

Learning Decision Trees

- A Decision Tree is a tree-structured plan of a set of attributes to test in order to predict the output.
- To decide which attribute should be tested first, simply find the one with the highest information gain.
- Then recurse...

Summary of Basic Decision Tree Building

From www.cs.cmu.edu/~awm/tutorials

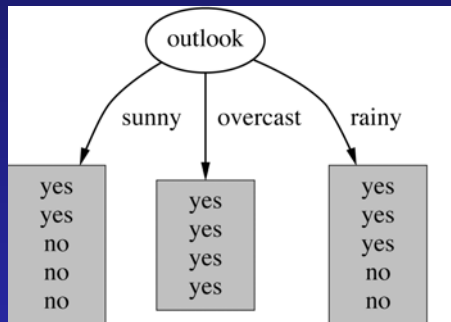
BuildTree(DataSet, Output)

- 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 .

Information gain (IG) of an attribute (outlook)



Info[2,3]=0.97 bits

Info[3,2]=0.97 bits

Info[4,0]=0 bits

$H(\text{playtennis}/\text{outlook}=\text{sunny})$

=Info[2,3]

=entropy(2/5, 3/5)

= $-2/5 \log 2/5 - 3/5 \log 3/5$

= **0.97 bits**

$H(\text{playtennis}/\text{outlook})=$

Average info of subtree(weighted)

= $0.97 \times 5/14 + 0 \times 4/14 + 0.97 \times 5/14$

= **0.693 bits**

$H(\text{playtennis}) = \text{Info of all training samples} = \text{Info}[9,5] = 0.94$

$\text{IG}(\text{playtennis}/\text{outlook}) = 0.94 - 0.693 = 0.247 \text{ bits}$

Decision Trees and Data mining

- Divide and conquer approach to decision tree induction
- Also called “top-down induction of decision trees”
- Information gain algorithm known as ID3
- Led to widely used system known as C4.5

About this Data set

- Tiny subset of the 1990 US census
- Publicly available from the UCI Machine Learning Repository

Used Attributes

age	edunum	race	hours_worked
employment	marital	gender	country
taxweighting	job	capitalgain	wealth
education	relation	capitalloss	agegroup

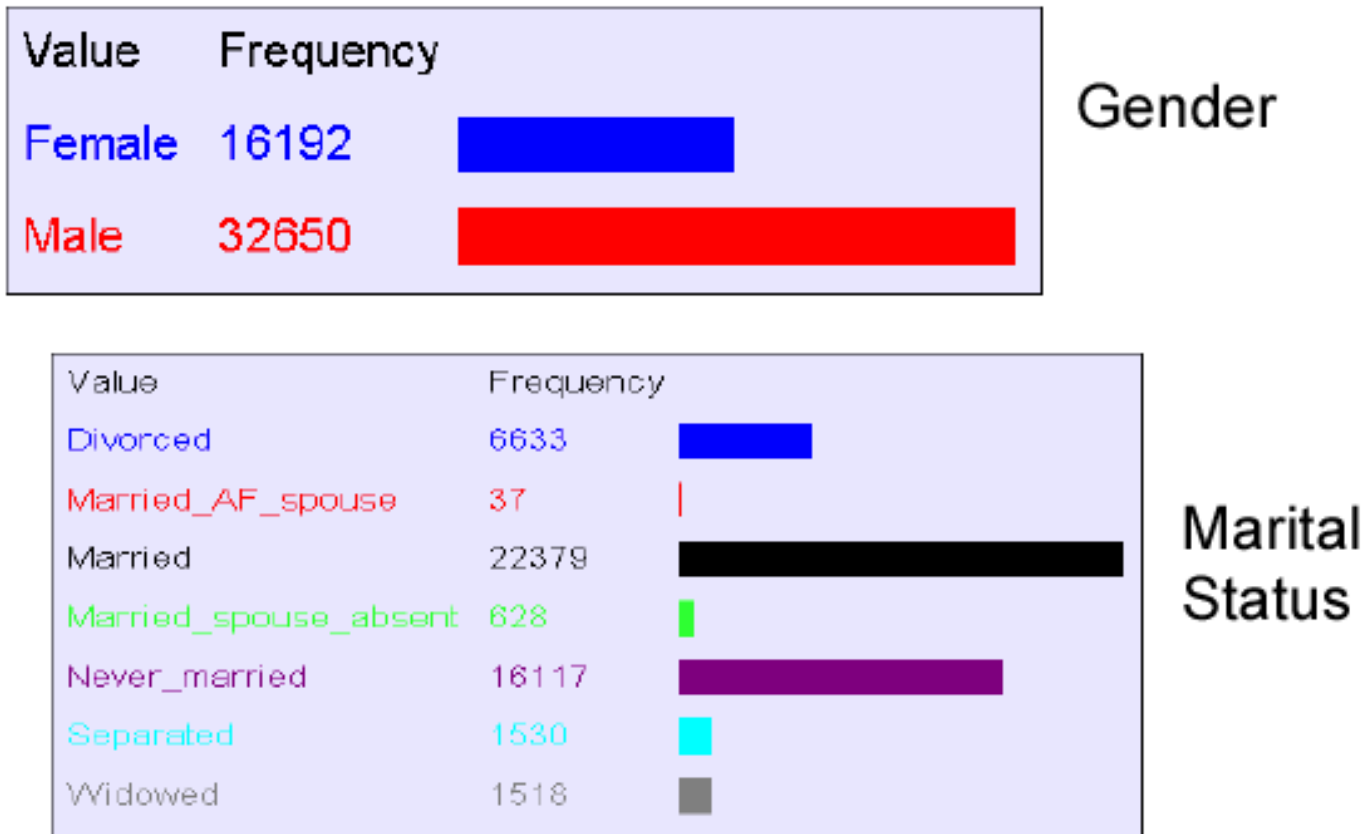
This color = Real-valued This color = Symbol-valued

Successfully loaded a new dataset from the file \tadult.fds. It has 16 attributes and 48842 records.

