

CS4063
NLP(CS/SE)

Serial No:

Final Exam

Total Time: 3 Hours

Total Marks: 100

Friday, January 5, 2024

Course Instructor

Dr. Mirza O Beg, Mr. Saad Salman


 Signature of Invigilator

Student Name

Roll No

Section

Signature

DO NOT OPEN THE QUESTION BOOK OR START UNTIL INSTRUCTED.**Instructions:**

1. Verify at the start of the exam that you have a total of nine (9) questions printed on ten (10) pages including this title page.
2. Attempt all questions on the question-book and in the given order.
3. This exam is **closed book**. Mobiles, Internet and note-sharing is not allowed. Please see that the area in your threshold is free of any material classified as *useful in the paper*, i.e. mobile/internet or else there may be a charge of cheating.
4. Read the questions carefully for clarity of context and understanding of meaning and make assumptions wherever required, for neither the invigilator will address your queries, nor the teacher/examiner will come to the examination hall for any assistance.
5. Fit in all your answers in the provided space. You may use extra space on the last page if required. If you do so, clearly mark question/part number on that page to avoid confusion.
6. Use only your own stationery and calculator. If you do not have your own calculator, do manual calculations.
7. Use only permanent ink-pens. Only the questions attempted with permanent ink-pens will be considered. Any part of paper done in lead pencil cannot be claimed for checking/rechecking.

| | Q1 | Q2 | Q3 | Q4 | Q5 | Q6 | Q7 | Q8 | Q9 | Total |
|----------------|----|----|----|----|----|----|----|----|----|-------|
| Marks Obtained | 10 | 6 | 9 | 25 | 10 | 10 | 12 | 8 | 6 | 79.5 |
| Total Marks | 10 | 10 | 10 | 10 | 10 | 10 | 15 | 15 | 10 | 100 |

Q1. $tf - idf$

(10 Marks) [4+2+4]

Term frequency - Inverse document frequency ($tf - idf$), is a numerical statistic that is intended to reflect how important a word is to a document in a corpus. Assuming that $tf(t, d) = \log(1 + f_{t,d})$ where $f_{t,d}$ is the raw count of a term t in a document d and $idf(t, D) = \log \frac{N}{n_t}$ where N is the total number of documents in the corpus D and n_t is the number of documents containing the term t , for the subsequent questions consider the following documents:

| ID | Document Text |
|-------|--|
| d_1 | annoyed by oysters devouring oysters |
| d_2 | happy at large <u>fish</u> devouring small <u>fish</u> |
| d_3 | what is the <u>fish</u> devouring once more? |
| d_4 | like crabs devouring oysters and fish |

Given the set of terms $T = \{crabs, devouring oysters, oysters, fish\}$ answer the following:

- (a) Compute the tf for the terms in T for each document.

| | Crabs | devouring oysters | Oysters | fish |
|-------|--------|-------------------|---------|--------|
| d_1 | 0 | 0.3010 | 0.4771 | 0 |
| d_2 | 0 | 0 | 0 | 0.4771 |
| d_3 | 0 | 0 | 0 | 0.3010 |
| d_4 | 0.3010 | 0.3010 | 0.3010 | 0.3010 |

- (b) Compute idf for the terms in T for the corpus.

$$df(crabs) = 1$$

$$df(devouring\ oysters) = 2$$

$$df(oysters) = 2$$

$$df(fish) = 3$$

$$idf(crabs) = \log \left[\frac{4}{1} \right] = 0.60205$$

$$idf(devouring\ oysters) = \log \left[\frac{4}{2} \right] = 0.301029$$

$$idf(oysters) = \log \left[\frac{4}{2} \right] = 0.301029$$

$$idf(fish) = \log \left[\frac{4}{3} \right] = 0.12493$$

- (c) Compute $tf-idf(t, D, d)$ for the terms in T for each document in the corpus.

| | Crabs | devouring oysters | Oysters | fish |
|-------|--------|-------------------|---------|--------|
| d_1 | 0 | 0.0906 | 0.1436 | 0 |
| d_2 | 0 | 0 | 0 | 0.0596 |
| d_3 | 0 | 0 | 0 | 0.0376 |
| d_4 | 0.1812 | 0.0906 | 0.0906 | 0.0376 |

Q2. Sentiment Analysis (10 Marks) [4+3+3]

The Boolean Naïve Bayes pseudocode given below (**Algorithm 1**) uses α -weighted Laplace smoothing to train a classification model. The algorithm uses words as features for classification.

