

## Contents

INF5820

Distributional Semantics: Extracting Meaning from Data

### Lecture 5

## Kings and queens, men and women: semantic relations between word embeddings

Andrey Kutuzov  
andreku@ifi.uio.no

23 November 2016

- 1 Semantic relations as geometrical directions
- 2 Why vector algebra works?
- 3 Possible applications
- 4 Projecting one model into another: models alignment
- 5 In the next week (last lecture)

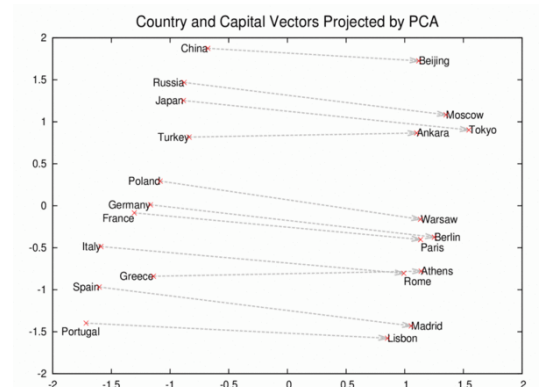
1

1

## Semantic relations as geometrical directions

## Semantic relations as geometrical directions

- ▶ Reducing high-dimensional representations to 2D or 3D can be used not only to create fancy pictures (or videos);
- ▶ Sometimes looking at the data in an understandable form can bring great insights...



'...ability of the model to automatically organize concepts and learn implicitly the relationships between them...' [Mikolov et al., 2013b]

2

3

## Semantic relations as geometrical directions

### A surprising discovery

- ▶ Algebraic operations on word embeddings reflect semantic relations
- ▶ simple operations with word vectors, like addition, subtraction, finding average, etc, produce intuitively meaningful results.

### Examples

- ▶ Element-wise adding phone vector to mobile vector results in a vector very close to cellphone;
- ▶ Element-wise subtracting France vector from Paris vector and adding Germany vector results in a vector very close to Berlin;

$$\vec{P} - \vec{F} + \vec{G} \approx \vec{B} \quad (1)$$

or

$$\vec{G} + \vec{P} - \vec{F} \approx \vec{B} \quad (2)$$

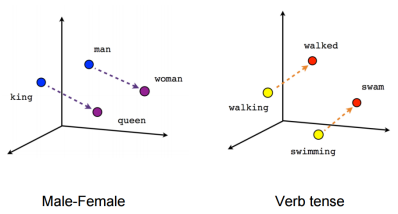
## Semantic relations as geometrical directions

- ▶ The latter operation can be expressed as solving proportions or performing analogical inference:
- ▶ 'France is to Paris as Germany is to what?'
- ▶ this is what Google Analogy dataset tries to evaluate.

Stanford folks call that 'directions of meaning' [Pennington et al., 2014]. The GloVe model to some extent was designed with this particular task in mind.

## Semantic relations as geometrical directions

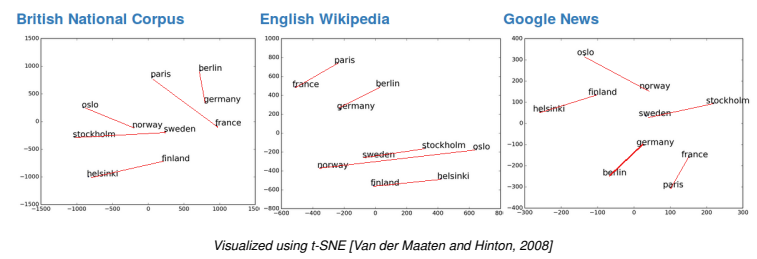
But isn't it simply finding 'geometrical' directions in the high-dimensional semantic space – from one word to another?



Can be any inference relation: plurality (dog → dogs), causality (life → experience), hypernymy (falcon → bird), meronymy (leg → body), antonymy (hot → cold), etc.