Classification

- A major data mining operation
- Try to predict an attribute, e.g., wealth by means of other attributes
- Applies to categorical outputs
 - Categorical attribute: an attribute which takes on two or more discrete values. Also known as a symbolic attribute
 - Real attribute: a column of real numbers

What can you do with a data set? Histograms

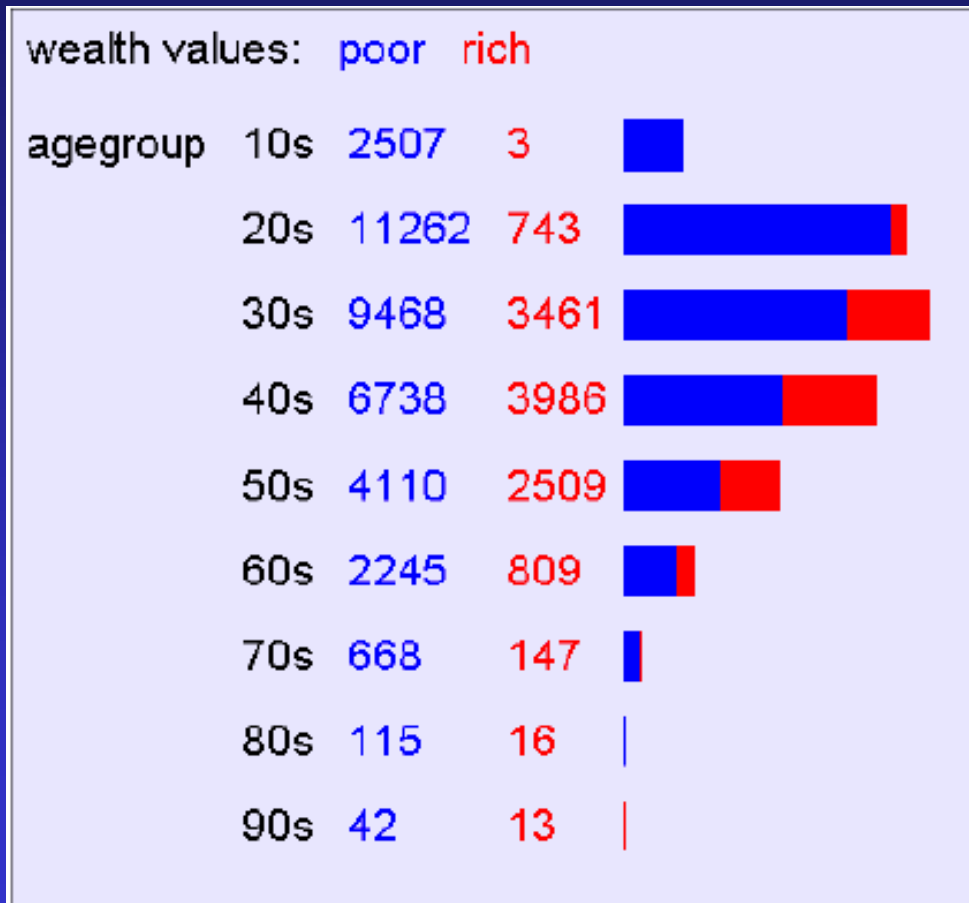


2-D Contingency Table

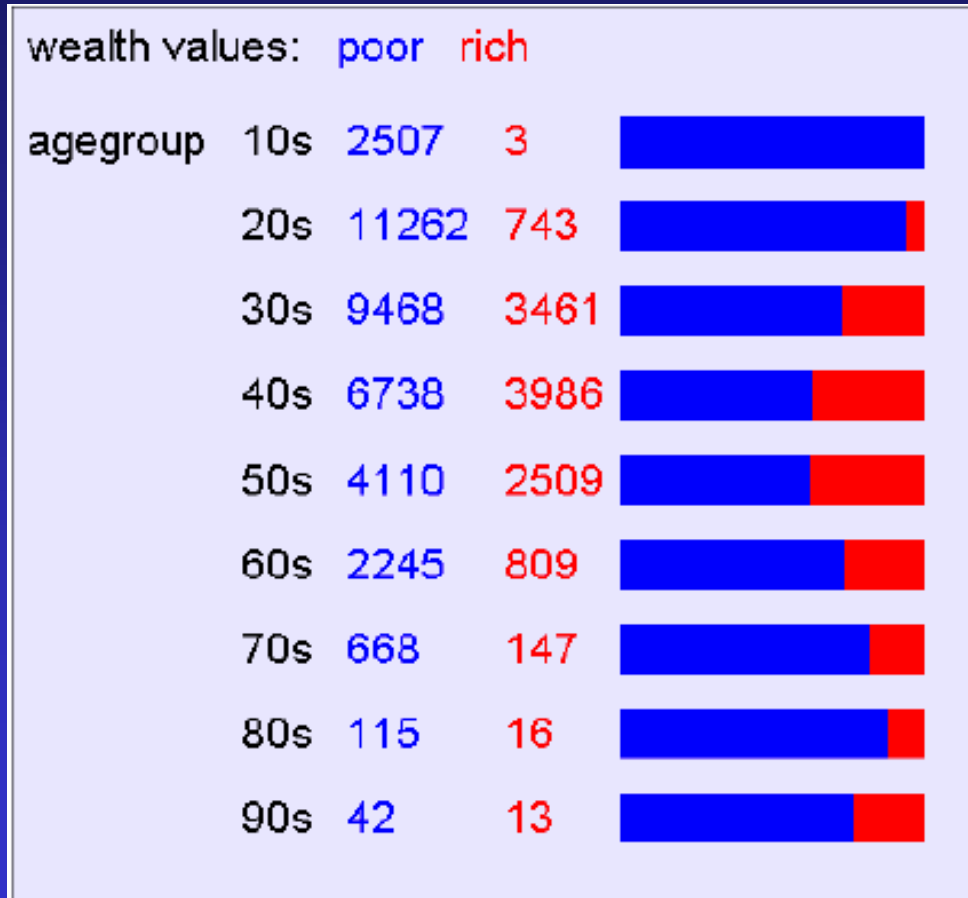
wealth values:		poor	rich
agegroup	10s	2507	3
	20s	11262	743
	30s	9468	3461
	40s	6738	3986
	50s	4110	2509
	60s	2245	809
	70s	668	147
	80s	115	16
	90s	42	13

