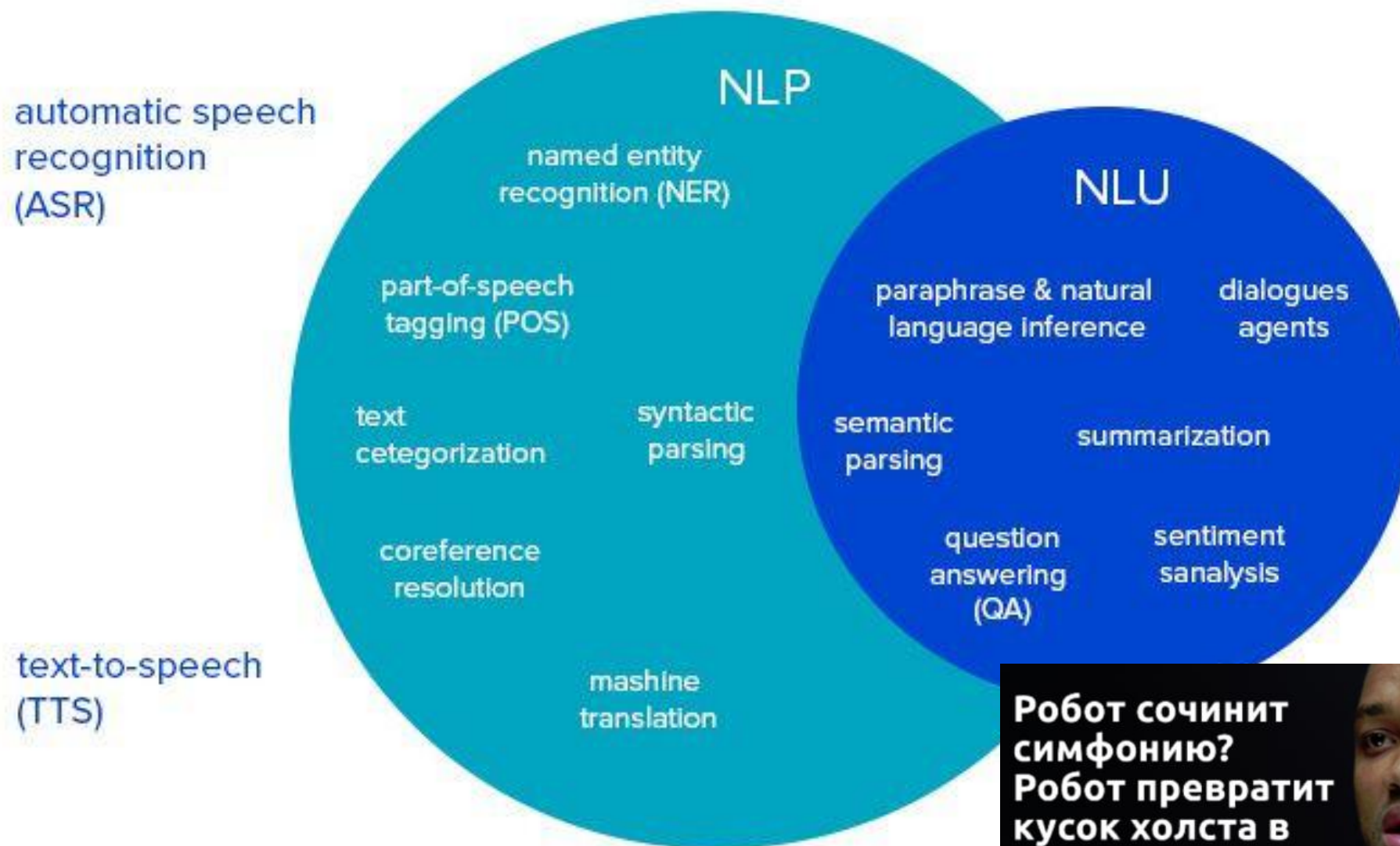


Lecture 4 – NLP. Part 1.

From Words to Vectors

Natural language processing



**Робот сочинит симфонию?
Робот превратит кусок холста в шедевр искусства?**





natasha

nae>al

raz|del

slovnet()

navec

corus.tar.gz

yargy

nerus.conllu

ipymarkup

Possible approaches to make vectors

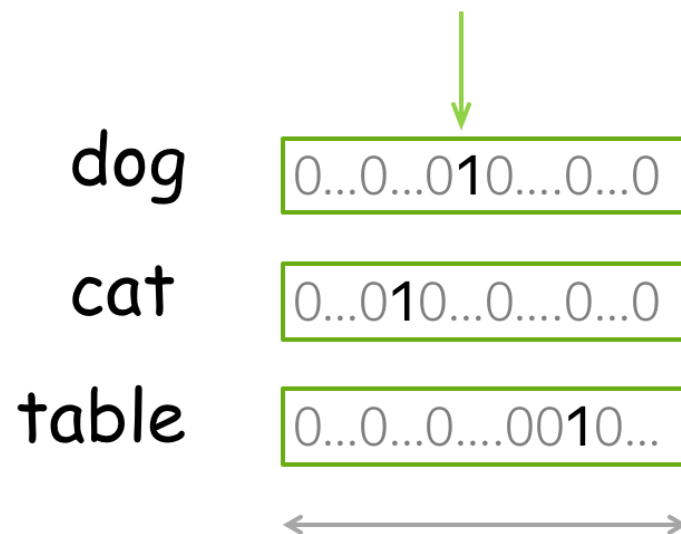
- Naive
 - One-hot
 - BoW
 - TF-IDF
- Count based
 - Co-occurency
 - PPMI & co
 - LSA/LSI
- Prediction-Based
 - Word2Vec
 - GloVe
 - FastText



Transformer-based

One-Hot Encoding

One is 1, the rest are 0



Embedding dimension =
vocabulary size

- What if on test we have word that is not in vocabulary?
- How to implement?
- Why not? What a problem?

Bag-of-Words (BoW)

What the difference with One-Hot to sentences?

Bag-of-Words (BoW)

- One-hot is originally for words, if you use it for sentences – it summarize one-hot vectors for words in sentence, but without repetition. BoW – update of One-Hot for sentences, that allow repetitions of words and display it on final vector.

*He is good player,
he is really good!*

	<i>he</i>	<i>good</i>	<i>work</i>	<i>really</i>
BoW	2	2	0	1
One-Hot	1	1	0	1

Term Frequency-Inverse Document Frequency

TF

Measures how often a term appears in a specific document

- Variants:
 - Raw count
 - Boolean (presence/absence)
 - Log-scale: $1 + \log(tf)$
 - Double-normalization:
 $0.5 + 0.5 * \max(tf \text{ in doc})$

IDF

Measures how rare/common a term is across the whole document corpus

- Variants:
 - Classic: $\log\left(\frac{N}{df_t}\right)$
 - Smooth: $\log\left(\frac{N+1}{df_t+1}\right) + 1$
 - Probabilistic: $\log\left(\frac{N-df_t}{df_t}\right)$
 - Max-IDF..

Overall problem

- How to use connections between words?

Use n-grams

Benefits:

- captures local context and word order (e.g. “not good” \neq “good”)
- helps with multi-word expressions / phrases (“New York”, “machine learning”)
- improves discrimination in tasks like sentiment, topic classification, document similarity

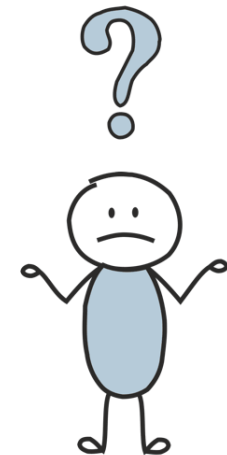
Trade-offs:

- increases feature space dimensionality \rightarrow more sparse matrices, more memory/time cost
- risk of overfitting on rare n-grams if data small
- need to tune n (2, 3, maybe more) and frequency thresholds

What is meaning?

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

(We hope you do not)



What is meaning?

Now look how this word is used in different contexts:

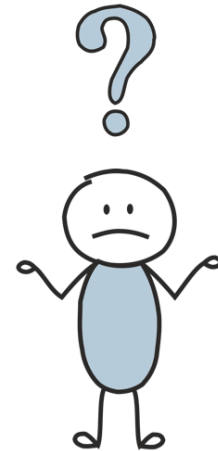
A bottle of **tezgüino** is on the table.

Everyone likes **tezgüino**.

Tezgüino makes you drunk.

We make **tezgüino** out of corn.

Can you understand what **tezgüino** means ?



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

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.



