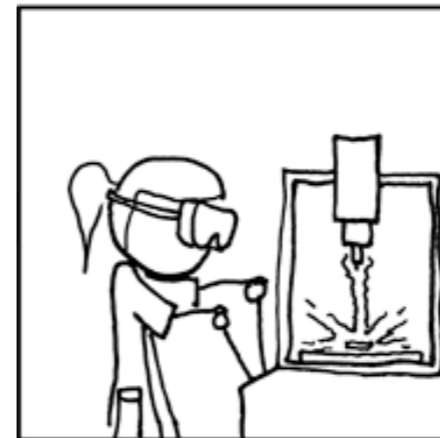
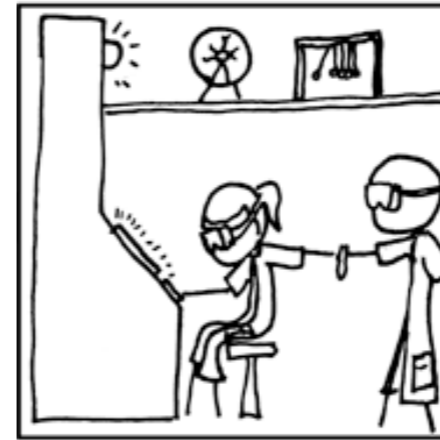


# Evaluation Basics

MOVIE SCIENCE MONTAGE



ACTUAL SCIENCE MONTAGE



# Reminders (Oct 17, 2019)

- Outdated slides on schedule: I leave older slides up, so its easy to see what could be coming
  - Don't necessarily trust slides for later days, as they might be from last year
- Class mini-project updates:
  - Can work in pairs or threes (if you want)
  - You can use packages, such as scikit, for the project (though not for assignments)
  - Your goal is to formalize your problem; you can ask me about it somewhat, but this is a part of the project

# Stepping back: What is machine learning?

- Central goal: obtain models that give good generalization performance (i.e., predict  $Y$  accurately from  $X$ )
  - Sometimes just want a good model for one problem
  - Typically want to identify approaches that work well across problems
  - Central to this theme: bias-variance trade-off
- Procedure:
  - Formalize the problem (e.g., as a maximum likelihood problem)
  - Propose a solution methodology (e.g., good optimization algorithm)
  - **Evaluate performance (either theoretically or empirically)**

# Crash course in evaluation

- Running a scientific experiment to compare machine learning algorithms necessary to draw conclusions
- Needs to be meticulous, even if sometimes tedious or need to re-run experiments
- Requires an experimental design and statistical significance tests
- See these slides: <http://pages.cs.wisc.edu/~dpage/cs760/evaluating.pdf>

# Hypotheses to test

- Algorithm A is better than Algorithm B
  - this is almost impossible to say
- Algorithm A is better than Algorithm B on this dataset
  - for all possible hyper-parameter settings?!
- With specific settings of hyperparameters, Algorithm A is better than Algorithm B on this dataset
  - but now is Algorithm B better with different hyperparameters?
  - want to ask a stronger question
- Specifying a testable hypothesis is part of the difficulty
- What do you mean by “better”? Why this dataset?

# Selecting the definition of better

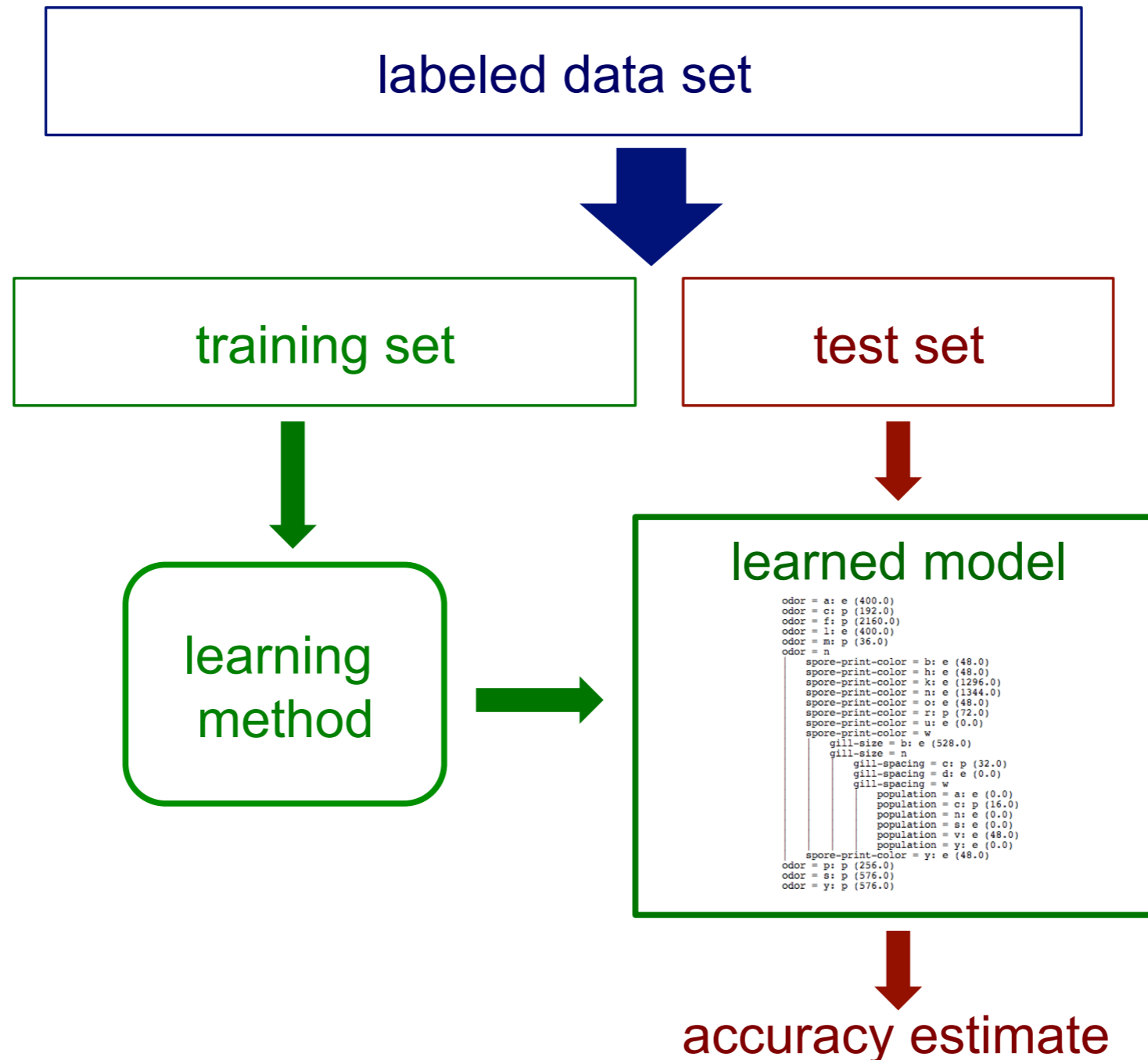
- Best classification accuracy
- Best classification accuracy on “most important” samples
- Robust classification accuracy that ignores outliers
- ROC curve and AUC
- Training time and space complexity
- Testing time and space complexity
- Ease of implementing the algorithm and interpretability
- We’ll talk about measures later; for now assume you’ve defined “better”

# Select dataset(s)

- Either you have some domain that you care about, and associated data
  - your goal is to predict well on future data, and generally better understand the data for your domain with whatever algorithm
- Or you care about exploring the properties of the algorithm and might find a diverse set of datasets
  - this includes potentially generating synthetic data for which you understand the properties
  - e.g. generate iid data from a Gaussian distribution
  - e.g., generate data from a simple neural network
- Again for now assume you've selected the data

