

Decision Trees

Machine Learning – CSE546
 Carlos Guestrin (by Sameer Singh)
 University of Washington
 October 16, 2014

©Carlos Guestrin 2005-2013

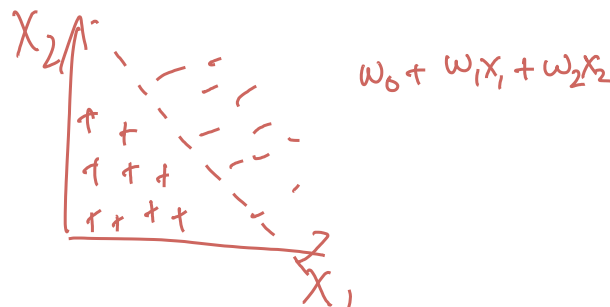
1

Linear separability

- A dataset is **linearly separable** iff there exists a **separating hyperplane**:

- Exists \mathbf{w} , such that:

- $w_0 + \sum_i w_i x_i > 0$; if $\mathbf{x}=\{x_1, \dots, x_k\}$ is a positive example
- $w_0 + \sum_i w_i x_i < 0$; if $\mathbf{x}=\{x_1, \dots, x_k\}$ is a negative example

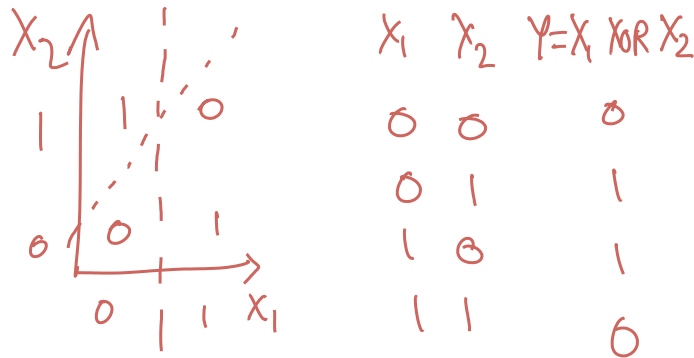


©Carlos Guestrin 2005-2013

2

Not linearly separable data

- Some datasets are **not linearly separable!**

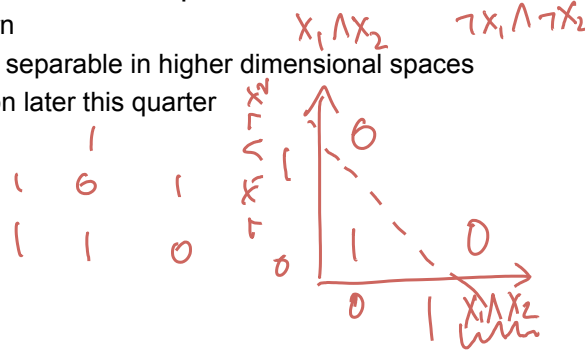


©Carlos Guestrin 2005-2013

3

Addressing non-linearly separable data – Option 1, non-linear features

- Choose non-linear features, e.g.,
 - Typical linear features: $w_0 + \sum_i w_i x_i$
 - Example of non-linear features:
 - Degree 2 polynomials, $w_0 + \sum_i w_i x_i + \sum_{ij} w_{ij} x_i x_j$
- Classifier $h_{\mathbf{w}}(\mathbf{x})$ still linear in parameters \mathbf{w}
 - As easy to learn
 - Data is linearly separable in higher dimensional spaces
 - More discussion later this quarter



©Carlos Guestrin 2005-2013

4

Addressing non-linearly separable data – Option 2, non-linear classifier

- Choose a classifier $h_w(\mathbf{x})$ that is non-linear in parameters \mathbf{w} , e.g.,
 - Decision trees, boosting, nearest neighbor, neural networks...
- More general than linear classifiers
- But, can often be harder to learn (non-convex/concave optimization required)
- But, but, often very useful
- (BTW. Later this quarter, we'll see that these options are not that different)

©Carlos Guestrin 2005-2013

5

A small dataset: Miles Per Gallon

Suppose we want to predict MPG

mpg	cylinders	displacement	horsepower	weight	acceleration	modelyear	maker
good	4	low	low	low	high	75to78	asia
bad	6	medium	medium	medium	medium	70to74	america
bad	4	medium	medium	medium	low	75to78	europa
bad	8	high	high	high	low	70to74	america
bad	6	medium	medium	medium	medium	70to74	america
bad	4	low	medium	low	medium	70to74	asia
bad	4	low	medium	low	low	70to74	asia
bad	8	high	high	high	low	75to78	america
.
.
.
.
bad	8	high	high	high	low	70to74	america
good	8	high	medium	high	high	79to83	america
bad	8	high	high	high	low	75to78	america
good	4	low	low	low	low	79to83	america
bad	6	medium	medium	medium	high	75to78	america
good	4	medium	low	low	low	79to83	america
good	4	low	low	medium	high	79to83	america
bad	8	high	high	high	low	70to74	america
good	4	low	medium	low	medium	75to78	europa
bad	5	medium	medium	medium	medium	75to78	europa

