Machine Learning, Artificial Intelligence, and Natural Language Processing

Marcelo C. Medeiros

Department of Economics The University of Illinois at Urbana-Champaign





Attention

Transformers





Searching for a better word embedding





Searching for a better word embedding





ECONOMETRICS LAR

Searching for a better word embedding



Attention is all you need

Attention Is All You Need

Ashish Vaswani* Google Brain avaswani@google.com Noam Shazeer* Google Brain noam@google.com Niki Parmar* Google Research nikip@google.com Jakob Uszkoreit* Google Research usz@google.com

Llion Jones* Google Research llion@google.com Aidan N. Gomez^{*}[†] University of Toronto aidan@cs.toronto.edu Łukasz Kaiser* Google Brain lukaszkaiser@google.com

Illia Polosukhin * ‡

illia.polosukhin@gmail.com





► For example,

The chicken didn't cross the road because it was too tired



► For example,

The chicken didn't cross the road because it was too tired

 \ast What is the meaning represented in the static embedding for "it"?



For example,

The chicken didn't cross the road because it was too tired

* What is the meaning represented in the static embedding for "it"?

Consider now

The chicken didn't cross the road because it was too wide



For example,

The chicken didn't cross the road because it was too tired

 ${\ensuremath{\ast}}$ What is the meaning represented in the static embedding for "it"?

Consider now

The chicken didn't cross the road because it was too wide

* Is the meaning the same?

Problem: The representation of the meaning of a word should be different in different contexts

Contextual Embedding: each word has a different vector that expresses different meanings depending on the surrounding words

Solution: Attention mechanism.



The chicken didn't cross the road because it ...

▶ What should be the properties of "it"?

The chicken didn't cross the road because it was too tired The chicken didn't cross the road because it was too wide

At this point in the sentence, it's probably referring to either the chicken or the street



Build up the contextual embedding from a word by selectively integrating information from all the neighboring words

► A word "attends to" some neighboring words more than others



Attention is a mechanism for helping compute the embedding for a token by selectively attending to and integrating information from surrounding tokens

Mathematically: a method for computing a ("smart") weighted sum of vectors









Simplified version

- ▶ Sequence of token embeddings: $x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_i$
- Measure distance between two tokens: inner product

$$\mathsf{score}(oldsymbol{x}_i,oldsymbol{x}_j) = oldsymbol{x}_i'oldsymbol{x}_j$$

► Transform the score into weights between 0 and 1:

$$\alpha_{ij} = \operatorname{softmax}(\operatorname{score}(\boldsymbol{x}_i, \boldsymbol{x}_j)) = \frac{e^{\operatorname{score}(\boldsymbol{x}_i, \boldsymbol{x}_j)}}{\sum_{k=1}^8 e^{\operatorname{score}(\boldsymbol{x}_i, \boldsymbol{x}_k)}}, \quad j \le n$$



$$oldsymbol{a}_i = \sum_{j \leq i} lpha_{ij} oldsymbol{x}_j$$



- High-level idea: instead of using embedding vectors directly, we will represent three separate roles each vector \boldsymbol{x}_i plays:
 - 1. query: element that is being compared to the preceding inputs
 - 2. key: preceding input that is being compared to the current element to determine a similarity
 - 3. value: a value of a preceding element that gets weighted and summed



- Matrices to project each vector x_i into a representation of its role as query, key, or value
 - st Query: Q
 - * Key: *K*
 - * Value: $oldsymbol{V}$



$$oldsymbol{q}_i = oldsymbol{Q} oldsymbol{x}_i \quad oldsymbol{k}_i = oldsymbol{K} oldsymbol{x}_i \quad oldsymbol{v}_i = oldsymbol{V} oldsymbol{x}_i$$

The three matrices are learned from data

▶ Take the representation in terms of query, key, and value matrices

$$oldsymbol{q}_i = oldsymbol{Q} oldsymbol{x}_i \quad oldsymbol{k}_i = oldsymbol{K} oldsymbol{x}_i \quad oldsymbol{v}_i = oldsymbol{V} oldsymbol{x}_i$$

 \blacktriangleright Similarity (score) between $m{x}_i$ and prior token $m{x}_j$: $m{q}_i'm{k}_j/\sqrt{d}$

 \blacktriangleright d is the dimension of the embeddings

► Therefore,

$$\alpha_{ij} = \operatorname{softmax}(\boldsymbol{q}'_i \boldsymbol{k}_j), \quad j \leq i$$

Attention:

$$\boldsymbol{a}_i = \sum_{j \leq i} lpha_{ij} \boldsymbol{v}_j$$



▶ Instead of one attention head, we'll have lots of them

Intuition: each head might be attending to the context for different purposes





Multi-head attention

For each head h:

$$oldsymbol{q}_i^h = oldsymbol{Q}^h oldsymbol{x}_i \quad oldsymbol{k}_i^h = oldsymbol{K}^h oldsymbol{x}_i \quad oldsymbol{v}_i^h = oldsymbol{V}^h oldsymbol{x}_i$$

Similarity (score) between x^h_i and prior token x_j: q^{h'}_ik^h_j/√d_h
d_h = d/H, H is the number of heads

► Therefore,

$$\alpha_{ij}^h = \mathsf{softmax}(\boldsymbol{q}^{h}{}'_i \boldsymbol{k}^h_j), \quad \mathsf{and} \quad \boldsymbol{a}_i^h = \sum_{j \leq i} \alpha_{ij}^h \boldsymbol{v}_j^h$$

Hence,

$$oldsymbol{a}_i = \left(oldsymbol{a}_i^1 \cdots oldsymbol{a}_i^H
ight)oldsymbol{O}$$



Attention = contextual embedding





Attention

Transformers





Goal: given a sequence of words, predict which word will come next

▶ It is based on NN architectures plus attention mechanism





Transformers



22