## Distributional Lexical Semantics (and word2vec)

## Yoav Goldberg

yogo@cs.biu.ac.il

Bar Ilan University

with some slides by Ido Dagan, Omer Levy, Marco Baroni, Kathrin Erk



#### A lot of text.



#### Need to understand what's being said.

this is where we come in.

▲□▶ ▲圖▶ ▲필▶ ▲필▶ \_ 필

 $\mathcal{O} \mathcal{Q} (\mathcal{P})$ 











#### What does it mean to **understand**?

▲□▶▲□▶▲≡▶▲≡▶ ● ● ● ●

The soup, which I expected to be good, was bad



#### This is called **Syntactic Parsing**.

▲□▶▲□▶▲≡▶▲≡▶ ≡ りへぐ





The gromp , which I furpled to be drogby , was spujky



🔀 u.cs.biu.ac.il/~yogo/ 🛛 🗙

C Index?text=The+gromp%2C+which+I+furpled+to+be+drogby%2C+was+spuj

#### **Greedy parsing to Stanford Dependencies**



#### Can understand structure without understanding words.

But the words are also important.



#### This is not a lecture about parsing.

Today we will focus on the words.

soup was bad

soup was bad soup was awful

soup was bad soup was awful soup was lousy

soup was bad soup was awful soup was lousy soup was abysmal

soup was bad soup was awful soup was lousy soup was abysmal soup was icky

soup was bad soup was awful soup was lousy soup was abysmal soup was icky

chowder was nasty



soup was bad soup was awful soup was lousy soup was abysmal soup was icky

chowder was nasty pudding was terrible

soup was bad soup was awful soup was lousy soup was abysmal soup was icky

chowder was nasty pudding was terrible cake was bad

soup was bad soup was awful soup was lousy soup was abysmal soup was icky

chowder was nasty pudding was terrible cake was bad hamburger was lousy

soup was bad soup was awful soup was lousy soup was abysmal soup was icky

chowder was nasty pudding was terrible cake was bad hamburger was lousy

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

service was poor

soup was bad soup was awful soup was lousy soup was abysmal soup was icky

chowder was nasty pudding was terrible cake was bad hamburger was lousy

service was poor atmosphere was shoddy

◆□▶ ◆□▶ ◆ □▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

soup was bad soup was awful soup was lousy soup was abysmal soup was icky

chowder was nasty pudding was terrible cake was bad hamburger was lousy

service was poor atmosphere was shoddy hammer was heavy

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

soup was bad soup was awful soup was lousy soup was abysmal soup was icky

chowder was nasty pudding was terrible cake was bad hamburger was lousy

service was poor atmosphere was shoddy hammer was heavy

- To the computer, each word is just a symbol, so these are all the same.
- But to us, some are more similar than others.
- We'd like a word representation that can capture that.

#### Use a dictionary?



▲□▶▲□▶▲≡▶▲≡▶ ≡ ∽♀♡

#### Use a dictionary?



#### Doesn't scale.

▲□▶▲□▶▲≡▶▲≡▶ ■ りへで

The distributional Hypothesis

Dr. Baroni saw a hairy little wampinuck sleeping behind a tree

The distributional Hypothesis

Dr. Baroni saw a hairy little wampinuck sleeping behind a tree

The Distributional Hypothesis – Haris, 1954 Words in similar contexts tend to have similar meanings

Firth, 1957 "you should know a word by the company it keeps"

## Distributional Representation of Word Meaning and Similarity

#### The distributional hypothesis in real life McDonald & Ramscar 2001

He filled the wampimuk, passed it around and we all drunk some

We found a little, hairy wampimuk sleeping behind the tree

# What word can appear in the context of all these words?

Word 1: drown, bathroom, shower, fill, fall, lie, electrocute, toilet, whirlpool, iron, gin

Word 2: eat, fall, pick, slice, peel, tree, throw, fruit, pie, bite, crab, grate

Word 3: advocate, overthrow, establish, citizen, ideal, representative, dictatorship, campaign, bastion, freedom

Word 4: spend, enjoy, remember, last, pass, end, die, happen, brighten, relive

# What word can appear in the context of all these words?

Word 1: drown, ba	bathtub
shower, fill, fall, lie	
electrocute, toilet,	2
whirlpool, iron, gin	l

