# Decision Trees

Machine Learning – CSE546

Carlos Guestrin (by Sameer Singh)

University of Washington

October 16, 2014

1

# Linear separability

- A dataset is **linearly separable** iff there exists a **separating hyperplane**:
  - ☐ Exists **w**, such that:
    - $w_0 + \sum_i w_i x_i > 0$; if $\mathbf{x}=\{x_1,\ldots,x_k\}$ is a positive example
    - $w_0 + \sum_i w_i x_i < 0$; if $\mathbf{x}=\{x_1,\ldots,x_k\}$ is a negative example

2

# Not linearly separable data

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

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

4

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

- Choose a classifier $h_w(x)$ that is non-linear in parameters $w$, e.g.,
  - □ Decision trees, boosting, nearest neighbor, neural networks…
- 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 quarter, we'll see that these options are not that different)

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 | 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 training examples

From the UCI repository (thanks to Ross Quinlan)

6

3

# 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 |

©Carlos Guestrin 2005-2013                                                                    7

# Recursion Step

mpg values:   bad   good

root
22   18
pchance = 0.001

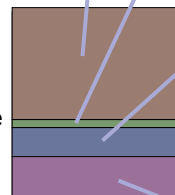| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0   0 | 4   17 | 1   0 | 8   0 | 9   1 |
| Predict bad | Predict good | Predict bad | Predict bad | Predict bad |

Take the Original Dataset..

