

NLP

Word embeddings

ОсноваНО на [NLP Course | For You](#)

План

- Особенности домена
 - Решаемые задачи
 - Предобработка: токенизация и эмбединги
-
- One-hot
 - Count based
 - Word2Vec, FastText, GloVe

Особенности текстового домена

- Слабая структурируемость
- На входе получаем не числа, а последовательность символов
- Длина текстов бывает разной
- Тексты сравнительно часто бывают с опечатками
- Сильная зависимость от контекста (например, кореференции)
- Грамматически и/или семантически связанные слова не всегда расположены рядом
- Неразмеченных данных очень много, даже с древних времён

Задачи решаемые в NLP

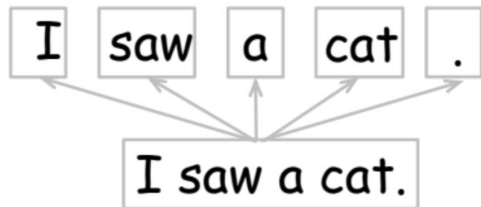
- Классификация текста
- Классификация слов
- Языковое моделирование (aka LM, language modeling)
- text-to-text
- Заполнение маскированных частей текста
- Схожесть текстов
- Выделение границ предложения, расстановка пунктуации

Токенизация

I saw a cat.

Text (your input)

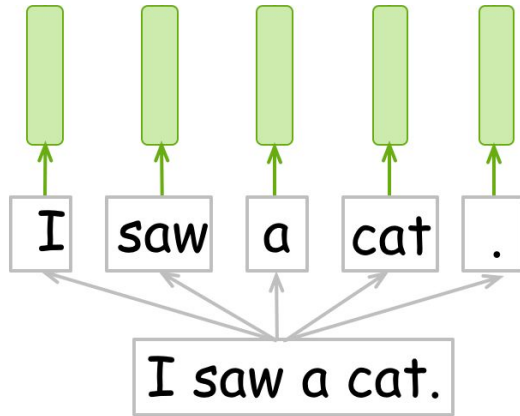
Токенизация



Sequence of tokens

Text (your input)

Embeddings

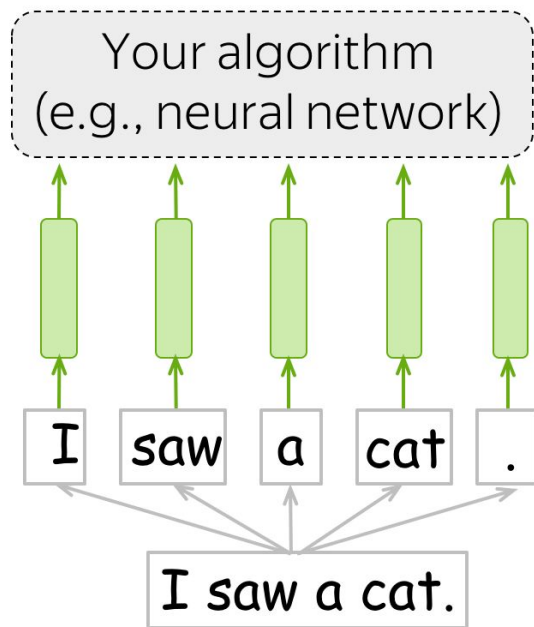


Word representation - vector
(input for your model/algorithm)

Sequence of tokens

Text (your input)

