

# DS 4400

## Machine Learning and Data Mining I

Alina Oprea  
Associate Professor  
Khoury College of Computer Science  
Northeastern University

October 27 2020

# Announcements

- HW 3 is due on Thu, Oct. 29
- Project proposal
  - Due on Monday, Nov. 2
  - Team of 2
  - Resources and example projects on Piazza

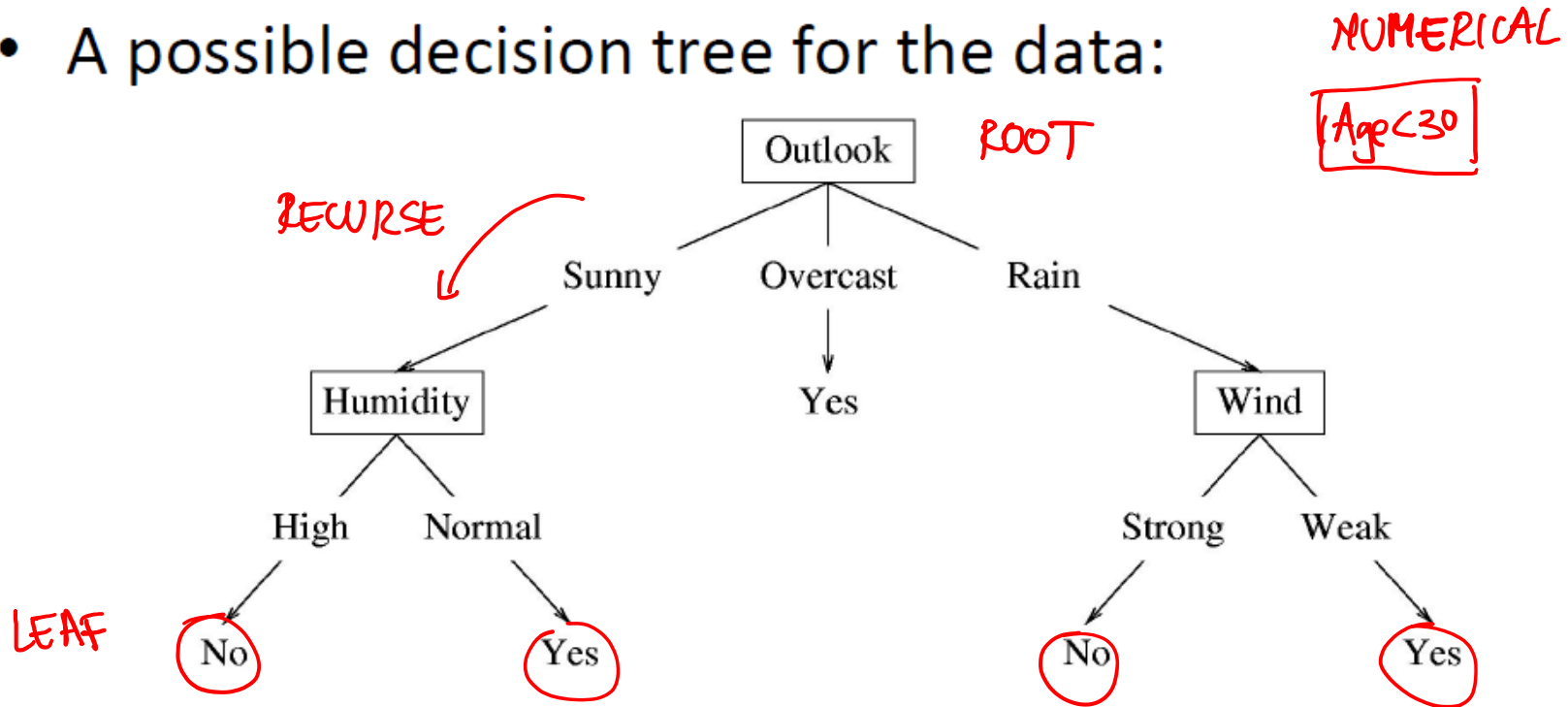
• EXAM- NOV. 19

# Outline

- Decision trees
  - Information gain
  - Learning decision trees
- Ensemble learning
  - Combine multiple classifiers to reduce model variance and improve accuracy
- Bagging
  - Bootstrap samples
  - Random Forests

# Decision Tree

- A possible decision tree for the data:



- Each internal node: test one attribute  $X_i$
- Each branch from a node: selects one value for  $X_i$
- Each leaf node: predict  $Y$  (or  $p(Y \mid x \in \text{leaf})$ )

# Information Gain

SHANNON

X = College Major

Y = Likes "Gladiator"

