

CS 536: Machine Learning

Decision Trees

Fall 2005

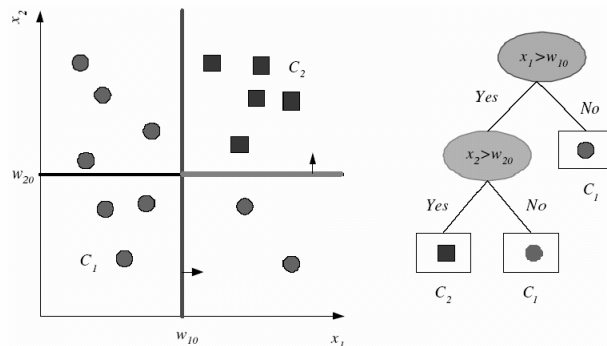
Ahmed Elgammal

Dept of Computer Science

Rutgers University

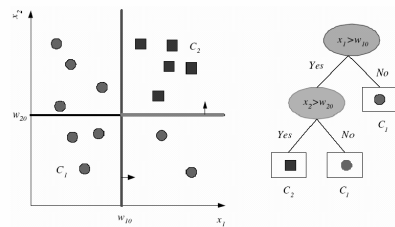
CS 536 - Fall 2005 - -

- Decision Tree: a hierarchical model for supervised learning whereby the local region is identified in a sequence of recursive splits.



CS 536 - Fall 2005 - -

- Internal Decision nodes: Each node m implement a test function $f_m(x)$ with discrete outcomes labeling the branches. Given an input, the test is applied and one of the branches is taken depending on the outcome.
- Terminal leaves: output: class code (for classification) or numeric value (for regression).
- Each $f_m(x)$ defines a discriminant in the d -dimensional input space dividing it into smaller regions which are further subdivided as we take a path from the root down.
- Each terminal leaf defines a localized region in the input space where instances falling in this region have the same label.



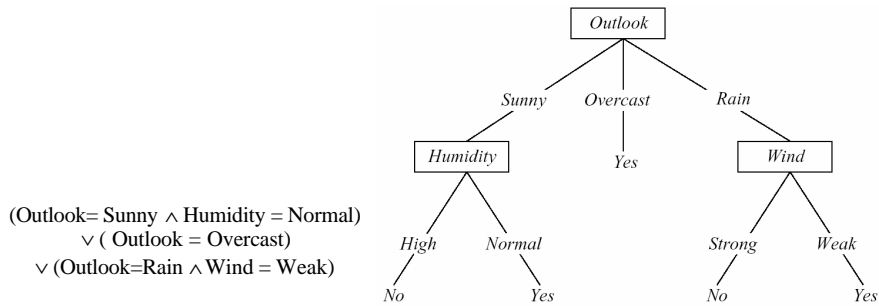
CS 536 – Fall 2005 - -

Training Examples for the target concept **PlayTennis**

Day	Outlook	Temp	Hum.	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Nml	Weak	Yes
D6	Rain	Cool	Nml	Strong	No
D7	Overcast	Cool	Nml	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Nml	Weak	Yes
D10	Rain	Mild	Nml	Weak	Yes
D11	Sunny	Mild	Nml	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Nml	Weak	Yes
D14	Rain	Mild	High	Strong	No

CS 536 – Fall 2005 - -

- Decision Tree represents a disjunction of conjunctions of the constraints on the attributes.
- Each path from the root to a leaf correspond to a conjunction of attribute tests.
- The tree itself is a disjunction of these conjunctions.