Обработка последовательности эмбедингов



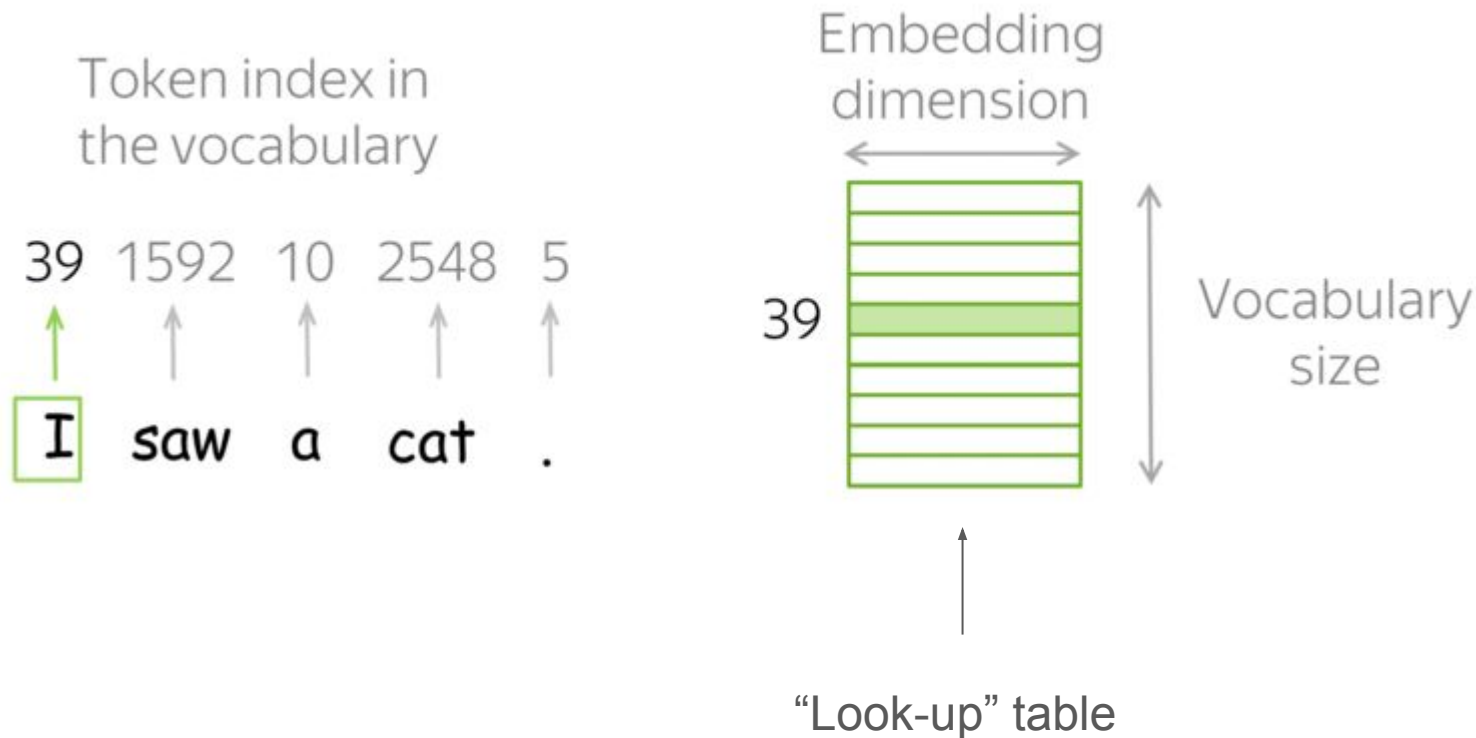
Any algorithm for solving a task

Word representation - vector
(input for your model/algorithm)

Sequence of tokens

Text (your input)

Как работаем с эмбедингами?



Что делать с незнакомыми словами?

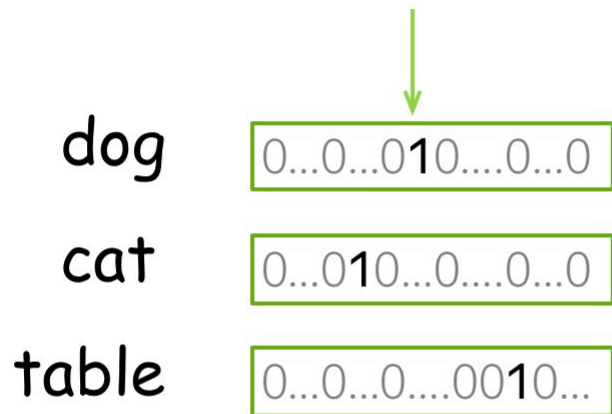
I saw a UNK .
↑ ↑ ↑ ↑ ↑
I saw a &%! . not in the
vocabulary

The diagram shows two rows of text. The top row is 'I saw a UNK .' and the bottom row is 'I saw a &%! .'. Vertical grey arrows point from each word in the bottom row to the corresponding word in the top row. The word '&%!' in the bottom row is enclosed in a red rectangular box. A red arrow points from the text 'not in the vocabulary' to this box.

Embeddings

One-hot encoding

One is 1, the rest are 0



В чём проблема?

Embedding dimension =
vocabulary size

One-hot encoding

One is 1, the rest are 0



dog

0...0...010...0...0

cat

0...010...0...0...0

table

0...0...0...0010...

