

## Q1. Generative Modeling(10 Marks) [10]

Recall your recent work in Assignment 2. Given the following sample text (from Eliot's *The Hollow Men*) train a model for poem generation.

Between the idea  
And the reality  
Between the motion  
And the act  
Falls the Shadow  
For Thine is the Kingdom

Between the conception  
And the creation  
Between the emotion  
And the response  
Falls the Shadow  
Life is long  
For Thine is the Kingdom

...

We want to generate poems such as the above. Given that we have a large corpus of Eliot's work  $T$ , write an algorithm that can learn a bigram model to generate a similar couplet. Be specific on how words and patterns are selected for generating each verse in the couplet.

### Algorithm 1 GENERATECOUPLET( $T$ )

- Create a dictionary named bigrams and tokenize corpus.
  - ~~start with~~ loop till the length of corpus.
    - If  $corpus[i]$  is not in dictionary add it and extend a dictionary from it ~~and~~ enter key  $corpus[i+1]$  with value  $corpus[i][corpus[i+1]] = 1$
    - if  $corpus[i]$  is in dictionary check if  $corpus[i+1]$  is in  $corpus[i]$  dictionary
    - if present add 1 ~~and~~ else enter it and assign value 1 to it.
  - once the bigram dictionary is created create a start words list consisting of the most used start words.
- 1: return C ▷ generated couplet

continued on page 5

9

Q2. Language Modeling

(15 Marks) [5+10]

- a) How many bigrams can be generated from the following sentence, after replacing punctuations by a single space. List them and calculate their probabilities in the bigram language model.

$s = \text{Big fish eat fish, eat fish, eat fish eat small fish.}$

5 ~~words~~ Bigrams

$$P(\text{fish} | \text{Big}) = \frac{1}{4} \quad (\text{eat} | \text{fish}) \quad (\text{fish} | \text{eat}) \quad (\text{small} | \text{eat})$$

$$P(\text{fish} | \text{small}) = 1$$

$$P(\text{eat} | \text{fish}) = \frac{4}{5}$$

$$P(\text{fish} | \text{eat}) = \frac{3}{4}$$

$$P(\text{small} | \text{eat}) = \frac{1}{4}$$

- (b) Perplexity measures how good a model is at predicting an unseen test set. Consider two bigram models that use Laplace Smoothing and are trained on the two corpora given below:

$c_1 = \text{oysters eat oysters eat oysters eat small oysters}$

$c_2 = \text{oysters eat oysters}$

Calculate the perplexity for the two models given that the test sentence is

$s = \text{oysters eat oysters eat oysters}$

bigrams in  $s$ : for  $c_1$

$P(\text{eat}|\text{oysters})$ , seeing how this occurs in  $c_1$

$$= \frac{3+1}{4+3} = \frac{4}{7}$$

$P(\text{oysters}|\text{eat})$

$$= \frac{2+1}{3+3} = \frac{3}{6} = \frac{1}{2}$$

$$P(\text{eat}|\text{oysters}) = \frac{3+1}{4+3} = \frac{4}{7}$$

$$P(\text{oysters}|\text{eat}) = \frac{2+1}{3+3} = \frac{1}{2}$$

for  $c_2$

$$P(\text{eat}|\text{oysters}) = \frac{1+1}{2+2} = \frac{2}{4} = \frac{1}{2}$$

$$P(\text{oysters}|\text{eat}) = \frac{1+1}{1+2} = \frac{2}{3}$$

$$P(\text{eat}|\text{oysters}) = \frac{1}{2}$$

$$P(\text{oysters}|\text{eat}) = \frac{2}{3}$$

$$P_{c_1} = \sqrt[5]{\frac{1}{\frac{4}{7} \times \frac{1}{2} \times \frac{4}{7} \times \frac{1}{2}}} = 1.65$$

$$P_{c_2} = \sqrt[5]{\frac{1}{\frac{1}{2} \times \frac{2}{3} \times \frac{1}{2} \times \frac{2}{3}}} = 1.55$$

~~$c_1$~~   $c_2$  has lesser perplexity



### Q3. Spelling Correction

(18 Marks) [10+8]

Given that the following table that shows the confusion matrix for character transposition in spelling errors and that the probability of transposition is defined as  $P(t|c) = \frac{\text{trans}[c, c+1]}{\text{count}[c, c+1]}$ .