# Evaluating on the data

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





# Our actual goal(s)

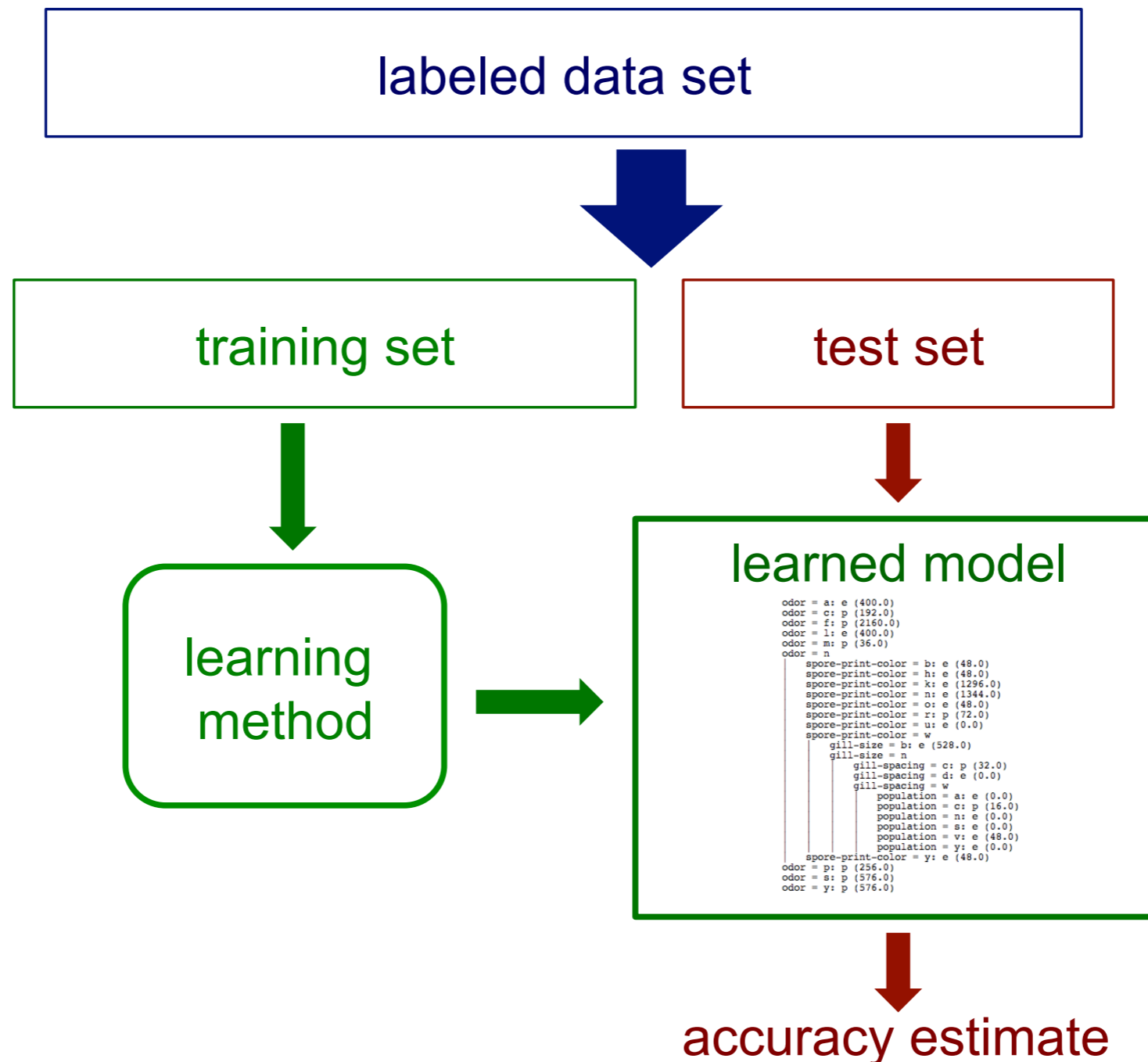
- **Goal 1:** train on our data, and then deploy in the world
- Before deployment, we'd like to have some idea of how our model might behave in the world
- We design experiments to mimic this deployment scenario
  - data is not actually separated into training/test — its just the data you have, and you decided to split it up that way

# Our actual goal(s)

- **Goal 2:** understand different properties of algorithms, by running them in different settings
- Here, you might have a dataset of 1 million samples, but you just want to use it to test how your algorithms might behave when learning on datasets of 10k samples
- Experimental design might be different compared to Goal 1
- You need to know your goal to design your experiment

# Back to evaluating on the data

- What do you think the goal is here? What is the question?

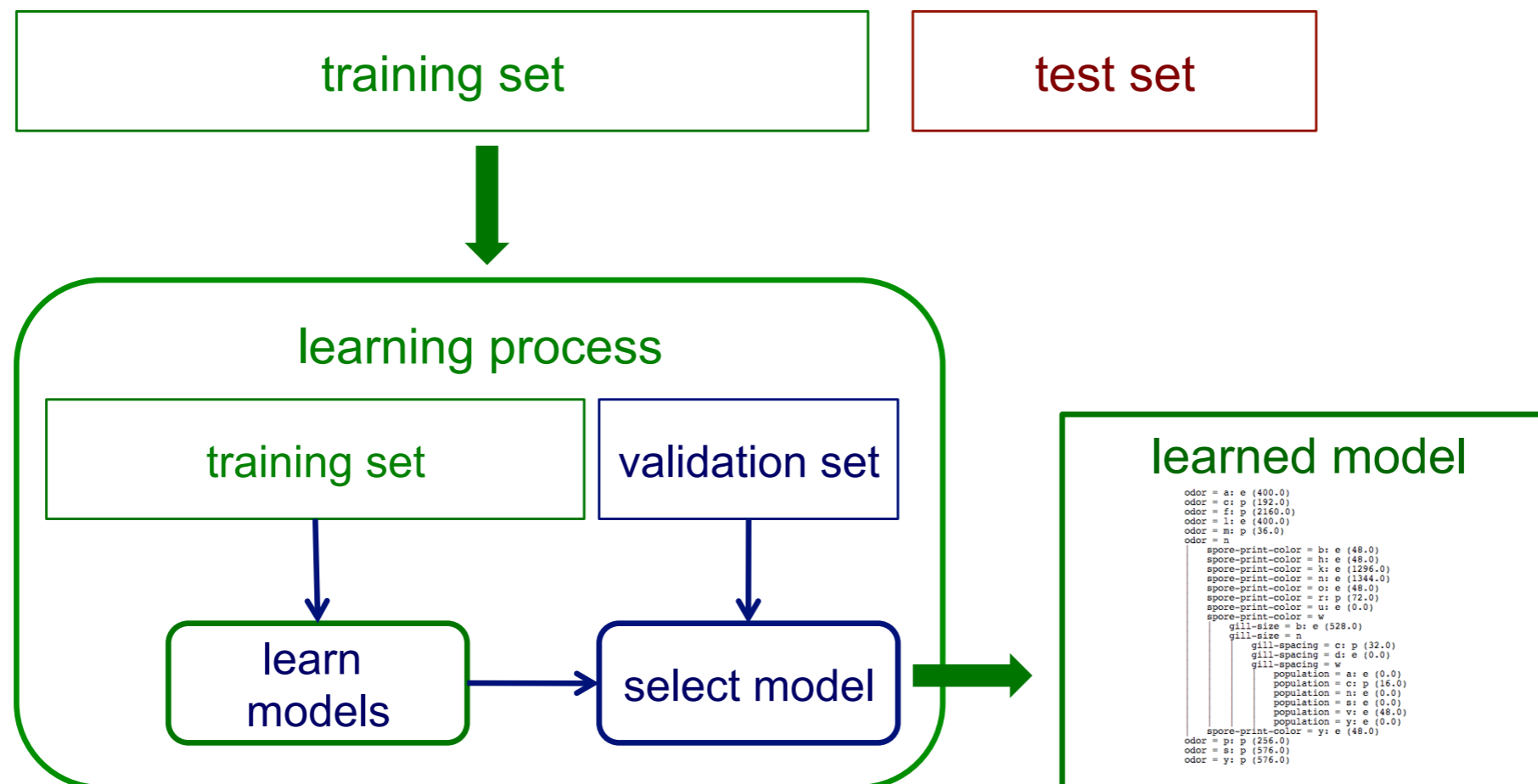


# Test sets revisited

- 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

# Validating (tuning) sets

- Suppose we want unbiased estimates of accuracy during the learning process (e.g. to choose the best regularization parameter for linear regression)?



