

Decision Trees

Machine Learning - 10601

Geoff Gordon, Miroslav Dudík

[partly based on slides of Carlos Guestrin and Andrew Moore]

<http://www.cs.cmu.edu/~ggordon/10601/>

October 21, 2009

Non-linear Classifiers

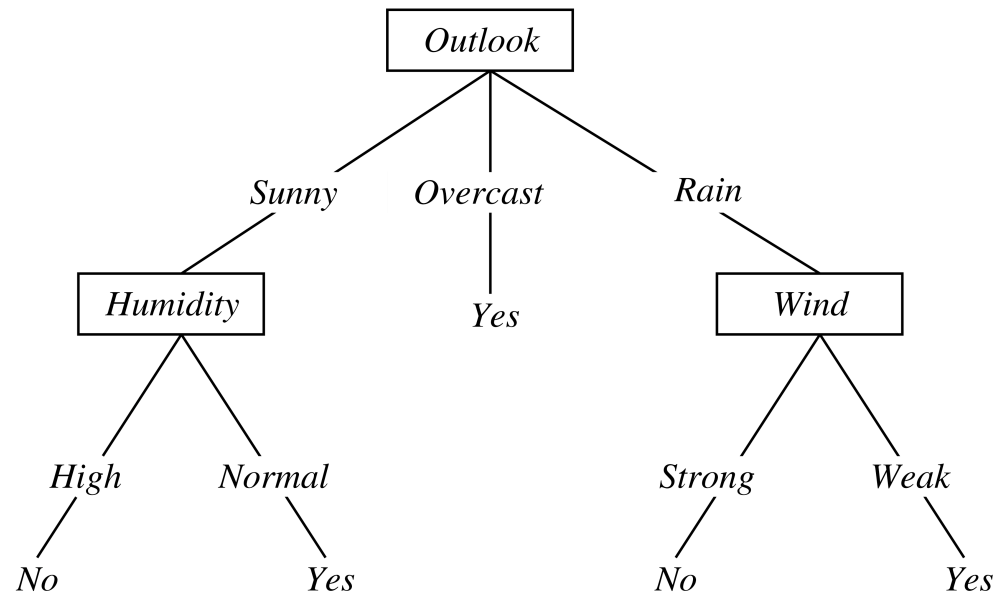
Dealing with non-linear decision boundary

1. add “non-linear” features
to a linear model (e.g., logistic regression)
2. **use non-linear learners**
(nearest neighbors, decision trees, artificial neural nets, ...)

k-Nearest Neighbor Classifier

- simple, often a good baseline
- can approximate arbitrary boundary: **non-parametric**
- **downside:** stores all the data

A Decision Tree for *PlayTennis*



Each internal node: **test one feature X_j**

Each branch from a node: **select one value for X_j**

Each leaf node node: **predict Y**

or **$P(Y \mid X \in \text{leaf})$**

Decision trees

How would you represent

$$Y = A \vee B \quad (A \text{ or } B)$$

Decision trees

How would you represent

$$Y = (A \wedge B) \vee (\neg A \wedge C) \quad ((A \text{ and } B) \text{ or } (\text{not } A \text{ and } C))$$

Optimal Learning of Decision Trees is Hard

- learning the smallest (simplest) decision tree is NP-complete (existing algorithms exponential)
- use “greedy” heuristics:
 - start with an empty tree
 - choose the **next best** attribute (feature)
 - recurse

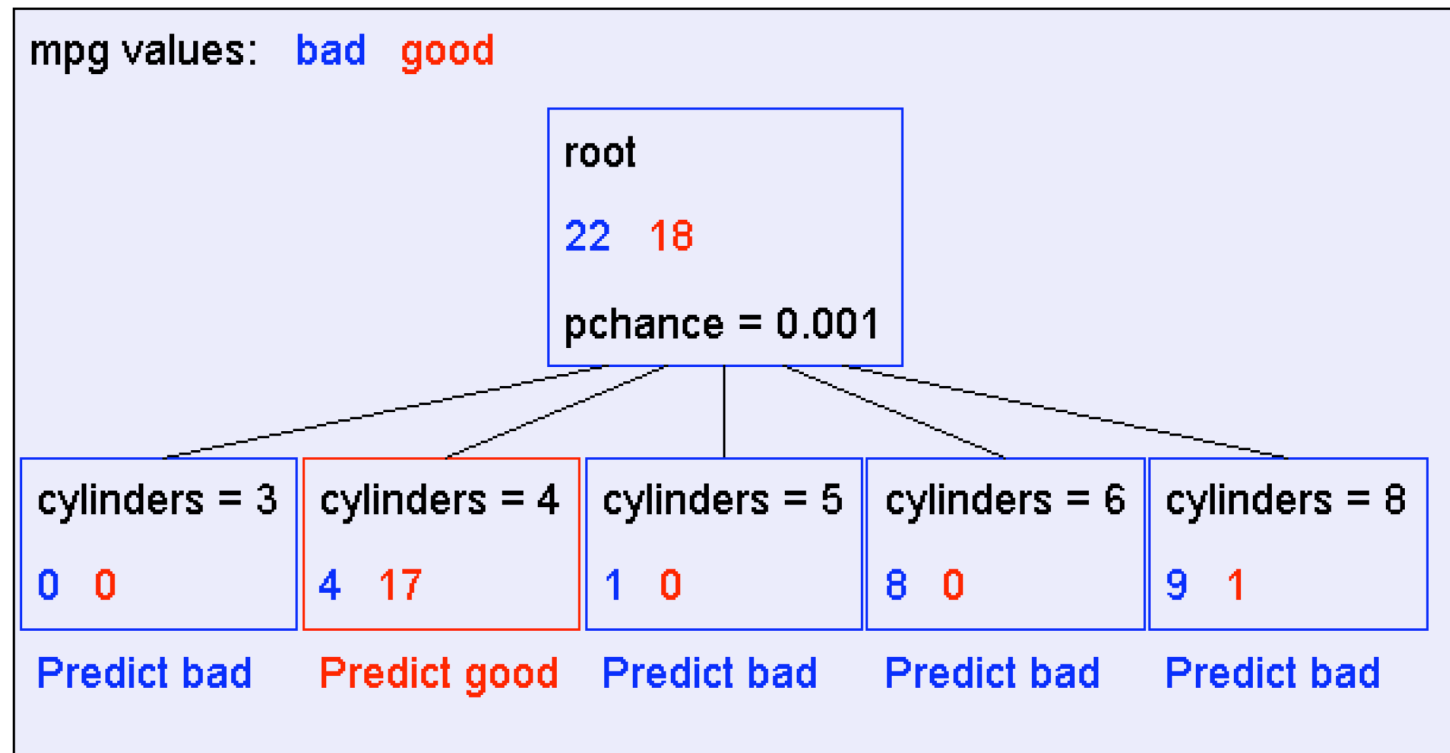
A small dataset: predict miles per gallon (mpg)

