

Lecture 25

[ RNN: backprop through time ]

Consider  $\begin{cases} \tilde{h}_t = W_{hx} \tilde{x}_t + W_{hh} \tilde{h}_{t-1}, \\ \tilde{o}_t = W_{ho} \tilde{h}_t \end{cases}$

no biases  
for simplicity

Loss function:  $L = \frac{1}{T} \sum_{t=1}^T \ell(y_t, \tilde{o}_t)$

$\uparrow$   
true target  
labels

Need to compute

$$\frac{\partial L}{\partial W_{hx}}, \frac{\partial L}{\partial W_{hh}}, \frac{\partial L}{\partial W_{ho}}$$

matrix  
derivatives

$$\frac{\partial L}{\partial W_{ho}} = \frac{1}{T} \sum_t \frac{\partial \ell(y_t, \tilde{o}_t)}{\partial W_{ho}} = \underbrace{\frac{1}{T} \sum_t \frac{\partial \ell(y_t, \tilde{o}_t)}{\partial \tilde{o}_t} \frac{\partial \tilde{o}_t}{\partial W_{ho}}}_{\text{local in time}}$$

Next, define  $\tilde{w}_h = \{W_{hh}, W_{hx}\}$  flattened weights  
 $\tilde{w}_o = \{W_{ho}\}$

Then  $\begin{cases} \tilde{h}_t = \tilde{f}(\tilde{x}_t, \tilde{h}_{t-1}, \tilde{w}_h), \\ \tilde{o}_t = \tilde{g}(\tilde{h}_t, \tilde{w}_o) \end{cases}$  give

$$\frac{\partial L}{\partial \tilde{w}_h} = \frac{1}{T} \sum_t \frac{\partial \ell(y_t, \tilde{o}_t)}{\partial \tilde{o}_t} \frac{\partial \tilde{g}(\tilde{h}_t, \tilde{w}_o)}{\partial \tilde{h}_t} \frac{\partial \tilde{h}_t}{\partial \tilde{w}_h}$$

Next,  ~~$\frac{\partial \tilde{h}_t}{\partial \tilde{w}_h}$~~   $\frac{\partial \tilde{h}_t}{\partial \tilde{w}_h} = \frac{\partial \tilde{f}(\tilde{x}_t, \tilde{h}_{t-1}, \tilde{w}_h)}{\partial \tilde{w}_h} +$   
 $+ \frac{\partial \tilde{f}(\tilde{x}_t, \tilde{h}_{t-1}, \tilde{w}_h)}{\partial \tilde{h}_{t-1}} \frac{\partial \tilde{h}_{t-1}}{\partial \tilde{w}_h} \leftarrow \text{recursive update}$

This yields

$$\frac{\partial \tilde{h}_t}{\partial \tilde{w}_h} = \frac{\partial \tilde{f}(\tilde{x}_t, \tilde{h}_{t-1}, \tilde{w}_h)}{\partial \tilde{w}_h} + \sum_{i=1}^{t-1} \left( \prod_{j=i+1}^t \frac{\partial \tilde{f}(\tilde{x}_j, \tilde{h}_{j-1}, \tilde{w}_h)}{\partial \tilde{h}_{j-1}} \right) \times \frac{\partial f(\tilde{x}_i, \tilde{h}_{i-1}, \tilde{w}_h)}{\partial \tilde{w}_h}.$$

$\Theta(T^2)$  complexity

Thus, the sum is truncated to the most recent K terms (K may be chosen adaptively)

---

[Gating and long-term memory]

RNNs "forget" inputs from the past due to vanishing gradients. Solutions:  
GRU + LSTM.

