# CSE 446: Week 1
# Decision Trees

# Administrative

- Everyone should have been enrolled into Gradescope, please contact Isaac Tian ([iytian@cs.washington.edu](mailto:iytian@cs.washington.edu)) if you did not receive anything about this

- Please check Piazza for news and announcements, now that everyone is (hopefully) signed up!

# Clarifications from Last Time

- "objective" is a synonym for "cost function"
  - later on, you'll hear me refer to it as a "loss function" – that's also the same thing

# Review

- Four parts of a machine learning problem [decision trees]
  - What is the data?
  - What is the hypothesis space?
    - It's big
  - What is the objective?
    - We're about to change that
  - What is the algorithm?

# Algorithm

- Four parts of a machine learning problem [decision trees]
  - What is the data?
  - What is the hypothesis space?
    - It's big
  - What is the objective?
    - We're about to change that
  - What is the algorithm?

# Decision Trees

[tutorial on the board]

[see lecture notes for details]

I.   Recap

II.  Splitting criterion: information gain

III. Entropy vs error rate and other costs

# Supplementary: 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 |
|---|---|---|---|

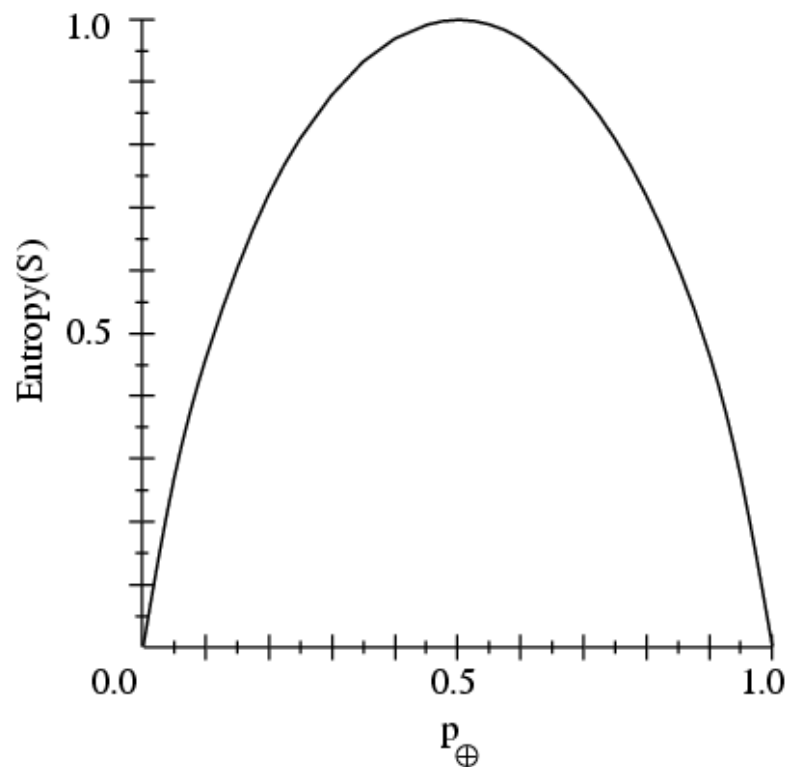| P(Y=A) = 1/4 | P(Y=B) = 1/4 | P(Y=C) = 1/4 | P(Y=D) = 1/4 |
|---|---|---|---|

# Supplementary: 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)
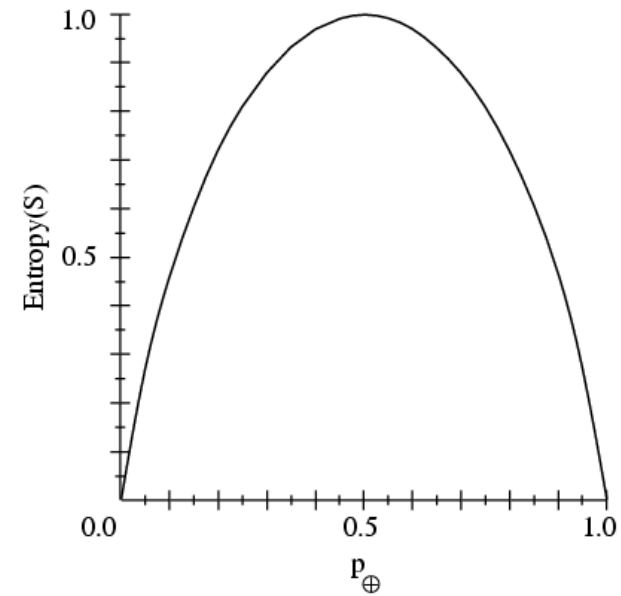
# Supplementary: Entropy Example

