Monolingual embeddings

Alignments

Applications and Results

Semi-supervised Learning for Multilingual Sentence Representation

Hedi Ben-younes - Alexandre Ramé

March 2016





・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ うへぐ

Monolingual embeddings

Alignments

Applications and Results

Outline



- 2 Monolingual embeddings
- 3 Alignments
- 4 Applications and Results





Introduction	Monolingual embeddings	Alignments	Applications and Results
0000	000000	000000000	0000000000

We need to know what is being said on the web

- \Rightarrow understand web pages.
 - Images,
 - Links to other pages,
 - Texts
 - ...

Understand text = being able to represent documents in a semantic space.



Introduction	
00000	

Problem: text comes in different languages We want to build a system that can compute similarity between texts:

- in different languages
- only based on semantics
- \Rightarrow We need a multilingual text embedding.



Monolingual embeddings	Align
0000000	0000

Alignments

Applications and Results

Two main problems are presented today:

How can we represent words and sentences using raw natural monolingual text ?



Figure: Monolingual embedding



Introduction

00000

Monolingual embeddings

Alignments

Applications and Results

e How can we align multiple embedding spaces and project them into a unique multilingual semantic space ?







Introduction
00000

Monolingual embeddings

Alignments

Applications and Results

How can we get a multilingual text embedding from

- a lot of raw unlabelled text,
- and a few aligned data ?







Monolingual embeddings

Alignments

Applications and Results

Outline



- 2 Monolingual embeddings
- 3 Alignments
- 4 Applications and Results





Monolingual embeddings

Alignments

Applications and Results

▲ロ▶ ▲周▶ ▲ヨ▶ ▲ヨ▶ ヨー のへで

Goal

Create as many monolingual embeddings as the number of languages.

How can we represent words and sentences in a semantic space ?



Monolingual embeddings

Alignments

Applications and Results

Word2Vec [Mikolov et al., 2013]

Word embeddings

- compact representation
- learnt from raw natural text
- close vectors = similar words
- linearize semantic relations





Figure: PCA of countries and their capital



Applications and Results

Word2Vec [Mikolov et al., 2013] - Skip-gram

Given a word, we try to predict which words are nearby

$$p(w_O|w_I) = \frac{\exp\left(v_{w_O}^{'\top} v_{w_I}\right)}{\sum_{w=1}^{W} \exp\left(v_w^{'\top} v_{w_I}\right)}$$



Figure: Skip-gram architecture

- The parameters of the model are the embeddings themselves
- Model trained with SGD
- Negative sampling for large vocabularies.



Monolingual embeddings

Alignments

Applications and Results

くしゃ 本語 を オヨ そ 山 や く しゃ

Sentence embeddings

How do we go from **word** embeddings to **sentence** embeddings ? First solution: bag-of-word2vec

$$S = [w_1, ..., w_n]$$

$$vec(S) = \sum_{i=1}^{n} word2vec(w_i)$$



Monolingual embeddings

Alignments

Applications and Results

Skip-thought vectors [Kiros et al., 2015]

Directly learn a sentence embedding

- learn how to compose word embeddings to get a task-independant sentence embedding,
- Applies the Skip-Gram idea to the sentence level: given a sentence, we try to generate the previous and the next sentence



- Uses a list of consecutive sentences
- One GRU-encoder and two GRU-decoders



Monolingual embeddings

Alignments

Applications and Results

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

FastSent [Hill et al., 2016]

Simplified version of Skip-thought vectors

• Learns source u_w and target v_w embeddings

• Encoder:
$$\mathbf{s_i} = \sum_{w \in S_i} u_w$$

Maximize

$$\sum_{w \in S_{i-1} \bigcup S_{i+1}} \Phi(\mathbf{s_i}, v_w)$$

where Φ is a softmax over the vocabulary



Monolingual embeddings

Alignments

Applications and Results

▲ロ▶ ▲周▶ ▲ヨ▶ ▲ヨ▶ ヨー のへで

Conclusion

Now we can:

- represent words
- represent sentences
 - sum of Word2Vecs
 - Skip-thought vectors
 - FastSent
- our current sentence embedding is based upon a sum of disambiguated Word2Vecs
- we are currently switching to FastSent



Monolingual embeddings

Alignments

Applications and Results

Outline

Introduction

2 Monolingual embeddings

3 Alignments







Introduction	
00000	

Monolingual embeddings

Alignments ••••••• Applications and Results

Up to now, we have one embedding model per language. We want to find projections from monolingual spaces to a multilingual space.

To do so, we need:

• aligned multi-lingual data

FR	EN	

• an alignment model



Monolingual embeddings

Alignments

Applications and Results

Alignment Data

European Languages

From Europarl (dataset of aligned sentences) to aligned words with GIZA++ ([Koehn et al., 2007])





Other Languages

Most common english words translated into Chinese, Russian, ...



Monolingual embeddings

Alignments

Applications and Results

Canonical Correlation Analysis

CCA:

- proposed in [Hotelling, 1936]
- very well explained in [Hardoon et al., 2004]

Problem

Given two matrices $X \in \mathbb{R}^{n \times d_x}$ and $Y \in \mathbb{R}^{n \times d_y}$ representing the observed values x and y of paired centered multivariate random variables,

we seek for two projections $W \in \mathbb{R}^{d_x imes p}$ and $V \in \mathbb{R}^{d_y imes p}$ such as

• $\forall i \in [1, p], < W_{:,i}, x > \text{and} < V_{:,i}, y > \text{are highly correlated}$

• the dimensions in the target space are uncorrelated



Monolingual embeddings

Alignments

Applications and Results

▲ロ▶ ▲周▶ ▲ヨ▶ ▲ヨ▶ ヨー のへで

Canonical Correlation Analysis

Let

$$egin{aligned} \Sigma_{xy} &= X'Y \in \mathbb{R}^{d_x imes d_y} \ \Sigma_{xx} &= X'X \in \mathbb{R}^{d_x imes d_x} \ \Sigma_{yy} &= Y'Y \in \mathbb{R}^{d_y imes d_y} \end{aligned}$$

We find these projections by solving:

$$\max_{W,V} \operatorname{Tr}(W\Sigma_{xy}V) \qquad s.c \begin{cases} W'\Sigma_{xx}W = I \\ V'\Sigma_{yy}V = I \\ W'\Sigma_{xy}V * (1-I) = \mathbf{0} \end{cases}$$
(1)

 \Rightarrow yields a close form



Monolingual embeddings

Alignments

Applications and Results

◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ○ ○ ○

Canonical Correlation Analysis

To go further on CCA:

- Kernel CCA [Hardoon et al., 2004]
- Deep CCA [Andrew et al., 2013]

In practice, CCA doesn't scale for large alignment datasets:

- Need to go through the whole dataset to compute Σ_{kl}
- Kernel version requires storing / inverting Graam matrices



Monolingual embeddings

Alignments

Applications and Results

Correlational Network [Chandar et al., 2015]

CCA + Auto Encoder => Correlational Network



Figure: Correlational Network Architecture

The loss function of Correlational Network is the sum of:

- correlation of projections in the common subspace
- self-reconstruction terms
- cross-reconstruction terms



Monolingual embeddings

Alignments

Applications and Results

Correlational Network [Chandar et al., 2015]



Figure: Correlational Network Architecture

$$\begin{aligned} \mathcal{J}_{\mathcal{Z}}(\theta) &= \sum_{i=1}^{N} L((x_i, y_i), g(h((x_i, y_i)))) + L((x_i, y_i), g(h(x_i))) \\ &+ L((x_i, y_i), g(h(y_i))) - \lambda corr(h(X), h(Y)) \\ &corr(h(X), h(Y)) = \frac{\sum_{i=1}^{N} (h(x_i) - \overline{h(X)})h(y_i) - \overline{h(Y)}}{\sqrt{\sum_{i=1}^{N} (h(x_i) - \overline{h(X)})^2 \sum_{i=1}^{N} (h(y_i) - \overline{h(Y)})^2}} \end{aligned}$$