Algorithm 1 TRAINNAÏVEBAYES(\mathcal{C}, \mathcal{T})

```

1: procedure NAÏVEBAYESTRAINING( $\mathcal{C}, \mathcal{T}$ )
2:    $V \leftarrow \text{EXTRACTVOCABULARY}(\mathcal{C})$ 
3:    $N \leftarrow \text{COUNTTEXTS}(\mathcal{T})$ 
4:   for each  $t \in \mathcal{T}$  do
5:     REMOVEDUPLICATEWORDS( $t$ )
6:   for each  $c \in \mathcal{C}$  do
7:      $N_c \leftarrow \text{COUNTTEXTSINCLASS}(\mathcal{T}, c)$ 
8:      $N_w \leftarrow \text{COUNTWORDSONALLTEXTSOFCORPUS}(\mathcal{T}, c)$ 
9:      $prior[c] \leftarrow \frac{N_c}{N}$   $\rightarrow prior[c] = \frac{\text{number of } c \text{ classes in corpus}}{\text{total number of classes in corpus}}$ 
10:     $doc_c \leftarrow \text{CONCATENATETEXTSINCLASS}(\mathcal{T}, c)$ 
11:    for each  $w_i \in V$  do
12:       $N_i \leftarrow \text{COUNTTOKENSOFWORDS}(doc_c, w_i)$ 
13:       $condprob[w_i][c] \leftarrow \frac{N_i + \alpha}{N_w + \alpha|V|}$ 
14:  return  $V, prior, condprob$ 

```

▷ THE NB MODEL

- (a) A variant of Naïve Bayes called Multinomial Naïve Bayes sometimes performs worse because it does not normalize text data. Modify the above code to convert it to the Multinomial version of the Naïve Bayes.
- (b) Some logical mistakes may have been intentionally added to the above code. Circle any error(s) and state their correction. If there are no errors, circle **None**.
- (c) **Algorithm 2** shows the testing for the Boolean Naïve Bayes classifier. The probability scores are computed using $P(c) \cdot \prod_i P(w_i|c)$ may result in an underflow because of small probability values being multiplied. Suggest a solution and modify the algorithm given below to implement your solution. To avoid the problem of underflow instead of multiplying the values we will take log of each probabilities and add them all.

