

Decision Trees

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

February 5th, 2007

©2005-2007 Carlos Guestrin

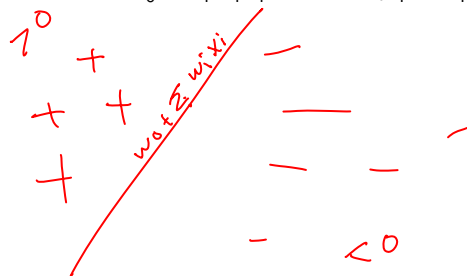
1

Linear separability

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

- $\exists \mathbf{w}$, such that:

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

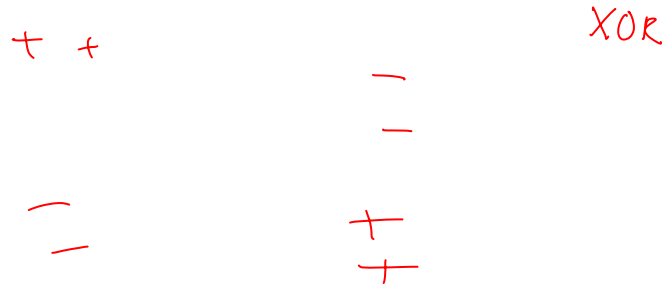


©2005-2007 Carlos Guestrin

2

Not linearly separable data

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

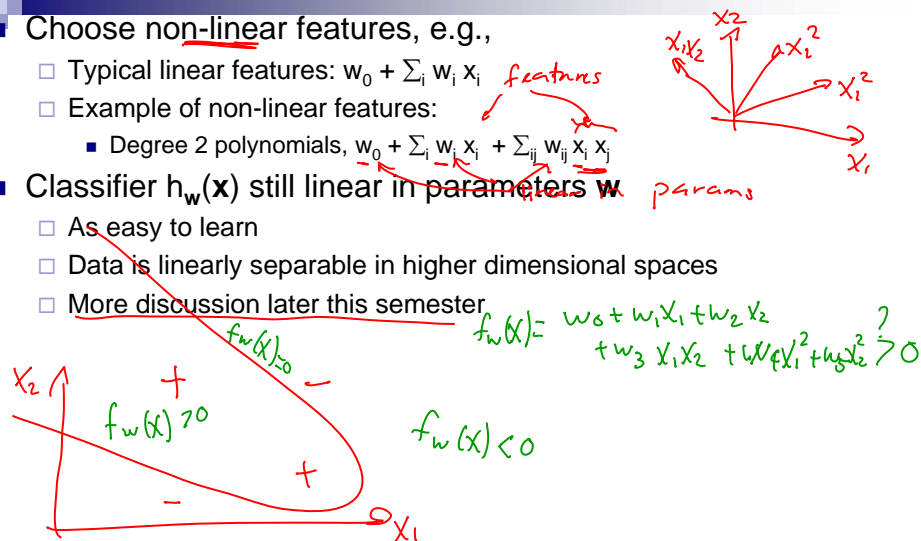


©2005-2007 Carlos Guestrin

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 $h_w(\mathbf{x})$ still linear in parameters w
 - As easy to learn
 - Data is linearly separable in higher dimensional spaces
 - More discussion later this semester



©2005-2007 Carlos Guestrin

4

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

- Choose a classifier $h_{\mathbf{w}}(\mathbf{x})$ that is non-linear in parameters \mathbf{w} , e.g.,
 - Decision trees, neural networks, nearest neighbor,...
- 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 semester, we'll see that these options are not that different)

©2005-2007 Carlos Guestrin

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 Records

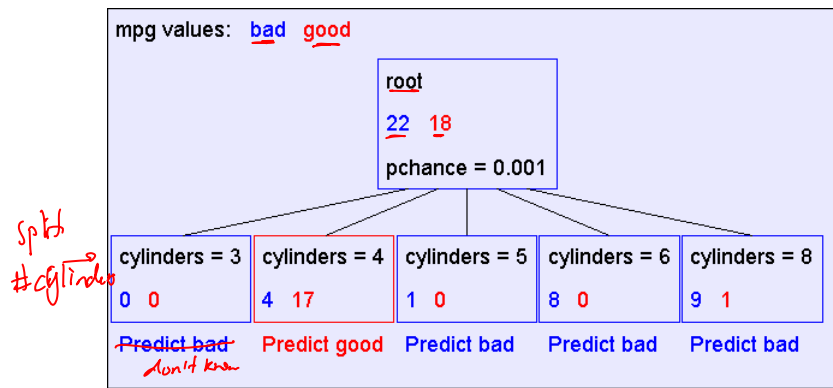
From the UCI repository (thanks to Ross Quinlan)

