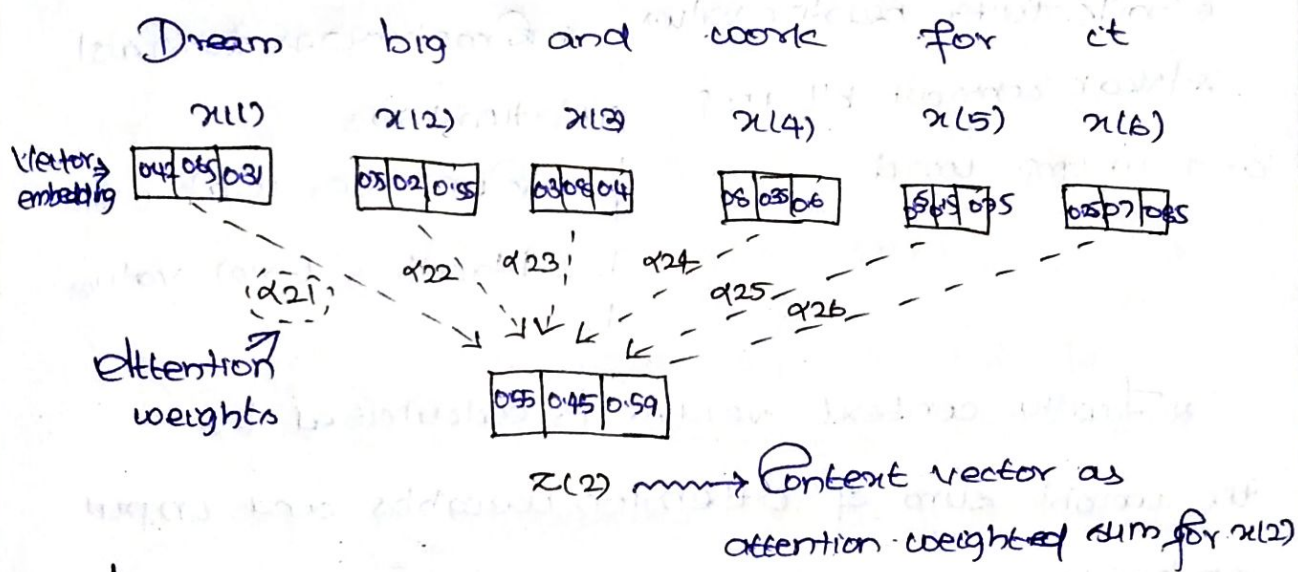


Lec-3

Introduction to Self Attention

- * Self attention generalizes Bahdanau attention.
- * The goal is to create context vector for every word in the sequence by taking into consideration of other words in the sequence.
- * Each word should know how much importance it has to attribute to itself and other words.

Calculation of context vector for each word in the sequence:



Attention score: (A metric to find similarity mathematically)

- * It is computed using dot product of two vectors.

$$W = \vec{v}_1 \cdot \vec{v}_2$$

Interpretation,

- High dot product \rightarrow high relevance
- Near zero \rightarrow weak relevance.
- Negative \rightarrow opposite direction.

* Query (Q) = Current word being processed.

* Key (K) = All words in the sequence.

* Attention score is calculated between query and every single input vector (key) by computing dot product b/w it.

* Attention weights = Normalized form of attention score where the sum is 1.

* We can normalize by 2 methods,

Simple Normalization

$$\alpha_c = \frac{w_c}{\sum w}$$

* Fails with negative values

* Weak contrast b/w imp and unimp. word

Softmax (Preferred)

$$\alpha_c = \frac{e^{w_c}}{\sum e^w}$$

* Emphasizes dominant relationships.

* Suppresses weak ones.

* Handles (-ve) values.

* Finally context vector is calculated by the weight sum of attention weights and input embedding.

* Consider we have 6 words in the input sequence, then the matrix representation can be,

→ Input embedding = 6×3

→ Attention = 6×6

→ Output context vector = 6×3

Attention weights

6×6

× Input embeddings

6×3

↓

Context Vector

6×3