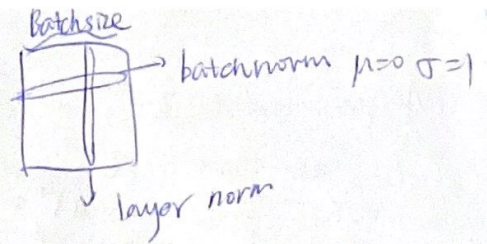
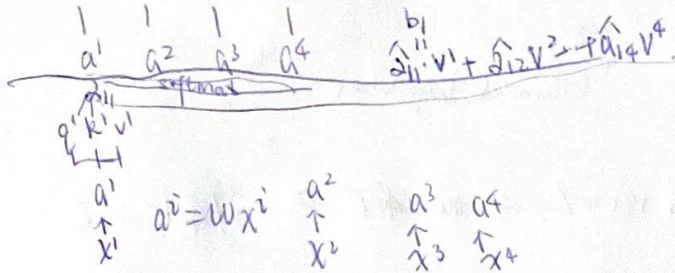
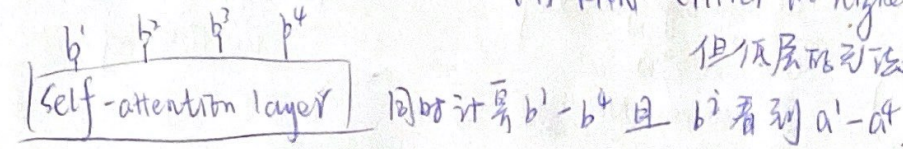


Transformer 作 summarizer
 seq2seq model with self-attention
 RNN 问题: hard to parallel



所有输入通过 CNN 代替 RNN (Filter in higher layer can consider longer sequence, 但低层只能考虑 short sequence)



query (to match others) : $q^i = W^q a^i$
 key (to be matched) : $k^i = W^k a^i$
 value (information to be extracted) : $v^i = W^v a^i$

对每个 q^i 对每个 k^j 做 attention

Scaled dot-product attention : $z_{ij} = \frac{q^i \cdot k^j}{\sqrt{d}}$
 (dim of q and k)

$A = K^T Q \rightarrow \hat{A}$: A softmax columnwise, $\rightarrow O = V \hat{A}$

input Z
 output O
 $Q = W^q Z$
 $K = W^k Z$
 $V = W^v Z$

$A = K^T Q \rightarrow \hat{A} \rightarrow O = V \hat{A}$

改进: multi-head self-attention
 2-head: $q^1, q^2, k^1, k^2, v^1, v^2$
 # heads

$q^1 = W^{q1} q^1$
 $q^2 = W^{q2} q^2$

每个 head 做独立 attention, 得到 b^1, b^2
 $b^i = W^b \begin{pmatrix} b^{i1} \\ b^{i2} \end{pmatrix}$

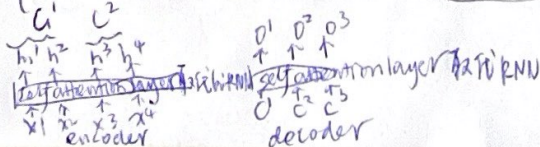
No position information in self-attention. 需要 positional encoding.

$e^i + a^i$. 对 x^i 要 append one-hot position vector.

positional vector e^i fixed not learn

$$(W^T W^P) \begin{pmatrix} x_i \\ p_i \end{pmatrix} = \frac{W^T x_i}{a_i} + \frac{W^P p_i}{e_i}$$

seq2seq with attention



Transformer 结构

