

Evaluating Machine Learning Methods

www.cs.wisc.edu/~dpage/cs760/

Goals for the lecture

you should understand the following concepts

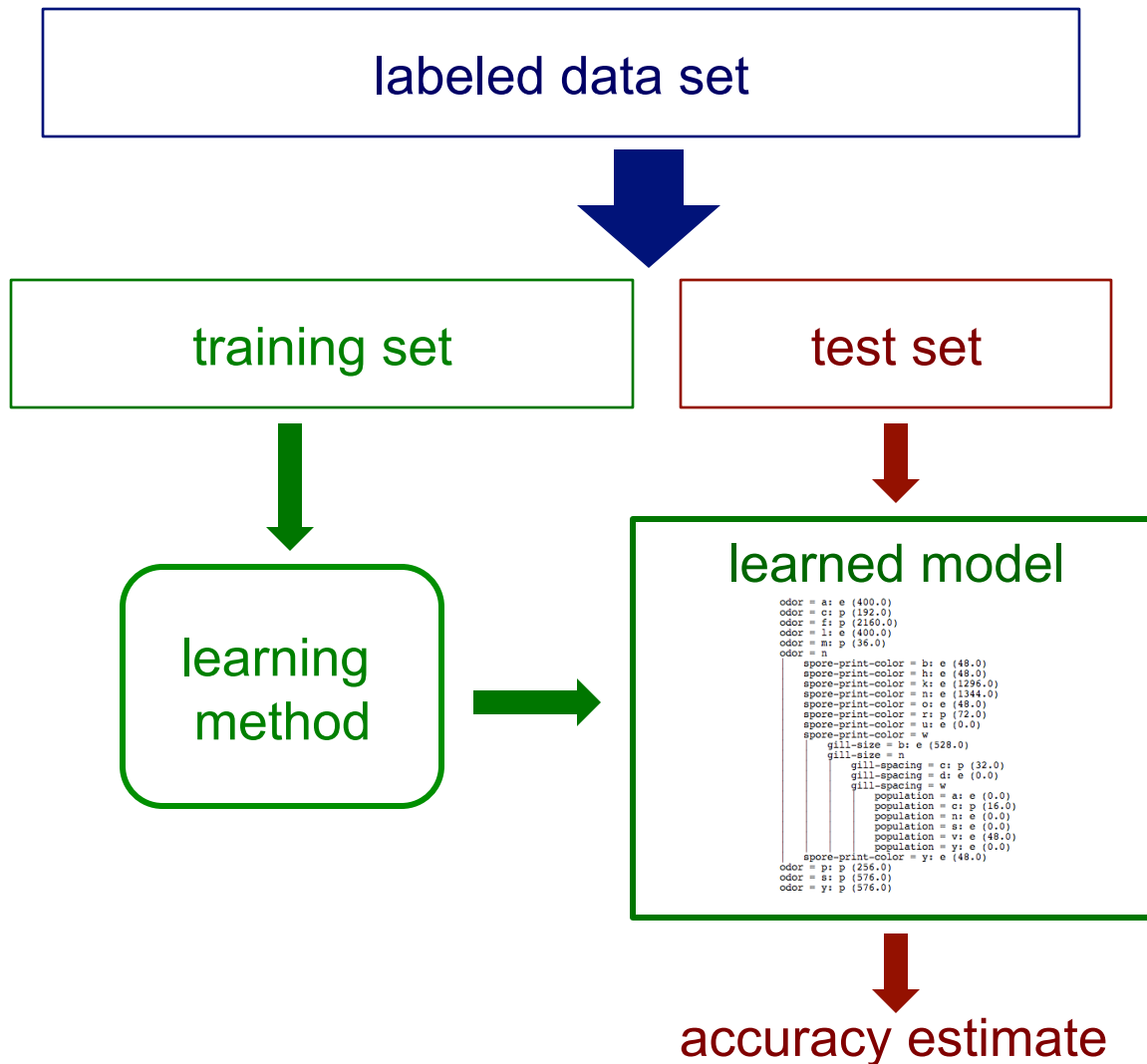
- test sets
- learning curves
- validation (tuning) sets
- stratified sampling
- cross validation
- internal cross validation
- confusion matrices
- TP, FP, TN, FN
- ROC curves
- confidence intervals for error
- pairwise t -tests for comparing learning systems
- scatter plots for comparing learning systems
- lesion studies

Goals for the lecture (continued)

- recall/sensitivity/true positive rate (TPR)
- precision/positive predictive value (PPV)
- specificity and false positive rate (FPR or 1-specificity)
- precision-recall (PR) curves

Test sets revisited

How can we get an unbiased estimate of the accuracy of a learned model?



Test sets revisited

How can we get an unbiased estimate of the accuracy of a learned model?

- when learning a model, you should pretend that you don't have the test data yet (it is "in the mail")*
- if the test-set labels influence the learned model in any way, accuracy estimates will be biased

* In some applications it is reasonable to assume that you have access to the feature vector (i.e. x) but not the y part of each test instance.

Learning curves

How does the accuracy of a learning method change as a function of the training-set size?

this can be assessed by plotting *learning curves*

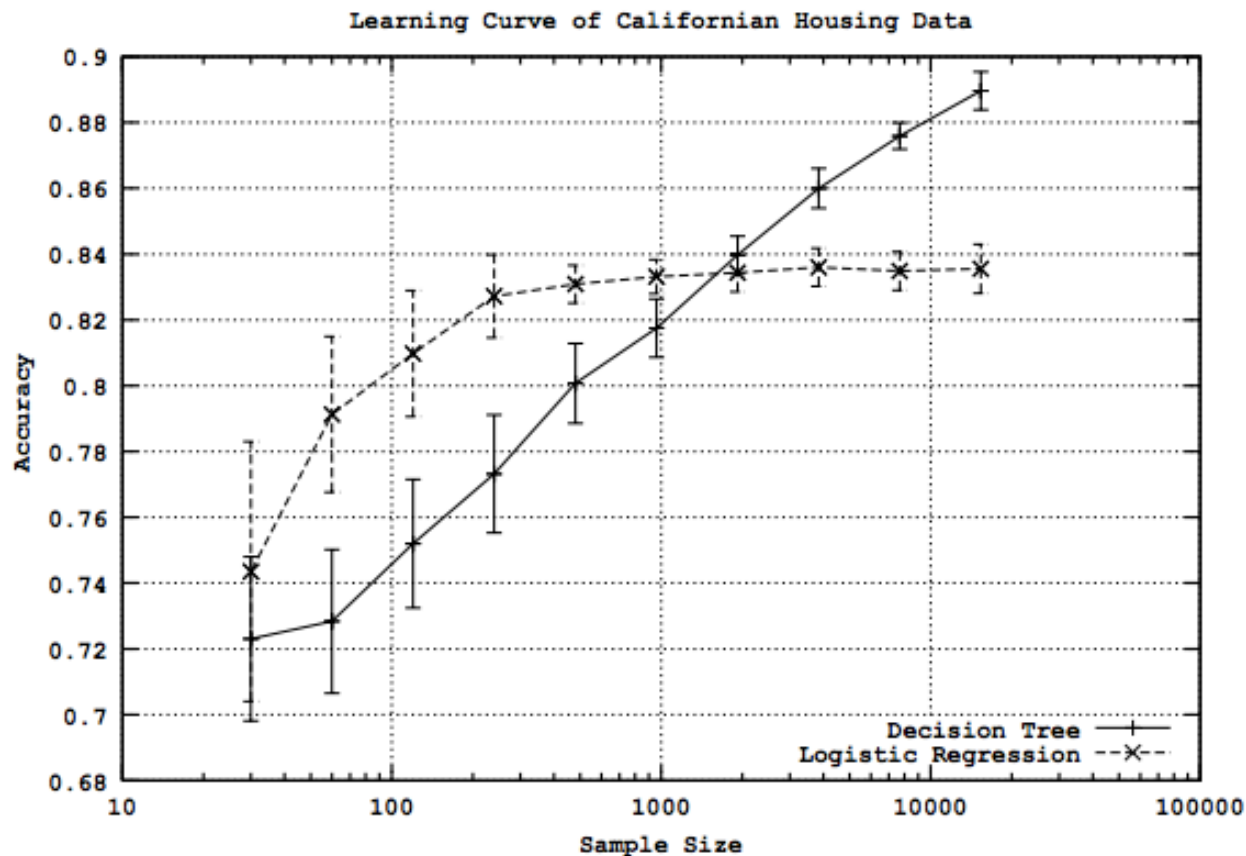
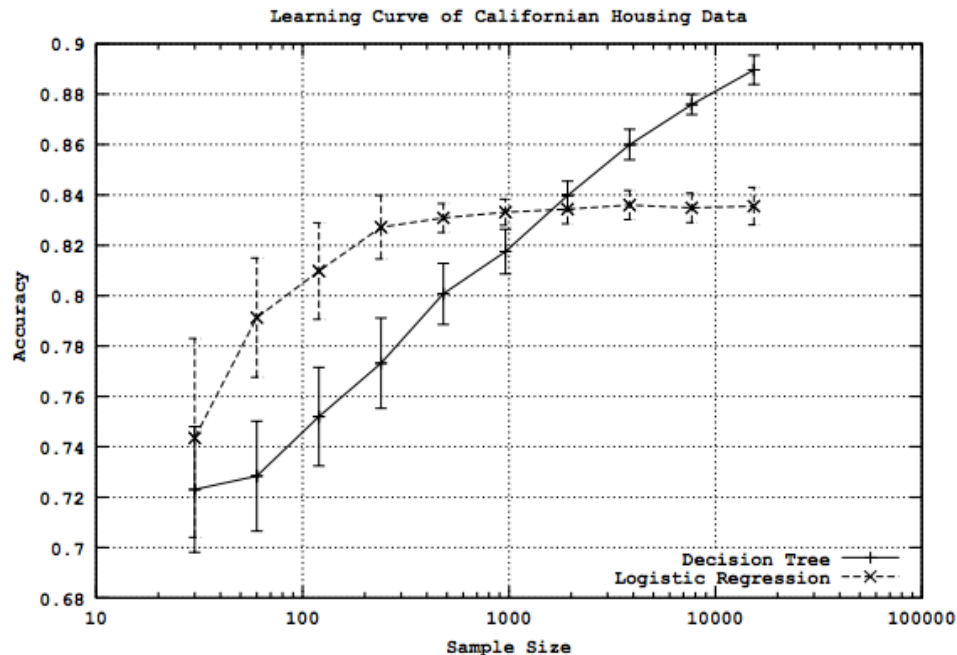


