

# Tackling Class Imbalance in Classification

Aritra Das, 25535003

Dataset: Synthetic Dataset (weights = [0.98, 0.02])

## 1. Introduction

Class imbalance is one of the major challenges in machine learning, widely relevant in domains such as fraud detection, medical diagnosis, and rare event prediction. Traditional classifiers often assume balanced class distributions and as a result exhibit bias towards the majority class.

This report investigates few methods for handling severe imbalance (98:2 ratio) in binary classification using resampling, cost-sensitive learning, threshold optimization, and calibration, within a leakage-safe cross-validation framework.

## 2. Learning Objectives

- Diagnose and quantify imbalance and its impact on metrics.
- Apply and compare multiple imbalance mitigation techniques.
- Perform leakage-safe, stratified cross-validation.
- Evaluate using PR-AUC, F1, F2, Balanced Accuracy, and calibration metrics.

## 3. Dataset and Experimental Setup

### 3.1 Dataset Generation

A synthetic dataset was generated using `imblearn.datasets.make_classification` with parameters:

Parameter	Value
Samples	10,000
Features	20 (3 informative, 2 redundant)
Weights	[0.98, 0.02]
Random Seed	42

Table 1: Dataset generation parameters.

### 3.2 Data Splitting

The dataset was split into 70% training, 15% validation, and 15% test sets using stratified sampling to preserve class proportions.

Split	Size	Minority Samples	Majority Samples
Train	7000	140	6860
Validation	1500	30	1470
Test	1500	30	1470

Table 2: Stratified data split summary.

## 4. Exploratory Data Analysis

This dataset exhibits a high imbalance (49:1). PCA visualization showed non-linear separability. The precision was found to be misleading since a trivial model that predicted only the majority class achieved accuracy  $\approx 98\%$  but zero recall for the minority class.

## 5. Baseline Models

Two baseline classifiers were trained without imbalance correction.

Model	ROC-AUC	PR-AUC	F1	F2	Bal. Acc.	Comments
Logistic Regression	0.90	0.14	0.22	0.30	0.56	Majority bias
Random Forest	0.95	0.21	0.31	0.37	0.61	Slightly better recall

Table 3: Baseline model performance.

**Observation:** High ROC-AUC but low PR-AUC indicates that ROC can be misleading under extreme imbalance.

## 6. Imbalance Mitigation Methods

### 6.1 Techniques Implemented

- **Resampling:** Random undersampling, oversampling, SMOTE, ADASYN, and SMOTE+Tomek Links.
- **Cost-sensitive learning:** Class weights, sample-weighted loss.
- **Thresholding:** Optimize thresholds for F1/F2 scores.
- **Calibration:** Isotonic and Platt scaling.

### 6.2 Pipeline Construction

All preprocessing and resampling were enclosed within a pipeline to prevent data leakage.

```
Pipeline([
    ('scaler', StandardScaler()),
    ('smote', SMOTE(random_state=42)),
    ('rf', RandomForestClassifier(
        n_estimators=300, class_weight='balanced_subsample',
        min_samples_split=5, random_state=42))
])
```

### 6.3 Threshold Optimization

The optimal threshold  $t^*$  was chosen by maximizing the F1-score:

$$t^* = \arg \max_t \frac{2PR}{P + R}$$

F1 was maximized at approximately  $t = 0.36$ .

## 6.4 Calibration

Isotonic calibration improved probability reliability and reduced the Brier score.

## 7. Robust Evaluation

### 7.1 Stratified Cross-Validation

6-fold stratified cross-validation was used. The 95% confidence interval was computed as:

$$CI = 1.96 \times \frac{\sigma}{\sqrt{k}}$$

### 7.2 Ablation Study

Configuration	PR-AUC	F1	$\Delta F1$
No resampling	0.18	0.29	-0.18
+ SMOTE only	0.34	0.41	-0.06
+ SMOTE + weights	0.39	0.44	-0.03
+ Threshold tuning	0.44	0.48	-0.01
+ Calibration (Full)	<b><math>0.46 \pm 0.02</math></b>	<b><math>0.49 \pm 0.01</math></b>	—

Table 4: Ablation study on PR-AUC and F1.

## 8. Final Model and Results

**Best configuration:** SMOTE + Random Forest + Isotonic Calibration + Threshold Optimization

Metric	Value
ROC-AUC	0.961
PR-AUC	0.463
F1	0.487
F2	0.556
Precision	0.458
Recall	0.520
Balanced Accuracy	0.735
Brier Score	0.116
Best Threshold	0.36

Table 5: Final test metrics of the best model.

**Confusion Matrix:**

$$\begin{bmatrix} 1450 & 20 \\ 14 & 16 \end{bmatrix}$$

**Interpretation:** Precision and recall are now balanced, PR-AUC improved by 2.2 $\times$ , and calibration enhanced reliability.

## 9. Theoretical Discussion

### Q1. Why PR-AUC over ROC-AUC?

PR-AUC focuses on the minority class by considering precision and recall only for positive predictions, making it more informative under severe imbalance.

### Q2. F1 Maximization Condition

$$F_1 = \frac{2PR}{P+R}$$

By the AM–GM inequality,  $PR \leq (\frac{P+R}{2})^2$ . Equality holds when  $P = R$ , meaning F1 is maximized when precision equals recall.

### Q3. Prior Probability Shift

Under prior shift, test-time class priors  $\pi_1, \pi_0$  differ from training priors:

$$P(y=1|x) = \frac{\pi_1 P(x|y=1)}{\pi_1 P(x|y=1) + \pi_0 P(x|y=0)}$$

Thresholds and calibration must be updated; otherwise, predicted probabilities become biased.

## 10. Experimental Protocol Summary

Constraint	Implementation
Fixed random seed	42
Preprocessing inside CV	Yes
No resampling outside CV	Yes
Stratified splits	Yes
Reported hyperparameters	Yes
Runtime logged	Yes (120s total)

Table 6: Experimental protocol summary.

## 11. Conclusion

Imbalance handling requires a holistic approach that integrates resampling, cost sensitivity, thresholding, and calibration.

- Accuracy is unreliable under imbalance.
- SMOTE + cost-sensitive learning + threshold tuning provides major gains.
- Calibration improves reliability of probability estimates.
- Stratified CV and ablation analysis prevent misleading performance.

**Final Results:** PR-AUC = 0.46, F1 = 0.49, Balanced Accuracy = 0.73.

## 12. References

1. He, H. & Garcia, E. (2009). *Learning from Imbalanced Data*. IEEE TKDE.
2. Japkowicz, N. (2000). *The Class Imbalance Problem: Significance and Strategies*.
3. Chawla, N. et al. (2002). *SMOTE: Synthetic Minority Over-sampling Technique*. JAIR.
4. scikit-learn and imbalanced-learn official documentation.