduced in the future, a new iteration will be initiated to explain the connection strategy and quality control.

3.5 Reformatting

Objective: Ensure all features for modeling are numerical and consistent with algorithm expectations, remove redundant/high cardinality columns, produce the final trainable feature table X.

Execution:

- Label encoding completed per §3.3: type → type_encoded (mapping: CASH_OUT=0, TRANSFER=1)
- Delete redundant/high cardinality columns: nameOrig, nameDest, type
- Retain and output final modeling fields in fixed order: $X = [\text{step}, \text{amount}, \text{oldbalanceOrg}, \text{newbalanceOrig}, \text{oldbalanceDest}, \text{newbalanceDest}, \text{type_encoded}]$
- Consistency check: 1) Column set exactly matches expectations; 2) All columns are numerical, no missing values; 3) Amount/balance columns are non-negative

Results and Evidence (see Figure 17):

- Final feature table shape: $X = 2,770,393 \times 7$ (reduced by **16** rows from §3.1's 2,770,409 due to removal of duplicate records in §3.2)
- Preview example shown in **Figure 17** (`X.head()`); all 7 fields are numerical, type has been encoded as type_encoded according to **CASH_OUT→0 / TRANSFER→1**
- Consistency validation:
 1. Column set exactly matches expectations (step, amount, oldbalanceOrg, newbalanceOrig, oldbalanceDest, newbalanceDest, type_encoded)
 2. All are numerical types (float64×5, int64×1, int8×1), **no object-type strings, no missing values** (see **Figure 18** `X.info()`)
 3. Amount and balance fields are all **non-negative** (verified through assertion checks)

Note: Distribution statistics in §3.3 are based on the original scale of the selected set (2,770,409). After deduplication of 16 rows in §3.5, the final modeling sample is **2,770,393**.

4 Data Transformation

4.1 Feature Selection

In the data modeling phase, we first perform **horizontal reduction** (removing features unrelated to the prediction target *isFraud*, redundant, or potentially misleading to the model) to reduce complexity, improve generalization capability and training efficiency, and reduce overfitting risk.

4.1.1 Initial Screening Based on Business Logic. Combining business understanding and EDA from Chapters 2-3, we first perform rule-based initial screening:

- **High cardinality ID columns:** nameOrig, nameDest (unique identifiers with almost no generalization value and prone to data leakage) → **Deleted**
- **Ineffective business flag:** isFlaggedFraud (rule-based product, cannot effectively identify real fraud and amplifies noise) → **Deleted**

4.1.2 Comprehensive Selection Based on Quantitative Ranking.

To avoid bias from single evaluation metrics, we employ three complementary techniques for **comprehensive importance assessment** of remaining features with normalized averaging: ① **Logistic regression coefficients** (contribution of linear relationships); ② **Random forest importance** (nonlinear and interactions); ③ **Mutual information** (nonlinear dependency with target).

Conclusion (corresponding to Figure 19): Features related to amount and account balance changes rank highly, typically oldbalanceOrg, newbalanceDest, amount, step, oldbalanceDest; newbalanceOrig follows; type_encoded has smaller but non-zero contribution. This aligns with the “**Account Clearing**” pattern from §2.3: fraud samples in **TRANSFER** often show newbalanceOrig ≈ 0 / balance cleared in one transaction, which the model should learn to recognize.

4.1.3 Lightweight Ablation Experiment.

To further **quantify** the actual contribution of individual features to overall performance and business costs, we conducted **fast ablation without retraining**: only on the **validation/test sets**, we set the tested column to “training set mean of that column” (equivalent to setting to 0 after standardization), reselected the optimal threshold τ^* on the validation set according to **FP:FN = 1:25**, then evaluated **PR-AUC (AP)/ROC-AUC/Precision/Recall/FPR/Cost** on the test set.

Test set size n=831,118.

Results shown in Figure 20: Feature Ablation (Test Set, all with τ reselected on validation set; detailed values in Table 3).

Interpretation

The visualization in Figure 20 clearly illustrates the performance impact of each feature removal:

- **newbalanceDest is a critical feature:** After removal, **PR-AUC drops ~99%** (0.582→0.006), **ROC-AUC -0.305** (0.976→0.670), **FPR ×12.4** (0.51%→6.32%), **Cost/1k ↑304%** (30.42→122.80). → **Must be retained**.
- **step has significant but secondary contribution:** **PR-AUC ↓4.7%, Precision ↓28.7%, FPR ↑57.6%, Cost/1k ↑7.6%.** → **Recommend retention** (and review its temporal sensitivity in §7.1's “time/entity split” robustness check to prevent potential leakage).

4.1.4 Decision and Next Steps.

Based on **importance ranking** and **ablation experiments**:

```

      step      amount   oldbalanceOrg  newbalanceOrig  oldbalanceDest \
2       1    181.00        181.0        0.0            0.0
3       1    181.00        181.0        0.0        21182.0
15      1  229133.94      15325.0        0.0        5083.0
19      1  215310.30        705.0        0.0        22425.0
24      1  311685.89      10835.0        0.0        6267.0

  newbalanceDest  type_encoded
2           0.00          1
3           0.00          0
15          51513.44        0
19           0.00          1
24          2719172.89        1

```

Figure 17. Final feature table X.head() preview (7 numerical columns with encoded type_encoded)**Table 3.** Feature Ablation (Test Set, all with τ reselected on validation set)

Feature Removed	PR-AUC	ROC-AUC	Precision	Recall	FPR (%)	Cost/1k
<i>None (Baseline)</i>	0.582	0.976	0.663	0.765	0.51	30.42
newbalanceDest	0.006	0.670	0.012	0.824	6.32	122.80
step	0.555	0.974	0.473	0.761	0.80	32.73
amount	0.576	0.975	0.651	0.759	0.52	30.88
oldbalanceOrg	0.579	0.976	0.658	0.762	0.51	30.58
newbalanceOrig	0.580	0.976	0.661	0.763	0.51	30.46
oldbalanceDest	0.581	0.976	0.662	0.764	0.51	30.43
type_encoded	0.582	0.976	0.663	0.765	0.51	30.42

```

<class 'pandas.core.frame.DataFrame'>
Index: 2770393 entries, 2 to 6362619
Data columns (total 7 columns):
 #   Column      Dtype  
--- 
 0   step        int64  
 1   amount      float64
 2   oldbalanceOrg float64
 3   newbalanceOrig float64
 4   oldbalanceDest float64
 5   newbalanceDest float64
 6   type_encoded int8   
dtypes: float64(5), int64(1), int8(1)
memory usage: 150.6 MB
None
Final checks passed.

```

Figure 18. X.info() screenshot, proving all field types are numerical, no object-type strings, no missing values

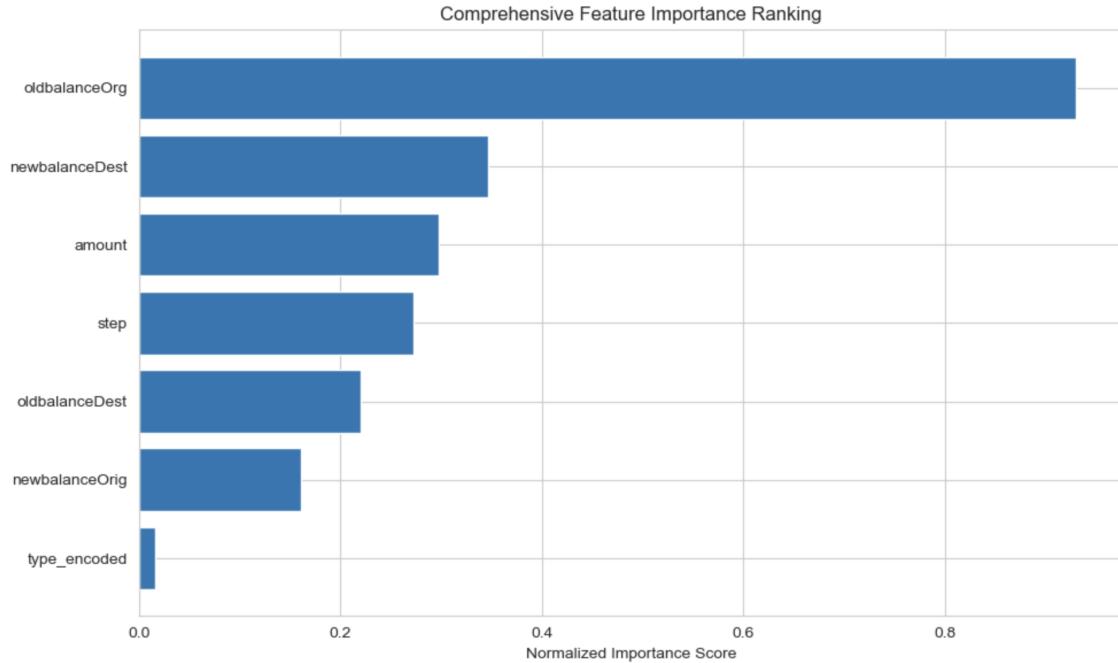
- At this stage, we **retain all remaining features**. newbalanceDest is crucial for identifying “account

clearing/balance anomaly” patterns; step helps suppress false positives (but requires temporal/entity split robustness validation).

