## CSE 446: Decision Trees
## Winter 2012

Slides adapted from Carlos Guestrin and Andrew Moore by
Luke Zettlemoyer & Dan Weld

---

## Overview of Learning

Type of Supervision
(eg, Experience, Feedback)

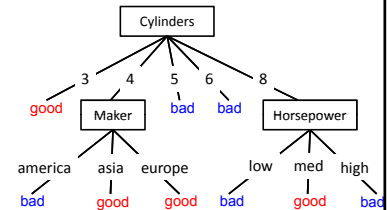| What is Being Learned? | | Labeled Examples | Reward | Nothing |
|---|---|---|---|---|
| | Discrete Function | Classification | | Clustering |
| | Continuous Function | Regression | | |
| | Policy | Apprenticeship Learning | Reinforcement Learning | |

5

---

## A learning problem: predict fuel efficiency

- 40 Records
- Discrete data (for now)

- 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 |

From the UCI repository (thanks to Ross Quinlan)

---

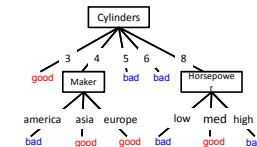## Hypotheses: decision trees $f : X \rightarrow Y$

- 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 $y$
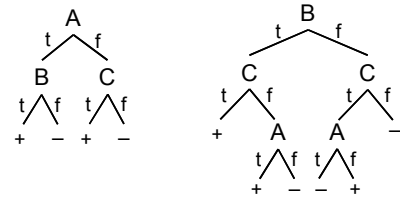


---

## 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?



---

## Are all decision trees equal?

- Many trees can represent the same concept
- But, not all trees will have the same size!

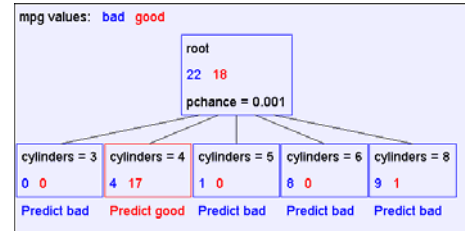  e.g., $\phi = (A \wedge B) \vee (\neg A \wedge C)$



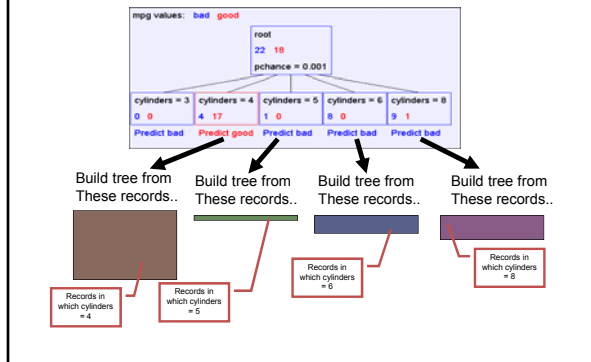- Which tree do we prefer?

## Learning decision trees is hard!!!

- Finding 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

## Improving Our Tree

predict
mpg=bad



## Recursive Step



## A full tree



## Two Questions

Greedy Algorithm:
- Start from empty decision tree
- Split on the **best attribute (feature)**
- Recurse

1. Which attribute gives the best split?
2. When to stop recursion?

## Which attribute gives the best split?

$A_1$: The one with the highest *information gain*
  Defined in terms of *entropy*

$A_2$: Actually many alternatives, eg, *accuracy*
  Seeks to reduce the *misclassification rate*

## 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 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 |

## 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)
            - 2/6 (1/2 $\log_2$ 1/2 + 1/2 $\log_2$ 1/2)
       = 2/6
       = 0.33

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

## Information Gain

*Advantage of attribute* – 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$_1$)
   = 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 |