**Definition of Information Gain:**

$IG(Y|X) =$  **I must transmit Y.**

**How many bits on average would it save me if both ends of the line knew X?**

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

$$IG(Y|X) = H(Y) - H(Y|X)$$

ENTROPY      COND. ENTROPY

**Example:**

- $H(Y) = 1$

- $H(Y|X) = 0.5$

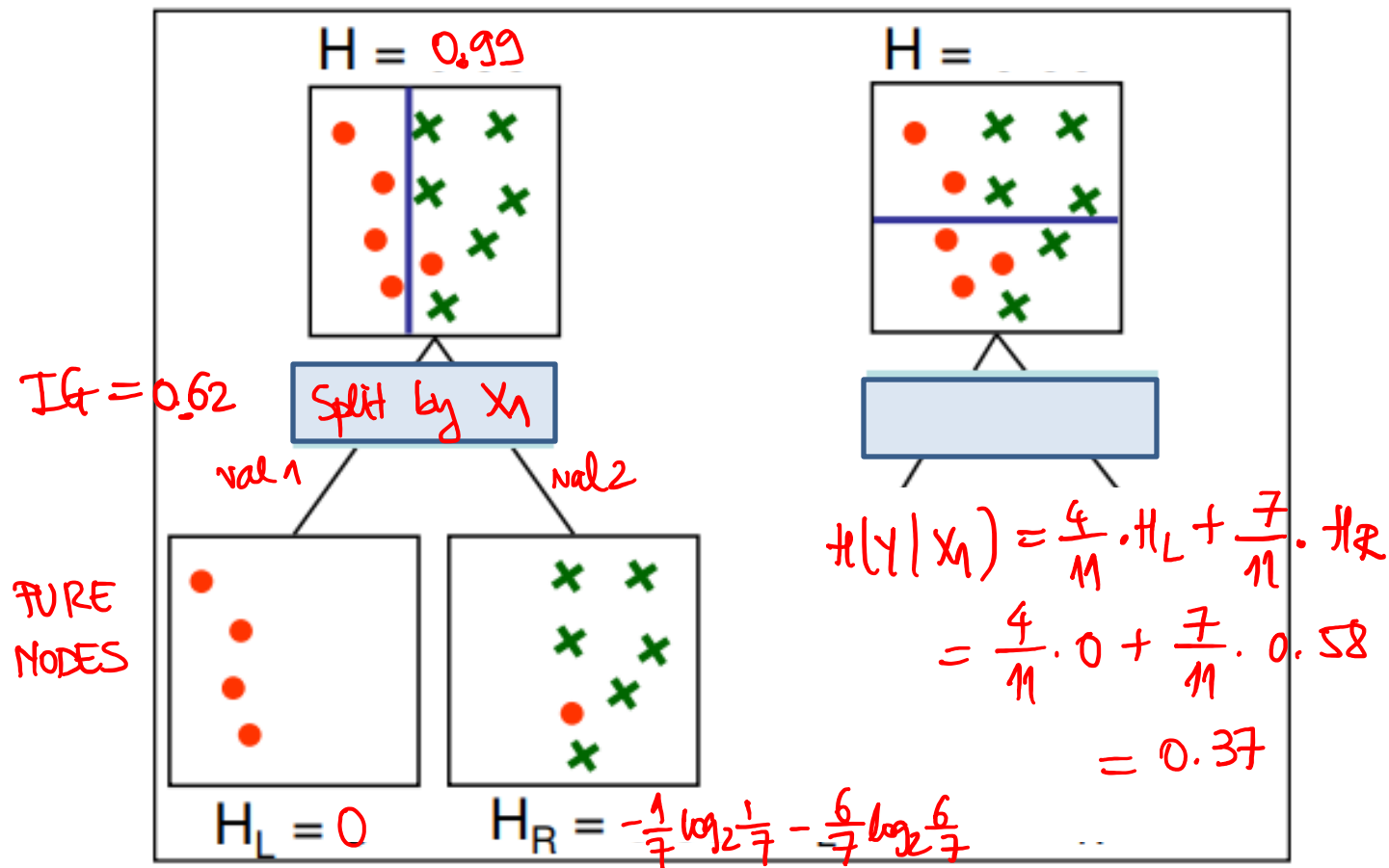
- Thus  $IG(Y|X) = 1 - 0.5 = 0.5$

$$H(Y|X) = \sum_{N \text{ values}} P[X=N] \cdot H(Y|X=N)$$

# Example Information Gain

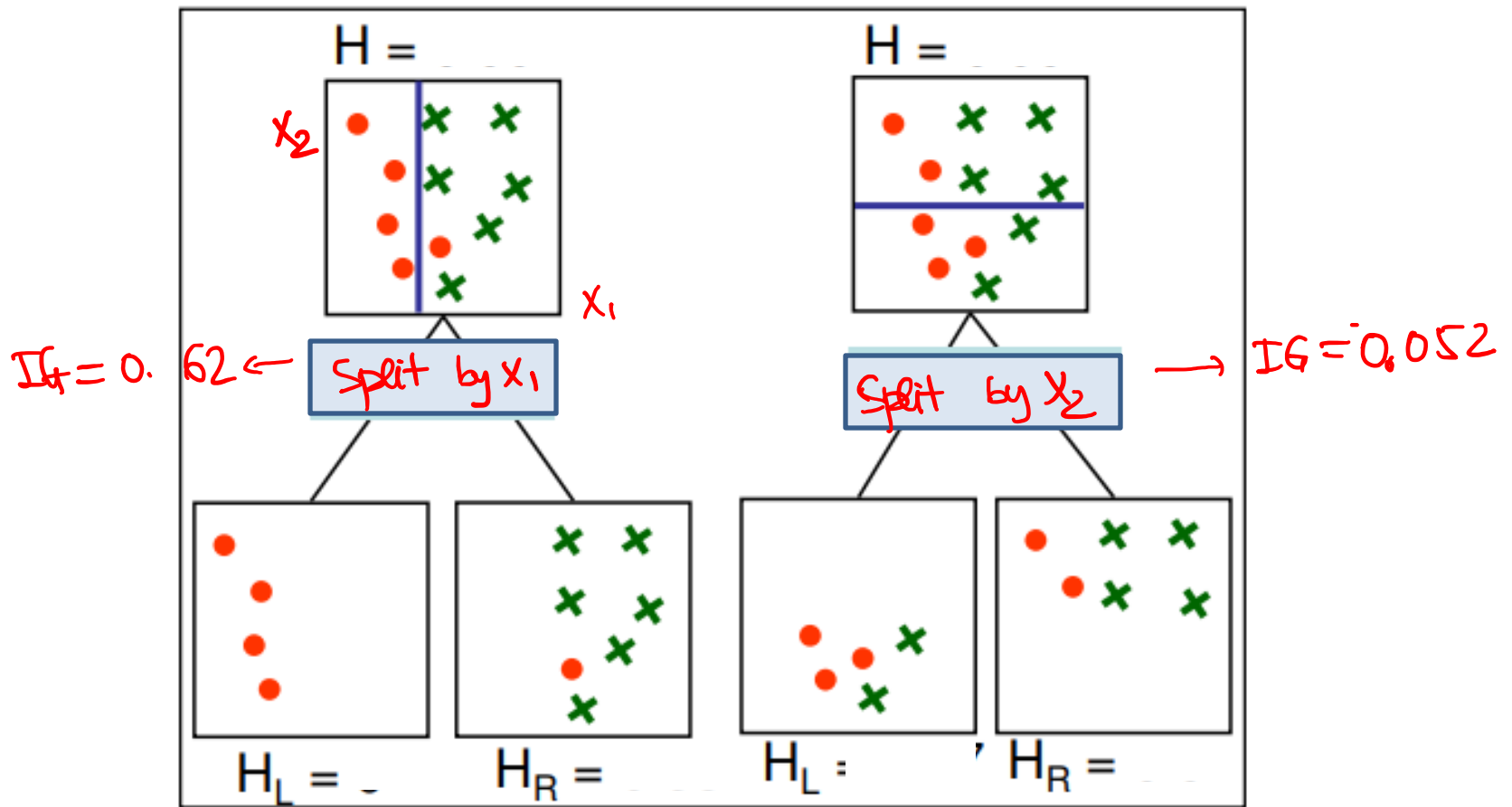
Y: (Red Green)  
5/11 6/11

$$H(Y) = -\frac{5}{11} \log_2 \frac{5}{11} - \frac{6}{11} \log_2 \frac{6}{11}$$



$$IG(Y|X_1) = H(Y) - H(Y|X_1) = 0.99 - 0.37 = 0.62$$

# Example Information Gain



PICK FEATURE max IG  
SELECT  $X_1$  - BETTER SPLIT

# Learning Decision Trees

ROOT: ALL FEATURES ARE AVAILABLE

- Start from empty decision tree
- Split on **next best attribute (feature)**
  - Use, for example, information gain to select attribute:

GINI INDEX

$$\boxed{\arg \max_i IG(X_i)} = \arg \max_i H(Y) - H(Y | X_i)$$

- Recurse : SUBSET OF FEATURES

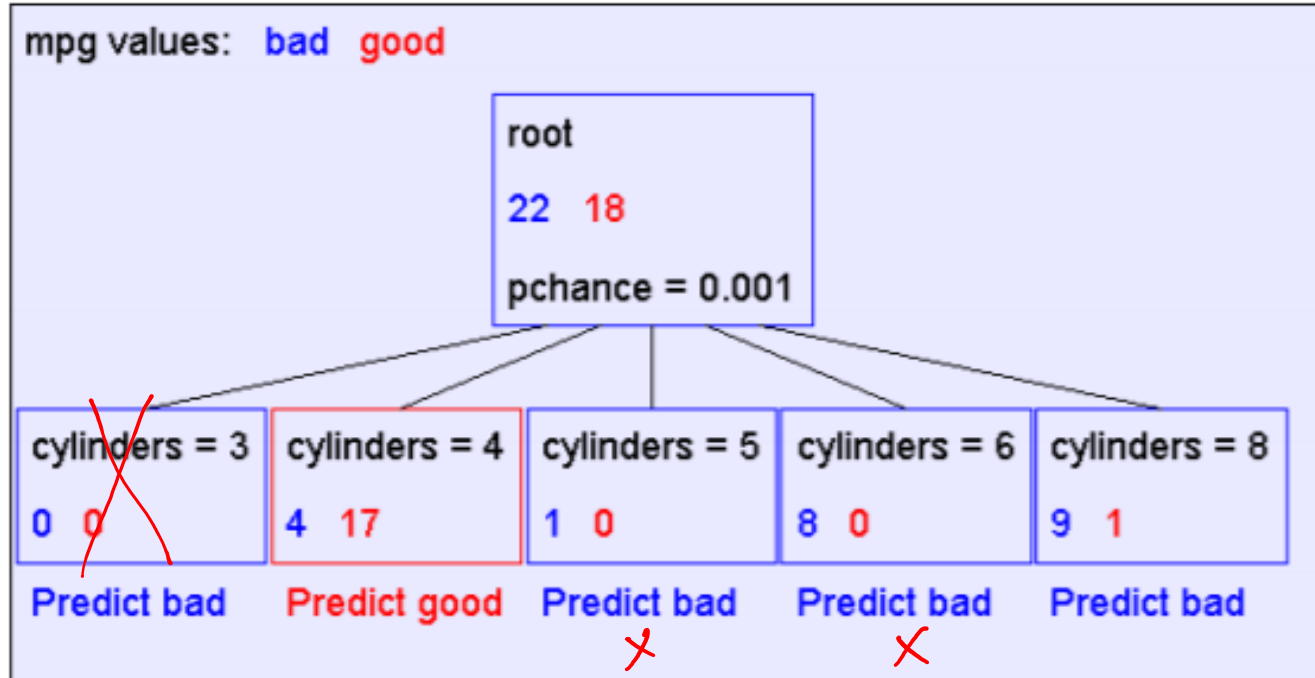
MAX. HEIGHT OF DECISION TREE: NUMBER OF FEATURES

ID3 algorithm uses Information Gain  
Information Gain reduces uncertainty on Y

OVERFITTING



# When to stop?

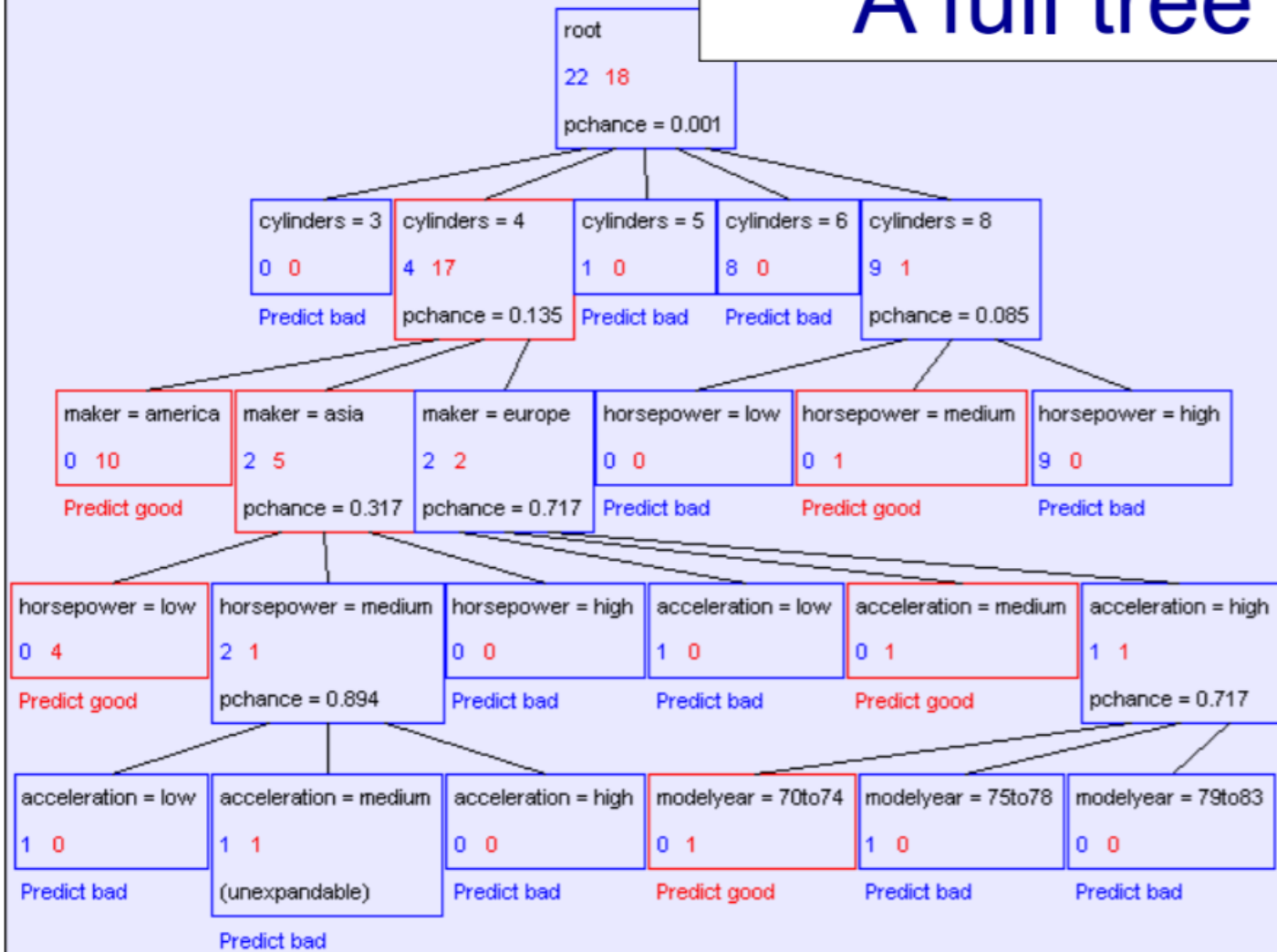


First split looks good! But, when do we stop?