Algorithm 2 TESTNAÏVEBAYES($\mathcal{C}, \mathcal{V}, prior, condprob, t$)

```

1: procedure NAÏVEBAYESINFERENCE( $\mathcal{C}, \mathcal{V}, prior, condprob, t$ )
2:    $W \leftarrow \text{EXTRACTWORDSFROMTEXT}(\mathcal{V}, t)$ 
3:   for each  $c \in \mathcal{C}$  do
4:      $score[c] \leftarrow \log prior[c]$ 
5:     for each  $w_i \in W$  do
6:        $score[c] *= condprob[w_i][c]$ 
7:  return  $\arg\max_{c \in \mathcal{C}} score[c]$ 

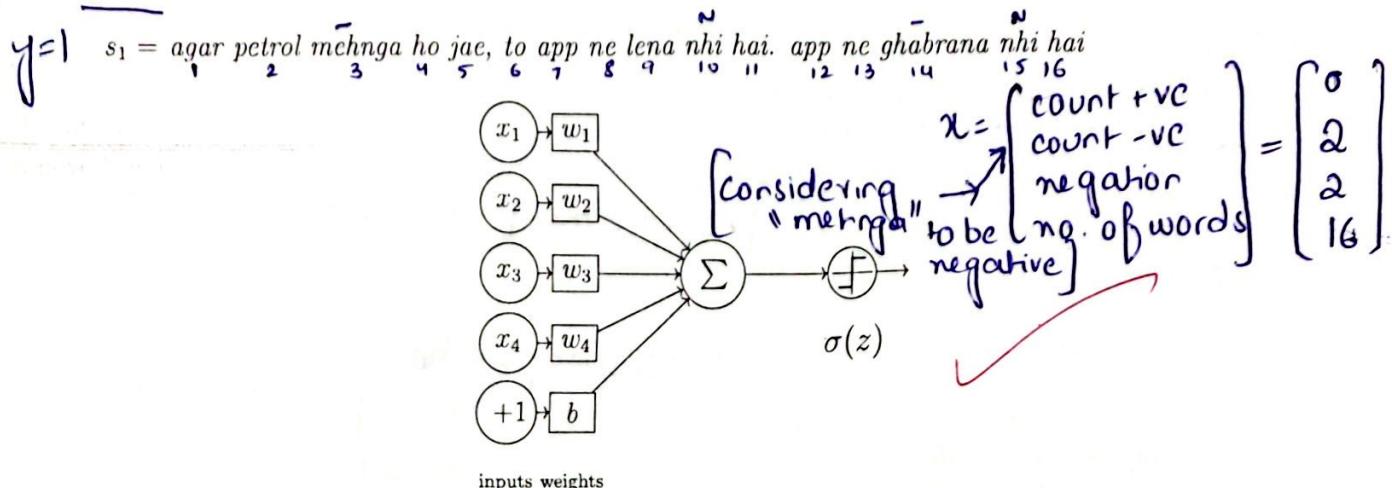
```

2
Score [c] += log[condprob[w_i]]
✓ ▷ THE PREDICTED CLASS

Q3. Logistic Regression

(10 Marks) [5+5]

Consider the following sentence used for training the following Logistic Regressor to detect sarcasm.



Given that s_1 is labeled positive and s_2 is labeled as negative, and features x_1 (count +ive lexicon), x_2 (count -ive lexicon), x_3 (negation words) and x_4 (number of words) and the initial weight and bias vector: $W, b = [0, 0, 0, 0, 0]$

Determine the weights and bias vector after training and backpropagating the logistic regressor on the above sentence, given that the learning rate $\alpha = 0.1$.

$$z = w^T x + b$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 & 2 & 2 & 16 \end{bmatrix} + 0$$

$$= 0 + (2)(0) + 0(2) + 16(0) + 0$$

$$z = 0$$

$$\hat{y} = \sigma(z) = \frac{1}{1+e^{-0}} = 0.5$$

$$w_{t+1} = w_t - \alpha \text{[loss function]}$$

$$w_{t+1} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + 0.1 \begin{bmatrix} \sigma(2) - y \\ \sigma(2) - y \\ \sigma(2) - y \\ \sigma(2) - y \\ \sigma(2) - y \end{bmatrix} \begin{bmatrix} x \end{bmatrix}$$

$$w_{t+1} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} - 0.1 \begin{bmatrix} 0.5 - 1 \\ 0.5 - 1 \\ 0.5 - 1 \\ 0.5 - 1 \\ 0.5 - 1 \end{bmatrix} \begin{bmatrix} x \end{bmatrix}$$

$$= \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} - 0.1 \begin{bmatrix} -1 \\ -1 \\ -8 \\ -0.5 \end{bmatrix}$$

$$w_{t+1} = \begin{bmatrix} 0 \\ 0.1 \\ 0.1 \\ 0.8 \\ 0.05 \end{bmatrix}$$

Fall 2023 after back propagation Page 4 of 10
 updated bias = 0.05

$$\text{updated weights} = \begin{bmatrix} 0 \\ 0.1 \\ 0.1 \\ 0.8 \end{bmatrix}$$