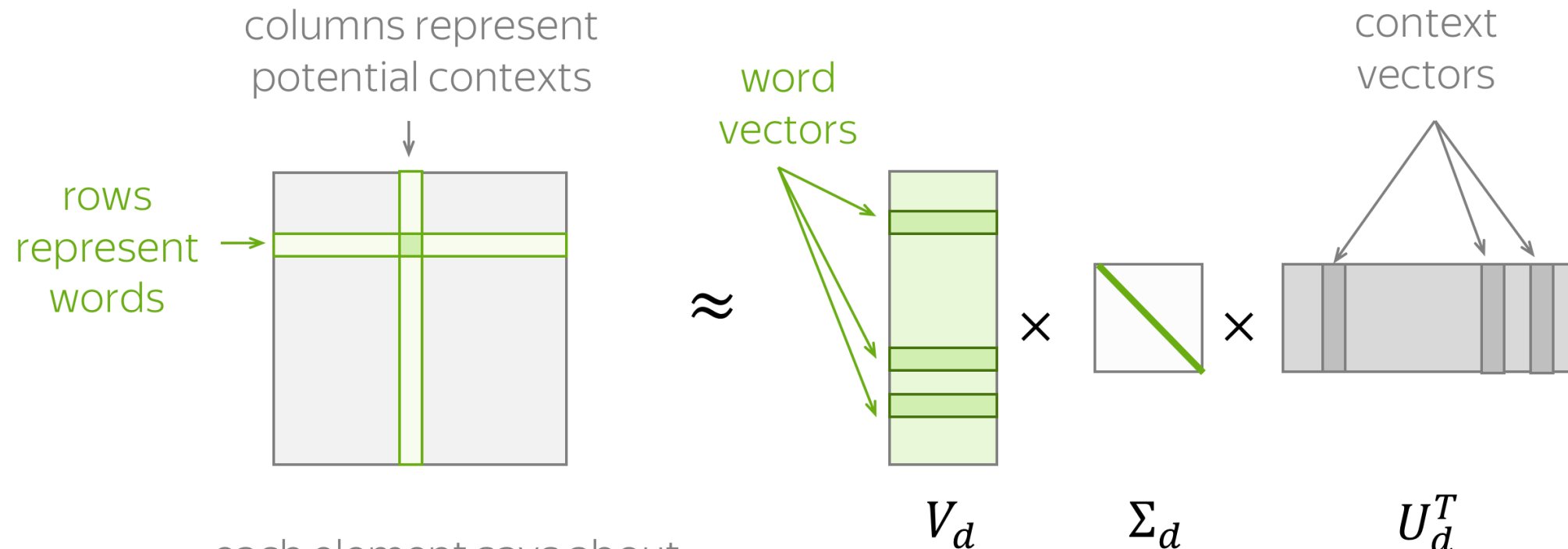
Tezgüino is a kind of alcoholic beverage made from corn.

With context, you can understand the meaning!



Count-based approaches

Put this information **manually**, based on global corpus statistics.



Reduce dimensionality:
Truncated Singular Value Decomposition (SVD)

What is context and matrix element

2-sized window for **cat**

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

Context:

- surrounding words in a L-sized window

contexts for **cat**

Matrix element:

- $N(w, c)$ – number of times **word** w appears in context c

About context

- Are context words at different distances equally important? If not, how can we modify co-occurrence counts?

[HAL paper](#)

- Context before interesting than after?
In [HAL paper](#) there are two context – left + right

- How to choose length of window?

[Redefining Context Windows for Word Embedding Models](#)

About matrix element

Context:

- surrounding words in a L-sized window

Matrix element:

- Raw count / conditional probability
- Positive Pointwise Mutual Information
- PPMI with Context Distribution Smoothing) – $P(c)^\alpha$, $\alpha < 1$ - increase probability rare context
- t-statistic - $\frac{P(w,c) - P(w)P(c)}{\sqrt{P(w,c)/N}}$, N – corpus size
- Shifted PPMI (SPPMI) – when PPMI is huge and sparse
- * SVD Weighting – additional weights for not null elements after PPMI

- $\text{PPMI}(w, c) = \max(0, \text{PMI}(w, c))$,
where
$$\text{PMI}(w, c) = \log \frac{P(w, c)}{P(w)P(c)} = \log \frac{N(w, c)|(w, c)|}{N(w)N(c)}$$

Latent Semantic Analysis

- While in the previous approaches contexts served only to get word vectors and were thrown away afterward, here we are also interested in context, or, in this case, document vectors.

Context:

- document d (from a collection D)

Matrix element:

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

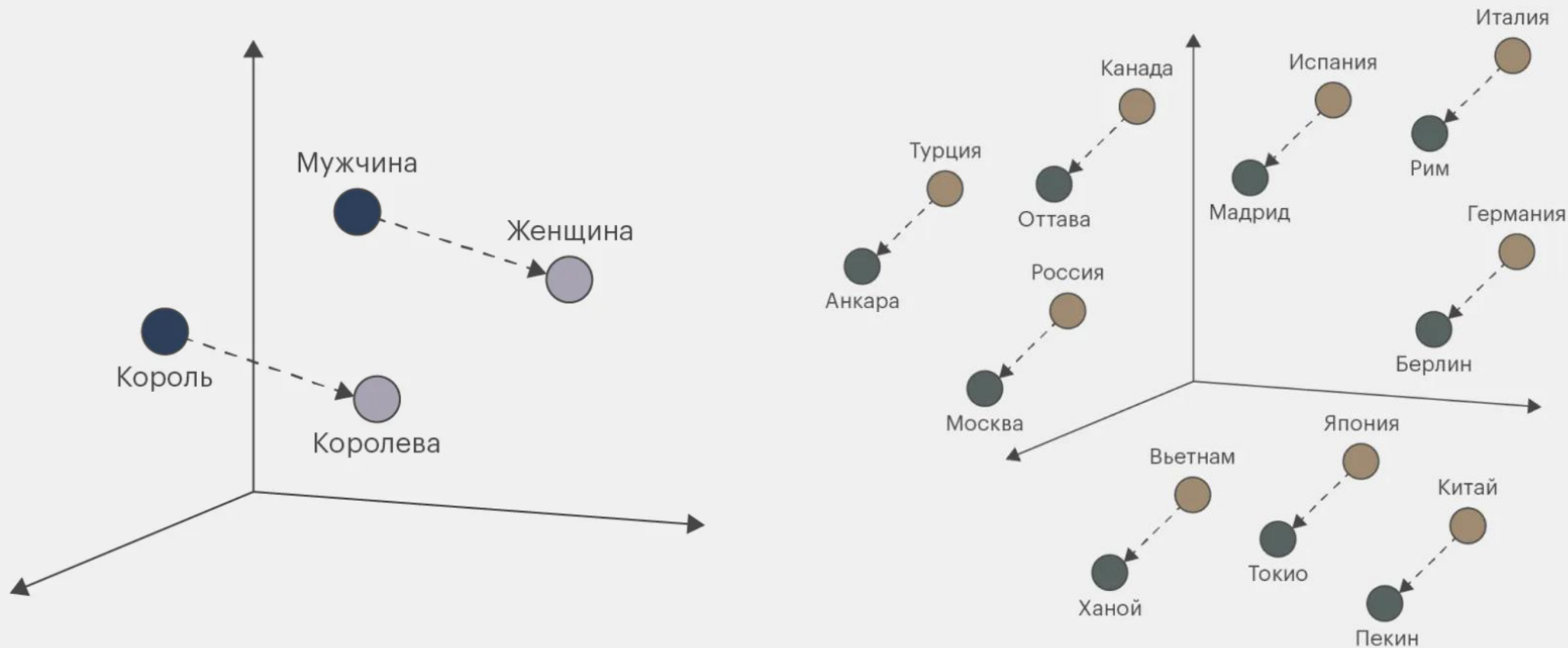
$N(\mathbf{w}, d)$

term frequency

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

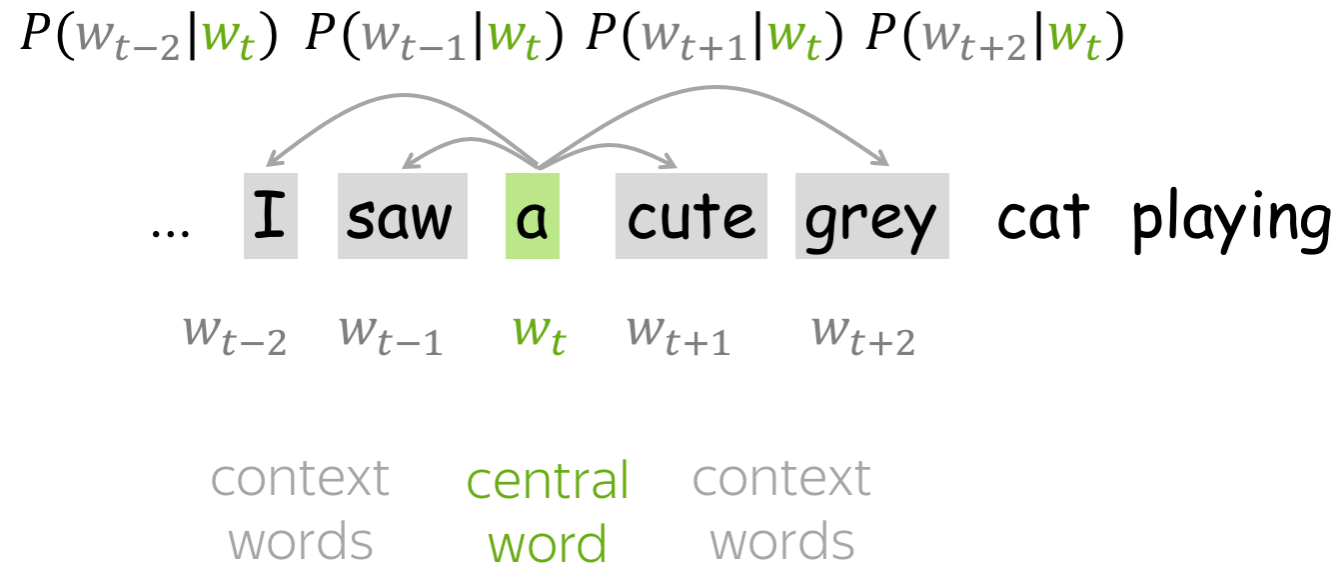
inverse document frequency

Word2Vec

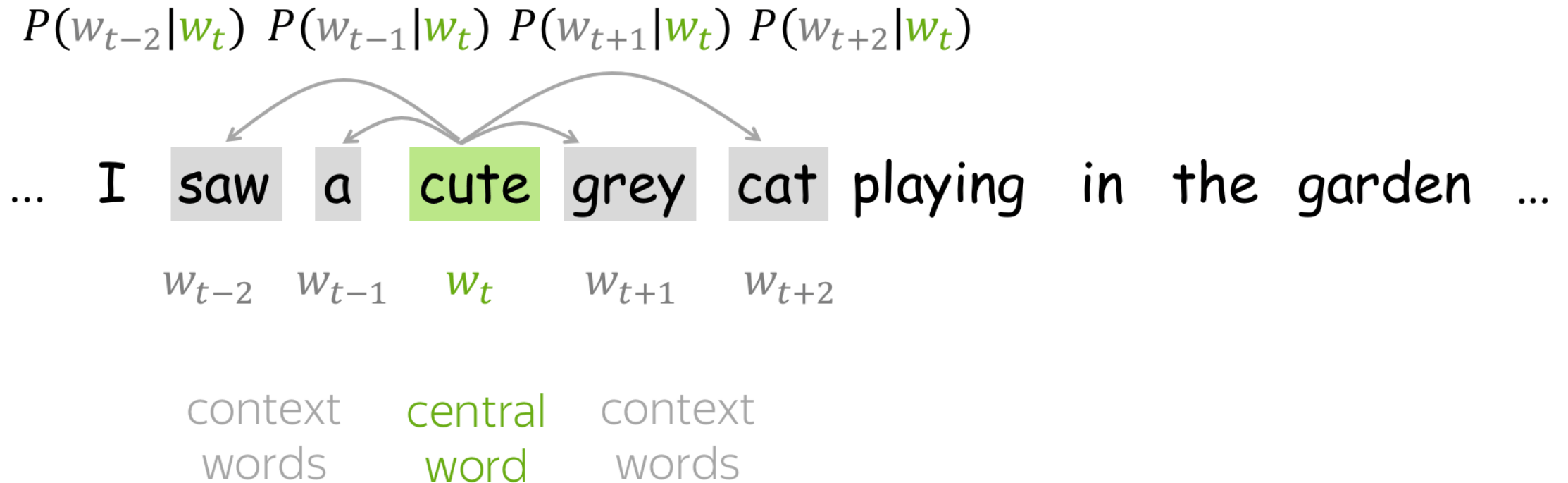


Word2Vec

- **Goal:** Learn word vectors by teaching them to predict contexts
 - take a huge text corpus
 - go over the text with a sliding window, moving one word at a time.
 - for the central word, compute probabilities of context words;
 - adjust the vectors to increase these probabilities.



Word2Vec



Objective Function: Negative Log-Likelihood

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:

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

How to compute this?

agrees with our plan above



go over text



with a sliding window



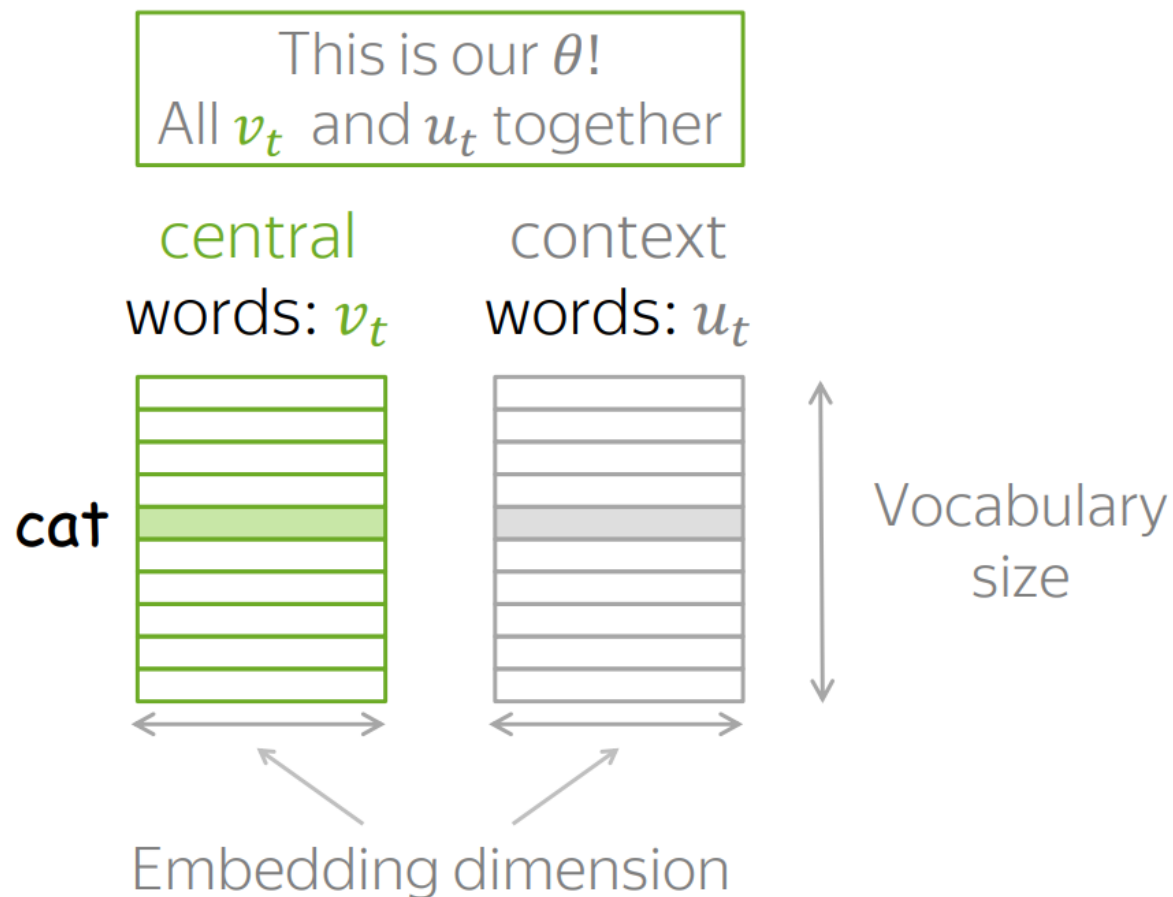
compute probability of the context word given the central

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):

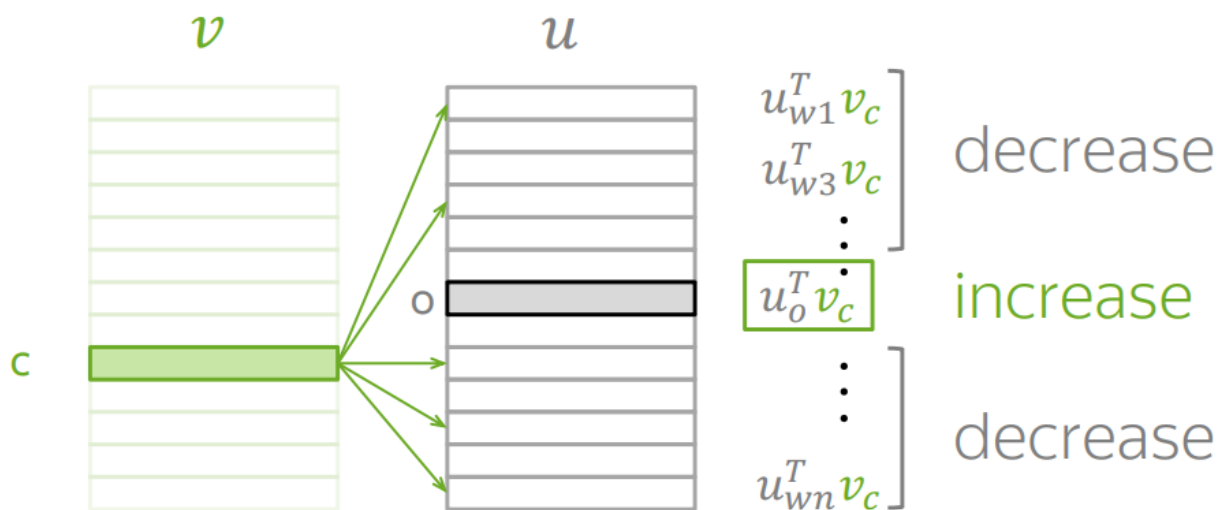
$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

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

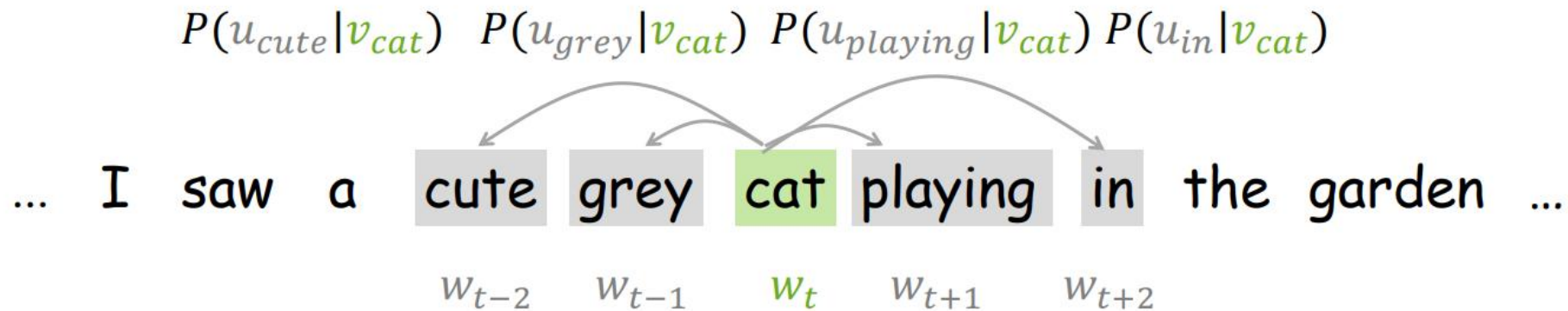
Normalize over entire vocabulary
to get probability distribution

