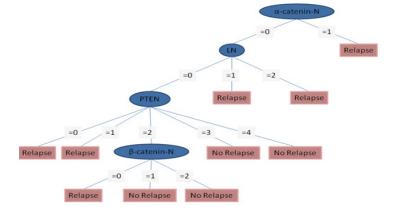
Cmput 466 / 551



Decision Trees

Covering: HTF Ch 9 (esp 9.2) + Induction of Decision Trees (Quinlan)

R Greiner
Dept of Computing Science
University of Alberta



Learning Decision Trees

- Def'n: Decision Trees
- Algorithm for Learning Decision Trees
- Overfitting
 - Def'n, χ², MDL, PostPruning
- Topics:
 - k-ary attribute values
 - Real attribute values
 - Other splitting criteria
 - Attribute Cost
 - Missing Values

- ...



Lessons

- Understanding Decision Tree
 - Def'n, Alg, Variants, ...
- Overview other topics:
 - Entropy
 - (Another) Example of Overfitting
 - χ^2 test
 - Minimal Description Language
 - Missing data ...
 - Importance of understanding data



DecisionTree Hypothesis Space

- Each internal nodes labeled with some feature x_i
- Each arc (from x_j) labeled with results of test x_j

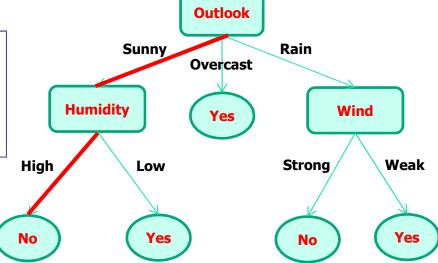
Each leaf node specifies class h(x)

Instance:

Outlook = Sunny
Temperature = Hot
Humidity = High
Wind = Strong

classified as "No"

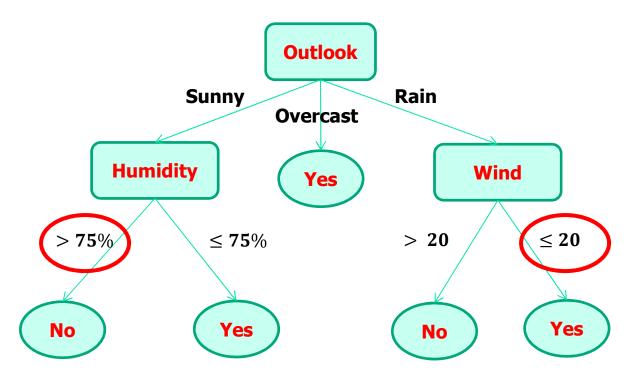
- (Temperature, Wind: irrelevant)
- Easy to use in Classification
 - Just need to answer short series of questions...





Continuous Features

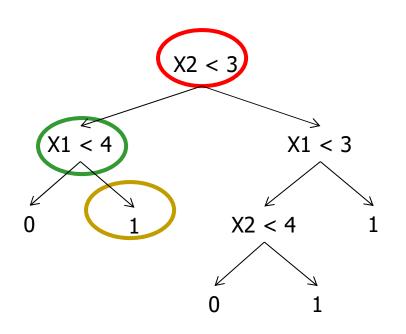
If feature is continuous:
 internal nodes may test value against threshold

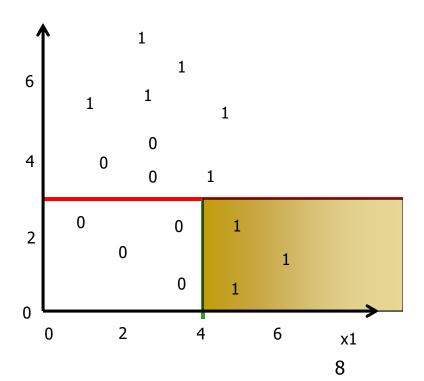




DecisionTree Decision Boundaries

 Decision trees divide feature space into axis-parallel rectangles,
 labeling each rectangle with one class

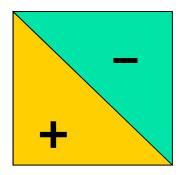


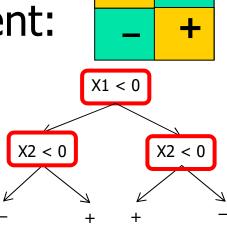


Properties

- İnstances represented by Attribute-Value pairs
 - "Bar = Yes", "Size = Large", "Type = French", "Temp = 82.6", ...
 - (Boolean, Discrete, Nominal, Continuous)
- Some concepts are easy to represent: using axis-parallel splits

Others are difficult :





4

Can Represent Any Boolean Fn

- V, &, ⊕, MofN
 (A ∨ B) & (C ∨ ¬D ∨ E)
- Arbitrary boolean functions (DNF)
 - ... but may require exponentially many nodes ...
 - "disjunctive concepts"

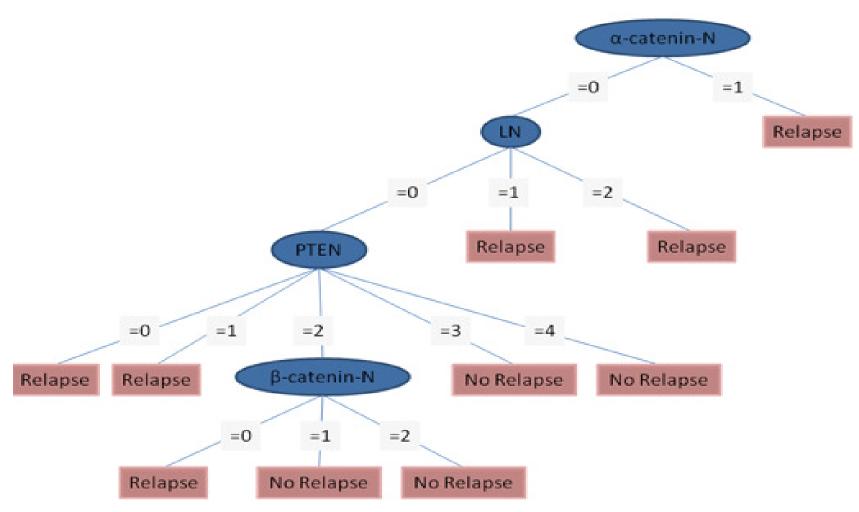
Variable-Size Hypothesis Space

- Can "grow" hypothesis space by increasing number of nodes
- depth 1 ("decision stump"): represent any boolean function of one feature
- **depth 2**: Any boolean function of two features; $(x_1 \lor x_2) \& (\neg x_1 \lor \neg x_3)$

...

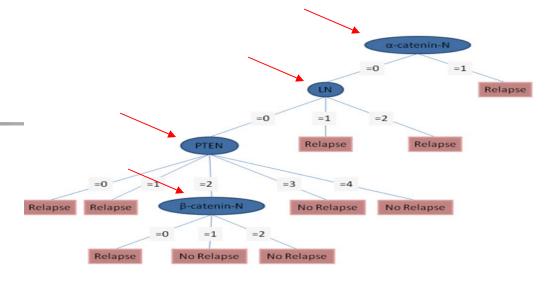


Learned Decision Tree





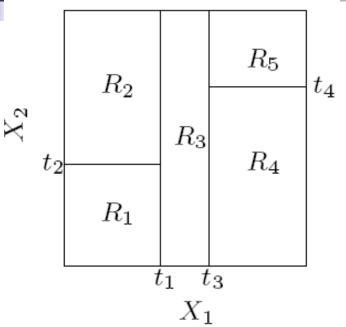
Meaning

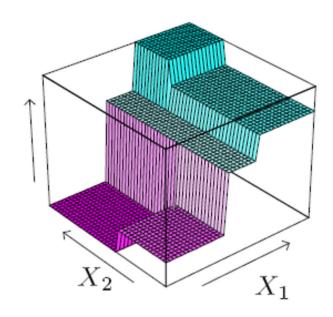


- Concentration of
 - α -catenin in nucleus is very important:
 - If >0, probably relapse
- If =0, then #lymph_nodes is important:
 - If >0, probably relapse
- If =0, then concentration of pten is important:
 - If <2, probably relapse
 - If >2, probably NO relapse
- If =2, then concentration of β -catenin in nucleus is important:
 - If =0, probably relapse
 - If >0, probably NO relapse



Regression (Constant) Tree





- Represent each region as CONSTANT
- Alternatively: could have LINEAR FUNCTION at each leaf node
 - ... over OTHER variables
 - ... or NON-linear ...

