# TRANSFORMER REVIEW

Decoder

→ Linear layer and softmax

→ Runs multiple times to produce a predicted word at each step.

FC

Add & norm

Positionwise FFN

$\in \mathbb{R}^{b \times n \times d}$

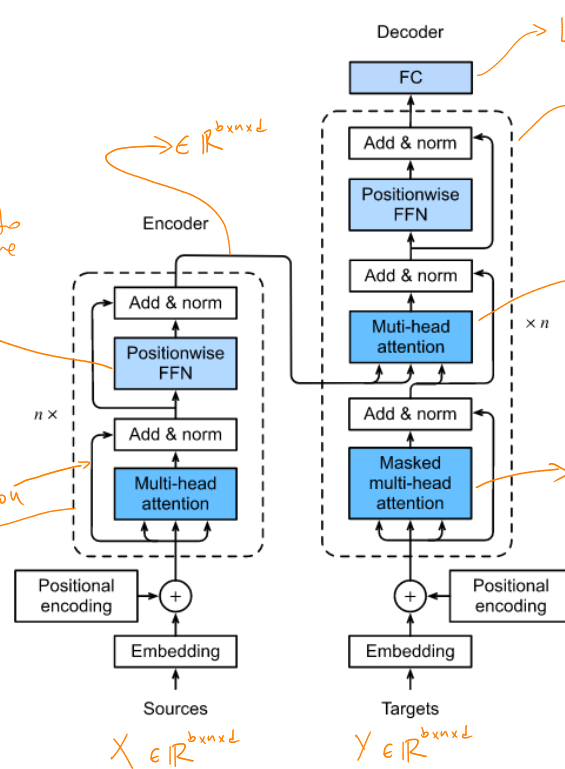Add & norm

Encoder

It seems this learns complex interactions of words.
- while wive
- keep up
- get angry

Applies transformations to the representation of the sequences. It is an MLP?

Add & norm

Positionwise FFN

$n \times$

Add & norm

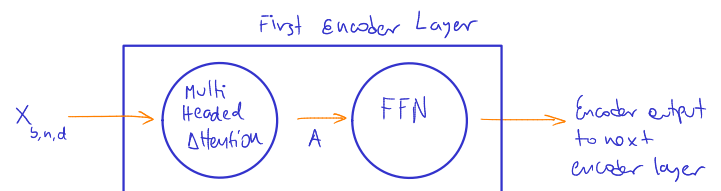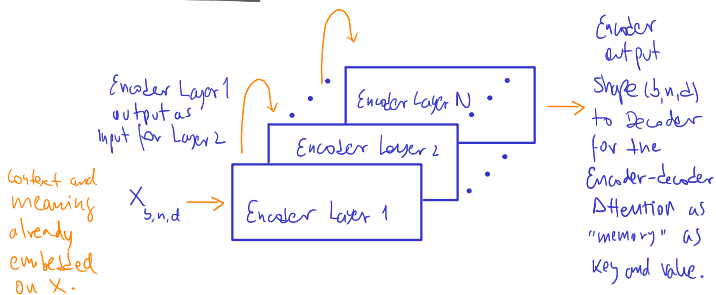Multi-head attention

Residual Connection

The encoder run once per batch

Multi-head attention

Inputs:
- x from prev. decoder layer (query)
- m (memory) from the output of the encoder. Twice as key and value.

Encoder-decoder attention

$\times n$

Add & norm

Masked multi-head attention

Queries, keys, values from outputs of prev. decoder layer except the first layer which is the target

Positional encoding  $+$

Positional encoding  $+$

Embedding

Embedding

Sources

Targets

$X \in \mathbb{R}^{b \times n \times d}$

$Y \in \mathbb{R}^{b \times n \times d}$

## LEGEND:

b: Batch size

n: sequence length

d: Embedding dimensionality

h: Number of heads

## THE ENCODER

Encoder Layer 1 outputs as Input for Layer 2

Context and meaning already embedded on X.

$X_{b,n,d}$

Encoder Layer 1

Encoder Layer 2

Encoder Layer N

Encoder output Shape (b,n,d) to Decoder for the Encoder-decoder Attention as "memory" as Key and Value.

First Encoder Layer

$X_{b,n,d}$

Multi Headed Attention  →  A  →  FFN

Encoder output to next encoder layer

## MULTIHEADED SELF-ATTENTION

Different on each encoder layer $\rightarrow$ $W_{d,d}^Q$

Note: We can also use one nn.Linear $(d, 3d)$ and then separate the output in three parts to get $Q$, $K$, and $V$

$X_{b,n,d}$

LINEAR → $Q_{b,n,d}$ —Reshape→ $Q_{b,n,h,d/h}$ —Swap dimensions→ $Q_{b,h,n,d/h}$

split into $h$ heads with embed. size of $d/h$

$W_{d,d}^K$

LINEAR → $K_{b,n,d}$ —Reshape→ $K_{b,n,h,d/h}$ —Swap dimensions→ $K_{b,h,n,d/h}$

$W_{d,d}^V$

LINEAR → $V_{b,n,d}$ —Reshape→ $V_{b,n,h,d/h}$ —Swap dimensions→ $V_{b,h,n,d/h}$

$\text{Attention}(Q,K,V)_{b,h,n,d/h} = A$

—Swap dimensions→ $A_{b,n,h,d/h}$

↓ Reshape (concat)

$A_{b,n,d}$

$W_{d,d}^A$ → LINEAR

↓ $A_{b,n,d}$ — Multi headed Attention output

We are feeding a tensor to a linear layer ⇒ In PyTorch the linear layer is applied only to the last dimension.

