

Data Mining

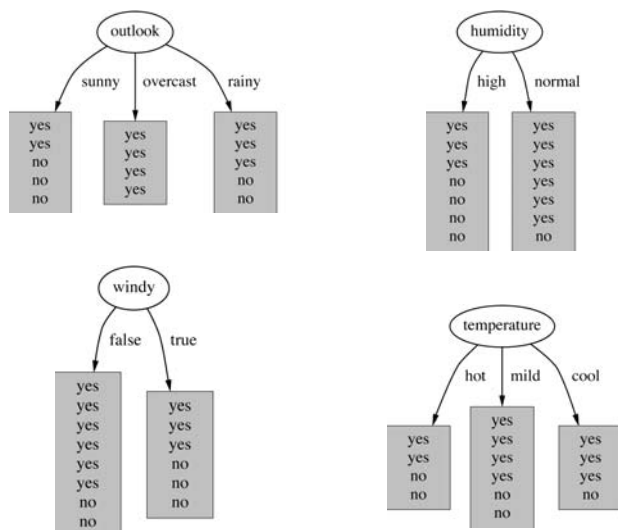
Practical Machine Learning Tools and Techniques

Slides for Chapter 4 of *Data Mining* by I. H. Witten and E. Frank
Decision Trees

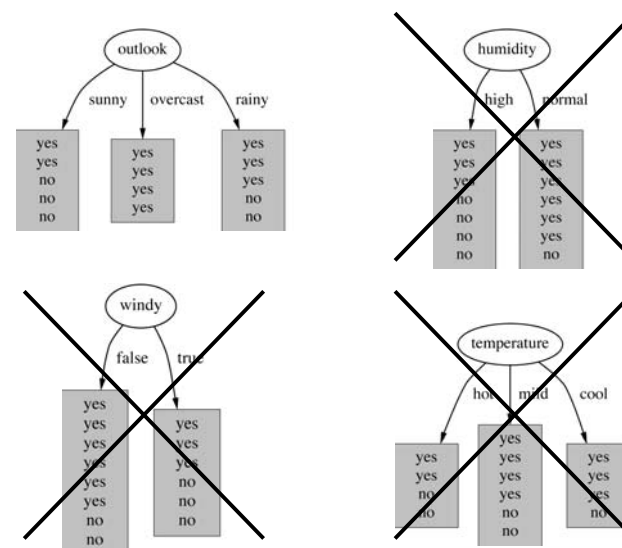
Constructing decision trees

- Strategy: top down
Recursive *divide-and-conquer* fashion
 - First: select attribute for root node
Create branch for each possible attribute value
 - Then: split instances into subsets
One for each branch extending from the node
 - Finally: repeat recursively for each branch, using only instances that reach the branch
- Stop if all instances have the same class

Which attribute to select?




Which attribute to select?



- Which is the best attribute?
 - ♦ Want to get the smallest tree
 - ♦ Heuristic: choose the attribute that produces the “purest” nodes
- Popular *impurity criterion: information gain*
 - ♦ Information gain increases with the average purity of the subsets
- Strategy: choose attribute that gives greatest information gain

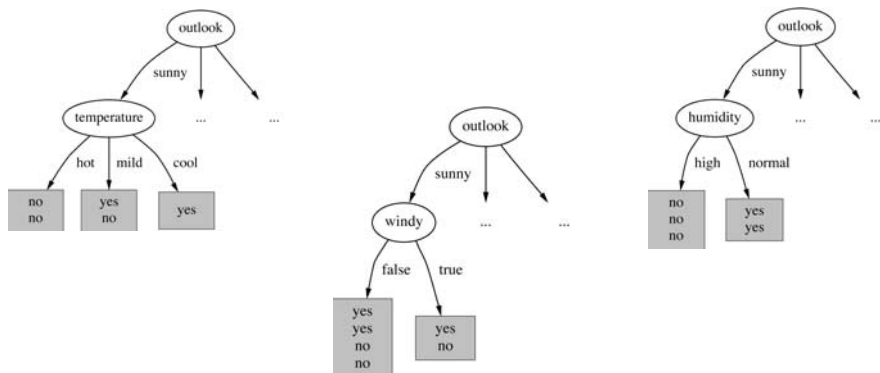
- Measure information in *bits*
 - ♦ Given a probability distribution, the info required to predict an event is the distribution’s *entropy*
 - ♦ Entropy gives the information required in bits (can involve fractions of bits!)
- Formula for computing the entropy:
$$\text{entropy}(p_1, p_2, \dots, p_n) = -p_1 \log p_1 - p_2 \log p_2 \dots - p_n \log p_n$$

