Foundations for Natural Language Processing Improving Encoder-Decoder, Attention

Ivan Titov (with graphics/materials from Elena Voita)



Plan for today

Today

- Improving Encoder-Decoder models
- Modeling Attention in Encoder-Decoder
- Decomposable Attention Model (Pre-transformer)

The key problem with this approach



The key problem with this approach



<u>Problem</u>: fixed source representation is suboptimal:

- for the encoder, it is hard to compress the sentence;
- for the decoder, different information may be relevant at different steps.

Solution: modeling "attention"

Attention: Intuition

At every step, the decoder decide on which input tokens to focus





• receives attention input: a decoder state h_t and all encoder states s_1 , s_2 , ..., s_m ;





At each decoder step, attention

- receives attention input: a decoder state h_t and all encoder states s_1 , s_2 , ..., s_m ;
- computes attention scores

For each encoder state s_k , attention computes its "relevance" for this decoder state h_t . Formally, it applies an attention function which receives one decoder state and one encoder state and returns a scalar value $score(h_t, s_k)$;





At each decoder step, attention

- receives attention input: a decoder state h_t and all encoder states s_1 , s_2 , ..., s_m ;
- computes attention scores

For each encoder state s_k , attention computes its "relevance" for this decoder state h_t . Formally, it applies an attention function which receives one decoder state and one encoder state and returns a scalar value $score(h_t, s_k)$;

 computes attention weights: a probability distribution - softmax applied to attention scores;



At each decoder step, attention

- receives attention input: a decoder state h_t and all encoder states s_1 , s_2 , ..., s_m ;
- computes attention scores

For each encoder state s_k , attention computes its "relevance" for this decoder state h_t . Formally, it applies an attention function which receives one decoder state and one encoder state and returns a scalar value $score(h_t, s_k)$;

- computes attention weights: a probability distribution softmax applied to attention scores;
- computes attention output: the weighted sum of encoder states with attention weights.



















How do we compute attention scores? Alternatives





Encoder-Decoder variants: Bahdanau



Source context c^(t) does not know what word y^(t-1) was chosen on the previous step

https://arxiv.org/pdf/1409.0473.pdf

Encoder-Decoder variants: Luong



The context vector $c^{(t)}$ is not affecting future steps

Is attention interpretable?

For Bahdanau – before even chosing the previous token

- The attention may be perceived as modelling *alignment* between input and output words, but
 - attention is computed **before** choosing the target token (i.e. the choice of the token does not influence attention)
 - sometimes states encode something unexpected (e.g., <eos> may capture the general topic of the sentence or mark than nothing is relevant)
 - attention is (?) not an explanation



Enoder-decoders and attention are a very general idea

For example, caption generation (encoder process the image, not a sentence)



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A <u>stop</u> sign is on a road with a mountain in the background.



A little <u>girl</u> sitting on a bed with a teddy bear.



A group of <u>people</u> sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

Attention is not necessarily faithful

Generally, you can make a change to model parameters such that attention score a_k for a state s_k decreases N-fold ($a'_k = s_k/N$) N) whereas s_k magnitude increases N-fold ($s'_k = s_k N$)

$$c_{\uparrow}^{(t)} = a_1^{(t)}s_1 + a_2^{(t)}s_2 + \dots + a_m^{(t)}s_m = \sum_{k=1}^m a_k^{(t)}s_k$$

Source context for decoder step t"

$$a_{k}^{(t)} = \frac{\exp(\text{score}(h_{t}, s_{k}))}{\sum_{i=1}^{m} \exp(\text{score}(h_{t}, s_{i}))}, k = 1..m$$

* "attention weight for source token k at decoder step t"

Since the attention output is the weighted sum of embedding encoders states, the attention output (c) will not be affected change

Attention is not necessarily faithful

Also, if we use a multilayer encoder (and we usually will), there is no guarantee that a state in a *n*-th layers (n > I) encodes the n-th token



There are attribution methods (e.g., norm x gradient, layer-wise relevant propagation, integrated gradients, attention flow, zero valuing, ...) which attempt to address these issues, but they also have their own pitfalls

Generally: attention cannot be trusted blindly but it can signal some issues with our models

Attribution can help us detect issues with our models

Even if not perfectly faithful attention (and attribution techniques) help us detect issues with our models and/or data

Detecting pneumonia from an x-ray, the model 'looks' at the corner of an image, rather than at lungs. Any thoughts on why?



