# Decision Trees
# Lecture 11

David Sontag

New York University

Slides adapted from Luke Zettlemoyer, Carlos Guestrin, and Andrew Moore

# 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 data points

- Goal: predict MPG

- Need to find:
  $f : X \rightarrow Y$
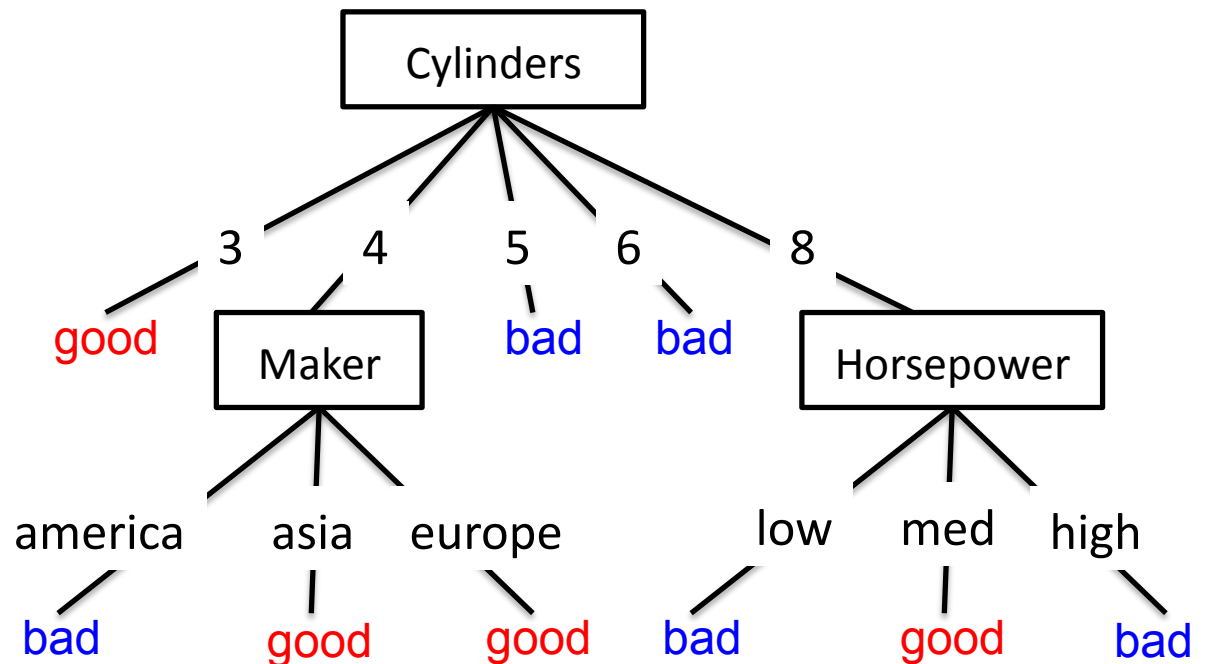
- Discrete data (for now)

$Y$         $X$

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 $y$
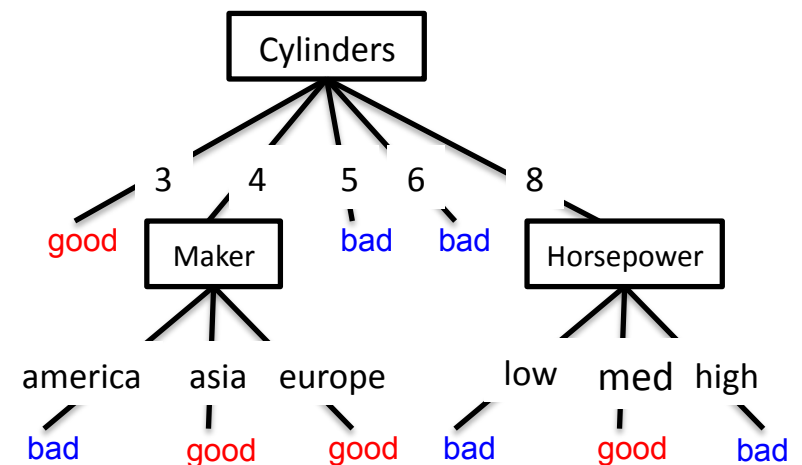
Cylinders

3    4    5    6    8

good    Maker    bad    bad    Horsepower

america    asia    europe    low    med    high

bad    good    good    bad    good    bad

Human interpretable!