$A_{b,n,d}$

FFN: $\text{ReLU}(AW_1 + b_1)W_2 + b_2$

$W_{1_{d,f}}$ → LINEAR

↓ $A'$

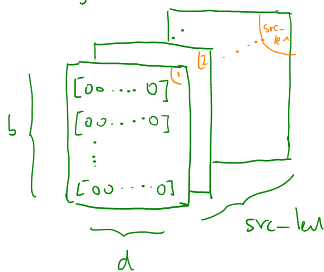$\text{ReLU}(A'_{b,n,d})$

$W_{2_{f,d}}$ → LINEAR

## A CLOSER LOOK TO THE ATTENTION MECHANISM

This will replace X as input in next encoder layer ⟨ Encoder output to next encoder layer : $(b,n,d)$
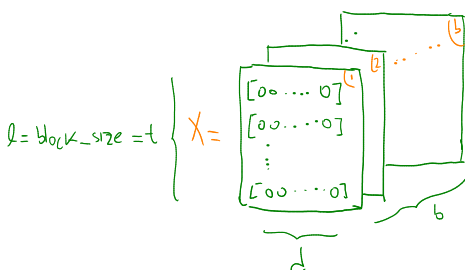
### THE INPUTS TO ATTENTION:

In RNNs we were interested in processing batches of the $i$-th word of each sentence at each timestep. Something like this:

where:
- $b \rightarrow$ A batch of all $i$-th words of all sentences
- $d \rightarrow$ The embedding size for each word
- $src\_len \rightarrow$ The len of the longest sentence in the input dataset.

$b \left\{ \begin{matrix} [00 \cdots 0] \\ [00 \cdots 0] \\ \vdots \\ [00 \cdots 0] \end{matrix} \right.$

src-len

$d$

With transformers we can pass a batch of entire sentences, so the inputs might look like this:

where:
- $b$: The # of sentences in the input dataset
- $\ell$: The # of words of sentence $i$ of the dataset
- $d$: The dimension of each word embedding in the sentence

$\ell = \text{block\_size} = t \left\{ X = \begin{matrix} [00 \cdots 0] \\ [00 \cdots 0] \\ \vdots \\ [00 \cdots 0] \end{matrix} \right.$

$b$

$d$

rows, each word

This is a tensor of shape $(b, \ell, d)$ when applying softmax over this tensor we do it row-wise, meaning, the second dimension

When transposing this tensors we need to forget the batches for a moment and just think about the dimensions independently of the batch dimension. It is not a good idea using the T operand in PyTorch since

$$(b, l, d)^T \rightarrow (d, l, b).$$ Better $(b, l, d).\text{transpose}(1,2) \Rightarrow (b, d, l)$ So we can do @ with torch.bmm:

switch

a single token?

$(b, l, d) \times (b, d, l) \rightarrow (b, l, l)$

bridge

As we are working with self attention each input vector $x_i$ assumes 3 roles:

- keys
- Queries    To make things easier we'll define 3 new vectors as linear
- Values     transformations of $x_i$ through 3 $d \times d$ matrices: $W_k, W_q, W_v$

$$q_i = W_q x_i \qquad K_i = W_k x_i \qquad V_i = W_r x_i$$
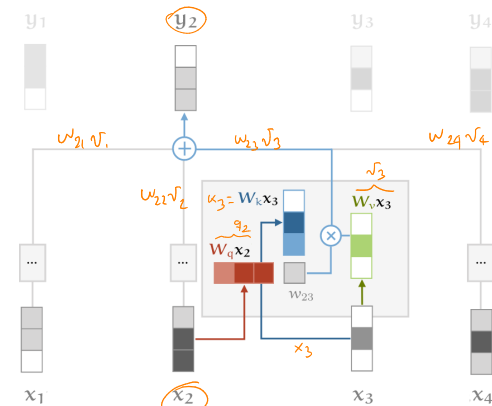
Shapes:
$(d,1) \quad (d,d) \quad (d,1)$
$K_i = W_k x_i$

**Attention Layer**

Attention output : $y_i = \sum_j w_{ij} \vec{v}_j$ $\longrightarrow$ $y_i$ of shape $(e \times 1)$