trans[X, Y] = number of times XY was written as YX

|       | a   | b | c  | d  | e   | f | g  | h  | i   | j | k | l  | m  | n  | o   | p  | q | r   | s  | t  | u   | v  | w | x | y  | z | Total |
|-------|-----|---|----|----|-----|---|----|----|-----|---|---|----|----|----|-----|----|---|-----|----|----|-----|----|---|---|----|---|-------|
| a     | 0   | 0 | 0  | 0  | 1   | 0 | 0  | 0  | 19  | 0 | 1 | 14 | 4  | 25 | 10  | 3  | 0 | 27  | 3  | 8  | 31  | 0  | 0 | 0 | 0  | 0 | 145   |
| b     | 0   | 0 | 0  | 0  | 2   | 0 | 0  | 0  | 0   | 0 | 0 | 1  | 1  | 0  | 2   | 0  | 0 | 0   | 2  | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 125   |
| c     | 0   | 0 | 0  | 0  | 1   | 0 | 0  | 1  | 85  | 0 | 0 | 15 | 0  | 0  | 13  | 0  | 0 | 0   | 3  | 0  | 0   | 7  | 0 | 0 | 0  | 0 | 10    |
| d     | 0   | 0 | 0  | 0  | 0   | 0 | 0  | 0  | 7   | 0 | 0 | 0  | 0  | 0  | 0   | 0  | 0 | 1   | 0  | 0  | 2   | 0  | 0 | 0 | 0  | 0 | 247   |
| e     | 1   | 0 | 4  | 5  | 0   | 0 | 0  | 0  | 0   | 0 | 0 | 21 | 0  | 16 | 0   | 2  | 0 | 29  | 6  | 0  | 85  | 0  | 0 | 0 | 2  | 0 | 13    |
| f     | 0   | 0 | 0  | 0  | 0   | 0 | 0  | 0  | 12  | 0 | 0 | 1  | 0  | 0  | 0   | 0  | 0 | 0   | 0  | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 28    |
| g     | 4   | 0 | 0  | 0  | 2   | 0 | 0  | 0  | 0   | 0 | 0 | 1  | 0  | 15 | 0   | 0  | 0 | 3   | 0  | 0  | 3   | 0  | 0 | 0 | 0  | 0 | 37    |
| h     | 12  | 0 | 0  | 0  | 15  | 0 | 0  | 0  | 0   | 0 | 0 | 0  | 0  | 0  | 0   | 0  | 0 | 0   | 0  | 10 | 0   | 0  | 0 | 0 | 0  | 0 | 213   |
| i     | 15  | 8 | 31 | 3  | 66  | 1 | 3  | 0  | 0   | 0 | 0 | 0  | 0  | 0  | 11  | 0  | 1 | 13  | 23 | 35 | 0   | 6  | 0 | 0 | 0  | 3 | 0     |
| j     | 0   | 0 | 0  | 0  | 0   | 0 | 0  | 0  | 0   | 0 | 0 | 0  | 0  | 0  | 0   | 0  | 0 | 0   | 0  | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 2     |
| k     | 0   | 0 | 0  | 0  | 2   | 0 | 0  | 0  | 0   | 0 | 0 | 0  | 0  | 0  | 0   | 0  | 0 | 0   | 0  | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 73    |
| l     | 11  | 0 | 0  | 12 | 20  | 0 | 1  | 0  | 4   | 0 | 0 | 0  | 0  | 0  | 1   | 3  | 0 | 0   | 1  | 1  | 3   | 9  | 0 | 0 | 7  | 0 | 35    |
| m     | 9   | 0 | 0  | 0  | 20  | 0 | 0  | 0  | 0   | 0 | 0 | 0  | 0  | 2  | 0   | 0  | 0 | 0   | 0  | 0  | 4   | 0  | 0 | 0 | 0  | 0 | 57    |
| n     | 15  | 0 | 6  | 2  | 12  | 0 | 8  | 0  | 1   | 0 | 0 | 0  | 3  | 0  | 0   | 0  | 0 | 0   | 6  | 4  | 0   | 0  | 0 | 0 | 0  | 0 | 44    |
| o     | 5   | 0 | 2  | 0  | 4   | 0 | 0  | 0  | 6   | 0 | 0 | 1  | 0  | 5  | 0   | 1  | 0 | 11  | 1  | 1  | 0   | 0  | 7 | 1 | 0  | 0 | 34    |
| p     | 17  | 0 | 0  | 0  | 4   | 0 | 0  | 1  | 0   | 0 | 0 | 0  | 0  | 0  | 1   | 0  | 0 | 5   | 3  | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 0     |
| q     | 0   | 0 | 0  | 0  | 0   | 0 | 0  | 0  | 0   | 0 | 0 | 0  | 0  | 0  | 0   | 0  | 0 | 0   | 0  | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 0     |
| r     | 12  | 0 | 0  | 0  | 24  | 0 | 3  | 0  | 14  | 0 | 2 | 2  | 0  | 7  | 30  | 1  | 0 | 0   | 0  | 2  | 10  | 0  | 0 | 0 | 2  | 0 | 109   |
| s     | 4   | 0 | 0  | 0  | 9   | 0 | 0  | 5  | 15  | 0 | 0 | 5  | 2  | 0  | 1   | 22 | 0 | 0   | 0  | 1  | 3   | 0  | 0 | 0 | 10 | 0 | 81    |
| t     | 4   | 0 | 2  | 0  | 4   | 0 | 0  | 21 | 47  | 0 | 0 | 4  | 0  | 0  | 3   | 0  | 0 | 5   | 0  | 0  | 11  | 0  | 2 | 0 | 0  | 0 | 106   |
| u     | 22  | 0 | 5  | 1  | 1   | 0 | 2  | 0  | 2   | 0 | 0 | 2  | 1  | 0  | 20  | 2  | 0 | 11  | 11 | 2  | 0   | 0  | 0 | 0 | 0  | 0 | 71    |
| v     | 0   | 0 | 0  | 0  | 1   | 0 | 0  | 0  | 4   | 0 | 0 | 0  | 0  | 0  | 0   | 0  | 0 | 0   | 0  | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 4     |
| w     | 0   | 0 | 0  | 0  | 0   | 0 | 0  | 4  | 0   | 0 | 0 | 0  | 1  | 1  | 1   | 0  | 0 | 0   | 0  | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 15    |
| x     | 0   | 0 | 0  | 0  | 0   | 0 | 0  | 0  | 0   | 0 | 0 | 0  | 0  | 0  | 0   | 1  | 0 | 0   | 0  | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 1     |
| y     | 0   | 1 | 2  | 0  | 0   | 0 | 1  | 0  | 0   | 0 | 0 | 3  | 0  | 0  | 0   | 2  | 0 | 1   | 10 | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 23    |
| z     | 0   | 0 | 0  | 0  | 0   | 0 | 0  | 0  | 0   | 0 | 0 | 0  | 0  | 0  | 0   | 0  | 0 | 0   | 0  | 0  | 0   | 0  | 0 | 0 | 0  | 0 | 0     |
| Total | 131 | 9 | 55 | 23 | 187 | 1 | 18 | 32 | 277 | 0 | 3 | 79 | 18 | 70 | 104 | 37 | 1 | 106 | 54 | 67 | 159 | 15 | 9 | 1 | 35 | 3 | 1586  |

