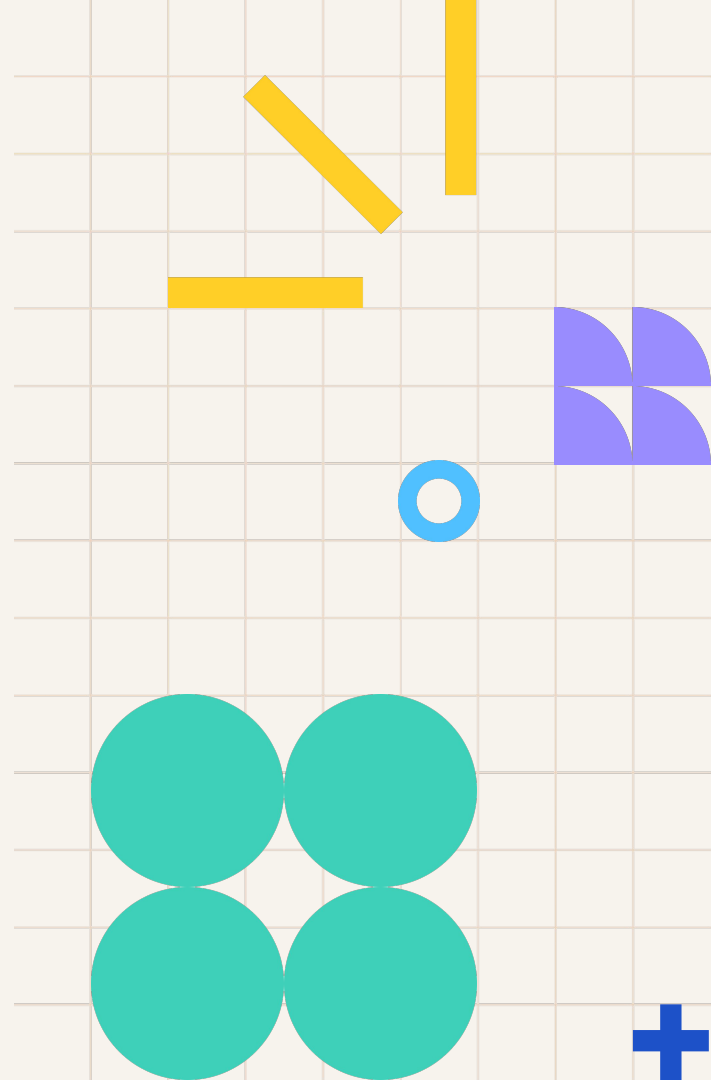




Introduction to NLP

Aydar Bulatov

DeepPavlov.ai
MIPT



Today

01 What is NLP

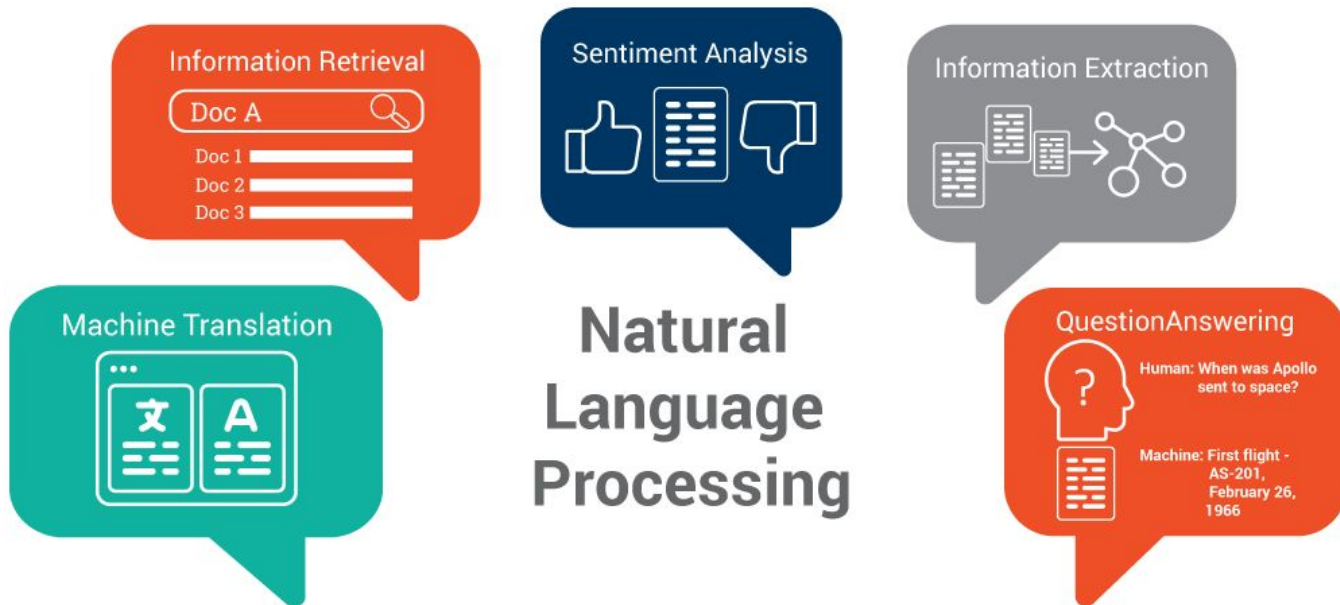
02 Text representation

03 Text processing



What is Natural Language Processing?





ML tasks in NLP

- **Text classification**
 - Sentiment, emotion, opinion analysis
 - Word sense disambiguation
- **Sequence labelling**
 - P-O-S tagging, named entity recognition, dialogue acts
- **Parsing dependencies**
- **Coreference resolution**
- **Text generation**
 - QA, summarization
 - data-to-text generation
 - dialogue systems
- **Learning without supervision**
 - clustering
 - matrix factorization
 - latent semantic indexing
- **Semi-supervised learning**
 - graph-based
 - clustering + classification

Word representations



Vector word representations

- dog [0.033, -0.009, 0.037, 0.068, ..., 0.053]
- cat [0.067, 0.069, -0.015, -0.0249, ... , 0.0317]
- home [-0.0004, 0.0199, 0.03, -0.04, ... , -0.031]
- computer [-0.0016 , 0.0872, -0.0314, -0.0523, ... , 0.0228]
- love [-0.169, 0.0382, 0.0084, -0.0182, ... , -0.0209]

Vector word representations

- dog [0.033, -0.009, 0.037, 0.068, ..., 0.053]
- cat [0.067, 0.069, -0.015, -0.0249, ... , 0.0317]
- home [-0.0004, 0.0199, 0.03, -0.04, ... , -0.031]
- computer [-0.0016 , 0.0872, -0.0314, -0.0523, ... , 0.0228]