Figure from Perlich et al. *Journal of Machine Learning Research*, 2003

Learning curves

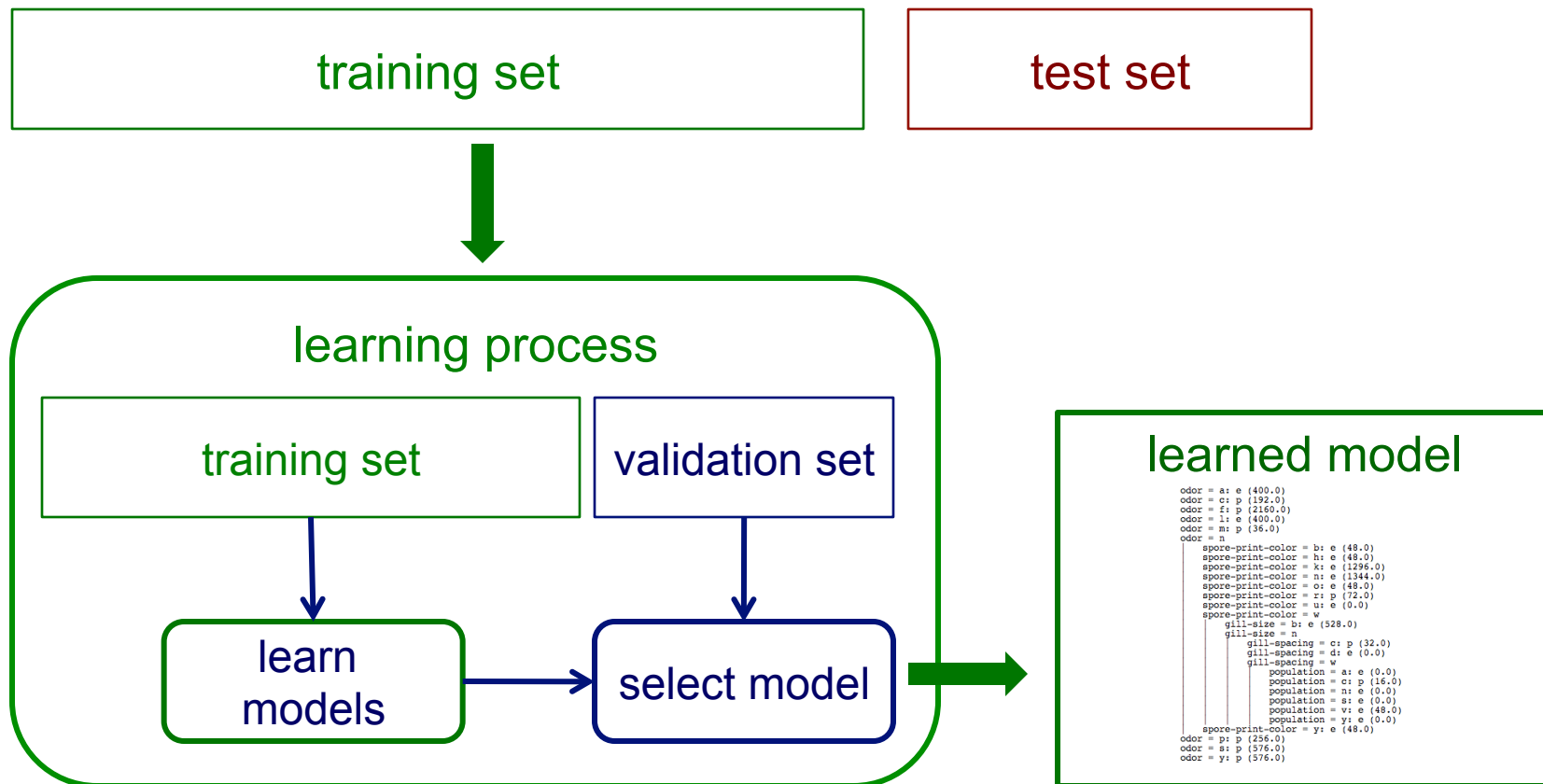
given training/test set partition

- for each sample size s on learning curve
 - (optionally) repeat n times
 - randomly select s instances from training set
 - learn model
 - evaluate model on test set to determine accuracy a
 - plot (s, a) or $(s, \text{avg. accuracy and error bars})$



Validation (tuning) sets revisited

Suppose we want unbiased estimates of accuracy during the learning process (e.g. to choose the best level of decision-tree pruning)?



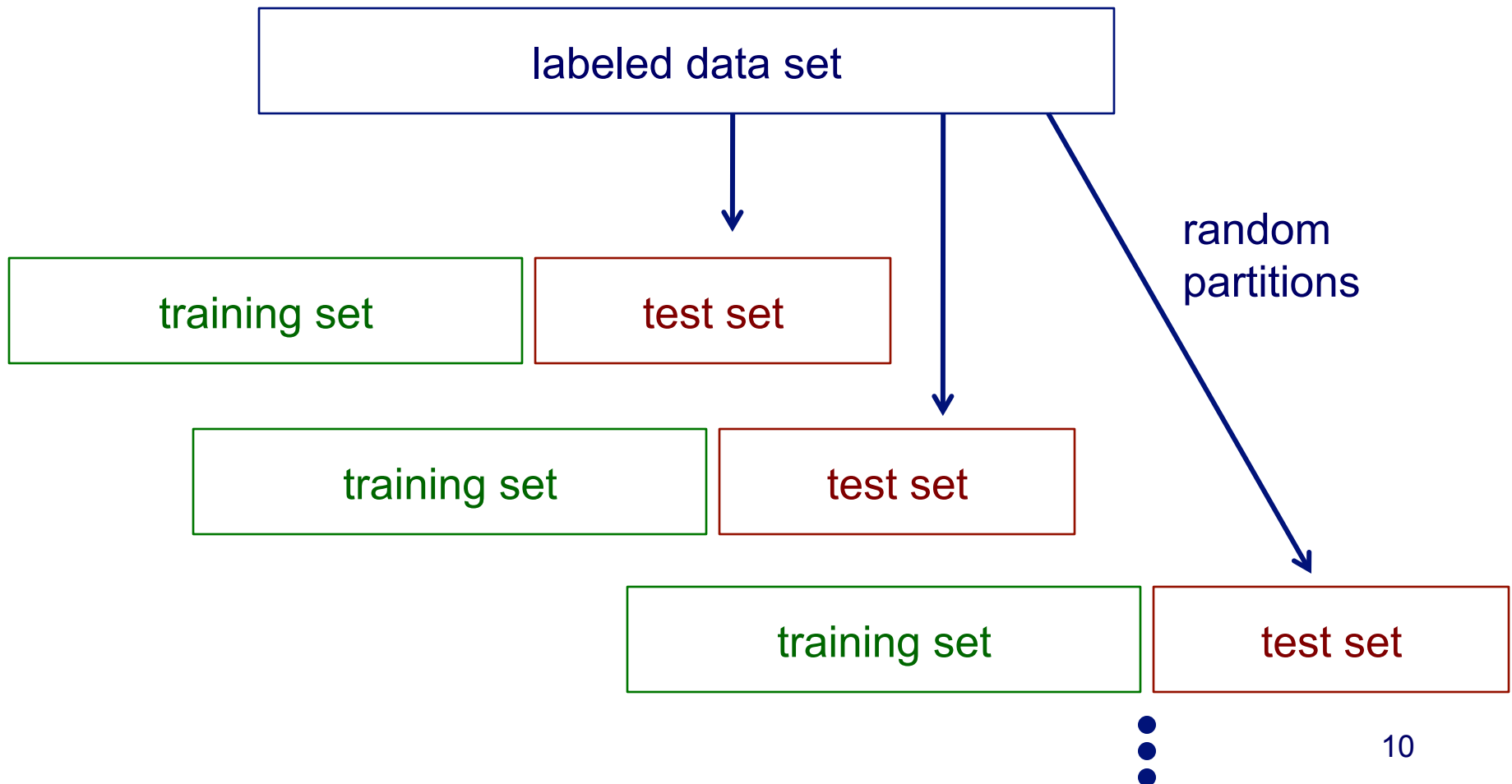
Partition training data into separate training/validation sets

Limitations of using a single training/test partition

- we may not have enough data to make sufficiently large training and test sets
 - a larger test set gives us more reliable estimate of accuracy (i.e. a lower variance estimate)
 - but... a larger training set will be more representative of how much data we actually have for learning process
- a single training set doesn't tell us how sensitive accuracy is to a particular training sample

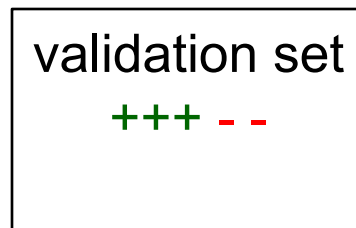
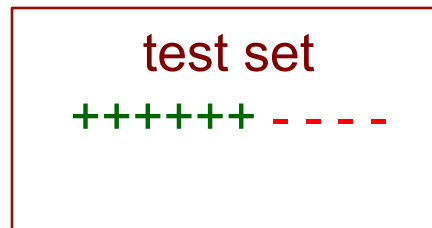
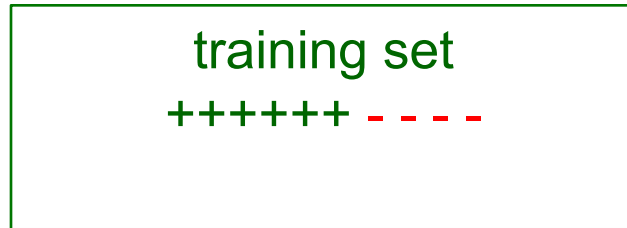
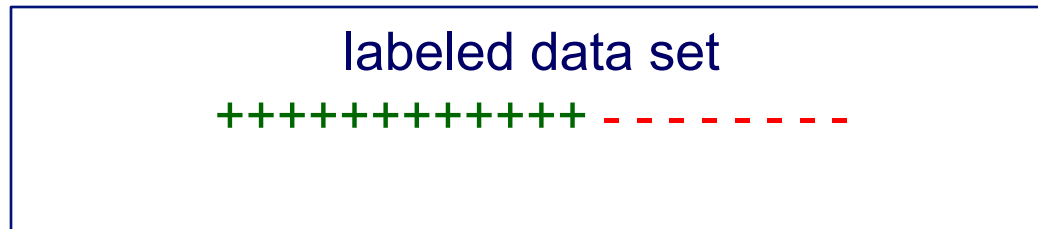
Random resampling

We can address the second issue by repeatedly randomly partitioning the available data into training and test sets.



