### CS378 Assignment 0: Linear Algebra, Probability, and Python Warmup

#### Due date: Friday, January 24 at 11:59pm CST

**Academic Honesty:** Reminder that assignments should be completed independently by each student. See the syllabus for more detailed discussion of academic honesty. Limit any discussion of assignments with other students to clarification of the requirements or definitions of the problems, or to understanding the existing code or general course material. Never directly discuss details of the problem solutions. Finally, you may not publish solutions to these assignments or consult solutions that exist in the wild.

**Goals** The main goal of this assignment is for you to **assess whether you have adequate preparation for the course.** It's fine to not be familiar with every concept here. However, if you find yourself struggling with much of this assignment, you should ask the course staff whether this course is appropriate for you given your background. This assignment is designed to take around 2 hours.

**Grading** The assignment is out of 100 points. Note that it is worth half as much as the other assignments, so your grade on Canvas will show up as your points divided by 2.

### 1 Linear Algebra (25 points)

**Q1 (15 points)** For each of the following matrices, give the answer or write "undefined" if the operation is invalid. You do not need to show work.

a) 
$$\begin{bmatrix} 1 & 2 & 4 \\ 3 & 4 & 2 \end{bmatrix} \begin{bmatrix} 4 \\ 2 \end{bmatrix}$$
 b)  $\begin{bmatrix} 1 & 2 & 4 \\ 3 & 4 & 2 \end{bmatrix} \begin{bmatrix} 4 \\ 5 \\ 2 \end{bmatrix}$  c)  $\begin{bmatrix} 1 & 2 & 4 \\ 3 & 4 & 2 \end{bmatrix} \begin{bmatrix} 4 & 2 \end{bmatrix}$  d)  $\begin{bmatrix} 6 \\ 2 \\ 4 \end{bmatrix}^{\top} \begin{bmatrix} 5 \\ 2 \\ 1 \end{bmatrix}$ 

**Q2** (10 points) Write a matrix operation capturing the following computation. Your answer should be a mathematical expression involving the vectors/matrices A, B, and C. Your math expression does not need to account for initialization of A, B, and C; it only needs to return the same value as sum given the same inputs.

```
A = np.rand(4)
B = np.rand(4,3)
C = np.rand(3)
sum = 0.0
for i in range(0,4):
   for j in range(0,3):
      sum += A[i] * B[i,j] * C[j]
return sum
```

#### 2 **Probability (35 points)**

Q3 (10 points) Consider the following joint distribution:

P(X,Y)	Y = 1	Y = 2	Y = 3
X = 1	0.1	0.2	0.2
X = 2	0.05	0.1	0.1
X = 3	0.1	0.1	0.05

- a) What is P(X|Y=2)?
- **b)** What is P(Y|X=1)?
- c) Are X and Y independent? Justify your answer.

**Q4 (10 points)** Suppose you have a distribution P(X, Y) where  $X \in \{0, 1\}$  and  $Y \in \{0, 1\}$ . You know that the marginal distribution P(X) = [0.5, 0.5] and P(Y) = [0.2, 0.8].

a) If X and Y are independent, what do we know about the value of the joint probability P(X = 0, Y = 0)? If it is not exactly knowable, give upper and lower bounds as precise as you can.

b) If X and Y are **not** independent, what do we know about the value of the joint probability P(X = 0, Y = 0)? If it is not exactly knowable, give upper and lower bounds as precise as you can.

Q5 (15 points) The binary entropy of a random variable X with discrete domain D is defined as:

$$H(X) = -\sum_{x \in D(X)} P(x) \log_2 P(x)$$

a) Compute the entropy of P(X) =Multinomial $(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n})$ , the uniform distribution over n variables. Your answer should be written symbolically.

**b**) When you have a joint distribution over X and Y, entropy is defined as:

$$H(X,Y) = -\sum_{x \in D(X)} \sum_{y \in D(Y)} P(x,y) \log P(x,y)$$

How does this relate to the entropy of the marginal distributions H(X) and H(Y) when X and Y are independent?

# 3 Language Basics / Coding Warmup (35 points)

In this part of the assignment, you will read in and do some basic manipulation of a text corpus. Included with the assignment is a file nyt.txt containing 8860 sentences taken from New York Times articles, one sentence per line. You should implement your solutions in a file called a0.py.

**Q6** (15 points) Here you will investigate tokenization schemes. Tokenization is the process of splitting raw text into words. In English, this involves splitting out punctuation and contractions (*shouldn't* becomes *should 'nt*) and is typically done with rules. In other languages like Chinese or Arabic, the process can be significantly more involved.

**a**) What are the 10 most frequent words in this dataset using whitespace tokenization? That is, split each sentence into words simply based on where the spaces are. List each word and its count and describe any patterns you see.

**b**) What are the 10 most frequent words in this dataset using smarter tokenization? You can either use the tokenizer in tokenizer.py or invoke another tokenizer like NLTK (nltk.word\_tokenize(sentence)) after importing NLTK) or spaCY:<sup>1</sup>

```
from spacy.lang.en import English
nlp = English()
tokenizer = nlp.Defaults.create_tokenizer(nlp)
tok_sent = tokenizer(sentence) # Returns a Doc object; see the spaCY docs
```

List each word and its count and describe any patterns you see.

c) Explain in a few sentences how these differences in tokenization could affect a downstream text processing system. Discuss at least two ways.

**Q7** (20 points) In this part, we are going to confirm a phenomenon known as Zipf's Law. A word has *rank* n if it is the nth most common word. Zipf's Law states that the frequency of a word in a corpus is inversely proportional to its rank. Roughly speaking, this means that the fifth most common word should be five times less frequent than the most common word, and the tenth most common word should occur half as much as the fifth most common word.

a) Make a plot of inverse rank vs. word count for the smart tokenization scheme. Inverse rank is the reciprocal of the rank of the word: 1 for the most frequently occurring word,  $\frac{1}{2}$  for the second most,  $\frac{1}{3}$  for the third most, etc. Include your plot in your submission. Matplotlib is a good tool to use, but Excel/Matlab/Gnuplot/others are okay too.

**b**) Based on the plot, where does Zipf's law appear to hold? Are there any outliers?

c) Look at your list of most frequent words. Identify **three words** out of the top 100 words in this dataset that you believe are unusually common in this data compared to written English text overall, and for each of these words, say why you think it is more common than expected here.

<sup>&</sup>lt;sup>1</sup>Instructions to install NLTK: https://www.nltk.org/install.html and spaCy: https://spacy.io/usage

# Submission

If you wish to complete either Parts 1 or 2 of this assignment on paper, you may. Submit it in person by 5pm Friday at Greg's office (slip it under his door if he's not there).

You should upload two files to Canvas. Please upload these as individual files; do not upload a zip/tgz.

- 1. A PDF of your answers/output; this might contain your answers to the whole assignment, or just Part 3
- 2. A single python file called a0.py. There are no specific requirements for what the code has to do, but it should provide sufficient evidence that you have completed this part of the assignment. For example, displaying the plot from part (a) is fine, or printing out the values that you used to create it is fine too.

Slip Days Slip days may be used on this assignment, per the policies described in the syllabus.