# Review

1. No sequentiality

2. Introducing Sequentiality.

3. Introducing Attention

# Part 1

Ler1 → lect3

XETRD

- binary classificate
- Multiclass classificat
- Regression
-
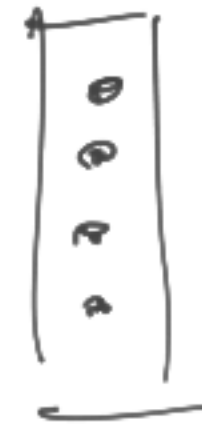
PS3

→ Preprocessing ← one hot encoding ✗
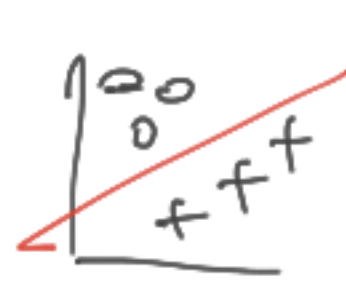
Scaling

→ Feature engineering

→ Feature Importance ⟶ prediction

~s __Models__

Tree based Models

$$\mathrel{\underset{\circ\,\,\circ\circ}{\diagup\!\!\!\!\!\diagdown}}_{+\,+\,+} \Rightarrow \mathrel{\underset{\circ\,\,\circ\,\,\circ}{\diagup\!\!\!\!\!\diagdown}}^{\circ\,\circ\,\circ}_{(+\!+\!+)\circ}$$

NN ( $\alpha$ )

---

__①F__

→ # trees

→ Depth of trees

→ # features / split

→ Impurity measure

⎫ HP

→ # Layers

→ # neurons / Layer

→ Act. function

→ Train →> lr
→ batch size
→ # epochs

⎫ HP

# HP optimization for RF (PS3)

Train | Val

Train ⟶ ⟶ θh

∴) h = (100 trees, Entropy, 200 depth)