©2005-2007 Carlos Guestrin

6

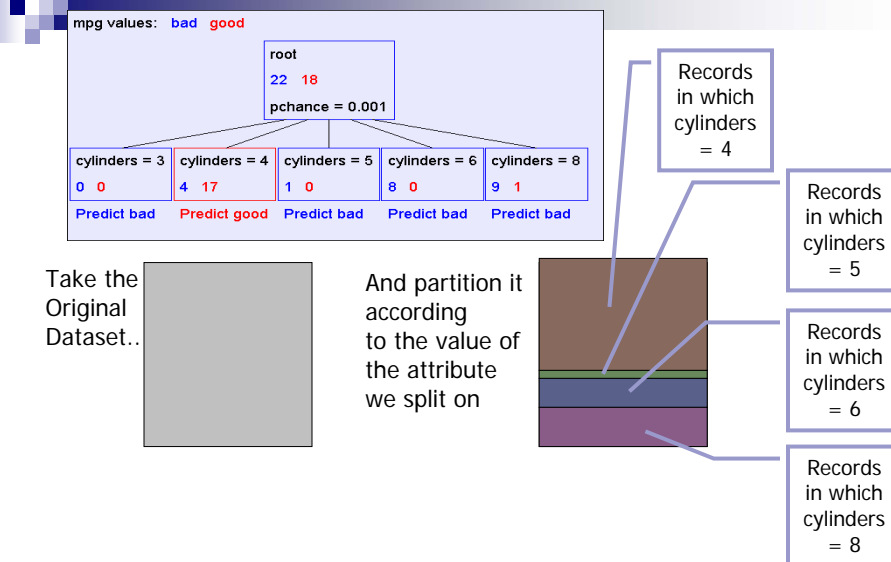
A Decision Stump

1 split



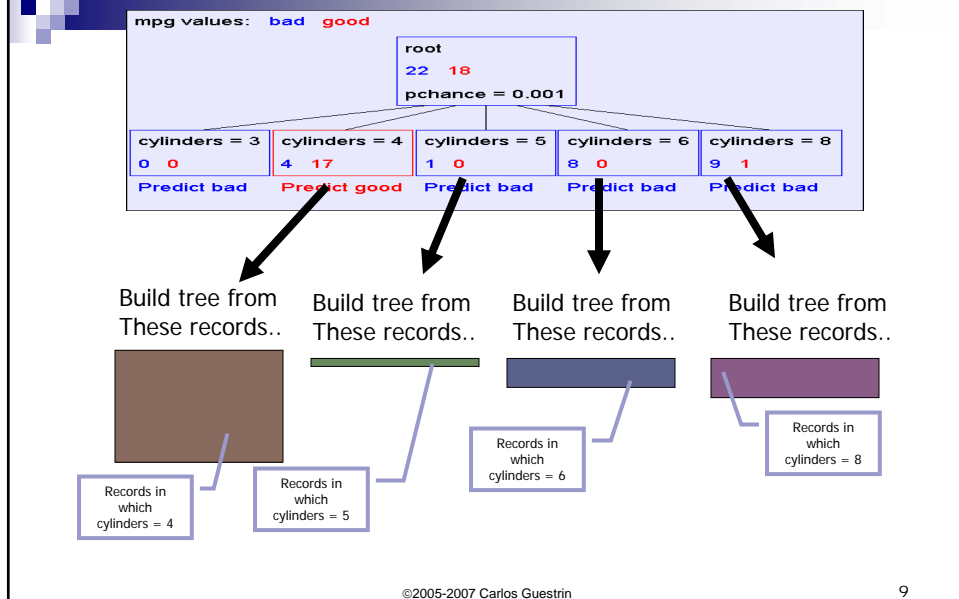
7

Recursion Step

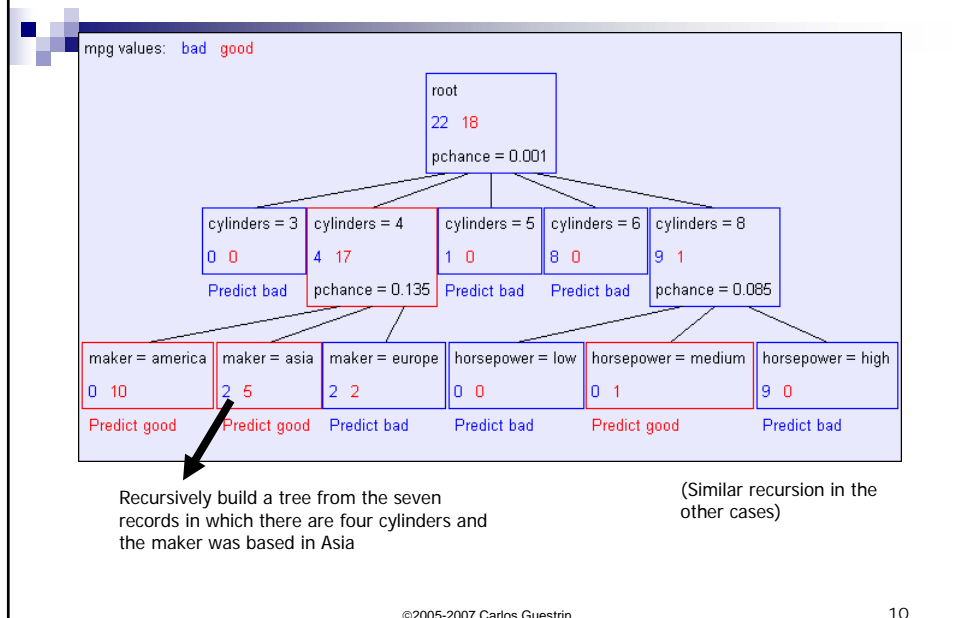


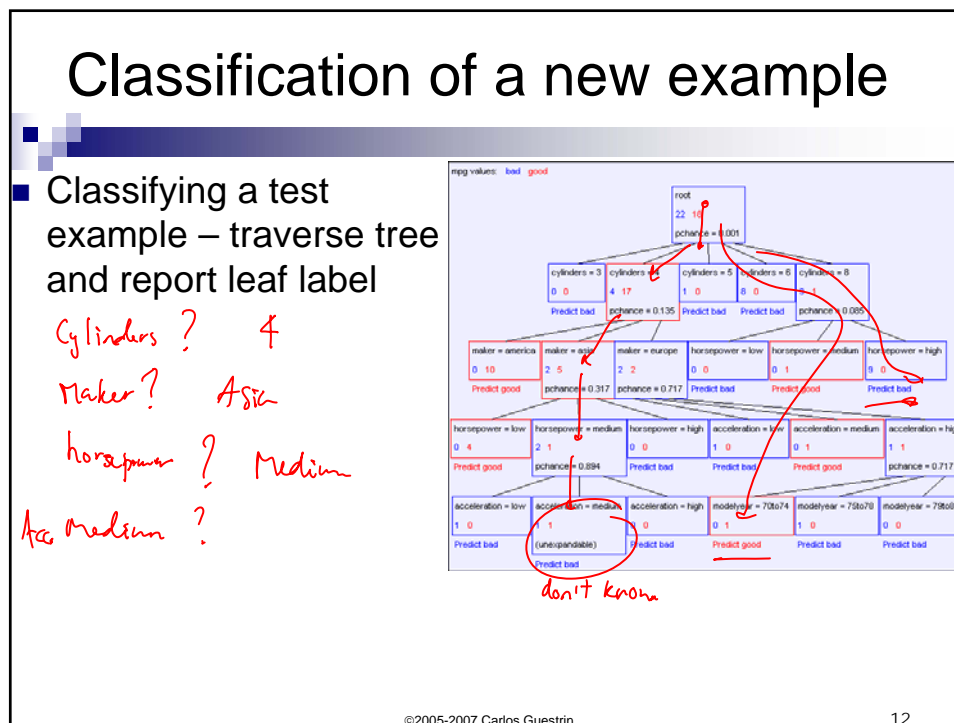
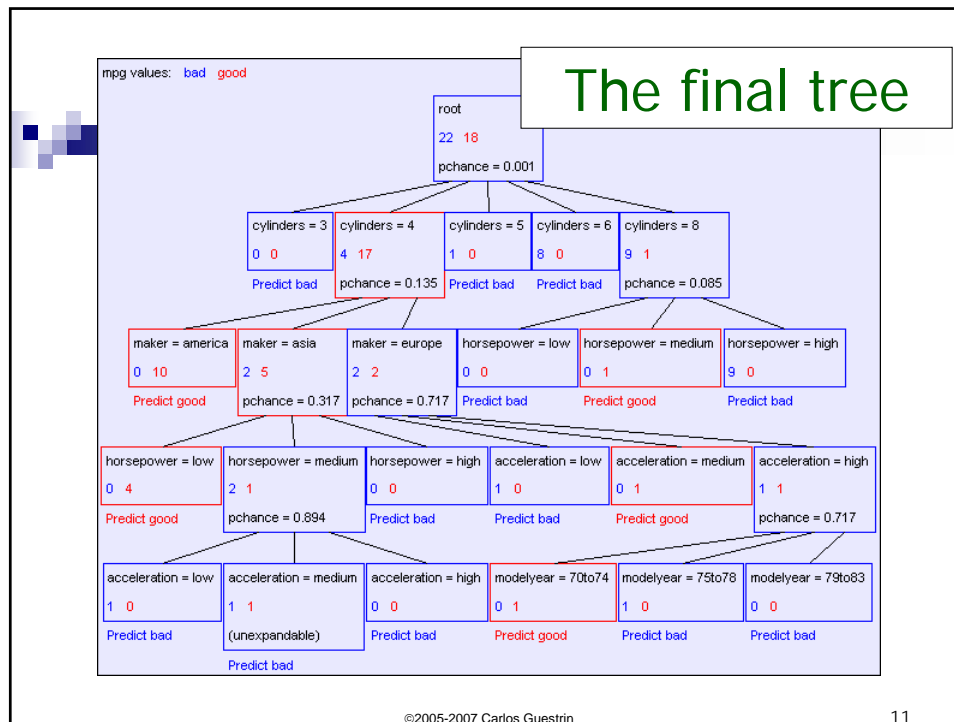
8

Recursion Step



