11413696

## Data Mining

**Assignment Four**        10 points
*Due: Thursday Oct. 19 @ 11:50 PM*

**Question 1 (6 points).** Consider the training examples shown in the Table below for a binary classification problem.

| Instance | A1 | A2 | A3 | Class |
|----------|----|----|-----|-------|
| 1 | T | T | 1.0 | Yes |
| 2 | T | T | 6.0 | Yes |
| 3 | T | F | 4.0 | No |
| 4 | F | F | 7.0 | Yes |
| 5 | F | T | 8.0 | Yes |
| 6 | F | F | 5.0 | No |
| 7 | F | F | 3.0 | No |
| 8 | F | F | 7.0 | No |
| 9 | F | T | 8.0 | No |

(a) What is the entropy of this collection of training examples with respect to the class attribute?
(b) What are the information gains of A1 and A2 relative to these training examples?
(c) For A3, which is a continuous attribute, compute the information gain for every possible split.
(d) What is the best split (among A1, A2, and A3) according to the information gain?
(e) What is the best split (between A1 and A2) according to the misclassification error rate?
(f) What is the best split (between A1 and A2) according to the Gini index?

**Question 2 (4 points).** Consider splitting a parent node P into two child nodes, C1 and C2, using some attribute test condition. The composition of labeled training instances at every node is summarized in the Table below.

| | P | C1 | C2 |
|--------|---|----|----|
| Class 0 | 7 | 3 | 4 |
| Class 1 | 3 | 0 | 3 |

(a) Calculate the Gini index and misclassification error rate of the parent node P.
(b) Calculate the weighted Gini index of the child nodes. Would you consider this attribute test condition if Gini is used as the impurity measure?
(c) Calculate the weighted misclassification rate of the child nodes. Would you consider this attribute test condition if misclassification rate is used as the impurity measure?

Good luck!

---

**1**

Entropy at a given node $t$

$$Entropy = -\sum_{i=0}^{c-1} p_i(t) log_2 p_i(t)$$

Where $p_i(t)$ is the frequency of class $i$ at node $t$, and $c$ is the total number of classes

**3**

Information Gain:

$$Gain_{split} = Entropy(p) - \sum_{i=1}^{k} \frac{n_i}{n} Entropy(i)$$

Parent Node, p is split into k partitions (children)
$n_i$ is number of records in child node i

– Choose the split that achieves most reduction (maximizes GAIN)

**5**

**Measure of Impurity: Classification Error**

Classification error at a node $t$

$$Error(t) = 1 - \max_i [p_i(t)]$$

– Maximum of $1 - 1/c$ when records are equally distributed among all classes, implying the least interesting situation
– Minimum of 0 when all records belong to one class, implying the most interesting situation

**7**

**Misclassification Error vs Gini Index**



**9**

Gini Index for a given node t :

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

– For 2-class problem (p, 1 – p):
  ◆ GINI = 1 – p² – (1 – p)² = 2p (1-p)



**Binary Attributes: Computing GINI Index**

▫ Splits into two partitions (child nodes)
▫ Effect of Weighing partitions:
  – Larger and purer partitions are sought



**2**



**Problem with large number of partitions**

▫ Node impurity measures tend to prefer splits that result in large number of partitions, each being small but pure

– Customer ID has highest information gain because entropy for all the children is zero

**Computing Error of a Single Node**

$$Error(t) = 1 - \max_i[p_i(t)]$$

| C1 | 0 | P(C1) = 0/6 = 0   P(C2) = 6/6 = 1 |
| C2 | 6 | Error = 1 – max (0, 1) = 1 – 1 = 0 |

| C1 | 1 | P(C1) = 1/6   P(C2) = 5/6 |
| C2 | 5 | Error = 1 – max (1/6, 5/6) = 1 – 5/6 = 1/6 |

| C1 | 2 | P(C1) = 2/6   P(C2) = 4/6 |
| C2 | 4 | Error = 1 – max (2/6, 4/6) = 1 – 4/6 = 1/3 |

▫ Gini Index for a given node $t$

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

Where $p_i(t)$ is the frequency of class $i$ at node $t$, and $c$ is the total number of classes

– Maximum of $1 - 1/c$ when records are equally distributed among all classes, implying the least beneficial situation for classification
– Minimum of 0 when all records belong to one class, implying the most beneficial situation for classification