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



P(Y=t) = 5/6

P(Y=f) = 1/6

H(Y) = - 5/6 log$_2$ 5/6 - 1/6 log$_2$ 1/6

= 0.65

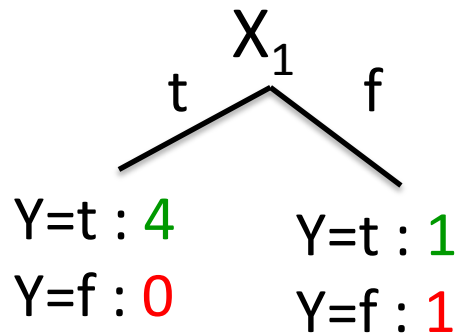| $X_1$ | $X_2$ | Y |
|-------|-------|---|
| T | T | T |
| T | F | T |
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |

# Supplementary: Conditional Entropy

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

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

Example:

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

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

$X_1$

t      f

Y=t : 4    Y=t : 1

Y=f : 0    Y=f : 1

$H(Y|X_1) = - 4/6 \ (1 \log_2 1 + 0 \log_2 0)$

$\qquad\qquad\quad - 2/6 \ (1/2 \log_2 1/2 + 1/2 \log_2 1/2)$

$\qquad = 2/6$

| $X_1$ | $X_2$ | Y |
|-------|-------|---|
| T | T | T |
| T | F | T |
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |

# Supplementary: Information gain

Decrease in entropy (uncertainty) after splitting

$$IG(X) = H(Y) - H(Y \mid X)$$

- IG(X) is non-negative (>=0)
- Prove by showing H(Y|X) <= H(X), with Jensen's inequality

In our running example:

IG($X_1$) = H(Y) − H(Y|$X_1$)

$\quad\quad$ =  0.65 − 0.33

IG($X_1$) > 0 → 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 |

# A learning problem: predict fuel efficiency

| 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 | europe |
| 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 | europe |
| bad | 5 | medium | medium | medium | medium | 75to78 | europe |

- 40 Records

- Discrete data (for now)

- Predict MPG

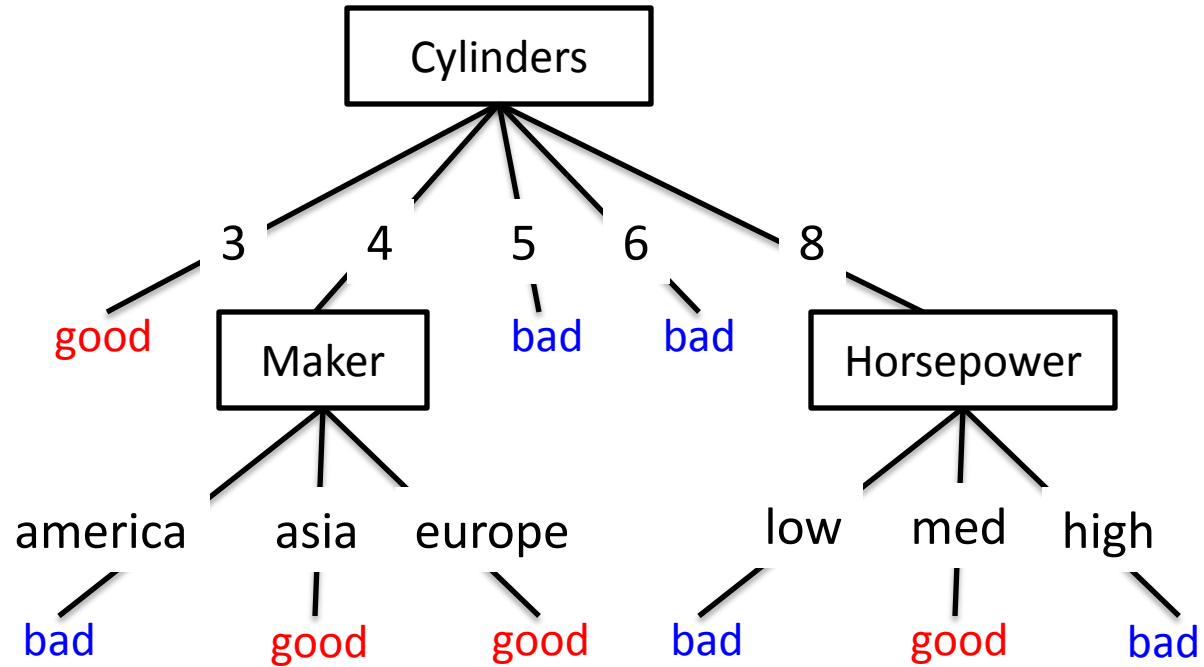- Need to find: $f : X \rightarrow Y$

$Y$          $X$

From the UCI repository (thanks to Ross Quinlan)