Attribution can help us detect issues with our models

The model decides that pronoun 'she' refers to 'Nurse' rather than Doctor (gender bias)



https://www.comet.com/site/blog/explainable-ai-for-transformers/

It is not an encoder-decoder attention, it is attention with a language model (i.e. what we will consider on Friday) Transformer

Attention is all you need

	Seq2seq without attention	Seq2seq with attention	Transformer
processing within <mark>encoder</mark>	RNN/CNN	RNN/CNN	attention
processing within <mark>decoder</mark>	RNN/CNN	RNN/CNN	attention
decoder-encoder interaction	static fixed- sized vector	attention	attention

Decomposable attention (pre-Transformer)



- It is sunny outside.

un 2017)3762v1 [cs.CL]

Attention Is All You Need

Ashish Vaswani* Noam Shazeer* Niki Parmar* Jakob Uszkoreit* Google Brain Google Brain Google Research Google Research avaswani@google.com noam@google.com nikip@google.com usz@google.com Llion Jones* Aidan N. Gomez*† Łukasz Kaiser* Google Research University of Toronto Google Brain llion@google.com aidan@cs.toronto.edu lukaszkaiser@google.com Illia Polosukhin* illia.polosukhin@gmail.com Abstract The dominant sequence transduction models are based on complex recurrent or

convolutional neural networks in an encoder-decoder configuration. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including

1 Introduction

Natural Language Inference (aka textual entailment)

Classification task:

<u>Premise</u>: Bob is in his room, but because of the thunder and lightning outside, he cannot sleep.

<u>Hypothesis</u>: It is sunny outside.

Is H entailed by P (i.e. logically follows), contradicts P or neither ?

Natural Language Inference (aka textual entailment)

Classification task:

<u>Premise</u>: Bob is in his room, but because of the thunder and lightning outside, he cannot sleep.

<u>Hypothesis</u>: It is sunny outside.

Is H entailed by P (i.e. logically follows), contradicts P or neither ?

The idea: the prediction can be done (or approximated) by matching / comparing parts of H and P

Classification task



Steps: (1) attend

- (1) attend / align (F),
- (2) compare (G)
- (3) aggregate and classify (H)



Representation of tokens in H aligned to i-th token in P:

$$\beta_i := \sum_{j=1}^{\ell_b} \frac{\exp(e_{ij})}{\sum_{k=1}^{\ell_b} \exp(e_{ik})} b_j$$

The same in the opposite direction:

$$\alpha_j := \sum_{i=1}^{\ell_a} \frac{\exp(e_{ij})}{\sum_{k=1}^{\ell_a} \exp(e_{kj})} \bar{a}_i$$

Step 2: Compare



$$\mathbf{v}_{1,i} := G([a_i, \beta_i]) \qquad \leftarrow \qquad \text{a representation of i-th token in H,} \\ \mathbf{v}_{2,j} := G([\overline{b}_j, \alpha_j]) \qquad \leftarrow \qquad \frac{\text{informed}}{\text{informed}} \text{ by all tokens in T}$$

G is just a feedforward neural network, v is a vector

Note: the input to G is concatenation, so not invariant to switching the arguments around

Step 2: Aggregate / classify





 $H([\mathbf{v}_1,\mathbf{v}_2])$

Just a linear classification layer, followed by softmax to get a distribution over 3 classes

All together now





All together now



This simple and very fast to train model without RNNs or CNNs components achieved very strong results on the NLP task

(a couple of small extra tricks were needed to inform the model about word position + to better handle unaligned tokens; we will discuss their counterparts in the Transformer)

Summary

- The general idea of attention
- Attention between encoder and decoder
- Does attention represent an alignment? / is it interpretable?
- Many applications in NLP and beyond, often as a crucial component of Transformers but not only
- We made the first step towards Transformers which we will discuss after the flexible learning week