Stratified sampling

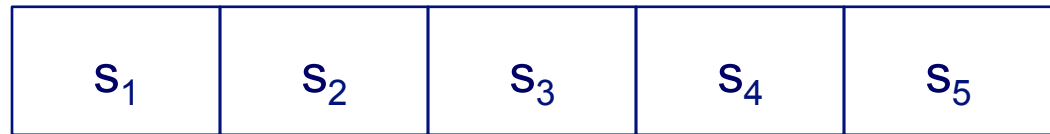
When randomly selecting training or validation sets, we may want to ensure that class proportions are maintained in each selected set



This can be done via stratified sampling: first stratify instances by class, then randomly select instances from each class proportionally.

Cross validation

labeled data set



partition data
into n subsamples

iteratively leave one
subsample out for
the test set, train on
the rest

iteration	train on	test on
1	S ₂ S ₃ S ₄ S ₅	S ₁
2	S ₁ S ₃ S ₄ S ₅	S ₂
3	S ₁ S ₂ S ₄ S ₅	S ₃
4	S ₁ S ₂ S ₃ S ₅	S ₄
5	S ₁ S ₂ S ₃ S ₄	S ₅

Cross validation example

Suppose we have 100 instances, and we want to estimate accuracy with cross validation

iteration	train on	test on	correct
1	s ₂ s ₃ s ₄ s ₅	s ₁	11 / 20
2	s ₁ s ₃ s ₄ s ₅	s ₂	17 / 20
3	s ₁ s ₂ s ₄ s ₅	s ₃	16 / 20
4	s ₁ s ₂ s ₃ s ₅	s ₄	13 / 20
5	s ₁ s ₂ s ₃ s ₄	s ₅	16 / 20

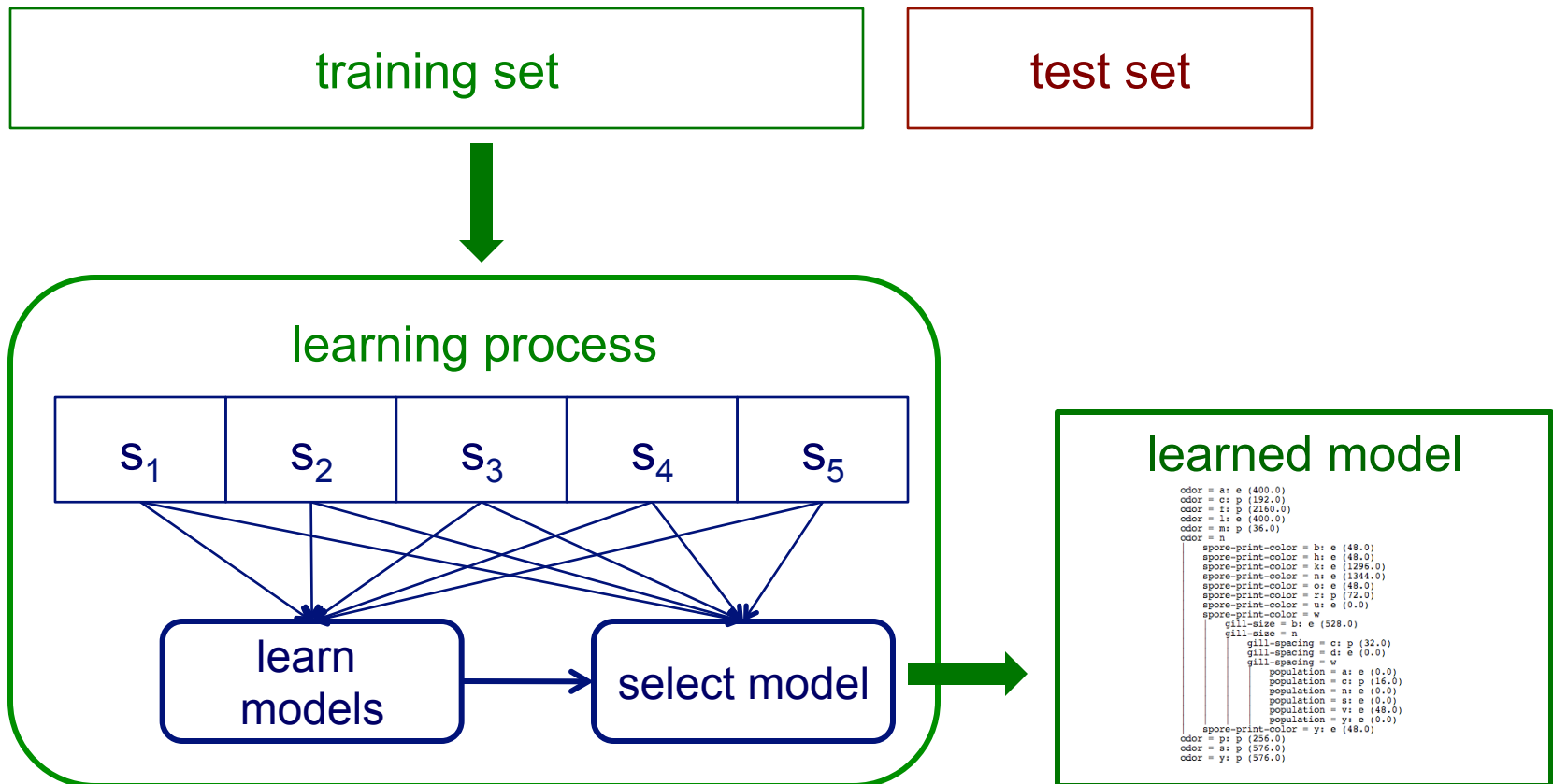
$$\text{accuracy} = 73/100 = 73\%$$

Cross validation

- 10-fold cross validation is common, but smaller values of n are often used when learning takes a lot of time
- in *leave-one-out* cross validation, $n = \#$ instances
- in *stratified* cross validation, stratified sampling is used when partitioning the data
- CV makes efficient use of the available data for testing
- note that whenever we use multiple training sets, as in CV and random resampling, we are evaluating a learning method as opposed to an individual learned model

Internal cross validation

Instead of a single validation set, we can use cross-validation within a training set to select a model (e.g. to choose the best level of decision-tree pruning)



Example: using internal cross validation to select k in k -NN

given a training set

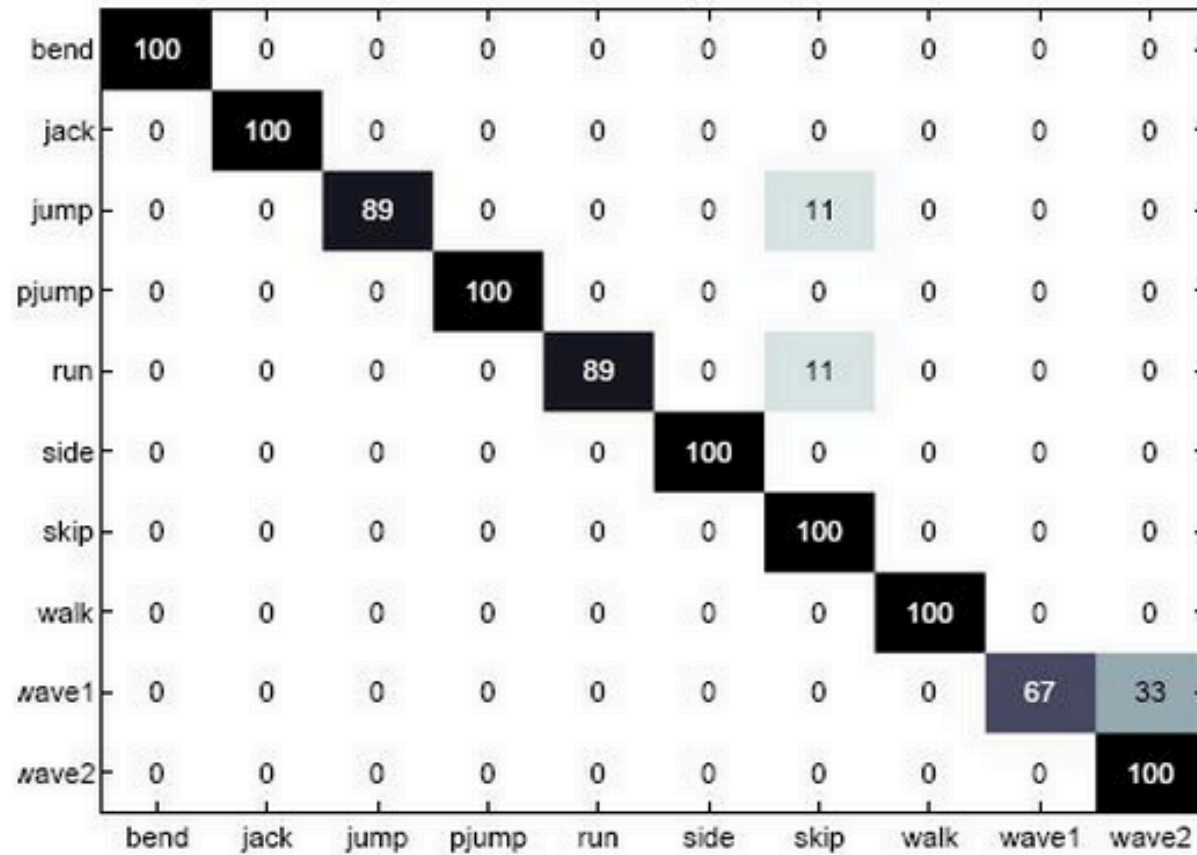
1. partition training set into n folds, $s_1 \dots s_n$
2. for each value of k considered
for $i = 1$ to n
learn k -NN model using all folds but s_i
evaluate accuracy on s_i
3. select k that resulted in best accuracy for $s_1 \dots s_n$
4. learn model using entire training set and selected k