Let us recall our plan:

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



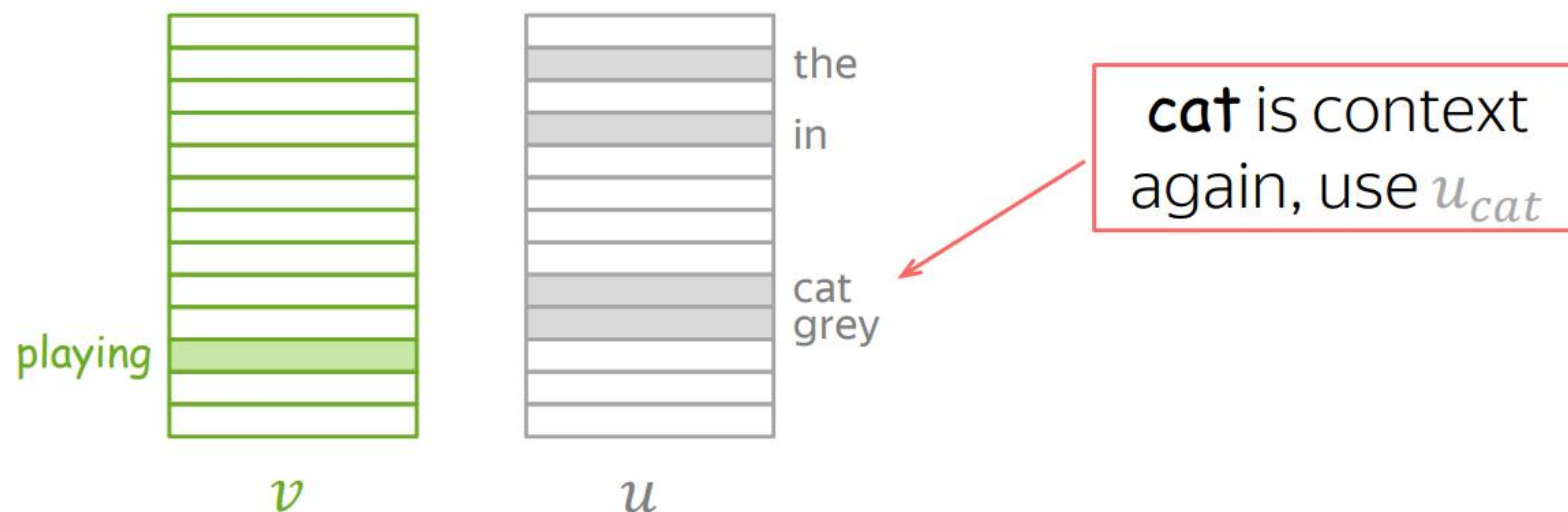
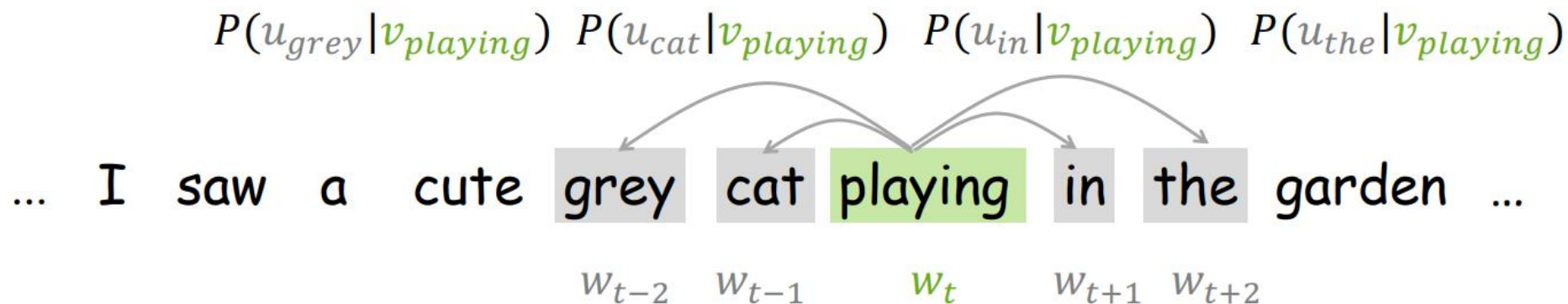
Two vectors for each word



cat is a central word, use v_{cat}



Two vectors for each word



One training step in detail

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

Loss for word j in window 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}

One training step in detail

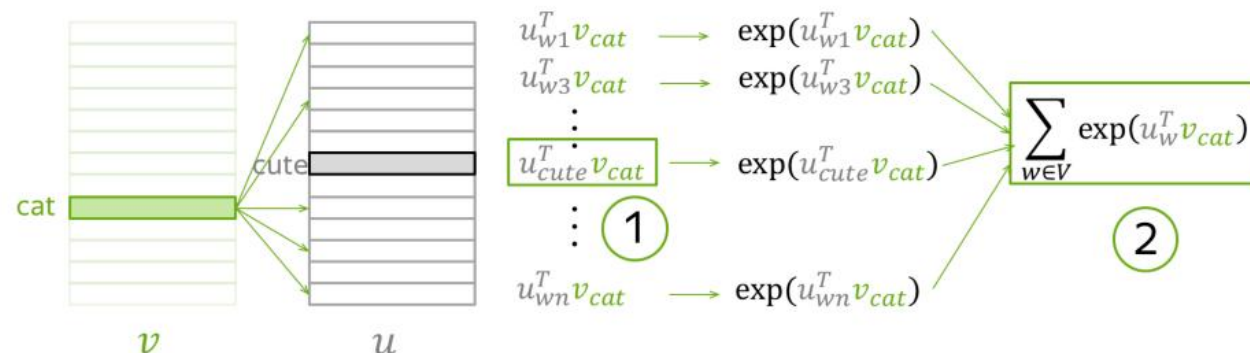
$$-\log P(\text{cute}|\text{cat})$$

$$= -u_{\text{cute}}^T v_{\text{cat}} + \log \sum_{w \in V} \exp(u_w^T v_{\text{cat}})$$

1. Take dot product of v_{cat} with all u

2. exp

3. sum all



4. get loss (for this one step)

$$J_{t,j}(\theta) = \underbrace{-u_{\text{cute}}^T v_{\text{cat}}}_{\text{1}} + \log \underbrace{\sum_{w \in V} \exp(u_w^T v_{\text{cat}})}_{\text{2}}$$

5. evaluate the gradient, make an update

$$v_{\text{cat}} := v_{\text{cat}} - \alpha \frac{\partial J_{t,j}(\theta)}{\partial v_{\text{cat}}}$$

$$u_w := u_w - \alpha \frac{\partial J_{t,j}(\theta)}{\partial u_w} \quad \forall w \in V$$

Let's do better: Negative Sampling

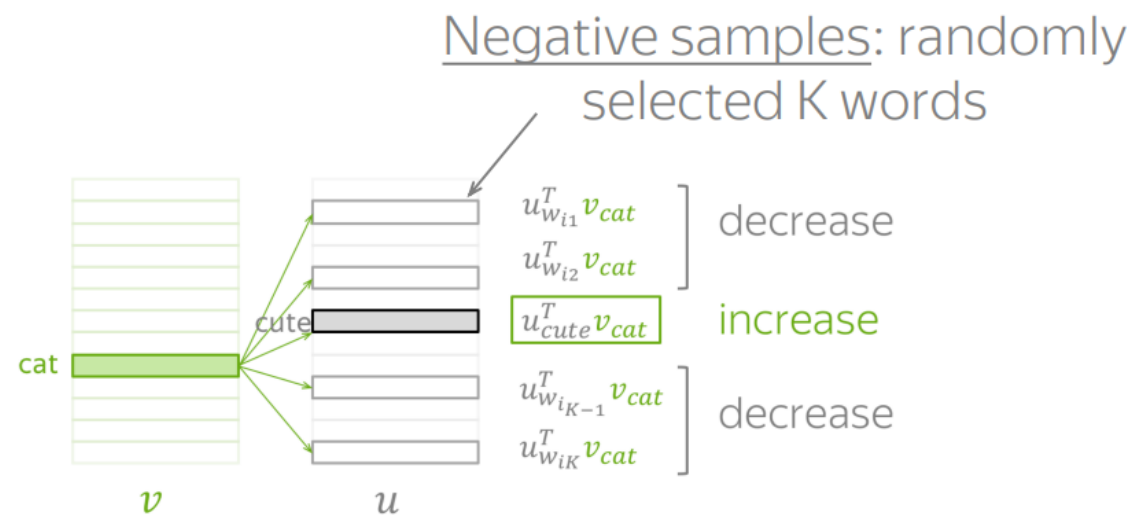
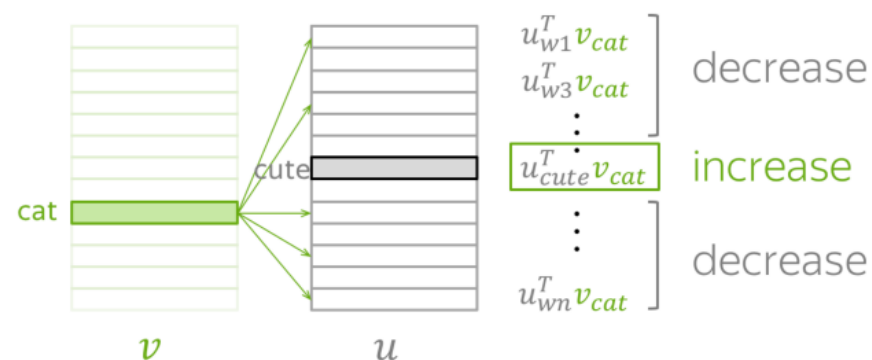
Dot product of v_{cat} :

