# Word embeddings

**Radoslav Neychev,** neychev@bigdatateam.org

ML Instructor, BigDataTeam

Assistant Professor, MIPT

LinkedIn profile

**09.12.2020,  online**

# Natural Language Processing:
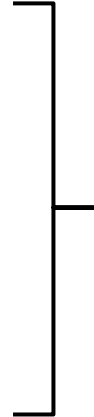
# Introduction

- Sentiment analysis
- Spam filtering
- Fake news detection
- Topic prediction
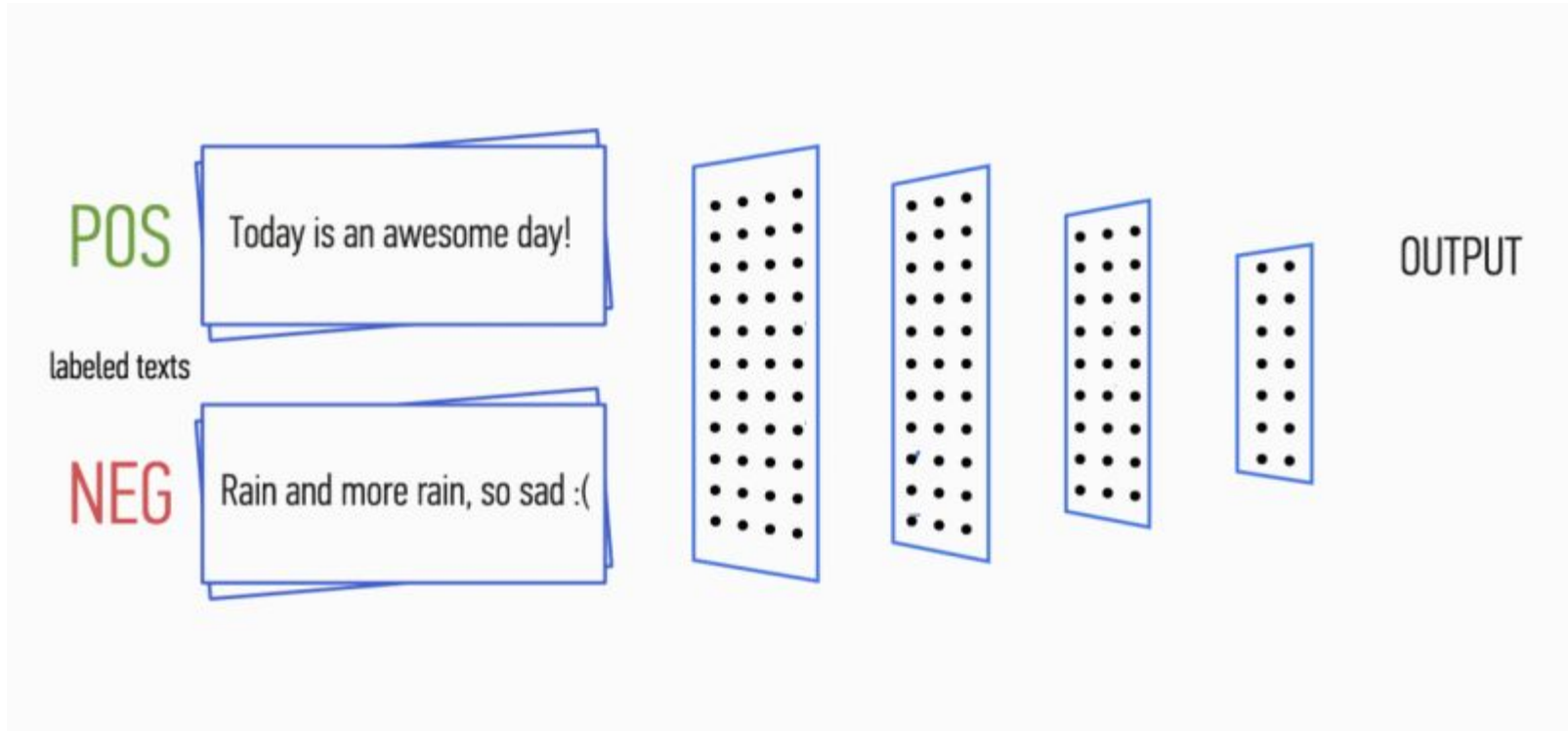- #hashtag prediction

Text classification tasks

Source: habr.com

Discrete labels:

- Binary
    - spam filtering, sentiment analysis
- Multi-class
    - categorization of items by its description
- Multi-label
    - #hashtag prediction

Continuous labels:

- Predict product price by its description

Fixed-size vector of
features

Class label /
Probability distribution

Fixed-size vector of features

$\blacktriangleright$ Linear Classifier

$\blacktriangleright$ Logistic Regression

$\blacktriangleright$ etc.

