



Word Embeddings

Mikhail Yurushkin
BroutonLab

October 1, 2020

Vector representations

Vector representation

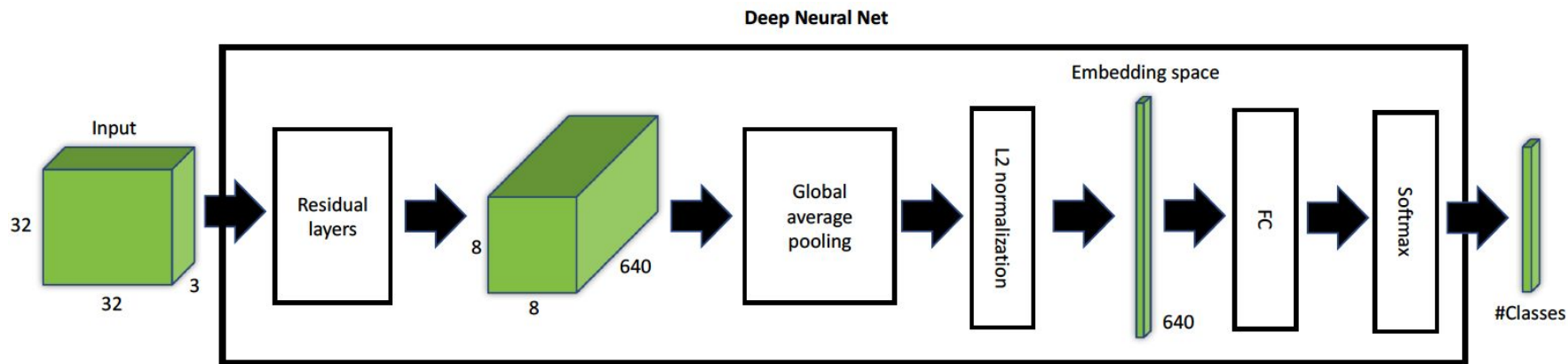
= embedding

= vector of fixed size

- typically gained in unsupervised way
- used in transfer learning
- speedup training of neural networks

