



Phrase2vec in Practice

Aerin Kim



At the end of this talk

```
AERINS-MacBook-Pro:bitbucket_BYOR aerin$ python WTB_phrase_similarity.py
Loading the data file... Please wait...
Successfully loaded 3.6 G bin file!
#####
#####
##### WELCOME TO THE PHRASE SIMILARITY CALCULATOR #####
#####
#####
Type the phrase1: How Trump Would Stimulate the U.S. Economy
Type the phrase2: Trump unveils plan to revitalize America's economy
#####
Similarity Score: 0.732695
#####
Type the phrase1: Let's make America great again
Type the phrase2: We will make America prosperous and powerful
#####
Similarity Score: 0.670172
#####
Type the phrase1: Let's make America great again
Type the phrase2: We are going to build a beautiful wall
#####
Similarity Score: 0.438079
#####
```



Download & Install

- <https://drive.google.com/file/d/0B7XkCwpI5KDYNINUTTISS21pQmM/edit>
- `easy_install -U gensim`
- `easy_install numpy`
- `easy_install scipy`
- `Pip install nltk`

Frustration

- Binary Vector (discrete representation)
- Apple = [0,0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0]
- AND
- Fruit = [0,0,0,0,0,0,0,1,0,0,0,0,0,0,0,0,0]
- = 0
- Dimensionality? – very sparse representation

Solution

- Statistical NLP → Co-occurrence matrix

[illegible]



Better Solution

- Predict surrounding words of every word
- You shall be judged by the company you keep

Since he announced his candidacy for the presidency, Trump has filed a number of lawsuits
Would a Trump presidency undo the UN climate change agreement?

↖ These words will represent “presidency” ↗



Word Embedding (Word2Vec)

- Objective function:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j} | w_t)$$

MAXIMIZE the log probability of any context word given the current center word.

One (very) Big Vector Θ

- Θ is the set of ALL parameters in one vector

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix}$$



OBJ function of single window

$$p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)}$$

- Trump announced his candidacy for president as a Republican
- Assuming window size = 1
- First element: $\text{Exp}(U^T(\text{his}) \cdot V(\text{candidacy}))$
- Second element: $\text{Exp}(U^T(\text{for}) \cdot V(\text{candidacy}))$



Result after the Optimization

- Semantically

Famous example:

$\text{Vec}(\text{King}) - \text{Vec}(\text{man}) = \text{Vec}(\text{Queen}) - \text{Vec}(\text{woman})$