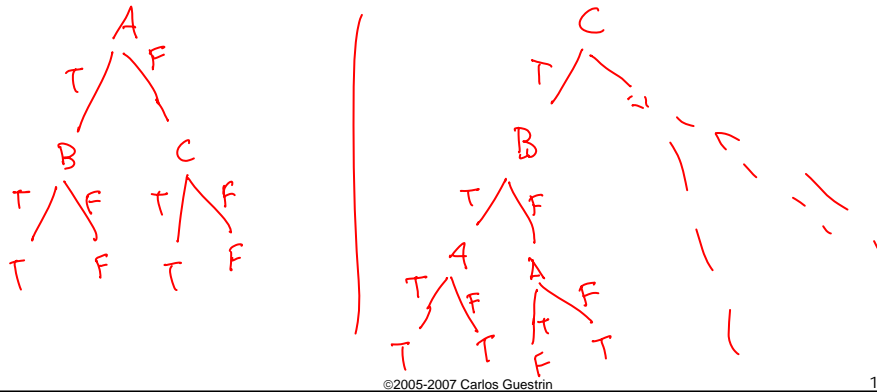
Second level of tree





Are all decision trees equal?

- Many trees can represent the same concept
- 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))



13

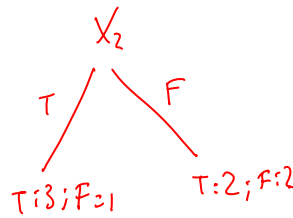
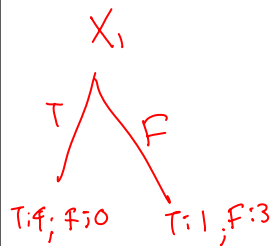
Learning decision trees is hard!!!