the steps inside the box are run independently for each training set (i.e. if we're using 10-fold CV to measure the overall accuracy of our k -NN approach, then the box would be executed 10 times)

Confusion matrices

How can we understand what types of mistakes a learned model makes?

activity recognition from video



actual class

predicted class

Confusion matrix for 2-class problems

		actual class	
		positive	negative
predicted class	positive	true positives (TP)	false positives (FP)
	negative	false negatives (FN)	true negatives (TN)

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$$

Is accuracy an adequate measure of predictive performance?

- accuracy may not be useful measure in cases where
 - there is a large class skew
 - Is 98% accuracy good if 97% of the instances are negative?
 - there are differential misclassification costs – say, getting a positive wrong costs more than getting a negative wrong
 - Consider a medical domain in which a false positive results in an extraneous test but a false negative results in a failure to treat a disease
- we are most interested in a subset of high-confidence predictions

Other accuracy metrics

actual class

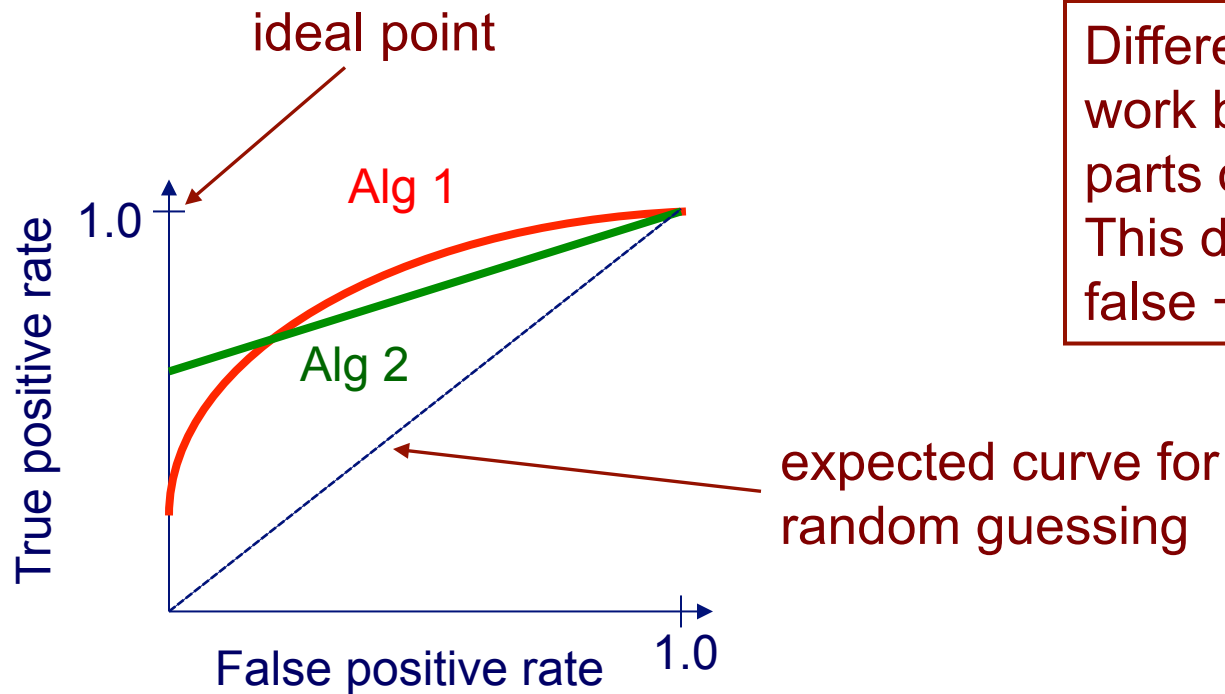
		actual class	
		positive	negative
predicted class	positive	true positives (TP)	false positives (FP)
	negative	false negatives (FN)	true negatives (TN)

$$\text{true positive rate (recall)} = \frac{\text{TP}}{\text{actual pos}} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{false positive rate} = \frac{\text{FP}}{\text{actual neg}} = \frac{\text{FP}}{\text{TN} + \text{FP}}$$

ROC curves

A Receiver Operating Characteristic (ROC) curve plots the TP-rate vs. the FP-rate as a threshold on the confidence of an instance being positive is varied



Different methods can work better in different parts of ROC space. This depends on cost of false + vs. false -

ROC curve example

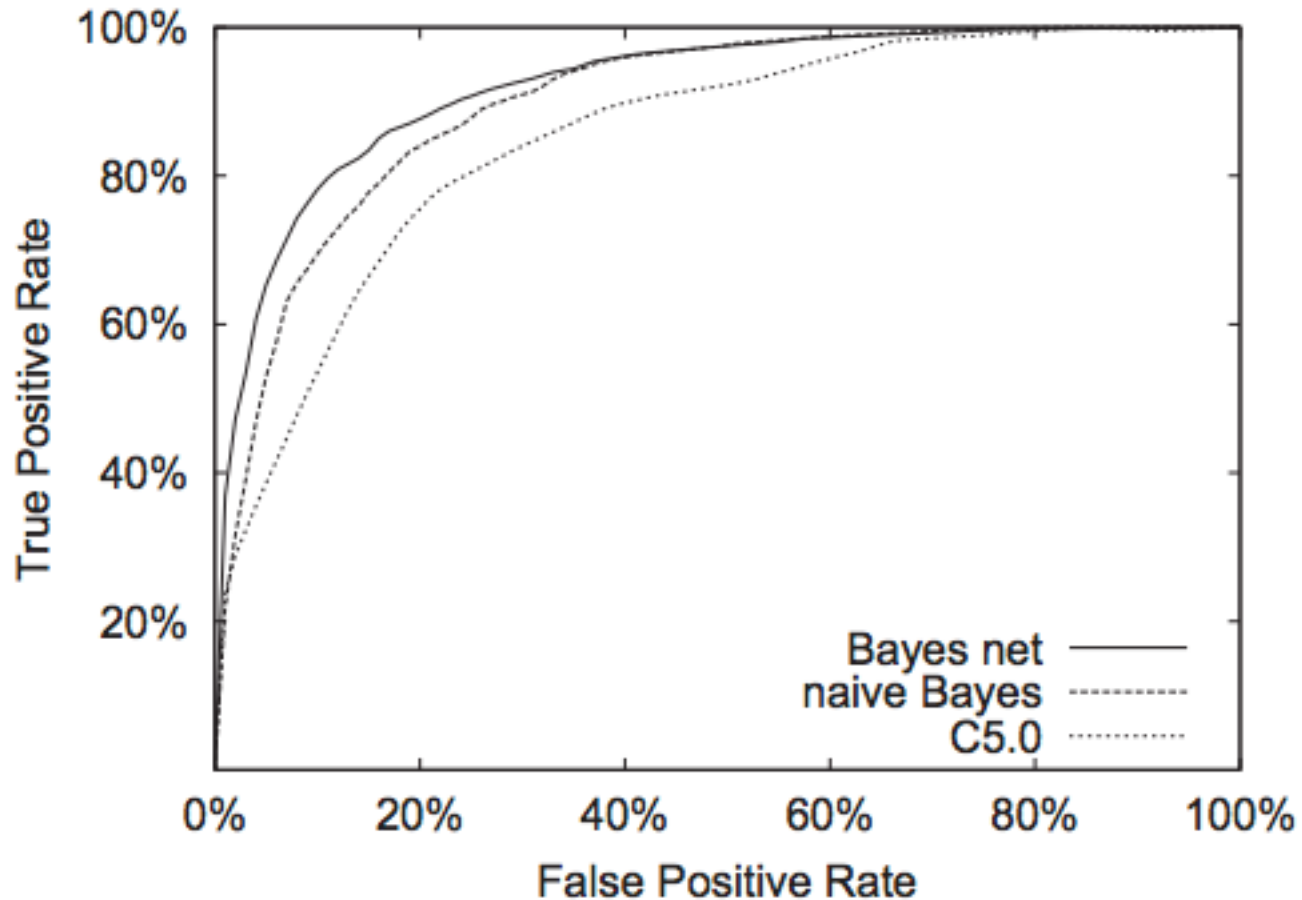
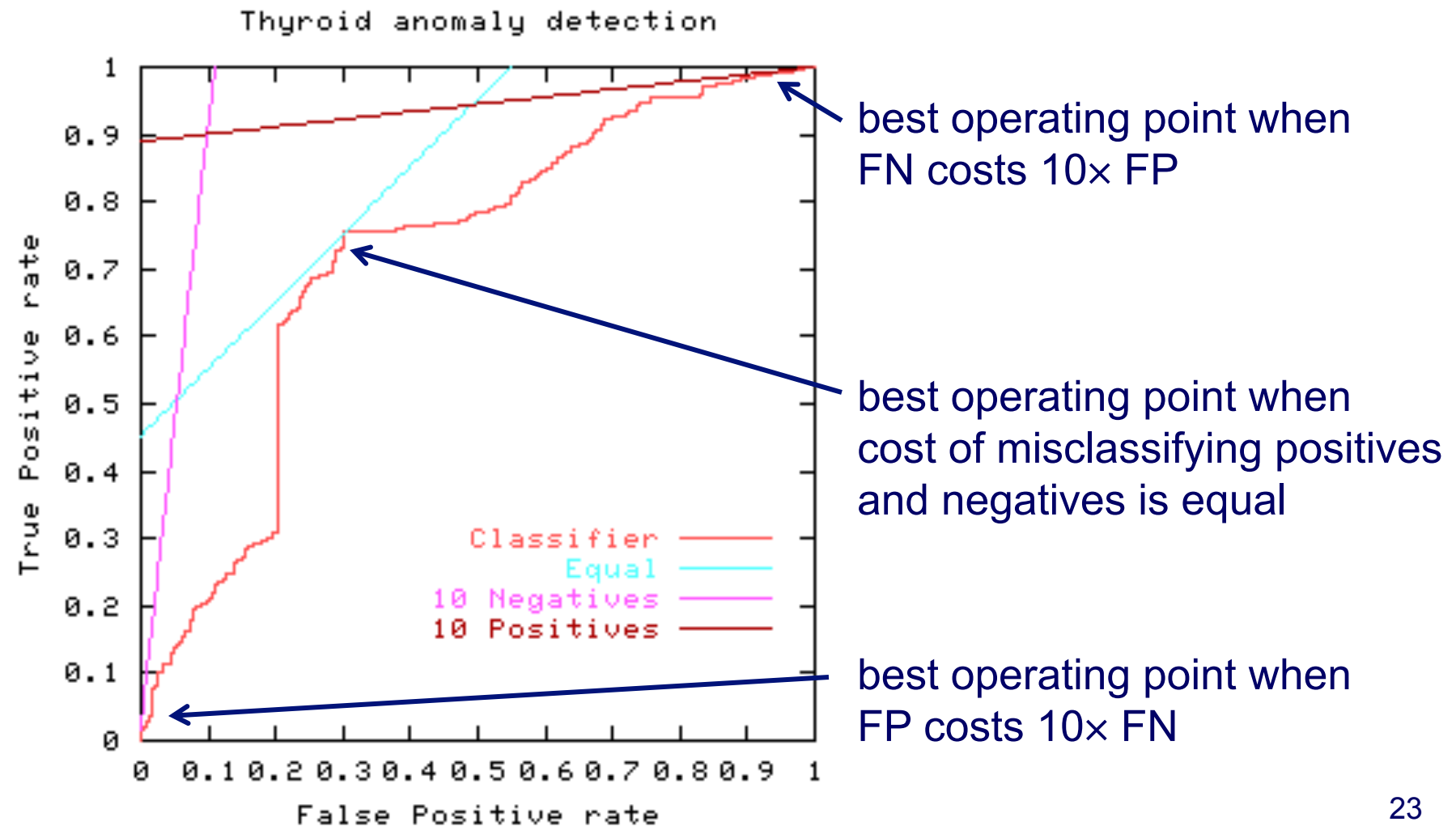


