

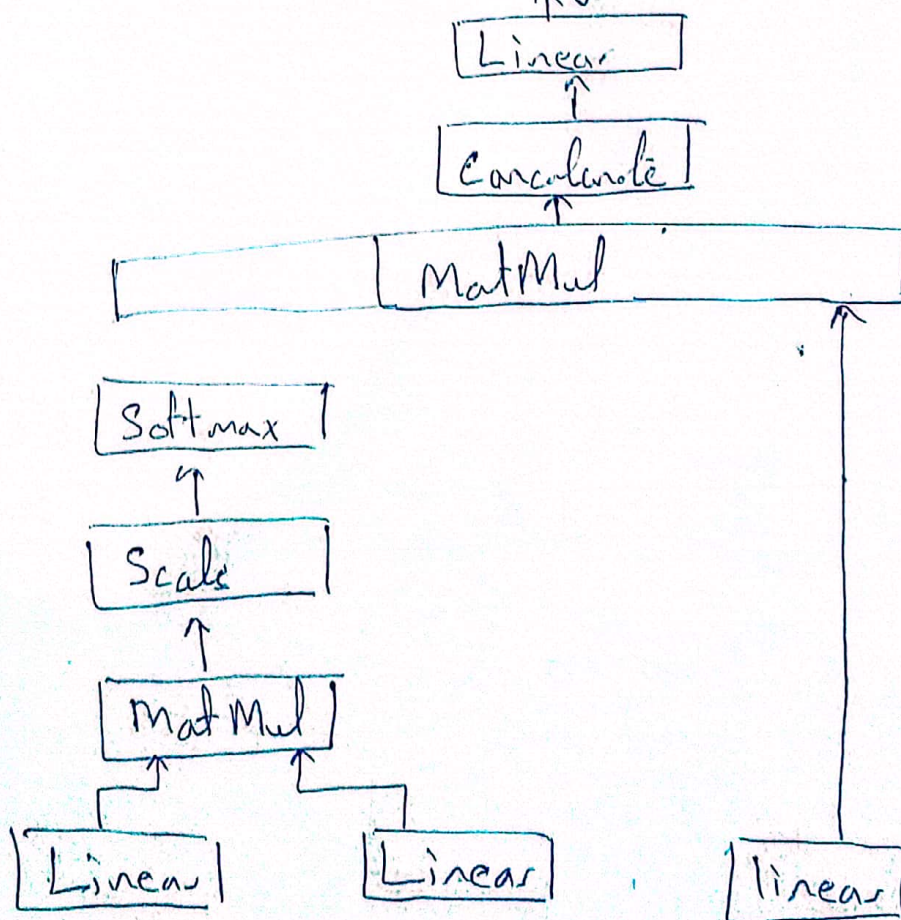
①
— Our mind will not pay equal attention to all the words in the sentence

Attention: focus on specific words in a sentence

"I went to a bank"

and

"I live by river bank"



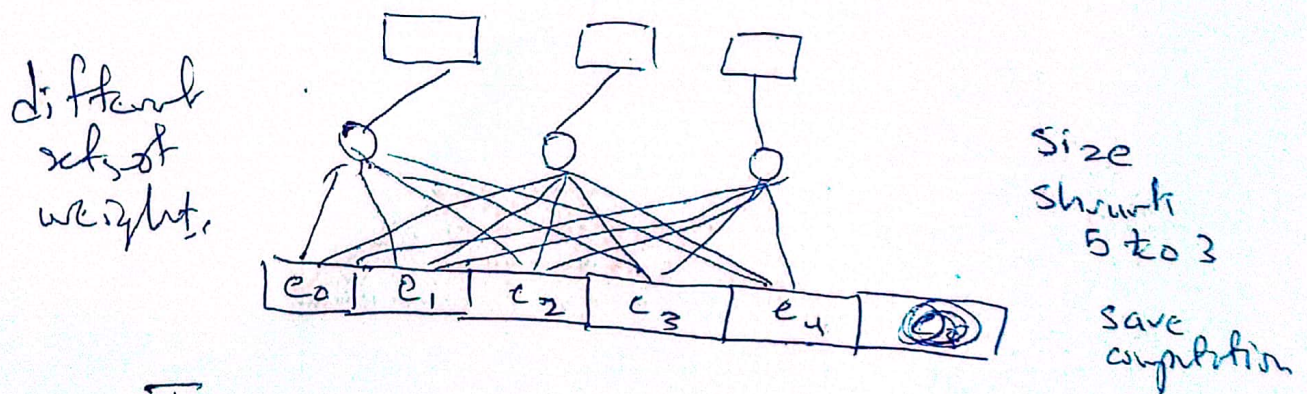
(2)

Three linear layers:

a combination of neurons without activation function

Purposes:

- 1) mapping inputs or to the outputs
- 2) Changing the dimensions.



