

# [CS224N Winter 2021] Assignment 5: Self-Attention, Transformers, and Pretraining

## 1. Attention exploration (21 points)

### 1. Attention Exploration

A. Suppose we want  $c = v_j$  for some  $j \in \{1, \dots, n\}$ . This means that  $v_j = \sum_{i=1}^n v_i \alpha_i$  which means

that  $v_j$  needs to be a linear combination of the value vectors  $\{v_1, \dots, v_n\} = V$ . The simplest case of this is when  $\alpha_j = 1$  and  $\alpha_i = 0$  for  $i \neq j$ . This is approximately achieved when

$k_j^T q \gg k_i^T q$  for  $i \neq j$ . This will often be the case when  $q$  is "similar" to  $k_j$  (both having large vector entries in the same direction):

$$v_{j,k} \approx q_k$$
$$|v_{j,k}| \approx |q_k| \gg 0$$

B. Two key vectors  $k_i$  and  $k_j$  are said to be perpendicular,  $k_i \perp k_j$ , if  $k_i^T k_j = 0$ . The length, norm, or magnitude of a key vector  $k_i$ ,  $\|k_i\| = \sqrt{k_{i,1}^2 + k_{i,2}^2 + \dots + k_{i,n}^2}$ .

Suppose we have key vectors  $K = \{k_1, \dots, k_n\}$  such that  $k_i \perp k_j$  for  $i \neq j$  and  $\|k_i\| = 1$  for all

i. Suppose we have  $V = \{v_1, \dots, v_n\}$  and our objective is to find a vector  $q$  such that

$$c = \frac{1}{2}(v_a + v_b)$$

for arbitrary  $a, b \in \{1, \dots, n\}$ . By the definition of  $c$  we get

$$\frac{1}{2}(v_a + v_b) = \sum_{i=1}^n v_i \alpha_i \quad \alpha_i = \frac{e^{k_i q}}{\sum_{j=1}^n e^{k_j^T q}}$$

Combining this all together we get that

$$\frac{1}{2}(v_a + v_b) = \sum_{i=1}^n v_i \frac{e^{k_i q}}{\sum_{j=1}^n e^{k_j^T q}}$$

Let  $q = c(k_a + k_b)$ . Then we have

$$k_i^T q = \frac{1}{2}(k_i^T k_a + k_i^T k_b) \quad \text{for all } i$$

when  $i \neq a, b$  then we have

$$k_i^T q = \frac{1}{2}(0 + 0) \quad \text{as } k_i \perp k_a \text{ and } k_i \perp k_b$$

when  $i = a$  we have

$$k_a^T q = \frac{1}{2}(k_a^T k_a + 0)$$

Because  $\|k_a\| = 1$ , we know that  $k_a^T k_a = 1$ :

$$1 = \sqrt{k_{i1}^2 + \dots + k_{in}^2}$$

$$1 = k_{i1}^2 + \dots + k_{in}^2$$

Therefore  $k_a^T q = \frac{1}{2} \cdot 1 = \frac{1}{2}$ .

Similarly  $k_b^T q = \frac{1}{2}$

$$\text{Therefore } \alpha_i = \left\{ \frac{e^0}{n-2 \cdot 2^{\frac{1}{2}}} \frac{1}{2^{\frac{1}{2}} + n-2} \right\}$$

Lemma 1) If  $\|v\| = 1$ , then that implies  $v^T v = 1$ .

Proof:

$$\|v\| = 1$$

$$\sqrt{v_1^2 + \dots + v_n^2} = 1$$

$$v_1^2 + \dots + v_n^2 = 1$$

$$v^T v = 1$$

Proof

Let  $q = \hat{c}(k_a + k_b)$  where  $\hat{c}$  is a very large constant:  $\hat{c} \Rightarrow \infty$ . First, let's compute the attention scores  $\alpha$  with this

choice of  $q$ . If  $i \neq a, b$  then we have

$$\begin{aligned} k_i^T q &= k_i^T \hat{C}(k_a + k_b) \\ &= \hat{C}(k_i^T k_a + k_i^T k_b) \\ &= \hat{C}(0 + 0) \quad \text{as } k_i \perp k_a \\ &= 0 \quad \text{as } k_i \perp k_b \end{aligned}$$

If  $i = a$  we have

$$\begin{aligned} k_a^T q &= k_a^T \hat{C}(k_a + k_b) \\ &= \hat{C}(k_a^T k_a + k_a^T k_b) \\ &= \hat{C}(1 + 0) \quad \text{as } k_a \perp k_b \\ &= \hat{C} \end{aligned}$$