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

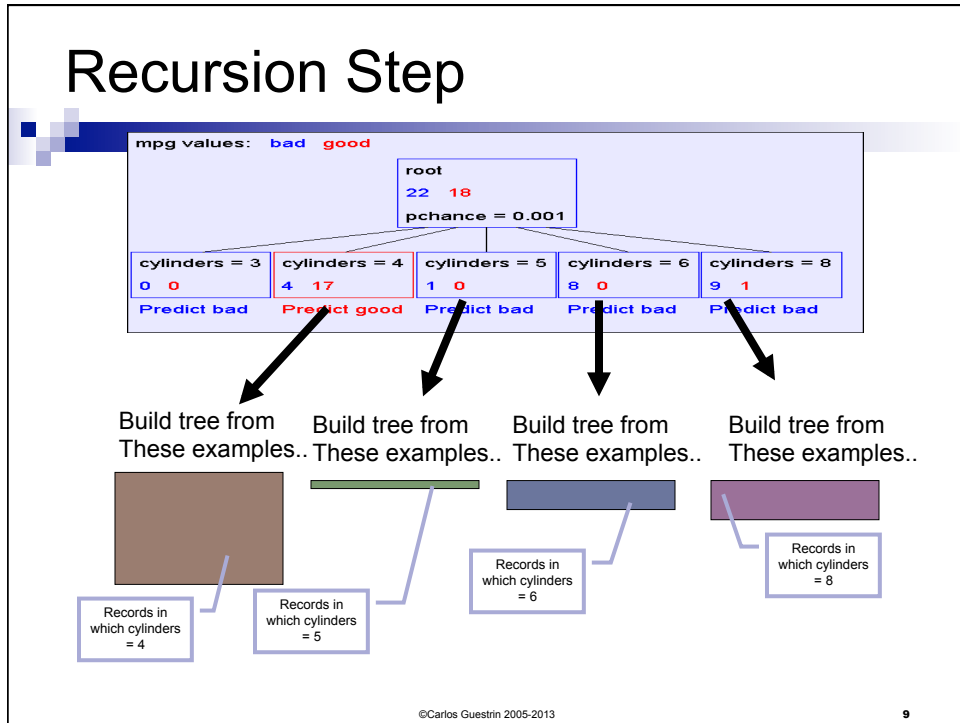Examples in which cylinders = 4

Examples in which cylinders = 5

Examples in which cylinders = 6

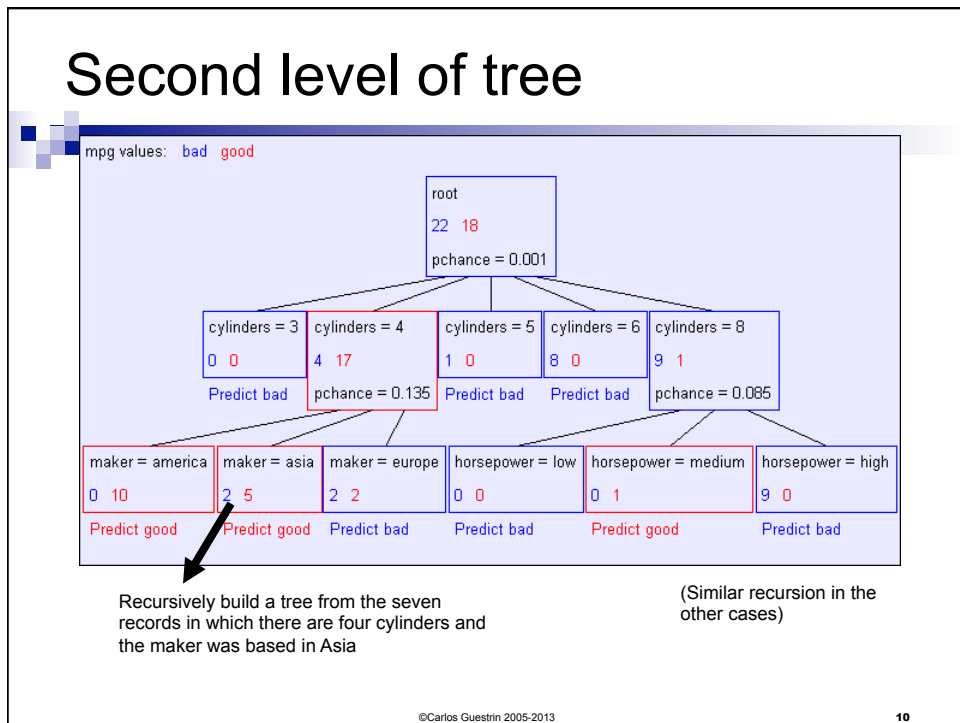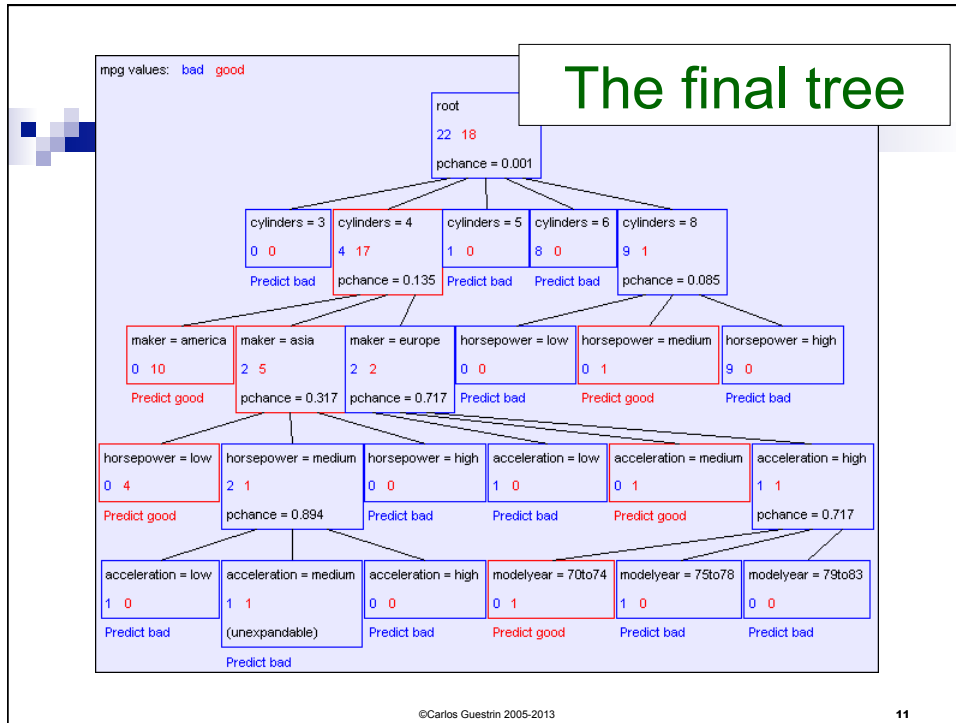Examples in which cylinders = 8

