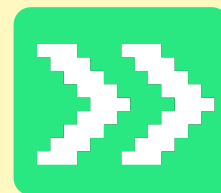


Информативные векторные представления в машинном обучении

YOUNG & YANDEX



Радослав Нейчев

Выпускник и преподаватель ШАД и МФТИ,
руководитель группы ML-разработки в Яндексе,
сооснователь girafe-ai



Содержание

01

Правдоподобие / Likelihood

02

Векторное представление текста

03

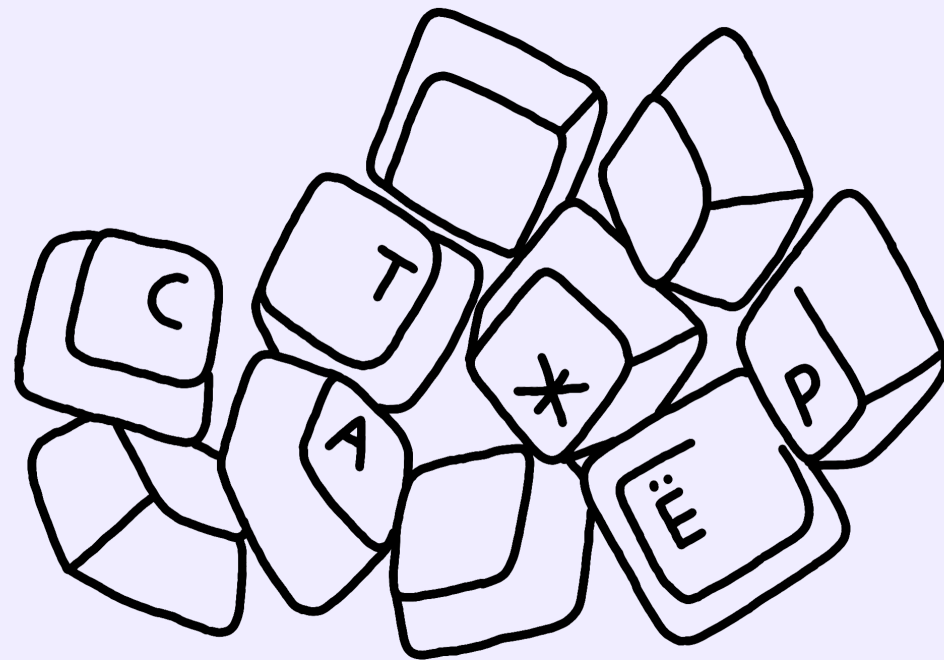
Информативные векторные представления –
эмбединги

04

Построение эмбедингов для слов
word2vec

Правдоподобие / Likelihood

01



Правдоподобие / Likelihood

Пусть задана выборка X, Y и модель с параметрами θ .

Правдоподобием будем называть

$$L(\theta|X, Y) = P(X, Y|\theta) \rightarrow \max_{\theta}$$

помним про i.i.d.

$$P(X, Y|\theta) = \prod_i P(x_i, y_i|\theta)$$

$$\log L(\theta|X, Y) = \log P(X, Y|\theta) = \sum_i \log P(x_i, y_i|\theta)$$

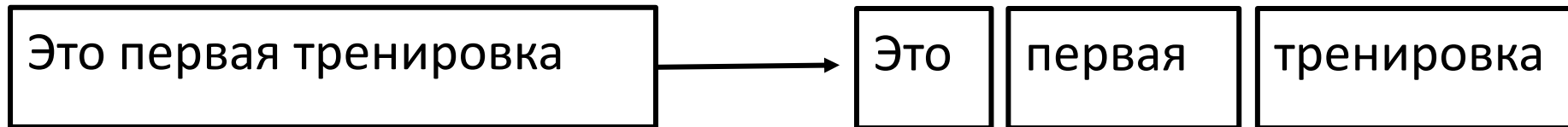
Векторные представления для текста

02



Токенизация

Токен – атомарный элемент последовательности.



Токеном может быть как слово, так и символ, и морфема – это вопрос договоренности в каждой задаче.

Мешок слов – Bag-of-Words

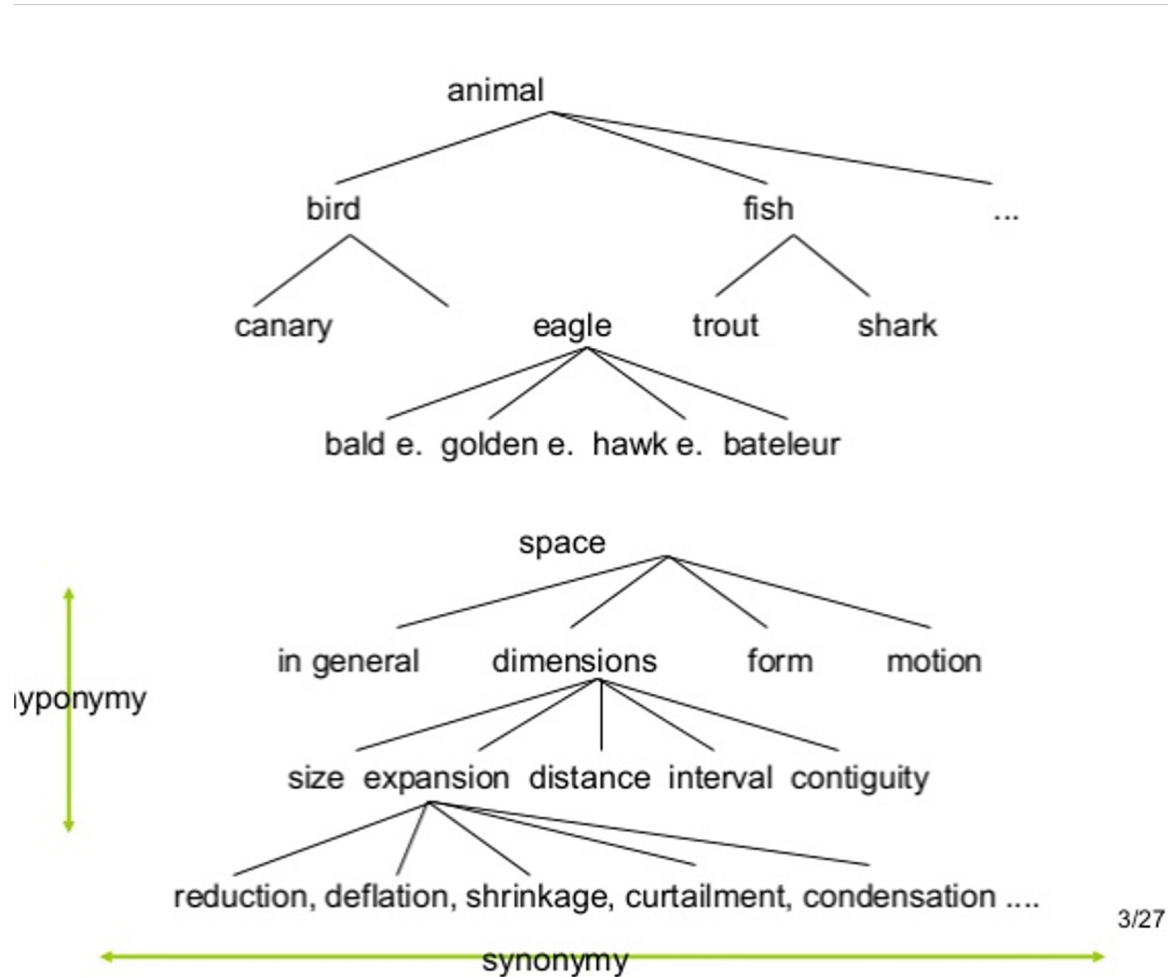
Как-то раз купил мужик шляпу – а она ему как раз