Gated recurrent units (GRU)

Consider  $X_t = \begin{matrix} N \times D \\ \uparrow \\ \text{batch size} \end{matrix}$  data matrix  
 $\begin{matrix} \uparrow \\ \text{vocab. size} \\ \text{(input data dim'n)} \end{matrix}$

$H_t = N \times H$   
 $\uparrow$   
# hidden states (i.e., units)

$N$   
 $N \times H$ , with  
the bias vector  
in  $b_r$  in each row

Compute

$\underbrace{D \times H}$

$\underbrace{H \times H}$

reset gate

$$R_t = \sigma(X_t W_{xr} + H_{t-1} W_{hr} + b_r) \quad N \times H$$

update gate

$$Z_t = \sigma(X_t W_{xz} + H_{t-1} W_{hz} + b_z) \quad N \times H$$

$R_t, Z_t \in [0, 1]$   
element-wise

Next, compute

$$\tilde{H}_t = \tanh(X_t W_{xh} + (R_t \otimes H_{t-1}) W_{hh} + b_h) \quad N \times H$$

If all entries in  $R_t$  are close to 1,

$\tilde{H}_t$  is approx.  $\sqrt{H_t}$  from the standard RNN.

If all entries in  $R_t$  are close to 0,

$\tilde{H}_t$  'forgets'  $H_{t-1}$  and acts as a simple MLP. Thus, reset gate captures short-term dependencies.

Finally, compute  $N \times H$  matrix of ones

$$H_t = Z_t \otimes H_{t-1} + (1 - Z_t) \otimes \tilde{H}_t$$

When  $Z_{ij} \approx 1 \Rightarrow H_{t,ij} \approx H_{t-1,ij}$ , hidden state passed unchanged  
*[i.e.,  $X_t$  is ignored]*

when  $Z_{ij} \approx 0 \Rightarrow H_{t,ij} \approx \tilde{H}_{t,ij}$

Overall, the update gate captures long-term dependencies

# Long short term memory (LSTM)

Compute

$$\begin{cases} O_t = \sigma(X_t W_{x_0} + H_{t-1} W_{h_0} + b_o), & \text{output gate} \\ I_t = \sigma(X_t W_{x_i} + H_{t-1} W_{h_i} + b_i), & \text{input gate} \\ F_t = \sigma(X_t W_{x_f} + H_{t-1} W_{h_f} + b_f) & \text{(not) forget gate} \end{cases}$$

Candidate cell state:

$$\tilde{C}_t = \tanh(X_t W_{x_c} + H_{t-1} W_{h_c} + b_c)$$

Cell state:

$$C_t = F_t \otimes C_{t-1} + I_t \otimes \tilde{C}_t$$

$\uparrow$   $\uparrow$   
 $C_{t-1}$  included       $\tilde{C}_t$  included if  
 if  $F_t$  is 'on'       $I_t$  is 'on'

$F_t \approx 1$  &  $I_t \approx 0$  over many  
 timesteps ~ long-term memories,  
 stored in cells  $C_t$

Finally,  $H_t = \underbrace{O_t \otimes \tanh(C_t)}_{\text{short-term memory}} \text{ influenced by long-term memory stored in } C_t$

## Attention

Consider  $m$  feature vectors of length  $v$  referenced by  $m$  keys of length  $k$ :

$$K \in \mathbb{R}^{m \times k} \Rightarrow V \in \mathbb{R}^{m \times v}$$

The model will produce a weighted combination of  $m$  feature vectors based on the input query vector  $\vec{q} \in \mathbb{R}^d$ .

'Soft' (differentiable) lookups:

$$\text{Attn}(\vec{q}; (\vec{k}_1, \vec{v}_1) \dots (\vec{k}_m, \vec{v}_m)) = \underbrace{\sum_{i=1}^m \alpha_i(\vec{q}, \vec{k}_{1:m}) \vec{v}_i}_{\text{attention weights}}$$

$$= \text{Attn}(\vec{q}; \vec{k}_{1:m}, \vec{v}_{1:m})$$

$$0 \leq \alpha_i \leq 1 ; \quad \sum_{i=1}^m \alpha_i = 1 .$$

$$\text{Explicitly, } \alpha_i(\vec{q}, \vec{k}_{1:m}) = \frac{e^{a(\vec{q}, \vec{k}_i)}}{\sum_{j=1}^m e^{a(\vec{q}, \vec{k}_j)}},$$

where  $a(\vec{q}, \vec{k}_i) \in \mathbb{R}$  = attention score

Parametric attention: assume that

$d \equiv q_f = k$ , compute

$$a(\vec{q}, \vec{k}) = \frac{\vec{q} \cdot \vec{k}}{\sqrt{d}} \text{ to scale out the variance}$$

