

# DS 4400

## Machine Learning and Data Mining I

Alina Oprea  
Associate Professor, CCIS  
Northeastern University

February 14 2019

# Logistics

- HW3 is due on Friday, February 22
- Project proposal due on Tuesday 02/26
  - 1 page description of your project, including problem statement, dataset, and ML algorithms
- Week of February 25
  - Lecture on 02/26 taught by Lisa Friedland
  - Lecture on 02/28 canceled

# Summary Decision Trees

- Representation: decision trees
- Bias: prefer small decision trees
- Search algorithm: greedy
- Heuristic function: information gain or information content or others
- Overfitting / pruning

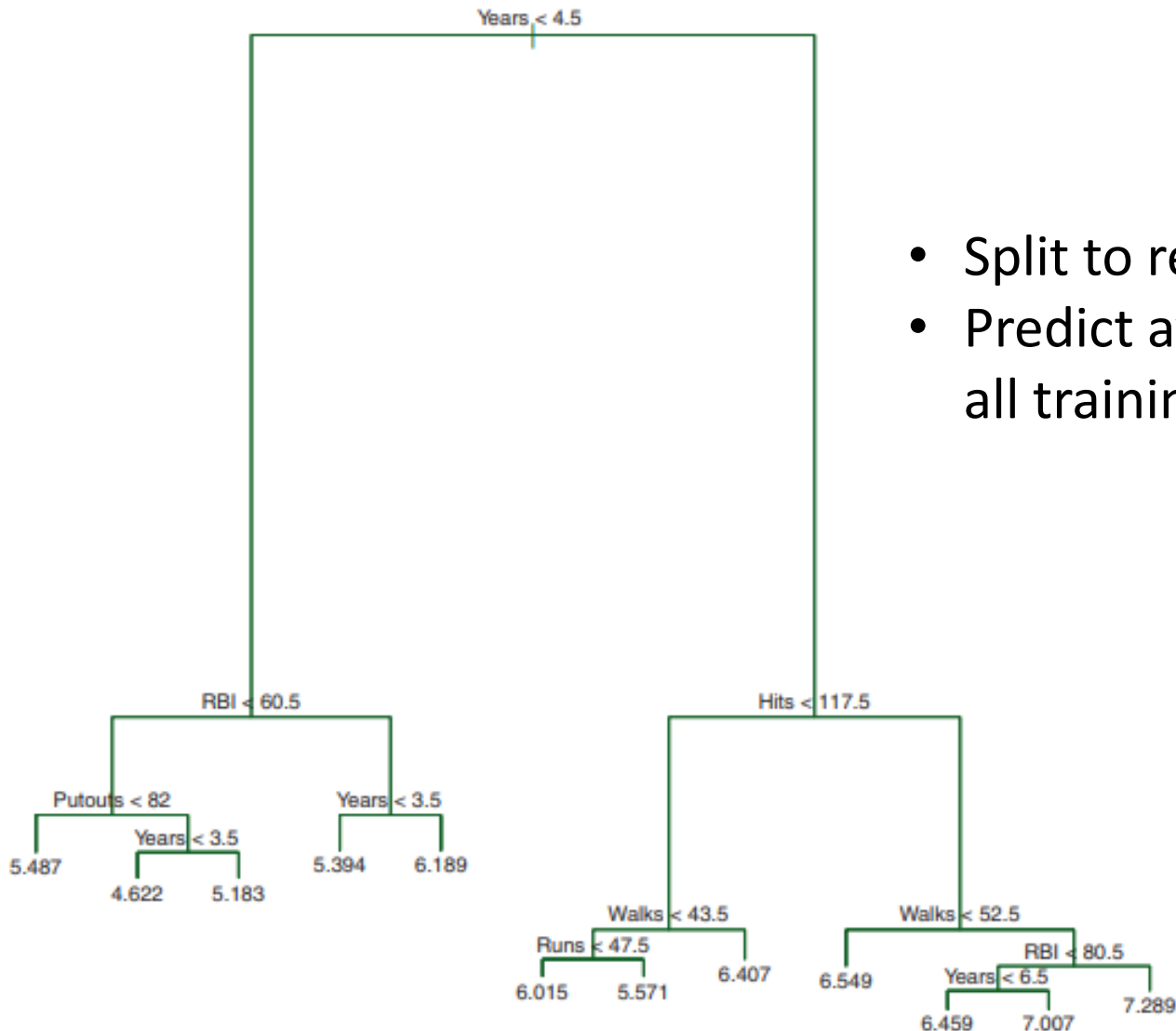
## Strengths

- Fast to evaluate
- Interpretable
- Generate rules
- Supports categorical and numerical data

## Weaknesses

- Overfitting
- Splitting method might not be optimal
- Accuracy is not always high
- Batch learning

