

# Evaluating Machine-Learning Methods (Part 2)

Mark Craven and David Page  
Computer Sciences 760  
Spring 2019

# Goals for the lecture

you should understand the following concepts

- confidence intervals for error
- pairwise  $t$ -tests for comparing learning systems
- scatter plots for comparing learning systems
- lesion studies
- model selection
- validation (tuning) sets
- internal cross validation

# 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

$$error_S(h) = \frac{r}{n}$$

# Confidence intervals on error

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

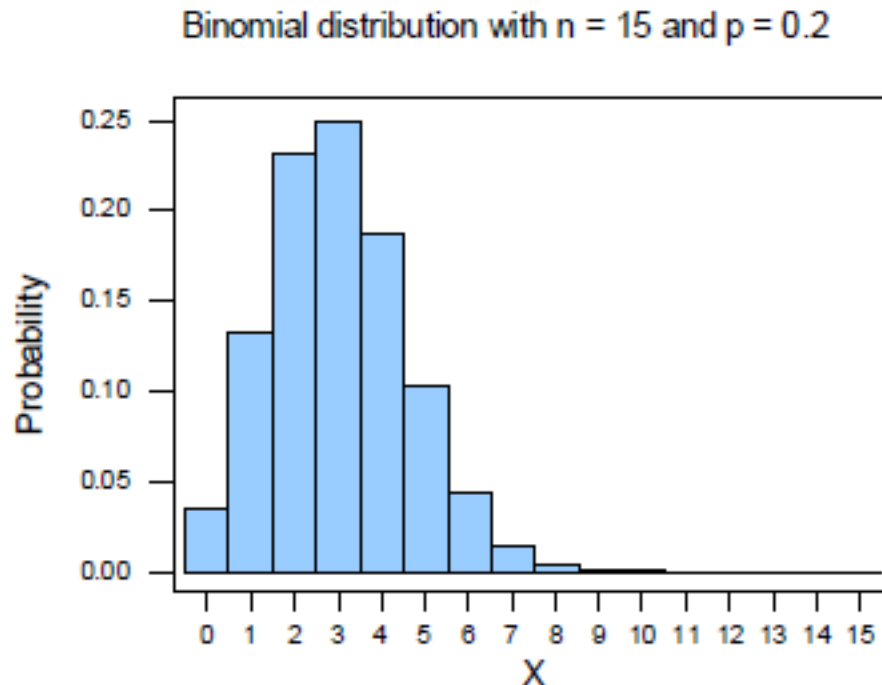
$$error_s(h) \pm z_C \sqrt{\frac{error_s(h)(1 - error_s(h))}{n}}$$

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

# Confidence intervals on error

How did we get this?

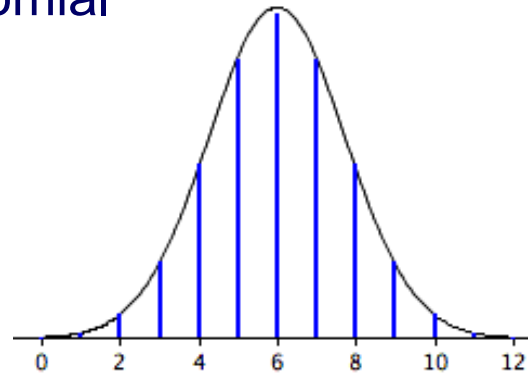
1. Our estimate of the error follows a binomial distribution given by  $n$  and  $p$  (the true error rate over the data distribution)



