

# BBM406

## Fundamentals of Machine Learning

Lecture 18:  
Decision Trees

# Today

- Decision Trees
- Tree construction
- Overfitting
- Pruning
- Real-valued inputs

# Machine Learning in the ER

Triage Information  
(Free text)



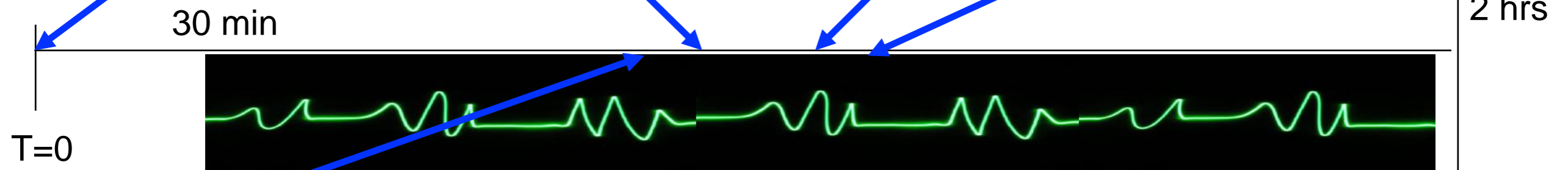
MD comments  
(free text)



Specialist consults



Physician  
documentation



T=0

30 min

2 hrs



Lab results  
(Continuous valued)

Repeated vital signs  
(continuous values)  
Measured every 30 s

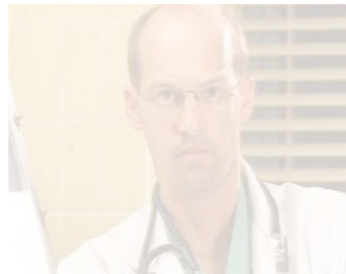
Disposition

# Can we predict infection?

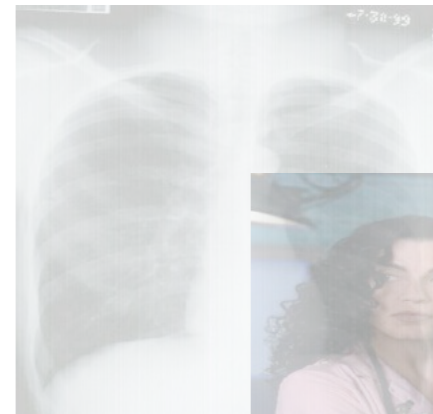
Triage Information  
(Free text)



MD comments  
(free text)



Specialist consults



Physician  
documentation

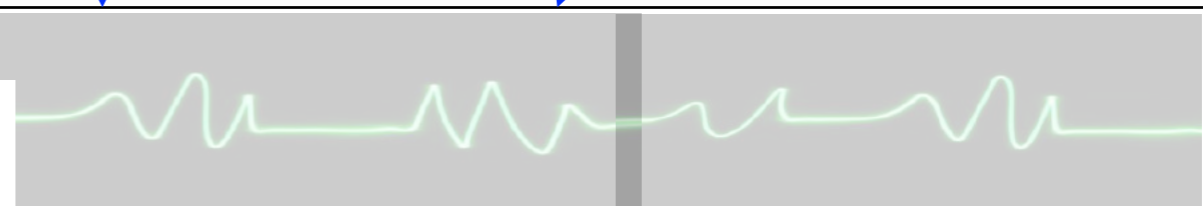


Many crucial decisions  
about a patient's care are  
made here!



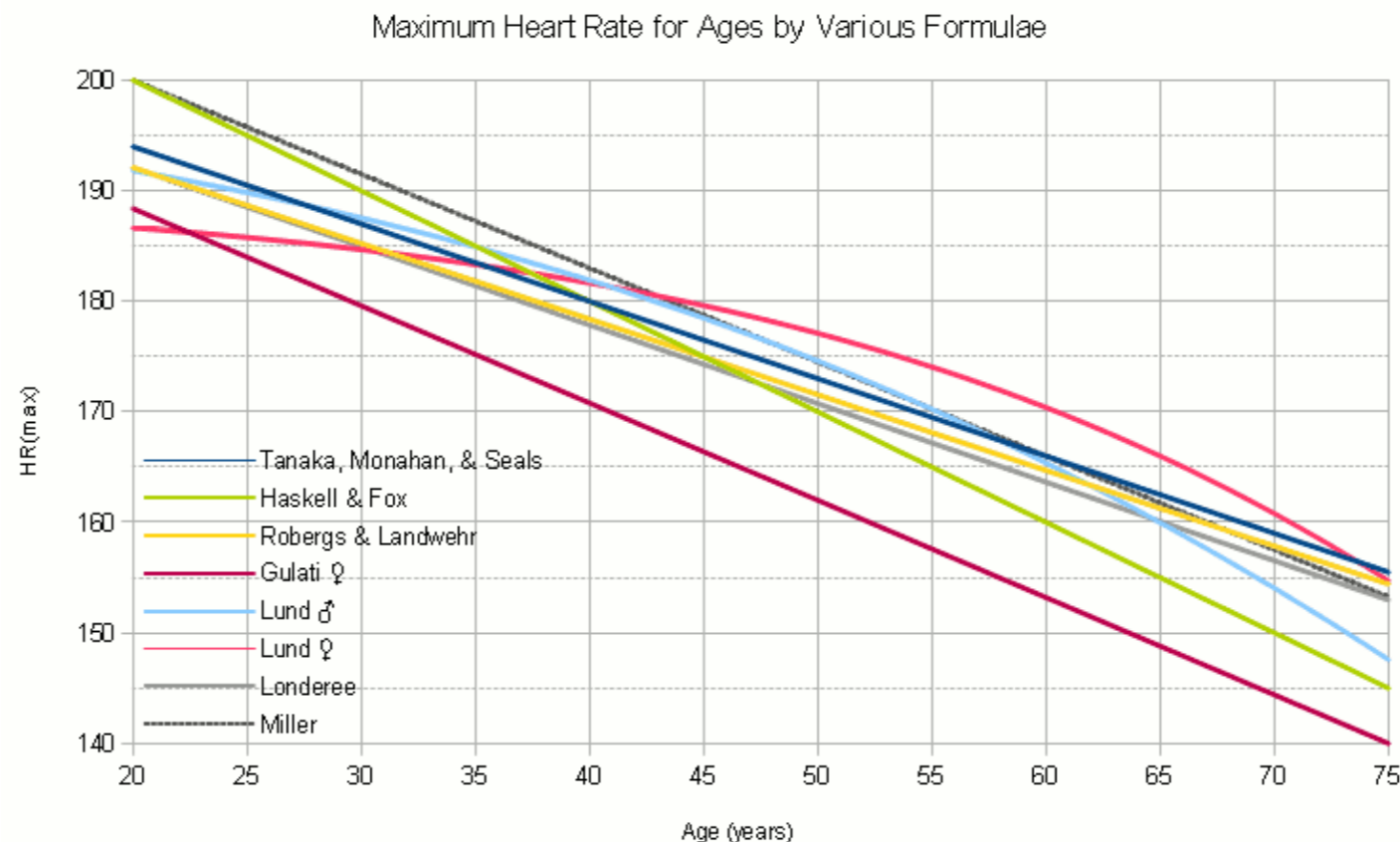
Lab results  
(Continuous valued)

Repeated vital signs  
(continuous values)  
Measured every 30 s



# Can we predict infection

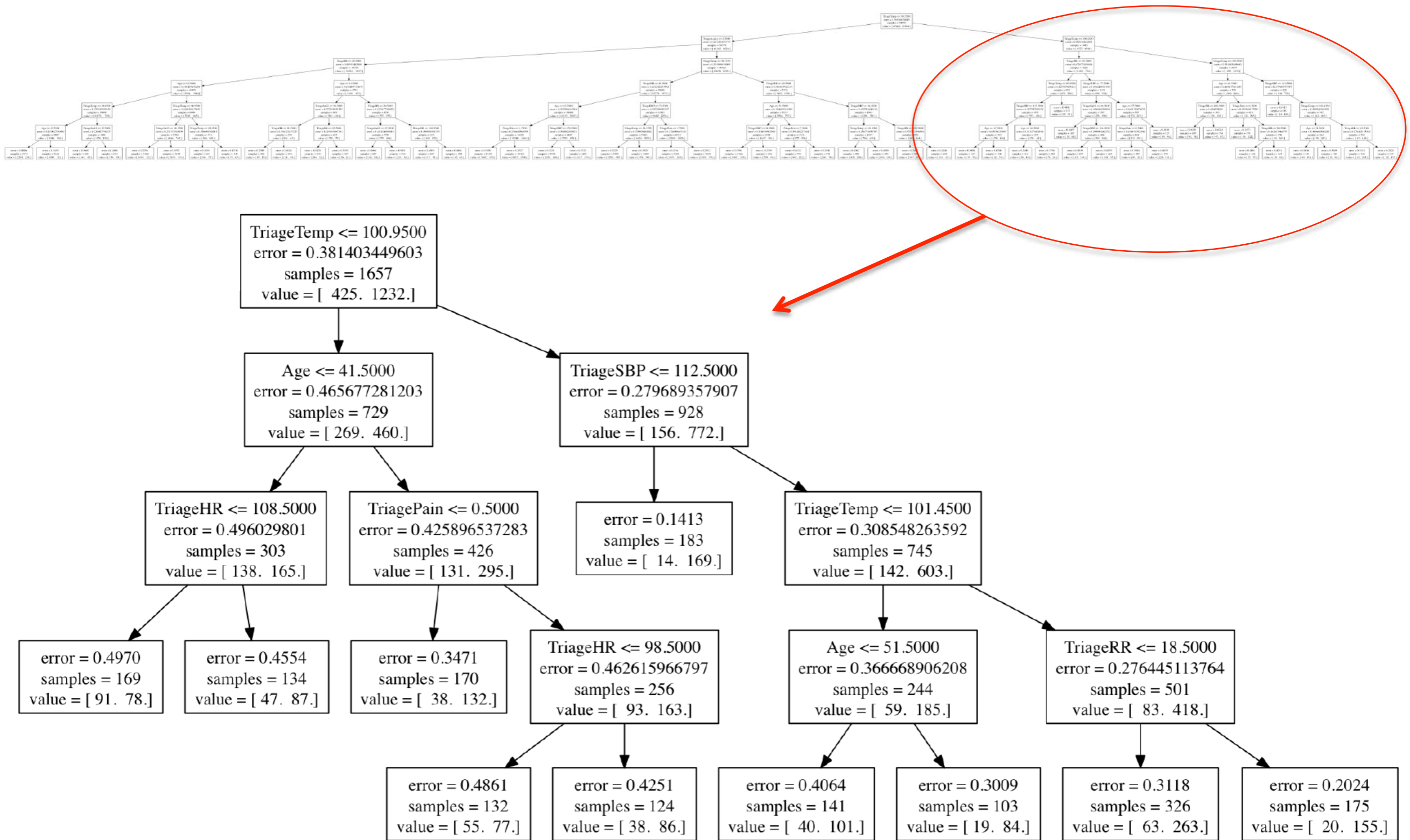
- Previous automatic approaches based on simple criteria:
  - Temperature  $< 96.8$  °F or  $> 100.4$  °F
  - Heart rate  $> 90$  beats/min
  - Respiratory rate  $> 20$  breaths/min
- Too simplified... e.g., heart rate depends on age!



