National University of Computer and Emerging Sciences, Lahore Campus

STICHAL UNIVERS
* * * * *
Woo as
SALIN S. EMERGINES

Course Name:	NLP	Course Code:	CS4063
Degree Program:	BS-CS	Semester:	Spring 2023
Exam Duration:	60 Minutes	Total Marks:	82
Paper Date:	28-02-2023	Weight	
Sections:	ALL	No of Page(s):	
Exam Type:	Midterm I		

Student : Name:_	Roll No	Section:
Instruction/Notes:	Attempt all questions. Programmable calculators are not allowed.	
Q1. You are give	n the following training corpus.	(1+2+2+5+2) 12
<s> i want to eat</s>	t thai food	

- <s> we ate pakistani food </s>
- <s> i ate apples </s>
- <s> they ate thai food </s>
- a) Calculate the probability of the following test sentence. Include </s> in your counts just like any other token.

<s> i ate thai food </s>

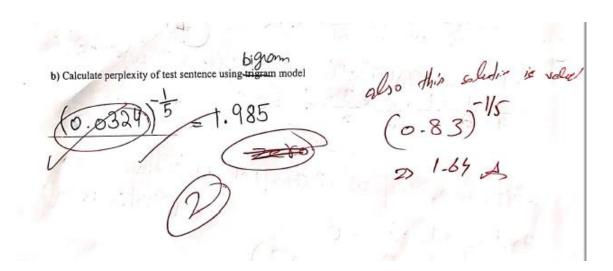
- i) **Unigram Model**
- ii) Bigram Model
- iii) Trigram Model
- Bigram model with linear interpolation Λ_1 = 0.7, Λ_2 = 0.3 iv)
- b) Calculate the perplexity of the test sentence using bigram model.

i. Unigram Model
$$= P(I) \times P(ate) \times P(chineste) \times P(food) \times P(
$$= \frac{2}{3} \times \frac{3}{3} \times \frac{2}{3} \times \frac{3}{3} \times \frac{4}{21} = 3.52$$
ii. Bigram Model
$$= P(I | < s >) \times P(ate | I) \times P(chineste) \times P(food | chinese) \times P(s | s) | food$$

$$= \frac{2}{4} \times \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3} \times \frac{3}{3} = \frac{1}{12} = 0.83$$

$$= \frac{2}{4} \times \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3} \times \frac{3}{3} = \frac{1}{12} = 0.83$$
iii. Trigram Model
$$= P(I | < s > < s >) \times P(ate | < s > I) \times P(chinese | I ate) + P(food | ate chinese) \times P(I | < s > < s >) \times P(ate | < s > I) \times P(chinese | I ate) + P(food | ate chinese) \times P(I | < s > < s >) \times P(ate | < s > I) \times P(a$$$$

iv. Bigram language model with linear interpolation.
$$0.\overline{7}$$
 0.3 0.3 0.7 $0.$



Q3: You are given two documents:

Doc 1 = the car is in the parking and the bike is in the garage (13)

Doc 2 = the truck is driven on the highway and the tractor is in the farm parking (15)

a) Compute the normalized term frequency and un-smoothed logarithmic inverse document frequency for the given corpus. You can use log of base 10 for the calculation of IDF. Fill the table based on your calculations: (15+10)

IDF: log(N/n)

Term	Count (doc1)	Count (doc2)	TF (doc1)	TF (doc2)	IDF	TF*IDF (doc1)	TF*IDF (doc2)

Term	Count	Count	Tf (doc 1)	Tf (doc 2)	IDF	TF*IDF (doc 1)	TF*IDF (doc 2)
	(doc 1)	(doc 2)					
the	4	4	0.30769231	0.26666667	0	0	0
car	1	0	0.07692308	0	0.301	0.02315385	0
is	2	2	0.15384615	0.13333333	0	0	0
in	2	1	0.15384615	0.06666667	0	0	0
parking	1	1	0.07692308	0.06666667	0	0	0
and	1	1	0.07692308	0.06666667	0	0	0
bike	1	0	0.07692308	0	0.301	0.02315385	0
garage	1	0	0.07692308	0	0.301	0.02315385	0
truck	0	1	0	0.06666667	0.301	0	0.02006667
driven	0	1	0	0.06666667	0.301	0	0.02006667
on	0	1	0	0.06666667	0.301	0	0.02006667
highway	0	1	0	0.06666667	0.301	0	0.02006667
tractor	0	1	0	0.06666667	0.301	0	0.02006667
farm	0	1	0	0.06666667	0.301	0	0.02006667

b) From TF*IDF vectors calculated in part a, compute Euclidean Distance and Cosine Similarity (write the formulae too).

	Euclidean Distance:	
	0.0634	
	Cosine Similarity:	
	cosme similarity.	
	0	
Q4: Aı	nswer the following.	(10)
I.	Which of the following is a type of Minkowski distance?	
II.	(a. Hamming , b. Levenshtein, c. Jaro) Damerau-Levenshtein allows4 edit operations, while	Hamming allows 1 operations, (0, 1, 2, 3, 4, 5)
III.	Root of the word "antidisestablishmentarianism" is establ	 :
IV.	What are the derivational morphemes in each of these wo	
17	"realism", and "higher"? (en-, re-, none, un-, -ism	•
V. blaa! b	Mention all of the following expressions that the regex /bla+. lat! black! bla?! bla.?! bla+?!	.?!/ Will Select?
a.	bla! b. blaa! c. blat! d. black! e. bla?! f. bla+?! g. bla?! h. bla	ı.?! i. bla +?! j. bla+.?!

Q5. Find the Levenshtein distance between PAYMENTS and APARTMENTS. Use the same algorithm and weights as discussed in the class (i.e. cost(Insertion)=1, cost(Deletion)=1, cost(Substitution)=2). (20)

	#	A	P	A	R	T	M	E	N	T	S
#	0	1	2	თ	4	5	6	7	8	9	10
P	1	2	1	2	3	4	5	6	7	8	9
A	2	1	2	1	2	3	4	5	6	7	8
Y	3	2	3	2	3	4	5	6	7	8	9
M	4	3	4	3	4	5	4	5	6	7	8
E	5	4	5	4	5	6	5	4	5	6	7
N	6	5	6	5	6	7	6	5	4	5	6
T	7	6	7	6	7	6	7	6	5	4	5
S	8	7	8	7	8	7	8	7	6	5	4

Minimum Edit Distance	Λ

• Show the optimal alignment between the sequences and one possible minimal edit sequence (a sequence of inserts *I*, deletes *D* and substitutions *S*) that would result in an optimal conversion from PAYMENTS to APARTMENTS.

	P	A		Y	M	Е	N	Т	S		
A	P								S		
I	S	S	I	S	S	S	S	S	S		