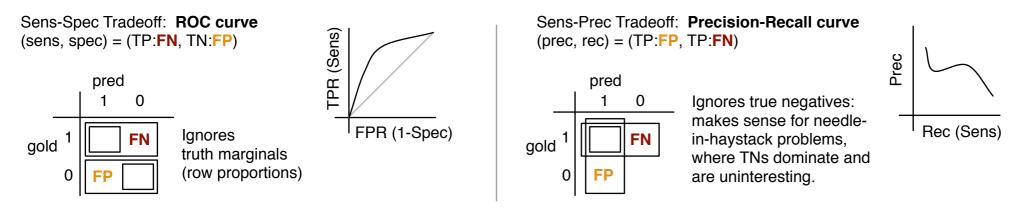


Confidence Tradeoffs

Classifier says YES if confidence in '1' is greater than a threshold t. Varying the threshold trades off between FPs versus FNs. For any one metric, you can make it arbitrarily good by being either very liberal or very conservative. Therefore you are often interested in a pair of metrics: one that cares about FPs, and one that cares about FNs. Here are two pairings that are popular. By changing the confidence threshold, you get a curve of them.



There are several summary statistics to make a **FP-FN** metric pairing invariant to the choice of threshold.

AUC area under the ROC curve. (a.k.a. ROCAUC, AUROC) Semantics: prob the classifier correctly orders a randomly chosen pos and neg example pair. (Similar to Wilcoxon-Mann-Whitney or Kendall's Tau). Neat analytic tricks with convex hulls.

PRAUC: area under the PR curve. No semantics that I know of. Similar to mean average precision. See also precision@K.

PR Breakeven: find threshold where prec=rec, and report that value.

Alternatively, directly use predicted probabilities without thresholding or confusion matrix: Cross-entropy (log-likelihood) or Brier score loss (squared error).

There are also summary stats that require a specific threshold -- or you can use if you don't have confidences -but have been designed to attempt to care about both FP and FN costs.

Balanced accuracy: arithmetic average of sens and spec. This is accuracy if you chose among pos and neg classes with equal prob each -- i.e. a balanced population. [In contrast, AUC is a balanced rank-ordering accuracy.] Corresponds to linear isoclines in ROC space.

F-score: harmonic average of prec and rec. Usually F1 is used: the unweighted harmonic mean. $F1 = 2^{*}P^{*}R / (P+R)$ No semantics that I know of. Corresponds to curvy isoclines in PR space.

Multiclass setting: each class has it own prec/rec/f1. Macroaveraged F1: mean of each class' F1. Cares about rare classes as much as common classes. **Microaveraged** F1: take total TP,FP,FN to calculate prec/rec/f1. Cares more about highly prevalent classes.

Balanced accuracy = macroaveraged recall

Accuracy = microavg rec = microavg prec = microavg f1

Classical Frequentist Hypothesis Testing

Task is to detect null hypothesis violations. Define "reject null" as a positive prediction.

- Type I error = false positive. Size (a) = 1-Specificity - Type II error = false negative. Power (β) = Sensitivity and **p-value** = 1-Specificity, kind of.

"Significance" refers to either p-value or size (i.e. FPR)

The point is to give non-Bayesian probabilistic semantics to a single test, so these relationships are a little more subtle. For example, size is a worst-case bound on specificity, and a p-value is for one example: prob the null could generate a test statistic at least as extreme (= size of the strictest test that would reject the null on that example). When applied to a single test, PPV and NPV are not frequentist-friendly concepts since they depend on priors. But in multiple testing, they can be frequentist again; e.g. Benjamini-Hochberg bounds expected PPV by assuming PriorNeg=0.9999999.

References

Diagrams based on William Press's slides: http://www.nr.com/CS395T/lectures2008/17-ROCPrecisionRecall.pdf Listings of many metrics on: http://en.wikipedia.org/wiki/Receiver_operating_characteristic Fawcett (2006), "An Introduction to ROC analysis": http://people.inf.elte.hu/kiss/13dwhdm/roc.pdf Manning et al (2008) "Intro to IR", e.g. http://nlp.stanford.edu/IR-book/html/htmledition/evaluation-of-ranked-retrieval-results-1.html Hypothesis testing: Casella and Berger, "Statistical Inference" (ch 8); also Wasserman, "All of Statistics" (ch 10)