©Carlos Guestrin 2005-2013                                                                    8
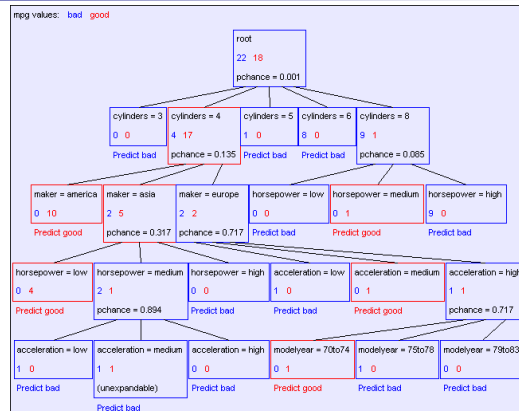
# Recursion Step

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 examples..

Build tree from These examples..

Build tree from These examples..

Build tree from These examples..

Records in which cylinders = 4

Records in which cylinders = 5

Records in which cylinders = 6

Records in which cylinders = 8

©Carlos Guestrin 2005-2013

9

# 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)

©Carlos Guestrin 2005-2013

10

5

The final tree

©Carlos Guestrin 2005-2013

11

# Classification of a new example

■ Classifying a test example – traverse tree and report leaf label



©Carlos Guestrin 2005-2013

12

# 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

14

# Choosing a good attribute

| X₁ | X₂ | 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 |

　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 |
|---|---|---|---|

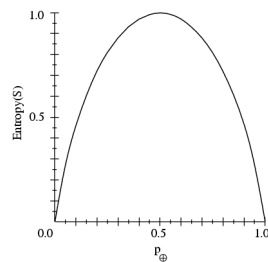| P(Y=A) = 1/4 | P(Y=B) = 1/4 | P(Y=C) = 1/4 | P(Y=D) = 1/4 |
|---|---|---|---|

　16

# 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)

17

---

# Information gain

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

- Advantage of attribute – decrease in uncertainty
  - □ Entropy of Y before you split

  - □ Entropy after split
    - Weight by probability of following each branch, i.e., normalized number of records

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

- Information gain is difference   $IG(X) = H(Y) - H(Y \mid X)$

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

**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 | | |

**22**

# 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 |

©Carlos Guestrin 2005-2013    23

---

## Base Case One

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

horsepower = low   horsepower = medium   horsepower = high
0  0                0  1                    9  0
Predict bad         Predict good            Predict bad

pchance = 0.717

medium   horsepower = high   acceleration = low   acceleration = medium   acceleration = high
         0  0                1  0                 0  1                    1  1
         Predict bad         Predict bad          Predict good            pchance = 0.717

medium   acceleration = high   modelyear = 70to74   modelyear = 75to78   modelyear = 79to83
         0  0                  0  1                 1  0                 0  0
         Predict bad           Predict good          Predict bad          Predict bad