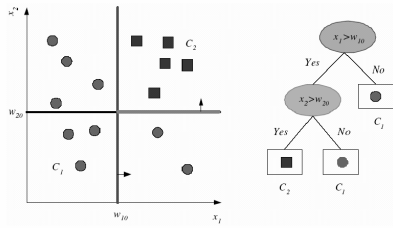


CS 536 – Fall 2005 - -

Different kinds of decision Trees:

- Internal decision nodes
 - Univariate: Uses a single attribute, x_i
 - Numeric x_i : Binary split : $x_i > w_m$
 - Discrete x_i : n -way split for n possible values
 - Multivariate: Uses all attributes, \mathbf{x}
- Leaves
 - Classification: Class labels, or proportions
 - Regression: Numeric; r average, or local fit
- Learning is greedy; find the best split recursively (Breiman et al, 1984; Quinlan, 1986, 1993)

CS 536 – Fall 2005 - -

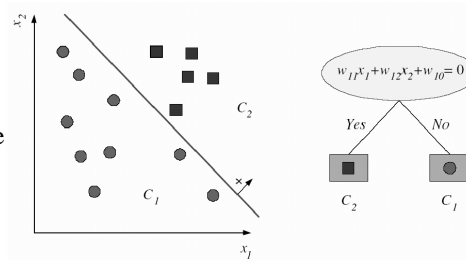


Univariate Tree:

- One attribute at a time:
- Each decision node defines axis-parallel hyperplane.
- Each leaf defines a hyper rectangular decision surface.

Multivariate decision tree:

- all attributes each time:
- each decision node defines a hyperplane
- Each leaf defines a polyhedral decision surface



CS 536 - Fall 2005 -

Whence Decision Trees?

Consider: Discrete Univariate Classification Trees

- Instances describable by attribute-value pairs
- Target function is discrete valued
- Disjunctive hypothesis may be required
- Possibly noisy training data or missing attribute values

Examples:

- Equipment or medical diagnosis
- Credit risk analysis
- Many successful applications that outperform human experts.

CS 536 - Fall 2005 -

Evolution of Decision Trees:

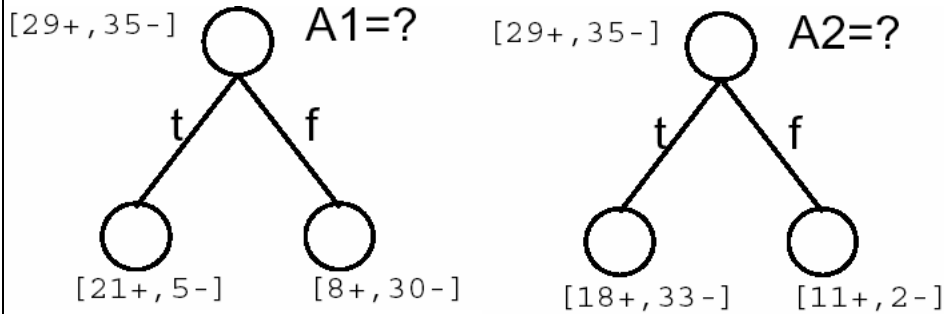
- CLS (Concept Learning System) Earl Hunt 1960's
- ID3 (Interactive Dichotomizer 3) Quinlan 70's and 80's
- C4.5 Quinlan 90's

Top-Down Induction

Main loop:

1. $A \leftarrow$ the "best" decision attribute for next *node*
2. Assign A as decision attribute for *node*
3. For each value of A , create new descendant of *node*
4. Sort training examples to leaf nodes
5. If training examples perfectly classified, Then STOP,
Else iterate over new leaf nodes

Which Attribute is Best?



What is a good quantitative measure of the worth of an attribute ?

CS 536 - Fall 2005 -

Measuring Entropy

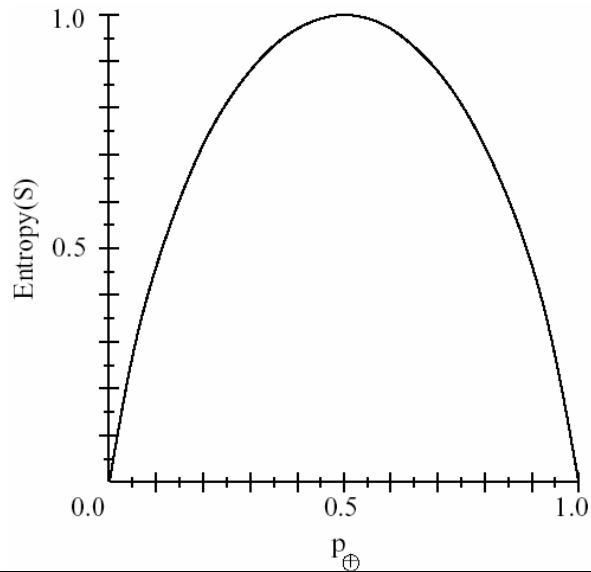
- S is a sample of training examples
- p_{\oplus} is the proportion of positive examples in S
- p_{\otimes} is the proportion of negative examples in S