# 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 $y$

# Learning decision trees

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

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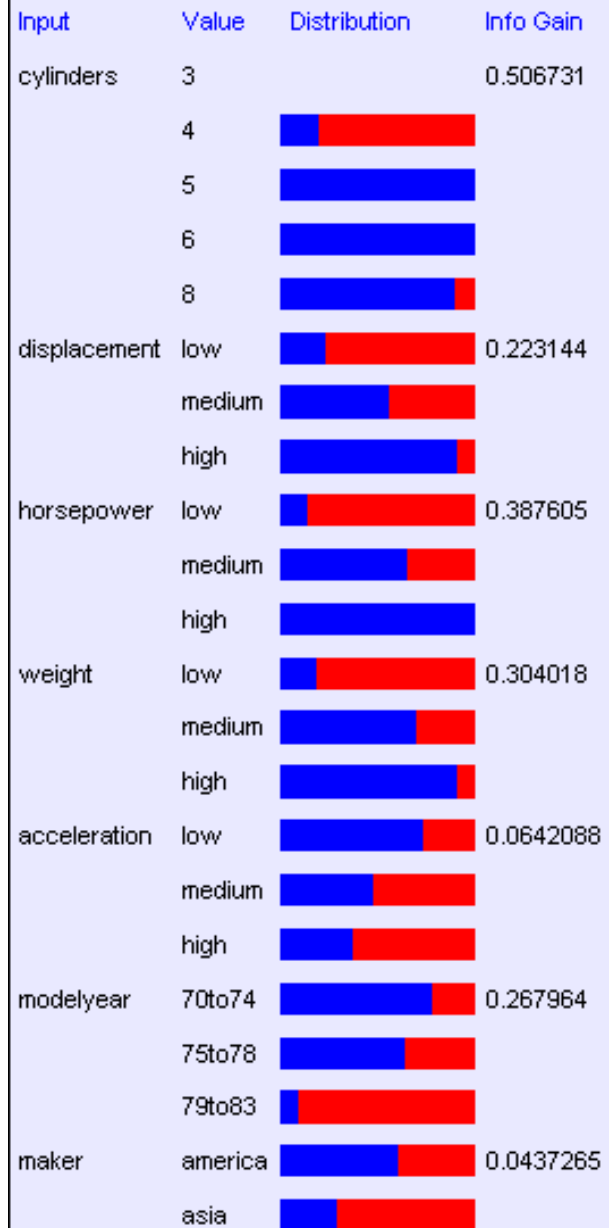
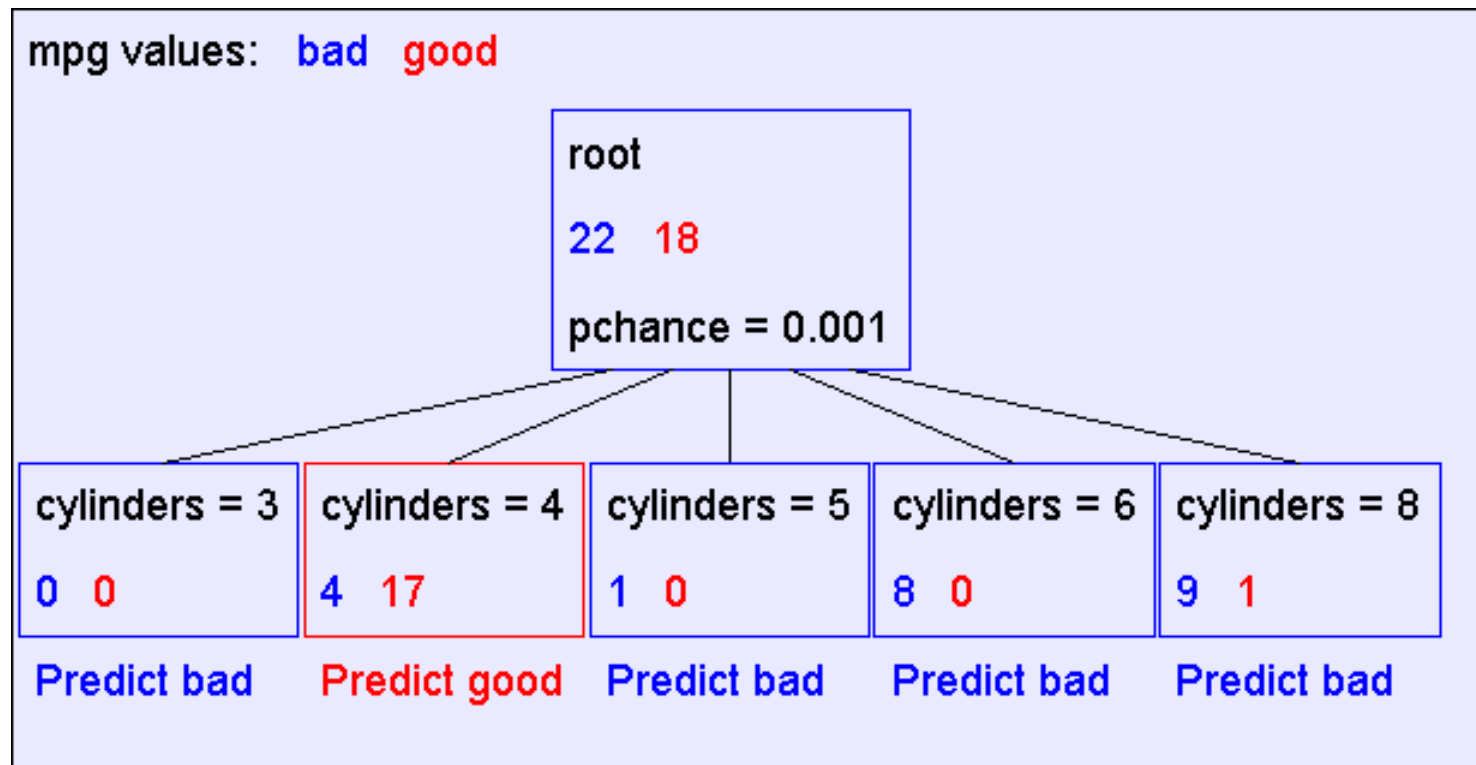- Recurse

