# Attention and Self-Attention

Tirtharaj Dash

Dept. of CS & IS and APPCAIR
BITS Pilani, Goa Campus

November 16, 2021

In the prev. lecture:

(r) Seq2Seq modelling → decoding

(L) attention mechanism.

greedy —
exhaustive —
beam search. $O\left(K |y| T_y\right)$

$K=2$

real vector
↓
one-hot encoder
output
↓

target sentence

Attention model:

(Machine translation) output



$h_0 \rightarrow$ ENC $\rightarrow$ summary $\in \mathbb{R}^d$ $\rightarrow$ DEC $\rightarrow$ output

input

"large . - - - Sentence"

↓ mimic

Which parts are relevant
(or should be given more
attention to ?)

↓

attention mechanism.

je          veux                    ⟨eos⟩

$S_0$
Zero
State          ⟨bos⟩

$S_1$

Content Vector

$C_1$          $\alpha_{2,1}$    $C_2$          $\alpha_{2,6}$    $\alpha_{1,6}$    probability
                                                                                    scores

First
word
other part    $\alpha_{1,1}$    $\alpha_{1,2}$    $\alpha_{1,3}$    $\alpha_{2,2}$    $\alpha_{2,3}$

concatenated
vector                                                                              2d-dimens
                                                                                    $\in \mathbb{R}^{2d}$

$x_1$  I    $x_2$ Want    $x_3$ to    $x_4$ live    $x_5$ in    $x_6$ Europe    . . . . . .

I          I  want          ↑ want to                      .

ENC

q

I      want                                    Europe .

let $\quad h_{t'} = \left( \overrightarrow{h}_{t'}, \overleftarrow{h}_{t'} \right)$     hidden state
(summary at
$b = t'$)
from the encoder

what we want

$$\sum_{t'} \alpha_{1,t'} = 1 \quad \text{and} \quad \alpha_{1,j} \geq 0$$

$\longrightarrow$ probability.

content vector

$$C_1 = \sum_{t'} \boxed{\alpha_{1,t'}} \; h_{t'}$$

attention score or
attention weight ✔

$$\alpha_{t,t'} = \text{softmax}\left(e_{t,t'}\right)$$

$$= \frac{\exp\left(e_{t,t'}\right)}{\sum_{t'=1}^{T_x} \exp\left(e_{t,t'}\right)}$$

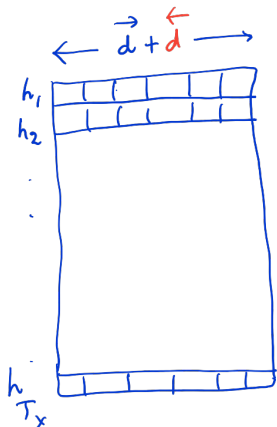attention
weight

attention
scores

$$\downarrow \quad +b_a$$

$$e_{t,t'} = \tanh\left([S_{t-1}, h_{t'}] \times W_a\right)$$

hidden state / summary of the decoder at
$t = t-1$

$h_{t'} = $ summary of the encoder
at $t = t'$

$S_{t-1} \rightarrow$

$O \rightarrow e_{t,t'}$

$h_{t'} \rightarrow$

$\xrightarrow{\vec{d}} + \xleftarrow{\vec{d}}$

concat

$\xleftarrow{} d_1 \xrightarrow{}$  summary of the decoder (prev. state)

$\oplus$

$S_{t-1}$

$h_1$
$h_2$

attention weight

$\uparrow$
$T_x$ $\boxed{\alpha}$
$\downarrow$

$\uparrow$

$h_{T_x}$

$\begin{bmatrix} T_x & \times & 2d \end{bmatrix}$

softmax$\left(\begin{array}{c} \text{tanh}(\boxed{}) \end{array}\right)$

bias?

$\begin{bmatrix} T_x \times (2d + d_1) \end{bmatrix}$ $\times$ $W_a \begin{bmatrix} T_x + 1 \end{bmatrix}$ $=$ attention scores.