- For each pair of values for attributes (agegroup, wealth)
- Count how many records match

Graphical Representation of 2-D contingency Table



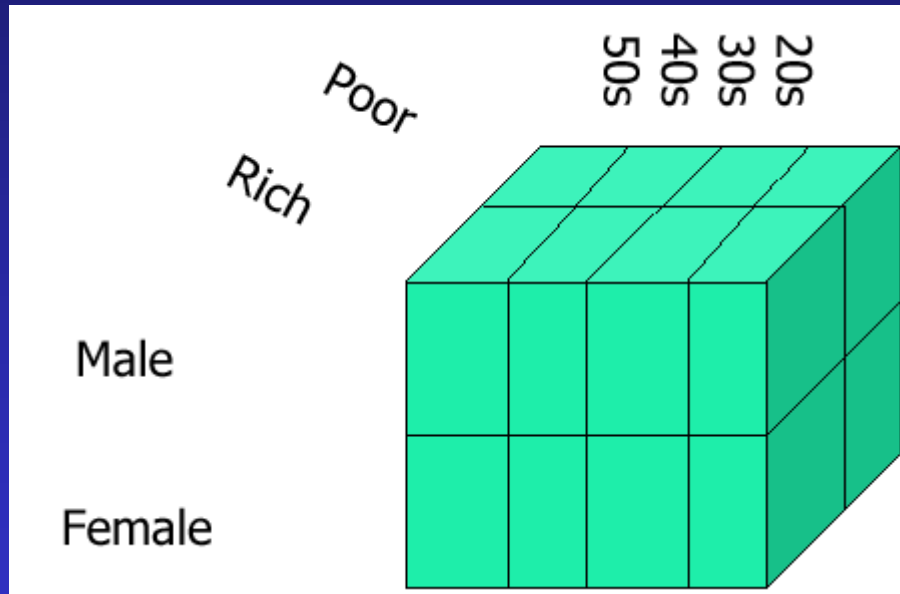
Easier to see more interesting things if we stretch out histogram bars



2-D Contingency Table showing all 16 attributes

job values:		Acn_clerical	Craft_repair	Farming_fishing	Machine_op_inspct	Priv_house_serv	Protective_serv	Tech_support									
MissingValue		Armed_Forces	Exec_managerial	Handlers_cleaners	Other_service	Prof_specialty	Sales	Transport_moving									
marital	Divorced	270	1192	0	679	890	90	197	434	762	46	795	121	664	239	254	
	Married_AF_spouse	5	6	0	4	3	1	1	1	5	0	4	1	5	0	1	
	Married	928	1495	7	3818	3600	889	724	1469	1088	27	3182	583	2491	609	1489	
	Married_spouse_absent	45	84	0	77	52	35	32	37	92	9	64	7	55	9	30	
	Never_married	1242	2360	8	1301	1260	434	1029	872	2442	99	1849	237	1992	506	486	
	Separated	97	224	0	160	126	23	63	123	275	21	145	23	146	48	56	
	Widowed	222	250	0	73	155	38	26	86	259	40	133	11	151	35	39	

3-D Contingency Table: much harder to look at



On Line Analytical Processing (OLAP)

- Software packages and database add-ons to do this are known as OLAP tools
- They usually include point and click navigation to view slices and aggregates of contingency tables
- They usually include nice histogram visualization
- Too many tables though
 - With 16 attributes there are:
 - 16 1-D contingency tables
 - $^{16}C_2 = 120$ 2-D contingency tables
 - 560 3-D contingency tables

Information Gain

wealth values: poor rich

relation Husband	10870	8846		$H(\text{wealth} \mid \text{relation} = \text{Husband}) = 0.992385$
Not_in_family	11307	1276		$H(\text{wealth} \mid \text{relation} = \text{Not_in_family}) = 0.473439$
Other_relative	1454	52		$H(\text{wealth} \mid \text{relation} = \text{Other_relative}) = 0.216617$
Own_child	7470	111		$H(\text{wealth} \mid \text{relation} = \text{Own_child}) = 0.110192$
Unmarried	4816	309		$H(\text{wealth} \mid \text{relation} = \text{Unmarried}) = 0.328606$
Wife	1238	1093		$H(\text{wealth} \mid \text{relation} = \text{Wife}) = 0.997207$

