

Συστήματα Ανακάλυψης Γνώσης από Βάσεις Δεδομένων

Εργασία 01 - HW1 - Decision Trees

Ονοματεπώνυμο: Παναγιώτης Γιαννουτάκης

AM: 12/38

(α)

$$\text{GINI}(\text{parent}) = 1 - (4/9)^2 - (5/9)^2 = 1 - 16/81 - 25/81 = 1 - 41/81 = 0.493$$

(β)

ID

$$\text{GINI}(\text{id1}) = 1 - (1/1)^2 = 1 - 1 = 0$$

$$\text{GINI}(\text{id2}) = 1 - (1/1)^2 = 1 - 1 = 0$$

...

$$\text{GINI}(\text{id9}) = 1 - (1/1)^2 = 1 - 1 = 0$$

$$\text{GINI}(\text{ID_split}) = 1*0 + 1*0 + \dots + 1*0 = 0$$

$$\text{GAIN}(\text{ID}) = 0.493 - 0 = 0.493$$

A1

$$\text{GINI}(T) = 1 - (3/4)^2 - (1/4)^2 = 1 - 9/16 - 1/16 = 1 - 10/16 = 0.375$$

$$\text{GINI}(F) = 1 - (1/5)^2 - (4/5)^2 = 1 - 1/25 - 16/25 = 1 - 17/25 = 0.32$$

$$\text{GINI}(a1_split) = (4/9) * 0.375 + (5/9) * 0.32 = 0.166 + 0.177 = 0.343$$

$$\text{GAIN}(a1) = 0.493 - 0.343 = 0.149$$

A2

$$\text{GINI}(T) = 1 - (2/5)^2 - (3/5)^2 = 1 - 4/25 - 9/25 = 1 - 13/25 = 0,48$$

$$\text{GINI}(F) = 1 - (2/4)^2 - (2/4)^2 = 1 - 4/16 - 4/16 = 1 - 8/16 = 0,5$$

$$\text{GINI}(a2_split) = 5/9 * 0,48 + (4/9) * 0,5 = 0,266 + 0,222 = 0,488$$

$$\text{GAIN}(a2) = 0,493 - 0,488 = 0,005$$

Δεν επιλέγουμε το $\text{GAIN}(\text{ID})$ γιατί σε μελλοντικές καταχωρίσεις καμία καταχώριση δεν θα μπει σε κάποιον από τους κόμβους που ήδη υπάρχουν αλλά θα δημιουργείται κάθε φορά ένας νέος.

(γ)

Ο πίνακας με τις διασπάσεις και τα GAIN είναι ο παρακάτω:

		+	-	+	-	+	+-	-	
		1	3	4	5	6	7	8	
		0,5	2	3,5	4,5	5,5	6,5	7,5	8,5
		<	>	<	>	<	>	<	>
+		0	4	1	3	1	3	2	2
-		0	5	0	5	1	4	1	4
GAIN		0	0,001	0,001	0,046	0,005	0,012	0,049	0

Ακολουθούν οι πράξεις που έγιναν στο χαρτί:

0.5

$$GINI(L) = 1 - \left(\frac{1}{6}\right)^2 - \left(\frac{2}{3}\right)^2 = 1$$

$$GINI(R) = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{3}\right)^2 = 1 - \frac{16}{36} - \frac{25}{36} = 1 - \frac{41}{36} = 1 - 0.506 \approx 0.493$$

$$GINI(0.5 \text{ split}) = 0.1 + \frac{2}{3} \cdot 0.493 = 0.493$$

$$GAIN(0.5) = 0.493 - 0.493 = 0$$

2

$$GINI(L) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 1 - 1 - 0 = 0$$

$$GINI(R) = 1 - \left(\frac{3}{8}\right)^2 - \left(\frac{5}{8}\right)^2 = 1 - \frac{9}{64} - \frac{25}{64} = 1 - \frac{34}{64} = 1 - 0.531 = 0.469$$

$$GINI(2 \text{ split}) = \frac{1}{2} \cdot 0 + \frac{5}{8} \cdot 0.469 = 0.293$$

$$GAIN(2) = 0.493 - 0.293 = 0.2$$

3.5

$$GINI(L) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 1 - 0.5 = 0.5$$

$$GINI(R) = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = 1 - \frac{9}{49} - \frac{16}{49} = 1 - \frac{25}{49} = 1 - 0.510 = 0.49$$

$$GINI(3.5 \text{ split}) = \frac{2}{7} \cdot 0.5 + \frac{5}{7} \cdot 0.49 = 0.143 + 0.381 = 0.524$$

$$GAIN(3.5) = 0.493 - 0.524 = -0.031$$

4.5

4,5

$$GINI(L) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 1 - \frac{4}{9} - \frac{1}{9} = 1 - \frac{5}{9} = 0,445$$