40
Records

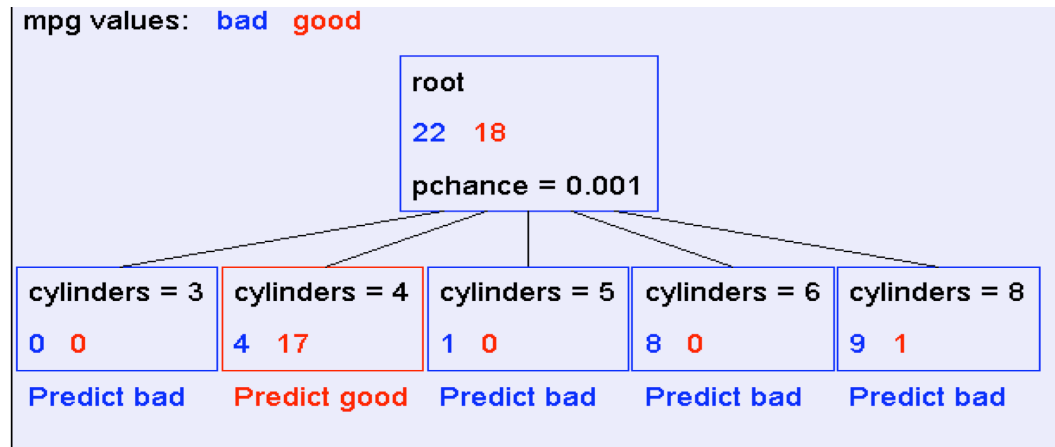
mpg	cylinders	displacement	horsepower	weight	acceleration	modelyear	maker
good	4	low	low	low	high	75to78	asia
bad	6	medium	medium	medium	medium	70to74	america
bad	4	medium	medium	medium	low	75to78	europa
bad	8	high	high	high	low	70to74	america
bad	6	medium	medium	medium	medium	70to74	america
bad	4	low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

From the UCI repository (thanks to Ross Quinlan)

A Decision Stump



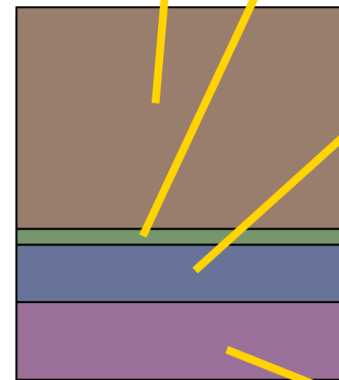
Recursion Step



Take the
Original
Dataset..



And partition it
according
to the value of
the attribute
we split on



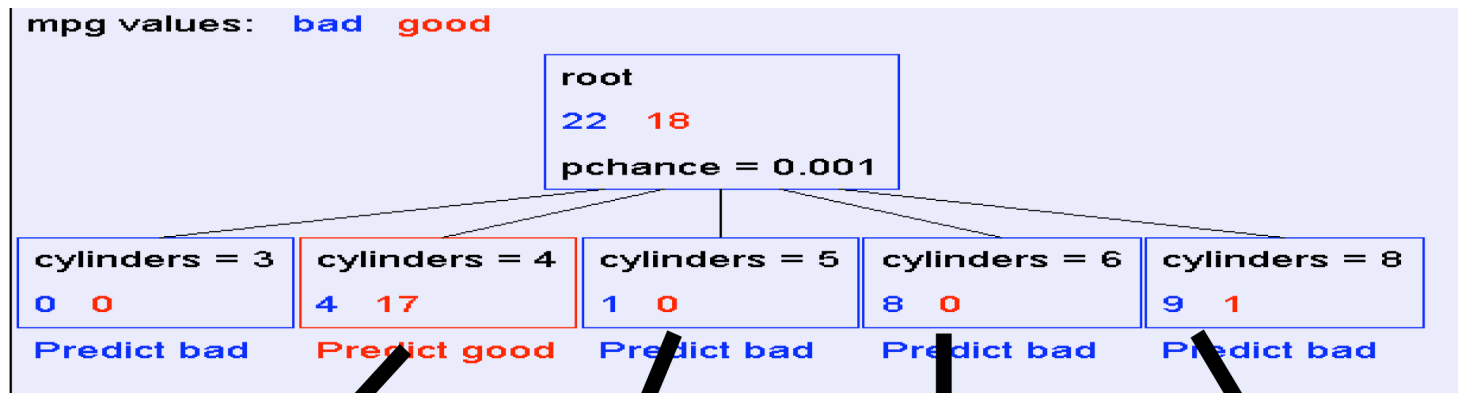
Records
in which
cylinders
= 4

Records
in which
cylinders
= 5

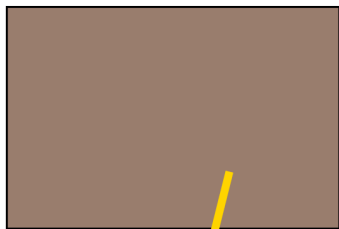
Records
in which
cylinders
= 6

Records
in which
cylinders
= 8

Recursion Step



Build tree from
These records..



Records in
which
cylinders = 4

Build tree from
These records..



Records in
which
cylinders = 5

Build tree from
These records..



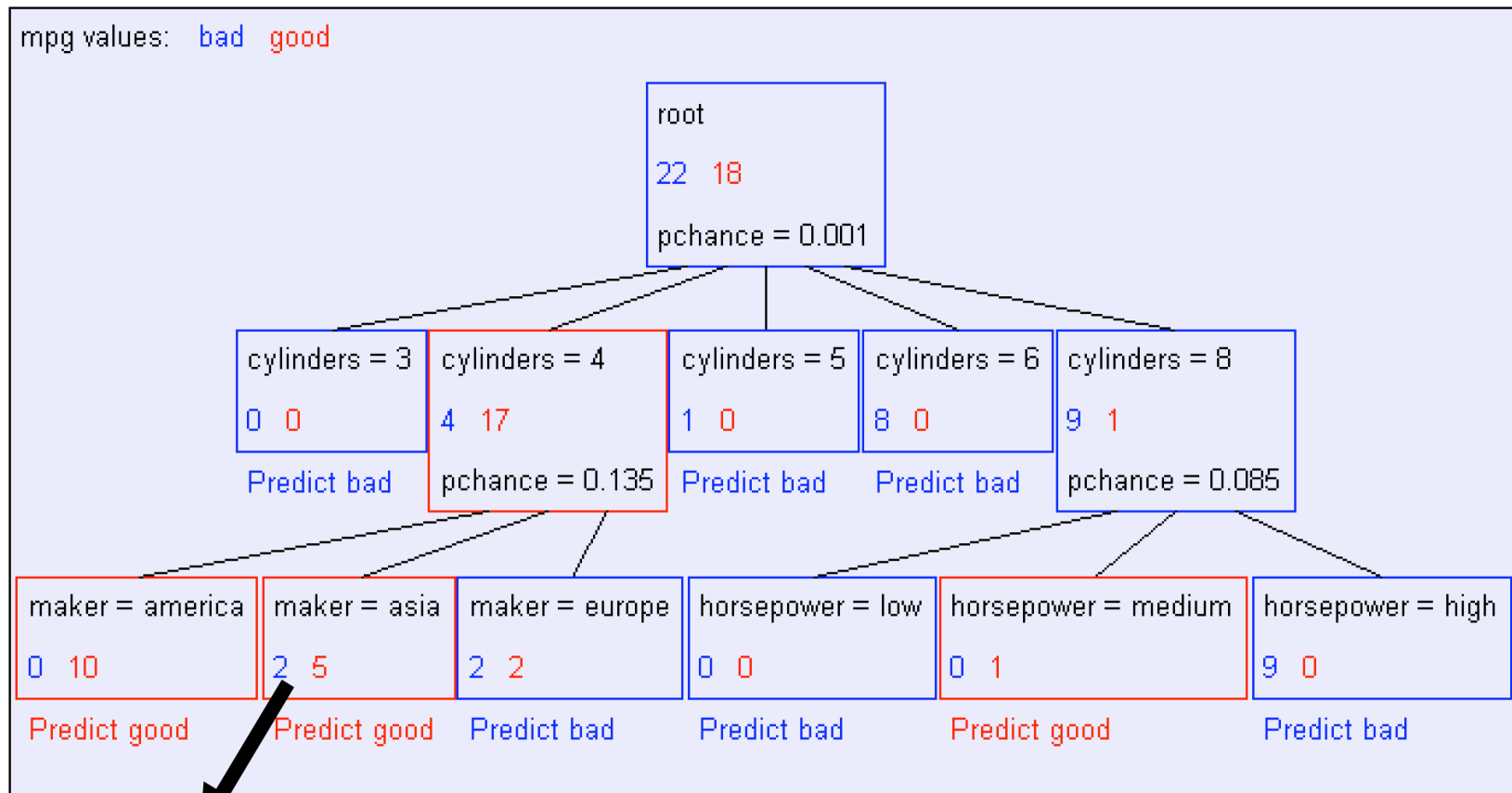
Records in
which
cylinders = 6

Build tree from
These records..