день	как	то	раз	купил	мужик	шляпу	а	она	ему	кот	гав
0	2	1	2	1	1	1	1	1	1	0	0

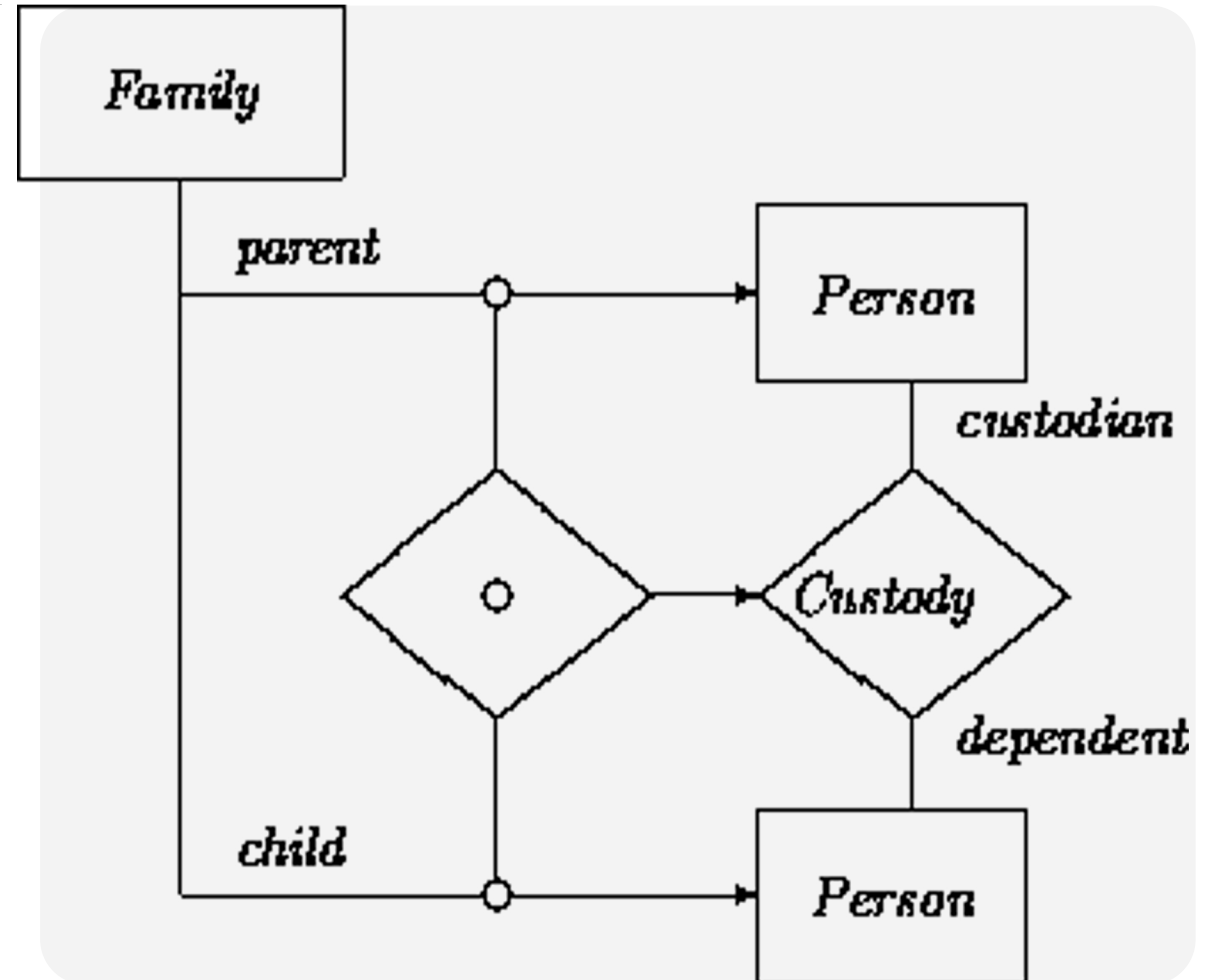
Минусы:

- Теряем информацию о порядке слов
- Векторы высокой размерности
- Векторы крайне разреженные
- Разные формы слов воспринимаются как разные слова

WordNet



3/27



Чем предобрабатывать тексты?

- NLTK
 - `nltk.stem.SnowballStemmer`
 - `nltk.stem.PorterStemmer`
 - `nltk.stem.WordNetLemmatizer`
 - `nltk.corpus.stopwords`
- BeautifulSoup (for parsing HTML)
- Regular Expressions (`import re`)
- Pymorphy2

TF-IDF

- **Term Frequency (tf)**: gives us the frequency of the word in each document in the corpus.

$$\text{tf}(t, d) = f_{t, d}$$

- **Inverse Document Frequency (idf)**: used to calculate the weight of rare words across all documents in the corpus. The words that occur rarely in the corpus have a high IDF score.

$$\text{idf}(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

N : total number of documents in the corpus $N = |D|$

$|\{d \in D : t \in d\}|$: number of documents where the term t

TF-IDF

- *Sentence A*: The car is driven on the road.
- *Sentence B*: The truck is driven on the highway.

(each sentence is a separate document)

TF-IDF

Word	TF		IDF	TF * IDF	
	A	B		A	B
The	1/7	1/7			
Car	1/7	0			
Truck	0	1/7			
Is	1/7	1/7			
Driven	1/7	1/7			
On	1/7	1/7			
The	1/7	1/7			
Road	1/7	0			
Highway	0	1/7			

