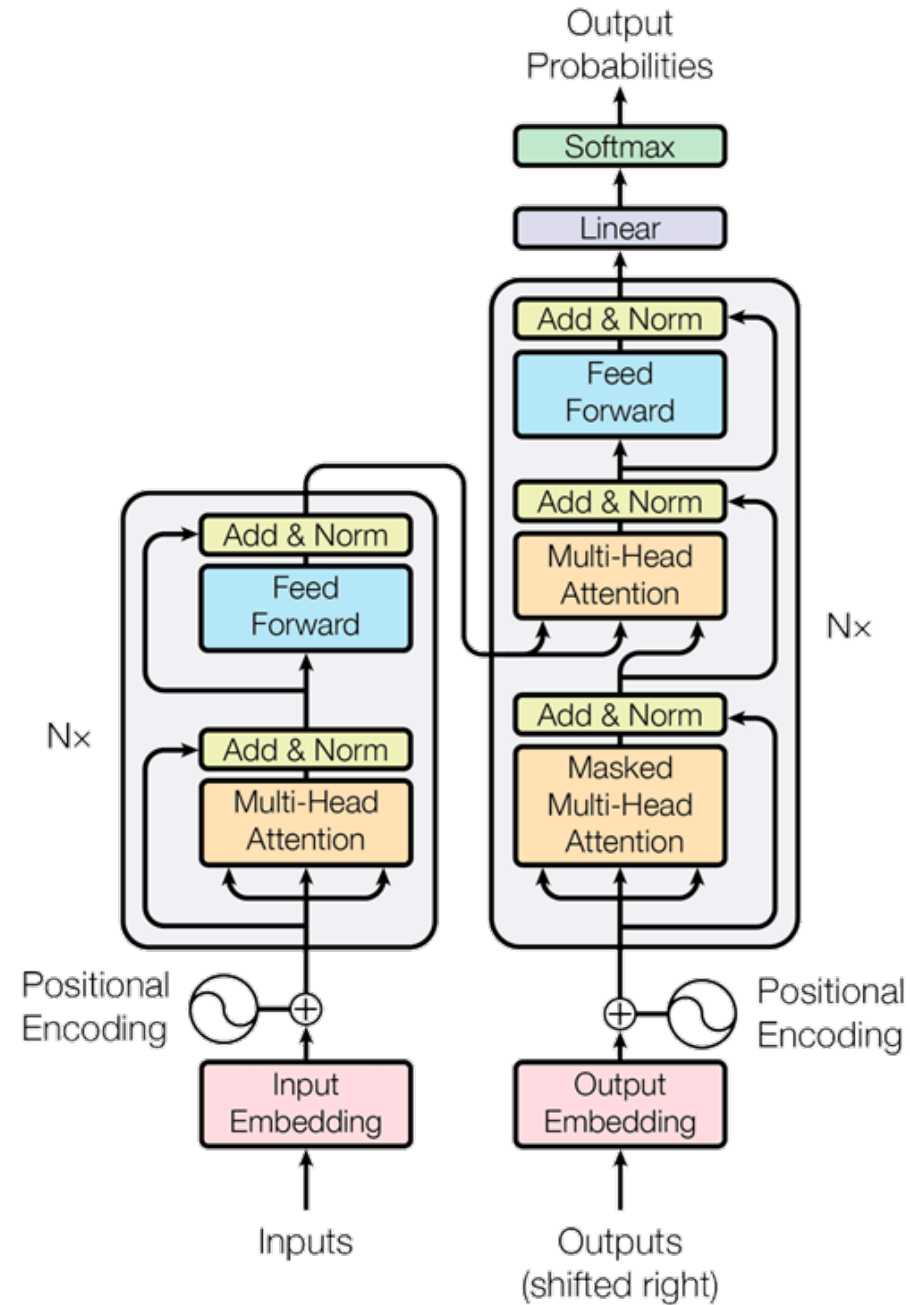


Transformer: Decoder Block with Masks

Ki Hyun Kim

nlp.with.deep.learning@gmail.com

Transformer



Equations

- Given Dataset,

$$\mathcal{D} = \{x^i, y^i\}_{i=1}^N$$
$$x^i = \{x_1^i, \dots, x_m^i\} \text{ and } y^i = \{y_0^i, y_1^i, \dots, y_n^i\},$$

