

Работа с текстовыми данными



Как связаны слово или предложения с алгеброй?

- Bag of words (мешок слов)
- Tf-idf
- Эмбединги (например, word2vec)

Raw Text

A dog in heat needs
more than shade

Bag of words vector

Dog	0
need	2
Cat	1
than	0
it	1
heat	2
needs	0

train_X

'This is good',

'This is bad'

'This is awesome'

Fit

CountVectorizer

word_index

{'this':0,
'is':1,
'good':2,
'bad':3,
'awesome':4}

Features

This	is	good	bad	awesome
1	1	1	0	0
1	1	0	1	0
1	1	0	0	1

TF*IDF

TF*IDF=**TERM FREQUENCY** * INVERSE
DOCUMENT FREQUENCY

TERM FREQUENCY=

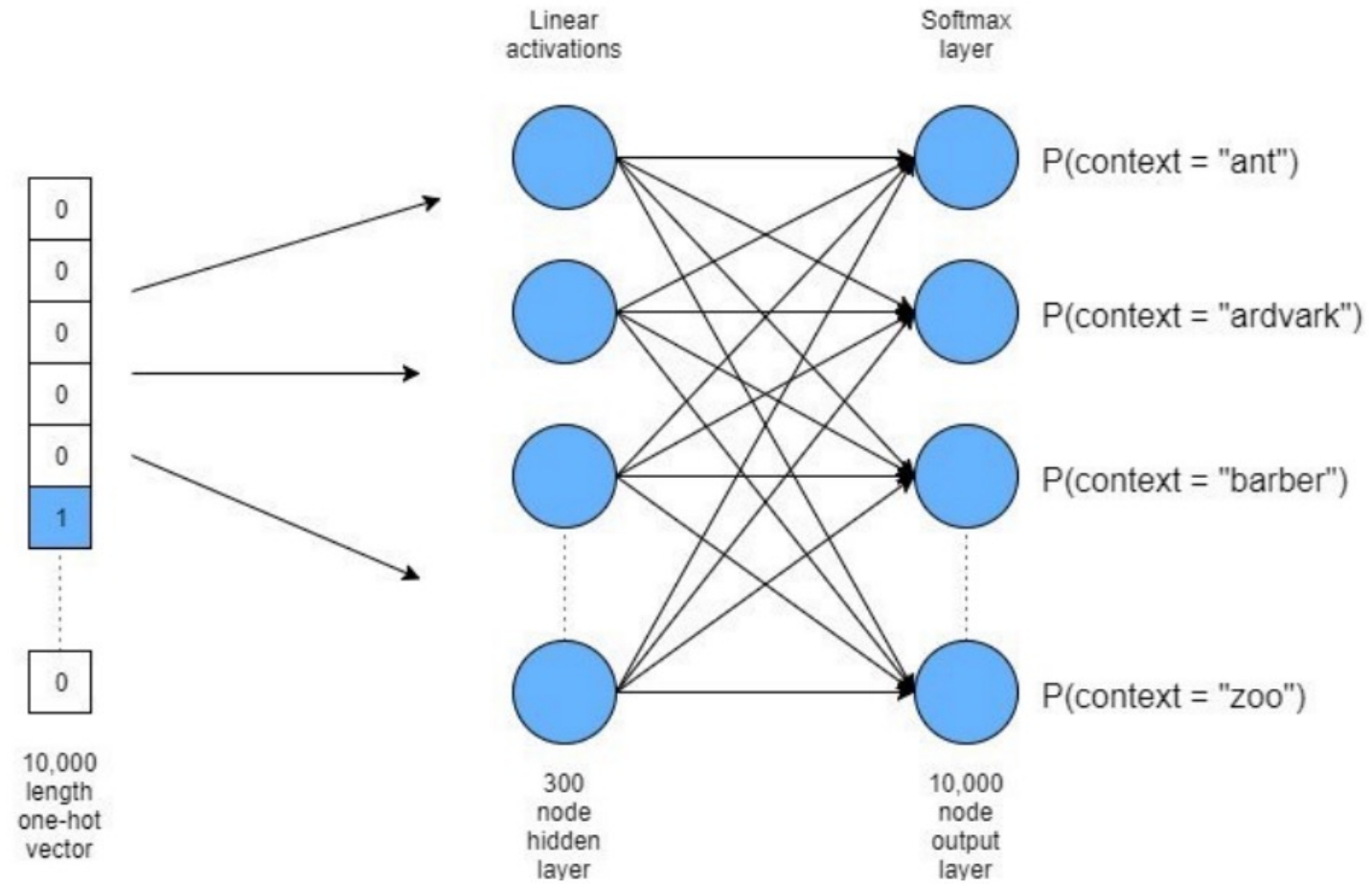
THE AMOUNT OF TIMES A
TERM APPEARS IN A
DOCUMENT