## Semantic relations as geometrical directions

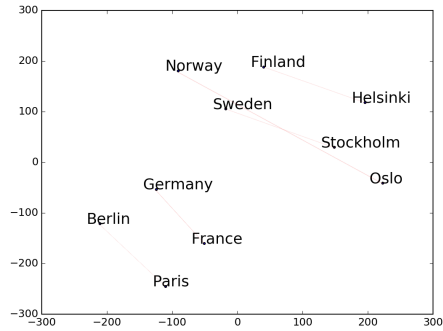
Semantic direction 'a capital of' in different models:



But why the Google News model looks so bad? Any ideas?

# Semantic relations as geometrical directions

...we used **lower-case** lemmas, and this model has **separate** embeddings for title-case and lower-case words (quite annoying).  
Let's try **title-case** countries and cities...



Looks better, doesn't it?

# Contents

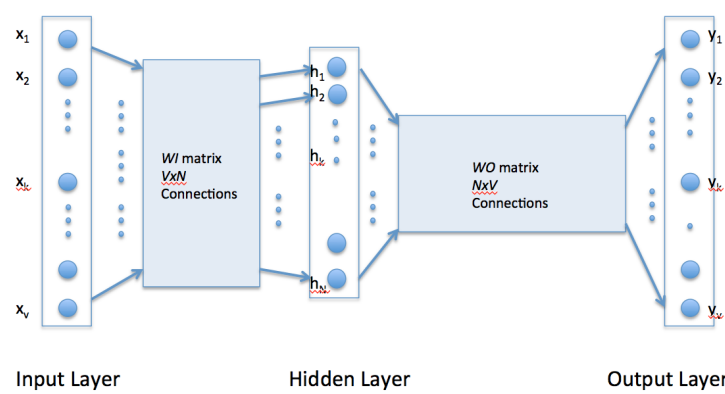
- 1 Semantic relations as geometrical directions
- 2 Why vector algebra works?
- 3 Possible applications
- 4 Projecting one model into another: models alignment
- 5 In the next week (last lecture)

# Why vector algebra works?

Why sum (element-wise addition) works?  
phone + mobile = cellphone

# Why vector algebra works?

Let's recall the prediction model architecture:



## Why vector algebra works?

### Why sum (element-wise addition) works?

- phone + mobile = cellphone
- ▶ input vectors represent context distribution;
  - ▶ they are related logarithmically to output probabilities;
  - ▶ the sum of 2 input vectors  $\vec{x}_1$  and  $\vec{x}_2$  is related to the product of 2 probability context distributions;
  - ▶ when multiplying 2 probability distributions, only words with high probabilities for both  $\vec{x}_1$  and  $\vec{x}_2$  will rank higher;
  - ▶ ...sort of AND function;
  - ▶ if cellphone appears frequently as a context word for phone and mobile, it will become the answer.

11

## Why vector algebra works?

### Why analogical inference works?

Germany + Paris - France = Berlin

- ▶ That's more difficult to explain.
- ▶ [Levy et al., 2015] are not sure that relational information exists in contextual features at all.
- ▶ In fact, Germany + Paris produces Berlin as well.
- ▶ So it might be that in many cases there are no 'relations' at all, only finding nearest associates of the positive input words.

But the idea is still attractive and produces empirical results.

12

## Contents

- 1 Semantic relations as geometrical directions
- 2 Why vector algebra works?
- 3 Possible applications
- 4 Projecting one model into another: models alignment
- 5 In the next week (last lecture)

## Possible applications

Geometric directions between word embeddings are used in many sense-related applications.

### Information extraction through semantic relations

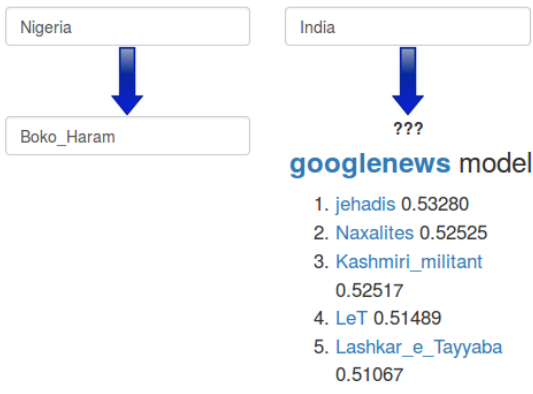
- ▶ It is possible to find semantic relations like **A is a capital of B**, or **A is located in B**.
- ▶ But it can be **A is an object of B's actions**, for that matter.
- ▶ We can even build full-fledged **ontologies** using this approach.
- ▶ ...or find phenomena analogical to other phenomena in complex ways.