- 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

Parameters to be updated:

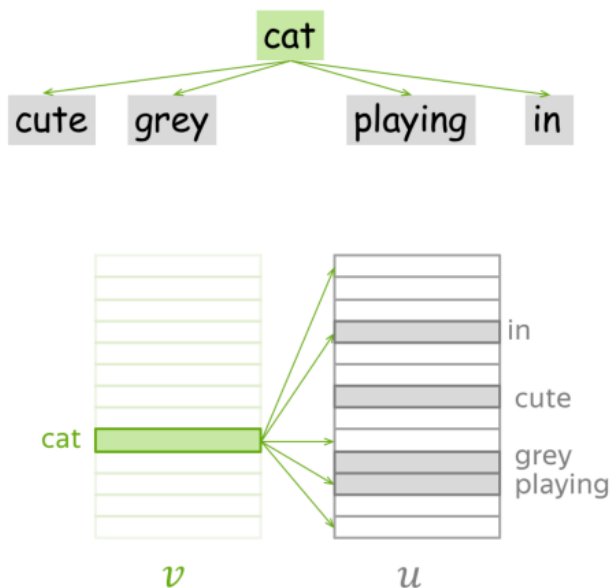
good

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

Word2Vec Variants: Skip-Gram and CBOW

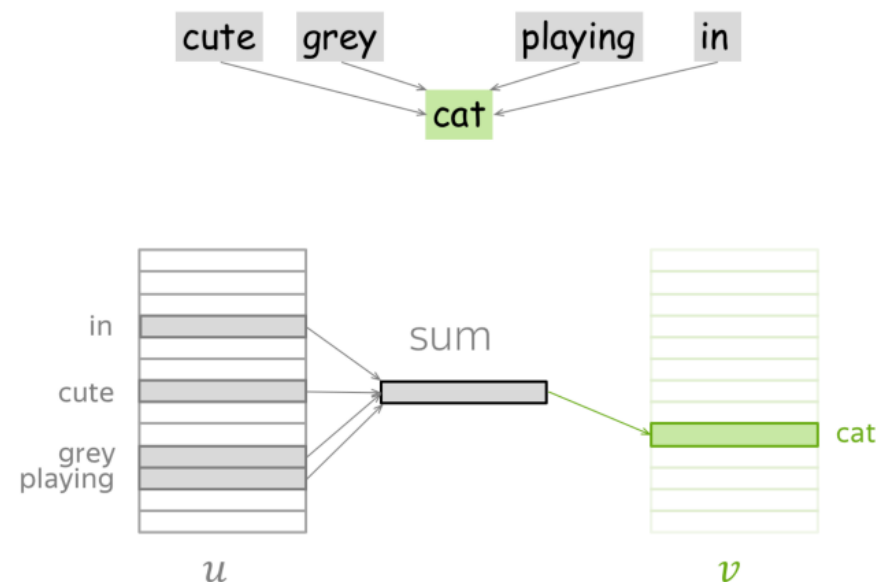
... I saw a cute grey cat playing in the garden ...

Skip-Gram: from **central** predict context
(one at a time)



(this is what we did so far)

CBOW: from sum of context predict **central**



(Continuous Bag of Words)

To read:

- Distributed Representations of Words and Phrases and their Compositionality. [Paper](#)
- Efficient Estimation of Word Representations in Vector Space. [Paper](#)
- word2vec Parameter Learning Explained. [Paper](#)
- [Habr](#): Russian translation of Jay Alammar blogpost

Relation to PMI Matrix Factorization



Related Papers

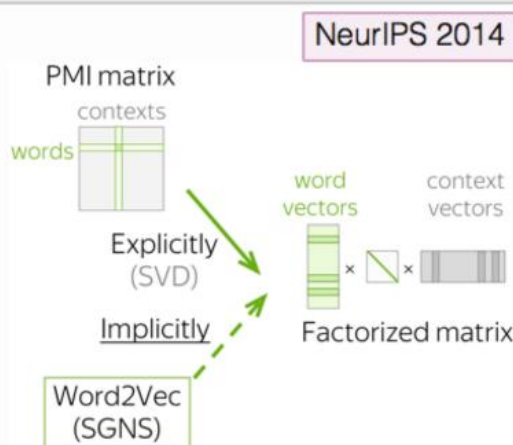
Good Old Classics

Neural Word Embedding as Implicit Matrix Factorization

Omer Levy, Yoav Goldberg

Theoretically, Word2Vec is not so different from matrix factorization approaches! Skip-gram with negative-sampling (SGNS) implicitly factorizes the shifted pointwise mutual information (PMI) matrix: $PMI(w, c) - \log k$, where k is the number of negative examples in negative sampling.

► More details



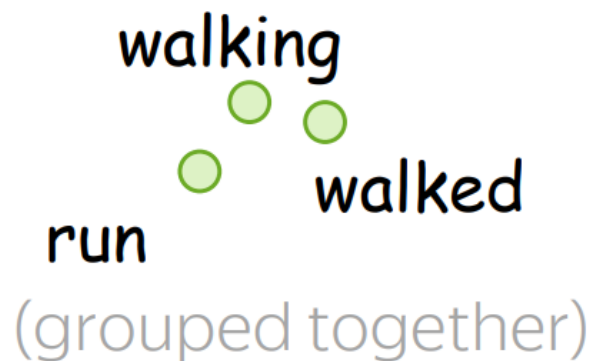
The Effect of Context Window Size

- **Larger windows** – more topical similarities



dog
bark leash
(grouped together)

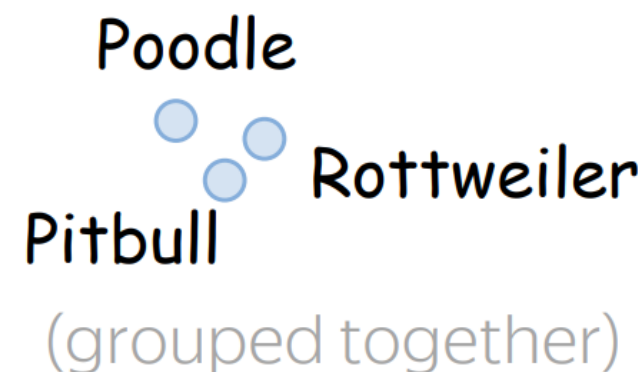
This diagram illustrates topical similarities using a large context window. It features four green circular markers. One marker is positioned near the word 'dog' at the top. Below it, two markers are placed near 'bark' and 'leash', which are positioned close to each other. A third marker is located further down and to the right, near the word 'run'. The words 'bark' and 'leash' are grouped together in a light gray font below the markers.



walking
run walked
(grouped together)

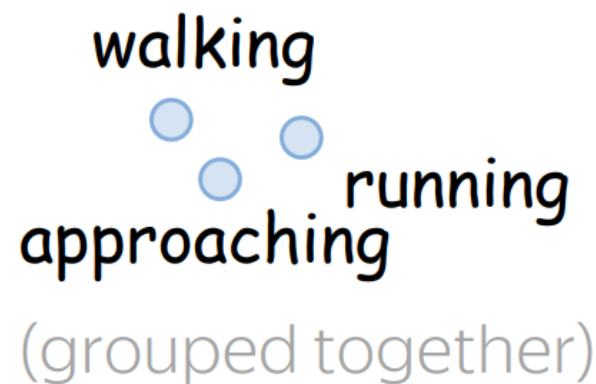
This diagram illustrates topical similarities using a large context window. It features four green circular markers. One marker is near 'walking' at the top. Below it, two markers are near 'run' and 'walked', which are positioned close to each other. A third marker is located further down and to the right, near the word 'leash'. The words 'run' and 'walked' are grouped together in a light gray font below the markers.