- love [-0.169, 0.0382, 0.0084, -0.0182, ... , -0.0209]



- Wikipedia
- Twitter
- НКРЯ - Национальный корпус русского языка
(собрание русских текстов в электронной форме)
- Тайга - корпус русского языка (в основе, художественная литература)

Vector word representations

Why not just use one-hot encoding?

- dog: [1, 0, 0, 0, ..., 0]
- cat: [0, 1, 0, 0, ... , 0]
- home: [0, 0, 1, 0, ... , 0]
- computer: [0, 0, 0, 1, ... , 0]
- love: [0, 0, 0, 0, ... , 1]

Vector word representations

Why not just use one-hot encoding?

- dog: [1, 0, 0, 0, ..., 0]
- cat: [0, 1, 0, 0, ... , 0]
- home: [0, 0, 1, 0, ... , 0]
- computer: [0, 0, 0, 1, ... , 0]
- love: [0, 0, 0, 0, ... , 1]

These vectors know nothing about the word!

Vector word representations

What properties should a representation share with its word?

Vector word representations

What properties should a representation share with its word?

- run ~ running
- walk ~ walked
- to describe ~ descriptor

Vector word representations

What properties should a representation share with its word?

- run ~ running
- walk ~ walked
- to describe ~ descriptor

- big ~ huge
- to buy ~ to purchase
- pineapple ~ coconut

Distributional semantics

You shall know a word by the company it keeps

Firth, 1957

Oculist and eye-doctor . . . occur in almost the same *environments*.