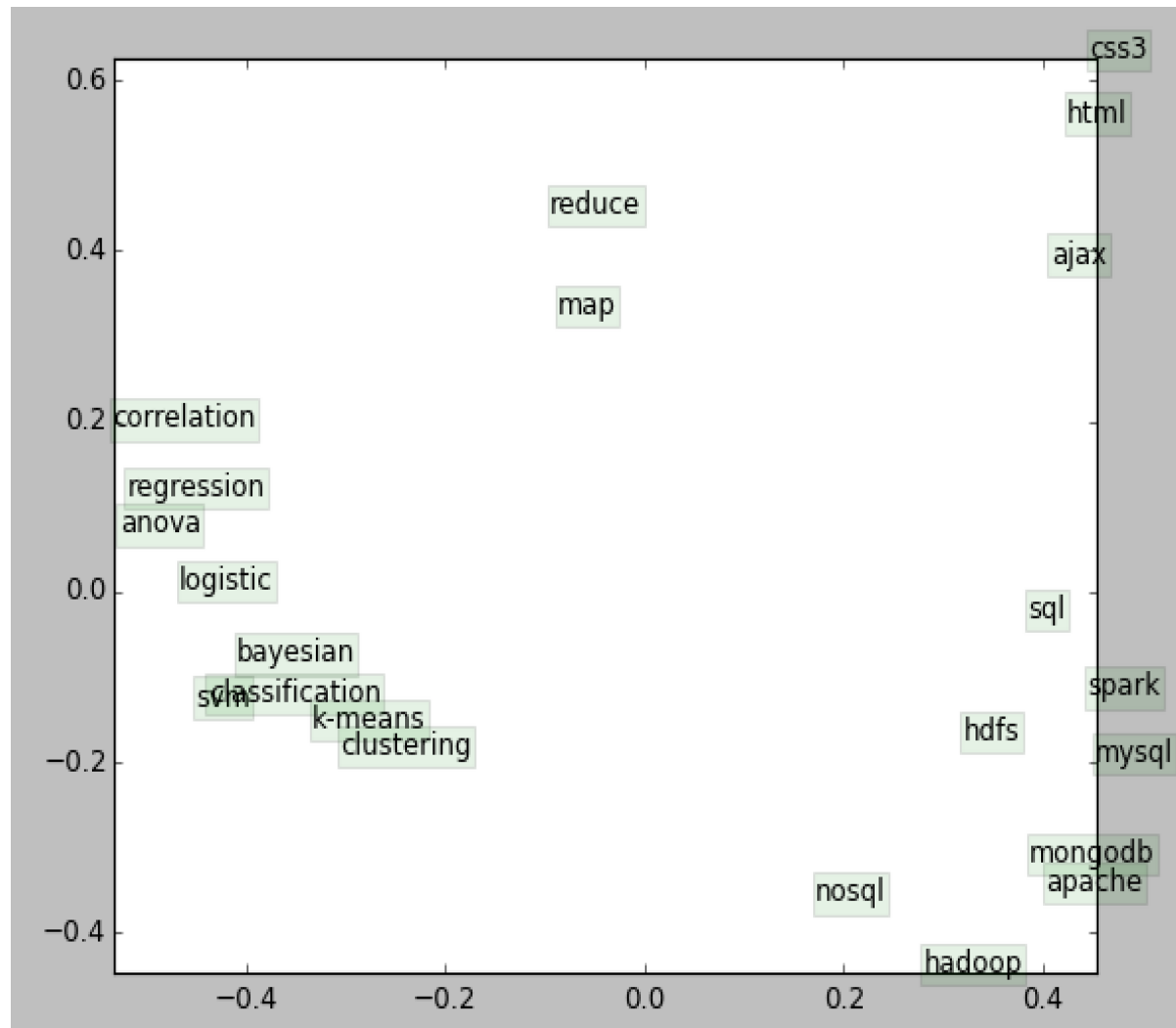
$\text{Vec}(\text{CSS}) - \text{Vec}(\text{Front-end}) = \text{Vec}(\text{Django}) - \text{Vec}(\text{Beck-end})$

- Syntactically

$\text{Vector}(\text{apple}) - \text{vector}(\text{apples}) = \text{vector}(\text{car}) - \text{vector}(\text{cars})$