Class label / Probability distribution

Fixed-size vector of features

Class label / Probability distribution

# Feature extraction

Tokenization: split the input into tokens

| This is a sentence |

→

| This | | is | | a | | sentence |

the dog is on the table

| 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|---|
| are | cat | dog | is | now | on | table | the |

Problems:
▷ No information about words order
▷ Word vectors are huge and very sparse
▷ Word vectors are not normalized
▷ Same words can take different forms

# Text Preprocessing

Token normalization

| Dog, dogs | → | dog |
| Bark, barks | → | bark |

Token normalization:

▷ **Stemming**: removing and replacing suffixes to get to the root of the word (**stem**)

Token normalization:

▷ **Stemming**: removing and replacing suffixes to get to the root of the word (**stem**)

▷ **Lemmatization**: to get base or dictionary form of a word (**lemma**)



| Playing | ⟶ | Play |
| Plays | ⟶ | Play |
| Played | ⟶ | Play |

Common root form 'play'

## Porter stemmer

▶ Published in 1979
▶ Base starting option

## Lancaster stemmer

▶ Published in 1990
▶ The most aggressive
▶ Easy adding of your own rules

## Snowball stemmer (Porter 2)

▶ Based on Porter
▶ More aggressive
▶ Most popular option now

Porter's stemmer:

▷ **Heuristics, applied one-by-one:**

■ SSES - SS (dresses - dress)

■ IES - I (ponies - poni)

■ S - <empty> (dogs - dog)

▷ **What's wrong?**

■ **Overstemming and understemming**

# Oversemming

University
Universal
Universities
Universe

} Univers

# Understemming

| Data | → | dat |

| Datum | → | datu |

Lemmatizer from NLTK:

- ▷ Tries to resolve word to its dictionary form
- ▷ Based on **WordNet** database
- ▷ For the best results feed part-of-speech tagger

**BIGDATA TEAM**



hyponymy

synonymy

3/27

22

Read more: https://en.wikipedia.org/wiki/WordNet

▶ NLTK

    ▷ nltk.stem.SnowballStemmer

    ▷ nltk.stem.PorterStemmer

    ▷ nltk.stem.WordNetLemmatizer

    ▷ nltk.corpus.stopwords

▶ BeautifulSoup (for parsing HTML)

▶ Regular Expressions (import re)

▶ Pymorphy2

- Capital Letters
- Punctuation
- Contractions (e.g, etc.)
- Numbers (dates, ids, page numbers)
- Stop-words ("the", "is", etc.)
- Tags

How to improve BOW?
▷ Use n-gramms instead of words!

The brown dog plays with a little cat

→

The brown

brown dog

dog plays

plays with

with a

a little

little

cat

The brown dog plays with a little cat

→

The brown

brown dog

dog plays

plays with

with a

a little

little

cat

Do we need all this bigramms?

The brown dog plays with a little cat → brown dog

dog plays
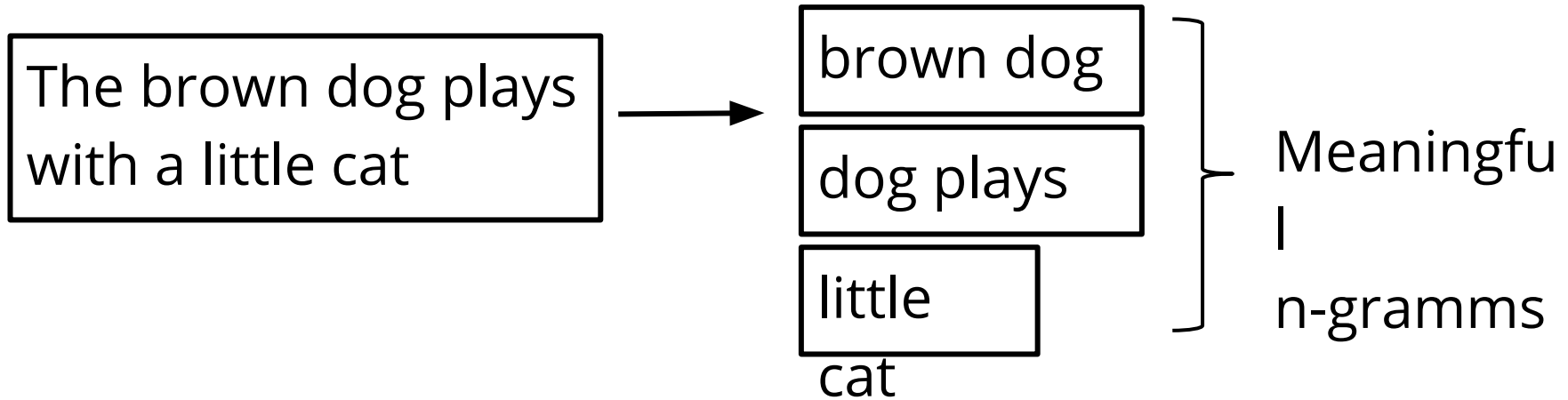
little cat

} Meaningful n-gramms

Meaningful n-gramms are often called **collocations**

How to detect meaningful n-gramms?

Delete:

▷ High-frequency n-gramms

  ■ Articles, prepositions

  ■ Auxiliary verbs (to be, to have, etc.)