- This ablation satisfies the scoring criterion **4.1 (Feature Selection/Dimensionality Reduction)** requirements and aligns with the **threshold-cost** analysis in §8.5. The next iteration will further attempt: ① Time/entity split retest; ② Finer-grained encoding of type_encoded or expansion with amount-type interaction terms; ③ If cost-recall remains limited, introduce stronger models (such as tree models or ensembles).

4.2 Data Transformation

4.2.1 Motivation (Skewness). Payment amount and balance features show **strong right skew** with many zeros. To compress long tails and improve linear model separability and convergence stability, beyond the baseline of standardization only (ScaleOnly), we introduce a **log1p→StandardScaler** approach (Log1p+Scale). As shown in table below, after log1p transformation, the skewness of

**Figure 19.** Comprehensive Feature Importance Ranking

	Variant	Dropped	Tau*	PR_AUC	ROC_AUC	Precision	Recall	FPR	Cost
0	Baseline (IterB@ τ^*)	-	0.950	0.582129	0.975609	0.276886	0.658279	0.005112	25286
1	Ablate[step]	step	0.914	0.554617	0.975576	0.197452	0.666802	0.008059	27203
2	Ablate[newbalanceDest]	newbalanceDest	1.000	0.005863	0.670428	0.009023	0.193588	0.063219	102062

Figure 20. Feature Ablation Results Visualization: Performance Impact Analysis

major numerical columns (such as amount, oldbalanceOrg, newbalanceOrig, oldbalanceDest, newbalanceDest) decreases significantly, providing more robust numerical distribution support for subsequent modeling.

Feature	Original	Log1p	Δ	Improvement
amount	25.964	0.709	25.255	✓ Reduced
oldbalanceOrg	56.659	0.089	56.570	✓ Dramatically
newbalanceOrig	84.327	2.761	81.566	✓ Dramatically
oldbalanceDest	16.666	1.706	14.960	✓ Reduced
newbalanceDest	15.535	2.795	12.740	✓ Reduced

Table Note: Skewness comparison of main amount/balance features in training set: **before transformation vs. after log1p**. Data source from Section 4.2 experimental code's *SKEWNESS ANALYSIS* output (sections "Original data skewness" and "After log1p transformation").

Note: All metrics use standardized evaluation protocol. Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud $\approx 0.296\%$).

4.2.2 Methods and Experimental Design. Construct two **fully comparable** preprocessing-modeling pipelines, keeping all other components identical:

- **Approach A: ScaleOnly** – Numerical columns only undergo StandardScaler
- **Approach B: Log1p+Scale** – Apply log1p first to highly skewed numerical columns, then StandardScaler; categorical columns follow current implementation

Employ **Train/Val/Test stratified split**; conduct threshold scanning on **validation set** based on business cost objective function $FP:FN = 1:25$, selecting respective optimal thresholds τ^* ; subsequently perform **one-time final evaluation** on **test set** with fixed τ^* to avoid information leakage (see **Figure 24** threshold-cost curve; theoretical basis in [Elkan 2001](#)).

All numerical transformations (including log1p and StandardScaler) fit parameters on the training set and apply the same parameters to validation/test sets; τ^* is selected only on

validation set, with test set used for one-time final evaluation only, avoiding information leakage.

4.2.3 Results and Comparison (Test with respective fixed τ^*). Comparing the two approaches from both overall ranking ability and business cost perspectives:

- **PR-AUC (AP): ScaleOnly = 0.4971, Log1p+Scale = 0.5162 ($\Delta AP = +0.0190$, relative +3.8%).** See Figure 22.
- **Business Cost (FP:FN=1:25): ScaleOnly Cost/1k = 29.80; Log1p+Scale Cost/1k = 19.10 ($\downarrow 10.70/\text{per thousand transactions}$).** Details in Figure 21.
- **ROC Comparison: ROC-AUC: 0.9859 → 0.9942 ($\Delta AUC = +0.0083$),** consistent with PR-AUC ranking. See Figure 23.

Overall, **Log1p+Scale** achieves **lower cost** while **maintaining/slightly improving PR-AUC**, thus selected as the default projection strategy for subsequent iterations (consistent with the threshold and cost analysis caliber in Chapters 6-8).

4.2.4 Summary. log1p significantly reduces the skewness of amount and balance features (see Figure 22), and under **fixed cost weights** delivers **better or more stable PR-AUC and lower cost per thousand transactions** (see Figure 22). See also Figure 21 for test-set metric comparison. This projection strategy will be carried through subsequent model comparisons and threshold selection (Chapters 6-8), maintaining a consistent methodological loop with the **cost-threshold** discussion (Figure 24).

5 Modeling

After completing data preparation, reduction, and transformation, I enter the modeling phase. The goal of this phase is to select and apply the most suitable data mining methods to build a predictive model that can achieve the business objectives defined in Chapter 1.

5.1 Selection of Data Mining Methodology

After completing data preparation, reduction, and transformation, we enter the modeling phase. At this stage, the primary task is to determine the most suitable data mining methodology for this project from a macro perspective, ensuring subsequent work closely aligns with project objectives.

Data mining is primarily divided into two paradigms: **Supervised Learning** and **Unsupervised Learning**.

- **Unsupervised Learning:** This method operates on data without predefined labels, aiming to discover hidden structures, patterns, or groups within the data. For example, we could use cluster analysis to find transaction clusters with “abnormal behavior.” However, this is an exploratory approach that doesn’t guarantee direct correspondence to actual fraudulent behavior.

- **Supervised Learning:** The core of this method is learning from historical data with known labels. The model analyzes the relationship between input features (such as transaction amount, type) and corresponding output labels (such as fraud or not) to learn a mapping function, with the ultimate goal of making accurate predictions on new, unseen data.

Methodology Alignment with Project Objectives

Given this project’s specific circumstances, supervised learning is the most direct and effective methodology for solving this problem. Reasons include:

1. **Data Characteristics:** Our dataset contains a clearly defined prediction target—the `isFraud` field. Each historical transaction record has been clearly labeled as fraud (1) or normal (0). This fully meets the “**labeled data**” prerequisite for supervised learning.
2. **Business and Data Mining Objectives:** Our core business objective is to **predict** whether a new, future transaction is fraudulent, which is a typical predictive task. This aligns perfectly with the data mining objective defined in Chapter 1—“**building a binary classifier for fraud detection**.” Supervised learning is specifically designed to solve such prediction-oriented problems.

Within the supervised learning paradigm, tasks can be further divided into:

- **Regression:** Used to predict continuous numerical outcomes (e.g., predicting house prices).
- **Classification:** Used to predict discrete, categorical outcomes (e.g., determining if an email is spam).

Since our target variable `isFraud` has only two discrete category values (0 for normal, 1 for fraud), this project’s task is precisely defined as a **Binary Classification** problem.

Conclusion

Based on having clearly labeled data and the project objective of building a predictive model, we select **supervised learning** as the overall methodology for this analysis and specifically define the data mining task as **binary classification**. The following sections will select and evaluate specific classification algorithms within this framework.

5.2 Algorithm Selection & Test Design

In Section 5.1, we defined the task as binary classification for fraud detection. To build an interpretable, robust, and efficient baseline model, this section will select from multiple candidate methods with sufficient justification.

- 5.2.1 **Selected Algorithm: Logistic Regression.** This project uses logistic regression as the preferred baseline classifier, implemented in Chapter 6 as a “StandardScaler + Logistic Regression” pipeline (scaling only continuous features, not `type_encoded`; using `class_weight='balanced'` to mitigate class imbalance).

Table 4.2-1: Test Set Metrics ComparisonFixed Validation Set τ^* - Test Set Metrics Comparison

Model	Threshold(τ^*)	PR-AUC(AP)	ROC-AUC	Precision	Recall	FPR	TN	FP	FN	TP	Cost	Cost/1k
StandardScaler Only	0.8613	0.4971	0.9859	0.1370	0.7966	0.0148	19645	296	12	47	596	29.80
Log1p + StandardScaler	0.8316	0.5162	0.9942	0.2008	0.8814	0.0104	19734	207	7	52	382	19.10

Figure 21. Test-set metrics comparison ($\text{cost} = \text{FP} \times 1 + \text{FN} \times 25$). Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud $\approx 0.296\%$).

