

# CSCI 5922 - 001 Neural Networks and Deep Learning

## Problem Set 4

### 1. Transformers vs RNNs vs MLPs

- a) Transformers and RNNs are similar because,
  - Both Transformers and RNNs handle the sequence to sequence data related problem statements such as language translation, speech recognition and text processing.
  - Both these models can be trained using a very large text datasets/corpus with a variety of word embeddings feeded into the model.
  - Both models understand the relationships between different words by passing the states through timestamps.
- b) Transformers and RNNs are different because,
  - RNNs struggle with the Long short term memory handling issue and in RNNs we cannot pass the input simultaneously but it is possible with Transformers with the help of self attention.
  - Transformers get the whole sentence at once but the RNNs will depend on the past words and sometimes it will miss important context required for predicting the next word.
  - RNNs handle different sequence lengths whereas the transformers use the positional encoding.
- c) Transformers are similar to MLPs because,
  - Both are fully connected architectures which contains weight vectors, activation functions and layers of neurons for learning the intricate patterns about the data.
  - Both of these models heavily rely on the mathematical computation of matrices so both of them are GPU intensive tasks.
  - Both of these models can be used for the classification task for example Transformers can be used in Classification of the text “Spam Detection” and MLP’s also classify ie.“Credit Card Fraud Detection”
- d) Transformers are different from MLPs, because
  - Transformers work with sequence to sequence text data (unstructured data) while the MLPs work well with the structured numerical data.
  - Transformers learns the contextual based relationships of generating a word / sequence of words but MLPs just maps the inputs to the outputs and they won’t even bother about the context.

- Transformers focus on the highly important words by means of self attention but in MLPs there is no ability like that.

## 2. Self attention Equation

a )

- i. When computing the  $QK^T$  the values will become very large if the  $d_k$  ie. the dimension of the key vectors is high. This may cause the softmax function to produce very small gradient values, which makes the learning process very slow and leads to vanishing gradient problem. Dividing by  $\sqrt{d_k}$  will scale down the values and keeps them in a stable range. This prevents the softmax function from producing overfitted probability values and it will also help in improving the training stability.
- ii. The softmax activation function will convert the computed attention scores into probabilities ie. from the range of 0 to 1. The output of the softmax function produces a vector which sums to 1. This activation function will solve the problem of which relevant tokens should be considered more priority ie. helps in deciding the most relevant tokens. It will also help in a level of differentiation in which it amplifies the larger scores by suppressing the smaller scores which helps the self attention mechanism easier to capture the correlation between the input sequence.

b) The term equation (1) is referred to as the attention map/attention matrix because

- It will show how much attention each word in a sentence should give to other words in a sentence.
- It is called the attention map because it marks the important words and reduces the attention on the less relevant words.
- Each row represents a Query (word) and the column represents Key (Word) which forms the Attention Scores matrix. ie. The matrix multiplication of Query and keys forms attention matrix.
- This will also help in the understanding of which words should get importance in the given context.

c) For example in an example of Machine Translation of English to Hindi, if we use single headed attention the model will focus only one aspect of word at a time and will miss the important contexts.

For example: I love apple. Here ‘apple’ is represented as a fruit / company?  
The single attention head will struggle to capture both of the meanings.

But the multi headed attentions will capture multiple meanings by means of different attention heads. And some heads will focus on the nearby words and other heads will capture the long range dependency words to help in improving the understanding.

Multiple heads will provide the wide range of perspectives, which prevents biasedness of relying on a single feature. It will also process the input simultaneously at the same time making the training process more efficient.

d) Transformer encoder will read the whole i/p sequence at once and it will understand how the different parts of the sentence relate to each other. It has layers that will help in learning the meaning and the context of the data.

The layers that it consists are self attention layers and feed forward neural networks which helps in learning the patterns of input data.