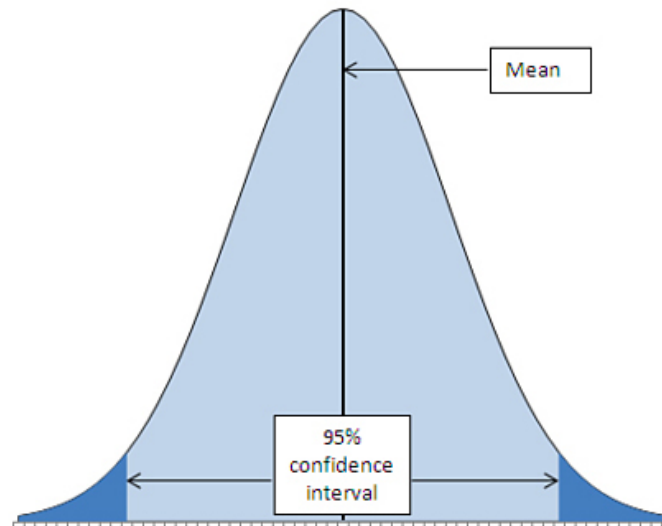
2. Most common way to determine a binomial confidence interval is to use the *normal approximation* (although can calculate exact intervals if  $n$  is not too large)

# 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  $C\%$  confidence interval by determining what bounds contain  $C\%$  of the probability mass under the normal



# Alternative approach: confidence intervals using bootstrapping

- *bootstrap sample*: given  $n$  examples in data set, randomly, uniformly, independently draw  $n$  examples with replacement
- repeat 1000 (or 10,000) times:
  - draw bootstrap sample
  - measure error on bootstrap sample
  - for 95% confidence interval, lower (upper) bound is set such that 2.5% of runs yield lower (higher) error

# 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 A:	80%	50	75	...	99
System B:	79	49	74	...	98
$\delta$ :	+1	+1	+1	...	+1

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

# Comparing systems using a paired $t$ test

- consider  $\delta$ '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  $\delta$ '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

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

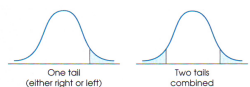
3. determine the corresponding  $p$ -value, by looking up  $t$  in a table of values for the Student's  $t$ -distribution with  $n-1$  degrees of freedom

APPENDIX B STATISTICAL TABLES 891

STATISTICAL TABLES

TABLE B.2 THE  $t$  DISTRIBUTION

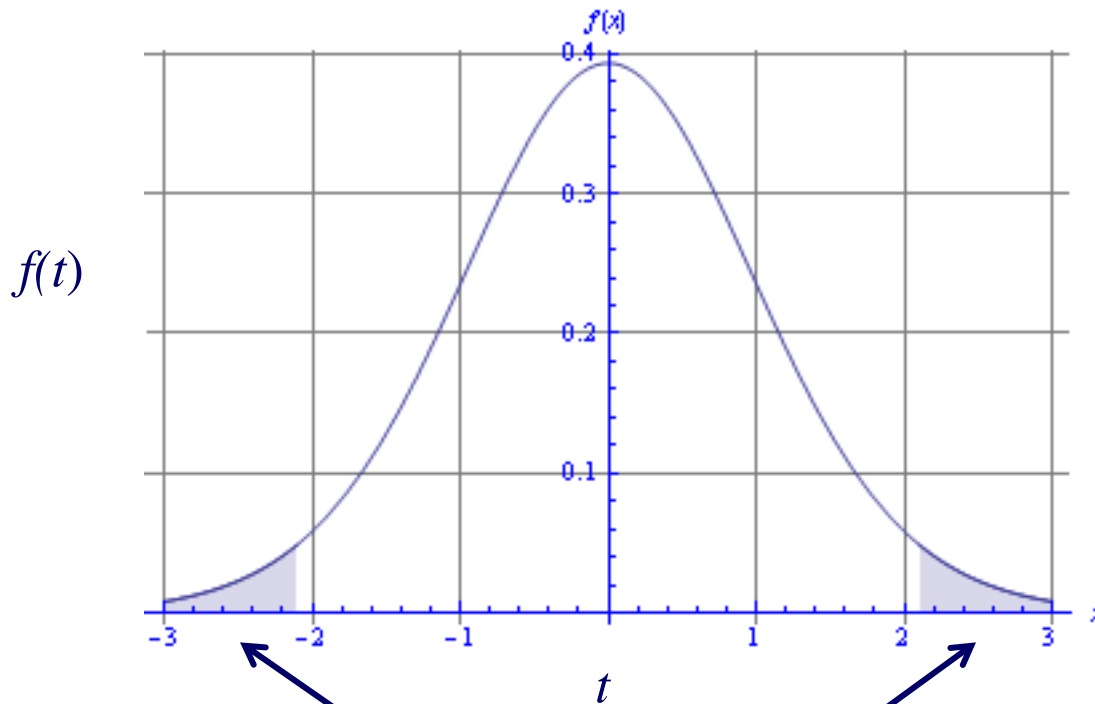
Table entries are values of  $t$  corresponding to proportions in one tail or in two tails combined.