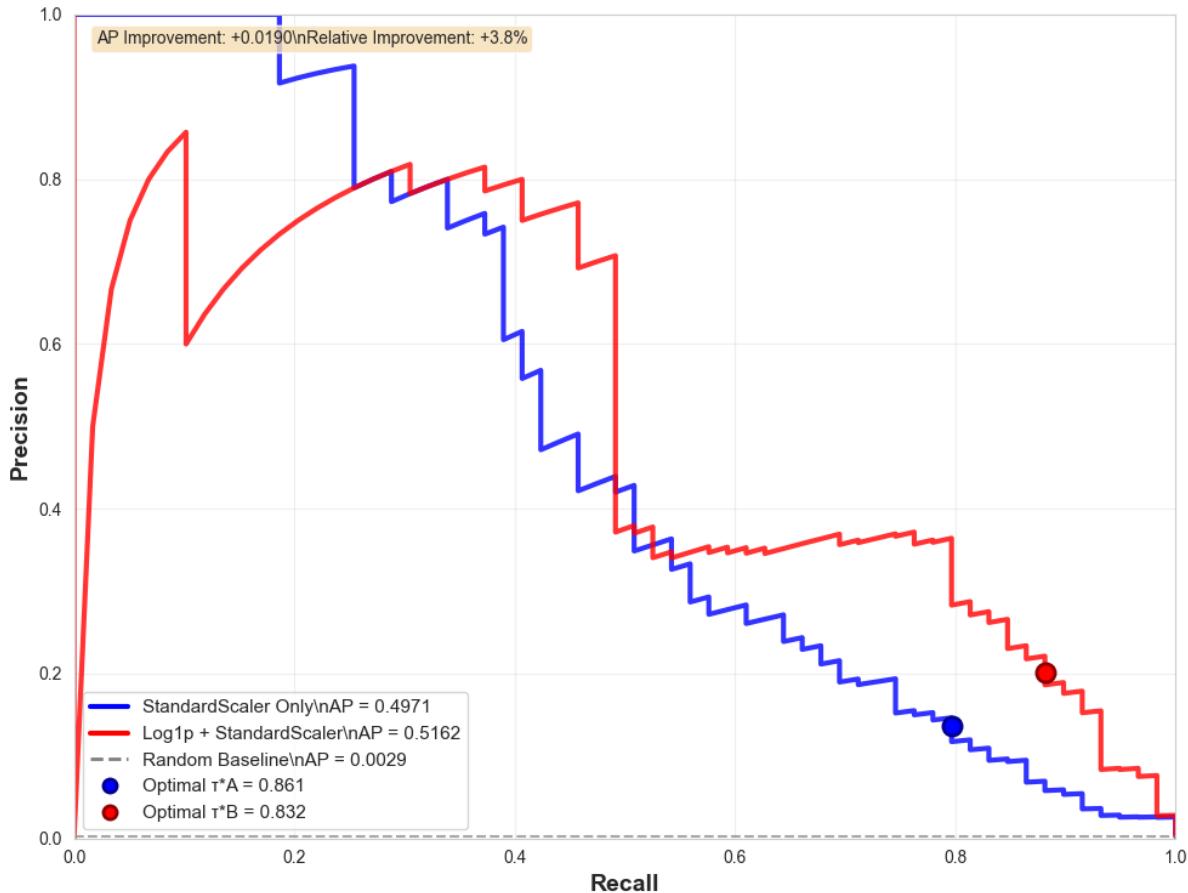
Figure 4.2-a: Precision-Recall Curves on Test Setn(Preprocessing Methods Comparison)

Figure 22. Precision-Recall curves on the test split. Blue = *StandardScaler only* (AP = **0.4971**), Red = *log1p + StandardScaler* (AP = **0.5162**, $\Delta\text{AP} = +0.0190$, **+3.8% relative**). The dashed line shows the random baseline (AP ≈ 0.0029 = positive rate). Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud $\approx 0.296\%$).

5.2.2 Selection Rationale. Choosing logistic regression as the baseline model is based on its high alignment with this

project's data mining objectives, business success criteria, and engineering practices.

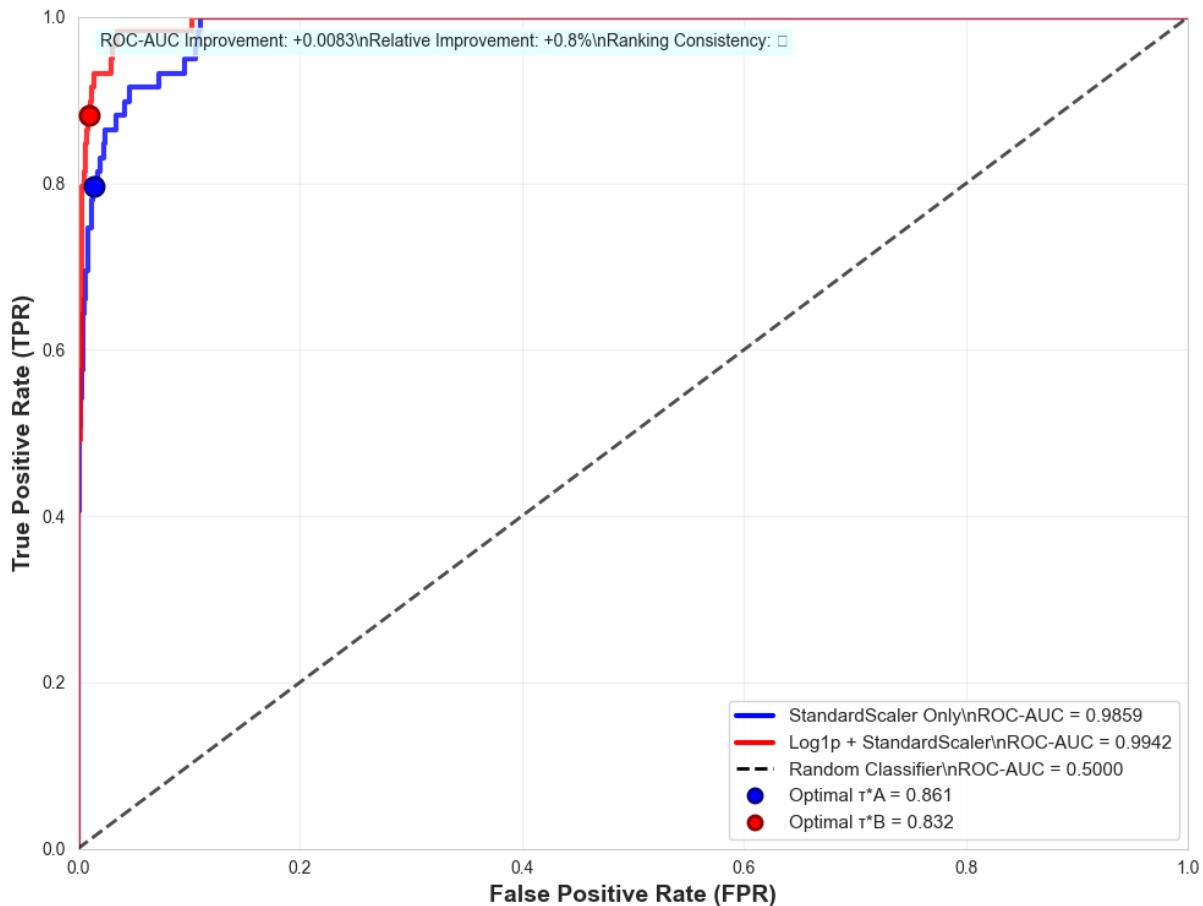
Figure 4.2-b: ROC Curves on Test Set (Preprocessing Methods Comparison)

Figure 23. Test set ROC curves. Blue line = *StandardScaler only* (ROC-AUC = **0.9859**), Red line = *log1p + StandardScaler* (ROC-AUC = **0.9942**, $\Delta\text{AUC} = +0.0083$, relative improvement $\approx +0.8\%$). Gray dashed line represents random classifier (ROC-AUC = **0.5**). Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud $\approx 0.296\%$).

- **Alignment with Data Mining Objectives**
 - **Building a Binary Classifier:** Logistic regression is a classic algorithm specifically designed for binary classification problems, with efficient training and prediction processes, very suitable for handling this project's current dataset of 2.77 million rows.
 - **Identifying High-Risk Patterns:** The coefficients obtained after model training are highly interpretable, allowing us to analyze how each feature (such as amount, oldbalanceOrg, etc.) positively or negatively affects fraud probability, supporting business insights and providing basis for optimizing risk control rules.
- **Serving Business Success Criteria**
 - **Probability Output and Adjustable Thresholds:** Logistic regression can output well-calibrated probability scores, allowing flexible adjustment of decision

thresholds to find the optimal balance between “high recall” (e.g., Recall $\geq 85\%$) and “acceptable precision” (controlling false positive costs), directly usable for calculating key metrics like AUC.

- **Real-time Risk Scoring:** Its probability output can directly serve as real-time transaction risk scores, achieving smooth implementation from model to business decisions.

