

# Machine Learning

CSCI 5622 Fall 2020

Prof. Claire Monteleoni

# Today

- Decision trees (continued)
- Overfitting vs. generalization
  - bias-variance tradeoff
- Model evaluation/selection
  - Validation, Cross-Validation
- [If time] Intro to ensemble methods
  - Random forests

# Decision tree issues

Very expressive family of classifiers:

Any type of data can be accommodated:  
real, Boolean, categorical, ...

Can perfectly fit any self-consistent training set

- no example pairs  $(x, y), (x, y')$

But this also means that there is serious danger of  
overfitting.

# Greedy decision tree building

Start with all points in a single node

Repeat

Pick the split with greatest benefit

Until ???

When to stop?

- (i) When each leaf is pure?
- (ii) When the tree is already pretty big?
- (iii) When each leaf has uncertainty < some threshold?

Common strategy: keep going until leaves are pure (recall: this didn't work too well for us earlier...)

Then, shorten the tree by pruning, to correct the overfitting problem.

# What is overfitting?

Data comes from some true underlying distribution  $D$  on  $X \times Y$ .  
[ $X$  = input space,  $Y$  = label space]

All we ever see are samples from  $D$ : training set, test set, etc.

When we choose a classifier  $h: X \rightarrow Y$ , we can talk about its error  
on the training set  $(x_1, y_1), \dots, (x_m, y_m)$ :

$$\hat{\epsilon}(h) = \frac{1}{m} \sum_{i=1}^m \mathbf{1}(h(x_i) \neq y_i)$$

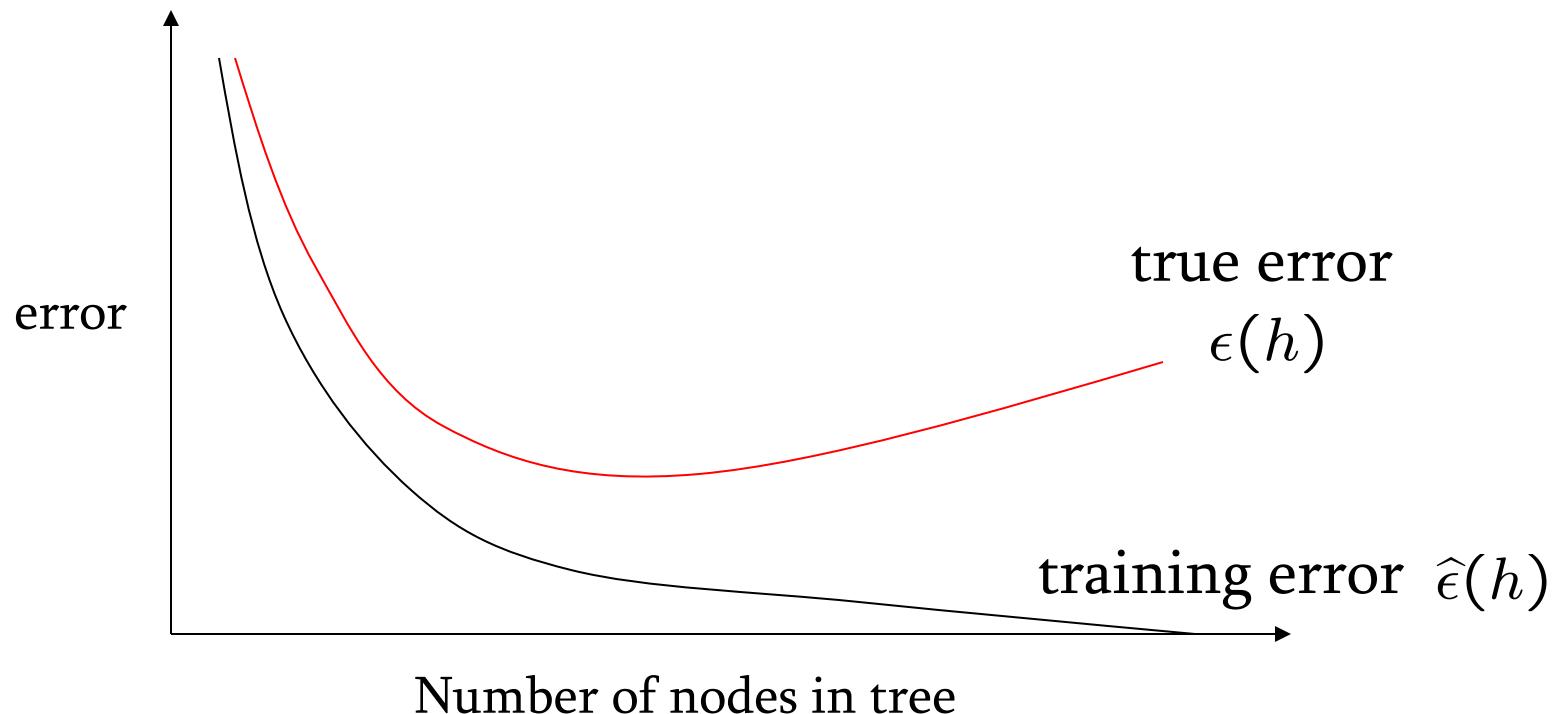
But we can also talk about its true error:

$$\epsilon(h) = \mathbf{P}_{(x,y) \sim D}(h(x) \neq y)$$

How are these two quantities related?

# Overfitting: picture

E.g., building a decision tree



As we make our tree more and more complex:

- The training error keeps going down
- but the true error stops improving and may even get worse!

# Overfitting: one perspective

1. The true underlying distribution D is the one whose structure we would like to capture
2. The training data reflects the structure of D, so it helps us.
3. But it also has **chance** structure of its own – we must try to avoid modeling this.



For instance:  $D = \text{uniform distribution over } \{1, 2, 3, \dots, 100\}$

Pick three training points: eg. 6, 12, 98.

They all happen to be even: but this is just chance structure.

It would be bad to build this into a classifier.

# Overfitting: another perspective

“Fit a line to a point”: absurd

“Fit a plane to two points”: likewise

Moral: It is not good to use a model which is so complex that there isn’t enough data to reliably estimate its parameters.

# Decision tree pruning

[1] Split the training data  $S_{full}$  into two parts:

- A smaller training set  $S$
- A validation set  $V$  (a model of reality, a surrogate test set)

[2] Build a full decision tree  $T$  using  $S$

[3] Then prune  $T$  using  $V$ :

Repeat

if there a node  $u$  in  $T$  such that removing the subtree rooted at  $u$  decreases the error on  $V$ :

$$T = T - \{\text{subtree rooted at } u\}$$

Of course,  $V$  has chance structure too, but its chance structure is unlikely to coincide with that of  $S$ .

# Example: SPAM data set

4601 points, each corresponding to an email message

39.4% are SPAM

Each point has 57 features:

48 check for specific words, eg. FREE

6 check for specific characters, eg. !