Example: Bidirectional Encoder Representations from Transformers [BERT] is used for the sentimental analysis of Food Reviews in Uber Eats. It helps the organisation to make investment and marketing decisions. It will process the entire text bidirectionally, it is well suitable for the understanding of context rich customer reviews.

Transformer decoder on the other end works step by step. It will predict based on the past word and the context from the encoder. ie. It will generate each and every word based on the previous words and the context is from the encoder.

The transformer decoder has self attention layer, encoder decoder attention layer and feed forward neural networks layers to understand both the previous knowledge as well as past knowledge from encoder.

Example: In Healthcare BioGPT which is created to write the medical reports. The decoder ensures the text flows such as summarizing the patient records and the test results which makes easier for the doctors to make decisions.

The encoders are best for analyzing the data and the decoders are best for generating new data.

## Neural Networks & Deep Learning

### Problem set 4

3) i/p tokens:

$$\begin{bmatrix} 1 & -1 & 0 & 2 \end{bmatrix}$$

Query weights

$$W_Q = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ -0.5 & 0 \\ 1 & 1 \end{bmatrix}$$

Key weights

$$W_K = \begin{bmatrix} 1 & -0.5 \\ 0 & 1 \\ 0.5 & -0.5 \\ 1 & 0.5 \end{bmatrix}$$

Value weights

$$W_V = \begin{bmatrix} 1 & 0.5 \\ -0.5 & 0 \\ -1 & 2 \\ 0.5 & -1 \end{bmatrix}$$

For i/p token 1:  $[1 \ -1 \ 0 \ 2]$

$$\text{Query } Q_1 : i/p_1 \times W_Q = [1 \ -1 \ 0 \ 2] \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ -0.5 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 4} \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ -0.5 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 2} \\ = [0 -1 + 0 + 2 \quad 1 + 0 + 0 + 2] \\ Q_1 = [1 \ 3]$$

$$\text{Key } K_1 : i/p_1 \times W_K = [1 \ -1 \ 0 \ 2] \begin{bmatrix} 1 & -0.5 \\ 0 & 1 \\ 0.5 & -0.5 \\ 0.5 & 0.5 \end{bmatrix}_{4 \times 4} \begin{bmatrix} 1 & -0.5 \\ 0 & 1 \\ 0.5 & -0.5 \\ 0.5 & 0.5 \end{bmatrix}_{4 \times 2} \\ = [1 + 0 + 0 + 2 \quad -0.5 - 1 + 0 + 1] \\ K_1 = [3 \ -0.5]$$

$$\text{Value } V_1 : i/p_1 \times W_V = [1 \ -1 \ 0 \ 2] \begin{bmatrix} 1 & 0.5 \\ -0.5 & 0 \\ -1 & 2 \\ 0.5 & -1 \end{bmatrix}_{4 \times 4} \begin{bmatrix} 1 & 0.5 \\ -0.5 & 0 \\ -1 & 2 \\ 0.5 & -1 \end{bmatrix}_{4 \times 2} \\ = [1 + 0.5 + 0 + 1 \quad 1/2 + 0 + 0 - 2] \\ V_1 = [2.5 \ -1.5]$$

for i/p token 2:  $[0 \ -1 \ 2 \ 1]$

$$\text{Query } Q_2 : i/p_2 \times W_Q = [0 \ -1 \ 2 \ 1] \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ -0.5 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 4} \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ -0.5 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 2} \\ = [0 -1 - 1 + 1 \quad 0 + 0 + 0 + 1] \\ Q_2 = [-1 \ 1]$$

$$\text{Key } K_2 : i/p_2 \times W_K = [0 \ -1 \ 2 \ 1] \begin{bmatrix} 1 & -0.5 \\ 0 & 1 \\ 0.5 & -0.5 \\ 0.5 & 0.5 \end{bmatrix}_{4 \times 4} \begin{bmatrix} 1 & -0.5 \\ 0 & 1 \\ 0.5 & -0.5 \\ 0.5 & 0.5 \end{bmatrix}_{4 \times 2} \\ = [0 + 0 + 1 + 1 \quad 0 - 1 - 1 + 1] \\ K_2 = [2 \ -1.5]$$