  ■ General vocabulary

▷ Low-frequency n-gramms

  ■ Typos

  ■ Combinations that occur 1-2 times in a text

▶ **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 $t$ appears

29

*Sentence A:* The car is driven on the road.

*Sentence B:* The truck is driven on the highway.

(each sentence is a separate document)

| 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 | | | |

| 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 | | |

| 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 |

33

```
from sklearn.feature_extraction.text

import TfidfVectorizer
```

# Word Embeddings

**BIGDATA TEAM**

**One-hot vectors:**

```
                 Paris
         Rome                                              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]
```

**Problems:**
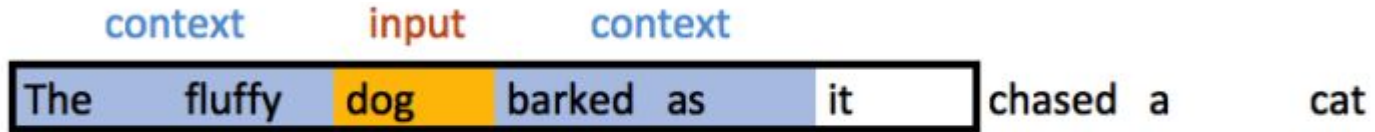
- Huge vectors
- VERY sparse
- No semantics or word similarity information included

Does vector similarity imply semantic similarity?

*"You shall know a word by the company it keeps"*
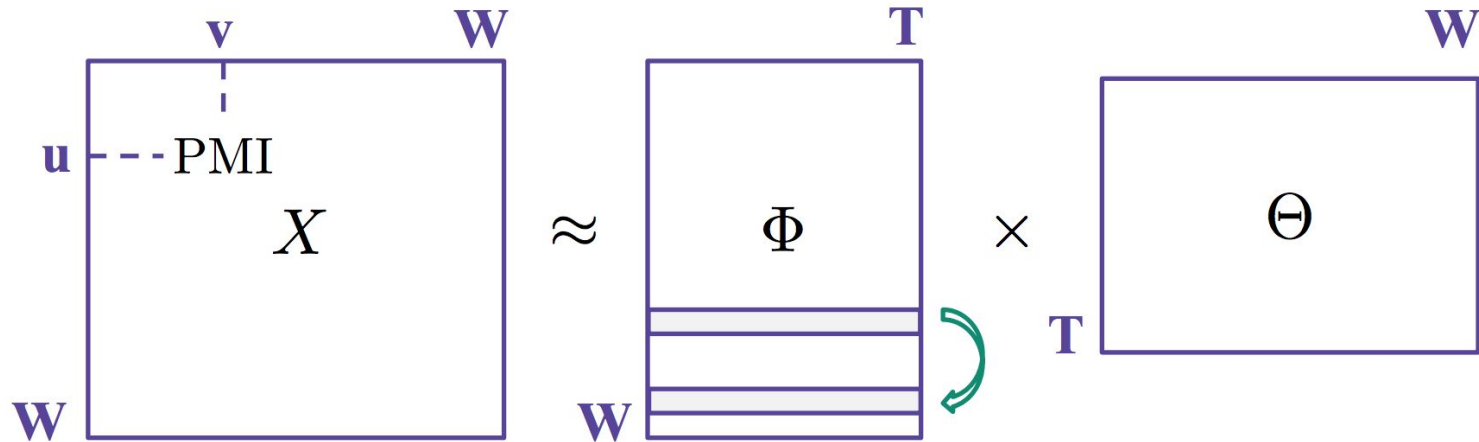
Firth, 1957

Input: **PMI, word coocurrences, etc.**

Method: **dimensionality reduction (SVD)**

Output: **word similarities**

Delete:

▷ High-frequency n-gramms

- ■ Articles, prepositions
- ■ Auxiliary verbs (to be, to have, etc.)
- ■ General vocabulary

▷ Low-frequency n-gramms

- ■ Typos
- ■ Combinations that occur 1-2 times in a text

Coocurrence counters in a window of fixed size

▷ $n_{uv}$ states for the number of times we've seen word $u$ and word $v$ together in the window

● Better solution: Pointwise Mutual Information (PMI)