Word 2: eat, fall, pick, slice, peel, tree, throw, fruit, pie, bite, crab, grate

Word 3: advocate, overthrow, establish, citizen, ideal, representative, dictatorship, campaign, bastion, freedom

dayWord 4: spend, enjoy,remember, last, pass, end,die, happen, brighten, relive

# What can you say about word number 5? *Distributional Similarity (2<sup>nd</sup>-order)*


# What can you say about word number 5? *Distributional Similarity (2<sup>nd</sup>-order)*



### Counting context words

- They picked up red
   apples that had fallen
   to the ground
- Eating apples is healthy

### Word count, 3-word

a	be	eat	fall	have	healthy	pick	red	that	up
2	1	2	1	1	1	2	2	1	1

- She ate a red apple)
- Pick an apple

### Distributional semantics

- Comparing two words:
  - Look at all context words for word1
  - Look at all context words for word2
  - How similar are those two context collections in their entirety?
- Compare distributional representations of two words

# How can we compare two context collections in their entirety?

Count how often "apple" occurs close to other words in a large text collection (corpus):

eat	fall	ripe	slice	peel	tree	throw	fruit	pie	bite	crab
794	244	47	221	208	160	145	156	109	104	88

Interpret counts as coordinates:



Every context word becomes a dimension.

# How can we compare two context collections in their entirety?

Count how often "apple" occurs close to other words in a large text collection (corpus):

eat	fall	ripe	slice	peel	tree	throw	fruit	pie	bite	crab
794	244	47	221	208	160	145	156	109	104	88

#### Do the same for "orange":

eat	fall	ripe	slice	peel	tree	throw	fruit	pie	bite	crab
265	22	25	62	220	64	74	111	4	4	8

# How can we compare two context collections in their entirety?

Then visualize both count tables as vectors in the same space:

eat	fall	ripe	slice	peel	tree	throw	fruit	pie	bite	crab
794	244	47	221	208	160	145	156	109	104	88
eat	fall	ripe	slice	peel	tree	throw	fruit	pie	bite	crab
265	22	25	62	220	64	74	111	4	4	8



Similarity between two words as proximity in space

#### Words as Vectors

We can arrange the words in a huge, sparse matrix, where each row is a word, and each column is a context.



now, each word is associated with a (sparse) vector of counts.

### Using distributional models

- Finding (near-)synonyms: automatically building a thesaurus
- Related: use distributional similarity of documents (containing similar words) in Information Retrieval

# Where can we find texts to use for making a distributional model?

- Text in electronic form!
- Newspaper articles
- Project Gutenberg: older books available for free
- Wikipedia
- Text collections prepared for language analysis:
  - Balanced corpora
  - WaC: Scrape all web pages in a particular domain
    - uk, fr, it, de (http://wacky.sslmit.unibo.it/doku.php?id=corpora)
  - ELRA, LDC hold corpus collections
    - For example, large amounts of newswire reports
  - Google n-grams, Google books

### How much text do we need?

- At least: British National Corpus, 100 million words
- Better: add
  - UKWaC (2 billion words)
  - Wikipedia (2 billion words)

### How much text do we need?

• At least:

British National Corpus, 100 million words

- Better: add
  - UKWaC (2 billion words)
  - Wikipedia (2 billion words)
  - Or go even larger with common-crawl (terrabytes)

## What do we mean by "similarity" of vectors?



## Problem with Euclidean distance: very sensitive to word frequency!



## What do we mean by "similarity" of vectors?



# Some counts for "letter" in "Pride and Prejudice". What do you notice?

the	to	of	and	l a	her	she	his	is	was	in	1	hat
102	75	72	56	52	50	41	36	35	34	34		33
had	i	fro	m	you	as	this	mr	for	not	on	be	he
32	28	28		25	23	23	22	21	21	20	18	17
but	e	lizab	eth	W	ith	him	$\mathbf{W}$	hich	by	when	jan	e
17	1	7		16	)	16	16	)	15	14	12	

### Some counts for "letter" in "Pride and Prejudice". What do you notice?

the	to	of	and	a	her	she	his	5	is	was		in		th	at
102	75	72	56	52	50	41	36		35	34		34		33	
had	i	fro	m	you	as	this	m	r	for	not	0	n	be		he
32	28	28		25	23	23	22	2	21	21	2	0	18		17
but	e	lizab	eth	wi	th	him		whi	ch	by	V	vhen	ja	ne	
17	1	7		16		16		16		15	1	4	12	,	

