

# Zero-Training Sentence Embedding via Orthogonal Basis

Ruslan Rakhimov   Alexey Bokhovkin  
Ivan Fursov   Eugenia Cheskidova

Skolkovo Institute of Science and Technology

December 2018

# Outline

---

## Introduction

## Approach

- Quantify new semantic meaning

- Novelty

- Significance

- Uniqueness

- Sentence vector

## Experiments

- Algorithm

- STS benchmark

- Paragraph embeddings

- Supervised tasks

## Improvements

- n-grams

## Summary

# Introduction

---

A **word embedding** is a learned representation for text where words that have the same meaning have a similar representation. Word embeddings are in fact a class of techniques where individual words are represented as real-valued vectors in a predefined vector space.

- ▶ Parameterized sentence embeddings
  - ▶ SkipThought (Kiros et al., 2015)
  - ▶ Sent2Vec (Pagliardini et al., 2018)
  - ▶ InferSent (Conneau et al. 2017)
  - ▶ Universal Sentence Encoder (Yang et al., 2018; Cer et al., 2018)
- ▶ Non-parameterized sentence embedding
  - ▶ SIF (Arora et al., 2017)
  - ▶ Concatenated  $p$ -mean (Ruckle et al., 2018)

# Problem statement

---

We study a new simple and robust non-parameterized approach for building sentence representations.

Inspired by the Gram-Schmidt Process in geometric theory, authors build an orthogonal basis of the subspace spanned by a word and its surrounding context in a sentence.

## Quantify new semantic meaning

---

Consider a word  $w_i$  in sequence and its  $m$ -neighbourhood window inside sentence. Then **Contextual Window Matrix**:

$$\mathbf{S}^i = [\mathbf{v}_{w_{i-m}}, \dots, \mathbf{v}_{w_{i-1}}, \mathbf{v}_{w_{i+1}}, \dots, \mathbf{v}_{w_{i+m}}, \mathbf{v}_{w_i}] \in \mathbb{R}^{d \times (2m+1)}$$

Then compute novel semantic information compared with its context:  $\mathbf{S}^i = \mathbf{Q}^i \mathbf{R}^i$ . In this way we generate new orthogonal word embedding vector:

$$\mathbf{Q}_{:,2m+1}^i = \mathbf{q}_i \rightarrow \{\mathbf{q}_1, \dots, \mathbf{q}_{i-1}\}$$

In order to generate the embedding for a sentence, weights of *novelty*, *significance* and *uniqueness* will be assigned to each of its words.

# Novelty

---

A word  $w_i$  is more important to a sentence if its novel orthogonal basis vector  $\mathbf{q}_i$  is a large component in  $\mathbf{v}_{w_i}$ .

$$\alpha_n = \exp\left(\frac{r_{-1}}{\|\mathbf{v}_{w_i}\|_2}\right) = \exp\left(\frac{r_{-1}}{\|\mathbf{r}\|_2}\right)$$

where  $\mathbf{r}$  is the last column of  $\mathbf{R}_i$ , and  $r_{-1}$  is the last element of  $\mathbf{r}$ .

$\alpha_n$  is the exponential of the normalized distance between  $\mathbf{v}_{w_i}$  and the subspace spanned by its context.

# Significance

---

Intuition: The significance of a word is related to how semantically aligned it is to the meaning of its context. SVD is used in to identify principal meanings of the context.

$$\mathbf{S}^i = \mathbf{U}^i \Sigma^i (\mathbf{V}^i)^T$$

A word is more important if its novel semantic meaning has a better alignment with more principal meanings.

$$\alpha_s = \frac{\|\sigma(\mathbf{S}^i) \odot (\mathbf{q}_i^T \mathbf{U}^i)\|_2}{2m + 1} = \frac{r_{-1}}{2m + 1}$$

$\alpha_s$  is essentially the distance between  $\mathbf{w}_i$  and the context hyper-plane, normalized by the context size.

## Uniqueness

---

A corpus-wise uniqueness of *stop words* (commonly present in the corpus) is small. We compute the principal directions of the corpus and then measure their alignment with the novel orthogonal basis vector  $\mathbf{q}_i$ . A high alignment means a relatively low corpus-wise uniqueness score, and vice versa.

We want to obtain an intermediate coarse-grained sentence embedding matrix  $\mathbf{X}^c = [\mathbf{g}_1, \dots, \mathbf{g}_N] \in \mathbb{R}^{d \times N}$

Suppose  $\mathbf{S} = [\mathbf{v}_{w_1}, \dots, \mathbf{v}_{w_n}] = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Then the coarse-grained embedding for the  $i$ -th sentence is defined as:

$$\mathbf{g}_i = \sum_{j=1}^n f(\sigma_j) \mathbf{U}_{:,j}$$

where  $f(\sigma_j)$  is a monotonically increasing function.