Entropy measures the impurity of S

$$Entropy(S) = -p_{\oplus} \log p_{\oplus} - p_{\otimes} \log p_{\otimes}$$

CS 536 - Fall 2005 -

Entropy Function



Entropy

Entropy(S) = expected number of bits needed to encode class \oplus or \otimes of a randomly drawn member of S (under the optimal, shortest-length code)

Why?

Information theory: optimal length code assigns $-\log_2 p$ bits to message having probability p .

So, expected number of bits to encode \oplus or \otimes of a random member of S :

$$p_{\oplus} (-\log p_{\oplus}) + p_{\otimes} (-\log p_{\otimes})$$

Information Gain

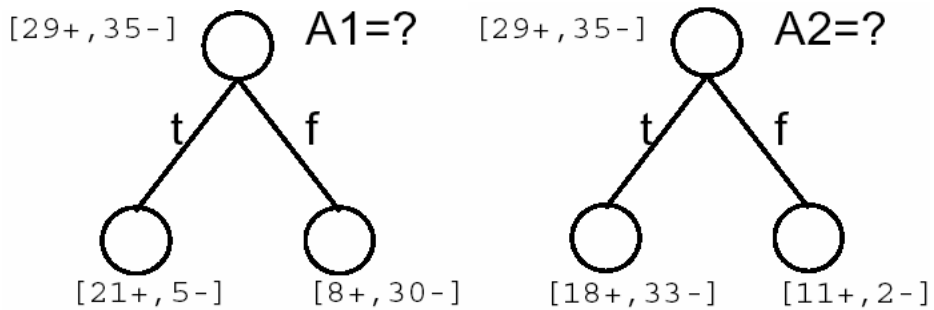
$Gain(S, A)$ = expected reduction in entropy due to sorting S on A

$Gain(S, A) \equiv$

$$Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

Here, S_v is the set of training instances remaining from S after restricting to those for which attribute A has value v .

Which Attribute is Best?

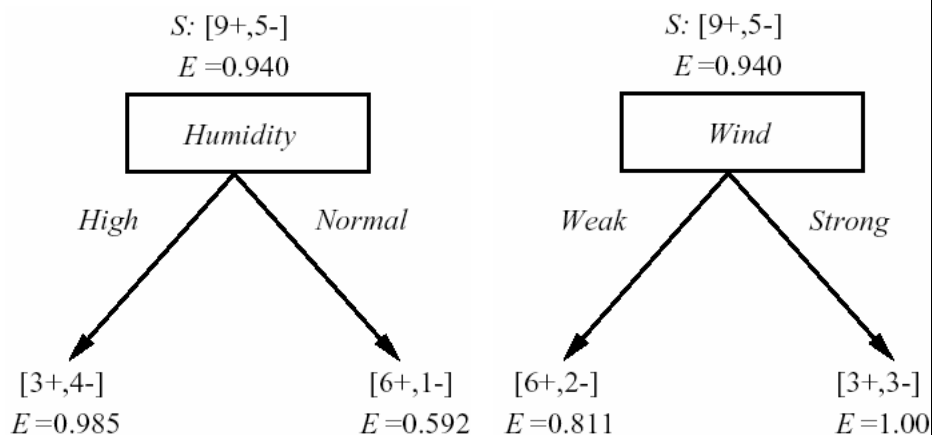


Training Examples

Day	Outlook	Temp	Hum.	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Nml	Weak	Yes
D6	Rain	Cool	Nml	Strong	No
D7	Overcast	Cool	Nml	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Nml	Weak	Yes
D10	Rain	Mild	Nml	Weak	Yes
D11	Sunny	Mild	Nml	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Nml	Weak	Yes
D14	Rain	Mild	High	Strong	No

