(a) if q = k; and $q \perp k$; (if)

then $c \approx v$;

(b)
$$q = M(Ma + Nb)$$
big[†] number

$$C = \frac{e \quad V_{0} + k_{0})k_{0}}{V_{0} + e \quad V_{0} + V_{0} + \dots}$$

$$e \quad W_{0} + e \quad V_{0} + k_{0} + k_{0$$



Intuitively, we can assume that the following statements are roughly true:

- If matrices A and B are both from a normal distributions with the mean μ and $\Sigma_i = \alpha I$ (like mentioned in the question) then their dot product will approximately equal to $\|\mu\|^2$. (Because A and B roughly point to the same direction and they have roughly the same norm.)
- If matrices A and B are from two normal distributions with the means μ_1 and μ_2 and $\mu_1 \perp \mu_2$ and $\Sigma_i = \alpha I$ for both (like mentioned in the question) then their dot product will approximately equal to 0. (Because A and B roughly point to the same direction as μ_1 and μ_2 hence A is roughly perpendicular to B)

$$Q = M(N_{a} + N_{b})$$

$$M(N_{a} + N_{b}) K_{a} M(N_{a} + N_{b}) K_{b} M(N_{a} + N_{b}) K_{c}$$

$$C = \frac{e}{M(N_{a} + N_{b}) K_{a}} M(N_{a} + N_{b}) K_{b} M(N_{a} + N_{b}) K_{c}$$

$$e + e + e$$

$$M(N_{a} + N_{b}) K_{a} M(N_{a} + N_{b}) K_{c}$$

$$e + e + e$$

$$M(N_{a} + N_{b}) K_{a} M(N_{a} + N_{b}) K_{c}$$

$$e + e + e$$

$$M(N_{a} + N_{b}) K_{a} M(N_{a} + N_{b}) K_{c}$$

$$e + e + e$$

$$M(N_{a} + N_{b}) K_{a} M(N_{a} + N_{b}) K_{c}$$

$$e + e + e$$

$$M(N_{a} + N_{b}) K_{a} M(N_{a} + N_{b}) K_{c}$$

$$e + e + e$$

$$M(N_{a} + N_{b}) K_{a} M(N_{a} + N_{b}) K_{c}$$

$$e + e + e$$

$$M(N_{a} + N_{b}) K_{a} M(N_{a} + N_{b}) K_{c}$$

$$e + e + e$$

$$e + e$$

$$e + e + e$$

$$e + e$$

$$e + e + e$$

$$e + e$$

The terms " $k_i^T.\mu_a$ ", " $k_i^T.\mu_b$ " and hence "exp(log(M)($k_i^T.\mu_a + k_i^T.\mu_b$)). v_i " are approximately zero according to the above statements.

ii. In this case the term " $k_a^T.\mu_a$ " will vary roughly between 1/2 and 3/2. Thus, c will have a value roughly between $(v_a+2v_b)/3$ and $(3v_a+v_b)/5$. In the first extreme case c will lean towards v_b and in the second case it will lean towards v_a . This phenomenon shows an inconsistency between different samplings of k_i .



$$q_1 = M/a$$

$$q_2 = M/b$$

$$C_{1} = \frac{e \quad V_{a} + e \quad V_{b} + \cdots + e \quad V_{i} + \cdots}{e \quad V_{a} + e \quad V_{b} + \cdots + e \quad V_{i} + \cdots}$$

$$e \quad H(N_{a}, N_{a}) \quad M(N_{b}, N_{a}) \quad M(N_{b}, N_{a})$$

$$e \quad + e \quad + \cdots + e \quad + \cdots$$

$$e \quad M(N_{a}, N_{a})$$

$$e \quad V_{a}$$

$$e \quad V_{a}$$

$$C_{2} = \frac{M(N_{a}N_{b})}{M(N_{a}N_{b})} \frac{M(N_{b}N_{b})}{M(N_{b}N_{b})} \frac{M(N_{b}N_{b})}{M(N_{b}N_{b})} \frac{M(N_{b}N_{b})}{M(N_{b}N_{b})} \frac{M(N_{b}N_{b})}{M(N_{b}N_{b})} \frac{M(N_{b}N_{b})}{M(N_{b}N_{b})} = \frac{M(N_{b}N_{b})}{M(N_{b}N_{b})} \frac{M(N_{b}N_{b})}{M(N_{b}N_{b})} \frac{M(N_{b}N_{b})}{M(N_{b}N_{b})}$$