Monolingual embeddings

Alignments

Applications and Results

Hyperparameters

Hyperparameters

- λ of the correlation term
- mini-batch size
- regularization
- size hidden layer
- activation function: tanh
- add terms in the loss function
- tied or not
- contrastive term



Monolingual embeddings

Alignments

Applications and Results

Deep Correlational Network



Figure: Deep Correlational Network Architecture

Training Procedure: similar to greedy layerwise pretraining of deep autoencoders



Introduction Monolingual embeddings 00000 0000000		Alignments	Applications and Results

Bridge Correlational Network [Rajendran et al., 2015]



Figure: Bridge Correlational Network Architecture

- one pivot language
- multiple languages



Monolingual embeddings

Alignments

Applications and Results

Outline

- Monolingual embeddings



4 Applications and Results





Monolingual embeddings

Alignments

Applications and Results

▲ロト ▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ ● 臣 ● のへで

Transliteration Equivalence

English	French	Deutsch	Spanish	Greek	Chinese
machine learning:1.0	apprentissage automatique:0.83	algorithmen:0.81	computacional:0.82	επιστήμη υπολογιστών:0.78	模式 :0.76
pattern recognition :0.89	algorithmique:0.80	datenanalyse:0.79	aprendizaje automático:0.79	δηλωτική:0.75	算机:0.75
algorithms:0.87	algorithmes:0.79	datenstrukturen:0.75	algoritmos:0.79	συνδυαστική:0.75	建模:0.73
data mining:0.87	modélisation:0.79	mustererkennung :0.75	redes neuronales:0.76	αλγορίθμων:0.74	化:0.73
natural language processing :0.86	heuristiques:0.78	numerische mathematik:0.75	aprendizaje automático:0.76	συνδυαστικής:0.74	形式化:0.71
Arab	Korean	Russian	Italian	Portuguese	Esperanto
0.78:كهۇرزمىّت	응용:0.75	алгоритмов:0.81	computazionale:0.80	computacional:0.81	komputiko :0.74
0.78:ئمدھج	데이터마이닝 :0.74	алгоритмы:0.79	computazionali:0.79	computacionais:0.80	komputado:0.73
0.78: ألهّسوبي	알고리즘과:0.73	моделирования:0.78	algoritmi:0.77	algoritmos:0.78	komputiko:0.72
0.74: ألهّسوبي	신경망:0.70	вычислений:0.77	statistica:0.77	otimização:0.76	algoritmoj:0.71
0.74:ھسو بي	최 적 화:0.70	многомерных:0.77	computazione:0.76	linguagens de programação :0.76	analitiko:0.71

The top 5 closest words to 'machine learning' in cosine similarity



Monolingual embeddings

Alignments

Applications and Results

T-SNE





◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Monolingual embeddings

Alignments

Applications and Results

English Word Similarity Evaluation

smart	stupid	5.81	OF
wood	forest	7.73	da
money	cash	9.15	dc
king	queen	8.58	co
king	rook	5.92	pl

OPEC	country	5.63
day	summer	3.94
day	dawn	7.53
country	citizen	7.31
planet	people	5.75

Figure: Wordsim353 similarity and relatedness

English Model	WS-Similarity	WS-Relatedness	WS
Our SGNS ¹ Model in the English space	0.742	0.621	0.67
Our SGNS ¹ Model in the multilingual space	0.765	0.627	0.685
SGNS ¹	0.737 ² -0.798 ³	$0.592^2 - 0.700^3$	
Glove	$0.651^2 - 0.746^3$	0.541 ² -0.643 ³	
Swivel	0.748 ²	0.616 ²	
Spearman's a correlation on English Word Similarity task			