Suppose we want to predict MPG

Look at all the information gains...



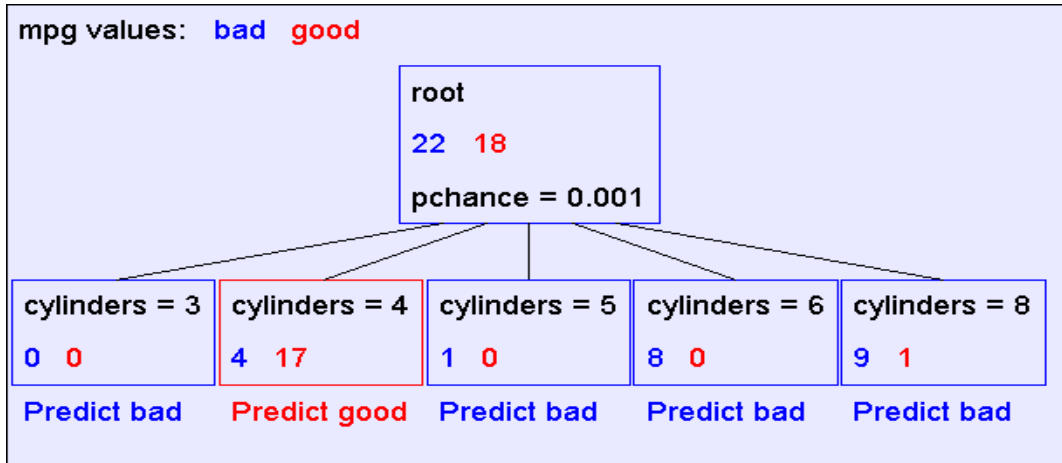Information gains using the training set (40 records)
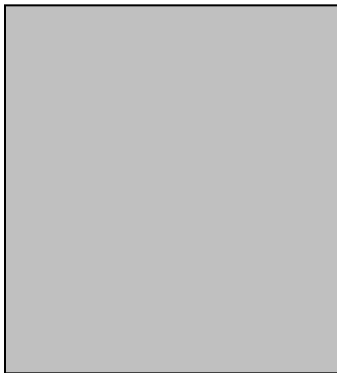
mpg values:  bad  good

| 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 | | |

# A Decision Stump

mpg values: bad good

root

22 18

pchance = 0.001

| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0 0 | 4 17 | 1 0 | 8 0 | 9 1 |
| Predict bad | Predict good | Predict bad | Predict bad | Predict bad |

# Recursive Step

mpg values:  bad   good

root
22   18
pchance = 0.001

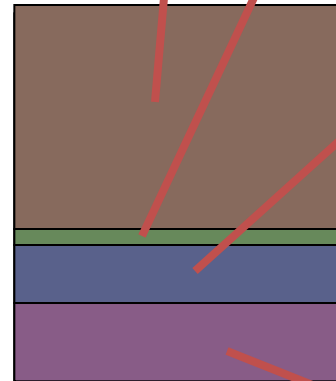| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0   0 | 4   17 | 1   0 | 8   0 | 9   1 |
| Predict bad | Predict good | Predict bad | Predict bad | Predict bad |

Take the Original Dataset..