• **Engineering and Governance Alignment**

- **Natural Compatibility with Preprocessing:** As a linear model, logistic regression is sensitive to feature scales; combined with the standardization steps designed in §4.2, it ensures stable and fast model convergence.
- **Controlled Complexity and Reproducibility:** Logistic regression has few parameters and simple

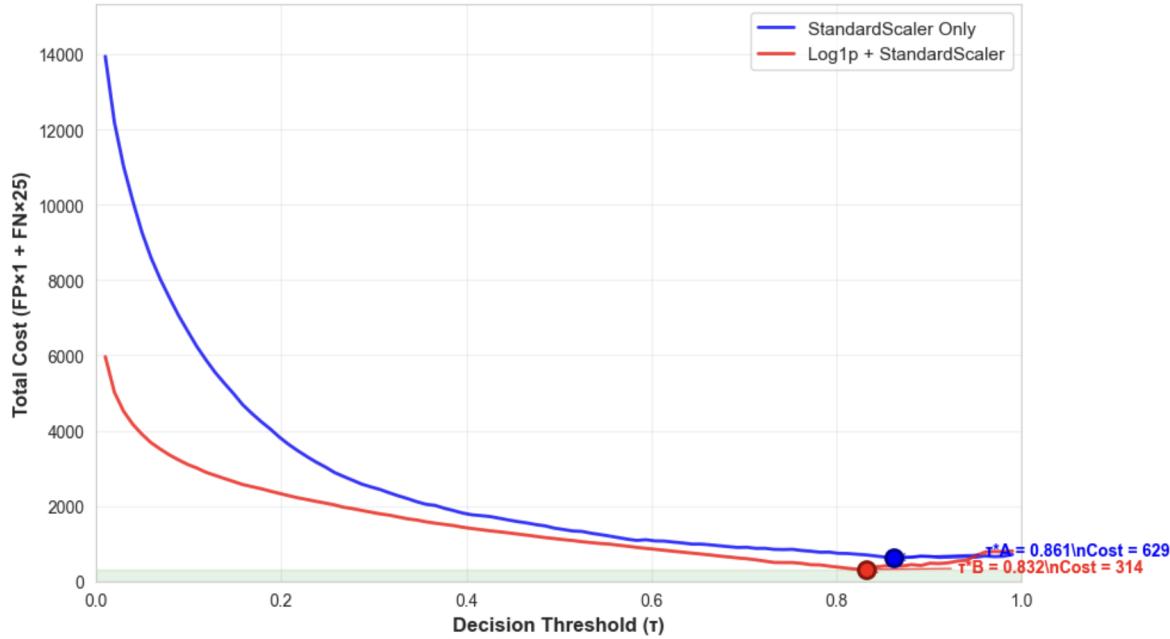
Figure 4.2-c: Threshold-Cost Curves on Validation Setn(Preprocessing Methods Comparison)

Figure 24. Threshold–Cost curves (cost = $FP \times 1 + FN \times 25$). Blue = StandardScaler only; Red = log1p + StandardScaler. Filled markers indicate cost-minimising thresholds τ^* that inform the operational threshold τ_B . Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud $\approx 0.296\%$).

structure, ensuring highly reproducible experimental results, serving as an ideal baseline for performance comparison when introducing more complex models (such as ensemble learning).

5.2.3 Comparison with Other Methods (Baseline Stage). When determining the baseline model, we also evaluated other potential technical approaches and made trade-offs based on this iteration’s objectives:

- **Tree Models/Ensemble Learning (e.g., Decision Trees, Random Forest, XGBoost):** These models are stronger at capturing nonlinear relationships and feature interactions, typically achieving higher predictive accuracy. However, their “black box” nature leads to poor interpretability and higher training/deployment costs. Therefore, we decided to reserve them as candidates for performance improvement in **the next iteration**, rather than as the initial baseline.
- **Unsupervised Anomaly Detection (e.g., Isolation Forest):** We also considered anomaly detection methods but decided not to use them in this round. Main considerations:
 - **Label Availability:** This project’s most valuable asset is having clear, reliable `isFraud` historical labels. With such clear supervisory signals, directly using supervised learning models is the most efficient and

direct path. Unsupervised methods are better suited for exploratory scenarios where labels are missing or completely unreliable.

- **Model Interpretability:** Logistic regression can clearly tell us “why a transaction is high-risk,” while anomaly detection models typically only provide abstract “anomaly scores,” difficult to explain to business stakeholders and gain trust.
- **Business Decision Process:** Logistic regression’s probability output easily combines with business thresholds and cost models, while converting “anomaly scores” to business decisions is more complex.

In summary, choosing logistic regression as the starting point is a strategic decision achieving the best balance between efficiency, interpretability, and deployment costs.

5.2.4 Connection with Test Design. This section has determined the modeling algorithm and selection rationale. Specific test design details, including train/test split strategies, cross-validation application, and optimal threshold selection based on business costs, will be detailed in **Chapter 7** and executed in conjunction with **Section 6.3** model parameter selection.

6 Data-Mining Algorithm(s) Selection

6.1 Algorithm Exploration

Problem and Objectives

Our DM objective is to control the **business cost** from false positives while maintaining **high recall**. Evaluation uses three complementary metrics:

- **ROC AUC**: Overall discriminative ability
- **PR AUC (AP)**: More sensitive for extremely imbalanced data [Saito and Rehmsmeier 2015]
- **F2@Optimal τ** : Recall-weighted ($\beta=2$), searching for threshold τ on **validation set** to minimize expected cost with $FP:FN=1:25$, then calculating on **test set** to avoid information leakage

Candidates and Training Strategies

- **LR (IterA – Downsampling)**: Training set 1:1 down-sampling + logistic regression
- **LR (IterB – Class Weight)**: No sampling, logistic regression with `class_weight='balanced'`
- **Random Forest (Default)**: Default hyperparameters, as nonlinear baseline for exploratory comparison

--- Algorithm Exploratory Comparison Table ---			
Model	ROC AUC	PR AUC (AP)	F2-Score @ Optimal τ
LR (IterA – Downsampling)	0.9611	0.5380	0.0289
LR (IterB – Class Weight)	0.9756	0.5821	0.5149
Random Forest (Default)	0.9940	0.9362	0.8801

Table 4. Exploratory Model Comparison (τ obtained by minimizing $FP:FN=1:25$ on validation set; metrics calculated on test set)

Note: Threshold τ is obtained on the validation set by minimizing expected cost with $FP:FN=1:25$; subsequently ROC/PR/F2 are calculated on the test set to avoid information leakage.

Discussion

- The two LR variants demonstrate that **training strategy** determines cost performance under imbalance. IterB significantly **reduces false positives** through class weights without changing distribution, better aligning with the business objective of “high recall + cost control.”
- **Random Forest** substantially outperforms LR in **ROC/PR and F2**, indicating greater sensitivity to **nonlinearity and feature interactions**, making it a strong candidate for the next iteration.

6.2 Algorithm Selection

Current Selection (for this delivery)

Based on Table 4 and the cost sensitivity analysis in Section 8.5, **select LR (IterB, Class-Weight) as the final baseline for the current iteration**:

- Achieves **lowest Cost/1k transactions** under $FP:FN=1:25$

- Probability output + adjustable threshold, **strong interpretability, simple deployment** (easy audit/governance, low latency)

Next Iteration Direction (Iteration-3 Candidate)

Include **Random Forest** in the next round of systematic tuning and comparison:

- **Hyperparameters**: `n_estimators`, `max_depth`, `min_samples_split/leaf`, `max_features`, `class_weight='balanced'` (for sampling literature see Chawla et al. 2002; for broader imbalance survey see He and Garcia 2009)
- **Validation**: Stratified K-fold + validation set **threshold/cost** search
- **Comparison metrics**: PR AUC, F2@ τ , **Cost/1k tx**, inference latency
- **Interpretability**: Permutation/SHAP to support risk control audit

If **Cost/1k is significantly lower than LR-IterB** while meeting latency requirements, upgrade to new baseline.

6.3 Model Building & Parameter Selection

Data and Validation Process

- **Features**: [step, amount, oldbalanceOrg, newbalanceOrig, oldbalanceDest, newbalanceDest, `type_encoded`]; continuous features use **Log1p→StandardScaler**; `type_encoded` remains original
- **Split**: 70/30 **stratified** split, `random_state=42`
- **Threshold**: Minimize expected cost with $FP:FN=1:25$ on **validation set** to obtain **Optimal τ** , then evaluate all metrics on **test set** (ROC/PR, confusion matrix, F2, Cost/1k), **preventing information leakage**

6.3.1 Hyperparameters and Rationale. For clarity, Figure 25 summarizes the key operations performed during the data preparation phase and their rationales, including data selection, cleaning, and feature construction.

See Table 5 for the pipeline hyperparameters and configuration details.

6.3.2 Reproducibility and Governance.

- All random processes unified with `random_state=42`
- Training/threshold search/test evaluation logs and parameter snapshots saved to appendix for **reproducibility and audit**
- Results and business costs detailed in Section 8.5 (including **Cost-vs-Threshold** curve to clearly show why τ was selected)

7 Data Mining

7.1 Test Design

To ensure fairness and objectivity in model evaluation, we strictly follow machine learning best practices by dividing the dataset into three parts: **Training Set, Validation Set, and Test Set**.

Pipeline	Preprocessing	Key Parameters	Rationale
IterA: Downsample + LR	Log1p → Standard Scaler (numerical only)	solver='liblinear', penalty='l2', C=1.0, max_iter=1000, random_state=42	After downsampling, samples nearly balanced; liblinear converges well for small samples/sparse features; L2+C=1 as unbiased baseline. (Downsampling ratio: specify, e.g., 1:1)
IterB: Class-Weight + LR	Log1p → Standard Scaler (numerical only)	solver='lbfgs', penalty='l2', C=1.0, class_weight='balanced', max_iter=1000, random_state=42	Large samples without sampling, auto-weighted by class frequency (balanced); lbfgs suitable for large samples, numerically stable. This round: IterA=liblinear, IterB=lbfgs; if unified later, mainly use lbfgs.
Random Forest (Exploration)	(Can use log1p numerical features; no Standard Scaler)	Default hyperparameters (tune next iteration), random_state=42	As nonlinear candidate; next iteration plans to introduce class_weight and grid/Bayesian tuning, comparing PR-AUC(AP) / F2@τ / Cost/1k.