Spearman's ρ correlation on English Word Similarity tasks



¹Skip Gram Negative Sampling ²According to [Shazeer et al., 2016] ³According to [Levy et al., 2015]

Monolingual embeddings

Alignments

Applications and Results

Other Languages' Word Similarity Evaluation

intelligent stupide 5.81 bois forêt 7.73 argent monnaie 9.15 roi reine 8.58 roi tour 5.92

聪明的 愚蠢的	5.81
木材 林 7.73	
钱 货币 9.15	
景 女王 8.58	
景圆 5.92	

Figure: French and Chinese Wordsim

Lang	WS	Lang	WS
en	0.685	fr	0.594
de	0.603	eo	0.584
ru	0.617	es	0.597
pt	0.587	ko	0.573
it	0.596	zh	0.563
el	0.562	ar	0.566

Spearman's ρ correlation on Word Similarity task for different



languages

Monolingual embeddings

Alignments

Applications and Results

Bilingual Word Similarity Evaluation

smart stupide 5.81 intelligent stupid 5.81 king reine 8.58 roi queen 8.58 聪明的 stupide 5.81 intelligent 笨 5.81 投资回报率 reine 8.58 roi 女王 8.58

Figure: English-French and Chinese-French Wordsim

Langs	WS	Langs	WS
en/fr	0.613	fr/fr	0.594
de/fr	0.554	eo/fr	0.526
ru/fr	0.561	es/fr	0.578
pt/fr	0.542	ko/fr	0.524
it/fr	0.574	zh/fr	0.513
el/fr	0.524	ar/fr	0.518

Spearman's ρ Correlation on Bilingual Word Similarity task for different languages with French



Monolingual embeddings

Alignments

Applications and Results

Sentence Similarity Evaluation: SICK

SICK(Sentences Involving Compositional Knowledge) dataset

- "A group of children is playing in the house and there is no man standing in the background"
- "A group of kids is playing in a yard and an old man is standing in the background"
- similarity=3.2

Model	Sick
Our English Model	0.59
Skip-thought	0.57 ¹
FastSent	0.61 ¹

Spearman's ρ correlation on Sentence Similarity task



Monolingual embeddings

Alignments

Applications and Results

▲ロ▶ ▲周▶ ▲ヨ▶ ▲ヨ▶ ヨー のへで

Sentence Multilingual Application

- Transfer Learning
 - Cross Language Topic Classification
 - Cross Language Moderation
- Multilingual Information Retrieval
- Named Entity Recognition: detection of persons, brands, topics, ...
- Multilingual Disambiguation



Monolingual embeddings

Alignments

Applications and Results

Multilingual Search Engine





Monolingual embeddings

Alignments

Applications and Results

◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ○ ○ ○

Conclusion

Summary

We are able to represent in a **unique** semantic space:

- words and sentences
- in multiple languages (12 up to now)
- => understanding of texts from the web

Ongoing Works

- Character-aware embeddings
- Add images in our multimodal space



Monolingual embeddings

Alignments

Applications and Results

Bibliographie I





Monolingual embeddings

Alignments

Applications and Results

・ロト ・ 四ト ・ 日ト ・ 日 ・

Bibliographie II



Levy, O., Goldberg, Y., and Dagan, I. (2015). Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics*, 3:211–225.



Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *NIPS 2013*, pages 3111-3119.



Rajendran, J., Khapra, M. M., Chandar, S., and Ravindran, B. (2015). Bridge correlational neural networks for multilingual multimodal representation learning. arXiv preprint arXiv:1510.03519.



Shazeer, N., Doherty, R., Evans, C., and Waterson, C. (2016). Swivel: Improving embeddings by noticing what's missing. *arXiv preprint arXiv:1602.02215.*