figure from Bockhorst et al., *Bioinformatics* 2003

ROC curves and misclassification costs

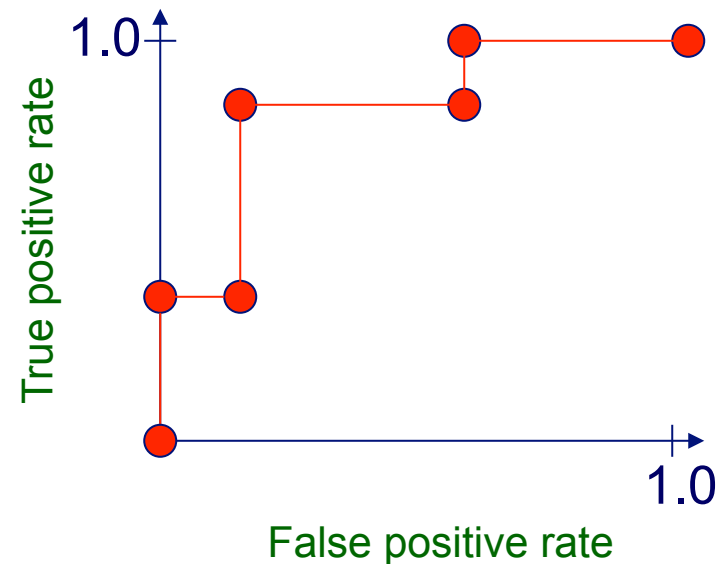


Algorithm for creating an ROC curve

1. sort test-set predictions according to confidence that each instance is positive
2. step through sorted list from high to low confidence
 - i. locate a *threshold* between instances with opposite classes (keeping instances with the same confidence value on the same side of threshold)
 - ii. compute TPR, FPR for instances above threshold
 - iii. output (FPR, TPR) coordinate

Plotting an ROC curve

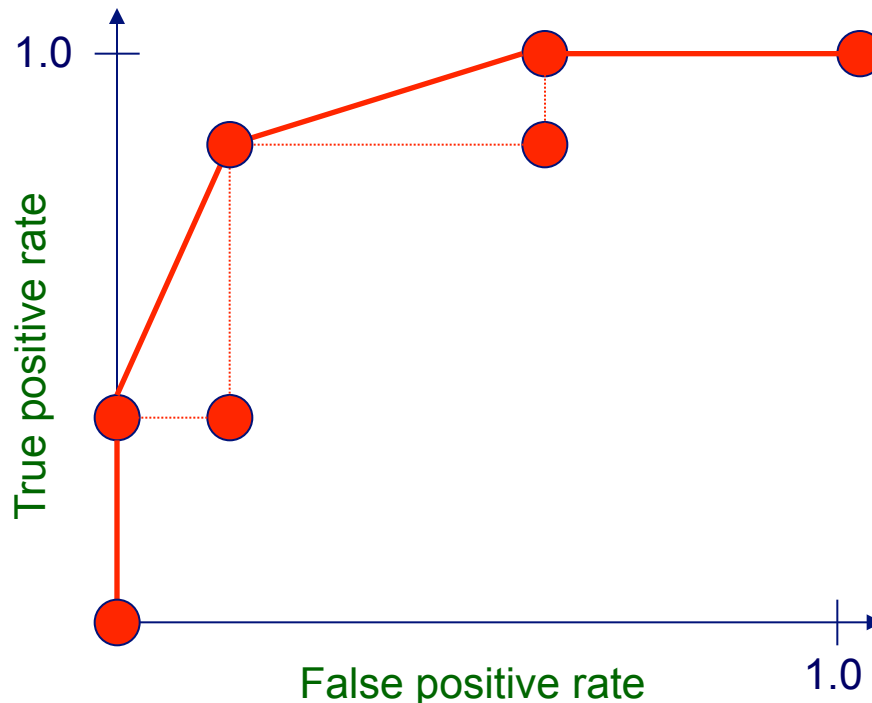
instance	confidence positive		correct class
Ex 9	.99		+
Ex 7	.98	TPR= 2/5, FPR= 0/5	+
Ex 1	.72	TPR= 2/5, FPR= 1/5	-
Ex 2	.70		+
Ex 6	.65	TPR= 4/5, FPR= 1/5	+
Ex 10	.51		-
Ex 3	.39	TPR= 4/5, FPR= 3/5	-
Ex 5	.24	TPR= 5/5, FPR= 3/5	+
Ex 4	.11		-
Ex 8	.01	TPR= 5/5, FPR= 5/5	-



Plotting an ROC curve

can interpolate between points to get *convex hull*

- convex hull: repeatedly, while possible, perform interpolations that skip one data point and discard any point that lies below a line
- interpolated points are achievable in theory: can flip weighted coin to choose between classifiers represented by plotted points



ROC curves

Does a low false-positive rate indicate that most positive predictions (i.e. predictions with confidence $>$ some threshold) are correct?

suppose our TPR is 0.9, and FPR is 0.01

fraction of instances that are positive	fraction of positive predictions that are correct
0.5	0.989
0.1	0.909
0.01	0.476
0.001	0.083