Attention distribution: $w_{ij} = \text{softmax}(w'_{ij})$

Attention scores: $w'_{ij} = q_i^T K_j$

$x_i$ as $W_q x_i = q_i$ of shape $(d \times 1)$

To avoid gradient explosion we clip the attention scores

$$w'_{ij} = \frac{q_i^T K_j}{\sqrt{e}}$$



$$w'_{23} = q_2^T K_3 = W_q x_2^T W_k x_3$$

$$w_{23} = \text{softmax}(W_q x_2^T W_k x_3)$$

$$y_2 = w_{21} \vec{v}_1 + w_{22} \vec{v}_2 + w_{23} \vec{v}_3 + w_{24} \vec{v}_4$$
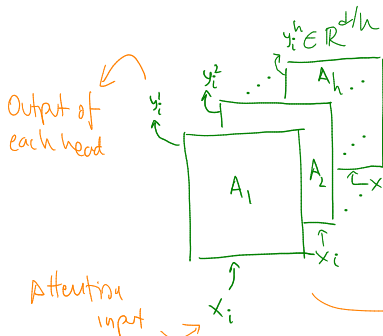
Scalar          vector

As human communication is complex we can't not capture the entire meaning with one attention pass. We need to stack more layers! $\Rightarrow$ Multi-headed Attention.

# MULTI - HEADED SELF ATTENTION

We actually split the d-dim. inputs in d/h heads

We just "stack" several Self-Attention mechanisms: we'll have a set of weight matrices per layer

$\Rightarrow$ $W_q^h \quad W_k^h \quad W_v^h$ where h is the # of layers.

Attention Heads

Output of each head

$y_i^h \in \mathbb{R}^{d/h}$

$y_i^2 \cdots A_h \cdots$
$A_1 \qquad A_2 \quad \uparrow x_i$
$\uparrow x_i$

concatenation
$\Rightarrow [y_i^1 \ y_i^2 \cdots y_i^h] \in \mathbb{R}^d$
$\rightarrow \boxed{\text{LINEAR TRANSF.}} \rightarrow y_i \in \mathbb{R}^d$
$W_{d,d}$

Attention input $\rightarrow x_i$

$X_i \in \mathbb{R}^d$

$X_{l \times d} = \begin{bmatrix} x_1 \\ \vdots \\ x_l \end{bmatrix}$

This is expensive $\Rightarrow$ It is better to:

Not sure about this representation

$W_q^1 \ W_q^2 \cdots W_q^h \Rightarrow W_q \in \mathbb{R}^{d \times d}$

$W_k^1 \ W_k^2 \cdots W_k^h \Rightarrow W_k \in \mathbb{R}^{d \times d}$

$W_v^1 \ W_v^2 \cdots W_v^h \Rightarrow W_v \in \mathbb{R}^{d \times d}$

concatenation     Just 3 matrices

Better ways to handle weights for multi headed Attention

$d$

$X_i = [0\ 0\ 0 \text{---------------------------------} 0]$

INPUT CHUNKS $\rightarrow$ Input for $A_1$ \quad Input for $A_2$ $\cdots$ Input for $A_h$
$[0 0 \cdots 0] \quad [0 0 \cdots 0] \cdots [0 0 \cdots 0] \in \mathbb{R}^{d/h}$
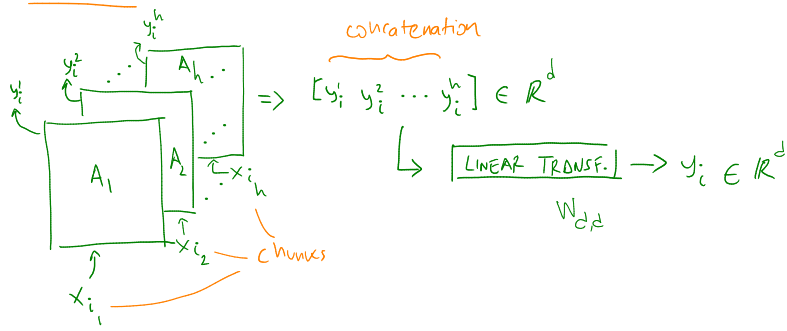
Each set contains $\leftarrow$ $W_q^1 \ W_k^1 \ W_v^1 \quad W_q^2 \ W_k^2 \ W_v^2 \cdots W_q^h \ W_k^h \ W_v^h \in \mathbb{R}^{d/h \times d/h}$
3 matrices.