**Computing Gini Index of a Single Node**

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

| C1 | 0 | P(C1) = 0/6 = 0   P(C2) = 6/6 = 1 |
| C2 | 6 | Gini = 1 – P(C1)² – P(C2)² = 1 – 0 – 1 = 0 |

| C1 | 1 | P(C1) = 1/6   P(C2) = 5/6 |
| C2 | 5 | Gini = 1 – (1/6)² – (5/6)² = 0.278 |

| C1 | 2 | P(C1) = 2/6   P(C2) = 4/6 |
| C2 | 4 | Gini = 1 – (2/6)² – (4/6)² = 0.444 |

---

## Question 1

**Part A:**

$Entropy(t) = -\Sigma \left[ P(t) * log2(P(t)) \right] = -\left[ \left(\frac{4}{9}\right) * log2\left(\frac{4}{9}\right) + \left(\frac{5}{9}\right) * log2\left(\frac{5}{9}\right) \right] = 0.99107605983$

**Part B:**
*"Yes" class labels*

$Entropy(A1) = -\Sigma \left[ P(t) * log2(P(t)) \right] = -\left[ \left(\frac{2}{4}\right) * log2\left(\frac{2}{4}\right) + \left(\frac{2}{4}\right) * log2\left(\frac{2}{4}\right) \right] = 1$

$Entropy(A2) = -\Sigma \left[ P(t) * log2(P(t)) \right] = -\left[ \left(\frac{4}{4}\right) * log2\left(\frac{4}{4}\right) + 0 \right] = 0$

*"No" class labels*

$Entropy(A1) = -\Sigma \left[ P(t) * log2(P(t)) \right] = -\left[ \left(\frac{1}{5}\right) * log2\left(\frac{1}{5}\right) + \left(\frac{4}{5}\right) * log2\left(\frac{4}{5}\right) \right] = 0.72192809488$

$Entropy(A2) = -\Sigma \left[ P(t) * log2(P(t)) \right] = -\left[ \left(\frac{2}{5}\right) * log2\left(\frac{2}{5}\right) + \left(\frac{3}{5}\right) * log2\left(\frac{3}{5}\right) \right] = 0.97095059445$

$Gainsplit(A1) = Entropy(T) - \Sigma\left(\frac{T_i}{T}\right) * Entropy(i) = 0.99107605983 - \left[ \left(\frac{4}{9}\right) * (1) + \left(\frac{5}{9}\right) * (0.72192809488) \right] = 0.14556045156$

$Gainsplit(A2) = Entropy(T) - \Sigma\left(\frac{T_i}{T}\right) * Entropy(i) = 0.99107605983 - \left[ \left(\frac{4}{9}\right) * (0) + \left(\frac{5}{9}\right) * (0.97095059445) \right] = 0.45165906291$