1  0     1  1
Predict bad  (unexpandable)
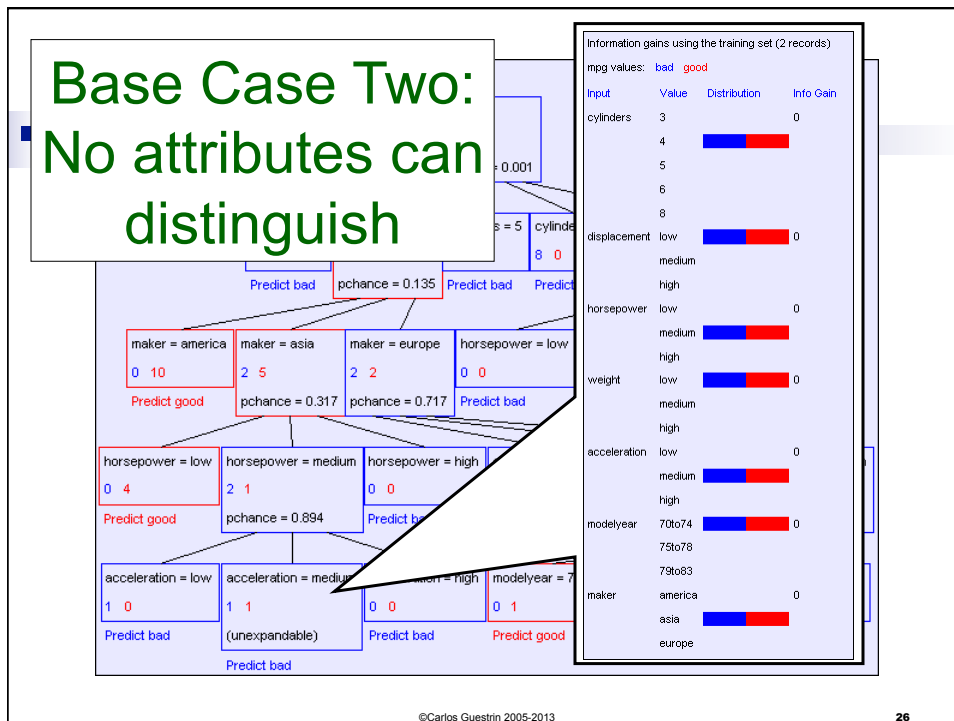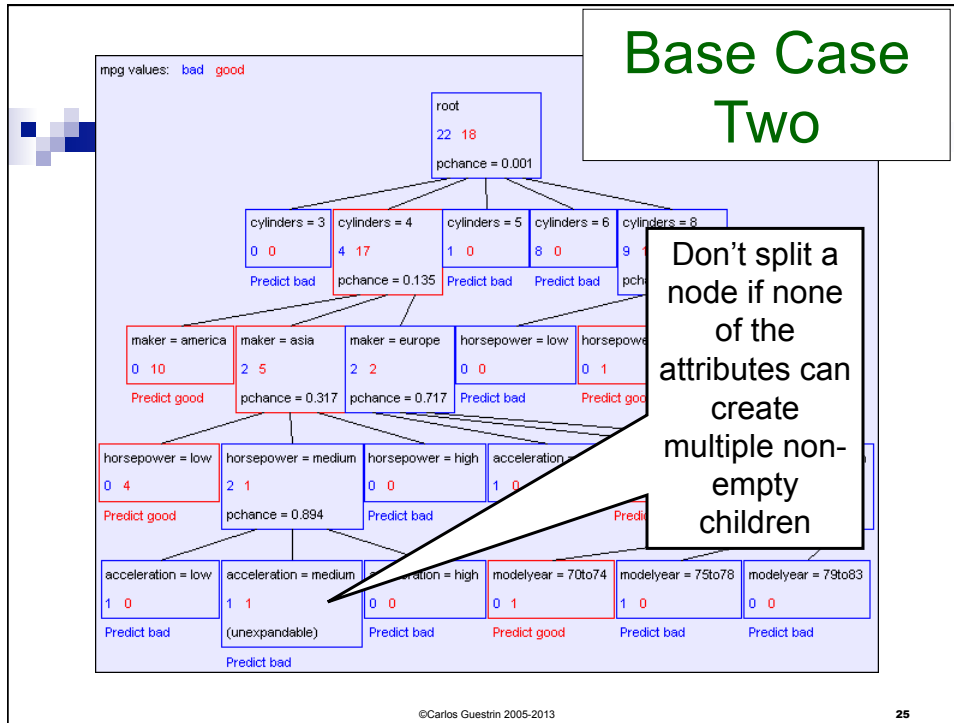Predict bad

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

©Carlos Guestrin 2005-2013    24

11

Base Case Two

mpg values: bad good

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

©Carlos Guestrin 2005-2013

25



Base Case Two: No attributes can distinguish

Information gains using the training set (2 records)
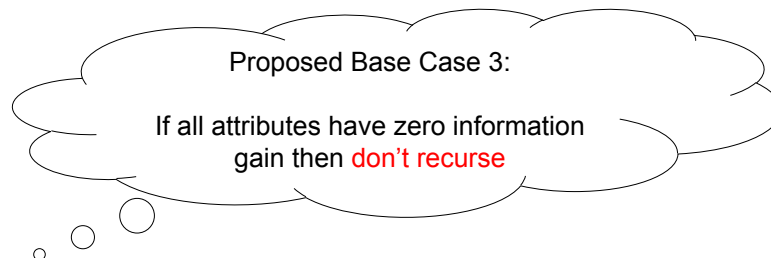
©Carlos Guestrin 2005-2013