### Some counts for "letter" in "Pride and Prejudice". What do you notice?

the	to	of	and	a	her	she	his	is	was	in		that
102	75	72	56	52	50	41	36	35	34	34		33
had	i	fro	m	you	as	this	mr	for	not	on	be	he
32	28	28		25	23	23	22	21	21	20	18	17
but	e	lizab	eth	) wi	th	him	W	hich	by	when	Jan	e
17	1	7		16		16	16	)	15	14	12	

All the most frequent co-occurring words are function words.

# Some words are more informative than others

• Function words co-occur frequently with <u>all</u> words

– That makes them less informative

- They have much higher co-occurrence counts than content words
  - They can "drown out" more informative contexts

# Using association rather than co-occurrence counts

- Degree of association between target and context:
  - High association: high co-occurrence with "letter", lower with everything else
  - Low association: lots of co-occurrence with all words
- Many ways of implementing this
- For example Pointwise Mutual Information between target a and context b:

$$PMI(a, b) = \log \frac{P(a, b)}{P(a) \cdot P(b)}$$

$$PMI(w,c) = \frac{P(w,c)}{P(w)P(c)}$$

$$P(w,c) = \frac{\#(w,c)}{\sum_{w,c\in D} \#(w,c)}$$

$$P(c) = \frac{\#(\Box, c)}{\sum_{w, c \in D} \#(w, c)}$$

$$P(w) = \frac{\#(w, \Box)}{\sum_{w,c\in D} \#(w, c)}$$

$$PMI(w,c) = \frac{P(w,c)}{P(w)P(c)}$$

matrix entry

$$P(w,c) = \frac{\#(w,c)}{\sum_{w,c\in D} \#(w,c)}$$

$$P(c) = \frac{\#(\Box, c)}{\sum_{w, c \in D} \#(w, c)}$$

$$P(w) = \frac{\#(w, \Box)}{\sum_{w,c\in D} \#(w, c)}$$

$$PMI(w,c) = \frac{P(w,c)}{P(w)P(c)}$$

matrix entry

sum of matrix entries 
$$P(w,c) = \frac{\#(w,c)}{\sum_{w,c\in D} \#(w,c)}$$

$$P(c) = \frac{\#(\Box, c)}{\sum_{w, c \in D} \#(w, c)}$$

$$P(w) = \frac{\#(w, \Box)}{\sum_{w,c\in D} \#(w, c)}$$

$$PMI(w,c) = \frac{P(w,c)}{P(w)P(c)}$$

matrix entry

sum of matrix entries 
$$P(w,c) = \frac{\#(w,c)}{\sum_{w,c\in D} \#(w,c)}$$

sum of matrix column

$$P(c) = \frac{\#(\Box, c)}{\sum_{w, c \in D} \#(w, c)}$$

$$P(w) = \frac{\#(w, \Box)}{\sum_{w,c\in D} \#(w, c)}$$

$$PMI(w,c) = \frac{P(w,c)}{P(w)P(c)}$$
sum of matrix entries 
$$P(w,c) = \frac{\#(w,c)}{\sum_{w,c\in D} \#(w,c)}$$
sum of matrix column 
$$P(c) = \frac{\#(\Box,c)}{\sum_{w,c\in D} \#(w,c)}$$

sum of matrix row  $P(w) = \frac{\#(w, \Box)}{\sum_{w,c\in D} \#(w, c)}$ 

#### We often move from PMI to PPMI (Positive PMI)

PPMI(a, b) = max(PMI(a, b), 0)

### Weighting

Adjusting raw co-occurrence counts:

 $\begin{array}{ccc} \text{bright} & \text{in} \\ \text{stars} & 385 & 10788 & \dots & \leftarrow \text{Counts} \end{array}$ 



### Weighting

Adjusting raw co-occurrence counts:





### Weighting

#### Adjusting raw co-occurrence counts:



Other weighting schemes:

- ► TF-IDF
- Local Mutual Information
- Dice

See Ch4 of J.R. Curran's thesis (2004) and S. Evert's thesis (2007) for surveys of weighting methods

#### Words as Vectors

We can arrange the words in a huge, sparse matrix, where each row is a word, and each column is a context.