If A and B have almost identical *environments* we say that they are *synonyms*

Zellig Harris (1954)

Distributional semantics

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

Distributional semantics

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

But what if we know the following contexts:

1. a bottle of *tezgüino* is on the table,
2. everybody likes *tezgüino*,
3. *tezgüino* makes you drunk,
4. we make *tezgüino* out of corn.

Distributional semantics

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

But what if we know the following contexts:

1. a bottle of *tezgüino* is on the table,
2. everybody likes *tezgüino*,
3. *tezgüino* makes you drunk,
4. we make *tezgüino* out of corn.

Now we know!

Tezgüino is a kind of alcoholic beverage made from corn!

Distributional semantics

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

But what if we know the following contexts:

1. a bottle of *tezgüino* is on the table,
2. everybody likes *tezgüino*,
3. *tezgüino* makes you drunk,
4. we make *tezgüino* out of corn.

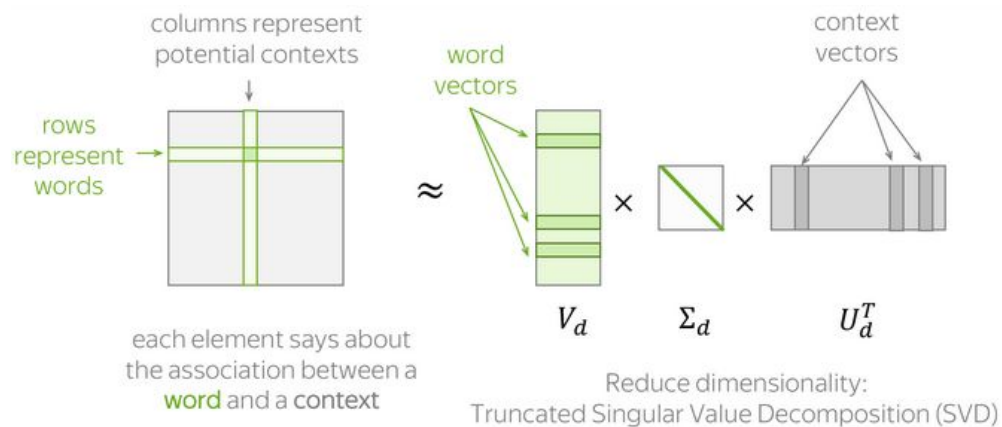
	context			
	1.	2.	3.	4.
tezgüino	1	1	1	1
loud	0	0	0	0
motor oil	1	0	0	1
tortillas	0	1	0	1
wine	1	1	1	0

Now we know!

Tezgüino is a kind of alcoholic beverage made from corn!

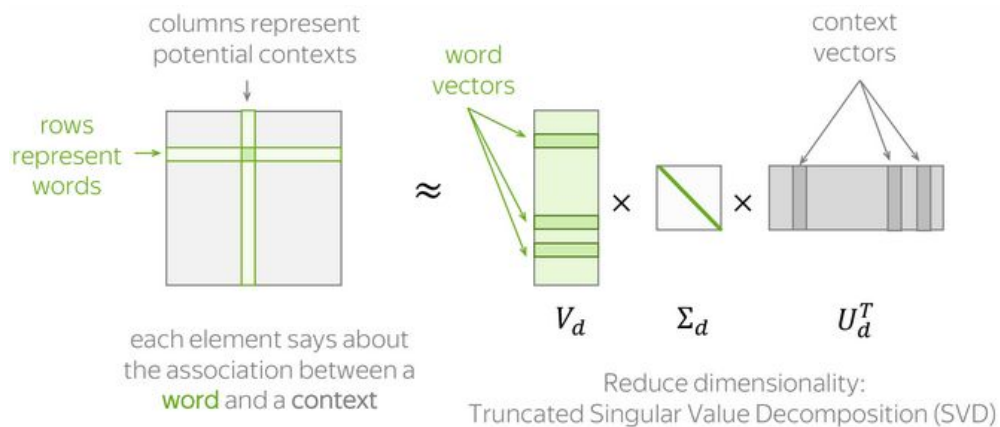
Count-based methods

Main idea: We have to put information about contexts into word vectors.



Count-based methods

Main idea: We have to put information about contexts into word vectors.