$H(\text{wealth}) = 0.793844$ $H(\text{wealth} \mid \text{relation}) = 0.628421$

$IG(\text{wealth} \mid \text{relation}) = 0.165423$

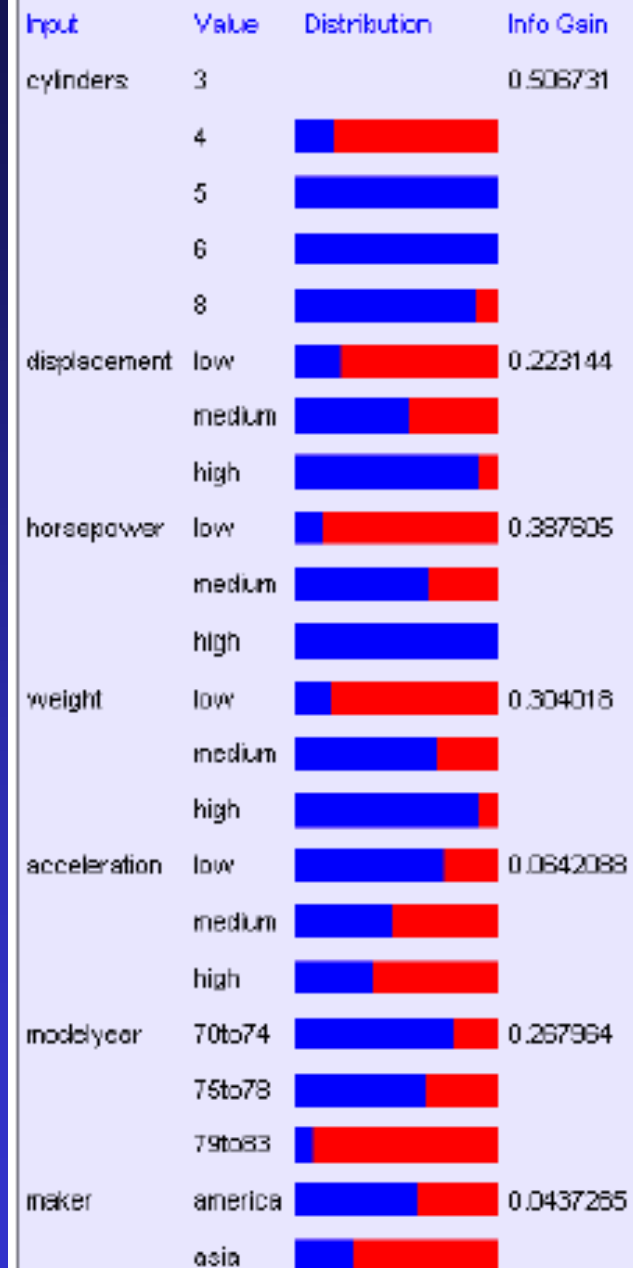
MPG Data Set (40 records from UCI repository)

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	europa
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	europa
bad	5	medium	medium	medium	medium	75to78	europa

- Predicting MPG
 - Has good/bad values
- Look at all the Information Gains
- Cylinders has the highest IG

Information gains using the training set (40 records)

mpg values: bad good



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

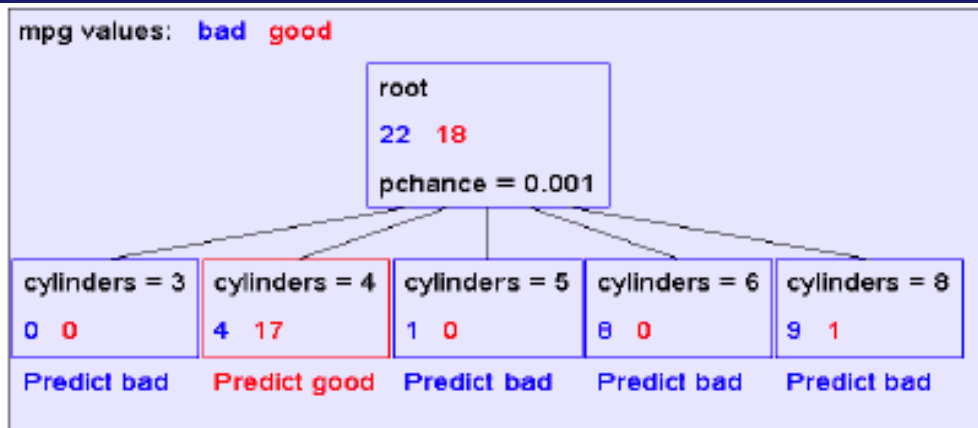
Predict good

Predict bad

Predict bad

Predict bad

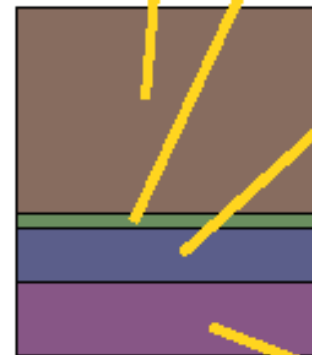
Recursion Step



Take the Original Dataset..