### **Dimensionality reduction**

- Vector spaces often range from tens of thousands to millions of context dimensions
- Some of the methods to reduce dimensionality:
  - Select context features based on various relevance criteria
  - Random indexing
  - Following claimed to also have a beneficial smoothing effect:
    - Singular Value Decomposition
    - Non-negative matrix factorization
    - Probabilistic Latent Semantic Analysis
    - Latent Dirichlet Allocation

Words as Vectors

We often apply SVD or similar technique of dimensionality reduction.



### Words as Vectors – It works

Nearest neighbours to dog

- cat
- horse
- ► fox
- pet
- rabbit
- pig
- animal
- mongrel
- sheep
- pigeon

### Words as Vectors – It works

Nearest neighbours to dog

#### 2-word window

- cat
- horse
- ► fox
- pet
- rabbit
- pig
- animal
- mongrel
- sheep
- pigeon

### Words as Vectors – It works

Nearest neighbours to dog

#### 2-word window

- cat
- horse
- ► fox
- pet
- rabbit
- pig
- animal
- mongrel
- sheep
- pigeon

#### 30-word window

- kennel
- puppy
- pet
- bitch
- terrier
- rottweiler
- canine
- cat
- to bark
- Alastian

▲□▶▲□▶▲≡▶▲≡▶ ● ● のへぐ

note that the choice of context has a very important effect.

so does the choice of underlying corpus.

#### word2vec

C C code.google.com/p/word2vec/	Q 🛇 🔂 🔘 🚝
	yoav.goldberg@gmail.com ▼   My favorites ▼   Profile   Sign ou
Word2vec         Tool for computing continuous distr         Project Home       Issues         Summary       People	ributed representations of words.
Project Information	Introduction
<ul> <li>Starred by 694 users <u>Project feeds</u></li> <li>Code license <u>Apache License 2.0</u></li> <li>Labels NeuralNetwork, MachineLearning, NaturalLanguageProcessing, WordVectors, Google</li> </ul>	<ul> <li>This tool provides an efficient implementation of the continuous bag-of-words and skip-gram architectures for computing vector representations of words. These representations can be subsequently used in many natural language processing applications and for further research.</li> <li>Quick start</li> <li>Download the code: svn checkout <a href="http://word2vec.googlecode.com/svn/trunk/">http://word2vec.googlecode.com/svn/trunk/</a></li> <li>Run 'make' to compile word2vec tool</li> <li>Run the demo scripts: ./demo-word.sh and ./demo-phrases.sh</li> </ul>
Members <u>tmiko@gmail.com</u> <u>6 contributors</u>	<ul> <li>For questions about the toolkit, see <a href="http://groups.google.com/group/word2vec-toolkit">http://groups.google.com/group/word2vec-toolkit</a></li> <li>How does it work</li> </ul>
l inks	The word?vec tool takes a text corpus as input and produces the word vectors as output. It first


About 384 results (0.56 seconds)

#### MLMU.cz - Radim Řehůřek - Word2vec & friends (7.1.2015 ...



www.youtube.com/watch?v=wTp3P2UnTfQ Jan 14, 2015 - Uploaded by Marek Modrý I'll go over a particular model published by Google, called word2vec, its optimizations, applications and ...

#### Word2Vec convergence on Vimeo



#### https://vimeo.com/112168934

Nov 18, 2014

This is "Word2Vec convergence" by MaciejLyst on Vimeo, the home for high quality videos and the people who ...

#### Statistical Semantic入門~分布仮説からword2vecまで #1 ...



www.ustream.tv/recorded/43497190 マ Statistical Semantic入門~分布仮説からword2vecまで #1. February 5, 2014 at 7:16pm ...