Using Decision Trees

- Our focus:
 - Label is discrete
- Advantages of Decision Trees?
 - Fast to learn
 - Fast(er) to evaluate
 - Uses only subset of features
 - ... relevant ones ?
 - Intuitive... easy to EXPLAIN ... kinda
- Disadvantages ...
 - Typically NOT most accurate
 - Some natural concepts are difficult
- Uses:
 - Credit risk analysis
 - Modeling calendar scheduling preferences
 - Equipment or medical diagnosis
 - ...



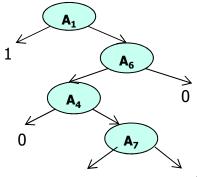
Learning Decision Trees

- Def'n: Decision Trees
- Algorithm for Learning Decision Trees
 - Greedy
 - Splitting Criterion: Entropy
 - Inductive Bias (Occam's Razor)
- Overfitting
- Topics

4

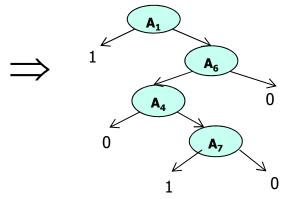
Learning / Using Decision Tree

- Classify: DecisionTree × Instance → ClassLabel
 - Given



- classify instance $[\mathring{A}_1 = -, A_4 = +, A_6 = +]$ as "0"
- Learn: Data → DecisionTree

A ₁	A ₂	A ₃	 A ₇	A ₈	Label
+	-	+	 +	+	1
_	+	_	 +	+	0
_	+	_	 -	+	0
_	_	_	 +	_	1
:	:	:	:	:	:
_	_	+	 +	_	0





Training Examples

Day	Outlook	Temp.	Humidity	Wind	P.Tennis
d ₁	Sunny	Hot	High	Weak	No
d ₂	Sunny	Hot	High	Strong	No
dз	Overcast	Hot	High	Weak	Yes
d ₄	Rain	Mild	High	Weak	Yes
d ₅	Rain	Cool	Normal	Weak	Yes
d ₆	Rain	Cool	Normal	Strong	No
d ₇	Overcast	Cool	Normal	Strong	Yes
de	Sunny	Mild	High	Weak	No
do	Sunny	Cool	Normal	Weak	Yes
d ₁₀	Rain	Mild	Normal	Weak	Yes
d ₁₁	Sunny	Mild	Normal	Strong	Yes
d ₁₂	Overcast	Mild	High	Strong	Yes
d ₁₃	Overcast	Hot	Normal	Weak	Yes
d ₁₄	Rain	Mild	High	Strong	No

- 4 discrete-valued attributes
- "Yes/No" classification

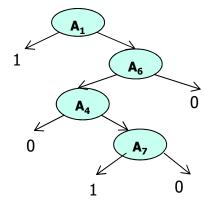
Want: Decision Tree

 DT_{PT} (Out, Temp, Humid, Wind) $\in \{ Yes, No \}$



Learn: Data → DecisionTree

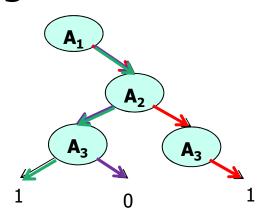
A ₁	A ₂	A ₃	 A ₇	A ₈	Label	
+	-	+	 +	+	1	
-	+	-	 +	+	0	\Longrightarrow
_	+	_	 -	+	0	
-	-	_	 +	-	1	
:	:	;	:	:	:	
_	_	+	 +	-	0	



Option 1: Just store training data

A ₁	A ₂	A ₃	Label
+	+	+	1
+	_	+	0
+	_	_	1

But ...





Learn ?Any? Decision Tree

- Just produce "path" for each example
 - May produce large tree
 - Any generalization? (what of other instances?)

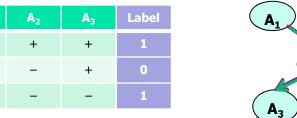
$$[-, +, +]$$
?

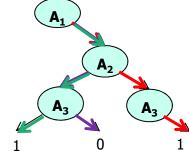
Noise in data

mis-recorded as

$$[[+,+,-], 1]$$

 $[[+,-,-], 0]$





Intuition:

Want SMALL tree

- ... to capture "regularities" in data ...
- ... easier to understand, faster to execute, ...

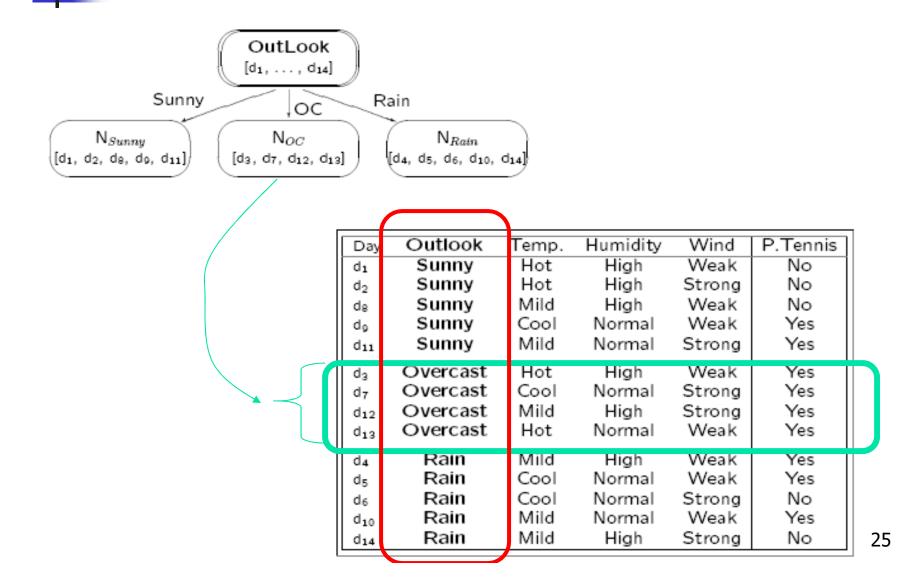


Grow a tree... first Split?

??

Day	Outlook	Temp.	lumidity	Wind	P.Tennis
d ₁	Sunny	Hot	High	Weak	No
d ₂	Sunny	Hot	High	Strong	No
de	Sunny	Mild	High	Weak	No
d٥	Sunny	Cool	Normal	Weak	Yes
d ₁₁	Sunny	Mild	Normal	Strong	Yes
d ₃	Overcast	Hot	High	Weak	Yes
d ₇	Overcast	Cool	Normal	Strong	Yes
d ₁₂	Overcast	Mild	High	Strong	Yes
d ₁₃	Overcast	Hot	Normal	Weak	Yes
d4	Rain	Mild	High	Weak	Yes
d ₅	Rain	Cool	Normal	Weak	Yes
d ₆	Rain	Cool	Normal	Strong	No
d ₁₀	Rain	Mild	Normal	Weak	Yes
d ₁₄	Rain	Mild	High	Strong	No

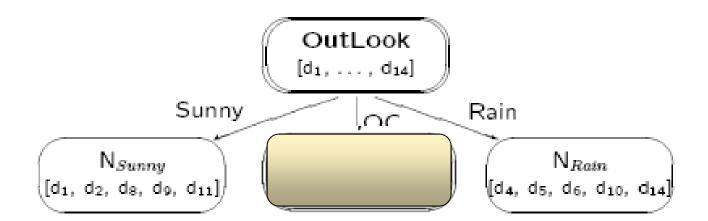
First Split: Outlook





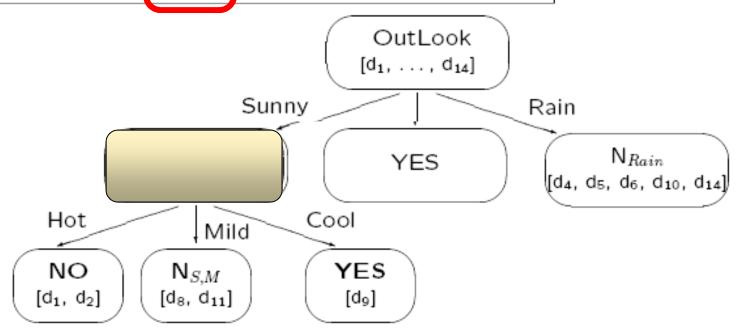
Onto N_{OC} ...