$h$

$\rightarrow$ # of Attention heads

Then we have:



$$\Rightarrow \quad [y_i^1 \; y_i^2 \; \cdots \; y_i^h] \in \mathbb{R}^d$$

concatenation

$$\hookrightarrow \boxed{\text{LINEAR TRANSF.}} \longrightarrow y_i \in \mathbb{R}^d$$

$$W_{d,d}$$

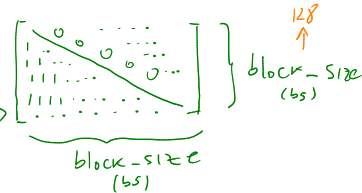$$x_i = [x_{i_1} \; x_{i_2} \; \cdots \; x_{i_h}] \; ; \; x_i \in \mathbb{R}^d$$

## CAUSAL MASK

A vanilla multi-head masked self attention layer with a projection at the end. This can also be achieved with torch.nn.MultiheadAttention()

Key, query, value $\in \mathbb{R}^{d \times d}$ where $d$ = embedding size. $\Rightarrow$ k, q, v Projections for all heads: nn.Linear (d, d) $\rightarrow 256$

Regularization layers $\rightarrow$ Attention dropout is nn.Dropout (attention_dropout) $\rightarrow 0.1$
$\hookrightarrow$ Residual connection dropout : nn.Dropout (residual_dropout)
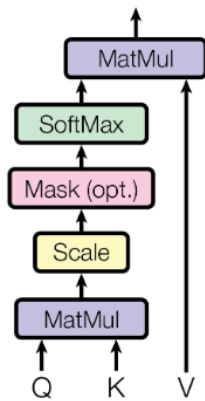
Output projection $\rightarrow$ nn.Linear (d,d)

Mask to hide "future" input for attention (right part) $\rightarrow$



$\}$ block_size (bs) $\leftarrow 128$
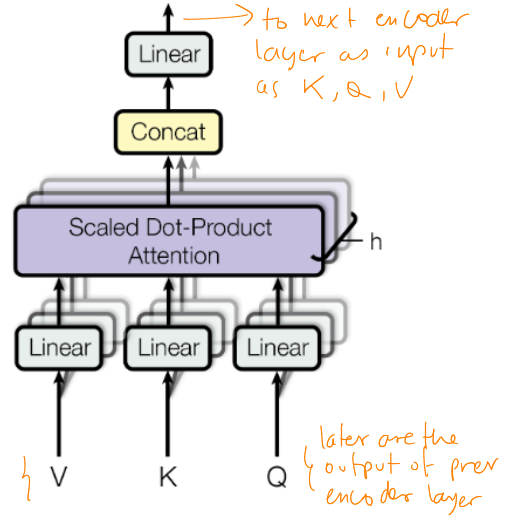
{ But we need a higher dimension tensor of shape (1,1,bs,bs) with the view function

$\underbrace{\phantom{xxxxxx}}_{\substack{\text{block-size} \\ (bs)}}$

n_head $\Rightarrow 8$

## Scaled Dot-Product Attention



Q  K  V

## Multi-Head Attention



V  K  Q

*(handwritten annotations on Multi-Head Attention figure)*

→ to next encoder layer as input as $K, Q, V$

— h

*(handwritten annotation near V on left figure)* on first encoder layer these are $X$

*(handwritten annotation near Q)* later are the output of prev encoder layer

*(handwritten title)* **CS224N – ASSIGNMENT 5**

---

## 1. Attention exploration (21 points)

(a) (2 points) **Copying in attention:** Recall that attention can be viewed as an operation on a query $q \in \mathbb{R}^d$, a set of value vectors $\{v_1, \ldots, v_n\}, v_i \in \mathbb{R}^d$, and a set of key vectors $\{k_1, \ldots, k_n\}, k_i \in \mathbb{R}^d$, specified as follows:

*(handwritten)* → attention output (vector)

$$c = \sum_{i=1}^{n} v_i \alpha_i \qquad (1)$$

*(handwritten)* $\alpha$ is attention distribution (scalar)

$$\alpha_i = \frac{\exp(k_i^\top q)}{\sum_{j=1}^{n} \exp(k_j^\top q)}. \qquad (2)$$

*(handwritten)* at time step $t$

where $\alpha_i$ are frequently called the "attention weights", and the output $c \in \mathbb{R}^d$ is a correspondingly weighted average over the value vectors.

We'll first show that it's particularly simple for attention to "copy" a value vector to the output $c$. Describe (in one sentence) what properties of the inputs to the attention operation would result in the output $c$ being approximately equal to $v_j$ for some $j \in \{1, \ldots, n\}$. Specifically, what must be true about the query $q$, the values $\{v_1, \ldots, v_n\}$ and/or the keys $\{k_1, \ldots, k_n\}$?

*(handwritten right margin)*

$q_t \in \mathbb{R}^d, \quad K_i \in \mathbb{R}^d, \quad v_i \in \mathbb{R}^d$

$i \in \{1, \ldots, n\}$

↳ # hidden cells in encoder

All calculations are done at each time step for all $i$.