$T_x \boxed{}$

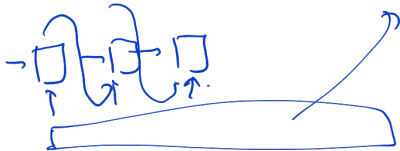# Transformer Network:

2017 : A. Vaswani et. al. "Attention is all you need"

Idea: → attention based representation ( content vectors )

+

→ convolution style of processing.



parallel processing.

2   mechanism :→

① Self-attention ←

② Multi-head attention. ← Next class

Self-attention:

$$A \left( q, K, V \right)$$

query key
value

attention-based representation.
"vector" for a word (token).

↑↓

word: word-embedding (one-hot vector)

I want to <u>live</u> in Europe.

$[0, 0, 0, 1, \ldots, 0]$    as a PostDoc ?

          as a businessman ?

$[$        as a husband ?

     rich embedding

"Self-attention" is used to construct a "rich" (very information) embedding of a word (or ⟨token⟩).

$$\text{attention} = \text{Softmax} \left( \quad \right)$$

$$= \frac{\exp(\square)}{\sum \exp(\square)}$$

(dot product)

attention for = word or token "$i$"

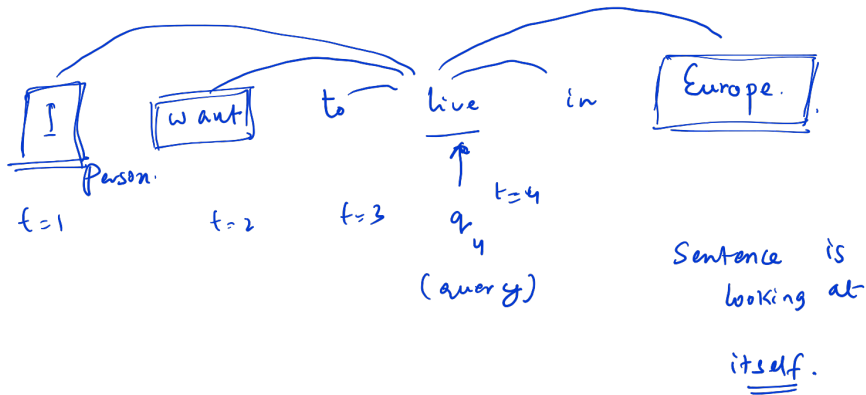$$\sum_i \frac{\exp(q \cdot k_i)}{\sum_j \exp(q \cdot k_j)} v_i$$

$A(q_i, K, V)$

query at $i$

keys values

$\leftarrow$ scalar

$b_i$ $v_i$ $\leftarrow$ vector

"self"

$$A(q_i, K, V) = \sum_i \left[ \frac{\exp(q_i \cdot k_i)}{\sum_j \exp(q_i \cdot k_j)} \right] v_i$$

Softmax

I — Person.
want
to
live ↑ $q_4$ (query) $t=4$
in
Europe.

$t=1$   $t=2$   $t=3$

Sentence is looking at

itself.

$$A(q_4, K, V)$$ vector

probability $\downarrow p$

$\otimes$

$V_i \leftarrow$ vector

scalar

Softmax

$q_4$

$\rightarrow K_1$

scalar

scalar

$\boxed{q_4 \cdot k_1}$  $\boxed{q_4 \cdot k_2}$  $\boxed{q_3 \cdot k_3}$  $\boxed{q_4 \cdot k_4}$  $\cdots$  $\boxed{q_4 \cdot k_6}$

$(q_1, K_1, v_i)$  $(q_2, k_3, v_2)$  $(q_3, k_6, v_4)$  $(q_6, K_6, v_6)$

$X_1$  $X_2$  $X_3$  $X_4$  $X_5$  $X_6$

I  want  to  live  in  Europe

In what context is "live" used?