In practice, we deal with  $n$  query vectors at the same time:

$$Q \in \mathbb{R}^{n \times d}, K \in \mathbb{R}^{m \times d}, V \in \mathbb{R}^{m \times v}$$

Then

$$\text{Attn}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \in \mathbb{R}^{n \times v}$$

softmax applied row-wise

seq2seq with attention (RNN)

Recall that in RNN decoder,

$$\tilde{h}_t^d = \tilde{f}_d(\tilde{h}_{t-1}^d, \tilde{y}_{t-1}, \tilde{c})$$

↑  
context vector,  
encodes input  $\tilde{x}_{1:T}$   
entire

Typically,  $\tilde{c} = \tilde{h}_{T+1}^e \Leftarrow$  final hidden state  
of RNN encoder

Idea: replace  $\tilde{c}$  with  $\tilde{h}_i^e$  feature  $i$

$$\tilde{c}_t = \sum_{i=1}^T \lambda_i (\tilde{h}_{t-1}^d, \tilde{h}_{1:T}^e) \tilde{h}_i^e$$

↑  
query      ↑  
set of  $T$   
keys

"decoder  
attends to  
encoder"

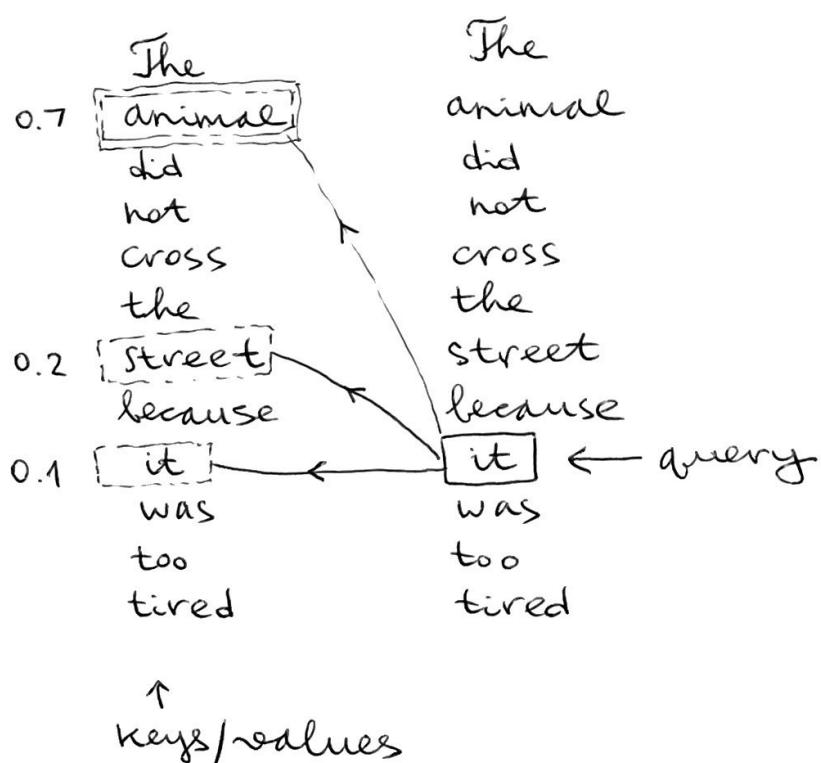
Then  $\tilde{h}_t^d = \tilde{f}_d(\tilde{h}_{t-1}^d, \tilde{y}_{t-1}, \tilde{c}_t)$

## Transformers:

seq2seq model which uses attention in both encoder and decoder.

Applications: machine translation, music generation, text summarization, image generation (treat image as seqn of pixels), etc.

Self-attention: encoder attends to itself given input tokens  $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^d$ , generate  $\tilde{y}_i = \text{Attn}(\tilde{x}_i; (\tilde{x}_1, \tilde{x}_1) \dots (\tilde{x}_n, \tilde{x}_n))$



Multi-headed attention: caption several notions of similarity at once.