Other accuracy metrics

actual class

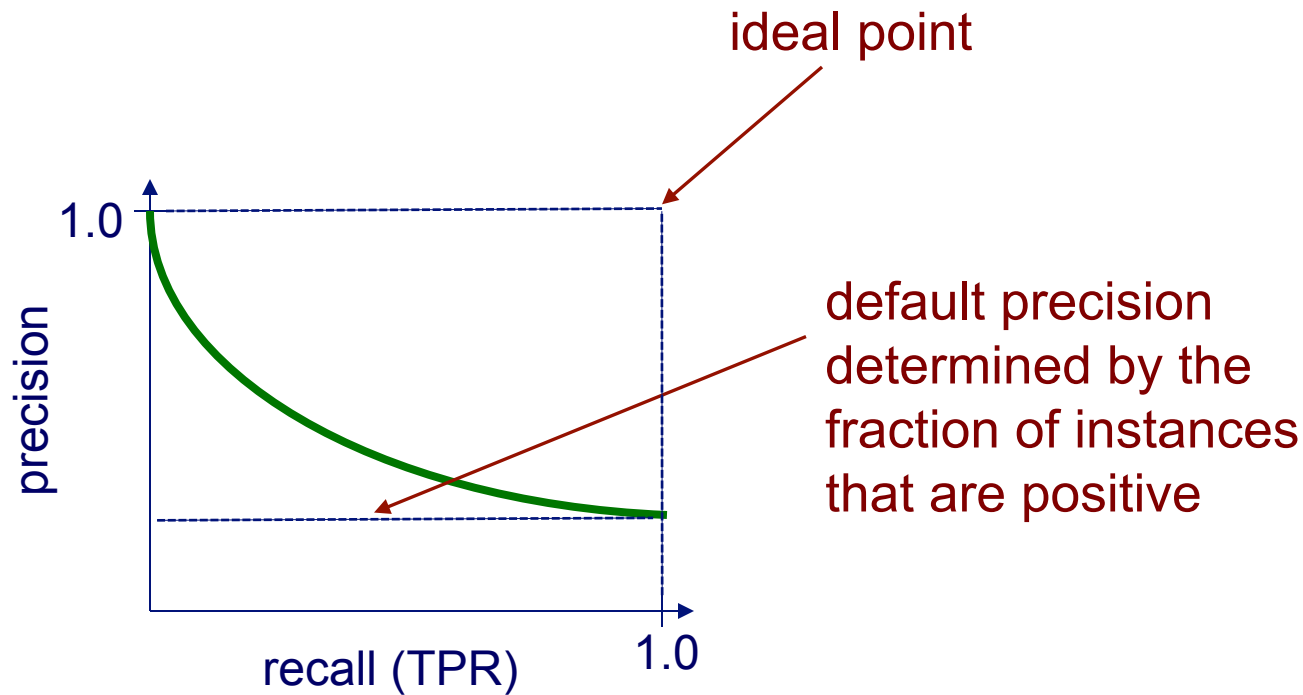
		actual class	
		positive	negative
predicted class	positive	true positives (TP)	false positives (FP)
	negative	false negatives (FN)	true negatives (TN)

$$\text{recall (TP rate)} = \frac{\text{TP}}{\text{actual pos}} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

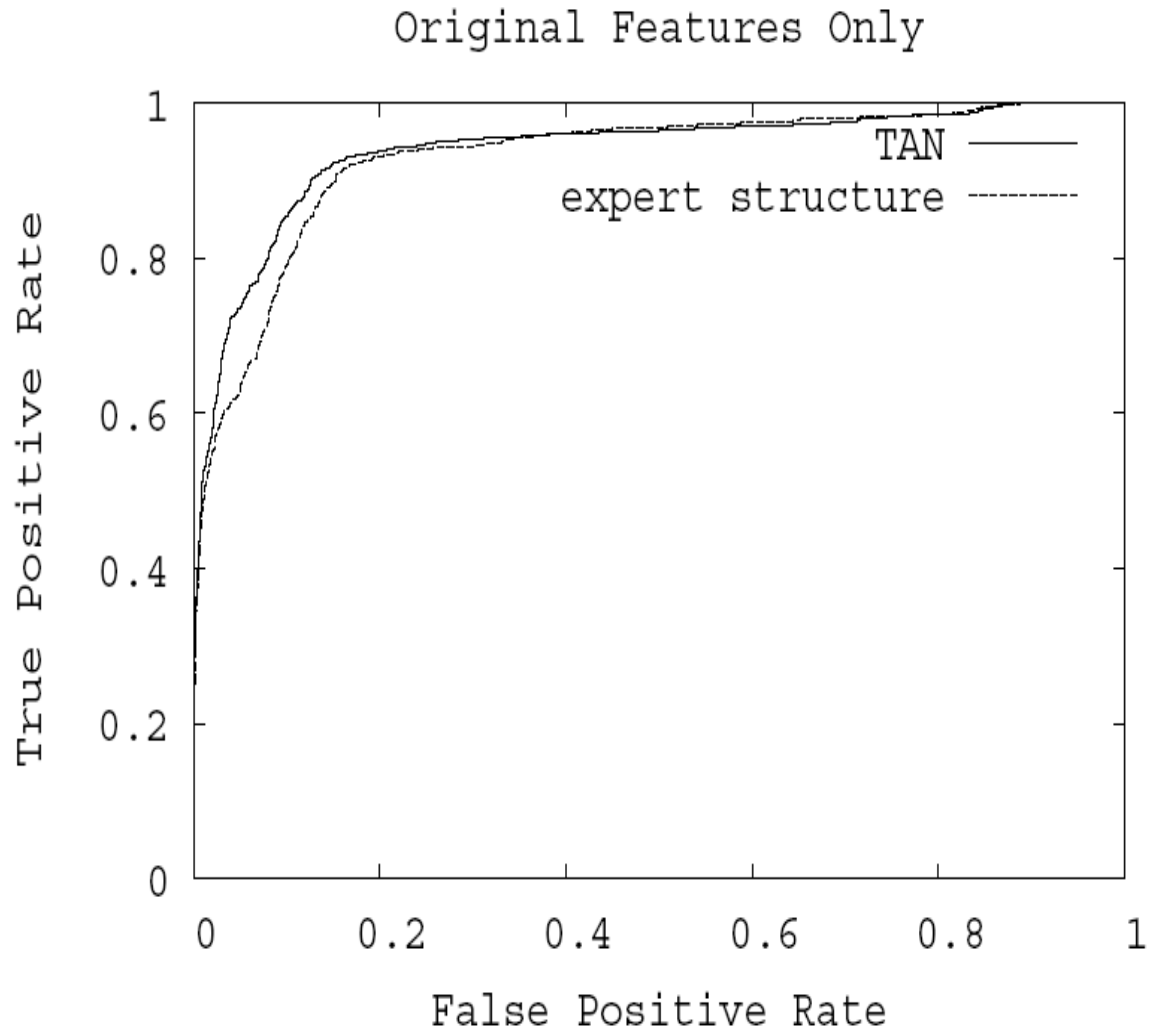
$$\text{precision} = \frac{\text{TP}}{\text{predicted pos}} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

Precision/recall curves

A *precision/recall curve* plots the precision vs. recall (TP-rate) as a threshold on the confidence of an instance being positive is varied

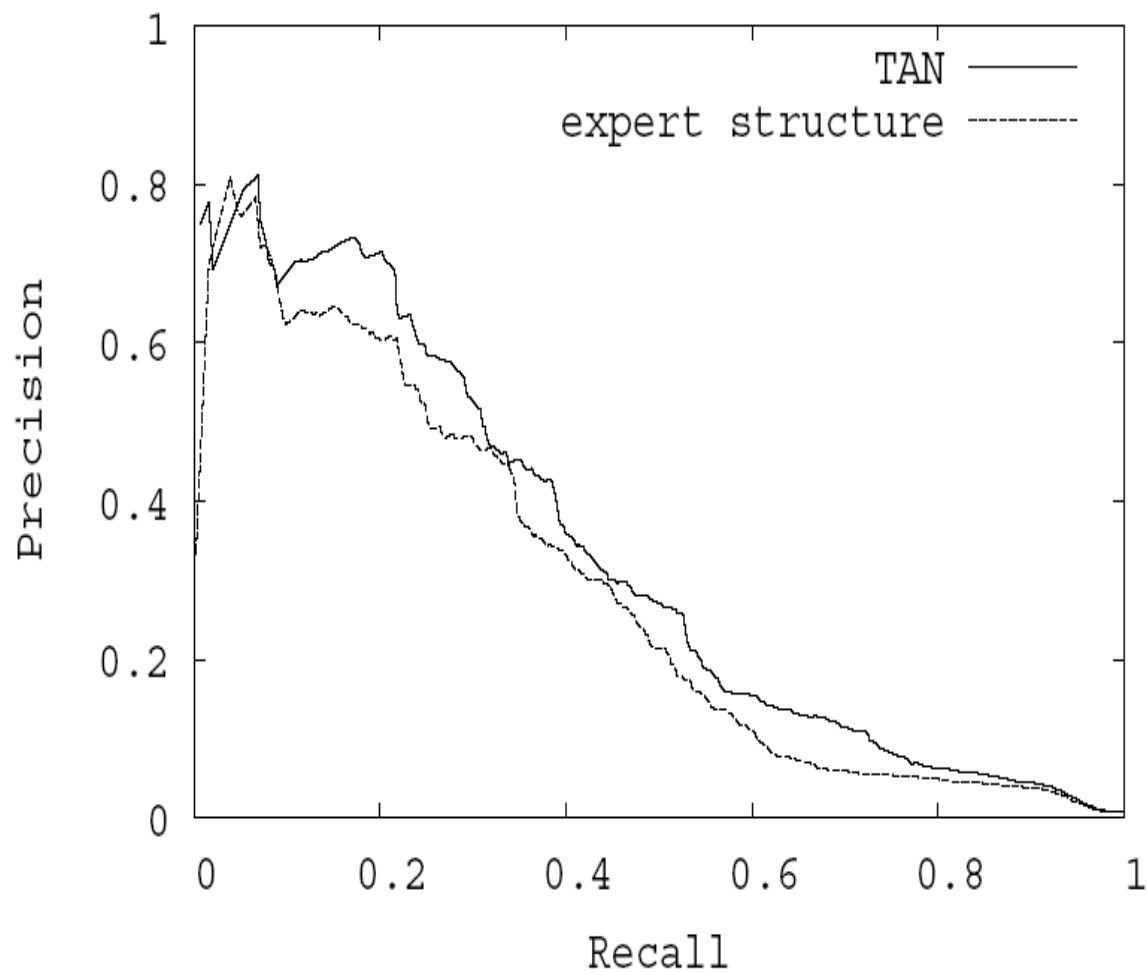


Mammography Example: ROC



Mammography Example: PR

Original Features Only



How do we get one ROC/PR curve when we do cross validation?

