

Q1

XOR is a function which isn't linear.

A single layer perceptron isn't able to classify the non-linear func.
Therefore, by introducing multiple layers of perceptrons, (MLP) we remove the linearity of the first model, helping us to classify non-linear func.

Q2

A linear relation would be of type : $\hat{y} = \beta_0 + \beta_1 x$

Multiple layers of same type would result in:

$$\hat{y} = (\beta_0 + \beta'_0 + \beta''_0) + (\beta_1 + \beta'_1 + \beta''_1)x$$

$$= \beta_0 + \beta_1 x + \beta'_0 + \beta'_1 x + \dots$$

$$= \hat{y}_1 + \hat{y}_2 + \dots$$

↳ This is result of just addition

of multiple linear func.

∴ linear layers stacked on linear layers remain linear.

In deep neural networks, ~~extensive~~ chain rule is used extensively to calculate error (Backprop.)

$$\text{Eq: } \frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a_3} \frac{\partial a_3}{\partial a_2} \dots \frac{\partial a_i}{\partial w_i}$$

Here the partial derivatives are usually of the less than 1, which when multiplied lot of times, result in a very small value, resulting in gradient shrink.

Sigmoid form' takes the value and scales it down to $(0, 1)$ resulting in even smaller derivative.

Hence is $\text{mse}(0, n)$ which gives this derivative as 1, solving the problem of vanishing grad.

Q3

When transforms use self attention, the tokens aren't given some number to preserve the original sequence. It doesn't know about what the correct order is. To solve this problem, we use positional encoding.

Absolute PE

Gives unique vectors to each index

Initialized randomly and optimized by backprop.

Can only handle seq. of having length less than an assigned val.

RoPE helps with long-context as it shifts how the model perceive distance from 'absolute slots' to relative rotation.

Unlike ~~some other~~ methods that add vectors together, RoPE rotates the Query & Key vectors.

Sinusoidal PE

Uses sine, cos fun' of different frequency

No training happens. Fixed with formula are used

Can be used on any length of sequence

- Query (Q) : It represents the "current" word looking for information
 Key (K) : This represents a "label" for all words in the sequence.
 Value (V) : This is the actual info tied to the word.

dk is used for scaling as it gives us stability numerically.

When dk is too large, Q.K' is very big. Hence by dividing by dk var of

$Q \cdot K^T$ is ~~≈~~ around 1.

In attention weight matrix, the diagonal represents ^a the word attending itself. ~~because~~ the word relates most to itself, the diagonal values are highest.

Q5 If we use single attention head, the model can only focus on one type of relationship at a time for a specific word. By using multiple heads our model can attend ^{to} different types of relationships at the same time.

d_{model}: Total size of embedding vector for each word.

h: No. of parallel attention "workers".

d_{head}: Size of vector processed by each individual head.

Q6

Greedy decoding can ~~lead~~ lead the model "into a "local optima" trap. By choosing a very likely word now, the model may be forced ~~to~~ into highly unlikely words later to finish the sentence.

Beam search is better as it explores multiple branches/choices simultaneously. If one branch starts with a slightly low-prob. word but leads to a much more logical and high-prob. conclusion, beam search will preserve that path while greedy would've discarded ~~one~~ immediately.

	<u>1st Word</u>	<u>2nd Word</u>
1 st	$\text{Word 1} \rightarrow 0.7 \text{ prob}$	$\Rightarrow \text{Word 2} \rightarrow 0.001 \text{ prob.}$
2 nd	$\text{Word 1} \rightarrow 0.3 \text{ prob}$	$\Rightarrow \text{Word 2} \rightarrow 0.9 \text{ prob.}$

$$1^{\text{st}} \text{ path prob} = 0.007$$

$$2^{\text{nd}} \text{ path prob} = 0.27$$

Q1

$$\text{a) } d_{\text{model}} = d_{\text{model}} = 64$$

$$\text{b) Paravari per Matrix} = d_{\text{model}} = (768)^2$$

$$\text{for 3 matrix: } 3 d_{\text{model}} = 3 \cdot (768)^2$$

Q2

$$\frac{e^{x_i}}{\sum e^{x_i}} \quad e^2 = 7.39$$

$$e^1 = 2.71$$

$$e^0 = 1$$

~~$$\sum e^{x_i} = 11.1$$~~

$$G_1 = \frac{7.39}{11.1} = 0.6657$$

$$G_2 = \frac{2.71}{11.1} = 0.244$$

$$G_3 = \frac{1}{11.1} = 0.09$$