

CS390DM0 Spring 2017 Midterm 1 *solutions*, 21 February, 2017

Prof. Chris Clifton

I would expect an A student to score at least 36, a B student at least 30, and a C student at least 19.

**Turn Off Your Cell Phone.** Use of any electronic device during the test is prohibited. As previously noted, you are allowed notes: Up to two sheets of 8.5x11 or A4 paper, single-sided (or one sheet double-sided).

Time will be tight. If you spend more than the recommended time on any question, **go on to the next one**. If you can't answer it in the recommended time, you are either giving too much detail or the question is material you don't know well. You can skip one or two parts and still demonstrate what I believe to be an A-level understanding of the material.

Note: It is okay to abbreviate in your answers, as long as the abbreviations are unambiguous and reasonably obvious.

In all cases, it is important that you give some idea of how you derived the answer, not simply give an answer. Setting up the derivation correctly, even if you don't carry out the calculations to get the final answer, is good for nearly full credit.

**The last page of the exam contains an example database that you will use for several of the questions. You may tear it off for easy reference.**

## 1 Probability (5 minutes, 4 points)

There are three cards. The first is green on both sides, the second is red on both sides, and the third is green on one side and red on the other. Consider the scenario where a card is chosen at random and one side is shown (also chosen at random). If the side shown is green, what is the probability that the other side is also green?

$$P(\text{2nd side green} \mid \text{1st side green}) = \frac{P(\text{1st and 2nd side green})}{P(\text{1st side green})}$$

$$= \frac{\frac{1}{3}}{\frac{3}{6}} = \frac{6}{9} = \boxed{\frac{2}{3}}$$

**Courtesy of Brian Lee.** 1 for setting it up as probability, 1 for noting conditional on first draw, 1 for noting that which card is the first draw is conditional on seeing green, 1 for correct

## 2 Conditional Probability (5 minutes, 4 points)

Prove that if  $A$  and  $B$  are independent events, then  $P(A|B) = P(A)$ .

A and B independent:

$$P(A|B) = P(A)P(B)$$

$$P(A|B) = \frac{P(AB)}{P(B)}$$

$$= \frac{P(A)P(B)}{P(B)}$$

$$= P(A) \quad \# \text{ (shown)}$$

4

Courtesy of a student who prefers to be anonymous. 1 for setting up as proof, 1 for definition of independent, 1 for some working it out, 1 for complete

### 3 Exploratory Data Analysis (10 minutes, 6 points)

These questions all refer to the dataset on the last page.

- A. Find two attributes that are highly correlated. Mathematically characterize this correlation relation.

The attributes A and B are highly correlated. We can take  $p=1$ ,  $q=2$  and  $r=3$  for attribute A and  $x=1$   $y=2$  for attribute B.  
 Mathematically correlation coefficient =  $\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sigma_x \sigma_y}$  where  $x_i = i^{\text{th}}$  value in column A and  $y_i = i^{\text{th}}$  value in column B.  
 Hence this correlation coefficient should be between 0.8 and 1, as they seem to be highly correlated.

Courtesy of Abhishek Shrinivasan. 1 for picking two that are correlated, 1 for a definition of correlation, 1 for working it out

- B. Are A and B conditionally independent with respect to the class? Explain.

$A = \{p, r, q\}$   $B = \{x, y\}$   $\text{Class} = \{+, -\}$  yes b/c the values do not depend on each other.  
 $\frac{P(A, B | \text{Class})}{P(\text{Class})} = \frac{P(A | \text{Class}) \cdot P(B | \text{Class})}{P(\text{Class})} = \frac{\{0\}}{2} \cdot \frac{\{0\}}{2} = 0$

Courtesy of Simone May. 1 for showing knowledge of conditional independence, 1 for correct answer, 1 for good explanation

### 4 Information Gain (8 minutes, 5 points)

Using the dataset on the last page,

- A. Calculate information gain for attribute C

$$IG = Entropy(\text{total}) - \left( \frac{1}{2} Entropy(C_1) + \frac{1}{2} Entropy(C_2) \right)$$

$$E_{\text{total}} = \frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \left( \frac{1}{2} \right) = 1$$

$$E_{C_1} = \frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \left( \frac{1}{2} \right) = 1$$

$$E_{C_2} = \frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \left( \frac{1}{2} \right) = 1$$

$$IG = 1 - \left( \frac{1}{2} E_{C_1} + \frac{1}{2} E_{C_2} \right)$$

$$IG = 0$$

Courtesy of a student who prefers to be anonymous. 1 for setting it up properly, 1 for working it out

- B. When choosing what attribute to use using information gain, do we choose the one with highest or lowest information gain? Explain why.

higher gain means lower entropy, which means that for that attribute, the values are better differentiated, hence making that attribute very suitable for us to use

Courtesy of Marshia Seto. 1 for "high", 1 for good explanation  
1 for showing some understanding of IG, can come from either part.