26

12

# 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

27

# Base Cases: An idea

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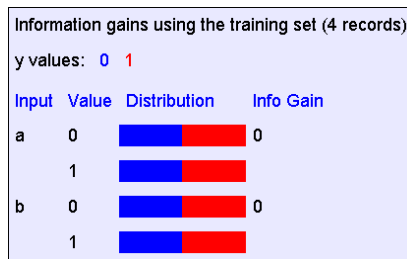- Base Case Two: If all records have exactly the same set of input attributes then don't recurse

Proposed Base Case 3:

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

- *Is this a good idea?*

28

# 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 bad decision tree:

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 | | |

y values: 0 1

root

2 2

Predict 0

©Carlos Guestrin 2005-2013    29

# If we omit Base Case 3:

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

y = a XOR b

The resulting decision tree:

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

©Carlos Guestrin 2005-2013    30

14

# 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"
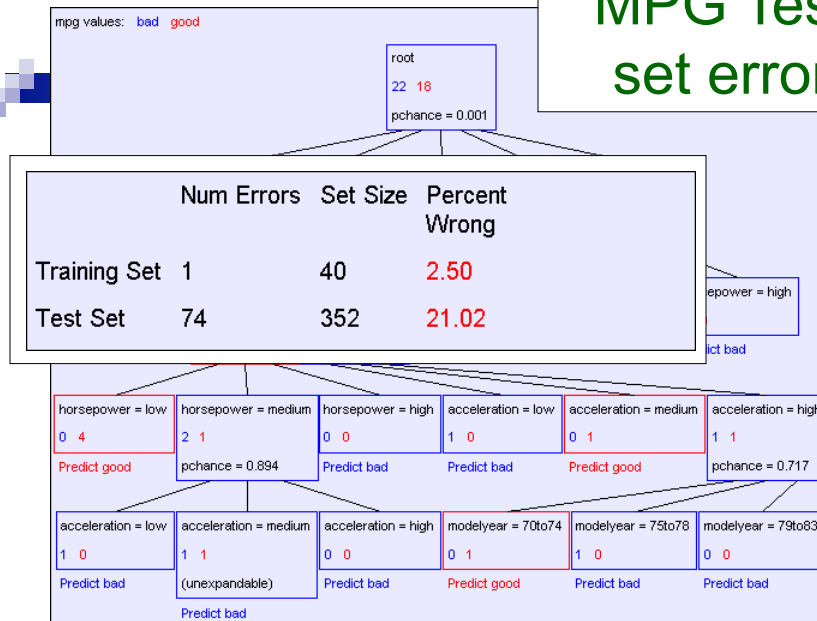- If all input values are the same, return a leaf node that says "predict the majority output"
- Else find attribute $X$ with highest Info Gain
- Suppose $X$ has $n_X$ distinct values (i.e. X has arity $n_X$).
  - ☐ Create and return a non-leaf node with $n_X$ children.
  - ☐ The $i$'th child should 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.

31

---

## 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 |

epower = high

ict bad

| horsepower = low | horsepower = medium | horsepower = high | acceleration = low | acceleration = medium | acceleration = high |
|---|---|---|---|---|---|
| 0  4 | 2  1 | 0  0 | 1  0 | 0  1 | 1  1 |
| Predict good | pchance = 0.894 | Predict bad | Predict bad | Predict good | pchance = 0.717 |

| acceleration = low | acceleration = medium | acceleration = high | modelyear = 70to74 | modelyear = 75to78 | modelyear = 79to83 |
|---|---|---|---|---|---|
| 1  0 | 1  1 | 0  0 | 0  1 | 1  0 | 0  0 |
| Predict bad | (unexpandable) | Predict bad | Predict good | Predict bad | Predict bad |
| | Predict bad | | | | |

32

15

## 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 |

...epower = high

...ict bad

| 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?