Day	Outlook	Temp.	Humidity	Wind	P. Tennis
dз	Overcast	Hot	High	Weak	Yes
d ₇	Overcast	Cool	Normal	Strong	Yes
d ₁₂	Overcast	Mild	High	Strong	Yes
d ₁₃	Overcast	Hot	Normal	Weak	Yes



What about N_{Sunny}?

	Day	Outlook	Temp.	Humidity	Wind	P.Tennis
	d ₁	Sunny	Hot	High	Weak	No
	d_2	Sunny	Hot	High	Strong	No
	d ₈	Sunny	Mild	High	Weak	No
	d ₁₁	Sunny	Mild	Normal	Strong	Yes
	d ₉	Sunny	Cool	Normal	Weak	Yes
1						





(Simplified) Algorithm for Learning Decision Tree

```
\label{eq:final_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_cont
```

- Many fields independently discovered this learning alg...
- Issues
 - > 2 labels> 2 feature values
 - ... continuous values ??
 - ? ChooseBestAttribute
 - to purity? ... used all attributes?pruning
 - oblique splits

Interesting issues...

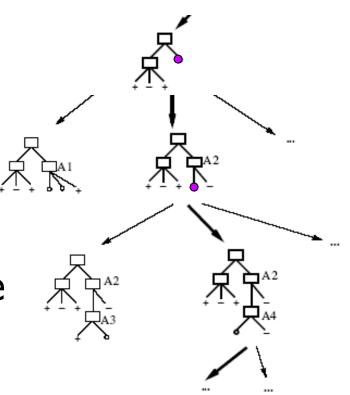
Alg for Learning Decision Trees

```
function DECISION-TREE-LEARNING(examples, attributes, default) returns a decision tree
  inputs: examples, set of examples
           attributes, set of attributes
           default, default value for the goal predicate
  if examples is empty then return default
  else if all examples have the same classification then return the classification
  else if attributes is empty then return MAJORITY-VALUE(examples)
  else
      best \leftarrow CHOOSE-ATTRIBUTE(attributes, examples)
      tree \leftarrow a new decision tree with root test best
      for each value v; of best do
          examples_i \leftarrow \{elements \text{ of } examples \text{ with } best = v_i\}
          subtree \leftarrow DECISION-TREE-LEARNING(examples<sub>i</sub>, attributes - best,
                                                    MAJORITY-VALUE(examples))
          add a branch to tree with label v; and subtree subtree
      end
    return tree
```



Search for Good Decision Tree

- Local search
 - expanding one leaf at-a-time
 - no backtracking
- Trivial to find a tree that perfectly "fits" training data*
- but... this is NOT necessarily best tree
- Prefer small tree
 - NP-hard to find smallest tree that fits data





Issues in Design of Decision Tree Learner



What attribute to split on?

- Avoid Overfitting
 - When to stop?
 - Any post-processing?
 - ... should tree by "pruned"?
- $\begin{aligned} \textit{GrowTree}(\ S: labeled_dataset) \\ & \text{if } (\ y ==0\) \ \text{for all } [\textbf{x},y] \in S: \ \text{return new_leaf}(\ 0\) \\ & \text{if } (\ y ==1\) \ \text{for all } [\textbf{x},y] \in S: \ \text{return new_leaf}(\ 1\) \end{aligned}$ $\begin{aligned} & /^* \ \text{else } ^*/ \\ & x_j = \textit{ChooseBestAttribute}(\ S\) \\ & S_0 = \{\ [\textbf{x},y] \in S \mid x_j ==0\ \} \\ & S_1 = \{\ [\textbf{x},y] \in S \mid x_j ==1\ \} \\ & \text{return new_node}(\ x_j, \textit{GrowTree}(\ S_0\), \textit{GrowTree}(\ S_1\)\) \end{aligned}$

- How to evaluate classifier (decision tree) ?
 - ... learner?



Choosing Best Splitting Test

How to choose best feature for split?



- After Gender split, still some uncertainty
 After Smoke split, no more Uncertainty
 - ⇒ NO MORE QUESTIONS!

(Here, Smoke is a great predictor for Cancer)

Want a "measure" that prefers
 Smoke over Gender



Statistics ...

If split on x_i , produce 2 children:

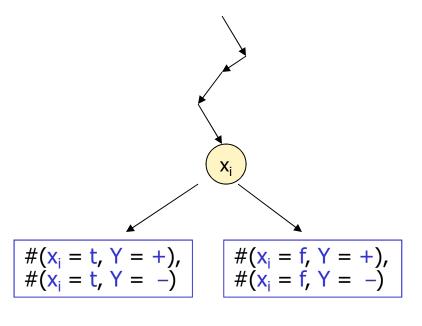
#(x_i = t) follow TRUE branch

$$\Rightarrow$$
 data: [$\#(x_i = t, Y = +), \#(x_i = t, Y = -)$]

• $\#(x_i = f)$ follow FALSE branch

$$\Rightarrow$$
 data: [$\#(x_i = f, Y = +), \#(x_i = f, Y = -)$]

Day	Outlook	Temp.	Humidity	Wind	P.Tennis
d ₁	Sunny	Hot	High	Weak	No
d ₂	Sunny	Hot	High	Strong	No
d3	Overcast	Hot	High	Weak	Yes
d ₄	Rain	Mild	High	Weak	Yes
d ₅	Rain	Cool	Normal	Weak	Yes
d ₆	Rain	Cool	Normal	Strong	No
d ₇	Overcast	Cool	Normal	Strong	Yes
de	Sunny	Mild	High	Weak	No
do	Sunny	Cool	Normal	Weak	Yes
d ₁₀	Rain	Mild	Normal	Weak	Yes
d ₁₁	Sunny	Mild	Normal	Strong	Yes
d ₁₂	Overcast	Mild	High	Strong	Yes
d ₁₃	Overcast	Hot	Normal	Weak	Yes
d ₁₄	Rain	Mild	High	Strong	No





Desired Properties

Score for split of data D, on x_i, related to

$$S \begin{pmatrix} \#(x_i = t, Y = +) \\ \#(x_i = t, Y = -) \end{pmatrix}$$
 $S \begin{pmatrix} \#(x_i = f, Y = +) \\ \#(x_i = f, Y = -) \end{pmatrix}$

- Score S(.) should be...
 - BEST for [+0, -200]
 - WORST for [+100, -100]
 - "symmetric"

Same for
$$[+19, -5]$$
 and $[+5, -19]$

Deals with any number of values

l

Play 20 Questions



- I'm thinking of integer ∈ {1, ..., 100 }
- Questions
 - Is it 22?
 - More than 90?
 - More than 50?
- Why?
- Q(r) = E[# of additional questions wrt set of size r]

$$= 22? \frac{1}{100} \times Q(1) + \frac{99}{100} \times Q(99)$$

■
$$\geq 90$$
? $\frac{11}{100} \times Q(11) + \frac{89}{100} \times Q(89)$

■
$$\geq 50$$
? $\frac{50}{100} \times Q(50) + \frac{50}{100} \times Q(50)$

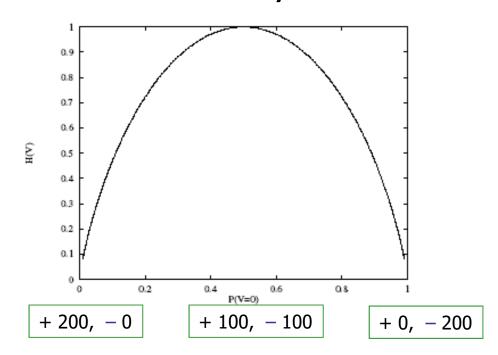
Want this to be **small** ...