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



Poodle
Pitbull Rottweiler
(grouped together)

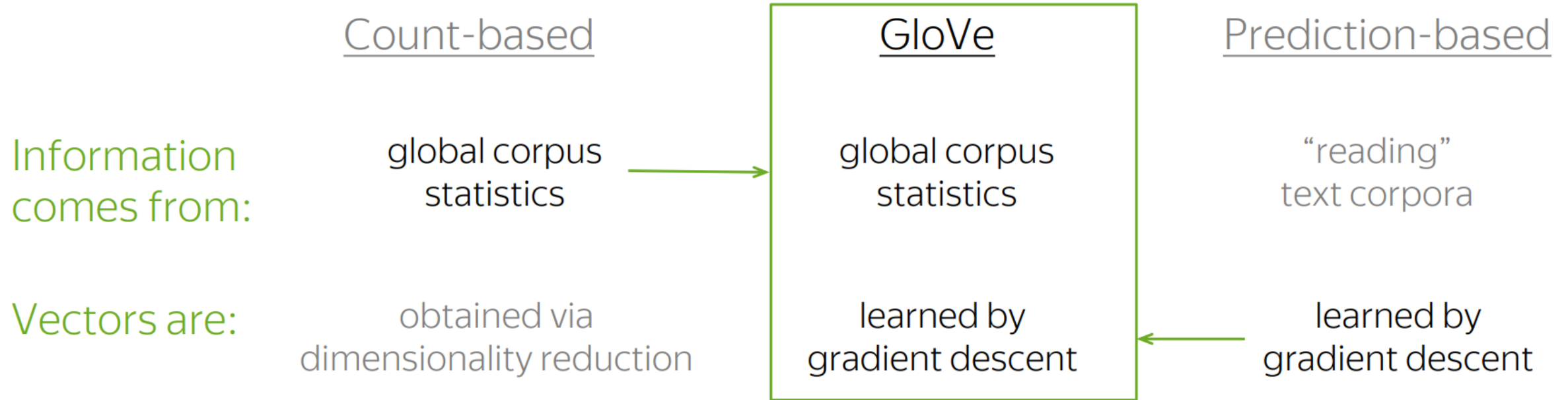
This diagram illustrates functional and syntactic similarities using a small context window. It features four blue circular markers. Three markers are clustered together near the words 'Poodle', 'Pitbull', and 'Rottweiler'. A fourth marker is located further down and to the right, near the word 'leash'. The words 'Pitbull' and 'Rottweiler' are grouped together in a light gray font below the markers.



walking
approaching running
(grouped together)

This diagram illustrates functional and syntactic similarities using a small context window. It features four blue circular markers. Three markers are clustered together near the words 'walking', 'approaching', and 'running'. A fourth marker is located further down and to the right, near the word 'leash'. The words 'approaching' and 'running' are grouped together in a light gray font below the markers.

GloVe: Global Vectors for Word Representation



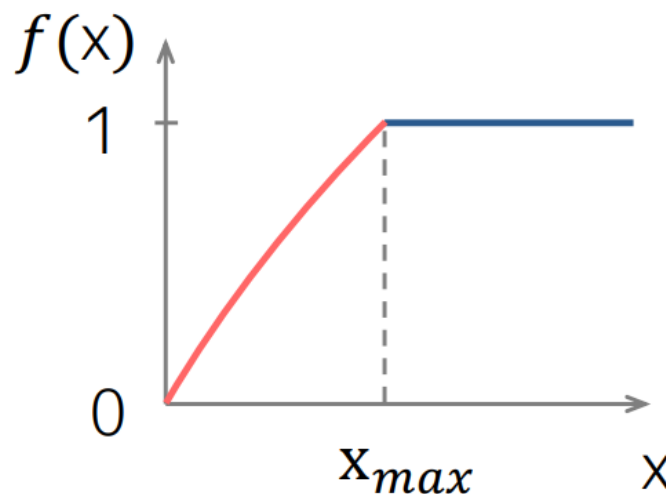
GloVe: Global Vectors for Word Representation

context vector word vector bias terms (also learned)

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

Weighting function to:

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



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

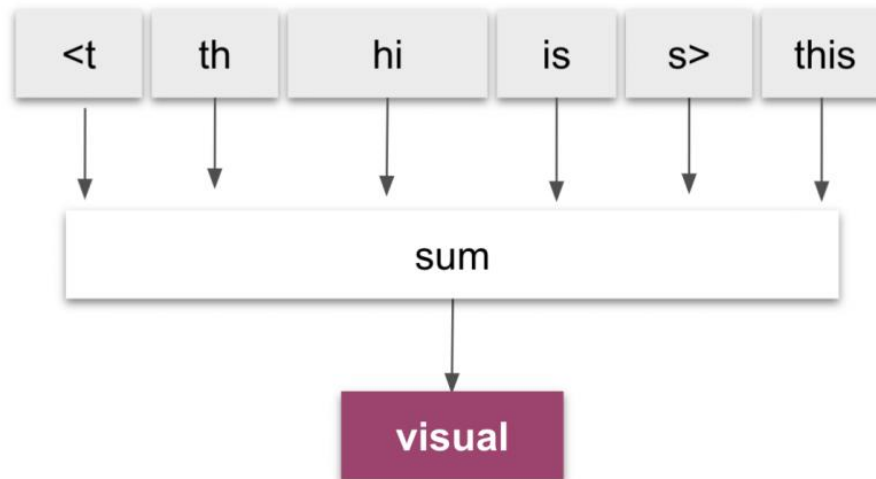
$$\alpha = 0.75, x_{max} = 100$$

FastText

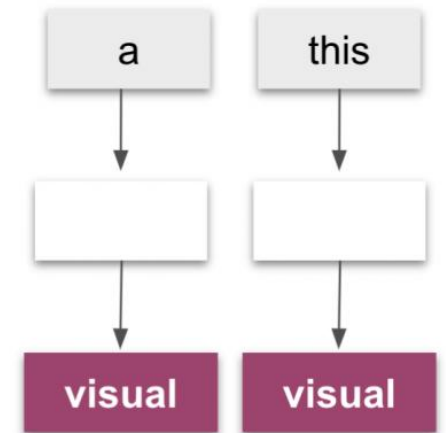
Let's analyse word as a set of symbolic n-grams

- Effective work with rare words and OOV
- Consideration of morphology and internal structure of words

