

Review of lecture 04.

D. Encoder-Decoder model



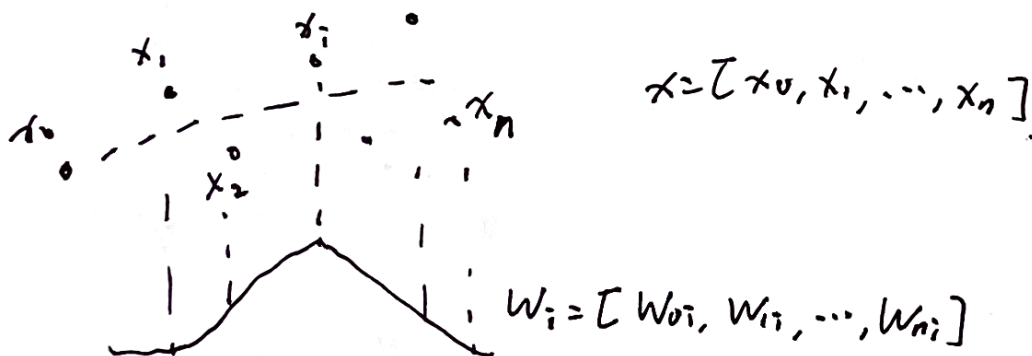
RNN: $x_1 \rightarrow h_1, x_2 \rightarrow h_2, \dots, x_n \rightarrow h_n,$
 $\tilde{c} = h_n \rightarrow \text{decoder}.$

RNN + Attention: $x_1 \rightarrow h_1, x_2 \rightarrow h_2, \dots, x_n \rightarrow h_n$

$w_1 h_1 + w_2 h_2 + \dots + w_n h_n \rightarrow \tilde{c}.$

new context: $[\tilde{c}, h_{n+1}]$ concatenate. \rightarrow output.

②. Intuition of Attention: Re-weighting : get more and better context.

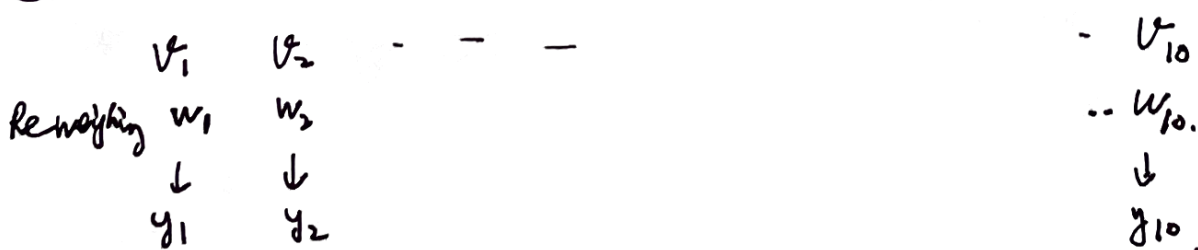


$y = [y_0, y_1, \dots, y_n].$

$y = W \cdot x, \quad W = [w_0, w_1, \dots, w_n].$

A good W should smooth the noise of x to get a good y .

③. Noa can be annoying but she is a great cut.