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



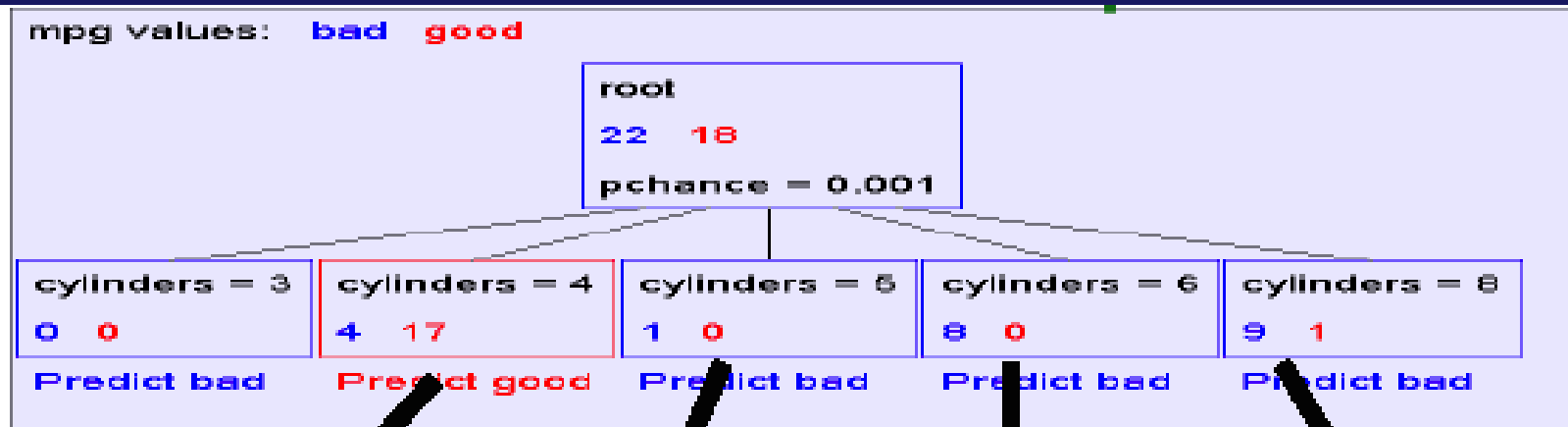
Records in which cylinders = 4

Records in which cylinders = 5

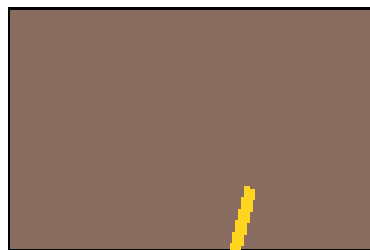
Records in which cylinders = 6

Records in which cylinders = 8

Recursion Step



Build tree from
These records..



Records in
which
cylinders = 4

Build tree from
These records..



Records in
which
cylinders = 5

Build tree from
These records..



Records in
which
cylinders = 6

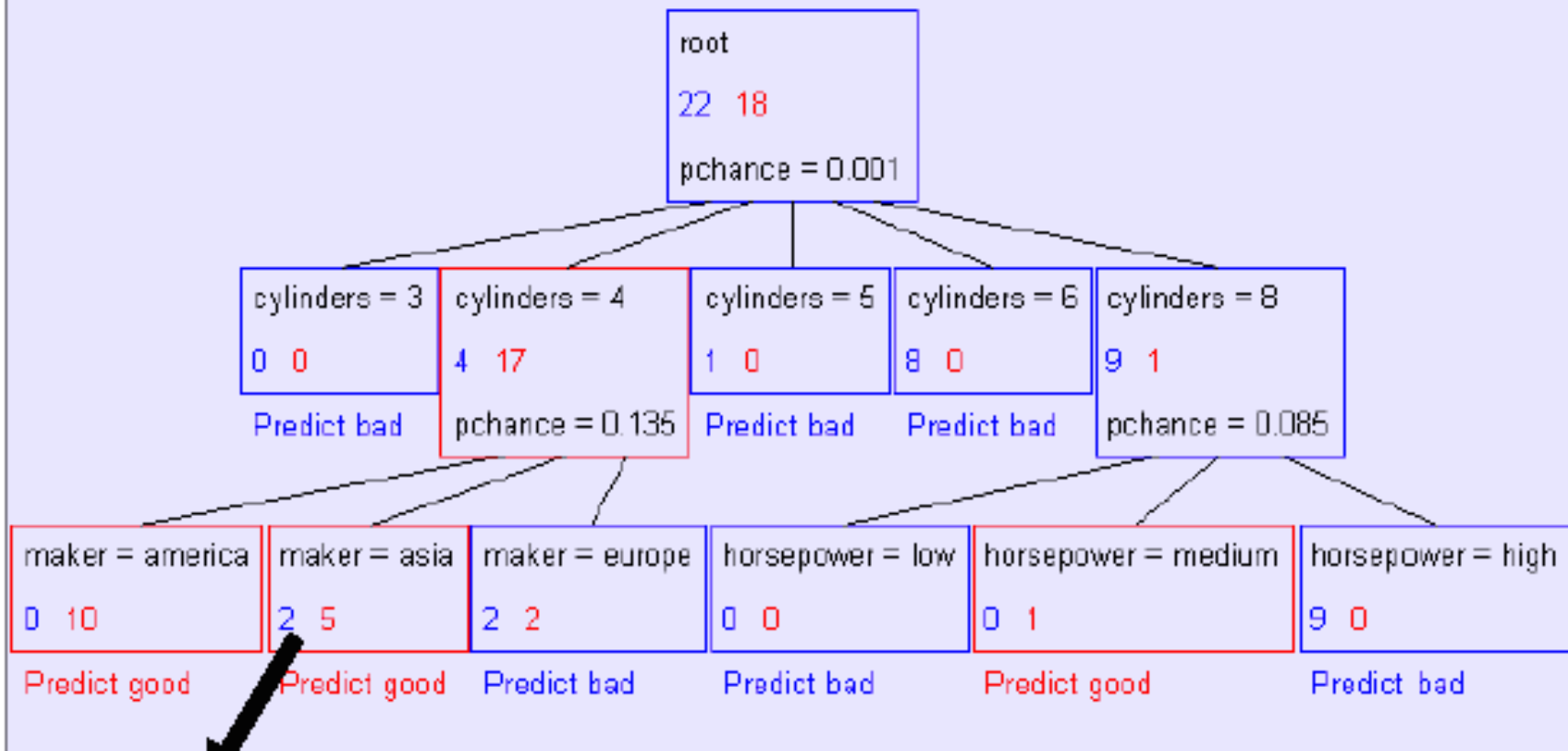
Build tree from
These records..