These weights are fed as a matrix

Why 3 layers :- Each has a special function Q, K, V

request titles contents

Similarity b/w Q and K

Cosine Similarity ? -1 to +1

↓
opposite direction ↓
Same direction

$$\cos(A, B) = \frac{A \cdot B}{\text{Scaling}}$$

$$\text{Similarity using matrix} = \frac{A \cdot B^T}{\text{Scaling}}$$

$$\text{Similarity}(Q, K) = \frac{Q \cdot K^T}{\text{Scaling}}!$$

~~Relat.~~

Query layer: Position-aware Embeddings

- two more copies to key and value layer.
- Rely on training - self attention for improvement
- Multiply with linear layers
transpose with embeddy matrix.
each linear layer has its own weights

1)

Focus on Q and K

Compute dot product of Q and K

o/p of dot product \rightarrow Attention filter

Once the training period is over

$7 \times 3 \quad \cdot \quad 3 \times 7$

weights are more meaningful.

Scale attention Scores:

- Divide the attention scores
by the dimension of key vector.

- Apply Softmax $\text{Softmax} \left(\frac{Q \cdot K^T}{\sqrt{d_k}} \right)$

We now have

- Value matrix.

- Attention filter (K, Q).

multiply them $\text{Softmax} \left(\frac{Q \cdot K^T}{\sqrt{d_k}} \right) \cdot V$

analogy with Computer Vision.

Filled value matrix o/p
of multi-head attention.

Why need multi-head

Focus on different linguistic phenomenon.

Original Paper: 8 heads.

- Concatenate the three heads together.
- Apply linear to shrink size

⑥

Residual Connection

- Knowledge Reservation
- Vanishing gradient

(information is not lost)

Add & Norm:

Simple addition which gets fed

- Normalization: layer normalization standardization along ^{axis} # features

$$x_i = \frac{x_i^d - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

↓
does not cd 0

Linear layer with RELU

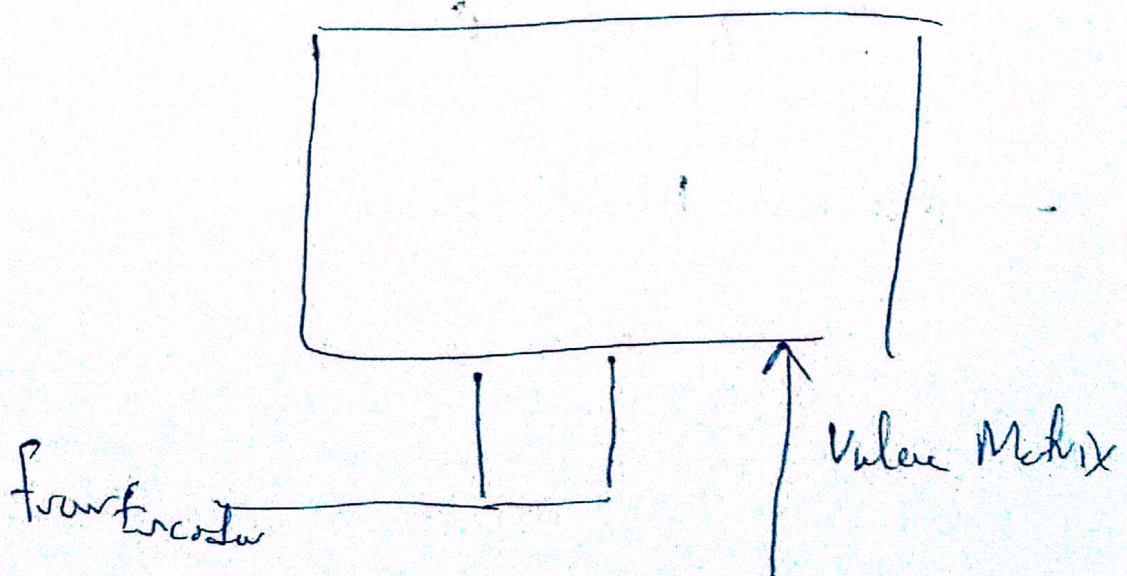
Encoder: Takes input and convert
to vectorized o/p

Decoder: two inputs

o/p of encoder is split in to
two copies K , & ^{2nd} input to decoder

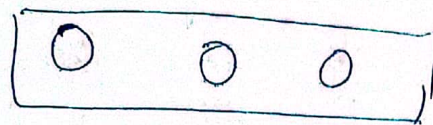
First input of Decoder:

- 1st word: $\langle \text{start} \rangle$
- Masked Multi-head
- Add & Norm
- Value Matrix



final linear layer

— fully connected Neurons
of classes.



dialogue generation
total vocabulary size

the final we need score
linear layer

— Flatten the ~~of~~ matrix

• We get single scores for each
word (logit)

— Apply Softmax

— Masked Multi-head

mask

(9)

During training
mask and show correct word
is shown teacher forcing

loss (prob(Pred), prob(Truth))

- Cross entropy loss
- Masking the input

Difference b/w training and testing

- We mask in both
- During training we apply actual words to decode
- During inference we input predicted words