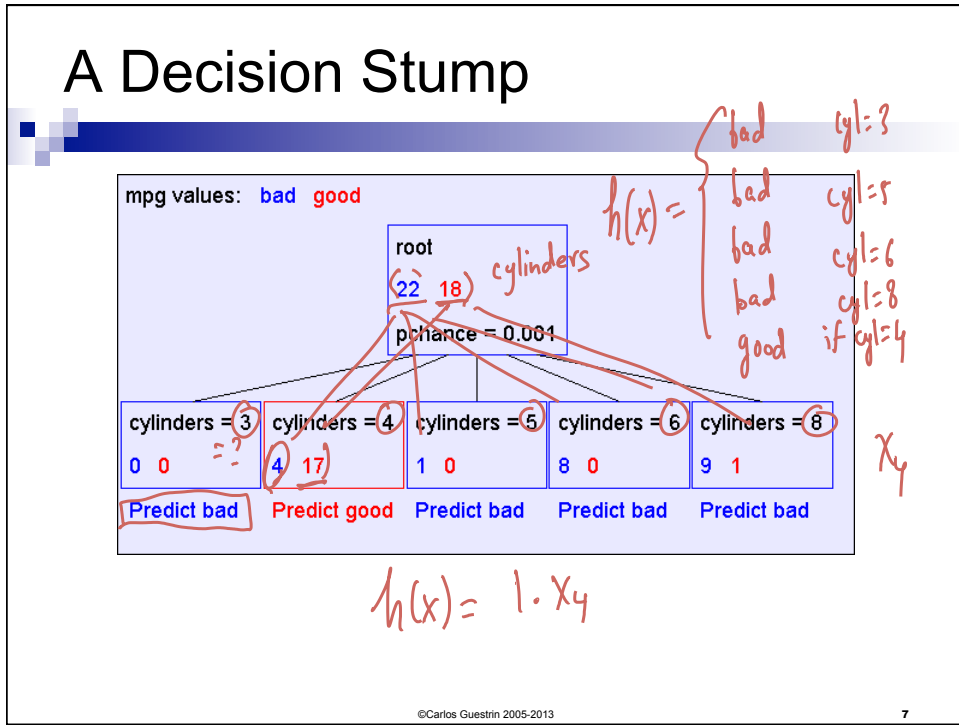
40 training examples

From the UCI repository (thanks to Ross Quinlan)

