Self attention          computer X

NLP → words → numbers → vectorization

OHE    1) mat cat mat
       2) cat rat rat

|       | mat | cat | rat |
|-------|-----|-----|-----|
| mat   | 1   | 0   | 0   |
| cat   | 0   | 1   | 0   |
| rat   | 0   | 0   | 1   |

num → [1 0 0] [0 1 0] [1 0 0]

BOW

|       | mat | rat | cat |
|-------|-----|-----|-----|
| S1    | [2  | 0   | 1]  |
| S2    | [0  | 2   | 1]  |

Tf Idf → Word embedding
         number → Semantic meaning

training wiki → → 5 dim

cricket [0.2  0.9  0.9  0.90 ...] n-dim vector    256, 512
                                                  64

king → [0.9  0.1  1  0  0.9]
                                → similar
queen → [0.9  0.2  0.4  1  0]
similar
royalty

meaning → vector  tech

Apple → vector

The problem of "Average Meaning"    → dataset

1) An (apple) a day keeps the doctor away
2) Apple is healthy
3) Apple is better than orange        10000
4) Apple makes great phones           9000   1000
   .                                  fruit  phone
   .        9000   1000
            tech   fru                [X  Y]
                                       ↑   ↑
                                    (taste) (technology)

→ [:] 2 dim

[0.8  0.2] → [0.9  0.3]

word embeddings
   create → use → static
                          smart
   static embedding → [contextual embedding]

NLP → trans eng-hindi

NLP → trans eng-hindi

fruit → [ Apple launched a new Phone
        while I was eating an orange ]

static → [ 0.9  0.3 ]
         [ 0.3  0.8 ]

word embedding
[static] → self attention
              ↓
           mechanism

$e_{apple}$
$e_{launch}$
$e_{phone}$
$e_{orange}$

Calculate ? How?

self attention → $Y_{apple}$
                 $Y_{am}$ → transformers
                 $Y_{man}$

## Embeddings

Hello → [ 0.3 | 0.9 | 1 | 0 ]  → OHE
                                  BOW
                                  Embeddings

↑ semantic

Money Bank                    static →        River Bank

[ 0.6 | 0.2 | 0.1 | 0.7 ]     [ 0.6 | 0.2 | 0.1 | 0.7 ]

Contextual embedding
  ↳ dynamic

Self Attention
embedding → dynamic
            contextual

Humans    Love    Smartphones

$e_1$    $e_2$    $e_3$

calculations    Self Attention    → ⌣ ?

$y_1$    $y_2$    $y_3$

$\boxed{\text{money} \quad \text{b}\underline{a}\text{nk} \quad \text{grows}}$

$\boxed{\text{ri}\underline{v}\text{er} \quad \text{ba}\underline{n}\text{k} \quad \text{flows}}$

$\text{bank} = \boxed{0.5 \text{river} + 9.4 \text{bank} + 0.1 \underline{\text{flows}}}$

$\underline{\text{bank}} \rightarrow \underline{\text{bank}}$

$\underline{\text{bank}} \rightarrow \boxed{0.3 \text{money} + 0.7 \underline{\text{bank}} + 0.1 \text{grows}}$

$S1$    Word embe

$\underline{\text{money}} = 0.7 \text{money} + 0.2 \text{bank} + 0.1 \text{grows}$

$\underline{\text{bank}} = 0.25 \text{money} + 0.7 \text{bank} + 0.05 \text{grows}$

$\underline{\text{grows}} = 0.1 \text{money} + 0.2 \text{bank} + 0.7 \text{grows}$

$S2$

$\text{river} = 0.8 \text{river} + 0.15 \text{bank} + 0.05 \text{flows}$

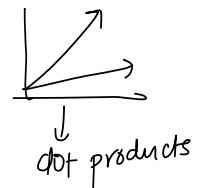$\text{bank} = 0.2 \text{river} + 0.78 \text{bank} + 0.02 \text{flows}$