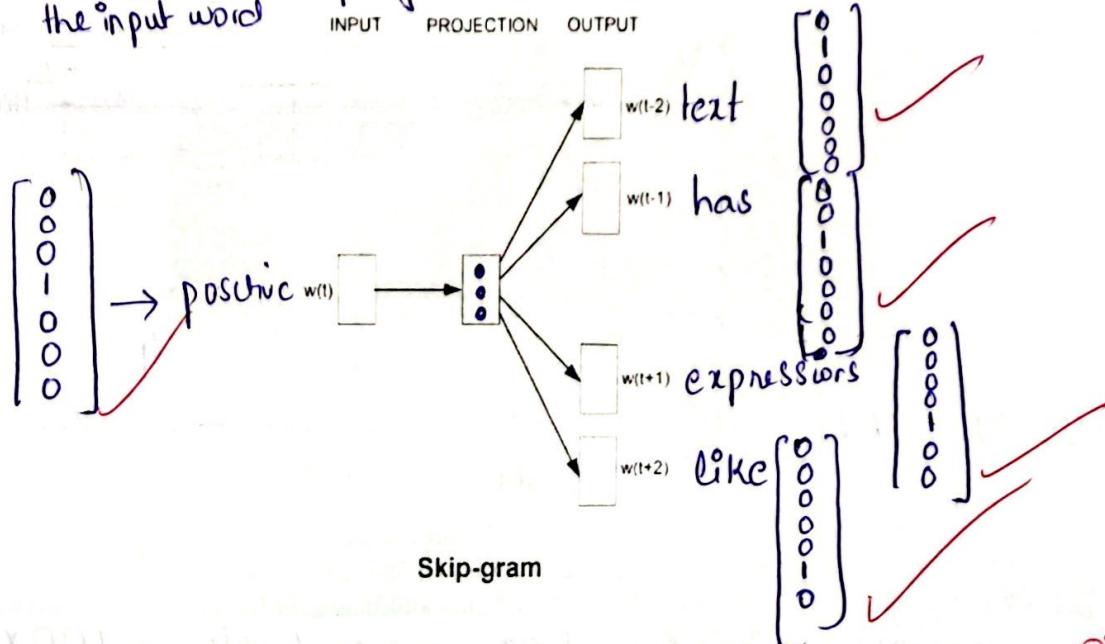
Q4. Vector Semantics

8/10

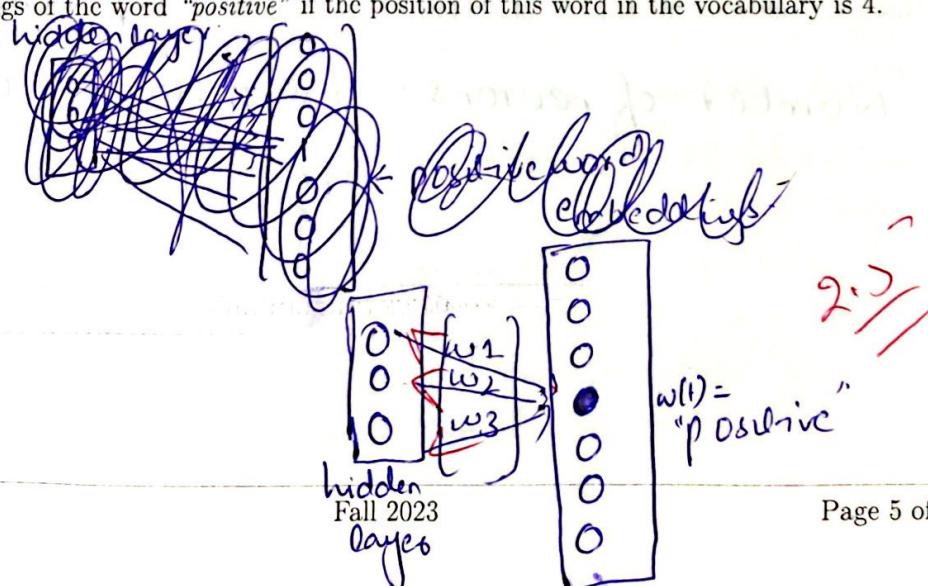
(10 Marks) [2+3+2+3]

Consider the test corpus ["the text has positive expressions like good"] for training the skip-gram model below for generating word embeddings, where the vocabulary size is 7.