**Part C:** Type equation here.

| Sorted Values | 1 | 3 | 4 | 5 | 6 | 7 | 8 | |
|---|---|---|---|---|---|---|---|---|
| Split Positions | 0.5 | 2 | 3.5 | 4.5 | 5.5 | 6.5 | 7.5 | 8.5 |
| | <= \| > | <= \| > | <= \| > | <= \| > | <= \| > | <= \| > | <= \| > | <= \| > |
| Yes | 0 \| 4 | 1 \| 3 | 1 \| 3 | 1 \| 3 | 1 \| 3 | 2 \| 2 | 3 \| 1 | 4 \| 0 |
| No | 0 \| 5 | 0 \| 5 | 1 \| 4 | 2 \| 3 | 3 \| 2 | 3 \| 2 | 4 \| 1 | 5 \| 0 |
| Gain | $Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{4}{9}\right) * log2\left(\frac{4}{9}\right) + \left(\frac{5}{9}\right) * log2\left(\frac{5}{9}\right)] = 0.99107605983$<br><br>$Gainsplit(t) = 0$ | $\leq: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{1}{1}\right) * log2\left(\frac{1}{1}\right)] (0) * log2(0)] = 0$<br><br>$>: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{3}{8}\right) * log2\left(\frac{3}{8}\right) + \left(\frac{5}{8}\right) * log2\left(\frac{5}{8}\right)] = 0.95443400292$<br><br>$Weighted\ Average: \left[\left(\frac{1}{9}\right) * (0)\right] + \left[\left(\frac{8}{9}\right) * (0.95443400292)\right] = 0.84838578037$<br><br>$Gainsplit(t) = 0.99107605983 - (0.84838578037) = 0.14269027946$ | $\leq: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{1}{1}\right) * log2\left(\frac{1}{1}\right)] = 0$<br><br>$>: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{3}{7}\right) * log2\left(\frac{3}{7}\right) + \left(\frac{4}{7}\right) * log2\left(\frac{4}{7}\right)] = 0.98522813603$<br><br>$Weighted\ Average: \left[\left(\frac{2}{9}\right) * (1)\right] + \left[\left(\frac{7}{9}\right) * (0.98522813603)\right] = 0.98851077246$<br><br>$Gainsplit(t) = 0.99107605983 - (0.98851077246) = 0.00256528737$ | $\leq: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{1}{2}\right) * log2\left(\frac{1}{2}\right) + \left(\frac{2}{2}\right) * log2\left(\frac{2}{2}\right)] = 1$<br><br>$>: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{3}{6}\right) * log2\left(\frac{3}{6}\right) + \left(\frac{3}{6}\right) * log2\left(\frac{3}{6}\right)] = 1$<br><br>$Weighted\ Average: \left[\left(\frac{3}{9}\right) * (0.91829583405)\right] + \left[\left(\frac{6}{9}\right) * (1)\right] = 0.97276527801$<br><br>$Gainsplit(t) = 0.99107605983 - (0.97276527801) = 0.01831078182$ | $\leq: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{1}{4}\right) * log2\left(\frac{1}{4}\right) + \left(\frac{3}{4}\right) * log2\left(\frac{3}{4}\right)] = 0.81127812445$<br><br>$>: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{3}{5}\right) * log2\left(\frac{3}{5}\right) + \left(\frac{2}{5}\right) * log2\left(\frac{2}{5}\right)] = 0.97095059445$<br><br>$Weighted\ Average: \left[\left(\frac{4}{9}\right) * (0.81127812445)\right] + \left[\left(\frac{5}{9}\right) * (0.97095059445)\right] = 0.89998505222$<br><br>$Gainsplit(t) = 0.99107605983 - (0.89998505222) = 0.09109100761$ | $\leq: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{2}{5}\right) * log2\left(\frac{2}{5}\right) + \left(\frac{3}{5}\right) * log2\left(\frac{3}{5}\right)] = 0.97095059445$<br><br>$>: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{2}{4}\right) * log2\left(\frac{2}{4}\right) + \left(\frac{2}{4}\right) * log2\left(\frac{2}{4}\right)] = 1$<br><br>$Weighted\ Average: \left[\left(\frac{5}{9}\right) * (0.97095059445)\right] + \left[\left(\frac{4}{9}\right) * (1)\right] = 0.98386144136$<br><br>$Gainsplit(t) = 0.99107605983 - (0.98386144136) = 0.00721461847$ | $\leq: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{3}{7}\right) * log2\left(\frac{3}{7}\right) + \left(\frac{4}{7}\right) * log2\left(\frac{4}{7}\right)] = 0.98522813603$<br><br>$>: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{1}{2}\right) * log2\left(\frac{1}{2}\right) + \left(\frac{1}{2}\right) * log2\left(\frac{1}{2}\right)] = 1$<br><br>$Weighted\ Average: \left[\left(\frac{7}{9}\right) * (0.98522813603)\right] + \left[\left(\frac{2}{9}\right) * (1)\right] = 0.98851077246$<br><br>$Gainsplit(t) = 0.99107605983 - (0.98851077246) = 0.00256528737$ | $\leq: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = -[\left(\frac{4}{9}\right) * log2\left(\frac{4}{9}\right) + \left(\frac{5}{9}\right) * log2\left(\frac{5}{9}\right)] = 0.99107605983$<br><br>$>: Entropy(t) = -\Sigma [P(t) * log2(P(t))] = 0$<br><br>$Weighted\ Average: \left[\left(\frac{9}{9}\right) * (0.99107605983)\right] + [0] = 0.99107605983$<br><br>$Gainsplit(t) = 0.99107605983 - (0.99107605983) = 0$ |

**Part D:**

According to the information gain, A2 is the best split between A1, A2, and A3. We prefer to have an information gain value that is close to 1. It is hard to tell with A3 because of the large number of partitions. Sure, the nodes are pure but they are small measures.

