

**Q1 (12 points): Chap 4 Exercise 3 (a) to (f) (2 points each)**

Given a table

Table 1 Dataset for Exercise 3

Instance	a <sub>1</sub>	a <sub>2</sub>	a <sub>3</sub>	Target Class
1	T	T	1.0	+
2	T	T	6.0	+
3	T	F	5.0	-
4	F	F	4.0	+
5	F	T	7.0	-
6	F	T	3.0	-
7	F	F	8.0	-
8	T	F	7.0	+
9	F	T	5.0	-

- a. What is the entropy of this collection of training examples with respect to the class attribute?

We count the number of (+) and (-) then we have:

C1 (+)	4
C2 (-)	5

Recap :

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

$$P(C1) = 4/9$$

$$P(C2) = 5/9$$

Thus:

$$Entropy = - (4/9) * \log_2 (4/9) - (5/9) * \log_2 (5/9) = \mathbf{0.99107606}$$

- b. What are the information gains of a<sub>1</sub>, a<sub>2</sub> relative to these training examples?

Recap:

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

First, we calculate entropies for a<sub>1</sub> and a<sub>2</sub>.

		Target Class	
		Positive (+)	Negative (-)
a <sub>1</sub>	T	3	1
	F	1	4

$$\text{Entropy}(a_1, \text{target class}) = P(T) * \text{Entropy}(3,1) + P(F) * \text{Entropy}(1,4) \\ = \frac{4}{9} \left[ -\frac{3}{4} * \log_2\left(\frac{3}{4}\right) - \frac{1}{4} * \log_2\left(\frac{1}{4}\right) \right] + \frac{5}{9} \left[ -\frac{1}{5} * \log_2\left(\frac{1}{5}\right) - \frac{4}{5} * \log_2\left(\frac{4}{5}\right) \right] = 0.7616$$

$$\text{GAIN } a_1 = 0.9912 - 0.7616 = 0.2296$$

		Target Class	
		Positive (+)	Negative (-)
a <sub>2</sub>	T	2	3
	F	2	2

$$\text{Entropy}(a_2, \text{target class}) = P(T) * \text{Entropy}(2,3) + P(F) * \text{Entropy}(2,2) = \\ \frac{5}{9} \left[ -\frac{2}{5} * \log_2\left(\frac{2}{5}\right) - \frac{3}{5} * \log_2\left(\frac{3}{5}\right) \right] + \frac{4}{9} \left[ -\frac{2}{4} * \log_2\left(\frac{2}{4}\right) - \frac{2}{4} * \log_2\left(\frac{2}{4}\right) \right] = 0.9839$$

$$\text{GAIN } a_2 = 0.9912 - 0.9839 = 0.0072$$

- c. For a<sub>3</sub>, which is a continuous attribute, compute the information gain for every possible split.

	1		3		4		5		6		7		8			
Split	0.5		2		3.5		4.5		5.5		6.5		7.5		8.5	
	≤	>	≤	>	≤	>	≤	>	≤	>	≤	>	≤	>	≤	>
P(+)	0	4	1	3	1	3	2	2	2	2	3	1	4	0	4	0
P(-)	0	5	0	5	1	4	1	4	3	2	3	2	4	1	5	0

**RECAP: 0\*log<sub>2</sub>0=0**

$$\text{Entropy}(1, \text{target class}) = P(\leq 0.5) * \text{Entropy}(0,0) + P(>0.5) * \text{Entropy}(4,5)$$

$$0 + \frac{9}{9} \left[ -\frac{4}{9} * \log_2\left(\frac{4}{9}\right) - \frac{5}{9} * \log_2\left(\frac{5}{9}\right) \right] = 0.9912$$

$$\text{GAIN}(1) = 0.9912 - 0.9912 = 0$$

$$\text{Entropy}(3, \text{target class}) = P(\leq 2) * \text{Entropy}(1,0) + P(>2) * \text{Entropy}(3,5)$$

$$\frac{1}{9} \left[ -\frac{1}{1} * \log_2\left(\frac{1}{1}\right) - \frac{0}{1} * \log_2\left(\frac{0}{1}\right) \right] + \frac{8}{9} \left[ -\frac{3}{8} * \log_2\left(\frac{3}{8}\right) - \frac{5}{8} * \log_2\left(\frac{5}{8}\right) \right] = 0.8484$$

$$\text{GAIN (3)} = 0.9912 - 0.8484 = 0.1428$$

$$\text{Entropy (4, target class)} = P(\leq 3.5) * \text{Entropy (1,1)} + P(>3.5) * \text{Entropy(3,4)}$$

$$\frac{2}{9} \left[ -\frac{1}{2} * \log_2\left(\frac{1}{2}\right) - \frac{1}{2} * \log_2\left(\frac{1}{2}\right) \right] + \frac{7}{9} \left[ -\frac{3}{7} * \log_2\left(\frac{3}{7}\right) - \frac{4}{7} * \log_2\left(\frac{4}{7}\right) \right] = 0.9885$$

$$\text{GAIN (4)} = 0.9912 - 0.9885 = 0.0027$$

$$\text{Entropy (5, target class)} = P(\leq 4.5) * \text{Entropy (2,1)} + P(>4.5) * \text{Entropy(2,4)}$$

$$\frac{3}{9} \left[ -\frac{2}{3} * \log_2\left(\frac{2}{3}\right) - \frac{1}{3} * \log_2\left(\frac{1}{3}\right) \right] + \frac{6}{9} \left[ -\frac{2}{6} * \log_2\left(\frac{2}{6}\right) - \frac{4}{6} * \log_2\left(\frac{4}{6}\right) \right] = 0.9183$$

