

Neural Word Embeddings as Matrix Factorization

Master's Thesis Mathematics

Presented by:
Theresa Klumpp

Supervisors:
Prof. P. Pfaffelhuber
Prof. H. Bast

January 15, 2020

Table of Contents

1 Problem

2 Solution

3 Evaluation

Problem

Goal: word vectors that reflect similarities and dissimilarities

Problem

Goal: word vectors that reflect similarities and dissimilarities

Underlying hypothesis: words in similar contexts have similar meanings

Problem

Goal: word vectors that reflect similarities and dissimilarities

Underlying hypothesis: words in similar contexts have similar meanings

- I get to work faster when I take the ***.

Problem

Goal: word vectors that reflect similarities and dissimilarities

Underlying hypothesis: words in similar contexts have similar meanings

- I get to work faster when I take the ***.
- This model has amazing acceleration for a *** of its size.

Problem

Goal: word vectors that reflect similarities and dissimilarities

Underlying hypothesis: words in similar contexts have similar meanings

- I get to work faster when I take the ***.
- This model has amazing acceleration for a *** of its size.
- I would never drive my *** into Paris if I could get there by train.

Problem

Goal: word vectors that reflect similarities and dissimilarities

Underlying hypothesis: words in similar contexts have similar meanings

- I get to work faster when I take the ***.
- This model has amazing acceleration for a *** of its size.
- I would never drive my *** into Paris if I could get there by train.

Demo

Contributions

- Gaining an understanding of the objective functions of skip-gram (with and without negative sampling) and the statistical models behind them.
- Finding a maximum for skip-gram's objective.
- Showing the connection between the neural networks and Singular Value Decomposition (SVD).
- Comparing different metrics on the sphere.
- Finding a formula for the expectation of the distance of the closest vector.
- An implementation of the SGNS neural network and the SVD variant for both skip-gram and SGNS.
- Evaluation of the models on word similarity and analogy tasks.

Contributions

- Gaining an understanding of the objective functions of skip-gram (with and without negative sampling) and the statistical models behind them.
- Finding a maximum for skip-gram's objective.
- Showing the connection between the neural networks and Singular Value Decomposition (SVD).
- Comparing different metrics on the sphere.
- Finding a formula for the expectation of the distance of the closest vector.
- An implementation of the SGNS neural network and the SVD variant for both skip-gram and SGNS.
- Evaluation of the models on word similarity and analogy tasks.

Questions?

Table of Contents

1 Problem

2 Solution

3 Evaluation

Definition: Context

	Text		Samples
	I get to work faster when I take the car.	⇒	(I, get) (I, to)
	I get to work faster when I take the car.	⇒	(get, I) (get, to) (get, work)
	I get to work faster when I take the car.	⇒	(to, I) (to, get) (to, work) (to, faster)
	I get to work faster when I take the car.	⇒	(work, get) (work, to) (work, faster) (work, when)

...

Notation

- V_W and V_C : word and context vocabulary (we have $V_W = V_C$)
- D : observed word context pairs
- $\#(\mathbf{w}, \mathbf{c})$: number of times the pair (w, c) appears in D
- $\#(\mathbf{w}) = \sum_{c' \in V_C} \#(w, c')$ and $\#(\mathbf{c}) = \sum_{w' \in V_W} \#(w', c)$

Mathematical Goal

Find embeddings such that $\vec{w} \cdot \vec{c}$ is

- high for pairs with large $\#(w, c)$ and
- small for pairs with low $\#(w, c)$

Mathematical Goal

Find embeddings such that $\vec{w} \cdot \vec{c}$ is

- high for pairs with large $\#(w, c)$ and
- small for pairs with low $\#(w, c)$

Why does this yield good embeddings?

Mathematical Goal

Find embeddings such that $\vec{w} \cdot \vec{c}$ is

- high for pairs with large $\#(w, c)$ and
- small for pairs with low $\#(w, c)$

Why does this yield good embeddings?