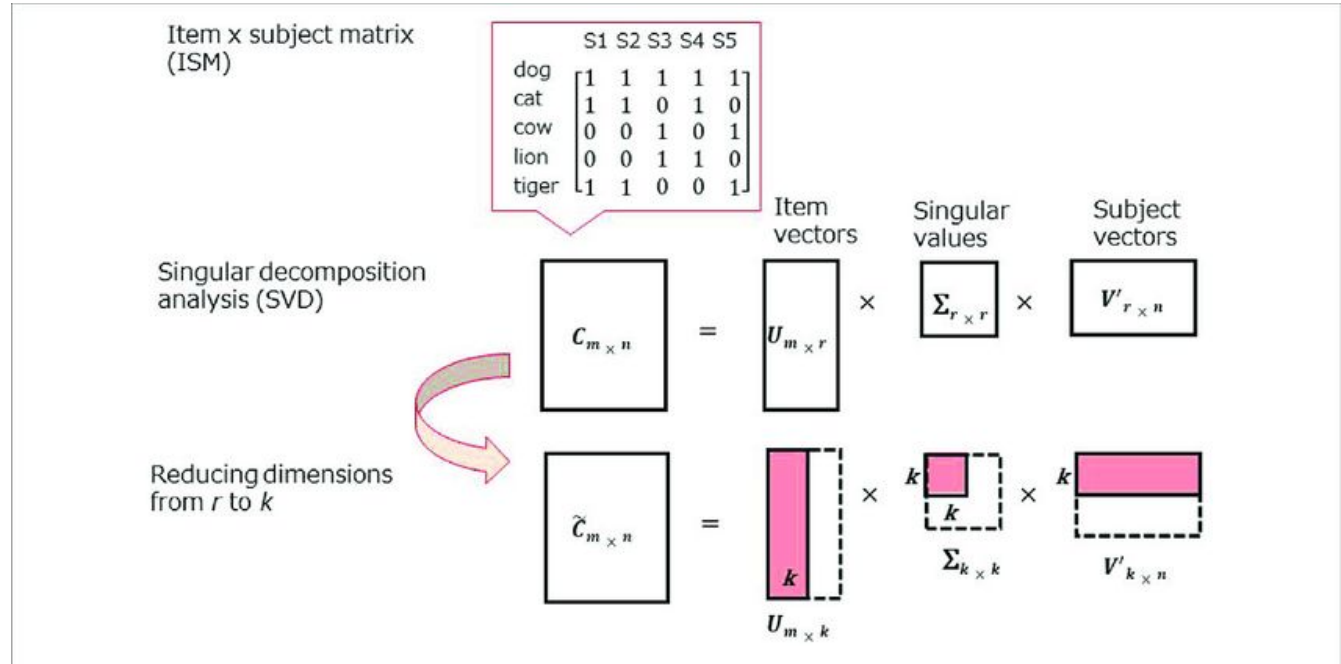
Classical example: one-hot encoding

Usage in Computer Vision



Classical approach: SVD

Not scalable :(



Word2Vec

Paper: **Efficient Estimation of Word Representations in Vector Space**. 2013, Jeffrey Dean et al.

- set of methods with the same idea
- C-BOW, Skip-Gram

$X = \text{vector}(\text{"biggest"}) - \text{vector}(\text{"big"}) + \text{vector}(\text{"small"})$

Result: **smallest**

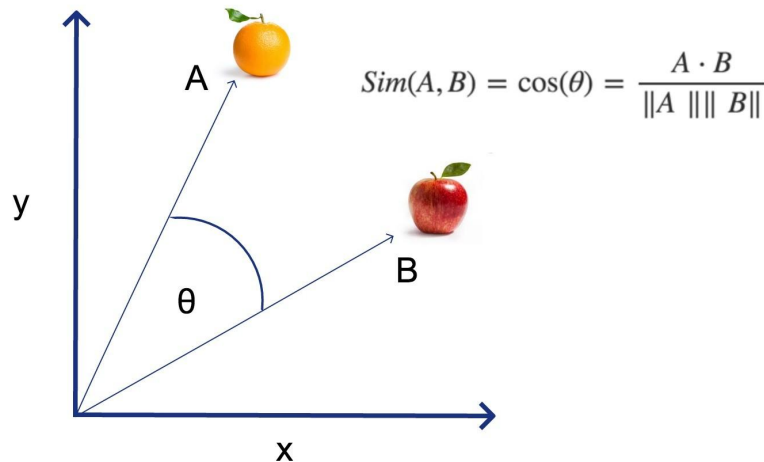
Measuring similarity of vectors

Vectors can be compared in different ways.

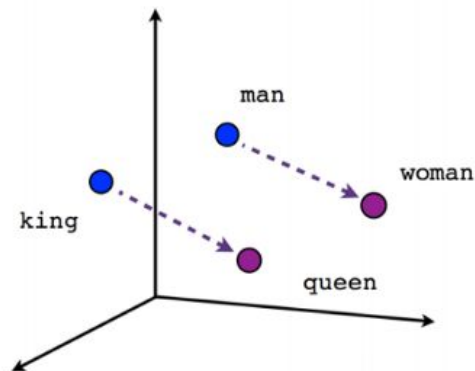
Most popular ways:

1. L2 distance
2. Cosine similarity

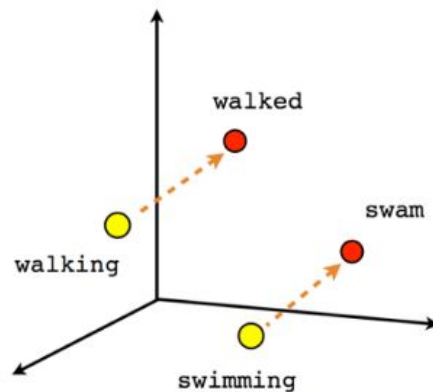
Cosine Similarity



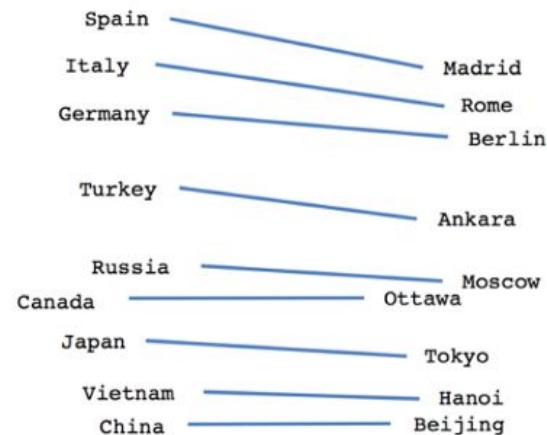
Semantic similarity



Male-Female



Verb tense



Country-Capital

Other semantic similarity examples

Type of relationship	Word Pair 1		Word Pair 2	
Common capital city	Athens	Greece	Oslo	Norway
All capital cities	Astana	Kazakhstan	Harare	Zimbabwe
Currency	Angola	kwanza	Iran	rial
City-in-state	Chicago	Illinois	Stockton	California
Man-Woman	brother	sister	grandson	granddaughter
Adjective to adverb	apparent	apparently	rapid	rapidly
Opposite	possibly	impossibly	ethical	unethical
Comparative	great	greater	tough	tougher
Superlative	easy	easiest	lucky	luckiest
Present Participle	think	thinking	read	reading
Nationality adjective	Switzerland	Swiss	Cambodia	Cambodian
Past tense	walking	walked	swimming	swam
Plural nouns	mouse	mice	dollar	dollars
Plural verbs	work	works	speak	speaks

Applications of word embeddings

- Search of typos or synonyms in search queries
- Can be used as a features in most of NLP tasks:
 - Named entity recognition
 - Word tagging
 - End2end machine translation
 - Sentiment analysis
 - Document classification
 - etc.

C-BOW:

- Idea: Using context words, we can predict center word

i.e. Probability("It is (?) to finish" \rightarrow **"time"**)
context words (window_size=2)

- Present word as distributed vector of probability \rightarrow Low dimension
- Goal: Train weight-matrix(W) satisfies below

$$\operatorname{argmax}_W \{ \text{Minimize}(|\mathbf{time} - \text{softmax}(pr(\mathbf{time} | \mathbf{it, is, to, finish}))|; W) \}$$

* Softmax(): K-dim vector of $x \in \mathbb{R} \rightarrow$ K-dim vector that has $(0,1) \in \mathbb{R}$

- Loss-function (using cross-entropy method)

$$E = -\log p(w_t | w_{t-c} \dots w_{t+c})$$

C-BOW:

$$\mathbf{h} = \mathbf{W}^T \mathbf{x}$$

$$\mathbf{u} = \mathbf{W}'^T \mathbf{h} = \mathbf{W}'^T \mathbf{W}^T \mathbf{x}$$

$$\mathbf{y} = \text{Softmax}(\mathbf{u}) = \text{Softmax}(\mathbf{W}'^T \mathbf{W}^T \mathbf{x})$$

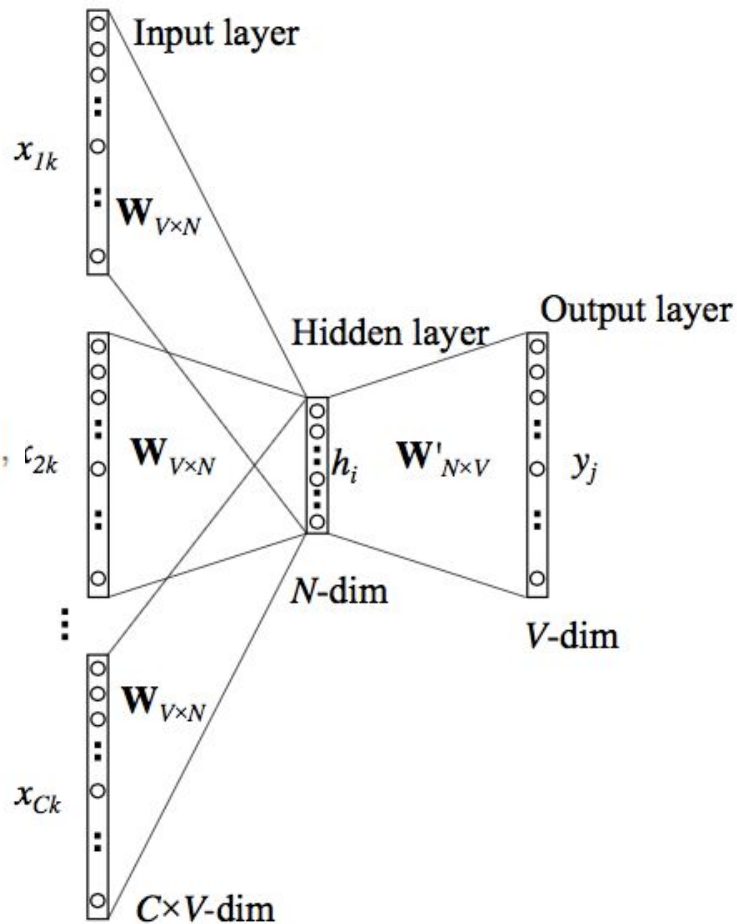
$$\mathcal{L} = -\log \mathbb{P}(w_t | w_c) = -\log y_{j^*} = -\log[\text{Softmax}(u_{j^*})] = -\log\left(\frac{\exp u_{j^*}}{\sum_i \exp u_i}\right), \ell_{2k}$$