= V_b

ii. In this case always c_1 equals to v_a approximately and c_2 approximately equals to v_b because all other terms are approximately zero and the remaining coefficients in the numerator and denominator cancel out. In other words, variance in the magnitude of k_a , unlike part (c), does not have much effect.

$$C_{q} = \frac{e^{-v_{1}q_{2}} + e^{-v_{2}q_{2}} + e^{-v_{3}q_{2}}}{e^{-v_{1}q_{2}} + e^{-v_{2}q_{2}} + e^{-v_{3}q_{2}}}$$

$$= \frac{e^{-v_{1}q_{2}} + e^{-v_{2}q_{2}} + e^{-v_{3}q_{2}}}{e^{-v_{1}q_{2}} + e^{-v_{3}q_{2}}}$$

$$= \frac{e^{-v_{1}q_{2}} + e^{-v_{2}q_{2}} + e^{-v_{3}q_{2}}}{2 + e^{-v_{3}q_{2}}}$$

$$= \frac{e^{-v_{1}q_{2}} + e^{-v_{3}q_{2}}}{2 + e^{-v_{3}q_{3}}}$$

$$= \frac{e^{-v_{1}q_{2}} + e^{-v_{3}q_{2}}}{2 + e^{-v_{3}q_{3}}}$$

$$= \frac{e^{-v_{1}q_{2}} + e^{-v_{3}q_{2}}}{2 + e^{-v_{3}q_{3}}}$$

$$= \frac{e^{-v_{1}q_{2}} + e^{-v_{3}q_{3}}}{2 + e^{-v_{3}q_{3}}}$$

$$= \frac{e^{-v_{1}q_{3}} + e^{-v$$

No, because x_1 and x_3 have the terms u_d or u_c and hence the softmax of $k_1^T q_2$ or $k_3^T q_2$ will not be near zero.

$$V_{1} = U_{5} = V_{1} = V(u_{5} + u_{d}) \Rightarrow V \text{ should have } \frac{1}{\|u_{b}\|^{2}} u_{b} u_{b}^{T}$$

$$V_{8} = U_{5} - u_{c} = V_{3} = V(u_{5} + u_{c}) \Rightarrow V \text{ should have } \frac{1}{\|u_{c}\|^{2}} u_{c} u_{c}^{T}$$

$$V_{8} = U_{5} - u_{c} = V_{3} = V(u_{5} + u_{c}) \Rightarrow V \text{ should have } \frac{1}{\|u_{c}\|^{2}} u_{c} u_{c}^{T}$$