#### Statistical Semantic入門~分布仮説からword2vecまで #2, PFI ...

www.ustream.tv/recorded/43497424



#### Feb 5, 2014 非常に説明がわかりやすいです!「ゲーミフィケーション入門」と「マー ケティングとスタートアップの話」を見ましたが、どちらも 非常に理解 しやすかった ...

#### GigaOM Show: Samsung watch secrets spilled! B&N's Nook ...





▲□▶▲□▶▲≡▶▲≡▶ ● ● ● ● ●

### From Distributional to Distributed Semantics

#### This part of the talk

- word2vec as a black box
- a peek inside the black box
- relation between word-embeddings and the distributional representation

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

tailoring word embeddings to your needs using word2vecf

#### word2vec

C C code.google.com/p/word2vec/	Q 🛇 🔂 🔘 🗮
	yoav.goldberg@gmail.com ▼   My favorites ▼   Profile   Sign ou
Word2vec         Tool for computing continuous distr         Project Home       Issues         Summary       People	ibuted representations of words.
Project Information	Introduction
<ul> <li>Starred by 694 users <u>Project feeds</u> </li> <li>Code license <u>Apache License 2.0</u> </li> <li>Labels NeuralNetwork, MachineLearning, NaturalLanguageProcessing, WordVectors, Google     </li> </ul>	<ul> <li>This tool provides an efficient implementation of the continuous bag-of-words and skip-gram architectures for computing vector representations of words. These representations can be subsequently used in many natural language processing applications and for further research.</li> <li>Quick start</li> <li>Download the code: svn checkout <a href="http://word2vec.googlecode.com/svn/trunk/">http://word2vec.googlecode.com/svn/trunk/</a></li> <li>Run 'make' to compile word2vec tool</li> </ul>
Members <u>tmiko@gmail.com</u> <u>6 contributors</u>	<ul> <li>Run the demo scripts: ./demo-word.sn and ./demo-phrases.sn</li> <li>For questions about the toolkit, see <u>http://groups.google.com/group/word2vec-toolkit</u></li> <li>How does it work</li> </ul>
Links	The word2vec tool takes a text corpus as input and produces the word vectors as output. It first

#### word2vec



#### word2vec

- dog
  - cat, dogs, dachshund, rabbit, puppy, poodle, rottweiler, mixed-breed, doberman, pig
- sheep
  - cattle, goats, cows, chickens, sheeps, hogs, donkeys, herds, shorthorn, livestock
- november
  - october, december, april, june, february, july, september, january, august, march
- jerusalem
  - tiberias, jaffa, haifa, israel, palestine, nablus, damascus katamon, ramla, safed
- teva
  - pfizer, schering-plough, novartis, astrazeneca, glaxosmithkline, sanofi-aventis, mylan, sanofi, genzyme, pharmacia

#### Word Similarity

Similarity is calculated using cosine similarity:

$$sim(d \vec{o} g, c \vec{a} t) = rac{d \vec{o} g \cdot c \vec{a} t}{||d \vec{o} g|| \, ||c \vec{a} t||}$$

For normalized vectors (||x|| = 1), this is equivalent to a dot product:

$$sim(d \vec{o} g, c \vec{a} t) = d \vec{o} g \cdot c \vec{a} t$$

Normalize the vectors when loading them.

Finding the most similar words to  $\vec{dog}$ 

• Compute the similarity from word  $\vec{v}$  to all other words.

Finding the most similar words to dog

- Compute the similarity from word  $\vec{v}$  to all other words.
- This is a single matrix-vector product:  $W \cdot \vec{v}^{\top}$



< □ ▶ < □ ▶ < 三 ▶ < 三 ▶ < 三 ∽ へ ○

Finding the most similar words to dog

- Compute the similarity from word  $\vec{v}$  to all other words.
- This is a single matrix-vector product:  $W \cdot \vec{v}^{\top}$



- Result is a |V| sized vector of similarities.
- Take the indices of the k-highest values.

Finding the most similar words to dog

- Compute the similarity from word  $\vec{v}$  to all other words.
- This is a single matrix-vector product:  $W \cdot \vec{v}^{\top}$



- Result is a |V| sized vector of similarities.
- Take the indices of the k-highest values.
- FAST! for 180k words, d=300: ~30ms

#### Most Similar Words, in python+numpy code

W,words = load\_and\_normalize\_vectors("vecs.txt")
# W and words are numpy arrays.
w2i = {w:i for i,w in enumerate(words)}

dog = W[w2i['dog']] # get the dog vector

sims = W.dot(dog) # compute similarities

most\_similar\_ids = sims.argsort()[-1:-10:-1]
sim\_words = words[most\_similar\_ids]

▲□▶▲□▶▲≡▶▲≡▶ ④�?

#### Similarity to a group of words

- "Find me words most similar to cat, dog and cow".
- Calculate the pairwise similarities and sum them:

$$W \cdot \vec{cat} + W \cdot \vec{dog} + W \cdot \vec{cow}$$

Now find the indices of the highest values as before.

