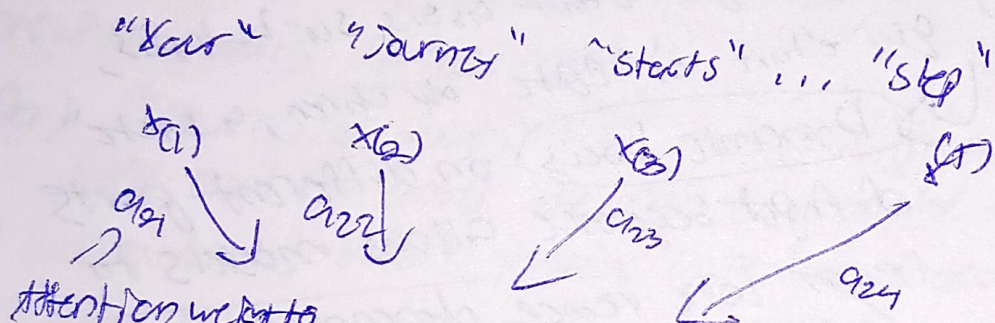


Simple variant of self-attention, free from any learnable weights

Sentence = "your Journey starts with one step"

How much attention do we need to pay to each of these words when we look at Journey

0 Embedding \rightarrow context (clips w/ predicting next word)



attention weights to weigh importance of input x_i

$\begin{bmatrix} 14 & 6 & 5 \end{bmatrix}$
 $z^{(2)}$

input sequence: x

consists of tokens

$x^{(1)}, x^{(2)}, \dots, x^{(T)}$

token embedding

token embedding

context vector $z^{(2)}$ is computed as a combo of all input vectors weighted w.r.t. input element $x^{(2)}$

3D embedding in example.

② "calculate a context vector $z^{(i)}$ for each element $x^{(i)}$ "
context vector: flattened embedding vector.

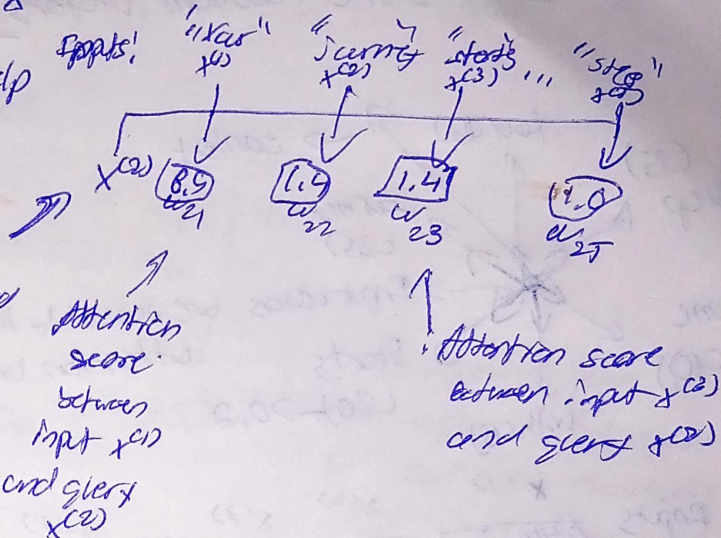
④ Focus on 2nd element $x^{(2)}$ (Journey) ← query

↳ corresponding context vector is $z^{(2)}$ which is an embedding which contains information about $x^{(1)}$ and all other input elements $x^{(1)}$ to $x^{(n)}$

