

Evaluating Hypotheses

What to measure?

How to measure it?

Two key difficulties when only limited data is available:

- Bias in the estimate
- Variance in the estimate

Estimating Hypothesis Accuracy

- Given a hypothesis h and a data sample containing n examples drawn at random from D , what is the best estimate of the accuracy of h over future instances?
- What is the probable error in this accuracy estimate?

Sample Error and True Error

- Sample error of hypothesis h with respect to target function f and data sample S is
$$error_S(h) = (1/n) \sum_{x \in S} \{\delta_K[f(x), h(x)]\}$$
- True error of hypothesis h with respect to target function f and distribution D is

$$error_D(h) = \Pr_{x \in D} [f(x) \neq h(x)]$$

How good an estimate of $error_D(h)$ is $error_S(h)$?

Confidence Intervals

For discrete-valued hypotheses, $n \geq 30$ and $error_s(h)=r/n$, statistical theory asserts:

- The most probable value of $error_D(h)$ is $error_s(h)$
- With approximately $N\%$ probability, $error_D(h)$ lies in the interval:

$$error_s(h) \pm z_N[error_s(h)(1 - error_s(h))/n]^{1/2}$$

Example: $N = 95$ $z_N \cong 1.96$

Valid when (Rule of Thumb): $n error_s(h) [1-error_s(h)] \geq 5$

In general

We can estimate error measures over finite samples

- With no bias
- With a variance that decreases with sample size

Regression: RMSE

- ▶ The Root-Mean Squared Error (RMSE) is a cost function for regression. The formula for the RMSE is:

$$RMSE(f) = \sqrt{\frac{1}{m} \sum_{i=1}^m (f(\mathbf{x}_i) - y_i)^2}$$

where m is the number of test examples, $f(\mathbf{x}_i)$, the predicted output on \mathbf{x}_i and y_i the actual values.

RMSE has no scale.

Normalized Root-Mean Squared Error (NRMSE), defined as RMSE over the variance of the data, is a better measure.

Confusion Matrix-Based Performance Measures

True class-> Hypothesized class	Pos	Neg
Yes	TP	FP
No	FN	TN
	P=TP+FN	N=FP+TN

Confusion Matrix

- ▶ **Multi-Class Focus:**
 - **Accuracy** = $(TP+TN)/(P+N)$
- ▶ **Single-Class Focus:**
 - **Precision** = $TP/(TP+FP)$
 - **Recall** = TP/P
 - **Fallout** = FP/N
 - **Sensitivity** = $TP/(TP+FN)$
 - **Specificity** = $TN/(FP+TN)$

Some issues with performance measures

True class->	Pos	Neg	True class ->	Pos	Neg
Yes	200	100	Yes	400	300
No	300	400	No	100	200
	P=500	N=500		P=500	N=500

- ▶ Both classifiers obtain **60% accuracy**
- ▶ They exhibit very different behaviours:
 - On the left: **weak** positive recognition rate/**strong** negative recognition rate
 - On the right: **strong** positive recognition rate/**weak** negative recognition rate

Some issues with performance measures (cont'd)

True class →	Pos	Neg	True class →	Pos	Neg
Yes	500	5	Yes	450	1
No	0	0	No	50	4
	P=500	N=5		P=500	N=5

- ▶ The classifier on the left obtains 99.01% accuracy while the classifier on the right obtains 89.9%
 - Yet, the classifier on the right is much more sophisticated than the classifier on the left, which just labels everything as positive and misses all the negative examples.