TF-IDF

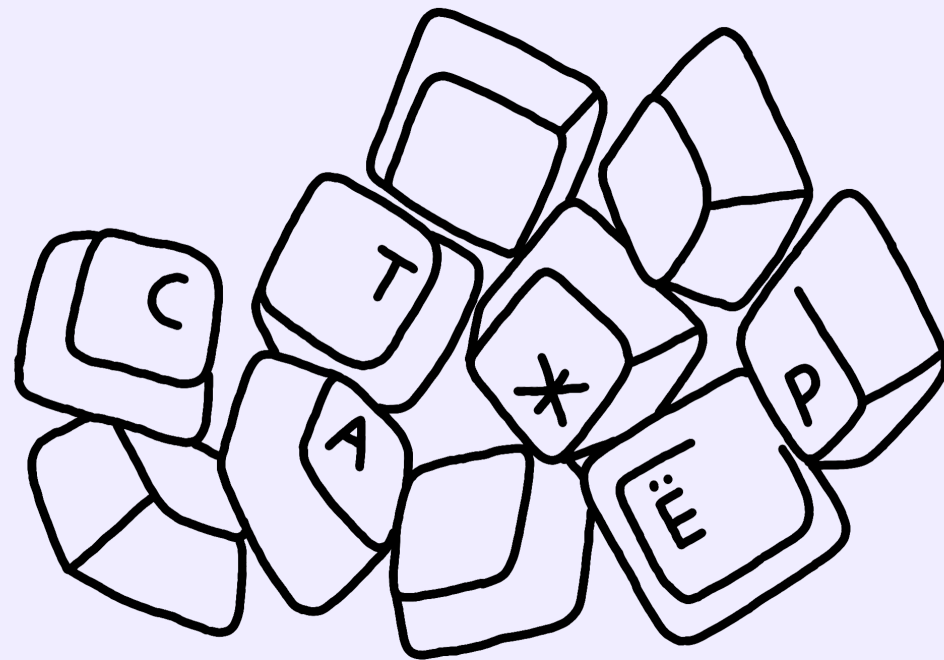
Word	TF		IDF	TF * IDF	
	A	B		A	B
The	1/7	1/7	$\log(2/2)=0$		
Car	1/7	0	$\log(2/1)=0.3$		
Truck	0	1/7	$\log(2/1)=0.3$		
Is	1/7	1/7	$\log(2/2)=0$		
Driven	1/7	1/7	$\log(2/2)=0$		
On	1/7	1/7	$\log(2/2)=0$		
The	1/7	1/7	$\log(2/2)=0$		
Road	1/7	0	$\log(2/1)=0.3$		
Highway	0	1/7	$\log(2/1)=0.3$		

TF-IDF

Word	TF		IDF	TF * IDF	
	A	B		A	B
The	1/7	1/7	$\log(2/2)=0$	0	0
Car	1/7	0	$\log(2/1)=0.3$	0.043	0
Truck	0	1/7	$\log(2/1)=0.3$	0	0.043
Is	1/7	1/7	$\log(2/2)=0$	0	0
Driven	1/7	1/7	$\log(2/2)=0$	0	0
On	1/7	1/7	$\log(2/2)=0$	0	0
The	1/7	1/7	$\log(2/2)=0$	0	0
Road	1/7	0	$\log(2/1)=0.3$	0.043	0
Highway	0	1/7	$\log(2/1)=0.3$	0	0.043

Информативные векторные представления – эмбединги

03



Простейший вариант – one-hot

Проблемы с one-hot:

- Высокая размерность
- Разреженность
- Все векторы взаимно ортогональны

One-hot vectors:

Rome

Paris

word V

Rome

=

[1,

0,

0,

0,

0,

0,

...

0]

Paris

=

[0,

1,

0,

0,

0,

0,

...

0]

Italy

=

[0,

0,

1,

0,

0,

0,

...

0]

France

=

[0,

0,

0,

1,

0,

0,

...

0]

Как определить значение слов, если
нельзя пользоваться словами?

“You shall know a word by the company it keeps”

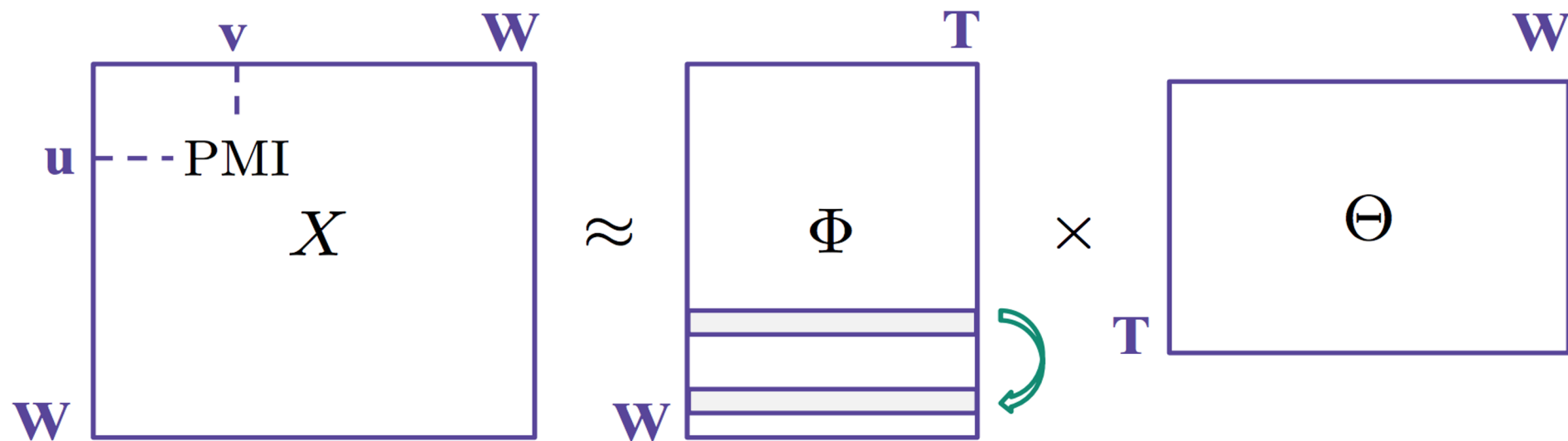
J. R. Firth, 1957: 11

...government debt problems turning into banking crises as happened in 2009...

...saying that Europe needs unified banking regulation to replace the hodgepodge...

...India has just given its banking system a shot in the arm...