# Hypothesis space

- How many possible hypotheses?

- What functions can be represented?

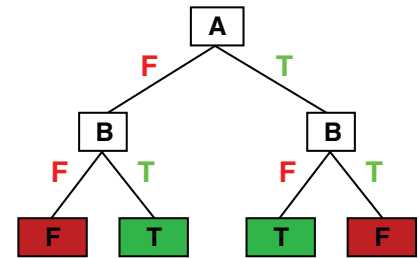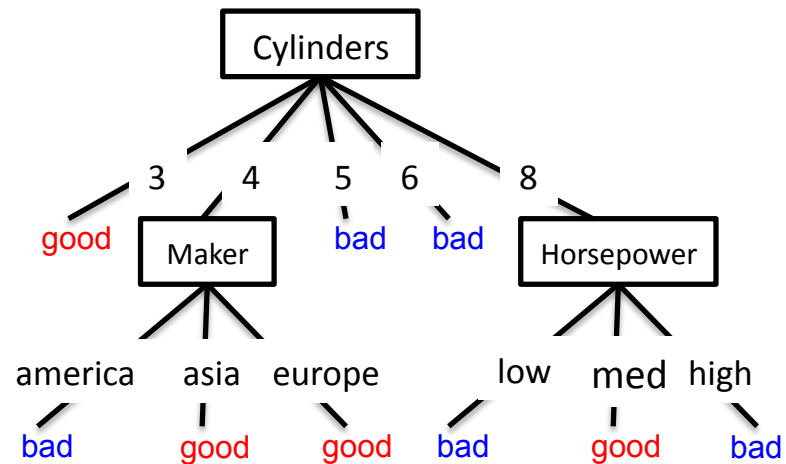| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|------|-----------|--------------|------------|--------|--------------|-----------|---------|
| | | | | | | | |
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium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

# 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)



cyl=3 ∨ (cyl=4 ∧ (maker=asia ∨ maker=europe)) ∨ …

# Hypothesis space

- ## How many possible hypotheses?

- ## What functions can be represented?

- ## How many will be consistent with a given dataset?

- ## How will we choose the best one?

  - ### Lets first look at how to split nodes, then consider how to find the best tree

| 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 |

# What is the Simplest Tree?

| 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 |

predict
mpg=bad

## Is this a good tree?

[22+, 18-]  ⬅  Means:
correct on 22 examples
incorrect on 18 examples

# A Decision Stump

# Recursive Step

mpg values: bad good

root

22 18

pchance = 0.001

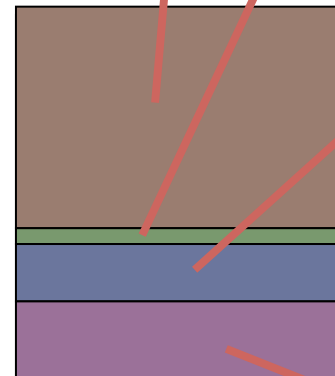| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0 0 | 4 17 | 1 0 | 8 0 | 9 1 |
| Predict bad | Predict good | Predict bad | Predict bad | Predict bad |

Take the Original Dataset..

