

# PART 3

2022-03-23

## Language Models

Example: What is the problem?

predict the next word

$i$	1	2	3	4
$x_i$	what	is	the	problem

$y_i$	is	the	problem	EOS
-------	----	-----	---------	-----

$$\star P(w_k | w_1, \dots, w_{k-1}) \quad \text{NLG} \quad P(w_i) = P(w_i / \text{Begin})$$

$$\star P(w_1, w_2, \dots, w_n) = P(w_1) P(w_2 | w_1) \dots P(w_n | w_1, \dots, w_{n-1})$$

prob of sentence

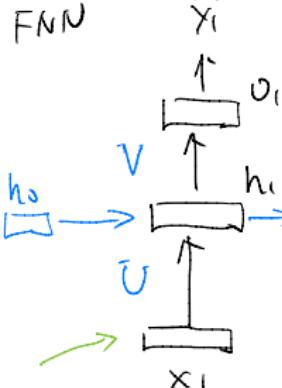
Dialogue:

- \* How are you?
- \* I am fine

i	1	2	3	4	5	6
X <sub>i</sub>	How	are	you	-	-	-

Y <sub>n</sub>	-	-	-	I	am	fine
----------------	---	---	---	---	----	------

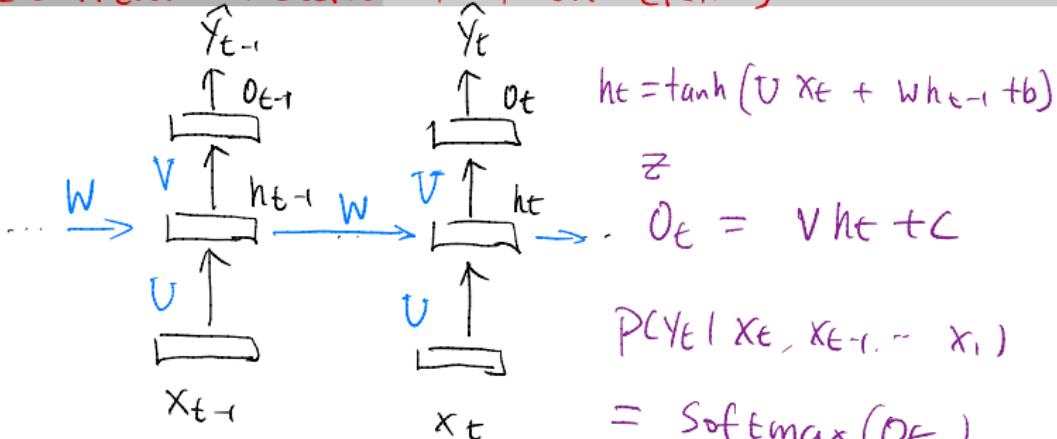
# Recurrent Neural Network (RNN)



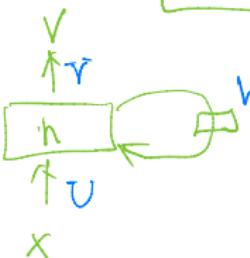
Embedding:

Vec  
Oem

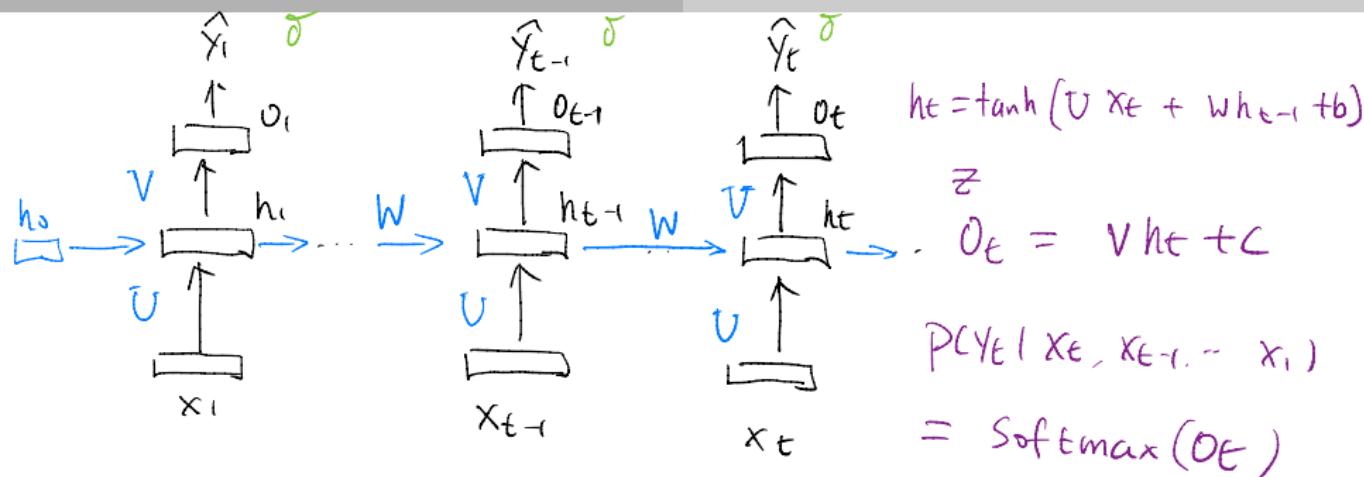
Circuit diagram



rolled-out model



Time delayed dependence



Parameters:  $W, U, V, b, c, \theta_{\text{em}}$  (embedding vectors)

$$L = E \left[ \frac{1}{T} \sum_{t=1}^T -\log P(Y_t | X_t, \dots, X_1) \right]$$

$\nabla_W L, \nabla_U L, \dots$  Average over pairs of training sequences

## LSTM

Standard RNN cell

$$a_t = Vx_t + Wh_{t-1} + b$$

$$h_t = \tanh(a_t)$$

LSTM cell

$$C_t = f_t \otimes C_{t-1} + i_t \otimes \tanh(c_{t-1})$$

$$\begin{bmatrix} 0.9 \\ 0.01 \\ \vdots \\ \vdots \end{bmatrix} \rightarrow \begin{bmatrix} - \\ - \\ \vdots \\ \vdots \end{bmatrix}$$

forget

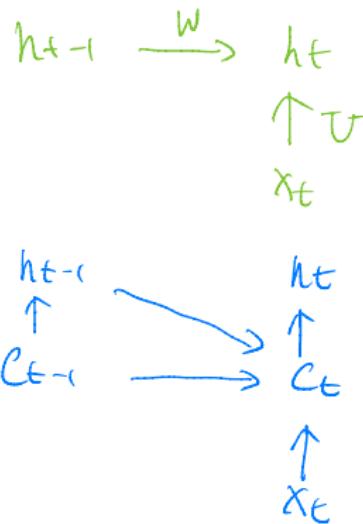
$$\begin{bmatrix} 0.01 \\ 0.9 \end{bmatrix} \rightarrow \begin{bmatrix} - \\ - \\ \vdots \\ \vdots \end{bmatrix}$$

Add to memory

Not add to memory

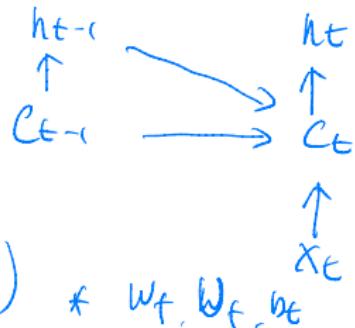
$f_t$ : forget gate

$i_t$ : input gate



## LSTM cell

$$C_t = f_t \otimes C_{t-1} + i_t \otimes \tanh(c_t)$$



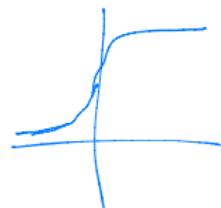
$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) * \text{learnable para}$$