fastText

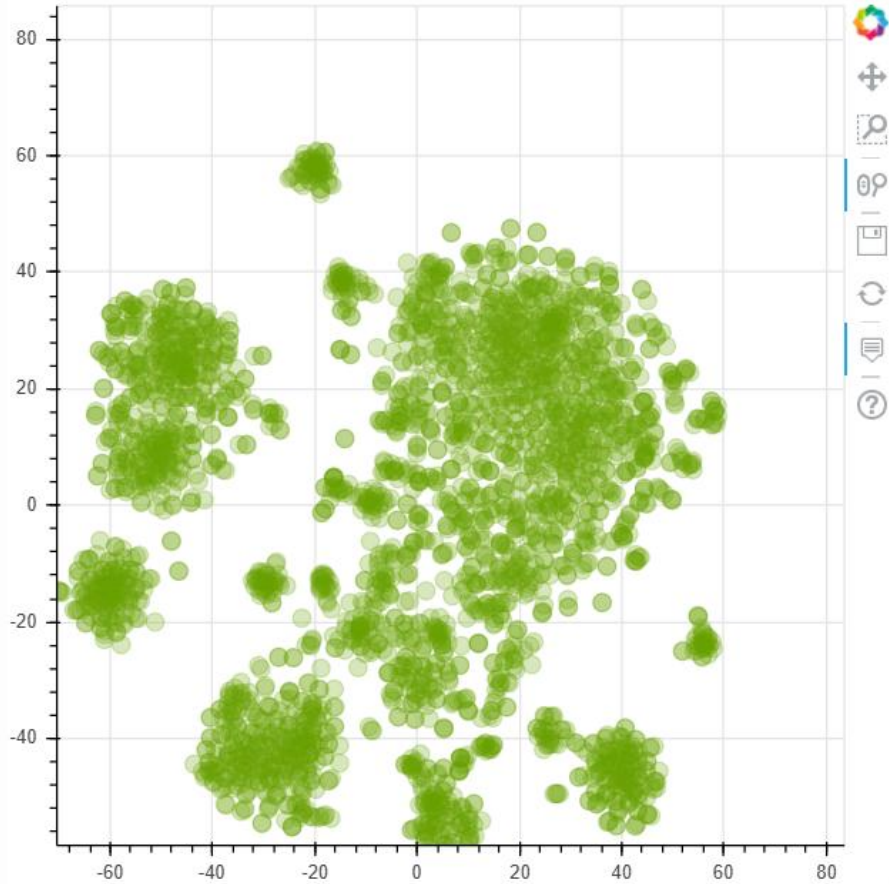


Word2Vec

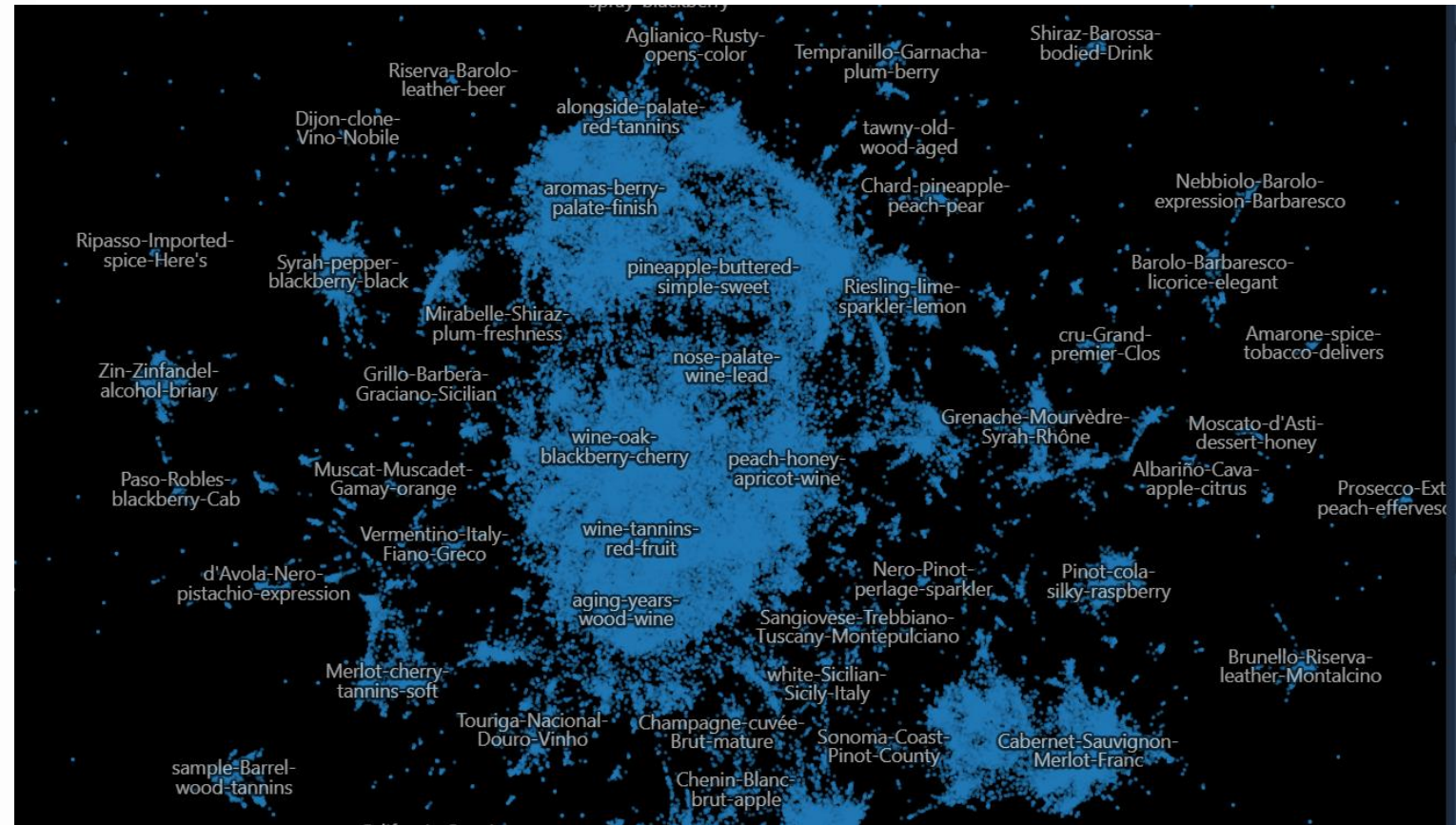


That's all! Thank you for your attention!

Visualization1



Visualization2



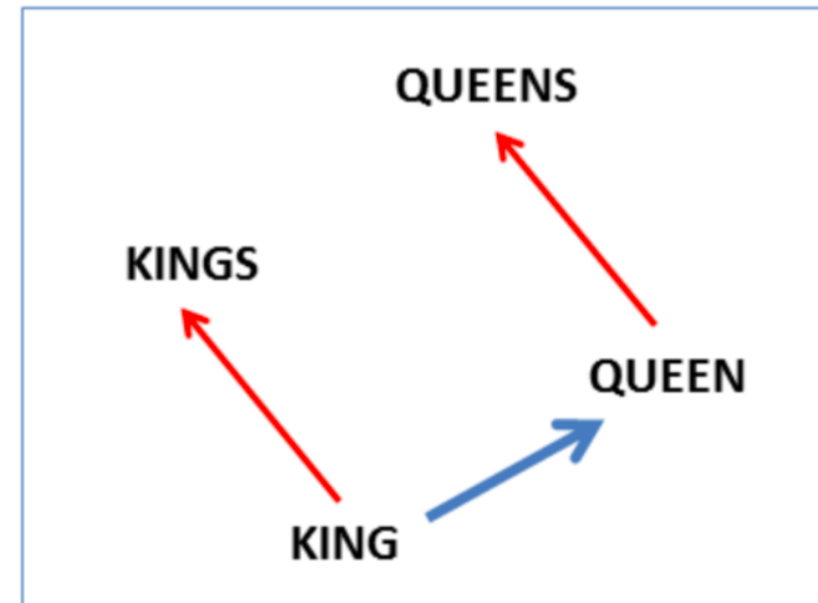
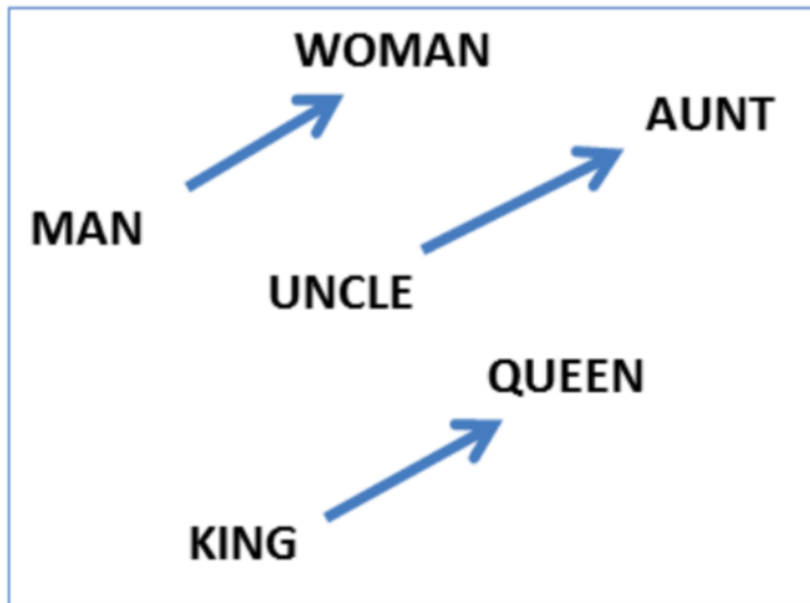
Word Similarity Benchmarks

<u>word pair</u>		<u>score</u>
vulgarism	profanity	9.62
subdividing	separate	8.67
friendships	brotherhood	7.5
exceedance	probability	5.0
assigned	allow	3.5
marginalize	interact	2.5
misleading	beat	1.25
radiators	beginning	0

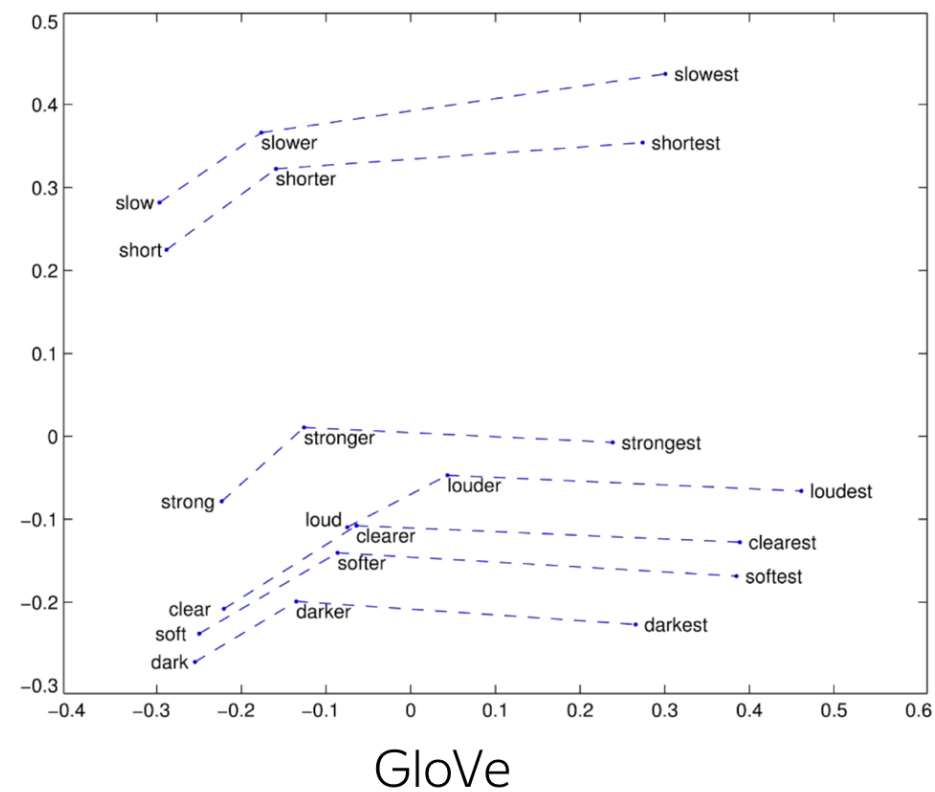
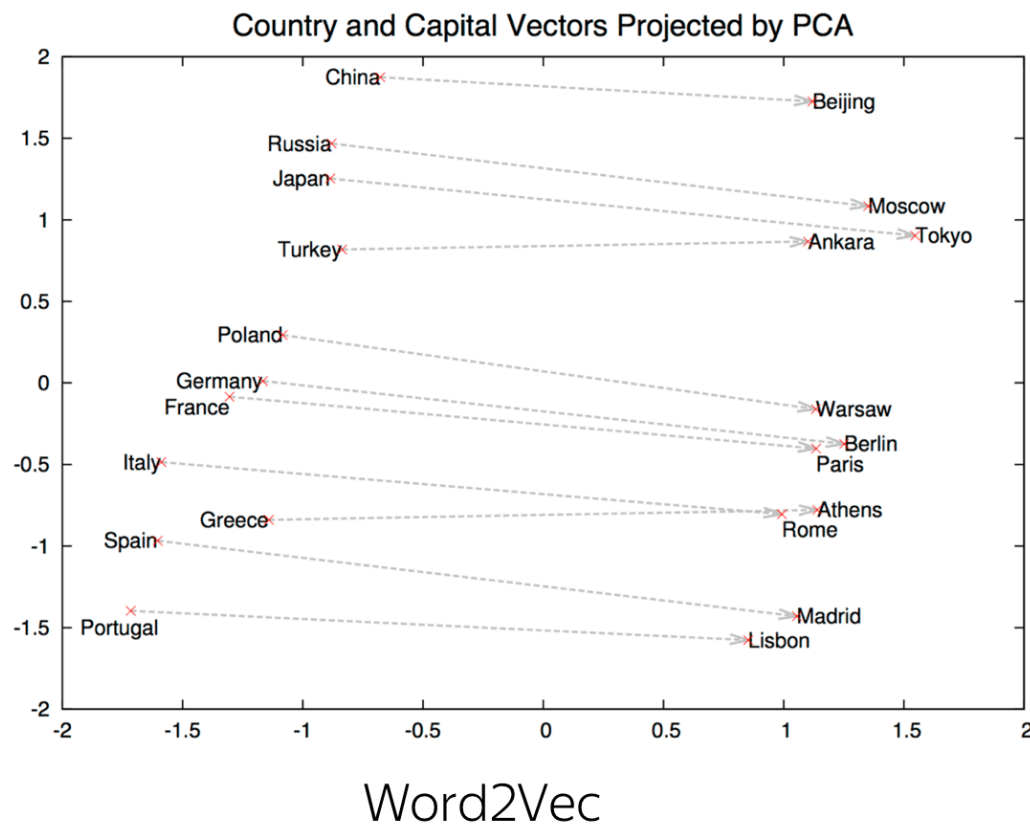
Linear Structure