$\text{flows} = 0.4 \text{river} + 0.01 \text{bank} + 0.59 \text{flows}$

$n \, dim$    $n \, dim \rightarrow$

$$e_{\underline{money}}^{(new)} = \boxed{0.7} e_{money} + \boxed{0.2} \, e_{bank} + \boxed{0.1} e_{grows}$$

$$e_{\underline{bank}}^{(new)} = \boxed{0.25} e_{money} + \boxed{0.7} e_{bank} + \boxed{0.05} e_{grows}$$

$$e_{\underline{grows}}^{(new)} = \boxed{0.1} e_{money} + \boxed{0.2} e_{bank} + \boxed{0.7} e_{grows}$$

$\textcircled{1}$    normalised

similarity

$e_{money} \quad e_{money}$

$e_{money} \quad e_{bank}$

$e_{money} \quad e_{grow}$

$e_{grows} \quad e_{money}$

$\downarrow$ dot products

$$e_{bank}^{(new)} = \left[ e_{bank} \cdot e_{money}^T \right] e_{money} + \left[ e_{bank} \cdot e_{bank}^T \right] e_{bank} + \left[ e_{bank} \cdot e_{grows}^T \right] e_{grows}$$

$S_{21}$    $e_{bank}$    $e_{money}$

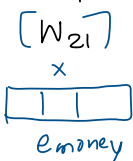$S_{22}$    $e_{bank}$    $e_{bank}$

$S_{23}$    $e_{bank}$    $e_{grows}$

$S_{21}$    $S_{22}$    $S_{23}$
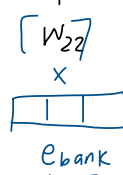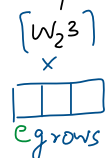
$$W_{22} = \frac{e^{S_{22}}}{e^{S_{21}} + e^{S_{22}} + e^{S_{23}}}$$

$$W_{21} = \frac{e^{S_{21}}}{e^{S_{21}} + e^{S_{22}} + e^{S_{23}}}$$

$\boxed{\text{Softmax}}$

$[W_{21}]$    $[W_{22}]$    $[W_{23}]$

$\times$    $\times$    $\times$

$e_{money}$    $+$    $e_{bank}$    $+$    $e_{grows}$

$e_{money}$ + $e_{bank}$ + $e_{grows}$

$$e_{bank}^{(new)} = \boxed{y_{bank}} \longrightarrow \text{contextual word embedding}$$

## Points to consider

→ This operation is a **parallel** operation

→ There are **no parameters** involved



W → data learn

suf model

nono → $e_1$ → $Y_1$

how are you

↓     ↓     ↓

$e_1$   $e_2$   $e_3$

Self mod

→ $y_1$   $y_2$   $y_3$

general contextual embedding

machine eng | hindi

$\boxed{\text{how are you}}$  कैसे हो

I am good  मैं वढ़िया हूँ

peia of cake → वेड़्द आसान काम

2 ढ़का

$\boxed{\text{piece of cake}}$ → $\boxed{\text{के का का इकड़ड़}}$

↓     ↓     ↓

$e_1$   $e_2$   $e_3$

↓     ↓     ↓

$y_1$   $y_2$   $y_3$

break a ug → दांगा तोड़ दो

X  task specific

Contextual embedding

# Progress

Words → machine

↓

embedding → semantic

↦ context

river <u>bank</u>    money <u>bank</u>

→ contextual embeddings

[ Simple self model ] → parameters

learnable parameters → [ general ]     task-specific

bank

$q_{bank}$
$k_{bank}$
$v_{bank}$

$e_{bank}$

✓ Query $e_{bank}$

✓ Key $e_{bank}$

✓ Value $e_{bank}$

Separation of Concerns

ebank    s21    ebank    s22    ebank    s23

emoney    ebank    egrows

s21    s22    s23

**Softmax**

w21    w22    w23

w21 × emoney    +    w22 × ebank    +    w23 × egrows

value

data

ybank

Author → (35) → jeevansaathi

→ 1) Profile → Key
→ 2) Search    query →        data
→ 3) match → value

$e_{bank}$

Profile    search    match

initial
embe  ⌐ quer
      ⌐ key
       value

$e_{bank}$ → $a_{bank}$
→ $k_{bank}$
→ $v_{bank}$        data

d magnitude (scaling)

