

# Natural Language Processing

## Assignment 4

### Type of Question: MCQ

Number of Questions: 8 Total Marks:  $(6 \times 1) + (2 \times 2) = 10$

=====

1. Baum-Welch algorithm is an example of - **[Marks 1]**
- a. Forward-backward algorithm
  - b. Special case of the Expectation-maximization algorithm
  - c. Both A and B
  - c. None

**Answer: C**

**Solution:** Theory.

=====

2. Once a day (e.g. at noon), the weather is observed as one of state 1: rainy state 2: cloudy state 3: sunny The state transition probabilities are :

0.4	0.3	0.3
0.2	0.6	0.2
0.1	0.1	0.8

Given that the weather on day 1 ( $t = 1$ ) is sunny (state 3), what is the probability that the weather for the next 7 days will be "sun-sun-rain-rain-sun-cloudy-sun"?

**[Marks 2]**

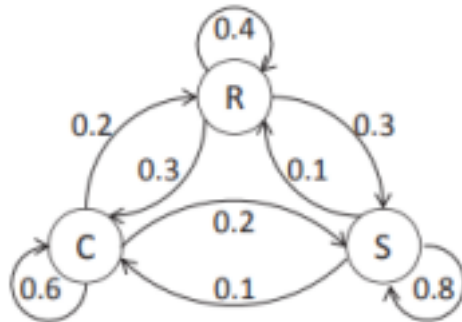
- a.  $1.54 \times 10^{-4}$
- b.  $8.9 \times 10^{-2}$
- c.  $7.1 \times 10^{-7}$
- d.  $2.5 \times 10^{-10}$

**Answer: A**

**Solution:**

$O = \{S3, S3, S3, S1, S1, S3, S2, S3\}$

$$\begin{aligned}
P(O \mid \text{Model}) &= P(S3, S3, S3, S1, S1, S3, S2, S3 \mid \text{Model}) \\
&= P(S3) P(S3|S3) P(S3|S3) P(S1|S3) P(S1|S1) P(S3|S1) P(S2|S3) \\
P(S3|S2) &= Q3 \cdot a_{33} \cdot a_{33} \cdot a_{31} \cdot a_{11} \cdot a_{13} \cdot a_{32} \cdot a_{23} \\
&= (1)(0.8)(0.8)(0.1)(0.4)(0.3)(0.1)(0.2) \\
&= 1.536 \times 10^{-4}
\end{aligned}$$



=====

3. Find the Viterbi Decoding of the sequence “ki fin yeni!”. Possible POS tags are {T1, T2, T3, T4}. Assume all POS tags are equally likely to be at the starting of a sequence. Emmission probabilities are, **[Marks 1]**

	ki	fin	yeni
T1	0.1	0.1	0.8
T2	0.8	0.1	0.1
T3	0.2	0.2	0.6
T4	0.8	0.1	0.1

Table 1: Output Symbol probabilities

Transition matrix is -

	T1	T2	T3	T4
T1	0.18	0.01	0.8	0.01
T2	0.9	0.0	0.05	0.05
T3	0.4	0.5	0.05	0.05
T4	0.4	0.5	0.05	0.05

Table 2: Hidden State transition matrix

Calculate  $P(x_1 = \text{“ki”}, x_2 = \text{“f in”}, y_1 = \text{“T1”}, y_2 = \text{“T2”})$ .

- 0.000025
- 0.0001
- 0.0025
- None of the above

**Answer: A**

**Solution:** Apply Markov Property and chain rule.

=====

4. Let us define an HMM Model with  $K$  classes for hidden states and  $T$  data points as observations. The dataset is defined as  $X = \{x_1, x_2, \dots, x_T\}$  and the corresponding hidden states are  $Z = \{z_1, z_2, \dots, z_T\}$ . Please note that each  $x_i$  is an observed variable and each  $z_i$  can belong to one of classes for hidden state. What will be the size of the state transition matrix, and the emission matrix, respectively for this example. **[Marks 1]**

- a.  $K \times K, K \times T$
- b.  $K \times T, K \times T$
- c.  $K \times K, K \times K$
- d.  $K \times T, K \times K$

**Answer: A**

**Solution:** Since there are  $K$  hidden states, the state transition matrix will be of size  $K \times K$ . The emission matrix will be of size  $K \times T$ , as it defines the probability of emitting an observed state from a hidden state.