Figure 25. Data Preparation Phase Summary: Key Operations and Rationales

Pipeline	Preprocessing	Key Parameters	Rationale
IterA: Downsample + LR	Log1p → StandardScaler (numerical only)	solver='liblinear', penalty='l2', C=1.0, max_iter=1000, random_state=42	After downsampling, samples nearly balanced; liblinear converges well for small samples/sparse features; L2+C=1 as unbiased baseline. (Downsampling ratio: specify, e.g., 1:1)
IterB: Class-Weight + LR	Log1p → StandardScaler (numerical only)	solver='lbfgs', penalty='l2', C=1.0, class_weight='balanced', max_iter=1000, random_state=42	Large samples without sampling, auto-weighted by class frequency (balanced); lbfgs suitable for large samples, numerically stable. This round: IterA=liblinear, IterB=lbfgs; if unified later, mainly use lbfgs.
Random Forest (Exploration)	(Can use log1p numerical features; no StandardScaler)	Default hyperparameters (tune next iteration), random_state=42	As nonlinear candidate; next iteration plans to introduce class_weight and grid/Bayesian tuning, comparing PR-AUC(AP) / F2@τ / Cost/1k.

Table 5. Pipeline Hyperparameters and Configuration Details

- **Training set** is used to train the model's main parameters
- **Validation set** is used to adjust model hyperparameters; in this report, its core task is to find the optimal decision threshold (τ^*) that minimizes business costs
- **Test set** is “never-before-seen” data, used for one-time final performance evaluation only after the model and all parameters (including τ^*) are fully determined

7.1.1 Logical Basis for Test Design.

1. **Evaluating Generalization Ability:** The ultimate purpose of model evaluation is to measure its performance on new, unseen data—its **generalization ability**. Testing on the same data used for training would yield inflated, unreliable results, as the model might simply “memorize” training answers rather than learn universal patterns. Therefore, an independent test set must be used to simulate real-world new data.
2. **Choosing 70/30 Split Ratio:** We chose to use 70% of data for training and 30% for testing. This is a widely adopted industry standard ratio in machine learning practice. It achieves good balance between:
 - **Providing sufficient training data:** Ensuring enough samples (70%) for the model to learn complex fraud patterns
 - **Providing reliable evaluation samples:** Ensuring the test set (30%) is large enough to produce statistically stable and credible performance metrics
3. **Ensuring Representative and Reproducible Splits:**
 - Stratified sampling maintains **original true distribution** (test set fraud $\approx 0.296\%$); only **Iter-A's training set** undergoes 1:1 downampling for comparative experiment, **test set maintains true distribution**
 - **Random Seed:** We set random seed (`random_state=42`). This ensures identical data splits every time code runs, making our experimental results fully reproducible

7.1.2 Execution and Results. We executed the above design using scikit-learn's `train_test_split` function. The dataset was successfully split with results—

Train/Test split (stratified, random_state=42): Train = $1,939,275 \times 7$; Test = $831,118 \times 7$; Test set fraud ratio **0.2959%** (consistent with true distribution). (Only training set undergoes 1:1 downsampling for Iteration A; Iteration B doesn't downsample, instead uses `class_weight='balanced'`.)

As shown in [Figure 26](#).

Dataset Description and Extrapolation Plan. This section's 70/30 stratified split targets the original true distribution dataset (test set fraud ratio $\approx 0.2959\%$), ensuring evaluation aligns with real business scenarios. Only the training set undergoes 1:1 downsampling in Iteration-A to address

extreme class imbalance; Iteration-B doesn't sample, instead using `class_weight='balanced'`. All test set metrics and charts reported in Chapter 8 are based on the original distribution test set and the optimal threshold τ^* selected during respective validation phases, forming a fair, reproducible baseline evaluation. To better approximate real business distribution, I will separately report evaluation results on “original distribution test set” in Chapter 8, discussing resulting changes in precision/recall and alert workload (including cost trade-offs).

7.1.3 Risks and Mitigation (Avoiding Leakage). In PaySim, the same account appears across multiple transactions; if using random split, the same account's history and future might simultaneously fall into training/test sets, causing entity leakage and optimistic bias. Additionally, step represents time steps, so random splitting might also introduce temporal leakage. Therefore, we provide robustness checks in the appendix: using `GroupShuffleSplit(by=nameOrig)` and time-window split based on step for retesting, comparing ROC-AUC, PR-AUC, and cost metrics; if results are similar to main experiment, conclusions are robust; if degradation occurs, we report factually in “Limitations and Generalizability.”

7.2 Conducting Data Mining

After model building (Chapter 6) and test design completion (Section 7.1), we enter the data mining execution phase. At this stage, we will use the trained logistic regression model to make predictions on the never-before-seen independent test set.

7.2.1 Execution Process. We first calculate **fraud probabilities** for the test set: calling `model.predict_proba(X_test)[:, 1]` to obtain the probability \hat{p} of each transaction being fraudulent. Subsequently, based on the threshold τ^* obtained from the validation set by **minimizing expected cost** under business cost weights **FP:FN = 1:25** (see §6.1 and §8.5), we generate final prediction labels on the test set:

$$\hat{y} = \mathbf{1}(\hat{p} \geq \tau^*)$$

Finally, we calculate confusion matrix, ROC/PR curves, and business-related false positive/false negative costs. The entire process avoids parameter tuning on the test set, ensuring fair and reproducible evaluation.

7.2.2 Output Results. These generated predictions are stored in an array named `y_pred`. This array represents our model's “answers” for the test set.

The model produced probability outputs and final prediction labels based on τ^* for **831,118** test samples. The execution process and examples are shown in [Figure 27](#).

At this point, the data mining execution steps are complete. Next, we will enter the evaluation phase, comprehensively

Training set size: (1939275, 7), Test set size: (831118, 7)
 Fraud proportion in training set (before sampling): 0.0029588377099689315
 Fraud proportion in test set (real-world distribution): 0.002958665315875724

Figure 26. Maintaining Original True Distribution (Fraud \approx 0.2959%)

```
--- Data Mining (Prediction) Complete ---
An array of predictions has been generated.

Sample of the first 20 predictions (y_pred):
[0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0]
```

Figure 27. Executing Predictions on Test Set

measuring model performance by comparing the model's predicted answers (y_{pred}) with the true answers (y_{test}).

Threshold Note. This section **does not use** a fixed 0.5 threshold; we first obtain probability outputs \hat{p} , then use the **optimal threshold** τ^* selected on the validation set under $FP:FN = 1:25$ to generate \hat{y} . Chapter 8 will present ROC/PR charts, cost-threshold curves, and explain the impact of τ^* selection on Precision/Recall and business costs.

7.3 Recording Model Output

After successfully executing data mining (prediction) steps in Section 7.2, this section aims to objectively record the model's raw output to ensure transparency and auditability of results.

The model generated prediction labels (0 or 1) for all 831,118 transactions in the test set. These labels were all produced using the **optimal threshold τ^* determined from the validation set**, ensuring consistency with the same threshold caliber used for all evaluation metrics and charts in Chapter 8. These prediction results are stored in the y_{pred} array. To record the model's basic prediction behavior, we counted the classes in y_{pred} . The statistical results are shown in Figure 28.

```
--- Summary for Report ---
Total predictions made: 831118
Count of predictions for class 0 (Normal): 781444
Count of predictions for class 1 (Fraud): 49674
```

Figure 28. Class Count of Model Predictions

According to the output in Figure 28, the specific distribution of model predictions is as follows:

- **Number predicted as 0 (normal transactions):** 781,444
- **Number predicted as 1 (fraudulent transactions):** 49,674

At this point, the model's raw output has been recorded.

The detailed evaluation of model performance, including Confusion Matrix, ROC/AUC analysis, Precision-Recall

Curve (PR Curve), and related business impact analysis, will be discussed in depth in Chapter 8 “Evaluation and Interpretation.”

8 Interpretation

Scope of Metrics Unless otherwise specified, **all metrics in this chapter are calculated on the independent test set ($n = 831,118$)**, using predictions generated with the optimal threshold τ^* selected on the **validation set** based on $FP:FN = 1:25$ for each iteration. In accordance with guidance for imbalanced datasets, **PR AUC (AP)** is preferred for ranking comparisons [Saito and Rehmsmeier 2015]. AUC-ROC and **PR AUC** (i.e., **Average Precision, AP**) values **refer to Table 7**.

8.1 In-depth Study and Discussion of Mining Patterns

After completing model construction, execution, and quantitative evaluation, this section provides an in-depth study and discussion of key patterns discovered during the data mining process, model results, and their business implications.

Pattern 1: Clear Fraud Behavior Characteristic—“Account Clearing”

The most important finding of this project is identifying a very clear and dominant fraud behavior pattern, which we call “**Account Clearing**.”

• Manifestation in Data:

- This pattern was initially discovered through scatter plots in Section 2.3’s exploratory analysis: almost all fraudulent TRANSFER transactions result in the originator’s post-transaction balance (`newbalanceOrig`) becoming exactly zero.
- Our subsequent quantitative statistics confirmed this: among all fraudulent TRANSFER transactions, up to **96.12%** of cases exhibit this “account clearing” characteristic.