∴) evaluation metrics (d1)
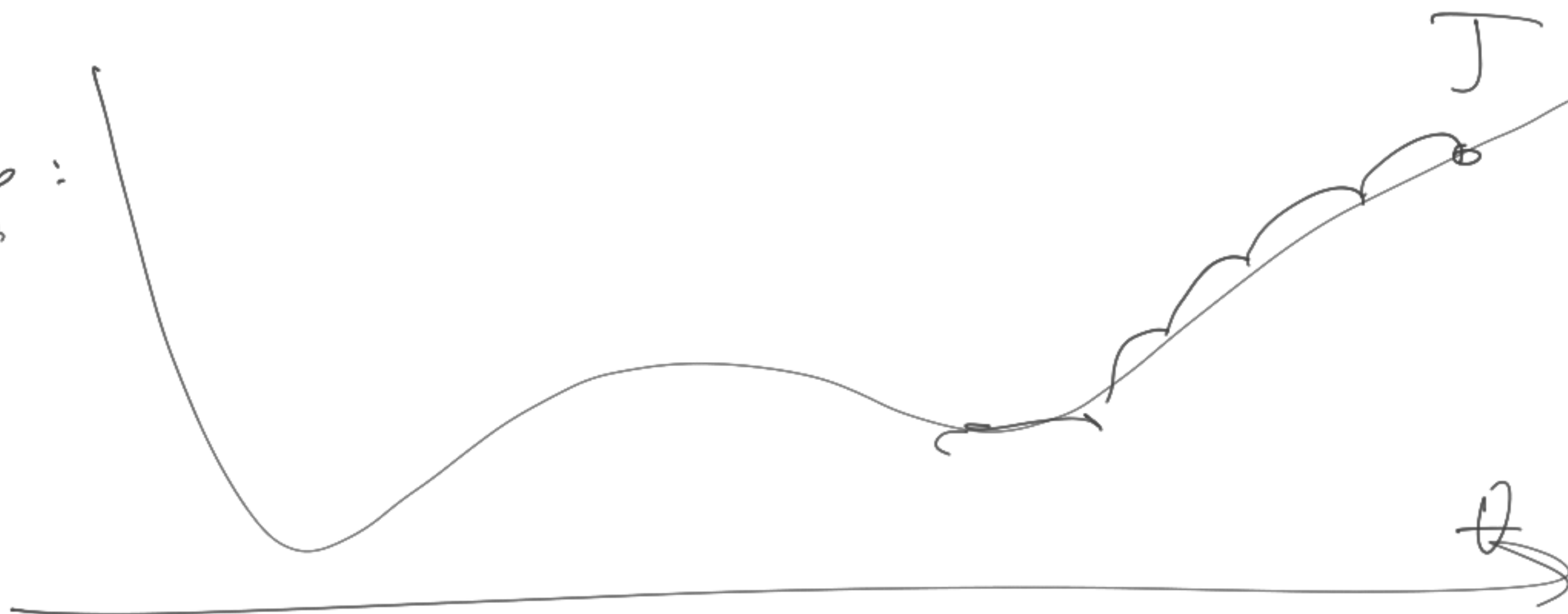
Training ⓛ2

based on

max IG

$$\boxed{NN}$$



$$z^2 = \left( W^T z^1 + b \right)$$

$$W \in \mathbb{R}^{d_1, d_2}, \quad b \in \mathbb{R}^{d_2}$$

Traing :



$J$

$\theta$

# Cross validation
## (small datasets)

$h = (100 \text{ trees}, \text{Entropy}, 200 \text{ depth})$

$\boxed{P_h} = \text{mean}_{1 \leq k \leq 5} P_h^k$

# Part 2 : Introducing Sequentiality:

(PSK) Introducing sequential data : Imdb dataset.

Dataset

doc$_1$ : $W_1^1$ ———— $W_1^{T_1}$ | $y_1$

$\vdots$

doc$_N$  $W_N^1$ $\vdots$ — — $W_N^{T_N}$ | $y_N$

$V$-dim vector

Dense

$K$ categories

2 limitations

$\rightarrow$ non sequentiality

$\rightarrow$ poor representation of words (more indices)

Lecture5

Dataset

$doc_1$ ————————————— } → word_index = $\{$ 'a' : 1

$\vdots$

$\vdots$

$doc_N$ —————————————

'zoo' : V $\}$

objective :  $E = \begin{bmatrix} \text{———— } W_1 \text{ ———} \\ \vdots \\ \text{———— } W_i \text{ ———} \\ \vdots \\ \text{— } W_V \text{ ———} \end{bmatrix}_{V \times D}$  : $W_i$ : D dim $\vec{r}$ representation word of index $i$

GloVe :

Dataset $\Longrightarrow$ $X$ : co-occurence $\Longrightarrow$ $\overset{\text{"target"}}{\log X_{ij}} \sim \overset{\text{prediction}}{W_i^T \tilde{W_j}}$
matrix

$X_{ij}$ : # times word of index $j$ is in the content of word of index $i$

Loss : $J_\theta = \sum_i \sum_j \left( \log X_{ij} - W_i^T \tilde{W_j} \right)^2$

$\theta = \begin{bmatrix} W_n \\ \vdots \\ \tilde{W} \\ \tilde{W} \\ \vdots \\ \tilde{W} \end{bmatrix}$

# Embedding Layer

Parameters : embedding matrix : $\mathcal{E} = \begin{bmatrix} \text{---} W_1 \text{---} \\ \vdots \\ \vdots \\ \text{---} W_V \text{---} \end{bmatrix}_{V \times D}$

$doc_i : [(17), (28) \cdots (47)] \longrightarrow \boxed{\begin{array}{c} \text{Embedd.} \\ \text{Layer} \end{array}} \longrightarrow$

$\underset{T}{\longleftrightarrow}$

$\begin{array}{c} \text{---} W_{17} \text{---} \\ \text{---} W_{28} \text{---} \\ \vdots \\ \text{---} W_{47} \text{---} \end{array}$

$T \quad \underset{D}{\longleftrightarrow}$

$(N, T) \longrightarrow (N, T, D)$

# Lecture 6  Introducing Seq Models : LSTM.

$h_i^1$

$h_i^0$ → [LSTM] → [LSTM] →

$X_i^1$

$h_i^2$

$X_i^2$ . . . . . . .