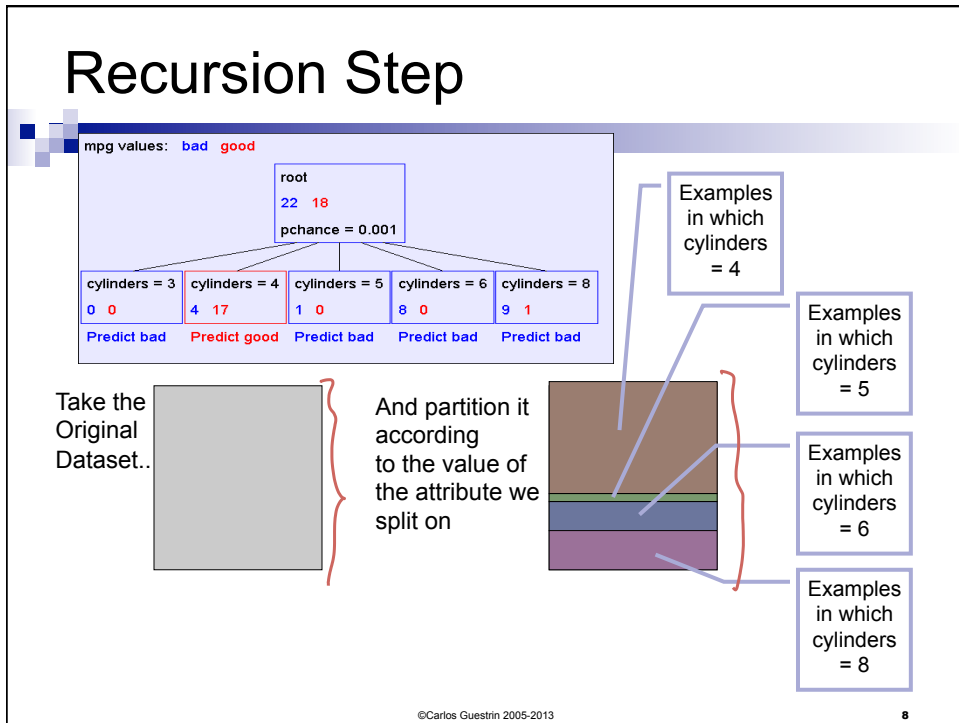
©Carlos Guestrin 2005-2013

6

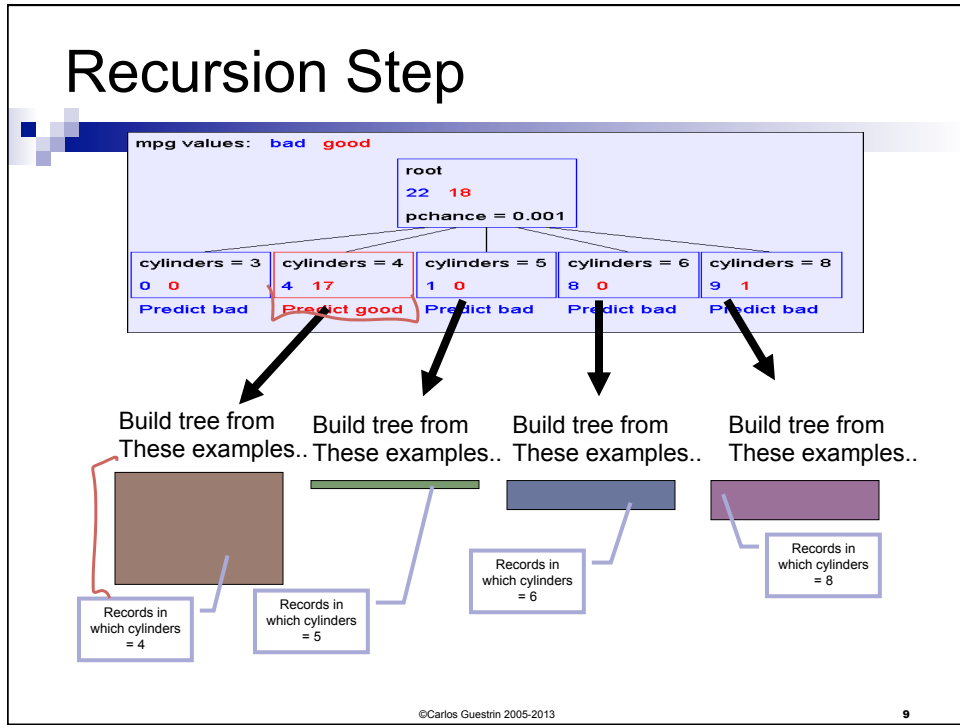
A Decision Stump



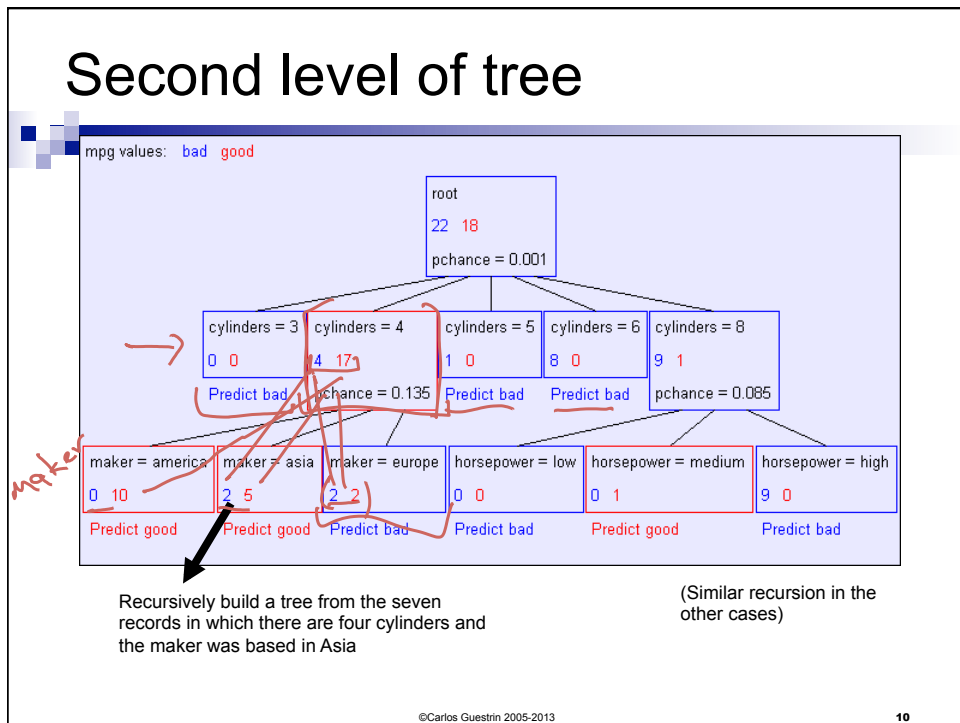
Recursion Step