CS 536 - Fall 2005 -

Selecting the Next Attribute

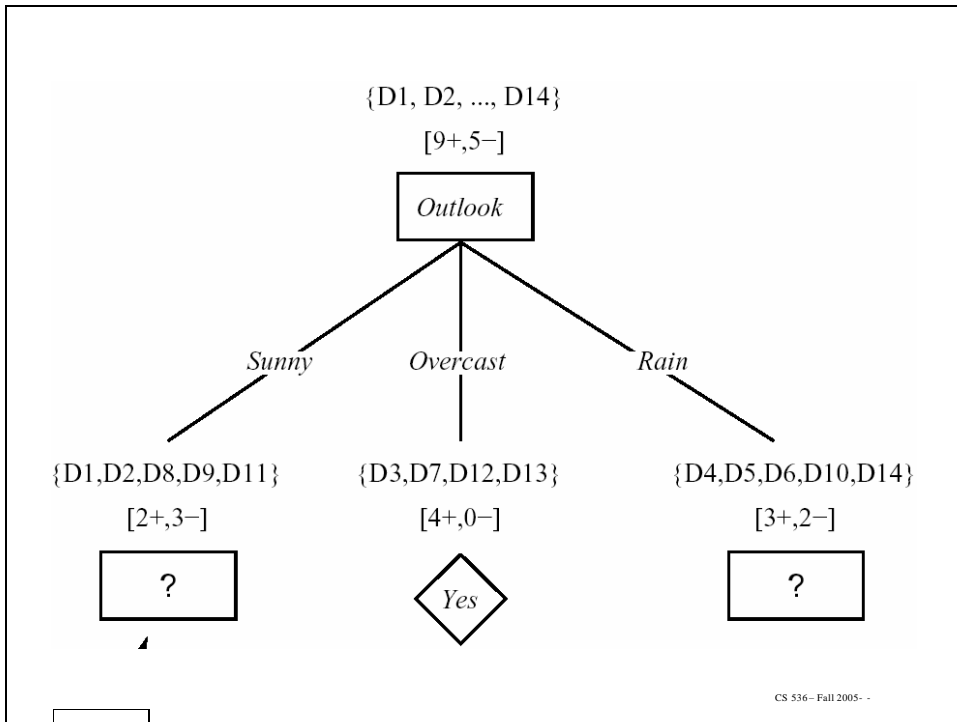


Which attribute is the best classifier?

$$Gain(S, Humidity) = .940 - (7/14).985 - (7/14).592 = .151$$

$$Gain(S, Wind) = .940 - (8/14).811 - (6/14)1.0 = .048$$

CS 536 - Fall 2005 -



Comparing Attributes

$$S_{\text{sunny}} = \{D1, D2, D8, D9, D11\}$$

- $\text{Gain}(S_{\text{sunny}}, \text{Humidity})$
 $= .970 - (3/5) 0.0 - (2/5) 0.0 = .970$
- $\text{Gain}(S_{\text{sunny}}, \text{Temp})$
 $= .970 - (2/5) 0.0 - (2/5) 1.0 - (1/5) 0.0 = .570$
- $\text{Gain}(S_{\text{sunny}}, \text{Wind})$
 $= .970 - (2/5) 1.0 - (3/5) .918 = .019$

What is ID3 Optimizing?

The hypothesis space searched by ID3 is the set of possible decision trees. Simple-to-Complex hill-climbing (greedy) search guided by information gain measure

How would you find a tree that minimizes:

- misclassified examples?
- expected entropy?
- expected number of tests?
- depth of tree given a fixed accuracy?
- etc.?