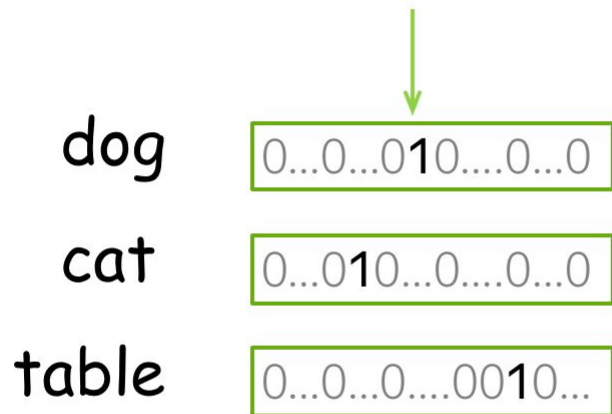


Embedding dimension =
vocabulary size

- размер эмбединга очень большой
- не отражают смысл слов

One-hot encoding

One is 1, the rest are 0



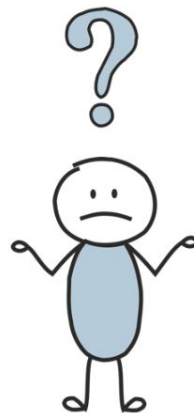
Embedding dimension =
vocabulary size

- размер эмбединга очень большой
- не отражают **СМЫСЛ** слов

Смысл

Do you know what the word **tezgüino** means ?

(We hope you do not)



СМЫСЛ

Now look how this word is used in different contexts:

A bottle of **tezgüino** is on the table.

Everyone likes **tezgüino**.

Tezgüino makes you drunk.

We make **tezgüino** out of corn.

Can you understand what **tezgüino** means ?



СМЫСЛ

Now look how this word is used in different contexts:

A bottle of **tezgüino** is on the table.

Everyone likes **tezgüino**.

Tezgüino makes you drunk.

We make **tezgüino** out of corn.



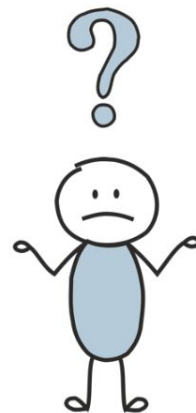
Tezgüino is a kind of alcoholic beverage made from corn.

With context, you can understand the meaning!



Смысл

How did you do this?



СМЫСЛ

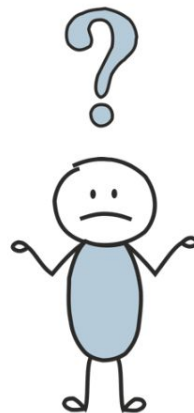
(1) A bottle of _____ is on the table.

(2) Everyone likes _____ .

(3) _____ makes you drunk.

(4) We make _____ out of corn.

What other words fit
into these contexts ?



СМЫСЛ

(1) A bottle of _____ is on the table.

(2) Everyone likes _____ .

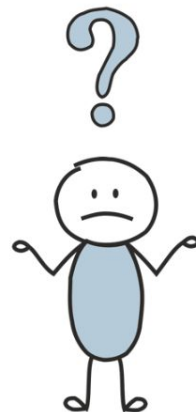
(3) _____ makes you drunk.

(4) We make _____ out of corn.

What other words fit
into these contexts ?

| | (1) | (2) | (3) | (4) | ... | ← contexts |
|-----------|-----|-----|-----|-----|-----|------------|
| tezgüino | 1 | 1 | 1 | 1 | | |
| loud | 0 | 0 | 0 | 0 | | |
| motor oil | 1 | 0 | 0 | 1 | | |
| tortillas | 0 | 1 | 0 | 1 | | |
| wine | 1 | 1 | 1 | 0 | | |

← rows show contextual
properties: 1 if a word can
appear in the context, 0 if not



СМЫСЛ

(1) A bottle of _____ is on the table.

(2) Everyone likes _____ .

(3) _____ makes you drunk.

(4) We make _____ out of corn.

| | (1) | (2) | (3) | (4) | ... |
|-----------|-----|-----|-----|-----|-----|
| tezgüino | 1 | 1 | 1 | 1 | |
| loud | 0 | 0 | 0 | 0 | |
| motor oil | 1 | 0 | 0 | 1 | |
| tortillas | 0 | 1 | 0 | 1 | |
| wine | 1 | 1 | 1 | 0 | |

rows are
similar

СМЫСЛ

(1) A bottle of _____ is on the table.

(2) Everyone likes _____ .

(3) _____ makes you drunk.

(4) We make _____ out of corn.

| | (1) | (2) | (3) | (4) | ... |
|-----------|-----|-----|-----|-----|-----|
| tezgüino | 1 | 1 | 1 | 1 | |
| loud | 0 | 0 | 0 | 0 | |
| motor oil | 1 | 0 | 0 | 1 | |
| tortillas | 0 | 1 | 0 | 1 | |
| wine | 1 | 1 | 1 | 0 | |

rows are
similar



Is this true?

meanings of the
words are similar

СМЫСЛ

(1) A bottle of _____ is on the table.