# Regression Trees

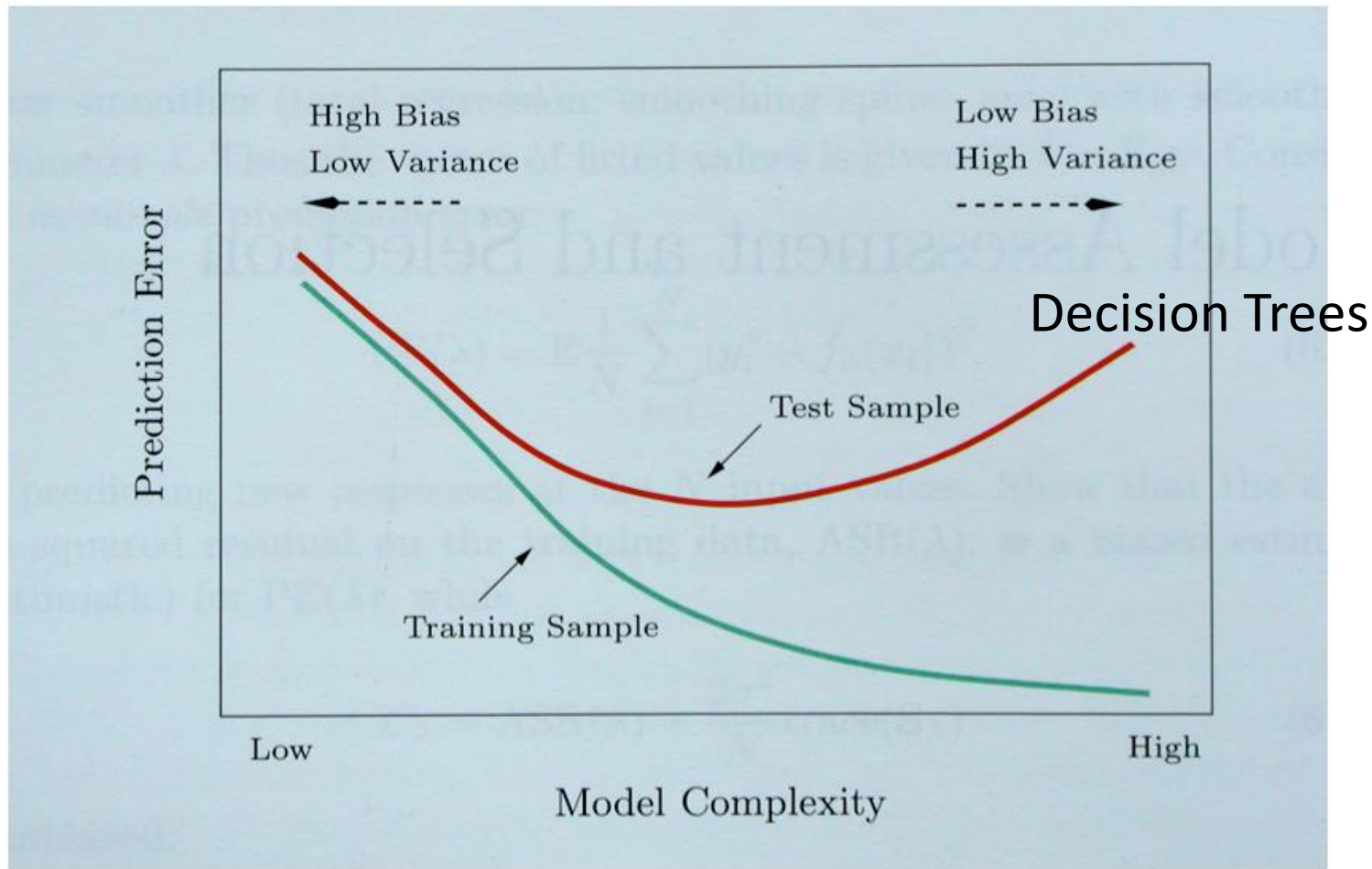


- Split to reduce MSE
- Predict average response of all training data at each leaf

# Outline

- Ensemble learning
  - Combine multiple classifiers to reduce model variance and improve accuracy
- Bagging
  - Bootstrap samples
  - Random Forests
- Boosting
  - AdaBoost

# Bias/Variance Tradeoff



Hastie, Tibshirani, Friedman "Elements of Statistical Learning" 2001

How to reduce variance of single decision tree?

# Ensemble Learning

Consider a set of classifiers  $h_1, \dots, h_L$

**Idea:** construct a classifier  $H(\mathbf{x})$  that combines the individual decisions of  $h_1, \dots, h_L$

- e.g., could have the member classifiers vote, or
- e.g., could use different members for different regions of the instance space

Successful ensembles require **diversity**

- Classifiers should make different mistakes
- Can have different types of base learners

# Build Ensemble Classifiers

- Basic idea
  - Build different “experts”, and let them vote
- Advantages
  - Improve predictive performance
  - Easy to implement
  - No too much parameter tuning
- Disadvantages
  - The combined classifier is not transparent and interpretable
  - Not a compact representation



# Practical Applications

**Goal:** predict how a user will rate a movie

- Based on the user's ratings for other movies
- and other peoples' ratings
- with no other information about the movies



This application is called “collaborative filtering”