$$PMI = log\frac{p(u,v)}{p(u)p(v)} = log\frac{n_{uv}n}{n_u n_v}$$

● Much better solution: **Positive PMI (pPMI)**

$$pPMI = \max(0, PMI)$$

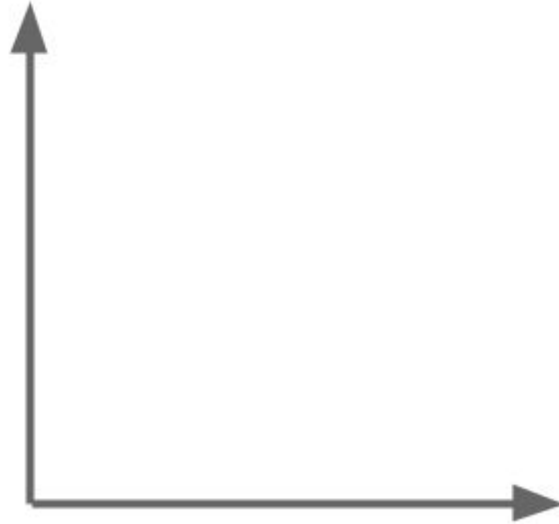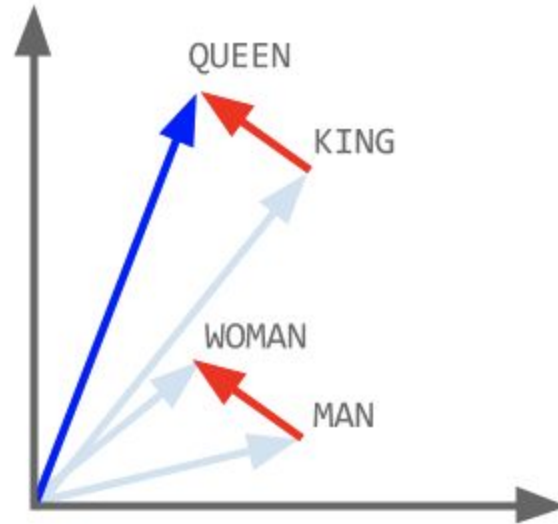| Frequency With Filter | PMI | T-test With Filter | Chi-Sq Test |
|---|---|---|---|
| (front, desk) | (universal, studios) | (front, desk) | (wi, fi) |
| (great, location) | (howard, johnson) | (great, location) | (cracker, barrel) |
| (friendly, staff) | (cracker, barrel) | (friendly, staff) | (howard, johnson) |
| (hot, tub) | (santa, barbara) | (hot, tub) | (la, quinta) |
| (clean, room) | (sub, par) | (continental, breakfast) | (front, desk) |
| (hotel, staff) | (santana, row) | (free, breakfast) | (universal, studios) |
| (continental, breakfast) | (e, g) | (great, place) | (santa, barbara) |
| (nice, hotel) | (elk, springs) | (parking, lot) | (santana, row) |
| (free, breakfast) | (times, square) | (customer, service) | (, more) |
| (great, place) | (ear, plug) | (desk, staff) | (flat, screen) |
| (desk, staff) | (la, quinta) | (walk, distance) | (french, quarter) |
| (parking, lot) | (fire, pit) | (comfortable, bed) | (elk, springs) |
| (customer, service) | (san, clemente) | (nice, hotel) | (walking, distance) |

# Why not to learn word vectors?

What is `king` - `man` + `woman`?

So king - man + woman = queen!
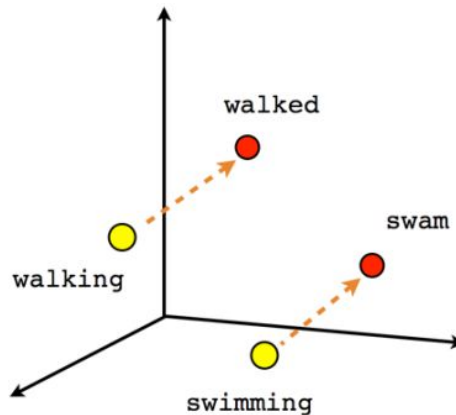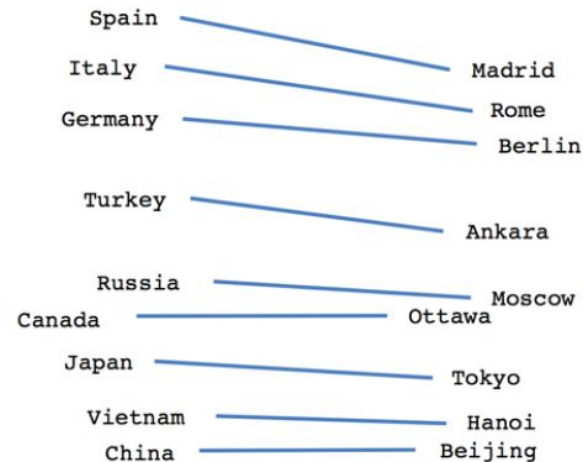
**BIGDATA TEAM**