# Full Tree

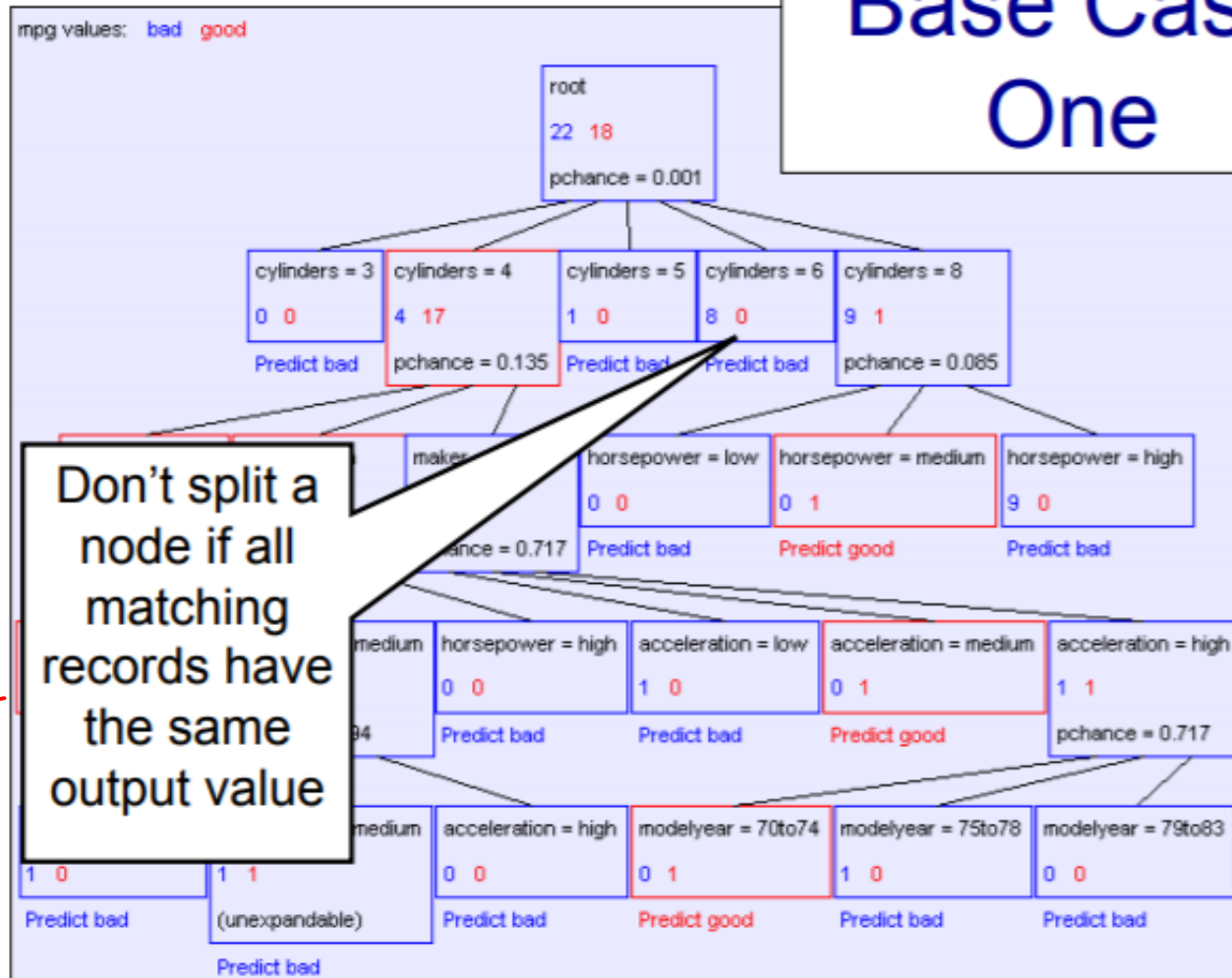
mpg values: bad good

## A full tree



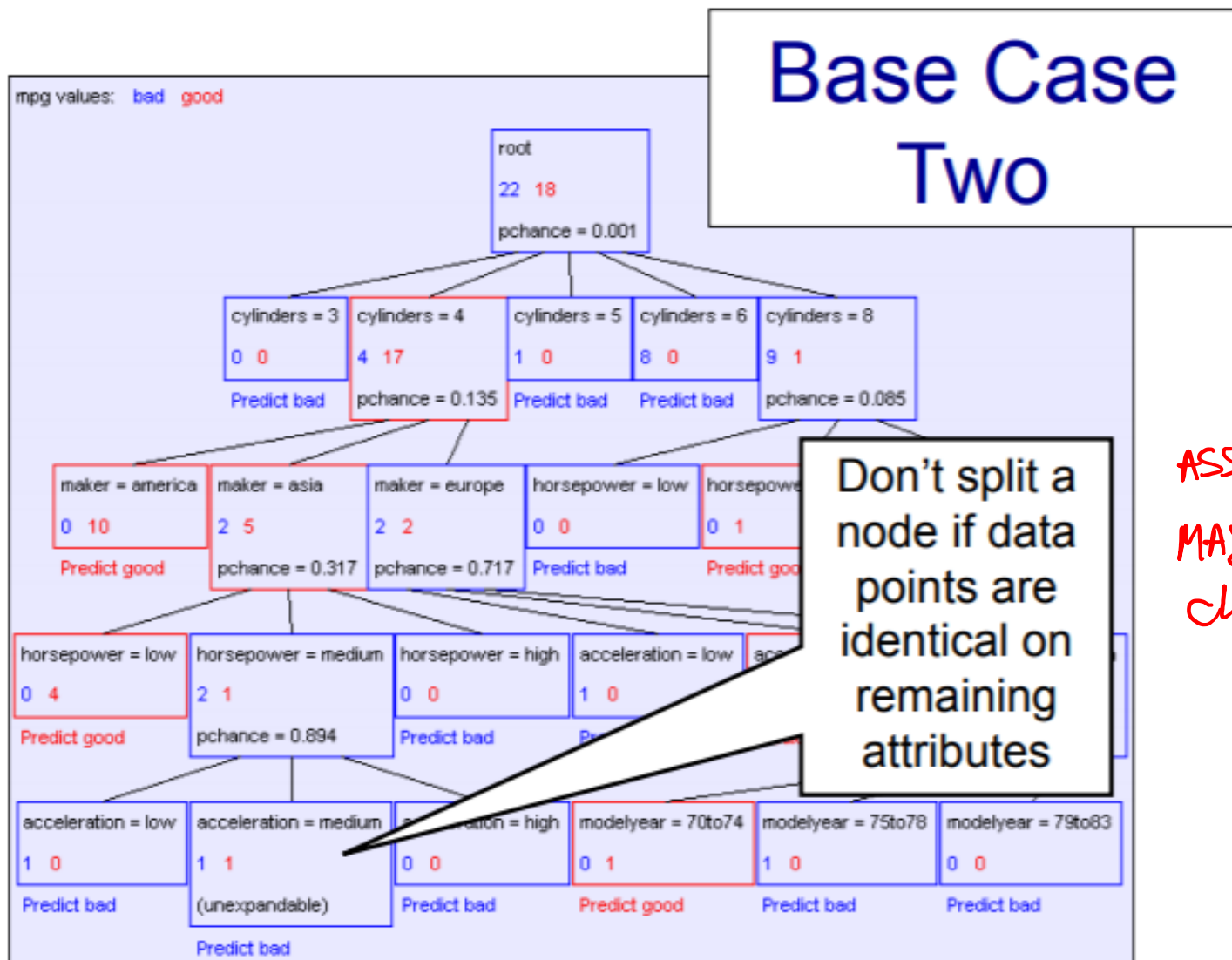
# Case 1

## Base Case One



ENTROPY = 0

# Case 2



# Decision Trees

TRAINING SET      LABELS  
↓                      ↓

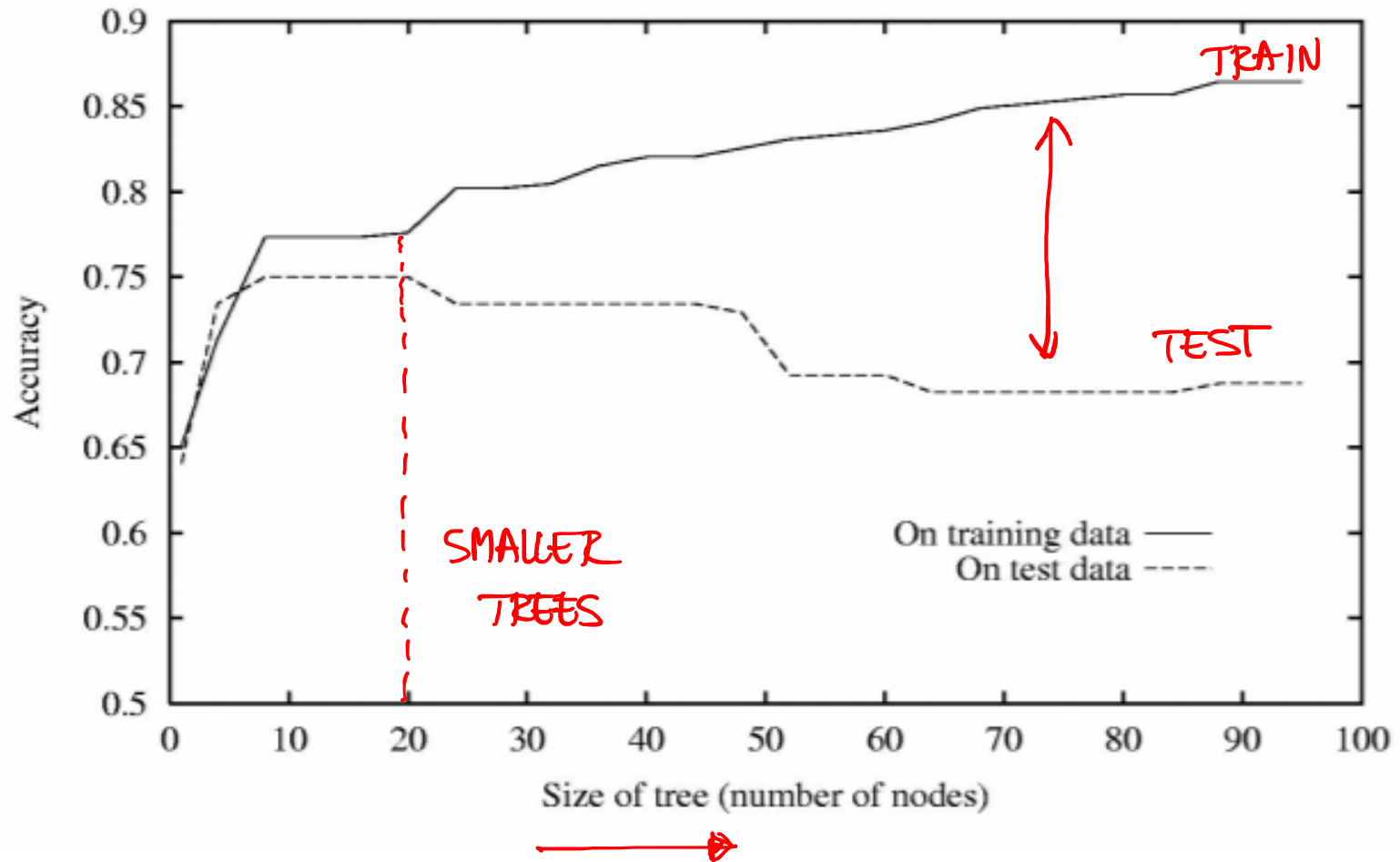
BuildTree(*DataSet*, *Output*)

- If all output values are the same in *DataSet*, return a leaf node that says “predict this unique output”      CASE 1
- If all input values are the same, return a leaf node that says “predict the majority output”      CASE 2
- Else find attribute  $X$  with highest Info Gain
- Suppose  $X$  has  $n_X$  distinct values (i.e.  $X$  has arity  $n_X$ ).
  - Create a non-leaf node with  $n_X$  children.
  - The  $i$ 'th child should be built by calling

BuildTree( $DS_i$ , *Output*)

