

Neural Nets (NN), глубокое обучение,
рекуррентные сети, LSTM и пр. Работа с
временными рядами (прогноз) и с текстами.
Трансформеры, BERT

Оленев Роман, Самарин Игорь

Санкт-Петербургский государственный университет
Кафедра статистического моделирования

Санкт-Петербург, 2024

Обработка текстов на естественном языке (Natural Language Processing, NLP) — общее направление искусственного интеллекта и математической лингвистики. Оно изучает проблемы компьютерного анализа и синтеза текстов на естественных языках.

Подходы к представлению текста:

1. Bag of words
2. Bag of Ngramms
3. TF-IDF
4. Word Embeddings

Word2Vec

Thou shalt not make a machine in the likeness of a human mind

thou	shalt	not	make	a	machine	in	the	...
------	-------	-----	------	---	---------	----	-----	-----

input word	target word
not	thou
not	shalt
not	make
not	a

Figure: SkipGram

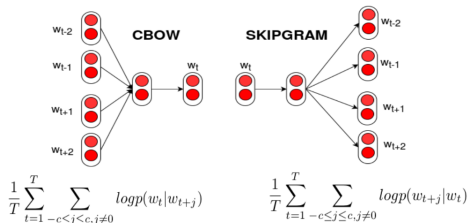
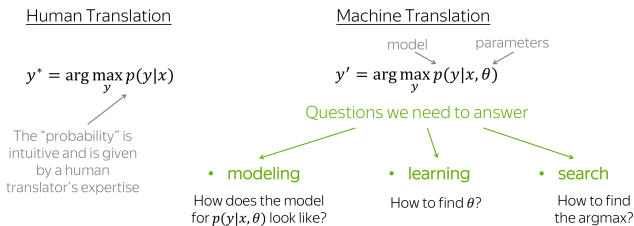


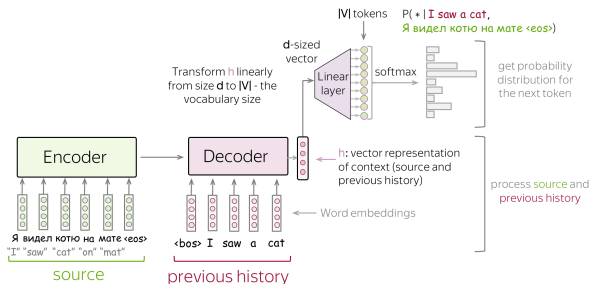
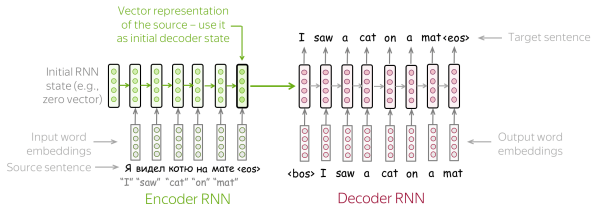
Figure: Word2Vec

Машинный перевод. Постановка задачи

Пусть есть входная последовательность x_1, x_2, \dots, x_m и выходная последовательность y_1, y_2, \dots, y_n . Тогда хотим максимизировать $p(y|x)$: $y^* = \arg \max_y p(y|x)$.



Машинный перевод. 2 RNNs



Машинный перевод. Обучение и предсказание

В качестве функции потерь используется кросс-энтропия. При обучении на вход случайным образом можно подавать не только предсказанное на предыдущем шаге слово, но и верное слово.

$$Loss(p^*, p) = -p^* \log(p) = - \sum_{i=1}^{|V|} p_i^* \log(p_i).$$

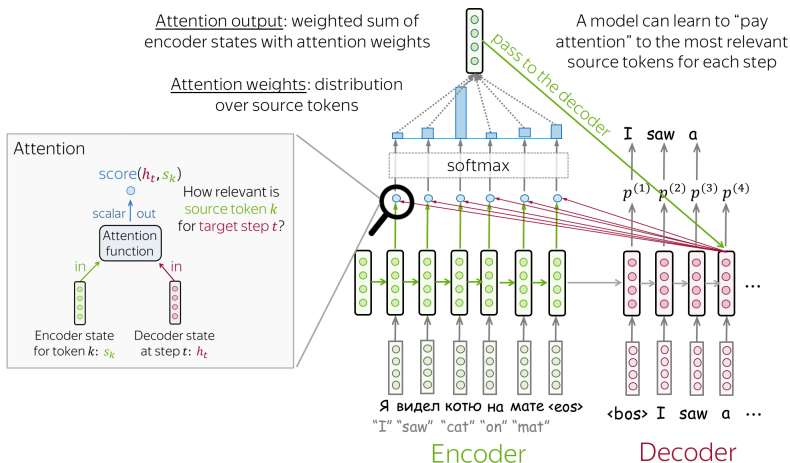
$$Loss(p^*, p) = -\log(p_{y_t}) = -\log(p(y_t | y_{< t}, x)).$$