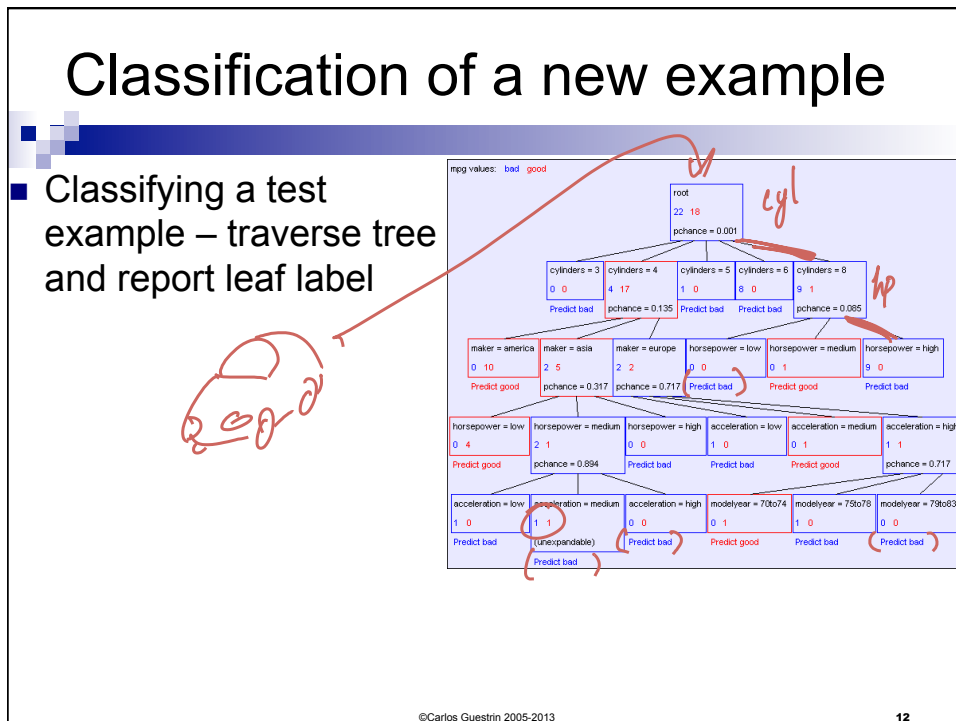
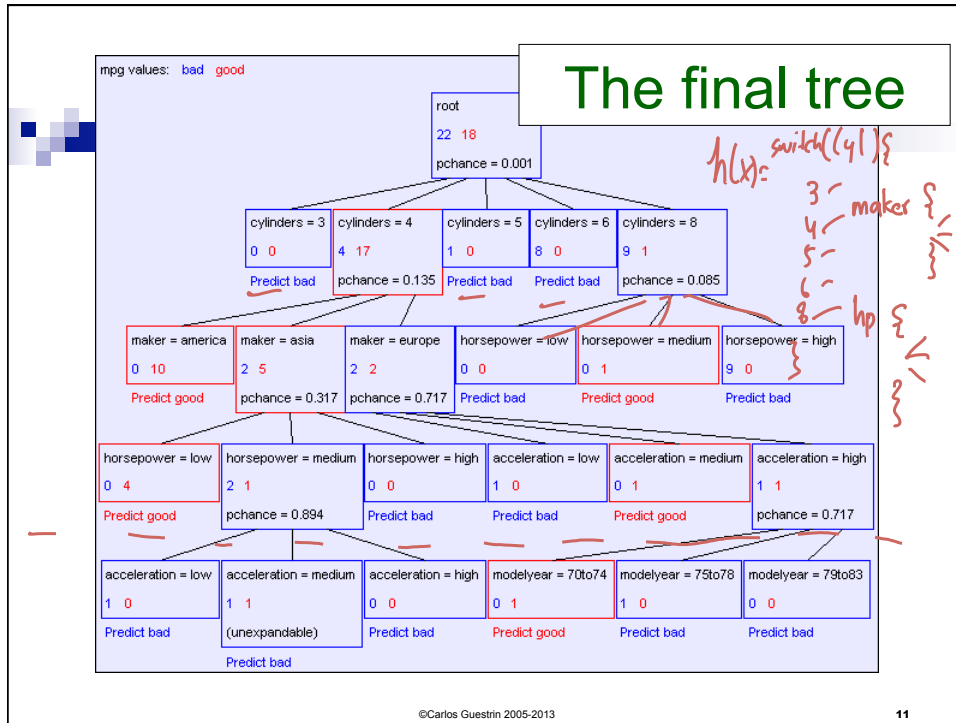


Recursion Step



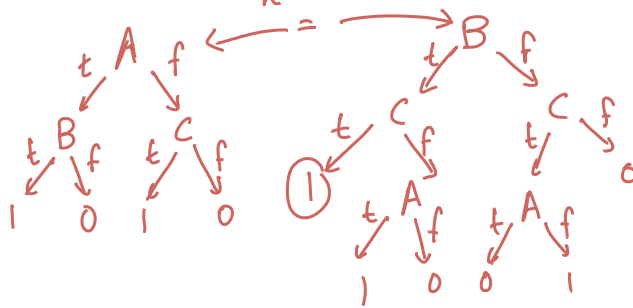
Second level of tree





Are all decision trees equal?