(2) Everyone likes _____ .

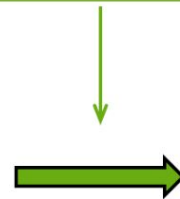
(3) _____ makes you drunk.

(4) We make _____ out of corn.

| | (1) | (2) | (3) | (4) | ... |
|-----------|-----|-----|-----|-----|-----|
| tezgüino | 1 | 1 | 1 | 1 | |
| loud | 0 | 0 | 0 | 0 | |
| motor oil | 1 | 0 | 0 | 1 | |
| tortillas | 0 | 1 | 0 | 1 | |
| wine | 1 | 1 | 1 | 0 | |

This is the **distributional hypothesis**

rows are
similar



meanings of the
words are similar

Дистрибутивная гипотеза

Words which frequently appear in **similar contexts** have **similar meaning**.

(Harris 1954, Firth 1957)

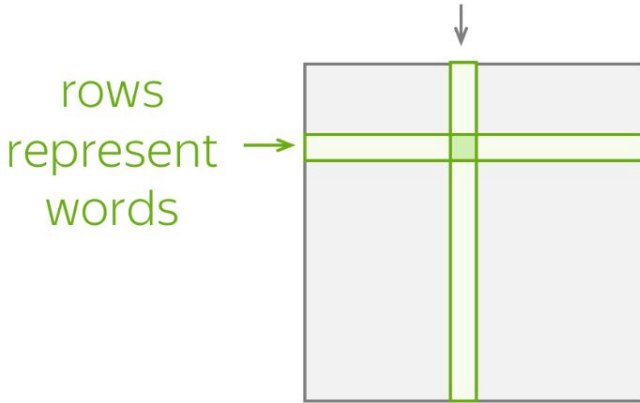
Main idea:

We have to put information about contexts into word vectors.