Прогноз:

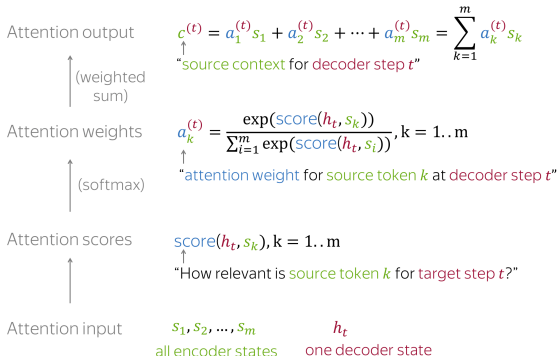
1. Greedy Decoding (предсказываем самое вероятное)
2. Beam Search (на каждом шаге храним несколько вероятных последовательностей)

Механизм Attention

Проблема: невозможно сохранить весь контекст в одном векторе. Решением проблемы является механизм внимания.



Механизм Attention



Dot-product

$$\begin{bmatrix} \text{red} & \text{red} & \text{red} & \text{red} \end{bmatrix} \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} s_k$$

$$\text{score}(h_t, s_k) = h_t^T s_k$$

Bilinear

$$\begin{bmatrix} \text{red} & \text{red} & \text{red} & \text{red} \end{bmatrix} \times \begin{bmatrix} \text{blue} \\ \text{blue} \\ \text{blue} \\ \text{blue} \end{bmatrix} W \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} s_k$$

$$\text{score}(h_t, s_k) = h_t^T W s_k$$

Multi-Layer Perceptron

$$\begin{bmatrix} \text{blue} & \text{blue} & \text{blue} & \text{blue} \end{bmatrix} w_2^T \times \tanh \left[\begin{bmatrix} \text{blue} \\ \text{blue} \\ \text{blue} \\ \text{blue} \end{bmatrix} W_1 \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} s_k \right] h_t$$

$$\text{score}(h_t, s_k) = w_2^T \cdot \tanh(W_1 [h_t, s_k])$$

Трансформеры

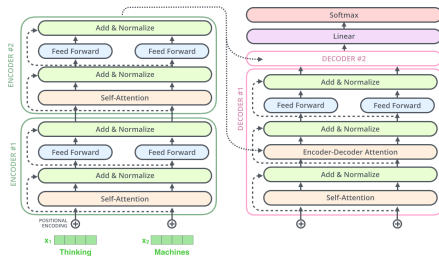
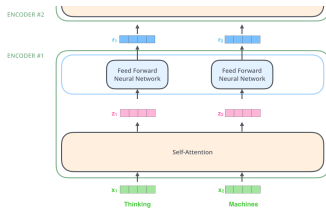
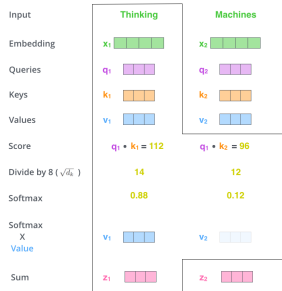
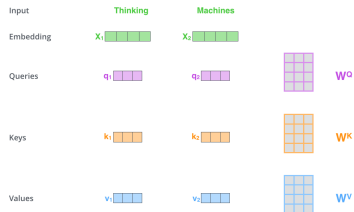


Figure: Архитектура Трансформер



Трансформеры. Self Attention



Трансформеры. Self Attention

$$X \times W^Q = Q$$

$$X \times W^K = K$$

$$X \times W^V = V$$

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^O that was trained jointly with the model