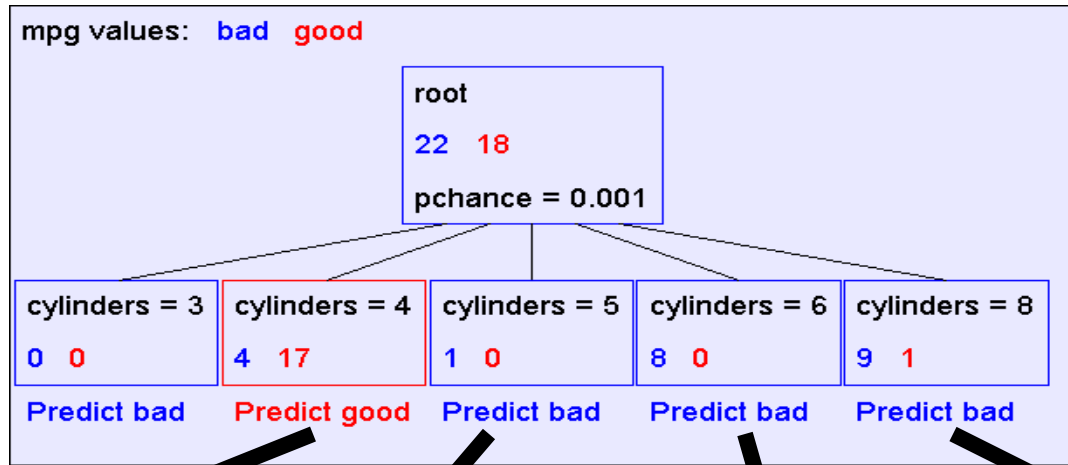And partition it according to the value of the attribute we split on

Records in which cylinders = 4
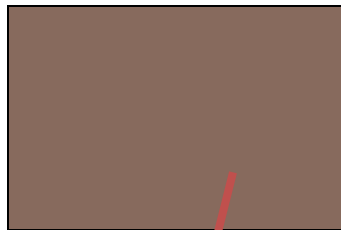
Records in which cylinders = 5

