

CSDS 440: Machine Learning

Soumya Ray (he/him, sray@case.edu)

Olin 516

Office hours T, Th 11:15-11:45 or by appointment

Today

- Decision Tree Induction (Ch 3 Mitchell)
- Overfitting and overfitting control

Decision Tree Induction

- Given a set of examples, produce a decision tree
- Decision tree induction works using the idea of **recursive partitioning**
 - At each step, the algorithm will **choose an attribute test**
 - If no attribute looks good, return
 - The chosen test will partition the examples into disjoint partitions
 - The algorithm will then recursively call itself on each partition until
 - a partition only has data from one class (**pure** node) OR
 - it runs out of attributes

Choosing an Attribute

- Which attribute should we choose to test first?
 - Ideally, the one that is “most predictive” of the class label
 - i.e., the one that gives us the “most information” about what the label should be
- This idea is captured by the “(Shannon) entropy” of a random variable

Entropy of a Random Variable

- Suppose a random variable X has density $p(x)$. Its (Shannon) “entropy” is defined by:

$$\begin{aligned} H(X) &= E(-\log_2(p(X))) \\ &= -\sum_x p(X = x) \log_2(p(X = x)) \end{aligned}$$

- Note: $0\log(0) = 0$.

What's the connection?

- Entropy measures the *information content* of a random variable
- Suppose we treat the class variable, Y , as a random variable and measure its entropy
- Then we measure its entropy after partitioning the examples with an attribute X

The Entropy Connection

- The difference will be a measure of the “information gained” about Y by partitioning the examples with X
- So if we can choose the attribute X that maximizes this “information gain”, we have found what we needed

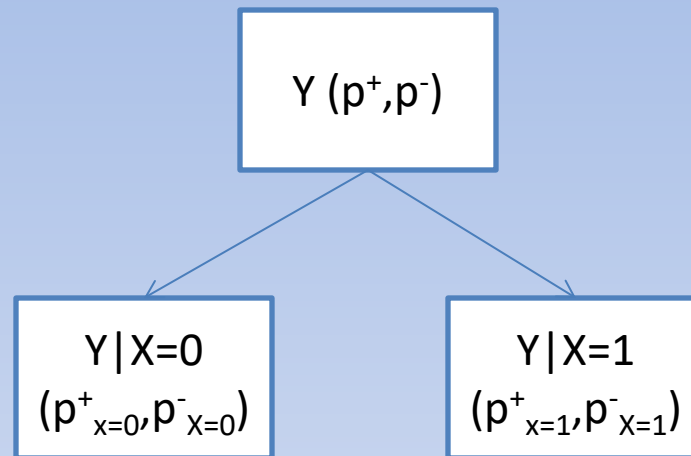
The class as a random variable

- Suppose at some point we have N training examples, of which pos are labeled “*positive*” and neg are labeled “*negative*” ($pos + neg = N$)
- We’ll treat the class label as a Bernoulli r.v. Y that takes value 1 with prob. $p^+ = pos/N$ and 0 with prob. $p^- = neg/N$
- Then $H(Y) = -p^+ \log_2(p^+) - p^- \log_2(p^-)$

Information Gain

- $IG(X)$ =reduction in entropy of the class label if the data is partitioned using X
- Suppose an attribute X takes two values 1 and 0. After partitioning, we get the quantities $p_{X=1}^+, p_{X=1}^-, p_{X=0}^+$ and $p_{X=0}^-$. Then,

Information Gain contd.



$$H(Y | X = 1) = -p_{X=1}^+ \log_2 p_{X=1}^+ - p_{X=1}^- \log_2 p_{X=1}^-$$

$$H(Y | X = 0) = -p_{X=0}^+ \log_2 p_{X=0}^+ - p_{X=0}^- \log_2 p_{X=0}^-$$

$$H(Y | X) = p(X = 1)H(Y | X = 1) + p(X = 0)H(Y | X = 0)$$

$$IG(X) = H(Y) - H(Y | X)$$

Nominal Attributes

- If X has v values:

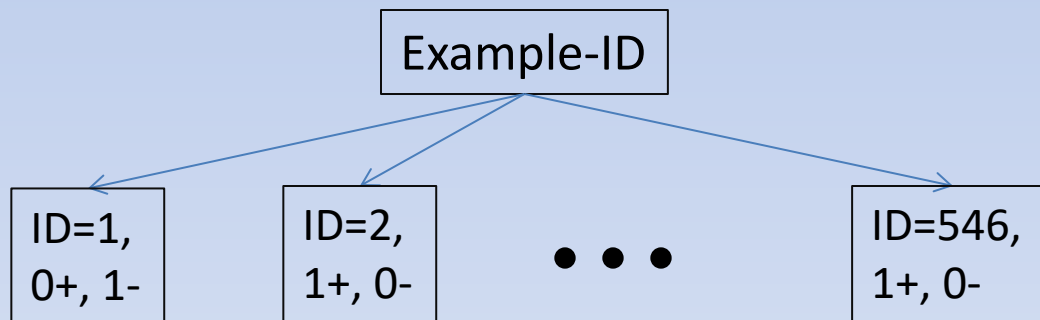
$$H(Y | X = v) = -p_{X=v}^+ \log_2 p_{X=v}^+ - p_{X=v}^- \log_2 p_{X=v}^-$$

$$H(Y | X) = \sum_v p(X = v) H(Y | X = v)$$

$$IG(X) = H(Y) - H(Y | X)$$

A Problem

- If an attribute has a lot of values, IG prefers it (resulting partitions tend to be pure)
- E.g., consider an “Example-ID” attribute



- This memorizes the data, so has perfect IG score

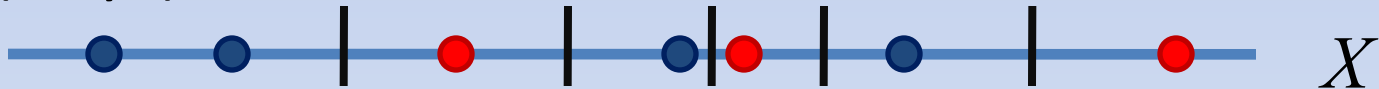
Fix: GainRatio

- Normalize IG with entropy of the attribute's distribution (computed from training data)