$$\text{Value } V_2 : i/p_2 \times W_V = [0 \ -1 \ 2 \ 1] \begin{bmatrix} 1 & 0.5 \\ -0.5 & 0 \\ -1 & 2 \\ 0.5 & -1 \end{bmatrix}_{4 \times 4} \begin{bmatrix} 1 & 0.5 \\ -0.5 & 0 \\ -1 & 2 \\ 0.5 & -1 \end{bmatrix}_{4 \times 2} \\ = [0 + 1/2 - 2 + 1/2 \quad 0 + 0 + 4 - 1] \\ V_2 = [-1 \ 3]$$

for i/p token 3:  $[-2 \ 2 \ -1 \ 0]$

$$\text{Query } Q_3 : i/p_3 \times W_Q = [-2 \ 2 \ -1 \ 0] \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ -0.5 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 4} \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ -0.5 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 2} \\ = [0 + 2 + k_3 + 0 \quad -2 + 0 + 0 + 0] \\ Q_3 = [-2.5 \ -2]$$

$$\text{Key } K_3 : i/p_3 \times W_K = [-2 \ 2 \ -1 \ 0] \begin{bmatrix} 1 & -0.5 \\ 0 & 1 \\ 0.5 & -0.5 \\ 0.5 & 0.5 \end{bmatrix}_{4 \times 4} \begin{bmatrix} 1 & -0.5 \\ 0 & 1 \\ 0.5 & -0.5 \\ 0.5 & 0.5 \end{bmatrix}_{4 \times 2} \\ = [-2 + 0 - 0.5 + 0 \quad 1 + 2 + 0.5 + 0] \\ K_3 = [-2.5 \ 3.5]$$

$$\text{Value } V_3 : i/p_3 \times W_V = [-2 \ 2 \ -1 \ 0] \begin{bmatrix} 1 & 0.5 \\ -0.5 & 0 \\ -1 & 2 \\ 0.5 & -1 \end{bmatrix}_{4 \times 4} \begin{bmatrix} 1 & 0.5 \\ -0.5 & 0 \\ -1 & 2 \\ 0.5 & -1 \end{bmatrix}_{4 \times 2} \\ = [-2 - 1 + 1 + 0 \quad -1 + 0 - 2 + 0] \\ V_3 = [-2 \ -3]$$

### Computing the attention scores

#### I/p token 1

We have to multiply the Query with all the keys.

$$Q_1 \times K_1 = [1 \ 3] \times [3 - 0.5]^T = 3 - 1.5 = 1.5$$

$$Q_1 \times K_2 = [1 \ 3] \times [2 - 1.5]^T = 2 - 4.5 = -2.5$$

$$Q_1 \times K_3 = [1 \ 3] \times [-2.5 \ 3.5]^T = -2.5 + 10.5 = 8$$

#### I/p token 2

$$Q_2 \times K_1 = [-1 \ 1] \times [1 - 0.5]^T = -3 - 0.5 = -3.5$$

$$Q_2 \times K_2 = [-1 \ 1] \times [2 - 1.5]^T = -2 - 1.5 = -3.5$$

$$Q_2 \times K_3 = [-1 \ 1] \times [-2.5 \ 3.5]^T = 2.5 + 3.5 = 6$$

#### I/p token 3:

$$Q_3 \times K_1 = [-2.5 \ -2] \times [3 - 0.5]^T = 7.5 + 1 = 8.5$$

$$Q_3 \times K_2 = [-2.5 \ -2] \times [2 - 1.5]^T = 5 + 3 = 8$$

$$Q_3 \times K_3 = [-2.5 \ -2] \times [-2.5 \ 3.5]^T = 6.25 - 7 = -13.25$$

### computing the attention weights & attention maps

#### I/p token 1:

$$\text{Attention weights} = \text{softmax}([1.5, -2.5, 8]) \rightarrow \max$$

$$\text{softmax}(z_i) = \frac{e^{z_i - \max(z)}}{\sum_{j=1}^n e^{z_j - \max(z)}}$$