- 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

©2005-2007 Carlos Guestrin

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

©2005-2007 Carlos Guestrin

15

Measuring uncertainty

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

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

← lower entropy

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

← High entropy

©2005-2007 Carlos Guestrin

16

Entropy

$$0 \log 0 = 0$$

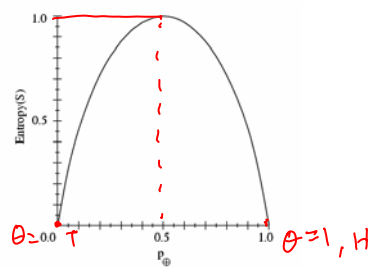
$$\lim_{x \rightarrow 0^+} x \log x = 0$$

Entropy $H(X)$ of a random variable Y

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

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)



coin flip $\theta = H$
 $1-\theta = T$

©2005-2007 Carlos Guestrin

17

Andrew Moore's Entropy in a nutshell



Low Entropy



High Entropy

©2005-2007 Carlos Guestrin

18

Andrew Moore's Entropy in a nutshell



Low Entropy



High Entropy

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

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

©2005-2007 Carlos Guestrin

19

Information gain

Advantage of attribute – decrease in uncertainty

□ Entropy of Y before you split

$$H(Y) = -p(Y=T) \cdot \log_2 P(Y=T) - p(Y=F) \cdot \log_2 P(Y=F)$$

□ Entropy after split

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

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

$P(X_1=T) = \frac{4}{6}$
 $H(Y|X_1=T) = -\frac{4}{6} \log_2 \frac{4}{6} - 0 \log_2 0$
 average over values of X_1

$H(Y|X_2=2)$

$P(X_2=F) = \frac{2}{6}$

$H(Y|X_2=F) = 1$

$H(Y|X_1) = P(X_1=T) \cdot H(Y|X_1=T) + P(X_1=F) \cdot H(Y|X_1=F)$

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

$H(Y) - H(Y|X_1) = H(Y) - \frac{1}{3}$

$H(Y|X_1) = \frac{4}{6} \cdot 0 + \frac{2}{6} \cdot 1 = \frac{1}{3}$

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

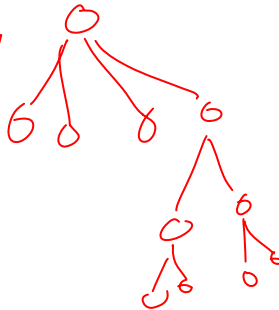
©2005-2007 Carlos Guestrin

20

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

if sure about
a leaf, $IG(x)=0$



©2005-2007 Carlos Guestrin

21

Suppose we want to
predict MPG

Look at all the
information
gains...

Information gains using the training set (40 records)

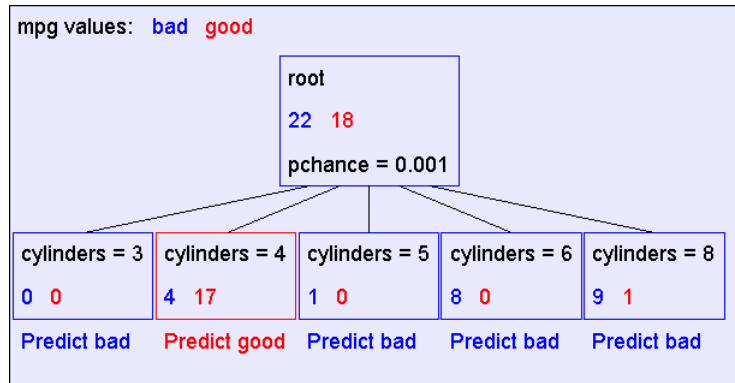
mpg values: bad good

Input	Value	Distribution	Info Gain
cylinders	3		0.506731
	4		
	5		
	6		
	8		
displacement	low		0.223144
	medium		
	high		
horsepower	low		0.387605
	medium		
	high		
weight	low		0.304018
	medium		
	high		
acceleration	low		0.0642088
	medium		
	high		
modelyear	70to74		0.267964
	75to78		
	79to83		
maker	america		0.0437265
	asia		