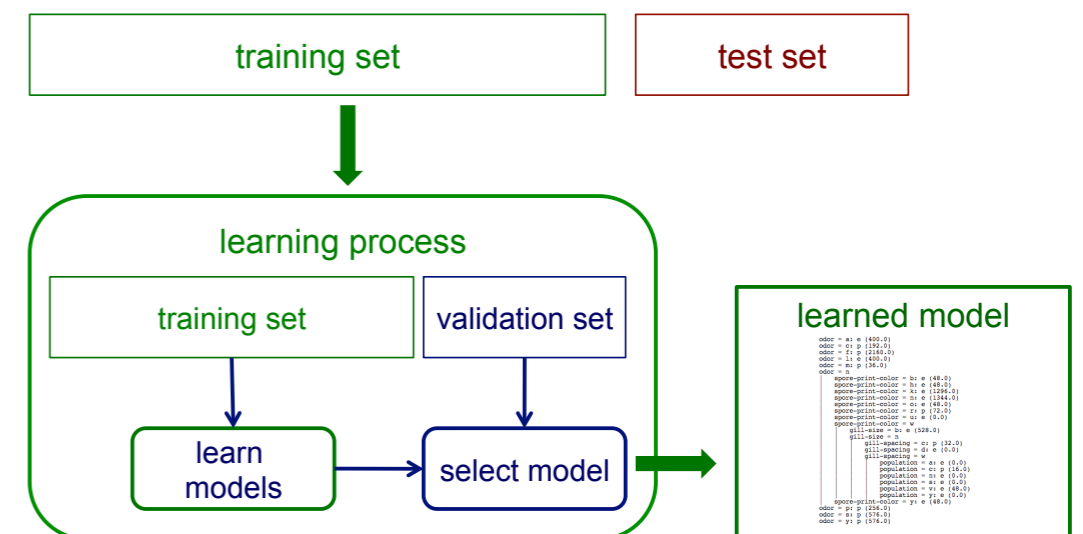
Partition training data into separate training/validation sets

# Exercise: L1 regression models

- Imagine you have a dataset with 5 observations (inputs)
- You expand up your inputs into a 9-th order polynomial
  - For  $d = 5$  (number of inputs),  $k = 9$  (the order), the total number of terms is  $(d+k)$  choose  $k$ . For  $d = 5, k = 9$ , this is 2000
- Now you are going to run L1 regression, to subselect features

• why? 
$$\min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

- How do you pick the regularization parameter, lambda?



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 the learning process
- A single training set doesn't tell us the variability of the algorithm across different training sets (if that is of interest)

# Exercise: How much hold-out test data would be enough?

- Hard to make the choice of having a single, hold-out test set, if not sure how many samples you lose for training
- Sufficient: enough test data to get an accurate estimate of the true expected error

- Average error on test data is a sample average of true error

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 \approx \mathbb{E}[(\mathbf{X}^\top \mathbf{w} - Y)^2]$$

- How quickly does this become correct?

Not the matrix  $X$ !  
This is the random vector  $X$



# How much hold-out test data is enough? Confidence intervals

- How quickly does this become correct?

$$Z = \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i^\top \mathbf{w} - Y_i)^2 \approx \mathbb{E}[(\mathbf{X}^\top \mathbf{w} - Y)^2] = \mu$$

- We want a confidence interval or upper bound on error:

$$\Pr(\mu \leq Z + u) = 1 - \alpha/2 \quad \text{for some } u$$

$$\Pr(Z - u \leq \mu \leq Z + u) = 1 - \alpha$$

- If squared-errors are Gaussian, have interval

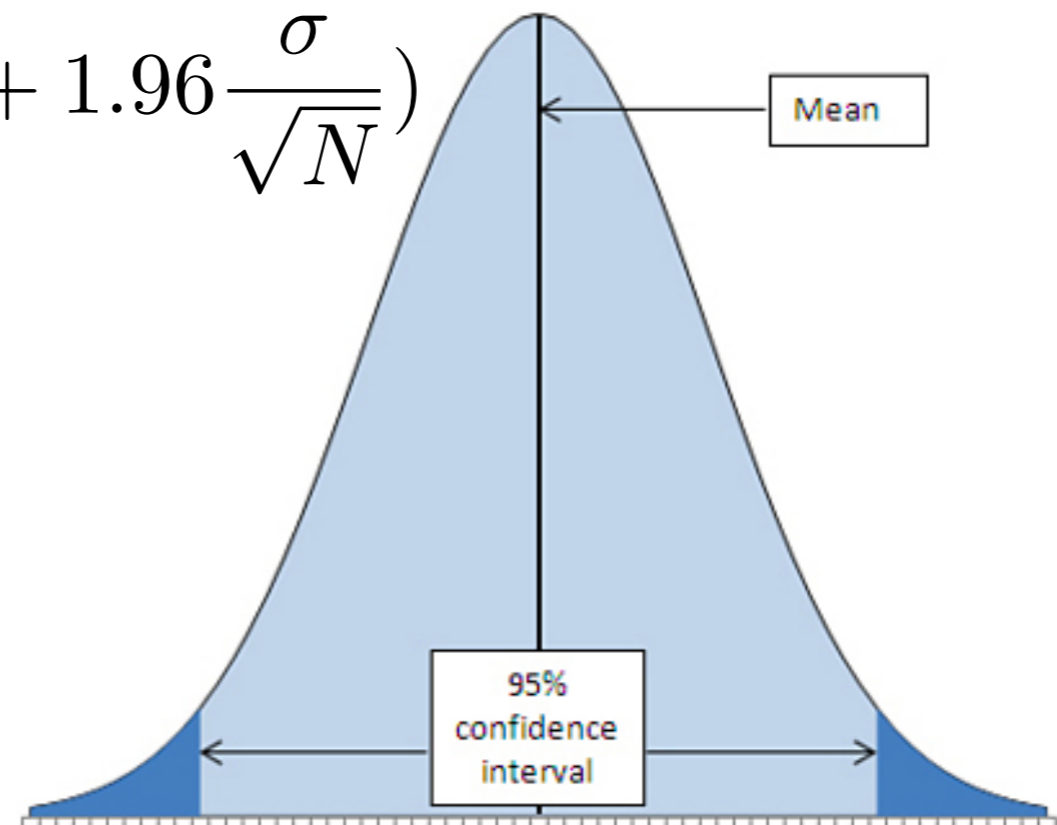
$$u = \frac{1.96\sigma}{\sqrt{n}}, \text{ giving } \Pr(\mu \leq Z + u) = 0.975$$

- If errors follow unknown distribution, need other approaches

# Normal confidence interval

- What if we plot our 100 errors and they do not look normally distributed? We'll talk about this later
- The normal 95% confidence interval determines the endpoints that contain 95% of the mass between them, under the curve

$$0.95 = P\left(\bar{X} - 1.96 \frac{\sigma}{\sqrt{N}} \leq \mu \leq \bar{X} + 1.96 \frac{\sigma}{\sqrt{N}}\right)$$



# Errors following unknown distribution

- How quickly does this become correct?

$$Z = \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i^\top \mathbf{w} - Y_i)^2 \approx \mathbb{E}[(\mathbf{X}^\top \mathbf{w} - Y)^2] = \mu$$

