

# Decision Trees: Information Gain

These slides were assembled by Byron Boots, with grateful acknowledgement to Eric Eaton and the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution.

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

### Last Time: Basic Algorithm for Top-Down Learning of Decision Trees [ID3, C4.5 by Quinlan]

*node* = root of decision tree

Main loop:

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

How do we choose which attribute is best?

### Entropy



*H*(*X*) is the expected number of bits needed to encode a randomly drawn value of *X* (under most efficient code)

$$\sum_{i=1}^{n} P(X = i)(-\log_2 P(X = i))$$

### Entropy



*H*(*X*) is the expected number of bits needed to encode a randomly drawn value of *X* (under most efficient code)

#### Why? Information theory:

- Most efficient code assigns -log<sub>2</sub>P(X=i) bits to encode the message X=i
- So, expected number of bits to code one random *X* is:

$$\sum_{i=1}^{n} P(X = i)(-\log_2 P(X = i))$$

### 2-Class Cases:

Entropy 
$$H(x) = -\sum_{i=1}^{n} P(x=i) \log_2 P(x=i)$$

• What is the entropy of a group in which all examples belong to the same class?

$$- \text{ entropy} = -1 \log_2 1 = 0$$



Minimum

- What is the entropy of a group with 50% in either class?
  - entropy = -0.5  $\log_2 0.5 0.5 \log_2 0.5 = 1$



Maximum

Based on slide by Pedro Domingos

# Sample Entropy



- $\bullet~S$  is a sample of training examples
- $p_{\oplus}$  is the proportion of positive examples in S
- $\bullet \; p_{\ominus}$  is the proportion of negative examples in S
- $\bullet$  Entropy measures the impurity of S

$$H(S)\equiv -p_\oplus \log_2 p_\oplus - p_\ominus \log_2 p_\ominus$$

# **Information Gain**

- We want to determine which attribute in a given set of training feature vectors is most useful for discriminating between the classes to be learned.
- Information gain tells us how important a given attribute of the feature vectors is.
- We will use it to decide the ordering of attributes in the nodes of a decision tree.

Entropy *H*(*X*) of a random variable *X* 

$$H(X) = -\sum_{i=1}^{n} P(X=i) \log_2 P(X=i)$$

$$H(X|Y = v) = -\sum_{i=1}^{n} P(X = i|Y = v) \log_2 P(X = i|Y = v)$$

$$H(X|Y) = \sum_{v \in values(Y)} P(Y = v)H(X|Y = v)$$

$$I(X,Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

Entropy *H*(*X*) of a random variable *X* 

$$H(X) = -\sum_{i=1}^{n} P(X = i) \log_2 P(X = i)$$

Specific conditional entropy H(X|Y=v) of X given Y=v:

$$H(X|Y = v) = -\sum_{i=1}^{n} P(X = i|Y = v) \log_2 P(X = i|Y = v)$$

$$H(X|Y) = \sum_{v \in values(Y)} P(Y = v)H(X|Y = v)$$

$$I(X,Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

Entropy *H*(*X*) of a random variable *X* 

$$H(X) = -\sum_{i=1}^{n} P(X=i) \log_2 P(X=i)$$

Specific conditional entropy H(X|Y=v) of X given Y=v:

$$H(X|Y = v) = -\sum_{i=1}^{n} P(X = i|Y = v) \log_2 P(X = i|Y = v)$$

Conditional entropy H(X|Y) of X given Y :

$$H(X|Y) = \sum_{v \in values(Y)} P(Y = v)H(X|Y = v)$$

$$I(X,Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

Entropy *H*(*X*) of a random variable *X* 

$$H(X) = -\sum_{i=1}^{n} P(X = i) \log_2 P(X = i)$$

Specific conditional entropy H(X|Y=v) of X given Y=v:

$$H(X|Y = v) = -\sum_{i=1}^{n} P(X = i|Y = v) \log_2 P(X = i|Y = v)$$

Conditional entropy H(X|Y) of X given Y :

$$H(X|Y) = \sum_{v \in values(Y)} P(Y = v) H(X|Y = v)$$

Mututal information (aka Information Gain) of *X* and *Y*: I(X,Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)

# Information Gain

Information Gain is the expected reduction in entropy of target variable Y for data sample S, due to sorting

#### **Calculating Information Gain**

**Information Gain** = entropy(parent) – [average entropy(children)]



Based on slide by Pedro Domingos

## Entropy-Based Automatic Decision Tree Construction



Node 1 What feature should be used? What values?

#### Quinlan suggested information gain in his ID3 system

#### Using Information Gain to Construct a Decision Tree



## Sample Dataset (was Tennis Played?)

- Columns denote features X<sub>i</sub>
- Rows denote labeled instances  $\langle {m x}_i, y_i 
  angle$
- Class label denotes whether a tennis game was played

	Predictors			Response
Outlook	Temperature	Humidity	Wind	Class
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rain	Mild	High	Strong	No

 $\langle \boldsymbol{x}_i, y_i \rangle$ 

#### Selecting the Next Attribute

Which attribute is the best classifier?



#### Selecting the Next Attribute

Which attribute is the best classifier?





Which attribute should be tested here?

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

 $Gain(S_{sunny}, Humidity) = .970 - (3/5) 0.0 - (2/5) 0.0 = .970$ 

 $Gain(S_{sunny}, Temperature) = .970 - (2/5) 0.0 - (2/5) 1.0 - (1/5) 0.0 = .570$ 

 $Gain(S_{sunny}, Wind) = .970 - (2/5) 1.0 - (3/5) .918 = .019$ 

### Which Tree Should We Output?



- ID3 performs heuristic search through space of decision trees
- It stops at smallest acceptable tree. Why?