Records in
which
cylinders = 8

Second Level of Tree

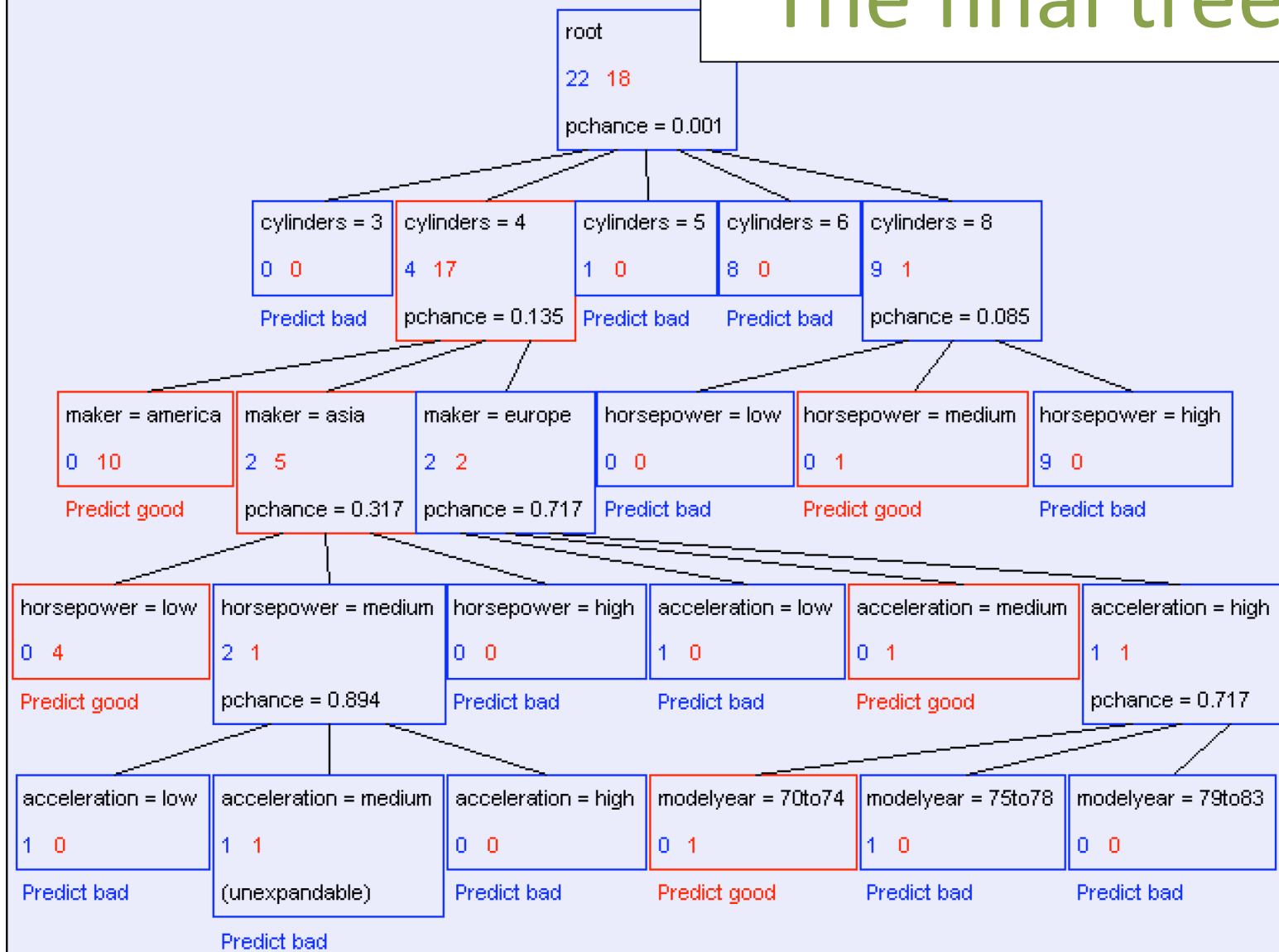


Recursively build a tree from the seven records in which there are four cylinders and the maker was based in Asia

(Similar recursion in the other cases)

The final tree

mpg values: bad good



Which attribute is the best?

A good split:

increases certainty about
classification **after split**

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

Entropy = measure of uncertainty

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

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

$H(Y)$ is the expected number of bits needed to encode a randomly drawn value of Y

Entropy = measure of uncertainty

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

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

$H(Y)$ is the expected number of bits needed to encode a randomly drawn value of Y

Why?

Entropy = measure of uncertainty

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

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

$H(Y)$ is the expected number of bits needed to encode a randomly drawn value of Y

Why?

Information Theory: most efficient code assigns
– $\log_2 P(Y=y_i)$ bits to message $Y=y_i$

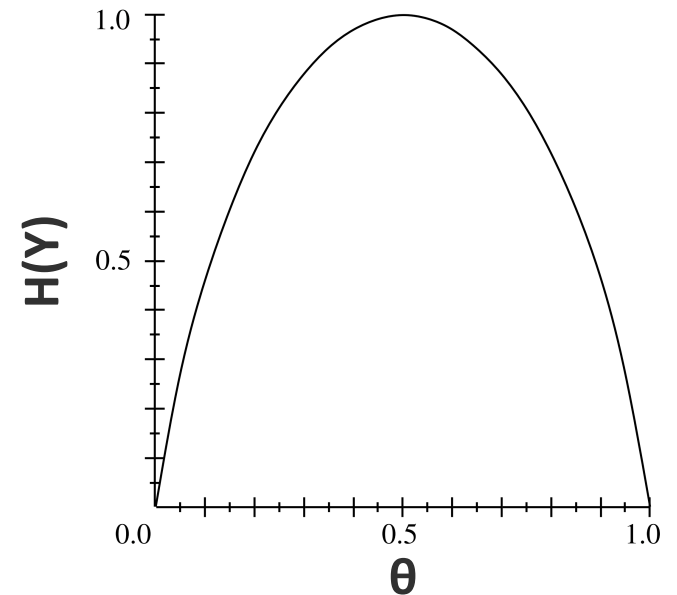
Entropy = measure of uncertainty

Y binary

$$P(Y=t) = \theta$$

$$P(Y=f) = 1 - \theta$$

$$H(Y) = \theta \log_2 \theta + (1 - \theta) \log_2 (1 - \theta)$$



Information Gain

= **reduction in uncertainty**

Entropy of Y before split:

$$H(Y)$$

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

Entropy of Y after split:

(weighted by probability of each branch)

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

Information gain = difference:

