

I still didnt understand the calibration part. Like I understand calibration makes it stable for minority class? But why ? I saw an example of 0.0001 becoming 80% true positive rate. How does this happen?

Let's break down calibration in imbalanced classification with a concrete example:

1. The Core Problem with Uncalibrated Models

In extreme class imbalance (1:60,000), your model becomes **overly conservative** with minority class predictions. Here's why:

Raw Output Example (Before Calibration):

Sample	True Class	Raw Probability (Minority)
A	Minority	0.0001
B	Minority	0.0015
C	Majority	0.0003

- **What you WANT:** If the model says "0.0001 probability" for minority class, you want this to mean "0.01% chance of being minority"
- **ACTUAL REALITY:** Due to imbalance, even 0.0001 might correspond to **80% chance** of being minority class in your dataset

2. Why This Happens - The Math

With 1:60,000 imbalance:

- Baseline probability (prior) of minority class = $1 / 60,001 \approx 0.0000167$
- If your model outputs 0.0001 (6× higher than baseline), this actually indicates:

$$\text{True Probability} = \frac{0.0001}{0.0001 + (1 - 0.0001)} \approx 0.01\%$$

But the actual conversion isn't linear due to:

- Model regularization
- Loss function biases

- Sampling artifacts

3. Calibration Fixes This Mapping

Calibration creates a **non-linear mapping** between raw outputs and true probabilities:

```
# Hypothetical calibration mapping (Platt Scaling)
def calibrate(raw_prob):
    return 1 / (1 + np.exp(-(2.5 + 15*raw_prob))) # Learned coefficients

# After calibration:
print(calibrate(0.0001)) # → 0.82 (82% true probability)
print(calibrate(0.0015)) # → 0.91
print(calibrate(0.0003)) # → 0.27
```

4. Real-World Analogy: Weather Forecasting

- **Uncalibrated Forecaster:** Always predicts "1% chance of rain" regardless of clouds. When it does rain, their 1% prediction actually corresponds to 50% historical rain days.
- **Calibrated Forecaster:** Adjusts predictions so when they say "50% chance", it truly rains 50% of the time.

5. Impact on Classification Thresholds

Threshold	Uncalibrated (Raw)	Calibrated
0.5	Misses ALL minority (requires probability > baseline)	Actually separates classes
0.001	99.9% FPR (too many false positives)	Balanced precision/recall

Before Calibration (Threshold=0.5):

```
Confusion Matrix:
      Predicted 0 Predicted 1
Actual 0      6,000,000         0
Actual 1         100           0
```

After Calibration (Threshold=0.5):

```
Confusion Matrix:
      Predicted 0 Predicted 1
Actual 0      5,800,000    200,000
Actual 1         20         80 # Now detects minorities
```

6. Implementation Visualization

```
from sklearn.calibration import calibration_curve

# Before calibration
prob_true_uncal, prob_pred_uncal = calibration_curve(y_test, raw_probs, n_bins=10)

# After calibration
prob_true_cal, prob_pred_cal = calibration_curve(y_test, calibrated_probs, n_bins=10)

plt.plot(prob_pred_uncal, prob_true_uncal, 's-', label='Uncalibrated')
plt.plot(prob_pred_cal, prob_true_cal, 'o-', label='Calibrated')
plt.plot([0,1], [0,1], 'k--', label='Perfect')
```

Calibration Curve

Key Takeaways

1. **Calibration \neq Balancing** - It fixes probability interpretations, not class frequencies
2. **Critical for Threshold Tuning** - Lets you use standard thresholds (0.5) meaningfully
3. **Preserves Ranking** - Good samples still have higher scores than bad ones, just rescaled
4. **Required for Meaningful Metrics** - AUC-PR, Brier Score, and Expected Calibration Error become trustworthy

For your 1:60,000 problem, calibration ensures that a predicted 1% probability actually means "1% chance of being minority", not "0.0001% chance" as raw outputs suggest.