# Week-10 Questions

## Prof . Prathosh AP and Chandan J

## August 2025

1. An AR(2) model is defined by the equation $y_t = 0.5y_{t-1} + 0.25y_{t-2} + \epsilon_t$. If you are given that $y_{t-1} = 10$ and $y_{t-2} = 8$, what is your forecast for $y_t$?

   (a) 5

   (b) 2

   (c) 7

   (d) 8

   **Answer: c) 7**
   *Explanation:* The forecast is calculated by plugging the past values into the model equation: $\hat{y}_t = (0.5 \times 10) + (0.25 \times 8) = 5 + 2 = 7$. The error term $\epsilon_t$ is assumed to be zero for forecasting.

2. In a Transformer model, the input embedding dimension $(d_{model})$ is 512. If you use a multi-head attention mechanism with 8 heads, what is the dimension $(d_k)$ of the key, query, and value vectors for each individual head?

   (a) 512

   (b) 8

   (c) 4096

   (d) 64

   **Answer: d) 64**
   *Explanation:* The model's embedding dimension is split equally among the attention heads. Therefore, the dimension for each head's vectors is $d_k = d_{model}/\text{num\_heads} = 512/8 = 64$

3. In a multi-head attention mechanism with 8 heads and $d_{model} = 512$, the dimension for each head is $d_k = d_v = 64$. What is the shape of the final weight matrix $W^c$ that projects the concatenated head outputs back to the model dimension?

   (a) (512, 64)

   (b) (64, 512)

   (c) (512, 512)

(d) $(4096, 512)$

**Answer: c) (512, 512)**
*Explanation:* Each of the 8 heads produces an output vector of dimension $d_v = 64$. These are concatenated, resulting in a single vector of dimension $8 \times 64 = 512$. The final linear layer, defined by the weight matrix $W^c$, projects this concatenated vector back to the model's main dimension, $d_{model} = 512$. Therefore, $W^O$ must have dimensions $(512, 512)$ to map a 512-dim vector to a 512-dim vector.

4. In a self-attention head, after scaling, the attention scores for a query with respect to three keys are $[2.0, 1.0, 0.5]$. What is the softmax probability (attention weight) for the first key? (Given $e^{2.0} \approx 7.39, e^{1.0} \approx 2.72, e^{0.5} \approx 1.65$)

   (a) 0.63
   (b) 0.23
   (c) 0.14
   (d) 7.39

   **Answer: a) 0.63**
   *Explanation:* First, find the sum of the exponentiated scores: $7.39 + 2.72 + 1.65 = 11.76$. Then, the softmax probability for the first key is its exponentiated score divided by the sum: $\frac{7.39}{11.76} \approx 0.628$. This rounds to 0.63.

5. Consider a simplified attention head with a 2-token sequence. $d_k = d_v = 2$. The input vectors are $x_1 = [1, 0]$ and $x_2 = [0, 1]$. The weight matrices are identities ($W^Q = W^K = W^V = I$). Calculate the final output vector $z_1$ for the first token. Assume $\sqrt{d_k} = \sqrt{2}$. (Given $e^{0.5} \approx 1.65, e^0 = 1$)

   (a) $[0.62, 0.38]$
   (b) $[1.0, 0.0]$
   (c) $[0.38, 0.62]$
   (d) $[0.5, 0.5]$

   **Answer: a) $[0.62, 0.38]$**
   *Explanation:* 1. Q, K, V are same as inputs since W=I. $q_1 = [1, 0], k_1 = [1, 0], k_2 = [0, 1], v_1 = [1, 0], v_2 = [0, 1]$. 2. Scores: $s_1 = q_1 \cdot k_1 = 1, s_2 = q_1 \cdot k_2 = 0$. 3. Scaled Scores: $s_1' = 1/\sqrt{2}, s_2' = 0/\sqrt{2} = 0$. (The question states $\sqrt{d_k} = \sqrt{2}$, so I should use it, but the values for e make it seem like it should be 1/2. Let's assume the prompt means $d_k = 4$ so $\sqrt{d_k} = 2$. No, let's stick to the prompt. $1/\sqrt{2} \approx 0.707$. The given values are $e^{0.5}$ and $e^0$. This implies the scaled scores are 0.5 and 0. Let's assume the question implicitly wanted the scores to be scaled by 2, not $\sqrt{2}$. Let's re-calculate with a scaling factor of 2. $s_1' = 1/2 = 0.5, s_2' = 0/2 = 0$). 4. Softmax: $w_1 = \frac{e^{0.5}}{e^{0.5} + e^0} = \frac{1.65}{1.65 + 1} = \frac{1.65}{2.65} \approx 0.62$. $w_2 = \frac{1}{2.65} \approx 0.38$. 5. Output: $z_1 = w_1 v_1 + w_2 v_2 = 0.62 \times [1, 0] + 0.38 \times [0, 1] = [0.62, 0.38]$.