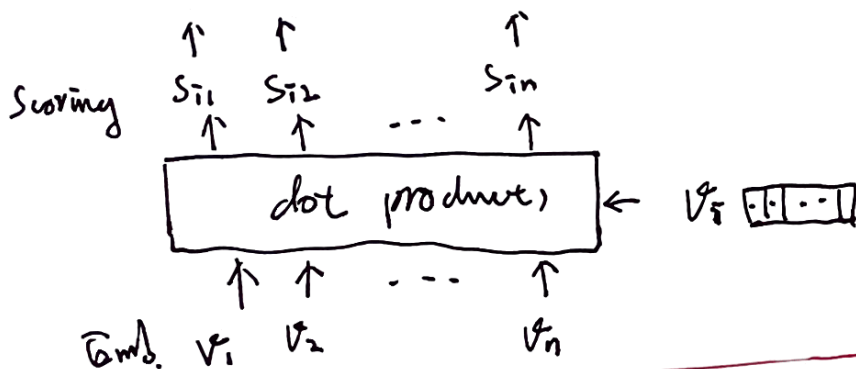
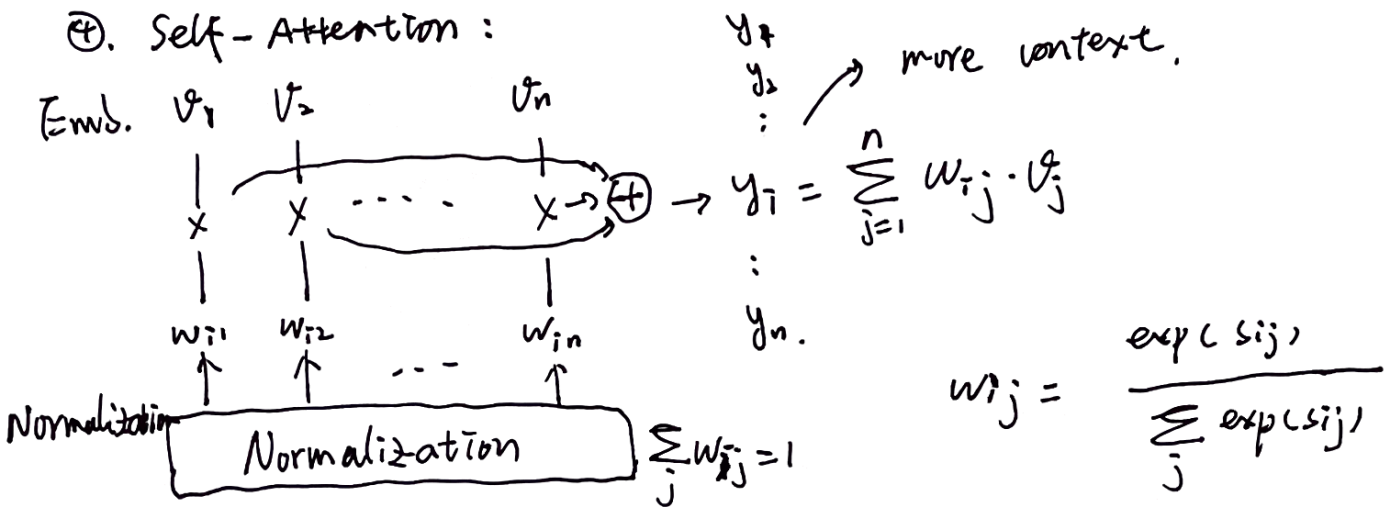


$y_{1:10}$: is more and better context compared with $v_{1:10}$.

$$y_i = w_{i1} \cdot v_1 + w_{i2} \cdot v_2 + \dots + w_{in} \cdot v_n.$$

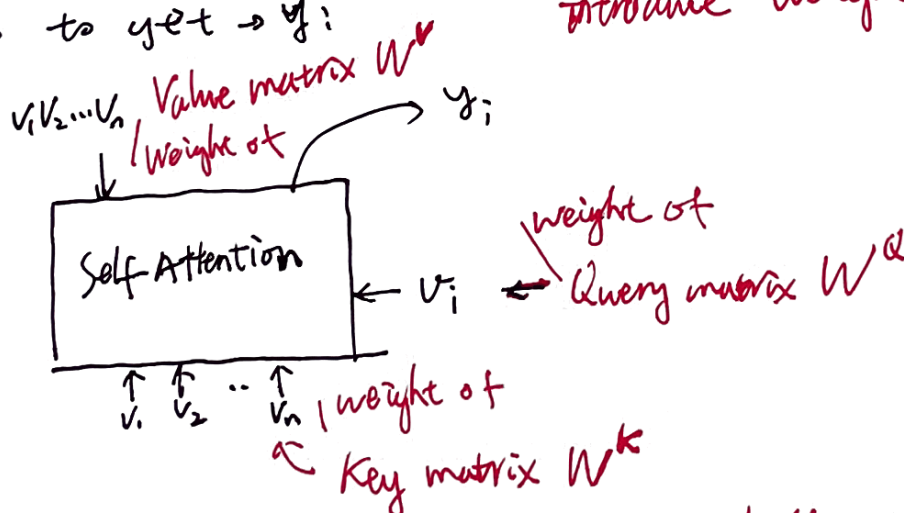
How to design w_{ij} ? How to introduce model weights?