## Uniqueness

---

1. We then compute the top  $K$  principal vectors  $\{\mathbf{d}_1, \dots, \mathbf{d}_K\}$  of  $\mathbf{X}^c$ , with singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_K$
2. For each sentence,  $\{\mathbf{d}_1, \dots, \mathbf{d}_K\}$  are re-ranked in descending order of their correlation with sentence matrix  $\mathbf{S}$ . The correlation is defined as  $o_i = \sigma_i \|\mathbf{S}^T \mathbf{d}_i\|_2$ ,  $1 \leq i \leq K$ .
3. Next, the top  $h$  principal vectors after re-ranking based on  $o_i$  are selected:  $\mathbf{D} = \{\mathbf{d}_{t_1}, \dots, \mathbf{d}_{t_h}\}$ , with  $o_{t_1} \geq o_{t_2} \geq \dots o_{t_h}$  and their singular values in  $\mathbf{X}^c$  are  $\boldsymbol{\sigma}_d = [\sigma_{t_1}, \dots, \sigma_{t_h}] \in \mathbb{R}^h$ .

Finally, a word  $w_i$  with new semantic meaning vector  $\mathbf{q}_i$  in this sentence will be assigned a corpus-wise uniqueness score:

$$\alpha_u = \exp(-\|\boldsymbol{\sigma}_d \odot (\mathbf{q}_i^T \mathbf{D})\|_2 / h)$$

This ensures that common stop words will have their effect diminished since their embeddings are closely aligned with the corpus' principal directions.

## Sentence vector

---

A sentence vector  $\mathbf{c}_s$  is computed as a weighted sum of its word embeddings, where the weights come from three scores: a novelty score ( $\alpha_n$ ), a significance score ( $\alpha_s$ ) and a corpus-wise uniqueness score ( $\alpha_u$ ).

$$\alpha_i = \alpha_n + \alpha_s + \alpha_u$$
$$\mathbf{c}_s = \sum_i \alpha_i \mathbf{v}_{w_i}$$

Given a set of sentence vectors, removing projections onto the principal components of the spanned subspace can significantly enhance the performance on semantic similarity task. However, as each sentence may have a different semantic meaning, it could be sub-optimal to remove the same set of principal components from all sentences.

$$\mathbf{c}_s \leftarrow \mathbf{c}_s - \sum_{j=1}^K (\mathbf{d}_{t_j}^T \mathbf{c}_s) \mathbf{d}_{t_j}$$

# Algorithm

---

## Algorithm 1 Geometric Embedding (GEM)

---

```

1: Inputs:
   A set of sentences  $\mathcal{S}$ , vocabulary  $\mathcal{V}$ , word embeddings  $\{v_w \in \mathbb{R}^d \mid w \in \mathcal{V}\}$ 
2: Outputs:
   Sentence embeddings  $\{c_s \in \mathbb{R}^d \mid s \in \mathcal{S}\}$ 
3: for ith sentence  $s$  in  $\mathcal{S}$  do
4:   Form matrix  $S \in \mathbb{R}^{d \times n}$ ,  $S_{i,j} = v_{w_j}$  and  $w_j$  is the  $j$ th word in  $s$ 
5:   The SVD is  $S = U \Sigma V^T$ 
6:   The  $i$ th column of the coarse-grained sentence embedding matrix  $X_{:,i}^c$  is  $U(\sigma(S))^3$ 
7: end for
8: Take first  $K$  singular vectors  $\{d_1, \dots, d_K\}$  and singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_K$  of  $X^c$ 
9: for sentence  $s$  in  $\mathcal{S}$  do
10:   Re-rank  $\{d_1, \dots, d_K\}$  in descending order by  $\sigma_i = \sigma_i \|S^T d_i\|_2$ ,  $1 \leq i \leq K$ .
11:   Select top  $h$  principal vectors as  $D = [d_{t_1}, \dots, d_{t_h}]$ , with singular values  $\sigma_d = [\sigma_{t_1}, \dots, \sigma_{t_h}]$ .
12:   for word  $w_i$  in  $s$  do
13:      $S^i = [v_{w_{i-m}}, \dots, v_{w_{i-1}}, v_{w_{i+1}}, \dots, v_{w_{i+m}}, v_{w_i}]$  is the contextual window matrix of  $w_i$ .
14:     Do QR decomposition  $S^i = Q^i R^i$ , let  $q_i$  and  $r$  denote the last column of  $Q^i$  and  $R^i$ 
15:      $\alpha_n = \exp(r_{-1}/\|r\|_2)$ ,  $\alpha_s = r_{-1}/(2m+1)$ ,  $\alpha_u = \exp(-\|\sigma_d \odot (q_i^T D)\|_2/h)$ 
16:      $\alpha_i = \alpha_n + \alpha_s + \alpha_u$ 
17:   end for
18:    $c_s = \sum_{v_i \in s} \alpha_i v_{w_i}$ 
19:   Principal vectors removal:  $c_s \leftarrow c_s - DD^T c_s$ 
20: end for