= 0.717

= 79to83

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

Predict bad

33

---

# Decision trees & Learning Bias

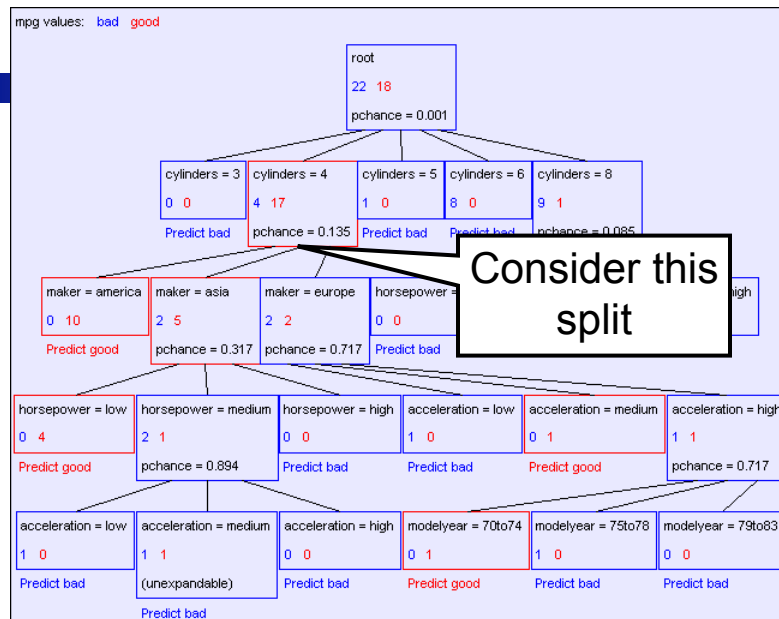| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|---|---|---|---|---|---|---|---|
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | medium | medium | medium | low | 75to78 |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

34

16

# 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
  - Will definitely overfit!!!
  - Must bias towards simpler trees
- Many strategies for picking simpler trees:
  - Fixed depth
  - Fixed number of leaves
  - Or something smarter…

35



36

17

# A chi-square test

mpg values:   bad   good

| maker | america | 0 | 10 | | H( mpg \| maker = america ) = 0 |
|-------|---------|---|----|---|------|
| | asia | 2 | 5 | | H( mpg \| maker = asia ) = 0.863121 |
| | europe | 2 | 2 | | H( mpg \| maker = europe ) = 1 |

H(mpg) = 0.702467   H(mpg|maker) = 0.478183

IG(mpg|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?

37

# A chi-square test

mpg values:   bad   good

| maker | america | 0 | 10 | | H( mpg \| maker = america ) = 0 |
|-------|---------|---|----|---|------|
| | asia | 2 | 5 | | H( mpg \| maker = asia ) = 0.863121 |
| | europe | 2 | 2 | | H( mpg \| maker = europe ) = 1 |

H(mpg) = 0.702467   H(mpg|maker) = 0.478183

IG(mpg|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 see readings…)

38

# 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 $p_{chance}$ > *MaxPchance*
  - ☐ Continue working you 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

39

# Pruning example

- With MaxPchance = 0.1, you will see the following MPG decision 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 | Predict good | Predict bad | Predict bad | Predict bad |

Note the improved test set accuracy compared with the unpruned tree

|  | Num Errors | Set Size | Percent Wrong |
|---|---|---|---|
| Training Set | 5 | 40 | 12.50 |
| Test Set | 56 | 352 | 15.91 |

40

# MaxPchance

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

Expected True Error

Decreasing ← → Increasing

MaxPchance

High Bias                    High Variance

41

---

# 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 | 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 |

Infinite number of possible split values!!!

Finite dataset, only finite number of relevant splits!

Idea One: Branch on each possible real value

42

20

# "One branch for each numeric value" idea:

mpg values:  bad   good

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | root | | | | | | | | |
| | | | | 22  18 | | | | | | | | |
| | | | | pchance = 0.222 | | | | | | | | |

| modelyear = 70 | modelyear = 71 | modelyear = 72 | modelyear = 73 | modelyear = 74 | modelyear = 75 | modelyear = 76 | modelyear = 77 | modelyear = 78 | modelyear = 79 | modelyear = 80 | modelyear = 81 | modelyear = 82 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 4  0 | 2  1 | 1  0 | 6  1 | 1  2 | 0  0 | 3  1 | 1  3 | 3  0 | 1  1 | 0  0 | 0  5 | 0  4 |
| Predict bad | Predict bad | Predict bad | Predict bad | Predict good | Predict bad | Predict bad | Predict good | Predict bad | Predict bad | Predict bad | Predict good | Predict good |

