

P.L.	A	A	A	A	A	A	B	B	B	B	B	B	B	B	B
X	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
					X				X	X				X	

Part II: Naïve Bayesian Classifier (20 points)

Consider the following training data from an employee database. The target feature (attribute) is salary.

department	status	age	salary
sales	senior	31-40	Medium
sales	junior	21-30	Low
sales	junior	31-40	Low
systems	junior	21-30	Medium
systems	senior	31-40	High
systems	junior	21-30	Medium
systems	senior	41-50	High
marketing	senior	31-40	Medium
marketing	junior	31-40	Medium
secretary	senior	41-50	Medium
secretary	junior	21-30	Low

If given a test instance with values: systems, senior, and 21-30 for the descriptive attributes department, status, and age, respectively, what would be a Naïve Bayesian classification for the salary of the test instance? (Show your calculation process)

Answer:

$$P(x_1 = \text{systems}, x_2 = \text{senior}, x_3 = 21-30 | A)$$

$$= \text{prob}(x_1 = a_1 | A) \cdot \text{prob}(x_2 = a_2 | A) \cdot \text{prob}(x_3 = a_3 | A) \cdot \text{prob}(A)$$

High: A=high	2/4	2/5	2/4	2/11
Med: A=med	2/4	3/5	1/4	6/11
Low: A=low	0	0	2/4	3/11

High: $2/4 \cdot 2/5 \cdot 2/4 \cdot 2/11 = 0.009090909$

Med: $2/4 \cdot 3/5 \cdot 1/4 \cdot 6/11 = 0.040909091$

Low: $0 \cdot 0 \cdot 2/4 \cdot 3/11 = 0$

Of this test instance of attributes the value would be a salary of category medium for this Naïve Bayes Classification.

Part III: Decision Tree (50 points).

Consider the following set of training data. The target feature (attribute) is *Class*, which can have values *text* (text file) or *.exe* (executable file) for different instances, is to be predicted based on other descriptive features (attributes) of the instance.

Instance	Writable	Updatable	Size	Class
1	yes	no	small	text
2	yes	yes	large	text
3	no	yes	medium	text
4	no	no	medium	.exe
5	yes	no	large	.exe
6	no	no	large	.exe

Note for Questions 1 and 2: Internal nodes of trees are annotated with descriptive feature names, leaf nodes are labeled with *Class* values. Sorted instances could be displayed within nodes of the trees.

Question 1 (10 points)

Construct the decision tree (called T1) generated by first considering the descriptive feature *Writable*, then *Updatable*, and then *Size*.

See Back of this page.

Question 2 (30 points)

Compute and construct the decision tree (called T2) using the ID3 algorithm. Show your work of calculations.

See Last Page Front : Back Sides

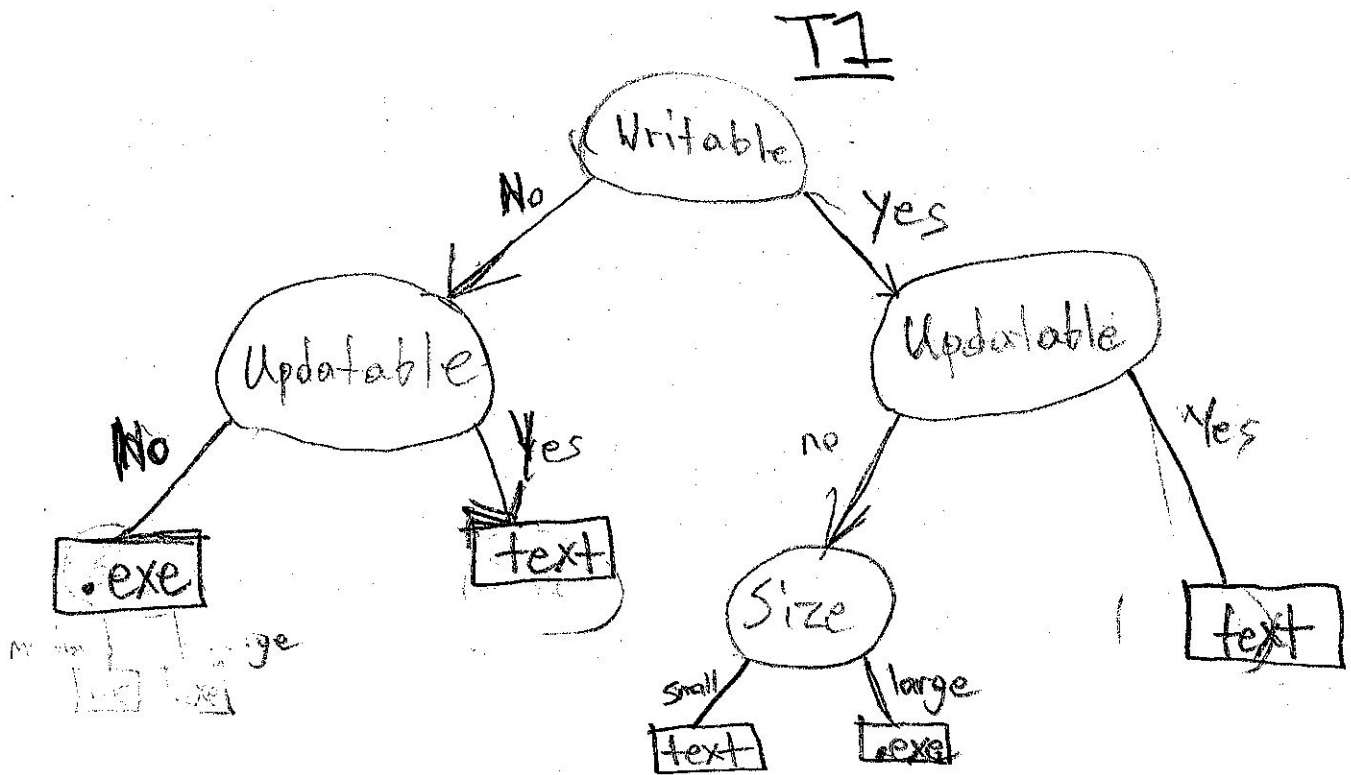
Question 3 (10 points)