Since it is skip-gram it will mask all other words here
 the input word

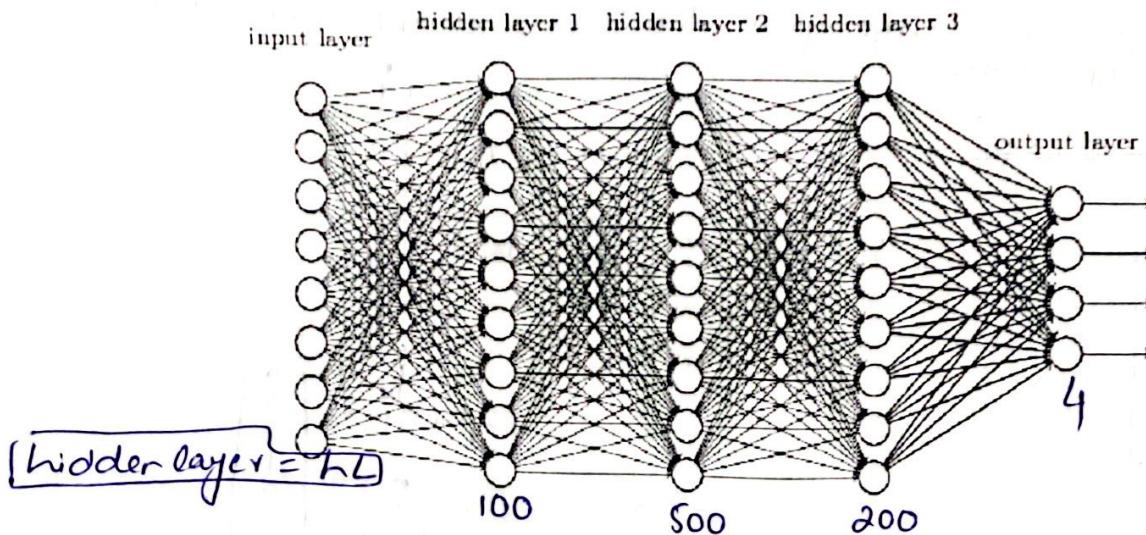


- (a) What is the format of the words in the output layer? Show an example on the diagram above.
 lets consider ~~no~~ positive is ~~needed~~ input we use one hot encoding for
- (b) Label the diagram above to show the predicted words if the input word is "positive".
~~Output-~~
- (c) What is the significance of the size of the hidden layer in the above network?
 size of hidden layer determines the dimensions of the output layer and fine tuning of the hyper parameters
- (d) Consider the size of the hidden layer to be 3. Now draw the weights that would represent the embeddings of the word "positive" if the position of this word in the vocabulary is 4.



Q5. Feed Forward Networks (10 Marks) [7+3]

Consider the Feed-Forward Deep Network architecture below:



- (a) If the network is used to classify datapoints into four classes and there are 100 input features to the network and the hidden layers 1, 2, 3 have 100, 500, 200 neurons respectively, how many weights are there in total, including both weights and biases, in the entire neural network?