- We want a confidence interval or upper bound on error:

$$\Pr(\mu \leq Z + u) = 1 - \alpha/2 \quad \text{for some } u$$

$$\Pr(Z - u \leq \mu \leq Z + u) = 1 - \alpha$$

- If errors follow unknown distribution, are i.i.d. and bounded in range  $[a,b]$ , could use Hoeffding inequality

$$u = \frac{(b - a)^2 \log(2/\alpha)}{\sqrt{2n}}, \quad \text{giving } \Pr(Z - u \leq \mu \leq Z + u) = 1 - \alpha$$

# Examples for specific $n$ and tolerance levels

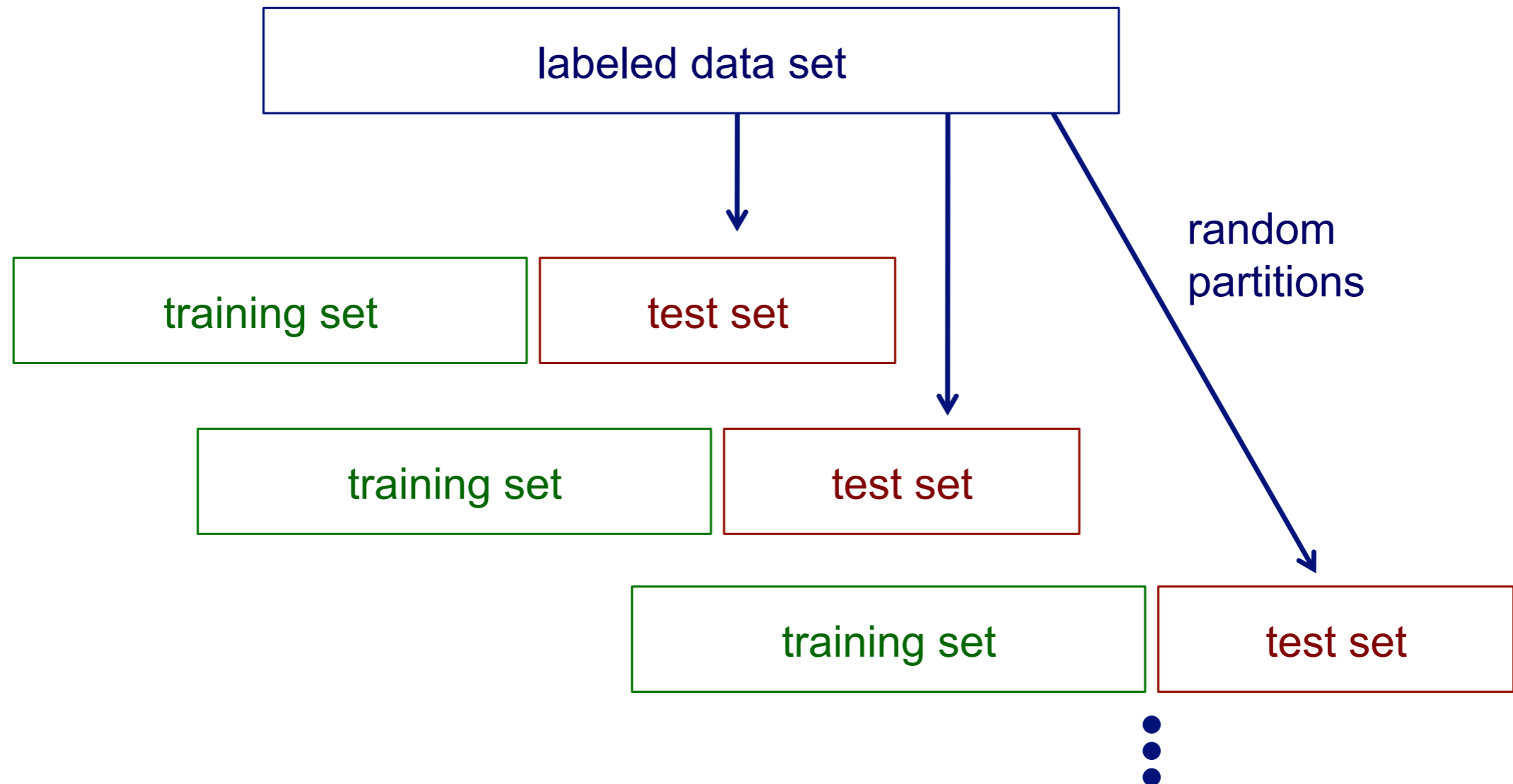
- If  $\alpha = 0.05$  and want  $Z$  to be within  $u = 10^{-2}$  of  $\mu$ :
- if errors Gaussian and  $\sigma = 1.0$ , then need about  $n = 40000$  test examples!
- This is only if want highly accurate estimates of error. If instead want to compare two models, just need to be able to say that one error is lower than the other
  - could even be quite off, and still be able to say one is better than the other with high-confidence
  - discuss this more later in the slides

# Recap: Using a single test set

- Even if have a hold-out test-set (i.e., enough data to do so), if the test-set labels influence the learned model in any way, the accuracy estimates will be biased
- It can be difficult to avoid peeking at test set, and so generally beneficial to consider other evaluation schemes
- One strategy: resampling approaches (random resampling and cross-validation)

# Random resampling

- Repeatedly randomly partitioning the available data into training and set sets.

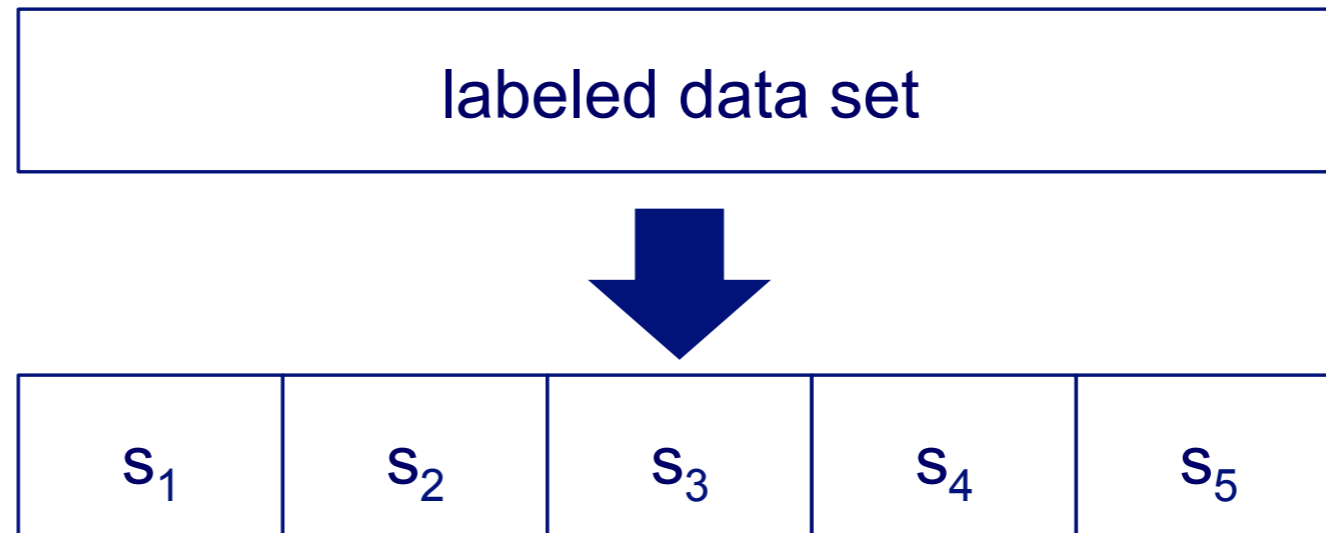


# Why might random sampling be more robust?

- Imagine a benchmark dataset
  - e.g., MNIST is a character recognition dataset, split into one training and test set
- One paper reports that a specific neural network did really well, with a narrow regularization range
- Now you build on this, adding say a few more nodes, increase the number of regularization parameters in that range and train and report test accuracy
- Lo-and-behold, you do better than the previous model!
- Is there an issue?