12

13

## Possible applications

Unsupervised information extraction: militant groups across countries



## Possible applications

Unsupervised information extraction: militant groups across countries



## Possible applications

Unsupervised information extraction: militant groups across countries



## Possible applications

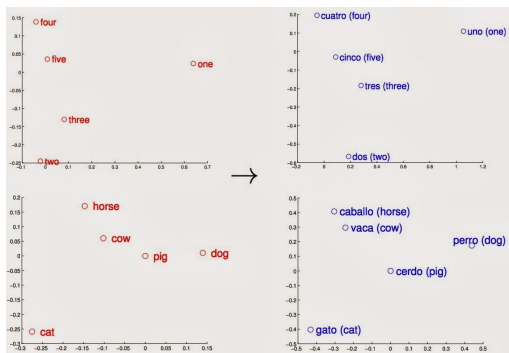
**What next?**

- ▶ We can **cross the boundaries of one model** and find links between languages (relation of 'being a translation of').
- ▶ **'Machine translation' for words**, relying only on **monolingual distributional models** and **small bilingual dictionaries**.
- ▶ Let's see how it can be done

Try yourself at <http://ltr.uio.no/semvec/calculator>

## Possible applications

### Inter-lingua?

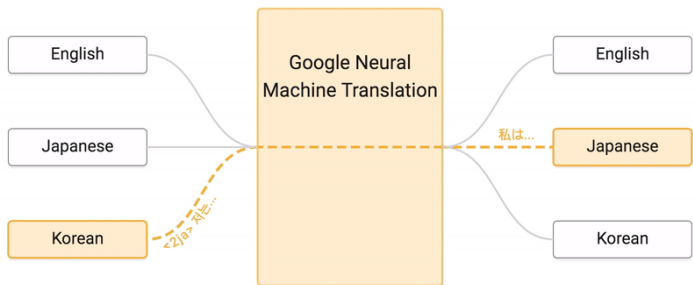


[Mikolov et al., 2013a]

Languages share concepts grounded in the real world → after proper rotation and scaling, models (semantic spaces) trained on comparable corpora from different languages should 'map' into each other.

## Possible applications

### Zero-shot

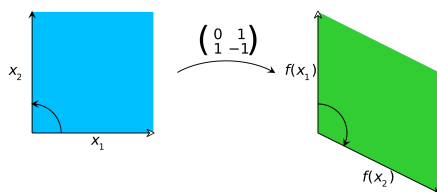


By the way, the recently announced Google's Zero-Shot neural machine translation model [Johnson et al., 2016] seems to rely on the same inter-lingua intuition.

## Possible applications

The Mikolov's idea of 'distributional translation' is to find regular correspondences between vector spaces representing each language:

1. train monolingual distributional models for languages L1 and L2;
2. take a small bilingual dictionary L1 → L2 with n frequent words present in both models;
3. learn linear transformation matrix M between embeddings of dictionary words in L1 and their translations in L2;



But how we do this?

## Contents

- 1 Semantic relations as geometrical directions
- 2 Why vector algebra works?
- 3 Possible applications
- 4 Projecting one model into another: models alignment
- 5 In the next week (last lecture)

Projecting one model into another: models alignment

- ▶ Mapping one model into another by finding a **linear transformation matrix  $M$**  is an **optimization problem** of minimizing prediction error – **linear regression**, to be exact;
- ▶ to find  $M$ , one either:
  1. learns the optimal weights **iteratively** (with gradient descent, etc.);
  2. or solves the quadratic minimization problem by finding **normal equation**.
- ▶ normal equation approach is **guaranteed to find the global optimum**.

$$\beta_i = (\mathbf{X}^T \times \mathbf{X})^{-1} \times \mathbf{X}^T \times y_i \tag{3}$$

$\mathbf{X}$  is the matrix of L1 dictionary word embeddings (input),  $y_i$  is the vector of the  $i^{th}$  components of corresponding L2 dictionary words (correct predictions), and  $\beta_i$  is our aim: the vector of optimal coefficients which transform an L1 embedding into the  $i^{th}$  component of the L2 embedding.