Some issues with performance measures (cont'd)

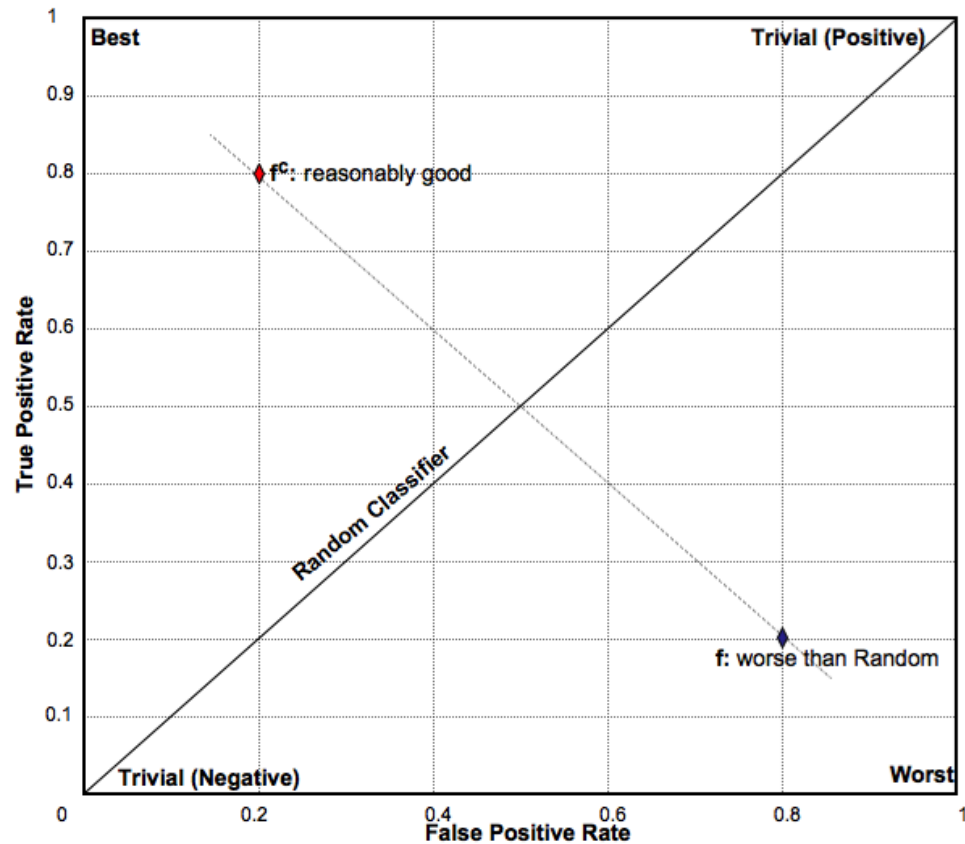
True class →	Pos	Neg	True class →	Pos	Neg
Yes	200	100	Yes	200	100
No	300	400	No	300	0
	P=500	N=500		P=500	N=100

- ▶ Both classifiers obtain the **same precision and recall** values of 66.7% and 40% (Note: the data sets are different)
- ▶ They exhibit very different behaviours:
 - Same positive recognition rate
 - Extremely different negative recognition rate: **strong** on the left / **nil** on the right
- ▶ Note: Accuracy has no problem catching this!