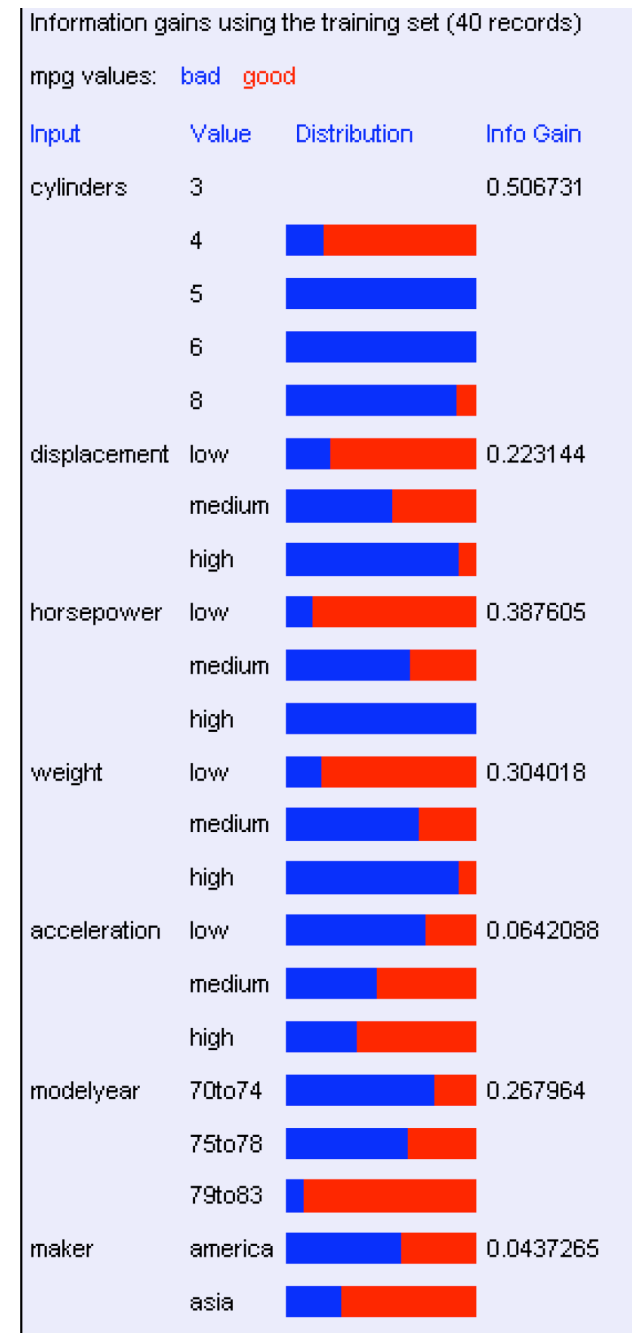
$$IG(X) = H(Y) - H(Y|X)$$

Learning decision trees

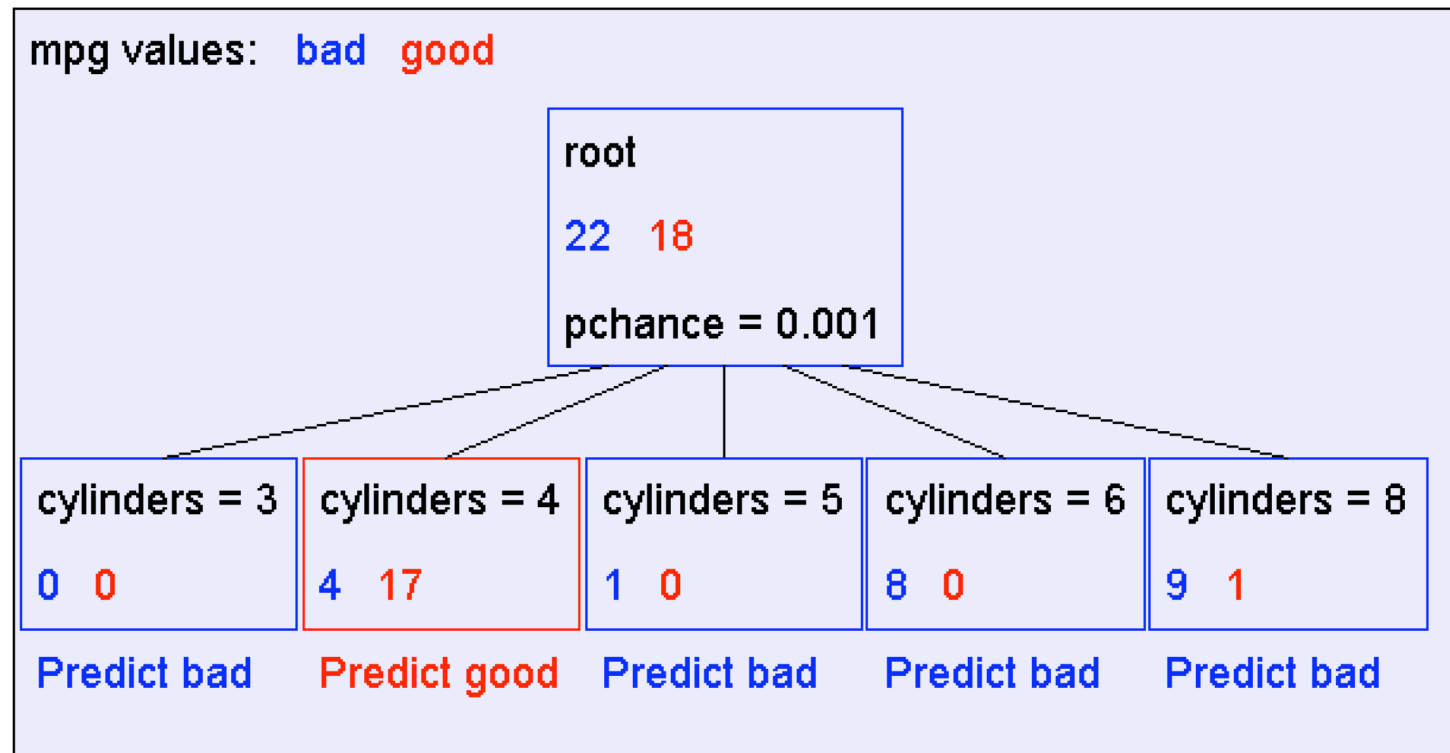
- start with an empty tree
- choose the **next best** attribute (feature)
 - for example, one that maximizes **information gain**
- split
- recurse

Suppose we want to
predict MPG.

Look at all
the
information
gains...