- Many trees can represent the same (concept) h
- But, not all trees will have the same size!
 - e.g., $\phi = A \wedge B \vee \neg A \wedge C$ ((A and B) or (not A and C))

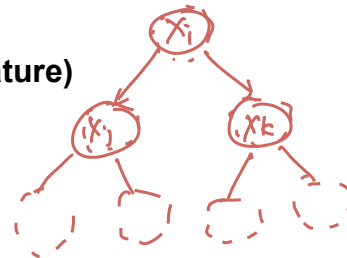


©Carlos Guestrin 2005-2013

13

Learning decision trees is hard!!!

- Learning the simplest (smallest) decision tree is an NP-complete problem [Hyafil & Rivest '76]
- Resort to a greedy heuristic:
 - Start from empty decision tree
 - Split on **next best attribute (feature)**
 - Recurse



©Carlos Guestrin 2005-2013

14

Choosing a good attribute

X_1	X_2	Y
T	T	T
T	F	T
T	T	T
T	F	T
F	T	T
F	F	F
F	T	F
F	F	F

"certainty" is good!

©Carlos Guestrin 2005-2013 15

Measuring uncertainty

- Good split if we are more certain about classification after split
 - Deterministic good (all true or all false)
 - Uniform distribution bad

$P(Y=A) = 1/2$	$P(Y=B) = 1/4$	$P(Y=C) = 1/8$	$P(Y=D) = 1/8$
----------------	----------------	----------------	----------------

$P(Y=A) = 1/4$	$P(Y=B) = 1/4$	$P(Y=C) = 1/4$	$P(Y=D) = 1/4$
----------------	----------------	----------------	----------------

©Carlos Guestrin 2005-2013 16

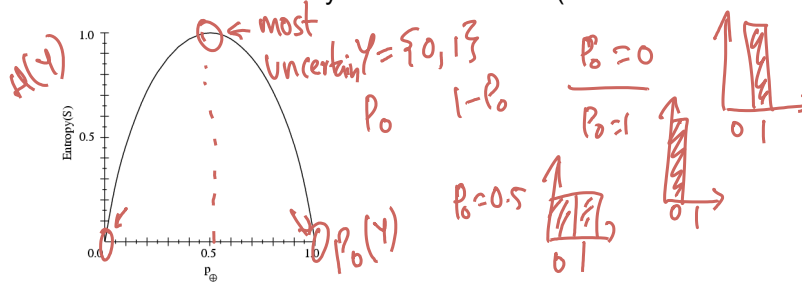
Entropy

Entropy $H(Y)$ of a random variable $Y = \{y_1, y_2, \dots, y_k\}$

$$H(Y) = - \sum_{i=1}^k P(Y = y_i) \log_2 P(Y = y_i)$$

More uncertainty, more entropy!

Information Theory interpretation: $H(Y)$ is the expected number of bits needed to encode a randomly drawn value of Y (under most efficient code)



©Carlos Guestrin 2005-2013

17

Andrew Moore's Entropy in a nutshell



Low Entropy



High Entropy

©Carlos Guestrin 2005-2013

18

Andrew Moore's Entropy in a nutshell



Low Entropy

..the values (locations of soup) sampled entirely from within the soup bowl



High Entropy

..the values (locations of soup) unpredictable... almost uniformly sampled throughout our dining room

Information gain

X ₁	X ₂	Y
T	T	T
T	F	T
T	T	T
T	F	T
F	T	T
F	F	F

Advantage of attribute – decrease in uncertainty

□ Entropy of Y before you split

$$H(Y) = \frac{5}{6} \log_2 \frac{5}{6} + \frac{1}{6} \log_2 \frac{1}{6} \approx 0.65$$

$$P(y=t) = \frac{5}{6}$$

$$P(y=f) = \frac{1}{6}$$

□ Entropy after split

■ Weight by probability of following each branch, i.e., normalized number of records

$$H(Y | X) = - \sum_{j=1}^v P(X_i = x_j) \sum_{i=1}^k P(Y = y_i | X_i = x_j) \log_2 P(Y = y_i | X_i = x_j)$$

$$H(Y | X_1) = \frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} = \frac{1}{3}$$

■ Information gain is difference $IG(X) = H(Y) - H(Y | X)$

$$IG(X_1) = 0.65 - \frac{1}{3} = [0.32]$$

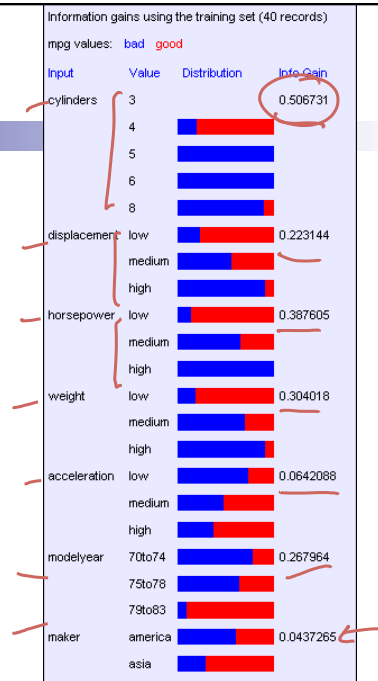
Learning decision trees

- Start from empty decision tree
- Split on **next best attribute (feature)**
 - Use, for example, information gain to select attribute
 - Split on $\arg \max_i IG(X_i) = \arg \max_i H(Y) - H(Y | X_i)$
- Recurse