# Can we predict infection?

- These are the attributes we have for each patient:
  - Temperature
  - Heart rate (HR)
  - Respiratory rate (RR)
  - Age
  - Acuity and pain level
  - Diastolic and systolic blood pressure (DBP, SBP)
  - Oxygen Saturation (SaO<sub>2</sub>)
- We have these attributes + label (infection) for 200,000 patients!
- Let's **learn** to classify infection

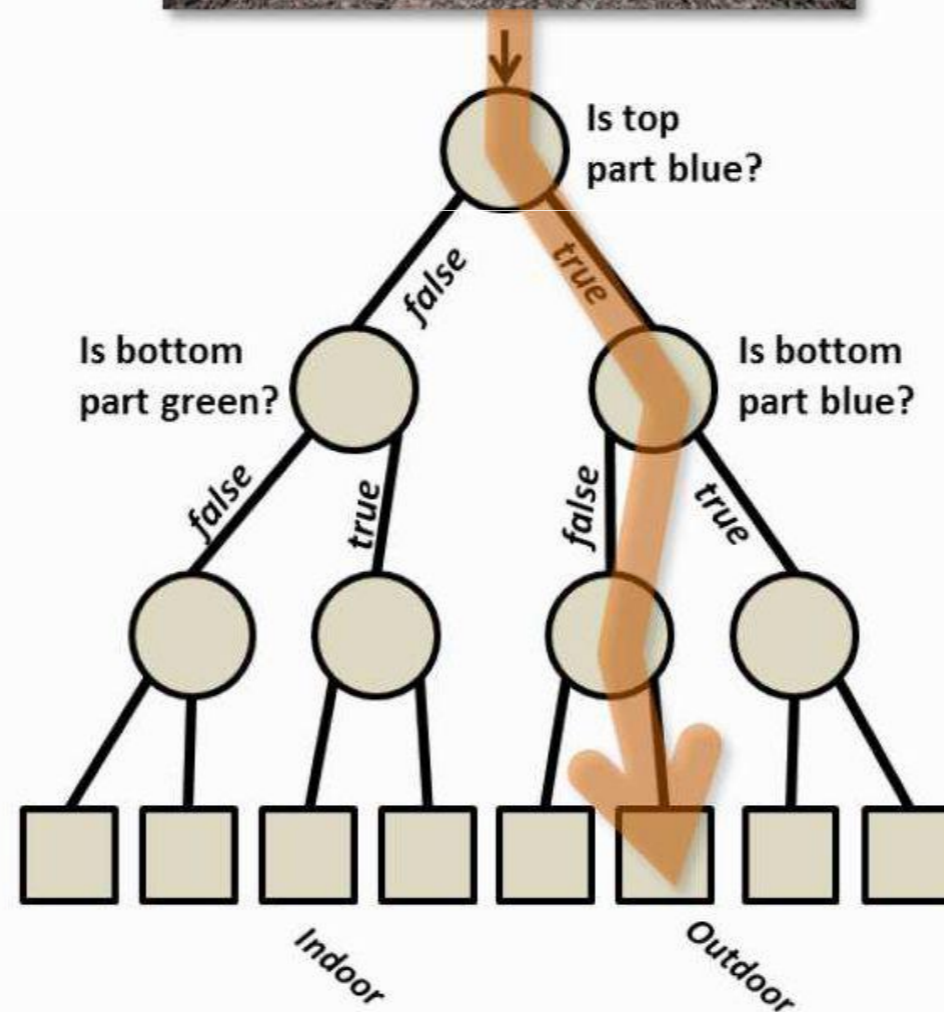
# Predicting infection using decision trees



# Example: Image Classification



A decision tree

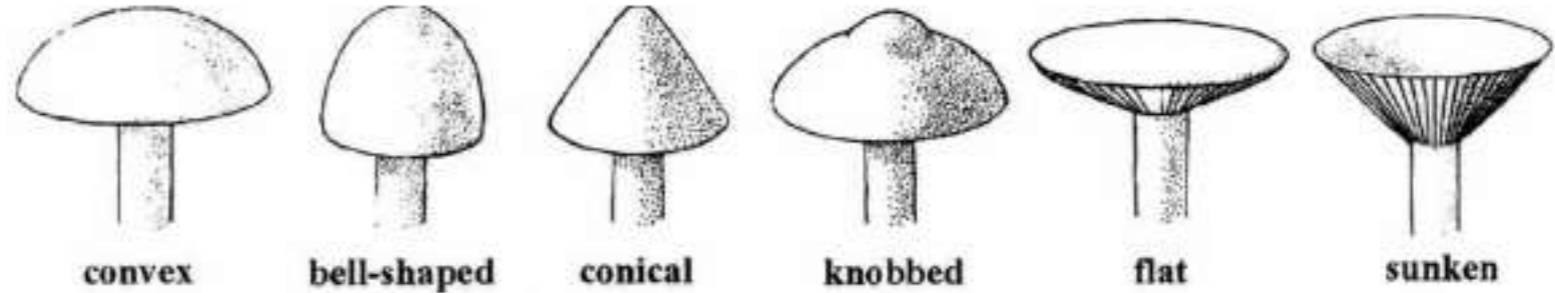




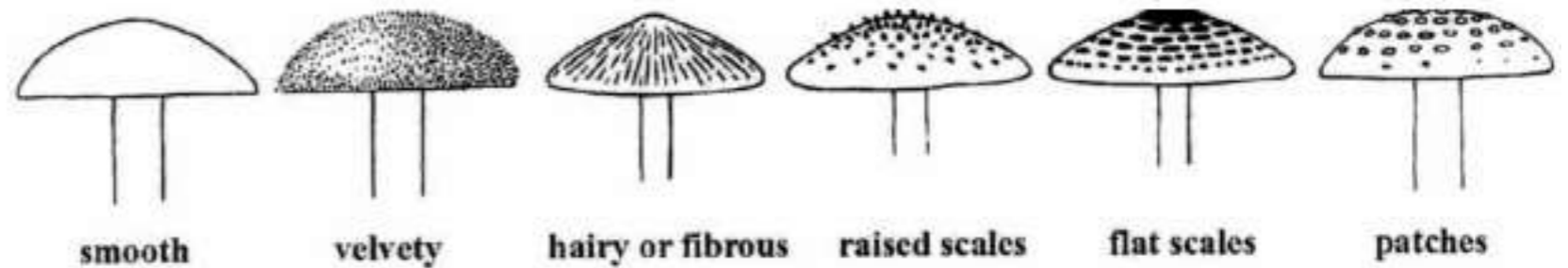
# Example: Mushrooms



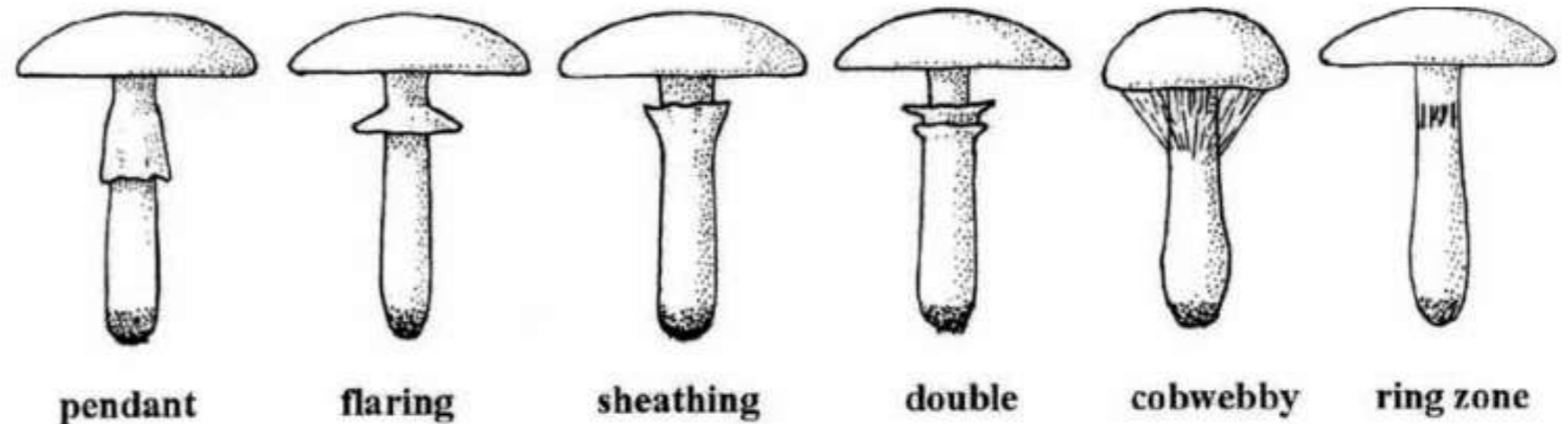
## Mushroom cap shapes



## Mushroom cap surfaces



## Annular rings



# Mushroom features

1. **cap-shape:** bell=b, conical=c, convex=x, flat=f, knobbed=k, sunken=s
2. **cap-surface:** fibrous=f, grooves=g, scaly=y, smooth=s
3. **cap-color:** brown=n, buff=b, cinnamon=c, gray=g, green=r, pink=p, purple=u, red=e, white=w, yellow=y
4. **bruises?:** bruises=t, no=f
5. **odor:** almond=a, anise=l, creosote=c, fishy=y, foul=f, musty=m, none=n, pungent=p, spicy=s
6. **gill-attachment:** attached=a, descending=d, free=f, notched=n
7. ...

# Two mushrooms

$x_1 = x, s, n, t, p, f, c, n, k, e, e, s, s, w, w, p, w, o, p, k, s, u$

$y_1 = p$

$x_2 = x, s, y, t, a, f, c, b, k, e, c, s, s, w, w, p, w, o, p, n, n, g$

$y_2 = e$

1. cap-shape: bell=b, conical=c, convex=x, flat=f, knobbed=k, sunken=s
2. cap-surface: fibrous=f, grooves=g, scaly=y, smooth=s
3. cap-color:  
brown=n, buff=b, cinnamon=c, gray=g, green=r,  
pink=p, purple=u, red=e, white=w, yellow=y
4. ...

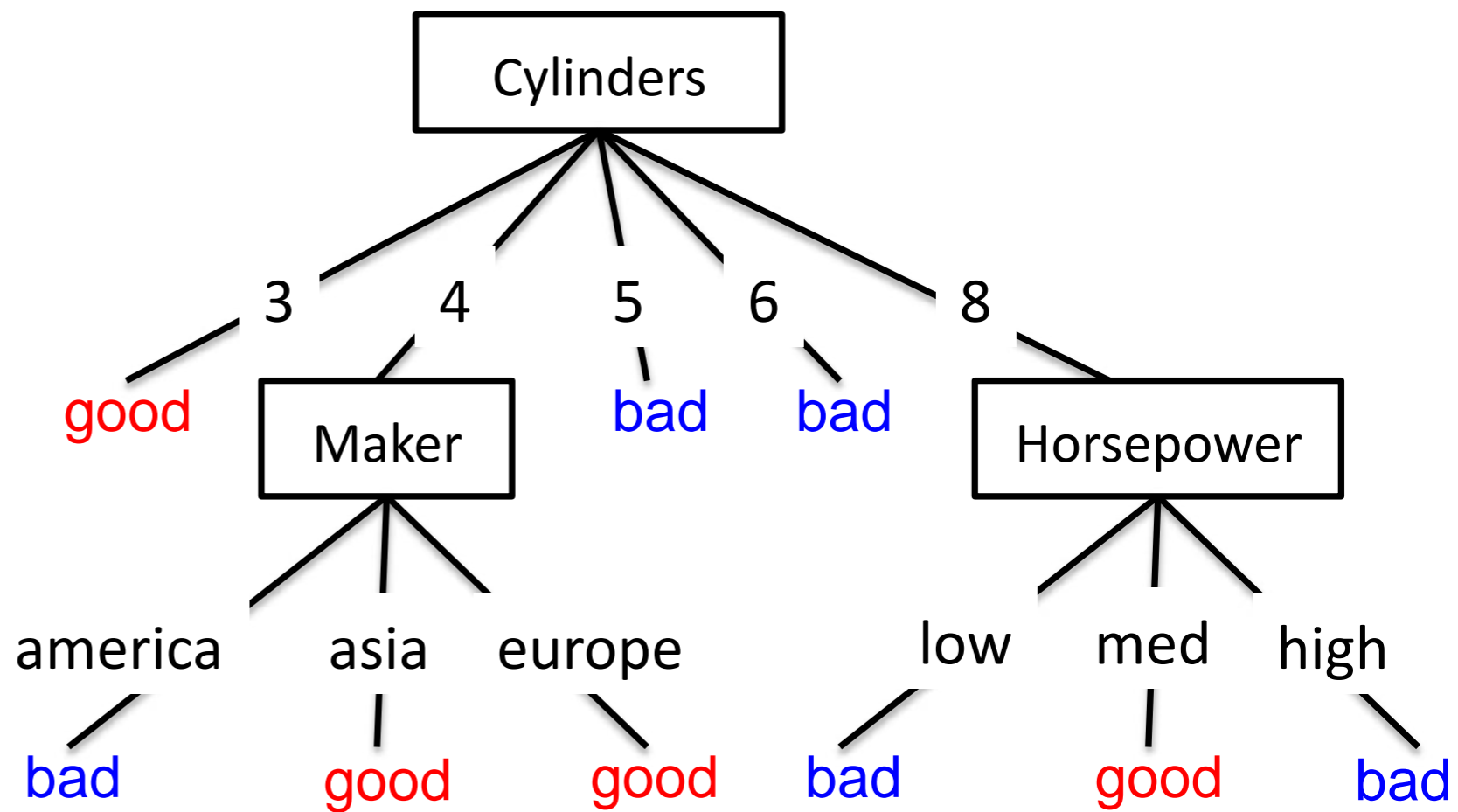
# Example: Automobile Miles-per-gallon prediction