```

---

$$\begin{aligned}
 &\text{Complexity: } \underbrace{O(Nd(\max \text{ length of sentence})^2)}_{\text{1st cycle}} + \underbrace{dN^2}_{\text{SVD of } X^c} \\
 &+ \underbrace{Nndm^2 + 3NdK}_{\text{2nd cycle (QR + matvec of } q_i \text{ and } D + \text{matvec } DD^T c_s)} )
 \end{aligned}$$

## STS benchmark

---

- ▶ To predict a cosine similarity score of two sentences given a sentence pair
- ▶ To evaluate the Pearson's coefficient  $r$  between human-labeled similarity (0 - 5 points) and predictions.

model	dev (article)	test (article)
Gem + LexVec	77.1 (81.9)	<b>63.9</b> (76.5)
Gem + Glove	<b>78.9</b> (—)	63.7 (—)
Mean + LexVec	66.6 (58.78)	28.8 (50.43)
Mean + Glove	51.8 (52.4)	23.8 (40.6)

Table: STS benchmark results

One can see that our results differ from the ones obtained in the article. Likely because of a bit different precomputed embeddings.

## Paragraph embeddings

---

Here we **developed ourselves** and compared three different techniques of a paragraph embedding computation:

1. Treat paragraph as one sentence (without any punctuation)
2. Compute paragraph embedding as average of sentence embeddings
3. Compute sentence embeddings, then treat them as words and paragraph as a sentence

First (full)	Second	Third
<b>65.30</b> (73.28)	61.75	56.50

**Table:** Accuracy of predicted sentiment with logistic regression; 2000 out of 25000 paragraphs (GEM embeddings on IMDb sentiments) are trained

## Supervised tasks

---

We tested GEM on text classification, sentiment analysis and textual semantic similarity tasks. We evaluate performance using cosine similarity and 0.5 threshold.

Task	GEM <sup>1</sup>	DIIN <sup>2</sup>	ULMFiT <sup>3</sup>
Quora Question Pairs	66.1	89.06	—
IMDB	73.28	—	95.4

**Table:** Comparison of GEM and SOTA models (accuracy)

**Supervised models are still superior to unsupervised methods.**

---

<sup>1</sup>Here we treat paragraphs as one sentence

<sup>2</sup>Gong et al., 2018

<sup>3</sup>Howard and Ruder, 2018

# n-grams

## Motivation

---

The idea is to implement  $n$ -grams and to average embeddings of nearest words that stay close to each other. Then see how the performance changes.

e.g. bigrams: A girl is Arya Stark of Winterfell NULL

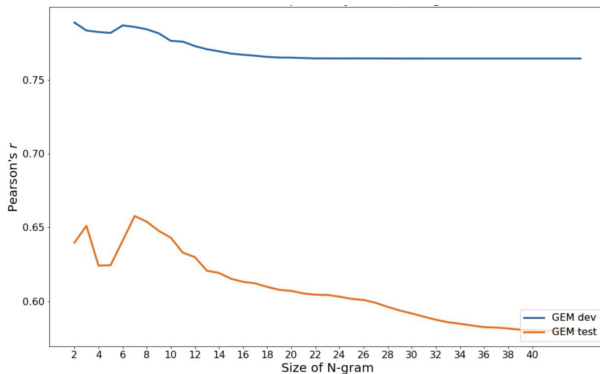


e.g. trigrams: A girl is Arya Stark of Winterfell NULL NULL



# n-grams

## Performance Dependence



**Figure:** Performance dependency on  $n$  for  $n$ -grams (STS dataset)



# n-grams

## Speed Dependence

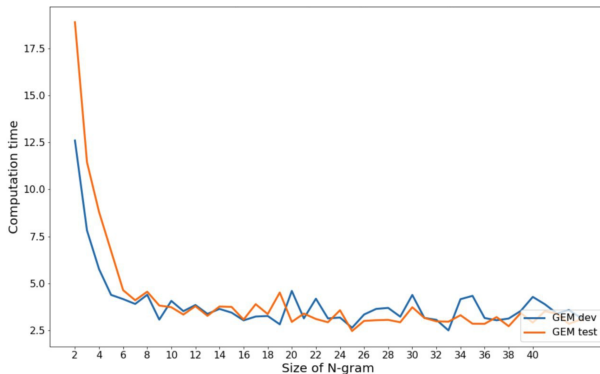


Figure: Speed dependency on  $n$  for  $n$ -grams (STS dataset)

# Summary

---

GEM has a potential use in model ensembling and fast prototyping.

## Advantages

1. Fast compared to supervised methods
2. No parameters
3. Unsupervised
4. No training needed

## Disadvantages

1. Time complexity grows as  $O(n^2)$  where  $n \doteq$  sentence length
2. Supervised methods are superior to GEM algorithm

# References

---



Yang, Z., Zhu, C., Chen, W. (2018). Zero-training Sentence Embedding via Orthogonal Basis.