One tail (either right or left)      Two tails combined

df	PROPORTION IN ONE TAIL				
	0.25	0.10	0.05	0.025	0.01
1	1.000	3.078	6.314	12.706	31.821
2	0.816	1.886	2.920	4.303	6.965
3	0.765	1.638	2.353	3.182	5.841
4	0.741	1.533	2.132	2.776	5.408
5	0.727	1.476	2.015	2.571	5.101
6	0.718	1.440	1.963	2.447	4.905
7	0.711	1.415	1.943	2.365	4.781
8	0.706	1.397	1.928	2.306	4.698
9	0.703	1.385	1.915	2.262	4.633
10	0.700	1.372	1.902	2.228	4.576
11	0.697	1.363	1.891	2.201	4.534
12	0.695	1.356	1.882	2.179	4.494
13	0.694	1.350	1.875	2.161	4.462
14	0.693	1.346	1.869	2.146	4.437
15	0.691	1.341	1.864	2.131	4.418
16	0.690	1.337	1.860	2.120	4.401
17	0.689	1.333	1.856	2.110	4.386
18	0.688	1.330	1.853	2.101	4.373
19	0.688	1.328	1.851	2.093	4.361
20	0.687	1.325	1.849	2.086	4.351
21	0.687	1.323	1.847	2.080	4.342
22	0.686	1.321	1.845	2.074	4.334
23	0.686	1.319	1.844	2.069	4.327
24	0.685	1.318	1.843	2.064	4.321
25	0.685	1.316	1.842	2.060	4.316
26	0.684	1.315	1.841	2.056	4.312
27	0.684	1.314	1.840	2.052	4.308
28	0.683	1.313	1.839	2.048	4.304
29	0.683	1.311	1.838	2.045	4.301
30	0.683	1.310	1.837	2.042	4.298
40	0.681	1.303	1.833	2.021	4.255
60	0.679	1.296	1.827	2.000	4.200
120	0.677	1.289	1.824	1.980	4.151
$\infty$	0.675	1.282	1.820	1.960	4.126

# Comparing systems using a paired $t$ test



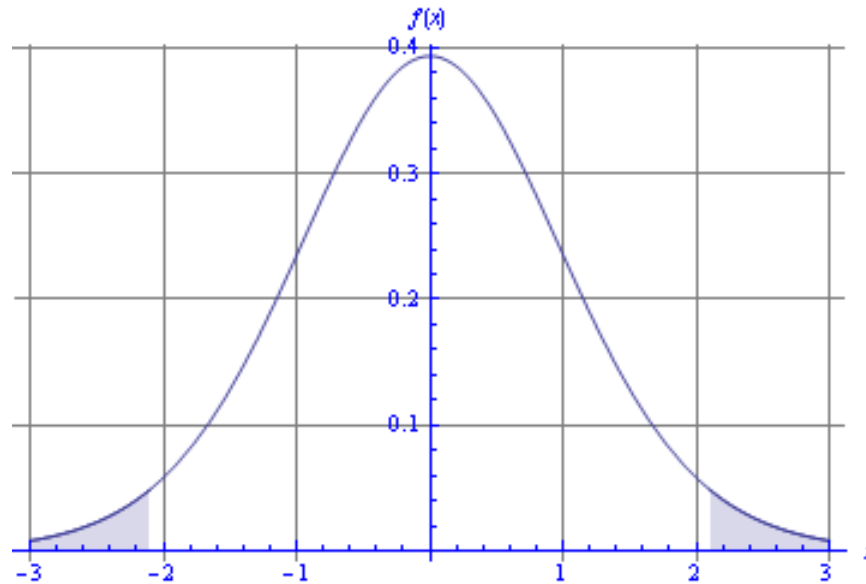
The null distribution of our  $t$  statistic looks like this