Where  $DS_i$  contains the records in *DataSet* where  $X = i$ th value of  $X$ .

# Overfitting

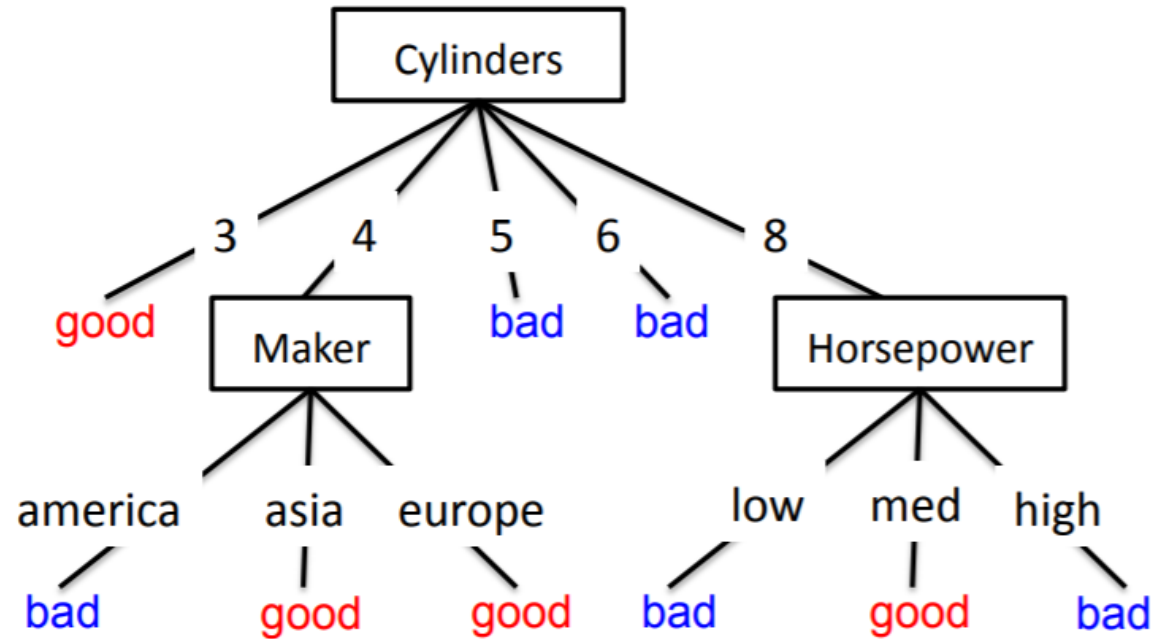


# Solutions against Overfitting

- Standard decision trees have no learning bias
  - Training set error is always zero!
    - (If there is no label noise)
  - Lots of variance
  - Must introduce some bias towards simpler trees
- Many strategies for picking simpler trees
  - Fixed depth
  - Minimum number of samples per leaf
- Pruning
  - Remove branches of the tree that increase error using cross-validation

# Interpretability

- Each internal node tests an attribute  $x_i$
- One branch for each possible attribute value  $x_i=v$
- Each leaf assigns a class  $y$
- To classify input  $x$ : traverse the tree from root to leaf, output the labeled  $y$