# Cross validation

partition data  
into  $n$  subsamples

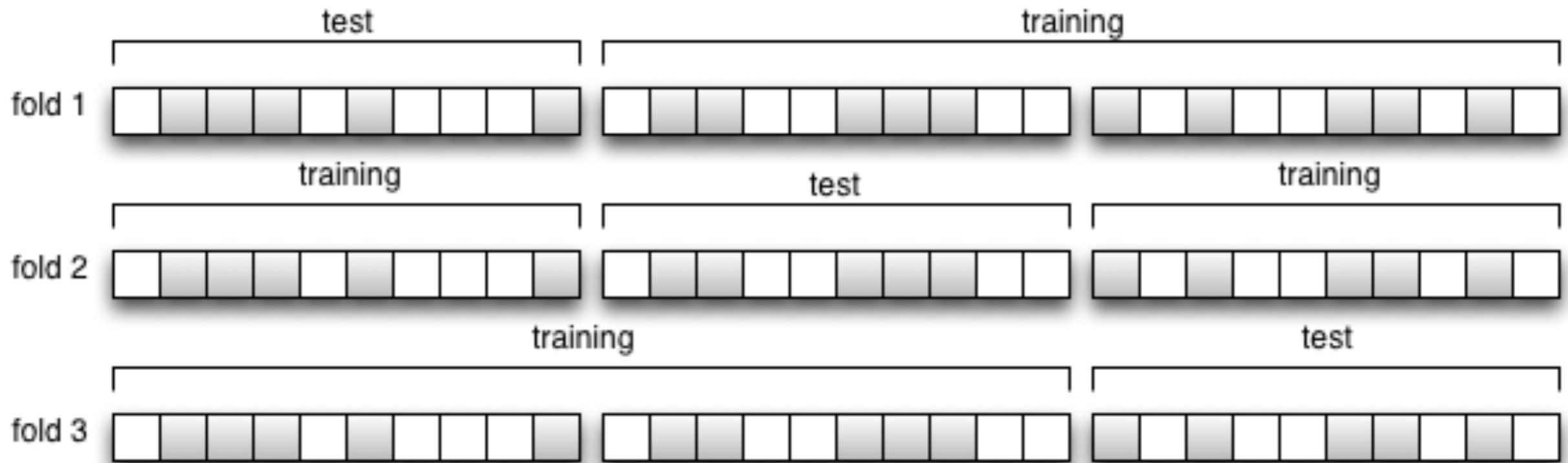


iteratively leave one  
subsample out for  
the test set, train on  
the rest

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	$S_1$ $S_2$ $S_3$ $S_4$	$S_5$



# Another view of cross validation



# 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

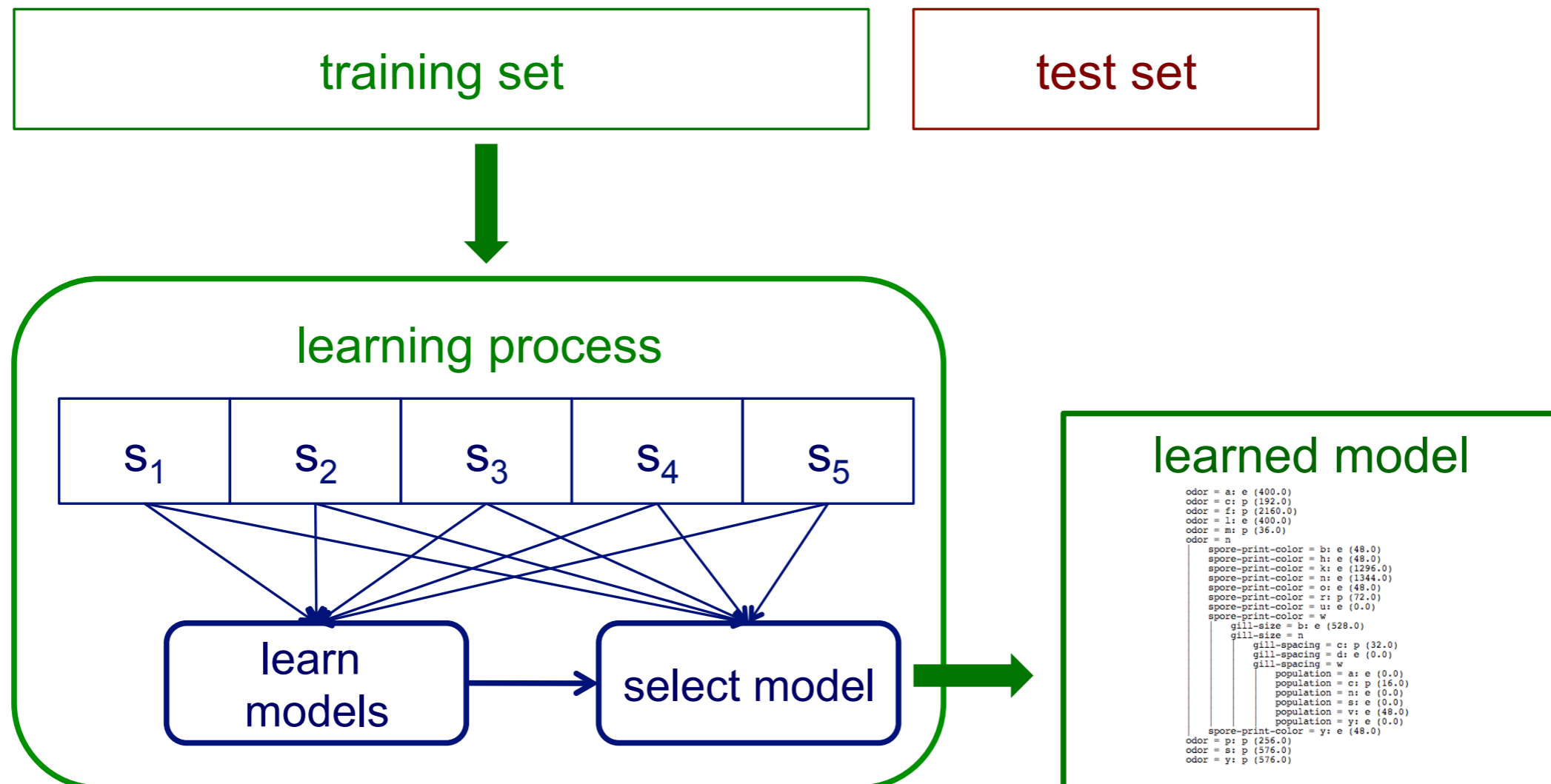
$$\text{accuracy} = 73/100 = 73\%$$

# 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
- 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 regularization parameter for linear regression)



# Example: selecting the regularization parameter with CV

- Given a training set
  1. partition the training set into  $n$  folds,  $s_1, \dots, s_n$
  2. for each value of  $\lambda$  considered
    - for  $i = 1$  to  $n$ 
      - learn regression model using all folds but  $s_i$
      - evaluate accuracy on  $s_i$
  3. select  $\lambda$  that resulted in the best accuracy for  $s_1, \dots, s_n$
  4. learn model using entire training set and selected  $\lambda$
- This is typically run separately for each training set
  - What would it mean to use this to pick hyperparameters across training sets?