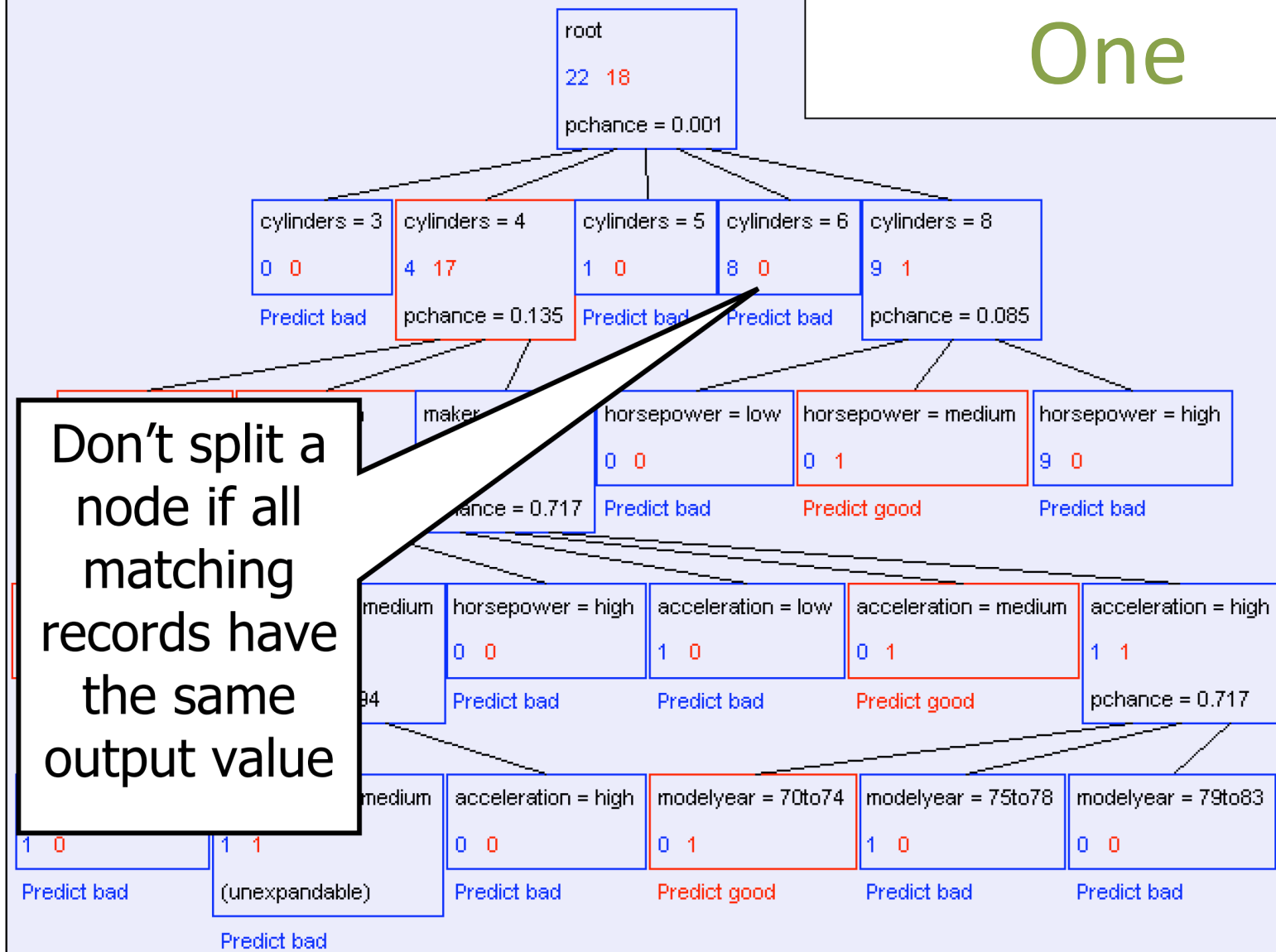


A Decision Stump

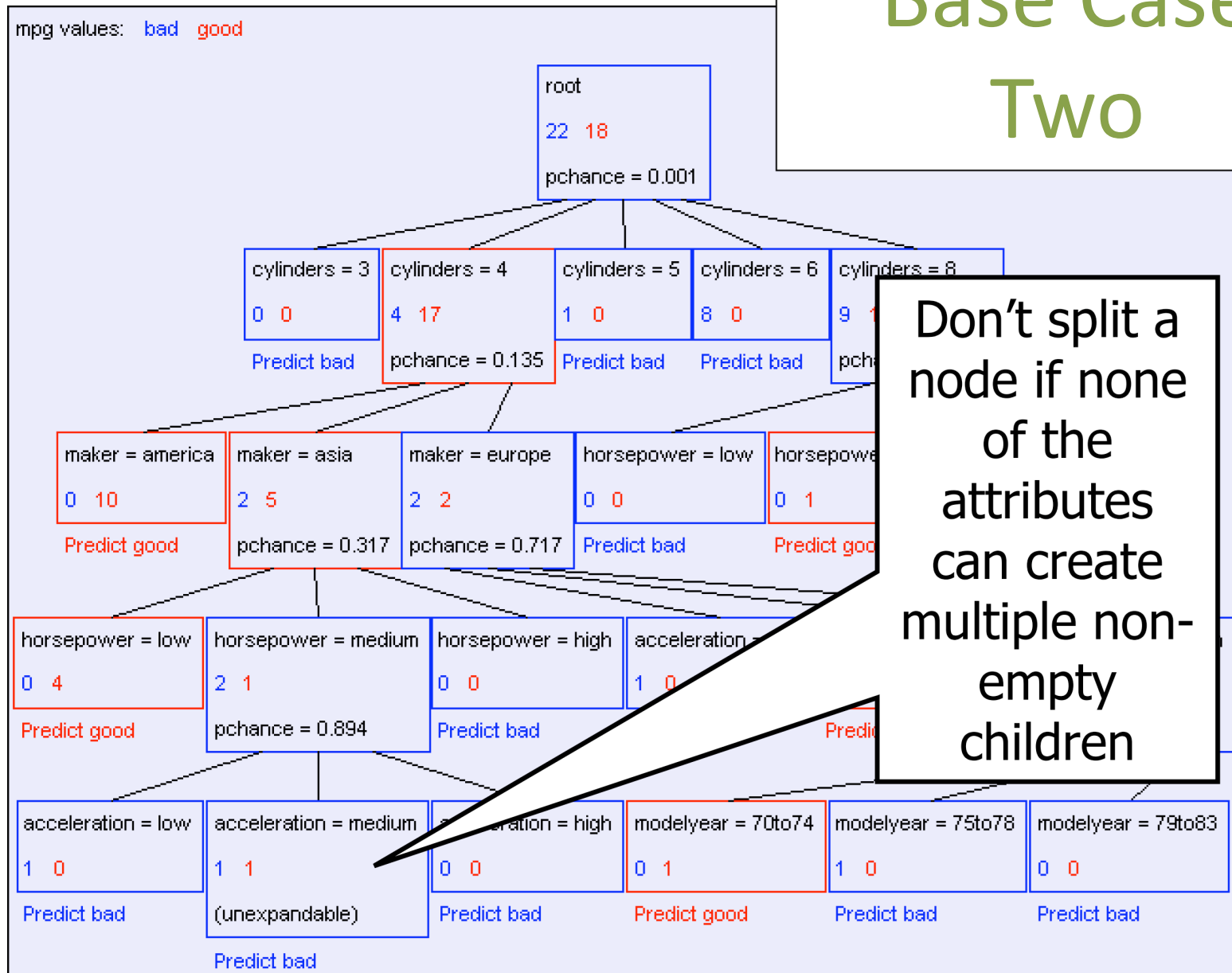


Base Case One

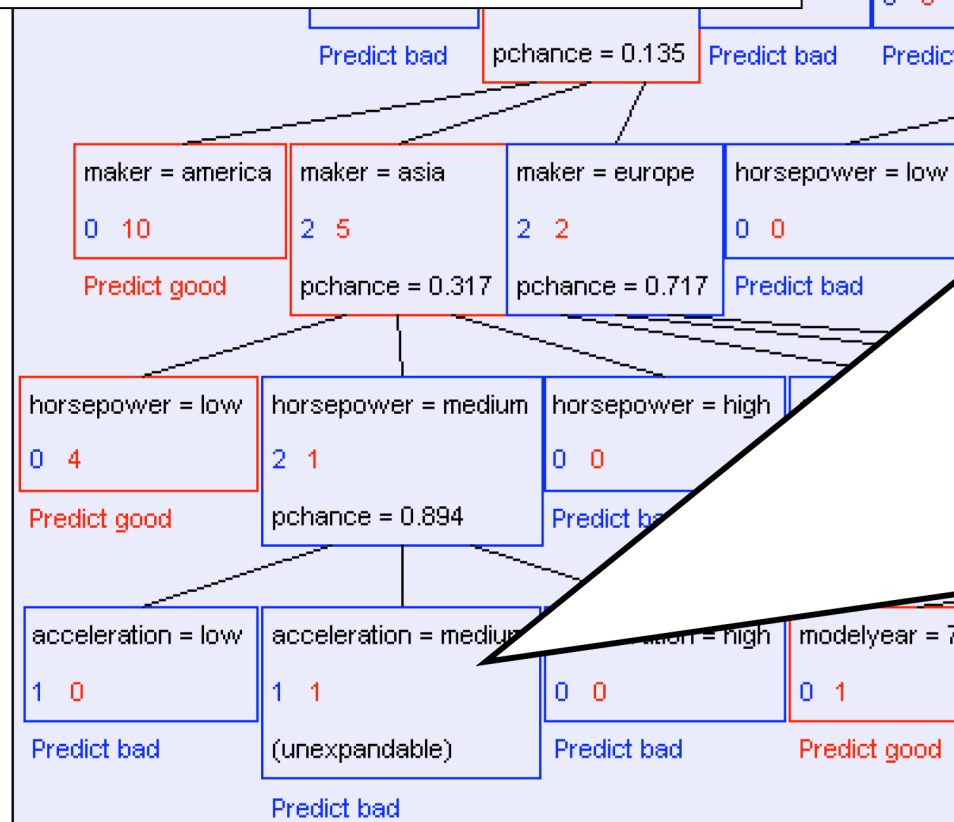
mpg values: bad good



Base Case Two



Base Case Two: attributes cannot distinguish classes



Information gains using the training set (2 records)

mpg values: bad good

Input	Value	Distribution	Info Gain
cylinders	3		0
	4		
	5		
	6		
	8		
displacement	low		0
	medium		
	high		
horsepower	low		0
	medium		
	high		
weight	low		0
	medium		
	high		
acceleration	low		0
	medium		
	high		
modelyear	70to74		0
	75to78		
	79to83		
maker	america		0
	asia		
	europe		

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**
- Base Case Two: If all records have exactly the same set of input attributes then **don't recurse**

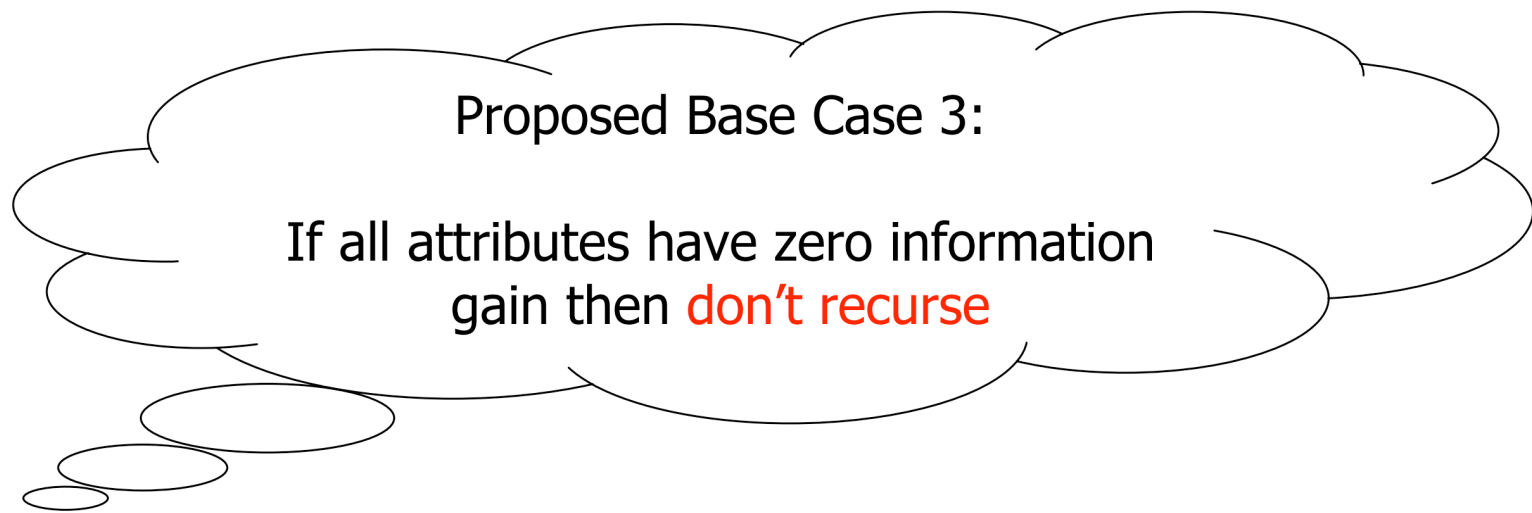