$$GINI(R) = 1 - \left(\frac{2}{6}\right)^2 - \left(\frac{4}{6}\right)^2 = 1 - \frac{4}{36} - \frac{16}{36} = 1 - \frac{20}{36} = 1 - 0,555 = 0,445$$

$$GINI(4,5-split) = 3 \cdot 0,445 + 6 \cdot 0,445 = 0,151 + 0,296 = 0,447$$

$$GAIN(4,5) = 0,493 - 0,447 = 0,046$$

5,5

$$GINI(L) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 1 - \frac{4}{25} - \frac{9}{25} = 1 - \frac{13}{25} = 1 - 0,52 = 0,48$$

$$GINI(R) = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 1 - \frac{4}{16} - \frac{4}{16} = 1 - \frac{8}{16} = 0,5$$

$$GINI(5,5-split) = 5 \cdot 0,48 + 4 \cdot 0,5 = 0,266 + 0,222 = 0,488$$

$$GAIN(5,5) = 0,493 - 0,488 = 0,005$$

6,5

$$GINI(L) = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = 1 - \frac{9}{36} - \frac{9}{36} = 1 - \frac{18}{36} = 0,5$$

$$GINI(R) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 1 - \frac{1}{9} - \frac{4}{9} = 1 - \frac{5}{9} = 0,445$$

$$GINI(6,5-split) = 6 \cdot 0,5 + 3 \cdot 0,445 = 0,333 + 0,148 = 0,481$$

$$GAIN(6,5) = 0,493 - 0,481 = 0,012$$

7,5

$$GINI(L) = 1 - \left(\frac{4}{8}\right)^2 - \left(\frac{4}{8}\right)^2 = 1 - \frac{16}{64} - \frac{16}{64} = 1 - \frac{32}{64} = 1 - 0,5 = 0,5$$

$$GINI(R) = 1 - \left(\frac{0}{1}\right)^2 - \left(\frac{1}{1}\right)^2 = 1 - 0 - 1 = 0$$

$$GINI(7,5-split) = 8 \cdot 0,5 + 1 \cdot 0 = 0,444$$

$$GAIN(7,5) = 0,493 - 0,444 = 0,049$$

8.5

$$GINI(L) = 1 - (40\%)^2 - (5\%)^2 = 1 - \frac{16}{81} - \frac{25}{81} = 1 - \frac{41}{81} = 1 - 0.506172839 = 0.493827161 \approx 0.493$$

$$GINI(R) = 1 - (0\%)^2 - (0\%)^2 = 1$$

$$GINI(8.5-split) = 9 \cdot 0.494 + 0 \cdot 1 = 0.494$$

$$GAIN(8.5) = 0.493 - 0.494 = 0$$

Η καλύτερη διάσπαση προκύπτει εκεί που υπάρχει το μεγαλύτερο GAIN, δηλαδή στην διάσπαση 2.

(δ)

Η ρίζα του δέντρου θα είναι η στήλη a1 γιατί έχει το μεγαλύτερο GAIN από όλες τις στήλες (εκτός από την στήλη ID που δεν την μετράμε). Το δέντρο είναι το παρακάτω:

|--- T (3,1) {+}

A1 ----- |

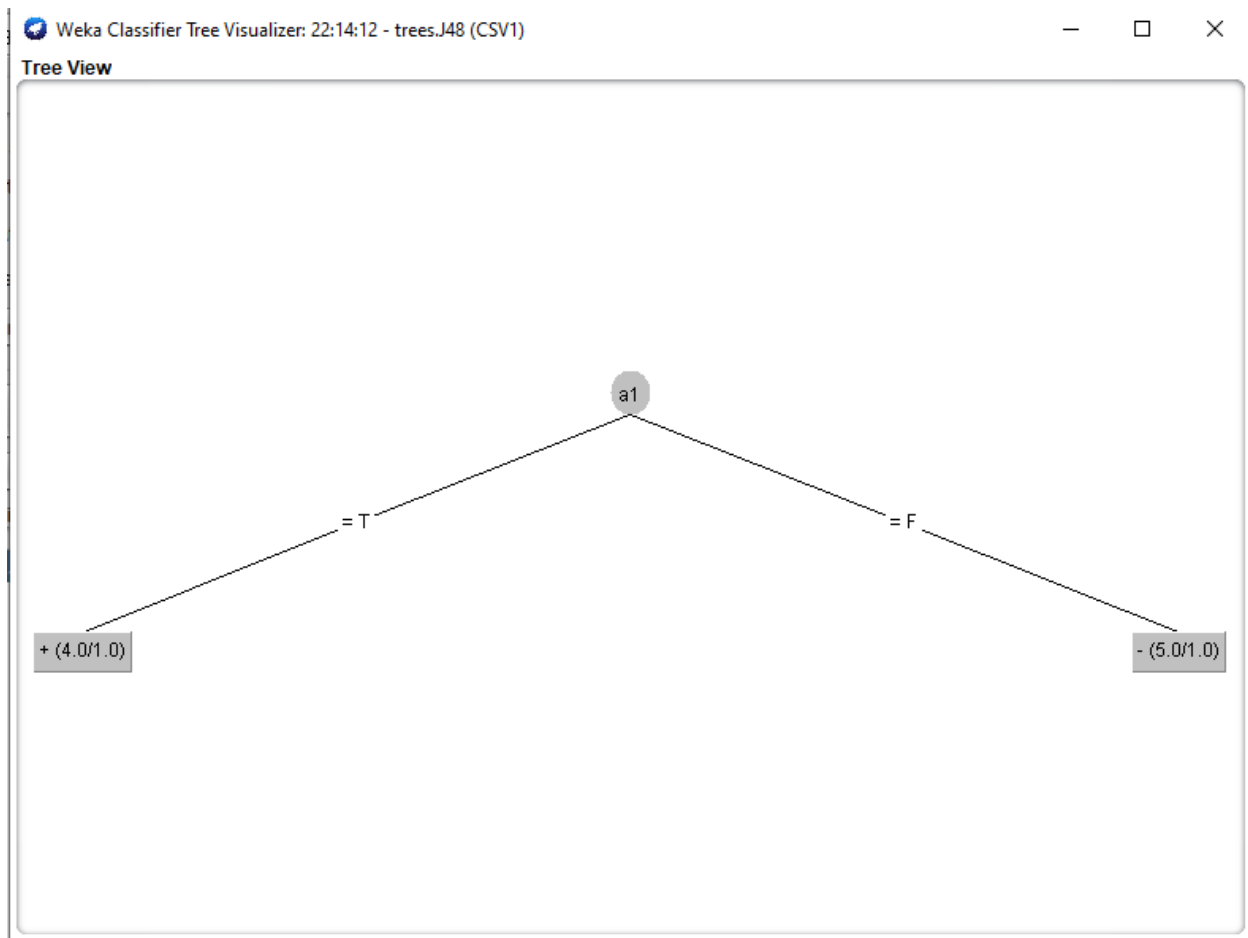
|--- F (1,4) {-}

Το ποσοστό των εγγραφών που κατηγοριοποιούνται σωστά είναι 7/9, δηλαδή 77,7 % και το υπόλοιπο 2/9 εγγραφές, δηλαδή 22,3% κατηγοριοποιούνται λάθος.

(ε)

Το WEKA διάλεξε την μεταβλητή a1 όπως υπολογίστηκε και προηγουμένως με το χέρι.

Η εικόνα του δέντρου είναι η παρακάτω:



(στ)

Με το WEKA πετυχαίνουμε 75% επιτυχία.

```
Classifier output
Scheme:      weka.classifiers.misc.InputMappedClassifier -I -trim -w weka.classifiers.trees.J48
Relation:    CSV1
Instances:   9
Attributes:  5
              ID
              a1
              a2
              a3
              class
Test mode:   user supplied test set: size unknown (reading incrementally)

=== Classifier model (full training set) ===

InputMappedClassifier:

J48 pruned tree
-----

a1 = T: + (4.0/1.0)
a1 = F: - (5.0/1.0)

Number of Leaves :    2
Size of the tree :    3

Attribute mappings:

Model attributes      Incoming attributes
-----
(nominal) ID          --> 1 (nominal) ID
(nominal) a1          --> 2 (nominal) a1
(nominal) a2          --> 3 (nominal) a2
(numeric) a3          --> 4 (numeric) a3
(nominal) class       --> 5 (nominal) class

Time taken to build model: 0 seconds

=== Evaluation on test set ===

Time taken to test model on supplied test set: 0 seconds

=== Summary ===

Correctly Classified Instances      3          75    %
Incorrectly Classified Instances    1          25    %
Kappa statistic                    0.5
Mean absolute error                 0.375
Root mean squared error            0.4486
Relative absolute error             71.7391 %
Root relative squared error        85.5785 %
Total Number of Instances          4

=== Detailed Accuracy By Class ===

              TP Rate  FP Rate  Precision  Recall  F-Measure  MCC      ROC Area  PRC Area  Cla
              0.667    0.000    1.000     0.667    0.800     0.577    0.833    0.917    +
              1.000    0.333    0.500     1.000    0.667    0.577    0.833    0.500    -
Weighted Avg.   0.750    0.083    0.875     0.750    0.767    0.577    0.833    0.813

=== Confusion Matrix ===

a b  <-- classified as
2 1 | a = +
0 1 | b = -
```

Με υπολογισμό στο χέρι λέμε πως πάλι είναι 75% γιατί είναι:

id10 (T = +) → Σωστό

id11 (T = +) → Σωστό

id12 (F = -) → Σωστό

id13(F = +) → Λάθος.

Οπότε καταλήγουμε ότι $\frac{3}{4}$ είναι σωστά και $\frac{1}{4}$ λάθος.