	$c_1 = \text{drive}$	$c_2 = \text{road}$	$c_3 = \text{space}$	$c_4 = \text{bottle}$
$w_1 = \text{car}$	0.9	0.8	0.2	0.1
$w_2 = \text{truck}$	0.8	0.7	0.2	0.2

Mathematical Goal

Find embeddings such that $\vec{w} \cdot \vec{c}$ is

- high for pairs with large $\#(w, c)$ and
- small for pairs with low $\#(w, c)$

$$W = \begin{pmatrix} \vec{w}_1 \\ \vdots \\ \vec{w}_{|V_W|} \end{pmatrix} \text{ and } C = \begin{pmatrix} \vec{c}_1 \\ \vdots \\ \vec{c}_{|V_C|} \end{pmatrix}$$

Mathematical Goal

Find embeddings such that $\vec{w} \cdot \vec{c}$ is

- high for pairs with large $\#(w, c)$ and
- small for pairs with low $\#(w, c)$

$$W = \begin{pmatrix} \vec{w}_1 \\ \vdots \\ \vec{w}_{|V_W|} \end{pmatrix} \text{ and } C = \begin{pmatrix} \vec{c}_1 \\ \vdots \\ \vec{c}_{|V_C|} \end{pmatrix}$$

\Rightarrow Find a function $\ell(W, C)$ that is maximized when the properties above hold.

Skip-Gram: Objective functions

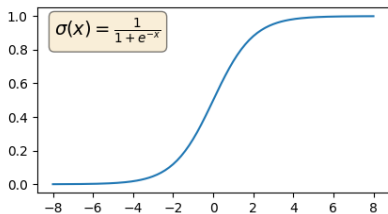
$$\ell_{SG}(W, C) = \sum_{(w, c) \in D} \left(\vec{w} \cdot \vec{c} - \log \left(\sum_{c' \in V_C} \exp(\vec{w} \cdot \vec{c}') \right) \right)$$

More

Skip-Gram: Objective functions

$$\ell_{SG}(W, C) = \sum_{(w, c) \in D} \left(\vec{w} \cdot \vec{c} - \log \left(\sum_{c' \in V_C} \exp(\vec{w} \cdot \vec{c}') \right) \right)$$

$$\ell_{SGNS}(W, C) = \sum_{(w, c) \in D} \left(\log \sigma(\vec{w} \cdot \vec{c}) + \sum_{j=1}^k \log \sigma(-\vec{w} \cdot \vec{c}_j) \right)$$



More

Optimal value for the dot products

- $\ell_{SGNS}(W, C)$ is maximized for

$$(\vec{w} \cdot \vec{c})^{\text{OPT}} = \log \left(\frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k$$

Optimal value for the dot products

- $\ell_{SGNS}(W, C)$ is maximized for

$$(\vec{w} \cdot \vec{c})^{\text{OPT}} = \log \left(\frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k$$

- Note that

$$(W \cdot C^T)_{ij} = \vec{w}_i \cdot \vec{c}_j$$

Optimal value for the dot products

- $\ell_{SGNS}(W, C)$ is maximized for

$$(\vec{w} \cdot \vec{c})^{\text{OPT}} = \log \left(\frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)} \right) - \log k$$

- Note that

$$(W \cdot C^T)_{ij} = \vec{w}_i \cdot \vec{c}_j$$

- Let M^{OPT} be the matrix containing the optimal dot products, that is

$$M_{ij}^{\text{OPT}} = (\vec{w}_i \cdot \vec{c}_j)^{\text{OPT}}$$

Singular Value Decomposition (SVD)

- $(W \cdot C^T)_{ij} = \vec{w}_i \cdot \vec{c}_j$ and $M_{ij}^{\text{OPT}} = (\vec{w}_i \cdot \vec{c}_j)^{\text{OPT}}$

Singular Value Decomposition (SVD)

- $(W \cdot C^T)_{ij} = \vec{w}_i \cdot \vec{c}_j$ and $M_{ij}^{\text{OPT}} = (\vec{w}_i \cdot \vec{c}_j)^{\text{OPT}}$

- Skip-gram with negative sampling is trying to find W and C such that

$$W \cdot C^T = M^{\text{OPT}}$$

Singular Value Decomposition (SVD)

- $(W \cdot C^T)_{ij} = \vec{w}_i \cdot \vec{c}_j$ and $M_{ij}^{\text{OPT}} = (\vec{w}_i \cdot \vec{c}_j)^{\text{OPT}}$

- Skip-gram with negative sampling is trying to find W and C such that

$$W \cdot C^T = M^{\text{OPT}}$$

- **Truncated SVD** gives us a factorization of the best rank d approximation of M^{OPT} :

$$W_{\text{SVD}} \cdot C_{\text{SVD}}^T = \arg \min_{M | \text{rk}(M)=d} \|M - M^{\text{OPT}}\|_F$$

Skip-Gram (without negative sampling)

Recall from previous slide:

$$\ell_{SG}(W, C) = \sum_{(w,c) \in D} \left(\vec{w} \cdot \vec{c} - \log \left(\sum_{c' \in V_C} \exp(\vec{w} \cdot \vec{c}') \right) \right)$$