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

# Comments on hypothesis testing to compare learning systems

- the paired  $t$ -test can be used to compare two learning systems
- other tests (e.g. McNemar's  $\chi^2$  test) can be used to compare two learned models
- a statistically significant difference is not necessarily a large-magnitude difference

# 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

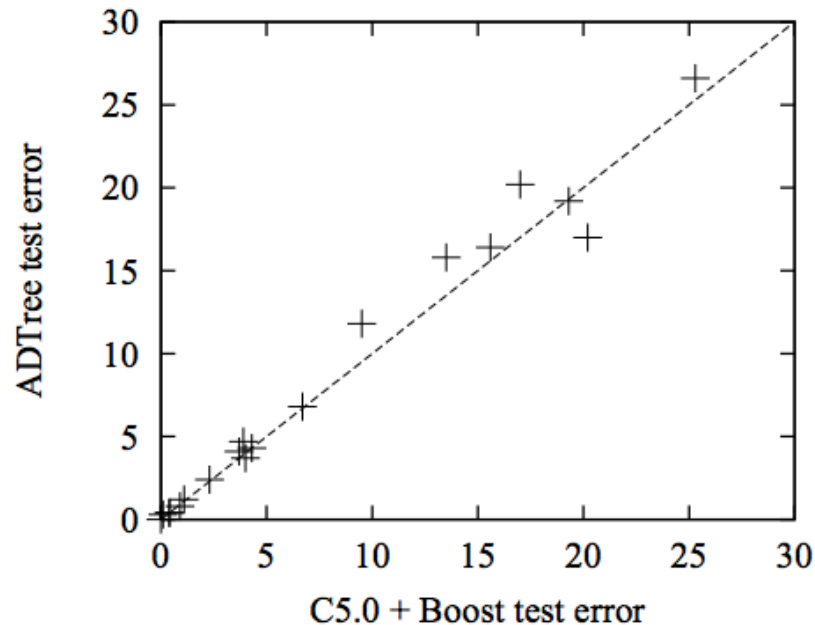