**Netflix Prize:** \$1M to the first team to do 10% better than Netflix' system (2007-2009)

**Winner:** BellKor's Pragmatic Chaos – an ensemble of more than 800 rating systems

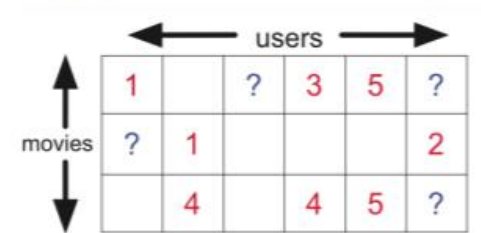
# Netflix Prize

Machine learning competition with a \$1 million prize

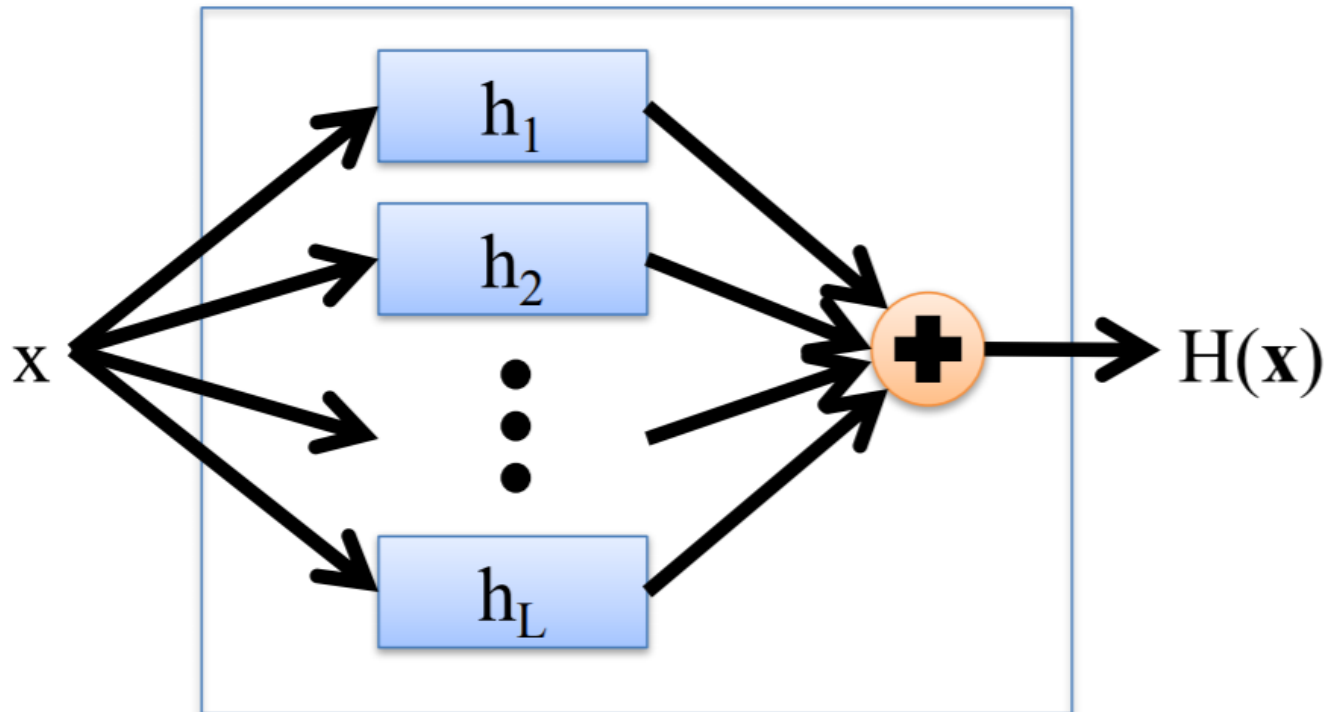
## Leaderboard

Display top 20 leaders.

Rank	Team Name	Best Score	% Improvement	Last Submit Time
1	<a href="#">The Ensemble</a>	0.8553	10.10	2009-07-26 18:38:22
2	<a href="#">BellKor in BigChaos</a>	0.8554	10.09	2009-07-26 18:18:28
<b>Grand Prize - RMSE &lt;= 0.8563</b>				
3	<a href="#">Grand Prize Team</a>	0.8571	9.91	2009-07-24 13:07:49
4	<a href="#">Opera Solutions and Vandelay United</a>	0.8573	9.89	2009-07-25 20:05:52
5	<a href="#">Vandelay Industries I</a>	0.8579	9.83	2009-07-26 02:49:53
6	<a href="#">PragmaticTheory</a>	0.8582	9.80	2009-07-12 15:09:53
7	<a href="#">BellKor in BigChaos</a>	0.8590	9.71	2009-07-26 12:57:25
8	<a href="#">Dace</a>	0.8603	9.58	2009-07-24 17:18:43
9	<a href="#">Opera Solutions</a>	0.8611	9.49	2009-07-26 18:02:08
10	<a href="#">BellKor</a>	0.8612	9.48	2009-07-26 17:19:11
11	<a href="#">BigChaos</a>	0.8613	9.47	2009-06-23 23:06:52
12	<a href="#">Feeds2</a>	0.8613	9.47	2009-07-24 20:06:46
<b>Progress Prize 2008 - RMSE = 0.8616 - Winning Team: BellKor in BigChaos</b>				
13	<a href="#">xianliang</a>	0.8633	9.26	2009-07-21 02:04:40
14	<a href="#">Gravity</a>	0.8634	9.25	2009-07-26 15:58:34
15	<a href="#">Ces</a>	0.8642	9.17	2009-07-25 17:42:38
16	<a href="#">Invisible Ideas</a>	0.8644	9.14	2009-07-20 03:26:12
17	<a href="#">Just a guy in a garage</a>	0.8650	9.08	2009-07-22 14:10:42
18	<a href="#">Craig Carmichael</a>	0.8656	9.02	2009-07-25 16:00:54
19	<a href="#">J.Dennis Su</a>	0.8658	9.00	2009-03-11 09:41:54
20	<a href="#">acmehill</a>	0.8659	8.99	2009-04-16 06:29:35
<b>Progress Prize 2007 - RMSE = 0.8712 - Winning Team: KorBell</b>				
<b>Cinematch score on quiz subset - RMSE = 0.9514</b>				