Вспомним матричные разложения



Построение эмбедингов для слов

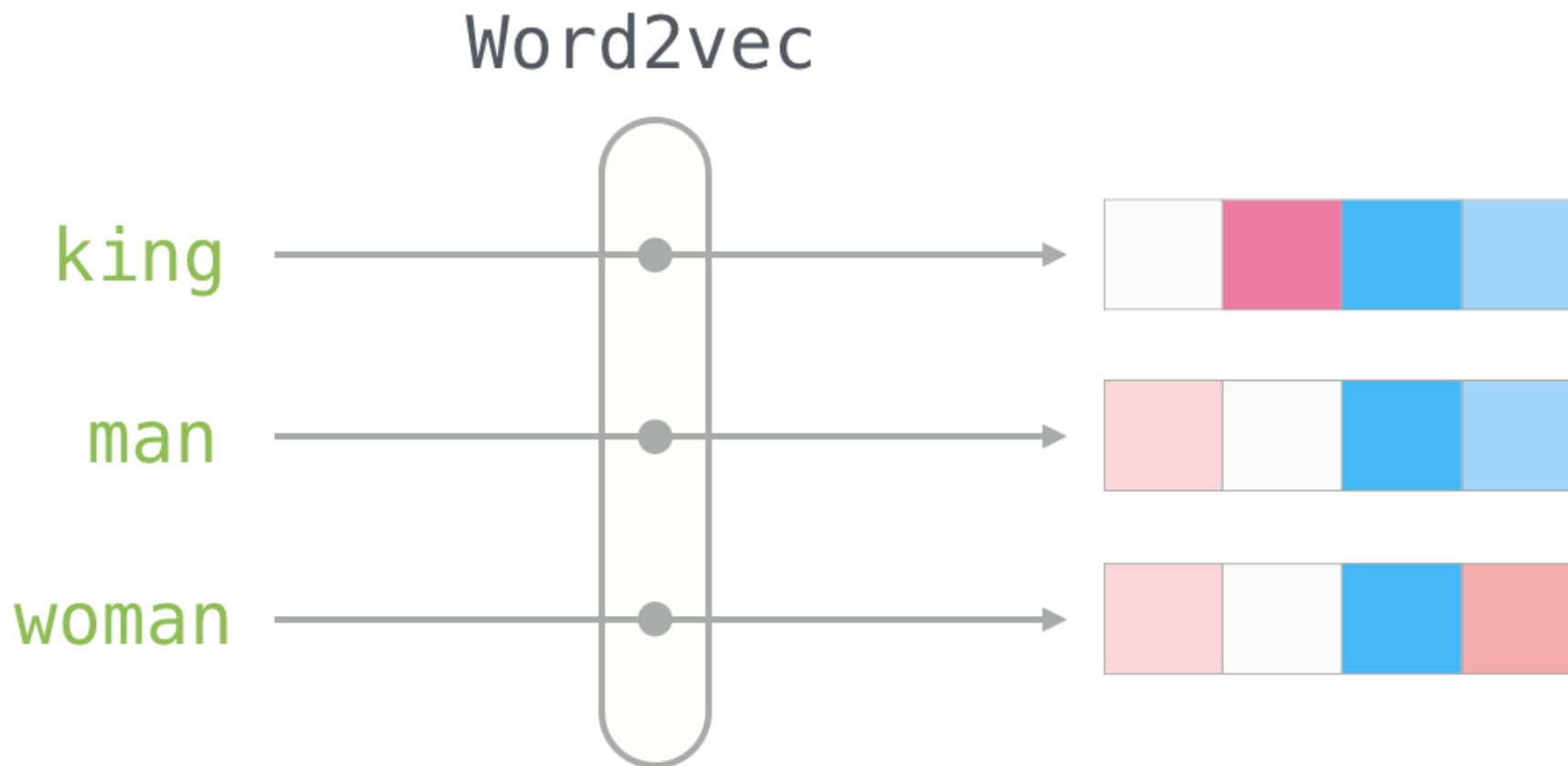
word2vec

04



word2vec

word2vec – метод построения информативных векторных представлений слов, представлен в работе 2013 за авторством Thomas Mikolov и его коллег

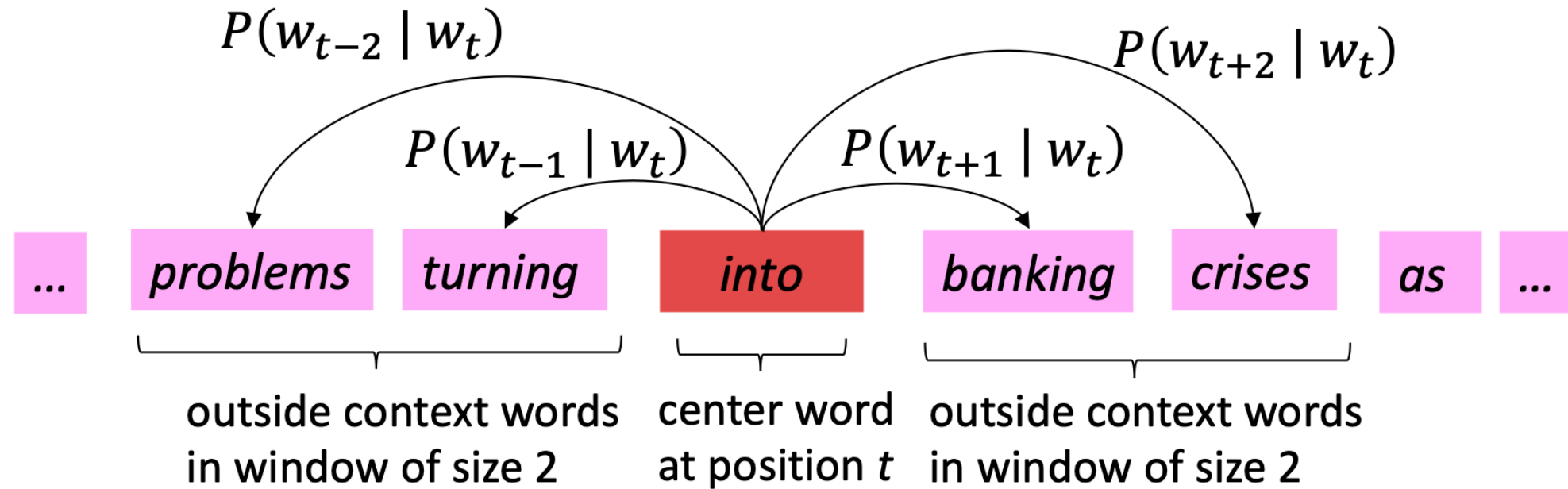


Source Text

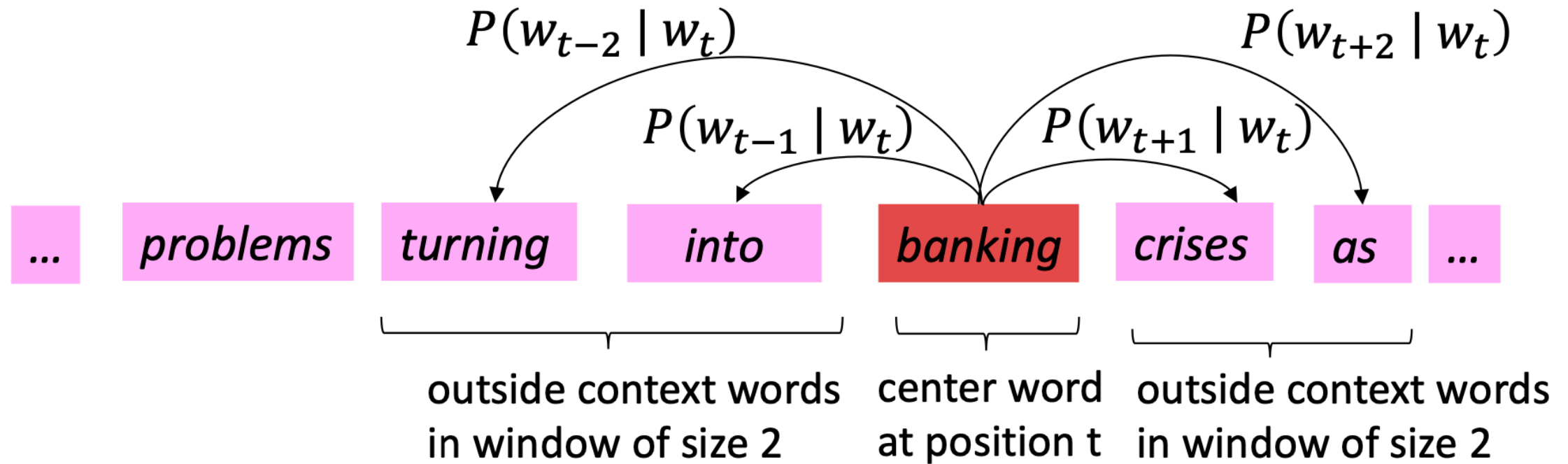
Training Samples

The quick brown fox jumps over the lazy dog. ➡	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. ➡	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. ➡	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. ➡	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