$$\text{GAIN (5)} = 0.9912 - 0.9183 = 0.0729$$

$$\text{Entropy (6, target class)} = P(\leq 5.5) * \text{Entropy (2,3)} + P(>5.5) * \text{Entropy(2,2)}$$

$$\frac{5}{9} \left[ -\frac{2}{5} * \log_2\left(\frac{2}{5}\right) - \frac{3}{5} * \log_2\left(\frac{3}{5}\right) \right] + \frac{4}{9} \left[ -\frac{2}{4} * \log_2\left(\frac{2}{4}\right) - \frac{2}{4} * \log_2\left(\frac{2}{4}\right) \right] = 0.9839$$

$$\text{GAIN (6)} = 0.9912 - 0.9839 = 0.0073$$

$$\text{Entropy (7, target class)} = P(\leq 6.5) * \text{Entropy (3,3)} + P(>6.5) * \text{Entropy(1,2)}$$

$$\frac{6}{9} \left[ -\frac{3}{6} * \log_2\left(\frac{3}{6}\right) - \frac{3}{6} * \log_2\left(\frac{3}{6}\right) \right] + \frac{3}{9} \left[ -\frac{1}{3} * \log_2\left(\frac{1}{3}\right) - \frac{2}{3} * \log_2\left(\frac{2}{3}\right) \right] = 0.9728$$

$$\text{GAIN (7)} = 0.9912 - 0.9728 = 0.0184$$

$$\text{Entropy (8, target class)} = P(\leq 7.5) * \text{Entropy (4,4)} + P(>7.5) * \text{Entropy(0,1)}$$

$$\frac{8}{9} \left[ -\frac{4}{8} * \log_2\left(\frac{4}{8}\right) - \frac{4}{8} * \log_2\left(\frac{4}{8}\right) \right] + \frac{1}{9} \left[ -\frac{0}{1} * \log_2\left(\frac{0}{1}\right) - \frac{1}{1} * \log_2\left(\frac{1}{1}\right) \right] = 0.8889$$

$$\text{GAIN (8)} = 0.9912 - 0.8889 = 0.1023$$

- d. What is the best split (among  $a_1, a_2, a_3$ ) according to the information gain?

Best of  $a_3$  at  $\text{GAIN (3)} = 0.1428$ ,

hence

$$\text{GAIN (a}_3\text{)} = 0.1428,$$

$$\text{GAIN (a}_2\text{)} = 0.0072,$$

$$\text{GAIN (a}_1\text{)} = 0.2296$$

We can see that,  $\text{MAX}(\text{GAIN}(a_3), \text{GAIN}(a_2), \text{GAIN}(a_1)) = \text{GAIN}(a_1) = 0.2296$ . Therefore,  $a_1$  provides the best split.

e. What is the best split (between  $a_1, a_2$ ) according to the classification error rate?

**Recap:**

Classification Error Rate:  $= (\text{FP} + \text{FN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$

Classification error rate  $a_1$ :  $(1+1) / (3+1+1+4) = 2/9 = 0.2222$

Classification error rate  $a_2$ :  $(2+2) / (2+3+2+2) = 4/9 = 0.4444$

Therefore,  **$a_1$  provides best split** because of lower classification error rate

f. What is the best split (between  $a_1, a_2$ ) according to Gini Index?

$$\text{GINI INDEX}(a_1) : \frac{4}{9} \left( 1 - \left( \frac{3}{4} \right)^2 - \left( \frac{1}{4} \right)^2 \right) + \frac{5}{9} \left( 1 - \left( \frac{1}{5} \right)^2 - \left( \frac{4}{5} \right)^2 \right) = 0.3444$$

$$\text{GINI INDEX}(a_2) : \frac{5}{9} \left( 1 - \left( \frac{2}{5} \right)^2 - \left( \frac{3}{5} \right)^2 \right) + \frac{4}{9} \left( 1 - \left( \frac{2}{4} \right)^2 - \left( \frac{2}{4} \right)^2 \right) = 0.4889$$

Therefore,  **$a_1$  provides best split** because of lower GINI Index

**Question 2: Chapter 4, Exercise 5**

Consider the following data set for a binary class 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 algorithm choose?

<b>Class Label (+)</b>	<b>4</b>
<b>Class Label (-)</b>	<b>6</b>

$$P(C1) = 4/10$$

$$P(C2) = 6/10$$

$$\text{Thus, Entropy} = -4/10 \log_2 4/10 - 6/10 \log_2 6/10 = 0.9710$$

		Class Label	
		Positive (+)	Negative (-)
A	T	4	3
	F	0	3

$$\begin{aligned} \text{Entropy (A, Class Label)} &= P(T) * \text{Entropy (4,3)} + P(F) * \text{Entropy (0,3)} \\ &= \frac{7}{10} \left[ -\frac{4}{7} * \log_2 \left( \frac{4}{7} \right) - \frac{3}{7} * \log_2 \left( \frac{3}{7} \right) \right] + \frac{3}{10} \left[ -\frac{0}{3} * \log_2 \left( \frac{0}{3} \right) - \frac{3}{3} * \log_2 \left( \frac{3}{3} \right) \right] = 0.6896 \\ \text{GAIN (A)} &= 0.9710 - 0.6896 = \mathbf{0.2813} \end{aligned}$$

		Class Label	
		Positive (+)	Negative (-)
B	T	3	1
	F	1	5

$$\begin{aligned} \text{Entropy (B, Class Label)} &= P(T) * \text{Entropy (3,1)} + P(F) * \text{Entropy (1,5)} \\ &= \frac{4}{10} \left[ -\frac{3}{4} * \log_2 \left( \frac{3}{4} \right) - \frac{1}{4} * \log_2 \left( \frac{1}{4} \right) \right] + \frac{6}{10} \left[ -\frac{1}{6} * \log_2 \left( \frac{1}{6} \right) - \frac{5}{6} * \log_2 \left( \frac{5}{6} \right) \right] = 0.7145 \\ \text{GAIN (B)} &= 0.9710 - 0.7145 = \mathbf{0.2565} \end{aligned}$$

Because  $\text{GAIN (A)} > \text{GAIN (B)}$ , **attribute A** will be chosen to split the node.

