

Assignment - 2

- 1) Perceptron calculates : $y = w_n + b$
 which is a straight line
 and XOR has output 0, 1.
 Hence, no straight line can separate
 0 and 1.

Multi-layer has hidden layer and activation functions. As these functions are non linear, multiple layers allow non linear decision making.

2) $y = f_1(f_2x) = (f_1f_2)x \rightarrow$ One linear transformation.

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a_n} \times \frac{\partial a_n}{\partial a_{n-1}} \times \dots \times \frac{\partial a_1}{\partial w}$$

During backpropagation, as multiple gradient terms gets multiplied continuously, it leads to exponential decay. Eg: Sigmoid outputs value b/w 0 and 1 so, it leads to vanishing gradient problem.

ReLU has no saturation, its gradient is equal to 1

- 3) Positional Embedding is necessary to keep the order of tokens correct which will make semantic meaning of sentence.

Sinusoidal PE has a fixed formula

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/12}}\right)$$

Absolute PE has fixed embeddings

Parity = what I am looking for

Key = what I offer

Value = what is actual content

$$\text{attention}(q, k, v) = \text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)v$$

→ Token is most similar to itself



5.) Each head gives attention to different relations.

Multi-head gives more broader and richer representation.

d_{model} = embedding size

h = number of heads

$d_{\text{head}} = \frac{d_{\text{model}}}{h}$

1

6.) Greedy decoding chooses locally best token.

→ Beam search keeps top-k partial sequences

→ It explores multiple futures.

Eg: let A → probability = 0.7

B → probability = 0.3

A → $P_{\text{end}} = 0.1$

, Greedy = 0.07

B → $P_{\text{end}} = 0.9$

, Beam search = 0.27

Greedy fails, Beam search succeeds

1)

$$(a) d_head = d_model \\ n$$

$$= \frac{768}{12} = 64$$

(b) Total parameters:

$$\text{Parameters} = 768 \times 768$$

$$= \underline{\underline{589,824}}$$

$$\text{For } \phi, k, v = 3 \times 589,824$$

$$= \underline{\underline{1,769,472}}$$

$$(2) \text{ Attention Scores} = [2.0, 1.0, 0.0]$$

$$\text{Softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$T_{\text{total}} = e^2 + e' + e^\circ = 11.107$$

$$\text{Normalise: } \left[\frac{e^2}{11.107}, \frac{e'}{11.107}, \frac{e^\circ}{11.107} \right]$$

$$= [0.665, 0.245, 0.09]$$