What comes next: different ways to do this

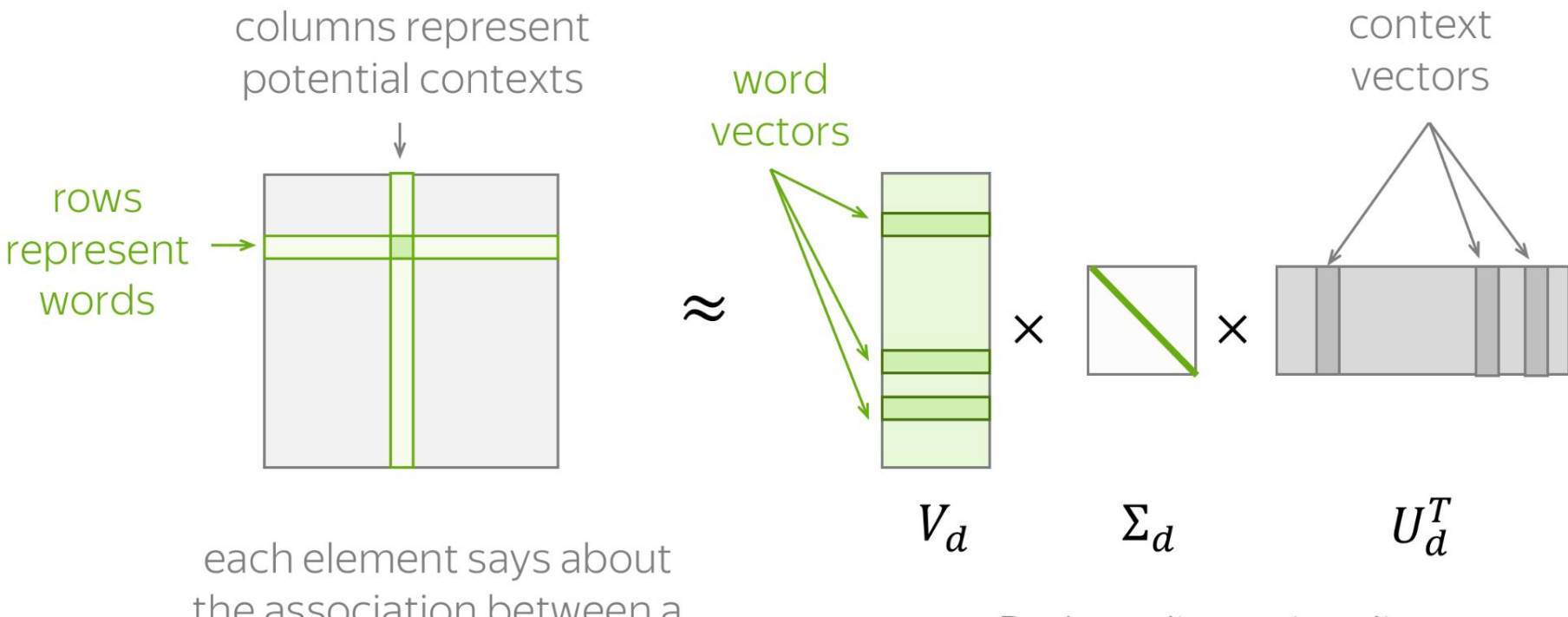
Count based

columns represent
potential contexts



each element says about
the association between a
word and a **context**

Count based




TF-IDF

Context:


- document d (from a collection D)

Matrix element:

- $\text{tf-idf}(w, d, D) = \text{tf}(w, d) \cdot \text{idf}(w, D)$


$$N(w, d)$$

term frequency


$$\log \frac{|D|}{|\{d \in D: w \in d\}|}$$

inverse document frequency

Word2Vec

Word2Vec

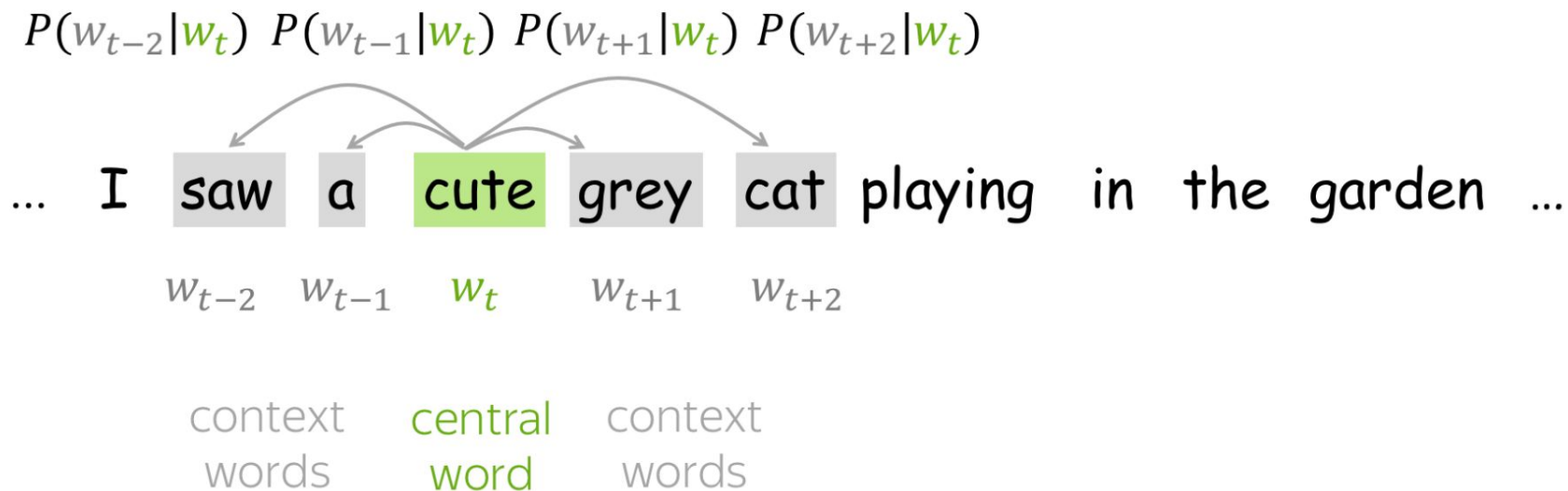
$$P(w_{t-2}|w_t) \quad P(w_{t-1}|w_t) \quad P(w_{t+1}|w_t) \quad P(w_{t+2}|w_t)$$

... I saw a cute grey cat playing in the garden ...