Desired Measure: Entropy

Entropy of V = [p(V = 1), p(V = 0)]:

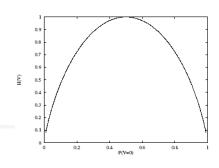
$$H(V) = -\sum_{v_i} P(V = v_i) \log_2 P(V = v_i)$$

- = # of bits needed to obtain full info
- ... average surprise of result of one "trial" of V
- Entropy ≈ measure of uncertainty





Examples of Entropy



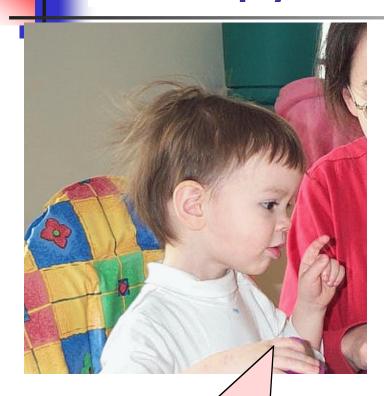
- Fair coin:
 - $H\left(\frac{1}{2}, \frac{1}{2}\right) = -\frac{1}{2}\log\frac{1}{2} \frac{1}{2}\log\frac{1}{2} = 1$ bit
 - ie, need 1 bit to convey the outcome of coin flip
- Biased coin:

$$H\left(\frac{1}{100}, \frac{99}{100}\right) = -\frac{1}{100}\log\frac{1}{100} - \frac{99}{100}\log\frac{99}{100} \approx 0.08 \text{ bit}$$

As P(heads) → 1, info of actual outcome → 0
 H(0, 1) = H(1, 0) = 0 bits
 ie, no uncertainty left in source

$$(0 \times \log_2(0) = 0)$$

Entropy in a Nut-shell





Low Entropy

High Entropy

...the values (locations of soup) sampled entirely from within the soup bowl ...the values

(locations of soup) unpredictable... almost uniformly sampled throughout dining room

Prefer Low Entropy Leaves

- Use decision tree h(.) to classify (unlabeled) test example x
 - ... Follow path down to leaf r
 - ... What classification?
- Consider training examples that reached r:
 - If all have same class C +200, -0 \Rightarrow label X as C (ie, entropy is 0)
 - If $\frac{1}{2}$ are +; $\frac{1}{2}$ are $\frac{100}{2}$ +100, -100 $\frac{1}{2}$ | (ie, entropy is 1)
- On reaching leaf r with entropy H_r, uncertainty w/label is H_r

(ie, need H_r more bits to decide on class)

⇒ prefer leaf with LOW entropy



Entropy of Set of Examples

- Don't have exact probabilities...
 - ... but training data provides estimates of probabilities:
- Given training set with $\begin{cases} p & \text{positive} \\ n & \text{negative} \end{cases}$ examples:

$$H\left(\frac{p}{p+n}, \frac{n}{p+n}\right) = -\frac{p}{p+n}\log_2\frac{p}{p+n} - \frac{n}{p+n}\log_2\frac{n}{p+n}$$

Eg: wrt 12 instances, D:

$$p = n = 6 \implies H(\frac{1}{2}, \frac{1}{2}) = 1 \text{ bit}$$

... so need 1 bit of info to classify example randomly picked from *D*

Remaining Uncertainty

Uncert(D, A) = remaining expected entropy after splitting on A

$$\equiv \sum_{v_i \in Values(A)} \frac{|S_{v_i}|}{|S|} \operatorname{Entropy}(S_{v_i})$$

$$\equiv \sum_{i=1}^{v} \frac{p_i^{(A)} + n_i^{(A)}}{p+n} \ H\left(\frac{p_i^{(A)}}{p_i^{(A)} + n_i^{(A)}}, \frac{n_i^{(A)}}{p_i^{(A)} + n_i^{(A)}}\right)$$

$$S: [9+, 5-]$$



$$[9+, 2-]$$

U = 0.683

$$[0+, 3-]$$

U = 0.0

Uncert(D, Rain)

$$= \frac{11}{14} 0.683 + \frac{3}{14} 0.0$$

$$= 0.536$$

$$[6+, 2-]$$

U = 0.811

$$[3+, 3-]$$

U = 1.00

Uncert(D, Wind)

$$= \frac{8}{14} 0.811 + \frac{6}{14} 1.00$$
$$= 0.892$$

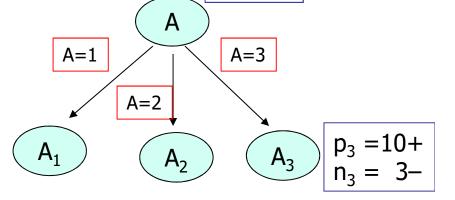
Entropy wrt Feature

- Assume [p,n] reach node
- Feature A splits into A₁, ..., A_v
 - A_i has { p_i(A) positive, n_i(A) negative }
- Entropy of each is:

$$H\left(\frac{p_i^{(A)}}{p_i^{(A)} + n_i^{(A)}}, \frac{n_i^{(A)}}{p_i^{(A)} + n_i^{(A)}}\right)$$

So for A₂: $H(\frac{28}{40}, \frac{12}{40})$

$$p_1 = 22 + 1 = 25 -$$



p = 60 +

n = 40 -

$$p_2 = 28 + n_2 = 12 - n_2$$

Uncert(A) =
$$\sum_{i=1}^{v} \frac{p_i^{(A)} + n_i^{(A)}}{p + n} H\left(\frac{p_i^{(A)}}{p_i^{(A)} + n_i^{(A)}}, \frac{n_i^{(A)}}{p_i^{(A)} + n_i^{(A)}}\right)$$



Minimize Remaining Uncertainty

- Greedy: Split on attribute that leaves least entropy wrt class
 ... over training examples that reach there
- Assume A divides training set E into E₁, ..., E_v
- E_i has $\{p_i^{(A)} \text{ positive, } n_i^{(A)} \text{ negative } \}$ examples
- Entropy of each E_i is $H\left(\frac{p_i^{(A)}}{p_i^{(A)} + n_i^{(A)}}, \frac{n_i^{(A)}}{p_i^{(A)} + n_i^{(A)}}\right)$
- Uncert(A) = expected information content
 ⇒ weighted contribution of each E_i

$$Uncert(A) = \sum_{i=1}^{v} \frac{p_i^{(A)} + n_i^{(A)}}{p + n} H\left(\frac{p_i^{(A)}}{p_i^{(A)} + n_i^{(A)}}, \frac{n_i^{(A)}}{p_i^{(A)} + n_i^{(A)}}\right)$$

Split on feature A* = argmin_A { Uncert(A) }



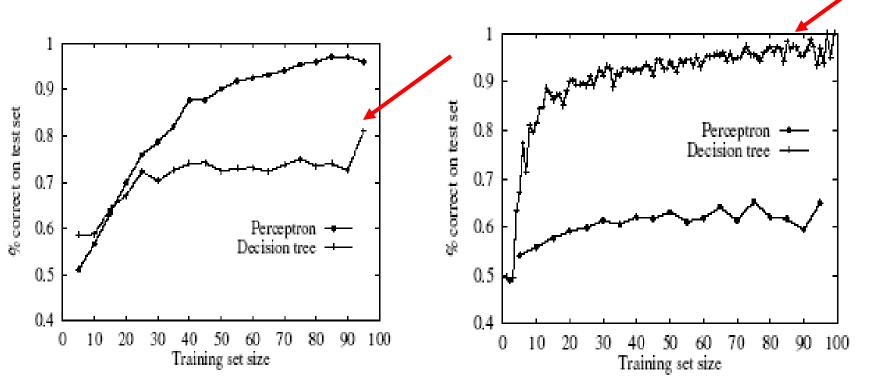


Occam's Razor

- Q: 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 *more likely to be* coincidence
- Argument opposed:
 - many ways to define small sets of hyps
 Eg, all trees with prime number of nodes
 whose attributes all begin with "Z"
 - What's so special about small sets based on size of hypothesis??