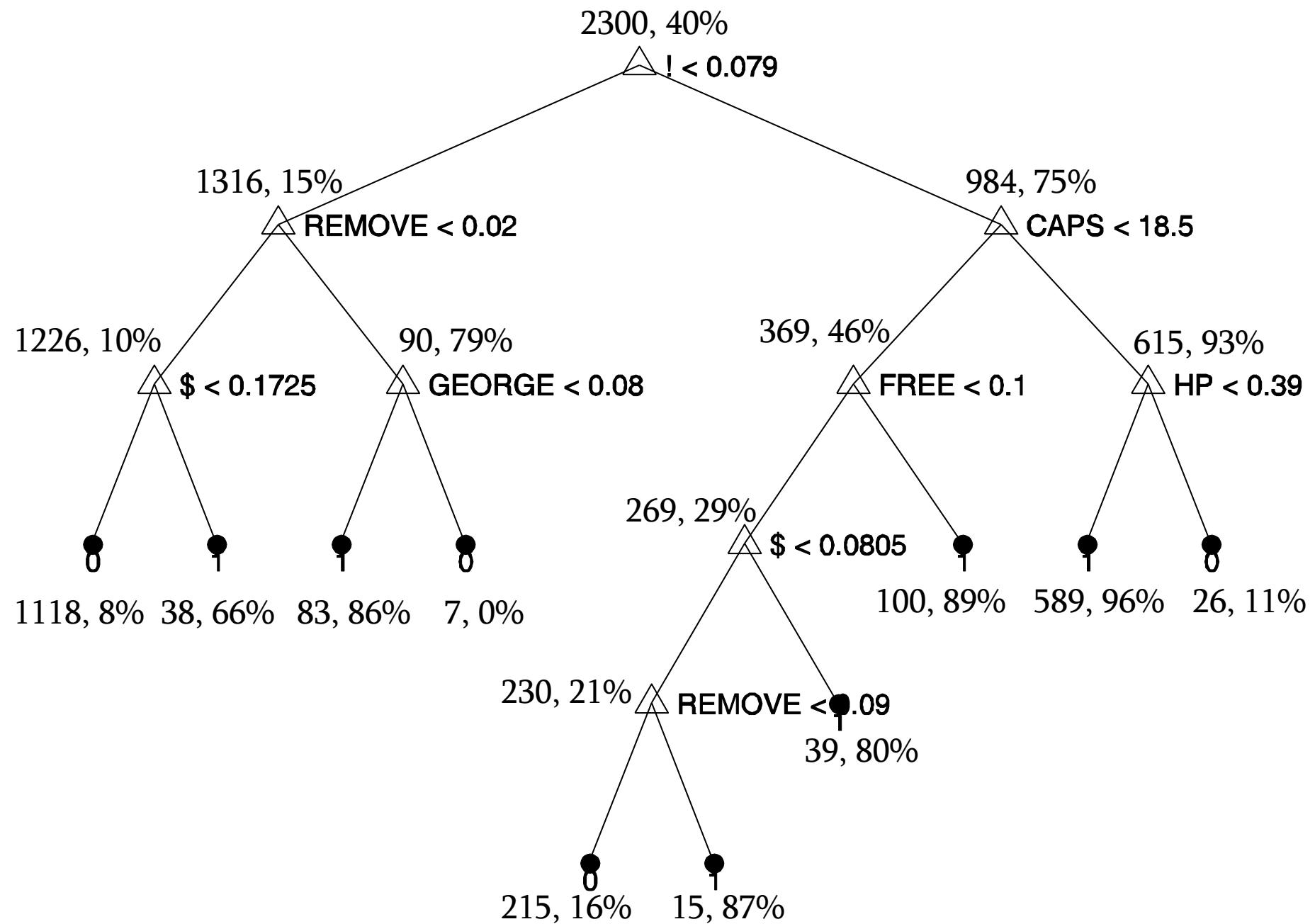
3 others, eg. longest run of capitals

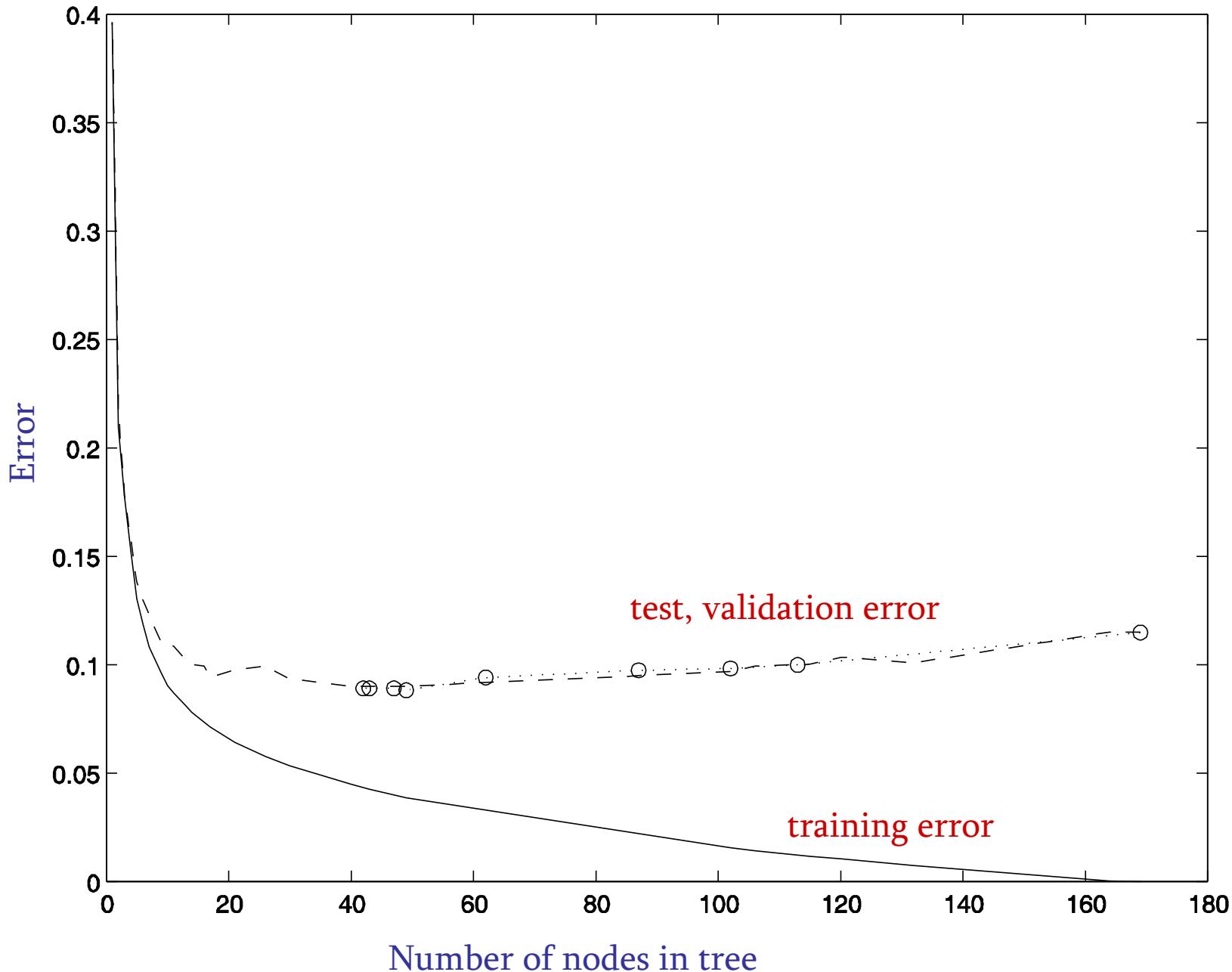
Randomly divide into three parts:

50% training data

25% validation

25% testing





# How accurate is the validation set?

How accurate are error estimates based on the validation set?

For any classifier  $h$  and underlying distribution  $D$  on  $X \times Y$ :

“true error”       $\text{err}(h) = \mathbf{P}_{(x,y) \sim D}(h(x) \neq y)$

“error on set  $S$ ”     $\text{err}(h, S) = \frac{1}{|S|} \sum_{(x,y) \in S} \mathbf{1}(h(x) \neq y)$

Suppose  $S$  is chosen i.i.d. (independent, identically distributed) from  $D$ . Then (over the random choices of  $S$ ),

$$\mathbf{E}[\text{err}(h, S)] = \text{err}(h)$$

And the standard deviation of  $\text{err}(h, S)$  is about  $1/\sqrt{|S|}$ .

# How accurate is the validation set?

Fix any  $h$ . Suppose  $S$  is chosen i.i.d. (independent, identically distributed) from  $D$ . Then (over the random choices of  $S$ ),

$$E[\text{err}(h, S)] = \text{err}(h)$$

And the standard deviation of  $\text{err}(h, S)$  is about  $1/\sqrt{|S|}$

- (i) In this scenario,  $S$  is used to assess the accuracy of a single, pre-specified classifier  $h$ .

Can the same  $S$  be used to check many classifiers simultaneously?

Answer: VC theory

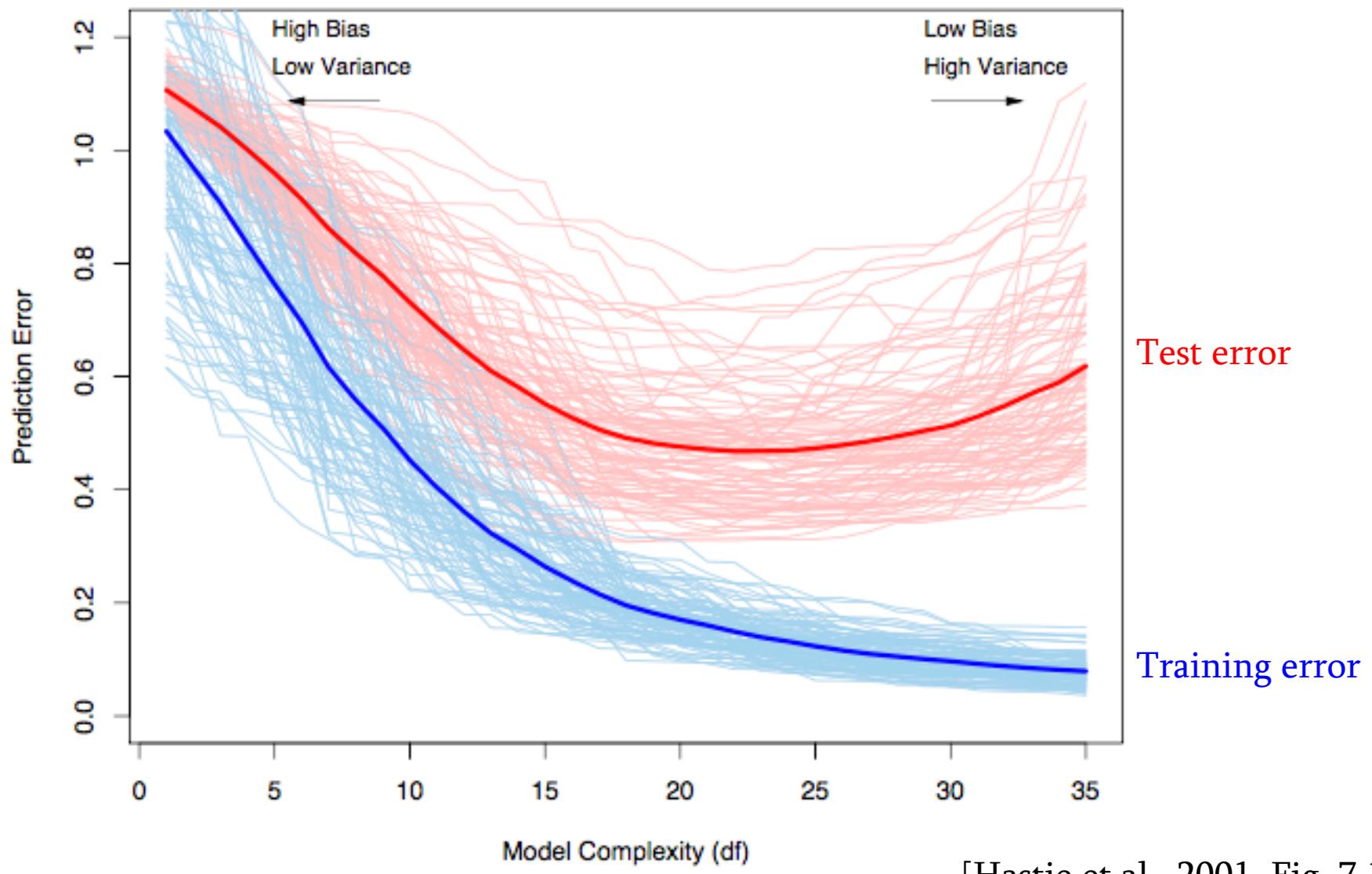
- (ii) Note however, if  $h$  was created using  $S$  as a training set, then the above scenario does not apply. In such situations,  $\text{err}(h, S)$  might be a very poor estimate of  $\text{err}(h)$ .

# Example

Predict flu trends using search data

- $X$ : search data,  $Y$ : fraction of population with flu
- $S_{\text{train}}$  = all data before 2012
- $S_{\text{test}}$  = all data in 2012

# Example



# Bias-Variance Tradeoff

Allowing the model (hypothesis class) to be more and more complex, the resulting classifier,  $h$ , will have a better and better fit to the training data, thus the bias of  $h$  reduces.