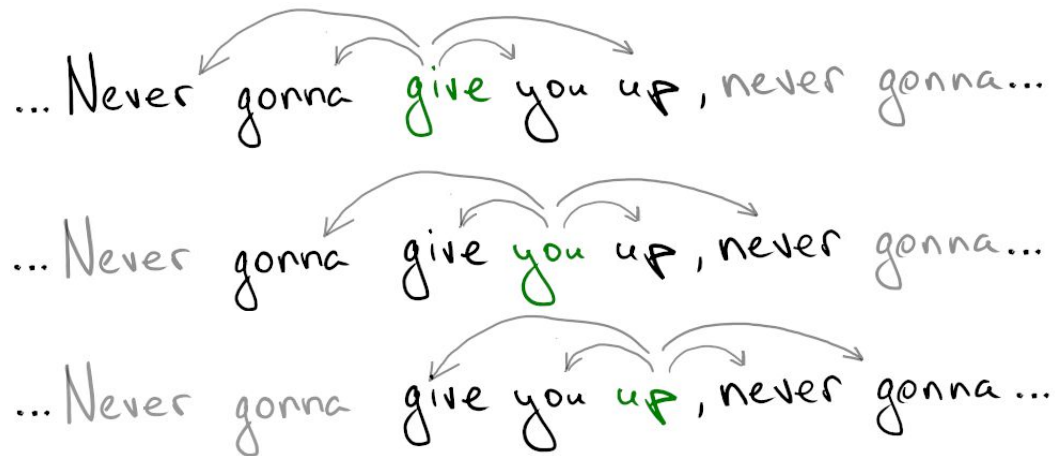


Ok but how do we represent context information?

Count-based methods

... Never gonna give you up, never gonna...

Count-based methods



Count-based methods

... Never gonna give you up, never gonna...

context

1. Co-Occurrence counts

$N(w, c)$ = number of times w appears in c

2. Tf-Idf

$$\text{tf-idf}(w, d, D) = \text{tf}(w, d) \cdot \text{idf}(w, D) = N(w, d) \cdot \log\left(\frac{|D|}{|\{d \in D : w \in d\}|}\right)$$

3. Positive PMI

$$\text{ppmi}(w, c) = \max(0, \text{PMI}(w, c))$$

$$\text{pmi}(w, c) = \log\left(\frac{P(w, c)}{P(w)P(c)}\right) = \log\left(\frac{N(w, c) \cdot I(w, c)}{N(w)N(c)}\right)$$

Word2Vec

Idea: teach embeddings to predict their contexts.

Algorithm

- 1) use a large text corpus
- 2) for each word start with 2 vectors: central \mathbf{u} + context \mathbf{v}
- 3) go over the text with a sliding window (*use negative sampling*)
- 4) compute probabilities for words from c based on w
- 5) adjust the vectors to increase probabilities

$$P(o|c) = \text{SM}(\mathbf{u}_o^T \cdot \mathbf{v}_c) = \frac{e^{\mathbf{u}_o^T \mathbf{v}_c}}{\sum_{w \in V} e^{\mathbf{u}_w^T \mathbf{v}_c}}$$

Word2Vec

Idea: teach embeddings to predict their contexts.

Algorithm

- 1) use a large text corpus
- 2) for each word start with 2 vectors: central \mathbf{u} + context \mathbf{v}
- 3) go over the text with a sliding window (*use negative sampling*)
- 4) compute probabilities for words from c based on w
- 5) adjust the vectors to increase probabilities

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

Objective function:

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

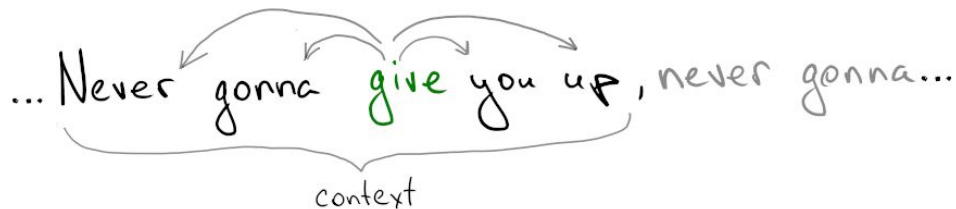
$$\text{Loss} = J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \log P(w_{t+j} | w_t, \theta)$$

agrees with our plan above \Rightarrow go over text \nearrow with a sliding window \nearrow compute probability of the context word given the central