**Word2vec** (Mikolov et al. 2013) - a framework for learning word embeddings



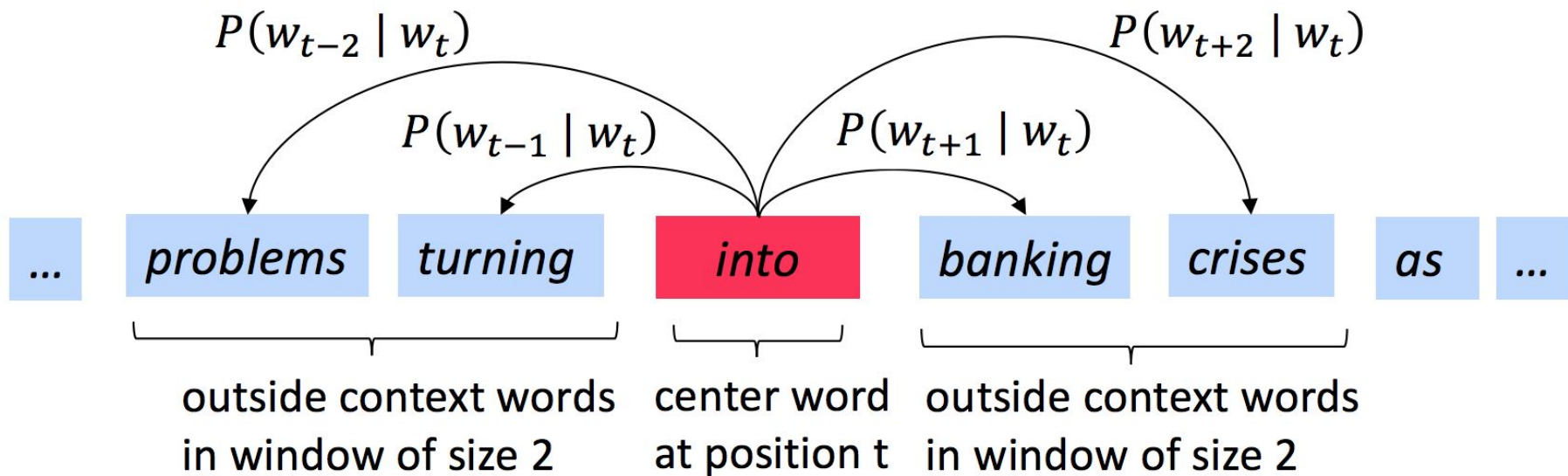Male-Female

Verb tense

Country-Capital

## Source Text

## Training Samples

| The | quick | brown | fox jumps over the lazy dog. ⟶ | (the, quick)<br>(the, brown) |

| The | quick | brown | fox | jumps over the lazy dog. ⟶ | (quick, the)<br>(quick, brown)<br>(quick, fox) |

| The | quick | brown | fox | jumps | over the lazy dog. ⟶ | (brown, the)<br>(brown, quick)<br>(brown, fox)<br>(brown, jumps) |

| The | quick | brown | fox | jumps | over | the lazy dog. ⟶ | (fox, quick)<br>(fox, brown)<br>(fox, jumps)<br>(fox, over) |

Source: http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/

$P(w_{t-2} \mid w_t)$ $\qquad$ $P(w_{t+2} \mid w_t)$

$P(w_{t-1} \mid w_t)$ $\qquad$ $P(w_{t+1} \mid w_t)$

... | problems | turning | into | banking | crises | as | ...

outside context words in window of size 2

center word at position t

outside context words in window of size 2

Source: CS224n: http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture01-wordvecs1.pdf

$$P(w_{t-2} \mid w_t)$$

$$P(w_{t+2} \mid w_t)$$

$$P(w_{t-1} \mid w_t)$$

$$P(w_{t+1} \mid w_t)$$

... | problems | turning | into | banking | crises | as | ...

outside context words in window of size 2

center word at position t

outside context words in window of size 2

48

Source: CS224n: http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture01-wordvecs1.pdf

**BIGDATA TEAM**

**Output Layer**
**Softmax Classifier**

**Hidden Layer**
**Linear Neurons**

**Input Vector**

| |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |

A '1' in the position corresponding to the word "ants"

| |
|---|
| 0 |

*10,000 positions*

Σ

Σ

Σ

*300 neurons*

Σ → Probability that the word at a randomly chosen, nearby position is "**abandon**"

Σ → ... "**ability**"

Σ → ... "**able**"

Σ → ... "**zone**"

*10,000 neurons*

Source: http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/

Hidden Layer Weight Matrix → Word Vector Lookup Table!

Source: http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/

- Word vectors with 300 components
- Vocabulary of 10,000 words.
- Weight matrix with 300 x 10,000 = 3 million weights each!