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

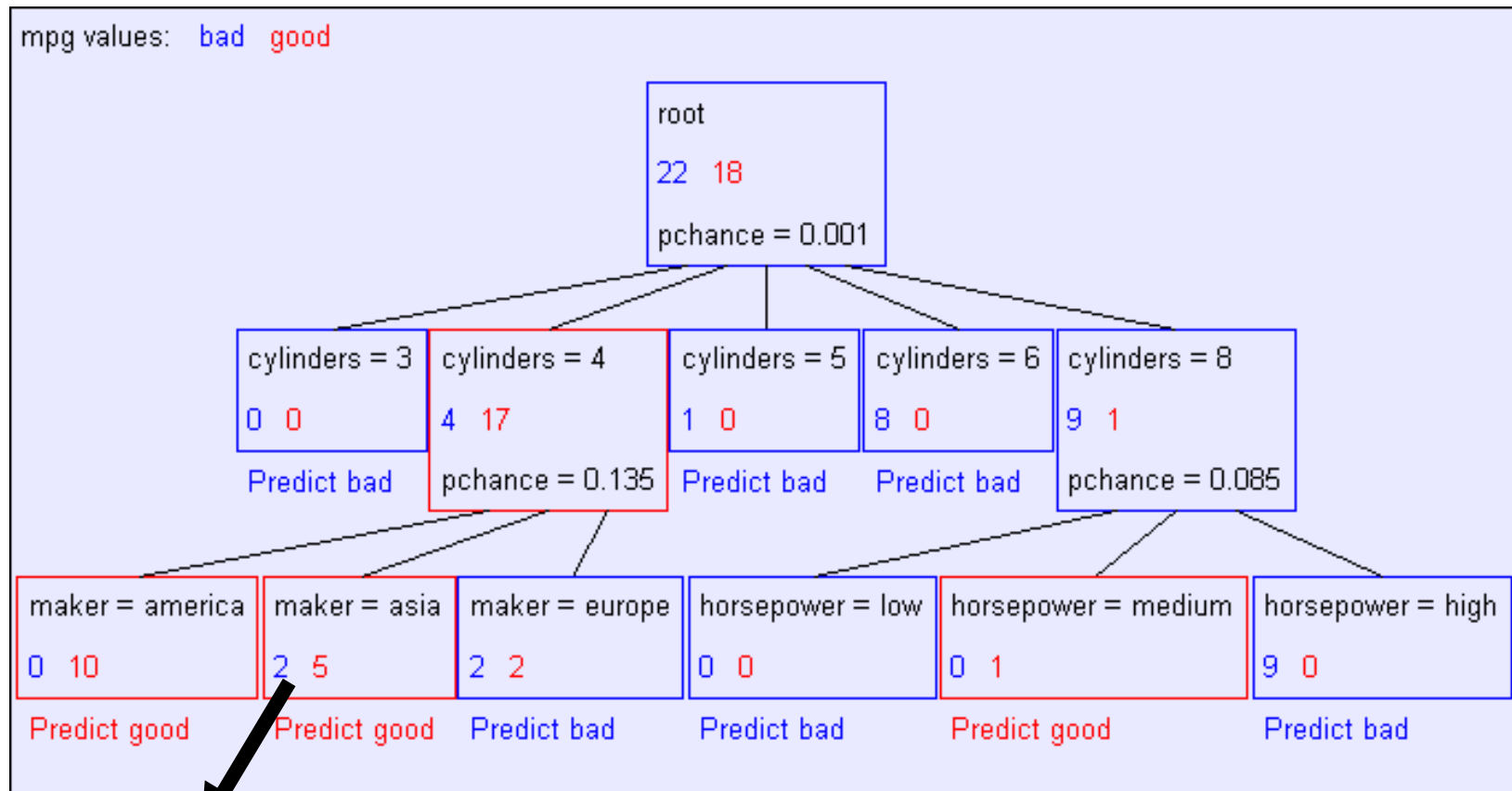Records in which cylinders = 4

Records in which cylinders = 5

Records in which cylinders = 6

Records in which cylinders = 8

# Recursive Step



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 |

Build tree from These records..

Build tree from These records..

Build tree from These records..

Build tree from These records..

Records in which cylinders = 4

Records in which cylinders = 5

Records in which cylinders = 6

Records in which cylinders = 8

# Second level of tree



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 | pchance = 0.135 | Predict bad | Predict bad | pchance = 0.085 |

| maker = america | maker = asia | maker = europe | horsepower = low | horsepower = medium | horsepower = high |
|---|---|---|---|---|---|
| 0  10 | 2  5 | 2  2 | 0  0 | 0  1 | 9  0 |
| Predict good | Predict good | Predict bad | Predict bad | Predict good | 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)

mpg values:  bad  good

A full tree

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

pchance = 0.317

maker = europe

2  2

pchance = 0.717

horsepower = low

0  0

Predict bad

horsepower = medium

0  1

Predict good

horsepower = high

9  0

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

# 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

# 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 |

Tree split on $X_1$:
- t: Y=t : 4, Y=f : 0
- f: Y=t : 1, Y=f : 3

Tree split on $X_2$:
- t: Y=t : 3, Y=f : 1
- f: Y=t : 2, Y=f : 2

**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)

Entropy of a coin flip



Entropy vs. Probability of heads

# 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

(Slide from Vibhav Gogate)

