

<b>Course Code: AI4001</b>	<b>Course Name: Fundamentals of Natural Language Processing</b>
<b>Instructor Name: Sumaiyah Zahid</b>	
<b>Student Roll No:</b>	<b>Section No:</b>

Instructions:

- Return the question paper.
- Read each question completely before answering it. There are **4 questions and 2 pages**.
- In case of any ambiguity, you may make assumptions. But your assumption should not contradict any statement in the question paper.
- Show all steps clearly.

**Time:** 60 minutes.

**Max Marks:** 30 points

**Question 1 [5 Points]:**

**[CLO 1]**

1. Compute the minimum edit distance (using insertion cost 1, deletion cost 1, substitution cost 2) of “fastian” to “fusion”. Show your work (using the edit distance grid). **[3 Points]**

**Solution:**

	f u s i o n						
0	1	2	3	4	5	6	
f	1	0	1	2	3	4	5
a	2	1	2	3	4	5	6
s	3	2	3	2	3	4	5
t	4	3	4	3	4	5	6
i	5	4	5	4	3	4	5
a	6	5	6	5	4	5	6
n	7	6	7	6	5	6	5

2. Write a Regex expression: **[2 Points]**

- a. To identify phone numbers in a simple format, such as "+92-123-1231234" or "0123-1231234."

**Solution:** `(\+\\d{2}-\\d{3}-\\d{7}|\\d{4}-\\d{7})`

- b. To validate usernames, allowing only alphanumeric characters and underscores.

**Solution:** `^\\b \\w+ \\b /`

**Question 2 [10 Points]:**

**[CLO 1]**

Consider the following training data:

<s> to be who or not to be who just </s>

<s> be who you want to be</s>

1. Compute probability of the test sentence “**be who you want**” using linear interpolation starting from trigram LM for the above training dataset with  $\lambda_1=0.5$ ,  $\lambda_2=0.3$  and  $\lambda_3=0.2$ . **[3 Points]**

**Solution:**

Tokenize the test sentence into trigrams:

	1

Trigrams: {<s> <s> be, <s> be who, be who you, who you want, you want <s>}

$$P(<s> <s> be) = \text{Count}(<s> <s> be) / \text{Count}(<s> <s>) = 1 / 2 = 0.5$$

$$P(<s> be who) = \text{Count}(<s> be who) / \text{Count}(<s> be) = 1 / 1 = 1$$

$$P(\text{be who you}) = \text{Count}(\text{be who you}) / \text{Count}(\text{be who}) = 1 / 3 = 0.3$$

$$P(\text{who you want}) = \text{Count}(\text{who you want}) / \text{Count}(\text{who you}) = 1 / 1 = 1$$

$$P(\text{you want } <s>) = \text{Count}(\text{you want } <s>) / \text{Count}(\text{you want}) = 0 / 1 = 0$$

Bigrams: {<s> be, be who, who you, you want, want <s>}

$$P(<s> be) = \text{Count}(<s> be) / \text{Count}(<s>) = 1 / 2 = 0.5$$

$$P(\text{be who}) = \text{Count}(\text{be who}) / \text{Count}(\text{be}) = 3 / 4 = 0.75$$

$$P(\text{who you}) = \text{Count}(\text{who you}) / \text{Count}(\text{who}) = 1 / 3 = 0.3$$

$$P(\text{you want}) = \text{Count}(\text{you want}) / \text{Count}(\text{you}) = 1 / 1 = 1$$

$$P(\text{want } <s>) = \text{Count}(\text{want } <s>) / \text{Count}(\text{want}) = 0 / 1 = 0$$

Unigrams: {<s>, be, who, you, want, <s>}

$$P(<s>) = 2 / 19$$

$$P(\text{be}) = 4 / 19$$

$$P(\text{who}) = 3 / 19$$

$$P(\text{you}) = 1 / 19$$

$$P(\text{want}) = 1 / 19$$

$$P(</s>) = 2 / 19$$

$P(\text{"be who you want"}) = \lambda_1 * \text{trigram} + \lambda_2 * \text{bigram} + \lambda_3 * \text{Unigram}$