## 5 Classification (12 minutes, 8 points)

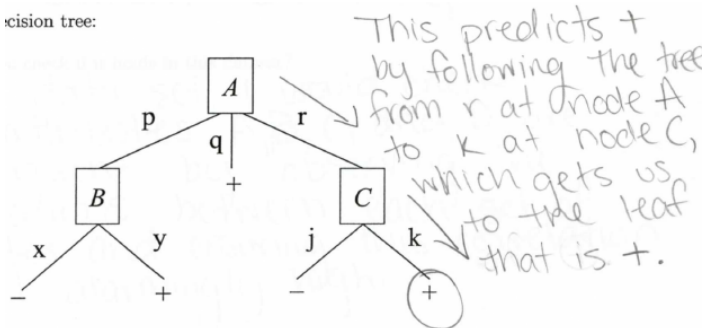
Given the dataset on the last page, what would we predict as the class of the instance:

**r x k 0.8**

if we used each of the following for classification. Give a brief discussion of why that is the prediction, and any assumptions you have to make (one or two sentences).

- A. The following decision tree:

Decision tree:



Courtesy of a student who prefers to be anonymous. 1 for correct answer, 1 for path

- B. 1-Nearest Neighbor

One closest neighbor would be r x j 0.8 therefore class is (-ve). Assumptions for B & C:  
 $d(x, y) = d(j, k)$ ;  $d(0.8, 0.9) > 0.00$

Courtesy of a student who prefers to be anonymous. 1 for correct answer, 1 for noting neighbor

C. 3-Nearest Neighbor

C. 3-nearest neighbor  
 — this is the prediction because the 3 nearest neighbors are  $-(r \times j \ 0.8)$ ,  $-(p \times k \ 0.7)$ , and  $+(r \ y \ k \ 0.9)$  and so the majority vote is  $-$ . I assumed all attributes are equally important, and since att. D is continuous, I used the Euclidean distance metric. (which is why  $(p \times k \ 0.3)$  is not a nearest neighbor.)

**Courtesy of Robyn Berkel.** 1 for correct answer, 1 for noting neighbor 1 point for discussion of distance function in either part

D. Is there a solid technical reason why I asked 1-NN and 3-NN, but not 2-NN?

Yes, KNN is majority vote, thus you must use odd values for k. With  $k=2$  there is chance for a stalemate.

**Courtesy of a student who prefers to be anonymous.** 1 for noting voting, need to break ties

## 6 Naïve Bayes (8 minutes, 5 points)

A. What is the key assumption made by Naïve Bayes?

what is the key assumption made by naive bayes:  
 Conditional independence of the attributes with respect to the class holds! That is the assumption.

**Courtesy of Avik Mikhija.** 1 for conditionally independent

B. How would you check if it holds in this dataset?

Calculate  
 $P(x_i, x_j | \text{class} = -) = P(x_i | \text{class} = -) P(x_j | \text{class} = -)$   
 and  $P(x_i, x_j | \text{class} = +) = P(x_i | \text{class} = +) P(x_j | \text{class} = +)$  for every pair of attributes. If any conditions fail, the assumption does not hold.

**Courtesy of Caleb Beth.** 1 for description. 1 for showing understanding of what CI means for NB

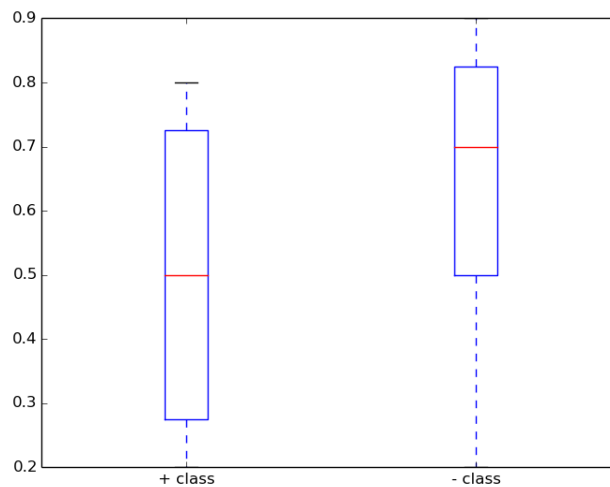
C. Is it possible for Naïve Bayes to have high accuracy even if this assumption is violated? Explain why or why not.

Yes. If the attributes are not mutually independent but are positively correlated, Naive Bayes still has robust performance most of the time. E

Courtesy of a student who prefers to be anonymous. 1 for yes, 1 for either example or reasonable discussion of why

## 7 Overlap of IQRs (5 minutes, 3 points)

Suppose that we plot the Inter-Quartile ranges for attribute  $D$  for the + and - classes, and we have the following plot:



Based on this plot, do you think attribute  $D$  would be an effective feature for distinguishing between the + and - class? Explain why or why not.

I don't believe  $D$  is a good choice since the entirety of the + class' third quartile is contained in the - class' IQR and the second quartile of the - class is contained in the + class' IQR. They overlap far too much to be effective.

