

Handling Class Imbalance in Classification Models

- Understanding Challenges & Mitigation Techniques

Understanding Class Imbalance

- Class imbalance occurs when one class is significantly more frequent than the other(s).
- Common in real-world applications like:
 - Fraud detection (fraudulent transactions $< 1\%$)
 - Medical diagnosis (rare diseases)
 - Spam detection (majority are non-spam emails)
- Why is it a problem?
 - High Accuracy is misleading.
 - Bias towards majority class.
 - Poor generalization.

Why Accuracy is Misleading?

- Example: Fraud Detection
 - Dataset: 95% Non-Fraud, 5% Fraud
 - Model predicts all transactions as Non-Fraud
 - Accuracy = 95%, but fails to detect fraud!
- Better Evaluation Metrics:
 - Precision & Recall
 - F1-Score
 - ROC-AUC

Strategies for Addressing Class Imbalance

- Data-Level Methods (Resampling Techniques):
 - Oversampling (SMOTE)
 - Undersampling
 - Hybrid methods
- Algorithmic Methods:
 - Cost-sensitive learning
 - Class-weighted loss functions

Resampling Methods for Balancing Data

- Oversampling:
 - Increases minority class instances.
 - SMOTE generates synthetic data.
- Undersampling:
 - Reduces majority class instances.
- Hybrid Approach:
 - Combination of both methods.

Cost-Sensitive Learning

- What is Cost-Sensitive Learning?
 - Adjusts learning algorithm to penalize errors on the minority class more.
 - Helps the model focus on rare classes by assigning a higher misclassification cost.
- Why Use Cost-Sensitive Learning?
 - Does not modify the dataset like resampling.
 - Effective when misclassification costs are known.
- Example Applications:
 - Fraud Detection: Missing fraud cases = financial loss.
 - Medical Diagnosis: Misdiagnosing a rare disease can be critical.

Implementing Cost-Sensitive Learning

- **Approach 1: Assign Class Weights**
 - Increase minority class sample weights.
 - Example in Scikit-Learn:

```
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier(class_weight={0: 1, 1: 5})
model.fit(X_train, y_train)
```

- **Approach 2: Custom Cost Matrices**
 - Define different penalties for different misclassifications.
 - Example cost matrix (higher penalty for misclassifying minority class):

Actual \ Predicted	Class 0	Class 1
Class 0	0	+1
Class 1	+5	0

Which Method Should You Use?

- If data is highly imbalanced: Use oversampling (SMOTE) or cost-sensitive learning.
- If dataset is large enough: Try undersampling or a hybrid approach.
- Dataset is small or fixed : Algorithmic solutions like cost-sensitive learning & weighted loss improve results.
- Evaluate models using: Precision-Recall AUC, F1-score, and Confusion Matrix.

Summary & Best Practices

- Class imbalance affects model performance.
- Accuracy is misleading; use Precision, Recall, F1-score, PR-AUC.
- Use resampling techniques: Oversampling, Undersampling, Hybrid methods.
- Algorithmic solutions like cost-sensitive learning & weighted loss improve results.
- Choose the right technique based on dataset size, model type, and class distribution.