$P(\text{"be who you want"}) = (0.5 * 0) + (0.3 * 0) + (0.2 * 48/19^6) = 2.04 \times 10^{-7}$

2. Compute the probability of the test sentence **“to be not who”** using stupid backoff starting from quadgram LM. [3 Points]

Quadgrams: {<s><s><s> to, <s><s> to be, <s> to be not, to be not who, be not who <s>} so many 0 probability then move to

Trigram: {<s> <s> to, <s> to be, to be not, be not who, not who <s>}

Bigram: {<s> to, to be, be not, not who, who <s>}

Unigram: {<s>, to, be, not, who, <s>}

$P(<s>) = 2 / 19$   
 $P(\text{to}) = 3 / 19$   
 $P(\text{be}) = 4 / 19$   
 $P(\text{not}) = 1 / 19$   
 $P(\text{who}) = 3 / 19$   
 $P(</s>) = 2 / 19$

$P(<s> \text{ to be not who } </s>) = 0.4 * 0.4 * 0.4 * 2/19 * 3/19 * 4/19 * 1/19 * 3/19 * 2/19$

3. Compute the perplexity of the test sentence **“to be who”** for bigram and trigram. [2 Points]

For Bigram:

$\text{Perplexity} = (1 / P(\text{to be who}))^{1/3}$

$P(\text{to be who}) = 1^{3/4} // \text{ skipping } <s> </s>$

$\text{Perplexity} = (4/3)^{1/3}$

For Trigram:

$P(\text{to be who}) = 2/3$

$\text{Perplexity} = (3/2)^{1/3}$

4. Generate 5 more words for the sentence starting from **“to...”** using bigram LM. [2 Points]

to be who you want

	2

**Question 3 [10 Points]:****[CLO 3]**

1. Write the Context Free Grammar rules for the following instructions of Assembly Language.

**[5 Points]**

- i. LOAD R0, 100;
- ii. STORE R1, 200;
- iii. ADD R0, R2, R3;
- iv. SUB R2, R0, R1;

$S \rightarrow \text{INS } \langle \text{REG} \rangle, \langle \text{CONS} \rangle; \mid \text{INS } \langle \text{REG} \rangle, \langle \text{REG} \rangle \langle \text{REG} \rangle;$

$\text{INS} \rightarrow \text{LOAD} \mid \text{STORE} \mid \text{ADD} \mid \text{SUB}$

$\text{REG} \rightarrow \text{R0} \mid \text{R1} \mid \text{R2} \mid \text{R3} \mid \dots \mid \text{Rn};$

$\text{CONS} \rightarrow \backslash d^+$

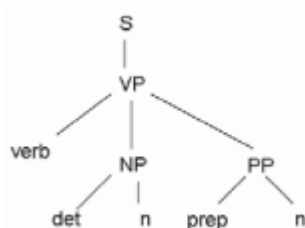
2. Draw the top-ranked parse tree for the sentence below by applying the given PCFG. Does the result seem reasonable to you? Why or why not?

**[5 Points]**

**Cut the envelope with scissors**

Consider the following PCFG:

production rule	probability
$S \rightarrow VP$	1.0
$VP \rightarrow \text{Verb NP}$	0.7
$VP \rightarrow \text{Verb NP PP}$	0.3
$NP \rightarrow NP PP$	0.3
$NP \rightarrow \text{Det Noun}$	0.7
$PP \rightarrow \text{Prep Noun}$	1.0
$\text{Det} \rightarrow \text{the}$	0.1
$\text{Verb} \rightarrow \text{Cut} \mid \text{Ask} \mid \text{Find} \mid \dots$	0.1
$\text{Prep} \rightarrow \text{with} \mid \text{in} \mid \dots$	0.1
$\text{Noun} \rightarrow \text{envelope} \mid \text{grandma} \mid \text{scissors} \mid \text{men} \mid \text{suits} \mid \text{summer} \mid \dots$	0.1



The top-ranked sentence structure is shown. (The leaf nodes representing words are omitted.) The probability of the resulting parse tree is  $1.0 * 0.3 * 0.7 * 1.0 * (0.1)^5$ , which is larger than  $1.0 * 0.7 * 0.3 * 0.7 * 1.0 * (0.1)^5$ , the probability of the alternative parse tree (with the  $[VP \rightarrow \text{Verb NP}]$  rule expansion).