$$GR(X) = \frac{IG(X)}{H(X)}$$

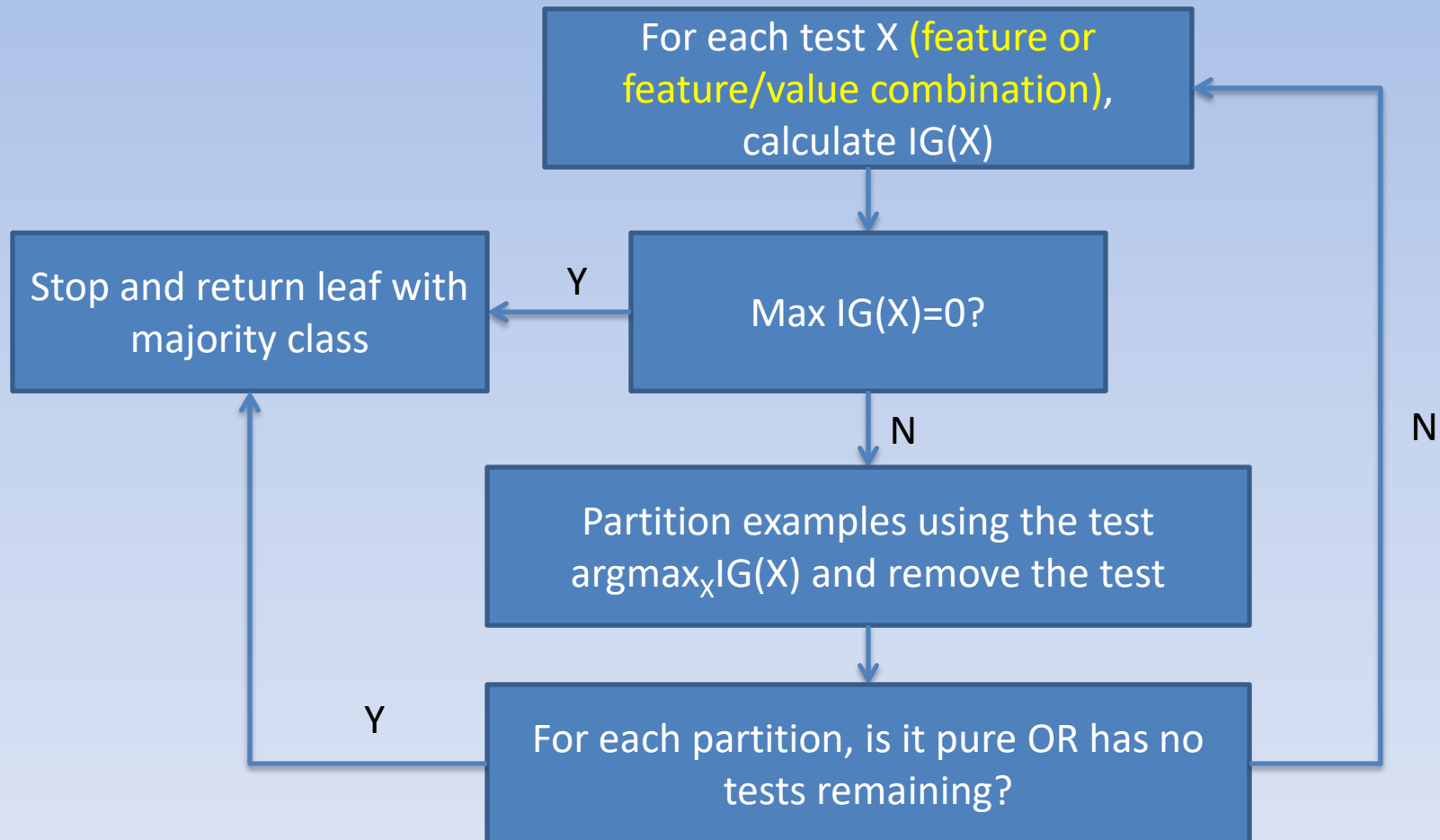
Continuous Attributes

- Cannot test for equality
- Consider all Boolean tests of the form $X \geq v$ (or $X \leq v$)
 - Only values of interest are those v that separate adjacent training examples with different classes (why?)



- Note: In this case, the attribute cannot be removed, though the test ((attribute, value) tuple) can be

ID3 Algorithm---Training phase



Example

| Color | Area | Shape | Class Label |
|-------|------|----------|-------------|
| red | 0.1 | circle | 1 |
| blue | 0.2 | triangle | 1 |
| green | 0.3 | triangle | 0 |
| green | 0.3 | circle | 0 |
| green | 0.4 | square | 0 |
| red | 0.4 | triangle | 1 |
| blue | 0.6 | circle | 0 |
| red | 0.7 | square | 0 |
| blue | 0.8 | square | 0 |

Example

| Color | Area | Shape | Class Label |
|-----------------|----------------|----------|--------------|
| red | 0.1 | circle | 1 |
| blue | 0.2 | triangle | 1 |
| green | 0.3 | triangle | 0 |
| green | 0.3 | circle | 0 |
| green | 0.4 | square | 0 |
| red | 0.4 | triangle | 1 |
| blue | 0.6 | circle | 0 |
| red | 0.7 | square | 0 |
| blue | 0.8 | square | 0 |

Example

| Color | Area | Shape | Class Label |
|-----------------|----------------|----------|--------------|
| red | 0.1 | circle | 1 |
| blue | 0.2 | triangle | 1 |
| green | 0.3 | triangle | 0 |
| green | 0.3 | circle | 0 |
| green | 0.4 | square | 0 |
| red | 0.4 | triangle | 1 |
| blue | 0.6 | circle | 0 |
| red | 0.7 | square | 0 |
| blue | 0.8 | square | 0 |