Computations for the skip-gram model (without negative sampling) yield a maximum for

$$(\vec{w} \cdot \vec{c})^{\text{OPT}} = \log \#(w, c)$$

Problems with SVD

$$M_{ij}^{\text{OPT}} = \log \left(\frac{\#(w_i, c_j) \cdot |D|}{\#(w_i) \cdot \#(c_j)} \right) - \log k$$

Problems with SVD

$$M_{ij}^{\text{OPT}} = \log \left(\frac{\#(w_i, c_j) \cdot |D|}{\#(w_i) \cdot \#(c_j)} \right) - \log k$$

- 1 What about pairs with $\#(w_i, c_j) = 0$?
(This is the case for more than 98% of our pairs!)
- 2 M^{OPT} is dense.

Problems with SVD

$$M_{ij}^{\text{OPT}} = \log \left(\frac{\#(w_i, c_j) \cdot |D|}{\#(w_i) \cdot \#(c_j)} \right) - \log k$$

- 1 What about pairs with $\#(w_i, c_j) = 0$?
(This is the case for more than 98% of our pairs!)
- 2 M^{OPT} is dense.

Solution: Factorize

$$M_{ij}^+ = \max \left(\log \left(\frac{\#(w_i, c_j) \cdot |D|}{\#(w_i) \cdot \#(c_j)} \right) - \log k, 0 \right)$$

Questions?

Table of Contents

1 Problem

2 Solution

3 Evaluation

Experiment Setup

data:	~ 4.6 million English Wikipedia articles
vocabulary size:	~ 160,000 (words that appeared at least 300 times)
window size:	2
word-context samples:	~ 9.7 billion
embedding dimension:	200

Table of Contents

- 3 Evaluation
 - Optimizing the objective
 - Word Similarity Tasks
 - Analogy Tasks

Optimizing the Objective

The following table shows the percentage of deviation from the optimal value, that is

$$\frac{l - l^{\text{OPT}}}{l^{\text{OPT}}}.$$

k	l^{OPT}	l^+	SVD	NN
0	0%	5.7%	25.1%	-
1	0%	29.3%	38.8%	22.7%
5	0%	120.9%	124.7%	9.5%
15	0%	309.0%	310.4%	8.9%

Table: Percentage of deviation from the optimal objective value.

Table of Contents

3 Evaluation

- Optimizing the objective
- **Word Similarity Tasks**
- Analogy Tasks

Word Similarity Tasks

Models were tested to two datasets:

- WordSim353: 353 word pairs
- MEN: 3000 word pairs

word pairs		human assigned similarity scores
stock	market	8.08
physics	chemistry	7.35
game	round	5.97
experience	music	3.47
stock	jaguar	0.92

Table: Examples from the WordSim353 dataset

Word Similarity Tasks

k	WordSim353		MEN	
	NN	SVD	NN	SVD
0	-	0.601	-	0.655
1	0.524	0.613	0.588	0.700
5	0.658	0.536	0.712	0.669
15	0.644	0.400	0.681	0.606

Table: Spearman's correlation between dataset similarity scores and similarity scores that different the models returned.

Note: Spearman's correlation $\rho_S \in [-1, 1]$, where negative (positive) numbers indicate negative (positive) correlation and zero indicates no correlation.

[More about Spearman's correlation](#)

Table of Contents

3 Evaluation

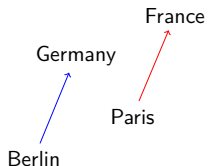
- Optimizing the objective
- Word Similarity Tasks
- **Analogy Tasks**

Analogy Tasks

Berlin is to **Germany** as **Paris** is to **France**.

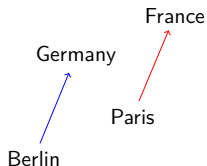
Analogy Tasks

Berlin is to **Germany** as **Paris** is to **France**.



Analogy Tasks

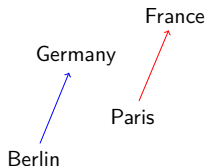
Berlin is to **Germany** as **Paris** is to **France**.



$$\Rightarrow \text{vec}(\mathbf{Germany}) - \text{vec}(\mathbf{Berlin}) = \text{vec}(\mathbf{France}) - \text{vec}(\mathbf{Paris})$$

Analogy Tasks

Berlin is to **Germany** as **Paris** is to **France**.



$$\Rightarrow \text{vec}(\mathbf{Germany}) - \text{vec}(\mathbf{Berlin}) = \text{vec}(\mathbf{France}) - \text{vec}(\mathbf{Paris})$$

in other words:

$$\text{vec}(\mathbf{France}) = \text{vec}(\mathbf{Germany}) - \text{vec}(\mathbf{Berlin}) + \text{vec}(\mathbf{Paris})$$

Analogy Tasks

k	Mixed dataset 19.500 analogies		Syntactic dataset 8.000 analogies	
	NN	SVD	NN	SVD
0	-	26.8%	-	28.7%
1	27.3%	30.6%	32.3%	19.6%
5	51.0%	12.0%	51.0%	5.7%
15	53.2%	5.9%	47.9%	1.4%

Table: Percentage of correct answers on two word analogy datasets.