Similarly if  $i = b$  we have

$$\begin{aligned} k_b^T q &= k_b^T \hat{C}(k_a + k_b) \\ &= \hat{C}(k_b^T k_a + k_b^T k_b) \quad \text{as } k_a \perp k_b \\ &= \hat{C} \end{aligned}$$

Therefore using this we get that the quantity

$$\sum_{i=1}^n e^{k_i^T q} = e^{\hat{C}} + e^{\hat{C}} = 2e^{\hat{C}}$$

Now we can compute the attention scores  $\alpha$ :

$$\alpha_i = \begin{cases} \frac{e^0}{2e^0} & \text{if } i \neq a, b \\ \frac{e^i}{2e^0} & \text{if } i = a, b \end{cases}$$

$$\alpha_i \approx \begin{cases} 0 & \text{if } i \neq a, b \\ 1/2 & \text{if } i = a, b \end{cases}$$

Therefore this gives us the following

$$C = \sum_{i=1}^n v_i \alpha_i \approx \frac{1}{2} v_a + \frac{1}{2} v_b = \frac{1}{2} (v_a + v_b)$$

as desired.

C

I. Suppose we have  $k = \{k_1, \dots, k_n\}$  where:  
 $k_i \sim N(u_i, \Sigma_i)$   $u_i \in \mathbb{R}^n$   $\Sigma_i \in \mathbb{R}^{n \times n}$

and  $u_i$  are known but  $\Sigma_i$  are not. Also suppose that:

$$u_i^T u_j = 0 \text{ if } i \neq j \text{ and } \|u_i\| = 1$$

meaning that all  $u_i$  are perpendicular to each other and have unit norm. Also assume that the covariance matrices are of the form:

$$\Sigma_i = \alpha I \quad \alpha \Rightarrow 0$$

As  $\alpha$  approaches 0,  $k_i$  approaches the mean vectors

$u_i$  as  $k_i \sim N(u_i, \Sigma_i \Rightarrow 0)$ . Therefore, if we want  $C = \frac{1}{2}(u_a + v_b)$ , then we can let:

$$q = C^* (k_a + k_b) \quad C^* \gg 0$$

using the results from above as:

$k_i \perp k_j \quad i \neq j$  because  $k_i \sim u_i$  and  $u_i \perp u_j \quad i \neq j$

$\|k_i\| = 1$  because  $k_i \sim u_i$  and  $\|u_i\| = 1$

Because  $k_a \sim u_a$  and  $k_b \sim u_b$  this reduces to

$$q = C^* (u_a + u_b)$$

as desired.

II Suppose that  $\Sigma_a = \alpha I + \frac{1}{n}(u_a u_a^T)$ . Because  $\|u_a\| = 1$ , we have

$$\|u_a\| = 1$$

$$\sqrt{u_{a1}^2 + \dots + u_{an}^2} = 1$$

$$u_{a1}^2 + \dots + u_{an}^2 = 1$$

$$u_a u_a^T = I$$

This means that as  $\alpha \Rightarrow 0$   $\Sigma_a = \frac{1}{2}I$ . Now

$$k_a \sim N(u_a, \frac{1}{2}I) \quad k_i \sim (u_i, 0) \quad i \neq a$$

That means that depending on the sample,

$\|k_a\| > 1$  or  $\|k_a\| < 1$  while  $\|k_i\| \sim 1 \quad i \in a$ .

Ext are equally likely. This means that

$$\begin{cases} \alpha_a > \alpha_b & \text{when } \|k_a\| > 1 \\ \alpha_a < \alpha_b & \text{when } \|k_a\| < 1 \\ \alpha_i = 0 & \text{when } i \neq a, b \end{cases}$$

This means that our output  $c$  will look more like  $v_a$  when  $\|k_a\| > 1$  and more like  $v_b$  when  $\|k_a\| < 1$ .