But as complexity increases,  $h$  will increasingly overfit to the chance structure of the particular training data set. So there will be higher variance of the true error of classifiers,  $h'$ , learned over different training sets.

# Bias-Variance Error Decomposition

If we assume all data points  $(x, y)$  obey:  $y = f(x) + \epsilon$

where:  $E[\epsilon] = 0$

$$\text{Var}(\epsilon) = \sigma^2$$

Then,  $\text{Err}_h(x) = E[(y - h(x))^2 | X = x]$

$$= \boxed{\sigma^2} + \text{Var}[h] + (\text{Bias}[h])^2$$

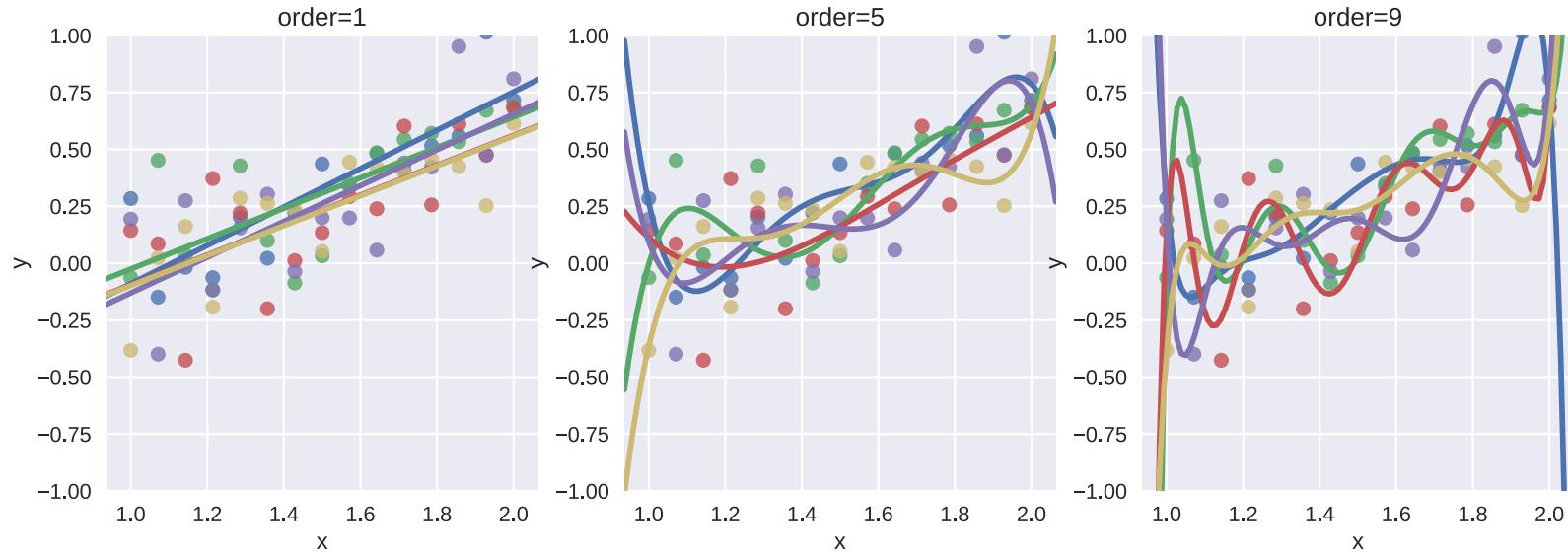
Irreducible error (can't be changed by choice of  $h$ )

where:  $\text{Var}[h] = E[h(x)^2] - E[h(x)]^2$

$$\text{Bias}[h] = (f(x) - E[h(x)])$$

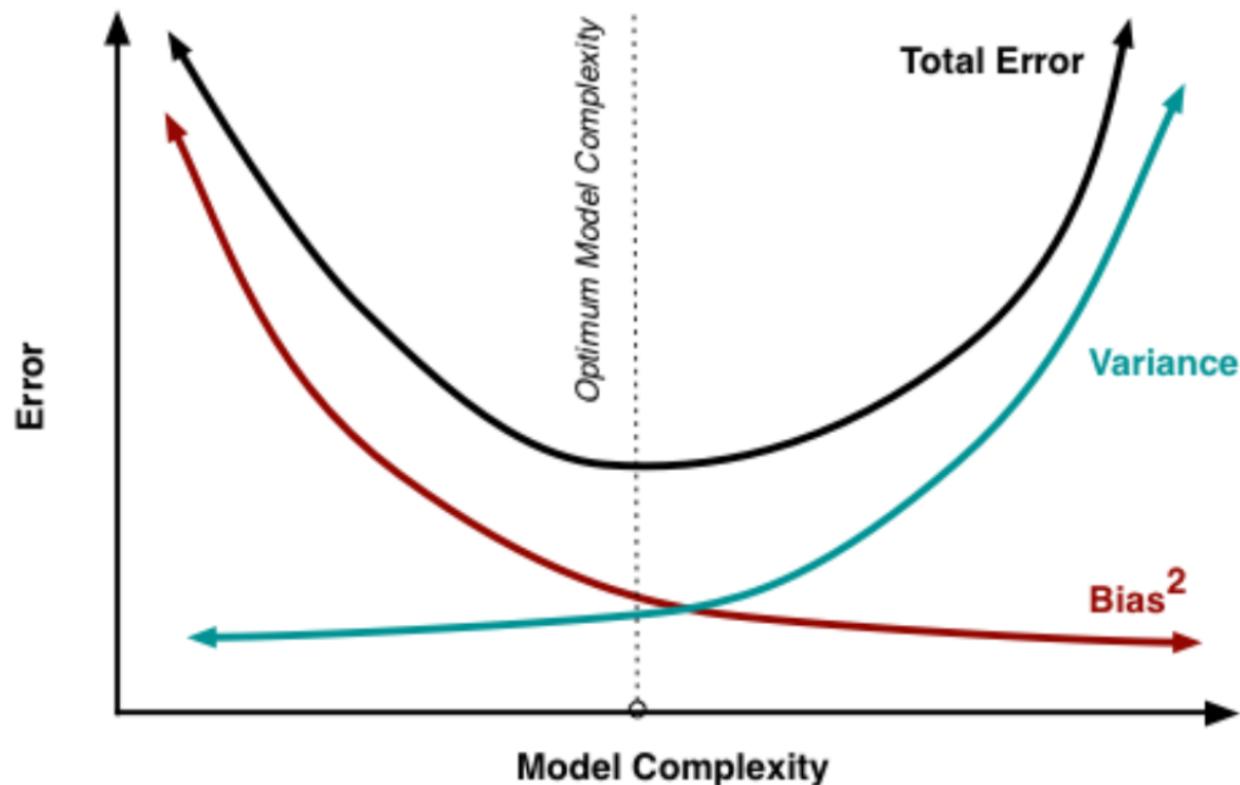
NOTE: All expectations are w.r.t. the random training set  $h$  is learned from, NOT  $x$ .