## Alternate Splitting Criteria

Misclassification Impurity

Minimum probability that a training pattern will be misclassified

$$M(Y) = 1 - \max_i P(Y = y_i)$$

Misclassification Gain

$$IG_M(X) = [1 - \max_i P(Y = y_i)] - [1 - (\max_j \; \max_i P(Y = y_i \mid x = x_j))]$$
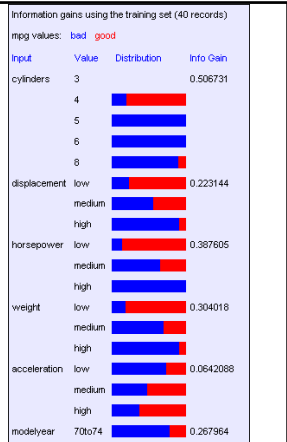
37

## Learning Decision Trees

- Start from empty decision tree
- Split on **next best attribute (feature)**
  - Use information gain (or…?) to select attribute:
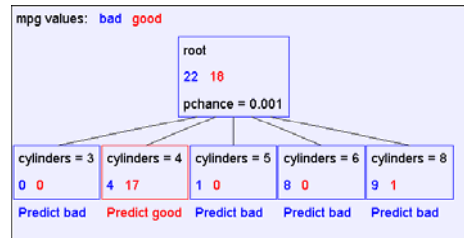  $$\arg\max_i IG(X_i) = \arg\max_i H(Y) - H(Y \mid X_i)$$
- Recurse

Suppose we want to predict MPG

predict mpg=bad

Now, Look at all the information gains…

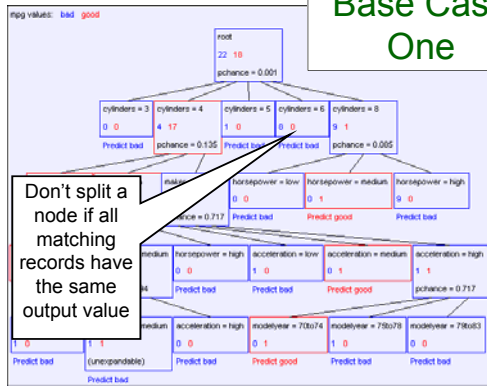Information gains using the training set (40 records)

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 |

---

## Tree After One Iteration

mpg values: bad good

| root |
|---|
| 22 18 |
| pchance = 0.001 |

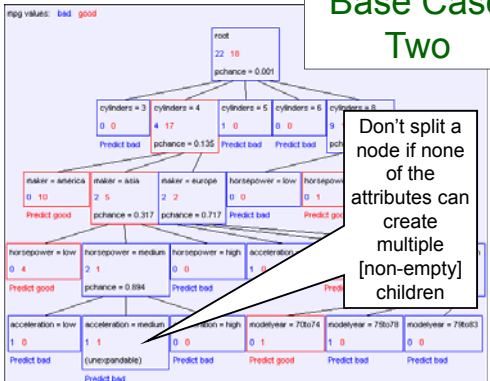| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0 0 | 4 17 | 1 0 | 8 0 | 9 1 |
| Predict bad | Predict good | Predict bad | Predict bad | Predict bad |

---

## When to Terminate?
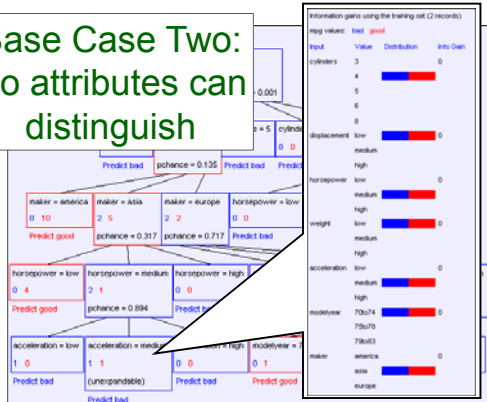
41

---



## Base Case One

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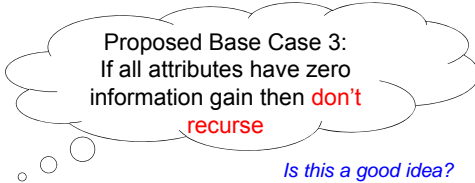
---