X

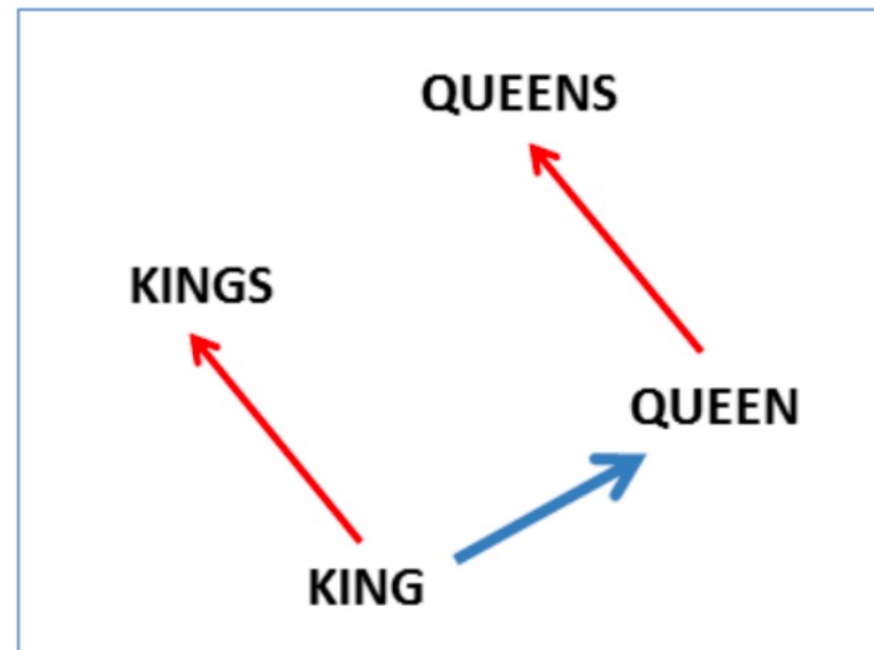
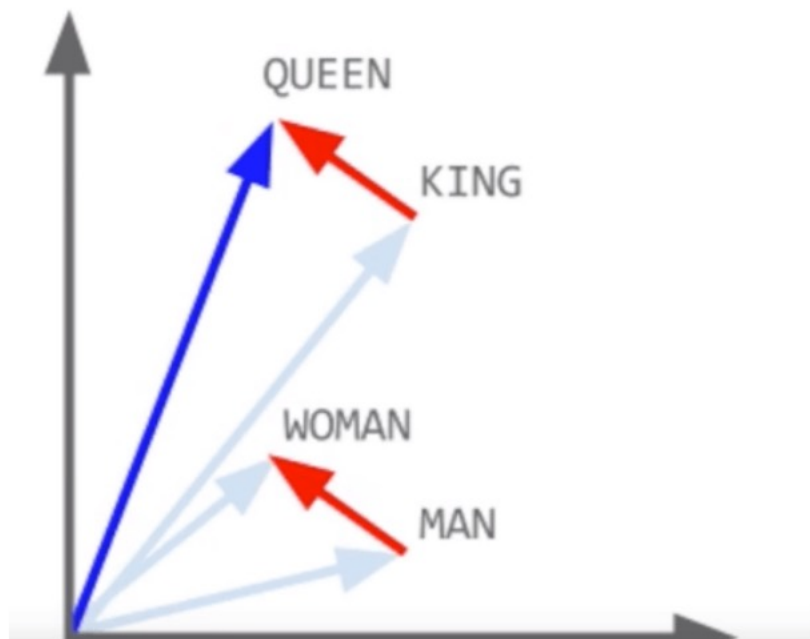
INVERSE DOCUMENT
FREQUENCY=

A MEASURE OF WHETHER A
TERM IS RARE OR COMMON
IN A COLLECTION OF
DOCUMENTS.