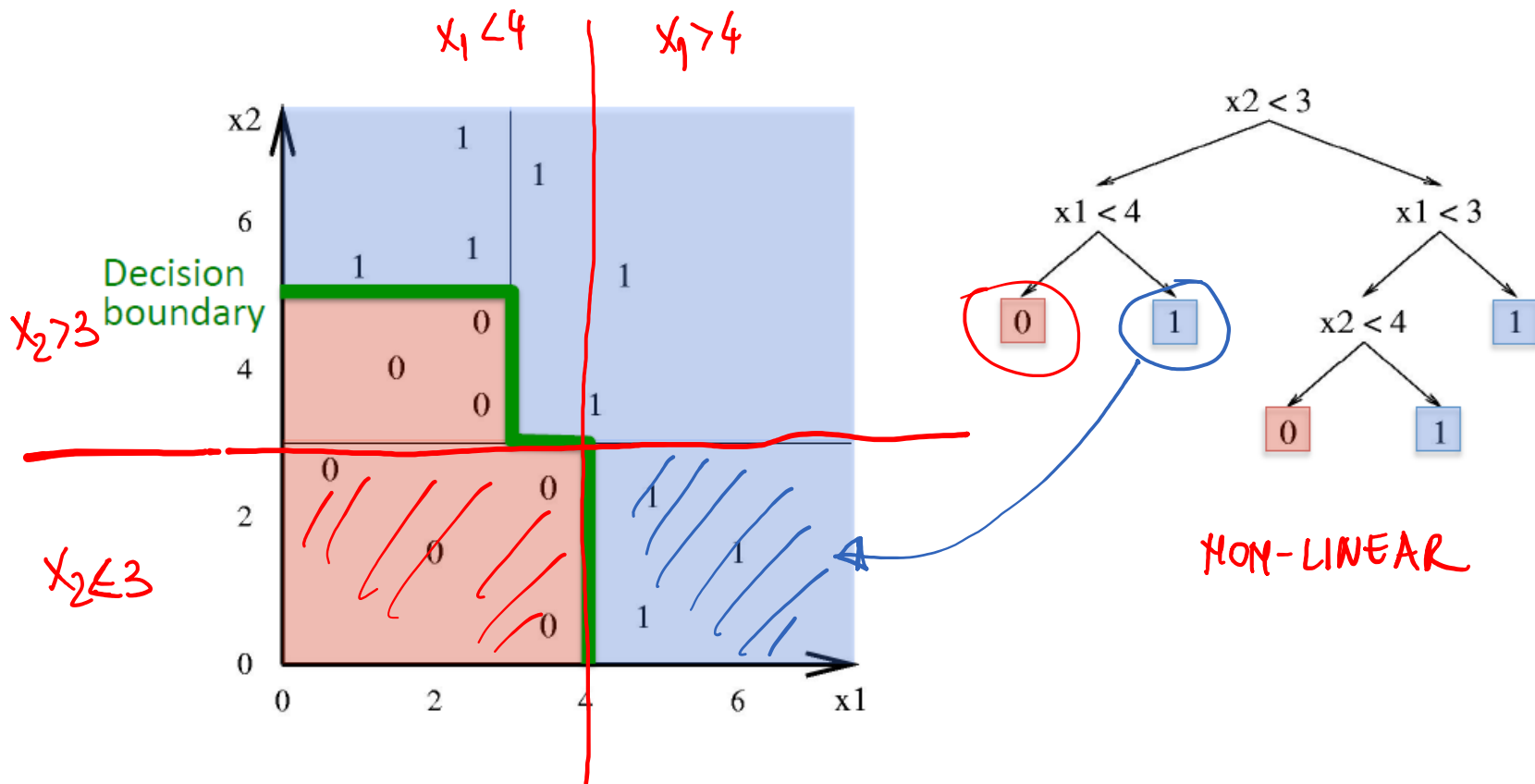


Human interpretable!

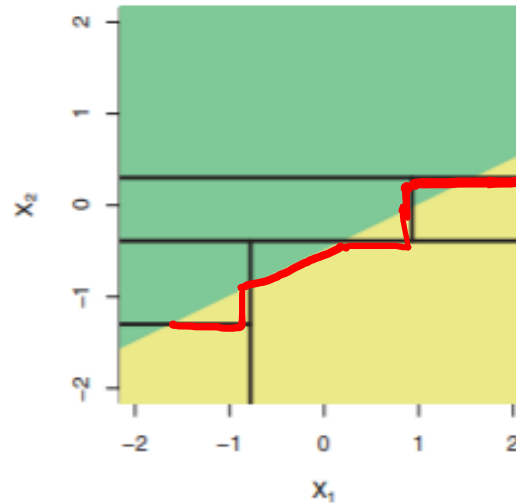
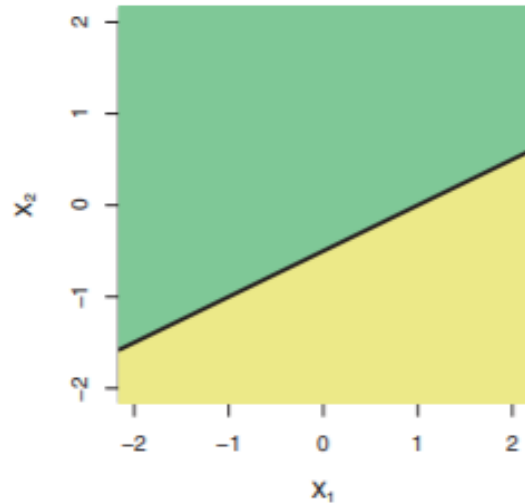


# Decision Boundary

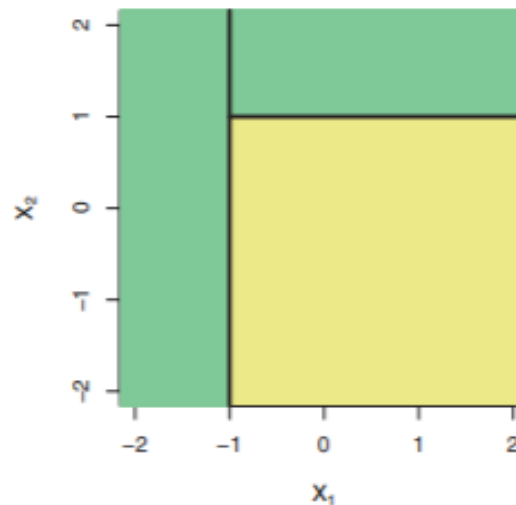
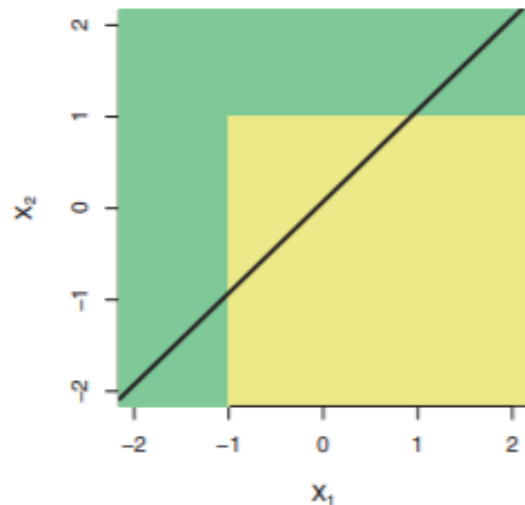
- Decision trees divide the feature space into axis-parallel (hyper-)rectangles
- Each rectangular region is labeled with one label