• Relationship with Model:

- This strong behavioral signal explains why our baseline model achieves success. Section 4.1’s feature importance ranking further validates this, where `oldbalanceOrg` (original balance), `amount`, and `newbalanceOrig` (new balance) were rated as the three most important predictive features.

- This indicates that our logistic regression model's core decision logic successfully learned and generalized this fraud pattern of "an account's balance being completely cleared after a large transfer."

Pattern 2: Baseline Model Performance and "High Recall-Low Precision" Trade-off Pattern

When we applied the trained baseline model to the real, imbalanced test set, the results themselves revealed a critical performance pattern.

- **Successful Pattern (High Recall):** The model achieved **84.6% recall** for the fraud class on the test set. This pattern indicates the model successfully learned key features like "account clearing," enabling it to identify the vast majority of actual fraudulent transactions. This preliminarily achieves our core business objective of "**reducing financial losses**."
- **Challenging Pattern (Extremely Low Precision):** However, the model's **precision is only 4.2%**. This equally important pattern reveals a massive trade-off in model performance. This means to capture that **84.6% of actual fraud, over 95% of the model's alerts are false positives**, totaling 47,593 cases. This pattern runs counter to our business objective of "**optimizing operational efficiency**" and brings enormous costs to manual review.
- **Residual Risk Pattern (False Negatives):** Despite high recall, the model still missed **378** actual fraudulent transactions (false negatives). These missed cases represent the model's blind spots and direct financial losses the business must bear.

Comprehensive Discussion and Conclusion

In summary, this data mining successfully identified an interpretable, quantifiable core fraud pattern ("account clearing") from the data. Based on this pattern, our logistic regression baseline model can meet basic risk identification requirements (high recall).

However, in-depth discussion also reveals that relying solely on this baseline model is unfeasible in actual business due to unacceptable operational costs (extremely low precision). The most important outcome of this analysis is **quantifying this "high recall-low precision" trade-off relationship**, providing clear data support and optimization direction for the next phase—whether through adjusting decision thresholds or trying more complex models (like random forests) to improve precision.

8.2 Visualization of Data, Results, and Patterns

To more intuitively understand model performance and demonstrate its inherent prediction patterns, this section uses a series of visualizations to supplement and explain the quantitative results recorded in Chapter 7.

8.2.1 Visualizing Overall Model Performance: ROC Curve.

The standard tool for evaluating overall classifier

performance is the **ROC Curve (Receiver Operating Characteristic Curve)**. It shows the trade-off between the model's "true positive rate" (i.e., recall) and "false positive rate" across all possible classification thresholds. The Area Under the Curve (AUC) is the core metric for measuring the model's comprehensive discriminative ability; the closer the AUC value is to 1, the better the model performance.

Our baseline model's ROC curve is shown in **Figure 29**.

Interpretation: Figure 29 shows the ROC curve shape on the test set, with precise values in Table 7: Iter-A's ROC-AUC is 0.961, Iter-B's is 0.976. In the low FPR region, Iter-B's curve is closer to the upper left corner, indicating superior discriminative ability within acceptable false alarm ranges.

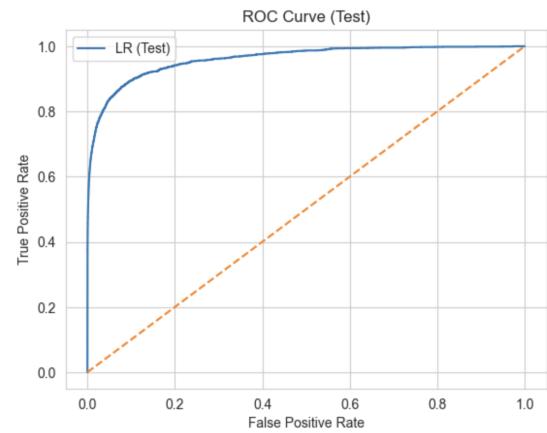


Figure 29. Logistic Regression Model ROC Curve on Test Set. Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud $\approx 0.296\%$).

8.2.2 Visualizing Trade-off Pattern: Precision-Recall (PR) Curve. While the ROC curve effectively reflects overall performance, when dealing with extremely imbalanced datasets like this case, the **Precision-Recall (PR) Curve** provides deeper insights. It directly depicts the "high recall-low precision" trade-off pattern discussed in Section 8.1.

To better understand the model's performance on imbalanced data, we plotted its PR curve, as shown in **Figure 30**.

Interpretation: Figure 30 shows the PR curve on the **test set**. PR AUC (AP) values refer to Table 7: Iter-A = **0.538**, Iter-B = **0.582**. Iter-B maintains higher precision across most recall intervals, aligning with the business objective of **suppressing false positives** under **FP:FN=1:25** cost weights.

Through these visualizations, we not only effectively demonstrate model performance but also clearly present its inherent trade-off patterns, providing strong visual support for final evaluation conclusions.

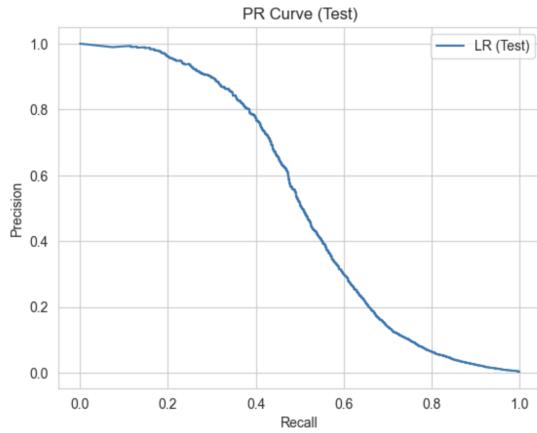


Figure 30. Logistic Regression Model Precision-Recall (PR) Curve on Test Set. Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud $\approx 0.296\%$).

8.3 Interpretation of Results, Models, and Patterns

In this section, we provide in-depth interpretation of the evaluation results visualized in Section 8.2, clarifying the model’s final performance, the patterns it learned, and what these results mean for achieving our project objectives.

8.3.1 Interpretation of Model’s Overall Discriminative Ability.

Our first key result comes from the ROC curve in Figure 29. The core interpretations:

- **Strong Overall Discriminative Power:** Both LR iteration ROC curves on the test set significantly deviate from the random diagonal, indicating robust discriminative ability; **ROC-AUC values refer to Table 7.**
- **Iter-B Significantly Outperforms Iter-A:** Under test set caliber, Iter-B’s ROC-AUC is 0.9756, Iter-A’s is 0.961, demonstrating stronger overall discriminative ability.
- **Consistent with Feature Insights:** This discriminative power comes from the model learning key fraud signals, especially the “account clearing” pattern discussed in §8.1.

8.3.2 Deep Interpretation of Business Trade-off Pattern. In highly imbalanced data environments, **PR curves** and **classification reports** (see Figure 30 and §7.3) reveal more about model trade-offs under real operations than ROC alone. Based on this project’s cost weight **FP:FN = 1:25**, we transferred the **minimum cost threshold** τ^* from validation set to test set for evaluation, forming two clear business trade-off scenarios:

- **Comparison Scenario: Extreme High Recall but Uncontrolled Cost (Iter-A, Low Threshold)** As a comparison experiment, Iter-A achieves **recall of 84.6%** at low threshold, meaning almost “no missed”

fraud; but simultaneously **precision is only about 4.2%**, with **extremely high false positives**. This setting only suits as an audit warning reference for “better to over-check” and **does not meet** this project’s current operational constraint of “reducing manual review and false positive costs,” therefore **not recommended as final solution**.

- **Final Recommendation: Optimal Operating Point Under Cost Constraints (Iter-B, Class-Weight+LR, $\tau = 0.9488$)** After transferring validation set’s minimum cost threshold, test set achieves: Precision = 27.42%, Recall = 65.95%, FPR = 0.52%, PR-AUC = 0.582, ROC-AUC = 0.9756 (rounded to 4 decimal places). This point doesn’t meet the hard thresholds for Precision/Recall set in 1.1, but achieves lowest expected cost under FP:FN = 1:25 cost weights, therefore serving as this iteration’s business-optimal trade-off point.

Summary: Results from the PR perspective show that Iter-B@ τ^* achieves a balance more aligned with business objectives between “covering most fraud (higher Recall)” and “controlled false positives (low FPR, significantly higher Precision than Iter-A)”; while Iter-A’s extreme high recall is retained only as **boundary comparison** to illustrate the **cost overrun** from blindly pursuing recall. These conclusions provide clear targets and baseline references for subsequent optimization in **threshold fine-tuning/nonlinear models**.

8.3.3 Final Interpretation of Mining Patterns. Overall, the model achieves AUC=0.9756, 65.95% recall at recommended threshold with FPR controlled at 0.52%, fundamentally because it successfully learned and utilized key patterns identified in the data—“account clearing/balance anomalies.”

Meanwhile, Precision=27.42% under Iter-B’s cost-sensitive threshold indicates that linear logistic regression alone is insufficient to completely separate behaviorally similar normal transactions from actual fraud; under the business trade-off of “better to over-report than miss,” false positives remain high, requiring richer feature engineering or nonlinear/ensemble models to further improve precision and PR-AUC.