$$\text{numer: } e^{z_i - \max(z)} = [1.5 - 8, -2.5 - 8, 8 - 8] = \left[ \frac{-6.5}{e^{-8}}, \frac{-10.5}{e^{-8}}, 0 \right] = [0, 0, 1]$$

$$\text{denom: } \sum_{j=1}^n e^{z_j - \max(z)} = e^{1.5 - 8} + e^{-2.5 - 8} + e^{8 - 8} = e^{-6.5} + e^{-10.5} + e^0 = 1$$

$$\text{here } e^{-6.5} \approx 0, e^{-10.5} \approx 0, e^0 = 1$$

$$\text{so, } \text{softmax}([1.5, -2.5, 8]) = (0, 0, 1)$$

$$\text{attention weights} = (0, 0, 1) \rightarrow (aw_1, aw_2, aw_3)$$

#### Attention maps $\rightarrow$ attention $\times$ values

$$aw_1 \times V_1 = 0 \times [2.5, -1.5] = [0 \ 0]$$

$$aw_2 \times V_2 = 0 \times [-1 \ 3] = [0 \ 0]$$

$$aw_3 \times V_3 = 1 \times [-2 \ -3] = [-2 \ -3]$$

#### I/p token 2:

$$\text{Attention weights} = \text{softmax}([-3.5, -3.5, 6]) \rightarrow \max$$

$$\text{softmax}(z_i) = \frac{e^{z_i - \max(z)}}{\sum_{j=1}^n e^{z_j - \max(z)}}$$

$$\text{numer: } e^{z_i - \max(z)} = [-3.5 - 6, -3.5 - 6, 6 - 6] = \left[ \frac{-9.5}{e^{-6}}, \frac{-9.5}{e^{-6}}, 0 \right] = [0, 0, 1]$$

$$\text{denom: } \sum_{j=1}^n e^{z_j - \max(z)} = e^{-9.5} + e^{-9.5} + e^0 = 1$$

$$\text{here } e^{-9.5} \approx 0, e^{-9.5} \approx 0, e^0 = 1$$

$$\text{so, } \text{softmax}([-3.5, -3.5, 6]) = (0, 0, 1)$$

$$\text{attention weights} = (0, 0, 1) \rightarrow (aw_1, aw_2, aw_3)$$

#### Attention maps $\rightarrow$ attention $\times$ values

$$aw_1 \times V_1 = 0 \times [2.5, -1.5] = [0 \ 0]$$

$$aw_2 \times V_2 = 0 \times [-1 \ 3] = [0 \ 0]$$

$$aw_3 \times V_3 = 1 \times [-2 \ -3] = [-2 \ -3]$$

#### I/p token 3:

$$\text{Attention weights} = \text{softmax}([8.5, 8, -13.25]) \rightarrow \max$$

$$\text{softmax}(z_i) = \frac{e^{z_i - \max(z)}}{\sum_{j=1}^n e^{z_j - \max(z)}}$$

$$\text{numer: } e^{z_i - \max(z)} = [8.5 - 8.5, 8 - 8.5, -13.25 - 8.5] = \left[ \frac{0}{e^{-8.5}}, \frac{0}{e^{-8.5}}, e^{-13.25} \right] = [0, 0, 1]$$

$$\text{denom: } \sum_{j=1}^n e^{z_j - \max(z)} = e^{8.5 - 8.5} + e^{8 - 8.5} + e^{-13.25 - 8.5} = 1 + 0.6 + 1.6 = 3.2$$