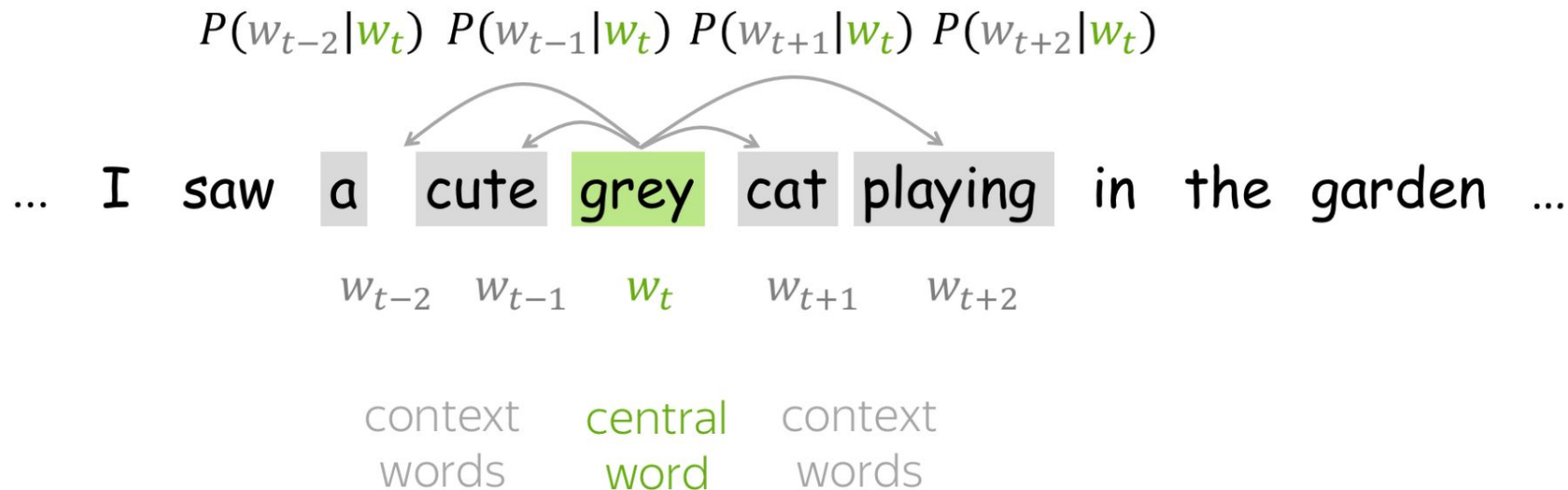
w_{t-2} w_{t-1} w_t w_{t+1} w_{t+2}

context central context
words word words

Word2Vec



Word2Vec



Objective

Word2Vec tries to find the parameters that maximize the data likelihood:

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m, \\ j \neq 0}} P(w_{t+j} | w_t, \theta)$$

We want our model to think that the training data is “likely”

To do this, it uses negative (log-)likelihood as its loss function:

$$\text{Loss} = J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \log P(w_{t+j} | w_t, \theta)$$

agrees with our plan above



go over text



with a sliding window



compute probability of the context word given the central

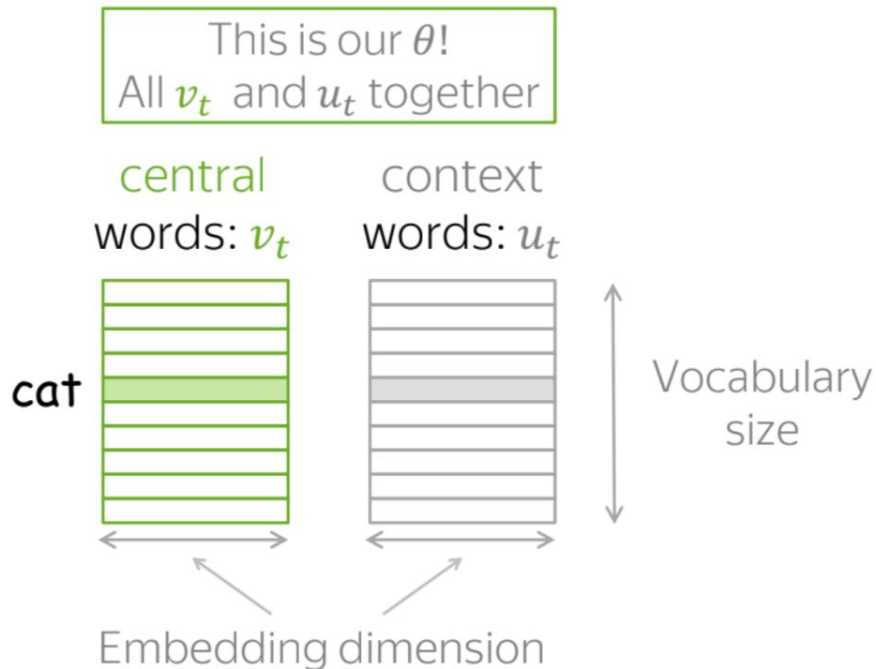


How to compute $P(w_{t+j} | w_t, \theta)$?