Consider  $\vec{q}_j \in \mathbb{R}^{d_q}$ ,  $\vec{k}_j \in \mathbb{R}^{d_k}$ ,  $\vec{v}_j \in \mathbb{R}^{d_v}$

Define embeddings:  $\vec{q}'_i = \underbrace{W_i^{(q)}}_{p \times d_q} \vec{q}_j$ ;  $i=1, \dots, h$

$$\vec{k}'_{i,j} = \underbrace{W_i^{(k)}}_{p \times d_k} \underbrace{\vec{k}_j}_{d_k}; \quad \vec{v}'_{i,j} = \underbrace{W_i^{(v)}}_{p_v \times d_v} \underbrace{\vec{v}_j}_{d_v}$$

Compute  $\underbrace{\vec{h}_i}_{p_v} = \text{Attn}(\vec{q}'_i; (\vec{k}'_{i,1}, \vec{v}'_{i,1}), (\vec{k}'_{i,2}, \vec{v}'_{i,2}), \dots)$

Finally,  $\vec{h} = \text{MHA}(\vec{q}; (\vec{k}_1, \vec{v}_1), (\vec{k}_2, \vec{v}_2), \dots) =$

$$= \underbrace{W_o}_{p_v \times (hp_v)} \begin{pmatrix} \vec{h}_1 \\ \vdots \\ \vec{h}_h \end{pmatrix} \in \mathbb{R}^{p_v}.$$

Positional embedding: use a set of basis functions to encode word positions in the sequence.

Consider a sequence of tokens labeled by  $i=1, \dots, n$ .

Compute

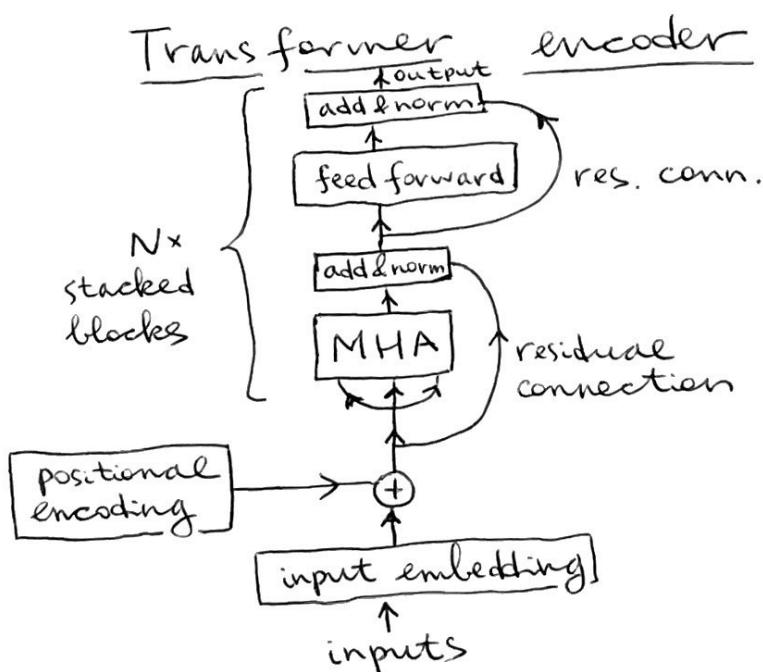
$$\begin{cases} p_{i,2j} = \sin\left(\frac{i}{C^{2j/d}}\right), \\ p_{i,2j+1} = \cos\left(\frac{i}{C^{2j/d}}\right). \end{cases}$$