Records in
which
cylinders = 8

Second Level of Tree

mpg values: bad good

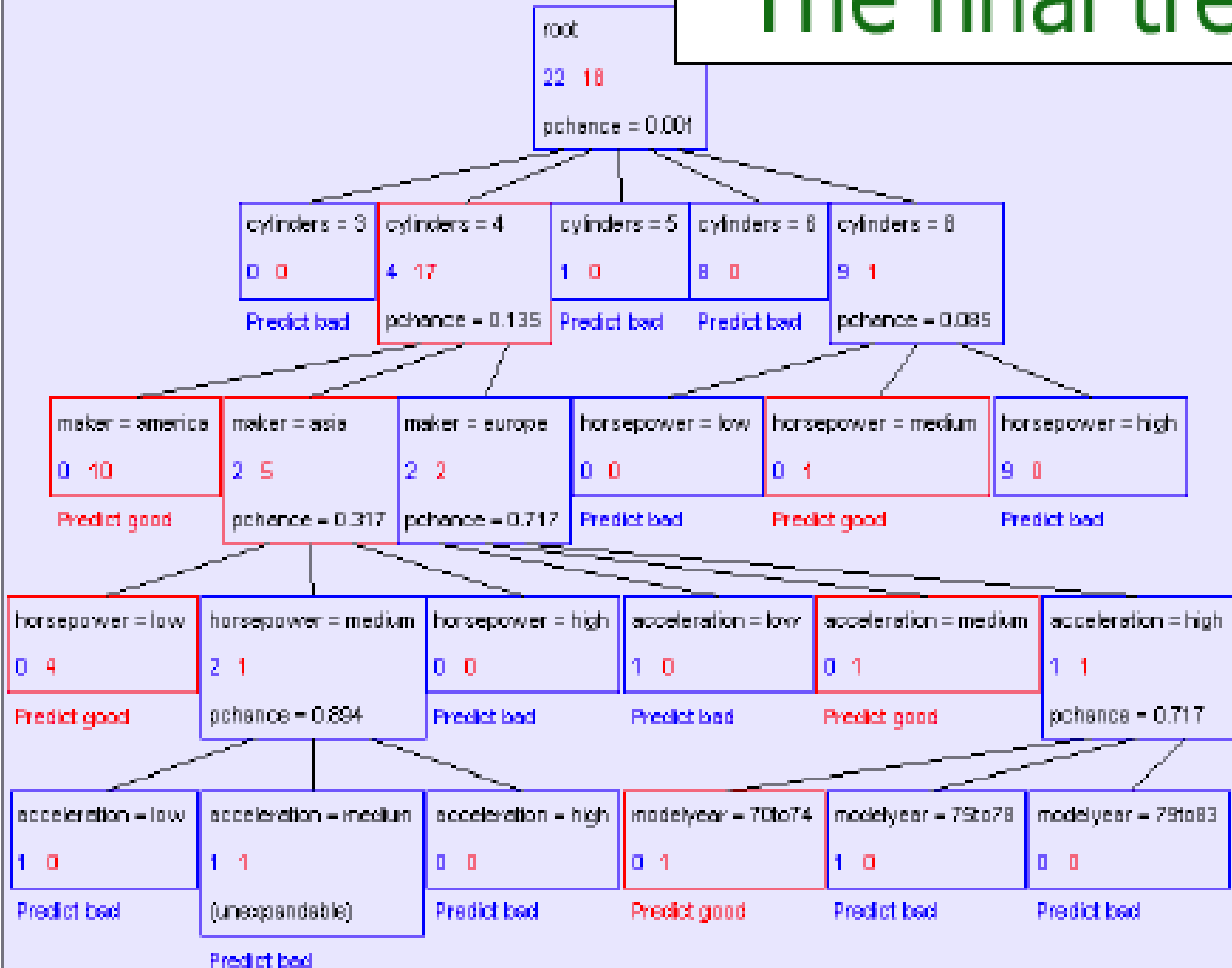


Recursively build a tree from the seven records in which there are four cylinders and the maker was based in Asia

(Similar recursion in the other cases)

The final tree

mpg values: bad good



Incorporating Continuous-Valued Attributes

- Initial definition of ID3 is restricted to attributes that take on discrete values
 - target attribute must be discrete valued
 - attributes tested at decision nodes must be discrete valued
- Example:
 - Convert a continuous-valued test attribute A into a Boolean valued attribute A_c
 - create a Boolean attribute A_c that is true if $A < c$, false otherwise

Example of Continuous-valued Features

	Feature	Minimum	Maximum
1	Ave Entropy	0.275	1.423216
2	Ave Grey-Threshold	170	226
3	NoBlackPixels	81105	873557
4	Ave NoExtContours	4	44
5	Ave NoIntContours	5	35
6	Ave SlopeHorizontal	0.094	0.5707
7	Ave SlopePositive	0.1	0.542
8	Ave SlopeVertical	0.066	0.546
9	Ave SlopeNegative	0.048	0.36
10	Ave StrokeWidth	4	11
11	Ave Slant	23.04	36.876
12	Ave Height	13	66

Thresholding a continuous value to obtain a Boolean Attribute

- Algorithm
 - Sort examples according to continuous attribute A
 - Identify adjacent examples that differ in target classification
 - Generate thresholds midway between corresponding values of A

Temperature	40	48	60	72	80	90
PlayTennis	No	No	Yes	Yes	Yes	No
			↓			↓

Thresholding Algorithm

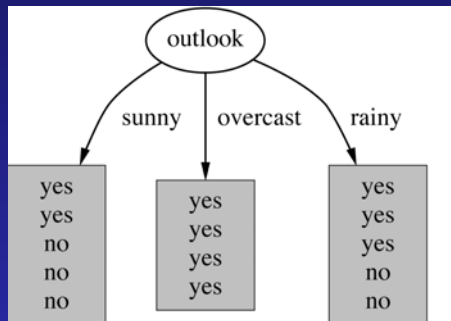
Temperature	40	48	60	72	80	90
PlayTennis	No	No	Yes	Yes	Yes	No

$$c = (48+60)/2 \\ = 54$$

$$c = (80+90)/2 \\ = 85$$

- There are two Boolean valued candidate attributes
 - $\text{Temperature}_{>54}$ and $\text{Temperature}_{>85}$
- Which one is better?
 - One that produces greatest information gain