→ linear transform

query
Key    $Q_{bank}$
value

→ [ ]

training

random /weigh value → $e_{bank}$ ←

$W^q$        $W^K$        $W^r$

data → machine transl

transl

$Q_{bank}$        $K_{bank}$        $V_{bank}$

$3 \times 3 \longrightarrow dot \longrightarrow (9)$ values

money bank grows
embeddings

$W_q$  $W_k$  $W_v$
parameter

word vec
512, 256, 64

dimn

money ✓  3 dim

(1,3)

(1,3)  (3,3)
(1,3)

$W_q$  $W_k$  $W_v$

[3×3]  [3×3]  [3×3]

(1,3)  (1,5)  (1,3)

qmoney  kmoney ✓  vmoney

3/10/5 e

$d_k \rightarrow$ dim of the
k vector
key

$d_k = d_q = d_v = 3$

bank ✓

$W_q$  $W_k$  $W_v$

qbank  kbank  vbank

grows ✓

$W_q$  $W_k$  $W_v$

qgrows  kgrows  vgrows

qmoney  kmoney  vmoney
qbank   kbank   vbank
qgrows  kgrows  vgrows

(3)  (2)  (9) dot → small
                    → high

(1,3)   9 vector-vector dot pro

(1,3)   Q   (3 vector)

512
dim

$K^T$   dot   9 numbers   Softmax   contextual embeddings

| w11 | w12 | w13 |
| w21 | w22 | w23 |
| w31 | w32 | w33 |

3×n

|| || ||

n×3  3 vectors

$\frac{1}{\sqrt{d_k}}$   $\mu$   $\sigma^2$

3×3   var(X)

| v11 | v12 | v13 |
| v21 | v22 | v23 |
| v31 | v32 | v33 |

3×3

V

3×n

money
bank
grows

3×n

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{Q K^T}{\sqrt{3}} \sqrt{d_k}\right) V \rightarrow \text{summary}$$

scaled dot product

$\frac{Q K^T}{\sqrt{d_k}} \rightarrow$ Unstable gradient $\rightarrow \sqrt{d_k}$

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

$[Q K^T]$
matrix  matrix
dot product
vectors dot product

$\text{Softmax}\left(\frac{Q K^T}{\sqrt{3}}\right) V$   $\frac{1}{\sqrt{d_k}}$ scale ?   Why ?

$\frac{1}{\sqrt{d_k}}$

Dot-product ka Nature

1) $[1, 2, 3] - [3, 2, 1]$   $\begin{bmatrix} f \\ g \\ h \\ i \\ j \end{bmatrix} \sigma_2^2$
2) [   ] - [   ]
3) [   ] - [   ]
4) [   ] - [   ]
5) [   ] - [   ]

low dimensional $\rightarrow$ dot product $\rightarrow$ low variance
high dim $\rightarrow$ dot product $\rightarrow$ high variance

1) $[1, 2] - [3, 2]$   $\begin{bmatrix} a \\ b \\ c \\ d \\ e \end{bmatrix} \sigma_1^2$
2) [   ] - [   ]
3) [   ] - [   ]
4) [   ] - [   ]
5) [   ] - [   ]

code

3 dim

1000 pair of vectors

1) — × — →
2) — × — →
⋮
1000) — × — →



100 dim

1000 pairs of vector

1) — × —
2) — × —
⋮
1000) — × —



1000 dim

1000 pair of vectors

1) — × —
2) — × —
⋮

vector

dim $\rightarrow$ variance
$\hookrightarrow$ high

high variance
$\downarrow$
it's a problem

variance
learn

prob
$\hookrightarrow$ 100/2
99%
1%

backprop

softmax

ex

train

variance

512 dim V
3 dim X

low dim

low

variance high

high

low

ignore

param
$\hookrightarrow$ vanish
$\downarrow$
upload X

grad

vanish

training stable

cach c

$$\begin{cases} X \rightarrow Var(x) \\ y \begin{array}{l} \rightarrow cX \\ \hookrightarrow c^2 Var(v) \end{array} \end{cases}$$