figure from Freund & Mason, *ICML* 1999

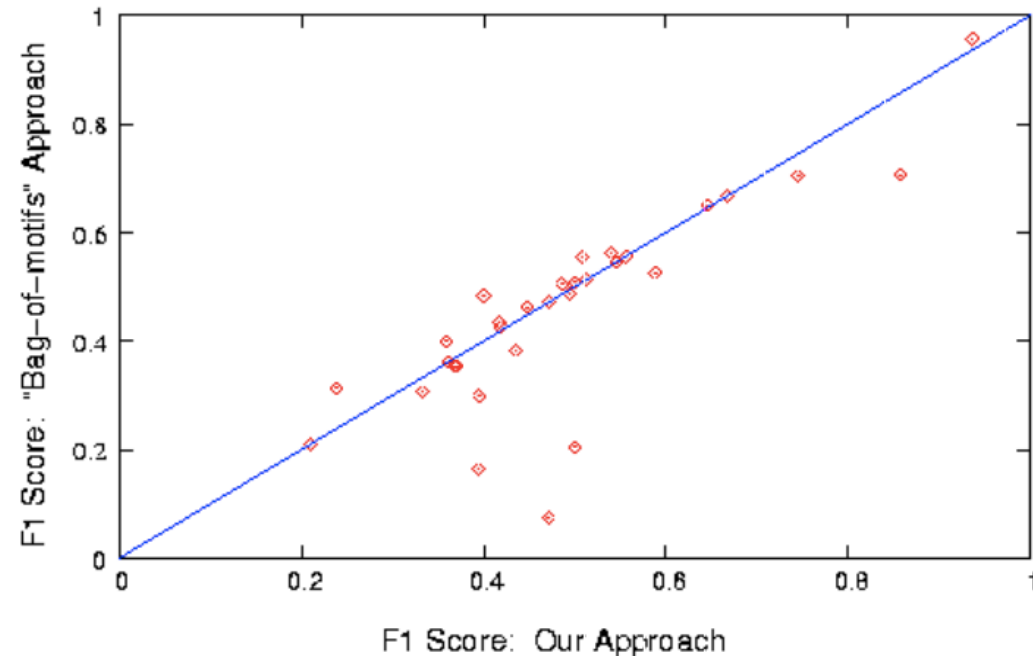
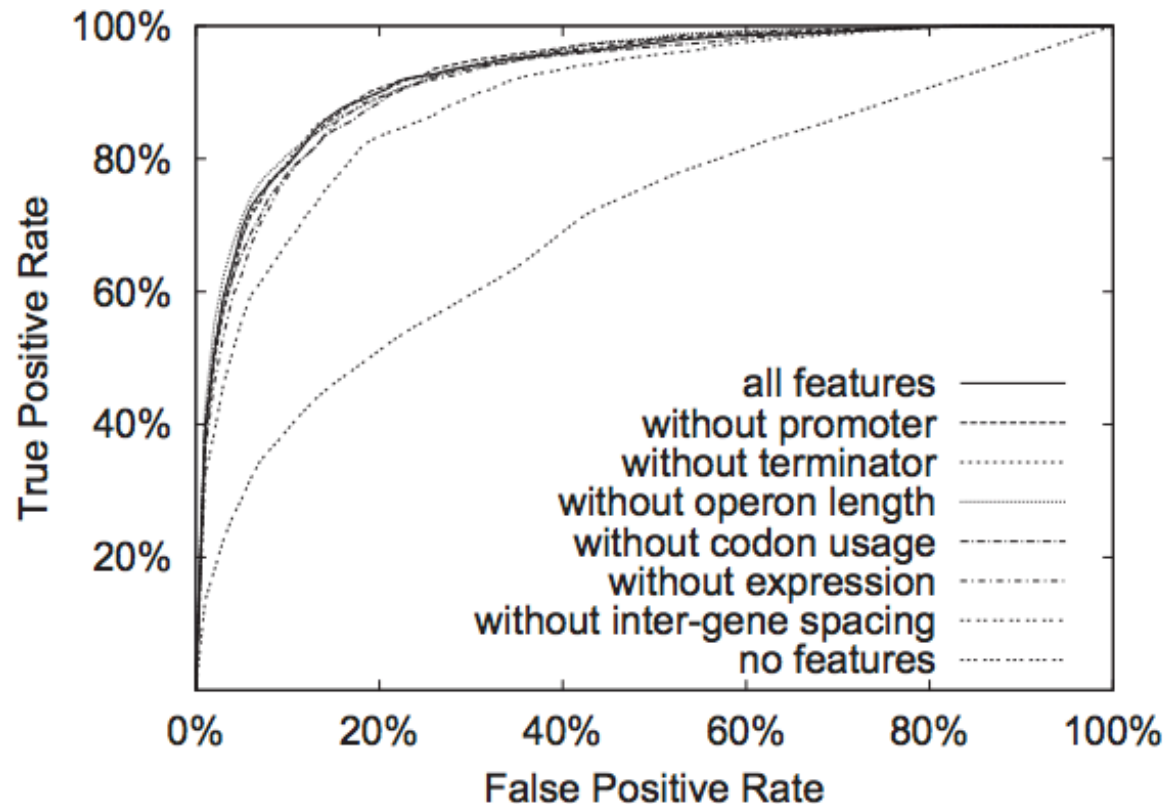


figure from Noto & Craven, *BMC Bioinformatics* 2006

# 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





# To avoid 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 instance?
  - Any preprocessing done over entire data set (feature selection, parameter tuning, threshold selection) must not use labels