- *Outlook = Sunny* :
 $\text{info}([2,3]) = \text{entropy}(2/5, 3/5) = -2/5 \log(2/5) - 3/5 \log(3/5) = 0.971 \text{ bits}$
- *Outlook = Overcast* :
 $\text{info}([4,0]) = \text{entropy}(1, 0) = -1 \log(1) - 0 \log(0) = 0 \text{ bits}$  **Note: this is normally undefined.**
- *Outlook = Rainy* :
 $\text{info}([2,3]) = \text{entropy}(3/5, 2/5) = -3/5 \log(3/5) - 2/5 \log(2/5) = 0.971 \text{ bits}$
- Expected information for attribute:
 $\text{info}([3,2],[4,0],[3,2]) = (5/14) \times 0.971 + (4/14) \times 0 + (5/14) \times 0.971 = 0.693 \text{ bits}$

- Information gain: information before splitting – information after splitting
$$\begin{aligned} \text{gain}(\text{Outlook}) &= \text{info}([9,5]) - \text{info}([2,3],[4,0],[3,2]) \\ &= 0.940 - 0.693 \\ &= 0.247 \text{ bits} \end{aligned}$$
- Information gain for attributes from weather data:

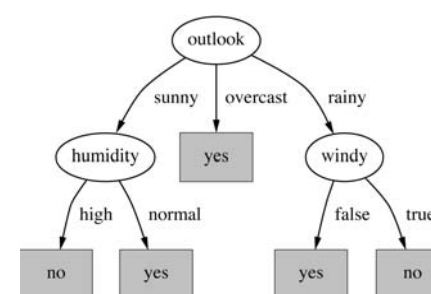
$\text{gain}(\text{Outlook})$	$= 0.247 \text{ bits}$
$\text{gain}(\text{Temperature})$	$= 0.029 \text{ bits}$
$\text{gain}(\text{Humidity})$	$= 0.152 \text{ bits}$
$\text{gain}(\text{Windy})$	$= 0.048 \text{ bits}$

Continuing to split



$$\begin{aligned} \text{gain}(\text{Temperature}) &= 0.571 \text{ bits} \\ \text{gain}(\text{Humidity}) &= 0.971 \text{ bits} \\ \text{gain}(\text{Windy}) &= 0.020 \text{ bits} \end{aligned}$$

Final decision tree



- Note: not all leaves need to be pure; sometimes identical instances have different classes
 \Rightarrow Splitting stops when data can't be split any further

Wishlist for a purity measure

- Properties we require from a purity measure:
 - When node is pure, measure should be zero
 - When impurity is maximal (i.e. all classes equally likely), measure should be maximal
 - Measure should obey *multistage property* (i.e. decisions can be made in several stages):

$$\text{measure}([2,3,4]) = \text{measure}([2,7]) + (7/9) \times \text{measure}([3,4])$$
- Entropy is the only function that satisfies all three properties!

Properties of the entropy

- The multistage property:

$$\text{entropy}(p, q, r) = \text{entropy}(p, q+r) + (q+r) \times \text{entropy}\left(\frac{q}{q+r}, \frac{r}{q+r}\right)$$
- Simplification of computation:

$$\begin{aligned} \text{info}([2,3,4]) &= -2/9 \times \log(2/9) - 3/9 \times \log(3/9) - 4/9 \times \log(4/9) \\ &= [-2 \times \log 2 - 3 \times \log 3 - 4 \times \log 4 + 9 \times \log 9] / 9 \end{aligned}$$
- Note: instead of maximizing info gain we could just minimize information

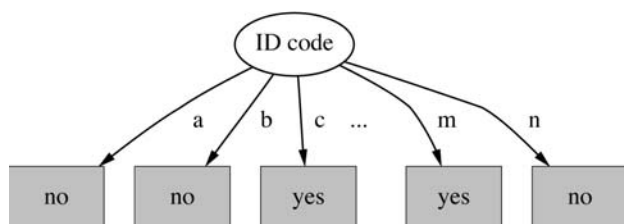
Highly-branching attributes

- Problematic: attributes with a large number of values (extreme case: ID code)
- Subsets are more likely to be pure if there is a large number of values
 - ⇒ Information gain is biased towards choosing attributes with a large number of values
 - ⇒ This may result in *overfitting* (selection of an attribute that is non-optimal for prediction)
- Another problem: *fragmentation*