$\frac{x \cdot}{c} \Rightarrow \left(\frac{x}{c}\right)$

$C \sim \Delta$ → num is muy pequeño $\Rightarrow \alpha_i \sim 0$

$C \sim B$ → num is muy grande $\Rightarrow \alpha_i \gg 1$

$\frac{x \cdot}{c} \Rightarrow \left(\frac{x}{c}\right) \sim 0$

*(handwritten left working)*

num $= \alpha_i^\top q$ ; $K_i^\top = [0\ 0 \cdots 0]$ ... $q = \begin{bmatrix} \Delta \\ \Delta \\ \vdots \\ \Delta \end{bmatrix} \} d$

num $= \exp(0\Delta + 0\Delta_2 + \cdots 0\Delta_d)$

$\underbrace{\qquad}_{\ll 0 \Rightarrow}$ num is muy pequeño

Esto significa $q \cdot K$ ó $q$ tienen sus vectores con componentes negativos

$K_i \uparrow$ and $q \downarrow \Rightarrow k_i \cdot q$ is negative $\quad [-1, -2, -3] \cdot \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = -14 \}$ This is the extreme case where $K_i = -q$, the absolute opposite

Mientras $K_i$ and $q$ sean más opuestos $\Rightarrow$ num $\ll 0 \Rightarrow \alpha_i \sim 0$

→ softmax$(k_i^\top q)$

*(number line diagram)* $n-q$ ... $0$ ... $1$ ... $2$ ... $3$
$\alpha_0 \sim 0 \quad \alpha_1 \sim 0 \quad \alpha_2 \sim 1 \quad \alpha_3 \sim 0$

$\Rightarrow \quad C = v_0 \alpha_0 + v_1 \alpha_1 + v_2 \alpha_2 + v_3 \alpha_3$

$\Rightarrow C \sim v_2$

The attention output will be a copy of one of the inputs when the all the keys except one of them oppose to the query.

$k_j^\top q < 0 ; \ j \neq i \qquad$ → Keys $K_j$ opposes to the query (low affinity between the query and keys)

**AND** $\quad k_i^\top q \gg k_j^\top q \qquad$ → Key $k_i$ is less opposed to the query (more affinity between $q$ and $k_i$)

(b) (4 points) **An average of two:** Consider a set of key vectors $\{k_1, \ldots, k_n\}$ where all key vectors are perpendicular, that is $k_i \perp k_j$ for all $i \neq j$. Let $\|k_i\| = 1$ for all $i$. Let $\{v_1, \ldots, v_n\}$ be a set of arbitrary value vectors. Let $v_a, v_b \in \{v_1, \ldots, v_n\}$ be two of the value vectors. Give an expression for a query vector $q$ such that the output $c$ is approximately equal to the average of $v_a$ and $v_b$, that is, $\frac{1}{2}(v_a + v_b)$.[1] Note that you can reference the corresponding key vector of $v_a$ and $v_b$ as $k_a$ and $k_b$.

when $\{k_1, k_2, k_3\}$    $n = 3$

$k_1 \perp k_2 \Rightarrow k_1 \cdot k_2 = 0$   $k_1 \perp k_3$   $k_2 \perp k_3$

$\{v_1, v_2, v_3\}$   $\|k_1\| = \|k_2\| = \|k_3\| = 1$

$k_1^T = [1, 0, 0]$
$k_2^T = [0, 1, 0]$
$k_3^T = [0, 0, 1]$

$\overrightarrow{v_{a=1}}, \overrightarrow{v_{b=3}}$

$q = ?$ such as $\boxed{c \sim \dfrac{v_a + v_b}{2}}$

$\vec{p} = p_1 \vec{i} + p_2 \vec{j} + p_3 \vec{k}$

$\|\vec{p}\| = \sqrt{p_1^2 + p_2^2 + p_3^2}$

To make this happen we need that

$\alpha_a = \alpha_b \approx \frac{1}{2}$ , and $\alpha_i \approx 0$ for $i \notin \{a, b\}$ ← no affinity between $q$ and $k_i$ !

$\Rightarrow$ softmax $(k_a^T q)$ = softmax $(k_b^T q) \approx \frac{1}{2}$

$\Rightarrow k_a^T q = k_b^T q$

[1]Hint: while the softmax function will never *exactly* average the two vectors, you can get close by using a large scalar multiple in the expression.

Taking the hint into account : Multiple of what? I guess the key vectors since $q$ interacts with them.

$\Rightarrow q = S(k_a + k_b)$

where $S$ is a big scalar

Now, $\alpha_i$ = softmax $(k_i^T q) = \dfrac{e^{k_i^T q}}{\sum e^{k_j^T q}} \Rightarrow \dfrac{e^{k_i^T S(k_a + k_b)}}{\sum_{i=1}^{n} e^{k_i^T S(k_a + k_b)}} = \dfrac{e^{S k_a k_i^T + S k_b k_i^T}}{\sum_{j=1}^{n} e^{S k_a k_j^T + S k_b k_j^T}}$

because $k_a k_a^T = 1$ and $k_b k_a^T = 0$

For $i \in \{a, b\}$   $\alpha_a = \dfrac{e^S}{2e^S + (n-2) \times 1} = \alpha_b \approx 0.5 = 0.5$ if $S = \infty$