Approach 1

- make assumption that confidence values are comparable across folds
- pool predictions from all test sets
- plot the curve from the pooled predictions

Approach 2 (for ROC curves)

- plot individual curves for all test sets
- view each curve as a function
- plot the average curve for this set of functions

Comments on ROC and PR curves

both

- allow predictive performance to be assessed at various levels of confidence
- assume binary classification tasks
- sometimes summarized by calculating *area under the curve*

ROC curves

- insensitive to changes in class distribution (ROC curve does not change if the proportion of positive and negative instances in the test set are varied)
- can identify optimal classification thresholds for tasks with differential misclassification costs

precision/recall curves

- show the fraction of predictions that are false positives
- well suited for tasks with lots of negative instances

To Avoid Cross-Validation Pitfalls, Ask:

- 1. Is my held-aside test data really representative of going out to collect new data?
 - Even if your methodology is fine, someone may have collected features for positive examples differently than for negatives – should be *randomized*
 - Example: samples from cancer processed by different people or on different days than samples for normal controls

To Avoid Pitfalls, Ask:

- 2. Did I repeat my entire data processing procedure on every fold of cross-validation, using only the training data for that fold?
 - On each fold of cross-validation, did I ever access in any way the label of a test case?
 - Any preprocessing done over *entire data set* (feature selection, parameter tuning, threshold selection) must *not* use labels

To Avoid Pitfalls, Ask:

- 3. Have I modified my algorithm so many times, or tried so many approaches, on this same data set that *I* (the human) am overfitting it?
 - Have I continually modified my preprocessing or learning algorithm until I got some improvement on this data set?
 - If so, I really need to get some additional data now to at least test on

Confidence intervals on error

Given the observed error (accuracy) of a model over a limited sample of data, how well does this error characterize its accuracy over additional instances?

Suppose we have

- a learned model h
- a test set S containing n instances drawn independently of one another and independent of h
- $n \geq 30$
- h makes r errors over the n instances

our best estimate of the error of h is

$$\text{error}_S(h) = \frac{r}{n}$$

Confidence intervals on error

With approximately $N\%$ probability, the true error lies in the interval

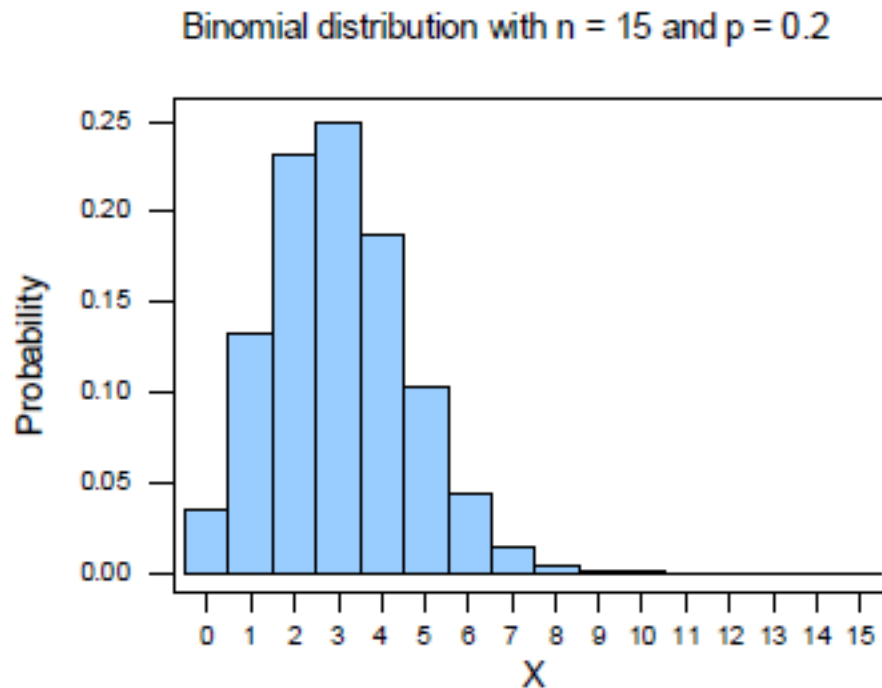
$$error_S(h) \pm z_N \sqrt{\frac{error_S(h)(1 - error_S(h))}{n}}$$

where z_N is a constant that depends on N (e.g. for 95% confidence, $z_N = 1.96$)

Confidence intervals on error

How did we get this?

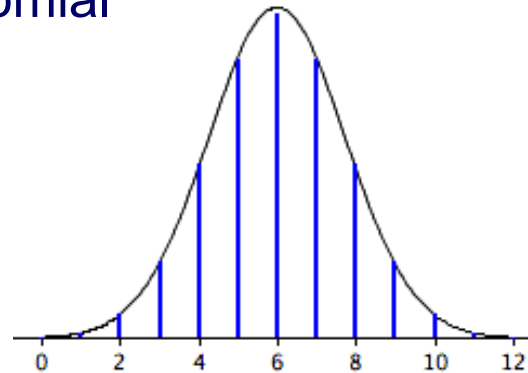
1. Our estimate of the error follows a binomial distribution given by n and p (the true error rate over the data distribution)



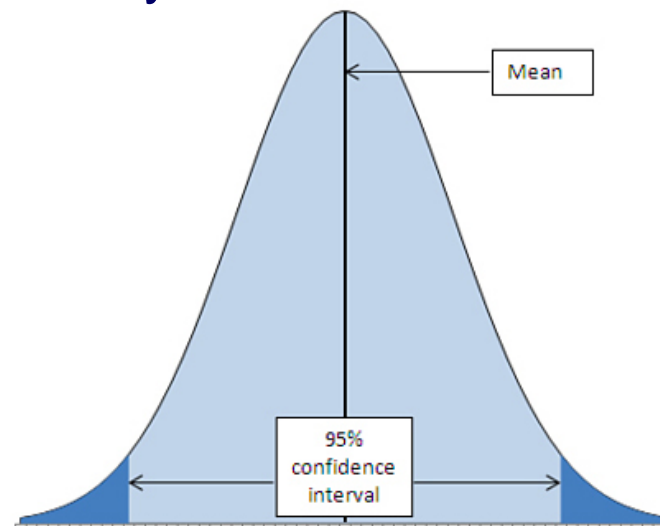
2. Simplest (and most common) way to determine a binomial confidence interval is to use the *normal approximation*

Confidence intervals on error

2. When $n \geq 30$, and p is not too extreme, the normal distribution is a good approximation to the binomial



3. We can determine the $N\%$ confidence interval by determining what bounds contain $N\%$ of the probability mass under the normal



Empirical Confidence Bounds

- Bootstrapping: Given n examples in data set, randomly, uniformly, independently (with replacement) draw n examples – bootstrap sample
- Repeat 1000 (or 10,000) times:
 - Draw bootstrap sample
 - Repeat entire cross-validation process
- Lower (upper) bound is result such that 2.5% of runs yield lower (higher)

Comparing learning systems

How can we determine if one learning system provides better performance than another

- for a particular task?
- across a set of tasks / data sets?

Motivating example

	<u>Accuracies on test sets</u>				
System 1:	80%	50	75	...	99
System 2:	79	49	74	...	98
δ :	+1	+1	+1	...	+1

- Mean accuracy for System 1 is better, but the standard deviations for the two clearly overlap
- Notice that System 1 is always better than System 2

Comparing systems using a paired t test

- consider δ 's as observed values of a set of i.i.d. random variables
- *null hypothesis*: the 2 learning systems have the same accuracy
- *alternative hypothesis*: one of the systems is more accurate than the other
- hypothesis test:
 - use paired t -test to determine probability p that mean of δ 's would arise from null hypothesis
 - if p is sufficiently small (typically < 0.05) then reject the null hypothesis

Comparing systems using a paired t test

1. calculate the sample mean

$$\bar{\delta} = \frac{1}{n} \sum_{i=1}^n \delta_i$$

2. calculate the t statistic

$$t = \frac{\bar{\delta}}{\sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n (\delta_i - \bar{\delta})^2}}$$

3. determine the corresponding p -value, by looking up t in a table of values for the Student's t -distribution with $n-1$ degrees of freedom

APPENDIX B STATISTICAL TABLES 691

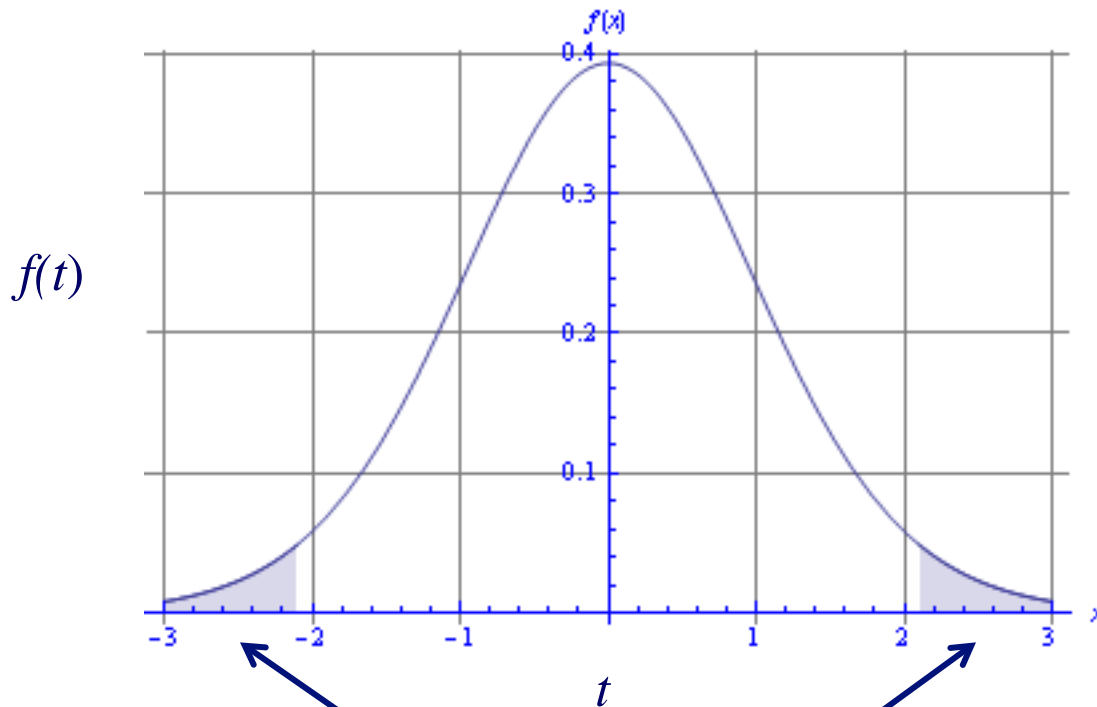
STATISTICAL TABLES

TABLE B.2 THE t DISTRIBUTION

Table entries are values of t corresponding to proportions in one tail or in two tails combined.