©2005-2007 Carlos Guestrin

22

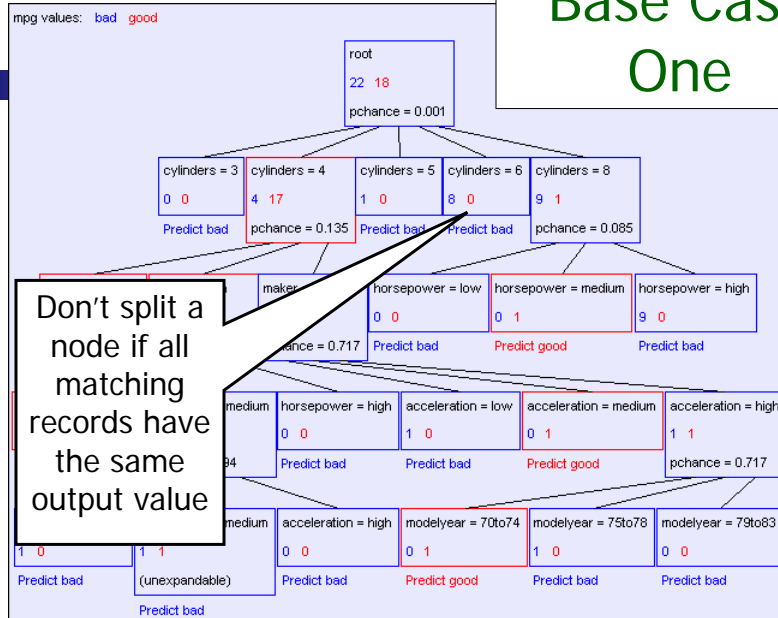
A Decision Stump



©2005-2007 Carlos Guestrin

23

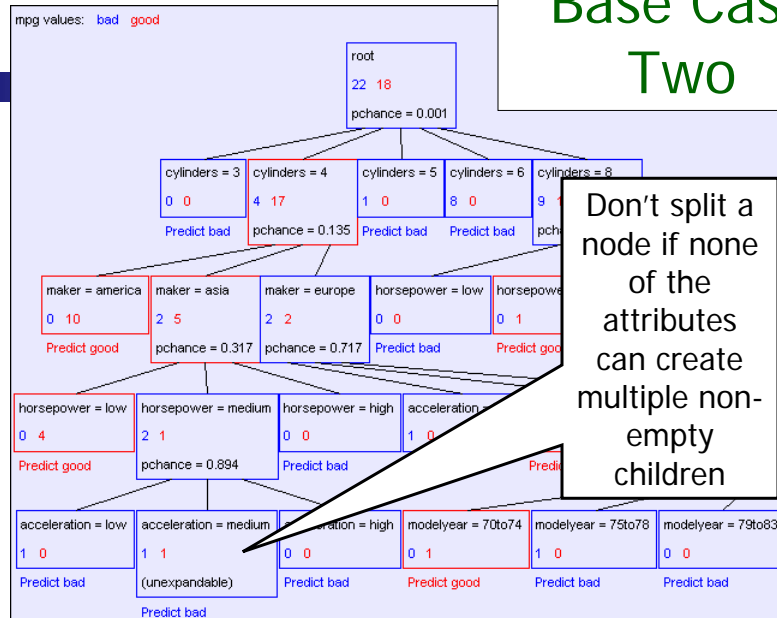
Base Case One



©2005-2007 Carlos Guestrin

24

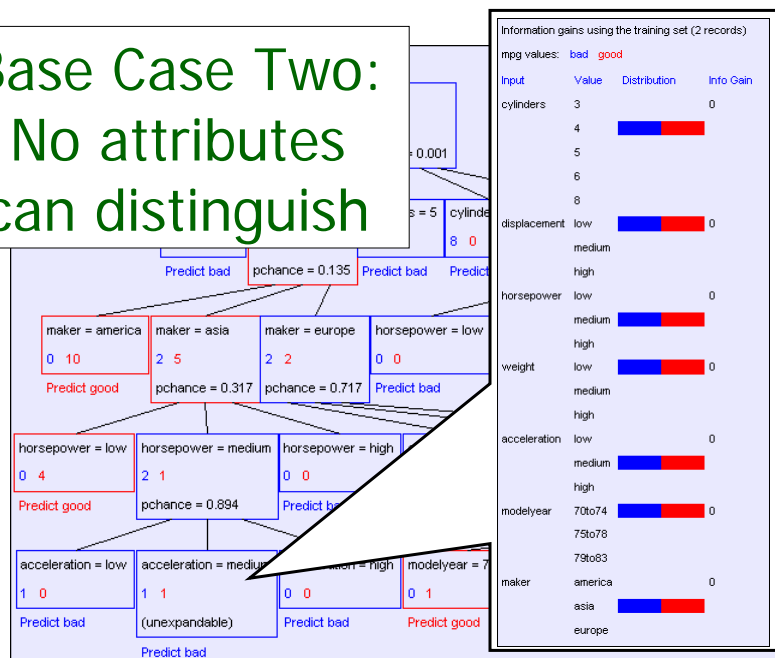
Base Case Two



©2005-2007 Carlos Guestrin

25

Base Case Two: No attributes can distinguish



©2005-2007 Carlos Guestrin

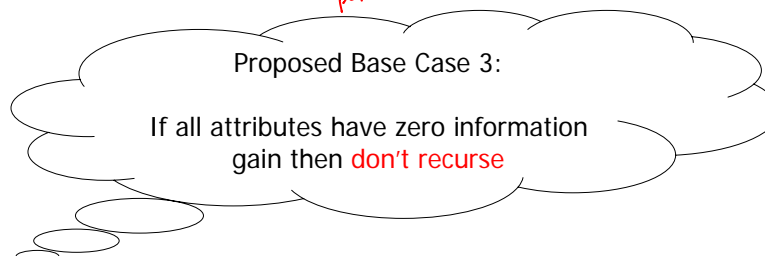
26

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

Base Cases: An idea

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

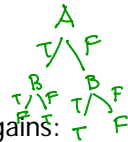


• *Is this a good idea?*

The problem with Base Case 3

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

$$y = a \text{ XOR } b$$

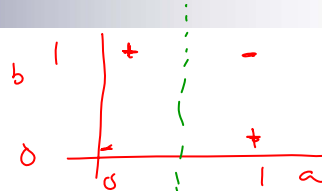


The information gains:

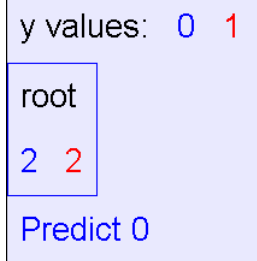
Information gains using the training set (4 records)

y values: 0 1

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



The resulting decision tree:



©2005-2007 Carlos Guestrin

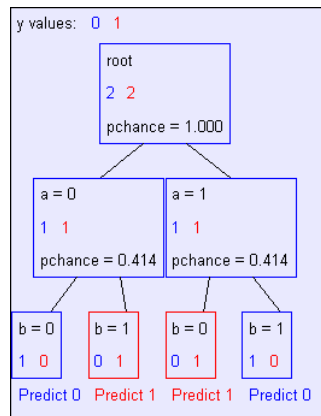
29

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:



©2005-2007 Carlos Guestrin

30

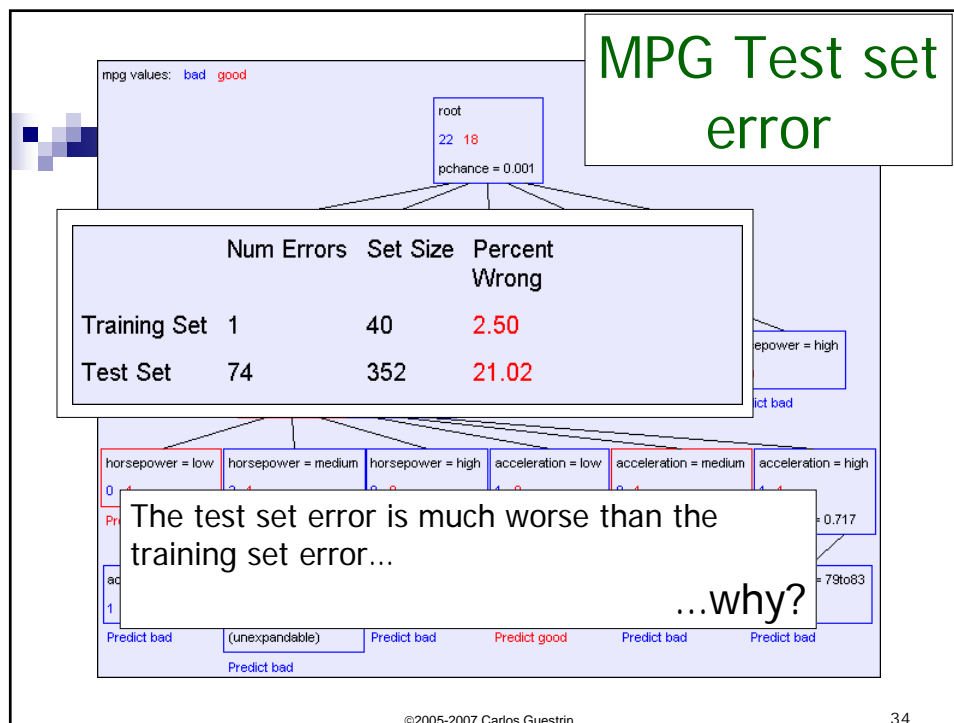
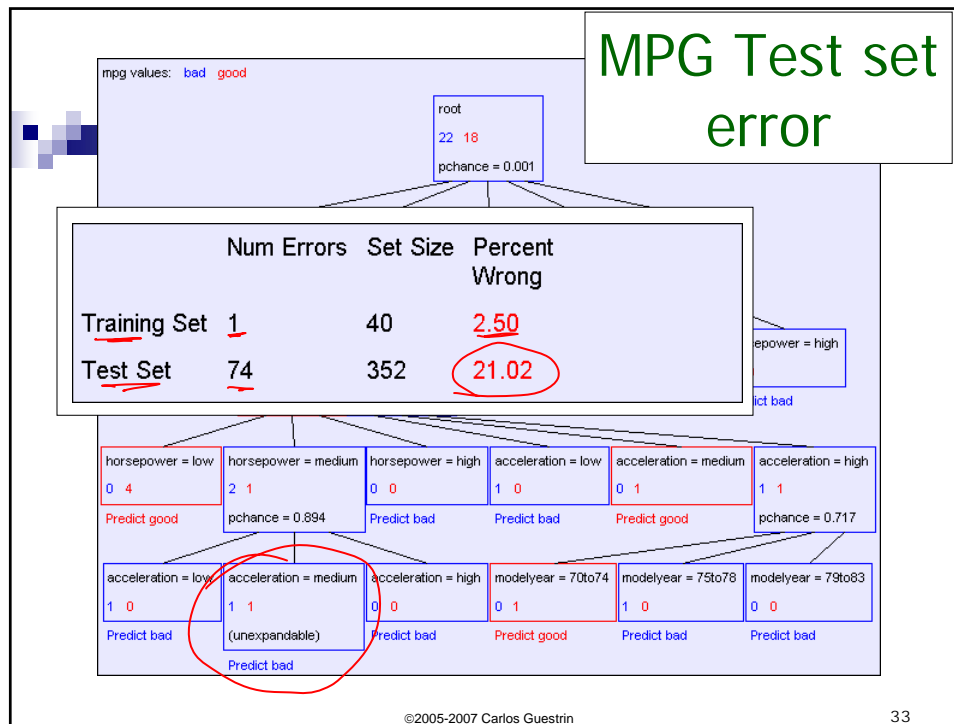
Basic Decision Tree Building Summarized

BuildTree(DataSet, Output)

- If all output values are the same in *DataSet*, return a leaf node that says “predict this unique output”
- 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*.

Announcements

- **Pittsburgh won the Super Bowl !!**
 - Last year...