More examples

Questions?

Expectation of the closest vector

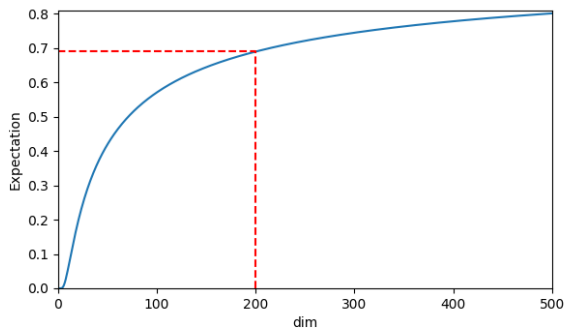


Figure: Expectation of the cosine distance to the nearest vector for 159,862 vectors depending on the embedding dimension.

Back

Expectation of the closest vector

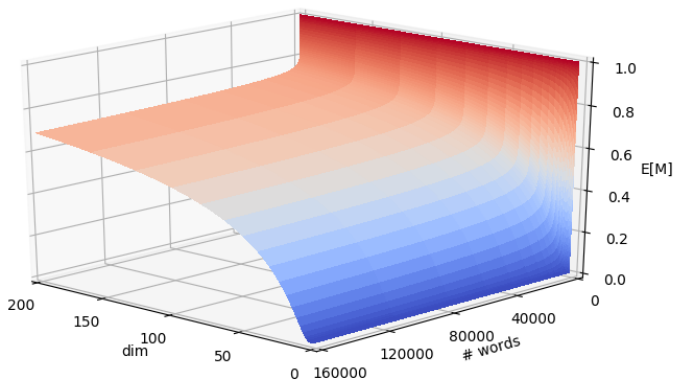
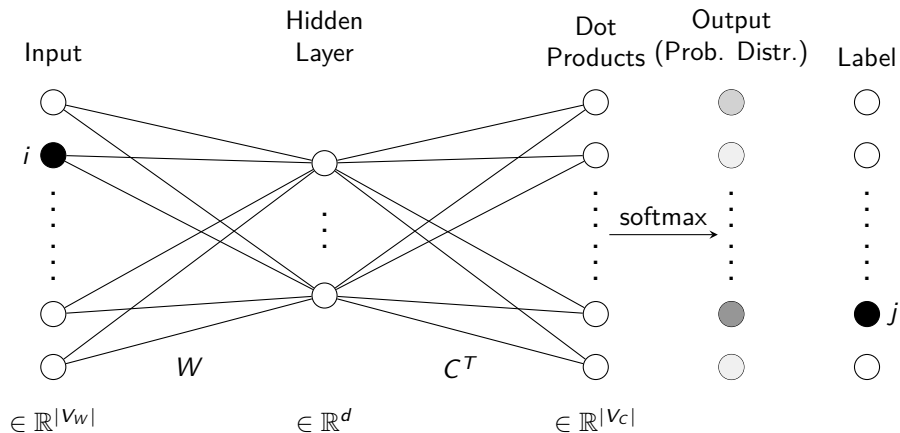


Figure: The expectation of the distance to the closest word depending on the embedding dimension and the number of words.

