

Knowledge Discovery and Data Mining

Unit # 4

Acknowledgement

- Most of the slides in this presentation are taken from course slides provided by
 - Han and Kimber (Data Mining Concepts and Techniques) and
 - Tan, Steinbach and Kumar (Introduction to Data Mining)

Alternative Splitting Criteria based on Entropy

- Entropy at a given node t :

$$Entropy(t) = -\sum_j p(j|t) \log p(j|t)$$

(NOTE: $p(j|t)$ is the relative frequency of class j at node t).

- Measures homogeneity of a node.
 - Maximum ($\log n_c$) when records are equally distributed among all classes implying least information
 - Minimum (0.0) when all records belong to one class, implying most information
- Entropy based computations are similar to the GINI index computations
- Hint: $\log_2 p = \ln p / \ln(2)$

Entropy in a nut-shell



Low Entropy



High Entropy

Examples for computing Entropy

$$Entropy(t) = -\sum_j p(j|t) \log_2 p(j|t)$$

| | |
|----|----------|
| C1 | 0 |
| C2 | 6 |

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Entropy = -0 \log 0 - 1 \log 1 = -0 - 0 = 0$$

| | |
|----|----------|
| C1 | 1 |
| C2 | 5 |

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$Entropy = -(1/6) \log_2 (1/6) - (5/6) \log_2 (5/6) = 0.65$$

| | |
|----|----------|
| C1 | 2 |
| C2 | 4 |

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$Entropy = -(2/6) \log_2 (2/6) - (4/6) \log_2 (4/6) = 0.92$$

Splitting Criteria based on Classification Error

- Classification error at a node t :

$$Error(t) = 1 - \max_i P(i|t)$$

- Measures misclassification error made by a node.
 - Maximum $(1 - 1/n_c)$ when records are equally distributed among all classes, implying least interesting information
 - Minimum (0.0) when all records belong to one class, implying most interesting information

Examples for Computing Error

$$Error(t) = 1 - \max_i P(i | t)$$

| | |
|----|----------|
| C1 | 0 |
| C2 | 6 |

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Error = 1 - \max(0, 1) = 1 - 1 = 0$$

| | |
|----|----------|
| C1 | 1 |
| C2 | 5 |

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$Error = 1 - \max(1/6, 5/6) = 1 - 5/6 = 1/6$$

| | |
|----|----------|
| C1 | 2 |
| C2 | 4 |

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$Error = 1 - \max(2/6, 4/6) = 1 - 4/6 = 1/3$$

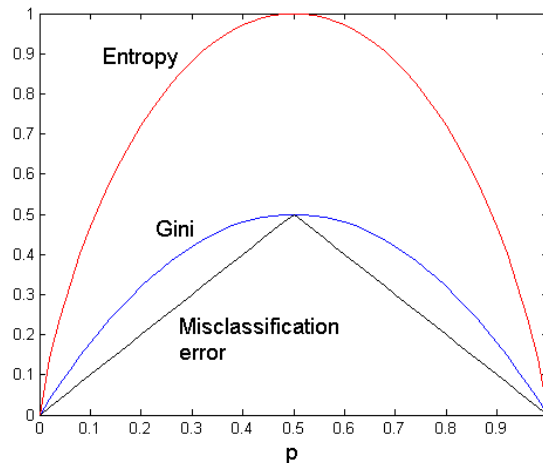
Sajjad Haider

Fall 2013

7

Comparison among Splitting Criteria

For a 2-class problem:



Sajjad Haider

Fall 2013

8

Example

| Attribute 1 | Attribute 2 | Attribute 3 | Class |
|-------------|-------------|-------------|-------|
| A | 70 | T | C1 |
| A | 90 | T | C2 |
| A | 85 | F | C2 |
| A | 95 | F | C2 |
| A | 70 | F | C1 |
| B | 90 | T | C1 |
| B | 78 | F | C1 |
| B | 65 | T | C1 |
| B | 75 | F | C1 |
| C | 80 | T | C2 |
| C | 70 | T | C2 |
| C | 80 | F | C1 |
| C | 80 | F | C1 |
| C | 96 | F | C1 |

Sajjad Haider

Fall 2013

9

Example II

| Height | Hair | Eyes | Class |
|--------|-------|-------|-------|
| Short | Blond | Blue | + |
| Tall | Blond | Brown | - |
| Tall | Red | Blue | + |
| Short | Dark | Blue | - |
| Tall | Dark | Blue | - |
| Tall | Blond | Blue | + |
| Tall | Dark | Brown | - |
| Short | Blond | Brown | - |

Sajjad Haider

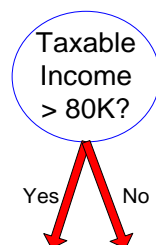
Fall 2013

10

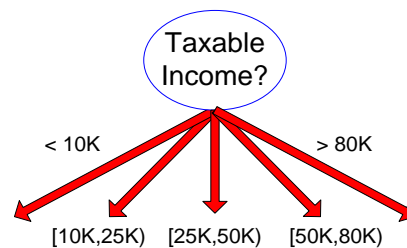
Splitting Based on Continuous Attributes

- Different ways of handling
 - **Discretization** to form an ordinal categorical attribute
 - Static – discretize once at the beginning
 - Dynamic – ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.
 - **Binary Decision**: ($A < v$) or ($A \geq v$)
 - consider all possible splits and finds the best cut
 - can be computationally intensive

Splitting Based on Continuous Attributes (Cont'd)



(i) Binary split



(ii) Multi-way split

Continuous Attributes: Computing GINI Index

- For efficient computation: for each attribute,
 - Sort the attribute on values
 - Linearly scan these values, each time updating the count matrix and computing gini index
 - Choose the split position that has the least gini index

| Cheat | No | No | No | Yes | Yes | Yes | No | No | No | No | |
|-----------------|----------------|-------|-------|-------|-------|-------|--------------|-------|-------|-------|-------|
| | Taxable Income | | | | | | | | | | |
| Sorted Values | 60 | 70 | 75 | 85 | 90 | 95 | 100 | 120 | 125 | 220 | |
| Split Positions | 55 | 65 | 72 | 80 | 87 | 92 | 97 | 110 | 122 | 172 | 230 |
| | <=> | <=> | <=> | <=> | <=> | <=> | <=> | <=> | <=> | <=> | <=> |
| Yes | 0 3 | 0 3 | 0 3 | 0 3 | 1 2 | 2 1 | 3 0 | 3 0 | 3 0 | 3 0 | 3 0 |
| No | 0 7 | 1 6 | 2 5 | 3 4 | 3 4 | 3 4 | 3 4 | 4 3 | 5 2 | 6 1 | 7 0 |
| Gini | 0.420 | 0.400 | 0.375 | 0.343 | 0.417 | 0.400 | <u>0.300</u> | 0.343 | 0.375 | 0.400 | 0.420 |