

# Data Mining

## Week 8 Exercise Sheet

1. Consider the following data set for a binary classification problem:

A	B	Class Label
T	F	+
T	T	+
T	T	+
T	F	-
T	T	+
F	F	-
F	F	-
F	F	-
T	T	-
T	F	-

- (a) Calculate the information gain when splitting on  $A$  and  $B$ . Which attribute would the decision tree induction choose?
- (b) Calculate the gain in the Gini index when splitting on  $A$  and  $B$ . Which attribute would the decision tree induction algorithm choose?

2. Consider the following set of training examples:

X	Y	Z	No. of Class C1 Examples	No. of Class C2 Examples
0	0	0	5	40
0	0	1	0	15
0	1	0	10	5
0	1	1	45	0
1	0	0	10	5
1	0	1	25	0
1	1	0	5	20
1	1	1	0	15

- (a) Compute a two-level decision tree using the greedy approach described in lectures. Use the classification error rate as the criterion for splitting. What is the overall error rate of the induced tree?

- (b) Repeat part (a) using  $X$  as the first splitting attribute and then choose the best remaining attribute for splitting at each of the two successor nodes. What is the error rate of the induced tree?
- (c) Compare the results of parts (a) and (b). Comment on the suitability of the greedy heuristic used for splitting attribute selection.