# Overall experiment design

- Hyperparameters to try for each algorithm
  - e.g., have to decide on the range of lambda to try
  - e.g., which optimizer to use in algorithm
- Depending on choices, answering a different question
  - e.g., one training/test split approximates how a learned model performs, when training on that much data
  - e.g., multiple training/test splits approximates performance of a learning method, with given hyperparameters
  - e.g., could report algorithm with parameter settings as a single learning method, understand parameter sensitivity
- Running a thorough experiment can be difficult, but rewarding

# Practical questions (and break)

- What are possible hyperparameters that you may have to deal with in machine learning?
- What does it mean if we get 10 measures of accuracy, and they are quite different from each other?
- What if we have more than one dataset?
- What might an experiment look like, that tests learning speeds?
  - measuring speed is straightforward. Can't we just test the learning speed on the training set, for our different algorithms?

# Reminders: Oct. 22, 2019

- Marks for Assignment 1 released today
- Delay Assignment 2 deadline until Friday night



# Testing for significance

- Imagine now that you have 100 error values from your experiment (obtained from 100 random training/test splits)
- Imagine you choose hyperparameters with CV each time on the training set, and so are evaluating Algorithm A and B
  - rather than just a specific learned model
  - assuming you tested a reasonably large range of hyperparameters
- The average error value for Algorithm A is smaller than Algorithm B: can you conclude that A is better?
- Need statistical significance tests, the mean not enough info

**Fun fact:** Your knowledge of ML models will help make statistical significance make much more sense

# Statistical significance tests

- Null hypothesis: A and B have the same generalization performance (i.e., A and B have the same expected error)
  - by running on multiple random test sets, obtaining unbiased estimates of expected error
- Alternative hypothesis: A and B have different generalization performance
- Confidence intervals and standard error
- Paired t-test

# Computing confidence intervals

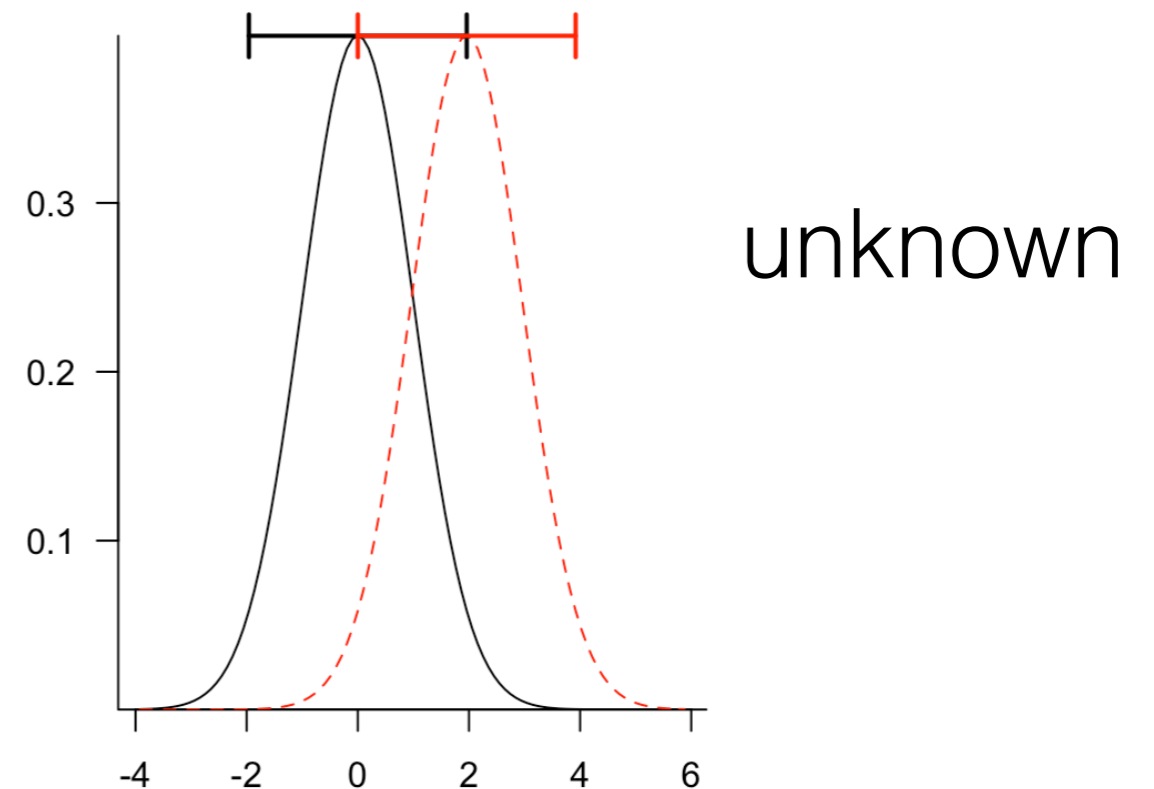
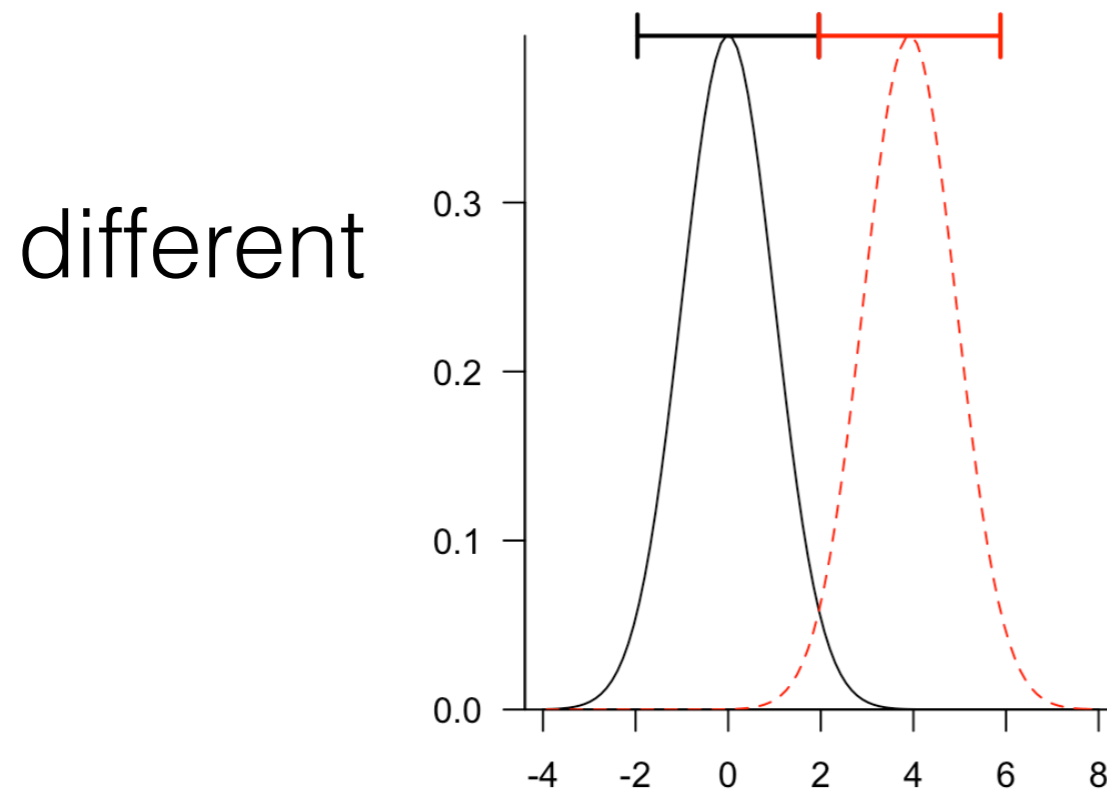
- Confidence interval around expected error  $\bar{X}$  of an algorithm
  - might make a normal assumption (because of central limit theorem)

$$0.95 = P\left(\bar{X} - 1.96 \frac{\sigma}{\sqrt{N}} \leq \mu \leq \bar{X} + 1.96 \frac{\sigma}{\sqrt{N}}\right)$$