=====

5. You are building a model distribution for an infinite stream of word tokens. You know that the source of this stream has a vocabulary of size 1000. Out of these 1000 words you know of 100 words to be stop words each of which has a probability of 0.0019. With only this knowledge what is the maximum possible entropy of the modelled distribution. (Use log base 10 for entropy calculation) **[Marks 2]**

- a. 5.079
- b. 0
- c. 2.984
- d. 12.871

**Answer: C**

**Solution:** There are 100 stopwords with each having an occurrence probability of 0.0019. Hence,

$$P(\text{Stopwords}) = 100 * 0.0019 = 0.19$$

$$P(\text{non - stopwords}) = 1 - 0.19 = 0.81$$

For maximum entropy, the remaining probability should be uniformly distributed.

For every non-stopword  $w$ ,  $P(w) = 0.81/(1000 - 100) = 0.81/900 = 0.0009$ .

Finally, the value of the entropy would be,

$$H = E(\log(1/p))$$

$$= -100(0.0019 * \log(0.0019)) - 900(0.0009 \log(0.0009))$$

$$= -2.9841$$

=====

6. For an HMM model with  $N$  hidden states,  $V$  observable states, what are the dimensions of parameter matrices  $A, B$  and  $\pi$ ?  $A$ : Transition matrix,  $B$ : Emission matrix,  $\pi$ : Initial Probability matrix. **[Marks 1]**

- a.  $N \times V, N \times V, N \times N$
- b.  $N \times N, N \times V, N \times 1$
- c.  $N \times N, V \times V, N \times 1$
- d.  $N \times V, V \times V, V \times 1$

**Answer: B**

**Solution:** Matrix  $A$  contains all the transition probabilities and have dimension  $N \times N$ . Similarly, matrix  $B$  contains all the emission probabilities and dimension  $N \times V$ . Similarly,  $\pi$  contains initial probability for all hidden states and have dimension  $N \times 1$ .

=====

7. Suppose you have the input sentence "Death Note is a great anime". And you know the possible tags each of the words in the sentence can take. • Death: NN, NNS, NNP, NNPS

• Note: VB, VBD, VBZ

• is: VB

• a: DT

• great: ADJ

• anime: NN, NNS, NNP

How many possible hidden state sequences are possible for the above sentence and States? **[Marks 1]**

- a.  $4 \times 3 \times 3$
- b.  $4^{3 \times 3}$
- c.  $2^4 \times 2^3 \times 2^3$
- d.  $2^{4 \times 3 \times 3}$

**Answer: A**

**Solution:** Each possible hidden sequence can take only one POS tag for each of the words. Hence the total possibility will be a product of the number of candidates for each word.

=====

8. In Hidden Markov Models or HMMs, the joint likelihood of an observed sequence  $O$  with a hidden state sequence  $Q$ , is written as  $P(O, Q; \theta)$ . In many applications, like POS tagging, one is interested in finding the hidden state sequence  $Q$ , for a given observation sequence, that maximizes  $P(O, Q; \theta)$ . What is the time required to compute the most likely  $Q$  using an exhaustive search? The required notations are,  $N$ : possible number of hidden states,  $T$ : length of the observed sequence. **[Marks 1]**

- a. Of the order of  $TN^T$
- b. Of the order of  $N^2T$
- c. Of the order of  $T^N$
- d. Of the order of  $N^2$

**Answer:** A

**Solution:** We will need to compute  $P(O, Q|\theta)$  for all possible  $Q$ . There are a total of  $N^T$  possible hidden sequences  $Q$  for a sequence of length  $T$ . Each individual probability calculation also requires  $T$  multiplications.