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

43

---

# Threshold splits

- Binary tree, split on attribute $X_i$
  - One branch: $X_i < t$
  - Other branch: $X_i \geq t$

44

# Choosing threshold split

- Binary tree, split on attribute $X_i$
  - One branch: $X_i < t$
  - Other branch: $X_i \geq t$
- Search through possible values of $t$
  - Seems hard!!!
- But only finite number of $t$'s are important
  - Sort data according to $X_i$ into $\{x_1,\ldots,x_m\}$
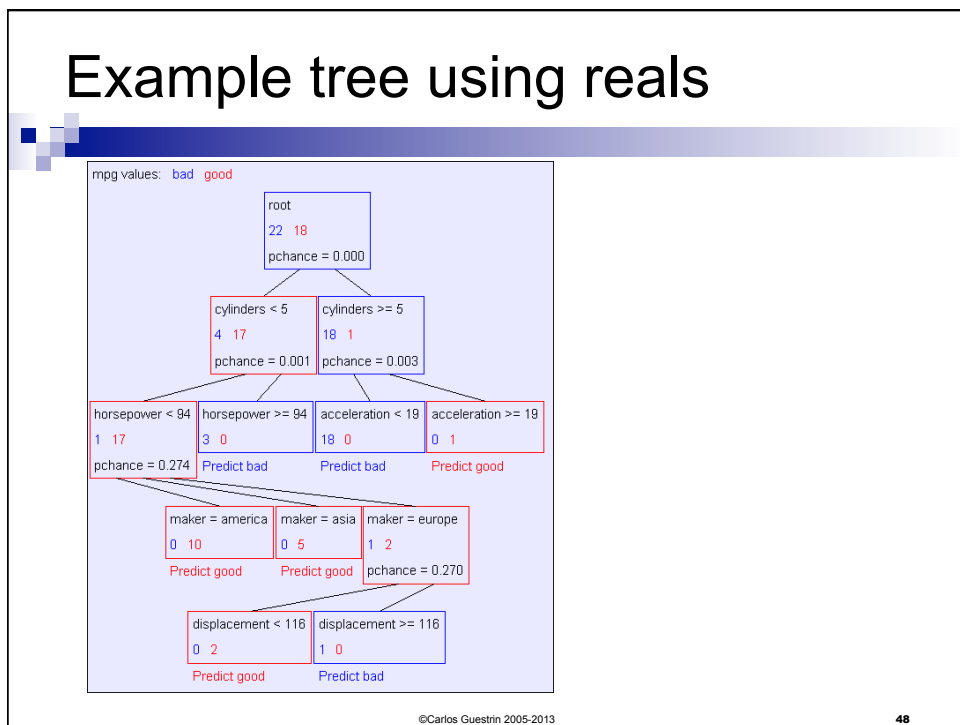  - Consider split points of the form $x_a + (x_{a+1} - x_a)/2$

45

# A better idea: thresholded splits

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

- Note, may split on an attribute multiple times, with different thresholds

46

## 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 | | |
| | europe | | |

©Carlos Guestrin 2005-2013                                    47

## Example tree using reals



mpg values:  bad   good

root
22  18
pchance = 0.000

cylinders < 5         cylinders >= 5
4  17                 18  1
pchance = 0.001       pchance = 0.003

horsepower < 94   horsepower >= 94   acceleration < 19   acceleration >= 19
1  17             3  0               18  0               0  1
pchance = 0.274   Predict bad        Predict bad         Predict good

maker = america   maker = asia   maker = europe
0  10             0  5           1  2
Predict good      Predict good   pchance = 0.270

displacement < 116   displacement >= 116
0  2                 1  0
Predict good         Predict bad

©Carlos Guestrin 2005-2013                                    48

23

## 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

©Carlos Guestrin 2005-2013                                          49

## Acknowledgements

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

©Carlos Guestrin 2005-2013                                          50