Information gain of an attribute



Info[2,3]=0.97 bits

Info[3,2]=0.97 bits

Info[4,0]=0 bits

$$\begin{aligned} \text{Info}[2,3] &= \text{entropy}(2/5, 3/5) \\ &= -2/5 \log 2/5 - 3/5 \log 3/5 \\ &= 0.97 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{Average info of subtree (weighted)} &= \\ &= 0.97 \times 5/14 + \\ &= 0 \times 4/14 + \\ &= 0.97 \times 5/14 = 0.693 \text{ bits} \end{aligned}$$

Info of all training samples, $\text{info}[9,5] = 0.94$

$$\text{gain}(\text{outlook}) = 0.94 - 0.693 = 0.247 \text{ bits}$$

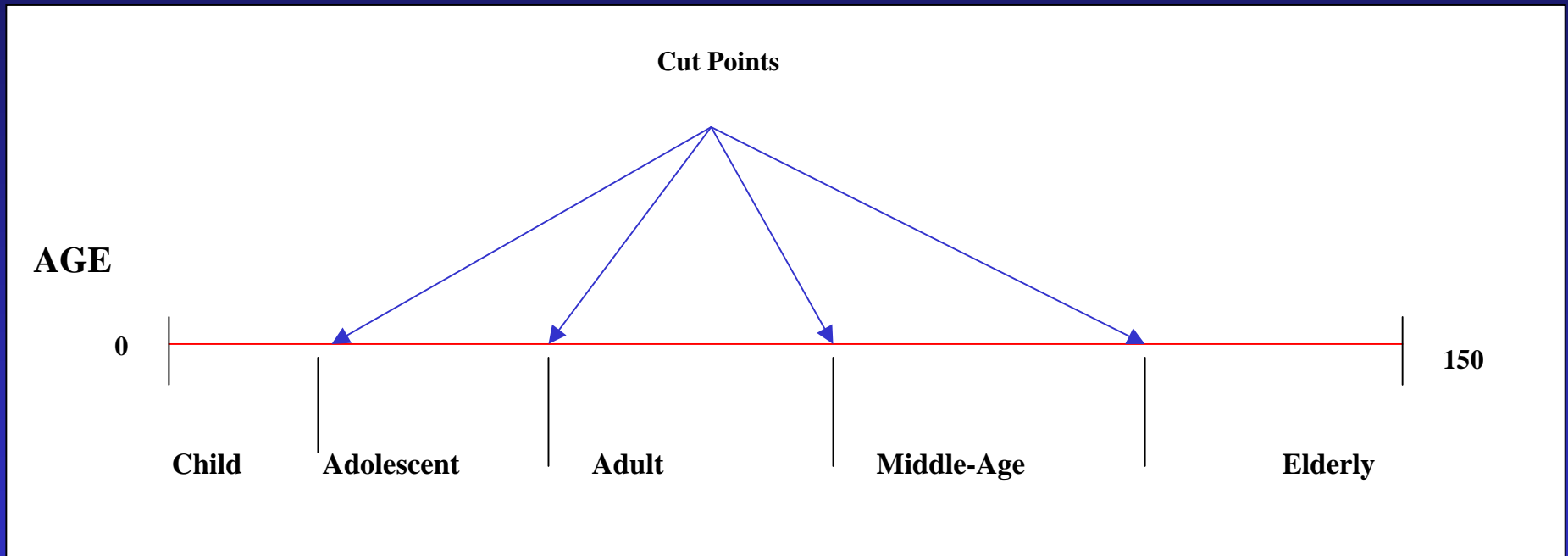
Algorithm for thresholding a continuous-valued feature

- At a given node
 - For a given continuous-valued feature
 - Obtain candidate thresholds for continuous-valued feature:
 - points where target attribute changes value
 - Determine information gain for each threshold
 - Choose threshold with highest information gain

Data Reduction

- Feature Discretization
 - Discretize values of continuous features into a small number of intervals
 - Each interval is mapped into a discrete symbol
 - Each feature handled independently of other features
- Benefits
 - Simplified data description
 - Easy to understand data
 - Easy to understand data mining results
 - Applicability of several machine learning algorithms
 - Quality of the analysis improved

Discretization of the age feature



Automated Discretization Techniques: not based on target attribute value

- Values are placed in bins
- Partitioning into groups with close values
- All values in a bin will be merged into a single concept represented by a single value – usually the mean or median of the bin's values

Clustering for binning

- Set of values for a given feature = {3,2,1,5,4,3,1,7,5,3}
- Sorted set = {1,1,2,3,3,3,4,5,5,7}
- Split into three bins of approx same size
 - {1,1,2, 3,3,3 4,5,5,7}
 - Bin1 Bin2 Bin3
- Bin representative selection
 - Modes:
 - {1,1,1 3,3,3 5,5,5,5}
 - Mean Values:
 - {1.33, 1.33, 1.33 3,3,3 5.25, 5.25, 5.25}
 - Closest Boundary Values:
 - {1,1,2 3,3,3 4,4,4,7}

K-means clustering

- Sort all values for a given feature
- Assign approximately equal numbers of sorted adjacent values to each of k bins
- Move a border element from one bin to next when that reduces the sum of all distances from each number to the mean or mode of the assigned bin
- {1,1,2,2,2,5,6,8,8,9}
- Gets assigned as {1,2,2,2,2, 5,6, 8,8,9}