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

# Hypotheses: decision trees $f : X \rightarrow Y$

- Each internal node tests an attribute  $x_i$
- Each branch assigns an 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  $y$

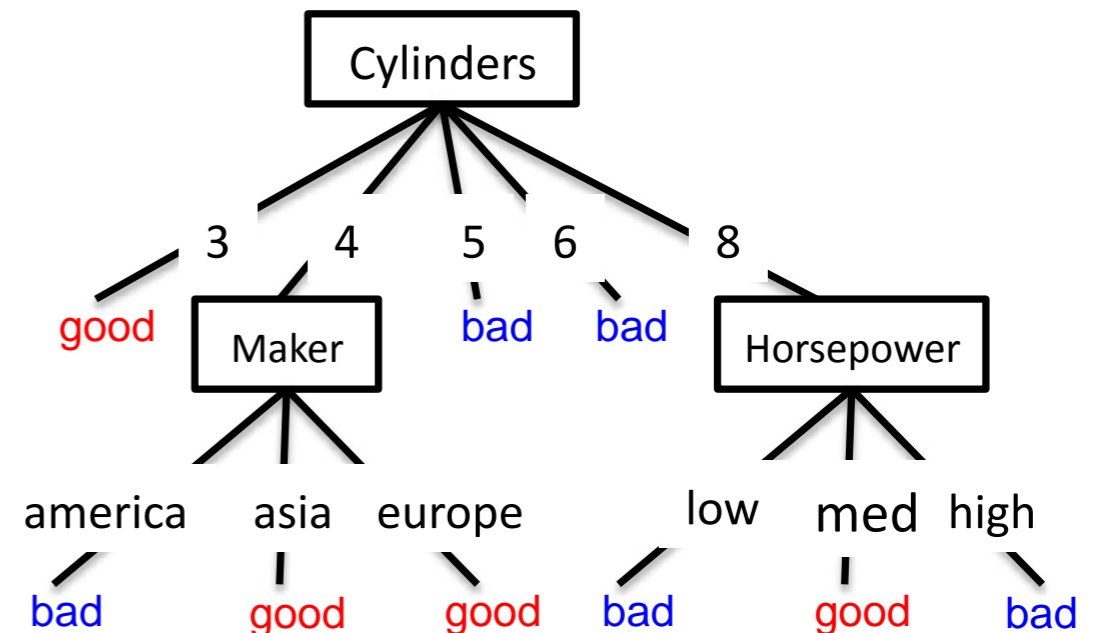


**Human interpretable!**

# Hypothesis space

- How many possible hypotheses?
- What functions can be represented?

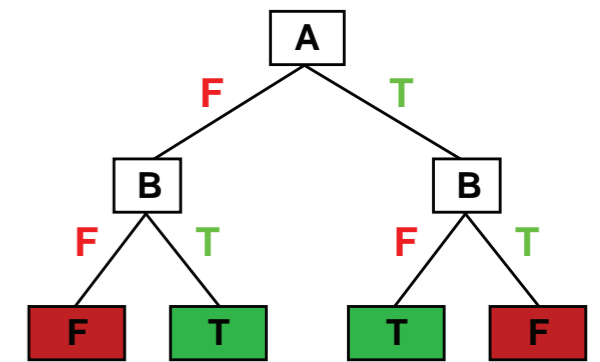
mpg	cylinders	displacement	horsepower	weight	acceleration	modelyear	maker
good	4	low	low	low	high	75to78	asia
bad	6	medium	medium	medium	medium	70to74	america
bad	4	medium	medium	medium	low	75to78	europa
bad	8	high	high	high	low	70to74	america
bad	6	medium	medium	medium	medium	70to74	america
bad	4	low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



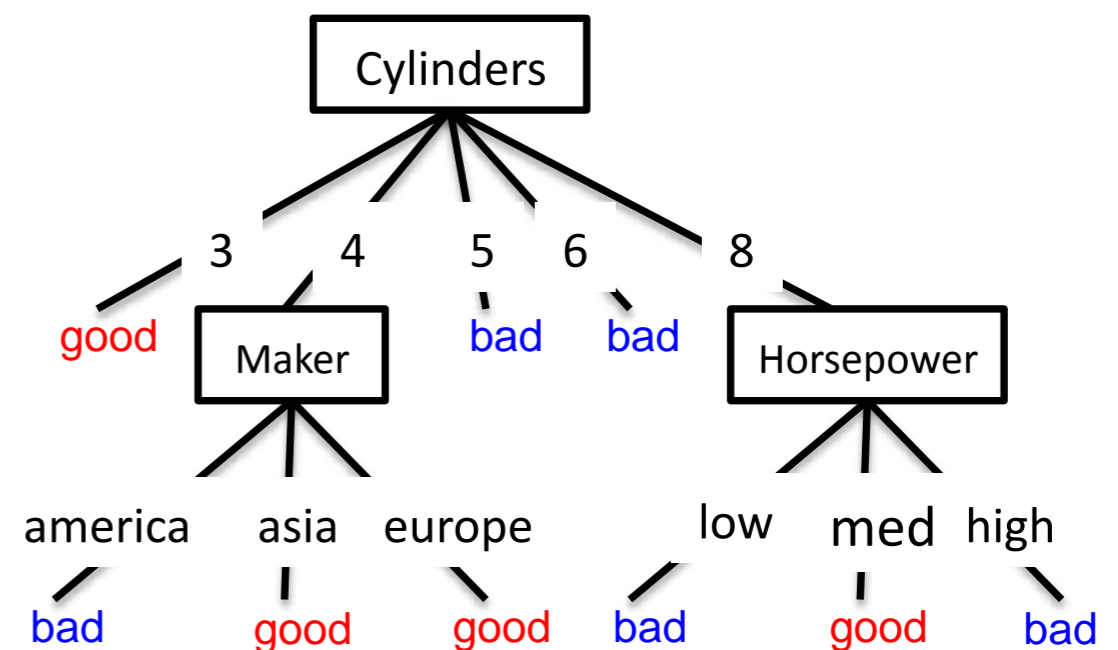
# What functions can be represented?

- Decision trees can represent any function of the input attributes!
- For Boolean functions, path to leaf gives truth table row
- But, could require exponentially many nodes...

A	B	A xor B
F	F	F
F	T	T
T	F	T
T	T	F



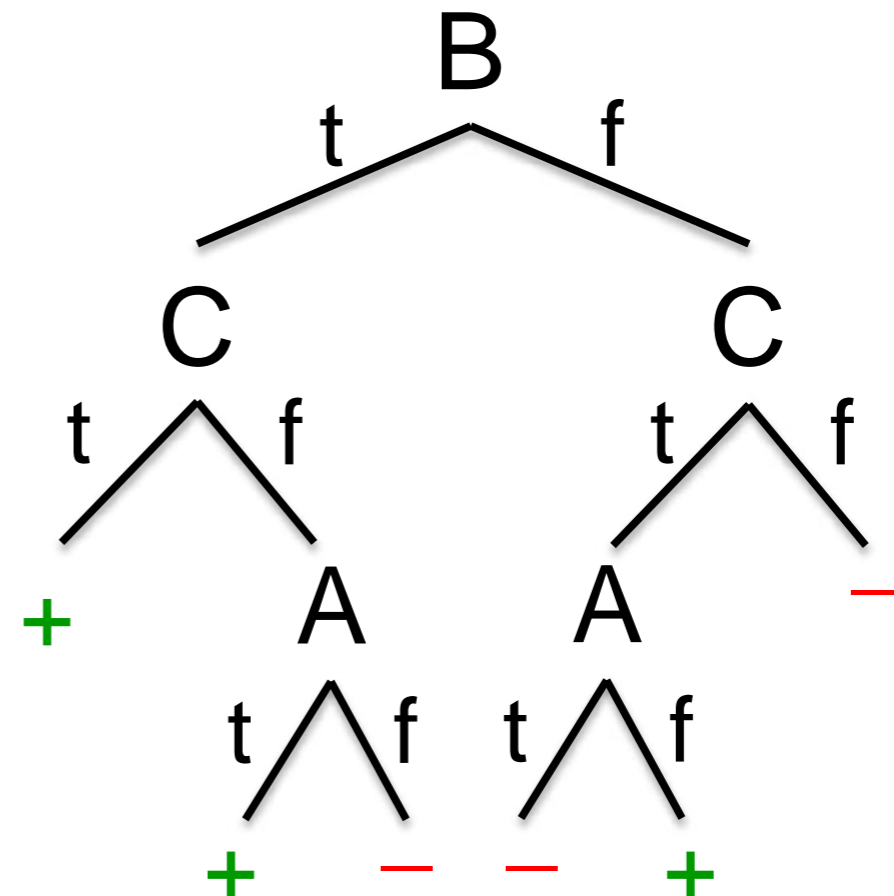
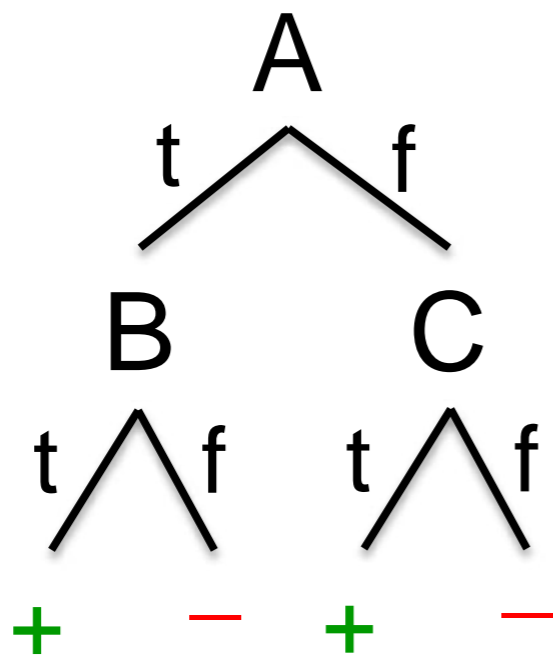
(Figure from Stuart Russell)



$$\text{cyl}=3 \vee (\text{cyl}=4 \wedge (\text{maker}=\text{asia} \vee \text{maker}=\text{europe})) \vee \dots$$

# 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))



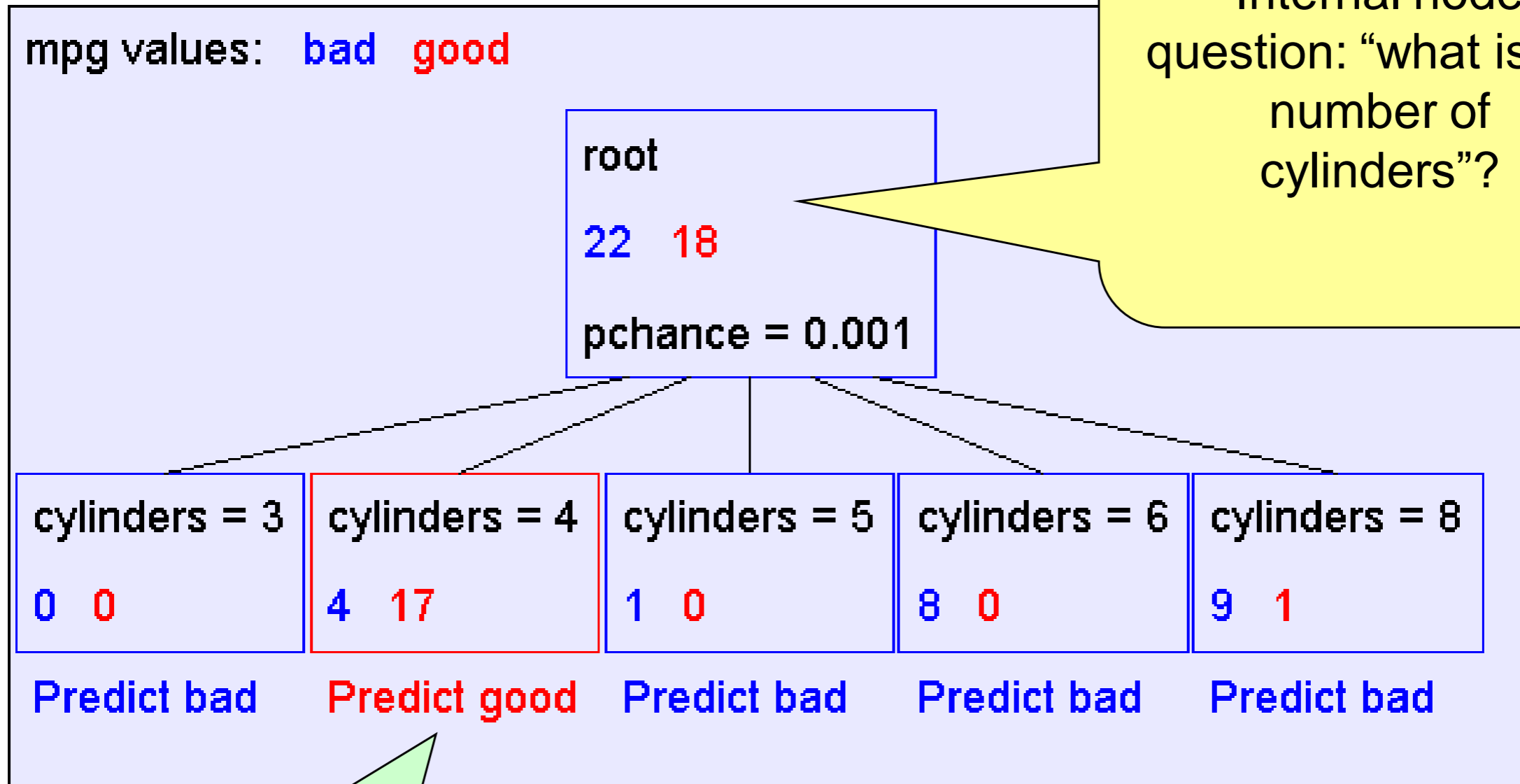
- Which tree do we prefer?