Perceptron vs Decision Tree



Majority Function (11 inputs)

WillWait predicate



Learning Decision Trees

- Defn: Decision Tree
- Algorithm for Learning Decision Trees
- Overfitting
 - Example, Def'n, Approaches
 - \mathbf{x}^2 -test
 - MDL
 - PostPruning
 - Topics

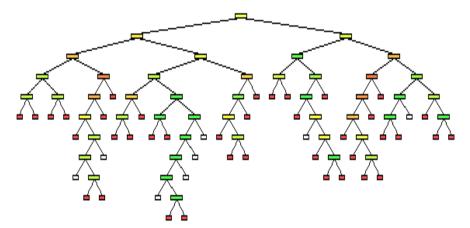
Example of Overfitting

- 25% have butterfly-itis
 - Random
- $\frac{1}{2}$ of patients have $F_1 = 1$
 - Eg: "odd birthday"
- $\frac{1}{2}$ of patients have $F_2 = 1$
 - Eg: "even SSN"
- ... for 10 features
- Note feature F_i is INDEPENDENT of class!!
- Decision Tree results
 - over 1000 patients (using these silly features) ...

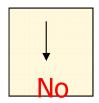




Standard decision tree learner:



Optimal decision tree:



Error Rate:

■ Train data: 0%

New data: (37.5%)

Error Rate:

■ Train data: 25%

New data:

25%

Example of Overfitting

- Spse 10 binary attributes (uniform), but class is random: $\begin{cases} N \text{ w/prob } p = 0.75 \\ Y \text{ w/prob } 1-p = 0.25 \end{cases}$
- (naïve) C4.5 builds nonsensical tree w/ 119 nodes!
 - ⇒ Should be SINGLE NODE!
- Error rate (hold-out): 37.5%
 - \Rightarrow Could be only 25% (just say "No")
- Why? Tree assigns leaf "N" w/prob 1-p, "Y" w/prob p
 - Tree sends instances to arbitrary leaves
 - Mistake if
 - Y-instance reaches N-leaf: p x (1-p)
 - N-instance reaches Y-leaf: (1-p) x p
 - Total prob of mistake = $2 \times p \times (1-p) = 0.375$
- Overfitting happens for EVERY learner ... not just DecTree !!



How to Avoid Overfitting (Decision Trees)

- When to act
 - Use more stringent STOPPING criterion while growing tree
 - ... only allow statistically significant splits ...
 - Grow full tree, then post-prune
- To evaluate tree, measure performance over ...
 - training data
 - separate validation data set
- How to represent classifier?
 - as Decision Tree
 - as Rule Set

Mod's?

```
\label{eq:growTree} \textit{GrowTree}(\ S: labeled\_dataset) \\ \textit{if } (\ y == 0\ ) \textit{ for all } [\textbf{x},y] \in S: \textit{ return new\_leaf}(\ 0\ ) \\ \textit{if } (\ y == 1\ ) \textit{ for all } [\textbf{x},y] \in S: \textit{ return new\_leaf}(\ 1\ ) \\ /* \textit{ else } */ \\ x_j = \textit{ChooseBestAttribute}(\ S\ ) \\ \textit{if } x_j \neq \textit{NULL} \\ S_0 = \{\ [\textbf{x},y] \in S \mid x_j == 0\ \} \\ S_1 = \{\ [\textbf{x},y] \in S \mid x_j == 1\ \} \\ \textit{tree} = \textit{new\_node}(\ x_j, \textit{GrowTree}(\ S_0\ ), \textit{GrowTree}(\ S_1\ )\ ) \\ \textit{else} \\ \textit{tree} = \textit{new\_leaf}(\ \textit{majority}(\ S\ )\ ) \\ \textit{return tree} \\ \end{aligned}
```

Currently

```
ChooseBestAttribute( S : labeled_dataset)
best = argmin<sub>x</sub> Uncert( S, x )
if NotGoodEnuf( S, x) return NULL
return best
```

- Maybe ...
 - Decide if this best features is good enough?





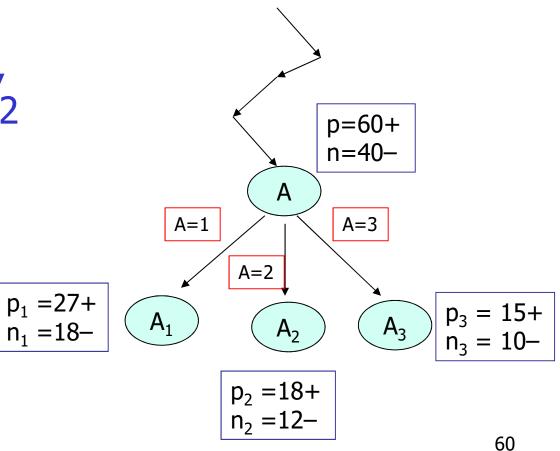
Avoid Overfitting #1

- (StopEARLY, Training-Data, DecTree)
- Add more stringent STOPPING criterion while growing tree
 - At leaf n_r (w/ instances S_r)
 spse proposed split is based on feature A
- A. Use χ^2 test, on data S_r
 - Apply statistical test to compare
 - T₀: leaf at r (majority label) vs
 - T_A: split using feature A
 - Is error of T_A statistically better than T₀?
- B. MDL: minimize size(tree) + size(misclassifications(tree))



Test for Significance

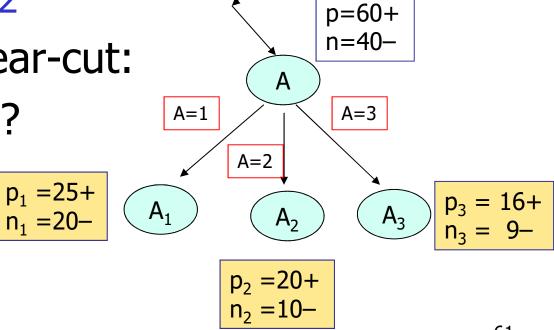
- Spse A is irrelevant
 - $[p_i, n_i] \propto [p, n]$
 - So if [p,n] = 3:2, then $[p_i, n_i] = 3:2$





Test for Significance

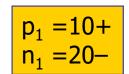
- Spse A is irrelevant
 - $[p_i, n_i] \propto [p, n]$
 - So if [p,n] = 3:2, then $[p_i, n_i] = 3:2$
- Not always so clear-cut:
- Is this significant?

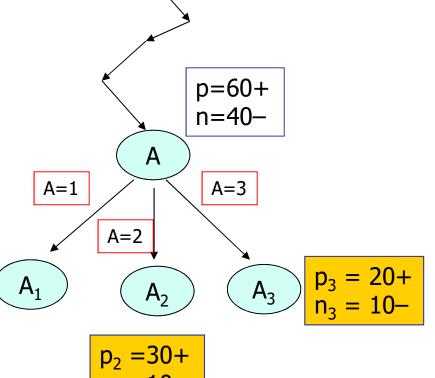




Test for Significance

- Spse A is irrelevant
 - $[p_i, n_i] \propto [p, n]$
 - So if [p,n] = 3:2, then $[p_i, n_i] = 3:2$
- Not always so clear-cut:
- Is this significant?
- Or this??