$\rightarrow$ associate each word with 3-tuple $\langle q, k, v \rangle$



Query      Key      Value

$t=1$ I    $q_1$      $k_1$      $v_1$

$t=2$ want    $q_2$      $k_2$      $v_2$

$q_4$      $k_4$      $v_4$ $= A(q_4, k, v)$

content specific questions

$t=6$ Europe    $q_6$      $k_6$      $v_6$

for every $t = 1 \ldots T_x$ :

$\qquad$ compute $A(q_t, K, V)$

$$\downarrow$$

$\qquad$ "rich" representations.

word $t=1$ $\qquad$ [ $\qquad$ $A(q_1, K, V)$ ]

$\qquad$ $t=2$ $\qquad$ [ $\qquad$ $A(q_2, K, V)$ ]

$\qquad$ $\vdots$

$\qquad$ $t=T_x$ $\qquad$ [ $\qquad$ $A(q_6, K, V)$ ]

Let $x_1$: be the word-embedding for the word

"live"

(one-hot representation)

then

$q_1$ : $W_Q$ $\times$ $x_1$

$k_1$ : $W_k$ $\times$ $x_1$

$v_1$ : $W_v$ $\times$ $x_v$

$\searrow$ parameters of the model

In vectorcsed form:

$$\text{attention}\,(Q, K, V) = \text{softmax}\left(\frac{Q K^T}{\sqrt{d_k}}\right) V$$

$q \cdot k_i$

used to
control explotion
of the numerator

(Scaling parameter)

Scaled dot-product
attention.

(Attention is all you
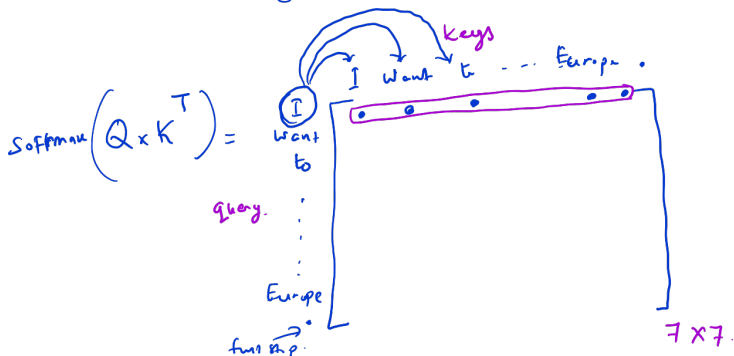need.)

example:    I    want    to    live    in    Europe.

sentence   length = 7

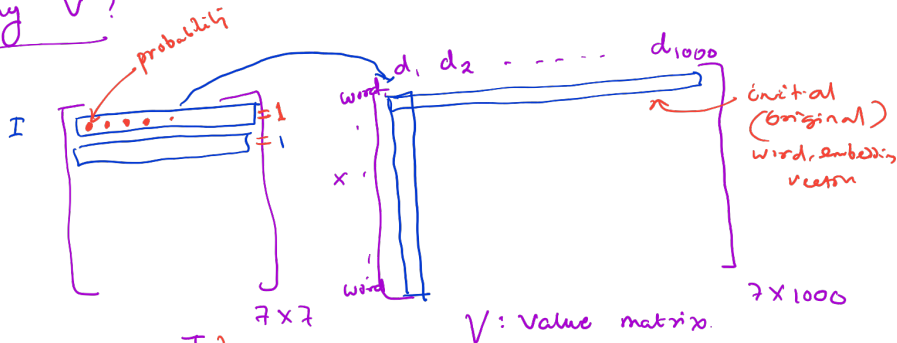word-emedding   dimension = 1000

Size  of  Q, K, V  = 7 × 1000

$$\text{softmax}\left(Q \times K^T\right) =$$



word dictionary

Keys

I  Want  to  --- Europe .

query.

Europe

first step.

probabilty matrix

(attention score matrix)

7 × 7.

why V ?

probability



$d_1$, $d_2$ - - - - - - $d_{1000}$

initial
(original)
word-embedding
vector

word

I

= 1
= 1

x

word

7 × 7

7 × 1000

V : Value matrix.

$\text{softmax}(Q \times K^T)$

self attention scores

=

I
want
.
.
Europe

new embedding
vector for
the word "I"
(rich)

7 × 1000

adjusted token/word
embeddings.

$\text{sotmax}(Q K^T) \cdot V$