Skip-Gram

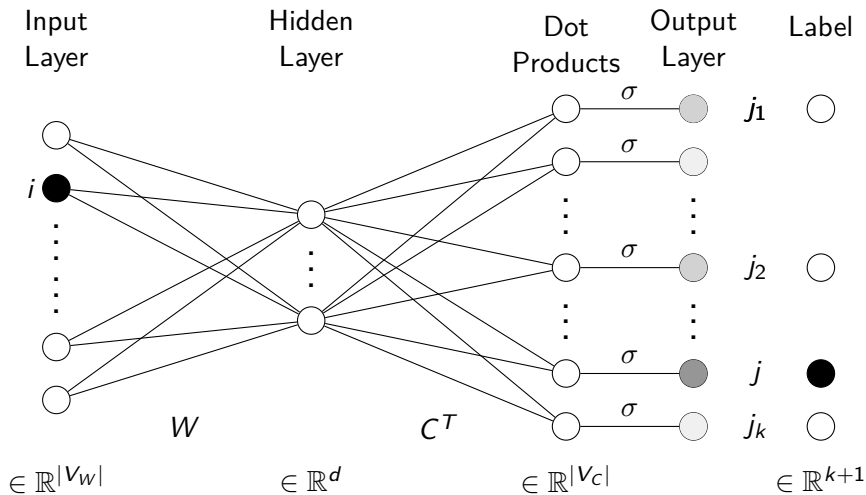


Objective function SG

$$\begin{aligned}\ell_{SG}(W, C) &= \sum_{(w,c) \in D} \log \frac{\exp(\vec{w} \cdot \vec{c})}{\sum_{c' \in V_C} \exp(\vec{w} \cdot \vec{c}')} \\ &= \sum_{(w,c) \in D} \left(\vec{w} \cdot \vec{c} - \log \left(\sum_{c' \in V_C} \exp(\vec{w} \cdot \vec{c}') \right) \right)\end{aligned}$$

Back

Skip-Gram with negative sampling

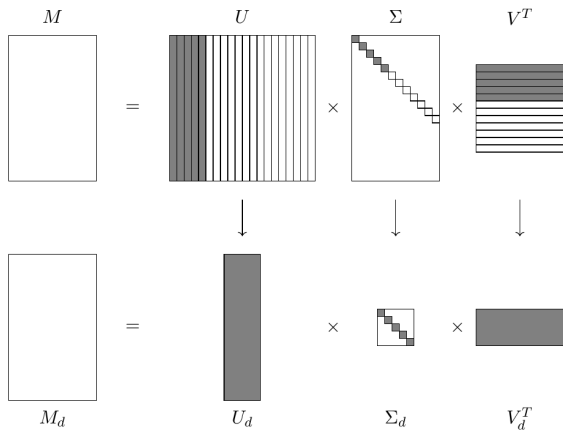


Objective function SGNS

$$\begin{aligned} \ell_{SGNS}(W, C) &= \sum_{(w_i, c_j) \in D} \left(\log \sigma(\vec{w}_i \cdot \vec{c}_j) + \sum_{l=1}^k \log(1 - \sigma(\vec{w}_i \cdot \vec{c}_{j_l})) \right) \\ &= \sum_{(w_i, c_j) \in D} \left(\log \sigma(\vec{w}_i \cdot \vec{c}_j) + \sum_{l=1}^k \log \sigma(-\vec{w}_i \cdot \vec{c}_{j_l}) \right) \\ &\approx \sum_{(w, c) \in D} \left(\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} [\log \sigma(-\vec{w} \cdot \vec{c}_N)] \right) \end{aligned}$$

Back

Truncated SVD

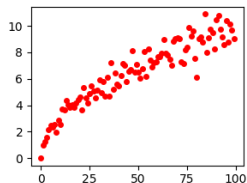


Back

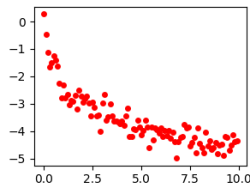
Spearman correlation

Let X_i be the human-assigned scores and Y_i be the cosine similarity of the vectors. Then, the Spearman correlation is defined as

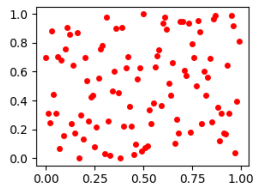
$$\rho_S = \frac{\text{cov}(\text{rg}(X), \text{rg}(Y))}{\sigma(\text{rg}(X))\sigma(\text{rg}(Y))} \in [-1, 1].$$



(a) positive



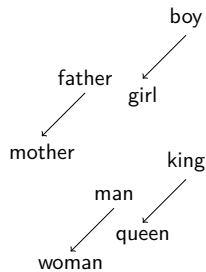
(b) negative



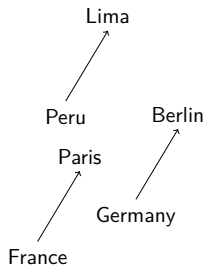
(c) around zero

Figure: Datasets with different Spearman correlation

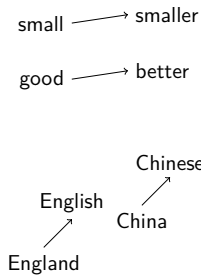
Analogy Tasks



(a) man-woman



(b) capitals



(c) syntactic relations

Figure: Examples of various relations between words

Analogy Tasks