Records in which cylinders = 6

Records in which cylinders = 8

# Recursive Step

# Second level of tree

mpg values:  bad  good

root
22 18
pchance = 0.001

cylinders = 3
0  0
Predict bad

cylinders = 4
4  17
pchance = 0.135

cylinders = 5
1  0
Predict bad

cylinders = 6
8  0
Predict bad

cylinders = 8
9  1
pchance = 0.085

maker = america
0  10
Predict good

maker = asia
2  5
Predict good

maker = europe
2  2
Predict bad

horsepower = low
0  0
Predict bad

horsepower = medium
0  1
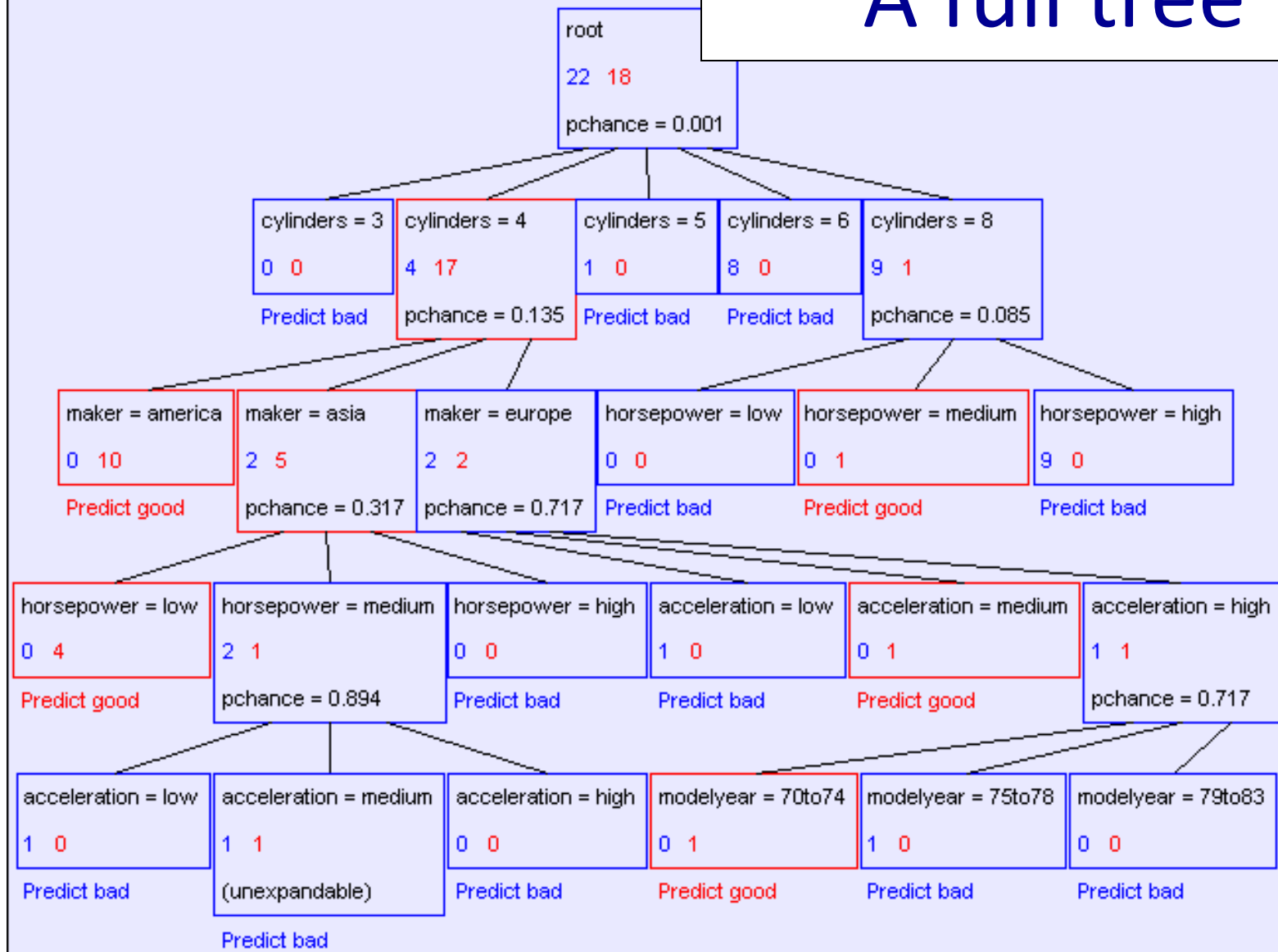Predict good

horsepower = high
9  0
Predict bad

Recursively build a tree from the seven records in which there are four cylinders and the maker was based in Asia
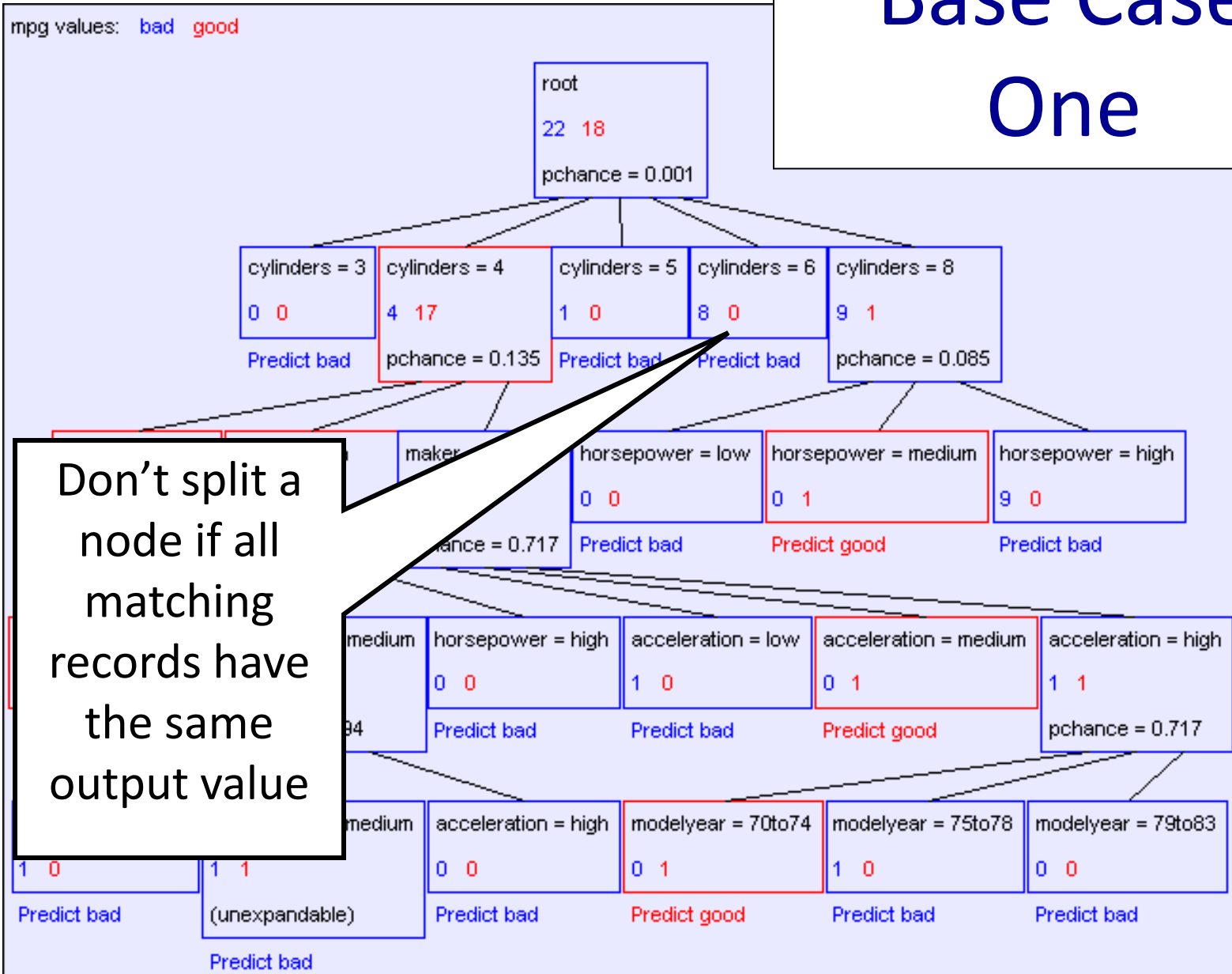
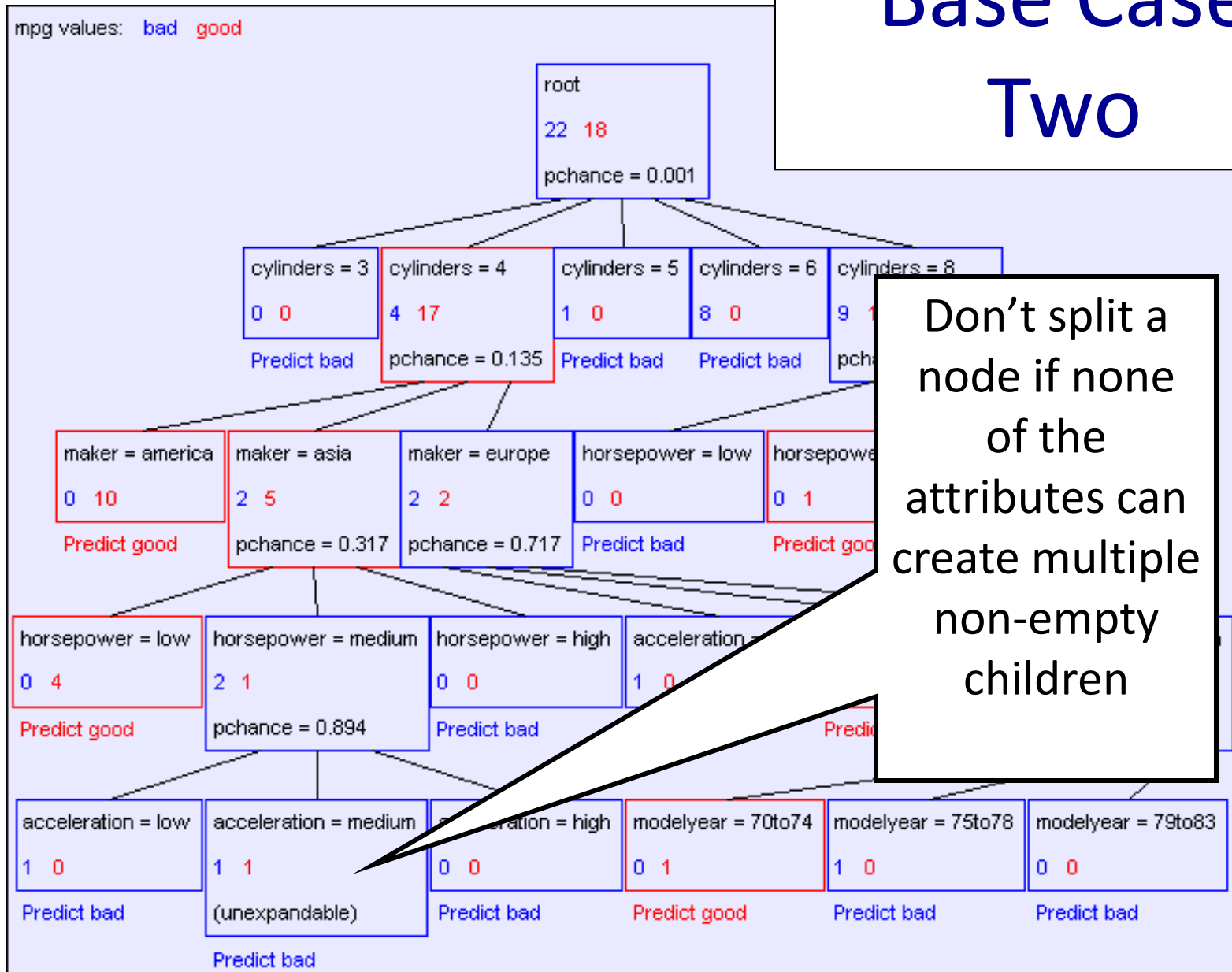(Similar recursion in the other cases)