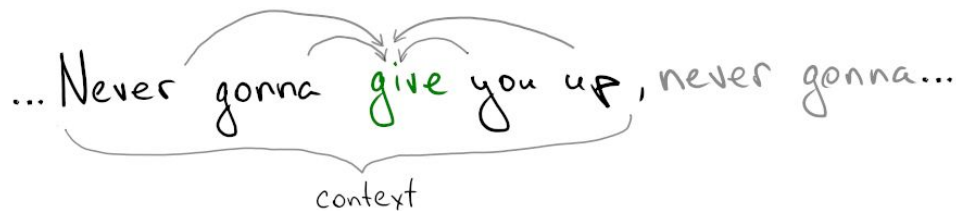
[Efficient Estimation of Word Representations in Vector Space](#) (2013)

Word2Vec

1. Skip-Gram
from central **word** predict context



2. CBOW
from sum of context predict **word**



Word2Vec



GloVe

Idea: use **global** corpus statistics to learn embedding **vectors**

W2V:
$$J_{t,j}(\theta) = -\log P(\text{cute}|\text{cat}) = -\log \frac{\exp u_{\text{cute}}^T v_{\text{cat}}}{\sum_{w \in \text{Voc}} \exp u_w^T v_{\text{cat}}} = -u_{\text{cute}}^T v_{\text{cat}} + \log \sum_{w \in \text{Voc}} \exp u_w^T v_{\text{cat}}.$$

GloVe

Idea: use **g**lobal corpus statistics to learn embedding **v**ectors

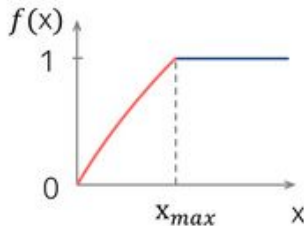
W2V:
$$J_{t,j}(\theta) = -\log P(\text{cute}|\text{cat}) = -\log \frac{\exp u_{\text{cute}}^T v_{\text{cat}}}{\sum_{w \in \text{Voc}} \exp u_w^T v_{\text{cat}}} = -u_{\text{cute}}^T v_{\text{cat}} + \log \sum_{w \in \text{Voc}} \exp u_w^T v_{\text{cat}}.$$

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

context vector \swarrow
word vector \swarrow
bias terms (also learned) \swarrow

Weighting function to:

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

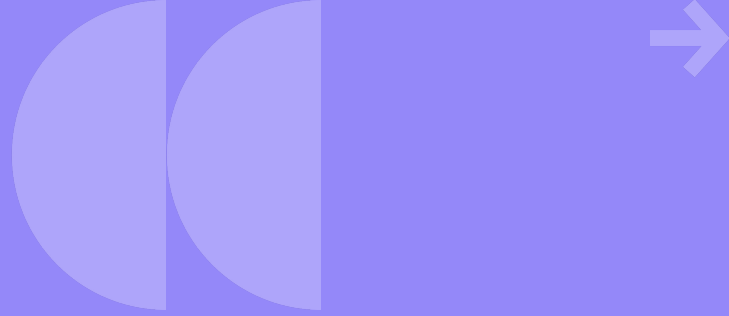


$$\begin{cases} (x/x_{\max})^\alpha & \text{if } x < x_{\max}, \\ 1 & \text{otherwise.} \end{cases}$$

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



Processing language



Problems

- OOV: what does “wher” mean?
- word ordering

Problems

- OOV: what does “wher” mean?
- word ordering

FastText (2016, Facebook)