For each word w , we will have two vectors:

- v_w when it is a central word
- u_w when it is a context word

Once the vectors are trained,
usually we throw away context
vectors and use only word vectors.



How to compute $P(w_{t+j} | w_t, \theta)$?

For the central word c and context word o (o - outside):

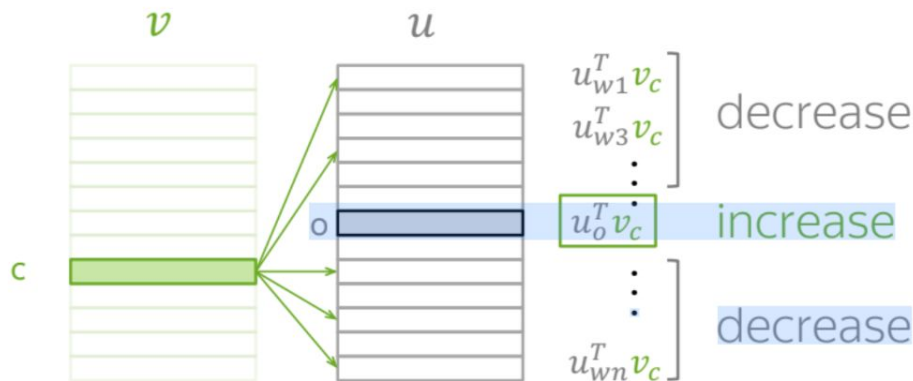
$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

Dot product: measures similarity of o and c
Larger dot product = larger probability

Normalize over entire vocabulary
to get probability distribution

Let us recall our plan:

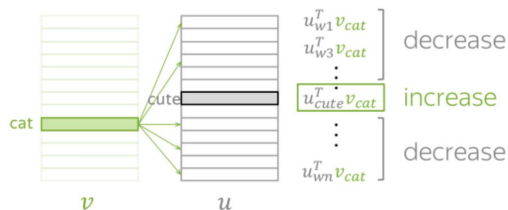
- ...
- adjust the vectors to increase these probabilities.



Negative sampling

Dot product of v_{cat} :

- with u_{cute} - increase,
- with all other u - decrease



Parameters to be updated:

- v_{cat}
 - u_w for all w in the vocabulary
- $|V| + 1$ vectors

Many parameters
at each step –
slow training

Negative sampling

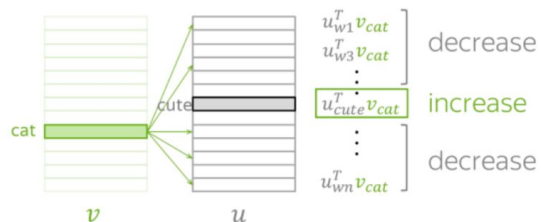
Dot product of v_{cat} :

- with u_{cute} - increase,
- with all other u - decrease



Dot product of v_{cat} :

- with u_{cute} - increase,
- with a subset of other u - decrease

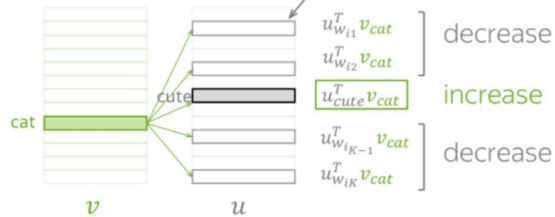


Parameters to be updated:

bad

- v_{cat}
 - u_w for all w in the vocabulary
- $|V| + 1$ vectors

Negative samples: randomly selected K words



Parameters to be updated:

good

- v_{cat}
 - u_{cute} and u_w for w in K negative examples
- $K + 2$ vectors

Window size

- **Larger windows** – more topical similarities

dog
bark leash
(grouped together)

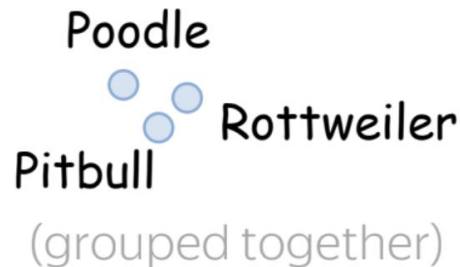


walking
run walked
(grouped together)



- **Smaller windows** – more functional and syntactic similarities

Poodle
Pitbull Rottweiler
(grouped together)

