

# Continuous Space Representation (PA153)

Pavel Rychlý

# Problems with statistical NLP

- ▶ many distinct words (items) (from Zipf)
- ▶ zero counts
  - ▶ MLE gives zero probability

$$p(w_3|w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)}$$

- ▶ not handling similarities
  - ▶ some words share some (important) features
  - ▶ *driver, teacher, butcher*
  - ▶ *small, little, tiny*

# Many distinct words

How to solve:

- ▶ use only most frequent ones (ignore outliers)
- ▶ use smaller units (subwords)
  - ▶ prefixes, suffixes
  - ▶ *-er, -less, pre-*

But:

- ▶ we want to add more words
- ▶ *black hole* is not *black* or *hole*
- ▶ even less frequent words are important
  - ▶ *deagrofertizace* from “*The deagrofertization of the state must come.*”
  - ▶ humans process them easily

# Zero counts

How to solve:

- ▶ complicated smoothing strategies
  - ▶ Good-Turing, Kneser–Ney, back-off, ...
- ▶ bigger corpora
- ▶ more data = better estimation

But:

- ▶ sometimes there is no more data
  - ▶ Shakespeare, new research field
- ▶ any size is not big enough

# How big corpus?

## Noun *test*

- ▶ British National Corpus
- ▶ 15789 hits, rank 918
- ▶ word sketches from the Sketch Engine
- ▶ object-of: *pass, undergo, satisfy, fail, devise, conduct, administer, perform, apply, boycott*
- ▶ modifier: *blood, driving, fitness, beta, nuclear, pregnancy*
- ▶ can we freely combine any two from that lists?

# How big corpus?

## Collocations of noun *test*

- ▶ *blood test* in BNC
  - ▶ object-of: *order* (3), *take* (12)
- ▶ *blood test* in enClueWeb16 (16 billion tokens)
  - ▶ object-of: *order* (708), *perform* (959), *undergo* (174), *administer* (123), *conduct* (229), *require* (676), *repeat* (80), *run* (347), *request* (105), *take* (1215)

# How big corpus?

Phrase *pregnancy test* in 16 billion corpus

**pregnancy test** (*noun*) enClueWeb - Sketches freq = 13677 (0.8 per million)  
(test-n filtered by pregnancy)

Constructions	PP_X	955	N_mod	13677	-1.6	and_or	1684	-4.2
wh	PP in-i	175	urine	314	3.07	ultrasound	65	2.25
that_0	PP at-i	150	home	2204	2.68	urine	39	1.31
Vinf_to	PP on-i	139	blood	248	1.36	counseling	44	0.9
	PP for-i	82	serum	53	0.56	condom	23	0.66
object_of	PP after-i	60	at-home	37	0.21	urinalysis	14	0.44
take	PP with-i	55				test	190	0.33
perform	PP from-i	37	AVP_post_mod	431	-2.8	smear	14	0.25
buy	PP within-i	32	prior	27	0.11			
administer	PP to-i	31	AJ_premod	3077	-3.0			
	PP as-i	26	positive	853	3.66	N_premod	1505	nan
	PP before-i	26				kit	317	2.48
						ept	54	1.15

Figure 1: pregnancy test word sketch

# How big corpus?

Phrase *black hole* in 16 billion corpus

## WORD SKETCH

enTenTen [2012]



black hole 30,327×



↩	☰	🔍	✕	↩	☰	🔍	✕	↩	☰	🔍	✕	↩	☰		
object_of				subject_of				modifier				modifie			
accrete	...			accrete	...			supermassive	...			quasar			
orbit	...			evaporate	...			super-massive	...			wormhole			
gape	...			orbit	...			stellar-mass	...			pulsar			
harbor	...			swallow	...			primordial	...			supernova			
collide	...			gobble	...			Supermassive	...			quark			
evaporate	...			collide	...			intermediate-mass	...			astronomer			
harbour	...			devour	...			stellar	...			comet			
yawn	...			lurk	...			massive	...			galaxy			
rotate	...			coalesce	...			Schwarzschild	...			remnant			
encircle	...			radiate	...			miniature	...			gravity			
form	...			emit	...			giant	...						
eject	...			suck	...			galactic	...						