Decision trees & Learning Bias

label noise

for (almost) any table \exists decision tree. exactly represents table \Rightarrow Bias low

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

©2005-2007 Carlos Guestrin

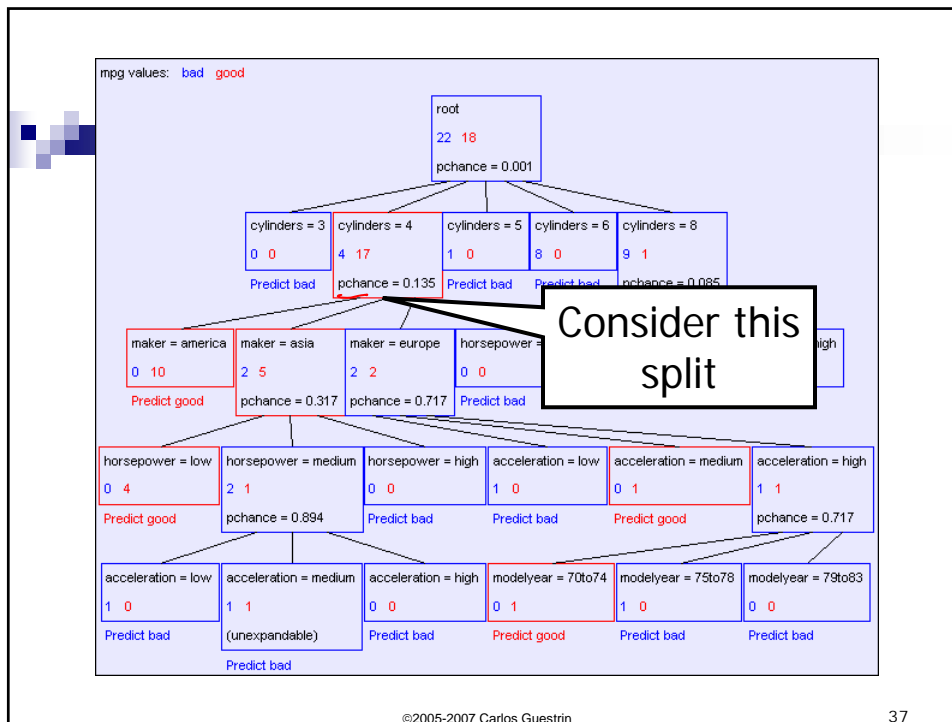
35

Decision trees will overfit

- Standard decision trees are have no learning biased
 - Training set error is always zero!
 - (If there is no label noise)
 - Lots of variance
 - Will definitely overfit!!!
 - Must bias towards simpler trees
- Many strategies for picking simpler trees:
 - Fixed depth
 - Fixed number of leaves
 - Or something smarter...

©2005-2007 Carlos Guestrin

36



A chi-square test

mpg values: bad good







maker	america	asia	europe	$H(\text{mpg} \text{maker})$
america	0 10			0
asia	2 5			0.863121
europe	2 2			1

$H(\text{mpg}) = 0.702467$ $H(\text{mpg}|\text{maker}) = 0.478183$
 $IG(\text{mpg}|\text{maker}) = 0.224284$

- Suppose that mpg was completely uncorrelated with maker.
- What is the chance we'd have seen data of at least this apparent level of association anyway?

A chi-square test

mpg values: bad good

maker	america	0	10			$H(\text{mpg} \mid \text{maker} = \text{america}) = 0$
	asia	2	5			$H(\text{mpg} \mid \text{maker} = \text{asia}) = 0.863121$
	europa	2	2			$H(\text{mpg} \mid \text{maker} = \text{europa}) = 1$

$H(\text{mpg}) = 0.702467$ $H(\text{mpg} \mid \text{maker}) = 0.478183$
 $IG(\text{mpg} \mid \text{maker}) = 0.224284$

- Suppose that mpg was completely uncorrelated with maker.
- What is the chance we'd have seen data of at least this apparent level of association anyway?

By using a particular kind of chi-square test, the answer is 7.2%

(Such simple hypothesis tests are very easy to compute, unfortunately, not enough time to cover in the lecture, but in your homework, you'll have fun! :))

©2005-2007 Carlos Guestrin

39

Using Chi-squared to avoid overfitting

- Build the full decision tree as before
- But when you can grow it no more, start to prune:
 - Beginning at the bottom of the tree, delete splits in which $p_{\text{chance}} > \text{MaxPchance}$
 - Continue working your way up until there are no more prunable nodes

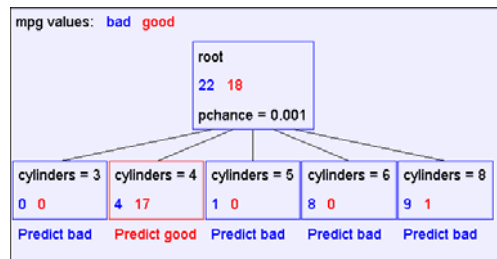
MaxPchance is a magic parameter you must specify to the decision tree, indicating your willingness to risk fitting noise

©2005-2007 Carlos Guestrin

40

Pruning example

- With MaxPchance = 0.1, you will see the following MPG decision tree:



Note the improved test set accuracy compared with the unpruned tree

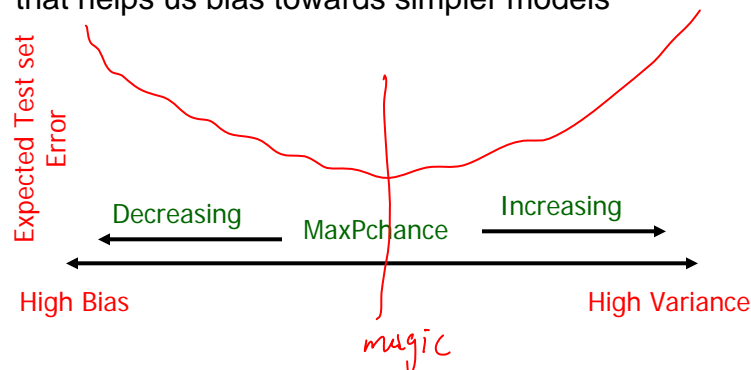
	Num Errors	Set Size	Percent Wrong
Training Set	5	40	<u>12.50</u>
Test Set	56	352	<u>15.91</u>

©2005-2007 Carlos Guestrin

41

MaxPchance

- Technical note MaxPchance is a regularization parameter that helps us bias towards simpler models



- We'll learn to choose the value of these magic parameters soon!