- b. Calculate the gain in the Gini Index when splitting on A and B. Which attribute would the decision tree induction algorithm choose?

$$\text{Overall GINI Index: } 1 - (4/10)^2 - (6/10)^2 = 0.48$$

$$\text{GINI Index (A)} : \frac{7}{10} \left( 1 - \left( \frac{4}{7} \right)^2 - \left( \frac{3}{7} \right)^2 \right) + \frac{3}{10} \left( 1 - \left( \frac{0}{3} \right)^2 - \left( \frac{3}{3} \right)^2 \right) = \mathbf{0.3428}$$

$$\text{GINI Index (B)} : \frac{4}{10} \left( 1 - \left( \frac{3}{4} \right)^2 - \left( \frac{1}{4} \right)^2 \right) + \frac{6}{10} \left( 1 - \left( \frac{1}{6} \right)^2 - \left( \frac{5}{6} \right)^2 \right) = \mathbf{0.3167}$$

Because  $\text{GINI (B)} < \text{GINI (A)}$ , **attribute B** will be chosen to split the node

- c. Figure 4.13 shows that entropy and the Gini Index are both monotonously increasing on the range  $[0, 0.5]$  and they are both monotonously decreasing on the range  $[0.5, 1]$ . Is it possible that information gain and the gain in the Gini Index favor different attribute? Explain

*Yes, as indicated in previous question a) and b) due to the use of different functions.*

### **Question 3. Chapter 4 Exercise 7**

The following table summarizes a data set with three attributes A, B, C and two class labels +, -. Build a two-level decision tree

A	B	C	Number of instances	
			+	-
T	T	T	5	0
F	T	T	0	20
T	F	T	20	0
F	F	T	0	5
T	T	F	0	0
F	T	F	25	0
T	F	F	0	0
F	F	F	0	25

- a. According to the classification error rate, which attribute would be chosen as the first splitting attribute? For each attribute, show the contingency table and the gains in classification error rate.

	Number of instances
+	50
-	50

Original classification error rate:  $1 - \text{MAX}(50/100, 50/100) = 0.5$

Entropy =  $-50/100 * \log_2 50/100 - 50/100 * \log_2 50/100 = 1$

$$\begin{aligned} \text{Entropy (A, Class Label)} &= P(T) * \text{Entropy}(25,0) + P(F) * \text{Entropy}(25,50) \\ &= \frac{25}{100} \left[ -\frac{25}{25} * \log_2 \left( \frac{25}{25} \right) - \frac{0}{25} * \log_2 \left( \frac{0}{25} \right) \right] + \frac{75}{100} \left[ -\frac{25}{75} * \log_2 \left( \frac{25}{75} \right) - \frac{50}{75} * \log_2 \left( \frac{50}{75} \right) \right] \\ &= 0.6887 \end{aligned}$$

$$\text{GAIN}(A) = 1 - 0.6887 = \mathbf{0.3113}$$

$$\begin{aligned} \text{Entropy (B, Class Label)} &= P(T) * \text{Entropy}(30,20) + P(F) * \text{Entropy}(20,30) \\ &= \frac{50}{100} \left[ -\frac{30}{50} * \log_2 \left( \frac{30}{50} \right) - \frac{20}{50} * \log_2 \left( \frac{20}{50} \right) \right] + \frac{50}{100} \left[ -\frac{20}{50} * \log_2 \left( \frac{20}{50} \right) - \frac{30}{50} * \log_2 \left( \frac{30}{50} \right) \right] \end{aligned}$$

$$= 0.9710$$

$$GAIN(B) = 1 - 0.9710 = 0.029$$

$$\begin{aligned} \text{Entropy}(C, \text{Class Label}) &= P(T) * \text{Entropy}(25, 25) + P(F) * \text{Entropy}(25, 25) \\ &= \frac{50}{100} \left[ -\frac{25}{50} * \log_2\left(\frac{25}{50}\right) - \frac{25}{50} * \log_2\left(\frac{25}{50}\right) \right] + \frac{50}{100} \left[ -\frac{25}{50} * \log_2\left(\frac{25}{50}\right) - \frac{25}{50} * \log_2\left(\frac{25}{50}\right) \right] \\ &= 1 \\ GAIN(C) &= 1 - 1 = 0 \end{aligned}$$

**Classification Error Rate: = (FP + FN) / (TP + TN + FP + FN)**

		+	-
<b>A</b>	<b>T</b>	25	0
	<b>F</b>	25	50

		+	-
<b>B</b>	<b>T</b>	30	20
	<b>F</b>	20	30

		+	-
<b>C</b>	<b>T</b>	25	25
	<b>F</b>	25	25

Classification Error Rate A:  $(25+0)/(25+0+25+50) = 25/100$

Classification Error Rate B:  $(20+20)/(30+20+20+30) = 40/100$

Classification Error Rate C:  $(25+25)/(25+25+25+25) = 50/100$

**Attribute A provides** the lowest classification error rate (or highest information gain); therefore, it will be chosen to split the node.

