

B.Tech (Computer Engineering) 8th Semester Examinations 2021
Natural Language Processing and Information Extraction

Paper Code: CEN 807

Maximum Marks: 60

Maximum Time: 3 hr

(Write your Roll No. on the top immediately on receipt of this question paper)

Note: Attempt any two parts from each question Assume suitable data, if necessary.

S.No.	Questions	Marks	CO																																														
1(a)	<p>Apply the Max-Match word segmentation algorithm on the following:</p> <p>“wecanonlyseeashortdistanceahead”</p> <p>Assume you have the English all vocabulary set.</p>	6	1																																														
1(b)	<p>Consider Count(w), Count(v,w) be unigram and bigram counts taken from a training corpus, where w is the single word, v,w is a bigram. Let N be te total number of words in the Corpus. What are the maximum likelihood estimates for unigram and bigram model?</p>	6	1																																														
1(c)	<p>Write the training and prediction algorithms for Bernoulli Naïve Bayes Text Classification. Consider the following dataset:</p> <table><tr><td>Feature</td><td>1</td><td>1</td><td>0</td><td>1</td><td>1</td><td>1</td><td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td><td>0</td><td>0</td><td>1</td></tr><tr><td>Label</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td><td>1</td><td>1</td><td>1</td><td>1</td><td>1</td><td>1</td><td>a</td><td>b</td></tr></table> <p>Predict the label values a and b using Bernoulli Naïve Bayes Text Classification</p>	Feature	1	1	0	1	1	1	0	0	0	1	0	1	1	1	0	1	1	1	1	0	0	1	Label	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	a	b	6	1
Feature	1	1	0	1	1	1	0	0	0	1	0	1	1	1	0	1	1	1	1	0	0	1																											
Label	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	a	b																											
2(a)	<table><tr><td></td><td>Doc</td><td>Words</td><td>Class</td></tr><tr><td rowspan="4">Training</td><td>1</td><td>Chinese, Beijing, Chinese</td><td>C</td></tr><tr><td>2</td><td>Chinese, Chinese, Shanghai</td><td>C</td></tr><tr><td>3</td><td>Chinese, Macao</td><td>C</td></tr><tr><td>4</td><td>Tokyo, Japan, Chinese</td><td>J</td></tr><tr><td>Test</td><td>5</td><td>Chinese, Chinese, Tokyo, Japan</td><td>?</td></tr></table> <p>Compute the most likely class for Doc 5. Assume a multinomial naive Bayes classifier and use add-α La Place smoothing for the likelihoods. Suitable value of α may be chosen.</p>		Doc	Words	Class	Training	1	Chinese, Beijing, Chinese	C	2	Chinese, Chinese, Shanghai	C	3	Chinese, Macao	C	4	Tokyo, Japan, Chinese	J	Test	5	Chinese, Chinese, Tokyo, Japan	?	6	2																									
	Doc	Words	Class																																														
Training	1	Chinese, Beijing, Chinese	C																																														
	2	Chinese, Chinese, Shanghai	C																																														
	3	Chinese, Macao	C																																														
	4	Tokyo, Japan, Chinese	J																																														
Test	5	Chinese, Chinese, Tokyo, Japan	?																																														

2(b)	<p>Given a corpus $w_1, w_2, w_3, \dots, w_N$ such that each word is labelled with three POS tags: NOUN, VERB, ADJECTIVE, OTHER.</p> <p>Construct the graphical sequence labelling model using HMM for POS tagging and define all the transition and emission probabilities.</p>	6	2
2(c)	Define Precision, Recall and F-measure.	6	2
3(a)	<p>Given the following Dictionary entry for line.</p> <p>line² a length of cord, rope, wire, or other material serving a particular purpose: <i>wring the clothes and hang them on the line a telephone line</i>.</p> <ul style="list-style-type: none"> • one of a vessel's mooring ropes. • a telephone connection: <i>she had a crank on the line</i>. • a railroad track. • a branch or route of a railroad system: <i>the Philadelphia to Baltimore line</i>. <p>line³ a horizontal row of written or printed words.</p> <ul style="list-style-type: none"> • a part of a poem forming one such row: <i>each stanza has eight lines</i>. • (lines) the words of an actor's part in a play or film. • a particularly noteworthy written or spoken sentence: <i>his speech ended with a line about the failure of justice</i>. <p>Which of these senses are related by homonymy, and which are related by polysemy? For any senses which are polysemous, give an argument as to how the senses are related.</p>	6	3
3(b)	<p>Assume the following sentence L in which the word line is in focus:</p> <p>L = you must wait in a long line at the checkout counter</p> <p>Give a collocation feature vector (including n-gram) for in the word line in L, given a window size of 3 words to the left and 3 words to the right.</p>	6	3
3(c)	<p>C = About three years ago, he nearly gave up because he nearly had nothing to sell; Now his shelves are full, and towels and clothes hang from a <u>line</u> overhead.</p> <p>For the word line in the above text L, generate the bag-of-words feature vector for window size = +-2, assume C as the whole corpus.</p>	6	3

4(a)	For the following term document matrix:	6	4																									
	<table><tr><td></td><td>Document 1</td><td>Document 2</td><td>Document 3</td><td>Document 3</td></tr><tr><td>digital</td><td>1</td><td>0</td><td>7</td><td>13</td></tr><tr><td>computer</td><td>114</td><td>80</td><td>62</td><td>89</td></tr><tr><td>information</td><td>36</td><td>58</td><td>1</td><td>4</td></tr><tr><td>data</td><td>20</td><td>15</td><td>2</td><td>3</td></tr></table>				Document 1	Document 2	Document 3	Document 3	digital	1	0	7	13	computer	114	80	62	89	information	36	58	1	4	data	20	15	2	3
				Document 1	Document 2	Document 3	Document 3																					
	digital			1	0	7	13																					
	computer			114	80	62	89																					
	information			36	58	1	4																					
data	20	15	2	3																								
Calculate the similarity between good and fool using (i) cosine similarity (ii) PPMI. Use add 2 smoothing if necessary.																												
4(b)	Describe the Logistic Regression based skip gram model with the help of suitable diagram.	6	4																									
4(c)	Construct the word co-occurrence matrix, window size+-1, for the following text: Document: “Roses are red. Sky is blue.	6	4																									
5(a)	Define Information Extraction, Named Entity Recognition and Relation Extraction. Why is Information Retrieval task not sufficient to perform information extraction tasks?	6	5																									
5(b)	What are the different encoding schemes for Named Entity Recognition? Illustrate with examples.	6	5																									
5(c)	Construct a rule based e-mail extractor which can distinguish between sender and receiver e-mail addresses.	6	5																									