χ^2 Te

χ² Test for Significance

Null hypothesis H₀:

Attribute A is irrelevant in context of r

Ie, dist'n of class labels at node $n_r \equiv$ dist'n after splitting on A

Observe some difference between these dist'ns.

What is prob (under H_0) of observing this difference, given $m = |S_r|$ iid samples?

- Defn: ${p \choose n}$ of S_r are ${positive \choose negative}$
- After splitting on A, get k subsets
 - wrt A = i : p_i positives, n_i negatives
- If H₀ (A irrelevant), would have

•
$$\widetilde{p_i} = p \times \frac{p_i + n_i}{p + n}$$
 positives

•
$$\widetilde{\mathbf{n_i}} = \mathbf{n} \times \frac{\mathbf{p_i} + \mathbf{n_i}}{\mathbf{p+n}}$$
 negatives



χ^2 Test – con't

 $\sum \frac{(expected - observed)^2}{expected}$

If H₀ (A irrelevant), would have

$$\widetilde{p_i} = p \times \frac{p_i + n_i}{p + n}$$
 positives $\widetilde{n_i} = n \times \frac{p_i + n_i}{p + n}$ negatives

Measure of deviation:

$$D(A) = \sum_{i=1}^{k} \frac{(p_i - \widehat{p_i})^2}{\widehat{p_i}} + \frac{(n_i - \widehat{n_i})^2}{\widehat{n_i}}$$

- Under H_0 , $D \sim \chi^2_{(k-1)(2-1)}$
 - k = size of A's domain; 2 = #classes
- If A irrelevant (H_0) , would you expect to see D?
- Extend to new internal node if A is not irrelevant... iff $D > T_{\alpha,d}$
 - $T_{\alpha,d}$ is threshold:
 - α = probability of making mistake (rejecting null hyp, when A is irrelevant)
 - \bullet d = degree of freedom

χ² Table

For χ^2 test, with $\alpha = 0.05$

various "degrees of freedom"

DOF	$T_{\alpha,DOF}$	DOF	$T_{\alpha,DOF}$
1	3.841	16	26.296
2	5.991	17	27.587
3	7.815	18	28.869
4	9.488	19	30.144
5	11.070	20	31.410
6	12.592	21	32.671
7	14.067	22	33.924
8	15.507	23	35.172
9	16.919	24	36.415
10	18.307	25	37.652
11	19.675	26	38.885
12	21.026	27	40.113
13	22.362	28	41.337
14	23.685	29	42.557
15	24.996	30	43.773

Minimum Description Length

A wants to transmit to B classification function $c(\cdot)$

- simplified to:
 - At and B agree on instances $\langle x_1, ..., x_M \rangle$
 - What should A send, to allow B to determine M bits:

$$\langle c(x_1), ..., c(x_M) \rangle$$

- Option#1: A can send M "bits"
- Option#2: A sends "perfect" decision tree d s.t. c(x_i) = d(x_i) for each x_i
- Option#3: A sends "imperfect" decision tree d'
 + set of indices of K exceptions F = { x_{i1} , ..., x_{iK} }

$$C(x_i) = \begin{cases} \neg d(x_i) & \text{if } x_i \in F \\ d(x_i) & \text{otherwise} \end{cases}$$

So... Increase tree-size
 IFF (significant) reduction in #exceptions

Mod's?

```
\begin{aligned} &\textit{GrowTree}(\ S: labeled\_dataset) \\ &\text{if } (\ y == 0\ ) \ \text{for all } [\textbf{x},y] \in S: \ \text{return new\_leaf}(\ 0\ ) \\ &\text{if } (\ y == 1\ ) \ \text{for all } [\textbf{x},y] \in S: \ \text{return new\_leaf}(\ 1\ ) \\ &/* \ \text{else } */ \\ &x_j = \textit{ChooseBestAttribute}(\ S\ ) \\ &\text{if } x_j \neq \text{NULL} \\ &S_0 = \{\ [\textbf{x},y] \in S \mid x_j == 0\ \} \\ &S_1 = \{\ [\textbf{x},y] \in S \mid x_j == 1\ \} \\ &\text{tree} = \text{new\_node}(\ x_j, \ \textit{GrowTree}(\ S_0\ ), \ \textit{GrowTree}(\ S_1\ )\ ) \\ &\text{else} \\ &\text{tree} = \text{new\_leaf}(\ \text{majority}(\ S\ )\ ) \\ &\text{return tree} \end{aligned}
```

```
ChooseBestAttribute( S : labeled_dataset)
best = argmin<sub>x</sub> Uncert( S, x )
if NotGoodEnuf( S, x) return NULL
return best
```

$NotGoodEnuf_{MDL}(S, x)$:

- I_{Tree} = bits to encode subtree rooted in x
- I_{Stop} = bits to encode the exceptions (wrt default)
- Return (I_{Tree} < I_{Stop})



Avoid overfitting#2: PostPruning

Grow tree, to "purity", then PRUNE it back!

```
Build "complete" decision tree h, from train For each penultimate node: n_i Let h_i be tree formed by "collapsing" subtree under n_i, into single node If h_i better than h Set h \leftarrow h_i
```

- How to decide if h_i better than h?
- 1. Test on Hold-Out data?
 - 3 sets: training, VALIDATION, testing
 Problematic if small total # of samples
- 2. Pessimistic Pruning
 ... re-use training samples ...

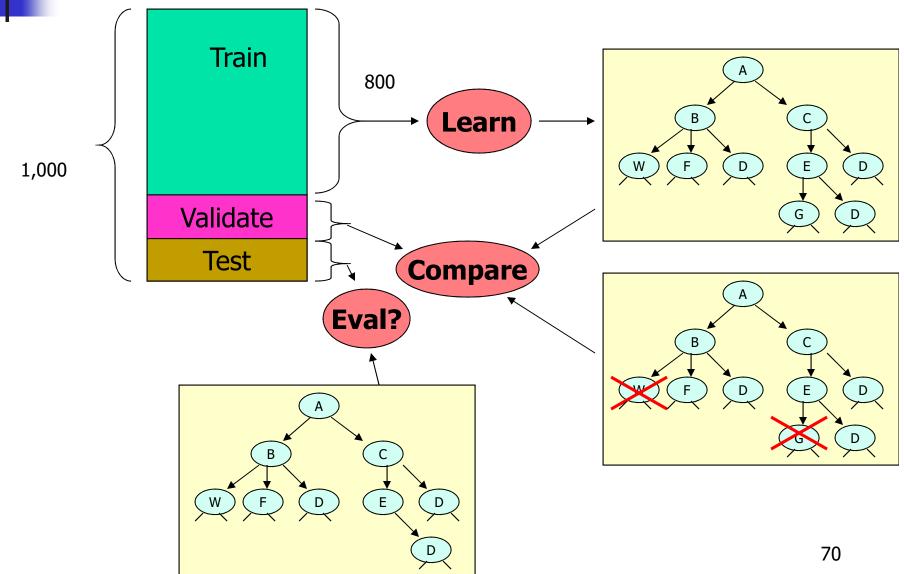
Train

Validate

Test

4

Using Validation Set



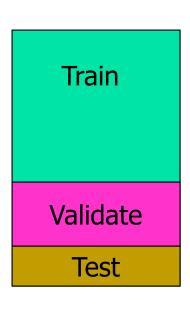


Avoid Overfitting#2.1 "Reduced-Error Pruning"

Split data into training and validation set

Alg: Do until further pruning is harmful:

- 1. Consider pruning each leaf node* Evaluate accuracy on *validation* set
- 2. Greedily remove the node that most improves accuracy on *validation* set



- Produces possibly smaller subtree... "regularization"
- What if data is limited?