- b. Repeat for the two children of the root node

Because  $A_T = 25$  (+) and 0 (-). It is pure for T's values; no further action is required!

For  $A_{T=F}$  we construct new table as follow by removing the shaded rows

A	B	C	Number of instances	
			+	-
T	T	T	5	0
F	T	T	0	20
T	F	T	20	0
F	F	T	0	5
T	T	F	0	0
F	T	F	25	0
T	F	F	0	0
F	F	F	0	25

We have:

A	B	C	Number of instances	
			+	-
F	T	T	0	20
F	F	T	0	5
F	T	F	25	0
F	F	F	0	25

	Number of instances
+	25
-	50

		+	-
B	T	25	20
	F	0	30

		+	-
C	T	0	25
	F	25	25

Classification Error Rate B:  $(20+0)/(25+20+0+30) = 20/75$

Classification Error Rate C:  $(25+25)/(0+25+25+25) = 50/75$

Entropy =  $-25/75 * \log_2 25/75 - 50/75 * \log_2 50/75 = 0.9183$

Entropy (B, Class Label) =  $P(T) * \text{Entropy}(25, 20) + P(F) * \text{Entropy}(0, 30)$   
 $= \frac{45}{75} \left[ -\frac{25}{45} * \log_2 \left( \frac{25}{45} \right) - \frac{20}{45} * \log_2 \left( \frac{20}{45} \right) \right] + \frac{30}{75} \left[ -\frac{0}{30} * \log_2 \left( \frac{0}{30} \right) - \frac{30}{30} * \log_2 \left( \frac{30}{30} \right) \right]$   
 $= 0.5946$

$GAIN(B) = 0.9183 - 0.5946 = 0.3237$



$$\begin{aligned}
 \text{Entropy (C, Class Label)} &= P(T) * \text{Entropy (0, 25)} + P(F) * \text{Entropy (25,25)} \\
 &= \frac{25}{75} \left[ -\frac{0}{25} * \log_2\left(\frac{0}{25}\right) - \frac{25}{25} * \log_2\left(\frac{25}{45}\right) \right] + \frac{50}{75} \left[ -\frac{25}{50} * \log_2\left(\frac{25}{50}\right) - \frac{25}{50} * \log_2\left(\frac{25}{50}\right) \right] \\
 &= 0.6667 \\
 \text{GAIN (B)} &= 0.9183 - 0.6667 = 0.2156
 \end{aligned}$$

From both classification error rate and information gain, we can see that, **attribute B** will be chosen to split node.

- c. How many instances are misclassified by the resulting decision tree?

		+	-
<b>B</b>	<b>T</b>	25	<b>20</b>
	<b>F</b>	0	30

Because  $B_F = 0$  (+) and 30 (-). It is pure for F's values; Thus, the resulting decision tree has 20 misclassified instances.

- d. Repeat parts a), b), c) using C as the first root/splitting attribute.

Recap:

		+	-
<b>C</b>	<b>T</b>	25	25
	<b>F</b>	25	25

$$\begin{aligned}
 \text{Entropy (C, Class Label)} &= P(T) * \text{Entropy (25,25)} + P(F) * \text{Entropy (25,25)} \\
 &= \frac{50}{100} \left[ -\frac{25}{50} * \log_2\left(\frac{25}{50}\right) - \frac{25}{50} * \log_2\left(\frac{25}{50}\right) \right] + \frac{50}{100} \left[ -\frac{25}{50} * \log_2\left(\frac{25}{50}\right) - \frac{25}{50} * \log_2\left(\frac{25}{50}\right) \right] \\
 &= 1
 \end{aligned}$$

$$\text{GAIN (C)} = 1 - 1 = 0$$

$$\text{Classification Error Rate C: } (25+25)/(25+25+25+25) = 50/100$$

After splitting node C, we have two sets (distinguished by shaded region)

<b>A</b>	<b>B</b>	<b>C</b>	<b>Number of instances</b>	
			+	-
T	T	T	5	0
F	T	T	0	20
T	F	T	20	0
F	F	T	0	5
T	T	F	0	0
F	T	F	25	0
T	F	F	0	0
F	F	F	0	25

For the 1<sup>st</sup> set

		+	-
A	T	25	0
	F	0	25

		+	-
B	T	5	20
	F	20	5

Classification error rate A:  $0/50 = 0$

Classification error rate B:  $40/50$

$$\begin{aligned} \text{Entropy (A, Class Label)} &= P(T) * \text{Entropy (25, 0)} + P(F) * \text{Entropy (0, 25)} \\ &= \frac{25}{50} \left[ -\frac{25}{25} * \log_2\left(\frac{25}{25}\right) - \frac{0}{25} * \log_2\left(\frac{0}{25}\right) \right] + \frac{25}{50} \left[ -\frac{0}{25} * \log_2\left(\frac{0}{25}\right) - \frac{25}{25} * \log_2\left(\frac{25}{25}\right) \right] \\ &= 0 \end{aligned}$$

$$GAIN(B) = 1 - 0 = 1.$$

$$\begin{aligned} \text{Entropy (B, Class Label)} &= P(T) * \text{Entropy (5, 20)} + P(F) * \text{Entropy (20, 5)} \\ &= \frac{25}{50} \left[ -\frac{5}{25} * \log_2\left(\frac{5}{25}\right) - \frac{20}{25} * \log_2\left(\frac{20}{25}\right) \right] + \frac{25}{50} \left[ -\frac{20}{25} * \log_2\left(\frac{20}{25}\right) - \frac{5}{25} * \log_2\left(\frac{5}{25}\right) \right] \\ &= 0.7220 \end{aligned}$$