A full tree

# What to stop?

Base Case One

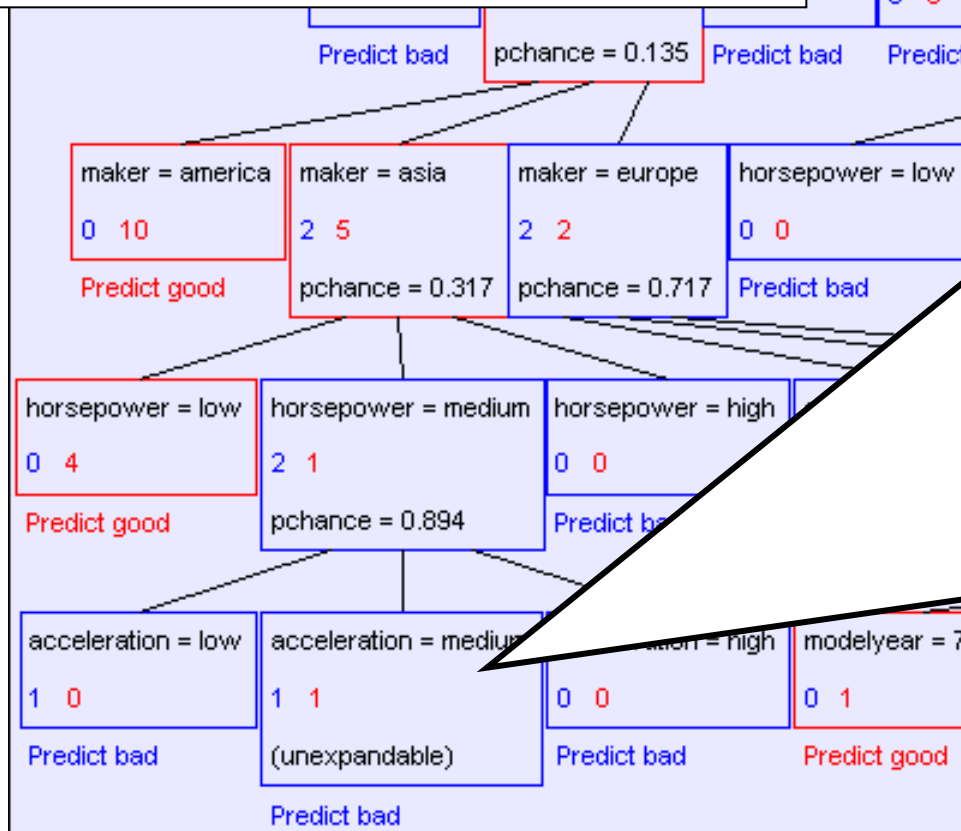Base Case Two

mpg values:  bad   good

Don't split a node if none of the attributes can create multiple non-empty children

# Base Case Two:
# No attributes can distinguish



Predict bad    pchance = 0.001

s = 5    cylinde

8   0

Predict bad    pchance = 0.135    Predict bad    Predict

| maker = america | maker = asia | maker = europe | horsepower = low |
|---|---|---|---|
| 0   10 | 2   5 | 2   2 | 0   0 |
| Predict good | pchance = 0.317 | pchance = 0.717 | Predict bad |

| horsepower = low | horsepower = medium | horsepower = high | |
|---|---|---|---|
| 0   4 | 2   1 | 0   0 | |
| Predict good | pchance = 0.894 | Predict b | |

| acceleration = low | acceleration = mediu | ion = high | modelyear = 7 |
|---|---|---|---|
| 1   0 | 1   1 | 0   0 | 0   1 |
| Predict bad | (unexpandable) | Predict bad | Predict good |

Predict bad

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: 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

- *Is this a good idea?*

# The problem with Base Case 3

y = a XOR b

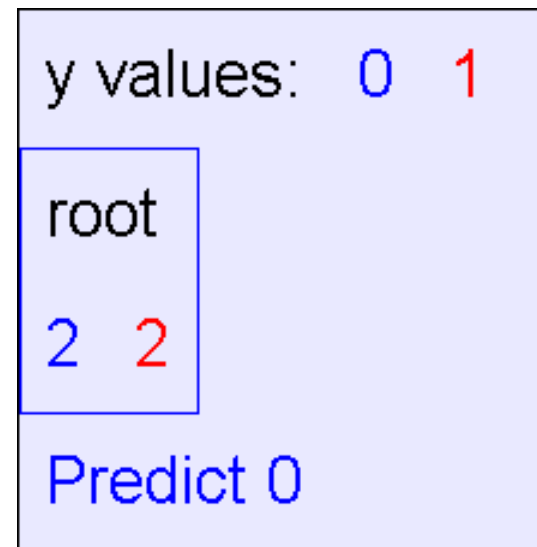| a | b | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

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:

y values: 0 1

root

2 2

Predict 0

# If we omit Base Case 3:

The resulting decision tree:

y = a XOR b

| a | b | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

# 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 = low
horsepower = medium
horsepower = high
acceleration = low
acceleration = medium
acceleration = high
epower = high

Predict bad
(unexpandable)
Predict bad
Predict good
Predict bad
Predict bad

Predict bad

= 0.717
= 79to83
ict bad

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

…why?

# 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
  - Or something smarter…

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | ? |
| 0 | 0 | 0 | 1 | ? |
| 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | ? |
| 1 | 0 | 0 | 0 | ? |
| 1 | 0 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | ? |
| 1 | 0 | 1 | 1 | ? |
| 1 | 1 | 0 | 0 | 0 |
| 1 | 1 | 0 | 1 | ? |
| 1 | 1 | 1 | 0 | ? |
| 1 | 1 | 1 | 1 | ? |

# Decision trees will overfit!!!