$h_i^T$  $\updownarrow d$

→ [LSTM] →

$X_i^T$  $\updownarrow D$

$(N, T, D)$ ──────────→ $(N, T, d)$

if I only keep $\boxed{h_i^T}$  $(N, T, D)$ ──────────→ $(N, d)$

$$\boxed{PSA}$$

$P_i \quad \bigsqcup^{\uparrow}_{\downarrow} K$

$\boxed{Dense}$

$\overbrace{\cancel{h_i^T}} \quad X \quad X \quad \left(\!h_n^T\!\right)$

$\boxed{LSTM}$

$X_i^n \quad\quad X_i^T$

$\boxed{Emb\ Layer} \longrightarrow \sim 2\ million\ params$

$: W_i \quad - - - - - \quad W_i$

<span style="color:red">⚠</span>

<span style="color:red">using preetrained
word vectors
reduced drastically
# parameters
to train</span>

(28) From may to one (PSJ) $\longrightarrow$ Many to many.

NLP : NMT

= France

Tom a été entarté

$V_y$ [ ] [ ] [ ] [ ]

[LSTM] [LSTM] [LSTM] [LSTM]

$h_i^1$

$h_i^{T_x}$

Tom was hit with a pie

$Y^1 \; - - - \; Y^{T_y}$

$\Uparrow$

$\overline{X^1}$ $- \; - \; -$ $X^{T_x}$

$T_y$

$T_x$

$t$

# Alignment matrix:

|  | Tom | was | hit | with | a | pie |
|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | $t$ | $T_x-1$ | $T_x$ |

$\longrightarrow$ Input Sequence $T_x$.

Tom ⊘

entarté'ty ⊘ ⊘ ⊘ ⊘

$\downarrow$

Output Sequence $T_y$

$\rightarrow T_x$

$\alpha_{i,y}$ } sum to 1

$T_y$

# Part 3 : Introducing the Attention Mechanism:

## Intuition :



title $k_1$
video $v_1$
$\vdots$
title $k_n$
video $v_n$

Database

query
Att. mec. $\rightarrow$ $q$

$A(q, \{k_i, v_i\}_{1 \leq i \leq n})$

$=$

$\sum_{i=1}^{n} \widehat{\alpha_i} \; v_i$

how relevant $v_i$ is to $q$

$\underline{\alpha_i ?} \quad \left( sim(q, k_i) \right)_{1 \leq i \leq n} \xrightarrow{softmax} \left( \alpha_i \right)_{1 \leq i \leq n}$

# 1st application:

Tom was hit with a pie $T_x$

$t_y = $ ...

$T_y$



$t$, $T_y$, $T_x$

Tom a été [entarté] (eos) $t_y$

what is the next word?

$\boxed{\phantom{x}} \xrightarrow{s_i^{t_y-1}}$ $LSTM_2$

$c_i^{t_y} = A(s_i^{t_y-1}, \{h_i^{t_x}\}_{t_x}^{t_x})$

$\langle t_y, 1 \rangle$
$\alpha_i$

$\langle t_y, T_x \rangle$
$\alpha_i$

Attention Layer

$h_i$ ... $h_i$ ... $h_i$ ... $h_i$ ... $h_i$ ... $h_i$
1    2    3    4    $T_x$

$LSTM$ $LSTM$ $LSTM$ $LSTM$ $LSTM$ $LSTM$

$x_i^1$ $x_i^2$ $x_i^3$ $x_i^4$ $x_i$ $x_i^{T_x}$
1    2    3    4

Tom was hit with a pie

⇨ The input to the Decoder LSTM in the sequence to sequence with Attention framework is:

$$c_i^{t_y} = \sum_{t_x=1}^{T_x} \underbrace{\alpha_i^{\langle t_y, t_x \rangle}}_{=} h_i^{t_x}$$

$$\frac{\exp\left(s_i^{t_y-1} \cdot h_i^{t_x}\right)}{\sum_{t_x'=1}^{T_x} \exp\left(s_i^{t_y-1} \cdot h_i^{t_x'}\right)}$$

2nd Application :   Self Attention layer.