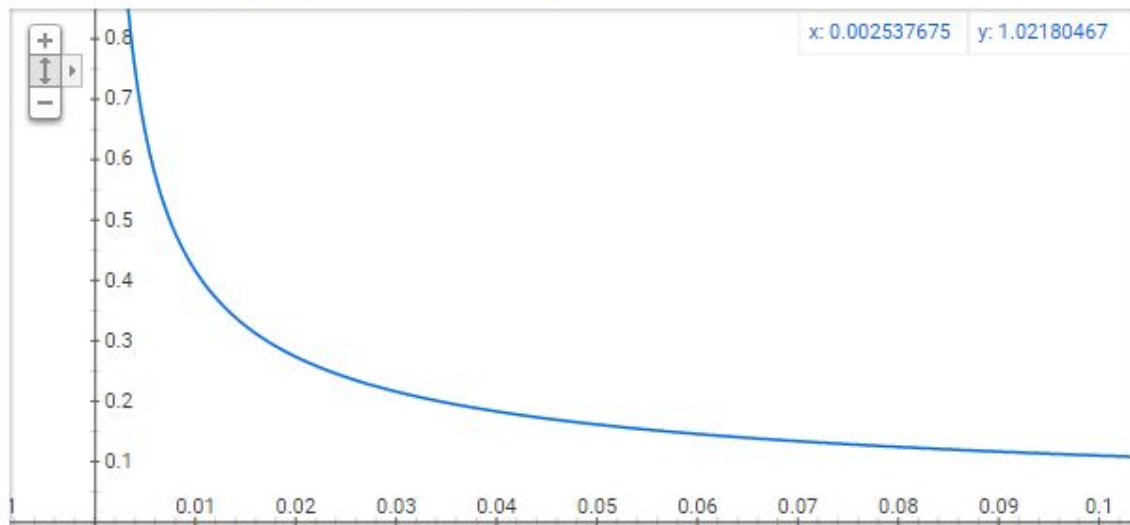Training is too long and computationally expensive
## How to fix this?

Basic approaches:

1. Treating common word pairs or phrases as single "words" in their model.
2. Subsampling frequent words to decrease the number of training examples.
3. Modifying the optimization objective with a technique they called "Negative Sampling", which causes each training sample to update only a small percentage of the model's weights.

## Subsampling frequent words.

$w_i$ is the word, $z(w_i)$ is the fraction of this word in the text

Graph for (sqrt(x/0.001)+1)*0.001/x

| | x: 0.002537675 | y: 1.02180467 |



$P(w_i)$ is the probability of *keeping* the word:

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

53

Source: http://mccormickml.com/2017/01/11/word2vec-tutorial-part-2-negative-sampling/

Negative Sampling idea: only few words error is computed. All other words have zero error, so no updates by the backprop mechanism.

More frequent words are selected to be negative samples more often.The probability for a selecting a word is just it's weight divided by the sum of weights for all words.

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

## Continuous BOW (CBOW)

$$p(w_i \mid w_{i-h} \ldots, w_{i+h})$$

Predict center word from (bag of) context words

- Predicting one word each time
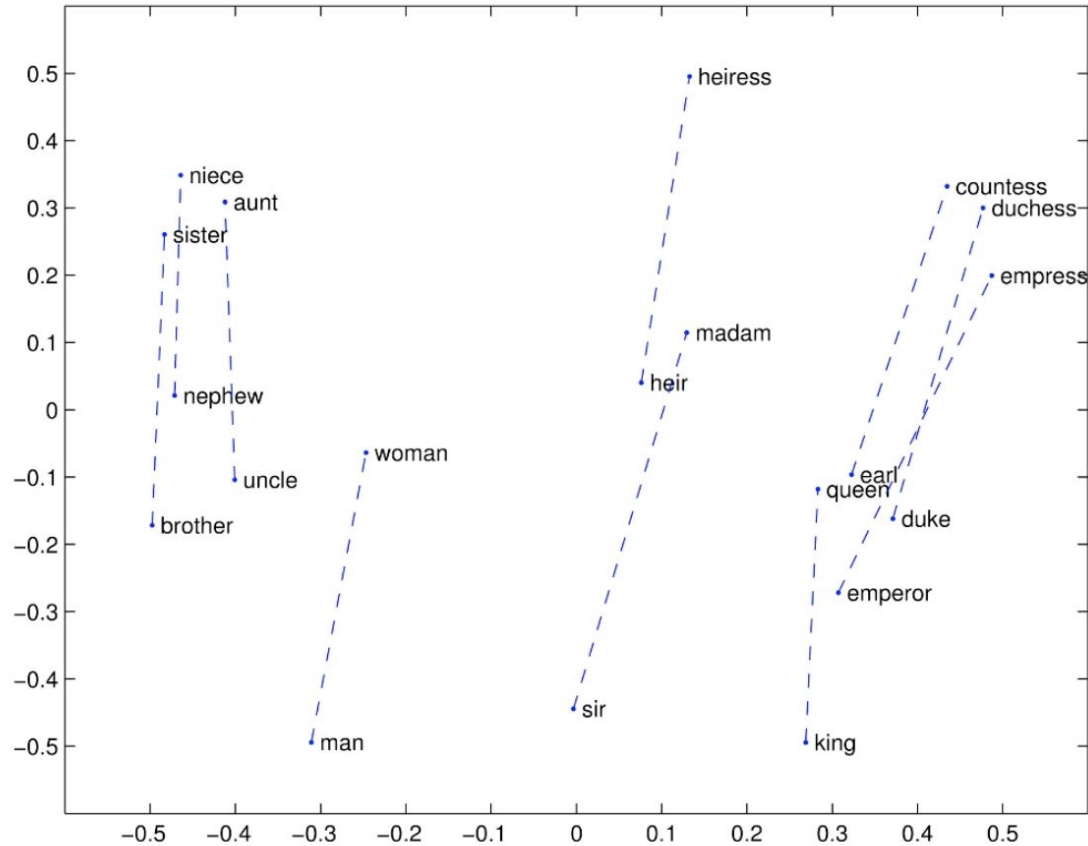- Relatively fast