Proposed Base Case 3:

If all attributes have zero information gain then **don't recurse**

A thought bubble graphic with a large central oval and several smaller circles trailing off to the bottom left. The text is contained within the large central oval.

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

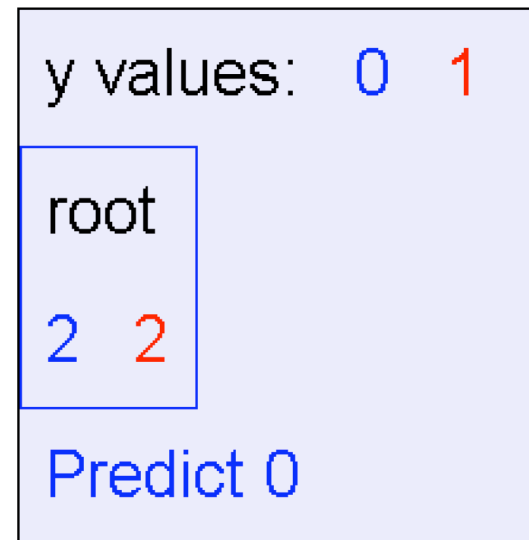
a	b	y
0	0	0
0	1	1
1	0	1
1	1	0

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

The information gains:

Information gains using the training set (4 records)			
y values: 0 1			
Input	Value	Distribution	Info Gain
a	0		0
	1		
b	0		0
	1		

The resulting decision tree:

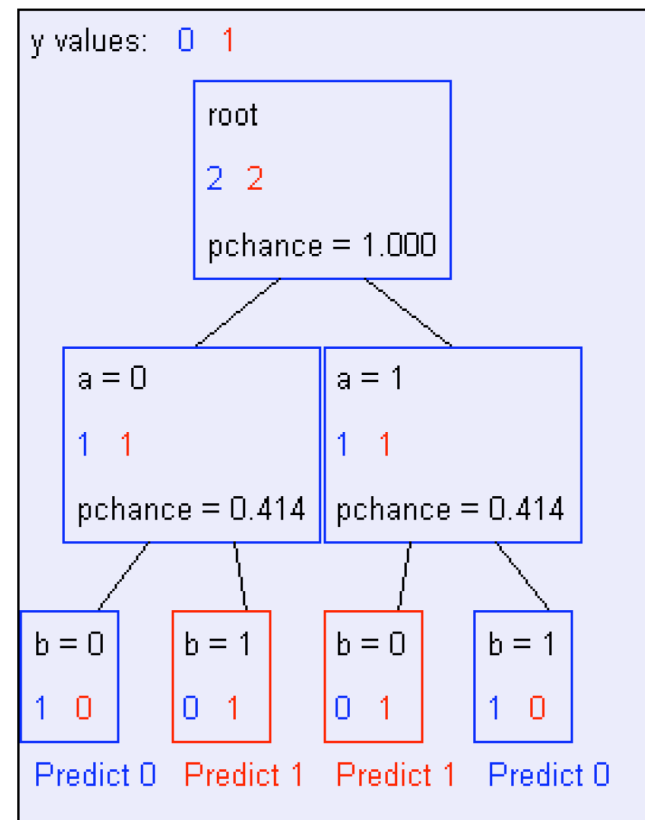


If we omit Base Case 3:

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*)

- 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 .

