

## Lesson 28 Summary: Class Imbalance

Data Science with Python – Key Concepts

Data Science Program

## Handling Class Imbalance

### The Problem

Minority ignored  
High accuracy misleading  
Model predicts majority

### Resampling

Oversample minority  
Undersample majority  
SMOTE: synthetic

### Class Weights

`class_weight=balanced`  
Penalize minority errors  
No data changes

### Better Metrics

Precision-Recall curve (not ROC)  
F1, balanced accuracy

### Finance Context

Default prediction: rare defaults  
Fraud detection: few frauds

`class_weight="balanced"` | `SMOTE()` | `RandomUnderSampler()`

Imbalanced data is common in fraud, default, and churn prediction

## Why imbalance hurts:

- **Majority bias:** Model predicts the common class
- **Misleading accuracy:** 99% majority = 99% accuracy
- **Lost minority:** Rare but important events missed

---

In fraud detection, the fraud class is typically less than 1%

## Add more minority samples:

- **Random:** Duplicate minority samples
- **SMOTE:** Generate synthetic samples
- **Caution:** May cause overfitting

---

Always oversample training data only, not test data

## Remove majority samples:

- **Random:** Randomly remove majority
- **Tomek links:** Remove ambiguous samples
- **Caution:** Loses potentially useful information

---

Works well when majority class has redundancy

## Algorithm-level adjustment:

- **Method:** Penalize errors on minority more heavily
- **sklearn:** `class_weight='balanced'`
- **Advantage:** No data modification needed

---

Class weights are often the simplest effective solution

## Synthetic Minority Oversampling:

- **How:** Create synthetic samples between neighbors
- **Benefit:** Adds diversity, not just duplicates
- **Library:** imbalanced-learn (imblearn)

---

SMOTE is the most popular oversampling technique

### Beyond accuracy:

- **Precision-Recall curve:** Better than ROC for imbalance
- **F1 score:** Balances precision and recall
- **Balanced accuracy:** Average of per-class recalls

---

Use PR-AUC instead of ROC-AUC for highly imbalanced data

# Threshold Adjustment

## Lower the decision threshold:

- **Default 0.5:** Too high for rare class
- **Lower threshold:** Catches more minority (more recall)
- **Tradeoff:** Also increases false positives

---

Optimize threshold based on business costs

## Recommended workflow:

- **Start simple:** Try `class_weight='balanced'` first
- **Then SMOTE:** If class weights don't work
- **Always:** Evaluate on original (imbalanced) test set

---

**Never resample your test set**

### Essential Techniques:

Method	Code
Class weights	<code>LogisticRegression(class_weight='balanced')</code>
Random oversample	<code>RandomOverSampler().fit_resample(X, y)</code>
SMOTE	<code>SMOTE().fit_resample(X, y)</code>
Undersample	<code>RandomUnderSampler().fit_resample(X, y)</code>
Balanced acc.	<code>balanced_accuracy_score(y, pred)</code>

---

Install imbalanced-learn: `pip install imbalanced-learn`