How will T1 and T2 classify the following two new instances, respectively?

Instance	Writable	Updatable	Size	Class
1	yes	yes	large	?
2	no	no	small	?

Instance	T1	T2
1	text	text
2	.exe	text

Part III Question 1



Part III Question 2 ID 3

Entropy of Dataset

Entropy(D)

$$= \left(\frac{3}{6}, \frac{3}{6} \right) = - \left[\frac{3}{6} \log_2 \left(\frac{3}{6} \right) + \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right]$$

Information Gain

$$= \boxed{1}$$

$$\text{Entropy(D)} = 1$$

$$H(\text{Writable}) = \left(\frac{3}{6} \left(- \left(\left[\frac{2}{6} \cdot \log_2 \left(\frac{2}{6} \right) \right] + \left[\frac{1}{6} \log_2 \left(\frac{1}{6} \right) \right] \right) \right) + \dots \right.$$

$$\left(\frac{3}{6} \left(- \left(\left[\frac{2}{6} \cdot \log_2 \left(\frac{2}{6} \right) \right] + \left[\frac{1}{6} \log_2 \left(\frac{1}{6} \right) \right] \right) \right) \right) = 0.959147917$$

Writable: Yes
Class: text

Writable: Yes
Class: .exe

Writable: No
Class: text

Writable: No
Class: .exe

$$\text{Info. Gain (Writable)} = 1 - 0.959147917 = \boxed{0.040852083}$$

$$H(\text{Updatable}) = \left(\frac{2}{6} \left(- \left(\left[\frac{2}{6} \cdot \log_2 \left(\frac{2}{6} \right) \right] + \left[0 \right] \right) \right) + \dots \right.$$

$$\left(\frac{4}{6} \left(- \left(\left[\frac{1}{6} \cdot \log_2 \left(\frac{1}{6} \right) \right] + \left[\frac{3}{6} \cdot \log_2 \left(\frac{3}{6} \right) \right] \right) \right) \right) = 0.750543056$$

Updatable: Yes
Class: Text

Updatable: Yes
Class: .exe

Updatable: No
Class: Text

Updatable: No
Class: .exe

$$\text{Info. Gain (Updatable)} = 1 - 0.750543056 = \boxed{0.249456944}$$

$$H(\text{Size}) = \left(\frac{1}{6} \left(- \left(\left[\frac{1}{6} \cdot \log_2 \left(\frac{1}{6} \right) \right] + \left[0 \right] \right) \right) + \left(\frac{2}{6} \left(- \left(\left[\frac{2}{6} \log_2 \left(\frac{2}{6} \right) \right] + \left[\frac{1}{6} \log_2 \left(\frac{1}{6} \right) \right] \right) \right) + \dots \right.$$

$$+ \left(\frac{3}{6} \left(- \left(\left[\frac{2}{6} \log_2 \left(\frac{2}{6} \right) \right] + \left[\frac{1}{6} \log_2 \left(\frac{1}{6} \right) \right] \right) \right) \right) = 0.789849653$$

Size: Small
Class: Text

Size: Small
Class: .exe

Size: Medium
Class: Text

Size: Medium
Class: .exe

Size: Large
Class: Text

Size: Large
Class: .exe

$$\text{Info Gain (Size)} = 1 - 0.789849653 = \boxed{0.210150347}$$

I03

T2

* Based off Info Geo.
Calculation