Projecting one model into another: models alignment

After solving such normal equations for all the components (their number is equal to the **L2** embeddings dimensionality  $i$ ), we have the matrix  $M$  which fits the data best:

$$\sum_{a=1}^n (M\vec{x}_a - \vec{y}_a)^2 \tag{4}$$

is **minimal over all training data** ( $n$  dictionary pairs).

- ▶ For any unseen word in **L1**, you can now multiply its vector  $\vec{x}$  in the **L1** model by  $M$ :  $\vec{z} = M\vec{x}$ .
- ▶  $\vec{z}$  is the **translation vector**;
- ▶ Find the word in **L2** with the embedding most similar to  $\vec{z}$
- ▶ Use that word as a **translation**.

Projecting one model into another: models alignment

Results					
[Mikolov et al., 2013a] report accuracy about <b>0.5 @1</b> and <b>0.75 @5</b> for <i>English</i> → <i>Spanish</i> translation.					
Our results for <i>Ukrainian</i> → <i>Russian</i> :					
	CBOW		SkipGram		Edit distance
	Training	Test	Training	Test	
@1	0.648	<b>0.57</b>	0.545	0.374	0.549
@5	0.764	<b>0.658</b>	0.644	0.486	0.619

By the way, **semantic fingerprints can be projected into another semantic space using the same transformation matrix!**  
 It means **we can 'translate' whole documents**.

Projecting one model into another: models alignment

- ▶ **Analogical inference** can be interpreted as **linear regression** as well:
- ▶ We learn a **regression matrix** from a pair of example instances ('*France* → '*Paris*') and apply it to a test instance ('*Germany*')
- ▶ Say, we have example vectors **[3,5]** and **[7,2]**, and a test instance **[4,4]**, which should map into **[8,1]**
- ▶ Linear regression will 'learn' weights **[2.3, 0.4]**, and for the test instance it will output **[9.2, 0.8]**, pretty close to the desired result **[8,1]**.

## Projecting one model into another: models alignment

### What else?

- ▶ **Cross-model projections** can be used not only for MT.
- ▶ We can study **models sharing the same vocabulary** (if there are reasons to think that the semantics might be different).

### Detecting semantic shifts

- ▶ Suppose we want to **trace how word meaning changed throughout time**;
- ▶ Can be useful for **diachronic linguistics**...
- ▶ ...or for **extracting new events from news texts**.
- ▶ We train distributional models on time-separated corpora:
  - ▶ XIX century
  - ▶ XX century
  - ▶ XXI century
- ▶ Another application is **comparing word meanings in different corpora** (for example, fiction texts and web pages).

25

## Projecting one model into another: models alignment

### Houston, we have a problem

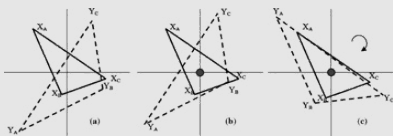
- ▶ But the **resulting dense embeddings are not directly comparable** (unlike traditional sparse count-based vectors)!
- ▶ Because of **stochastic nature** of our modeling, we cannot simply measure cosine distance between  $w_{xix}$ ,  $w_{xx}$  and  $w_{xxi}$  to find whether the meaning has changed.
- ▶ What can be done?