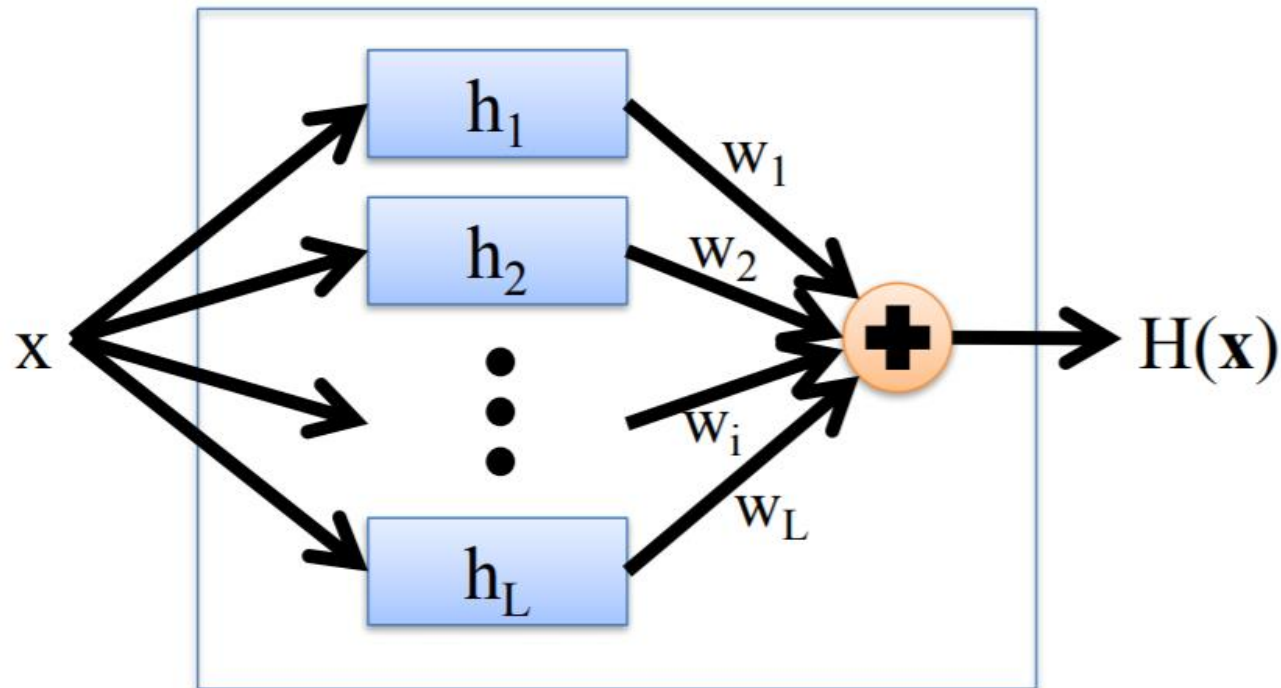


# Combining Classifiers: Averaging



- Final hypothesis is a simple vote of the members

# Combining Classifiers: Weighted Averaging



- Coefficients of individual members are trained using a validation set

# Reduce error

- Suppose there are 25 base classifiers
- Each classifier has error rate,  $\varepsilon = 0.35$
- Assume independence among classifiers
- Probability that the ensemble classifier makes a wrong prediction:

# Reduce Variance

- **Averaging** reduces variance:

$$\text{Var}(\bar{X}) = \frac{\text{Var}(X)}{N} \quad (\text{when predictions are independent})$$

Average models to reduce model variance

One problem:

only one training set

where do multiple models come from?

**Assuming models are independent!**

# How to Achieve Diversity

- Avoid overfitting
  - Vary the training data
- Features are noisy
  - Vary the set of features

Two main ensemble learning methods

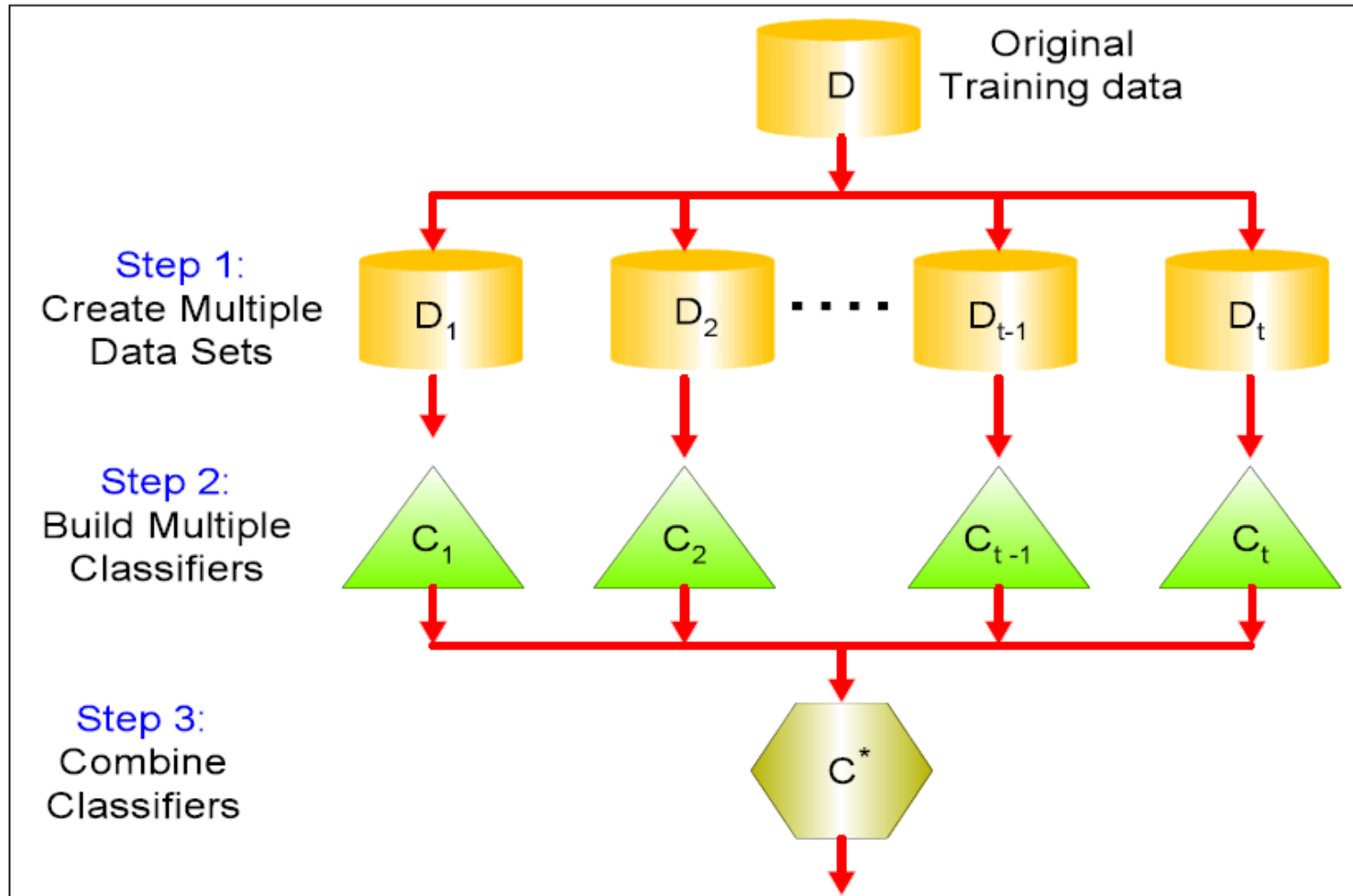
- **Bagging** (e.g., Random Forests)
- **Boosting** (e.g., AdaBoost)

# Bagging

- Leo Breiman (1994)
- Take repeated **bootstrap samples** from training set  $D$
- *Bootstrap sampling*: Given set  $D$  containing  $N$  training examples, create  $D'$  by drawing  $N$  examples at random **with replacement** from  $D$ .
- Bagging:
  - Create  $k$  bootstrap samples  $D_1 \dots D_k$ .
  - Train distinct classifier on each  $D_i$ .
  - Classify new instance by majority vote / average.