# To avoid pitfalls, ask

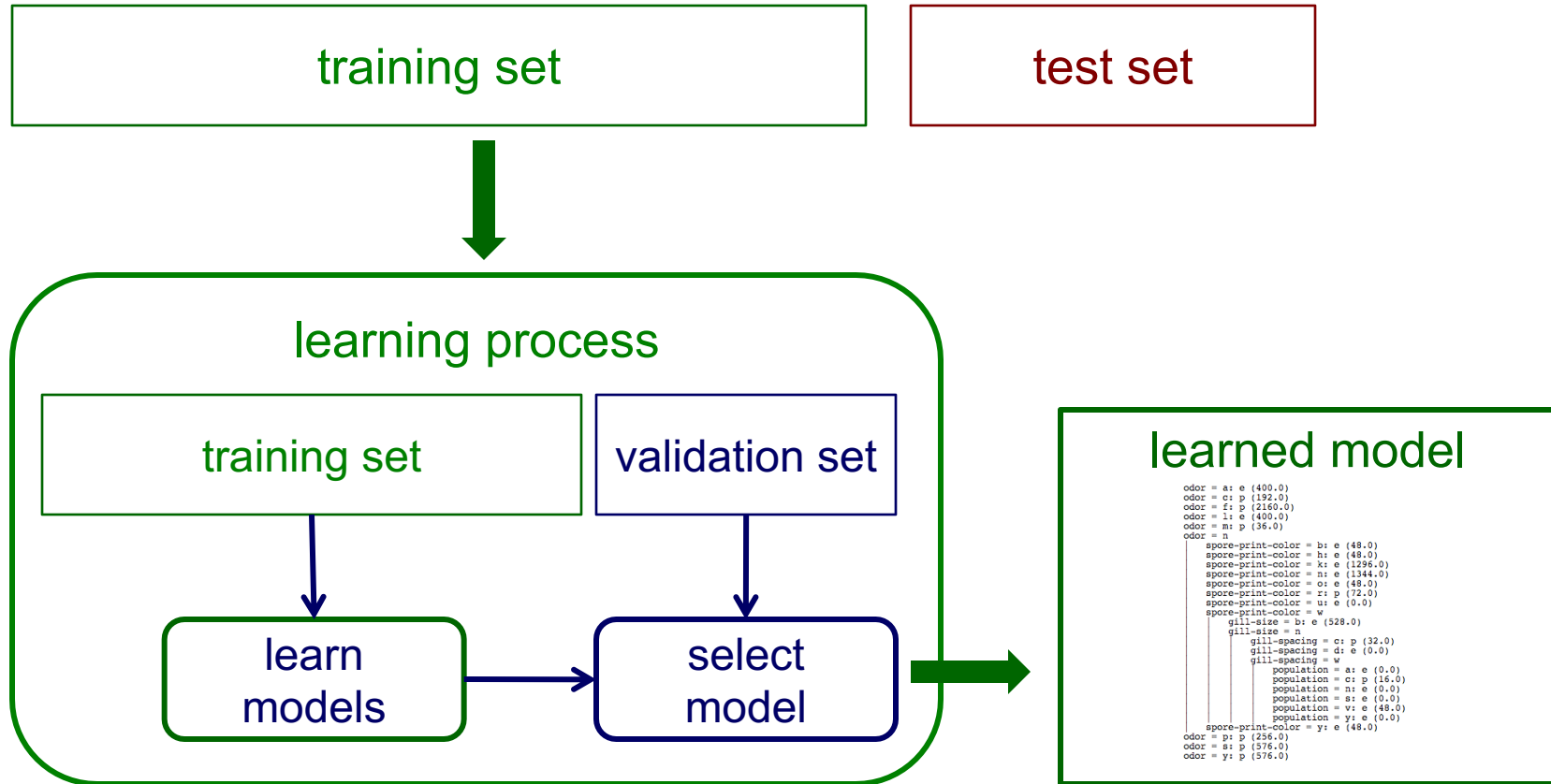
3. Have I modified my algorithm so many times, or tried so many approaches, on this same data set that I (the human) am overfitting it?
  - Have I continually modified my preprocessing or learning algorithm until I got some improvement on this data set?
  - If so, I really need to get some additional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

# Model selection

- *model selection* is the task of selecting a model from a set of candidate models
  - selecting among decision trees with various levels of pruning
  - selecting  $k$  in  $k$ -NN
  - etc.
- one approach to model selection is to use a tuning set or *internal* cross validation