- (a) What would the noisy channel model return for the misspelled word mdcol given the following candidate words and their priors? Circle the values that you use in the matrix given above.

| w     | P(w)   | P(x w)            | P(x w).P(w)   |
|-------|--------|-------------------|---|
| mdcol | 0.0009 | $\frac{11}{1500}$ | $\frac{9}{10000} \times \frac{11}{1500} = \frac{99}{1500000}$   |
| model | 0.157  | 0                 | 0   |
| mdole | 0.077  | $\frac{20}{1500}$ | $\frac{77}{10000} \times \frac{20}{1500} = \frac{154}{1500000}$ |
| dmole | 0.0045 | 0                 | 0   |

arg max?  $\hat{w} = ?$  8

- (b) Consider the following text where the ordering of letters between words is reversed by transposition. Use the above matrix to find the probability of each spelling error.

s = fi yuo ena rzed tihs, yuo hvae a srtagne mnid too

$P(\text{trans}[i, f]) \rightarrow \frac{1}{1500}$

$P(\text{trans}[o, y]) \rightarrow 0$

$P(\text{trans}[e, a]) \rightarrow \frac{1}{1500}$

$P(\text{trans}[h, i]) \rightarrow 0$

$P(\text{trans}[a, v]) \rightarrow 0$

$P(\text{trans}[t, r]) \rightarrow \frac{5}{1500}$

$P(\text{trans}[g, o]) \rightarrow \frac{8}{1500}$



continued

This page has been intentionally left blank

Q3  $P(\text{trans}(i,n)) \rightarrow \frac{5}{1500}$

$P(\text{if}|\text{fi}) \rightarrow \frac{1}{1500} \checkmark$

$P(\text{you}|\text{you}) \rightarrow 0 \checkmark$

$P(\text{read}|\text{read}) \rightarrow \frac{1}{1500} \checkmark$

$P(\text{have}|\text{have}) \rightarrow 0 \checkmark$

$P(\text{this}|\text{this}) \rightarrow 0 \checkmark$

~~$P(\text{start}|\text{start})$~~

$P(\text{strange}|\text{strange}) = \frac{5}{140} + \frac{2}{1500} = \frac{40}{2250000}$

can?