$$GAIN(B) = 1 - 0.7220 = 0.278$$

Attribute A will be chosen to split the node with 0 misclassified instances

**For the 2<sup>nd</sup> set**

		+	-
A	T	0	0
	F	25	25

		+	-
B	T	25	0
	F	0	25

Classification error rate A:  $25/50 = 0.5$

Classification error rate B:  $0/50 = 0$

$$\begin{aligned} \text{Entropy (A, Class Label)} &= P(T) * \text{Entropy (0, 0)} + P(F) * \text{Entropy (25, 25)} \\ &= \left[ -\frac{25}{50} * \log_2\left(\frac{25}{50}\right) - \frac{25}{50} * \log_2\left(\frac{25}{50}\right) \right] \\ &= 1 \end{aligned}$$

$$GAIN(A) = 1 - 1 = 0.$$

$$\begin{aligned} \text{Entropy (B, Class Label)} &= P(T) * \text{Entropy}(25, 0) + P(F) * \text{Entropy}(0, 25) \\ &= \frac{25}{50} \left[ -\frac{25}{25} * \log_2\left(\frac{25}{25}\right) - \frac{0}{25} * \log_2\left(\frac{0}{25}\right) \right] + \frac{25}{50} \left[ -\frac{0}{25} * \log_2\left(\frac{0}{25}\right) - \frac{25}{25} * \log_2\left(\frac{25}{25}\right) \right] \\ &= 0 \\ GAIN(B) &= 1 - 0 = 1 \end{aligned}$$

Attribute B will be chosen to split the node with 0 misclassified instance

- e. Use the results in parts (c) and (d) to conclude about the greedy nature of decision tree induction algorithm.

It can be seen that, starting from attribute A can lead to 20 misclassified instances but starting from attribute C, overall classification error rate is 0. This concludes that greedy heuristic does not always lead to the best decision tree.

#### **Question 4: Chapter 5 Exercise 7**

Consider the data set shown in Table 5.10

Record	A	B	C	Class
1	0	0	0	+
2	0	0	1	-
3	0	1	1	-
4	0	1	1	-
5	0	0	1	+
6	1	0	1	+
7	1	0	1	-
8	1	0	1	-
9	1	1	1	+
10	1	0	1	+

- a. Estimate the conditional probabilities for  $P(A|+)$ ,  $P(B|+)$ ,  $P(C|+)$ ,  $P(A|-)$ ,  $P(B|-)$ ,  $P(C|-)$ ?

$$P(A_0|+) = 2/5 = 0.4 ; \quad P(A_0|-) = 3/5 = 0.6$$

$$P(A_1|+) = 3/5 = 0.6 ; \quad P(A_1|-) = 2/5 = 0.4$$

$$P(B_0|+) = 4/5 ; \quad P(B_0|-) = 3/5$$

$$P(B_1|+) = 1/5 ; \quad P(B_1|-) = 2/5$$

$$P(C_0|+) = 1/5 ; \quad P(C_0|-) = 0/5$$

$$P(C_1|+) = 4/5 ; \quad P(C_1|-) = 5/5$$

- b. Use the estimate of conditional probabilities given in the previous question to predict the class label for a test sample ( $A=0, B=1, C=0$ ) using the naïve Bayes approach.

To predict the class label of a test record  $X = (A=0, B=1, C=0)$  we need to compute the posterior probabilities  $P(+|X)$  and  $P(-|X)$

$$P(+) = P(-) = 0.5$$

$$P(X|+) = P(A=0|+) * P(B=1|+) * P(C=0|+) = 2/5 * 1/5 * 1/5 = 2/75$$

$$P(X|-) = P(A=0|-) * P(B=1|-) * P(C=0|-) = 3/5 * 2/5 * 0/5 = 0$$

$$P(+|X) = P(X|+) * P(+) / P(X) = a * 0.5 * 2/75 \text{ where } a = 1/P(X) \text{ is a constant term}$$

$$P(-|X) = P(X|-) * P(-) / P(X) = 0$$

Since  $P(+|X) > P(-|X)$ , the record is classified as ‘+’

- c. Estimate the conditional probabilities using the m-estimate approach, with  $p=1/2$ , and  $m=4$

Recap:

$$P(x_i|y_i) = (n_c + mp)/(n+m)$$

$$P(A_0|+) = (2 + 4*1/2)/5+4 = 4/9$$

$$P(A_0|-) = (3 + 4*1/2)/5+4 = 5/9$$

$$P(A_1|+) = (3 + 4*1/2)/5+4 = 5/9$$

$$P(A_1|-) = (2 + 4*1/2)/5+4 = 4/9$$

$$P(B_0|+) = (4 + 4*1/2)/5+4 = 6/9$$

$$P(B_0|-) = (3 + 4*1/2)/5+4 = 5/9$$

$$P(B_1|+) = (1 + 4*1/2)/5+4 = 3/9$$

$$P(B_1|-) = (2 + 4*1/2)/5+4 = 4/9$$

$$P(C_0|+) = (1 + 4*1/2)/5+4 = 3/9$$

$$P(C_0|-) = (0 + 4*1/2)/5+4 = 2/9$$

$$P(C_1|+) = (4 + 4*1/2)/5+4 = 6/9$$

$$P(C_1|-) = (5 + 4*1/2)/5+4 = 7/9$$

- d. Repeat part (b) using the conditional probabilities given in part (c)

$$P(X|+) = P(A=0|+) * P(B=1|+) * P(C=0|+) = 4/9 * 3/9 * 3/9 = 4/81 = 0.04938$$

$$P(X|-) = P(A=0|-) * P(B=1|-) * P(C=0|-) = 5/9 * 4/9 * 2/9 = 40/729 = 0.05487$$

$$P(+|X) = P(X|+) * P(+) / P(X) = a * 0.5 * 0.04938$$

where  $a = 1/P(X)$  is a constant term

$$P(-|X) = P(X|-) * P(-) / P(X) = a * 0.5 * 0.05487$$

Since  $P(+|X) < P(-|X)$ , the record is classified as ‘-’

- e. Compare the two methods for estimating probabilities. Which method is better and why?