## Base Case Two

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

---



## Base Case Two: No attributes can distinguish

4

## 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

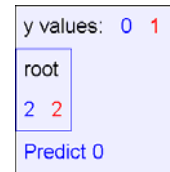| a | b | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

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

---

## But *Without* Base Case 3:

y = a XOR b

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

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

**So: Base Case 3? Include or Omit?**

---

## Building Decision Trees

**BuildTree(*DataSet,Output*)**

**If** all output values are the same in *DataSet*,
    **Then** return a leaf node that says "predict this unique output"
**If** all input values are the same,
    **Then** return a leaf node that says "predict the majority output"
**Else** find attribute *X* with highest Info Gain
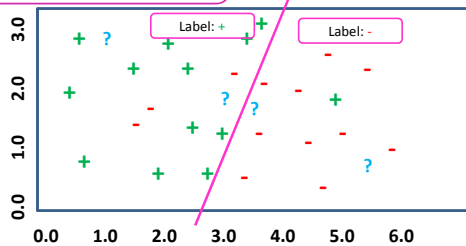    Suppose *X* has $n_X$ distinct values (i.e. X has arity $n_X$).
    Create and return a non-leaf node with $n_X$ children.
    The $i^{th}$ child is built by calling **BuildTree(*DS$_i$,Output*)**
      Where *DS$_i$* consists of all those records in DataSet
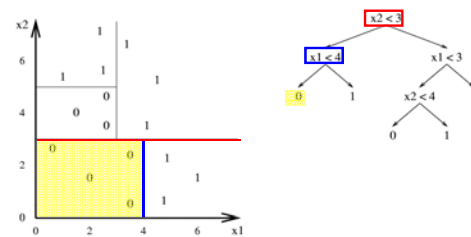        for which X = $i^{th}$ distinct value of X.

---

## General View of a Classifier

Hypothesis:
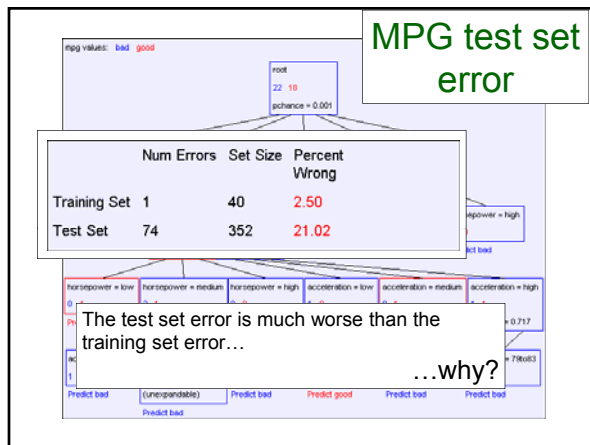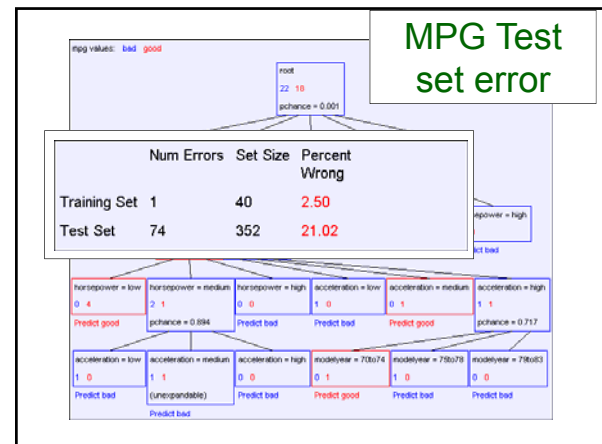Decision Boundary for labeling function

Label: +
Label: -

---

**Decision Tree Decision Boundaries**

Decision trees divide the feature space into axis-parallel rectangles, and label each rectangle with one of the $K$ classes.

x2 < 3
x1 < 4
x1 < 3
0
1
x2 < 4
0
1

## Ok, so how does it perform?

51

---

### MPG Test set error



| | Num Errors | Set Size | Percent Wrong |
|---|---|---|---|
| Training Set | 1 | 40 | 2.50 |
| Test Set | 74 | 352 | 21.02 |

---

### MPG test set error



| | Num Errors | Set Size | Percent Wrong |
|---|---|---|---|
| Training Set | 1 | 40 | 2.50 |
| Test Set | 74 | 352 | 21.02 |

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

…why?