⑤ Compute Attention weights w , also referred to as attention scores

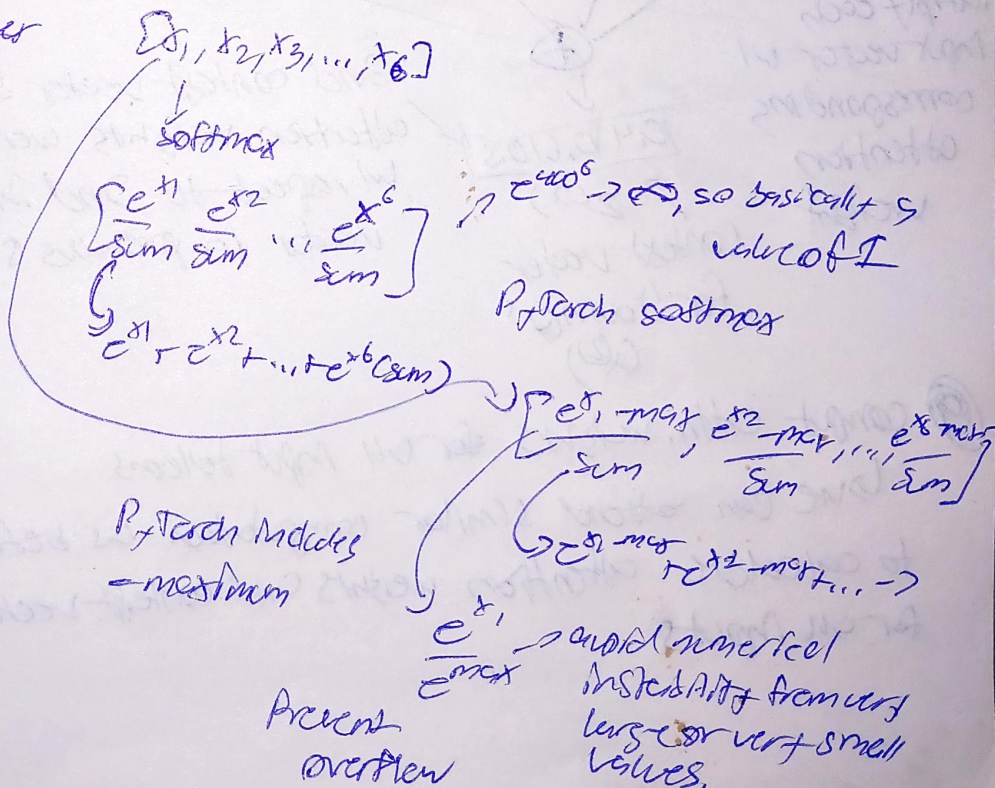
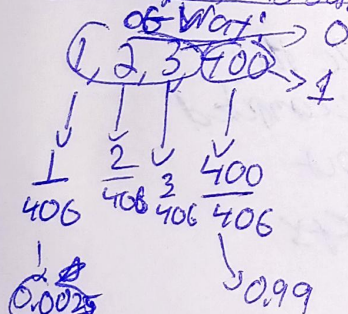
⑥ Later, add trainable weights to help

LMS learn to construct these context vectors, so that they are relevant for LLM to generate next token