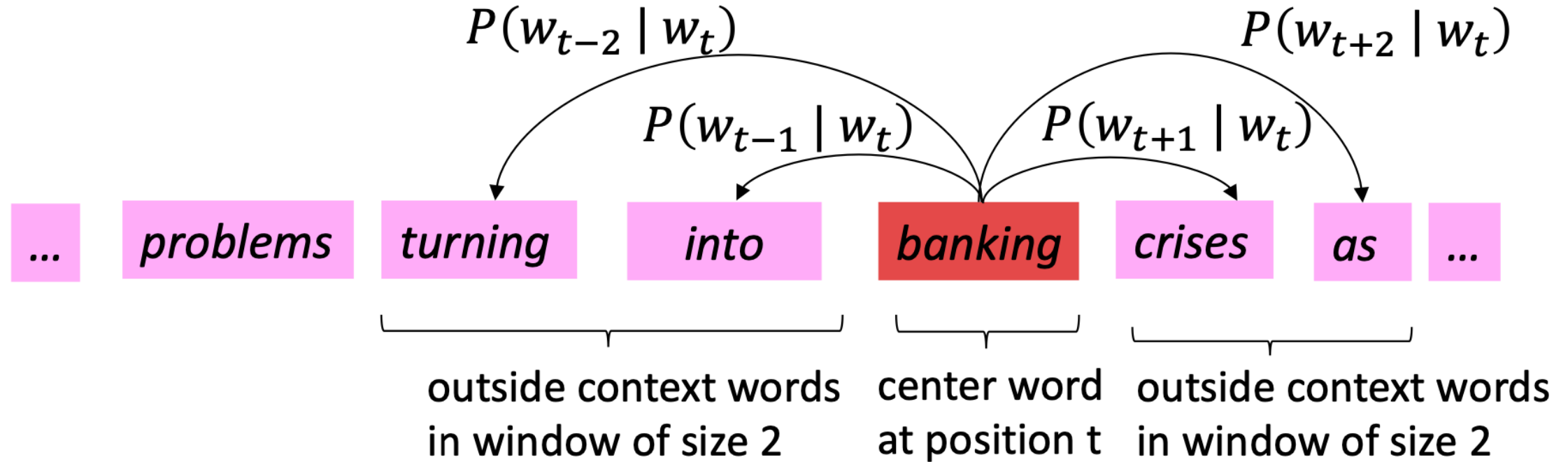
word2vec



word2vec



word2vec



Максимизируемый функционал:
(логарифм правдоподобия)

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$

word2vec

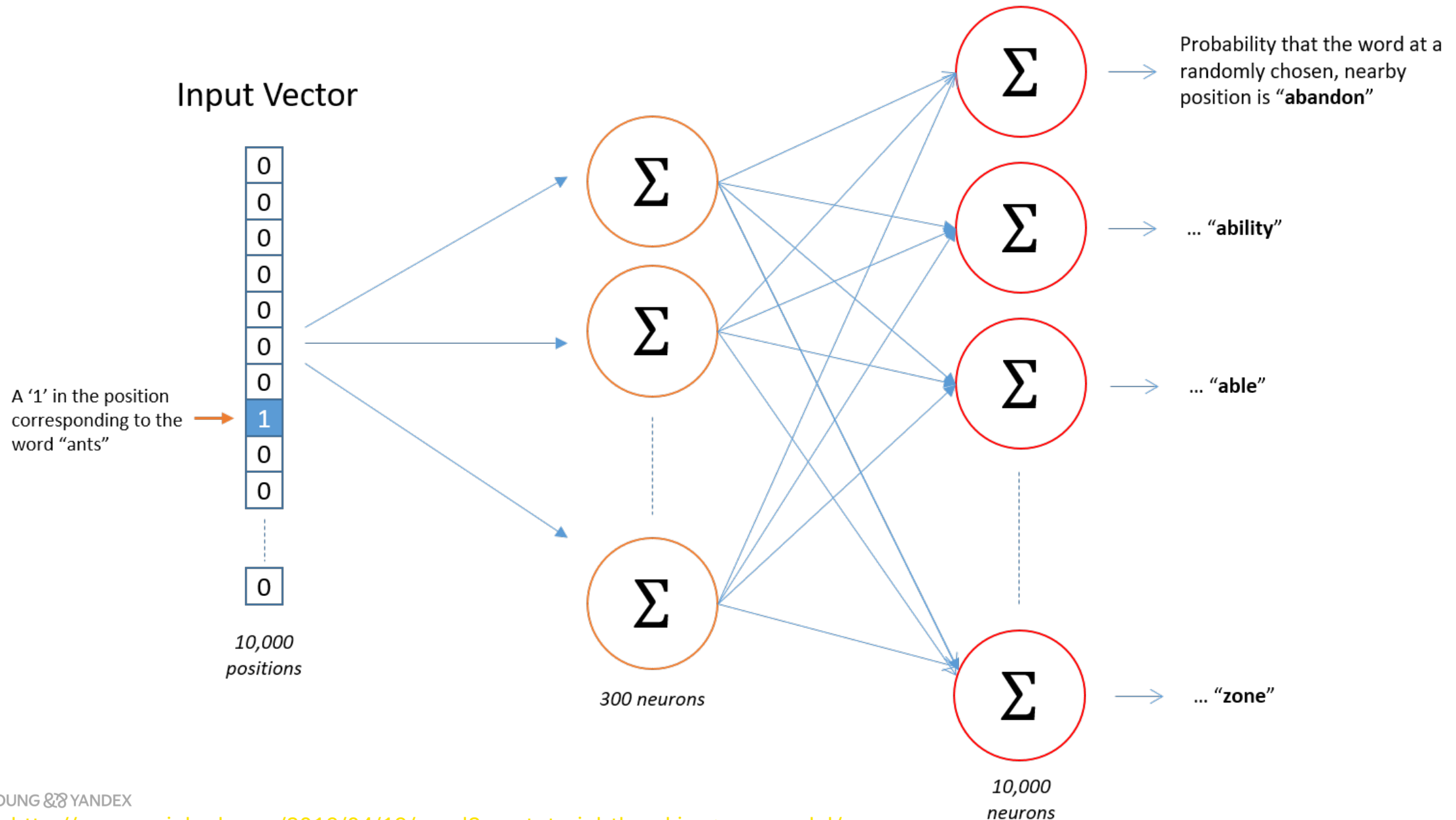
Максимизируемый функционал:
(логарифм правдоподобия)