where = <wh, whe, her, ere, re> and <where>

wher = <wh, er>

Problems

- OOV: what does “wher” mean?
- word ordering

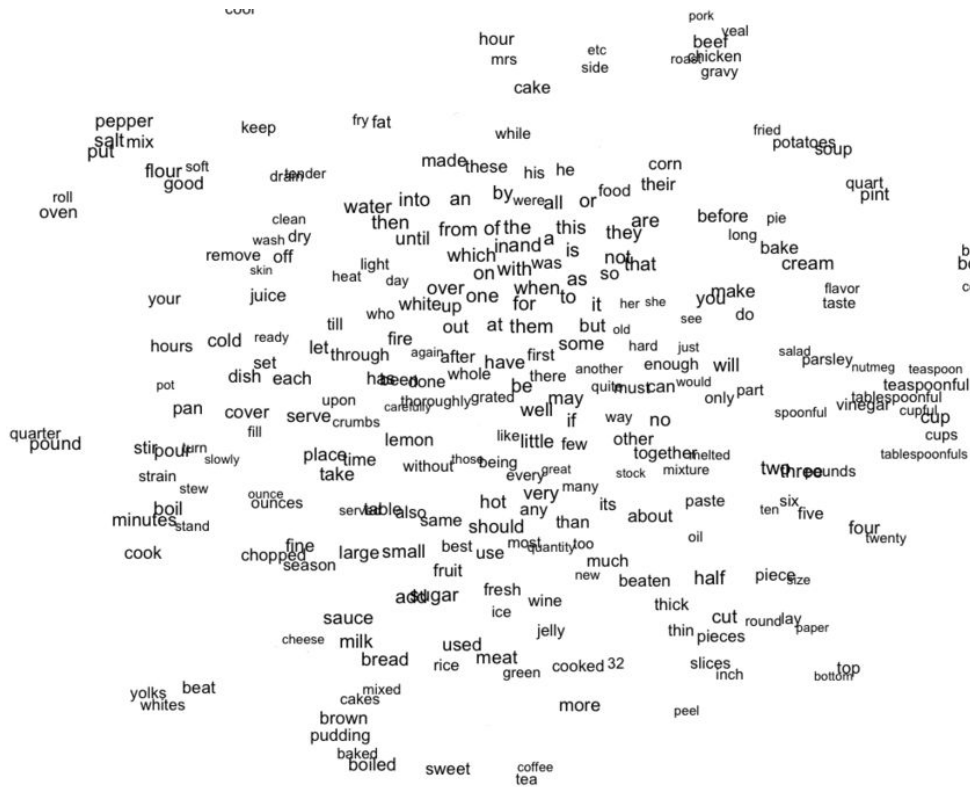
FastText (2016, Facebook)

where = <wh, whe, her, ere, re> and <where>

wher = <wh, er>

- ambiguity (context-dependency): what does “plant” mean?
- Los Angeles, Barack Obama, Doctor House

What now?