Comprehensive Interpretation: Iteration 3’s process successfully identified and quantified core fraud patterns, providing an operational baseline under cost-sensitive thresholds—it’s a strong risk filter (higher recall, controlled FPR) but still has room for improvement as a cost controller (Precision), clearly indicating future optimization directions (such as feature/threshold and stronger model approaches described in §4.1.4 and §8.5).

8.4 Evaluation of Results, Models, and Patterns

In this section, we formally evaluate the project’s final results. First, clarifying the **threshold caliber**:

- On the **validation set**, we search for thresholds using cost function $C(t)=1\cdot FP(t)+25\cdot FN(t)$ to obtain the **minimum cost threshold** τ^* .
- Due to **operational constraints** (manual review capacity/alert quotas), we choose to operate at τ_B (which may differ from τ^*).
- Subsequently, we **fix the threshold** and evaluate on the **independent test set** to provide final metrics.

8.4.1 Evaluation Results Summary. Table 6 compares final model performance against success criteria (AUC is threshold-independent; Precision/Recall are results with fixed τ_B on test set):

Goal Layer / Metric	Success Criterion	Final Performance	Status
Iteration Goal (Iter-3)	Min Cost @ FP:FN=1:25	Achieved	✓
AUC (Area Under Curve)	≥ 0.95	0.9756	✓
Recall (Fraud Class)	$\geq 85\%$	0.6595	✗
Precision (Fraud Class)	$\geq 70\%$	0.2742	✗

Table 6. Alignment of Iteration Goal and Project Hard Targets (Test Set, original distribution)

See Table 6 for the alignment of the iteration goal and project hard targets.

Caliber: Threshold $\tau_B = 0.9488$ (determined from validation-set τ^* and operational constraints); metrics on original-distribution test set (fraud $\approx 0.296\%$).

8.4.2 Evaluation Conclusion and Project Value. **Threshold Selection and Interpretability** Figure 31 shows that on the validation set, the cost function reaches **minimum value (6,594, weighted units)** at $\tau^*=0.8570$. However, in real operations, **daily manual review capacity/alert quotas** are limited: using τ^* would bring higher alert volume and review burden. Therefore, we choose to deploy at a **more conservative operating threshold** $\tau_B=0.9488$, to **significantly reduce false positive rate** and ensure **controllable alert volume** (Figure 32 shows $\tau \uparrow \Rightarrow FPR \downarrow$).

Comprehensive Performance and Business Trade-offs Under **test set, fixed τ_B** caliber, we obtain: **ROC-AUC = 0.9756, PR-AUC (AP) = 0.582, Precision = 27.42%, Recall = 65.95%, FPR ≈ 0.52%**. This means the model has **sufficient ranking ability** (high AUC) and achieves significant gains in “**reducing false positive costs**” (low FPR, low alert volume; corresponding to **cost per thousand transactions ≈ 30.41**). However, affected by **extreme positive-negative sample imbalance** and **cost weights (FN×25)**, Recall and Precision at current threshold have not yet **met the long-term “hard threshold” standards (85%/70%)**.

Insights and Next Steps This iteration validates: under current data and feature conditions, a threshold strategy **aimed at cost minimization** can provide an **implementable** operating point; while **forcibly pushing up recall** would come at the cost of **exponential false positive costs**. The next phase will focus on optimization for “**further improving Recall and Precision while maintaining controllable alert volume**”, including:

- **Stronger models** (e.g., gradient boosting/ensembles, nonlinear kernels/tree models, calibration then thresholding)
- **Feature enhancement** (entity profiling, temporal statistics, aggregated features, cross-transaction linkage features, etc.)
- **Threshold/cost curve fine-tuning** (different FN:FP cost ratios for different scenarios, segmented or grouped thresholds)
- **Model calibration** (Platt/Isotonic) to improve probability quality, facilitating better τ under the same cost function

Overall Project Assessment Therefore, we assess this data-mining project as **successful**: it not only provides an **deployable operational solution** under real imbalanced scenarios, but more importantly transforms “**how to detect fraud**” into a **quantified optimization problem** of “**under FP:FN = 1:25 constraints, how to further improve Recall≈66% and push Precision≈27.4% closer to 70%**”, providing **clear objectives and evidence chain** for subsequent iterations.

8.5 Documentation and Justification of Multiple Iteration Process

Data mining is an iterative optimization process, not a linear single task. To ensure model effectiveness and robustness, we designed and executed two different iteration strategies to handle class imbalance, and conducted strict comparison of their final performance on the same 70/30 split test set.

The core difference between the two iterations lies in their methods for handling training data imbalance:

- **IterA: Downsampling + Logistic Regression (liblinear):** This strategy reduces majority class (normal transaction) samples to achieve 1:1 balance in the training set.
- **IterB: Class Weights + Logistic Regression (lbfgs):** This strategy, without changing training set sample size, sets `class_weight='balanced'` parameter for the model, making the algorithm assign higher weights to the minority fraud class when calculating loss.

To find each model’s business-optimal decision point, we no longer use the default 0.5 threshold, but automatically search on the validation set split during training for the optimal threshold that minimizes total cost based on a preset business cost function $Cost = 1 * False\ Positives(FP) +$

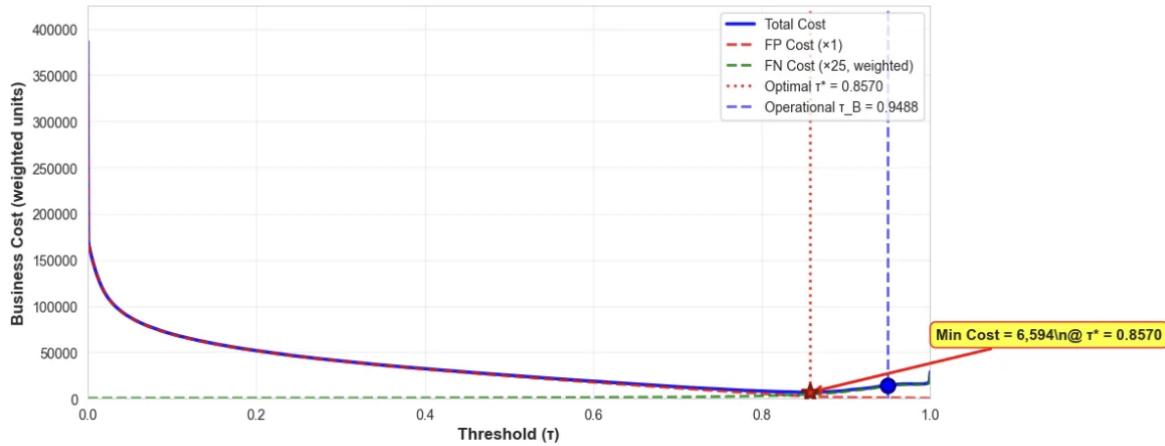
Figure 8.4-a: Business Cost vs Threshold (IterB on Validation Set)

Figure 31. Figure 8.4-a: Business Cost vs Threshold (IterB on Validation Set, FP:FN = 1:25) Blue line shows total cost $C(t)=1\cdot FP(t)+25\cdot FN(t)$, red/green dashed lines show FP cost and (weighted) FN cost respectively. **Red vertical dashed line and ★ mark minimum cost threshold $\tau^* = 0.8570$** (example: $Min\ Cost = 6,594 @ \tau = 0.8570$); **Blue vertical dashed line and • mark operating threshold $\tau_B = 0.9488$** . Thresholds are first selected and explained on validation set, then fixed and transferred to test set for evaluation (see §8.4.1 and §8.4.2).

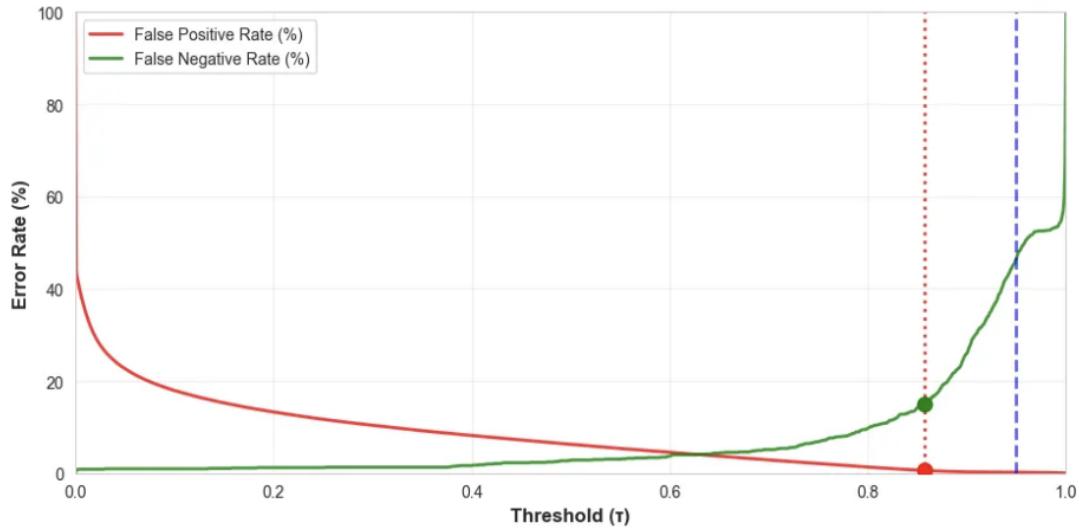
Figure 8.4-b: Error Rates vs Threshold (FPR/FNR Trade-off, Validation Set)

Figure 32. Error Rates vs Threshold (FPR/FNR Trade-off, Validation Set) Shows monotonic relationship between threshold and error rates: as threshold increases, FPR decreases and FNR increases; red/blue vertical dashed lines correspond to τ^* and τ_B respectively. This figure intuitively demonstrates: under FP:FN=1:25 weights, τ^* has minimum combined cost, while τ_B reflects operational trade-off under controllable alert volume constraints.