$$\text{weights + bias from input to hidden layer 1} = 100 \times 100 + 100$$

$$\text{weights + bias from hidden L1 to L2} = 100 \times 500 + 500 = 50500$$

$$\text{" " " } hL2 - hL3 = 500 \times 200 + 200 = 100200$$

$$\text{total} = 10100 + 50500 + 100200 + 804 = 161604$$

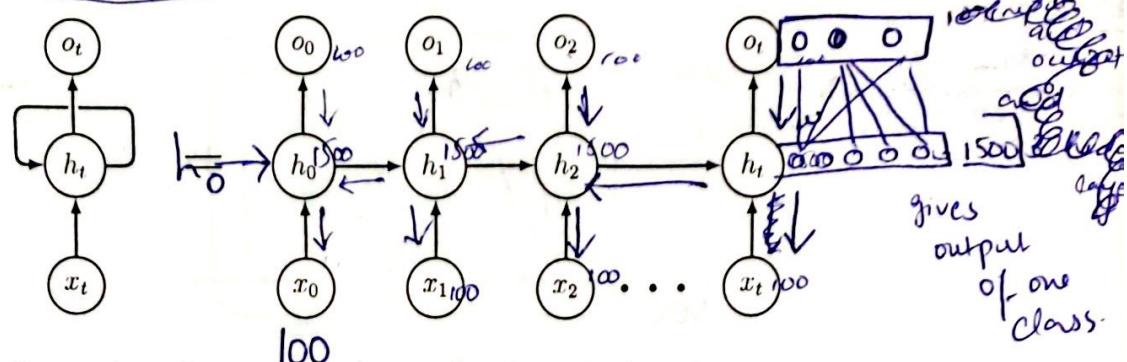
- (b) How many neurons are there in total, in the neural network above?

$$\text{Number of neurons} = 100 + 100 + 500 + 200 + 4 = 904$$

✓ (10)

Q6. Recurrent Neural Networks (10 Marks) [4+3+3]

Consider the RNN given below and its representation unrolled in time. Assume that the network is being used for classifying sentiments as $C = \{+, -, =\}$, the inputs are vectors of size 100 and the hidden layer contains 1500 neurons.



- (a) Determine the number of parameters that need to be trained in this network. Show your work.

bias are 1500

Let's say calculate number of parameters.
 we have input size of 100 and neurons in the hidden layer are 1500 so $[100 \times 1500] + 1500 = 2[15000]$ multiply with 2 because $(100 - 1500) - 100$
 then this is a recurrent process so we will add $1500 \times 1500 = 2250000$ added in them is now at final layer output we have 3 classes so $= 100 \times 3$
 so total we have $= 2[15000] + 2250000 + 300 + \text{bias}$

- (b) Modify the unrolled RNN diagram above to show what the output layer would look like for the given classification task?

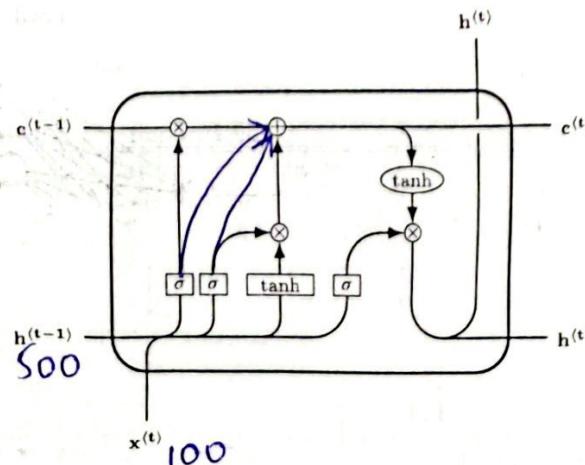
- (c) How do vanishing and exploding gradient problems affect the training of RNNs? What are the consequences of these gradient issues during the backpropagation process?

Vanishing gradient problem occurs after the value becomes nearer to 0 or 0 while back propagating. as we move far back in the sequence so the value for the earlier vectors in sequence we loose them thus context of far away or earlier words we forget and in exploding gradient problem the values number of bits become too big higher than they are hard to store. and they begin to overflow.

Q7. Long-Short Term Memory

(15 Marks) [5+5+5]

Consider the LSTM cell given below.



- (a) Given that the LSTM cell is being used in a network for classifying text into 5 classes, the input $x^{(t)}$ is a vector of size 100 and the hidden state vector $h^{(t)}$ is of size 500, find the total number of parameters to be trained for the classification network containing the LSTM cell above.

$$4[(100+500) \times 500 + 500] + 500 \times 5 + 5 \\ = 1202000 + 12500$$

Total number of parameters to be trained. $\rightarrow = 1214500$

- (b) We want to modify the LSTM cell by directly adding whatever we have learnt from the input and hidden state at timestep t to update the cell state. Make changes to the diagram above to reflect this modification.

whatever learnt from input

- (c) What would be the updated equation for the new cell state with the modification made in part (b) above.