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

horsepower = high

Predict bad

horsepower = low

0 4

Predict good

horsepower = medium

2 1

pchance = 0.894

horsepower = high

0 0

Predict bad

acceleration = low

1 0

Predict bad

acceleration = medium

0 1

Predict good

acceleration = high

1 1

pchance = 0.717

acceleration = low

1 0

Predict bad

acceleration = medium

1 1

(unexpandable)

Predict bad

acceleration = high

0 0

Predict bad

modelyear = 70to74

0 1

Predict good

modelyear = 75to78

1 0

Predict bad

modelyear = 79to83

0 0

Predict bad

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

horsepower = high

Predict bad

horsepower = low

horsepower = medium

horsepower = high

acceleration = low

acceleration = medium

acceleration = high

0 4

0 4

0 0

4 0

0 4

4 4

Predict

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

...why?

= 0.717

= 79to83

Predict bad

(unexpandable)

Predict bad

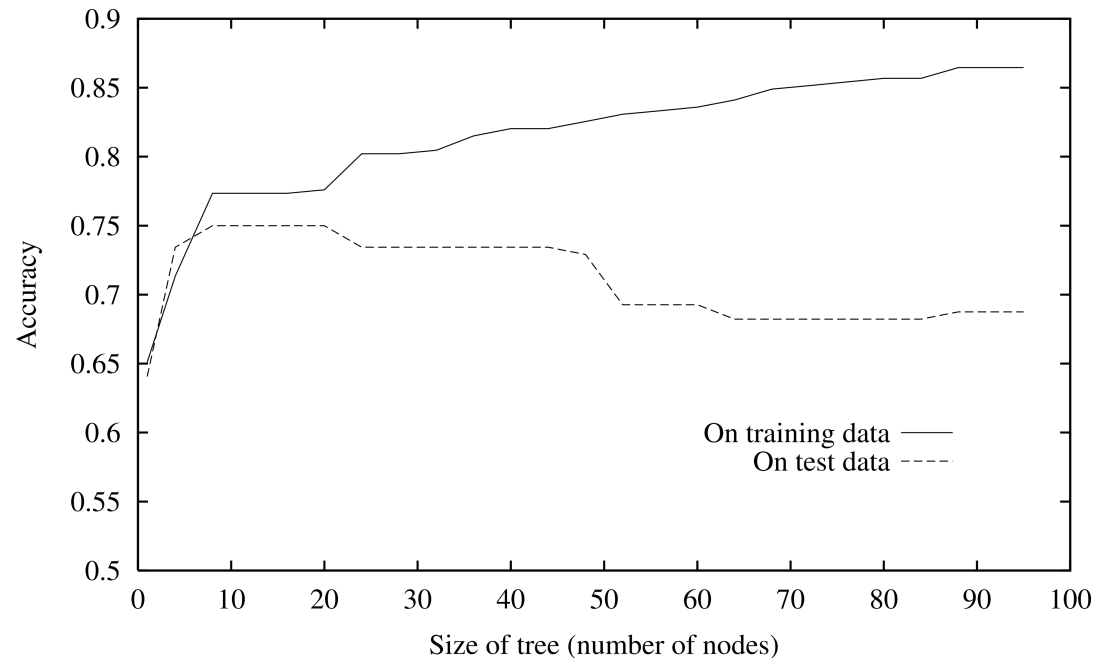
Predict good

Predict bad

Predict bad

Predict bad

Decision trees overfit!



Standard decision trees:

- **training error always zero** (if no label noise)
- **lots of variance**

Avoiding overfitting

- fixed depth
- fixed number of leaves
- stop when splits not statistically significant

Avoiding overfitting

- fixed depth
- fixed number of leaves
- stop when splits not statistically significant

OR:

- grow the full tree,
then **prune**
(collapse some subtrees)

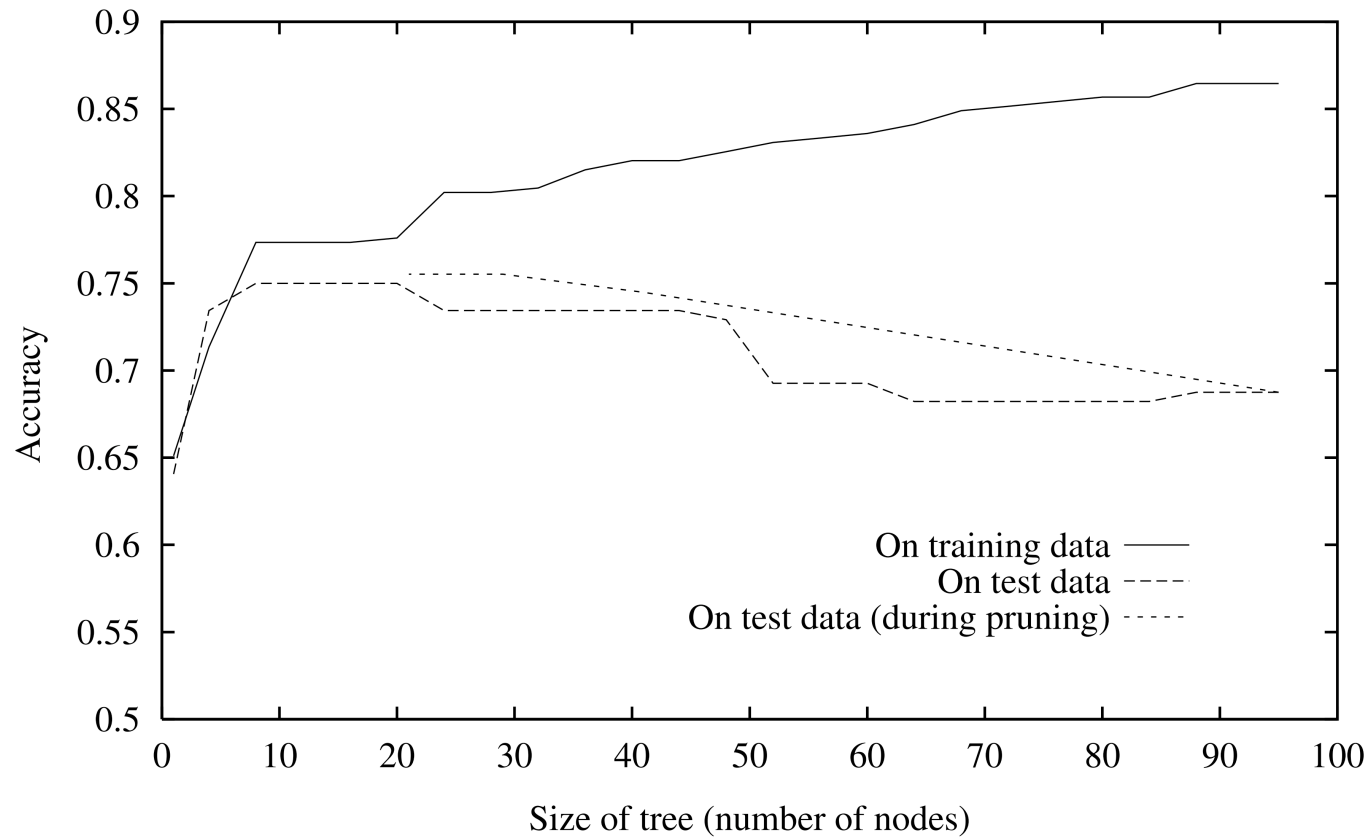
Reduced Error Pruning

Split available data into **training** and **pruning** sets

1. Learn tree that classifies **training set perfectly**
2. Do **until further pruning** is **harmful** over pruning set
 - consider pruning each node
 - collapse the node that best **improves pruning set accuracy**