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





Similarities Across Languages

The recipe for building large dictionaries from small ones

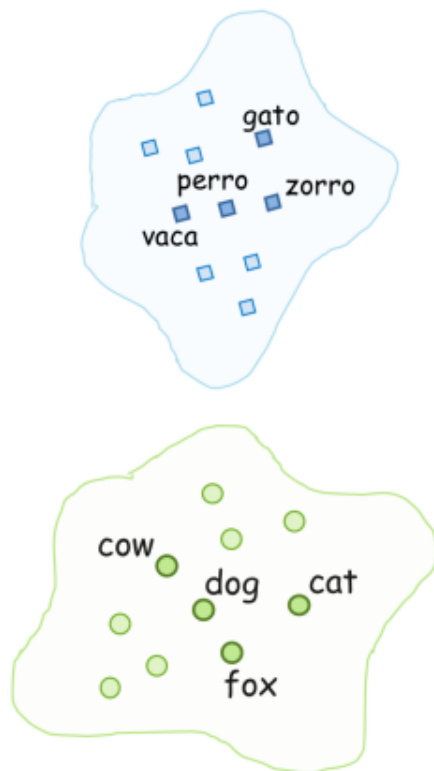
Ingredients:

- corpus in one language (e.g., **English**)
- corpus in another language (e.g., **Spanish**)
- very small dictionary

cat ↔ gato
cow ↔ vaca
dog ↔ perro
fox ↔ zorro
...

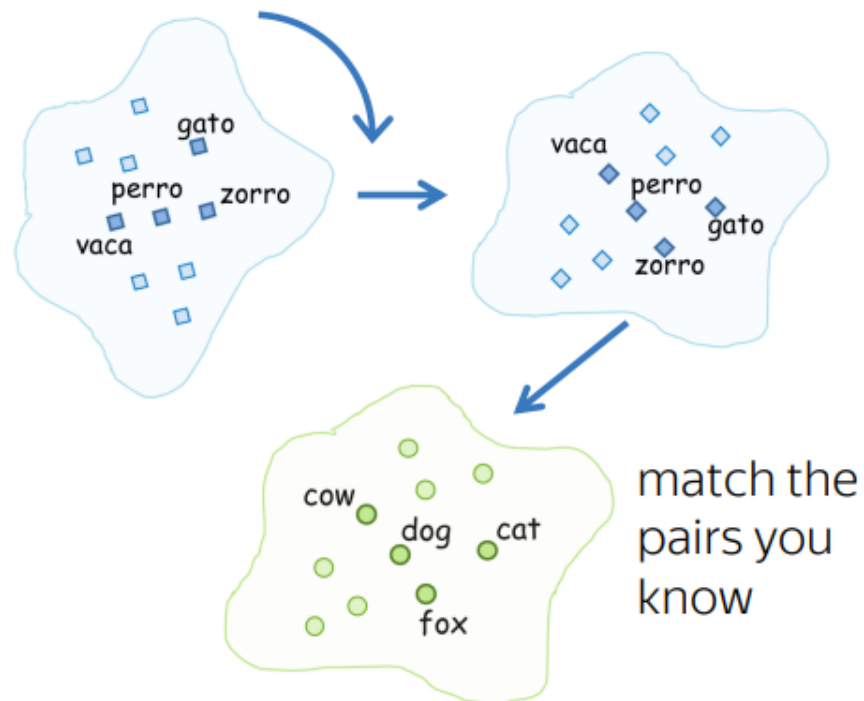
Step 1:

- train embeddings for each language



Step 2:

- linearly map one embeddings to the other to match words from the dictionary





Similarities Across Languages

The recipe for building large dictionaries from small ones

Ingredients:

- corpus in one language (e.g., **English**)
- corpus in another language (e.g., **Spanish**)
- very small dictionary

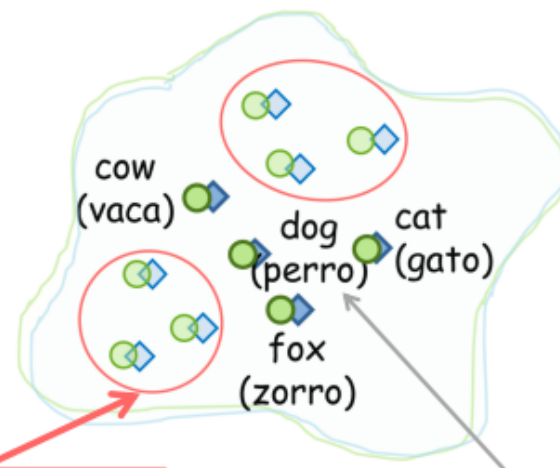
cat ↔ gato
cow ↔ vaca
dog ↔ perro
fox ↔ zorro
...

Steps 1-2:

- match words from the vocabulary

Step 3:

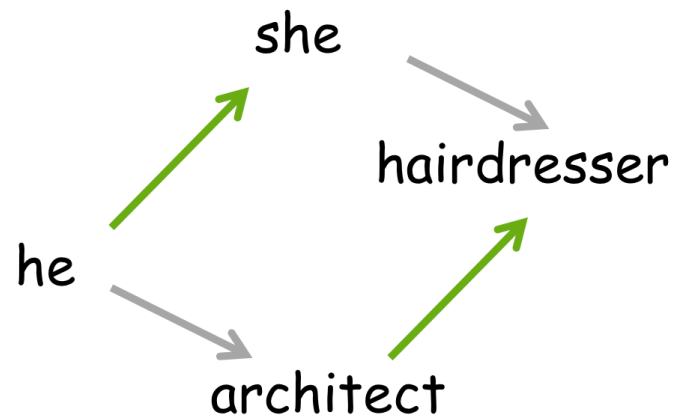
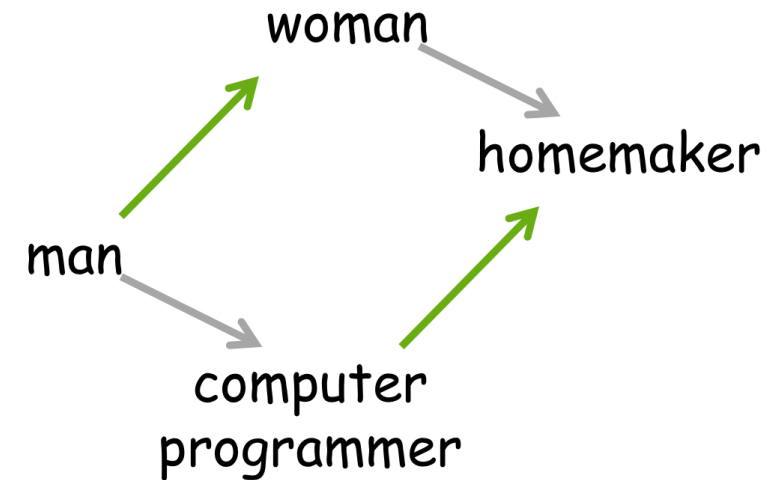
- after matching the two spaces, get new pairs from the new matches



new translations –
you learned them!

old translations –
the ones you knew

Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings



Gender stereotype **she-he** analogies

sewing-carpentry
nurse-surgeon
blond-burly
giggle-chuckle
sassy-snappy
volleyball-football

registered nurse-physician
interior designer-architect
feminism-conservatism
vocalist-guitarist
diva-superstar
cupcakes-pizzas

housewife-shopkeeper
softball-baseball
cosmetics-pharmaceuticals
petite-lanky
charming-affable
lovely-brilliant

Gender appropriate **she-he** analogies

queen-king
waitress-waiter

ovarian cancer-prostate cancer
convent-monastery

sister-brother
mother-father