$$\mathcal{K}_{i} = \mathcal{N}_{i} \longrightarrow \mathbf{K} = \mathbf{T}$$

$$Q_{i} = \mathcal{Q}_{i} = \mathcal{Q}_{i}$$



```
8/8 [00:04<00:00,
                                                                                         1.68it/s]
epoch 65 iter 7: train loss 0.32717. lr 5.315115e-04: 100%
                                                                      8/8 [00:04<00:00, 1.66it/s]
epoch 66 iter 7: train loss 0.30116. lr 5.294740e-04: 100%
                                                                      8/8 [00:04<00:00, 1.67it/s]
epoch 67 iter 7: train loss 0.30572. lr 5.274107e-04: 100%
                                                                      8/8 [00:04<00:00, 1.68it/s]
epoch 68 iter 7: train loss 0.29130. lr 5.253217e-04: 100%
                                                                      8/8 [00:04<00:00, 1.68it/s]
epoch 69 iter 7: train loss 0.27051. lr 5.232074e-04: 100%
                                                                      8/8 [00:04<00:00, 1.67it/s]
epoch 70 iter 7: train loss 0.25012. lr 5.210680e-04: 100%
                                                                      8/8 [00:04<00:00, 1.68it/s]
epoch 71 iter 7: train loss 0.23997. lr 5.189037e-04: 100%
                                                                      8/8 [00:04<00:00, 1.67it/s]
epoch 72 iter 7: train loss 0.23392. lr 5.167147e-04: 100%
                                                                     8/8 [00:04<00:00, 1.68it/s]
epoch 73 iter 7: train loss 0.22486. lr 5.145014e-04: 100%
                                                                     8/8 [00:04<00:00, 1.66it/s]
epoch 74 iter 7: train loss 0.21690. lr 5.122639e-04: 100%
                                                                     8/8 [00:04<00:00, 1.67it/s]
epoch 75 iter 7: train loss 0.20826. lr 5.100024e-04: 100%
                                                                     8/8 [00:04<00:00, 1.67it/s]
```

data has 418352 characters, 256 unique.

number of parameters: 3323392 number of parameters: 3323392

500it [00:41, 12.01it/s]

Correct: 9.0 out of 500.0: 1.79999999999998%

der ace

data has 418352 characters, 256 unique.

number of parameters: 3323392 number of parameters: 3323392

437it [00:42, 10.31it/s]

No gold birth places provided; returning (0,0)

Predictions written to vanilla.nopretrain.test.predictions; no targets provided

data has 418352 characters, 256 unique.

number of parameters: 3323392

437it [00:35, 12.27it/s]

Correct: 63 out of 437: 14.416475972540047%

landon baseline



```
2 f (1)
  epoch 643 iter 22: train loss 0.50243. lr 6.000000e-04: 100%
                                                                   23/23 [00:07<00:00,
                                                                                     3.04it/s
 epoch 644 iter 22: train loss 0.47090. lr 6.000000e-04: 100%
                                                                   23/23 [00:07<00:00, 3.04it/s]
 epoch 645 iter 22: train loss 0.45273. lr 6.000000e-04: 100%
                                                                   23/23 [00:07<00:00, 3.04it/s]
                                                                  23/23 [00:07<00:00, 3.05it/s]
 epoch 647 iter 22: train loss 0.49273. lr 6.296935e-04: 100%
                                                                   23/23 [00:07<00:00,
                                                                                     3.04it/s
 epoch 648 iter 22: train loss 0.46037. lr 6.586443e-04: 100%
                                                                  23/23 [00:07<00:00, 3.04it/s]
                                                                                     3.04it/s
 epoch 650 iter 22: train loss 0.48776. lr 7.182453e-04: 100%
                                                                  23/23 [00:07<00:00,
                                                                                     3.04it/s]
```

```
data has 418352 characters, 256 unique.

number of parameters: 3323392

number of parameters: 3323392

epoch 1 iter 7: train loss 0.71179. lr 5.999844e-04: 100%| | 8/8 [00:04<00:00, 1.70it/s]

epoch 2 iter 7: train loss 0.56339. lr 5.999351e-04: 100%| | 8/8 [00:04<00:00, 1.80it/s]

epoch 3 iter 7: train loss 0.45983. lr 5.998521e-04: 100%| | 8/8 [00:04<00:00, 1.79it/s]

epoch 4 iter 7: train loss 0.40272. lr 5.997352e-04: 100%| | 8/8 [00:04<00:00, 1.79it/s]

epoch 5 iter 7: train loss 0.33338. lr 5.995847e-04: 100%| | 8/8 [00:04<00:00, 1.79it/s]

epoch 6 iter 7: train loss 0.29572. lr 5.994004e-04: 100%| | 8/8 [00:04<00:00, 1.74it/s]

epoch 7 iter 7: train loss 0.22406. lr 5.991823e-04: 100%| | 8/8 [00:04<00:00, 1.77it/s]

epoch 8 iter 7: train loss 0.18956. lr 5.989306e-04: 100%| | 8/8 [00:04<00:00, 1.76it/s]

epoch 9 iter 7: train loss 0.12842. lr 5.983263e-04: 100%| | 8/8 [00:04<00:00, 1.76it/s]

epoch 10 iter 7: train loss 0.12842. lr 5.983263e-04: 100%| | 8/8 [00:04<00:00, 1.76it/s]
```

data has 418352 characters, 256 unique.