D We will define  $q_1$  and  $q_2$  such that  $c = \frac{1}{2}(v_a + v_b)$  with the added condition that  $q_1 \neq q_2$ . In order to do this we can do the following:

$$q_1 = c^{\otimes} u_a \quad c^{\otimes} \gg 0$$

$$q_2 = c^{\otimes} u_b \quad c^{\otimes} \gg 0$$

If we do that then we have:

$$\alpha_i^1 = \begin{cases} 1 & i=a \\ 0 & i \neq a \end{cases} \quad \alpha_i^2 = \begin{cases} 1 & i=b \\ 0 & i \neq b \end{cases}$$

This gives us

$$C_1 = V_a \quad C_2 = V_b$$

$$C = \frac{1}{2}(C_1 + C_2) = \frac{1}{2}(V_a + V_b)$$

II Suppose we have a sample of  $k = \{k_1, \dots, k_n\}$

Then we have:

$$k_a = N(u_1, \frac{1}{2}) \quad k_i = N(u_1, 0) \quad i \neq a$$

This means that

$$\alpha_i^1 = \begin{cases} 1 & \text{when } i = a \\ 0 & i \neq a \end{cases}$$

$$\alpha_i^2 = \begin{cases} 1 & \text{when } i = b \\ 0 & i \neq b \end{cases}$$

These are regardless of the variation of the norm  $\|k_a\|$  from sampling. Thus

$$C_1 = V_a$$

$$C_2 = V_b$$

And  $C = \frac{1}{2}(C_1 + C_2) = \frac{1}{2}(V_a + V_b)$ . The key here is that the perturbation on  $k_a$  had no effect on the resulting attention scores since we had multiple queries.



Extended elaboration: The takeaway here is that multihaded attention allows the network to construct independent queries and take a even average of those individual attention output. This helps the network out when the scale of keys is slightly different.

[E] If we perform self-attention then we get that

$$C_2 = \sum_{j=1}^n \alpha_{2j} V_j$$

In order to compute this we first compute the required quantities:

$$q_2 = x_2 = u_a$$

$$k_1 q_2 = (u_a + u_b)(u_a) = 0 \quad k_2 q_2 = (u_a)(u_a) = B$$

$$k_3 q_2 = (u_c + u_d)u_a = 0$$

$$\alpha_{21} = \frac{e^0}{\sum_{j=1}^3 e^{k_j q_2}} = \frac{1}{e^B} = 0$$

$$\alpha_{22} = \frac{e^B}{e^B} = 1 \quad \alpha_{23} = \frac{e^0}{e^0} = 0$$

Then we have  $C^2$  to be

$$C_2 = \sum_{j=1}^n \alpha_{2j} V_j = V_2 = x_2 = u_a$$

So  $C^2$  will approximate  $u_a$ . It would not be

possible for  $c_2$  to approximate  $u_b$  by adding either  $u_a$  or  $u_c$ . This is because:

1) If we set  $x_2$  to be  $u_a + u_b$ , then  $\alpha_2 = (0, 1, 0)$  as  $e^{2B} \gg e^B$ . So  $c^2$  will be  $u_a + u_b$

2) If we set  $x_2$  to be  $u_a + u_c$ , then again  $\alpha_2 = (0, 1, 0)$  as  $e^{2B} \gg e^0$ . So  $c^2$  will be  $u_a + u_c$

Neither of these are  $u_b$

II First let's try to find  $V$  as recommended by the author. Note that  $u_a, u_b, u_c, u_d \in \mathbb{R}^{n+1}$

$$V_1 = Vx_1 = V(u_a + u_b) = Vu_a + Vu_b$$

So we want  $Vu_a = 0$  and  $Vu_b = u_b$

$$V_3 = Vx_3 = V(u_c + u_b) = Vu_c + Vu_b$$

So we want  $Vu_c = -u_c$  and  $Vu_b = u_b$

Let's focus on the first condition  $Vu_b = u_b$ .

Using the hint from the assignment, note that

$$(u_b u_b^T) u_b = u_b \|u_b\|^2$$

So let's set  $V = \frac{1}{B} u_b u_b^T$ . Then we have

$$\begin{aligned} V_1 = Vx_1 &= \frac{1}{B} u_b u_b^T (u_a + u_b) \\ &= \frac{1}{B} u_b u_b^T u_b \quad \text{as } u_b u_b^T u_a = 0 \end{aligned}$$

$$= \frac{1}{B} u_b B = u_b \quad \checkmark$$

Let's do something similar by adding  $-\frac{1}{B} u_c u_c^T$  to our choice of  $V$ . Then

$$V = \frac{1}{B} u_b u_b^T - \frac{1}{B} u_c u_c^T$$

Then

$$\begin{aligned} V_3 &= V(u_c + u_b) = Vu_c + Vu_b \\ &= (0 + -\frac{1}{B} u_c u_c^T u) + (\frac{1}{B} u_b u_b^T u_b + 0) \quad \text{by orthogonality} \\ &= -u_c + u_b \quad \checkmark \end{aligned}$$

$$V_1 = V(u_d + u_b) = 0 + \frac{1}{B} u_b u_b^T = u_b \quad \checkmark$$

So our choice in  $V$  is good. Now let's try to find  $K$  and  $Q$ . We want  $c_2 \approx u_b$ . That will only occur when  $\alpha_2 = (1, 0, 0)$  given our choice of  $V$ . Similarly  $c_1 \approx u_a - u_c$  will only occur when  $\alpha_1 = (0, 0, 1)$ .