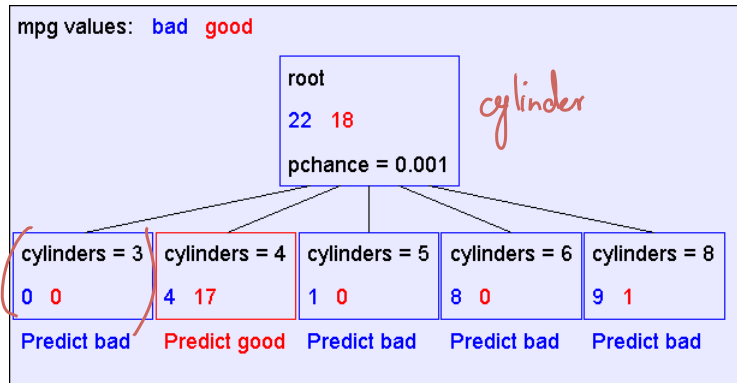
when do you stop?
 1) entropy is 0
 2) cannot split
 3) when IG is 0

Suppose we want to predict MPG

Look at all the information gains...



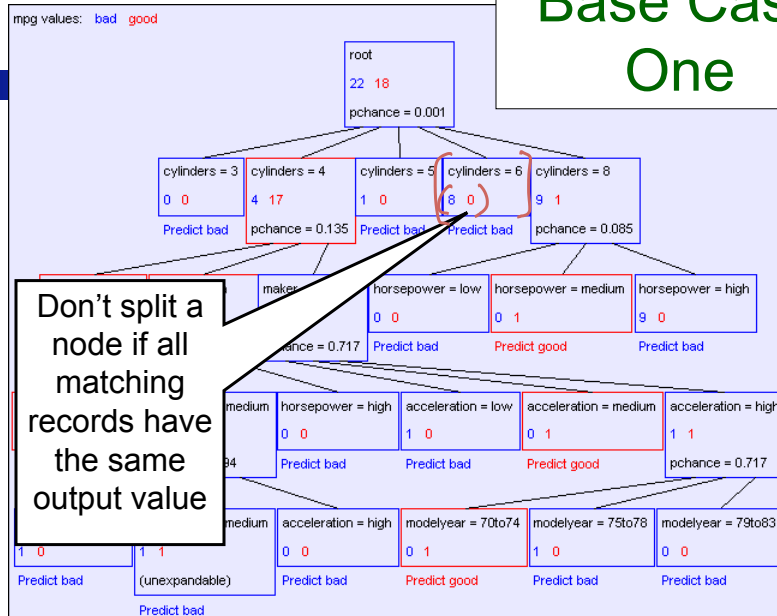
A Decision Stump



©Carlos Guestrin 2005-2013

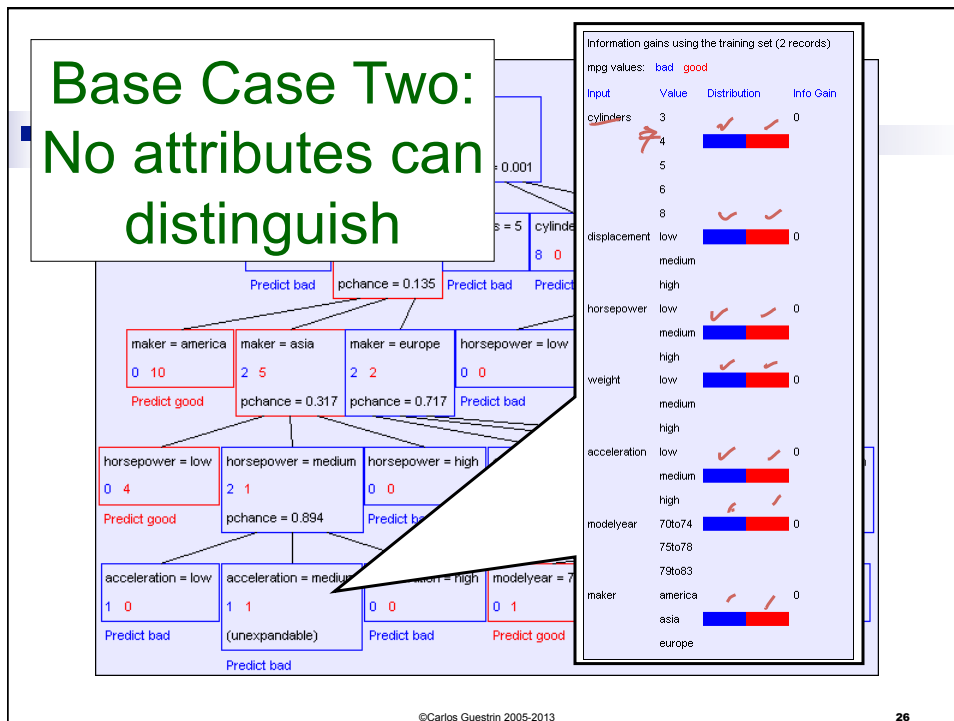
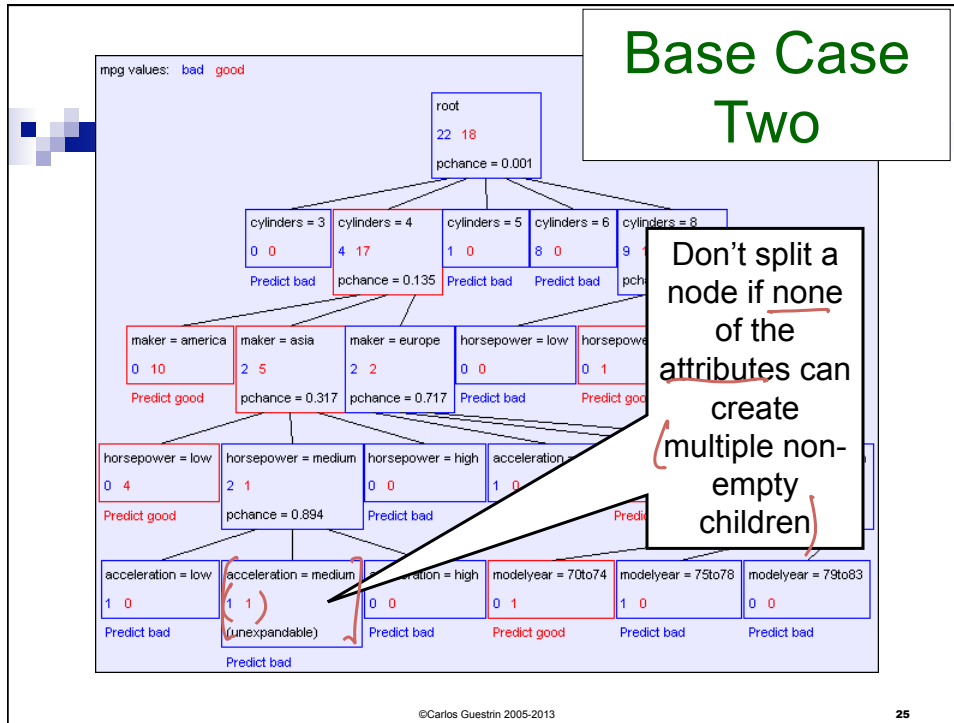
23

Base Case One



©Carlos Guestrin 2005-2013

24



Base Cases

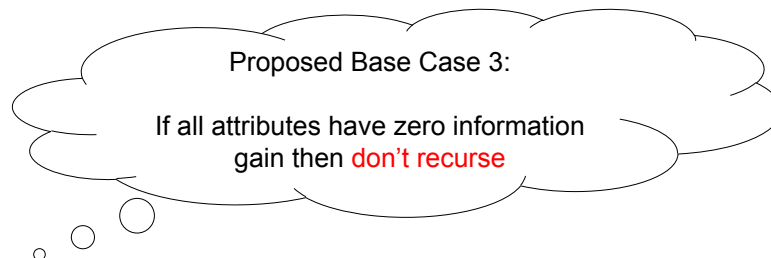
- Base Case One: If all records in current data subset have the same output then don't recurse
- Base Case Two: If all records have exactly the same set of input attributes then don't recurse

3rd case: 0 info gain

Does this work?

Base Cases: An idea

- Base Case One: If all records in current data subset have the same output then **don't recurse**
- Base Case Two: If all records have exactly the same set of input attributes then **don't recurse**



•Is this a good idea?