**Question 4 [5 Points]:****[CLO 1]**

Consider the following short article reviews each labeled with a type, either political or scientific. Use a naive Bayes classifier to classify the test data.

vote, election, policy, policy  $\rightarrow$  Political  
 experiments, artificial, intelligence  $\rightarrow$  Scientific  
 debate, vote, budget, election, policy  $\rightarrow$  Political

	3

vote, budget, policy, vote, gun —> Political  
intelligence, research, artificial —> Scientific

D: **budget, policy, intelligence, vote.**

1. Compute the most likely class for D using Naive Bayes and add-1 smoothing. [2.5 points]

$$P(\text{Political}) = 3/5$$

$$P(\text{Scientific}) = 2/5$$

$$P(\text{budget} | \text{Political}) = (2 + 1) / (14 + 10) = 3/24 = 1/8$$

$$P(\text{policy} | \text{Political}) = (4 + 1) / (14 + 10) = 5/24$$

$$P(\text{intelligence} | \text{Political}) = (0 + 1) / (14 + 10) = 1/24$$

$$P(\text{vote} | \text{Political}) = (4 + 1) / (14 + 10) = 5/24$$

$$P(\text{budget} | \text{Scientific}) = (0 + 1) / (6 + 10) = 1/16$$

$$P(\text{policy} | \text{Scientific}) = (0 + 1) / (6 + 10) = 1/16$$

$$P(\text{intelligence} | \text{Scientific}) = (2 + 1) / (6 + 10) = 3/16$$

$$P(\text{vote} | \text{Scientific}) = (0 + 1) / (6 + 10) = 1/16$$

$$\begin{aligned} P(D | \text{Political}) &= P(\text{Political}) * P(\text{budget} | \text{Political}) * P(\text{policy} | \text{Political}) * P(\text{intelligence} | \text{Political}) \\ &* P(\text{vote} | \text{Political}) * P(\text{research} | \text{Political}) \\ &= (3/5) * (1/8) * (5/24) * (1/24) * (5/24) \end{aligned}$$

$$\begin{aligned} P(D | \text{Scientific}) &= P(\text{Scientific}) * P(\text{budget} | \text{Scientific}) * P(\text{policy} | \text{Scientific}) * P(\text{intelligence} | \text{Scientific}) \\ &* P(\text{vote} | \text{Scientific}) * P(\text{research} | \text{Scientific}) \\ &= (2/5) * (1/16) * (1/16) * (3/16) * (1/16) \end{aligned}$$

Document D is Political.

2. What is the class of the test sentence if you use Binary Naive Bayes? [2.5 points]

$$P(\text{Political}) = 3/5$$

$$P(\text{Scientific}) = 2/5$$

$$P(\text{budget} | \text{Political}) = 2 / 12$$

$$P(\text{policy} | \text{Political}) = 3/12$$

$$P(\text{intelligence} | \text{Political}) = 0$$

$$P(\text{vote} | \text{Political}) = 3/12$$

$$P(\text{budget} | \text{Scientific}) = 0$$

$$P(\text{policy} | \text{Scientific}) = 0$$

$$P(\text{intelligence} | \text{Scientific}) = 2/6$$

$$P(\text{vote} | \text{Scientific}) = 0$$

$$P(D | \text{Political}) = (3/5) * 2/12 * 3/12 * 0 * 3/12$$

$$P(D | \text{Scientific}) = (2/5) * 0 * 0 * 2/6 * 0$$

\*\*\*Best of luck while transforming Text into Machine Intelligence.\*\*\*

	4