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) * \text{sigmoid good ac gate}$$

$$h_t = O_t \otimes \tanh(C_t)$$

$O_t$ : output gate  $\begin{bmatrix} 0.9 \\ 0.01 \\ 0.7 \end{bmatrix}$

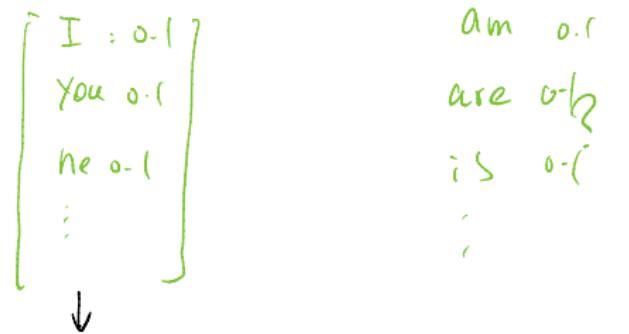
$$O_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$



# Encoder - Decoder / Seq 2 Seq Model Architecture

## Decoder

$s_1 \longrightarrow s_2$



$y_1 : I$

$\begin{bmatrix} \text{Am } 0.8 \\ \text{are } ? \\ \text{is } ? \\ \vdots \end{bmatrix}$

## Attention

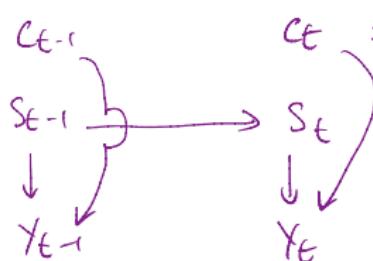
Encoder

$$x_1 \quad x_2 \cdots \quad x_j \cdots \\ \downarrow \quad \downarrow \quad \quad \quad \downarrow \\ h_1 \rightarrow h_2 \cdots \rightarrow h_j \cdots$$

$$\alpha_{tj} \leftarrow C_{T-t}^T \cdot h_j$$

$$\alpha_{tj} \leftarrow \frac{e^{\alpha_{tj}}}{\sum_j e^{\alpha_{tj}}}$$

Decoder



$$c_t = \sum_j \alpha_{tj} h_j : \alpha_{tj} - \text{how much attention to } h_j \text{ at } t \text{ (decoder)}$$

①  $c_{t-1}$  (what focused on at  $t-1$ )

Query

{ held  
drink

$c_t$  (what to focus on at  $t$ )

talk, meeting, cup, ...  
water, juice - -

② what info at  $j$ :  $h_j$

key

2022-03-25

$$z_1 \dots \underline{z_i} \dots$$

At  $z_i$ , how much attention to pay to  $x_j$ ?

$$x_1 \dots x_i \dots x_j \dots$$

$\uparrow$   
1x $d_m$   
512

Query:  $q_i = x_i W^Q \rightarrow d_m \times d_k$

$1 \times d_k$        $1 \times d_m$       ↑ learnable

key:  $k_j = x_j W^K \uparrow$

$$\alpha_{ij} = q_i k_j^T / \sqrt{d_k}$$

scaled ...

↑ Dot product attention

$$\alpha_{ij} \leftarrow \frac{e^{\alpha_{ij}}}{\sum_j e^{\alpha_{ij}}}$$

$$z_i = \sum_j \alpha_{ij} v_j$$

1x $d_v$ 

64

Value:  $v_j = x_j W^V \rightarrow d_m \times d_v \rightarrow 64$

1x $d_v$

Je

suis

etudiant(e)

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \quad Q = \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_N \end{bmatrix} \quad K = \begin{bmatrix} k_1 \\ k_2 \\ \vdots \\ k_N \end{bmatrix} \quad V = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_N \end{bmatrix} \quad Z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_N \end{bmatrix}$$