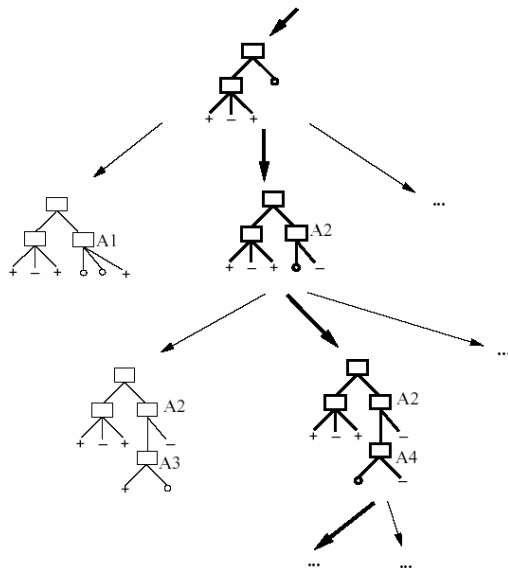
How decide if one tree beats another?

CS 536 - Fall 2005 -

Hypothesis Space Search by ID3

ID3:

- representation: trees
- scoring : entropy
- search : greedy



Hypothesis Space Search by ID3

- Hypothesis space is complete!
 - Target function surely in there...
- Outputs a single hypothesis (which one?)
 - Can't generate all consistent hypotheses... Single Concept.
- No back tracking
 - Local minima... Not globally optimal.
- Statistically-based search choices
 - Robust to noisy data...
- Inductive bias ~ “prefer shortest tree”

Inductive Bias in ID3

Note H is the power set of instances X

- Unbiased?

Not really...

- Preference for short trees, and for those with high information gain attributes near the root
- Bias is a *preference* for some hypotheses, rather than a *restriction* of hypothesis space H
- Occam's razor: prefer the shortest hypothesis that fits the data

Occam's Razor

Why prefer short hypotheses?

Argument in favor:

- Fewer short hyps. than long hyps.
 - a short hyp that fits data unlikely to be coincidence
 - a long hyp that fits data might be coincidence

Argument opposed:

- There are many ways to define small sets of hyps
- e.g., all trees with a prime number of nodes that use attributes beginning with "Z"
- What's so special about small sets based on size of hypothesis??

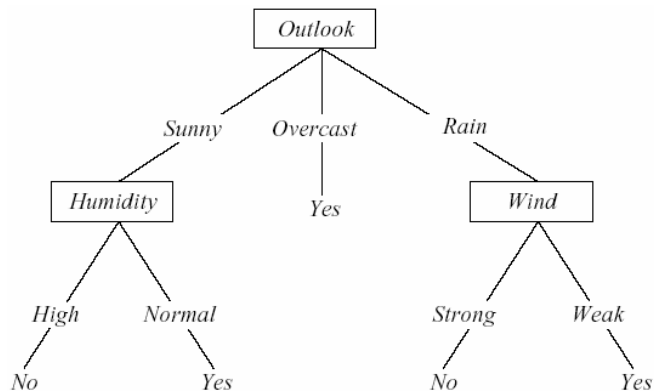
CS 536 - Fall 2005 -

Overfitting

Consider adding noisy training example #15:

Sunny, Hot, Normal, Strong, PlayTennis = No

What effect on earlier tree?



CS 536 - Fall 2005 -

Overfitting

Consider error of hypothesis h over

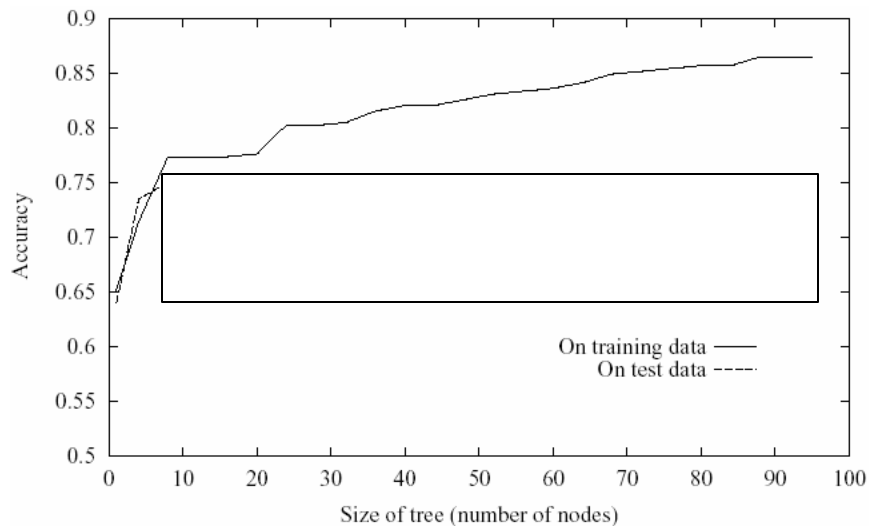
- training data: $error_{train}(h)$
- entire distribution D of data: $error_D(h)$

