

# 10-601: Homework 1

Due: 18 September 2014 11:59pm (Autolab)

TAs: Abhinav Maurya, Jingwei Shen

Name: Alvin Chou

Andrew ID: alvincho

Please answer to the point, and do not spend time/space giving irrelevant details. You should not require more space than is provided for each question. If you do, please think whether you can make your argument more pithy, an exercise that can often lead to more insight into the problem. Please state any additional assumptions you make while answering the questions. You need to submit a single PDF file on autolab. Please make sure you write legibly for grading.

You can work in groups. However, no written notes can be shared, or taken during group discussions. You may ask clarifying questions on Piazza. However, under no circumstances should you reveal any part of the answer publicly on Piazza or any other public website. The intention of this policy is to facilitate learning, not circumvent it. Any incidents of plagiarism will be handled in accordance with CMU's Policy on Academic Integrity.

---

## \*: Code of Conduct Declaration

---

- Did you receive any help whatsoever from anyone in solving this assignment? (Yes) / No.
- If you answered *yes*, give full details: Graeme Rock & Davei Wang explained (e.g. Jane explained to me what is asked in Question 3.4) a few questions to me.
- Did you give any help whatsoever to anyone in solving this assignment? (Yes) / No.
- If you answered *yes*, give full details: I provided pointers to Graeme (e.g. I pointed Joe to section 2.3 to help him with Question 2). Rock & Davei Wang on a couple questions.

---

## 1: The truth will set you free. (TA:- Abhinav Maurya)

---

State whether true (with a brief reason) or false (with a contradictory example). Credit will be granted only if your reasons/counterexamples are correct.

(a) During decision tree construction, if you reach a node where the maximum information gain for a node split using any attribute is zero, then all training examples at that node have the same label.

False, if the remaining attributes provides the same amount of entropy as the first node, the information gain would be 0 even though the data may have different labels. [2 points]

(b) Whenever a set  $S$  of labeled instances is split into two sets  $S_1$  and  $S_2$ , the average entropy will not increase, irrespective of the split attribute or the split point.

(b) We have two coins - an unbiased one with probability  $p_1 = 1/2$  of showing heads on a toss, and a biased one with probability  $p_2 = 1/3$  for showing heads. We do 100 tosses. Each time we choose one of the two coins. With an unknown probability  $p$ , we choose the biased coin, and with probability  $1 - p$ , we choose the unbiased one. And we observe 40 heads during the 100 tosses. Write down the MLE estimate of parameter  $p$  and explain it. (You do not have to derive it.)

$$p(\text{heads}) = p(p_1) + (1-p)(p_2) = \frac{40}{100} \quad [3 \text{ points}]$$

$$p(0.5) + (1-p)(0.33) = 0.4$$

$$\boxed{p = 0.41}$$

This means that a given probability of ~~getting heads~~ using the unbiased coin of 41% best matches the observed 40/100 heads scenario with the information provided.

### 3: Three Prisoners and a Warden (TA:- Jingwei Shen)

Three prisoners - A, B, and C - are on death row. The governor decides to pardon one of the three and chooses the prisoner to pardon at random. He informs the warden of his choice but requests the name to be kept as a secret.

Having heard of the pardon rumor through grapevine, A tries to get the warden to tell him his fate. The warden refuses. Then A asks which of B or C will be executed. The warden thinks a while and tells A that B is to be executed. (Assume that the warden picks a random legal answer for A's question).

(a) Let  $A, B, C$  denote the event that A, B, C will be pardoned respectively. Let  $!B$  denote the event that the warden says B will die. Compute  $P(A | !B)$ . Does the chance of A's survival increase with the additional information about B's death? (Hint: compare  $P(A | !B)$  and  $P(A)$ ).

$$P(A | !B) = \frac{P(!B | A)P(A)}{P(!B)} \quad P(A) = \frac{1}{3} = 33\% \quad [3 \text{ points}]$$

$$P(A | !B) = \boxed{33\%}$$

(b) Suppose A reveals all of the above to C. Show the probability of C surviving at this time is  $2/3$ . (Hint: Prove  $P(C | !B) = 2/3$ ).

$$\begin{aligned} P(C | !B) &= \frac{P(!B | C)P(C)}{P(!B)} \\ &= \frac{P(!B | C)P(C)}{P(!B | A)P(A) + P(!B | B)P(B) + P(!B | C)P(C)} \\ &= \frac{1(0.33)}{\frac{1}{2}(0.33) + 0 + 1(0.33)} = \boxed{\frac{2}{3}} \end{aligned} \quad [3 \text{ points}]$$

$$\begin{aligned}
 \int_0^{\frac{1}{2}} f(x) &= \int_0^{\frac{1}{2}} (x + \frac{1}{2})^2 \\
 &= \left[ \frac{(x + \frac{1}{2})^3}{3} \left( \frac{2}{x^2} \right) \right]_0^{\frac{1}{2}} \\
 &= \frac{1}{12} \left( \frac{1}{2} \right) \left( (1) + 3 + 3 \right) \\
 &= \boxed{\frac{7}{24}}
 \end{aligned}$$

### 5: Nearest neighbors to the rescue. (TA:- Jingwei Shen)

(a) Consider two classes  $C_1, C_2$  in the two-dimensional space. The data from class  $C_1$  are uniformly distributed in a circle of radius  $r$ . The data from class  $C_2$  are uniformly distributed in another circle of radius  $r$ . The centers of two circles are at a distance greater than  $4r$ . Show that the accuracy of 1-NN is greater than or equal to the accuracy of  $k$ -NN, where  $k$  is an odd integer and  $k \geq 3$ .