Now we will let  $Q$  and  $K$  be the following

$$Q = u_a u_a^T + u_c u_c^T \quad K = I$$

$$q_1 = (u_a u_a^T + u_c u_b^T)(u_a + u_b) = 0 + u_c u_b^T u_a = u_c$$

$$k_1 = kx_1 = u_a + u_b$$

$$k_2 = kx_2 = u_a$$

$$k_3 = kx_3 = u_c + u_b$$

$$\alpha_{11} = \frac{(u_a + u_b)^T u_c}{e^0} = 0 \quad \alpha_{12} = \frac{u_a^T u_c}{e^0} = 0 \quad \alpha_{13} = \frac{(u_a + u_b)^T u_b}{e^0} = 1$$

$$\alpha_1 = (0, 0, 1)$$

$$c_1 = v_3 = V(u_c + u_b) = u_b - u_c \quad \checkmark$$

$$q_2 = (u_a u_a^T + u_c u_b^T) u_a = u_a$$

$$k_1 = u_a + u_b \quad k_2 = u_a \quad k_3 = u_c + u_b \quad \text{this is from the result above}$$

$$\alpha_2 = (1, 0, 0)$$

$$c_2 = v_1 = V(u_a + u_b) = u_b \quad \checkmark$$

## 2. Pretrained Transformer models and knowledge access (35 points)

D.

On the development test set, I got a 1.401% accuracy score from a model trained from scratch versus a 5% accuracy score from a baseline model that just predicts "London" for each test row.

```
(local_nmt) sh-4.2$ python src/london_baseline.py
500 500
Result if just predicted London: (500.0, 25.0)
(local_nmt) sh-4.2$ python src/london_baseline.py
Result if just predicted London: correct:25.0, total:500.0: percentage:0.05
(local_nmt) sh-4.2$ python src/run.py evaluate vanilla.wiki.txt --reading_params_path vanilla.model.params --eval_corpus_path birth_dev.tsv --outputs_path vanilla.nopretrain.dev.predictions
data has 418352 characters, 256 unique.
number of parameters: 3323392
500it [01:03, 7.83it/s]
Correct: 7.0 out of 500.0: 1.4000000000000001%
```

F.

I got a 25.8% accuracy score on the development test set after pretraining and finetuning.

```
epoch 1 iter 7: train loss 0.73257. lr 5.999844e-04: 100%|██████████| 8/8 [00:13<00:00, 1.71s/it]
epoch 2 iter 7: train loss 0.53881. lr 5.999351e-04: 100%|██████████| 8/8 [00:16<00:00, 2.03s/it]
epoch 3 iter 7: train loss 0.44053. lr 5.998521e-04: 100%|██████████| 8/8 [00:17<00:00, 2.18s/it]
epoch 4 iter 7: train loss 0.35198. lr 5.997352e-04: 100%|██████████| 8/8 [00:16<00:00, 2.05s/it]
epoch 5 iter 7: train loss 0.30953. lr 5.995847e-04: 100%|██████████| 8/8 [00:16<00:00, 2.05s/it]
epoch 6 iter 7: train loss 0.26981. lr 5.994004e-04: 100%|██████████| 8/8 [00:16<00:00, 2.04s/it]
epoch 7 iter 7: train loss 0.21797. lr 5.991823e-04: 100%|██████████| 8/8 [00:16<00:00, 2.03s/it]
epoch 8 iter 7: train loss 0.18256. lr 5.989306e-04: 100%|██████████| 8/8 [00:16<00:00, 2.03s/it]
epoch 9 iter 7: train loss 0.14337. lr 5.986453e-04: 100%|██████████| 8/8 [00:16<00:00, 2.03s/it]
epoch 10 iter 7: train loss 0.13564. lr 5.983263e-04: 100%|██████████| 8/8 [00:16<00:00, 2.03s/it]
(local_nmt) sh-4.2$ python src/run.py evaluate vanilla.wiki.txt --reading_params_path vanilla.finetune.params --eval_corpus_path birth_dev.tsv --outputs_path vanilla.pretrain.dev.predictions
data has 418352 characters, 256 unique.
number of parameters: 3323392
500it [01:05, 7.68it/s]
Correct: 129.0 out of 500.0: 25.8%
(local_nmt) sh-4.2$
```

G.