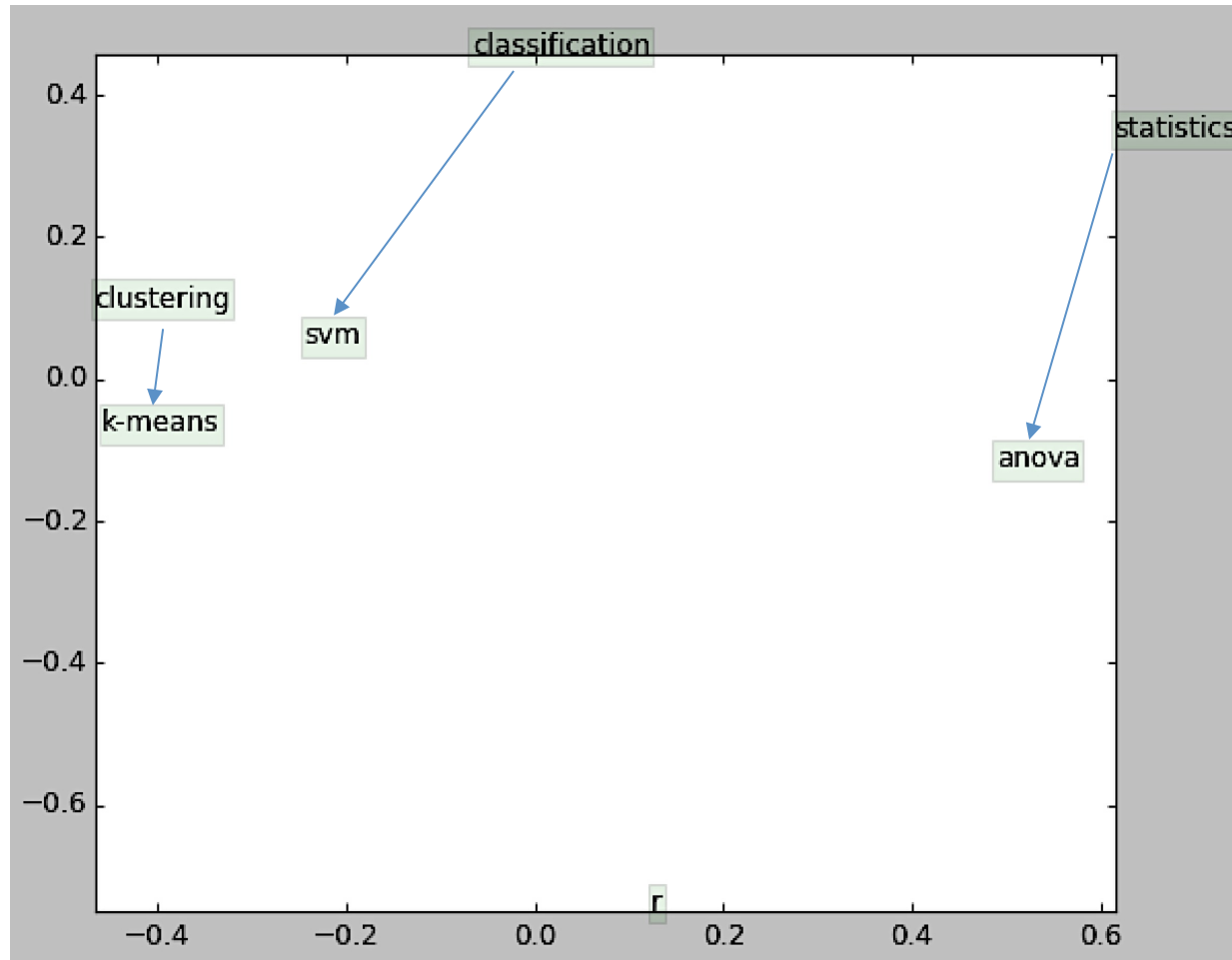
$\text{Vec}(\text{built}) - \text{Vec}(\text{build}) = \text{Vec}(\text{developed}) - \text{Vec}(\text{develop})$

Classifying Data Science Keywords





More sophisticated relationships





Let's make the Phrase Vectors!

- https://bitbucket.org/yunazzang/aiwiththebest_byor

```
AERINs-MacBook-Pro:bitbucket_BYOR aerin$ python WTB_phrase_similarity.py
Loading the data file... Please wait...
Successfully loaded 3.6 G bin file!
#####
##### WELCOME TO THE PHRASE SIMILARITY CALCULATOR #####
#####
Type the phrase1: How Trump Would Stimulate the U.S. Economy
Type the phrase2: Trump unveils plan to revitalize America's economy
#####
Similarity Score: 0.732695
#####
Type the phrase1: Let's make America great again
Type the phrase2: We will make America prosperous and powerful
#####
Similarity Score: 0.670172
#####
Type the phrase1: Let's make America great again
Type the phrase2: We are going to build a beautiful wall
#####
Similarity Score: 0.438079
#####
```