The problem with Base Case 3

a	b	y
0	0	0
0	1	1
1	0	1
1	1	0

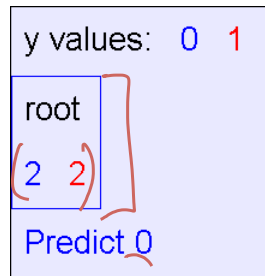
$Y = A \text{ XOR } B$

The information gains:

Information gains using the training set (4 records)
y values: 0 1

Input	Value	Distribution	Info Gain
(a)	0		0
	1		0
(b)	0		0
	1		0

The resulting bad decision tree:

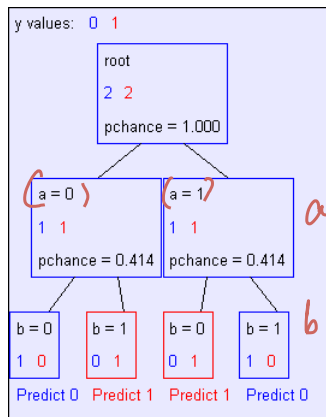


If we omit Base Case 3:

a	b	y
0	0	0
0	1	1
1	0	1
1	1	0

$y = a \text{ XOR } b$

The resulting decision tree:



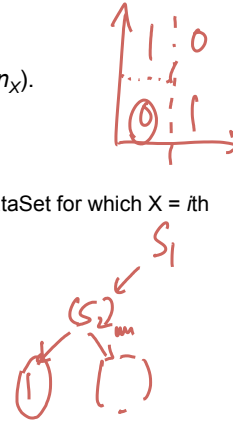
Basic Decision Tree Building Summarized

BuildTree(DataSet, Output)

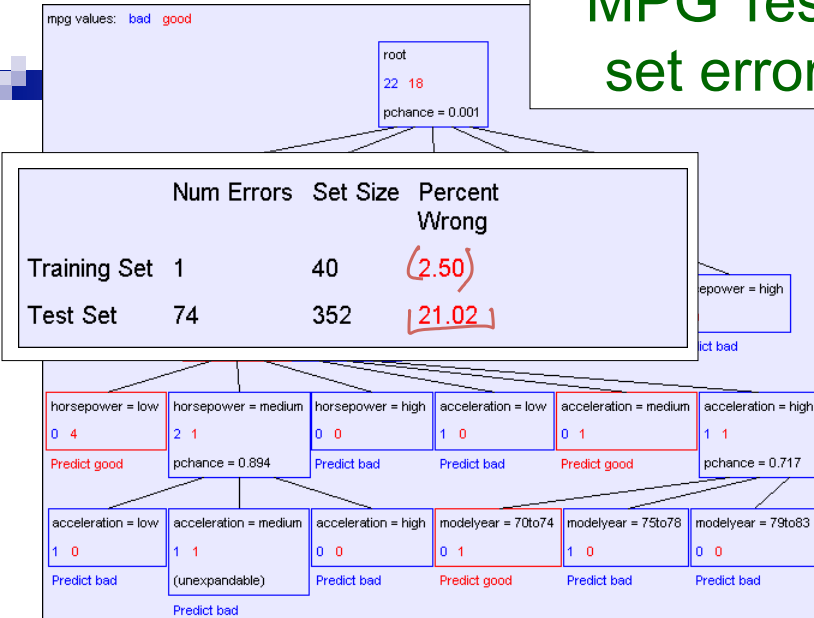
- bc1 ■ If all output values are the same in DataSet, return a leaf node that says "predict this unique output"
- bc2 ■ If all input values are the same, return a leaf node that says "predict the majority output"
- Else find attribute X with highest Info Gain
- Suppose X has n_x distinct values (i.e. X has arity n_x).
 - Create and return a non-leaf node with n_x children.
 - The i 'th child should be built by calling BuildTree(DS_i , Output)

Where DS_i built consists of all those records in DataSet for which $X = i$ th distinct value of X .

$DS_i \subseteq DS$



MPG Test set error



MPG Test set error

mpg values: bad good

root
22 18
pchance = 0.001

	Num Errors	Set Size	Percent Wrong
Training Set	1	40	2.50
Test Set	74	352	21.02

overfitting

The test set error is much worse than the training set error...
...why?

horsepower = low

horsepower = medium

horsepower = high

acceleration = low

acceleration = medium

acceleration = high

Predict bad

(unexpandable)

Predict bad

Predict good

Predict bad

Predict bad

©Carlos Guestrin 2005-2013 33

Decision trees & Learning Bias

$h(x)$ bias variance

linear high low
decision high high

mpg	cylinders	displacement	horsepower	weight	acceleration	modelyear	maker
good	4	low	low	high	high	75to78	asia
bad	6	medium	medium	medium	medium	70to74	america
bad	4	medium	medium	medium	low	75to78	europa
bad	8	high	high	high	low	70to74	america
bad	6	medium	medium	medium	medium	70to74	america
bad	4	low	medium	low	medium	70to74	asia
bad	4	low	medium	low	low	70to74	asia
bad	8	high	high	high	low	75to78	america
...
bad	8	high	high	high	low	70to74	america
good	8	high	medium	high	high	75to83	america
bad	8	high	high	high	low	75to78	america
good	4	low	low	low	low	75to83	america
bad	6	medium	medium	medium	high	75to78	america
good	4	medium	low	low	low	75to83	america
good	4	low	low	medium	high	75to83	america
bad	8	high	high	high	low	70to74	america
good	4	low	medium	low	medium	75to78	europa
bad	5	medium	medium	medium	medium	75to78	europa

©Carlos Guestrin 2005-2013 34