The reason to have m-estimate of conditional probability is to overcome the vanish of the overall posterior probability if the class-conditional probability for one of the attribute is zero. This problem is shown in part (b). So **it is obvious that the m-estimate approach is better.**

### Question 5: Chapter 5 Exercise 12

Given the Bayesian network shown in Figure 5.48, compute the following probabilities.

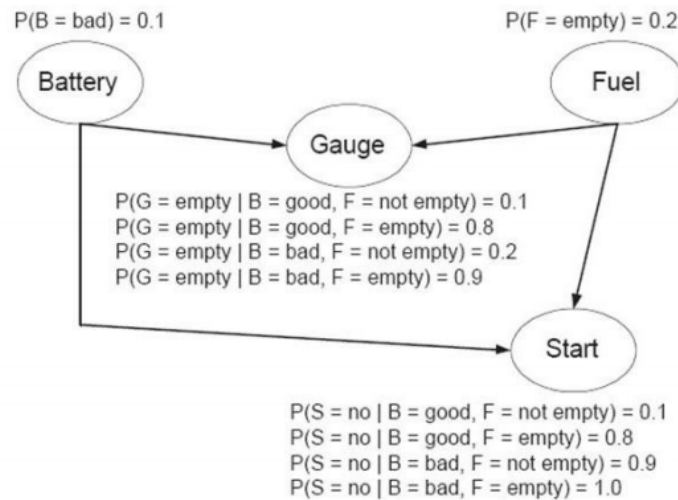


Figure 5.48. Bayesian belief network for Exercise 12.

- $P(B=\text{good}, F=\text{empty}, G=\text{empty}, S=\text{yes})$   
 $= P(G,S|B,F) \cdot P(B,F) = P(G|BF) \cdot P(S|BF) \cdot P(B) \cdot P(F)$   
 $= 0.8 * (1-0.8) * (1-0.1) * 0.2 = 0.8 * 0.2 * 0.9 * 0.2 = \mathbf{0.0288}$
- $P(B=\text{bad}, F=\text{empty}, G=\text{not empty}, S = \text{no})$   
 $= P(G,S|B,F) \cdot P(B,F) = P(G|BF) \cdot P(S|BF) \cdot P(B) \cdot P(F)$   
 $= (1-0.9) * 1.0 * 0.1 * 0.2 = \mathbf{0.002}$
- Given that the battery is bad, compute the probability that the cart will start.  
 Let  $\alpha \in \{\text{empty}, \text{not empty}\}$ . Given the conditions above we need to calculate  $P(S=\text{start}|B=\text{bad}, F=\alpha)$ .

First we have

$$P(S=\text{start}|B=\text{bad}, F=\text{empty}) = 1 - P(S=\text{no}|B=\text{bad}, F=\text{empty}) = 1 - 1 = 0$$

Then

$$\begin{aligned}
 &P(S=\text{start}|B=\text{bad}, F=\alpha) \\
 &= \sum_{\alpha} P(S = \text{start} | B = \text{bad}, F = \alpha) P(B = \text{bad}) P(F = \alpha) \\
 &= P(S = \text{start} | B = \text{bad}, F = \text{empty}) P(B = \text{bad}) P(F = \text{empty}) + \\
 &\quad P(S = \text{start} | B = \text{bad}, F = \text{not empty}) P(B = \text{bad}) P(F = \text{not empty})
 \end{aligned}$$

$$\begin{aligned}
 &= 0 + (1 - P(S = \text{no} | B = \text{bad}, F = \text{not empty})) * P(B = \text{bad}) * P(F = \text{not empty}) \\
 &= (1 - 0.9) * 0.1 * (1 - 0.2) = 0.1 * 0.1 * 0.8 = \mathbf{0.008}
 \end{aligned}$$

### **Question 6: Chapter 5 Exercise 17**

You are asked to evaluate the performance of two classification models, M1 and M2. The test set you have chosen contains 26 binary attributes, labeled as A through Z.

Table 5.14 shows the posterior probabilities obtained by applying the models to the test set. (Only the posterior probabilities for the positive class are shown). As this is a two-class problem,  $P(-) = 1 - P(+)$  and  $P(-|A, \dots, Z) = 1 - P(+|A, \dots, Z)$ . Assume that we are mostly interested in detecting instances from the positive class.

- Plotting the ROC curve for both M1 and M2 (You should plot them on the same graph). Which model do you think is better? Explain your reasons?