df	PROPORTION IN ONE TAIL				PROPORTION IN TWO TAILS COMBINED	
	0.25	0.10	0.05	0.025	0.01	0.005
1	1.000	3.078	6.314	12.506	31.821	63.657
2	0.816	1.886	2.920	4.303	6.965	9.925
3	0.766	1.638	2.353	3.182	4.541	5.841
4	0.741	1.533	2.132	2.776	3.747	4.604
5	0.727	1.476	2.015	2.571	3.363	4.032
6	0.718	1.440	1.963	2.447	3.143	3.707
7	0.711	1.415	1.945	2.365	2.998	3.499
8	0.706	1.397	1.930	2.306	2.896	3.355
9	0.703	1.383	1.913	2.262	2.821	3.250
10	0.700	1.372	1.902	2.228	2.764	3.169
11	0.697	1.363	1.896	2.201	2.718	3.106
12	0.695	1.356	1.892	2.179	2.681	3.055
13	0.694	1.350	1.889	2.160	2.650	3.012
14	0.693	1.346	1.887	2.145	2.624	2.977
15	0.691	1.341	1.885	2.131	2.602	2.947
16	0.689	1.337	1.884	2.120	2.583	2.921
17	0.689	1.333	1.883	2.110	2.567	2.898
18	0.688	1.330	1.882	2.101	2.552	2.878
19	0.688	1.328	1.882	2.093	2.539	2.861
20	0.687	1.325	1.881	2.086	2.528	2.845
21	0.686	1.323	1.881	2.080	2.518	2.831
22	0.686	1.321	1.881	2.074	2.508	2.819
23	0.685	1.319	1.881	2.069	2.500	2.807
24	0.685	1.318	1.881	2.064	2.492	2.797
25	0.684	1.316	1.880	2.060	2.485	2.787
26	0.684	1.315	1.880	2.056	2.479	2.779
27	0.684	1.314	1.880	2.052	2.473	2.771
28	0.683	1.313	1.880	2.048	2.467	2.763
29	0.683	1.311	1.879	2.042	2.462	2.756
30	0.683	1.310	1.879	2.042	2.457	2.750
40	0.681	1.303	1.878	2.021	2.423	2.704
60	0.679	1.296	1.877	2.000	2.390	2.660
120	0.677	1.289	1.876	1.980	2.358	2.617
∞	0.675	1.282	1.875	1.960	2.326	2.576

Comparing systems using a paired t test



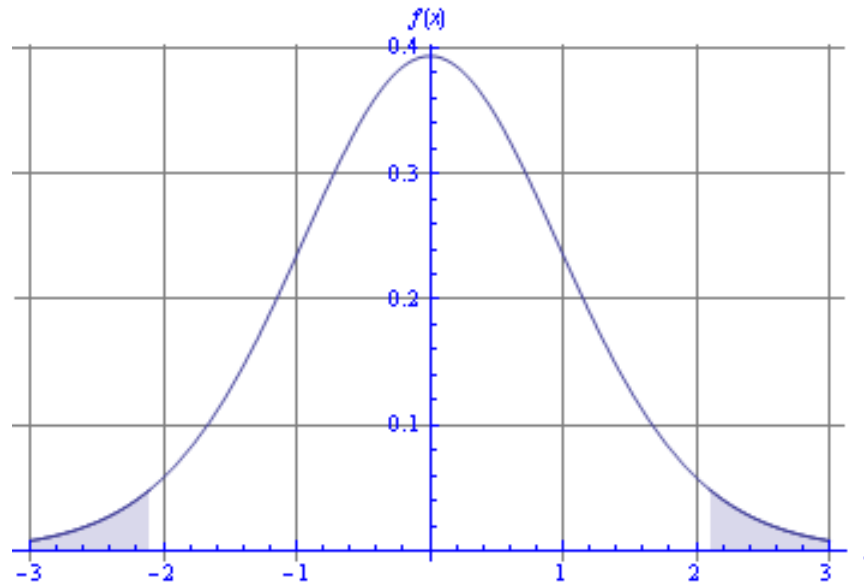
The null distribution of our t statistic looks like this

The p -value indicates how far out in a tail our t statistic is

If the p -value is sufficiently small, we reject the null hypothesis, since it is unlikely we'd get such a t by chance

for a two-tailed test, the p -value represents the probability mass in these two regions

Why do we use a two-tailed test?



- a two-tailed test asks the question: is the accuracy of the two systems different
- a one-tailed test asks the question: is system A better than system B
- a priori, we don't know which learning system will be more accurate (if there is a difference) – we want to allow that either one might be

Sign Test

- If less than 300 examples, we won't have 30 test examples per fold
- Prefer leave-one-out cross-validation
- Count “wins” for Algorithm A and B over the N test examples on which they disagree
- Let M be the larger of these counts
- What is probability under $b(N, 0.5)$ that either A or B would win *at least* M times

Scatter plots for pairwise method comparison

We can compare the performance of two methods *A* and *B* by plotting (*A performance*, *B performance*) across numerous data sets

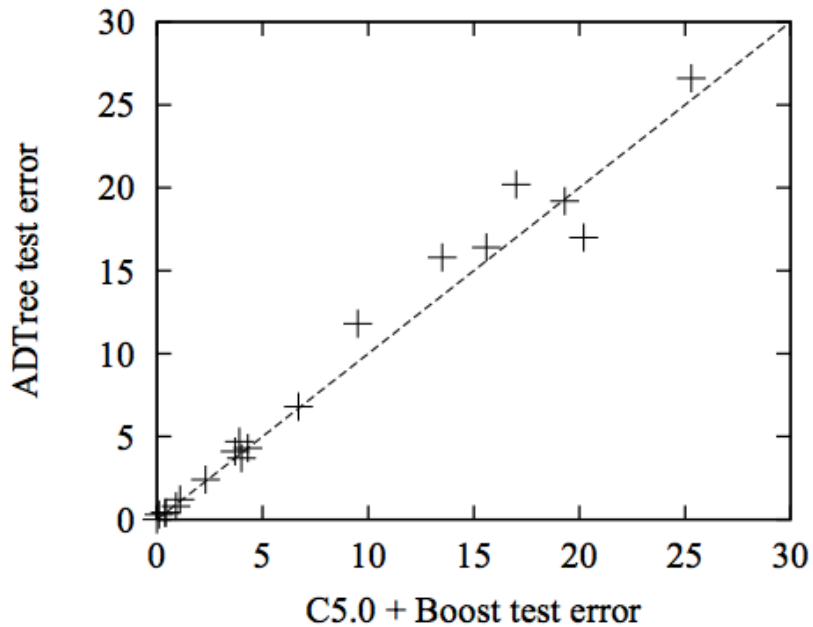


figure from Freund & Mason, *ICML* 1999

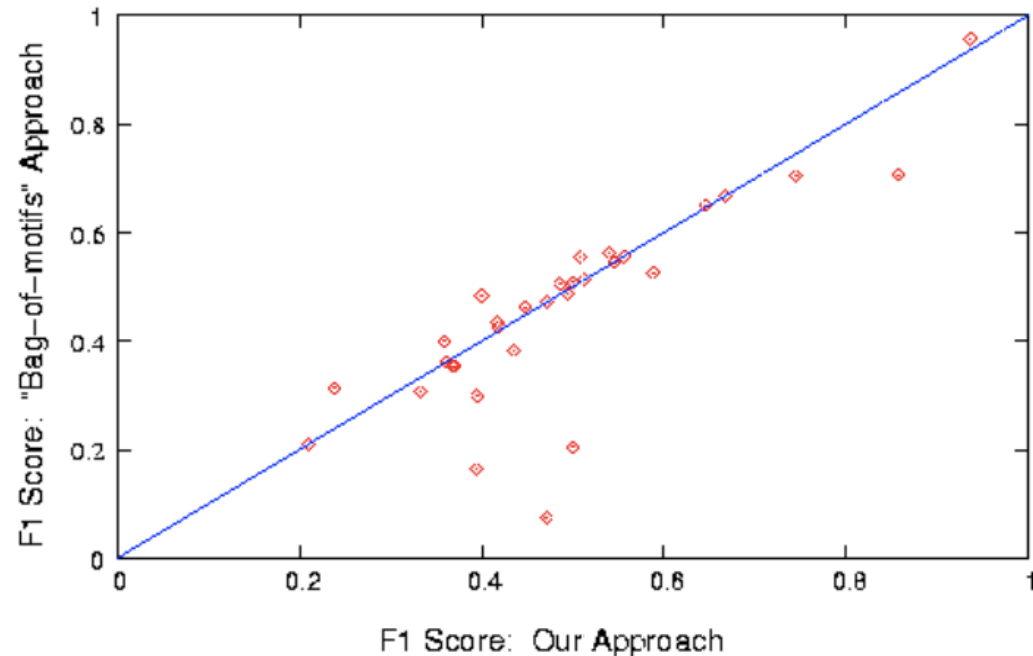


figure from Noto & Craven, *BMC Bioinformatics* 2006

Lesion studies

We can gain insight into what contributes to a learning system's performance by removing (lesioning) components of it

The ROC curves here show how performance is affected when various feature types are removed from the learning representation