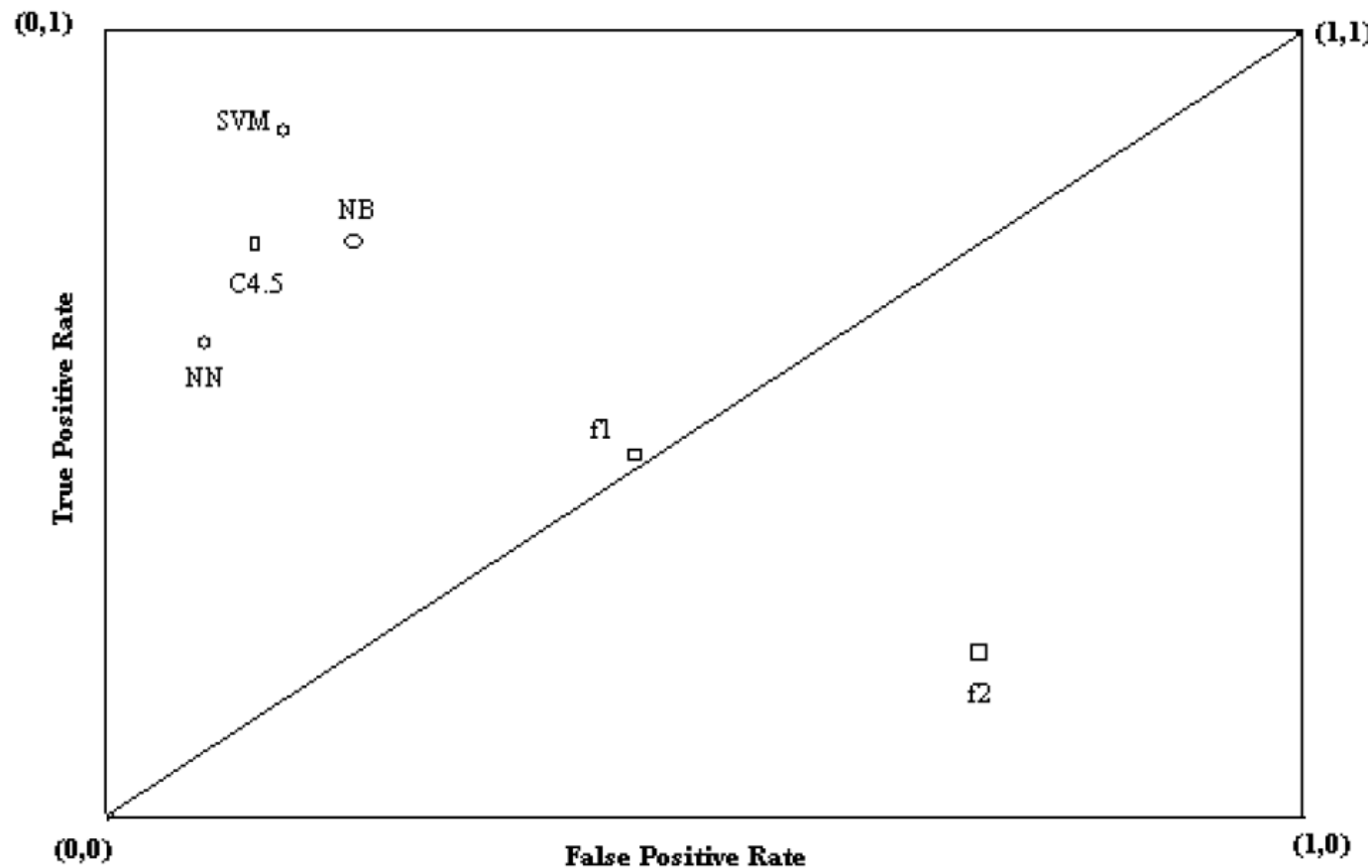
ROC Curves

- ▶ Performance measures for scoring classifiers
 - most classifiers are in fact scoring classifiers
 - Scores needn't be in pre-defined intervals, or even likelihood or probabilities
- ▶ ROC analysis has origins in Signal detection theory to set an operating point for desired signal detection rate
 - Signal assumed to be corrupted by noise (Normally distributed)
- ▶ ROC maps FPR on horizontal axis and TPR on the vertical axis; Recall that
 - $\text{FPR} = \text{FP}/(\text{FP} + \text{TN}) = 1 - \text{Specificity}$
 - $\text{TPR} = \text{TP}/(\text{TP} + \text{FN}) = \text{Sensitivity}$

