M T W T F S

# Assignment # 4

## Question # 1

Explain concept of self-attention of Transformer mode? How it differ...? mathematical

## SELF - Attention :

Self - Attention is a key component of the transformer model. Which was introduced by 'Vaswani et al.' in Paper 'Attention is All you Need'

It is differ from traditional attention mechanism in its ability to capture relationships between different positions in a sequence, allowing the model to focus on different parts of input sequence when producing output e.g. Sequence-to-Sequence tasks like machine translation.

## Traditional Vs Self Attention Mechanism :

CamScanner

In traditional attention mechanism, a context vector is computed based on the weighted sum of input sequence elements, where weights are determined by similarity b/w current decoding position & each position in the input sequence.

CV (used to compute output at current decoding position)

In Self - attention, input sequence transformed into three vectors.

Query (Q)
Key (K)  } → Compute attention scores
Value (V)

& output → weighted sum of value where weights are determined by attention scores

# Mathmatically :

Give input sequence

$$X = (x_1, x_2, \ldots x_n)$$

we compute Query (Q), key (K), value (V) matrices

$$\left.\begin{array}{l} Q = X W_Q \\ K = X W_K \\ V = X W_V \end{array}\right\} \begin{array}{l} \text{learnable} \\ \text{weights matrices} \end{array}$$

Self-attention scores are computed using dot product of Query & key

$$\text{Attention } (Q, K) = \frac{Q K^T}{\sqrt{d_K}}$$

$$\downarrow$$
dimension of key vector

Scores are scaled & passed through softmax function

$$\underline{\text{Attention\_weights}} = \text{softmax}\left(\frac{\text{Attention}(Q, K)}{\sqrt{d_K}}\right)$$

final output is a weighted sum of values vectors

Output = Attlention_weights V

## Question # 02

Summarize architecture of BERt ? discuss how it differs from ~~transfomdise~~ orginal Focus on theoretical & key innovations?

# BERT :

( Bidi-rectional Encoder Representations From transformers) is a pre-trained language model introduced. by google H builds upon orginal transformer model but introduces key innovations to address the limitations of uni-directional language models.

## ~ARCHITECTURE :

1) Bidirectional context
2) Masked language Model (MLM)
3) Pre-training tasks
4) Transformer Encoder layers.

# 1) Bidirectional Context :-

Bert Processes input text bi-directionally. It considers both left & right during training.

## 2) MLM :

Bert utilizes a masked language model objective during pre-training Randomly selected words in the input sequence are masked. & model is trained to predict these masked words based on their context.

## 3) Pre-training Tasks :

Bert is pre-trained on 2 Unsupervised tasks

- MLM
- NSP ( task involves predicting whether 2 sentences follow each other in training data .
  - ↳ These help Bert Capture deeper Understanding of language semantics.

# Key Difference From

## Original Transformer:
- Bi-directional Processing
- Task Specific Pre-training
- NO Cequence to Cequence task.

## Bi-directional (Bert introduces

Bi-directional Processing, capturing Contextual info from both directions. Orginal transformer which processes input sequentially.

## Task-Specific Pre training:
Bert pre-teams on specific task like MLM & NSP, enable model to capture broader understanding of language semantics & relationships.

## No Sts : Unlike orginal transformer, Bert is not designed for sequence-to sequence tasks.

## Q#03 :-

Consider a scenario where transformer model is struggling with long-range dependencies in input sequence ...... address this issue & explain why your approach could be effective.

To address the issue of long-range dependencies in a transformer model. One effective, modification is to incorporate additional positional information in the form of 'Relative Positional encoding'

## Modification Proposel : RPE

In the original transformer model, Positional encoding are added to the input embeddings to provide information about the positions of each tokens in the sequence.

However, these absolute positional encoding may not effectively capture the relative distances b/w tokens, especially in long sequences.

The proposed modification involves introducing relative positional embeddings that explicitly model the relative distance b/w tokens.

This can be achieved by incorporating trainable parameters for encoding relative positions, allowing model to learn the importance of different positional relationships during training.

## Why it Could be Effective :

† Improved capture of LRDependencies
† Adaptability to sequence lengths
† Enhanced training Dynamics.

1) Relative Positional encodings enable the model to better capture the relationships between tokens at varying distances

This added information helps model understand the context more effectively, especially in situations where long-range dependencies are crucial.

2) Unlike absolute positional encodings, which are fixed regardless of sequence length, relative positional encoding adapt dynamically to difference sequence length.

3) The introduction of trainable parameters for relative positional encoding introduces additional flexibility during training.