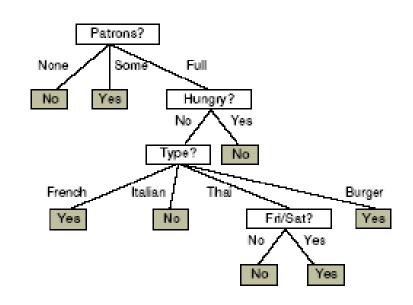
Avoid Overfitting #3 Using Rule Post-Pruning

- Grow decision tree.
 Fit data as well as possible.
 Allow overfitting.
- Convert tree to equivalent set of rules:
 - One rule for each path from root to leaf.
- 3. Prune each rule independently of others
 - ie, delete preconditions that improve its accuracy
- Sort final rules into best sequence-for-use, depending on accuracy
- Use ordered sequence for classification



Converting Trees to Rules

- Every decision tree corresponds to set of rules:
 - IF (Patrons = None)THEN WillWait = No
 - IF (Patrons = Full)
 & (Hungry = No)
 &(Type = French)
 THEN WillWait = Yes



 Why? (Small) RuleSet MORE expressive small DecTree ≈ small RuleSet (DecTree is subclass of ORTHOGONAL DNF)



Learning Decision Trees

- Def'n: Decision Trees
- Algorithm for Learning Decision Trees
- Overfitting
- Topics:
 - k-ary attribute values
 - Real attribute values
 - Other splitting criteria
 - Attribute Cost
 - Missing Values

<u>Skip</u>

...

Attributes with Many Values

- Problem: Uncert prefers attribute with many values, as Entropy ≈ ln(k)
 - Problem: what about using

```
Date = Jun 3 1996 {as LOTS of Date's \Rightarrow k is large }
PatientId = 217331 {as LOTS of PatientId's \Rightarrow k is large }
```

 If using a multi-way split: perhaps use *EntropyRatio* instead of Entropy

EntropyRatio(S, A) =
$$\frac{\text{Entropy}(S, A)}{\text{SplitInfo}(S, A)}$$
SplitInfo(S, A) =
$$-\sum_{i=1}^{K} \frac{|S_i|}{|S|} \log \frac{|S_i|}{|S|}$$

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

- Other Approaches:
 - Test one value versus all of the others?
 - Group values into two disjoint subsets?

Finding Split for Real-Valued Features

- Best threshold θ_j for attribute X_j
 - Sort training instances by X_{i,i}

y i	A	A	В	A	В	В	A	В	В
X _{i,j}	0.2	0.4	0.7	1.1	1.3	1.7	1.9	2.4	2.9

• For each possible threshold $\theta \in \mathfrak{R}$,

	< θ	$\geq \theta$
Α	$n_{A,I}(\theta)$	$n_{A,r}(\theta)$
В	$n_{B,I}(\theta)$	$n_{B,r}(\theta)$
(both)	$n_{l}(\theta)$	$n_r(\theta)$

Compute mutual information MI(θ) wrt label {A, B}:

$$\frac{n_{l}(\theta)}{n_{l}(\theta)+n_{r}(\theta)} \left[\frac{n_{A,l}(\theta)}{n_{A,l}(\theta)+n_{B,l}(\theta)} \ln \frac{n_{A,l}(\theta)}{n_{A,l}(\theta)+n_{B,l}(\theta)} + \frac{n_{B,l}(\theta)}{n_{A,l}(\theta)+n_{B,l}(\theta)} \ln \frac{n_{B,l}(\theta)}{n_{A,l}(\theta)+n_{B,l}(\theta)} \right] + \frac{n_{B,l}(\theta)}{n_{A,l}(\theta)+n_{B,l}(\theta)} \ln \frac{n_{B,l}(\theta)}{n_{A,l}(\theta)+n_{B,l}(\theta)} + \frac{n_{B,l}(\theta)}{n_{A,l}(\theta)+n_{B,l}(\theta)} \ln \frac{n_{B,l}(\theta)}{n_{A,l}(\theta)+n_{B,l}(\theta)} \right]$$

• Use $\theta^* = \operatorname{argmax}_{\theta} \{ MI(\theta) \}$



Finding Split for

For quality of resulting DecTree: find θ_i within each CrossValidation fold

Real-Valued Features

- Best threshold θ_i for attribute X_i
 - Sort training instances by X_{i,i}

y i	A	A	В	A	В	В	A	В	В
X _{i,j}	0.2	0.4	0.7	1.1	1.3	1.7	1.9	2.4	2.9

• Eq. for $\theta = 1.2$:

< 1.2	≥ 1.2
$n_{A,I}=3$	$n_{A,r}=1$
$n_{B,I}=1$	$n_{B,r}=4$
$n_I = 4$	$n_r=5$

- Mutual information = 0.2294
- Sequentially set $\theta = \frac{x_{i,j} + x_{i,j+1}}{2}$ over all j
- Note: just need to consider j s.t. $y_i \neq y_{i+1}$!



Learning Decision Trees

- Def'n: Decision Trees
- Algorithm for Learning Decision Trees
- Overfitting
- Topics:
 - k-ary attribute values
 - Real attribute values
 - Other splitting criteria
 - Attribute Cost
 - Missing Values

<u>Skip</u>

...



Desired Properties



Score for split M(D, x_i) related to

$$S \begin{pmatrix} \#(x_i = t, Y = +) \\ \#(x_i = t, Y = -) \end{pmatrix} S \begin{pmatrix} \#(x_i = f, Y = +) \\ \#(x_i = f, Y = -) \end{pmatrix}$$

- Score S(.) should be
 - Score is BEST for [+0, -200]
 - Score is WORST for [+100, -100]
 - Score is "symmetric"

 Same for [+19, -5] and [+5, -19]
 - Deals with any number of values

Other Splitting Criteria

- Why use *Uncert* as splitting criterion?
- Want:
 - Large "use me" value if split is [85, 0, 0, ..., 0]
 - Small "avoid me" value if split is [5, 5, 5, ..., 5]
- True of *Uncert*, *Uncert_Ratio*... and also for...
- χ^2 statistical tests:

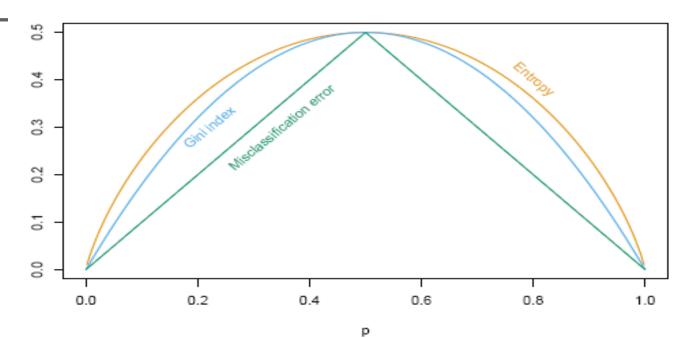
For each attr A, compute deviation:

$$\chi^{2}(A) = \sum_{i=1}^{k} \frac{(p_{i} - \widehat{p_{i}})^{2}}{\widehat{p_{i}}} + \frac{(n_{i} - \widehat{n_{i}})^{2}}{\widehat{n_{i}}}$$

- GINI index: GINI(A) = $\sum_{i} \sum_{j \neq i} p_i p_j = 1 \sum_{i} p_i^2$
- Others: "Marshall Correction" "G" statistic



Node Impurity Measures



- Node impurity measures for 2-class classification
 - function of the proportion p in class 2
 - Scaled cross-entropy has been scaled to pass through (0.5, 0.5)



Example of χ^2

Attributes T₁, T₂ class c

• As 4.29 $(T_1) > 0.533 (T_2)$, ... use T_1 (less likely to be irrelevant)