If you have a random variable $X$ with a variance of $Var(X)$, and you create a new variable $Y$ by scaling $X$ with a constant $c$, so that $Y = cX$, the variance of $Y$ ($Var(Y)$) is related to the variance of $X$ by the square of the scaling factor $c$. Mathematically, this relationship is expressed as:

$$Var(Y) = c^2 Var(X)$$

$$\begin{bmatrix} V_1 \odot V_4 \\ V_1 \odot V_5 \end{bmatrix} \quad \dim \uparrow \quad \rightarrow \boxed{\text{main quantify}}$$
$\hookrightarrow$ variance $\uparrow$

var

softmax

1 dim
$[a]$ ─ 3

$[b]$ 2 ── $\boxed{ab|6}$ 6 $\rightarrow$ random   var $\rightarrow$ X
$[c]$ 1 ── $\boxed{ac(3)}$ 3 $\rightarrow$ variance (sample)
$[d]$ 4 ── $\boxed{ad|12}$ 12 $\rightarrow$ expected var    $\rightarrow \boxed{Var(y)}$

$V_1$ ─ $V_4$ / $V_5$ / $V_6$

2 dim

$(X) \rightarrow var(x)$
$cX \rightarrow c^2 Var(x)$

$\begin{bmatrix} a & b \end{bmatrix} \quad \begin{bmatrix} c & d \end{bmatrix} \rightarrow \boxed{\dfrac{ag+bd}{\sqrt{2}}}$
$\begin{bmatrix} e & f \end{bmatrix} \rightarrow \boxed{\dfrac{ae+bf}{\sqrt{2}}}$  $\rightarrow Y \rightarrow Var(v)$
$\begin{bmatrix} g & h \end{bmatrix} \rightarrow \boxed{\dfrac{ag+bh}{\sqrt{2}}} \rightarrow Var(X)$

$y \rightarrow Var(Y) \rightarrow 2Var(x)$
$\dfrac{y}{\sqrt{2}} \rightarrow \dfrac{1}{2} Var(y) \rightarrow \dfrac{1}{2} 2 Var(x) \rightarrow Var(x)$

2 Var(X)

$Var(y) \simeq 2 Var(x)$

$Var(y) > Var(X)$

Yes   No

3 dim

$\begin{bmatrix} a & b & c \end{bmatrix} \quad \begin{bmatrix} d & e & f \end{bmatrix} \rightarrow \boxed{\dfrac{ad+be+cf}{\sqrt{3}}}$
$\begin{bmatrix} g & h & i \end{bmatrix} \rightarrow \boxed{\dfrac{ag+bh+ci}{\sqrt{3}}}$  $\rightarrow Z$
$\begin{bmatrix} K & l & m \end{bmatrix} \rightarrow \boxed{\dfrac{ac+bl+cm}{\sqrt{3}}} \rightarrow Var(X)$

$\dfrac{1}{\sqrt{K}}$

$Var(z) > Var(y) > Var(x)$

$Var(z) \simeq 3 Var(x)$

$z \rightarrow Var(z) \simeq 3 Var(y)$
$\dfrac{1}{\sqrt{3}} z \rightarrow (\dfrac{1}{\sqrt{3}})^2 Var(z) = \dfrac{1}{3} 3 Var$

d dim $\rightarrow$ $d \underline{Var(x)}$
$\hookrightarrow$ varina dim = 1

dim $\rightarrow$ variance
linear

$$\begin{cases} 1 \text{ dim} \rightarrow Var(x) \longrightarrow Var(x) \\ 2 \text{ dim} \rightarrow 2 Var(x) \longrightarrow Var(x) \\ 3 \text{ dim} \rightarrow 3 Var(x) \longrightarrow Var(x) \\ \vdots \\ d \text{ dim} \rightarrow d Var(x) \longrightarrow Var(x) \end{cases}$$

