

# When Class Imbalance Shifts the Decision Boundary: A Short Proof

Boammani Aser Lompo

October 7, 2025

## 1 Setup

Let  $Y \in \{0, 1\}$  with true prior  $\pi = P(Y = 1)$  and feature vector  $X \in \mathcal{X}$ . Let  $p(x | y)$  denote class-conditional densities and  $p(x) = \sum_y p(x | y)P(Y = y)$ . Suppose we train on an imbalanced sample with empirical prior  $\hat{\pi}$  (typically  $\hat{\pi} \neq \pi$ ).

## 2 Bayes decision rule

Under 0–1 loss with equal costs, the Bayes classifier predicts 1 iff

$$\log \frac{p(x | 1)}{p(x | 0)} \geq \log \frac{1 - \pi}{\pi}.$$

Equivalently,  $p(Y = 1 | x) \geq \frac{1}{2}$ .

**Proposition 1** (Imbalanced-sample boundary shift). *If a learner minimizes empirical 0–1 risk on a sample whose class prior is  $\hat{\pi}$ , then as  $n \rightarrow \infty$  it converges to the decision rule*

$$\log \frac{p(x | 1)}{p(x | 0)} \geq \log \frac{1 - \hat{\pi}}{\hat{\pi}}.$$

*If  $\hat{\pi} \neq \pi$ , the decision threshold differs from the true Bayes threshold by*

$$\Delta = \log \frac{\pi}{1 - \pi} - \log \frac{\hat{\pi}}{1 - \hat{\pi}},$$

*biasing decisions toward the majority class of the training sample.*

*Proof.* Minimizing empirical 0–1 risk on the sample is equivalent to using the sample distribution as the data-generating law; the Bayes optimal under that law uses prior  $\hat{\pi}$ . Therefore the induced threshold is  $\log \frac{1 - \hat{\pi}}{\hat{\pi}}$ . Subtract the true Bayes threshold  $\log \frac{1 - \pi}{\pi}$  to obtain  $\Delta$ .  $\square$

## 3 Cross-entropy and prior correction

Under cross-entropy, logistic models learn  $p_{\text{train}}(Y = 1 | x)$  with prior  $\hat{\pi}$ . If deployment prior  $\pi$  differs, correct via

$$\text{logit } p_{\text{test}}(1 | x) = \text{logit } p_{\text{train}}(1 | x) + \log \frac{\pi/(1 - \pi)}{\hat{\pi}/(1 - \hat{\pi})}.$$

Alternatively, use class-weighted losses with  $w_y \propto 1/\hat{\pi}_y$  or resampling to neutralize the sample prior.

## 4 Practical notes

State assumptions (no covariate shift, equal costs), clarify what “bias” means (boundary shift vs calibration vs minority metrics), and mention regularization/data scarcity effects.

## 5 Conclusion

Imbalance per se shifts the decision rule when the sample prior differs from deployment prior. The effect can be neutralized with weighting, resampling, calibrated thresholds, or prior-shift correction.