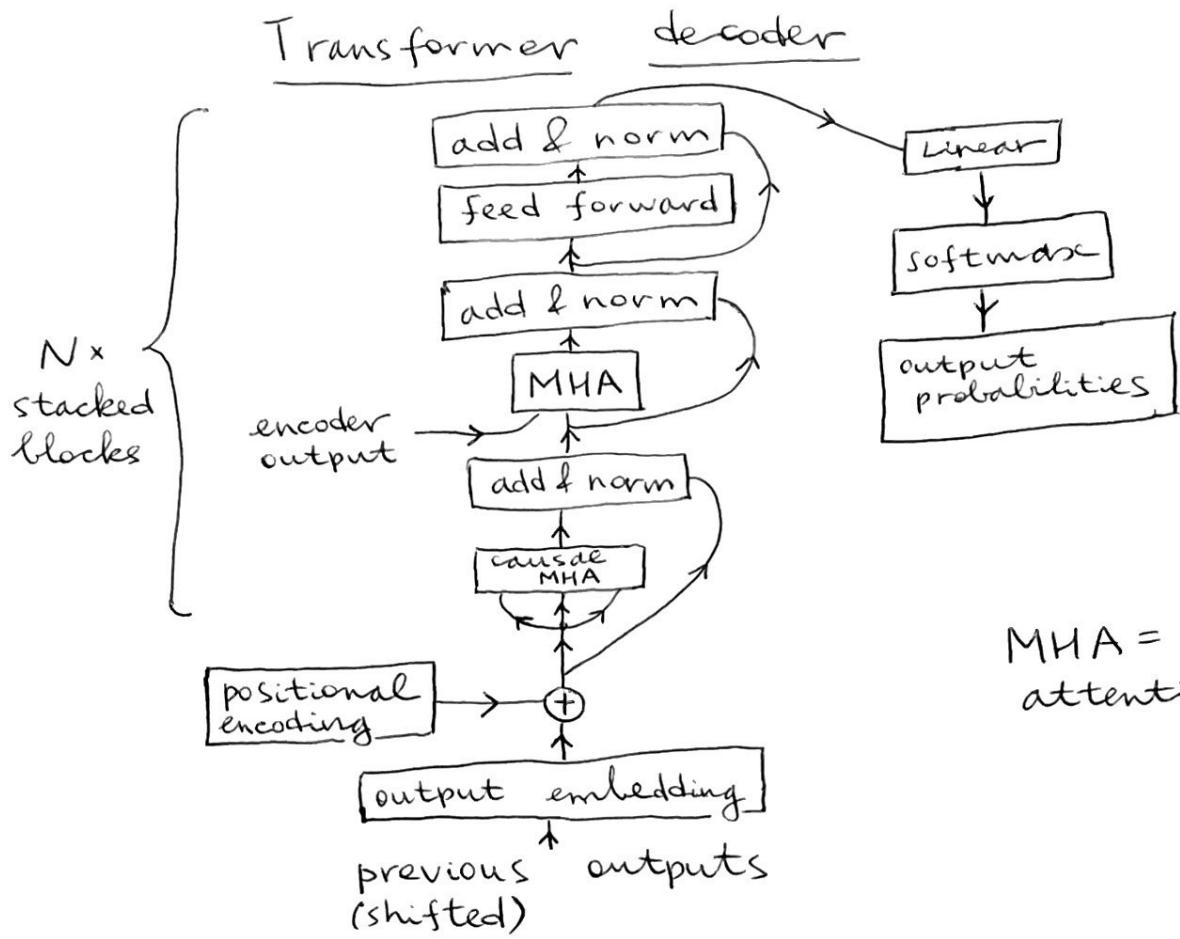
Here,  $C=10^4$  = max sequence length;  
 $d$  = dim of embedding;  
 $j=1, \dots, d$ .

For example, if  $d=4$ , token  $i$  is encoded as  $\vec{p}_i = (\sin\frac{i}{C^{0/4}}, \cos\frac{i}{C^{0/4}}, \sin\frac{i}{C^{2/4}}, \cos\frac{i}{C^{2/4}})$ .

If a token is represented by a (learned) embedding  $\vec{x}_i$  of length  $d$ , we compute  $\vec{p}_i$  of length  $d$  and add them:

$$\text{Pos}(\text{Embed}(\text{token}_i)) = \vec{x}_i + \vec{p}_i.$$





Causal (or masked) MHA : future tokens  
masked / unavailable