# Decision Trees vs Linear Models



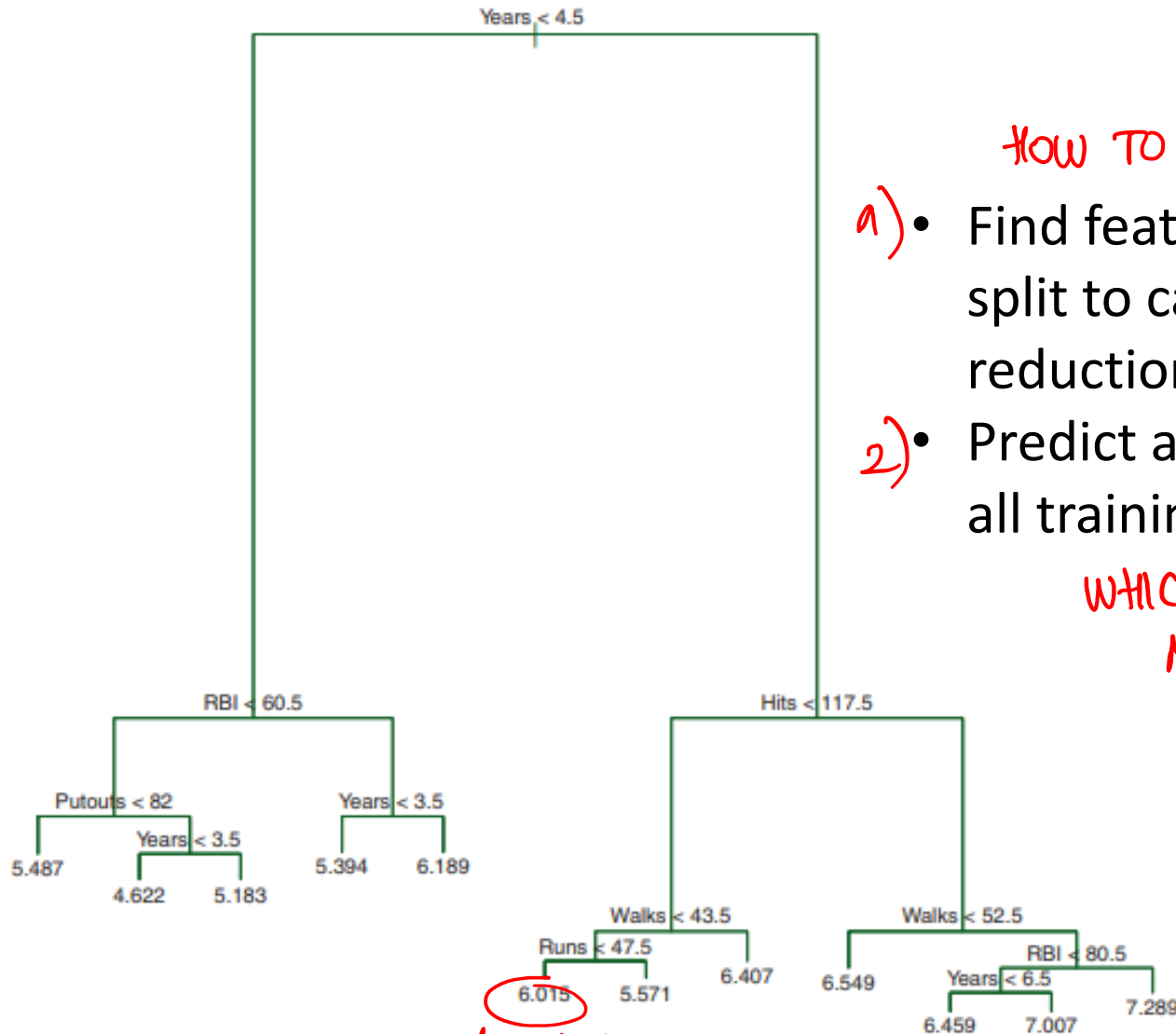
LINEARLY  
SEPARABLE



Linear model

Decision tree

# Regression Trees



## HOW TO SPLIT

- 1) Find feature and value to split to cause the maximum reduction in MSE
- 2) Predict average response of all training data at each leaf