:) Introduced to handle the polysemy problem

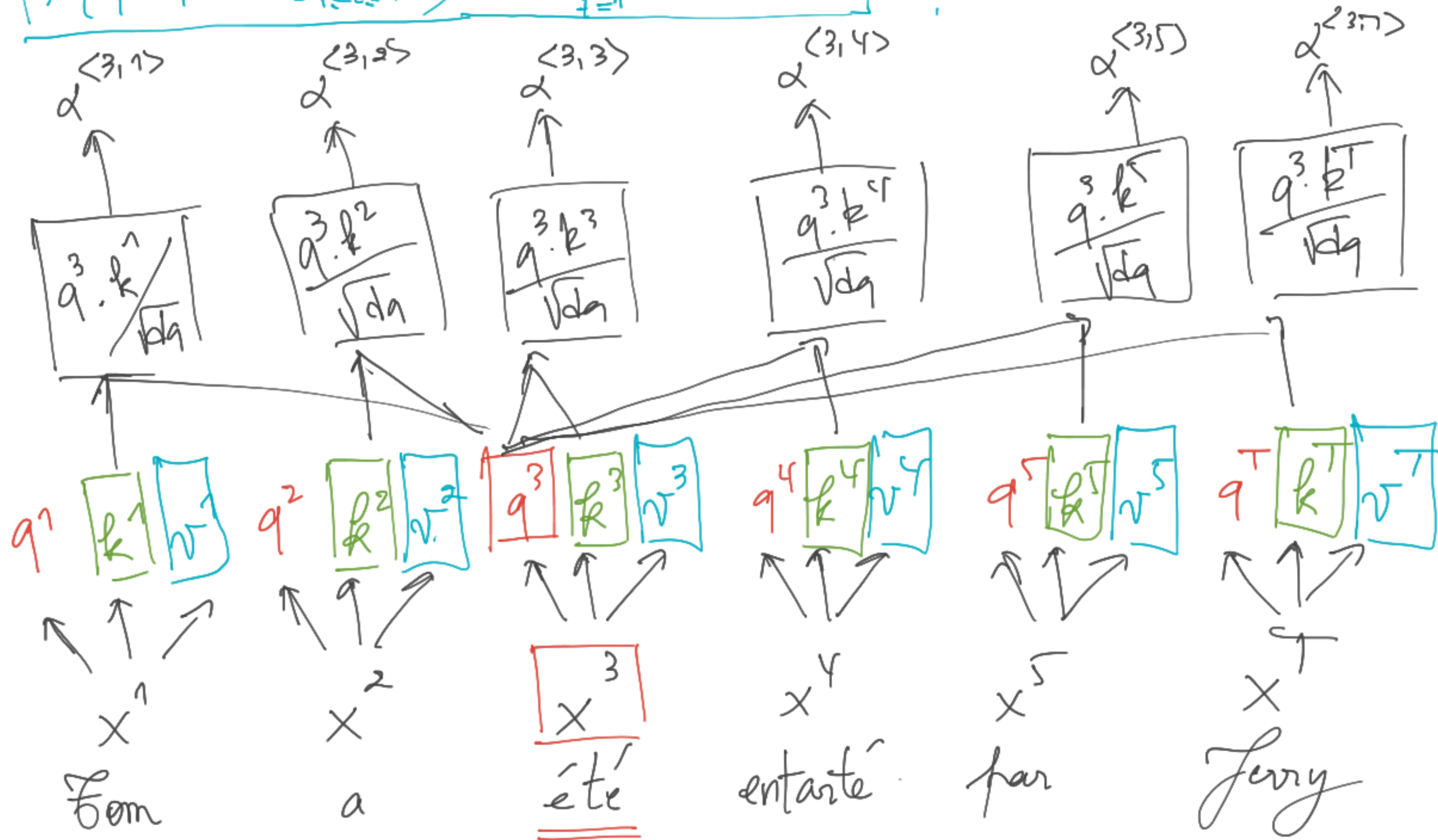→ Tom a été entarté par Jerry .

(i.e, Tom WAS hit with a pie by Jerry )

→ Cet été il fera horriblement chaud

(i.e, This summer will be unbearably hot )

:) GloVe approach : Static embedding vectors

⌐→ Need to create contextual embedding vectors

$$A\left(q^3, \{k^t, v^t\}_{1 \le t \le T}\right) = \sum_{t'=1}^{T} \alpha^{\langle 3, t' \rangle} v^{t'}$$

$\alpha^{\langle 3, 1 \rangle}$   $\alpha^{\langle 3, 2 \rangle}$   $\alpha^{\langle 3, 3 \rangle}$   $\alpha^{\langle 3, 4 \rangle}$   $\alpha^{\langle 3, 5 \rangle}$   $\alpha^{\langle 3, T \rangle}$