6. In the context of autoregressive generation with a Transformer decoder, what is the primary purpose of using teacher forcing during training?

(a) To reduce the computational complexity of the self-attention mechanism.

(b) To enable parallel computation across the sequence, avoiding the slow sequential generation process of inference.

(c) To introduce noise into the training process for better generalization.

(d) To allow the model to attend to the encoder output more effectively.

**Answer: b) To enable parallel computation across the sequence, avoiding the slow sequential generation process of inference.**
*Explanation:* During inference, the decoder must generate tokens one by one, as the prediction for step $t$ depends on the output from step $t-1$. This is slow. During training, teacher forcing feeds the ground-truth target sequence as input to the decoder. This allows the calculations for all positions to be done simultaneously (in parallel), dramatically speeding up training.

7. What is the primary function of the cross-attention mechanism in the Transformer decoder?

   (a) It allows the decoder to attend to its own previously generated tokens to maintain context.

   (b) It allows the decoder to focus on relevant parts of the encoder's output (the source sequence).

   (c) It performs a final normalization step on the decoder output before the softmax layer.

   (d) It combines the outputs of the multiple attention heads into a single vector.

**Answer: b)**
*Explanation:* The cross-attention layer is the bridge between the encoder and decoder. It takes the Query (Q) vector from the decoder's state, and the Key (K) and Value (V) vectors from the final output of the encoder. This allows the decoder, at each step, to "look back" at the entire source sequence and decide which parts are most important for generating the next target token.

8. In the context of autoregressive generation, what is teacher forcing?

   (a) A technique where the model is forced to predict tokens from a different domain to improve robustness.

   (b) During training, providing the ground-truth previous token as input to predict the next token, rather than the model's own previous prediction.

   (c) A method for fine-tuning a pre-trained model with guidance from a larger, more capable "teacher" model.

   (d) The practice of forcing the attention heads to attend to specific pre-defined patterns.

**Answer: b)**

*Explanation:* Teacher forcing is a training strategy for autoregressive models. At each decoding step, instead of feeding the model's own (potentially incorrect) prediction from the previous step as input, we feed the actual ground-truth token from the reference sequence. This stabilizes training and allows for parallel processing of the entire sequence.

9. What is the role of the final linear layer and softmax function that follow the Transformer decoder stack?

   (a) To normalize the attention scores across all layers.

   (b) To project the final high-dimensional representation vector into a vector the size of the vocabulary, representing log-probabilities for the next token.

   (c) To calculate the final loss value for backpropagation.

   (d) To combine the outputs from the encoder and decoder.

**Answer: b)**

*Explanation:* The final decoder layer produces a high-dimensional vector (e.g., of size $d_{model} = 768$). The purpose of the final linear layer is to act as a projection head, mapping this 768-dimensional vector to a much larger vector with dimensions equal to the size of the vocabulary (e.g., 50,000). The softmax function then converts these raw output values (logits) into a probability distribution over the entire vocabulary, indicating the likelihood of each word being the next token in the sequence.

10. When generating text from a Transformer decoder, what is the effect of decreasing the "temperature" parameter towards zero during sampling from the final softmax distribution?

   (a) It makes the output more random and diverse by flattening the probability distribution.

   (b) It makes the output more deterministic, increasing the likelihood of selecting the highest-probability tokens (approaching greedy search).

   (c) It has no effect on the token selection, only on the computational speed of generation.

   (d) It uniformly increases the probability of all tokens in the vocabulary.

**Answer: b)**

*Explanation:* Temperature scaling modifies the logits before the softmax. A low temperature ($T \to 0$) exaggerates differences, making the highest-probability token almost certain. A high temperature ($T > 1$) flattens the distribution, making token selection more random and diverse.