⊕. Self-Attention:



⊖-1. to get $\rightarrow y_i$:

introduce weights.



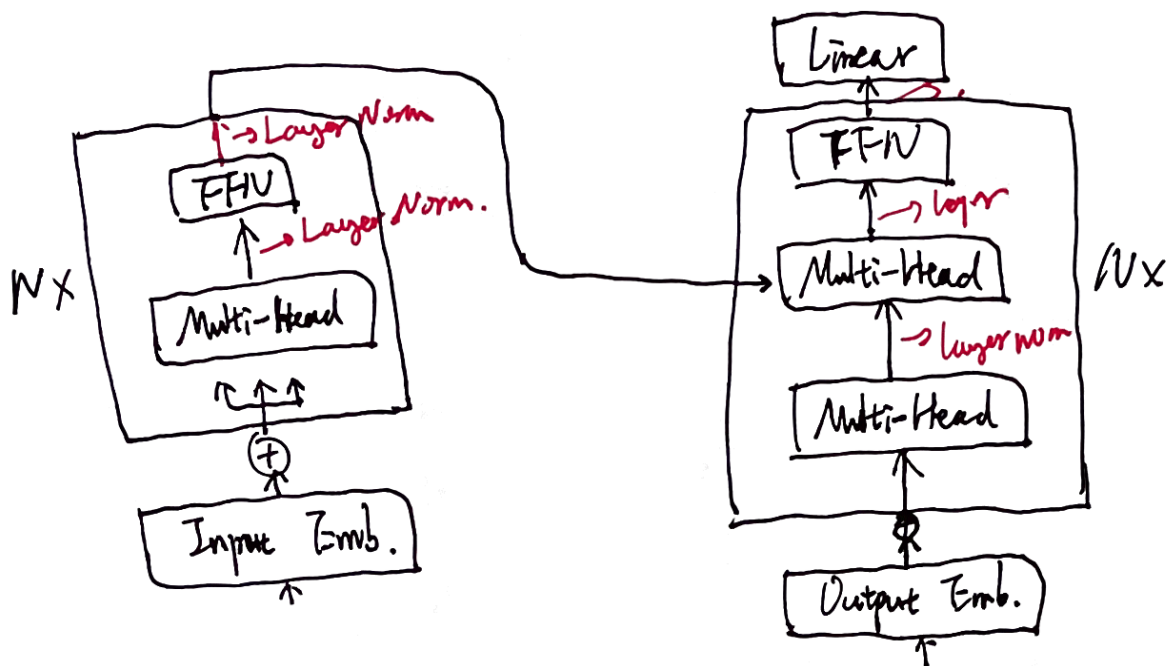
Let $X = [v_1, v_2, \dots, v_n]$ is the input embedding matrix.

Query matrix: $Q = X \cdot W^Q$

Key matrix: $K = X \cdot W^K$

Value matrix: $V = X \cdot W^V$

- Each row is a query/key/value vector.
- Time complexity: $O(n^2 \cdot d)$.



Transformer

Model configuration:

- N : number of layers
- d_{model} : dimension of output layers
- d_{ff} : inner layer dimension of FFN
- h : # heads
- d_k : queries and keys
- d_v : values

Total model size:

$$\text{Total} = A1 + A2$$

$$A1 = \text{Input Emb.} + \text{Output Emb.} + \text{Linear}$$

$$A2 = N \times \text{Encoder} + N \times \text{Decoder.}$$

$A1$: all three weight matrix share same weight. (See 3.4 Embeddings and Softmax).
 $W^{\text{emb}} \in \mathbb{R}^{V \times d_{\text{model}}}$

$$A2 = N \times \text{Encoder} + N \times \text{Decoder.}$$

$$\text{Encoder} = \text{Multihead} + \text{FFN.}, \quad \text{Decoder} = 2 \times \text{Multihead} + \text{FFN.}$$

$$\text{Multihead} = (|W_i^Q| + |W_i^K| + |W_i^V|) \cdot h + |W^O| + 2 \cdot \text{layer norm}$$

$$W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}, \quad W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}, \quad W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$$

$$W^O \in \mathbb{R}^{h \times d_v \times d_{\text{model}}}$$

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

$$W_1 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}, \quad b_1 \in \mathbb{R}^{d_{\text{ff}}}, \quad W_2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}, \quad b_2 \in \mathbb{R}^{d_{\text{model}}}$$

