# ملیکا محمدی فخار ۹۹۵۲۲۰۸۶

## منابع مورد استفاده:

https://stats.stackexchange.com/questions/386535/why-cant-we-use-back-propagation-in-hard-attention-but-we-can-use-it-in-relu

### /https://chat.openai.com

1.

بخش a) گزینههای ب و ج صحیح هستند.

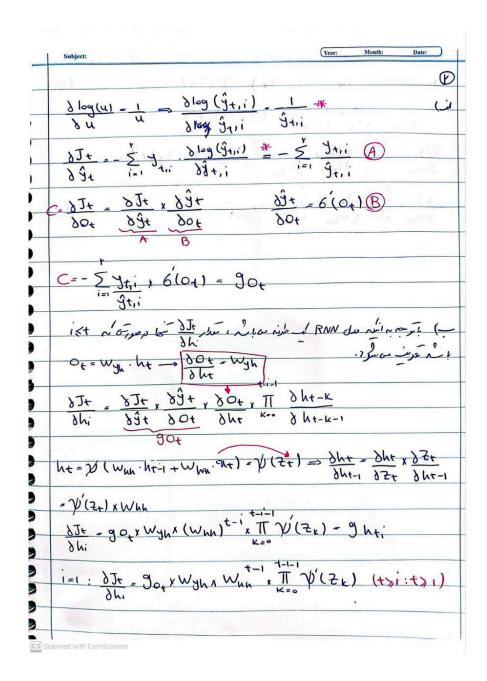
معماری many to one RNN برای تسکهایی مناسب هستند که ورودی آنها فرم sequential دارد و خروجی واحد نیست انتخاب نشده است.

بخش b) گزینه ج صحیح است.

زیرا اخلاق گربه به آب و هوای گذشته و امروز بستگی دارد اما به آب و هوای آینده بستگی ندارد و ما آن را نمیدانیم زیرا هنوز اتفاق نیفتاده است.

بخش c) گزینه ج صحیح است.

در یک مدل RNN هدف پیشبینی یک مرحله با توجه به مراحل قبلی آن است.



| ht= 2 (Whh ht= + What a                | 1 = 1                     | (7+)     |          | (2                                      |
|--|---------------------------|----------|----------|---|
|  |                           |          | 197      | - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 |
| => She She x BZ+ - V'                  | (5+)x                     | h t-1    | , b      | x. Z                                    |
| dyn i=1 dhi dwh                        | I gh                      | · W 12   | +) xhi-  | 1-16                                    |
| dun i=1 thi twhn                       | i=1 41                    | + 15     |          | 196                                     |
| 8J - 5 8Jt - 5 9,                      | A-1                       | +58.     | +E3      |   |
| DWAN t=1 8Whn t=1                      | uh, t                     | y 0 6    | 760      | ,96                                     |
|  |                           | - 3      | 2/10-    |   |
|  | a                         | L.A.     |          |   |
|  | 700                       |          | 77       |   |
|  | (8)                       |          | - 1      |   |
|  | -                         | 1.12     |          |   |
|  | 1500                      |          | AN IC    | V = 30                                  |
| Standard In                            |                           | 100      | 27       | 47.2                                    |
| (-ab-17) 6                             | 9 440                     | 107      | 168      | ,,,,,,                                  |
|  |                           | 1        |          |   |
| 11 6 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 | 1.55                      | es Vital | SA Jan L | 130 - 30                                |
| F-24 E - F-7 - MMC                     | TO NOT THE REAL PROPERTY. |          |          |   |
|  |                           |          |          |   |
|  |                           |          | 110      |   |
|  |                           |          | 7.       | 101.3                                   |
|  |                           |          |          |   |

الف. در گام نخست ضرب داخلی پرسوجو و کلیدها را محاسبه میکنیم:

for key0 = q.key0 = 3 - 2 - 3 = -2

for key1 = 9, key1 = 6 - 2 - 1 = 3

for key2 => q.key2 = 0 - 1 + 1 = 0

for key3 = q.key3 = 0 + 2 + 4 = 6

سپس اگر argmax را اعمال کنیم مشابهترین کلید اندیس ۳ را دارد.

حال با استفاده از این اندیس، مقدار مربوطه را به دست می آوریم:

Values[3] = [6, 1, 2]

ب.

استفاده از argmax در مکانیسمهای توجه، هم مزایا و هم معایبی دارد. توجه argmax شامل انتخاب عنصر با بالاترین مقدار از توزیع توجه است، که به طور مؤثر بر روی یک مکان یا ویژگی خاص تمرکز میکند. بیایید تأثیر آن بر آموزش مدلها را بررسی کنیم:

#### مزایا:

کارآیی: argmax attention به لحاظ محاسباتی بهینهتر از از soft attention است. این گونه توجه هزینه محاسباتی کمتری حین forward/backward pass ایجاد میکند.

تفسیرپذیری: وزنهای حاصل از توجه argmax دودویی (0 یا 1) هستند، که تفسیر قسمتهای مهم ورودی توسط مدل را اَسان تر میکند.

#### معایب:

عدم مشتق پذیری: چالش اصلی با توجه argmax این است که قابل مشتق گیری نیست. این یک مشکل حین فرآیند backpropagation ایجاد میکند که در آن گرادیانها محاسبه و برای بهروزرسانی پارامترهای مدل استفاده میشوند.

از بین رفتن اطلاعات: عدم مشتقپذیری توجه argmax منجر به از دست رفتن اطلاعات گرادیان میشود، که باعث سختی یادگیری و سازگاری مدل در طول آموزش میشود. این میتواند به همگرایی غیر بهینه و کارایی کمتر منجر شود.