Figure 2: black hole word sketch



# Similarities of words

Distinct words?:

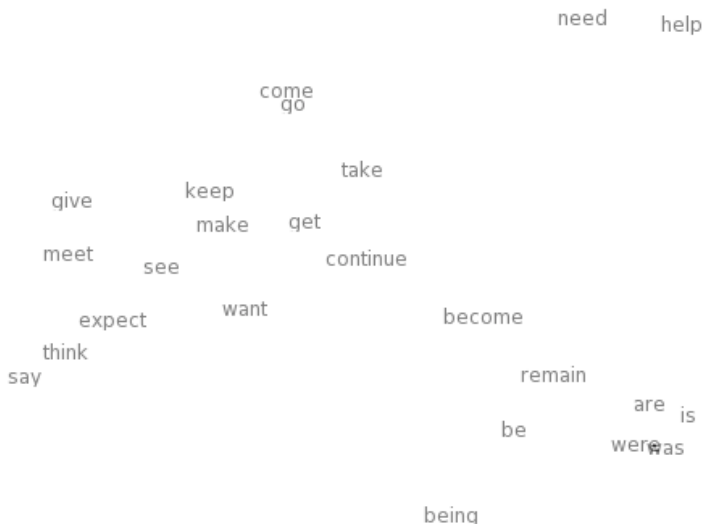
- ▶ *supermassive, super-massive, Supermassive*
- ▶ *small, little, tiny*
- ▶ *black hole, star*
- ▶ *apple, banana, orange*
- ▶ *red, green, orange*
- ▶ *auburn, burgundy, mahogany, ruby*

# Continuous space representation

- ▶ words are not distinct
- ▶ represented by a vector of numbers
- ▶ similar words are *closer* each other
- ▶ more dimensions = more features
  - ▶ tens to hundreds, up to 1000

## Words as vectors

*continue* = [0.286, 0.792, -0.177, -0.107, 0.109, -0.542, 0.349]



# How to create a vector representation

From co-occurrence counts:

- ▶ Singular value decomposition (SVD)
  - ▶ each word one dimension
  - ▶ select/combine important dimensions
  - ▶ factorization of co-occurrence matrix
- ▶ Principal component analysis (PCA)
- ▶ Latent Dirichlet Allocation (LDA)
  - ▶ learning probabilities of hidden variables
- ▶ Neural Networks

# Neural Networks

- ▶ training from examples = supervised training
- ▶ sometimes negative examples
- ▶ generating examples from texts
- ▶ from very simple (one layer) to deep ones (many layers)

## NN training method

- ▶ one training example = (input, expected output) =  $(x, y)$
- ▶ random initialization of parameters
- ▶ for each example:
  - ▶ get output for input:  $y' = NN(x)$
  - ▶ compute loss = difference between expected output and real output:  $loss = y - y'$
  - ▶ update parameters to decrease loss

## Are vectors better than IDs

- ▶ even one hit could provide useful information
- ▶ Little Prince corpus (21,000 tokens)
- ▶ modifiers of “planet”
  - ▶ *seventh, stately, sixth, wrong, tine, fifth, ordinary, next, little, whole*
  - ▶ each with 1 hit
  - ▶ many are *close* together, share a feature

## Simple vector learning

- ▶ each word has two vectors
  - ▶ node vector ( $node_w$ )
  - ▶ context vector ( $ctx_w$ )
- ▶ generate ( $node, context$ ) pairs from text
  - ▶ for example from bigrams:  $w_1, w_2$
  - ▶  $w_1$  is *context*,  $w_2$  is *node*
- ▶ move closer  $ctx_{w_1}$  and  $node_{w_2}$



## Simple vector learning

```
node_vec = np.random.rand(len(vocab), dim) * 2 - 1
ctx_vec = np.zeros((len(vocab), dim))

def train_pair(nodeid, ctxid, alpha):
    global node_vec, ctx_vec
    Nd = node_vec[nodeid]
    Ct = ctx_vec[ctxid]
    corr = (1 - expit(np.dot(Nd, Ct))) * alpha
    Nd += corr * (Ct - Nd)
    Ct += corr * (Nd - Ct)
```

## Simple vector learning

```
for e in range(epochs):
    last = tokIDs[0]
    for wid in tokIDs[1:]:
        train_pair(wid, last, alpha)
        last = wid
        # update alpha
```

# Embeddings advantages

- ▶ no problem in number of parameters
- ▶ similarity in many different directions
- ▶ good estimations of scores
- ▶ **generalization**
  - ▶ learnig for some words generalize to similar words

# Embeddings of other items

- ▶ lemmata
- ▶ part of speech
- ▶ topics
- ▶ any list of items with some structure