Weather data with *ID* code

ID code	Outlook	Temp.	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 *ID* code attribute



- Entropy of split:

$$\text{info}(\text{ID code}) = \text{info}([0,1]) + \text{info}([0,1]) + \dots + \text{info}([0,1]) = 0 \text{ bits}$$
 ⇒ Information gain is maximal for ID code (namely 0.940 bits)

Gain ratio

- *Gain ratio*: a modification of the information gain that reduces its bias
- Gain ratio takes number and size of branches into account when choosing an attribute
 - ♦ It corrects the information gain by taking the *intrinsic information* of a split into account
- Intrinsic information: entropy of distribution of instances into branches (i.e. how much info do we need to tell which branch an instance belongs to)

Computing the gain ratio

- Example: intrinsic information for ID code

$$\text{info}([1,1,\dots,1]) = 14 \times (-1/14 \times \log(1/14)) = 3.807 \text{ bits}$$

- Value of attribute decreases as intrinsic information gets larger
- Definition of gain ratio:

$$\text{gain_ratio}(\text{attribute}) = \frac{\text{gain}(\text{attribute})}{\text{intrinsic_info}(\text{attribute})}$$

- Example:

$$\text{gain_ratio}(\text{ID code}) = \frac{0.940 \text{ bits}}{3.807 \text{ bits}} = 0.246$$

Gain ratios for weather data

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])	1.557
Gain ratio: 0.247/1.577	0.157	Gain ratio: 0.029/1.557	0.019
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

More on the gain ratio

- “Outlook” still comes out top
- However: “ID code” has greater gain ratio
 - ♦ Standard fix: *ad hoc* test to prevent splitting on that type of attribute
- Problem with gain ratio: it may overcompensate
 - ♦ May choose an attribute just because its intrinsic information is very low
 - ♦ Standard fix: only consider attributes with greater than average information gain

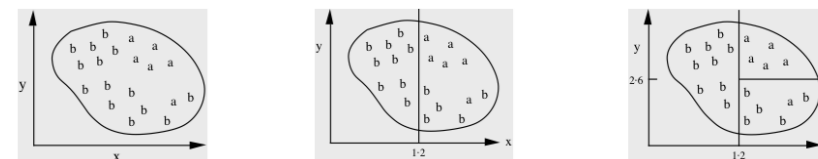
Discussion

- Top-down induction of decision trees: ID3, algorithm developed by Ross Quinlan
 - ♦ Gain ratio just one modification of this basic algorithm
 - ♦ \Rightarrow C4.5: deals with numeric attributes, missing values, noisy data
- Similar approach: CART
- There are many other attribute selection criteria!
(But little difference in accuracy of result)

Covering algorithms

- Convert decision tree into a rule set
 - Straightforward, but rule set overly complex
 - More effective conversions are not trivial
- Instead, can generate rule set directly
 - for each class in turn find rule set that covers all instances in it (excluding instances not in the class)
- Called a *covering* approach:
 - at each stage a rule is identified that “covers” some of the instances

Example: generating a rule



If true
then class = a

If $x > 1.2$
then class = a

If $x > 1.2$ and $y > 2.6$
then class = a

- Possible rule set for class “b”:

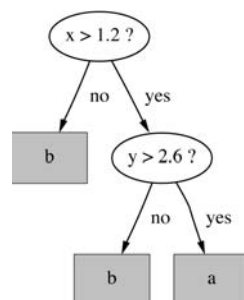
If $x \leq 1.2$ then class = b

If $x > 1.2$ and $y \leq 2.6$ then class = b

- Could add more rules, get “perfect” rule set

Rules vs. trees

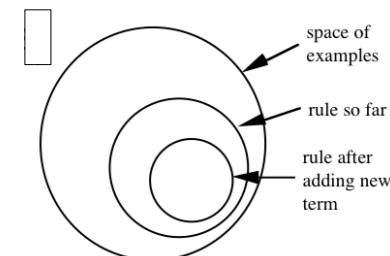
Corresponding decision tree:
(produces exactly the same predictions)



- But: rule sets *can* be more perspicuous when decision trees suffer from replicated subtrees
- Also: in multiclass situations, covering algorithm concentrates on one class at a time whereas decision tree learner takes all classes into account

Simple covering algorithm

- Generates a rule by adding tests that maximize rule’s accuracy
- Similar to situation in decision trees: problem of selecting an attribute to split on
 - But: decision tree inducer maximizes overall purity
- Each new test reduces rule’s coverage:



Selecting a test

- Goal: maximize accuracy
 - t total number of instances covered by rule
 - p positive examples of the class covered by rule
 - $t - p$ number of errors made by rule \Rightarrow Select test that maximizes the ratio p/t
- We are finished when $p/t = 1$ or the set of instances can't be split any further

Example: contact lens data

- Rule we seek: $\text{If ? then recommendation = hard}$
- Possible tests:

Age = Young	2/8
Age = Pre-presbyopic	1/8
Age = Presbyopic	1/8
Spectacle prescription = Myope	3/12
Spectacle prescription = Hypermetrope	1/12
Astigmatism = no	0/12
Astigmatism = yes	4/12
Tear production rate = Reduced	0/12
Tear production rate = Normal	4/12

Modified rule and resulting data

- Rule with best test added:
 $\text{If astigmatism = yes then recommendation = hard}$
- Instances covered by modified rule:

Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	Yes	Reduced	None
Young	Myope	Yes	Normal	Hard
Young	Hypermetrope	Yes	Reduced	None
Young	Hypermetrope	Yes	Normal	hard
Pre-presbyopic	Myope	Yes	Reduced	None
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetrope	Yes	Reduced	None
Pre-presbyopic	Hypermetrope	Yes	Normal	None
Presbyopic	Myope	Yes	Reduced	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetrope	Yes	Reduced	None
Presbyopic	Hypermetrope	Yes	Normal	None

Further refinement

- Current state: $\text{If astigmatism = yes and ? then recommendation = hard}$
- Possible tests:

Age = Young	2/4
Age = Pre-presbyopic	1/4
Age = Presbyopic	1/4
Spectacle prescription = Myope	3/6
Spectacle prescription = Hypermetrope	1/6
Tear production rate = Reduced	0/6
Tear production rate = Normal	4/6

- Rule with best test added:

```
If astigmatism = yes
    and tear production rate = normal
then recommendation = hard
```

- Instances covered by modified rule:

Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	Yes	Normal	Hard
Young	Hypermetrope	Yes	Normal	hard
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetrope	Yes	Normal	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetrope	Yes	Normal	None

- Current state:

```
If astigmatism = yes
    and tear production rate = normal
    and ?
then recommendation = hard
```

- Possible tests:

Age = Young	2/2
Age = Pre-presbyopic	1/2
Age = Presbyopic	1/2
Spectacle prescription = Myope	3/3
Spectacle prescription = Hypermetrope	1/3

- Tie between the first and the fourth test

- We choose the one with greater coverage

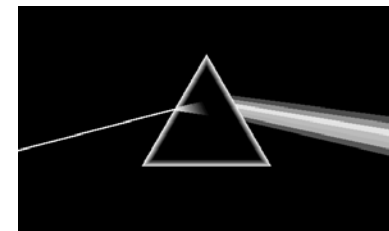
- Final rule:

```
If astigmatism = yes
    and tear production rate = normal
    and spectacle prescription = myope
then recommendation = hard
```
- Second rule for recommending “hard lenses”:
(built from instances not covered by first rule)

```
If age = young and astigmatism = yes
    and tear production rate = normal
then recommendation = hard
```

- These two rules cover all “hard lenses”:
 - Process is repeated with other two classes

```
For each class C
    Initialize E to the instance set
    While E contains instances in class C
        Create a rule R with an empty left-hand side that predicts class C
        Until R is perfect (or there are no more attributes to use) do
            For each attribute A not mentioned in R, and each value v,
                Consider adding the condition A = v to the left-hand side of R
                Select A and v to maximize the accuracy p/t
                (break ties by choosing the condition with the largest p)
            Add A = v to R
        Remove the instances covered by R from E
```



- PRISM with outer loop removed generates a decision list for one class
 - ♦ Subsequent rules are designed for rules that are not covered by previous rules
 - ♦ But: order doesn't matter because all rules predict the same class
- Outer loop considers all classes separately
 - ♦ No order dependence implied
- Problems: overlapping rules, default rule required

- Methods like PRISM (for dealing with one class) are *separate-and-conquer* algorithms:
 - ♦ First, identify a useful rule
 - ♦ Then, separate out all the instances it covers
 - ♦ Finally, “conquer” the remaining instances
- Difference to divide-and-conquer methods:
 - ♦ Subset covered by rule doesn't need to be explored any further