©2005-2007 Carlos Guestrin

42

Real-Valued inputs

- What should we do if some of the inputs are real-valued?

mpg	cylinders	displacemen	horsepower	weight	acceleration	modelyear	maker
good	4	97	75	2265	18.2	77	asia
bad	6	199	90	2648	15	70	america
bad	4	121	110	2600	12.8	77	europa
bad	8	350	175	4100	13	73	america
bad	6	198	95	3102	16.5	74	america
bad	4	108	94	2379	16.5	73	asia
bad	4	113	95	2228	14	71	asia
bad	8	302	139	3570	12.8	78	america
.
.
.
.
good	4	120	79	2625	18.6	82	america
bad	8	455	225	4425	10	70	america
good	4	107	86	2464	15.5	76	europa
bad	5	131	103	2830	15.9	78	europa

Infinite number of possible split values!!!

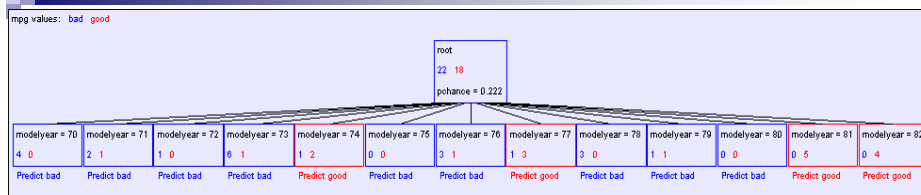
Finite dataset, only finite number of relevant splits!

Idea One: Branch on each possible real value

©2005-2007 Carlos Guestrin

43

“One branch for each numeric value” idea:



Hopeless: with such high branching factor will shatter the dataset and overfit

©2005-2007 Carlos Guestrin

44

Threshold splits

- Binary tree, split on attribute X

- One branch: $X < t$

t ← threshold

- Other branch: $X \geq t$



©2005-2007 Carlos Guestrin

45

Choosing threshold split

- Binary tree, split on attribute X

- One branch: $X < t$

- Other branch: $X \geq t$

- Search through possible values of t

- Seems hard!!!

- But only finite number of t 's are important

- Sort data according to X into $\{x_1, \dots, x_m\}$

- Consider split points of the form $x_i + (x_{i+1} - x_i)/2$

x_i:
10

12
- - - -
14
- - - -
15
- - - -
17

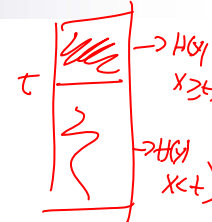
23

©2005-2007 Carlos Guestrin

46

A better idea: thresholded splits

- Suppose X is real valued
- Define $IG(Y|X:t)$ as $H(Y) - H(Y|X:t)$
- Define $H(Y|X:t) = H(Y|X < t)P(X < t) + H(Y|X \geq t)P(X \geq t)$
 - $IG(Y|X:t)$ is the information gain for predicting Y if all you know is whether X is greater than or less than t
- Then define $IG^*(Y|X) = \max_t IG(Y|X:t)$
- For each real-valued attribute, use $IG^*(Y|X)$ for assessing its suitability as a split



©2005-2007 Carlos Guestrin

47

Information gains using the training set (40 records)

mpg values: bad good

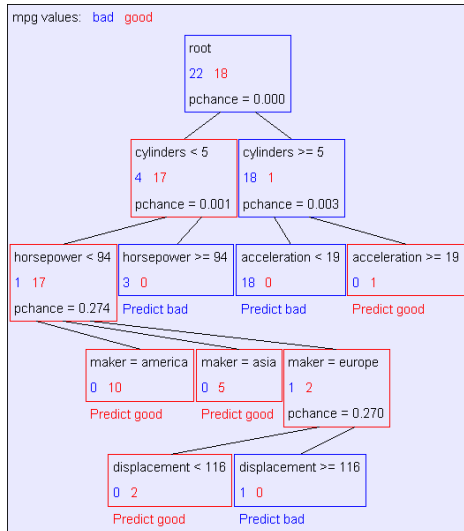
Input	Value	Distribution	Info Gain
cylinders	< 5		0.48268
	>= 5		
displacement	< 198		0.428205
	>= 198		
horsepower	< 94		0.48268
	>= 94		
weight	< 2789		0.379471
	>= 2789		
acceleration	< 18.2		0.159982
	>= 18.2		
modelyear	< 81		0.319193
	>= 81		
maker	america		0.0437265
	asia		
	europa		

Example with MPG

©2005-2007 Carlos Guestrin

48

Example tree using reals



allowed to split
on an x_i
multiple times

©2005-2007 Carlos Guestrin

49

What you need to know about decision trees

- Decision trees are one of the most popular data mining tools
 - Easy to understand
 - Easy to implement
 - Easy to use
 - Computationally cheap (to solve heuristically)
- Information gain to select attributes (ID3, C4.5,...)
- Presented for classification, can be used for regression and density estimation too
- Decision trees will overfit!!!
 - Zero bias classifier → Lots of variance
 - Must use tricks to find “simple trees”, e.g.,
 - Fixed depth/Early stopping
 - Pruning
 - Hypothesis testing

©2005-2007 Carlos Guestrin

50

Acknowledgements

- Some of the material in the presentation is courtesy of Andrew Moore, from his excellent collection of ML tutorials:
 - <http://www.cs.cmu.edu/~awm/tutorials>