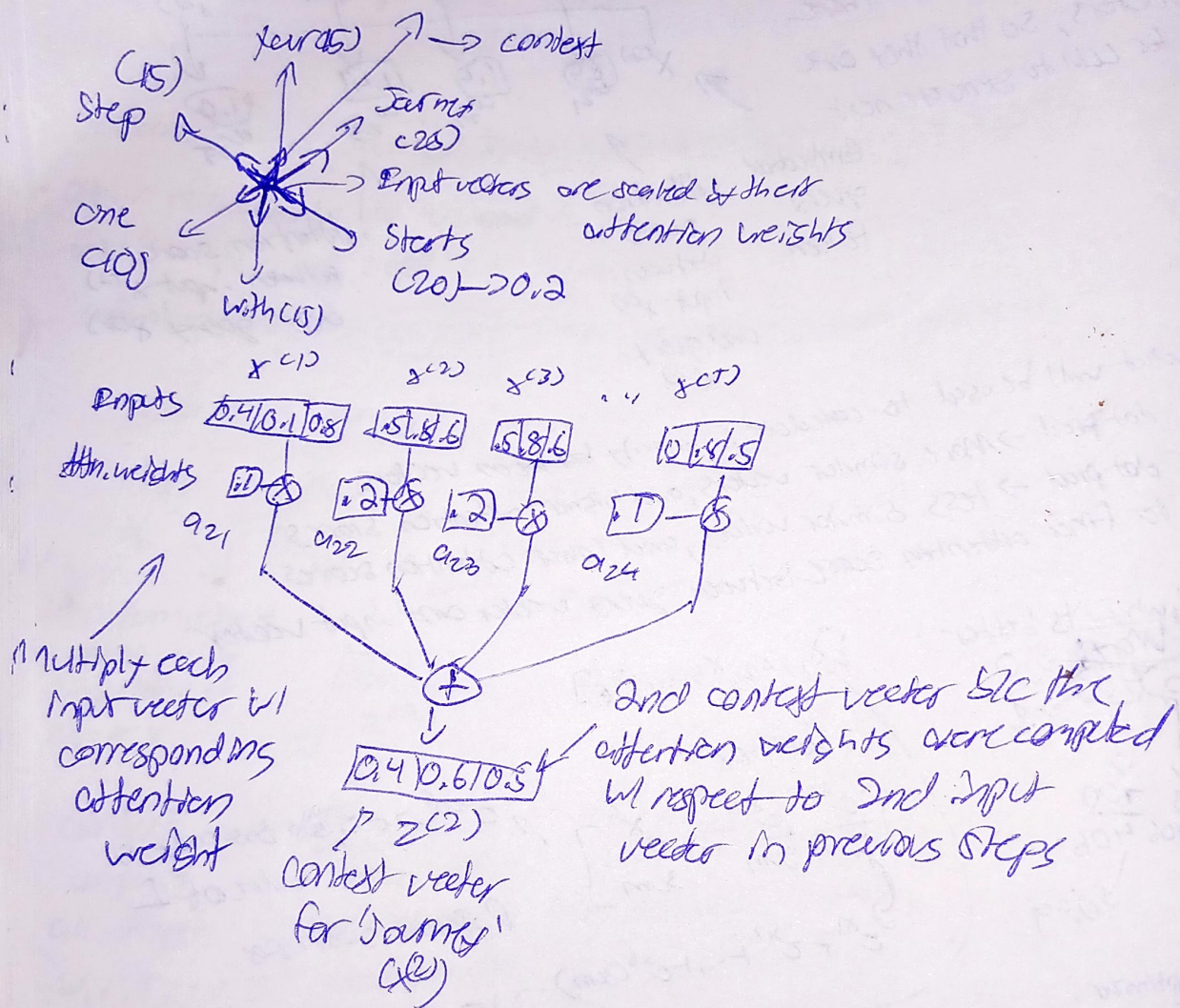
* Dot product will be used to calculate similarity between vectors, High dot prod → More similar vectors, and higher attention scores
Lower dot prod → less similar values, and lower attention scores
↳ use to find attention score between query vector and input vector

softmax, which is better



causes optimizer
s.t. during
backpropagation,
it still gives same
weights to these
values, when
we shouldn't
shrink it

H ⑧ After computing normalized attention weights, we calculate the context vector $z^{(2)}$ by multiplying the embedded input tokens $x^{(i)}$ with the corresponding attention weights and then [summing the resultant vectors]



⑨ Compute attn. weights for all input tokens

Since we can extend similar computation as before to calculate attention weights and context vectors for all inputs

⑩ Follow same 3 steps as before

- 1) compute attention scores \leftarrow dot products between inputs
- 2) compute attention weights \leftarrow normalized version of attention scores
- 3) compute context vectors \leftarrow weighted sum over inputs

Attention weights

Inputs

$$\begin{bmatrix} 0.20 & 0.18 & 0.20 & 0.14 & 0.12 & 0.11 & 0.14 \\ 0.138 & 0.23 & 0.23 & 0.124 & 0.108 & 0.158 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

6×6

$$\begin{bmatrix} 0.43 & 0.15 & 0.69 \\ 0.55 & 0.87 & 0.66 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

6×3

0.138 represents "car"
attn-weight, so multiply it
by "car" vector embeddings
and we have essentially
scaled it by "car"
attn. weights, what
the sum over
"Downing"

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

6×3

r_{c1}, r_{c2}, r_{c3}

$$= 0.138 [0.43 \ 0.15 \ 0.69] + 0.23 [0.55 \ 0.87 \ 0.66]$$

$$r_{c1} \cdot r_{c1} + r_{c2} \cdot r_{c2} + r_{c3} \cdot r_{c3} + \dots + r_{c6} \cdot r_{c6}$$