$$|V| = 37000,$$

$$N = 6, d_{\text{model}} = 512, d_{\text{ff}} = 2048, h = 8, d_k = 64, d_v = 64$$

$$Total = A1 + A2.$$

$$h \times d_v = 8 \times 64 = 512$$

$$A1 = |V| \times d_{\text{model}} = 37000 \times 512$$

$$A2 = N \times \text{Encoder} + N \times \text{Decoder}$$

$$\begin{aligned} \text{Encoder} &= h \cdot (|W_i^Q| + |W_i^K| + |W_i^V|) + |W^O| + \text{FFN} \\ &= 8 \cdot (512 \times 64 + 512 \times 64 + 512 \times 64) + 512 \times 512 + \text{FFN} \\ &= 4 \cdot 512 \times 512 + \text{FFN} \end{aligned}$$

$$\begin{aligned} \text{FFN} &= 2 \cdot d_{\text{model}} \times d_{\text{ff}} + d_{\text{ff}} + d_{\text{model}} \\ &= 2 \cdot 512 \times 2048 + 2048 + 512 \\ &= 2 \cdot 1024^2 + 2 \cdot 1024 + 512 \end{aligned}$$

$$\text{Encoder} = 1024^2 + 2 \cdot 1024^2 + 2 \cdot 1024 + 512$$

$$\text{Decoder} = 2 \cdot 1024^2 + 2 \cdot 1024^2 + 2 \cdot 1024 + 512$$

$$\begin{aligned} Total &= 37000 \times 512 + 6 \cdot (1024^2 \times 3 + 2 \cdot 1024 + 512) \\ &\quad + 6 \cdot (1024^2 \times 4 + 2 \cdot 1024 + 512) \end{aligned}$$

$$37000 \approx 18.06640625 \times 1024 \times 2$$

$$18 + 18 + 24 \approx 60.$$

$$= 37000 \times 512 + 42 \times 1024^2 + 30 \times 1024$$

$$= 63,014,912 \quad 63M.$$

Use the following:

$$Total = |V| \cdot d_{\text{model}} + N \cdot [B_1 + B_2]$$


$$\begin{aligned} B_1 &= (2 \cdot d_{\text{model}} \times d_k + d_{\text{model}} \cdot d_v) \cdot h + h \cdot d_v \cdot d_{\text{model}} + 2 \cdot d_{\text{model}} \cdot d_{\text{ff}} + d_{\text{model}} + d_{\text{ff}} \\ &= (2 \cdot d_{\text{model}} \cdot d_k + d_{\text{model}} \cdot d_v) \cdot 2h + 2h \cdot d_v \cdot d_{\text{model}} + 2 \cdot d_{\text{model}} \cdot d_{\text{ff}} + d_{\text{model}} + d_{\text{ff}} \end{aligned}$$

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

For each q : $\underline{q \cdot K^T}$. \rightarrow we only need to consider small

subset of closest keys. 32, 64.

Locality sensitive hashing:

Idea of random projection:  \rightarrow false.

Given n points $P \subset \mathbb{R}^d$, for a given $q \in \mathbb{R}^d$, find the point $p \in P$ that is closest to q . $O(n^p)$. $p < 1$.

Linformer:

$$X \cdot W_i^Q = Q,$$

$$\Pr(\| \tilde{p} W^T - p W^T \| \leq \epsilon \cdot \| p \cdot W^T \|) \geq 1 - O(\epsilon).$$

$$\text{rank}(B) = \Theta(\log(n)).$$

Remark: \oplus layer normalization:

$$y = r \cdot \frac{x - \mathbb{E}(x)}{\sqrt{\text{Var}(x) + \epsilon}} + \beta.$$

r : weights.

β : bias.

$$r \in \mathbb{R}^{d_{\text{model}}}, \beta \in \mathbb{R}^{d_{\text{model}}}$$

Encoder = $N \times 2 \times d_{\text{model}}$

Decoder = $N \times 3 \times d_{\text{model}}$.

Time complexity: $O(n^2 \cdot d + n \cdot d^2)$

w_1, w_2, \dots, w_n :

weighted sum of all words.