## Skip-gram

$$p(w_{i-h}, \ldots w_{i+h} \mid w_i)$$

Predict context ("outside") words (position independent) given center word

- Predicting context by one word
- Much slower
- Better with infrequent words

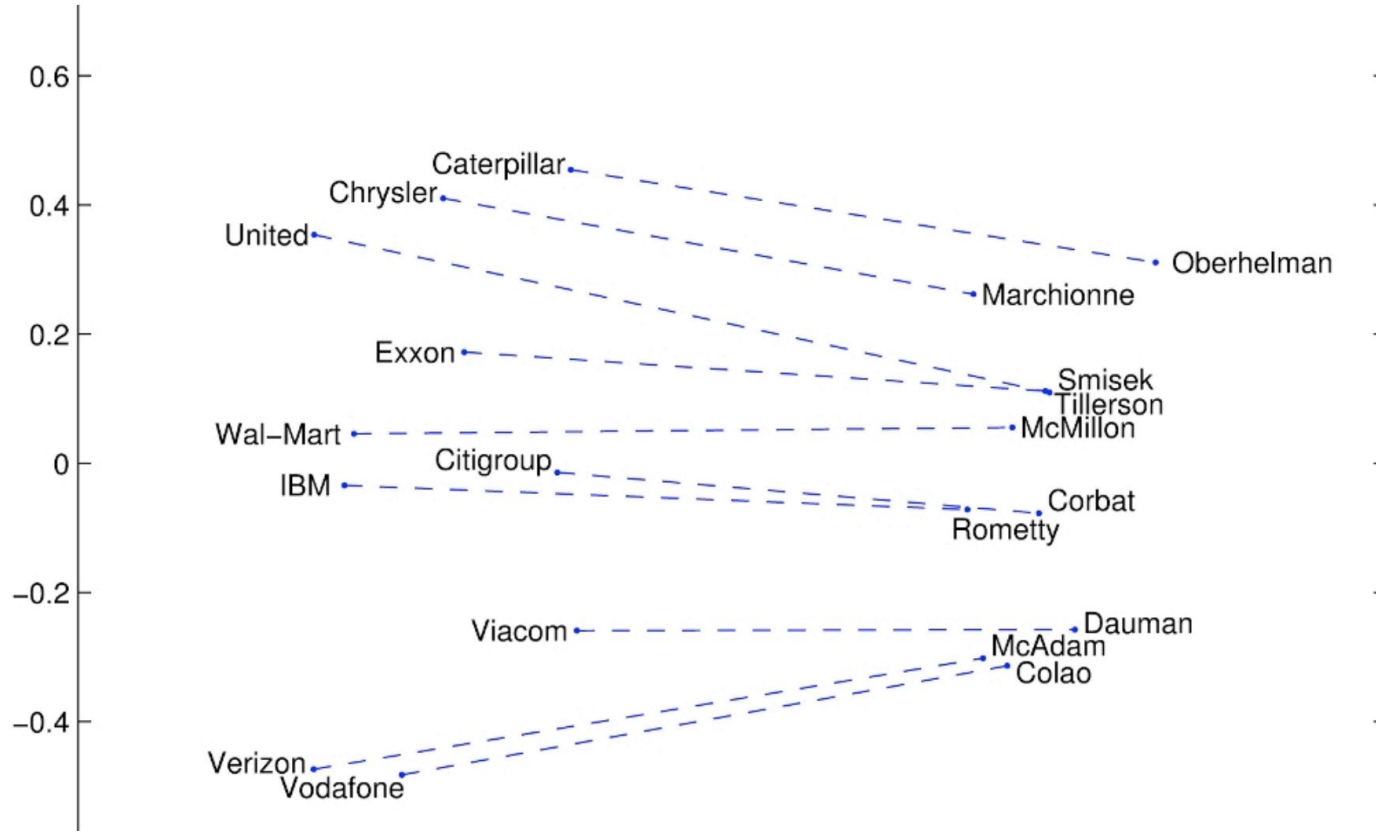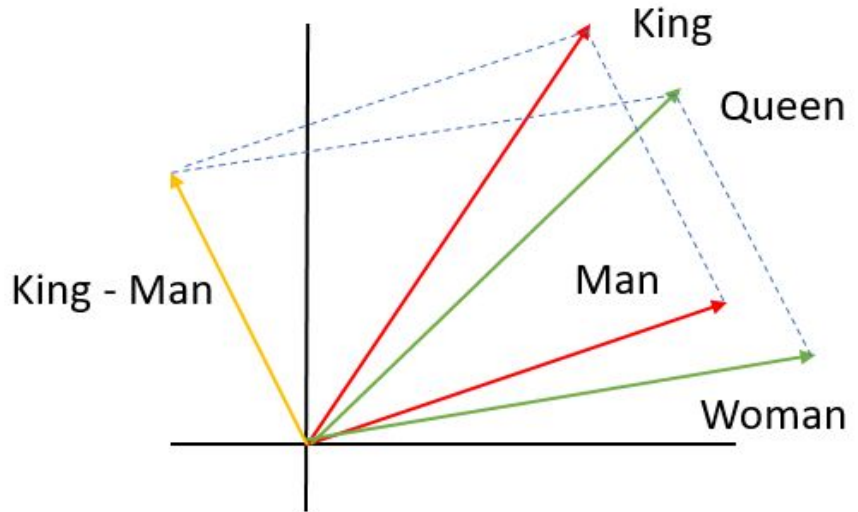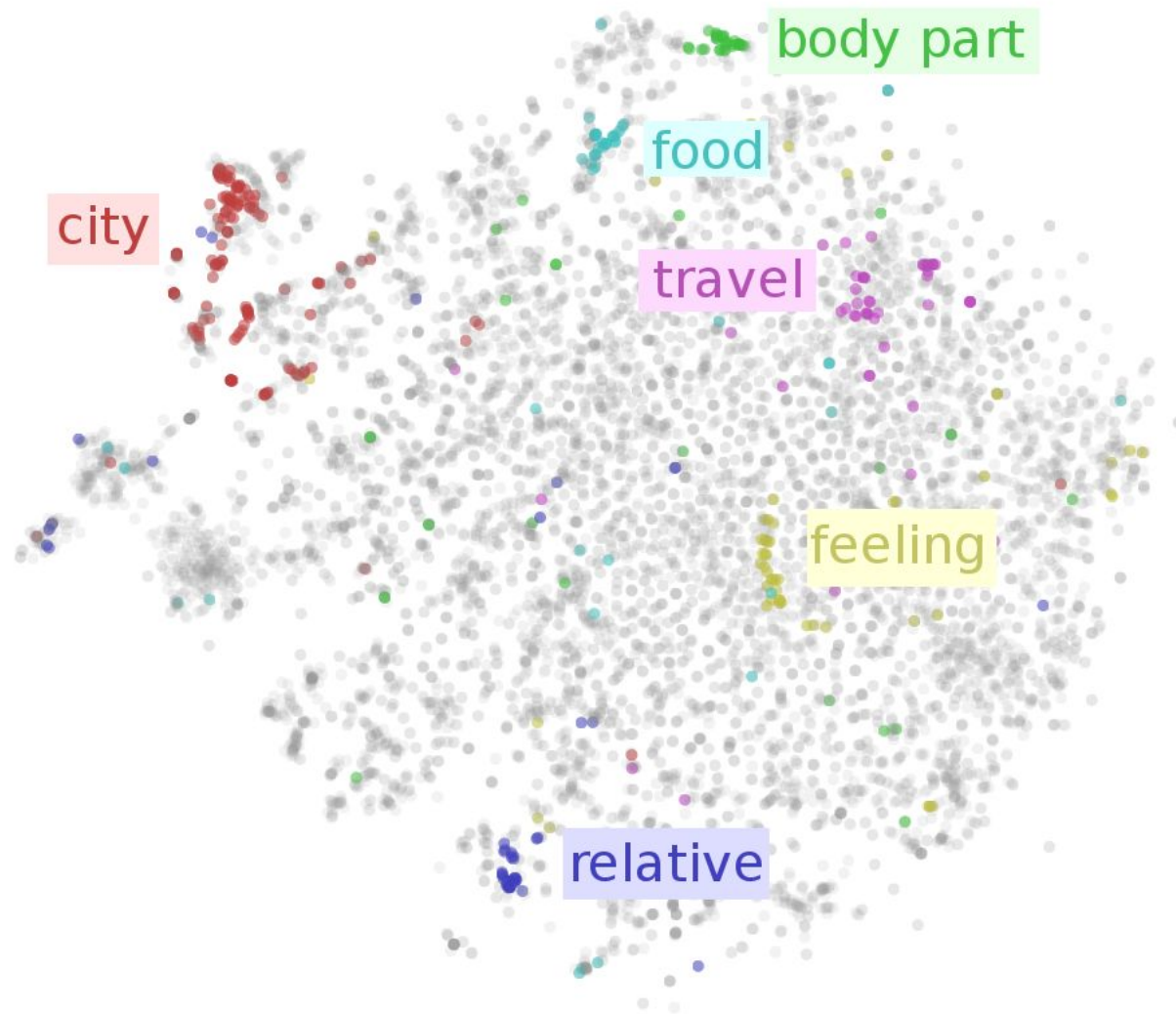King - man + woman = queen

$$x \qquad y \qquad y' \qquad target$$

$$\cos(x-y+y', target) \rightarrow \max_{target}$$

body part

food

city

travel

feeling

relative

Word vectors are simply vectors of numbers that represent the meaning of a word

Approaches:

- One-hot encoding
- Bag-of-words models
- Counts of word / context co-occurrences
- TF-IDF
- Predictions of context given word (skip-gram neural network models, e.g. word2vec)