## Attention

- key mechanism for accessing relevant information in a context to make predictions

* How can attention impact language modeling?

**Ex.** Fixed length sequences of A and B

AAAAAAA
ABAAAA B      All A's : Last letter is A
ABAABA B      Any B : Last letter is B
AAAABA B

- Keys and Query

$$A \quad A \quad B \quad A \quad \underline{\quad}^{(B)}$$

Keys : embeddings of sequence
Query : what we want to find

Assume : $A = [1, 0]$ $\Big\}$ one-hot encoding
           $B = [0, 1]$ $\Big\}$ embeddings $e_i$

- **Step 1** : Compute score for each key given query

$$[1\ 0] \quad [1\ 0] \quad [0\ 1] \quad [1\ 0]$$
$$A \qquad A \qquad B \qquad A$$

Dot product    0       0       1       0

score $s_i = K_i^T q$

set $q = [0\ 1]$ to find B's

- **step 2** : Softmax

$$A \qquad A \qquad B \qquad A$$
$$0 \qquad 0 \qquad 1 \qquad 0 \longrightarrow \left[ \tfrac{1}{6} \quad \tfrac{1}{6} \quad \tfrac{1}{2} \quad \tfrac{1}{6} \right] \alpha_i$$

softmax : (assume $e = 3$)

- **Step 3** : Compute output as weighted sum of input

$$\text{result} = \sum_{i=1}^{4} \alpha_i \, e_i = \left( \tfrac{1}{6}[1\ 0] + \tfrac{1}{6}[1\ 0] + \tfrac{1}{2}[0\ 1] + \tfrac{1}{6}[1\ 0] \right)$$
$$= [\tfrac{1}{2} \quad \tfrac{1}{2}]$$

\* Can make attention more peaked by amplifying the embeddings

$$k_i = W_k \, e_i \qquad W_k = \begin{matrix} 10 & 0 \\ 0 & 10 \end{matrix}$$

$$\overset{A}{[10 \ 0]} \ \overset{A}{[10 \ 0]} \ \overset{B}{[0 \ 10]} \ \overset{A}{[10 \ 0]}$$
$$\quad 0 \qquad \quad 0 \qquad \quad 1 \qquad \quad 0$$

- Original Dot product attention : $\quad s_i = k_i^T q$

- Scaled Dot product attention : $\quad s_i = k_i^T W q$

- Equivalent to having two weight matrices: $\quad s_i = (W^k k_i)^T (W^q q)$

## Self Attention

* Mechanism in transformer that processes entire sequence at once
* Every word in sentence is both a key and query simultaneously

$Q$ :  sequence length $\times$ $d$  matrix    ($d$ = embedding dimensions = 2 for this example)

$K$ :  sequence length $\times$ $d$  matrix

$$W^Q = \begin{matrix} 0 & 1 \\ 0 & 1 \end{matrix} \qquad \left( \text{no matter the value, we're going to look for B's} \right)$$

$$W^k = \begin{matrix} 10 & 0 \\ 0 & 10 \end{matrix} \qquad \text{"Booster" as before}$$

$$E = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{pmatrix} \qquad Q = E W^Q = \begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}$$

$$k = E W^k = \begin{pmatrix} 10 & 0 \\ 10 & 0 \\ 0 & 10 \\ 10 & 0 \end{pmatrix}$$

Scores :  $S = Q K^T \qquad S_{ij} = q_i \cdot k_j$

$$S = \begin{pmatrix} 0 & 0 & 10 & 0 \\ 0 & 0 & 10 & 0 \\ 0 & 0 & 10 & 0 \\ 0 & 0 & 10 & 0 \end{pmatrix}$$

* rows represent attention scores
(first row: I care about word 3)

$\longrightarrow$ row-wise softmax turns it into a distribution per row  (A)

Output :  $AE$

$\left( \text{actually } A(E W^v) \right)$

- Vasuani et al.

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

$$Q = EW^Q, \quad K = EW^k, \quad V = EW^v$$

✱ Normalizing by $\sqrt{d_k}$ helps control scale of softmax, makes it less peaked
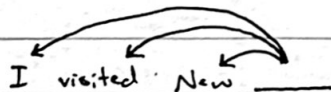
- What does self attention produce?
  - Square attention matrix (A) × input = same dimension as input
  - Computes a contextualized encoding for each word, preserving length of sequence

* "The Illustrated Transformer" by Alammar
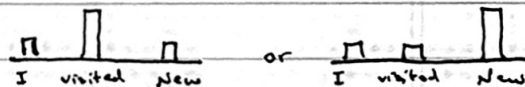  ↳ Visualizes all of this

- <u>Multihead Self Attention</u>

  * Attention can theoretically learn to attend to multiple tokens, but in practice, softmax distributions become peaked:

  I visited New ____          we want:

  we get:          or

  * <u>Solution</u> : Multiple heads to do independent copies of attention
    ↳ <u>Alammar for visualization</u>