$$\text{here } e^{-13.25} \approx 0, e^{-8.5} \approx 0, e^0 = 1$$

$$\text{so, } \text{softmax}([8.5, 8, -13.25]) = (0.62, 0.37, 0)$$

$$\text{attention weights} = (0.62, 0.37, 0) \rightarrow (aw_1, aw_2, aw_3)$$

#### Attention maps $\rightarrow$ attention $\times$ values

$$aw_1 \times V_1 = 0.62 \times [2.5, -1.5] = [1.55, -0.92]$$

$$aw_2 \times V_2 = 0.37 \times [-1 \ 3] = [-0.37, 1.11]$$

$$aw_3 \times V_3 = 0 \times [-2 \ -3] = [0 \ 0]$$

b) By using the attention maps got from (a) The i/p tokens are :

$$\text{for i/p token 1: } [0 \ 0] + [0 \ 0] + [-2 \ -3] = [-2 \ -3]$$

$$\text{for i/p token 2: } [0 \ 0] + [0 \ 0] + [-2 \ -3] = [-2 \ -3]$$

$$\text{for i/p token 3: } [1.55 \ -0.92] + [-0.37 \ 1.11] + [0 \ 0] = [1.18 \ 0.18]$$

c) Visualization of attention map from (a)



d) Observations :

\* Output 1 (o/p 1) & output 2 (o/p 2) are strongly influenced by the Input 1 (i/p 1) which indicates that i/p 1 has the highest attention weight for both o/p 1 & o/p 2 that is indicated by the solid line

\* Output 3 (o/p 3) is strongly influenced by the Input 1 (i/p 1) which indicates that i/p 1 has the highest attention weight for o/p 3 that is indicated by the solid line

\* Output 2 (o/p 2) is slightly influenced by the Input 2 (i/p 2) which indicates that i/p 2 has the second highest attention weight in the attention weight vector that is indicated by the dashed line

f)

a) Computing the cross attention map.

Query

$$Q = X_1 W_Q$$

$$= \begin{bmatrix} -1 & 1 & 0 \\ 0 & 1 & 0 \\ 2 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & -1 \\ 1 & 1 \\ -1 & 0 \end{bmatrix}_{3 \times 2}$$

$$= \begin{bmatrix} 0+1+0 & 1+1+0 \\ 0+1+0 & 0+1+0 \end{bmatrix}$$

$$Q = \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix}$$

Key

$$K = X_2 W_K$$

$$= \begin{bmatrix} -1 & 1 \\ 0 & 1 \\ 1 & 0 \\ 2 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -1 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix}_{4 \times 2}$$

$$K = \begin{bmatrix} -2 & 1 \\ -1 & 1 \\ 0 & 1 \\ 2 & 0 \end{bmatrix}$$

Value

$$V = X_2 W_V$$

$$= \begin{bmatrix} -1 & 1 \\ 0 & 1 \\ 1 & 1 \\ 2 & 0 \end{bmatrix} \begin{bmatrix} -1 & 0 & 1 & 2 \\ 1 & -1 & 0 & -2 \end{bmatrix}$$

$$V = \begin{bmatrix} 2 & -1 & -1 & -4 \\ 1 & -1 & 0 & -2 \\ 0 & -1 & 1 & 0 \\ -2 & 0 & 2 & 4 \end{bmatrix}$$

$$\text{Cross attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V$$

$$QK^T = \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} -2 & -1 & 0 & 2 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

$$QK^T = \begin{bmatrix} 0 & 1 & 2 & 2 \\ -1 & 0 & 1 & 2 \end{bmatrix}$$

softmax

softmax for row 1 of  $QK^T$ 

$$\text{softmax}([0 \ 1 \ 2 \ 2]) \rightarrow \max(z^1)$$

$$\text{softmax}(z_i) = \frac{e^{z_i - \max(z)}}{\sum_{j=1}^4 e^{z_j - \max(z)}}$$