$[\ \underline{\quad} \ - - - - ] \rightarrow Var(X) \rightarrow$
  $\uparrow$    d dim

a

Q

K

| s11 | s12 | s13 |
| s21 | s22 | s23 |
| s31 | s32 | s33 |

3×3

Scale

1/sqrt(dk)

| s11 | s12 | s13 |
| s21 | s22 | s23 |
| s31 | s32 | s33 |

3×3

Softmax

| w11 | w12 | w13 |
| w21 | w22 | w23 |
| w31 | w32 | w33 |

3×3

3×n

3×n

money
phone
grass

3×n

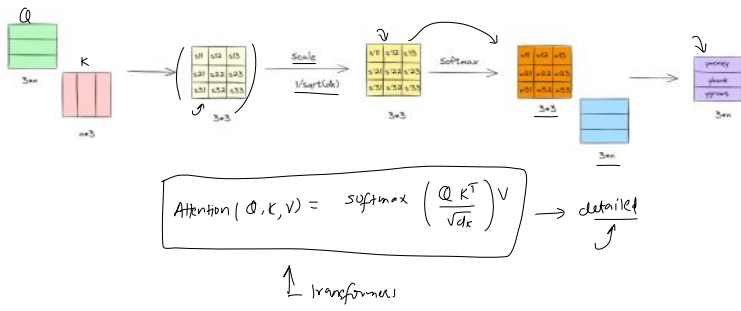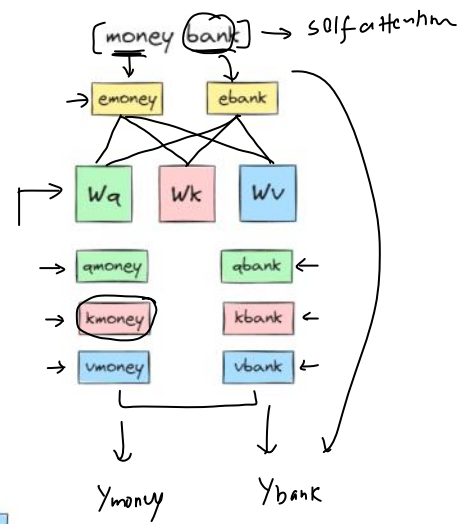$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q\,K^T}{\sqrt{d_k}}\right)V \longrightarrow \text{detailed}$$

↑ transformers

# What is $d_k$

28 February 2024      16:59

# Recep

To move canvas, hold mouse wheel or spacebar while dragging, or use the hand tool

bank → contextual

qmoney
s11

qmoney
s12

qbank
s21

qbank
s22

[money bank] → self-attention

kmoney
s11

kbank
s12

kmoney
s21

kbank
s22

emoney    ebank

Scaling (1/sqrt(dk))

Scaling (1/sqrt(dk))

Wq    Wk    Wv

s11'  s12'

s21'  s22'

qmoney    qbank

Softmax

Softmax

kmoney    kbank

w11    w12

w21    w22

vmoney    vbank

w11 × vmoney    +    w12 × vbank

w21 × vmoney    +    w22 × vbank

Ymoney    Ybank

ymoney

ybank

money bank

embedding
word2vec

emoney    ebank

$[7,7]$    $[9,3]$

Wq → Wk → Wv

vector
·
matrix
linear
transform

qmoney    qbank ✓
kmoney    kbank ✓
vmoney    vbank ✓

3 new
vector

$$Wq = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad 2\times2$$

$$W_K = \begin{bmatrix} 3 & 4 \\ 5 & 1 \end{bmatrix} \quad 2\times2$$

$$W_V = \begin{bmatrix} 4 & 1 \\ 2 & 1 \end{bmatrix} \quad 2\times2$$

✗ All values are hypothetical



qmoney    s11        qmoney    s12
kmoney              kbank
s11                 s12