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$

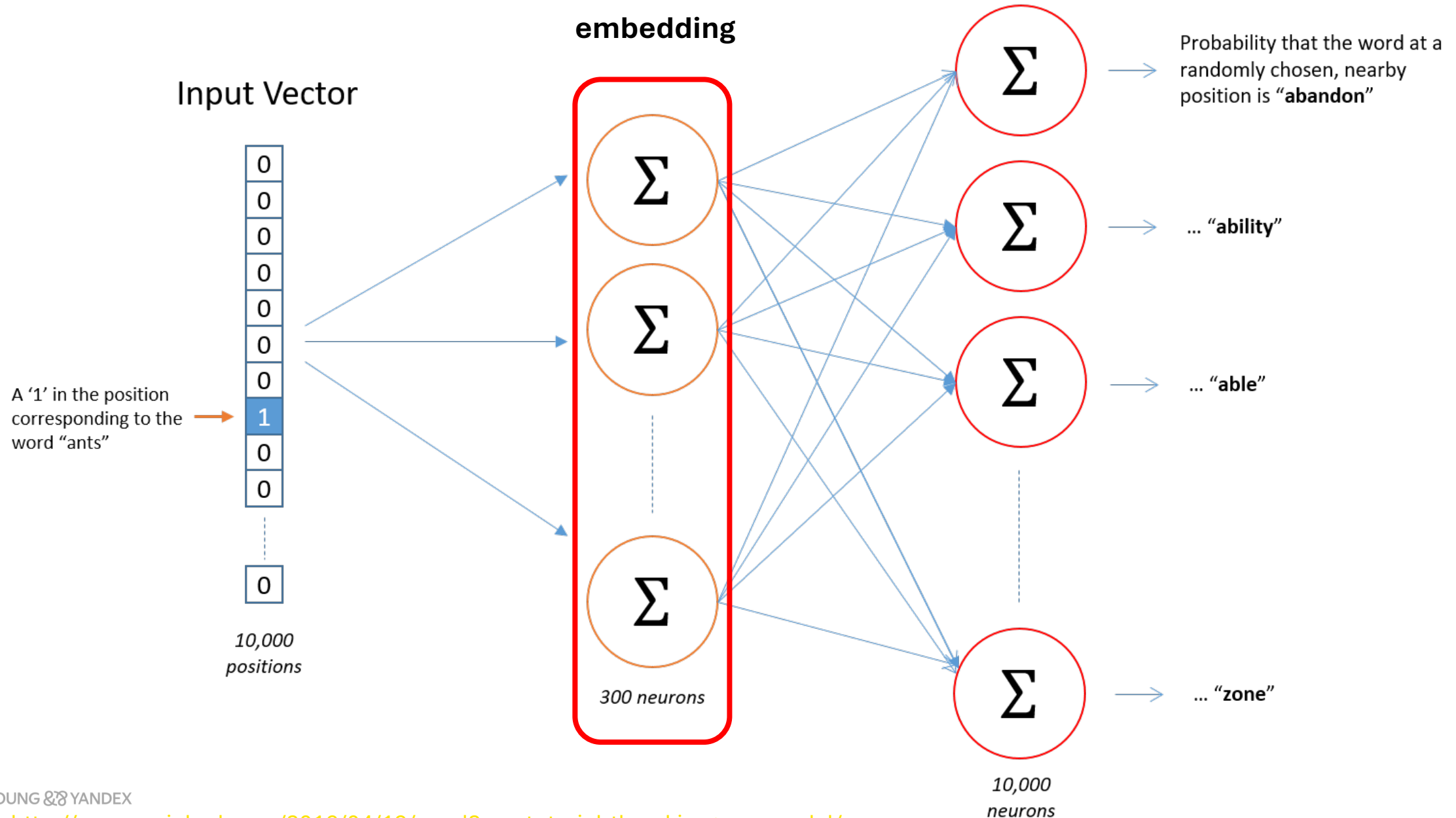
Оценка вероятности через softmax:
(очевидный способ)

$$p(w_O | w_I) = \frac{\exp \left(v'_{w_O}{}^\top v_{w_I} \right)}{\sum_{w=1}^W \exp \left(v'_w{}^\top v_{w_I} \right)}$$