- If two confidence intervals do not overlap, can say that the two algorithms have statistically significantly different expected errors (and so that one has statistically significantly lower expected error than the other)
- If the two confidence intervals do overlap, need to use a statistical significance test
  - see <http://www.cs.iastate.edu/~honavar/dietterich98approximate.pdf> for a comparison of the performance of several statistical tests

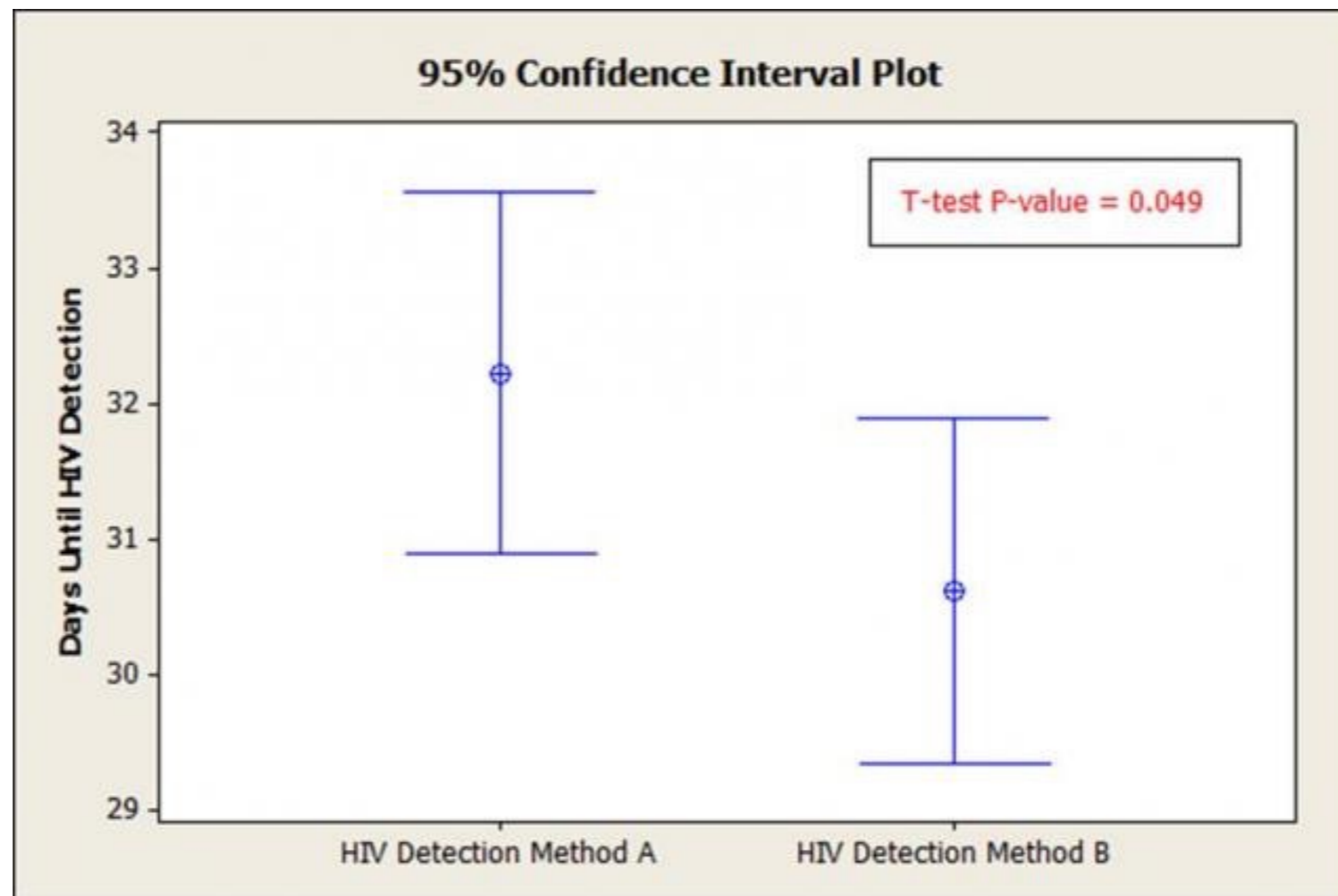
# Comparing two algorithms on a single domain

- Imagine you have  $N$  independent test-samples, giving  $N$  paired measures of error for the two algorithms
- Simplest (not very powerful) strategy:
  - compute two (95%) confidence intervals for the means
  - if the two intervals do not overlap, means are significantly different



# More powerful strategy: t-test

- Confidence intervals may overlap, but the means may still be statistically different
- Paired t-test enables a more powerful comparison
  - more ability to reject the null hypothesis



# Paired t-test

- 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

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

Now look at variability in delta, and compute confidence on average delta

# 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

TABLE B.2 THE  $t$  DISTRIBUTION

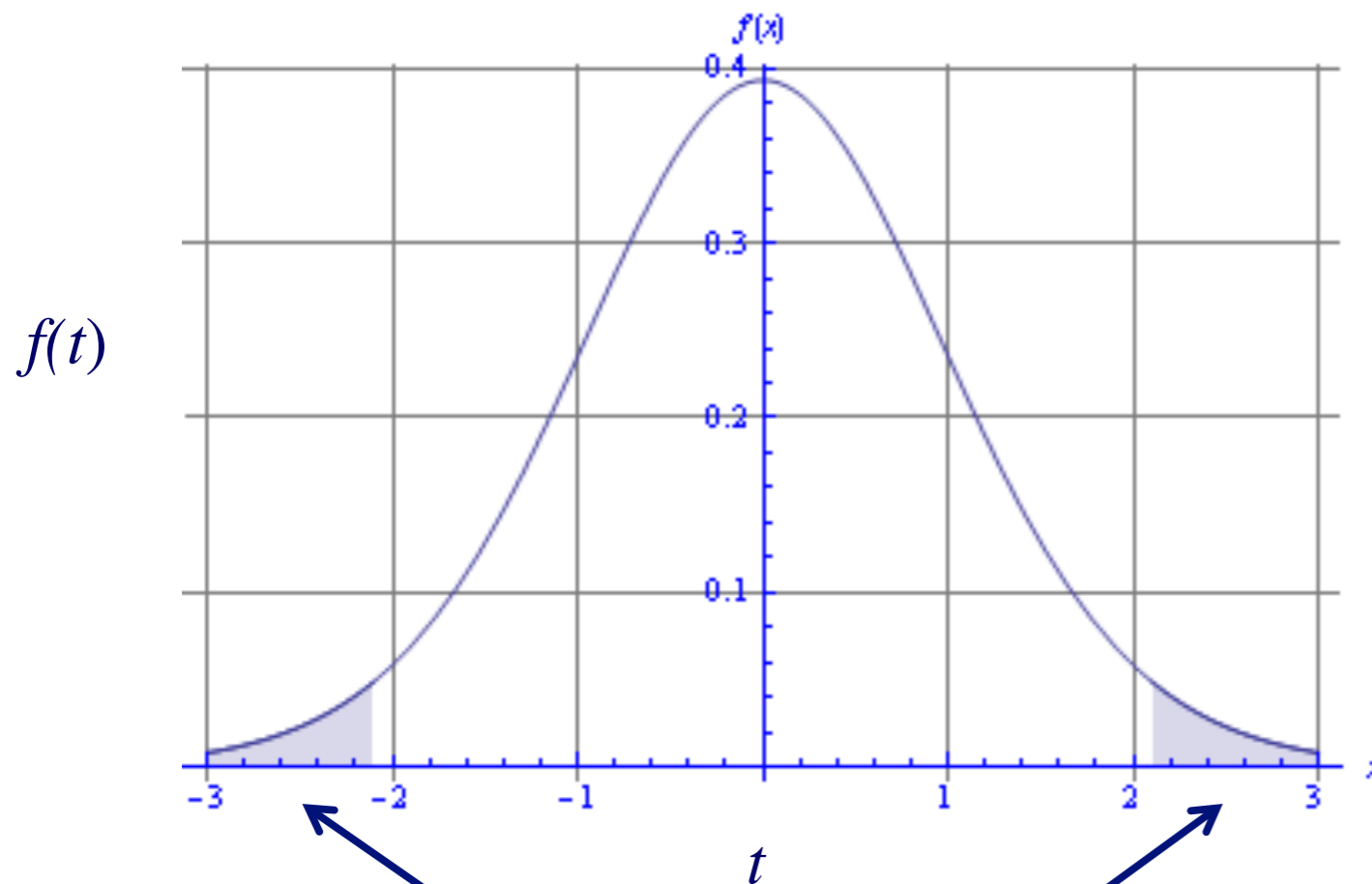
Table entries are values of  $t$  corresponding to proportions in one tail or in two tails combined.