number of parameters: 3323392

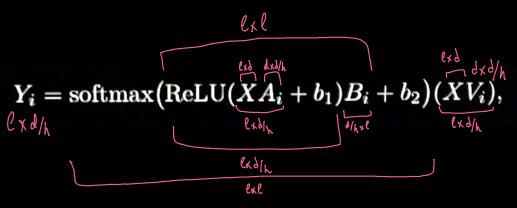
number of parameters: 3323392

437it [00:40, 10.77it/s]

No gold birth places provided; returning (0,0)

Predictions written to vanilla.pretrain.test.predictions; no targets provided





```
epoch 637 iter 22: train loss 0.58724. Lr 0.0000000e-04: 100% 23/23 [00:07<00:00, 3.05it/s] epoch 638 iter 22: train loss 0.58698. Lr 6.000000e-04: 100% 23/23 [00:07<00:00, 3.05it/s] epoch 639 iter 22: train loss 0.60296. Lr 6.000000e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 640 iter 22: train loss 0.56730. Lr 6.000000e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 641 iter 22: train loss 0.58239. Lr 6.000000e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 642 iter 22: train loss 0.57320. Lr 6.000000e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 643 iter 22: train loss 0.57320. Lr 6.000000e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 643 iter 22: train loss 0.58172. Lr 6.000000e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 644 iter 22: train loss 0.55105. Lr 6.000000e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 645 iter 22: train loss 0.56613. Lr 6.000000e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 646 iter 22: train loss 0.56613. Lr 6.000000e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 647 iter 22: train loss 0.56178. Lr 6.013186e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 648 iter 22: train loss 0.58965. Lr 6.296935e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 649 iter 22: train loss 0.56178. Lr 6.586443e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 649 iter 22: train loss 0.57799. Lr 6.881640e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 650 iter 22: train loss 0.57799. Lr 6.881640e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 650 iter 22: train loss 0.57799. Lr 6.881640e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 650 iter 22: train loss 0.57799. Lr 6.881640e-04: 100% 23/23 [00:07<00:00, 3.06it/s] epoch 650 iter 22: train loss 0.57799. Lr 6.881640e-04: 100% 23/23 [00:07<00:00, 3.06it/s]
```

🧼 2 g (2) :

🥏 2 g (3)

data has 418352 characters, 256 unique. number of parameters: 3076988

500it [00:39, 12.50it/s]

Correct: 40.0 out of 500.0: 8.0%

4 der acc

🥏 2 g (4)

```
data has 418352 characters, 256 unique.
number of parameters: 3076988
437it [00:35, 12.28it/s]
No gold birth places provided; returning (0,0)
Predictions written to synthesizer.pretrain.test.predictions; no targets provided
```



The synthesizer method does not take contextual and relative similarity of tokens into consideration. (Does not have a mechanism similar to Q and K)





The pretrained model captured low level information about the data using the corrupted span strategy.



- 1. Models of this kind which can output some wrong results that seem very similar to a true result are sometimes really hard to debug.
- 2. They can be misleading and sometimes dangerous if they are used in an important decision making process.
- 3. They may learn unwanted biases.



The model may learn a correlation between a name (and it's similar forms) and a location which is not always true. In such cases the model can use this similarity to make up a somewhat realistic birthplace if a name similar (not identical) to a previously learnt name is given.

In some use cases it may raise a subtle issue: information leakage. For example consider an application in which a user X has the authority to query about a name A but does not have permission to query about a name B which is very similar to A. This type of model can leak some info about B to the user X.