

Intro: Finding Four-Leaf Clovers in a Field of Grass

Rare event prediction is like finding four-leaf clovers * # in a massive field of grass.

Most of what we see is just grass - ordinary, uneventful, normal. The real challenge isn't getting the model to say "everything is fine" all the time. It's training a model that can identify just a few special blades hiding in the field and not by chance.

That's what rare event detection is about.



X What is Breed and Battle?

- Models evolve over generations
- Top performers breed children with slightly mutated hyperparameters
- O New models enter the arena and compete
- The strongest survive to the next round

Each round is a survival-of-the-fittest trial.

Only the models best at finding rare signals move forward.

Most ML systems optimize for what's frequent We're optimizing for what's rare — and often, what's most meaningful.

▲ Why This Matters

Most systems are metric-driven but context-blind

- .

 Hand sanitizer dispensers failing to detect dark skin.
- Skin cancer models trained mostly on healthy images.
 Financial models that crumble during rare stress events.

We're not just trying to "get high scores."

- ※ Financial shocks, where early detection changes everything
- Social equity, where the rare is often most vulneral

Seed & Battle Royale for Rare Event Detection

This notebook runs a generational tournament to find the best resampling strategy for rare event classification.

```
[16]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import sys, os
sys.path.insert(0, os.path.abspath("../src"))
                from model_eval import evaluate_model
from breed_and_battle import breed_and_battle
                from resampling import (
    manual_upsampling, smote, adasyn, borderline_smote,
    smote_tomek, smote_enn, random_undersample, cluster_centroids
                , rfrom preprocessing import load_and_prepare_data, split_and_scale, select_top_features from visualization import print_champion_summary, plot_pr_auc_tracking from resamplian_registry import get_resamplers from battle_logger import BattleLogger.
```

Load & Prepare Data

```
[19]: data_path = "../data/synth_rare_event_data.csv"
   target_col = "rare_event"
            X, y = load_and_prepare_data(data_path, target_col)
X_train_scaled, X_test_scaled, y_train, y_test = split_and_scale(X, y)
X_train_X_test__ = select_top_features(X_train_scaled, y_train, X_test_scaled, return_features=True)
```

[21]: df = pd.read_csv(data_path)

@ Compute baseline (random guessing PR AUC)
positive_rate = df["rare_event"].mean()
print(f"\(\sum_{\text{N}} \) Baseline PR AUC (random guessing): \(\{ \text{positive_rate:.3f}} \)")

── Baseline PR AUC (random guessing): 0.018

Resampling Strategies

[24]: # M Load resampling strategies via registry resamplers = get_resamplers(X_train, y_train, target_col) resampled_datasets = {name: fn() for name, fn in resamplers.items()}