# Example



Increasing the order of the polynomial model used to fit the data, increases model complexity eventually leading to overfitting.

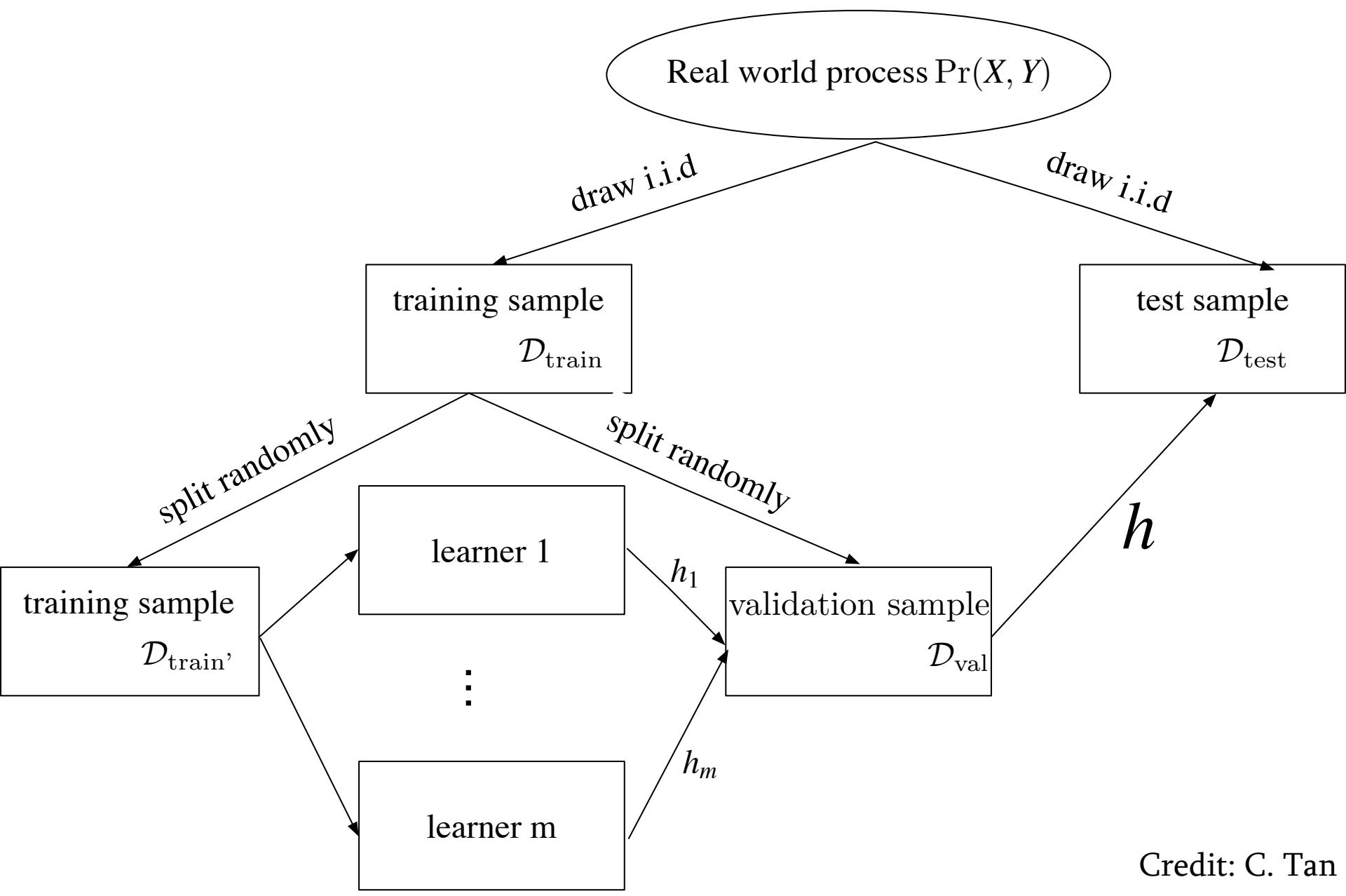
# Bias-Variance Tradeoff



Credit:

<http://scott.fortmann-roe.com/docs/BiasVariance.html>

# Model Selection



Credit: C. Tan

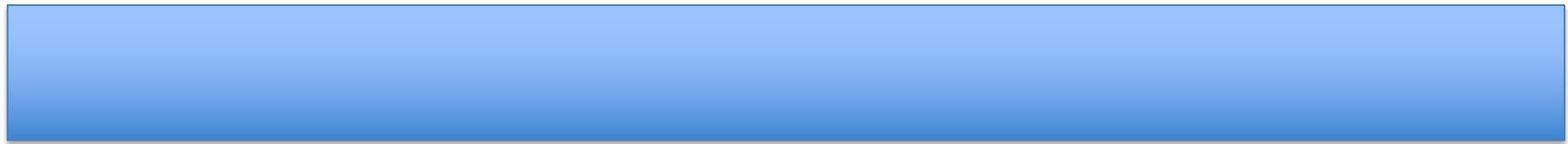
# Managing data sets: Train v. Test

- In general, you should have a **separate** training set and test set.
- Ideally they will be drawn from the same distribution.
- E.g. randomly permute your entire input data set, then fix  $T$ , and use all  $(x_t, y_t)$  for  $t < T$  for Train, and the remaining data for Test.

# Managing data sets: Train v. Test

- First create the holdout data set, the Test set, for computing the final test error. For best performance it will be a (labeled) dataset from the same distribution as the Training set.
  - E.g. technique on previous slide.
- With the remaining data set, you can either divide it into two parts for Train and Validation, as specified on the previous slide, or you can run **Cross-Validation** on it.
  - This is for model selection / parameter tuning

- Whole data set:



- Separate the test set, and hold it out:



- From the remaining data, use part for Validation:



- Or run cross validation to train and fit parameters:



# Cross validation: motivation

- To get a more **robust** estimate of the error rate, we should run more experiments, and average the error rates over the experiments.
- If we have a **huge** amount of labeled data, we can run several **disjoint** experiments (e.g. further split Train and Test into more data sets) and average the test error over the results. This is usually not the case.
  - Instead we do **Cross-Validation**.

# N-fold Cross-validation

- Cross validation data set:



- Fold 1:



- Fold 2:



- ...

- Fold n



# N-fold cross-validation

- Fix N, for example 5 or 10.
- Run N experiments. In each experiment,  $1/N$  fraction of the data is used for Test and the remaining data is used for Train. The Test sets over the N experiments are **disjoint**, and the Train sets are overlapping.
- Average the Test error over the N experiments.
- Note: different classifiers will likely be learned in each experiment, as the Train sets differ slightly.