26

## Projecting one model into another: models alignment

### How to compare incomparable?

- ▶ One can do the same: **learn transformation matrix**.
- ▶ One can stick to **comparing ranked lists of nearest associates**, abstracting away from the vector values.
- ▶ One can transform the semantic space **A** so that it matches the semantic space **B** as closely as possible for the shared words, at the same time retaining pair-wise similarities.
- ▶ For example, using the **orthogonal Procrustes** transformation:
  - ▶ basically applying **SVD** to the dot product ( $B \cdot A^T$ )

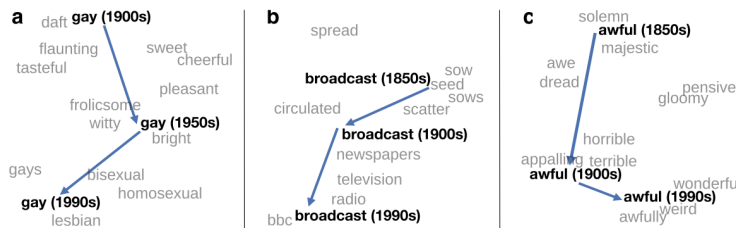


Orthogonal Procrustes (from [Schneider and Borlund, 2007])

27

## Projecting one model into another: models alignment

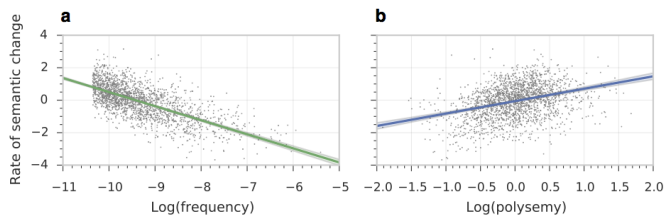
Once the models are aligned, **cosine distances between words from different models (across models boundaries) become meaningful**. Semantic shifts can be traced:



28



Such results can be produced even by comparing nearest associates lists. But models alignment allows uncovering **statistical laws of semantic change on large scale:**



See more at <http://nlp.stanford.edu/projects/histwords/> and in [Hamilton et al., 2016].

**Event extraction**

- ▶ While **updating a model with new data**, continue aligning the previous model with the current one;
- ▶ monitor **changes in semantic relations**;
- ▶ detect current events:
- ▶ 'The country **A** is no more **in adversary relations** to the country **B**'.
- ▶ Lots of opportunities for automatizing the handling of large textual data.



Questions?




INF5820  
 Distributional Semantics: Extracting Meaning from Data  
**Lecture 5**  
**Kings and queens, men and women:**  
**semantic relations between word embeddings**




- 1 Semantic relations as geometrical directions
- 2 Why vector algebra works?
- 3 Possible applications
- 4 Projecting one model into another: models alignment
- 5 In the next week (last lecture)**

What's going on: recent advances and trends in the word embeddings world

- ▶ Discussion on the results of the obligatory assignment.
- ▶ The exam: what to expect?
- ▶ Multilingual word embeddings.
- ▶ Language and vision embeddings.
- ▶ etc...

-  Hamilton, L. W., Leskovec, J., and Jurafsky, D. (2016). Diachronic word embeddings reveal statistical laws of semantic change. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1489–1501. Association for Computational Linguistics.
-  Johnson, M., Schuster, M., Le, Q. V., Krikun, M., Wu, Y., Chen, Z., Thorat, N., Viégas, F., Wattenberg, M., Corrado, G., et al. (2016). Google's multilingual neural machine translation system: Enabling zero-shot translation. *arXiv preprint arXiv:1611.04558*.

-  Levy, O., Remus, S., Biemann, C., and Dagan, I. (2015). Do supervised distributional methods really learn lexical inference relations? In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 970–976. Association for Computational Linguistics.
-  Mikolov, T., Le, Q., and Sutskever, I. (2013a). Exploiting similarities among languages for machine translation. *arXiv preprint arXiv:1309.4168*.
-  Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013b). Distributed representations of words and phrases and their compositionality. *Advances in Neural Information Processing Systems 26*.

-  Pennington, J., Socher, R., and Manning, C. D. (2014). GloVe: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
-  Schneider, J. W. and Borlund, P. (2007). Matrix comparison, part 2: Measuring the resemblance between proximity measures or ordination results by use of the Mantel and Procrustes statistics. *Journal of the American Society for Information Science and Technology*, 58(11):1596–1609.
-  Van der Maaten, L. and Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(2579-2605):85.