Hypothesis h in H **overfits** training data if there is an alternative hypothesis h' in H such that

- $error_{train}(h) < error_{train}(h')$, and
- $error_D(h) > error_D(h')$

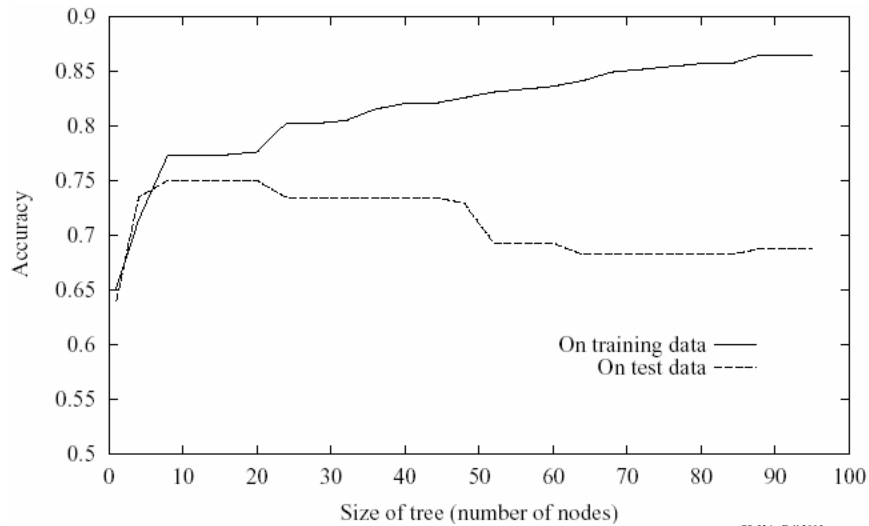
CS 536 – Fall 2005 - -

Overfitting in Learning



CS 536 – Fall 2005 - -

Overfitting in Learning



Avoiding Overfitting

How can we avoid overfitting?

- stop growing when data split not statistically significant
- grow full tree, then post-prune (DP alg!)

How to select “best” tree:

- Measure performance over training data
- Measure performance over separate validation data set
- MDL: minimize

$$size(tree) + size(misclassifications(tree))$$

Reduced-Error Pruning

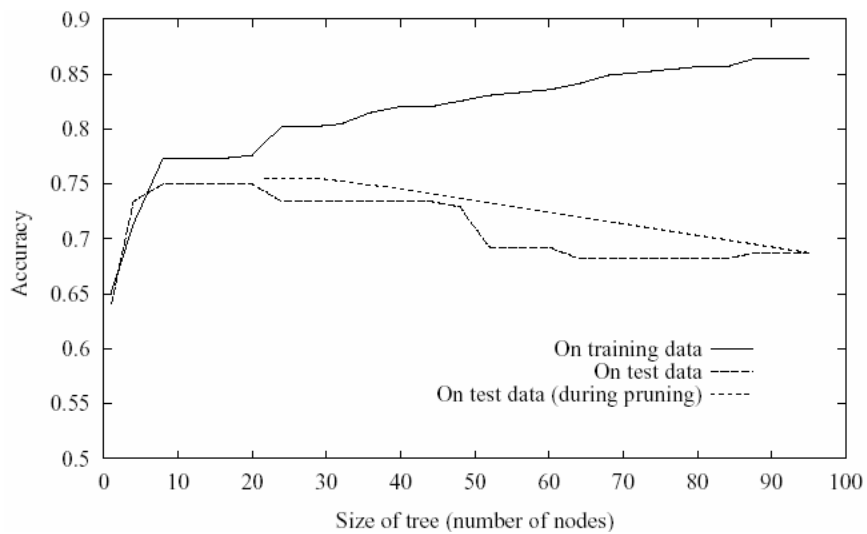
Split data into *training* and *validation* set