# Learning decision trees is hard!!!

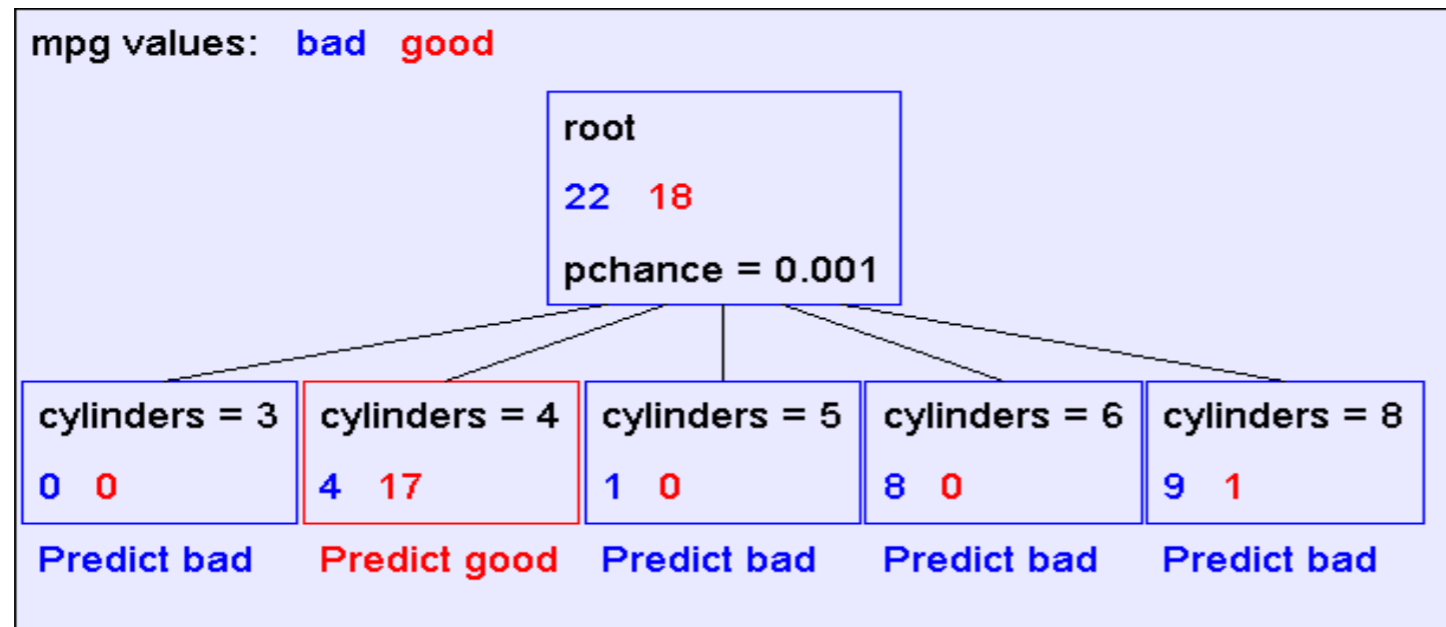
- Learning the simplest (smallest) decision tree is an NP-complete problem [Hyafil & Rivest '76]
- Resort to a greedy heuristic:
  - Start from empty decision tree
  - Split on next best attribute (feature)
  - Recurse

# A Decision Stump



Leaves: classify by majority vote

# Key idea: Greedily learn trees using recursion



Records in which cylinders = 4

Records in which cylinders = 5

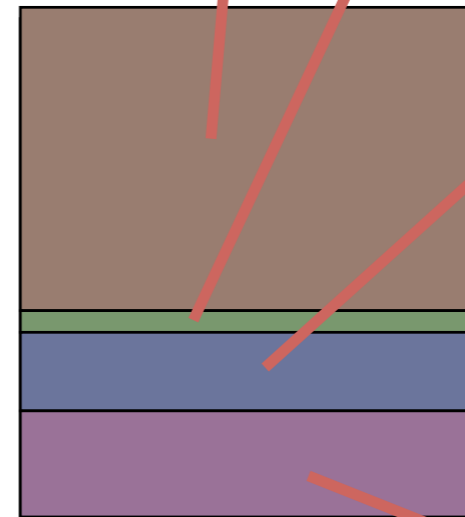
Records in which cylinders = 6

Records in which cylinders = 8

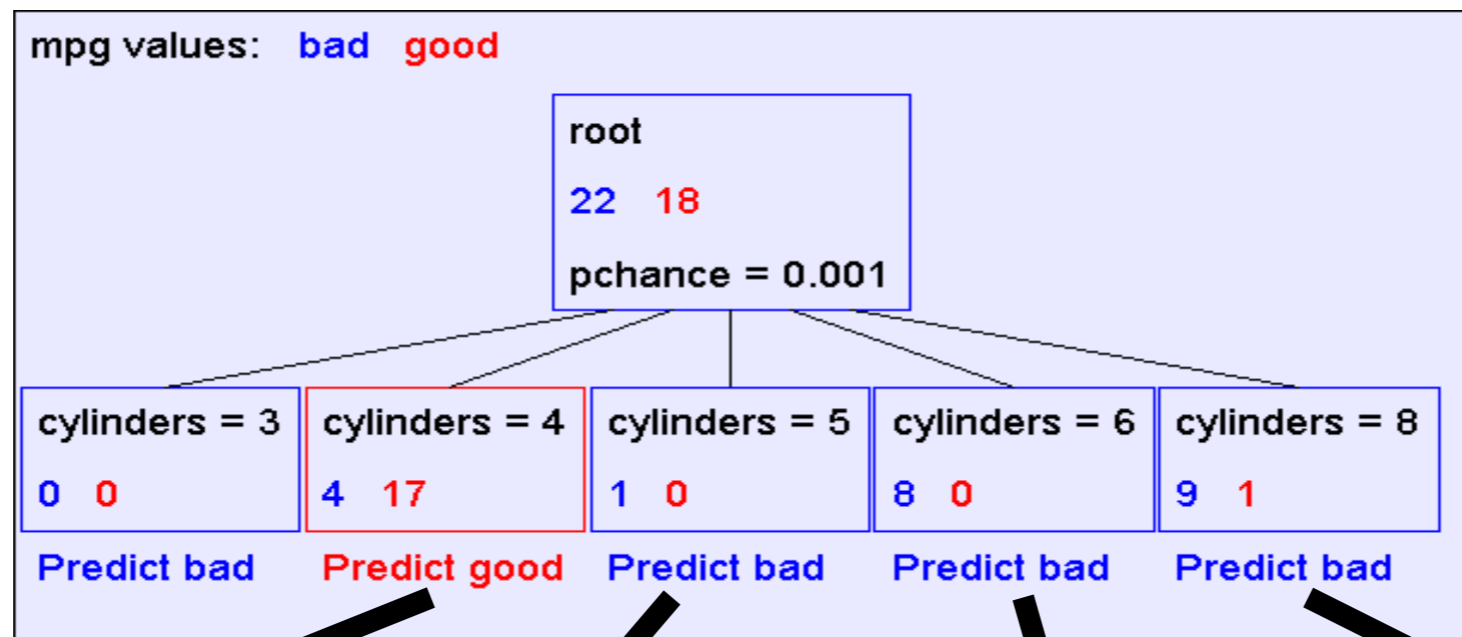
Take the Original Dataset..



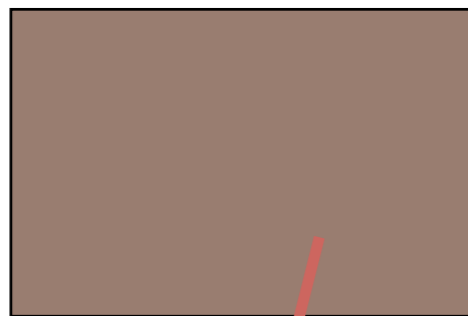
And partition it according to the value of the attribute we split on



# Recursive 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..



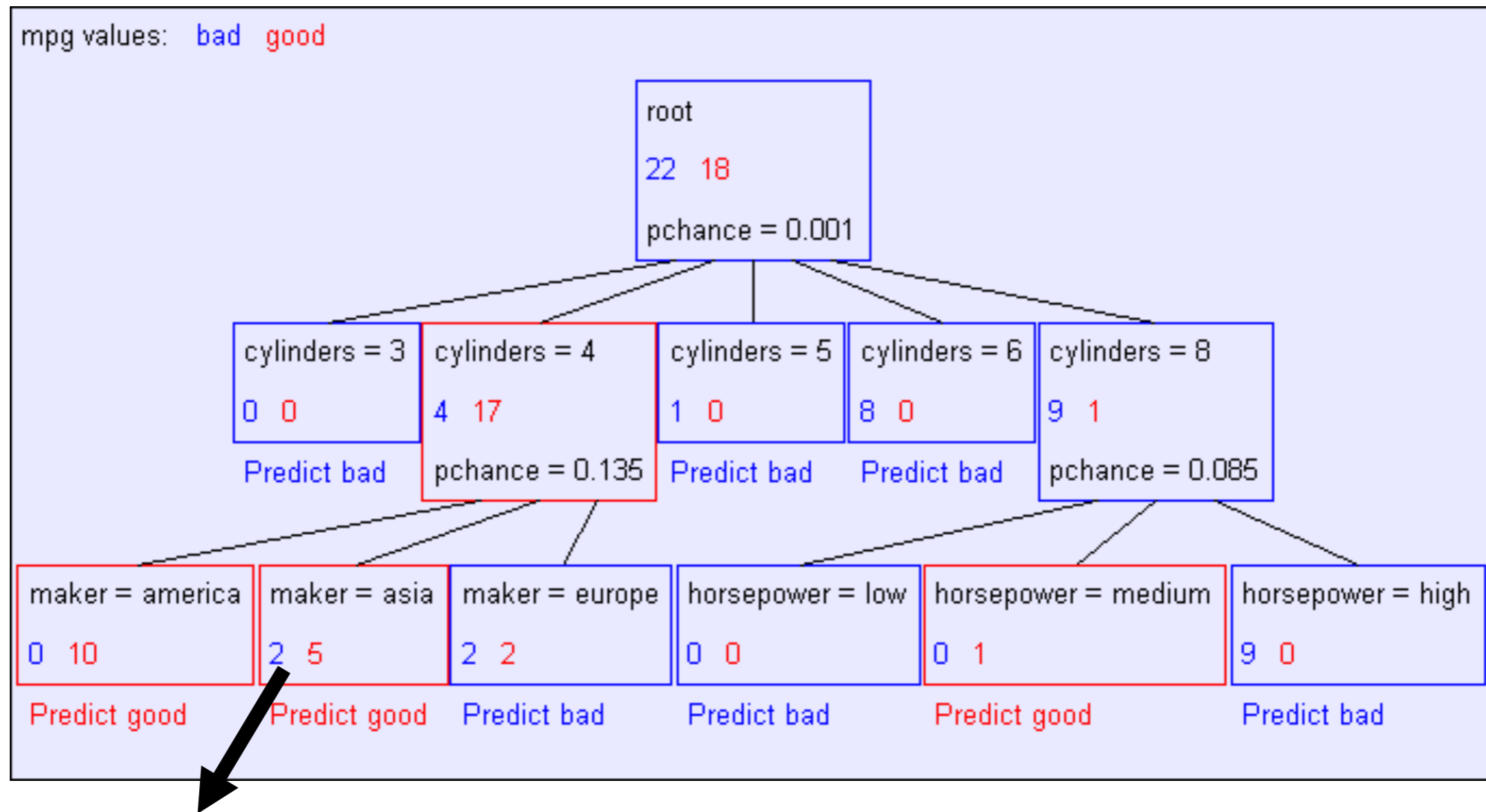
Records in  
which cylinders  
= 6

Build tree from  
These records..