before $C_t = (f_t \times C_{t-1}) + (I_t \times \tilde{C}_t)$
 "since we are directly adding the input and hidden state
 after $C_t = f_t + I_t$

$C_t = \text{forget gate} + \text{Input gate}$

Final Exam

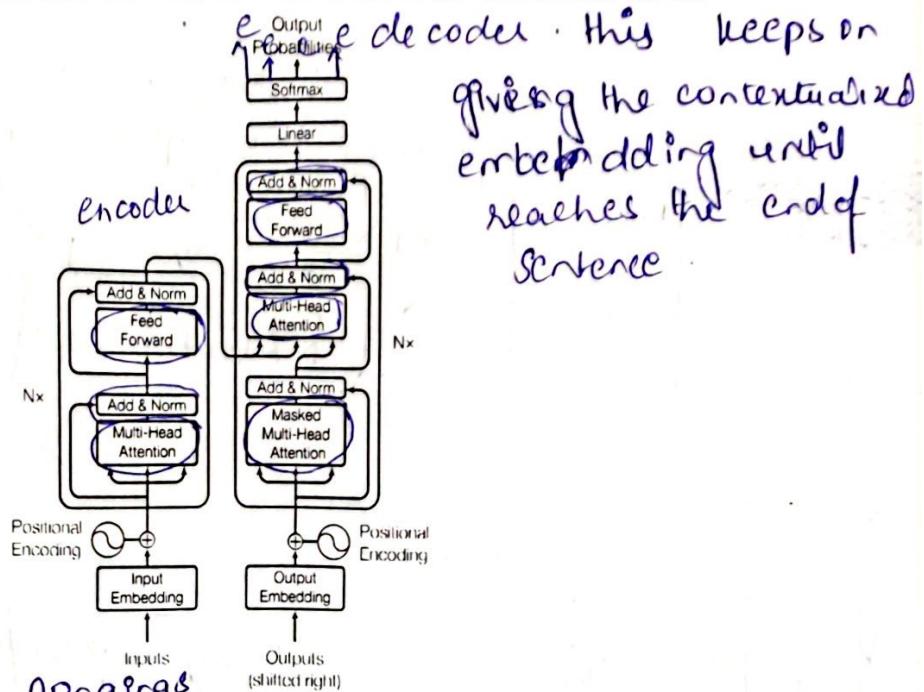
Fall 2023

Page 8 of 10

updated $C_t = \sigma \left[W_f^T [h_{t-1} | x_t] + b_f \right] + \sigma \left[W_i^T [h_{t-1} | x_t] + b_i \right]$

Q8. Attention and Transformers (15 Marks) [4+3+3+5]

The following diagram shows the Encoder-Decoder architecture of the Transformer model.



Token embeddings are forwarded

Figure 1: Transformer Model Architecture.

- (a) Consider that the model is trained on English-to-Urdu translation. Given that the input is "Tell these longings to go dwell elsewhere", use the diagram above to label and describe the output of the Encoder for this input. Assume that the transformer embeddings are of size 3.

