

Lecture 10. Attention Models and Transformers

Lecture 10. Attention Models and Transformers

Alex Avdyushenko

Kazakh-British Technical University

November 12, 2022



Good afternoon, dear students. Today we are going to talk about attention and transformers — probably the most breakthrough architecture of the last maybe five years of deep learning. This is a very interesting approach, and again it is quite similar to human information processing.

Lecture 10. Attention Models and Transformers

└ Five-minutes block

- List the disadvantages of convolutional neural networks
- Write the formula for the simplest (vanilla) recurrent network module
- What kind of filters (gates) does LSTM have?

Ok, great, and as always we start with traditional five-minutes questions on the previous lecture. Please, write answers or send photos with them directly to me in private messages here in Teams or may be in Telegram, but choose only one option please :)

Lecture 10. Attention Models and Transformers

└ Disadvantages of vanilla Recurrent Neural Network

Now, let's remember vanilla recurrent neural network and its disadvantages.

We use one hidden vector

$$h_t = f_W(h_{t-1}, x_t)$$

As a function f_W we set a linear transformation with a non-linear component-wise "sigmoid".

$$h_t = \tanh(W_{\text{in}}h_{t-1} + W_{\text{in}}x_t)$$

$$y_t = W_{\text{out}}h_t$$

Disadvantages

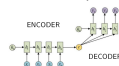
- ❑ input and output sequence lengths must match
- ❑ "reads" the input only from left to right, does not look ahead
- ❑ therefore it is not suitable for machine translation, question answering tasks and others

Lecture 10. Attention Models and Transformers

└ RNN for sequence synthesis (seq2seq)

RNN for sequence synthesis (seq2seq)

Recall

 $X = (x_1, \dots, x_n)$ — input sequence $Y = (y_1, \dots, y_m)$ — output sequence $c := h_n$ encodes all information about X to synthesize Y 

$$h_i = f_{\phi}(x_i, h_{i-1})$$

$$h'_i = f_{\omega}(h'_{i-1}, y_{i-1}, c)$$

$$y_i = f_{\psi}(h'_i, y_{i-1})$$

Ok, next step. How can we solve this problems? For example, we can use seq2seq architecture.

Lecture 10. Attention Models and Transformers

$$h_i = f_h(x_i, h_{i-1})$$

$$h'_i = f_{att}(h'_{i-1}, y_{i-1}, c)$$

$$y_i = f_y(h'_i, y_{i-1})$$

Disadvantages

- c remembers the end (h_n) better than the start
- the more n , the more difficult to pack all the information into vector c
- we should control the vanishing and explosions of the gradient
- RNN is difficult to parallelize

Question

How can you fix some of the problems above?

Hint

How do people perceive information?

Now, as is often the case we have new, more difficult problems with seq2seq architecture.

Lecture 10. Attention Models and Transformers

└ Let's count the number of passes

Now let's have some fun and count the number of passes in the video.

Lecture 10. Attention Models and Transformers

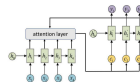
└ RNN with attention mechanism

RNN with attention mechanism

$a(h, h')$ is the similarity function of input h and output h' (for example, dot product or $\exp(h^T h')$ and others)

α_i = importance of input i for output t (attention score), $\sum_i \alpha_i = 1$

c_t = input context vector for output t



So, how can we use this mechanism for improving our network architecture?

We use function a , it is the similarity function of two vectors. And so on...

Lecture 10. Attention Models and Transformers

└ The Main Areas for Attention Models

Converting one sequence to another, i.e. seq2seq

- Machine translation
- Question answering
- Text summarization
- Annotation of images, videos (multimedia description)
- Speech recognition
- Speech synthesis

Sequence processing

- Classification of text documents
- Document Sentiment Analysis

Attention models allow us to solve a lot of problems from different areas.

Lecture 10. Attention Models and Transformers

└ Vector Similarity Functions

Vector Similarity Functions

$a(h, h') = h^T h'$ is the scalar (inner) product
 $a(h, h') = \exp(h^T h') \rightarrow$ norm becomes SoftMax
 $a(h, h') = h^T W h'$ with the learning parameter matrix W
 $a(h, h') = w^T \tanh(Uh + V h')$ is additive attention with w, U, V