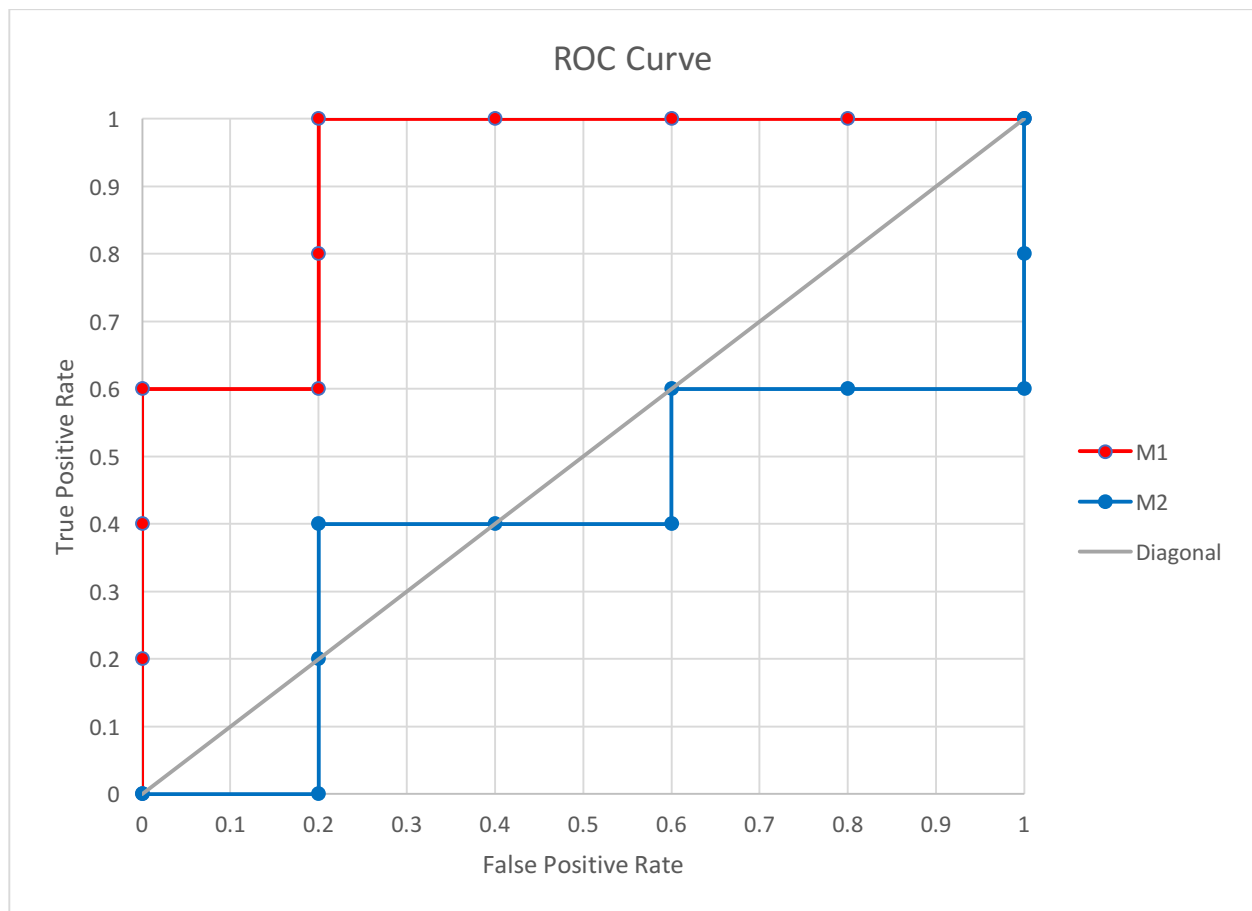
**Table 5.14.** Posterior probabilities for Exercise 17.

Instance	True Class	$P(+ A, \dots, Z, M_1)$	$P(+ A, \dots, Z, M_2)$
1	+	0.73	0.61
2	+	0.69	0.03
3	-	0.44	0.68
4	-	0.55	0.31
5	+	0.67	0.45
6	+	0.47	0.09
7	-	0.08	0.38
8	-	0.15	0.05
9	+	0.45	0.01
10	-	0.35	0.04

We construct two new tables as follow

Class	-	-	-	-	+	+	-	+	+	+	
M1	0.08	0.15	0.35	0.44	0.45	0.47	0.55	0.6	0.69	0.73	1
TP	5	5	5	5	5	4	3	3	2	1	0
FP	5	4	3	2	1	1	1	0	0	0	0
TN	0	1	2	3	4	4	4	5	5	5	5
FN	0	0	0	0	0	1	2	2	3	4	5
TPR	1	1	1	1	1	0.8	0.6	0.6	0.4	0.2	0
FPR	1	0.8	0.6	0.4	0.2	0.2	0.2	0	0	0	0

Class	+	+	-	-	+	-	-	+	+	-	
M2	0.01	0.03	0.04	0.05	0.09	0.31	0.38	0.45	0.61	0.68	1
TP	5	4	3	3	3	2	2	2	1	0	0
FP	5	5	5	4	3	3	2	1	1	1	0
TN	0	0	0	1	2	2	3	4	4	4	5
FN	0	1	2	2	2	3	3	3	4	5	5
TPR	1	0.8	0.6	0.6	0.6	0.4	0.4	0.4	0.2	0	0
FPR	1	1	1	0.8	0.6	0.6	0.4	0.2	0.2	0.2	0



From the graph above, it is clearly **that M1 performs better than M2** because its area under ROC curve is larger than that of M2's area.

- b. For model M1, suppose you choose the cutoff threshold to be  $t=0.5$ . In other words, any test instances whose posterior probability is greater than  $t$  will be classified as a positive example. Compute the precision, recall, and F-measure for the model at this threshold value

At threshold  $t = 0.5$  we have the following table.