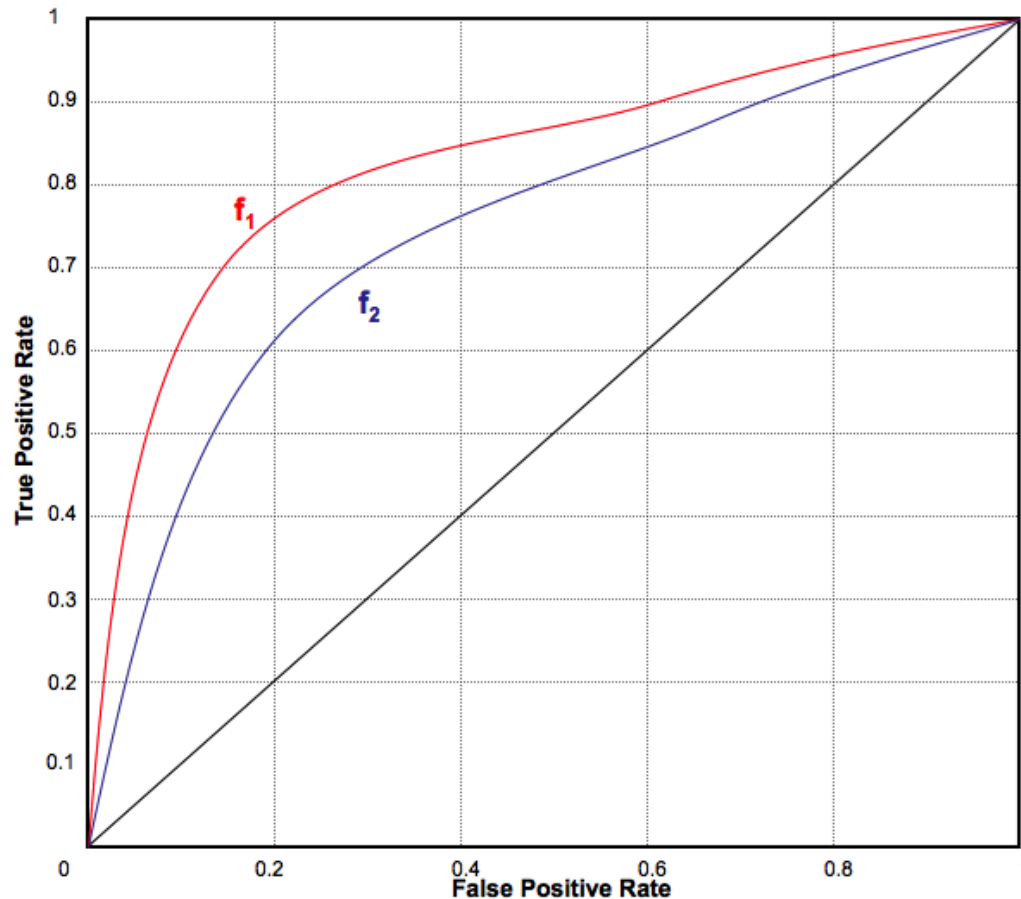
ROC Space



ROC Plot for discrete classifiers



ROC plot for two hypothetical scoring classifiers



AUC

- ▶ ROC Analysis allows a user to visualize the performance of classifiers over their operating ranges.
- ▶ However, it does not allow the quantification of this analysis, which would make comparisons between classifiers easier.
- ▶ The Area Under the ROC Curve (AUC) allows such a quantification: it represents the performance of the classifier averaged over all the possible cost ratios.

Se we decided on a performance measure.

- ▶ How do we estimate it in an unbiased manner?
- ▶ What if we used all the data?
 - Re-substitution: Overly optimistic (best performance achieved with complete over-fit)

Hold-out approach

- ▶ Set aside a separate test set T . Evaluate your measure on this set
- ▶ Pros
 - independence from training set
 - Generalization behavior can be characterized
 - Estimates can be obtained for *any* classifier
- Cons
 - We lose data for learning