# General Idea



Majority Votes

# Example of Bagging

- Sampling with replacement

Training Data  
↙

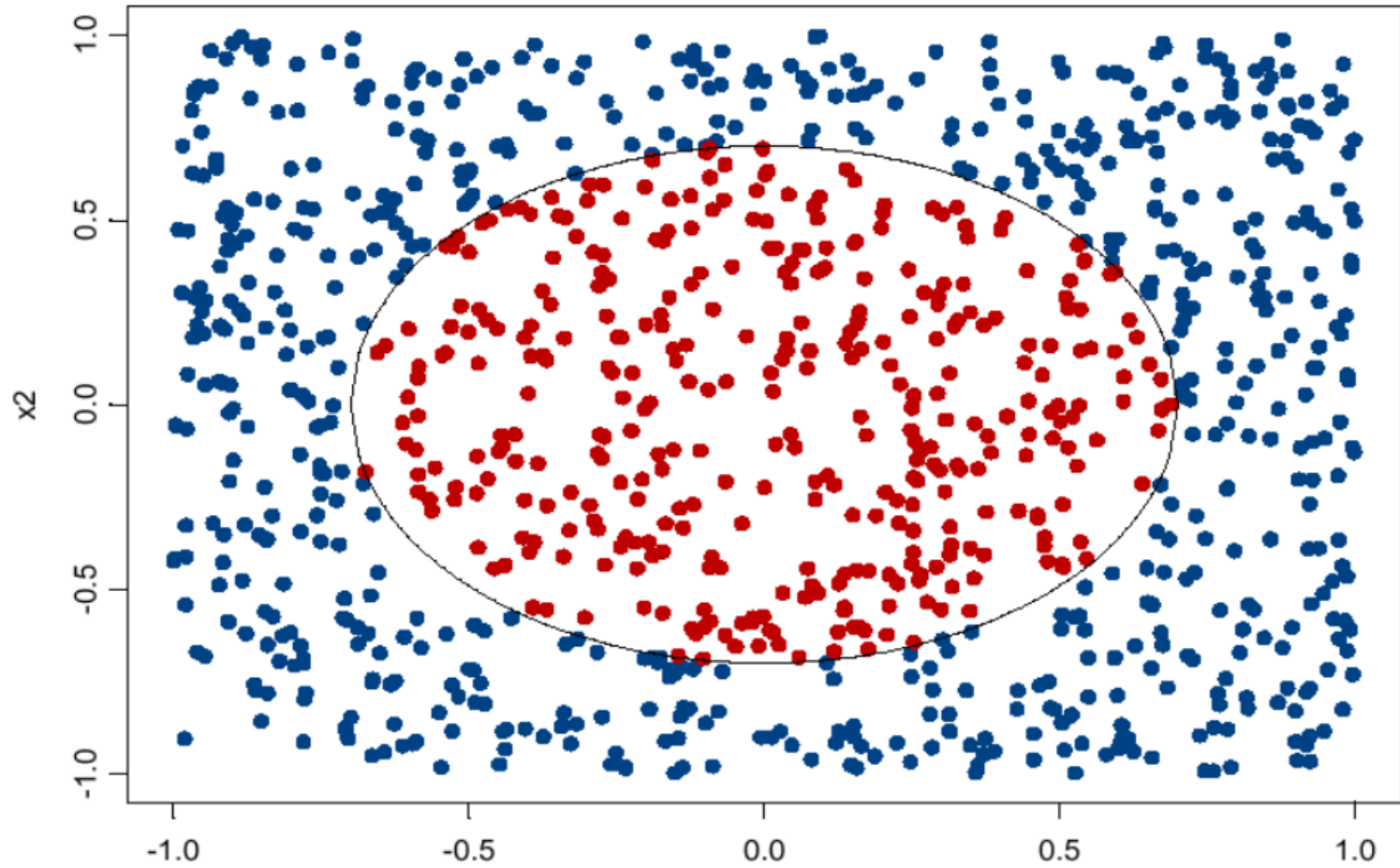
Data ID	1	2	3	4	5	6	7	8	9	10
Original Data	1	2	3	4	5	6	7	8	9	10
Bagging (Round 1)	7	8	10	8	2	5	10	10	5	9
Bagging (Round 2)	1	4	9	1	2	3	2	7	3	2
Bagging (Round 3)	1	8	5	10	5	5	9	6	3	7

- Sample each training point with probability  $1/n$
- **Out-Of-Bag (OOB) observation**: point not in sample
  - For each point: prob  $(1-1/n)^n$
  - About  $1/3$  of data
  - OOB error: error on OOB samples
- **OOB average error**
  - Compute across all models in Ensemble
  - Use instead of Cross-Validation error