Records in  
which cylinders  
= 8

# Second level of tree

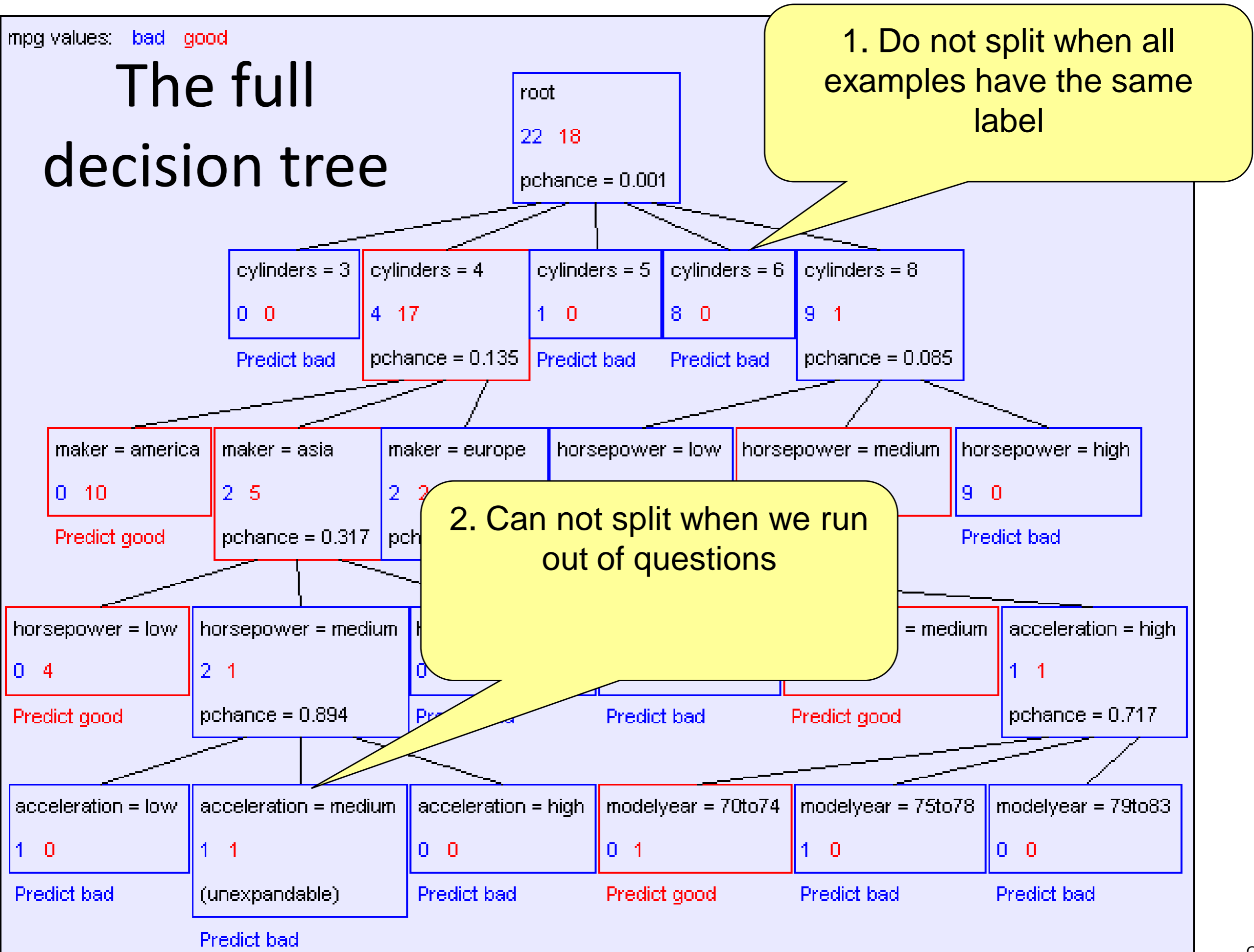


Recursively build a tree from the seven records in which there are four cylinders and the maker was based in Asia

(Similar recursion in the other cases)

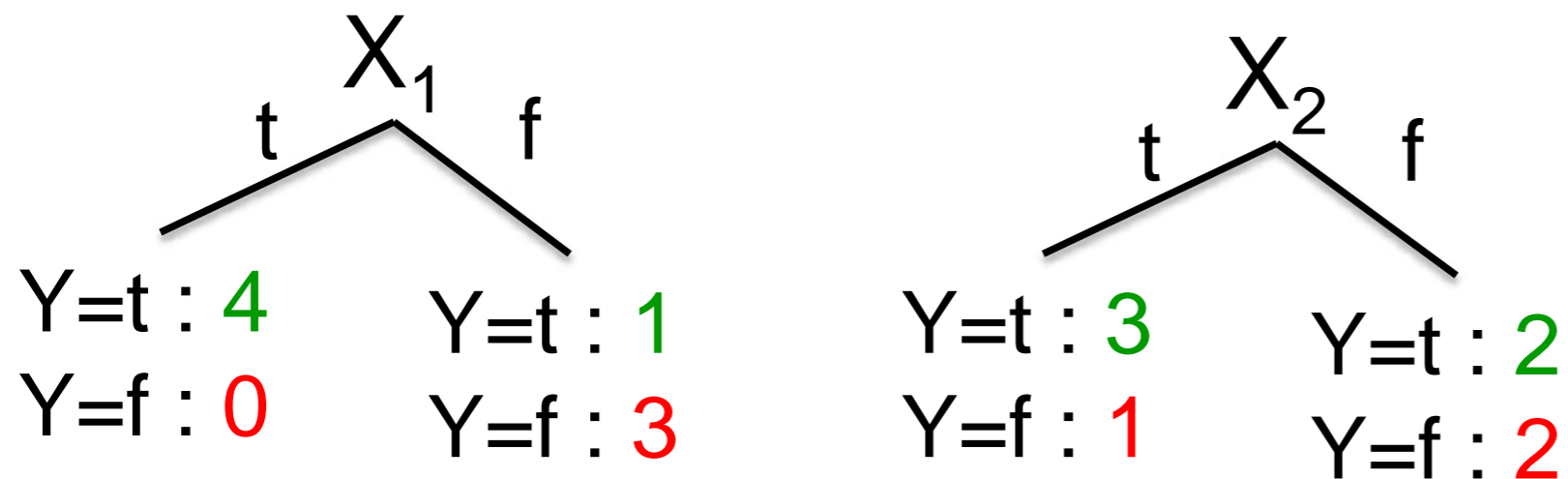
mpg values: bad good

# The full decision tree



# Splitting: Choosing a good attribute

- Would we prefer to split on  $X_1$  or  $X_2$ ?



$X_1$	$X_2$	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

**Idea:** use counts at leaves to define probability distributions, so we can measure uncertainty!

# Measuring uncertainty

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

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

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

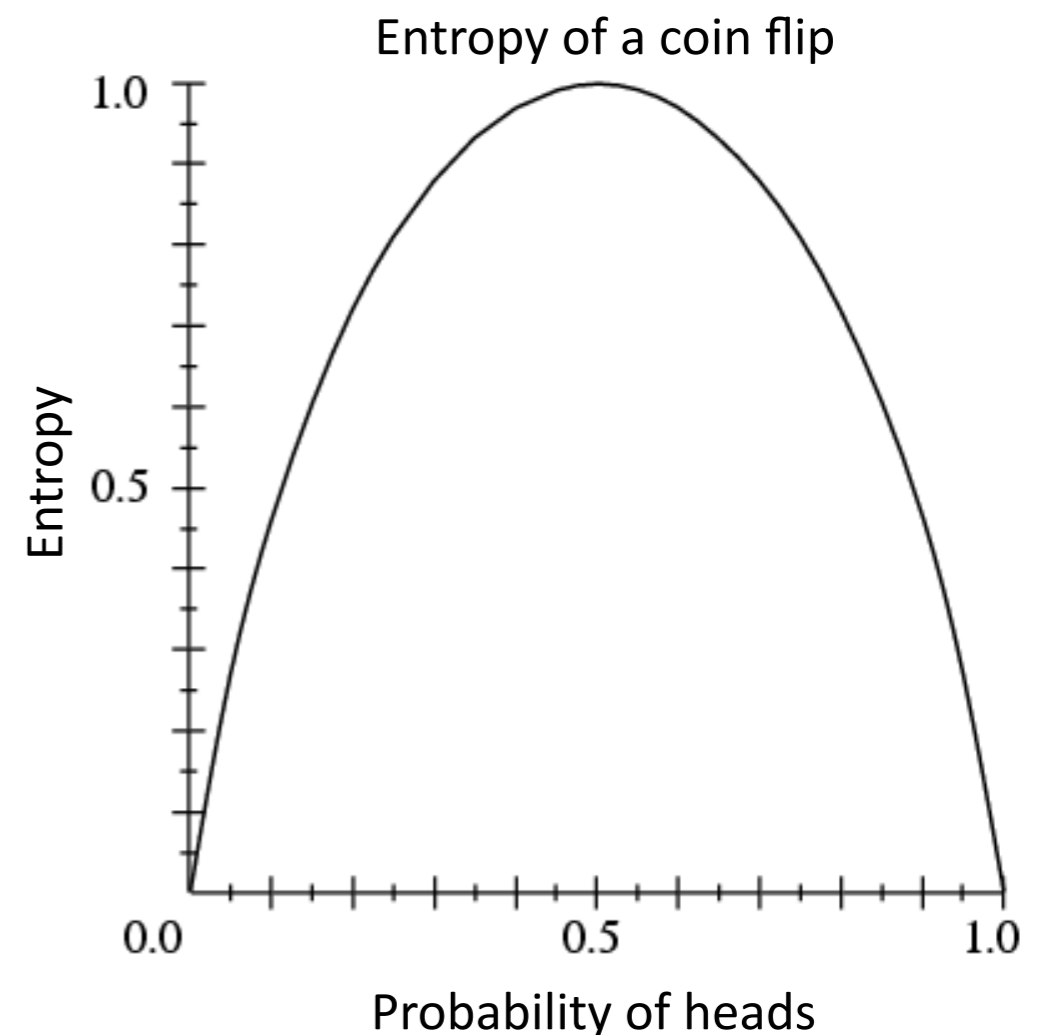


# Entropy

- Entropy  $H(Y)$  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)



# High, Low Entropy

- “High Entropy”
  - Y is from a uniform like distribution
  - Flat histogram
  - Values sampled from it are less predictable
- “Low Entropy”
  - Y is from a varied (peaks and valleys) distribution
  - Histogram has many lows and highs
  - Values sampled from it are more predictable

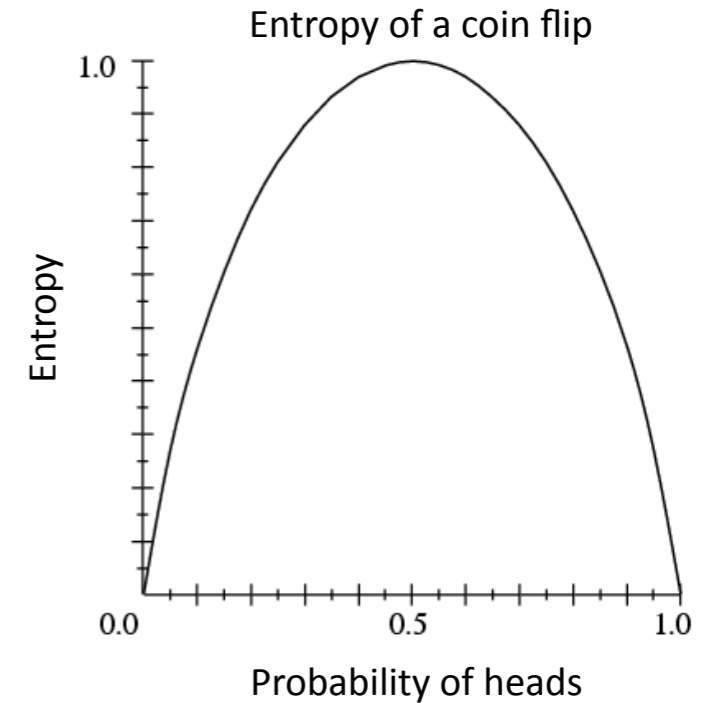
# Entropy Example

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

$$P(Y=t) = 5/6$$

$$P(Y=f) = 1/6$$

$$\begin{aligned} H(Y) &= - 5/6 \log_2 5/6 - 1/6 \log_2 1/6 \\ &= 0.65 \end{aligned}$$



