

## Natural Language Processing

### Assignment- 2

TYPE OF QUESTION: MCQ

Number of questions: 10

Total mark: 10 X 1 = 10

---

#### **QUESTION 1:**

According to Zipf's law which statement(s) is/are correct?

- (i) A small number of words occur with high frequency.
- (ii) A large number of words occur with low frequency.
- a. Both (i) and (ii) are correct
- b. Only (ii) is correct
- c. Only (i) is correct
- d. Neither (i) nor (ii) is correct

**Correct Answer: a**

Solution:

---

#### **QUESTION 2:**

Consider the following corpus  $C_1$  of 4 sentences. What is the total count of unique bi-grams for which the likelihood will be estimated? Assume we do not perform any pre-processing.

**today is Nayan's birthday**  
**she loves ice cream**  
**she is also fond of cream cake**  
**we will celebrate her birthday with ice cream cake**

- a. 24
- b. 28
- c. 27
- d. 23

**Correct Answer: a**

**Detailed Solution:**

Unique bi-grams are:

<s> today	today is	is Nayan's	Nayan's birthday	birthday <\s>
<s> she	she loves	loves ice	ice cream	cream <\s>
She is	is also	also fond	fond of	of cream
cake <\s>	<s> we	we will	will celebrate	celebrate her
her birthday	birthday with	with ice	cream cake	

---

### **QUESTION 3:**

A 4-gram model is a \_\_\_\_\_ order Markov Model.

- a. Constant
- b. Five
- c. Four
- d. Three

**Correct Answer: d**

**Detailed Solution:**

---

### **QUESTION 4:**

Which one of these is a valid Markov assumption?

- a. The probability of a word depends only on the current word.
- b. The probability of a word depends only on the previous word.
- c. The probability of a word depends only on the next word.
- d. The probability of a word depends only on the current and the previous word.

**Correct Answer: b**

Solution:

---

**QUESTION 5:**

For the string '**mash**', identify which of the following set of strings have a Levenshtein distance of 1.

- a. smash, mas, lash, mushy, hash
- b. bash, stash, lush, flash, dash
- c. smash, mas, lash, mush, ash
- d. None of the above

**Correct Answer: c**

**Detailed Solution:**

---

**QUESTION 6:**

Assume that we modify the costs incurred for operations in calculating Levenshtein distance, such that both the insertion and deletion operations incur a cost of 1 each, while substitution incurs a cost of 2. Now, for the string '**lash**' which of the following set of strings will have an edit distance of 1?

- a. ash, slash, clash, flush
- b. flash, stash, lush, blush,
- c. slash, last, bash, ash
- d. None of the above

**Correct Answer: d**

**Detailed Solution:**

---

**QUESTION 7:**

Given a corpus  $C_2$ , the Maximum Likelihood Estimation (MLE) for the bigram "dried berries" is 0.3 and the count of occurrence of the word "dried" is 580. for the same corpus  $C_2$ , the likelihood of "dried berries" after applying add-one smoothing is 0.04. What is the vocabulary size of  $C_2$ ?

- a. 3585
- b. 3795
- c. 4955
- d. 3995

**Correct Answer: b**

**Detailed Solution:**

$$P_{MLE}(\text{berries} | \text{dried}) = \frac{C(\text{dried}, \text{berries})}{C(\text{dried})}$$

$$0.3 = \frac{C(\text{dried}, \text{berries})}{580}$$

$$C(\text{dried}, \text{berries}) = 174$$

$$P_{Add-1}(\text{berries} | \text{dried}) = \frac{C(\text{dried}, \text{berries}) + 1}{C(\text{dried}) + V}$$

$$0.04 = \frac{174 + 1}{580 + V}$$

$$V = 3795$$

---

**For Question 8 to 10, consider the following corpus  $C_3$  of 3 sentences.**

**there is a big garden**

**children play in a garden**

**they play inside beautiful garden**

**QUESTION 8:**

Calculate **P(they play in a big garden)** assuming a bi-gram language model.

- a.  $1/8$
- b.  $1/12$
- c.  $1/24$
- d. None of the above

**Correct Answer: b**

**Detailed Solution:**

$$P(\text{they} | \langle s \rangle) = 1/3$$

$$P(\text{play} | \text{they}) = 1/1$$

$$P(\text{in} | \text{play}) = 1/2$$

$$P(\text{a} | \text{in}) = 1/1$$

$$P(\text{big} | \text{a}) = 1/2$$

$$P(\text{garden} | \text{big}) = 1/1$$

$$P(\langle s \rangle | \text{garden}) = 3/3$$

$$P(\text{they play in a big garden}) = 1/3 \times 1/1 \times 1/2 \times 1/1 \times 1/2 \times 1/1 \times 3/3 = 1/12$$

---

**QUESTION 9:**

Considering the same model as in Question 7, calculate the perplexity of **<s> they play in a big garden <\s>**.

- a. 2.289
- b. 1.426
- c. 1.574
- d. 2.178

**Correct Answer: b**

**Detailed Solution:**

$$\text{perplexity} = \sqrt[3]{12} = 1.426$$


---

**QUESTION 10:**

Assume that you are using a bi-gram language model with add one smoothing. Calculate **P(they play in a beautiful garden)**.

- a.  $4.472 \times 10^{-6}$
- b.  $2.236 \times 10^{-6}$
- c.  $3.135 \times 10^{-6}$
- d. None of the above

**Correct Answer: b**

**Detailed Solution:**

$$|V|=11$$

$$P(\text{they} \mid \text{<s>}) = (1+1)/(3+11)$$

$$P(\text{play} \mid \text{they}) = (1+1)/(1+11)$$

$$P(\text{in} \mid \text{play}) = (1+1)/(2+11)$$

$$P(\text{a} \mid \text{in}) = (1+1)/(1+11)$$

$$P(\text{beautiful} \mid \text{a}) = (0+1)/(2+11)$$

$$P(\text{garden} \mid \text{beautiful}) = (1+1)/(1+11)$$

$$P(\text{<\s>} \mid \text{garden}) = (3+1)/(3+11)$$

$$\begin{aligned} P(\text{they play in a beautiful garden}) &= 2/14 \times 2/12 \times 2/13 \times 2/12 \times 1/13 \times 2/12 \times 4/14 \\ &= 2.236 \times 10^{-6} \end{aligned}$$


---

\*\*\*\*\*END\*\*\*\*\*