Operations with embeddings

$$p = (p_1, p_2, \dots, p_n)$$

$$q = (q_1, q_2, \dots, q_n)$$

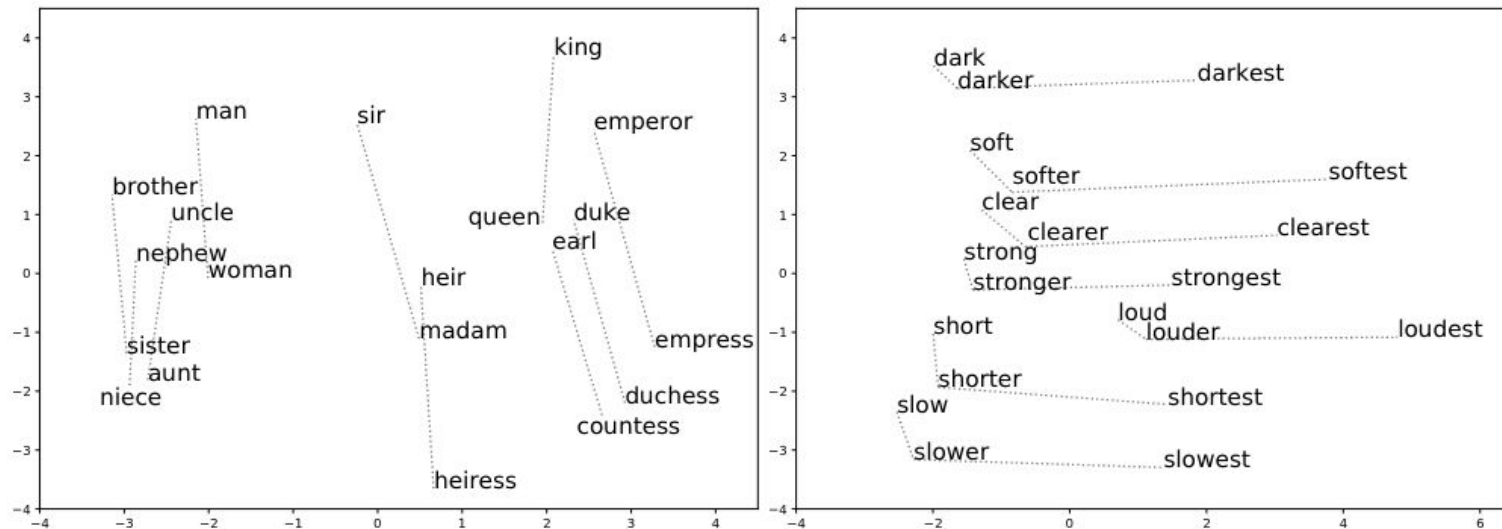
Euclidean distance

$$d(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

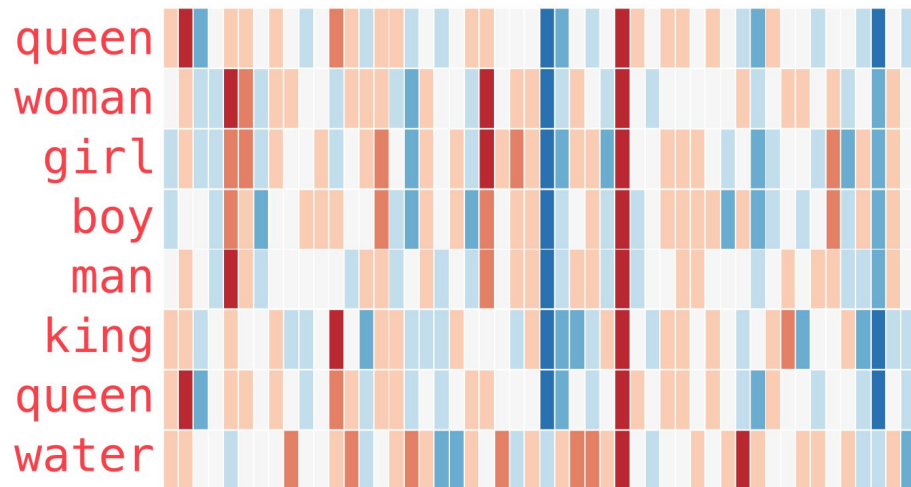
Cosine distance

$$\cos(p, q) = 1 - \frac{P \times Q}{\|P\|_2 \|Q\|_2} = 1 - \frac{\sum_{i=1}^n p_i \times q_i}{\sqrt{\sum_{i=1}^n p_i^2} \times \sqrt{\sum_{i=1}^n q_i^2}}$$

Operations with embeddings



Operations with embeddings



king – queen \approx man – woman

man – boy \approx ? – girl

japan – sushi \approx ? – pasta

euro – france \approx ? – russia

Centroid(january, february, march) \approx ?

Centroid(moscow, siberia, caucasus) \approx ?

Centroid(sunday, monday, tuesday) \approx ?

How to process sequences?

...Never gonna give you up, never gonna...