X Run the Evolutionary Tournament

[54]: with BattleLogger(
 to_file="battle_log_data.js",
 js_file="battle_log_data.js",
 inject_html=True, # # this enables auto HTML generation

```
html_template="battle_template.html", # your base template
html_output="battle_arena.html" # final version to view
                population = breed_and_battle(
    resampled_datasets=resampled_datasets,
    X_test=x^t_lest,
    y_test=y_lest,
    generations=10,
    top_k=3,
    debug=True
          /Users/sophiaboettcher/anaconda3/lib/python3.11/site-packages/sklearn/linear_model/_sag.py:349: ConvergenceWarning: The max_iter was reached which means the coef_ did not converge warnings.warn(
[58]: with open("battle_log.txt") as f:
    print(f.read()[:1000]) # First 1000 chars
           Initial Fighters Enter the Arena!
          M Generation 1 Begins!
          Top Performers:
1. Cluster Centroids_G0: PR AUC = 0.043
2. Manual Upsampling_G0: PR AUC = 0.035
3. Random Undersample_G0: PR AUC = 0.032

▼ Cluster Centroids breeds 2 children

          Cluster Centroids_child_1_G1 throws their gloves into the arena!

9 Style: Cluster Centroids

$\tilde{G}$ Genome: C = 0.97, 11_ratio = 0.54

$\tilde{G}$ Generation: 1

$\tilde{L}$ Last Score: PR AUC = 0.040

$\tilde{J}$ Lineage: child of Cluster Centroids_60
          Cluster Centroids_child_2_61 throws their gloves into the arena!

Style: Cluster Centroids

Genome: C = 1.09, 11_ratio = 0.40

Generation: 1

Last Score: PR AUC = 0.042

Lineage: child of Cluster Centroids_60

■ Manual Upsampling breeds 2 children

          Manual Upsampling_child_1.G1 throws their gloves into the arena!

9 Style: Manual Upsampling

$\precedup \text{censure}$ can be $\text{disc}$ conserved.

$\precedup \text{censure}$ censure $\text{censure}$ can be $\text{disc}$ conserved.

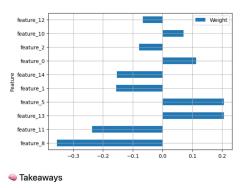
$\text{disc}$ Lineage: child of Manual Upsampling_60
           Manual Upsampling_child_2_G1 throws their
          M Final Champion Summary
TGRAND CHAMPION GAME: Cluster Centroids_G5
Generation: 5
GF Genome: C = 1.02, ll_ratio = 0.52
FR AUC: 0.043
ROCA AUC: 0.638
Lineage: child of Cluster Centroids_G0
          II PR AUC Over Generations
[39]: plot_pr_auc_tracking(population)
          /Users/sophiaboettcher/Param_IndianMutualFunds/A3_FeatureEngineering/rare_event_project/src/visualization.py:36: UserWarning: Glyph 128200 (\W.CKMRT WITH UPWARDS TREND}) missing from current font.
plt.tight_layout()
/Users/sophiaboettcher/anaconda3/lib/python3.11/site-packages/IPython/core/pylabtools.py:170: UserWarning: Glyph 128200 (\W.CHART WITH UPWARD 5 TREND)) missing from current font
fig.carwas.print_figure(bytes_10, **kw)
                                                                                        PR AUC Evolution Over Generations
               0.0425
               0.0400
               0.0375
               0.0350
               0.0325

    Manual Upsampling

               0.0300
                                                                                                                                                                             SMOTE
                                                                                                                                                                            - ADASYN
                                                                                                                                                                             Borderline SMOTE
SMOTETomek
SMOTEENN
               0.0275
               0.0250

    Random Undersample

                                                                                                                                                                           - Cluster Centroids
                                                                                                                                                                            No Resampling
                                                                                                                 Generation
          Baseline vs. Best
[42]: from sklearn.metrics import precision_recall_curve, auc
          # Calculate baseline PR AUC based on class imbalance
y_baseline = np.full_like(y_test, fitl\value=1, dtype=int)  # pretend all positives
_ recall, _ precision_recall.curve(y_test, y_baseline)
baseline.pr_auc = auc(recall, np.full_like(recall, y_test.mean()))
          Baseline PR AUC (random guessing): 0.018 Thampion PR AUC: 0.043 The Improvement: 2.4x
           5 Feature Signals: What Makes an Event 'Rare'?
top_feats.plot(kind="barh", x="Feature", y="Weight", title="Top Predictive Features")
```



- The Breed-and-Battle approach evolves better models across generations.
 The Cluster Centroids strategy consistently produced strong contenders.
 Our final model beats random guessing by -2.4x in terms of PR AUC.
 This shows that even with limited signal, structured evolution can surface meaningful patterns in rare event detection.

Next steps could include trying tree-based models, SHAP value interpretation, or PCA/UMAP visualization.

```
[48]: import webbrowser import os
          def launch_battle_arena(html_filename="battle_arena.html"):
                Opens the generated battle arena HTML file in the default web browser. \\
               html_path = os.path.abspath(html_filename)
if os.path.exists(html_path):
print("" Launching Battle Arena from:\n(html_path)")
webbrowser.open(f"file://{html_path}")
else:
print(f"X Could not find (html_filename). Make sure it has been generated.")
[50]: launch_battle_arena()
          |+ Launching Battle Arena from:
/Users/sophiaboettcher/Param_IndianMutualFunds/A3_FeatureEngineering/rare_event_project/notebooks/battle_arena.html
```