4

Example of GINI

Attributes T₁, T₂ class c

Scores:
$$GINI(T_1) = \begin{bmatrix} 1 - (\frac{5}{10})^2 - (\frac{5}{10})^2 \end{bmatrix} - \begin{bmatrix} \frac{3}{10} \left(1 - (\frac{0}{3})^2 - (\frac{3}{3})^2\right) + \frac{7}{10} \left(1 - (\frac{5}{7})^2 - (\frac{2}{7})^2\right) \end{bmatrix}$$

$$= \frac{1}{10} \left((\frac{0^2}{3} + \frac{3^2}{3} + \frac{5^2}{7} + \frac{2^2}{7}) - (\frac{5^2}{10} + \frac{5^2}{10}) \right)$$

$$= 0.21429$$

$$GINI(T_2) = \frac{1}{10} \left((\frac{3^2}{5} + \frac{2^2}{5} + \frac{1^2}{3} + \frac{2^2}{3} + \frac{1^2}{2} + \frac{1^2}{2} \right) - (\frac{5^2}{10} + \frac{5^2}{10}) \right)$$

$$= 0.026667$$

- As $GINI(T_1) > GINI(T_2)$, split on T_1
- As " $(\frac{5^2}{10} + \frac{5^2}{10})$ " is subtracted from both, can just use first term. . . and ignore $\frac{1}{10}$ multiplier

Cost-Sensitive Classification ... Learning

- So far, only considered ACCURACY
 In gen'l, may want to consider COST as well
 - medical diagnosis: BloodTest costs \$150
 - robotics: Width_from_1ft costs 23 sec
- Learn a consistent tree with low expected cost?
 - ... perhaps replace *Uncert(S,A)* by
 - Gain²(S,A) / Cost(A) [Tan/Schlimmer'90]
 - $[2^{Gain(S,A)} 1] / [Cost(A) + 1]^{w}$

where $w \in [0, 1]$ determines importance of cost [Nunez'88]

General utility (arb rep'n)

```
E[\sum_{i} cost(A_{i}) + Misclass penalty] [Greiner/Grove/Roth'96]
```

Unknown Attribute Values

- Q: What if some examples are incomplete
 - ... missing values of some attributes?

When learning:

- A1: Throw out all incomplete examples?
 - ... may throw out too many...
- A2: Fill in most common value ("imputation")
 - May miss correlations with other values
 - If impute wrt attributes: may require high-order statistics
- A3: Follow all paths, w/ appropriate weights
 - Huge computational cost if missing MANY values

When classifying

- Similar ideas ...
- ISSUE: Why are values missing?
 - Transmission Noise
 - "Bald men wear hats"
 - "You don't care"

See [Schuurmans/Greiner'94]

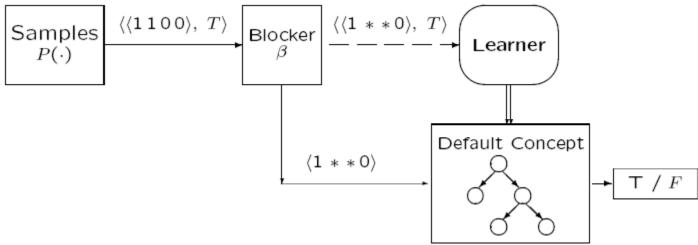


Handling Missing Values: Proportional Distribution

- Associate weight w_i with example [x_i, y_i]
 At root, each example has weight 1.0
- Modify mutual information computations: use weights instead of counts
- When considering test on attribute j, only consider examples that include x_{ii}
- When splitting examples on attribute j:
 - p_L = prob. non-missing example sent left p_R = prob. non-missing example sent right
 - For each example [x_i, y_i]missing attribute j: send it to both children;
 - to left w/ $w_i := w_i \times p_L$
 - to right w/ $w_i := w_i \times p_R$
 - To classify example missing attribute j:
 - Send it down left subtree; $P(\tilde{y}_L \mid x) = \text{resulting prediction}$
 - Send it down left subtree; $P(\tilde{y}_R \mid x) = \text{resulting prediction}$
 - Return $p_L \times P(\tilde{y}_L \mid x) + p_R \times P(\tilde{y}_R \mid x)$



Questions



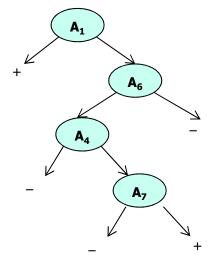
- 1. How to represent default concept?
- 2. When is best default concept learnable?
- If so, how many instances are required?
- 4. Is it better to learn from ...
 - Complete Samples, or
 - Incomplete Samples?



Learning Task

Input:

Output: (small) decision tree consistent with data



Motivation

- Most learning systems work best when
 - few attribute values are missing
 - missing values randomly distributed
- but... [Porter, Bareiss, Holte'90]
 - many datasets missing > ½ values!
 - not randomly missing but ...
 - "[missing] when they are known to be irrelevant for classification or redundant with features already present in the case description"
 - ⇒ Our Situation!!
- Why Learn?
 - A: ... when experts are ...
 - not available, or
 - unable to articulate classification process



Learning Decision Trees ... with "You Don't Care" Omissions

- No known algorithm for PAC-learning gen'l Decision Trees given all attribute values
- ... but Decision Trees are TRIVIAL to learn,
 if superfluous values are omitted:

```
Algorithm GrowDT

Collect "enough" labeled (blocked) instances

Let root = never-blocked feature x_i

Split instances by x_i = 1 vs x_i = 0,

and recur (until purity)
```

4

Conclusions

- Decision trees are extremely popular ML mining tool
 - Easy to understand
 - Easy to implement
 - Easy to use
 - Computationally cheap
- Software Systems:
 - C5.0 (from ID3, C4.5) [Quinlan'93]
 - CART
 - **...**
- Many many applications
- Need to avoid overfitting



What we haven't discussed...

- Decision Stumps" (1-level DT) seem to work surprisingly well
- Efficient alg's for learning optimal "depth-k decision trees" ... even if continuous variables
- Oblique Decision Trees Not just " $x_3 > 5$ ", but " $7 x_4 - 3 x_8 > 91$ "
- Use of prior knowledge
 - Incremental Learners ("Theory Revision")
 - "Relevance" info
- Real-valued outputs Regression Trees
- Bayesian Decision Trees
 - a different approach to preventing overfitting
- Boosting: a simple way to improve accuracy
 - ... ensemble methods



For more information

Two nice books

- Classification and Regression Trees. L. Breiman, J. H. Friedman,
 R. A. Olshen, and C. J. Stone. Wadsworth, Belmont, CA, 1984.
- C4.5: Programs for Machine Learning (Morgan Kaufmann Series in Machine Learning) by J. Ross Quinlan

Both started ≈ 1983, in Bay Area...done independently -- CS vs Stat

Dozens of nice papers, including

- Learning Classification Trees, Wray Buntine, Statistics and Computation (1992), Vol 2, pages 63-73
- On the Boosting Ability of Top-Down Decision Tree Learning Algorithms. Kearns and Mansour, STOC: ACM Symposium on Theory of Computing, 1996
- Dozens of software implementations available on the web for free and commercially for prices ranging between \$50 - \$300,000



Nice Web Sites

Check out...

http://webdocs.cs.ualberta.ca/~aixplore /learning/DecisionTrees/index.html

http://www.r2d3.us/visual-intro-tomachine-learning-part-1/



What you should know

- Information gain:
 - What is it? Why use it?
- Recursive algorithm for building an unpruned decision tree
- Why pruning can reduce test set error
- How to include real-valued inputs
- Computational complexity
 - straightforward, cheap
- Coping with Missing Data
- Alternatives to Information Gain for splitting nodes

Why include these slides?

- Decision trees are often used as they are "intuitive"
 - Suggest explanations for decisions
 - ... kinda ...
- But typically not as accurate as other models
- Include this material as:
 - another approach to learning (beyond linear-ish models)
 - other intuitions about learning process
 - Entropy, information content, ...
 - χ^2 , MDL, ...
 - Missing data ("You don't care", not "MCAR", ...)
 - alternative tweaks...