$\dfrac{q^3 \cdot k^1}{\sqrt{d_q}}$   $\dfrac{q^3 \cdot k^2}{\sqrt{d_q}}$   $\dfrac{q^3 \cdot k^3}{\sqrt{d_q}}$   $\dfrac{q^3 \cdot k^4}{\sqrt{d_q}}$   $\dfrac{q^3 \cdot k^5}{\sqrt{d_q}}$   $\dfrac{q^3 \cdot k^T}{\sqrt{d_q}}$

$q^1$ $k^1$ $v^1$    $q^2$ $k^2$ $v^2$    $q^3$ $k^3$ $v^3$    $q^4$ $k^4$ $v^4$    $q^5$ $k^5$ $v^5$    $q^T$ $k^T$ $v^T$

$x^1$   $x^2$   $x^3$   $x^4$   $x^5$   $x$

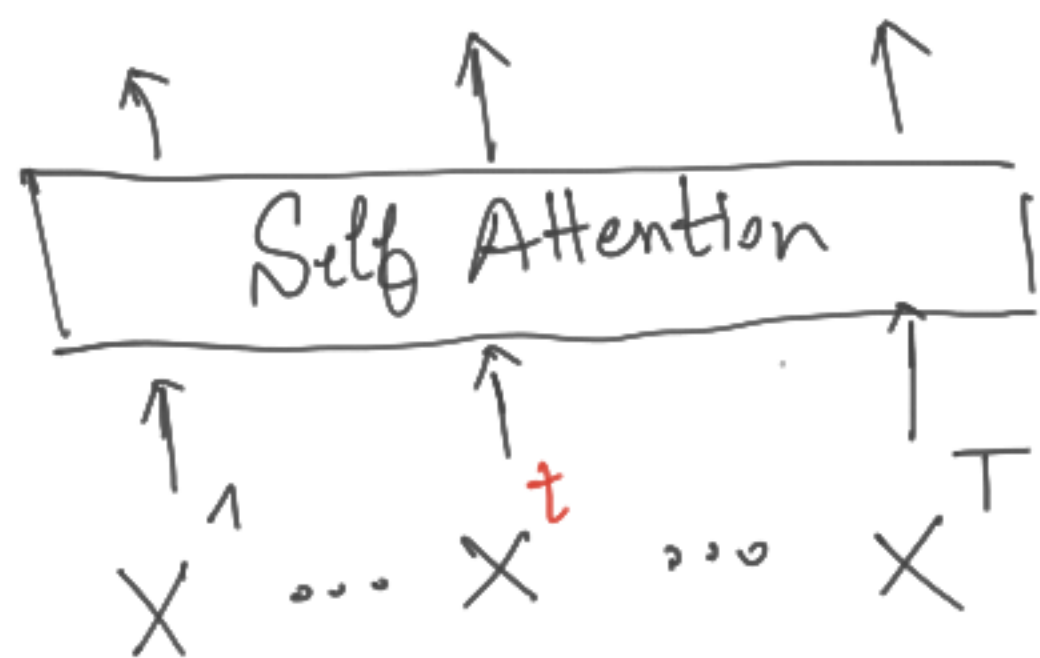Tom   a   été   entarté   par   Jerry

# The Self Attention Layer

**Parameters** $W_Q \in \mathbb{R}^{D \times d_q}$ ; $W_K \in \mathbb{R}^{D \times d_k}$ , $W_v \in \mathbb{R}^{D \times d_v}$

$$A(q^t, \{k^{t'}, v^{t'}\}_{1 \leq t' \leq T})$$

$$A(q^t, \{k^{t'}, v^{t'}\}_{1 \leq t' \leq T}) \text{ is the contextual representation of } x^t$$

$$A(q^t, \{k^{t'}, v^{t'}\}_{1 \leq t' \leq T})$$

$(N, T, d)$

$(N, T, D)$

Forward Propagation

$$\sum_{t'=1}^{T} \frac{\exp\left(\frac{q^t_{\cdot} k^{t'}}{\sqrt{d_q}}\right)}{\sum_{t''=1}^{T} \exp\left(\frac{q^t_{\cdot} k^{t''}}{\sqrt{d_q}}\right)} v^{t'}$$

Self Attention

$$X^1 \; \ldots \; X^t \; \ldots \; X^T$$

the end .