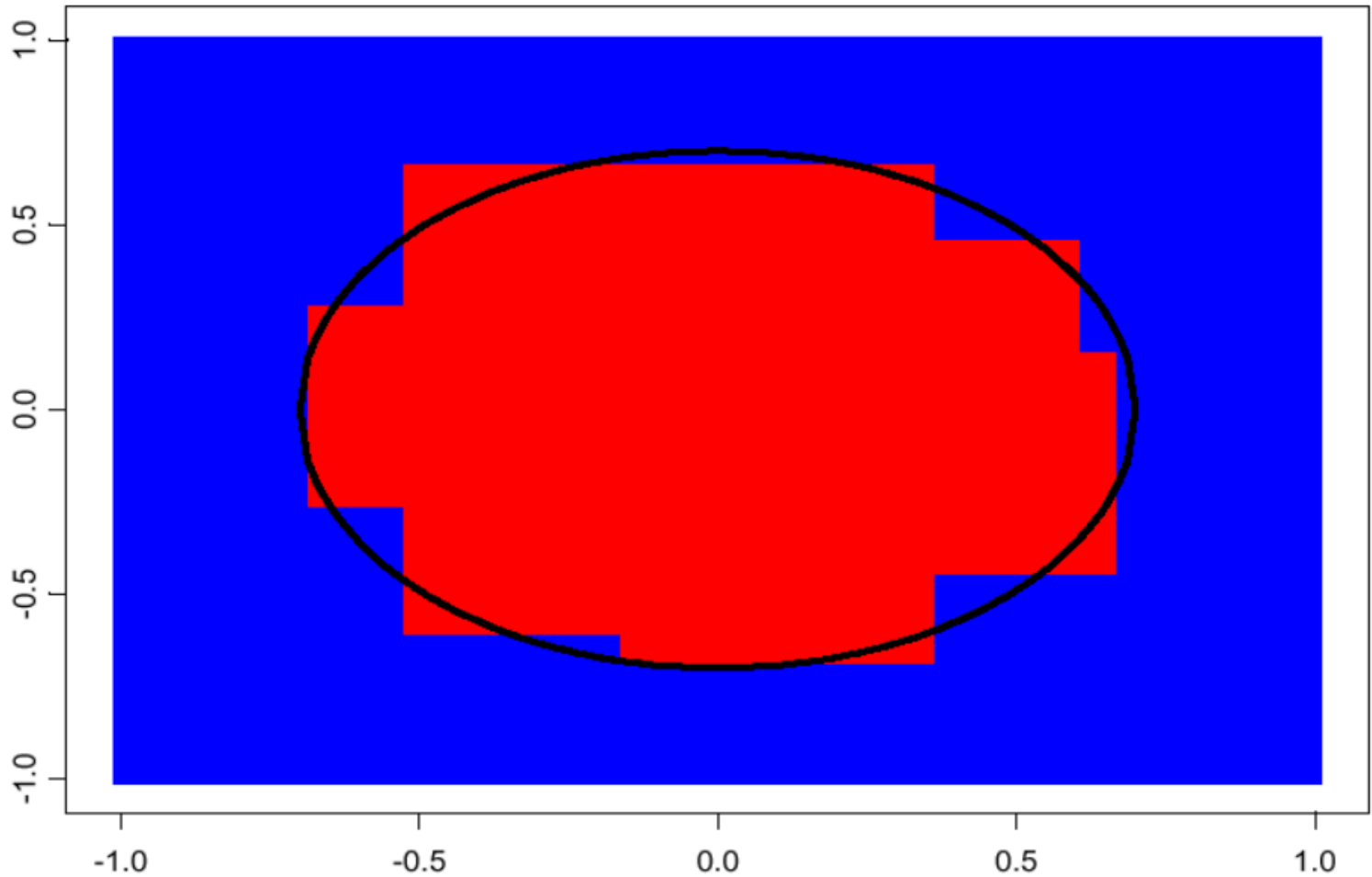
# Bagging

- Can be applied to multiple classification models
- Very successful for decision trees
  - Decision trees have high variance
  - Don't prune the individual trees, but grow trees to full extent
  - Precision accuracy of decision trees improved substantially
- OOB average error used instead of Cross Validation