---

### Decision trees will overfit

- Our 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 introduce some bias towards *simpler* trees

- Why might one pick simpler trees?

---

### Occam's Razor

- Why Favor Short Hypotheses?
- Arguments for:
  - Fewer short hypotheses than long ones
    - →A short hyp. less likely to fit data by coincidence
    - →Longer hyp. that fit data may might be coincidence
- Arguments against:
  - Argument above on really uses the fact that hypothesis *space* is small!!!
  - What is so special about small sets based on the *complexity* of each *hypothesis*?

---

### How to Build Small Trees

Several reasonable approaches:
- **Stop growing tree before overfit**
  - Bound depth or # leaves
  - Base Case 3
  - *Doesn't work well in practice*

- **Grow full tree; then prune**
  - **Optimize on a held-out (development set)**
    - If growing the tree hurts performance, then cut back
    - Con: Requires a larger amount of data…
  - **Use statistical significance testing**
    - Test if the improvement for any split is likely due to noise
    - If so, then prune the split!
  - **Convert to logical rules**
    - Then simplify rules

6

## Reduced Error Pruning
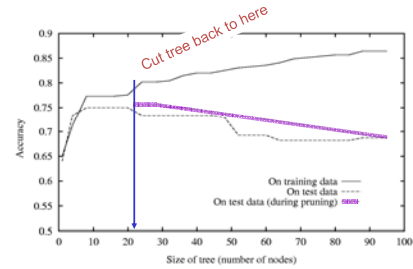
Split data into *training* & *validation* sets (10-33%)

Train on training set (overfitting)

Do until further pruning is harmful:

1) Evaluate effect on validation set of pruning *each* possible node (and tree below it)

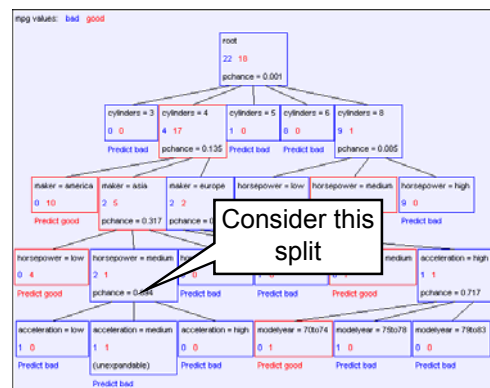2) Greedily remove the node that *most improves accuracy of validation set*

57

---

## Effect of Reduced-Error Pruning



---

## Alternatively

- Chi-squared pruning
  - Grow tree fully
  - Consider leaves in turn
    - Is parent split worth it?

- Compared to Base-Case 3?

59

---



mpg values: bad good

Consider this split

---

## A chi-square test



mpg values: bad good

| maker | | | | H(mpg \| maker = america) = 0 |
|---|---|---|---|---|
| america | 0 | 10 | | |
| 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 13.5%

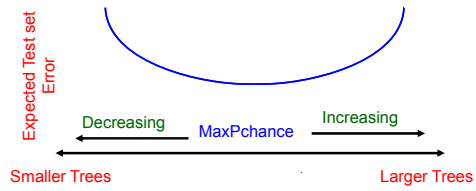Such hypothesis tests are relatively easy to compute, but involved

---

## 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

---

7

## Regularization

- Note for Future: MaxPchance is a regularization parameter that helps us bias towards simpler models



We'll learn to choose the value of magic parameters like this one later!

## 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
- Improved by
  - Carlos Guestrin &
  - Luke Zettlemoyer