Do until further pruning is harmful:

1. Evaluate impact on *validation* set of pruning each possible node (plus those below it)
2. Greedily remove the one that most improves *validation* set accuracy
 - produces smallest version of most accurate subtree
 - What if data is limited?

CS 536 - Fall 2005 -

Effect of Pruning



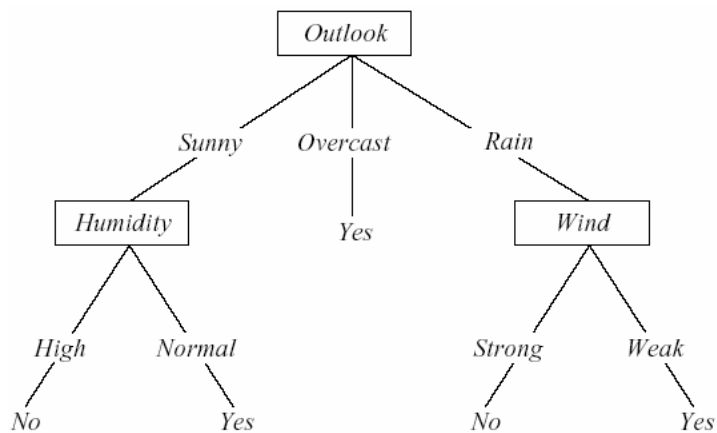
CS 536 - Fall 2005 -

Rule Post-Pruning

1. Convert tree to equivalent set of rules
 2. Prune each rule independently of others
 3. Sort final rules into desired sequence for use
- Perhaps most frequently used method (e.g., C4.5)

CS 536 - Fall 2005 -

Converting Tree to Rules



CS 536 - Fall 2005 -

The Rules

IF (Outlook = Sunny) ^ (Humidity = High)
THEN PlayTennis = No
IF (Outlook = Sunny) ^ (Humidity = Normal)
THEN PlayTennis = Yes
...

Attributes with Many Values - C4.5

Problem:

- If one attribute has many values compared to the others, *Gain* will select it
- Imagine using *Date = Jun_3_1996* as attribute

One approach: use *GainRatio* instead

$$\text{GainRatio}(S,A) \equiv \text{Gain}(S,A) / \text{SplitInfo}(S,A)$$

$$\text{SplitInfo}(S,A) \equiv -\sum_{i=1}^c |S_i|/|S| \log_2 |S_i|/|S|$$

where S_i is subset of S for which A has value v_i

Attributes with Costs

Consider

- medical diagnosis, *BloodTest* has cost \$150
- robotics, *Width_from_1ft* has cost 23 sec.

How to learn a consistent tree with low expected cost? Find min cost tree.

Another approach: replace gain by

- Tan and Schlimmer (1990)
 $Gain^2(S,A)/Cost(A)$
- Nunez (1988) [w in $[0,1]$: importance]
 $(2^{Gain(S,A)}-1)/(Cost(A)+1)^w$

Unknown Attribute Values

Some examples missing values of A ?

Use training example anyway, sort it

- If node n tests A , assign most common value of A among other examples sorted to node n
- assign most common value of A among other examples with same target value
- assign probability p_i to each possible value v_i of A (perhaps as above)
 - assign fraction p_i of example to each descendant in tree
- Classify new examples in same fashion

Sources

- ML: Chapter 3
- i2ML: Chapter 9
- Slides by Ethem Alpaydin
- Slides by Tom Mitchell as provided by Michael Littman