Linear vector transformations query, key, value:

$$a(h_i, h_{j-1}) = (W_Q h_i)^T (W_K h_{j-1}) / \sqrt{d}$$

$$\alpha_i = \text{SoftMax}_j a(h_i, h_{j-1})$$

$$c_i = \sum_j \alpha_{ij} W_V h_j$$

$W_Q, d \times \text{dim}(h); W_K, d \times \text{dim}(h); W_V, d \times \text{dim}(h) \rightarrow$ linear neuron weight matrices, possible simplification
 $W_K \equiv W_V$

Zichao Hu. An introductory survey on attention mechanisms in NLP problems. 2018



Let's discuss in detail the math inside attention. The first one is Similarity Functions of two vectors. Of course in the simplest case it could be the scalar (inner) product. Also you could use the trainable parameters and our favorite linear transformation layer with non-linear sigmoid function and weights vector. But the most common procedure is to have query, key, value parameters matrices.

Lecture 10. Attention Models and Transformers

└ Attention formula

In general case we could say, that it is three-layer network, and attention mechanism chooses the best way to store information from the input to the one output vector c . And one important case needs to be highlighted here. It is called self-attention.

Attention formula

q is the query vector for which we want to calculate the context
 $K = (k_1, \dots, k_n)$ — key vectors compared with the query
 $V = (v_1, \dots, v_n)$ — value vectors forming the context
 $a(k_i, q)$ — score of relevance (similarity) of key k_i to query q
 c is the desired context vector relevant to the query

Attention Model

This is a 3-layer network that computes a convex combination of v_i values relevant to the query q :

$$c = \text{Attn}(q, K, V) = \sum_i v_i \text{SoftMax}_i(a(k_i, q))$$

$c_i = \text{Attn}(W_q h'_i, W_k H, W_v H)$ is the example from the previous slide, where $H = (h_1, \dots, h_n)$ are input vectors, h'_i is output

Self-attention:

$c_i = \text{Attn}(W_q h_i, W_k H, W_v H)$ is a special case when $h' \in H$

Lecture 10. Attention Models and Transformers

└ Multi-Head Attention

Multi-Head Attention

Idea: J different attention models are jointly trained to highlight various aspects of the input information (for example, parts of speech, syntax, idioms)

$$c^j = \text{Attn}(W_q^j q, W_k^j H, W_v^j H), j = 1, \dots, J$$

Variants of aggregating the output vector:

$$c = \frac{1}{J} \sum_{j=1}^J c^j \quad \text{-- averaging}$$

$$c = [c^1 \dots c^J] \quad \text{-- concatenation}$$

$$c = [c^1 \dots c^J] W \quad \text{-- to return to the desired dimension}$$

Regularization: to make aspects of attention as different as possible, rows $J \times n$ of matrices $A, v_j = \text{SoftMax}(W_k^j h; W_q^j q)$, decorrelated ($v_i^T v_j \rightarrow 0$) and sparse ($v_j^T v_j \rightarrow 1$):

$$\|AA^T - I\|^2 \rightarrow \min_{(W_k^j, W_q^j)}$$

Zhouhan Lin, Y Bengio et al. A structured self-attentive sentence embedding. 2017

Ok, and the last generalization in this lecture — is multi-head attention.
Here you do all the same, but having J different attention models :)



Ok, finally, what is it, who knows?

Lecture 10. Attention Models and Transformers

└ Additions and remarks

Additions and remarks

- a lot of such blocks $N = 6, h_1 \rightarrow \dots \rightarrow x_i$ are connected in series
- calculations can be easily parallelized in x_i
- it is possible to use pre-trained x_i embeddings
- it is possible to learn embeddings x_i of words $w_i \in V$
- Layer Normalization (LN). $x, \mu, \sigma \in \mathbb{R}^d$

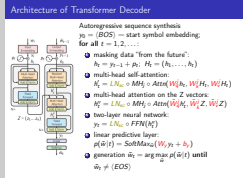
$$LN_L(x; \mu, \sigma) = \mu_L + \sigma_L \frac{x_L - \bar{x}}{\sigma_L}$$

$$\bar{x} = \frac{1}{d} \sum_{i=1}^d x_i, \sigma_L^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \bar{x})^2$$

Layer Normalization is quite easy — you train the mean and variance values of components of your vectors

Lecture 10. Attention Models and Transformers

└ Architecture of Transformer Decoder



Let's move on to the second part of transformer architecture — the decoder.

Temporary page!

\LaTeX was unable to guess the total number of pages correctly. As there was some unprocessed data that should have been added to the document, this extra page has been added to receive it.

If you rerun the document (without altering it) this surplus page will be removed, because \LaTeX now knows how many pages to expect for the final document.