**Part E:**

According to the misclassification error rate, $Misclassification\ Error(t) = 1 - \max(p_i(t))$, A2 is the best split between A1 and A2 since it has a lower classification error when compared to its counterpart. The lower the classification error, the more accurate of a sample set.

$Misclassification\ Error(A1) = 1 - \max(p_i(t)) = 1 - \max\left[\left(\frac{3}{9}\right),\left(\frac{6}{9}\right)\right] = 1 - \left(\frac{6}{9}\right) = 0.3333$

$Misclassification\ Error(A2) = 1 - \max(p_i(t)) = 1 - \max\left[\left(\frac{2}{9}\right),\left(\frac{7}{9}\right)\right] = 1 - \left(\frac{7}{9}\right) = 0.2222$

**Part F:**

According to the gini index, $Gini(t) = 1 - \sum p_i(t)^2$, A2 is the best split between A1 and A2 since it has a lower gini index when compared to its counterpart. The lower the gini index, the closer a node is to purity.

$Gini(A1) = 1 - \sum p_i(t)^2 = 1 - \left[\left(\frac{3}{9}\right)^2 + \left(\frac{6}{9}\right)^2\right] = 0.44444444444$

$Gini(A2) = 1 - \sum p_i(t)^2 = 1 - \left[\left(\frac{2}{9}\right)^2 + \left(\frac{7}{9}\right)^2\right] = 0.34567901234$

*Question 2*

**Part A:**

$Gini(t) = 1 - \sum p_i(t)^2,\ Gini(P) = 1 - \sum p_i(t)^2 = 1 - \left[\left(\frac{7}{10}\right)^2 + \left(\frac{3}{10}\right)^2\right] = 0.42$

$Misclassification\ Error\ Rate(t) = 1 - \max(p_i(t)),\ Misclassification\ Error\ Rate(P) = 1 - \max(p_i(t)) = 1 - \max\left[\left(\frac{7}{10}\right),\left(\frac{3}{10}\right)\right] = 1 - \left(\frac{7}{10}\right) = 0.30$

**Part B:**

$Gini(t) = 1 - \sum p_i(t)^2,\ Gini(C1) = 1 - \sum p_i(t)^2 = 1 - \left[\left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2\right] = 0$

$Gini(t) = 1 - \sum p_i(t)^2,\ Gini(C2) = 1 - \sum p_i(t)^2 = 1 - \left[\left(\frac{4}{7}\right)^2 + \left(\frac{3}{7}\right)^2\right] = 0.48979591836$

$Weighted\ Gini\ Index: \left[\left(\frac{3}{10}\right)*(0)\right] + \left[\left(\frac{7}{10}\right)*(0.48979591836)\right] = 0.34285714285$

If gini is used as the impurity measure, we would consider this attribute test condition. Consider looking at the gini index of the parent node P before the split. The gini impurity of each child node after the split then calculating the weighted gini gives a smaller output (which is what we want) in comparison to its parent P. That is the reason we would consider using this attribute test condition.

**Part C:**

$Misclassification\ Error\ Rate(t) = 1 - \max(p_i(t)),\ Misclassification\ Error\ Rate(C1) = 1 - \max(p_i(t)) = 1 - \max\left[\left(\frac{3}{3}\right),\left(\frac{0}{3}\right)\right] = 1 - \left(\frac{3}{3}\right) = 0$

$Misclassification\ Error\ Rate(t) = 1 - \max(p_i(t)),\ Misclassification\ Error\ Rate(C2) = 1 - \max(p_i(t)) = 1 - \max\left[\left(\frac{4}{7}\right),\left(\frac{3}{7}\right)\right] = 1 - \left(\frac{4}{7}\right) = 0.42857142857$

$Weighted\ Misclassification\ Error\ Rate: \left[\left(\frac{3}{10}\right)*(0)\right] + \left[\left(\frac{7}{10}\right)*(0.42857142857)\right] = 0.29999999999$

If misclassification error rate is used as the impurity measure, we would NOT consider this attribute test condition. Consider looking at the misclassification error rate of the parent node P before the split. The misclassification impurity of each child node after the split then calculating the weighted misclassification error gives an output equal (there is no drop in the error rate) to its parent P. Therefore, it would bring us no benefit to consider this attribute test condition.