df	PROPORTION IN ONE TAIL				
	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	4.541
4	0.741	1.533	2.132	2.776	3.747
5	0.727	1.476	2.015	2.571	3.365
6	0.718	1.440	1.943	2.447	3.143
7	0.711	1.415	1.895	2.365	2.998
8	0.706	1.397	1.860	2.306	2.896
9	0.703	1.385	1.833	2.262	2.821
10	0.700	1.372	1.812	2.228	2.764
11	0.697	1.363	1.796	2.201	2.718
12	0.695	1.356	1.782	2.179	2.681
13	0.694	1.350	1.771	2.160	2.650
14	0.692	1.345	1.761	2.145	2.624
15	0.691	1.341	1.753	2.131	2.602
16	0.690	1.337	1.746	2.120	2.583
17	0.689	1.333	1.740	2.110	2.567
18	0.688	1.330	1.734	2.101	2.552
19	0.688	1.328	1.729	2.093	2.539
20	0.687	1.325	1.725	2.086	2.528
21	0.686	1.323	1.721	2.080	2.518
22	0.686	1.321	1.717	2.074	2.508
23	0.685	1.319	1.714	2.069	2.500
24	0.685	1.318	1.711	2.064	2.492
25	0.684	1.316	1.708	2.060	2.485
26	0.684	1.315	1.706	2.056	2.479
27	0.684	1.314	1.703	2.052	2.473
28	0.683	1.313	1.701	2.048	2.467
29	0.683	1.311	1.699	2.045	2.462
30	0.683	1.310	1.697	2.042	2.457
40	0.681	1.303	1.684	2.021	2.423
60	0.679	1.296	1.671	2.000	2.390
120	0.677	1.289	1.658	1.980	2.358
$\infty$	0.674	1.282	1.645	1.960	2.326

STATISTICAL TABLES

© H. B. Fisher and F. Yates, *Statistical Tables for Biological, Agricultural and Medical Research*, 6th ed. London: Longman Group Ltd., 1963 (previously published by Oliver and Boyd Ltd., Edinburgh). Adapted and reprinted with permission of the Addison Wesley Longman Publishing Co.

# 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

$p$  is the probability of seeing  $t$ , under our null hypothesis

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

$p(t \mid \text{means are the same})$



# Evaluation summary

- Define “better” and your hypothesis upfront, to direct experiment to answer that hypothesis
- Be cognizant of your choices
  - e.g., hyperparameters, optimizers
  - e.g., number of folds
- Do not cheat by looking at test data
  - then you’ll just do poorly on some new data, outside the test data
  - Random sampling approaches are more robust to this cheating
- To avoid overfitting hyperparameter selection on the given training data, systematically sweep parameters rather than using human guessing and testing
  - this could bias you to training set, and cause bad performance on test

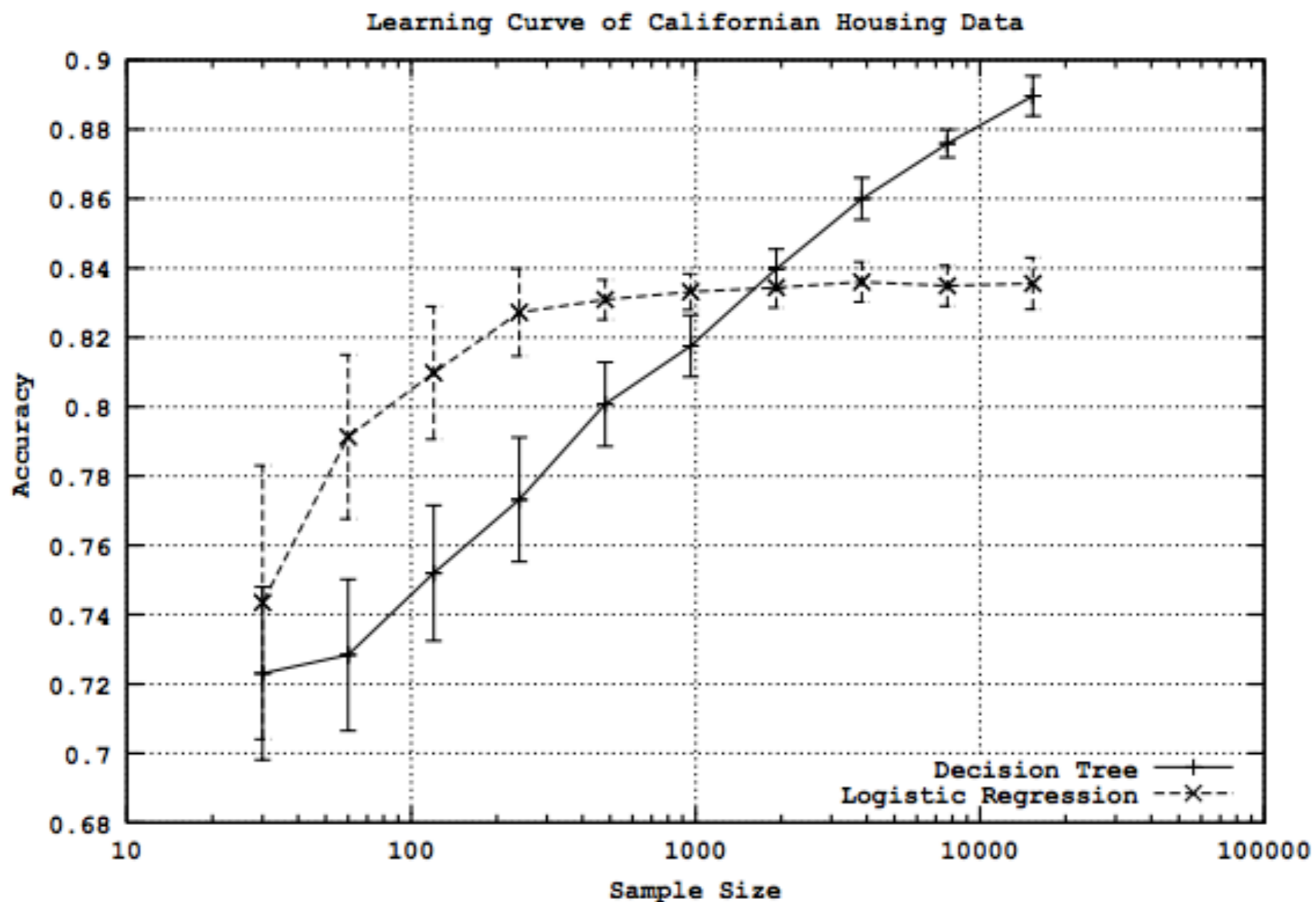
# Exercise: Generating learning curves

- Imagine you want to test the sample efficiency of logistic regression and decision trees
  - An algorithm is sample efficient if it can get good generalization performance using a small number of samples
  - For example, Algorithm 1 needs at least 1000 samples whereas Algorithm 2 only needs about 100 samples, to reach a similar level of performance on test data
  - Which models might be more or less sample efficient?
- How would you do this?

# 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*



# Learning curves

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