Scaling (1/sqrt(dk))

s11'    s12'

Softmax

w11        w12

w11 × vmoney    +    w12 × vbank

ymoney

[Dot Product]

$S_{21} = \boxed{10}$       $S_{22} = \dfrac{\boxed{32}}{\sqrt{2}}$

[Scaling]

$S_{21}' = \boxed{\dfrac{10}{\sqrt{2}}} = 7.09$    $S_{22}' = 22.69$

[Softmax]       $\boxed{1}$

$W_{21} = \boxed{0.2}$    $W_{22} = \boxed{0.8}$

qbank    s21        qbank    s22
kmoney              kbank
s21                 s22

$\dfrac{1}{\sqrt{2}}$

$\dfrac{1}{\sqrt{d_K}}$ →   Scaling (1/sqrt(dk))

$d_K = 2$           s21'    s22'

scaling            Softmax

w21        w22

$V_{money}$  $V_{bank}$

$V_{money}$  $V_{bank}$

$0.8\ V_{bank}$

$0.2\ V_{money}$

scaling

$w21$  $w22$

$0.2$  $w21 \times$  $V_{money}$  $+$  $0.8$  $w22 \times$  $V_{bank}$

$Y_{bank}$

money    bank

self attention
$\rightarrow$ gravity

money

bank

river

$V_{money}$  $V_{bank}$
$Y_{bank}$
$0.8\ V_{bank}$
$0.2\ V_{money}$

self attention

$e_{river}$  $Y_{bank}$

$e_{bank}$

$e_{money}$  $Y_{bank}$

$\rightarrow$ dataset

pull

$e_{bank}$

river bank

long trans → NN

encoder-decoder

Attention

decoder  लाइट  बंद  करो  <end>



<start>  लाइट  बंद  करो

h0  h1  h2  h3  h4

turn  off  the  lights  → 30 words

Set of numbers

Context vector

encoder

30 words

$S_1 h_1$

$Q_{11}$    $S_1 h_2$

$e_{12}$    $S_1 h_3$    $e_{14}$  $S_1 h_4$

$e_{13}$

$t_i = i$

hindi

$C_1$  बंद  $C_2$  <end>

$S_0$  1  $S_1$  2  $S_2$  3  $S_3$  4  $S_4$

<start>  लाइट  बंद  करो

Attention

inter sequence attention

$\alpha_{11}$  $\alpha_{12}$  $\alpha_{13}$  $\alpha_{14}$

$t = j$

h0  h1  h2  h3  h4

english  turn  off  the  lights

$c_i = \sum \alpha_{ij} h_j$ ✓

$\alpha_{ij} = softmax(e_{ij})$

$\boxed{e_{ij} = S_i^+ h_j}$

$S_i \longrightarrow$ query

$h_j \longrightarrow$ key

$h_j \longrightarrow$ value

self attention

$y_{turn}$  Turn  off  the  lights

$e_1$    $e_2$  $e_3$  $e_4$

$softmax(s_{11})$

$s_{12}/s_{13}, s_{14}$

$q_{turn}$      $v_{turn}$    $y_{turn} = W_{11} v_{turn} + W_{12} v_{off} +$
$q_{turn}$  $k_{turn}$  $W_{13} v_{the} + W_{14} v_{lights}$

query  Turn  off  the  lights
      $q_{turn}$  $q_{off}$  $q_{the}$  $q_{lights}$

$y_{off}$

$s_{21}$    $s_{22}$  $s_{23}$    $s_{24}$

Turn  off  the  lights
$k_{turn}$  $k_{off}$  $k_{the}$  $k_{lights}$

self

$y_{off}$

$\longrightarrow softmax$ $w_{21}$ $w_{22}$ $w_{23}$ $w_{24}$

intra sequence
↓
self

$c_i = \sum \alpha_{ij} h_j$    → weights

$\alpha_{ij} = softmax(e_{ij})$    → alignment score    → Luong attention

$\boxed{e_{ij} = S_i^T h_j}$

self attention