WHICH PREDICTION TO MAKE AT EACH LEAF

10 TRAINING EXAMPLE

# 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

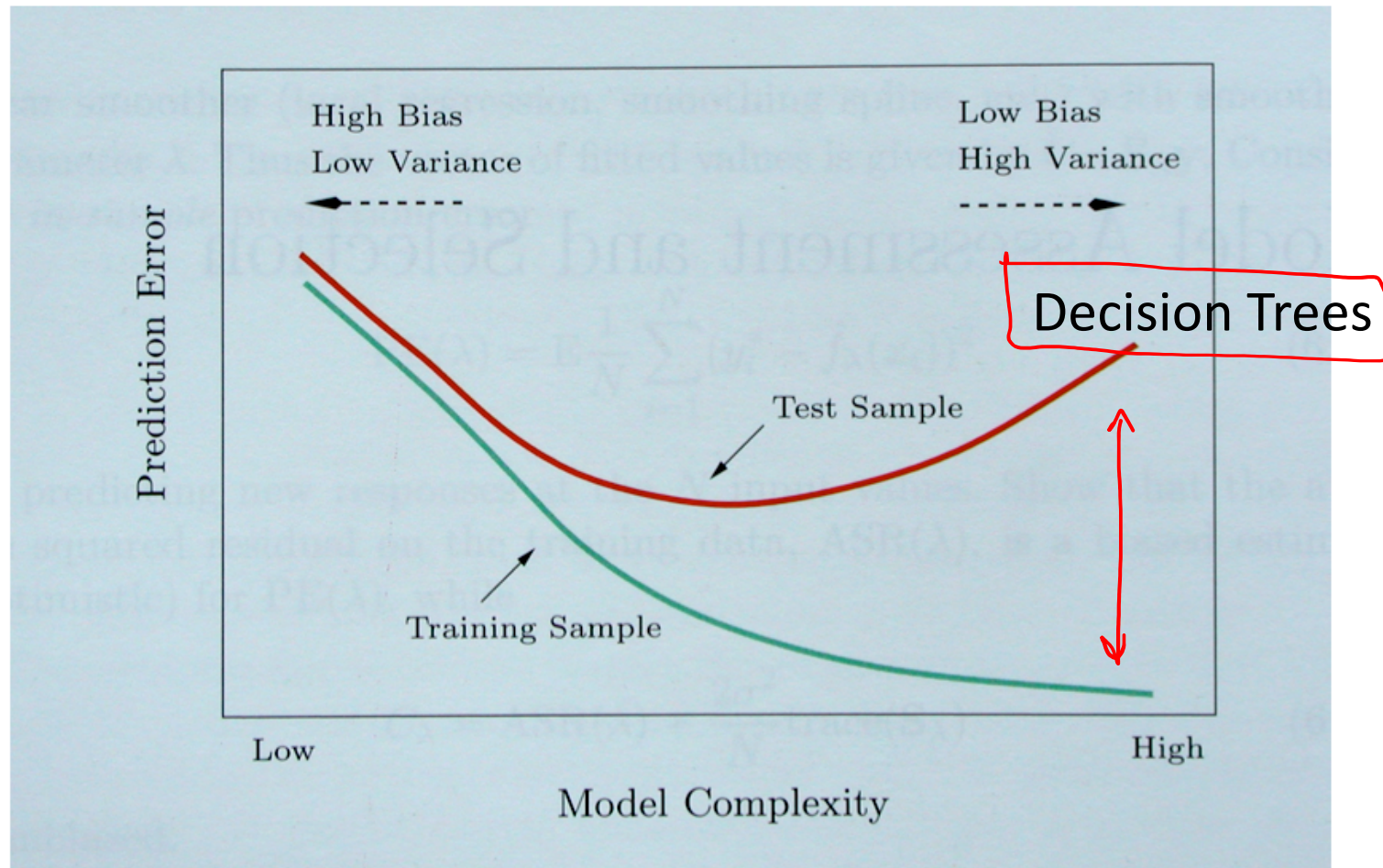
## Strengths

- 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

# 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



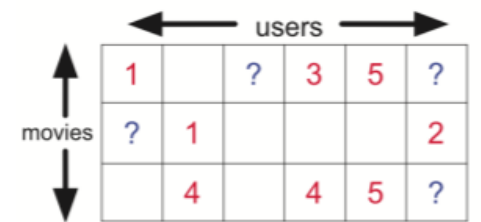
# Netflix Prize

Machine learning competition with a \$1 million prize

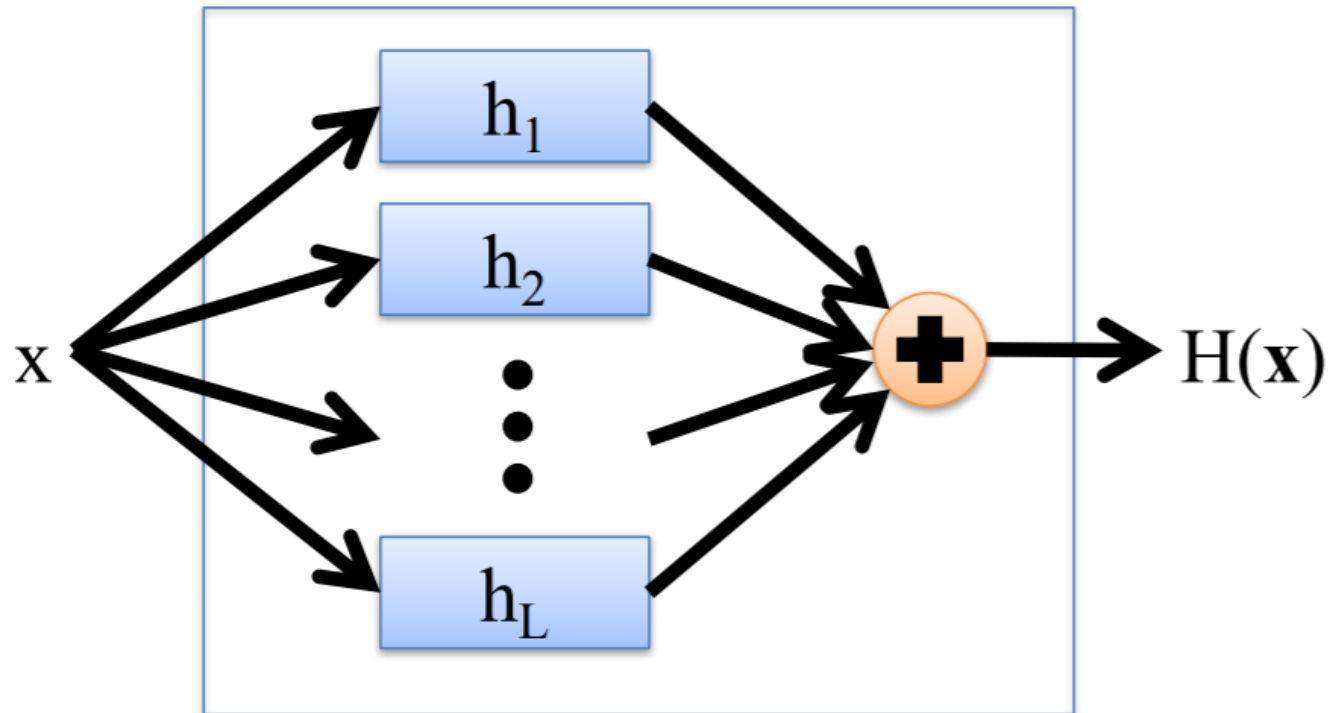
## Leaderboard

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>				

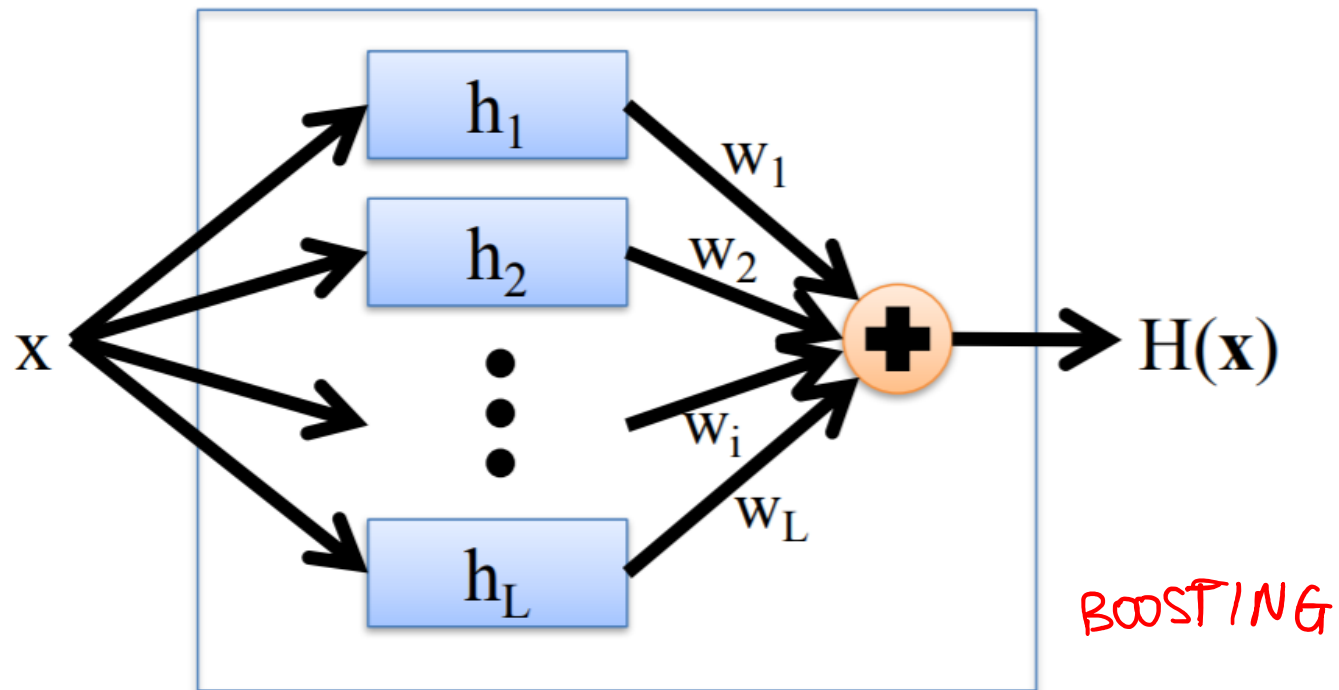


# Combining Classifiers: Averaging



- 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

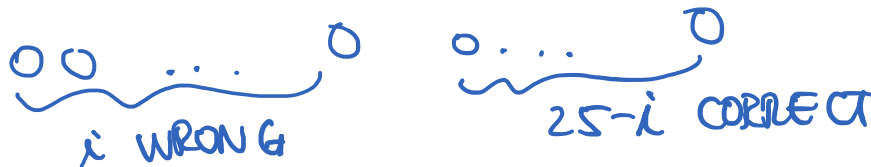
- Suppose there are 25 base classifiers
- Each classifier has error rate,  $\epsilon = 0.35$
- Assume independence among classifiers
- Probability that the ensemble classifier makes a wrong prediction:

$\geq 13$  MODELS MAKE AN ERROR

$$P[\text{ENSEMBLE MAKES ERROR}] = \sum_{i=13}^{25} \binom{25}{i} \epsilon^i \cdot (1-\epsilon)^{25-i} \approx 0.06$$

$i = \#$  CLASSIFIERS WRONG

$\binom{25}{i}$  SETS



# Reduce Variance

ENSEMBLES:

- REDUCE ERROR
- REDUCE VAR.

ASSUMPTION: MODELS ARE INDEPENDENT

REALITY: ONE TRAINING SET

HOW TO ACHIEVE DIVERSITY?

# Acknowledgements

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