Courtesy of Ian Spooner. 1 for "no", 1 for some reasonable discussion, 1 for noting overlap of distributions not working well

## 8 Probability (10 minutes, 4 points)

Suppose that 30 percent of computer owners use an Apple machine, 50 percent use a Windows machine, and 20 percent use Linux. Suppose that 65 percent of Apple users have succumbed to a computer virus, 82 percent of Windows users get the virus, and 50 percent of Linux users get the

virus. We select a person at random and learn that their system was infected with the virus. What is the probability that the person is a Windows user?

$V$ : system infected with virus.  
 $W$ : windows user  
 $A$ : apple user  
 $L$ : Linux user

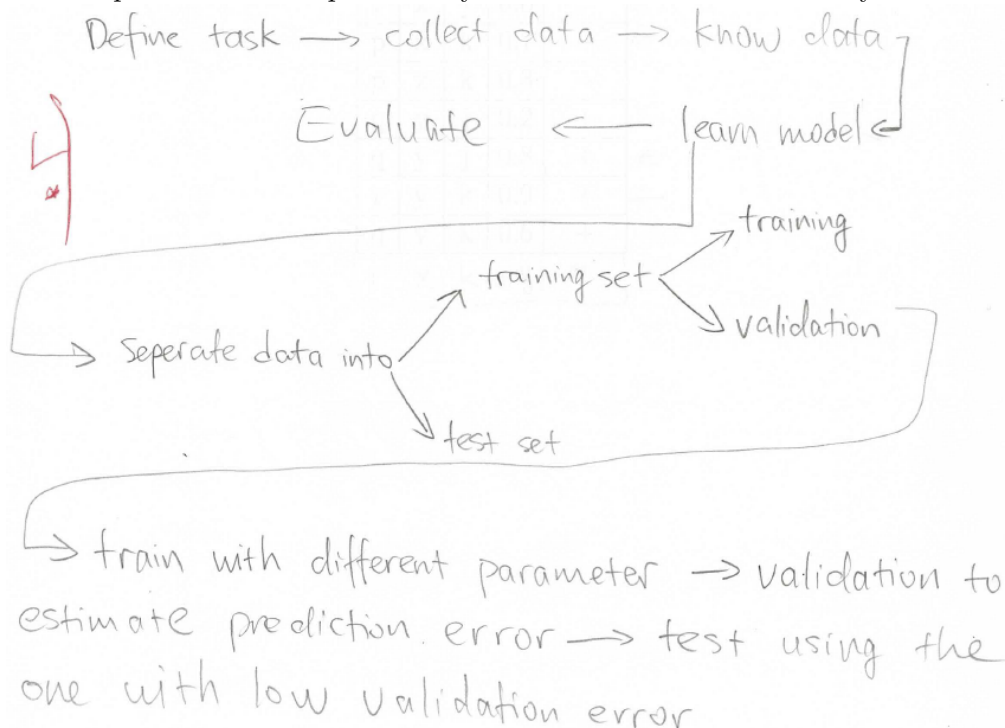
$$P(W|V) = \frac{P(V|W)P(W)}{P(V|W)P(W) + P(V|A)P(A) + P(V|L)P(L)}$$

$$= \frac{0.82 \times 0.5}{0.82 \times 0.5 + 0.65 \times 0.3 + 0.5 \times 0.2} = \frac{0.41}{0.705} = 0.58$$

**Courtesy of Yixuan Ding.** 1 for setting up problem, 1 for noting that having virus changes things, 1 for correct setup, 1 for answer

### 9 Data Mining process (10 minutes, 4 points)

Given the dataset on the last page, I would like you to build a classifier. Discuss the process you would follow. This should be independent of the type of classifier. You want to build as accurate a classifier as possible. Also explain how you would evaluate the accuracy of the resulting classifier.



**Courtesy of a student who prefers to be anonymous.**

There were many points that could have been covered in answering this question, and no one answer covered everything that could have been included. However, the answer above covers several of the most important points in a particularly clear and succinct manner.

1 for division of existing data into test and training (not enough to say "I'll get more test data")

- how do you know you can do so?), 1 for discussion of computing accuracy. In addition, a point was given for each of several other relevant things, including discussion of steps, cross-validation (which counts as discussing splitting into training and test data), data exploration, model selection (discussing how or why you pick a particular model, not just picking one), data cleaning/handling missing values, or giving a formula for any of this. To get full credit, you had to mention both test/training division and discuss computing accuracy, and you couldn't state anything that was actually incorrect. E.g., a common error was saying that you would use the test set for both model selection and reporting accuracy; these would need to be two different sets (validation and test).

**Dataset to be used for questions 3, 4, 5, 6, and 9.**

| <i>A</i> | <i>B</i> | <i>C</i> | <i>D</i> | <i>class</i> |
|----------|----------|----------|----------|--------------|
| p        | x        | j        | 0.2      | −            |
| r        | x        | j        | 0.8      | −            |
| p        | x        | k        | 0.7      | −            |
| p        | x        | k        | 0.3      | −            |
| q        | y        | j        | 0.2      | +            |
| q        | y        | j        | 0.8      | +            |
| r        | y        | k        | 0.9      | +            |
| q        | y        | k        | 0.6      | +            |