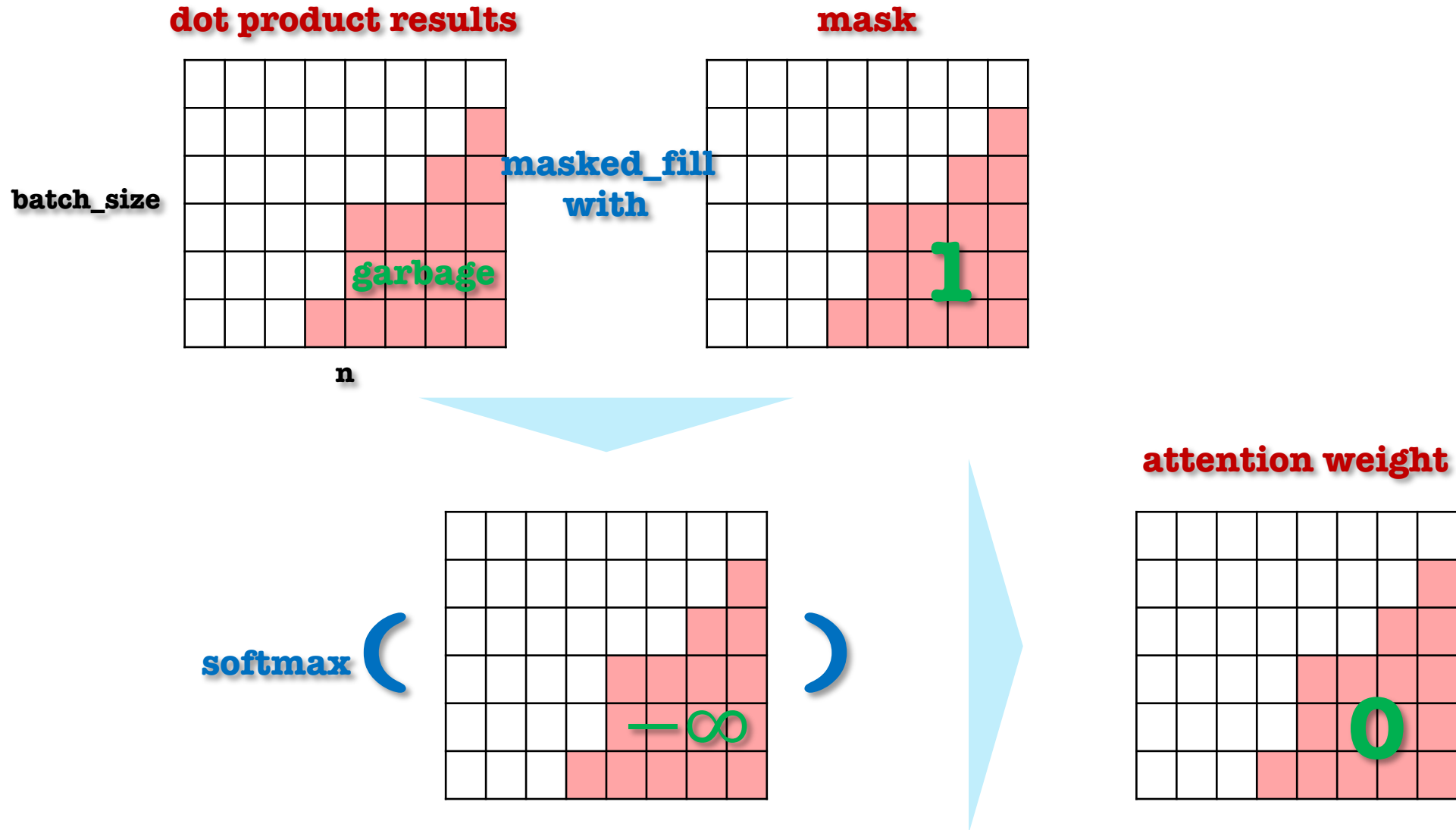
where $y_0 = \langle \text{BOS} \rangle$ and $y_n = \langle \text{EOS} \rangle$.

- What we want is

$$\hat{y}_{1:n} = f(x_{1:m} : \theta)$$

Before we start,

- Using mask, assign $-\infty$ to make 0s for softmax results.



Equations

- Self-attention with mask

$$\begin{aligned}h_{0,1:n} &= \text{emb}(y_{0:n-1}) + \text{pos}(0, n - 1) \\ \tilde{h}_{i,1:n}^{\text{dec}} &= \text{LayerNorm}(\text{Multihead}_i(Q, K, V) + h_{i-1,1:n}^{\text{dec}}), \\ &\text{where } Q = K = V = h_{i-1,1:n}^{\text{dec}}.\end{aligned}$$

Equations

- Attention from encoder with mask for <pad>

$$\tilde{h}_{i,1:n}^{\text{dec}} = \text{LayerNorm}(\text{Multihead}_i(Q, K, V) + h_{i-1,1:n}^{\text{dec}}),$$

where $Q = \tilde{h}_{i,1:n}^{\text{dec}}$ and $K = V = h_{\ell,1:m}^{\text{dec}}$.

Equations

- FC layers

$$\text{FFN}(h_{i,t}) = \text{ReLU}(h_{i,t} \cdot W_i^1) \cdot W_i^2$$

where $W_i^1 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}$ and $W_i^2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}$.

$$h_{i,1:m}^{\text{dec}} = \text{LayerNorm}([\text{FFN}(\tilde{h}_{i,1}^{\text{dec}}); \dots; \text{FFN}(\tilde{h}_{i,m}^{\text{dec}})] + \tilde{h}_{i,1:m}^{\text{dec}})$$

Equations

- Decoder is stack of decoder blocks.

$$\begin{aligned}h_{\ell_{\text{dec}},1:m}^{\text{dec}} &= \text{Block}_{\text{dec}}(h_{\ell_{\text{dec}}-1,1:m}^{\text{dec}}) \\&\dots \\h_{1,1:m}^{\text{dec}} &= \text{Block}_{\text{dec}}(h_{0,1:m}^{\text{dec}})\end{aligned}$$

- Generator:

$$\begin{aligned}\hat{y}_{1:n} &= \text{softmax}(h_{\ell_{\text{dec}},1:m}^{\text{dec}} \cdot W_{\text{gen}}), \\ \text{where } h_{\ell_{\text{dec}},1:m}^{\text{dec}} &\in \mathbb{R}^{\text{batch_size} \times n \times \text{hidden_size}} \text{ and } W_{\text{gen}} \in \mathbb{R}^{\text{hidden_size} \times |V|}.\end{aligned}$$

Summary

- Decoder는 2가지의 attention으로 구성됨
 - Attention from encoder:
 - K 와 V 는 encoder의 최종 출력 값, Q 는 이전 레이어의 출력 값
 - Self-attention with mask:
 - Q, K, V 는 이전 레이어의 출력 값
 - Attention weight 계산 시, softmax 연산 이전에 masking을 통해 음의 무한대를 주어, 미래 time-step을 보는 것을 방지
- 추론 때에는 self-attention의 mask는 필요 없으나, 모든 layer의 t 시점 이전의 모든 time-step($< t$)의 hidden state가 필요