This produces smallest version of
most accurate tree (over the pruning set)

Impact of Pruning



A Generic Tree-Learning Algorithm

Need to specify:

- an objective to select **splits**
- a criterion for **pruning** (or **stopping**)
- **parameters** for pruning/stopping
(usually determined by **cross-validation**)

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

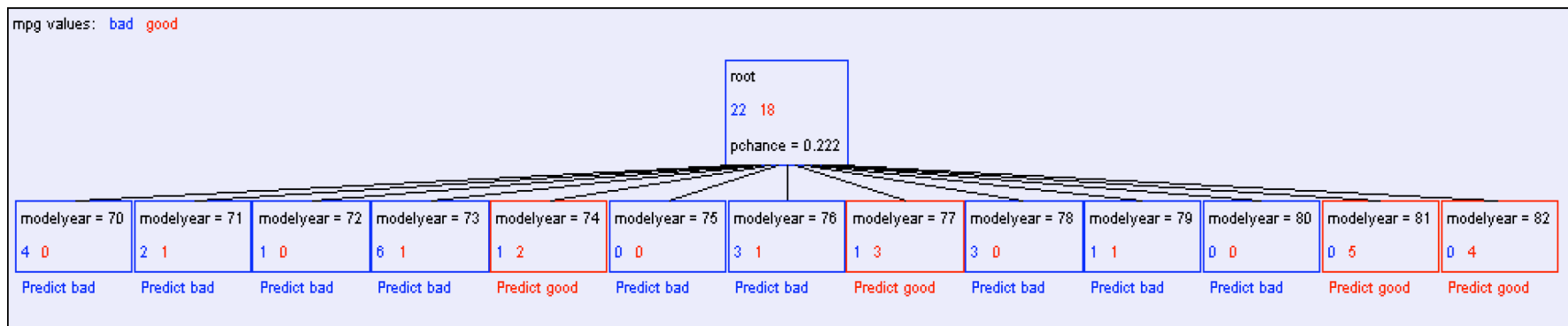
mpg	cylinders	displacemen	horsepower	weight	acceleration	modelyear	maker
good	4	97	75	2265	18.2	77	asia
bad	6	199	90	2648	15	70	america
bad	4	121	110	2600	12.8	77	europa
bad	8	350	175	4100	13	73	america
bad	6	198	95	3102	16.5	74	america
bad	4	108	94	2379	16.5	73	asia
bad	4	113	95	2228	14	71	asia
bad	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

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

Idea One: Branch on each possible real value

“One branch for each numeric value” idea:



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
















A better idea:

thresholded splits

- Binary tree, split on attribute X :
 - one branch: $X < t$
 - other branch: $X \geq t$
- Search through **all possible values of t**
 - seems hard, but only finite set relevant
 - sort values of X : $\{x_1, \dots, x_m\}$
 - consider splits at $t = (x_i + x_{i+1})/2$
- Information gain for **each split**
as if a **binary variable**: “true” for $X < t$
“false” for $X \geq t$

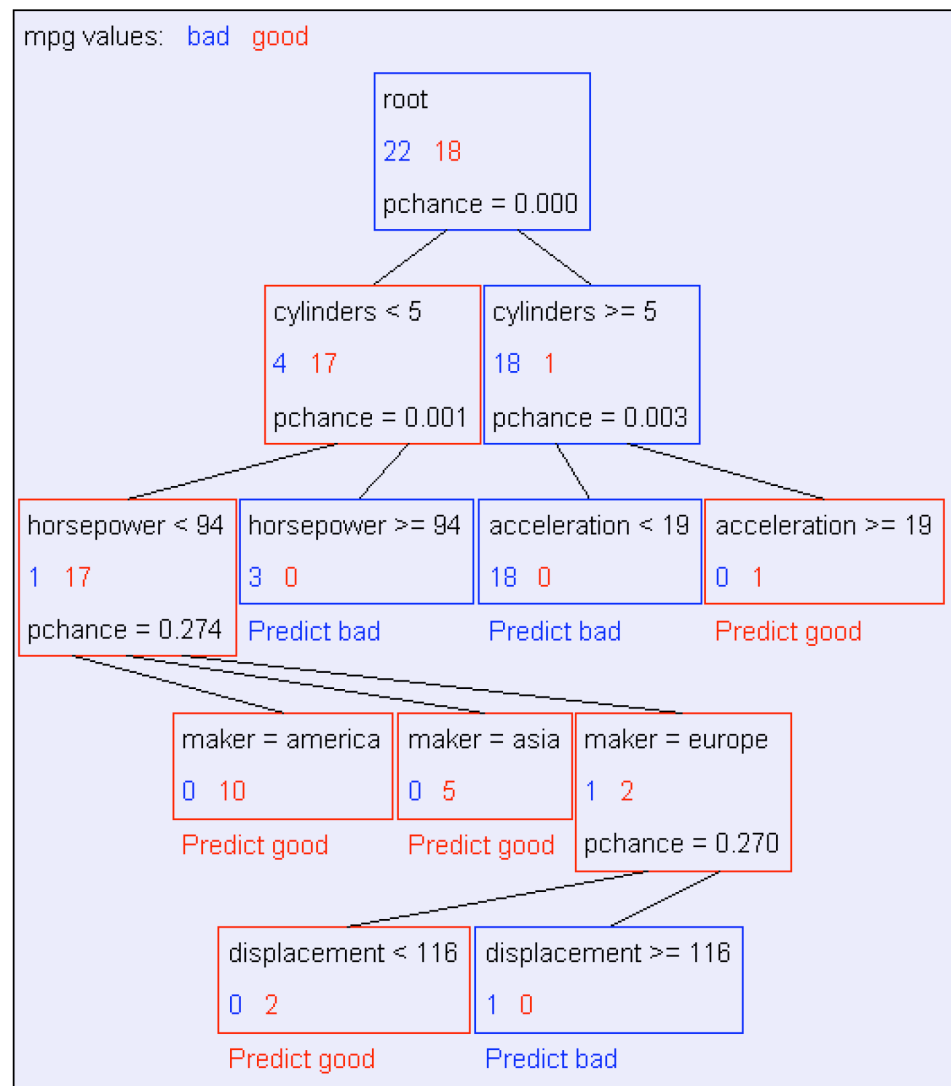
Information gains using the training set (40 records)

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		
	europe		

Example with MPG

Example tree using reals



What you should know about decision trees

- among **most popular** data mining tools:
 - easy to understand
 - easy to implement
 - easy to use
 - computationally fast (but only a greedy heuristic!)
- not only **classification**, also **regression**, **density estimation**
- meaning of information gain
- decision trees overfit!
 - many pruning/stopping strategies

Acknowledgements

Some material in this presentation is courtesy of
Andrew Moore, from his collection of ML tutorials:
<http://www.autonlab.org/tutorials/>

LEARNING THEORY

Computational Learning Theory

What general laws constrain “learning”?

- how many **examples** needed to learn a **target concept** to a given **precision**?
- what is the impact of:
 - complexity of the **target concept**?
 - complexity of our **hypothesis space**?
 - **manner** in which examples presented?
 - random samples—what we mostly consider in this course
 - learner can make queries
 - examples come from an “adversary”
(worst-case analysis, no statistical assumptions)