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

# Conditional Entropy

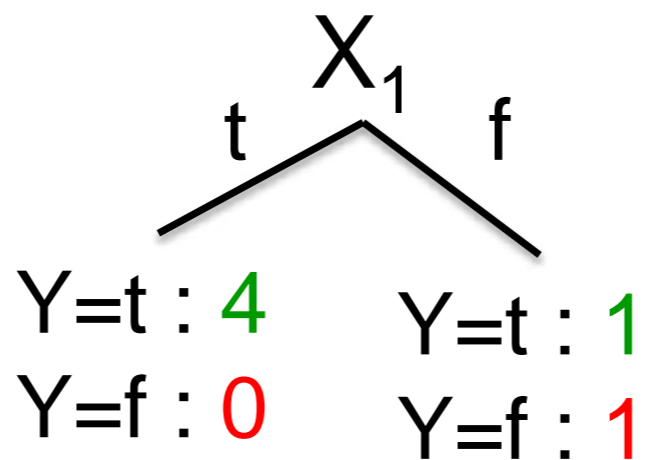
Conditional Entropy  $H(Y|X)$  of a random variable  $Y$  conditioned on a random variable  $X$

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

Example:

$$P(X_1=t) = 4/6$$

$$P(X_1=f) = 2/6$$



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

$$\begin{aligned}
 H(Y|X_1) &= - 4/6 (1 \log_2 1 + 0 \log_2 0) \\
 &\quad - 2/6 (1/2 \log_2 1/2 + 1/2 \log_2 1/2) \\
 &= 2/6
 \end{aligned}$$

# Information gain

- Decrease in entropy (uncertainty) after splitting

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

In our running example:

$$\begin{aligned} IG(X_1) &= H(Y) - H(Y|X_1) \\ &= 0.65 - 0.33 \end{aligned}$$

$IG(X_1) > 0 \rightarrow$  we prefer the split!