$\quad i \in \{a, b\}$    $j \notin \{a, b\}$

$j \in \{a, b\}$ : $e^S$     $j \notin \{a, b\}$ : $1$ because the key vectors are orthogonal then if $j \notin \{a, b\}$ $S k_i k_j^T = 0$ for $i \in \{a, b\}$ and $e^0 = 1$

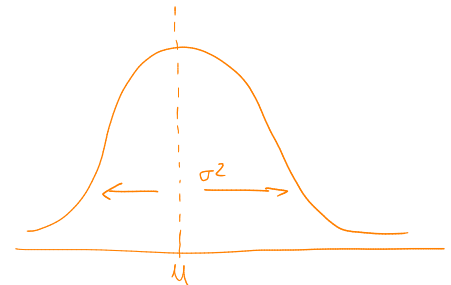$c \approx$ 

For $i \notin \{a, b\}$ $\alpha_i \approx 0$ if there is no affinity between the query and the keys to the values for $i$ in these cases. This means $q$ and $k$ are opposite.

(c) (5 points) **Drawbacks of single-headed attention:** In the previous part, we saw how it was *possible* for a single-headed attention to focus equally on two values. The same concept could easily be extended to any subset of values. In this question we'll see why it's not a *practical* solution. Consider a set of key vectors $\{k_1, \ldots, k_n\}$ that are now randomly sampled, $k_i \sim \mathcal{N}(\mu_i, \Sigma_i)$, where the means $\mu_i$ are known to you, but the covariances $\Sigma_i$ are unknown. Further, assume that the means $\mu_i$ are all perpendicular; $\mu_i^T \mu_j = 0$ if $i \neq j$, and unit norm, $\|\mu_i\| = 1$.

i. (2 points) Assume that the covariance matrices are $\Sigma_i = \alpha I$, for vanishingly small $\alpha$. Design a query $q$ in terms of the $\mu_i$ such that as before, $c \approx \frac{1}{2}(v_a + v_b)$, and provide a brief argument as to why it works.



$\Sigma_i = \begin{bmatrix} \alpha & 0 & 0 & \cdots & 0 \\ 0 & \alpha & 0 & \cdots & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \cdots & 0 & \alpha \end{bmatrix}$ vanishing small $\alpha \Rightarrow$ The variance $\approx 0 \Rightarrow k_i \approx \mu_i$ $\therefore q = S(\mu_a + \mu_b)$

where $S$ is a big scalar.

ii. (3 points) Though single-headed attention is resistant to small perturbations in the keys, some types of larger perturbations may pose a bigger issue. Specifically, in some cases, one key vector $k_a$ may be larger or smaller in norm than the others, while still pointing in the same direction as $\pi_a$. As an example, let us consider a covariance for item $a$ as $\Sigma_a = \alpha I + \frac{1}{2}(\mu_a \mu_a^T)$ for vanishingly small $\alpha$ (as shown in figure 1). Further, let $\Sigma_i = \alpha I$ for all $i \neq a$. When you sample $\{k_1, \ldots, k_n\}$ multiple times, and use the $q$ vector that you defined in part i., what qualitatively do you expect the vector $c$ will look like for different samples?

$\mu_a \mu_a^T = 1$

Not sure about this but I'd say sampling with $q$ as defined in the part i will → output big results comparing with the cases where $i \neq a$. Why? Because the variance $\Sigma_a$ is greater than $\Sigma_i$. Then the $c$ vector will replicate the large perturbations observed in the keys.

(d) (3 points) **Benefits of multi-headed attention:** Now we'll see some of the power of multi-headed attention. We'll consider a simple version of multi-headed attention which is identical to single-headed self-attention as we've presented it in this homework, except two query vectors ($q_1$ and $q_2$) are defined, which leads to a pair of vectors ($c_1$ and $c_2$), each the output of single-headed attention given its respective query vector. The final output of the multi-headed attention is their average, $\frac{1}{2}(c_1+c_2)$. As in question 1(c), consider a set of key vectors $\{k_1,\ldots,k_n\}$ that are randomly sampled, $k_i \sim \mathcal{N}(\mu_i, \Sigma_i)$, where the means $\mu_i$ are known to you, but the covariances $\Sigma_i$ are unknown. Also as before, assume that the means $\mu_i$ are mutually orthogonal; $\mu_i^\top \mu_j = 0$ if $i \neq j$, and unit norm, $\|\mu_i\| = 1$.

  i. (1 point) Assume that the covariance matrices are $\Sigma_i = \alpha I$, for vanishingly small $\alpha$. Design $q_1$ and $q_2$ such that $c$ is approximately equal to $\frac{1}{2}(v_a + v_b)$.

For this to happen $c_1 = v_a$ and $c_2 = v_b$ and as per question a) that will be true if there is equal affinity between the query and the keys $k_a$ and $k_b$; and no affinity with $k_i$ for $i \notin \{a,b\}$.

