

# XLNet: Generalized Autoregressive Pretraining for Language Understanding

Zhilin Yang et al.  
Speaker: Fedor Noskov

MIPT, 792

October, 2020

# AR modeling vs. autoencoding

В случае с авторегрессионной языковой моделью логарифм правдоподобия:

$$\log p_{\theta}(\mathbf{x}) = \sum_{t=1}^T \log p_{\theta}(x_t | \mathbf{x}_{<t}), \quad (1)$$

т.е. модель для предсказания нового слова смотрит только на предыдущие:

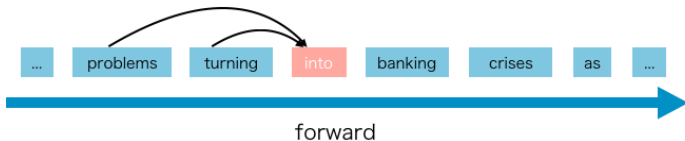


Рис. 1: (source)

# AR modeling vs. autoencoding

В случае BERT предполагаем, что при условии некоторого контекста  $\hat{\mathbf{x}}$  остальные токены появляются независимо:

$$\log p_{\theta}(\mathbf{x}_{masked}|\hat{\mathbf{x}}) = \sum_{t=1}^T \mathbb{I}\{x_t \text{ is unknown}\} \log p_{\theta}(x_t|\hat{\mathbf{x}}) \quad (2)$$



Рис. 2: (source)

## AR modeling vs. autoencoding

Но в реальности это предположение неверно! Решение: заметим, что порядок обуславливания неважен.

- |    |     |            |       |
|----|-----|------------|-------|
| 1) | --- | ???        | ---   |
| 2) | --- | немедленно | ???   |
| 3) | ??? | немедленно | выпил |
| 4) | И   | немедленно | выпил |

Т.е. для любой перестановки  $\mathbf{z}$ :

$$\log p_{\theta}(\mathbf{x}) = \sum_{t=1}^T \log p_{\theta}(x_{z_t} | \mathbf{x}_{<z_t}) \quad (3)$$

## AR modeling vs. autoencoding

Будем максимизировать следующую величину:

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{Z}_T} \sum_{t=1}^T \log p_{\theta}(x_t | \mathbf{x}_{<_z t}) \quad (4)$$

А параллельность?..

# Two-Stream Self-Attention

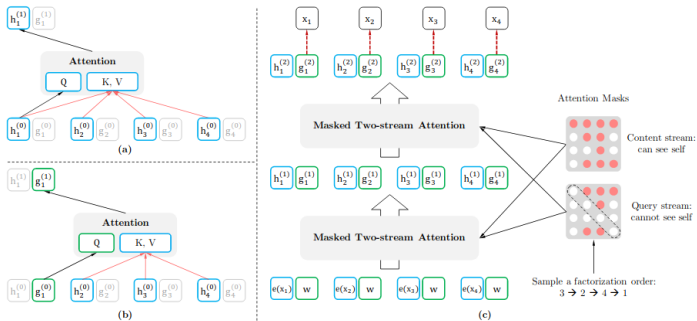


Рис. 3: (source)

## Тонкости обучения

- ▶ Обучать с учетом всех перестановок сложно. По этой причине будем в последовательности предсказывать только два последних токена
- ▶ Переставлять длинную последовательность сложно. Будем разбивать ее на несколько частей и для каждой части передавать в слой трансформера вектор контекста предыдущих. Перестановки будем совершать только внутри этих частей
- ▶ Стоимость обучения около 245 тыс. долларов

# Результаты

Model	IMDB	Yelp-2	Yelp-5	DBpedia	AG	Amazon-2	Amazon-5
CNN [15]	-	2.90	32.39	0.84	6.57	3.79	36.24
DPCNN [15]	-	2.64	30.58	0.88	6.87	3.32	34.81
Mixed VAT [31, 23]	4.32	-	-	0.70	4.95	-	-
ULMFIT [14]	4.6	2.16	29.98	0.80	5.01	-	-
BERT [35]	4.51	1.89	29.32	0.64	-	2.63	34.17
XLNet	<b>3.20</b>	<b>1.37</b>	<b>27.05</b>	<b>0.60</b>	<b>4.45</b>	<b>2.11</b>	<b>31.67</b>

Рис. 4: Классификационный датасеты