# Entropy Example

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

**Entropy of a coin flip**



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 |

# 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:

$X_1$

t     f

Y=t : 4    Y=t : 1
Y=f : 0    Y=f : 1

$P(X_1=t) = 4/6$
$P(X_1=f) = 2/6$

$H(Y|X_1) = -\ 4/6\ (1\ \log_2 1 + 0\ \log_2 0)$
$\qquad\qquad -\ 2/6\ (1/2\ \log_2 1/2 + 1/2\ \log_2 1/2)$
$\qquad = 2/6$

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

# Information gain

- Decrease in entropy (uncertainty) after splitting

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

In our running example:

IG($X_1$) = H(Y) – H(Y|$X_1$)

$\quad$ = 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 |

# 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...

# A Decision Stump



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

Base Case One

Base Case Two

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
pchance = 0.317

maker = europe
2  2
pchance = 0.717

horsepower = low
0  0
Predict bad

horsepow
0  1
Predict goo

horsepower = low
0  4
Predict good

horsepower = medium
2  1
pchance = 0.894

horsepower = high
0  0
Predict bad

acceleration = low
1  0

ac

Don't split a node if data points are identical on remaining attributes

acceleration = low
1  0
Predict bad

acceleration = medium
1  1
(unexpandable)
Predict bad

ation = high
0  0
Predict bad

modelyear = 70to74
0  1
Predict good

modelyear = 75to78
1  0
Predict bad

modelyear = 79to83
0  0
Predict bad

Base Case Two: No attributes can distinguish

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 | | |

= 0.001

s = 5   cylinde
        8  0

Predict bad   pchance = 0.135   Predict bad   Predict

maker = america
0  10
Predict good

maker = asia
2  5
pchance = 0.317

maker = europe
2  2
pchance = 0.717

horsepower = low
0  0
Predict bad

horsepower = low
0  4
Predict good

horsepower = medium
2  1
pchance = 0.894

horsepower = high
0  0
Predict b

acceleration = low
1  0
Predict bad

acceleration = medium
1  1
(unexpandable)
Predict bad

...ion = high
0  0
Predict bad

modelyear = 7
0  1
Predict good

# 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

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

$y = a$ XOR $b$

The information gains:

The resulting decision tree:

# 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 |

Is it OK to omit Base Case 3?



y values: 0  1

root
2  2
pchance = 1.000

a = 0
1  1
pchance = 0.414

a = 1
1  1
pchance = 0.414

b = 0
1  0
Predict 0

b = 1
0  1
Predict 1

b = 0
0  1
Predict 1

b = 1
1  0
Predict 0

# Summary: Building Decision Trees

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 a non-leaf node with $n_X$ children.

  – The *i*'th child should be built by calling

  BuildTree($DS_i,Output$)

  Where $DS_i$ contains the records in DataSet where X = *i*th 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 = low    horsepower = medium    horsepower = high    acceleration = low    acceleration = medium    acceleration = high

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

…why?

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

Predict bad

# 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…

# Decision trees will overfit!!!

# How to Build Small Trees

Two reasonable approaches:

- Optimize on the held-out (development) set
  - If growing the tree larger hurts performance, then stop growing
  - Requires a larger amount of data…
- Use statistical significance testing
  - Test if the improvement for any split it likely due to noise
  - If so, don't do the split!
  - Can also use this to prune the tree bottom-up

# Real-Valued inputs

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

Infinite number of possible split values!!!

Finite dataset, only finite number of relevant splits!

| 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 | europe |
| 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 | europe |
| bad | 5 | 131 | 103 | 2830 | 15.9 | 78 | europe |
| | | | | | | | |

# "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 ≥ t

- **Requires small change**
  - Allow repeated splits on same variable
  - How does this compare to "branch on each value" approach?

# The set of possible thresholds

- Binary tree, split on attribute X
  - One branch: X < t
  - Other branch: X ≥ 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 $\{x_1,...,x_m\}$
  - Consider split points of the form $x_i + (x_{i+1} - x_i)/2$
  - Morever, only splits between examples of different classes matter!



(Figures from Stuart Russell)

# 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:

  - $H(Y|X:t) = p(X < t) H(Y|X < t) + p(X >= t) H(Y|X >= t)$

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

  - $IG^*(Y|X) = \max_t IG(Y|X:t)$

- Use: $IG^*(Y|X)$ for continuous variables

# Example with MPG

Example tree for our continuous dataset

# 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
    - Hypothesis testing