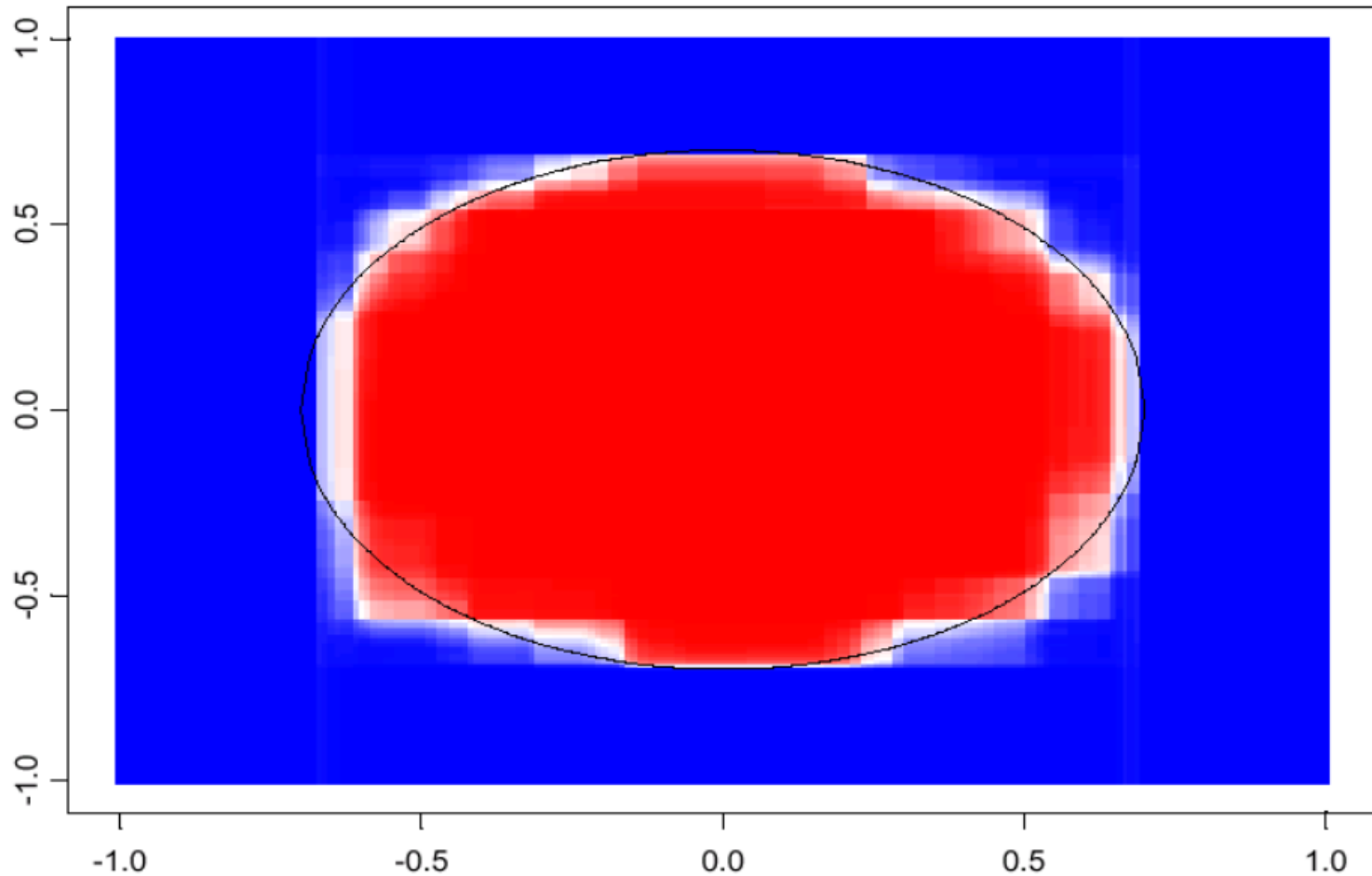
# Example Distribution



# Decision Tree Decision Boundary



# 100 Bagged Trees



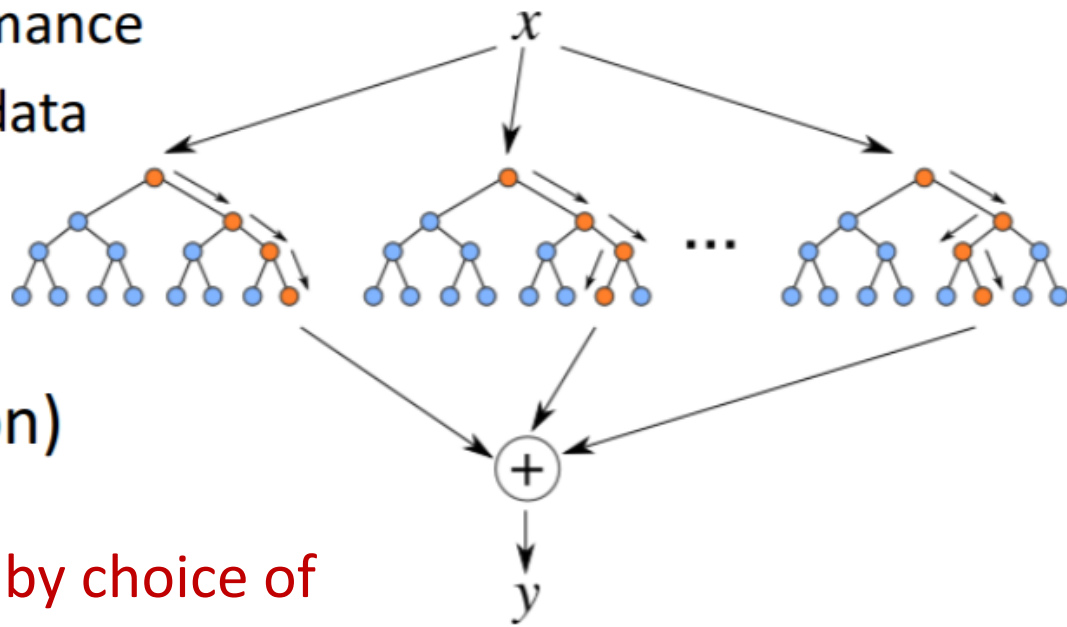
shades of blue/red indicate strength of vote for particular classification

# Random Forests

- Ensemble method specifically designed for decision tree classifiers
- Introduce two sources of randomness: “Bagging” and “Random input vectors”
  - **Bagging method**: each tree is grown using a bootstrap sample of training data
  - **Random vector method**: **At each node**, best split is chosen from a random sample of  $m$  attributes instead of all attributes

# Random Forests

- Construct decision trees on bootstrap replicas
  - Restrict the node decisions to a small subset of features picked randomly for each node
- Do not prune the trees
  - Estimate tree performance on out-of-bootstrap data
- Average the output of all trees (or choose mode decision)



Trees are de-correlated by choice of random subset of features



# Random Forest Algorithm

1. For  $b = 1$  to  $B$ :
  - (a) Draw a **bootstrap sample**  $\mathbf{Z}^*$  of size  $N$  from the training data.
  - (b) Grow a random-forest tree  $T_b$  to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size  $n_{min}$  is reached.
    - i. Select  **$m$  variables at random** from the  $p$  variables.
    - ii. Pick the best variable/split-point among the  $m$ .
    - iii. Split the node into two daughter nodes.
2. Output the ensemble of trees  $\{T_b\}_1^B$ .