$$\mathcal{L} = -u_{j^*} + \log \sum_i \exp(u_i).$$

$$\mathbf{h} = \frac{1}{C} \mathbf{W}^T \sum_{c=1}^C \mathbf{x}^{(c)} = \mathbf{W}^T \bar{\mathbf{x}}$$

$$\mathbf{u} = \mathbf{W}'^T \mathbf{h} = \frac{1}{C} \sum_{c=1}^C \mathbf{W}'^T \mathbf{W}^T \mathbf{x}^{(c)} = \mathbf{W}'^T \mathbf{W}^T \bar{\mathbf{x}}$$

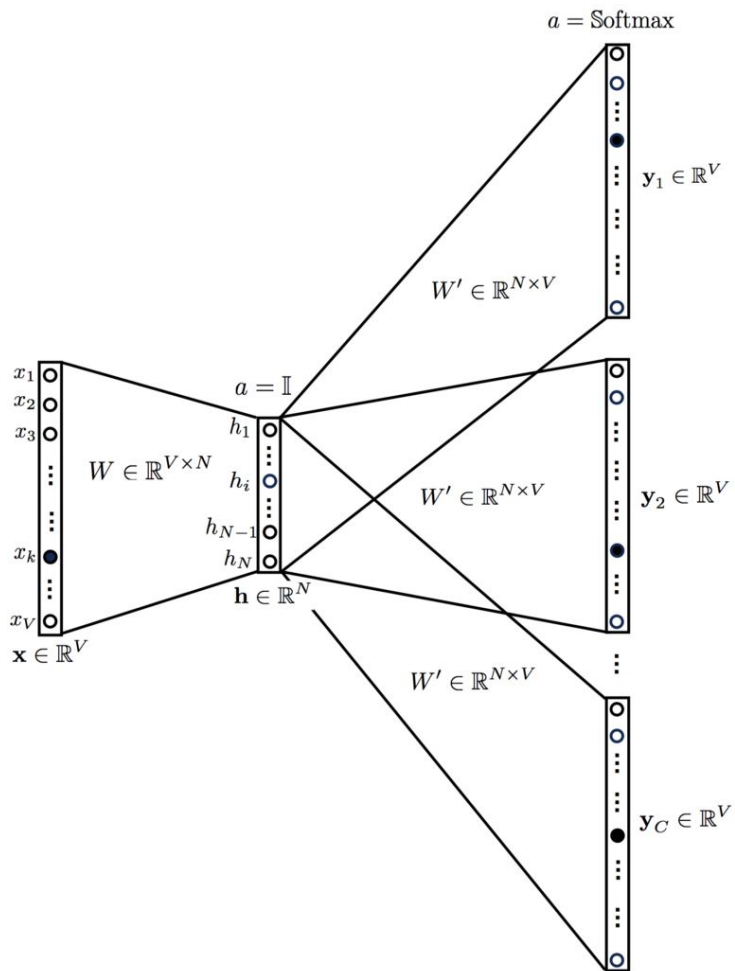
$$\mathbf{y} = \text{Softmax}(\mathbf{u}) = \text{Softmax}(\mathbf{W}'^T \mathbf{W}^T \bar{\mathbf{x}})$$



Backpropagation [details](#)

Skip-Gram:

Inversed C-BOW

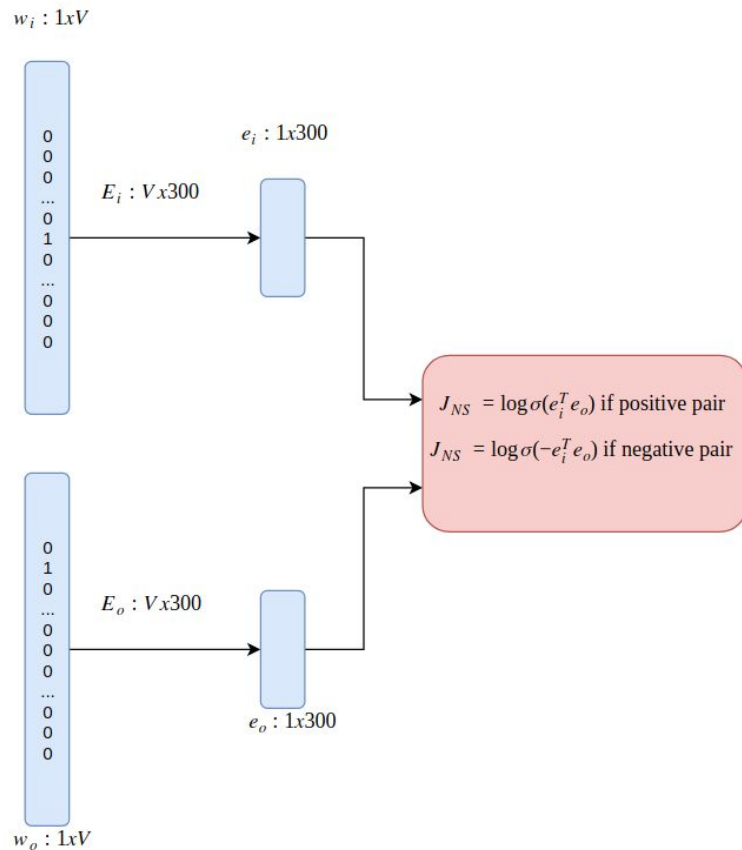


Negative sampling

- Difficult to calculate big softmax?
- let's solve binary classification problem!

$$\mathcal{L} = \sum_{(w,s) \in D_1} \log \sigma(v_w^T v_s) + \sum_{(w,s) \in D_2} \log \sigma(-v_w^T v_s),$$

$$D_1 = \{(w, s) : s \in c(w)\}, D_2 = \{(w, s) : s \notin c(w)\}$$



Negative sampling

■ : Center Word
■ : Context Word

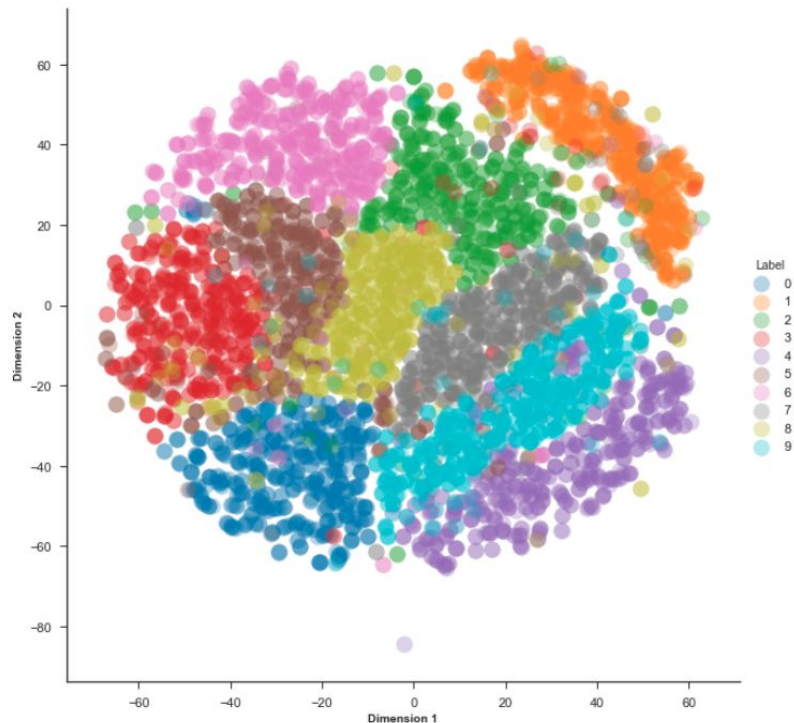
c=0 The cute **cat** jumps over the lazy dog.

c=1 The **cute** **cat** **jumps** over the lazy dog.

c=2 **The** **cute** **cat** **jumps** **over** the lazy dog.

Source Text	Training Samples
<div>The quick brown fox jumps over the lazy dog. →</div>	(the, quick) (the, brown)
<div>The quick brown fox jumps over the lazy dog. →</div>	(quick, the) (quick, brown) (quick, fox)
<div>The quick brown fox jumps over the lazy dog. →</div>	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
<div>The quick brown fox jumps over the lazy dog. →</div>	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

How to check your embeddings are good?



2-D t-SNE

- Unsupervised
- clusterization
- non-linear technique primarily used for data exploration and visualizing high-dimensional data.