(b) Label the residual connections of the Decoder on the diagram above. X ①

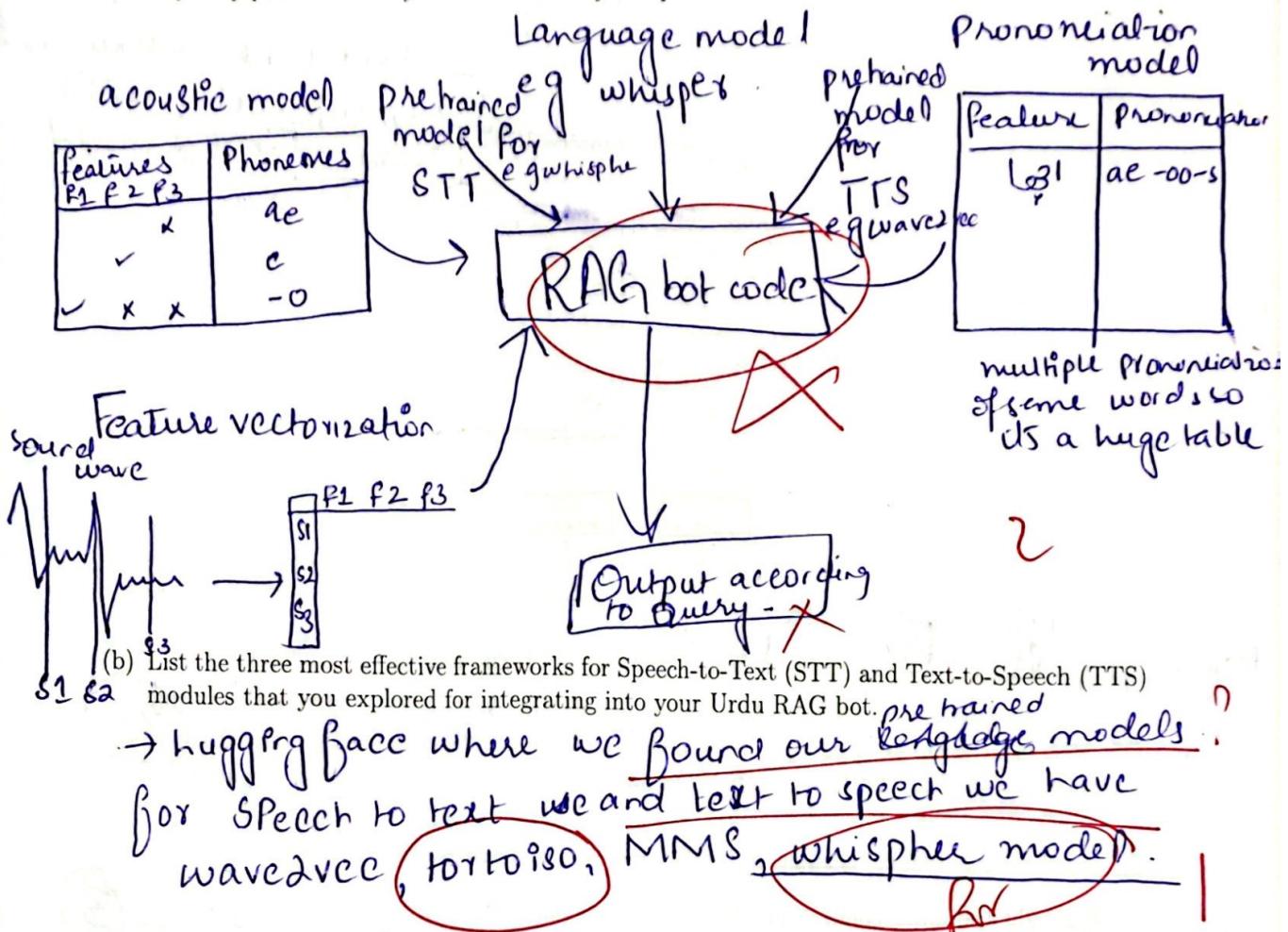
(c) Highlight the units/connections that contain the parameters in the transformer encoder. ✓ ③

(d) Describe using the example sentence "Tell these longings to go dwell elsewhere" how *Positional Encoding* is determined? Assume that the input embeddings are of size 3.

Positional encodings are formed by the position by ~~a one-hot encoding~~ in the 768 dimensions of a word eg we take Tell $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ it's similar to that as in one hot encoding except for in this we add the sequence at which the word is present - ③

Q9. Retrieval Augmented Generation (10 Marks) [5+2+3]

- (a) Recall your Urdu Retrieval Augmented Generation (RAG) bot code. Draw the architecture of your pipeline: identify and label the key components.



- (c) List any observed limitations or issues in the voice integration process and propose potential improvements. Your response should showcase a thoughtful analysis of the impact of voice integration on the overall user experience.

The voice generation for all sorts of texts whether it be a normal reading, some sort of story, or poetry it's monotone. All read the same way so to improve the user experience the voice generated should do a semantic analysis and read poetry or generate poetry ~~in~~ voice in a way that it is read with emphasis laid on words that are required by the context and tone.