

CLASSIFICATION

"A or Not A" — Pick prob. threshold and that's your hypothesis.

But need new eval metrics! For linear we used RMSE (measuring dist from prediction to label). For logistic:

— Accuracy? Fraction of predictions correct
↳ Imbalance problems when positive/negative is very rare (in the distribution)
Consider {true, false} {pos., neg.}

— Precision and Recall

$$\left(\frac{TP}{TP+FP}\right) \quad \left(\begin{array}{l} \text{Predicted true} \\ \text{was it true?} \end{array}\right) \quad \left(\begin{array}{l} \text{How many trues} \\ \text{did we predict true?} \end{array}\right) \quad \left(\frac{TP}{TP+FN}\right)$$

These are in tension:

- too much precision can lower recall
- vice versa

ROC Curve Evaluate every possible class. threshold
(Receiver Operating Characteristic)

For each threshold $\epsilon \in [0, 1]$, let

x = false pos. rate @ threshold and

y = true pos. rate @ threshold. Then

the ROC curve is defined by $\{x(t), y(t) \mid t \in [0, 1]\}$

$$\text{AUC (area under ROC curve)} = \int_{x_0}^{x_1} y(t) dx(t)$$

Scale invariant
(not comparing abs. values)

Threshold invariant

(For random pos. + random neg., what is $P(\text{model ranks correctly})$?
— this is AUC)

Prediction bias: Want predictions to match observations. $\text{Avg. (pred.)} = \text{Avg. (obs.)}$

→ Easy to fool but is a good canary.

→ Too much bias?? Consider:

- incomplete feature set
- too much regularization
- buggy pipeline
- biased training sample

→ Fix in the model, not "calibration layer"