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

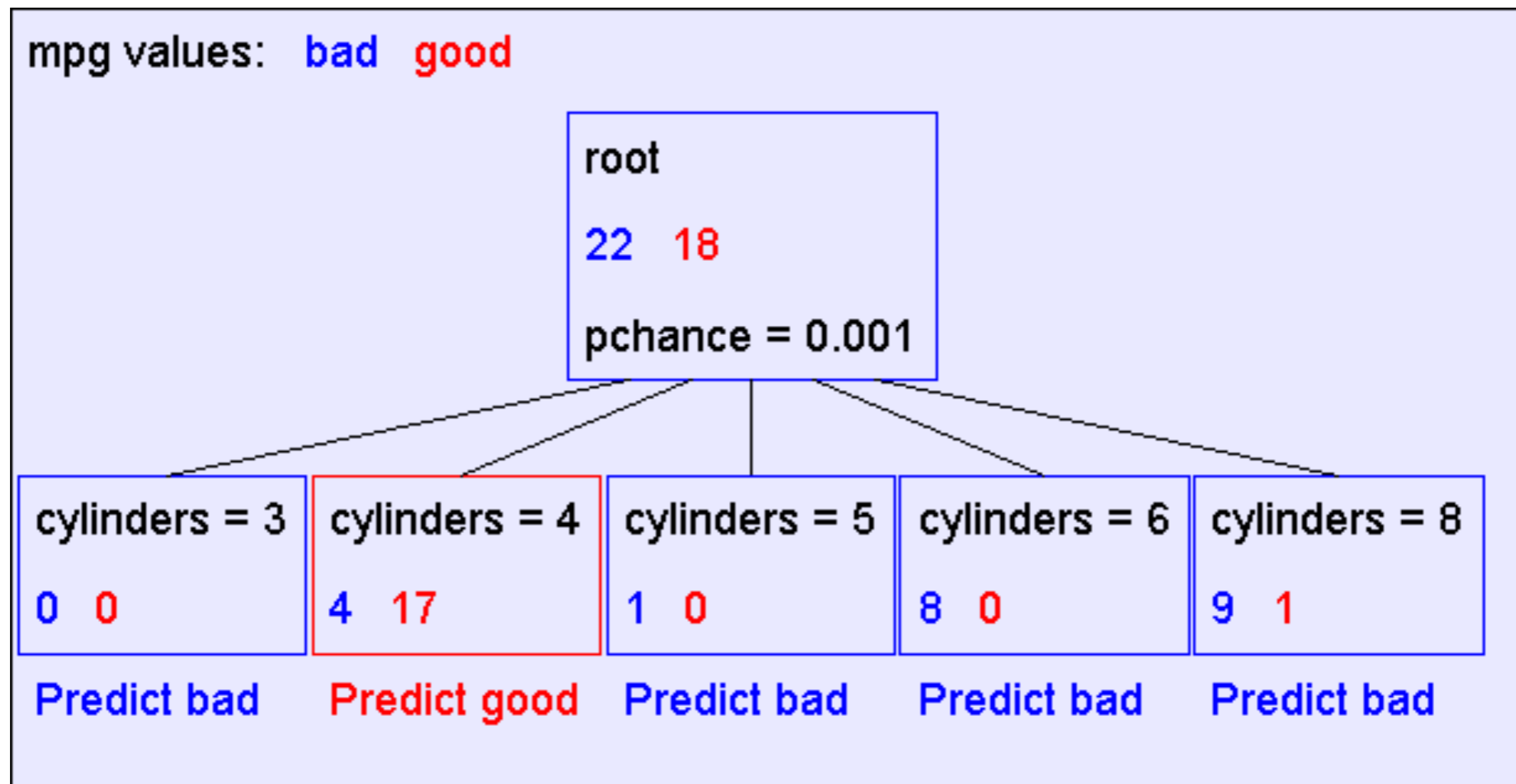
# Learning decision trees

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

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

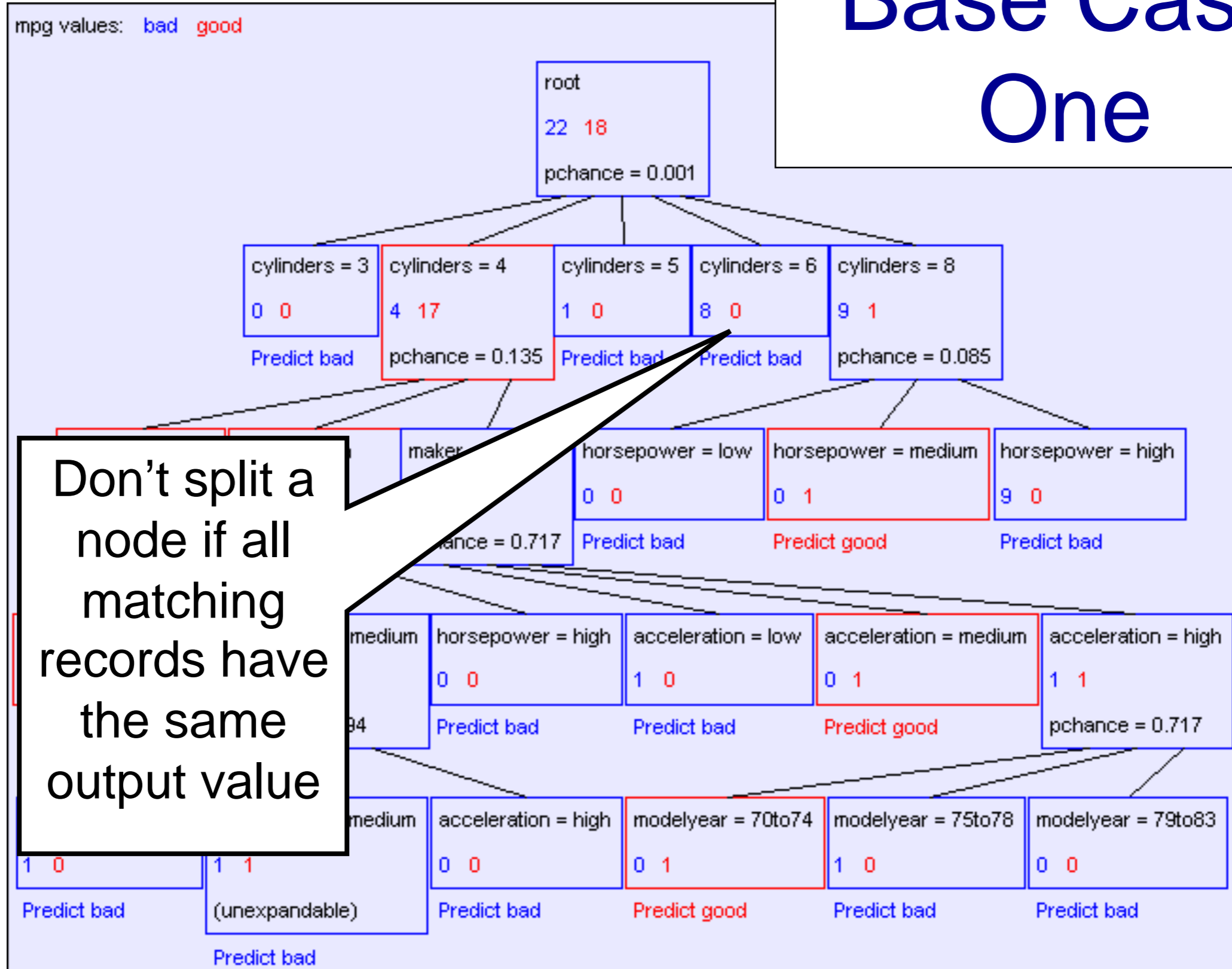
- Recurse

# When to stop?



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

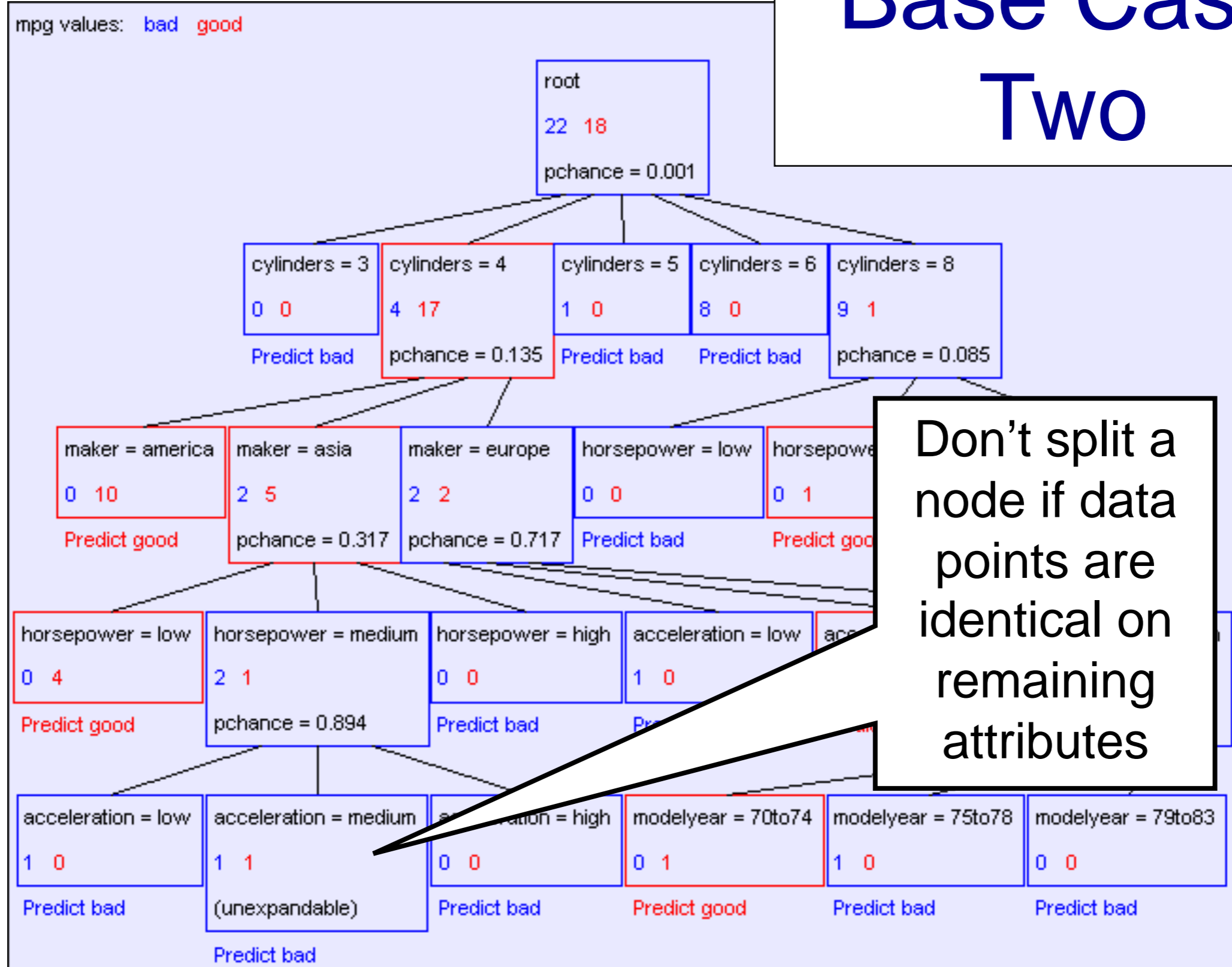
# Base Case One



Don't split a node if all matching records have the same output value

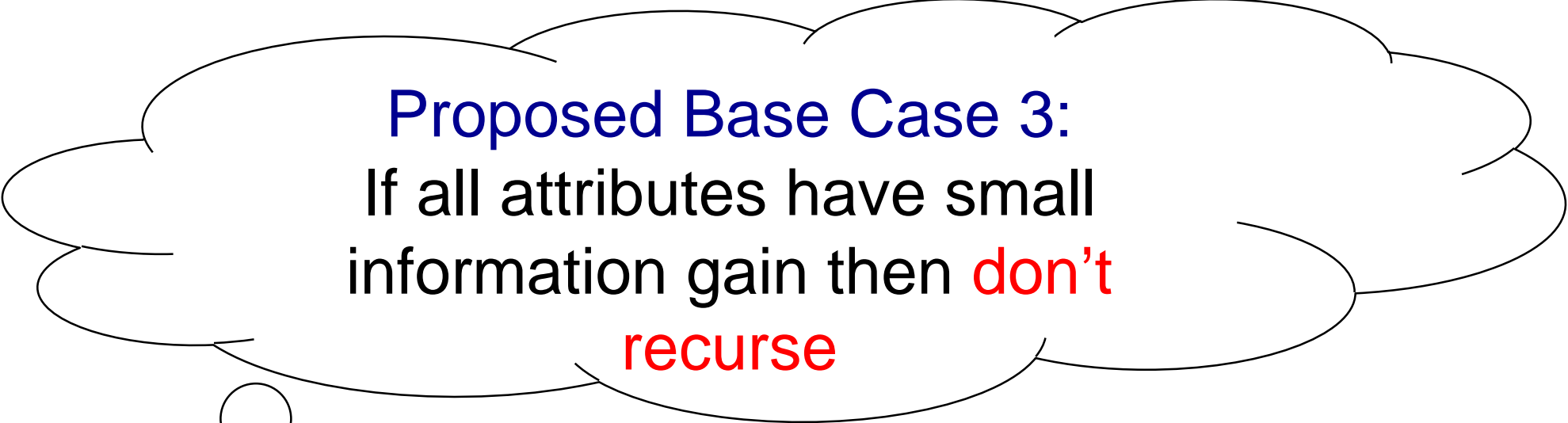


# Base Case Two



# 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 small  
information gain then **don't  
recurse**

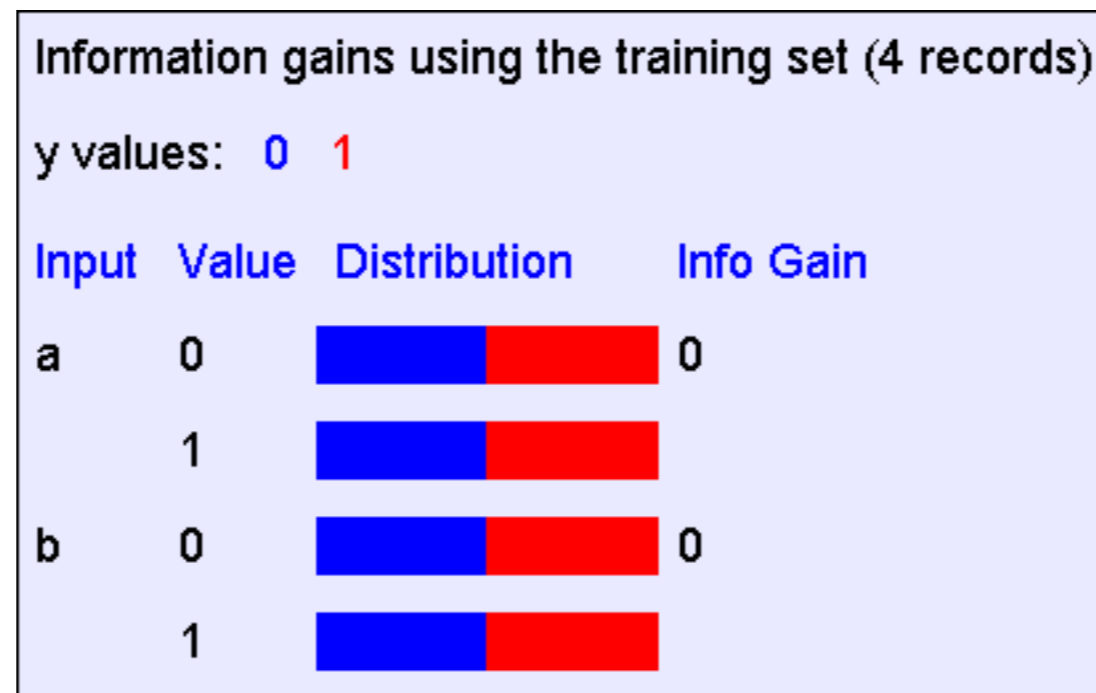
- *This is not a good idea*

# The problem with proposed case 3

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

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

The information gains:



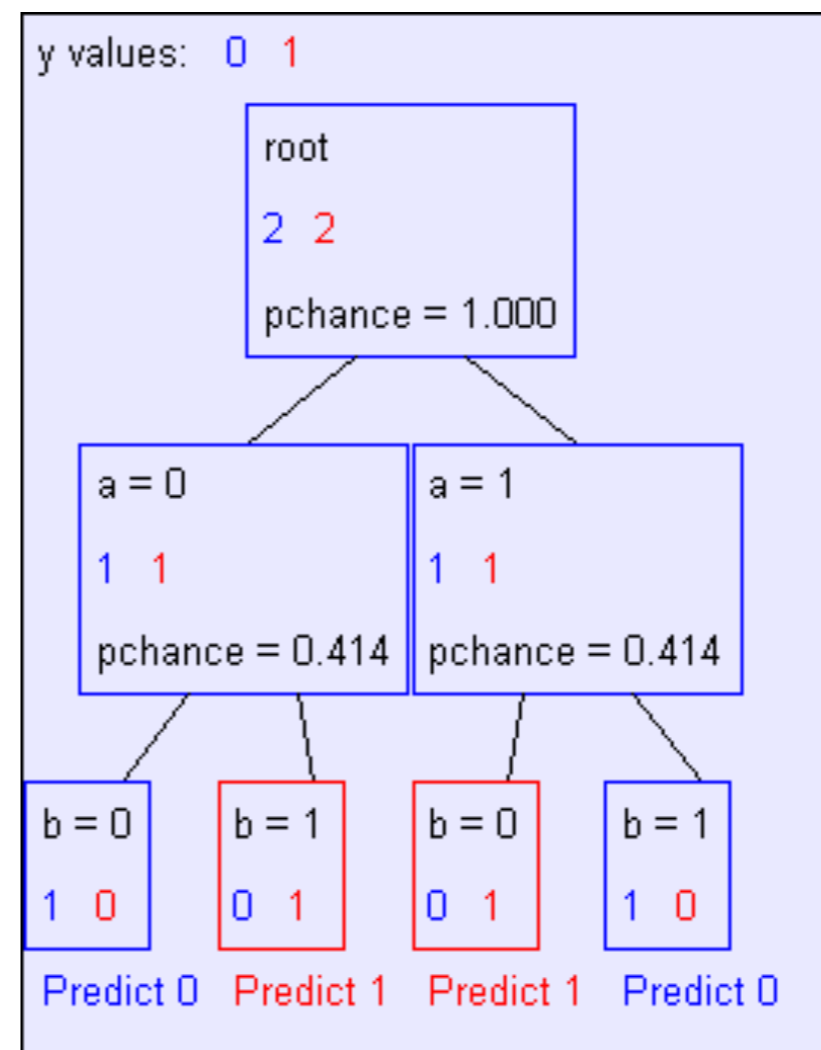
# If we omit proposed case 3:

$y = a \text{ XOR } b$

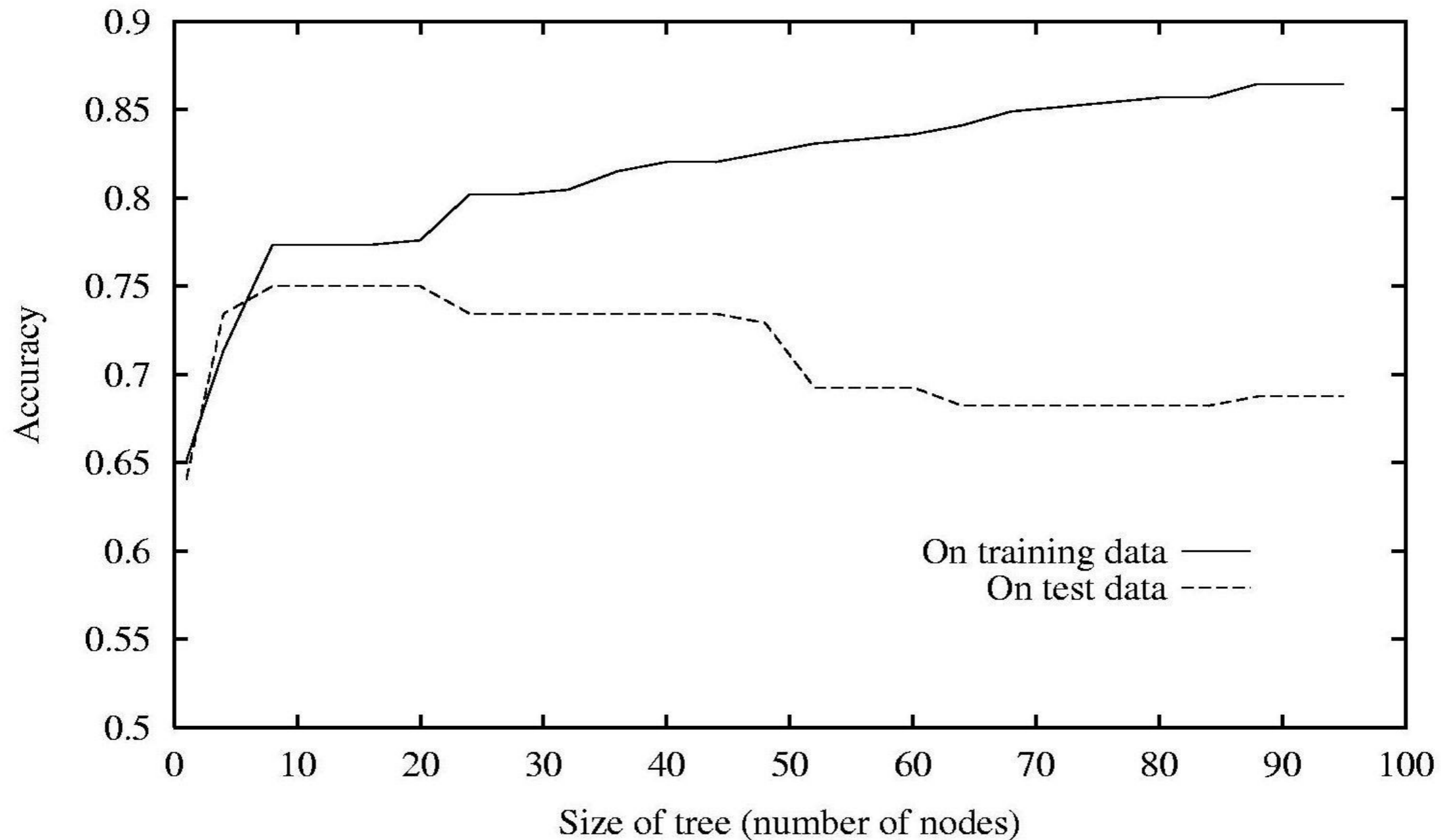
a	b	y
0	0	0
0	1	1
1	0	1
1	1	0

Instead, perform **pruning** after building a tree

The resulting decision tree:



# Decision trees will overfit



# Decision trees will overfit

- 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
  - Fixed number of leaves
- Random forests

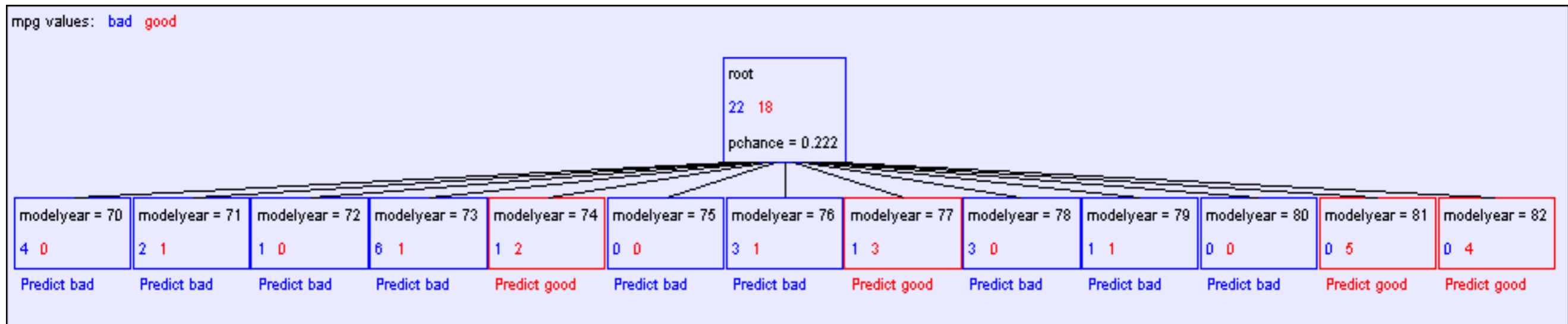
# Real-valued inputs

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

Infinite  
number of  
possible split  
values!!!

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

# “One branch for each numeric value” idea:

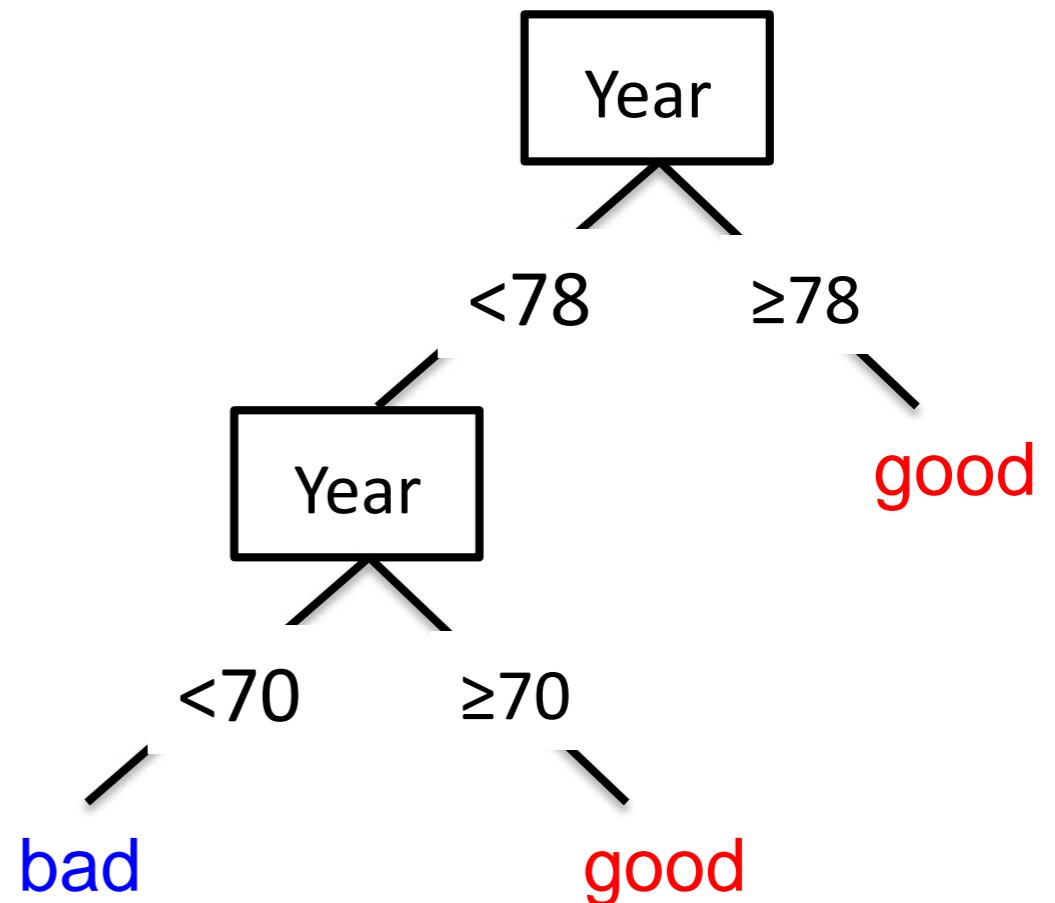


**Hopeless:** hypothesis with such a high branching factor will shatter any dataset and overfit



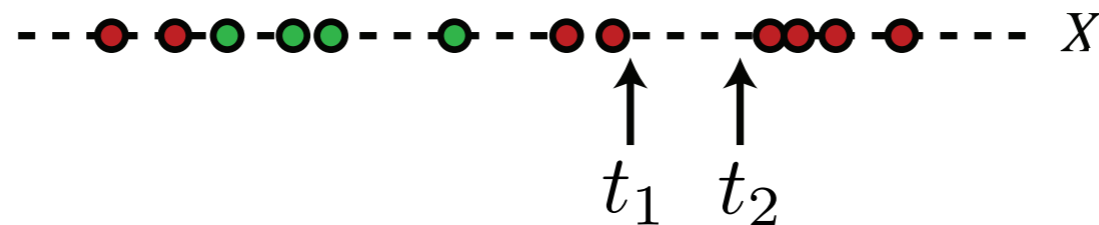
# Threshold splits

- Binary tree: split on attribute  $X$  at value  $t$ 
  - One branch:  $X < t$
  - Other branch:  $X \geq t$
- Requires small change
  - Allow repeated splits on same variable **along a path**

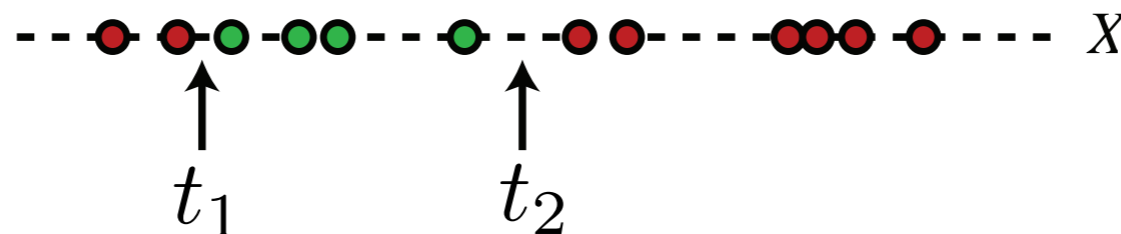


# The set of possible thresholds

- 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 a finite number of  $t$ 's are important:



- Sort data according to  $X$  into  $\{x_1, \dots, x_m\}$
- Consider split points of the form  $x_i + (x_{i+1} - x_i)/2$
- Moreover, only splits between examples from different classes matter!


















# Picking the best threshold

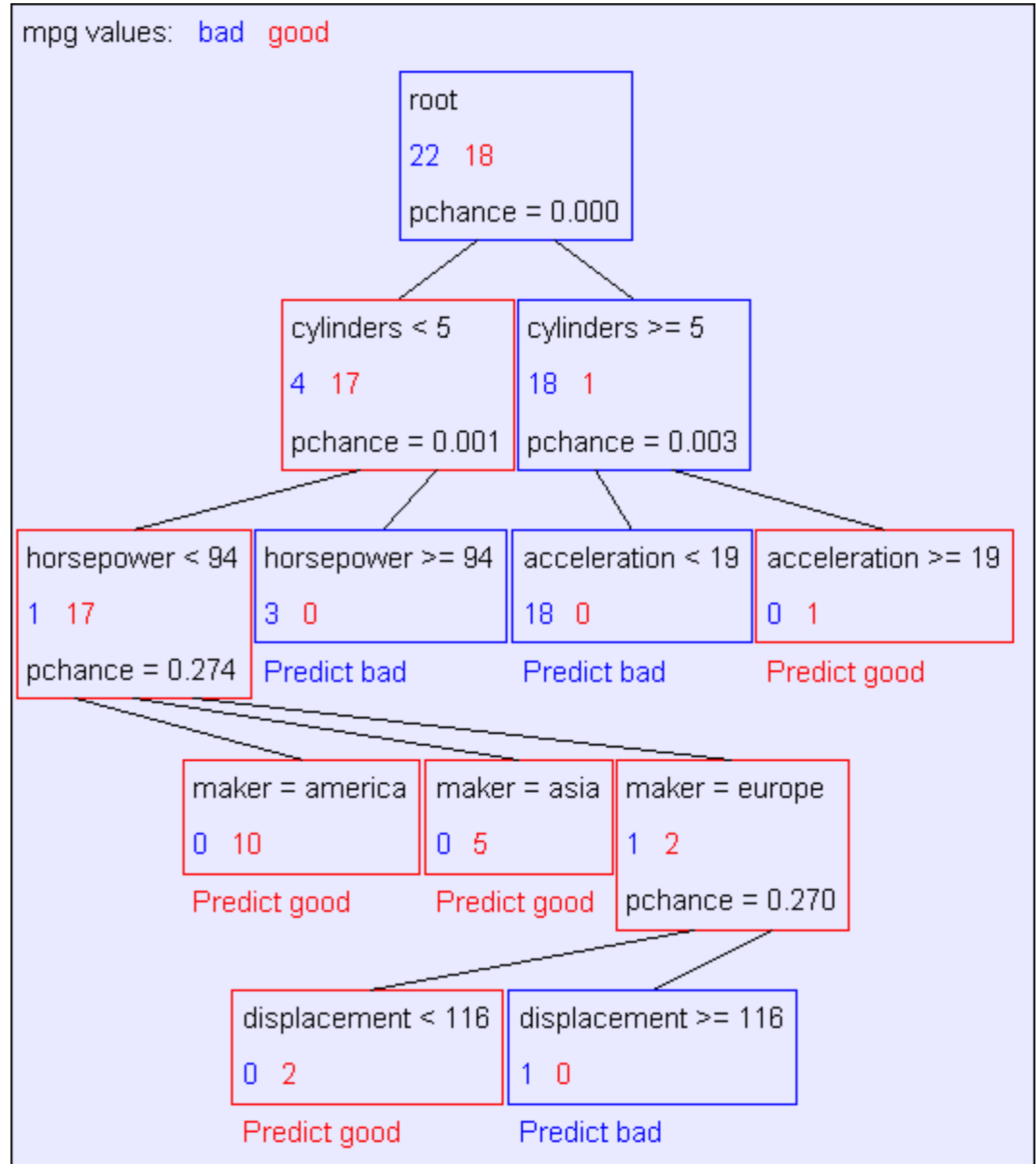
- Suppose  $X$  is real valued with threshold  $t$
- Want  $IG(Y | X:t)$ , the information gain for  $Y$  when testing if  $X$  is greater than or less than  $t$
- Define:
  - $H(Y | X:t) = p(X < t)H(Y | X < t) + p(X \geq t)H(Y | X \geq t)$
  - $IG(Y | X:t) = H(Y) - H(Y | X:t)$
  - $IG^*(Y | X) = \max_t IG(Y | X:t)$
- Use:  $IG^*(Y | X)$  for continuous variables

# Example with MPG

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 tree for our continuous dataset



Demo time...

# What you need to know about decision trees

- Decision trees are one of the most popular ML tools
  - Easy to understand, implement, and 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!!!
  - Must use tricks to find “simple trees”, e.g.,
    - Fixed depth/Early stopping
    - Pruning
  - Or, use ensembles of different trees (random forests)

# Decision Trees vs SVM

Characteristic	SVM	Trees
Natural handling of data of “mixed” type	▼	▲
Handling of missing values	▼	▲
Robustness to outliers in input space	▼	▲
Insensitive to monotone transformations of inputs	▼	▲
Computational scalability (large $N$ )	▼	▲
Ability to deal with irrelevant inputs	▼	▲
Ability to extract linear combinations of features	▲	▼
Interpretability	▼	◆
Predictive power	▲	▼