# Leave-one-out cross-validation

- $N$  is set to the size of the data set. For each experiment, Test consists of the  $i$ -th point, and Train consists of all points but the  $i$ -th point.
- Average the test error over the  $N$  experiments.

# Model selection / Parameter tuning

- For each meta-parameter, e.g.  $k$  in  $k$ -NN, explore a range of parameter values. Log-search can be performed by doubling, or halving the value of the parameter.
- On the tuning data set (a.k.a. Validation set, or CV on the training data), try to find the “knee” in the curve: the parameter setting such that the tuning test error is minimized, before increasing again.

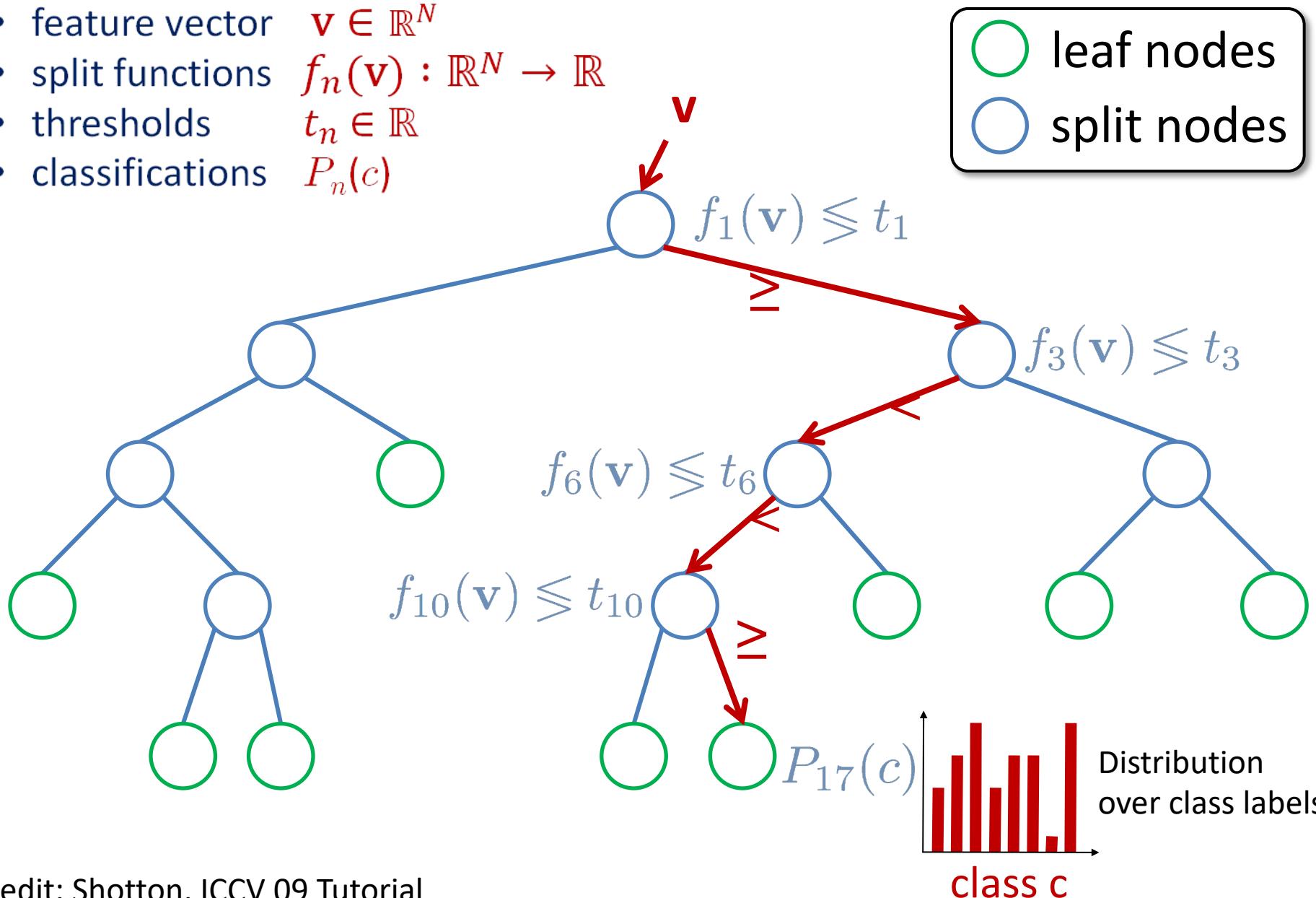
# Ensemble methods

An **ensemble** classifier combines a set of weak “base” classifiers into a “strong” ensemble classifier.

- “boosted” performance
- more robust against overfitting
- Decision Forests, Random Forests [Breiman ‘01], Bagging
- Voted-Perceptron
- Boosting
- Learning with expert advice
- ....

# Review: (Binary) Decision Trees

- feature vector  $\mathbf{v} \in \mathbb{R}^N$
- split functions  $f_n(\mathbf{v}) : \mathbb{R}^N \rightarrow \mathbb{R}$
- thresholds  $t_n \in \mathbb{R}$
- classifications  $P_n(c)$