I

```
(local_nmt) sh-4.2$ python src/run.py evaluate synthesizer.wiki.txt --reading_params_path synthesizer.finetune.params --eval_corpus_path birth_dev.tsv --outputs_path synthesizer.pretrain.dev.prediction
data has 418352 characters, 256 unique.
number of parameters: 3076988
500it [01:13, 6.82it/s]
Correct: 47.0 out of 500.0: 9.4%
(local_nmt) sh-4.2$ python src/run.py evaluate synthesizer.wiki.txt \ --reading_params_path synthesizer.finetune.params \ --eval_corpus_path birth_test_inputs.tsv \r.pretrain.dev.predictions
usage: run.py [-h] [--reading_params_path READING_PARAMS_PATH] [--writing_params_path WRITING_PARAMS_PATH] [--finetune_corpus_path FINETUNE_CORPUS_PATH] [--eval_corpus_path EVAL_CORPUS_PATH]
               [--outputs_path OUTPUTS_PATH]
               {pretrain,finetune,evaluate} {vanilla,synthesizer} pretrain_corpus_path
run.py: error: unrecognized arguments: --reading_params_path synthesizer.finetune.params --eval_corpus_path birth_test_inputs.tsv --outputs_path r.pretrain.dev.predictions
(local_nmt) sh-4.2$ python src/run.py evaluate synthesizer.wiki.txt --reading_params_path synthesizer.finetune.params --eval_corpus_path birth_test_inputs.tsv --outputs_path synthesizer.pretrain.dev.predictions
data has 418352 characters, 256 unique.
number of parameters: 3076988
437it [01:02, 7.01it/s]
No gold birth places provided; returning (0,0)
Predictions written to synthesizer.pretrain.dev.predictions; no targets provided
```

I got a 9.4% accuracy on the development test set (which hopefully translates to > 5% accuracy on the test set, if these are both held out random samples (no hyperparameter tuning was done)).

II

CausalSelfAttention implements attention by first generating a query representation for each token from the input, generating a key representation for each token from the input, and then comparing the two to decide how to weight the value representations in the output. This in the

synthesizer paper is referred to as token-token interactions, as the input is involved in a dot product with itself.

The primary purpose of attention is learning self-alignment: figuring out which individual tokens in the input sequence are relevant to the current output. This is especially important if the input sequence is long.

SynthesizerAttention skips token-token interactions to learn self alignment by instead learning a more powerful neural network that ingests the input once to learn pairwise attention weights. This is much faster as it requires 3 matrix multiplications instead of 4 (and compensates a little bit using a relu on top of the query representation).

So if a problem we are working on has self-alignment properties that need to compare token representations exactly to themselves (to see how similar they are), synthesizer attention will more likely struggle with this.

Another way to summarize this is that synthesizer attention can't flexibly create custom query and key representations based on the input sequence. It has to create one general representation for the keys, and then create ideal queries using the input sequence to compute attention from those keys.

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

A.

The pretrained model was able to achieve an accuracy score much higher than 10% because it was able to learn a lot about language and places associated to people beforehand from the character deletion task. This allowed it much better accuracy to predict people's birthplace after training on the birth-place tasks because both of these learnings are relevant for that task.

B.

The predictions in birth place prediction from the pre-trained model mostly look plausible (regardless of whether they are correct or not). For a user using this system, this might cause a problem of misinformation - where because the prediction looks plausible, the user is most likely to use and share that prediction. This also is less desirable for the user because it makes it difficult for the user to filter for or create confidence intervals around the prediction. If instead the predictions were either 1) correct or 2) gibberish, it would be easy for the user of the system to just throw out the gibberish incorrect predictions and just use the correct ones in their downstream applications.

C.

If the model didn't see the birthplace in either the pretraining or training task, it's impossible for the model to have learned the correct birthplace. In this case, the model will most likely default to "average" distributions based on the name origin of the person. For instance, if the model has learned that a name is frequent in the middle east (e.g. Muhammad), it might predict a high population middle eastern city like Mecca. Additionally, if the model has learned that the name John is frequent in the West, it might predict a high population western city like New York City. This can cause a bias in error based on name origin, where people with names that originate from different locations often get a higher error than people with names common to their actual birth-place.