Now taking question b) into account we can say that $q_1 \approx S k_a$ and $q_2 \approx S k_b$. Finally, as $\Sigma_i = \alpha I$ for a vanishing small $\alpha$ we know that the variance $\sigma_i^2 \approx 0$ ∴

$$q_1 = S \mu_a , \quad \text{where } S \text{ is a big scalar.}$$
$$q_2 = S \mu_b$$

  ii. (2 points) Assume that the covariance matrices are $\Sigma_a = \alpha I + \frac{1}{2}(\mu_a \mu_a^\top)$ for vanishingly small $\alpha$, and $\Sigma_i = \alpha I$ for all $i \neq a$. Take the query vectors $q_1$ and $q_2$ that you designed in part i. What, qualitatively, do you expect the output $c$ to look like across different samples of the key vectors? Please briefly explain why. You can ignore cases in which $q_i^\top k_a < 0$.

$$C = \frac{c_1 + c_2}{2} \quad ; \quad \sigma_a^2 > \sigma_{ita}^2$$

If we have the same $q_1$ and $q_2$ then $c$ should be $\approx (v_a + v_b)/2$ then when sampling keys there would be a tendency of having bigger norm for the $c$ vector around $v_a$, but as we are averaging $c_1$ and $c_2$ the perturbations might be reduced.

(e) (7 points) **Key-Query-Value self-attention in neural networks:** So far, we've discussed attention as a function on a set of key vectors, a set of value vectors, and a query vector. In Transformers, we perform *self-attention*, which roughly means that we draw the keys, values, and queries from the same data. More precisely, let $\{x_1,\ldots,x_n\}$ be a sequence of vectors in $\mathbb{R}^d$. Think of each $x_i$ as representing word $i$ in a sentence. One form of self-attention defines keys, queries, and values as follows. Let $V, K, Q \in \mathbb{R}^{d\times d}$ be parameter matrices. Then

$$v_i = V x_i \quad i \in \{1,\ldots,n\} \tag{3}$$
$$k_i = K x_i \quad i \in \{1,\ldots,n\} \tag{4}$$
$$q_i = Q x_i \quad i \in \{1,\ldots,n\} \tag{5}$$

Then we get a context vector for each input $i$; we have $c_i = \sum_{j=1}^{n} \alpha_{ij} v_j$, where $\alpha_{ij}$ is defined as $\alpha_{ij} = \frac{\exp(k_j^\top q_i)}{\sum_{\ell=1}^{n} \exp(k_\ell^\top q_i)}$. Note that this is single-headed self-attention.

In this question, we'll show how key-value-query attention like this allows the network to use different aspects of the input vectors $x_i$ in how it defines keys, queries, and values. Intuitively, this allows networks to choose different aspects of $x_i$ to be the "content" (value vector) versus what it uses to determine "where to look" for content (keys and queries.)

  i. (3 points) First, consider if we didn't have key-query-value attention. For keys, queries, and values we'll just use $x_i$; that is, $v_i = q_i = k_i = x_i$. We'll consider a specific set of $x_i$. In particular, let $u_a, u_b, u_c, u_d$ be mutually orthogonal vectors in $\mathbb{R}^d$, each with equal norm $\|u_a\| = \|u_b\| = \|u_c\| = \|u_d\| = \beta$, where $\beta$ is very large. Now, let our $x_i$ be:

$$x_1 = u_d + u_b \tag{6}$$
$$x_2 = u_a \tag{7}$$
$$x_3 = u_c + u_b \tag{8}$$

If we perform self-attention with these vectors, what vector does $c_2$ approximate? Would it be possible for $c_2$ to approximate $u_b$ by adding either $u_d$ or $u_c$ to $x_2$? Explain why or why not (either math or English is fine).

$$c_2 = \sum_{j=1}^{n} \alpha_{2j} v_j ; \quad \text{where}$$
$$\alpha_{2j} = \text{Softmax}(K_j^\top q_2)$$
↳ the keys showing more affinity with $q_2$

A vector highlighting the characteristics of the input values with higher affinity w.r.t. $q_2$

Since $v_i = q_i = k_i = x_i$ then

$$c_2 = \sum_{j=1}^{n} \text{softmax}(x_j^\top x_2) x_j = \frac{x_1 e^{x_1^\top x_2} + x_2 e^{x_2^\top x_2} + x_3 e^{x_3^\top x_2}}{e^{x_1^\top x_2} + e^{x_2^\top x_2} + e^{x_3^\top x_2}}$$
$$\Downarrow$$
$$c_2 = x_1 \text{softmax}(x_1^\top x_2) + x_2 \text{softmax}(x_2^\top x_2)$$
$$+ x_3 \text{softmax}(x_3^\top x_2)$$
↙

Let's take a look to $x_1^\top x_2$. $x_1^\top x_2 = (u_d + u_b)^\top u_a$. Since the dot product is distributive with addition:

$$\underbrace{u_a u_d^\top + u_a u_b^\top}$$ As $u_a, u_b$ and $u_d$ are orthogonal amongst them the dot product is $0$. Same happens with $x_3^\top x_2$.