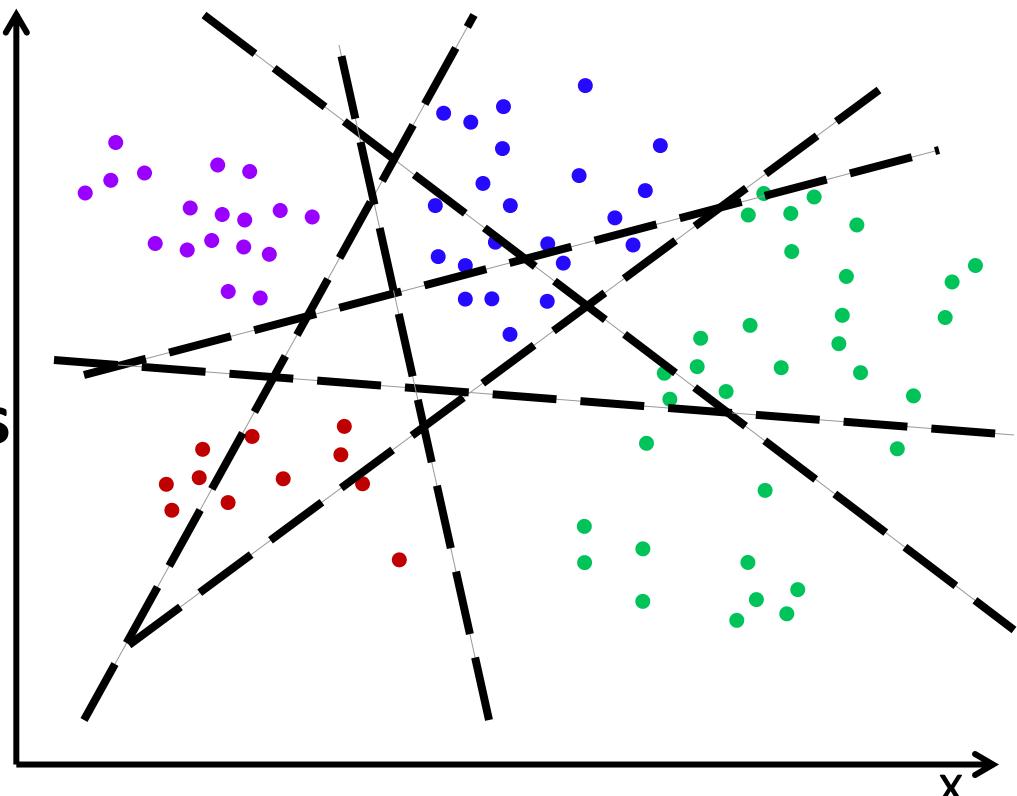
# Speeding up Decision Tree Learning

It is cumbersome to test all **possible** splits

Try several **random** splits

Keep the split that best separates data

- Reduces uncertainty



Recurse

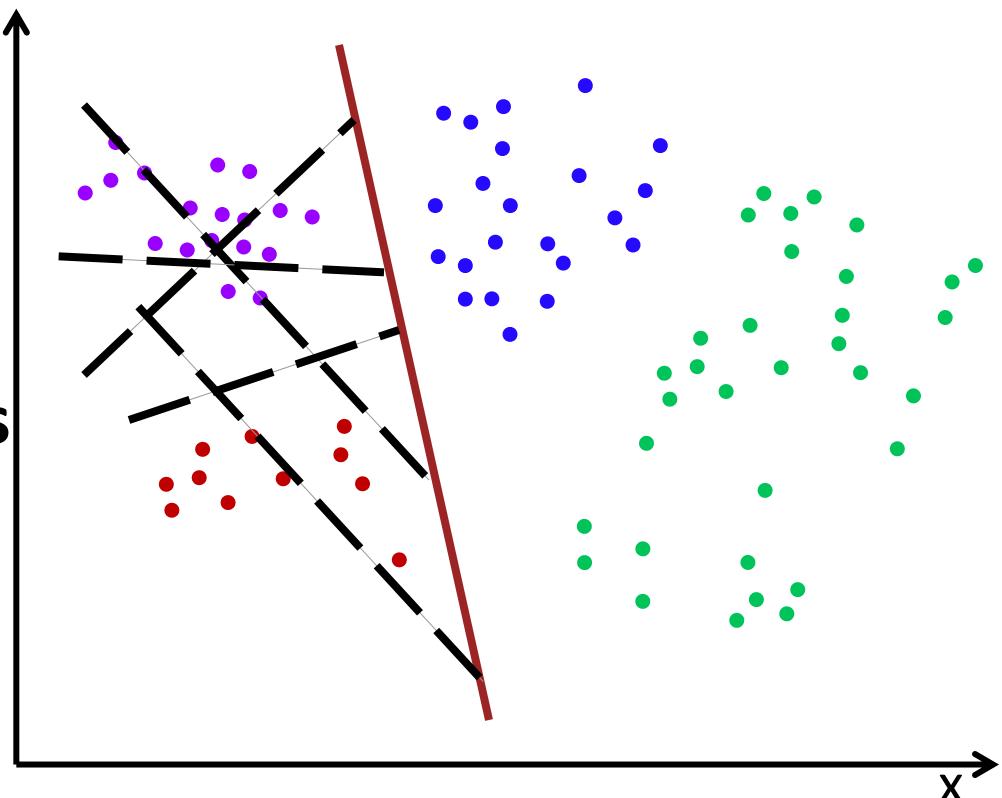
# Speeding up Decision Tree Learning

It is cumbersome to test all **possible** splits

Try several **random** splits

Keep the split that best separates data

- Reduces uncertainty



Recurse

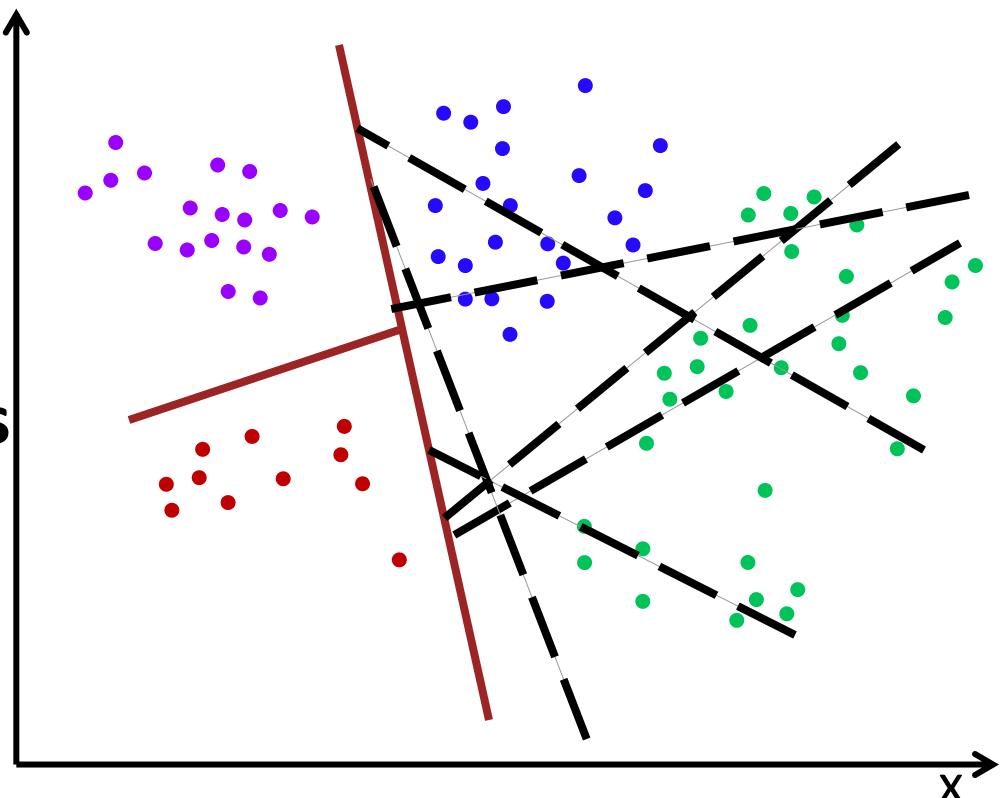
# Speeding up Decision Tree Learning

It is cumbersome to test all **possible** splits

Try several **random** splits

Keep the split that best separates data

- Reduces uncertainty



Recurse

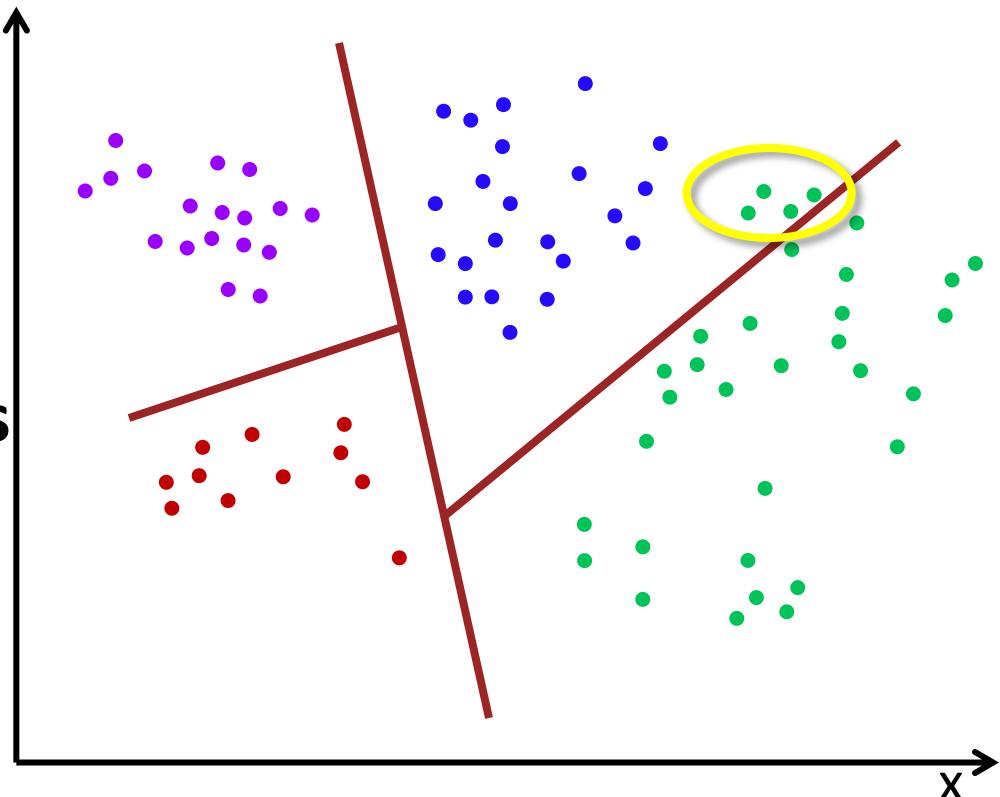
# Speeding up Decision Tree Learning

It is cumbersome to test all **possible** splits

Try several **random** splits

Keep the split that best separates data