<b>M1</b>	<b>Actual Class</b>	<b>Predicted Class</b>
0.08	-	-
0.15	-	-
0.35	-	-
0.44	-	-
0.45	+	-
0.47	+	-
0.55	-	+
0.67	+	+
0.69	+	+
0.73	+	+

We have confusion table as follow

<b>Actual Class (M1)</b>	<b>Predicted Class</b>		
		Class = +	Class = -
	Class =+	a = 3	b = 2
	Class =-	c = 1	d = 4

Precision (p) =  $a/(a+c) = 3/4 = 0.75$

Recall (r) =  $a/(a+b) = 3/(3+2) = 3/5 = 0.6$

F-measure (F) =  $2rp/(r+p) = 2a/(2a+b+c) = 2*3/(2*3+ 2+1) = 6/9 = 0.667$

- c. Repeat the analysis for part (c) using the same cutoff threshold on model M2. Compare the F-measure results for both models. Which model is better? Are the results consistent with what you expect from the ROC curve?

<b>M2</b>	<b>Actual Class</b>	<b>Predicted Class</b>
0.01	+	-
0.03	+	-
0.04	-	-
0.05	-	-
0.09	+	-
0.31	-	-
0.38	-	-
0.45	+	-
0.61	+	+
0.68	-	+

<b>Actual Class (M2)</b>	<b>Predicted Class</b>		
		Class = +	Class = -
	Class =+	a = 1	b = 4
	Class =-	c = 1	d = 4



$$\text{Precision (p)} = a/(a+c) = 1/2 = 0.5$$

$$\text{Recall (r)} = a/(a+b) = 1/(1+4) = 1/5 = 0.2$$

$$\text{F-measure (F)} = 2rp/(r+p) = 2a/(2a+b+c) = 2*1/(2*1+4+1) = 2/7 = 0.2857$$

Because F-measure (M1=0.667) > F-measure (M2=0.2857) so M1 is better than M2. Without looking at the ROC curve, we expect that M1 is closer to (0,1) than M2 at t=0.5.

We have M1 = (0.2, 0.6), M2 = (0.2, 0.2).

From the ROC curve, this result is consistent. Indeed

$$d(M1) = \sqrt{((1-0.6)^2 + 0.2^2)} = \sqrt{0.2} = 0.447$$

$$d(M2) = \sqrt{((1-0.2)^2 + 0.2^2)} = \sqrt{0.68} = 0.4624$$

Therefore  $d(M1) < d(M2)$  meaning that M1 is closer

- d. Repeat part (c) for model M1 using the threshold t=0.1. Which threshold do you prefer, t=0.5 or t=0.1? Are the results consistent with what you expect from the ROC curve?

M1	Actual Class	Predicted Class	
0.08	-	-	
0.15	-	+	
0.35	-	+	
0.44	-	+	
0.45	+	+	
0.47	+	+	
0.55	-	+	
0.67	+	+	
0.69	+	+	
0.73	+	+	
	Predicted Class		
Actual Class (M1)		Class = +	Class = -
	Class =+	a = 5	b = 0
	Class =-	c = 4	d = 1

$$\text{Precision (p)} = a/(a+c) = 5/9 = 0.556$$

$$\text{Recall (r)} = a/(a+b) = 5/(5+0) = 5/5 = 1$$

$$\text{F-measure (F)} = 2rp/(r+p) = 2a/(2a+b+c) = 2*5/(2*5+0+4) = 10/14 = 0.7143$$

In terms of F measure, t=0.1 gives the highest F value so it is preferable. We would expect that M1 at t=0.1 is closer to the point (0,1) than M1 at t=0.5 or  $d(M1=0.1) < d(M1=0.5)$

We have

$$M1(0.1) = (0.8, 1), M1(0.5) = (0.2, 0.6)$$

$$d(M1=0.1) = 0.8 \text{ and } d(M1=0.5) = 0.447.$$

Here  $d(M1) > d(M2)$  that is contradict to what we expect. Therefore, this result is NOT consistent with the ROC curve.

Since F-measure and ROC curve are NOT consistent. Area under the curve is a good indicator to select threshold  $t$

$$\text{Area } M1 (t=0.1) = 1 * (1 - 0.8) = 0.2$$

$$\text{Area } M1 (t=0.5) = 0.6 * (1 - 0.2) = 0.48$$

Since  $\text{Area } M1 (t=0.5) > M1 (t=0.1)$  so  **$t=0.5$  is preferable.**

**Question 7. Given a Twitter dataset, how do you build a classifier which can identify Harvey rescue tweets automatically (you can define your target if your data does not contain Harvey rescue data)**

To classify Harvey rescue tweets automatically, we use “bag-of-words” approach. That is, we define some words that have similar meaning to rescue such as “save, recover, release, relief” and compute the distance between each document with “bag-of-words” vector. If the distance is non-zero value, we classify this document as “Rescue”, otherwise “None-Rescue”.

- a. What is your distance function? (3 points)

Since we are working with document, each document will be represented as a vector where each attribute denotes the frequency of a term that occurs in the document. The distance function we are going to use is “**Euclidean distance**”

- b. Can you run KNN on your dataset? How much training set do you have? (7 points)

Yes, since our documents are labelled with “Rescue” and “None-Rescue” we can run KNN algorithm.

Our original data contains approximately 5,300 records, after preprocessing step (remove duplicate tweet or retweets), our final data contains 1700 records. As a rule of thumb, 70% of data will be used as training data 1190 documents.

- c. Can you evaluate your KNN classifier using precision, recall and F-measure? (4 points).

Yes, we can evaluate KNN classifier based on the prediction output vs testing data

- d. How to determine the best K? Using different K, compare their recall, precision and F-measure? (4 points)

As a rule of thumb, the best K is the **square root** of the total records in the dataset.

Note that we will select K as an **odd number**.

- e. Visualize classification results? (2 points)