Then we have: $c_2 \approx x_2 \text{softmax}(x_2^\top x_2) \approx \frac{x_2 e^{x_2^\top x_2}}{e^{x_1^\top x_2} + e^{x_2^\top x_2} + e^{x_3^\top x_2}} \approx x_2 \frac{e^{x_2^\top x}}{2 + e^{x_2^\top x_2}} \approx x_2 \approx u_a$ //

Notice: $\text{softmax}(x_1^\top x_2) = \frac{e^0}{2 + e^{x_2^\top x_2}} = \frac{1}{2 + e^{\beta^2}} \approx 0$

**Part g)** Synthesizer Attention

$X_{(l,d)}$

Linear transf. with $W_1$

$\rightarrow W_2$ in the implementation code. Not a linear transf.
The same in the implementation code

$Y_i = \text{softmax}(\text{ReLU}(XA_i + b_1)B_i + b_2)(XV_i) \rightarrow$ A linear transformation with self.value()

- $l \times d$   $d \times d$   $d/h \times l$   $l \times d$   $d \times d$

LINEAR  $A_i$ $(d,d)$

$X'_{(l,d)}$ but we need $(l, d/h)$ not to lose data we reshape as $(l, h, d/h)$

$l \times d = l \times h \times d/h$   $d/h \times l$   $l \times d = l \times h \times d/h$

RELU

$X''_{(h, l, d/h)}$  — swap —

softmax ( $l \times l$ )   $l \times d/h$

@ ← $B_i$ $(d/h, l)$ : $B_i$ $(1, d/h, l)$ reshape

$X'''_{(h, l, l)}$   $l \times l$   $l \times d/h$

+ ← $b_2$ $(l)$ $[\cdots, :l]$

$X'''_{(h, l, l)}$   $l \times l/h$

Softmax

$X^{iv}_{(h, l, l)}$

Then the output of the synthesizer mechanism is a linear transformation of the concatenation of each Attention head output:

@ ← $(h, l, d/h)$ ← $X_{(l,d)}$   LINEAR

$Y = [Y_1, Y_2 \cdots Y_n]W$
$\rightarrow d \times d$

$Y_i = X^v_{(h, l, d/h)}$   $V_{i\,(d,d)}$

$l \times \frac{d}{h}$ times $n \Rightarrow l \times d$

$l \times d \rightarrow$ Same shape as the input $X$!

---

ii. (2 points) Why might the *synthesizer* self-attention not be able to do, in a single layer, what the key-query-value self-attention can do?

I'd say it is because the synthesizer is not inspecting the potential relationships between each possible combination in the input. It is basically transforming the input into another representation and comparing it with the value representation. It lacks the power of the $Q, K, V$ system.

## 3. Considerations in pretrained knowledge (5 points)

(a) (1 point) Succinctly explain why the pretrained (vanilla) model was able to achieve an accuracy of above 10%, whereas the non-pretrained model was not.

Because the Vanilla model w/o pretraining was trained as a language model without an "understanding" of the complex relationships in the language (person $\xrightarrow{\text{born}}$ place). The pretraining activity was able to capture some of that language dependencies boosting the final objective thanks to transfer learning. Also the pretraining task had better chances to learn helpful dependencies for the downstream tasks because its dataset was bigger (68%) than the finetuning dataset

(b) (2 points) Take a look at some of the correct predictions of the pretrain+finetuned vanilla model, as well as some of the errors. We think you'll find that it's impossible to tell, just looking at the output, whether the model *retrieved* the correct birth place, or *made up* an incorrect birth place. Consider the implications of this for user-facing systems that involve pretrained NLP components. Come up with two reasons why this indeterminacy of model behavior may cause concern for such applications.

I think the model has learned to generate a place name after a question like "Where was xxx born?" When it gives a correct answer is probably because a relationship between that person name and a place was captured and embedded during pretraining. Otherwise outputs any location name. The implications are big, focus on trust and confidence of the predictions.

Low confidence → can't be applied to critical processes
Poor prediction → Low usability rate. Undesirable consequences when acting blindly based on erroneous model output.

(c) (2 points) If your model didn't see a person's name at pretraining time, and that person was not seen at fine-tuning time either, it is not possible for it to have "learned" where they lived. Yet, your model will produce *something* as a predicted birth place for that person's name if asked. Concisely describe a strategy your model might take for predicting a birth place for that person's name, and one reason why this should cause concern for the use of such applications.

If the prediction score is below a threshold the model can query a database to retrieve the real answer. This of course defeats the purpose of the model and might wrongly bias the user about the real power of AI. Perhaps it is better to take advantage of this situation to interact with the user asking for the missing info and/or creating a new datapair from the database for later finetuning.