word2vec



word2vec



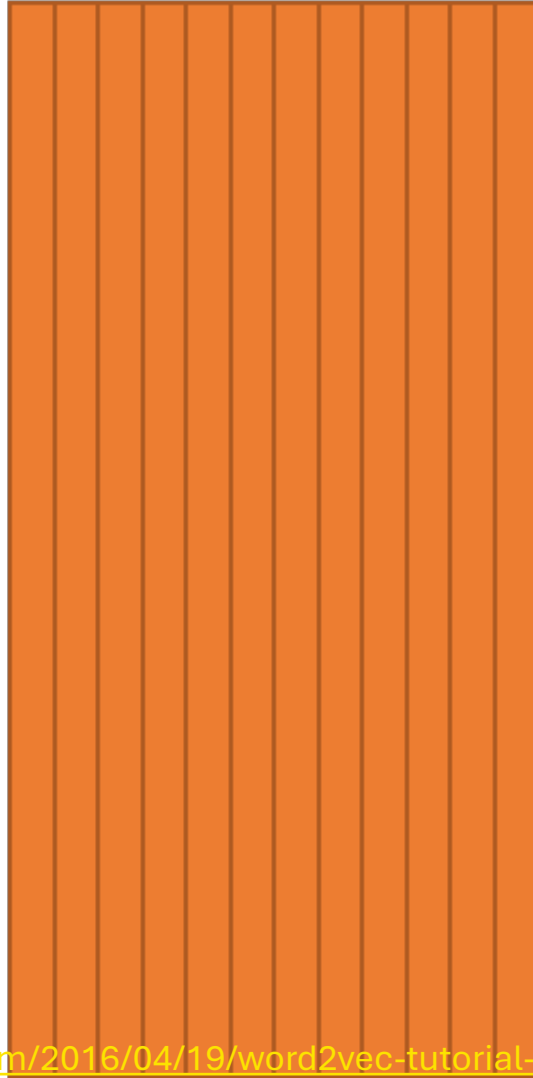
Hidden Layer
Weight Matrix



*Word Vector
Lookup Table!*

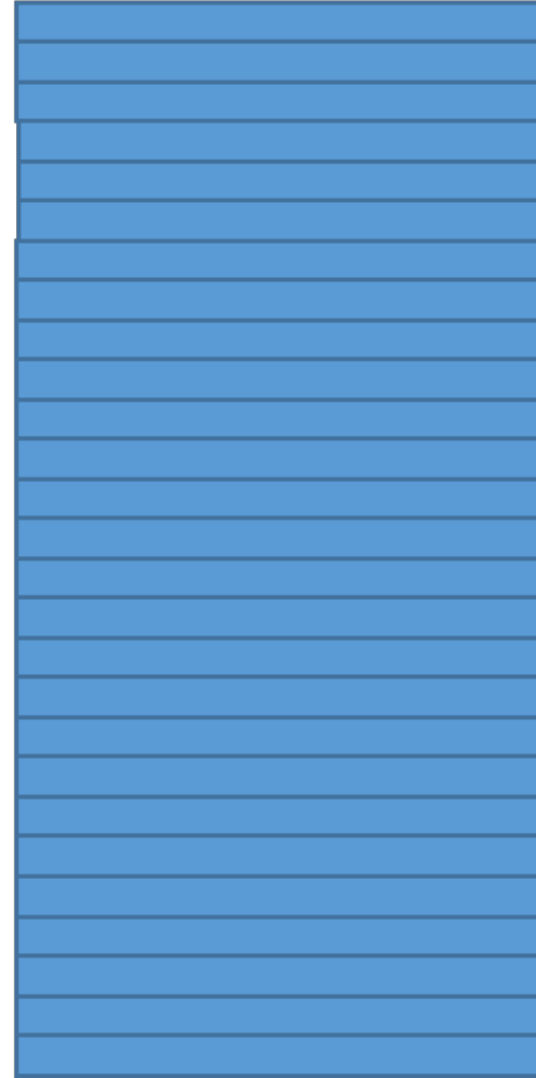
300 neurons

10,000 words



300 features

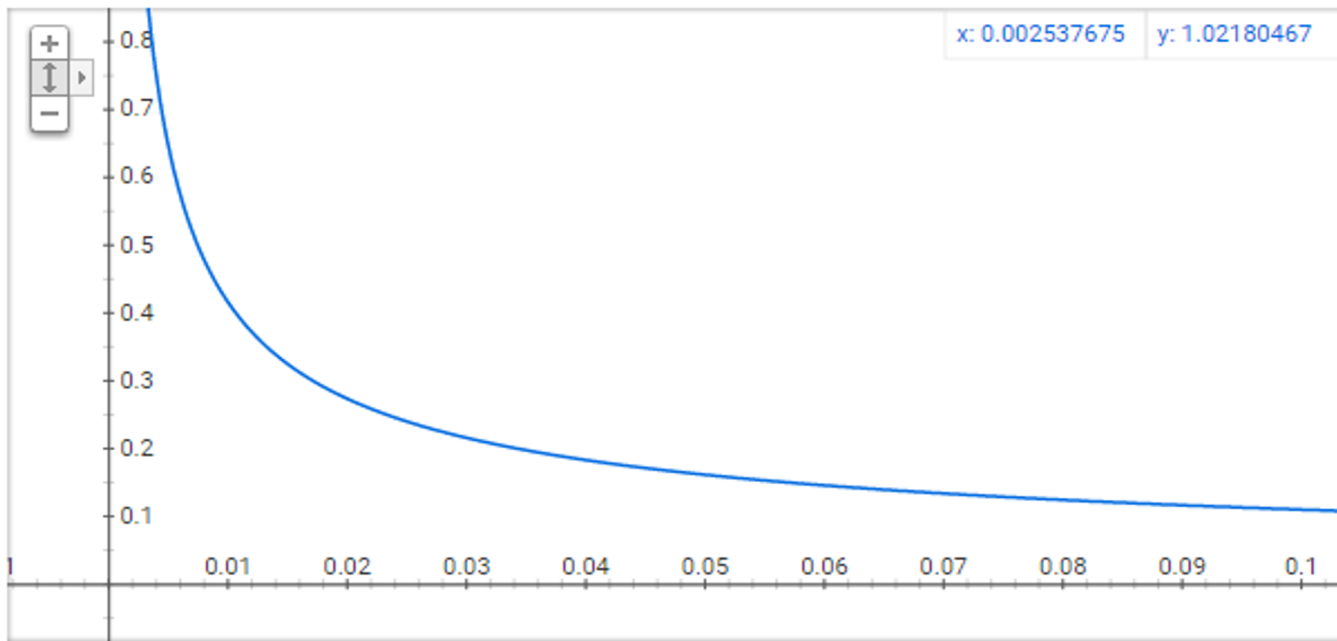
10,000 words



word2vec: subsampling

Часто встречающиеся слова чаще выступают в роли объектов – будем пропускать их случайным образом в зависимости от частоты слова $z(w_i)$

Graph for $(\sqrt{x/0.001}+1)*0.001/x$



$P(w_i)$ – вероятность
использовать пару в обучении

$$P(w_i) = \left(\sqrt{\frac{z(w_i)}{0.001}} + 1 \right) \cdot \frac{0.001}{z(w_i)}$$

word2vec: negative sampling

Имеет смысл не только “сближать” похожие (близкие по контексту) слова, но и “отдалять” непохожие. Для этого воспользуемся механизмом negative sampling.

Чем чаще встречается слово в обучающем корпусе, тем больше вероятность использовать его в качестве negative sample.

$$P(w_i) = \frac{f(w_i)}{\sum_{j=0}^n (f(w_j))}$$

word2vec: negative sampling

Имеет смысл не только “сближать” похожие (близкие по контексту) слова, но и “отдалять” непохожие. Для этого воспользуемся механизмом negative sampling.

Чем чаще встречается слово в обучающем корпусе, тем больше вероятность использовать его в качестве negative sample.

$$P(w_i) = \frac{f(w_i)}{\sum_{j=0}^n (f(w_j))} \longrightarrow P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=0}^n (f(w_j)^{3/4})}$$