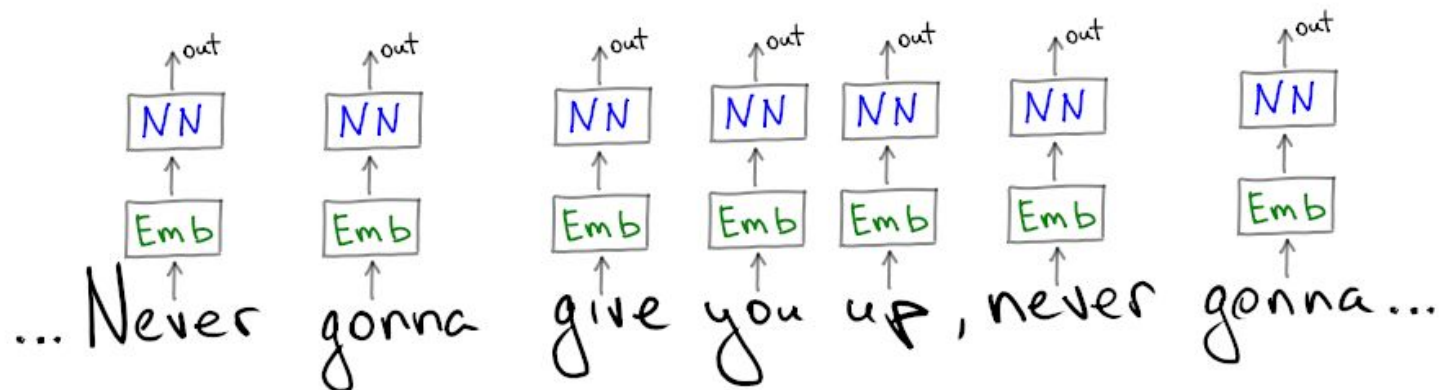
How to process sequences?

... Never gonna give you up, never gonna ...

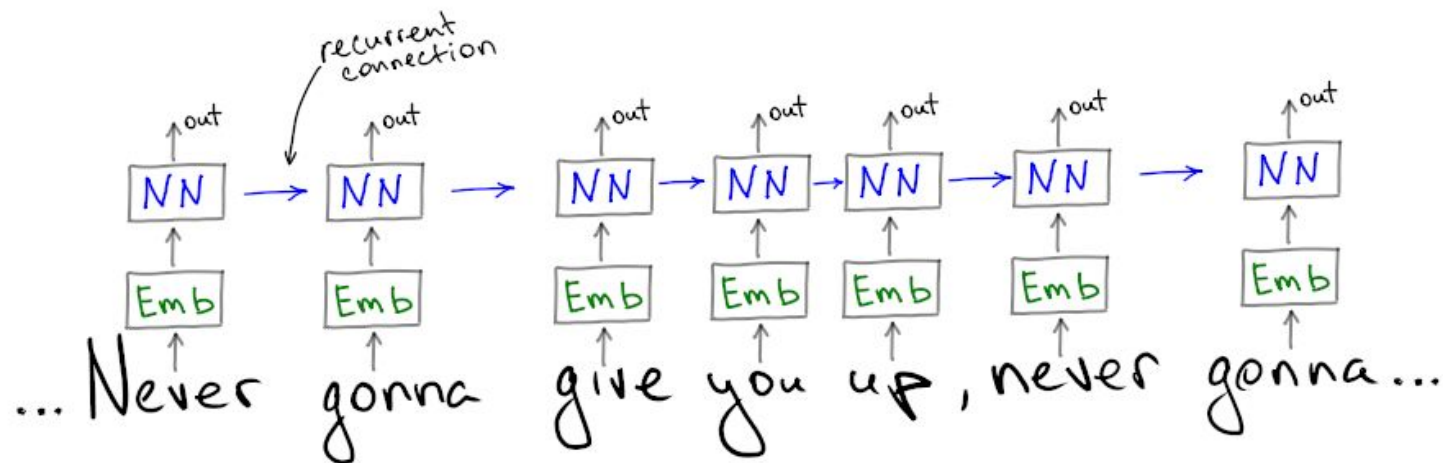


A diagram illustrating the sequential processing of a sentence. The sentence is written in a cursive, handwritten style: "... Never gonna give you up, never gonna ...". Above the text, a series of seven curved arrows point from left to right, connecting each word to the next word in the sequence. The arrows start above 'Never', point to 'gonna', then to 'give', then to 'you', then to 'up,', then to 'never', and finally to 'gonna...'. This visualizes the step-by-step processing of the sequence of words.

How to process sequences?



How to process sequences?



Next: RNN, LSTM



Artificial Intelligence Research Institute

airi.net



[airi_research_institute](https://t.me/airi_research_institute)



[AIRI Institute](https://vk.com/AIRI_Institute)



[AIRI Institute](https://www.youtube.com/AIRI_Institute)



[AIRI_inst](https://twitter.com/AIRI_inst)



[artificial-intelligence-research-institute](https://www.linkedin.com/company/artificial-intelligence-research-institute)