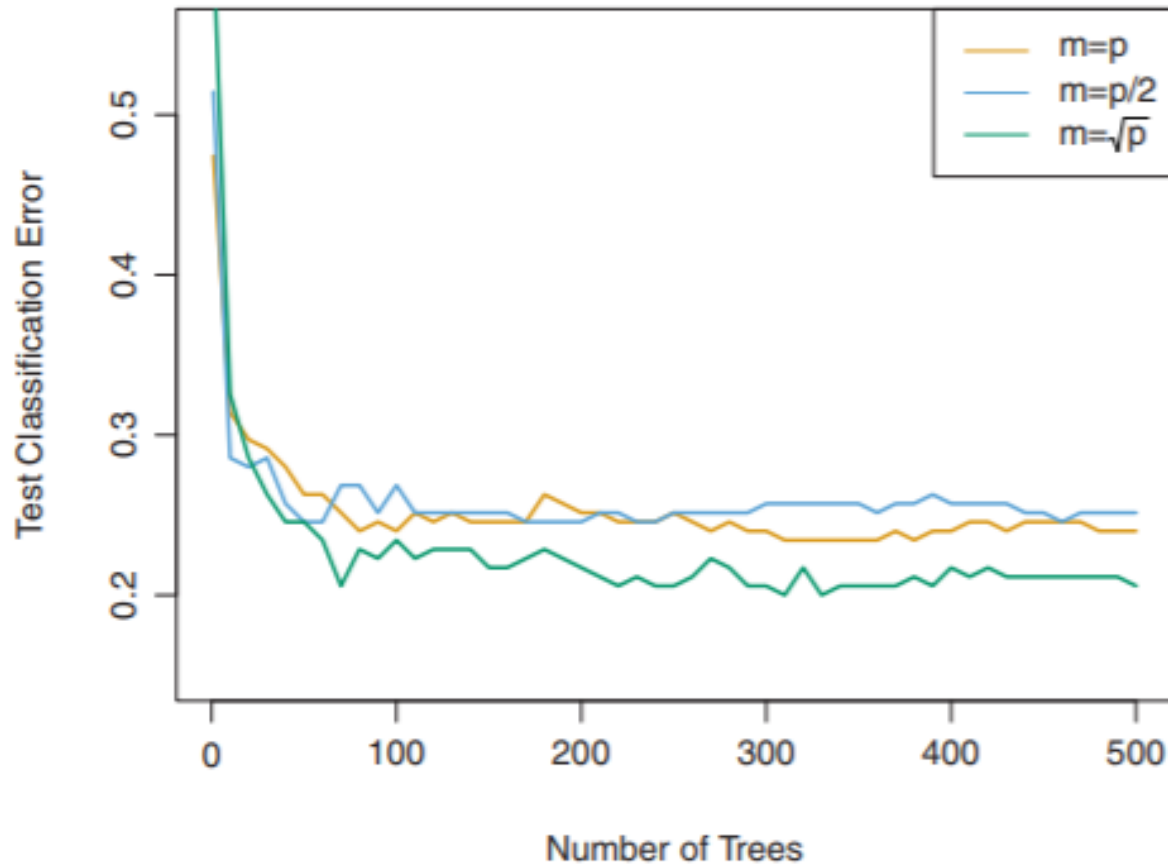
To make a prediction at a new point  $x$ :

*Regression:*  $\hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$ .

*Classification:* Let  $\hat{C}_b(x)$  be the class prediction of the  $b$ th random-forest tree. Then  $\hat{C}_{\text{rf}}^B(x) = \text{majority vote } \{\hat{C}_b(x)\}_1^B$ .

**If  $m=p$ , this is equivalent to Bagging**

# Effect of Number of Predictors



- $p$  = total number of predictors;  $m$  = predictors chosen in each split
- Random Forests uses  $m = \sqrt{p}$

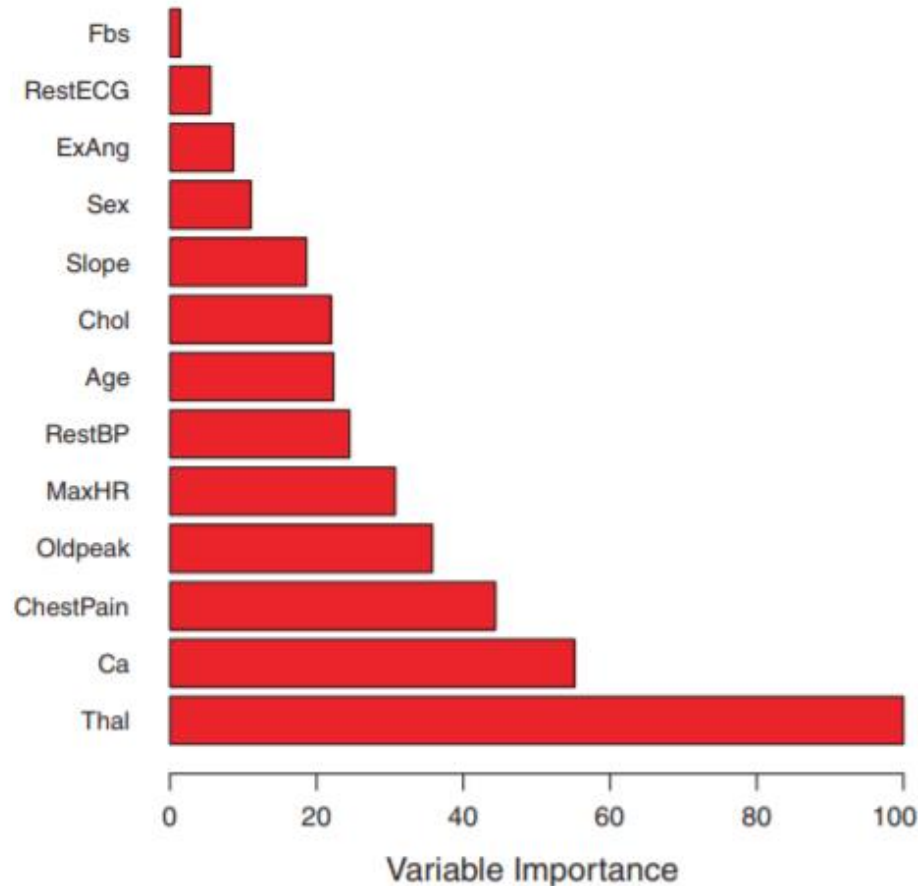
# Variable Importance

- Ensemble of trees loses somewhat interpretability of decision trees
- Which variables contribute mostly to prediction?
- Random Forests computes a Variable Importance metric

# Gini index

- Take a node of decision tree
- Let  $p_i$  be the fraction of examples from class  $i$
- Measures the “purity” of the node
  - If node has most examples from one class, Gini index is low
- What is the probability that a random example is mis-classified at that node?
  - $\sum_{i=1}^k p_i(1 - p_i)$
- Close to Information Gain

# Variable Importance Plots



**FIGURE 8.9.** A variable importance plot for the **Heart** data. Variable importance is computed using the mean decrease in Gini index, and expressed relative to the maximum.

# How to Achieve Diversity

- Avoid overfitting
  - Vary the training data
- Features are noisy
  - Vary the set of features

Two main ensemble learning methods

- **Bagging** (e.g., Random Forests)
- **Boosting** (e.g., AdaBoost)

# Acknowledgements

- Slides made using resources from:
  - Andrew Ng
  - Eric Eaton
  - David Sontag
- Thanks!