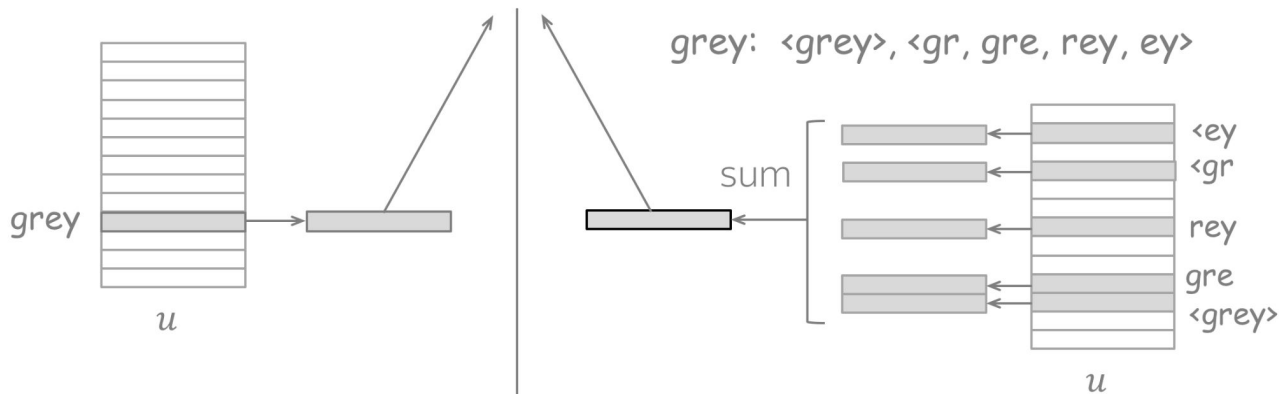


walking
approaching running
(grouped together)



FastText

... I saw a cute grey cat playing in the garden ...



Word2Vec

Vocabulary consists of:

- words

Word vector is:

- one vector from the look-up table

FastText

Vocabulary consists of:

- words and character n-grams

Word vector is:

- sum of word vector and vectors for its n-grams

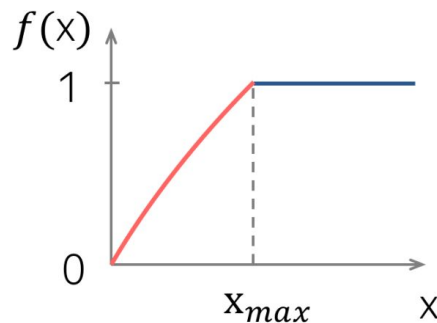
GloVe

context vector word vector bias terms (also learned)

$$J(\theta) = \sum_{w,c \in V} \underbrace{f(N(w, c))}_{\downarrow} \cdot (u_c^T v_w + b_c + \overline{b_w} - \log N(w, c))^2$$

Weighting function to:

- penalize rare events
- not to over-weight frequent events



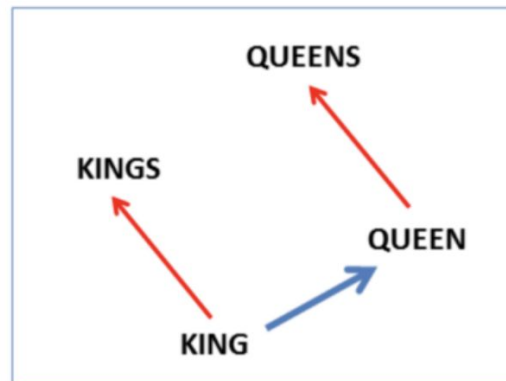
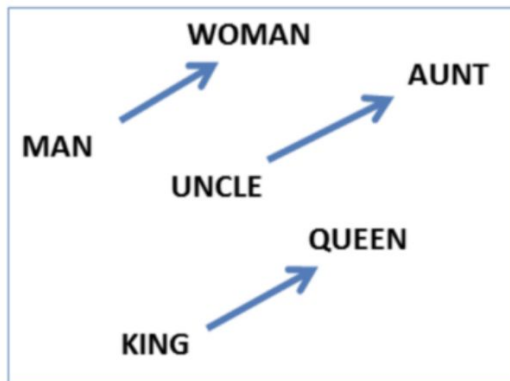
$$\begin{cases} (x/x_{max})^\alpha & \text{if } x < x_{max}, \\ 1 & \text{otherwise.} \end{cases}$$

$$\alpha = 0.75, x_{max} = 100$$

Linear structure

semantic: $v(\text{king}) - v(\text{man}) + v(\text{woman}) \approx v(\text{queen})$

syntactic: $v(\text{kings}) - v(\text{king}) + v(\text{queen}) \approx v(\text{queens})$



Linear structure