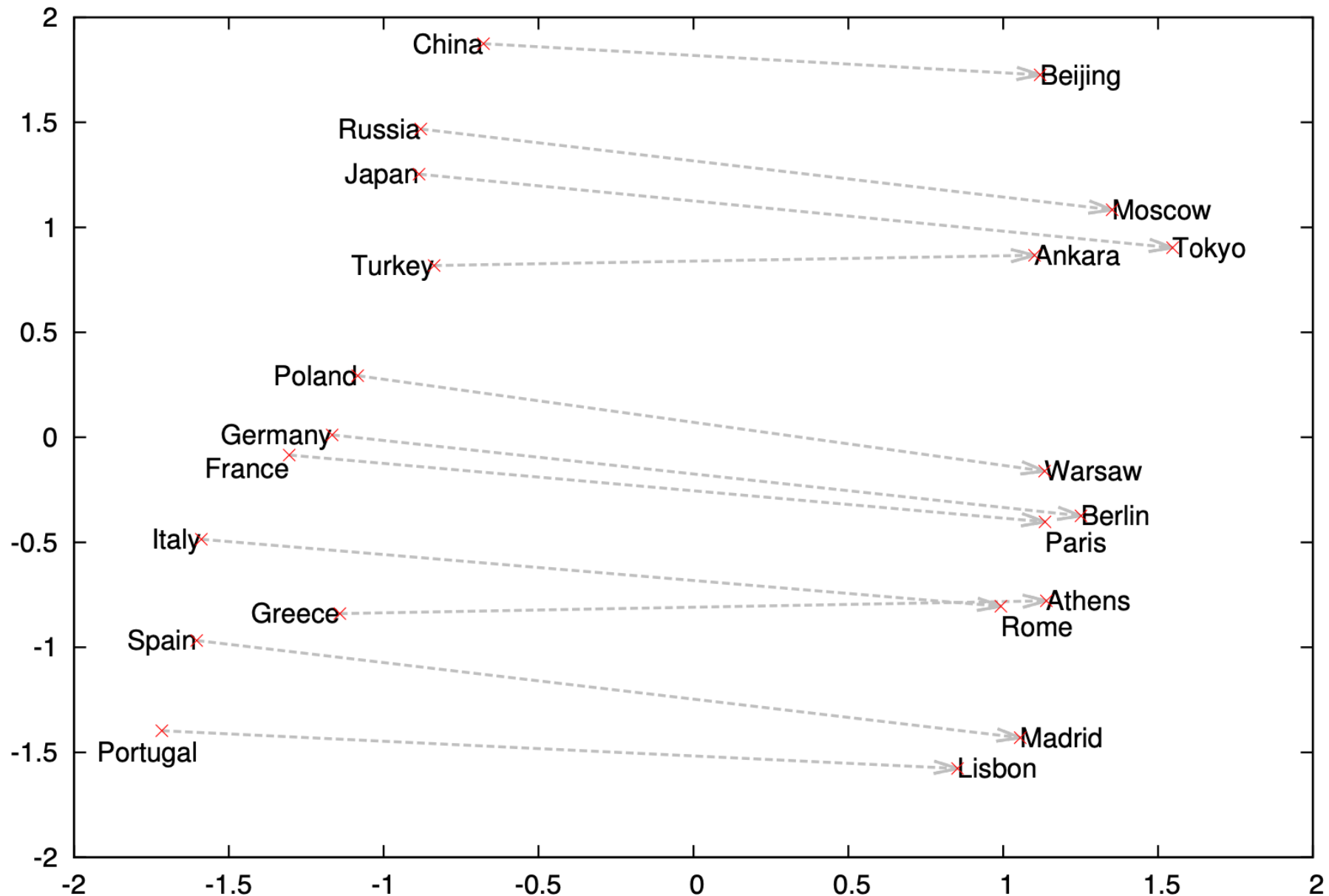
word2vec

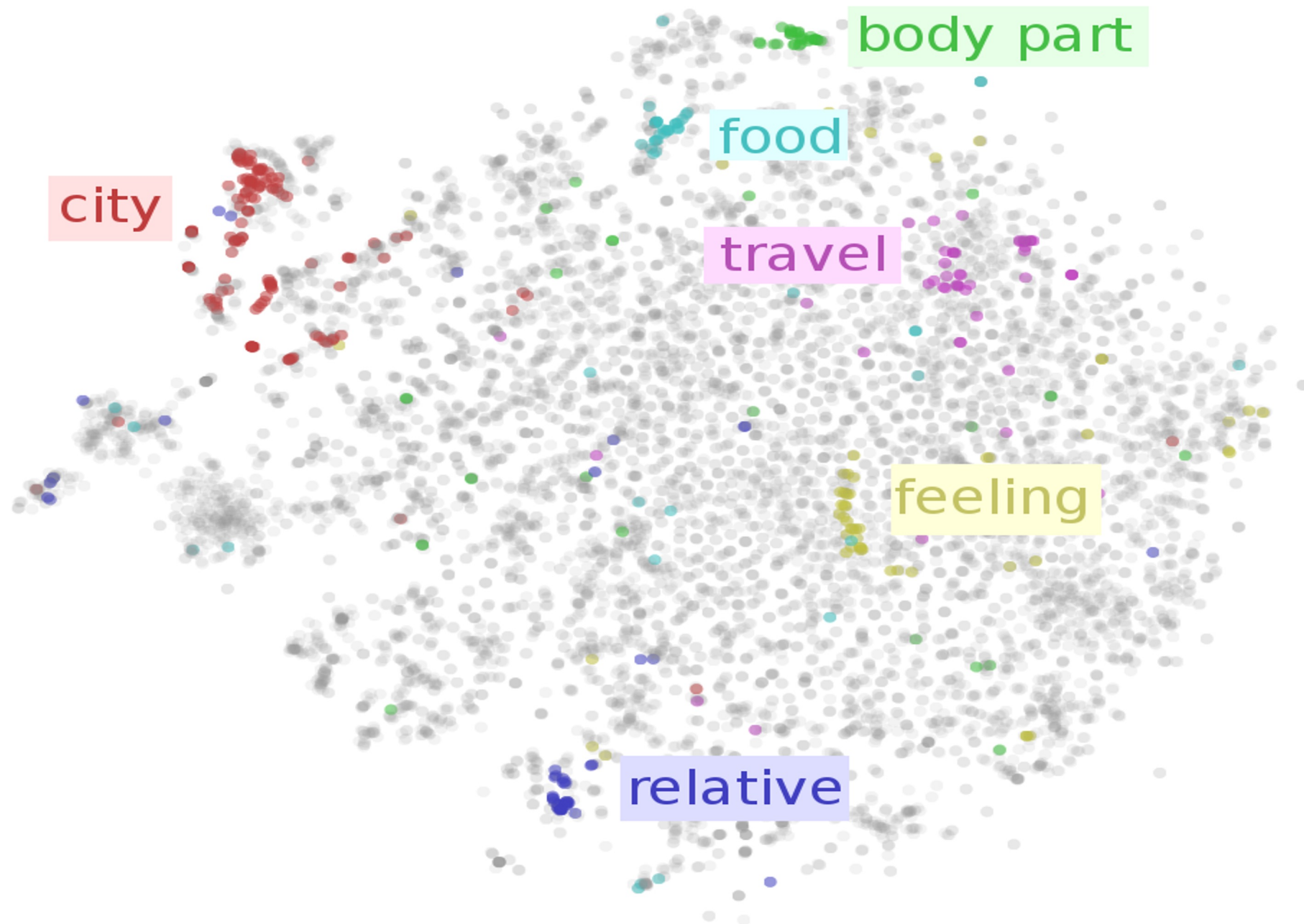
Обновленный оптимизируемый функционал: рассматриваем лишь **положительный** пример и **несколько отрицательных**:

$$\log \sigma(v'_{w_O}{}^\top v_{w_I}) + \sum_{i=1}^k \mathbb{E}_{w_i \sim P_n(w)} \left[\log \sigma(-v'_{w_i}{}^\top v_{w_I}) \right]$$

Перед нами один из многих примеров контрастного обучения – contrastive learning. В дальнейшем мы еще не раз столкнемся с ним.

Country and Capital Vectors Projected by PCA





Спасибо за внимание



Y&OY