Test data is in  $C_1$

a) If training data is such that  $C_1$  - 1 data &  $C_2$  - 2 data  
 $1NN \rightarrow C_1$        $3NN \rightarrow C_2$

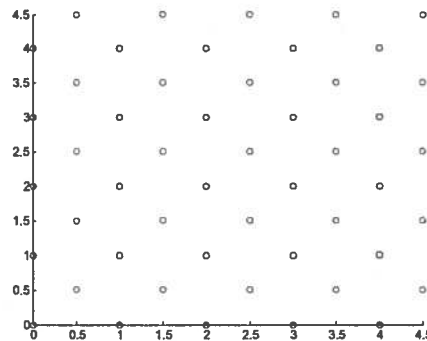
b) If training data is such that  $C_1$  - 2 data &  $C_2$  - 1 data  
 $1NN \rightarrow C_1$        $3NN \rightarrow C_1$

c) If training data is such that  $C_1$  - 3 data &  $C_2$  - 0 data  
 $1NN \rightarrow C_1$        $3NN \rightarrow C_1$

d) If training data is such that  $C_1$  - 0 data &  $C_2$  - 3 data  
 $1NN \rightarrow C_2$        $3NN \rightarrow C_2$

[3 points]	
1NN	3NN
✓	X
✓	✓
✓	✓
X	X

Figure 1: Q4 Dataset

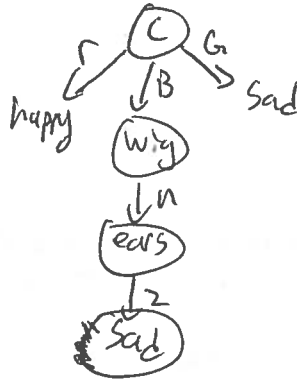


(b) In the dataset shown in figure 1, what is the leave-one-out accuracy of the  $k$ -NN method when  $k = 2$ ? Remember that a data point cannot be considered its own neighbor since it is left out. (Ignore the datapoints that have an output tie for  $k = 2$  nearest neighbors.)

0%.  $k=2$  will always yield the wrong label.

[2 points]

(c) Draw the full decision tree that would be learned for this data (assume no pruning).



[3 points]

(d) What would be the training set error for this dataset? Express your answer as the percentage of records that would be misclassified.

$\boxed{1/9}$

[2 points]

### 7: Digging up the dense binary tree. (TA:- Abhinav Maurya)

Consider the following data with three binary attributes, where  $x^i$  denotes the  $i^{\text{th}}$  datapoint,  $x_j$  denotes the  $j^{\text{th}}$  feature of the datapoint, and  $y$  denotes the class label:-

	$x_1$	$x_2$	$x_3$	$y$
$x^0$	0	0	0	0
$x^1$	0	0	1	1
$x^2$	0	1	0	1
$x^3$	0	1	1	0
$x^4$	1	0	0	1
$x^5$	1	0	1	0
$x^6$	1	1	0	0
$x^7$	1	1	1	1

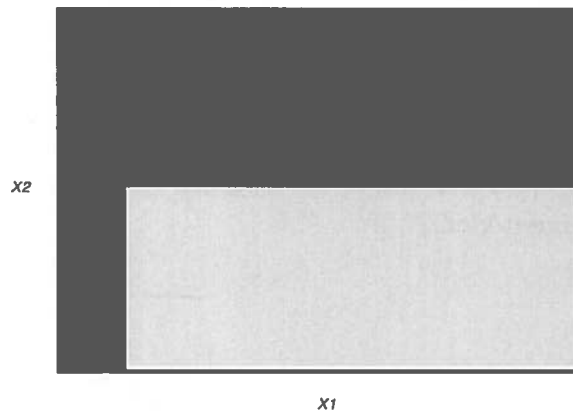
(a) Draw the decision tree for the above dataset using the entropy criterion to decide node splits (assume no pruning).

[3 points]

**8: On the hardness of learning optimal binary decision trees (TA:- Abhinav Maurya)**

In figure 2, assume that the rectangular region consisting of two features  $x_1$  and  $x_2$  is densely packed with points. The red, green, and yellow subrectangles represent the three classes  $C_1$ ,  $C_2$ , and  $C_3$  of datapoints. The  $x_1 \times x_2$  dimensions of the red, green, and yellow rectangles are  $1 \times 6$ ,  $7 \times 3$ , and  $7 \times 3$  respectively. The red rectangle is uniformly populated with 6,000 datapoints of class  $C_1$ . The green rectangle is uniformly populated with 42,000 datapoints of class  $C_2$ . The yellow rectangle is uniformly populated with 42,000 datapoints of class  $C_3$ .

Figure 2: A 2D dataset with three classes



(a) What is the minimum number of nodes that a decision tree needs to have in order to classify the above dataset correctly?

2  $x_1$  then  $x_2$

[2 points]

(b) What is the number of nodes in the decision tree trained on the above dataset using the entropy criterion?

3  $x_2$  then  $x_1$

[2 points]

(c) Are the number of nodes in the two cases identical or different? Why do you think that is?

Different. Since the entropy case will always split along the line where it could categorize the largest amount of datapoints,

[3 points]