- Reduces uncertainty



Random Decision Tree

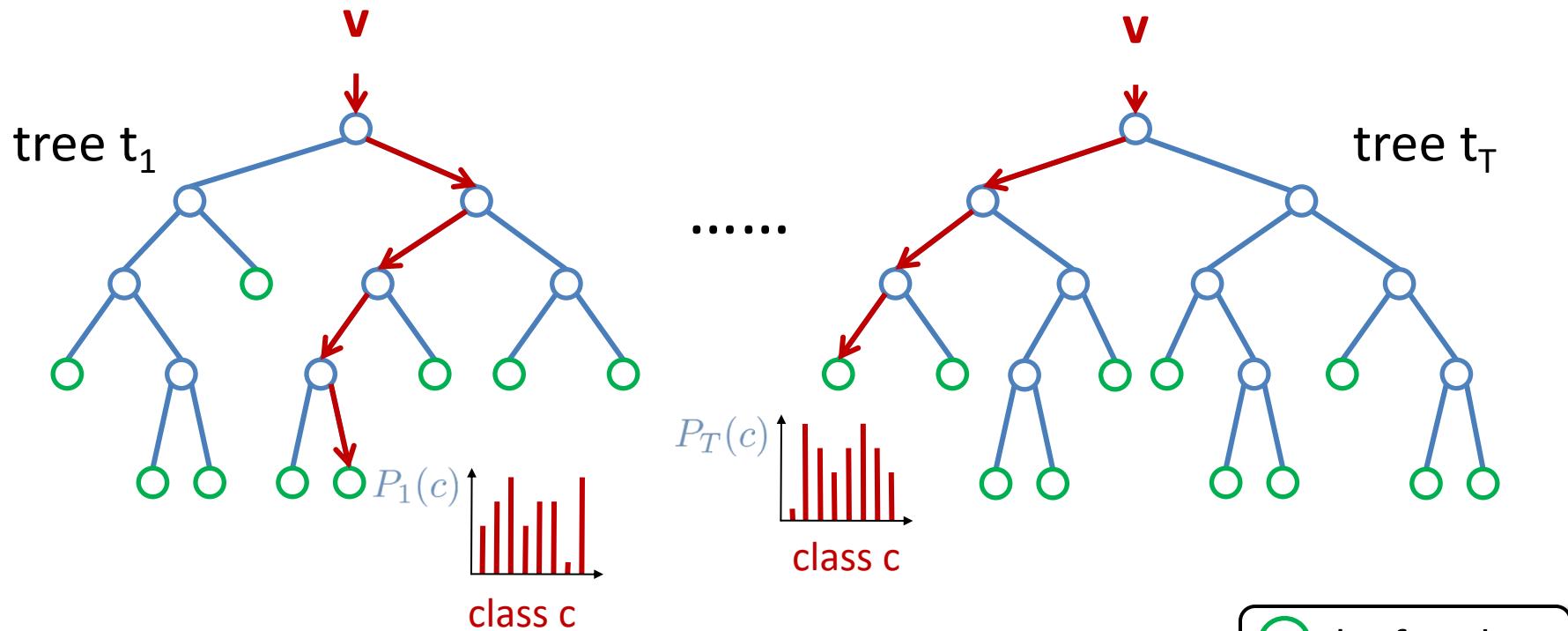
Recurse

# Random Decision Trees

- How many features and thresholds to try?
  - just one: “extremely randomized” [Geurts *et al.* 06]
  - few: fast training, may under-fit
  - many: slower training, may over-fit
- When to stop growing the tree?
  - maximum depth
  - minimum entropy gain
  - threshold changes in class distribution
  - pruning

# Decision Forests

- A forest is an ensemble of several decision trees



Classification: 
$$P(c|\mathbf{v}) = \frac{1}{T} \sum_{t=1}^T P_t(c|\mathbf{v})$$

○	leaf nodes
○	split nodes

# Learning a Forest

- Divide training examples into  $T$  subsets  $S_t$ 
  - improves generalization
  - reduces memory requirements & training time
- Train each decision tree,  $t$ , on subset  $S_t$ 
  - same decision tree learning as before
- Easy to parallelize