\times



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$= Z$$

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) \times V = Z$$

Трансформеры. Self Attention

1) This is our input sentence*

2) We embed each word*

3) Split into 8 heads.
We multiply X or R with weight matrices

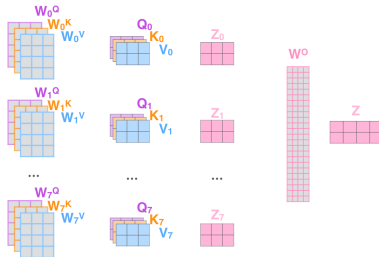
4) Calculate attention using the resulting $Q/K/V$ matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

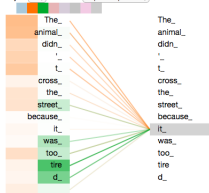
Thinking
Machines



* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



Layer: 5 Attention: Input - Input



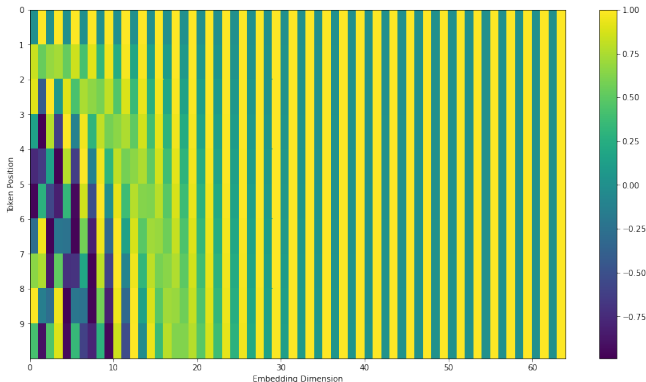
Трансформеры. Positional encoding

Можем закодировать информацию о позиции слова следующим образом:

$$PE_{pos,2i} = \sin(pos/10000^{2i/d_{model}}),$$

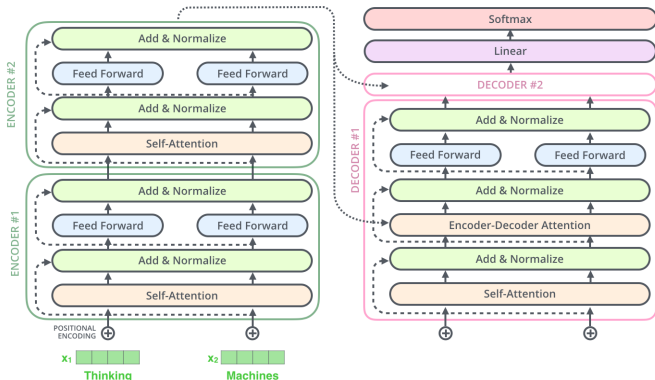
$$PE_{pos,2i+1} = \cos(pos/10000^{2i/d_{model}}),$$

после чего можем добавить получившийся энкодинг к эмбедингам.



Трансформеры. Masked Attention

Входная последовательность декодера маскируется таким образом, что Attention Score для каждого токена считается только с токенами, которые находятся перед ним. Для encoder-decoder attention матрица Q берется из декодера, а матрицы K и V — из энкодера.



Трансформеры. Train & Inference

Последний линейный слой преобразовывает выход декодера из размерности (batch size, sequence length, embedding size) в размерность (batch size, sequence length, vocabulary size), после чего считается

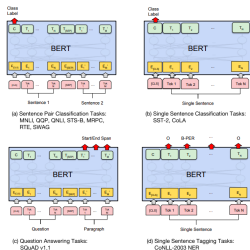
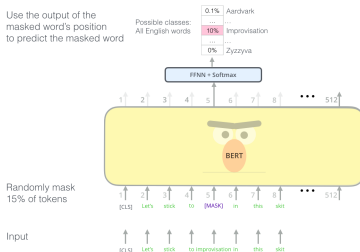
$$Loss = - \sum_{t=1}^T \log(p_{y_t}).$$

Для прогнозирования подаем декодеру <SOS> (start of sequence) токен, после чего добавляем ко входу декодера предсказанный токен и повторяем процесс до того, как декодер выведет <EOS> (end of sequence) токен.

Преимущества архитектуры:

1. Работает быстрее, так как нет рекуррентных блоков.
2. Много вычислений, которые можно распараллелить.

BERT (Bidirectional Encoder Representations from Transformers)



Transfer Learning