$$\text{numer: } e^{z_i - \max(z)} = e^{0-2}, e^{1-2}, e^{2-2}, e^{2-2}$$

$$= e^{-2}, e^{-1}, e^0, e^0$$

$$= \left( \frac{1}{e^2}, \frac{1}{e^1}, 1, 1 \right) = \left( \frac{1}{7.34}, \frac{1}{2.71}, 1, 1 \right)$$

$$\text{denom: } \sum_{j=1}^4 e^{z_j - \max(z)} = \frac{1}{e^2} + \frac{1}{e^1} + 1 + 1 = \frac{1}{7.34} + \frac{1}{2.71} + 1 + 1 = 2.50$$

$$\text{here } e^0 = 1, e^1 = 2.71, e^2 = 7.34$$

$$\text{so, softmax}([0 \ 1 \ 2 \ 2]) = \left[ \frac{1}{7.34 \times 2.5}, \frac{1}{2.71 \times 2.5}, \frac{1}{2.5}, \frac{1}{2.5} \right]$$

$$= [0.05, 0.15, 0.40, 0.40]$$

softmax for row 2 of  $QK^T$ softmax  $(-1 \ 0 \ 1 \ 2) \rightarrow \max(z^2)$ 

$$\text{numer: } e^{z_i - \max(z)} = (e^{-1-2}, e^{0-2}, e^{1-2}, e^{2-2})$$

$$= (e^{-3}, e^{-2}, e^{-1}, e^0) = \left( \frac{1}{19.90}, \frac{1}{7.34}, \frac{1}{2.71} \right)$$

$$\text{denominator: } e^{-3} + e^{-2} + e^{-1} + e^0 = \frac{1}{19.90} + \frac{1}{7.34} + \frac{1}{2.71} + 1 = 1.55$$

$$\text{softmax}(-1 \ 0 \ 1 \ 2) = \left[ \frac{1}{19.90 \times 1.55}, \frac{1}{7.34 \times 1.55}, \frac{1}{2.71 \times 1.55}, \frac{1}{1.55} \right]$$

$$= [0.03, 0.08, 0.24, 0.64]$$

b) Output

O/p 1

$$= [0.05 \ 0.15 \ 0.40 \ 0.40] \times V$$

$$= [0.05 \ 0.15 \ 0.40 \ 0.40] \times \begin{bmatrix} 2 & -1 & -1 & -4 \\ 1 & -1 & 0 & -2 \\ 0 & -1 & 1 & 0 \\ -2 & 0 & 2 & 4 \end{bmatrix}_{4 \times 4}$$

$$= [-0.55, -0.6, 1.15, 1.1]$$

O/p 2

$$= [0.03 \ 0.08 \ 0.24 \ 0.64] \times \begin{bmatrix} 2 & -1 & -1 & -4 \\ 1 & -1 & 0 & -2 \\ 0 & -1 & 1 & 0 \\ -2 & 0 & 2 & 4 \end{bmatrix}_{4 \times 4}$$

$$= [-1.14, -0.35, 1.49, 2.28]$$

$$\text{cross attention map} = \begin{bmatrix} -0.55 & -0.6 & 1.15 & 1.1 \\ -1.14 & -0.35 & 1.49 & 2.28 \end{bmatrix}$$

from the cross attention map

$$O/p 1 = -0.55 - 1.14 = -1.69$$

$$O/p 2 = -0.6 - 0.35 = -0.95$$

$$O/p 3 = 1.15 + 1.49 = 2.64$$

$$O/p 4 = 1.1 + 2.28 = 3.38$$

c) visualization of attention map



d) \* O/p 1 is strongly influenced by the i/p x2(3) &amp; x3(4) indicated by the solid line &amp; is slightly influenced by the i/p x2(2) indicated by the dashed line.

\* O/p 2 is strongly influenced by the i/p x2(4) indicated by the solid line &amp; is slightly influenced by the i/p x2(3) indicated by the dashed line.