The need for re-sampling

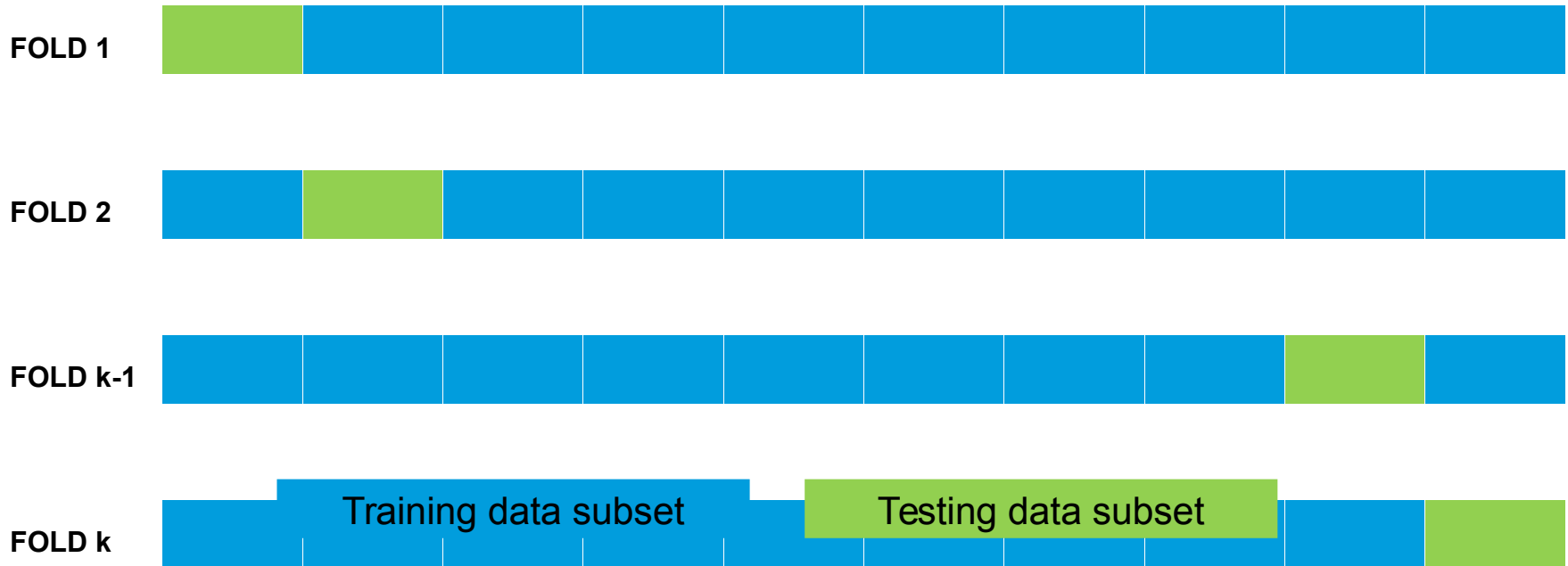
- ▶ Too few training examples -> learning a poor classifier
 - Having too few examples in the training set affects the bias of algorithm by making its average prediction unreliable
- ▶ Too few test examples -> bogus error estimates
 - Having too few examples in test set results in high variance in the estimation
- ▶ Hence: **Resampling**
 - Delivers accurate performance estimates while allowing the algorithm to train on most data examples

Simple Resampling:

Repeated hold-out

- ▶ Set aside a separate random test set T . Evaluate your measure on this set
- ▶ Repeat n times
- ▶ Average your estimations
- ▶ Test and training sets have partial overlap.
No independence.

Simple Resampling: K-fold Cross Validation



In Cross-Validation, the data set is divided into k folds and at each iteration, a different fold is reserved for testing while all the others are used for training the classifiers.

Observations

- ▶ k-fold CV is arguable the best known and most commonly used resampling technique
 - With k of reasonable size, less computer intensive than Leave-One-Out
 - Easy to apply
- ▶ In all the variants, the testing sets are independent of one another, as required, by many statistical testing methods
 - but the training sets are highly overlapping. This can affect the bias of the error estimates (generally mitigated for large dataset sizes).
- ▶ Results in an averaged estimate over k different classifiers