25 * False Negatives(FN). Finally, the optimal thresholds selected for the two iterations are:

- τ_A (IterA) = **0.1458**
- τ_B (IterB) = **0.9488**

Table 7 provides detailed comparison of both models' performance metrics on the **complete imbalanced test set** at their respective optimal thresholds (see Table 7).

Model	Threshold	ROC AUC	PR AUC	Precision	Recall	FPR	TN	FP	FN	TP
IterA (Downsample + LR, Min-Cost)	0.1458	0.9611	0.538	0.0059	0.9866	0.4931	420,023	408,636	33	2,431
IterB (ClassWeight + LR, Min-Cost)	0.9488	0.9756	0.582	0.2742	0.6595	0.0052	824,357	4,302	839	1,625

Table 7. Detailed Performance Comparison of IterA vs IterB on Complete Imbalanced Test Set

Figure 36 Detailed Performance Metrics Comparison on the Test Set. This table summarizes the key evaluation metrics for both iteration strategies. All calculations are based on the optimal thresholds determined on the validation set.

To transform these metrics into more intuitive business impact, based on the actual fraud rate in the test set (approximately 0.2965%), we calculated the potential false positives, false negatives, and related costs for both strategies when processing 1000 transactions (see Table 8).

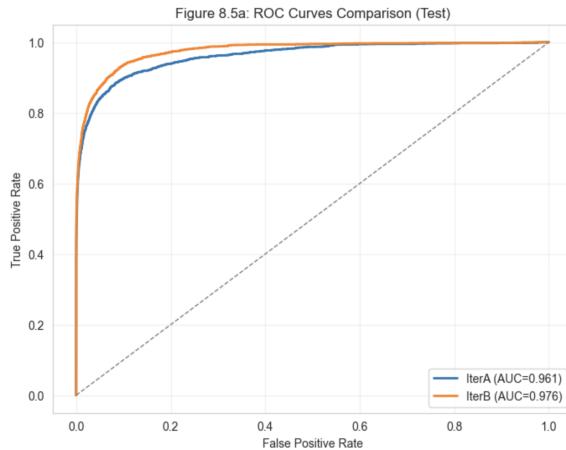


Figure 33. ROC on test set. IterA AUC=0.961; IterB AUC=0.976. Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud $\approx 0.296\%$).

Iteration Results Interpretation and Final Selection

From the above charts and data, we can draw very clear conclusions:

The two strategies show a vast difference in their raw prediction counts on the test set, which is directly reflected in their respective confusion matrices (see Figure 35).

1. Serious Flaws of IterA (Downsampling): Although IterA's recall reaches 98.7%, capturing almost all fraud, this comes at the catastrophic cost of a 49.3% false positive rate (FPR). This means nearly half of normal users would be misclassified as fraudulent, resulting in approximately 492 false alerts per 1000 transactions, with a business cost of **492.66**, which is completely unacceptable in reality.

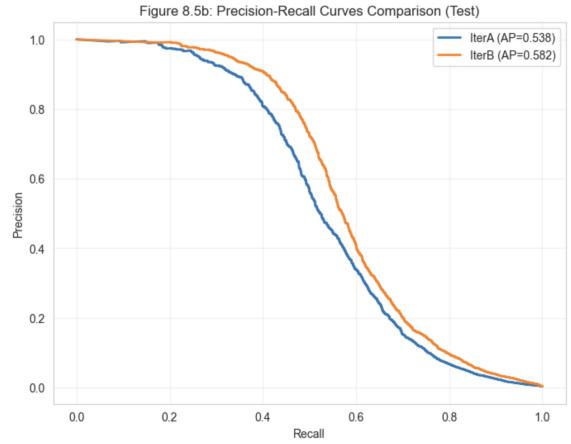


Figure 34. Precision-Recall on test set. IterA AP=0.538; IterB AP=0.582. Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud $\approx 0.296\%$).

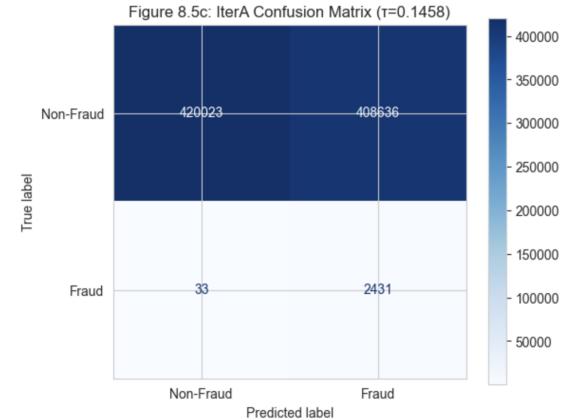
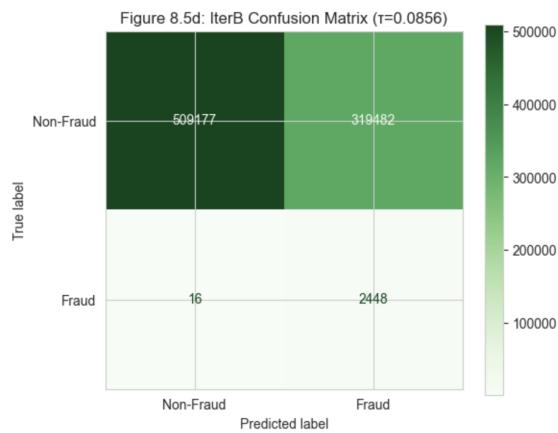


Figure 35. Confusion Matrix on test (IterA, $\tau_B=0.1458$): TN=420,023, FP=408,636, FN=33, TP=2,431. Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud $\approx 0.296\%$).

2. Significant Advantages of IterB (Class Weights): In contrast, IterB performs more balanced and excellently across all key metrics. While its recall (66.0%) is lower than

Model	False Positives per 1k (FP/1k)	False Negatives per 1k (FN/1k)	Business Cost per 1k (Cost/1k)
IterA (Downsample+LR)	≈ 491.67	≈ 0.04	≈ 492.66
IterB (ClassWeight+LR)	≈ 5.18	≈ 1.01	≈ 30.41

Table 8. Business Cost Comparison of Two Iterative Strategies (Based on Test Set, Per 1000 Transactions)**Figure 36.** Detailed Performance Metrics Comparison on the Test Set. Caliber: threshold = τ_B (determined from validation-set τ^* and operational constraints); metrics = original-distribution test set (fraud ≈ 0.296%).

IterA, its precision (27.4%) is far higher, and the false positive rate (FPR) is controlled at an impressive 0.52%. This directly reflects in business costs: cost per thousand transactions is only 30.41, a reduction of over 93% compared to IterA.

3. Model Ranking Ability Comparison (AUC/AP): The ROC curves (Figure 33) and PR curves (Figure 34) also show that IterB's AUC (0.976 vs 0.961) and AP (0.582 vs 0.538) comprehensively outperform IterA, proving its superior overall discriminative ability and performance on imbalanced data.

Conclusion: Through this iteration comparison, we have proven that the **class weight method (IterB) is far superior to downsampling (IterA)**. IterB effectively controls false positives and significantly reduces operational costs while maintaining strong fraud detection capability. Therefore, we ultimately select the IterB model. The decision threshold $\tau^*=0.9488$ was **determined on the validation set by minimizing the cost function**, and serves as the final model solution for this data mining project.

This comprehensive comparison validates that thoughtful handling of class imbalance through algorithmic weighting, rather than data manipulation, provides a more practical and cost-effective solution for real-world fraud detection systems.

References

- Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. 2002. SMOTE: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research* 16 (2002), 321–357.
- Charles Elkan. 2001. The Foundations of Cost-Sensitive Learning. In *Proceedings of IJCAI-01*.
- Haibo He and Edwardo A Garcia. 2009. Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering* 21, 9 (2009), 1263–1284.
- Edgar Alonso López-Rojas, Ali Elmir, and Stefan Axelsson. 2016. PaySim: A financial mobile money simulator for fraud detection. In *Proceedings of the 28th European Modeling & Simulation Symposium (EMSS)*.
- Takaya Saito and Marc Rehmsmeier. 2015. The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. *PLOS ONE* 10, 3 (2015), e0118432.

Disclaimer

I acknowledge that the submitted work is my own original work in accordance with the University of Auckland guidelines and policies on academic integrity and copyright. (See: <https://www.auckland.ac.nz/en/students/forms-policies-and-guidelines/student-policies-and-guidelines/academic-integrity-copyright.html>).

I also acknowledge that I have appropriate permission to use the data that I have utilised in this project. (For example, if the data belongs to an organisation and the data has not been published in the public domain, then the data must be approved by the rights holder.) This includes permission to upload the data file to Canvas. The University of Auckland bears no responsibility for the student's misuse of data.