25/2

7

Q1 continued: the Bigram's from high to low  
Sort all the start words heap<sup>0</sup> fill rand(2)  
Once we have start words  
Select a start word from start word list randomly and use  
Sent ← start word: word-start word  
heap<sup>1</sup> ← rand(3, 6) word  
if heap<sup>1</sup> Sent ← Sent + word + Bigram [randomly select a key from top]  
Printed line C ← sentence + "\n"  
if C ← sentence + "\n" integer is a multiple of 5  
Printed line C ← sentence + "\n"  
Return C



(7 Marks) [3+4]

#### Q4. Parts-of-Speech

Consider the following POS annotated story.

Little/JJ Red/NNP Riding/NNP Hood/NNP lived/VBD in/IN the/DT woods/NN with/IN her/PRP mother/NN .  
One/CD day/NN Little/NNP Red/NNP Riding/NNP Hood/NNP went/VBD to/TO visit/VB her/PRP granny/NN .  
She/PRP had/VBD a/DT nice/JJ cake/NN in/IN her/PRP basket/NN .  
On/IN her/PRP way/NN Little/JJ Red/NNP Riding/NNP Hood/NNP met/VBD a/DT wolf/NN .  
Hello/UH ! the/DT wolf/NN said/VBD . Where/WRB are/VBP you/PRP going/VBG ?  
I/PRP a/RS going/VBG to/TO see/VB my/PRP grandmother/NN . She/PRP lives/VBZ in/IN a/DT house/NN  
behind/IN those/DT trees/NN .  
The/DT wolf/NN ran/VBD to/TO Granny/NNP 's/NNP a/NN house/NN and/CC ate/NN Granny/NNP up/AB . He/PRP  
got/VBD into/IN Granny/NNP 's/NNP a/NN bed/NN . A/DT little/JJ later/AB , Little/NNP Red/NNP Riding/NNP  
Hood/NNP reached/VBD the/DT house/NN . She/PRP looked/VBD at/IN the/DT wolf/NN .  
Granny/NNP , what/VP big/JJ eyes/NN you/PRP have/VBP !  
All/PDT the/DT better/JJR to/TO see/VB you/PRP with/IN ! said/VBD the/DT wolf/NN .  
Granny/NNP , what/VP big/JJ ears/NN you/PRP have/VBP !  
All/PDT the/DT better/JJR to/TO hear/VB you/PRP with/IN ! said/VBD the/DT wolf/NN .  
Granny/NNP , what/VP a/DT big/JJ nose/NN you/PRP have/VBP !  
All/PDT the/DT better/JJR to/TO smell/VB you/PRP with/IN ! said/VBD the/DT wolf/NN .  
Granny/NNP , what/VP big/JJ teeth/NN you/PRP have/VBP !  
All/PDT the/DT better/JJR to/TO eat/VB you/PRP with/IN ! shouted/VBD the/DT wolf/NN .  
A/DT woodcutter/NN was/VBD in/IN the/DT wood/NN . He/PRP heard/VBD a/DT loud/JJ scream/NN and/CC  
ran/NN to/TO the/DT house/NN . The/DT woodcutter/NN hit/VBD the/DT wolf/NN over/IN the/DT head/NN  
 . The/DT wolf/NN opened/VBD his/PRP mouth/NN wide/JJ , gave/VBD a/DT cry/NN and/CC Granny/NNP  
jumped/VBD out/AB . The/DT wolf/NN ran/VBD away/AB and/CC Little/JJ Red/NNP Riding/NNP Hood/NNP  
never/AB saw/VBD the/DT wolf/NN again/AB .

- (a) The similarity() method in NLTK can be used for semi-supervised training. It uses the context of surrounding word annotations to find similar words i.e. words used in the same context. Highlight or underline the words in the above text that are used in a similar context as the word 'eyes'.

3

- (b) How do you plan to use POS tagging in your course project?

haven't decided on the course project yet  
but will use it to understand the  
content and infer implicit knowledge

3