# Validation (tuning) sets revisited

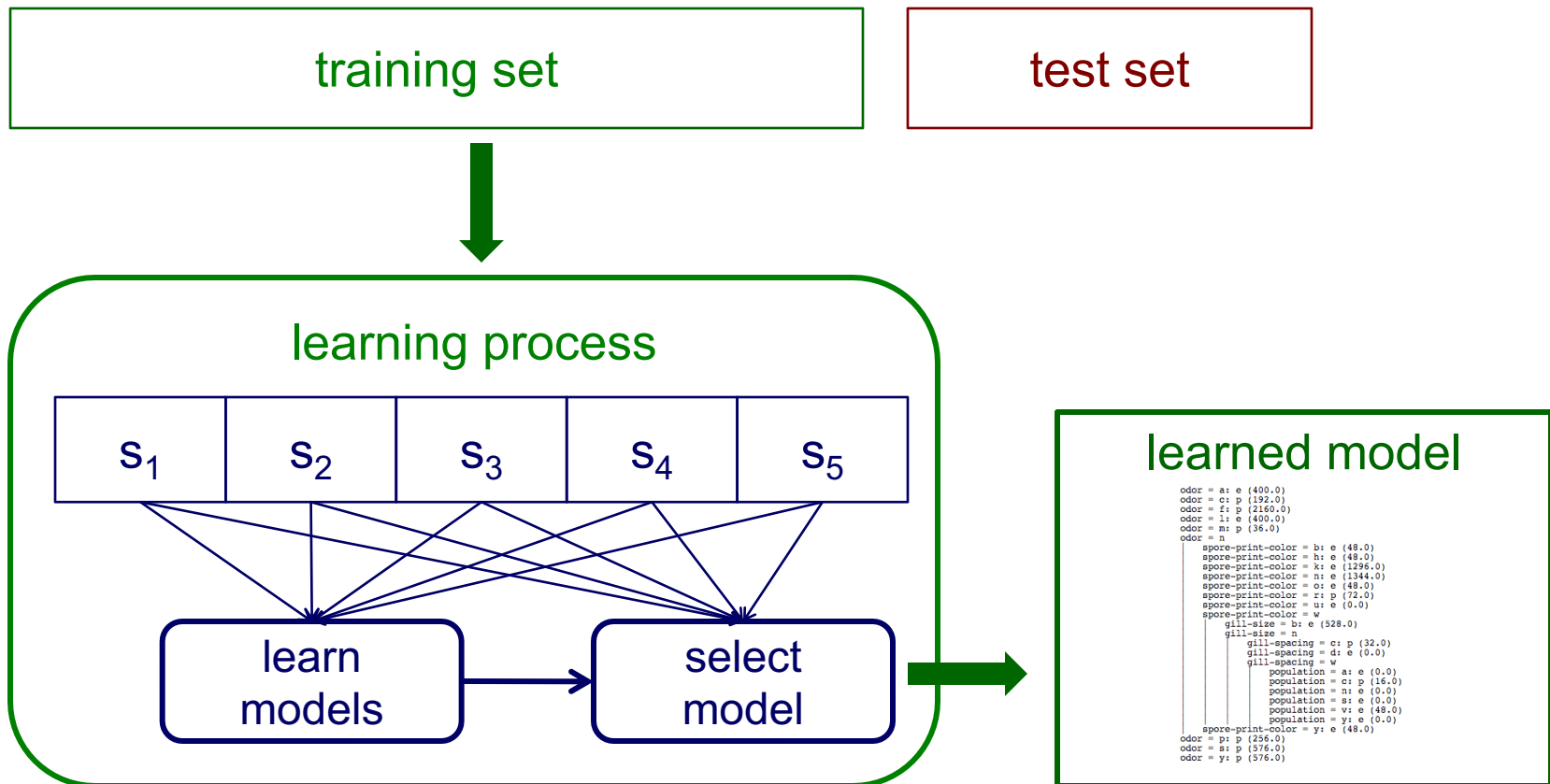
Suppose we want estimates of accuracy during the learning process (e.g. to choose the best level of decision-tree pruning)?



Partition training data into separate training/validation sets

# 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)