$N \times d_m \quad N \times d_K \quad N \times d_K \quad N \times d_V \quad N \times d_V$

$$Q = X W^Q$$

$$K = X W^K$$

$$V = X W^V$$

$$Z = \text{softmax}\left(\frac{Q K^T}{\sqrt{d_K}}\right) V$$

$N \times N$

Softmax



single-head Self-attention

$$X \xrightarrow{W_i^Q, W_i^K, W_i^V} Z$$

$N \times dm$                            $N \times d_v$

Multi-head self-attention

$$X \xrightarrow{W_1^Q, W_1^K, W_1^V} Z_1$$

:

$$Z \leftarrow Z W^O$$

$N \times dm$      $N \times h \times d_v$      $d_v \times dm$

$$h = 8$$

$$h = 16$$

$$\xrightarrow{W_h^Q, W_h^K, W_h^V} Z_h$$

$$Z = \text{cat}(Z_1, \dots, Z_h)$$

$$N \times \underline{h \times d_v}$$

$$dm = \underline{512}$$

$$\underline{512} = 8 \times 64$$

$$\underline{1024} = 16 \times 64$$

$\square$     $\square$

$\square$     $t$

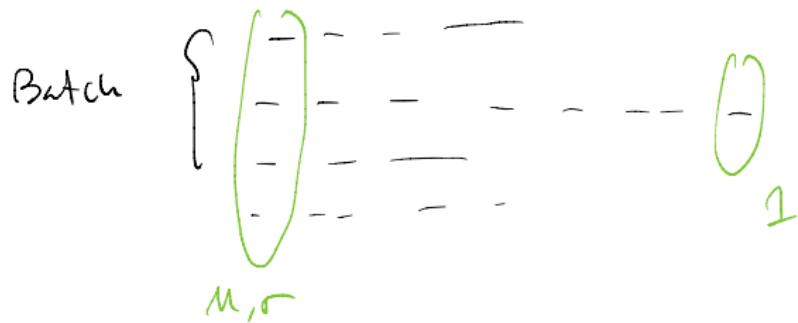
$\square$     $\square$  - - -    $\square$     $t$

Complexity :  $t^2$

512  $\times$  512

# Why Not Batch Normalization in NLP (Transformer)

- \* BN destroys sequential dependence
- \* Sentence length varies



- \* Memory

$z$

$\square$

$\square$

$\square$

$\square$

$z'$

$\square$

$\square$

$\square$

$\square$

$\square$

$\square$

$\square$

John

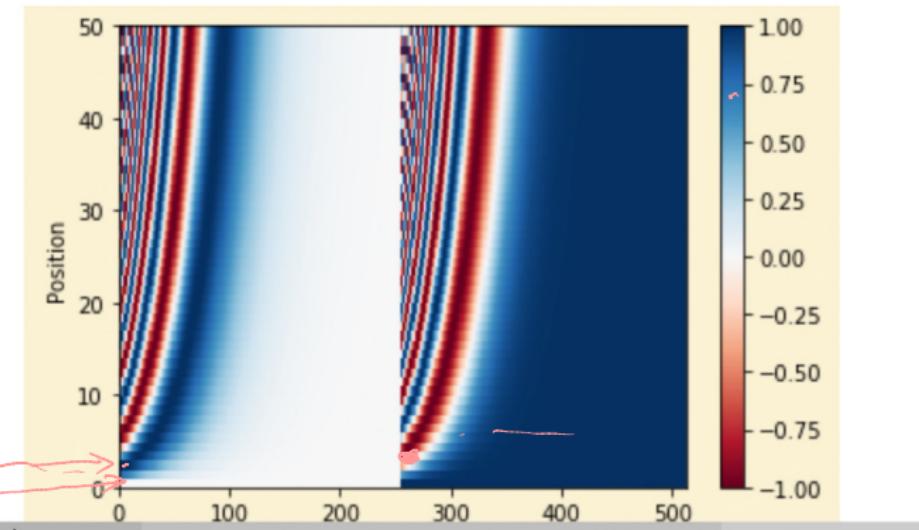
loves Mary

Mary

loves

John

$$z = z'$$



Je suis étudiant [0, ..., 0, 1, 1, ..., 1]

Suis étudiant [0.7, ..., 0.7, -0.7, ..., 1]

# Transformer

- \* Encoder + Decoder
- \* Seq2seq model

BERT: Bidirectional Encoder Representation from Transformer

- \* Encoder only
- \* Representation learning  
→ Downstream NLU tasks

GPT: Generative Pre-trained Transformer

- \* Decoder only
- \* For NLG tasks