Chimerge: an automatic discretization algorithm

- Analyzes the quality of multiple intervals for a given feature by using chi-squared statistics

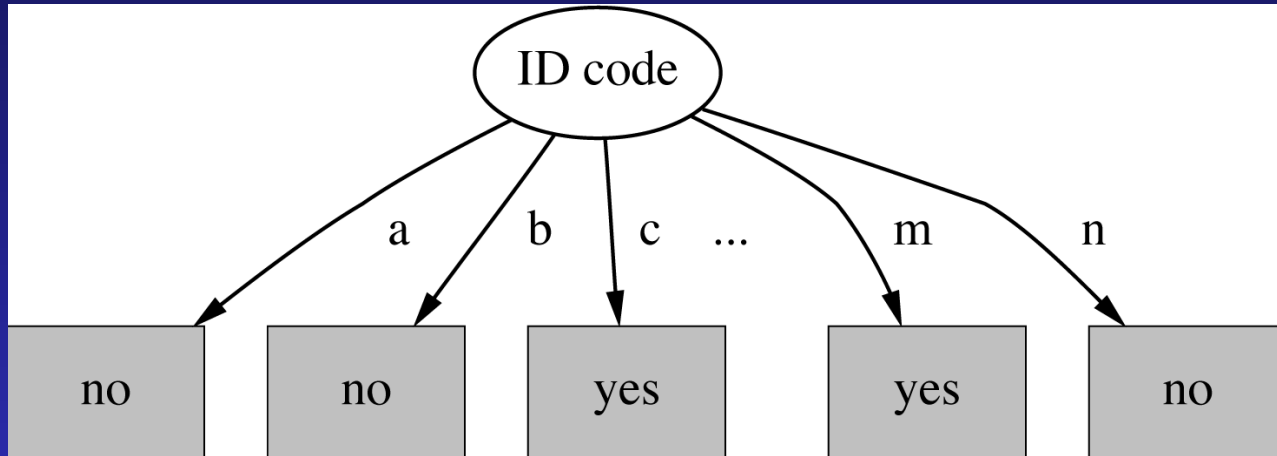
Highly Branching Attributes

- When attribute has a large number of possible values, information gain is large
 - example attributes are Date (March 4, 1979), or ID code
- Consider case when each instance has an Identification Code attribute

Play Tennis Data with Identification Codes

ID Code	Outlook	Temperature	Humidity	Windy	Play
a	sunny	hot	high	false	no
b	sunny	hot	high	true	no
c	overcast	hot	high	false	yes
d	rainy	mild	high	false	yes
e	rainy	cool	normal	false	yes
f	rainy	cool	normal	true	no
g	overcast	cool	normal	true	yes
h	sunny	mild	high	false	no
i	sunny	cool	normal	false	yes
j	rainy	mild	normal	false	yes
k	sunny	mild	normal	true	yes
l	overcast	mild	high	true	yes
m	overcast	hot	normal	false	yes
n	rainy	mild	high	true	no

Tree Stump for the ID Code Attribute



Information gain of each node is just the information at the root (0.94)

Therefore chosen as the splitting attribute!

Alternative Measures for Selecting Attributes

- Split information is sensitive to how broadly and uniformly the attribute splits the data

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

- S_1 through S_c are c subsets of examples resulting from partitioning the c -values of attribute A
- It is the entropy of S with respect to the values of attribute A
- (not wrt target attribute, as considered before)

Alternative Measures for Selecting Attributes, continued

$$\textit{GainRatio}(S, A) \equiv \frac{\textit{Gain}(S, A)}{\textit{SplitInformation}(S, A)}$$

S is the collection of data and A is the attribute

The Splitinformation term discourages the selection of attributes with many uniformly distributed values

Gain Ratio

- Take into account number and size of child nodes
- For ID code attribute:
 - Split Information[1,1,...,1] = $-1/14 \times \log 1/14 \times 14$
 - $= \log 14 = 3.807$ bits
 - Gain ratio = $0.940/3.807$
 - $= 0.246$
 - Gain Ratio reduces the effect of high branching

Gain Ratio Calculations for the Tree Stumps

Outlook		Temperature	
info:	0.693	info:	0.911
gain:0.940-0.693	0.247	gain:0.940-0.911	0.029
split info: info([5,4,5])	1.577	split info: info([4,6,4])	1.362
gain ratio:0.247/1.577	0.156	0.029/1.362	0.021
Humidity		Windy	
info:	0.788	info:	0.892
gain:0.940-0.788	0.152	gain:0.940-0.892	0.048
split info: info([7,7])	1.000	split info: info([8,6])	0.985
gain ratio:0.152/1	0.152	gain ratio:0.048/0.985	0.049

Handling Training Examples with Missing Attribute Values (3.7.4)

- Medical domain
 - Predict patient outcome based on lab tests
 - Lab test available only for subset of patients
- Estimate missing attribute using examples with known values

Handling Training Examples with Missing Attribute Values

- $Gain(S, A)$ is to be calculated at node n to test whether attribute A is the best attribute to test at this decision node
- Suppose $\langle x, c(x) \rangle$ is a training example in S and the value of $A(x)$ is unknown

Handling Training Examples with Missing Attribute Values

- Strategies
 - assign missing attribute the value that is most common among training examples at node n .
 - assign it the most common value among examples at node n that have classification $c(x)$
 - More complex procedure:
 - assign a probability to each of the possible values of A rather than simply assigning the most common value to $A(x)$
 - e.g., $P(A(x)=1) = 0.6$ and $P(A(x)=0)=0.4$
 - use fractional samples for computing information Gain