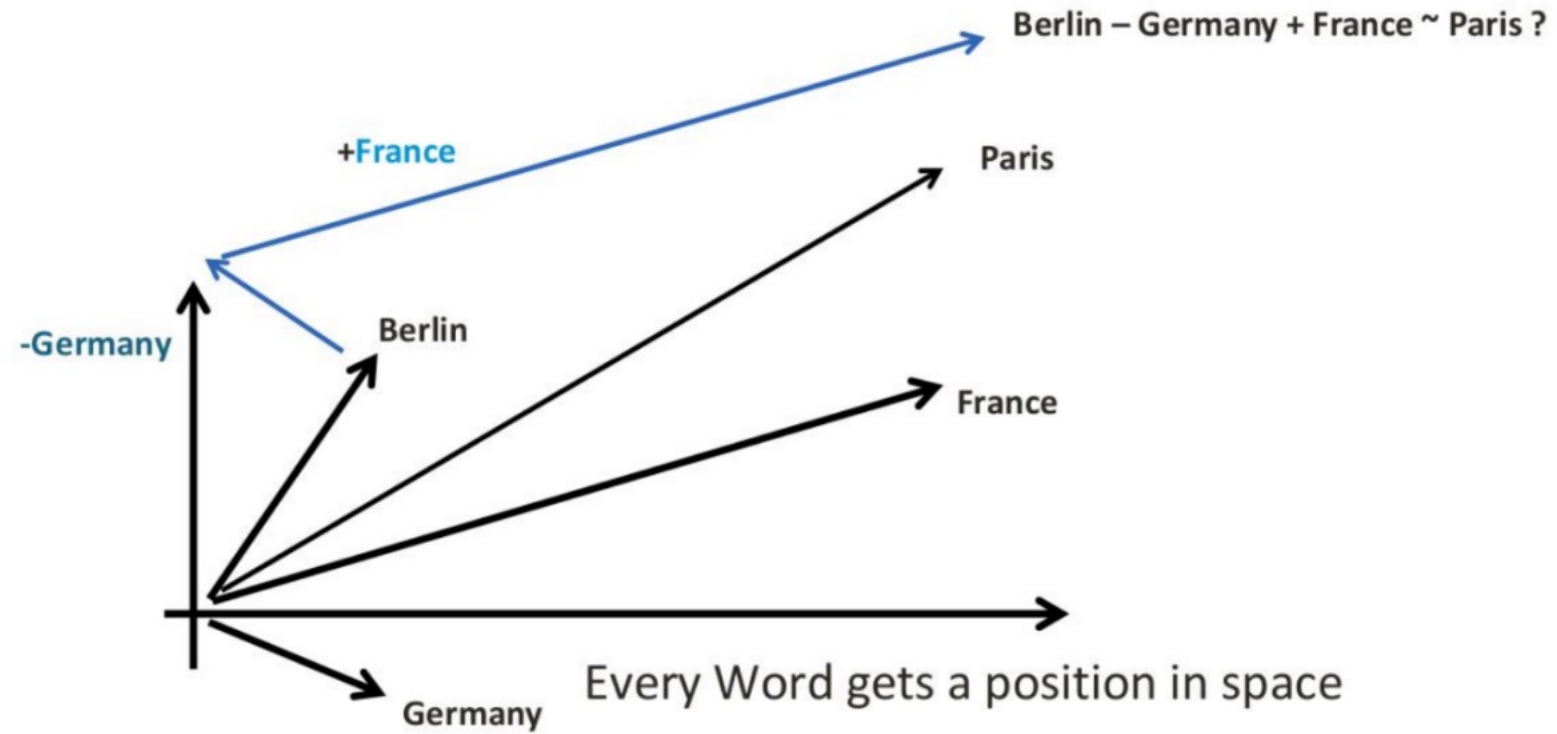
Word2Vec



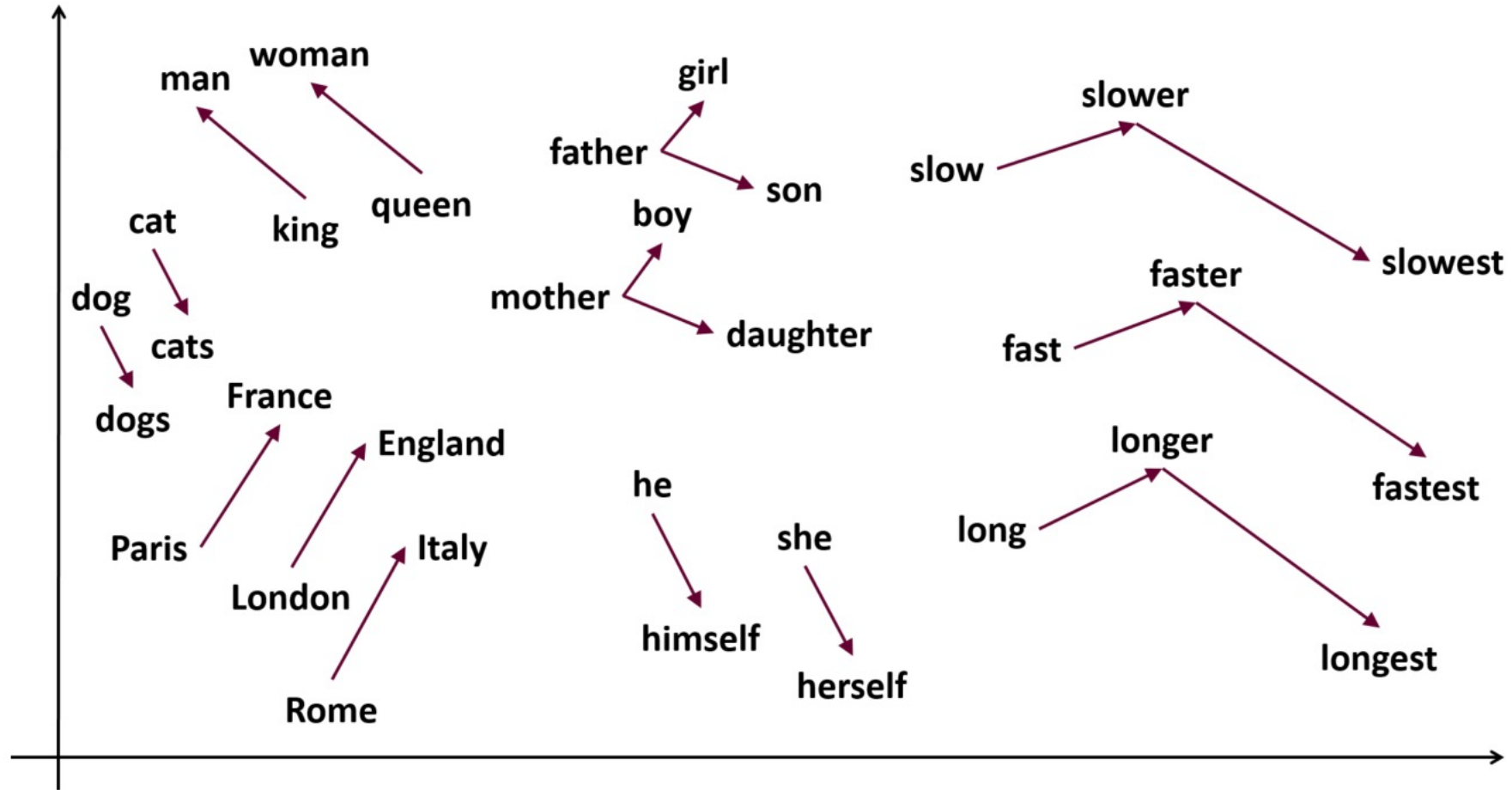
Word2Vec



Word2Vec



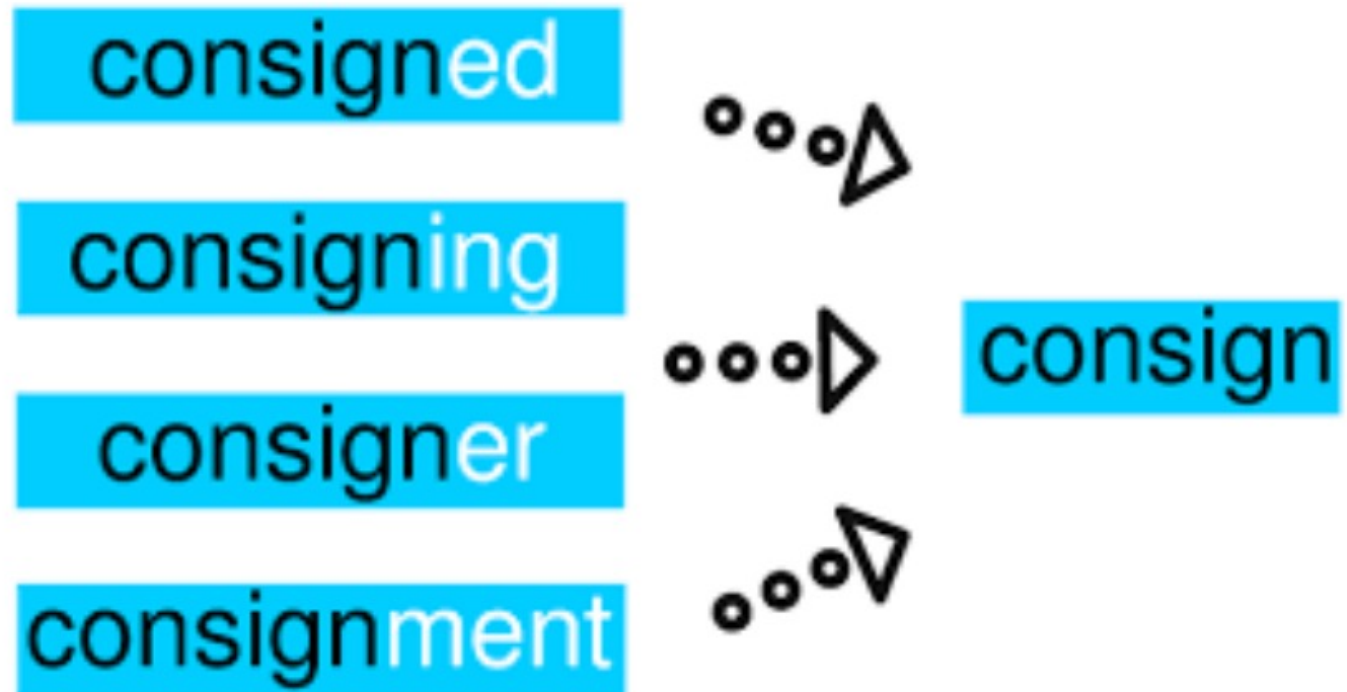
Word2Vec



Стемминг



Стемминг



Лемматизация

Пушистая	кошка	очень	спадко	спит	на	мягком	кожаном	диване	,	поймав	и	съев	маленькую	противную	мышку	с	длинным	хвостиком
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
пушистый	кошка	очень	спадко спадкий	спать	на	мягкий	кожаный	диван	,	поймав	и	съев	маленький	противный	мышка	с	длинный	хвостик

Лемматизация

обезьяны ➡ обезьяна

искал ➡ искать

любезных ➡ любезный

СТОП-слова!

- Союзы
- Предлоги
- Частицы
- Высокочастотные слова
- И т.д. и т.п.

N-gramms