#### Similarity to a group of words

- "Find me words most similar to cat, dog and cow".
- Calculate the pairwise similarities and sum them:

$$W \cdot \vec{cat} + W \cdot \vec{dog} + W \cdot \vec{cow}$$

- Now find the indices of the highest values as before.
- Matrix-vector products are wasteful. Better option:

$$W \cdot (\vec{cat} + \vec{dog} + \vec{cow})$$

Working with dense word vectors can be very efficient.

▲□▶▲□▶▲■▶▲■▶ ■ りへぐ

Working with dense word vectors can be very efficient.

But where do these vectors come from?



word2vec implements several different algorithms:

#### Two training methods

- Negative Sampling
- Hierarchical Softmax

#### Two context representations

Continuous Bag of Words (CBOW)

< □ ▶ < □ ▶ < 三 ▶ < 三 ▶ < 三 りへ ○

Skip-grams

word2vec implements several different algorithms:

#### Two training methods

- Negative Sampling
- Hierarchical Softmax

#### Two context representations

- Continuous Bag of Words (CBOW)
- Skip-grams

We'll focus on skip-grams with negative sampling.

intuitions apply for other models as well.

- Represent each word as a *d* dimensional vector.
- Represent each context as a d dimensional vector.
- Initalize all vectors to random weights.
- ► Arrange vectors in two matrices, *W* and *C*.



While more text:

#### Extract a word window:

A springer is [ a cow or heifer close to calving ].  $C_1 \quad C_2 \quad C_3 \quad W \quad C_4 \quad C_5 \quad C_6$ 

<ロ> < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

▶ *w* is the focus word vector (row in *W*).

•  $c_i$  are the context word vectors (rows in C).

While more text:

#### Extract a word window:

A springer is [ a cow or heifer close to calving ].  $C_1 \quad C_2 \quad C_3 \quad W \quad C_4 \quad C_5 \quad C_6$ 

Try setting the vector values such that:

 $\sigma(\mathbf{W} \cdot \mathbf{C}_1) + \sigma(\mathbf{W} \cdot \mathbf{C}_2) + \sigma(\mathbf{W} \cdot \mathbf{C}_3) + \sigma(\mathbf{W} \cdot \mathbf{C}_4) + \sigma(\mathbf{W} \cdot \mathbf{C}_5) + \sigma(\mathbf{W} \cdot \mathbf{C}_6)$ 

< □ ▶ < □ ▶ < 三 ▶ < 三 ▶ < 三 りへぐ

is high

While more text:

#### Extract a word window:

A springer is [ a cow or heifer close to calving ].  $C_1 \quad C_2 \quad C_3 \quad W \quad C_4 \quad C_5 \quad C_6$ 

Try setting the vector values such that:

 $\sigma(\mathbf{W} \cdot \mathbf{C}_1) + \sigma(\mathbf{W} \cdot \mathbf{C}_2) + \sigma(\mathbf{W} \cdot \mathbf{C}_3) + \sigma(\mathbf{W} \cdot \mathbf{C}_4) + \sigma(\mathbf{W} \cdot \mathbf{C}_5) + \sigma(\mathbf{W} \cdot \mathbf{C}_6)$ 

#### is high

- Create a corrupt example by choosing a random word w'

   [ a cow or comet close to calving ]

   c1
   c2
   c3
   w'
   c4
   c5
   c6
- Try setting the vector values such that:

 $\sigma(\mathbf{W}' \cdot \mathbf{C}_1) + \sigma(\mathbf{W}' \cdot \mathbf{C}_2) + \sigma(\mathbf{W}' \cdot \mathbf{C}_3) + \sigma(\mathbf{W}' \cdot \mathbf{C}_4) + \sigma(\mathbf{W}' \cdot \mathbf{C}_5) + \sigma(\mathbf{W}' \cdot \mathbf{C}_6)$ 

#### is **low**

# Word2vec Formulation

- *D* : a set of original (correct) word-context pairs
- $\overline{D}$ : a set of corrupt (incorrect) word-context pairs
- *k negative samples* are generated for each correct sample

The model needs to estimate:

P(D=1|w,c): the probability that (w,c) is from D - Should be high for pairs from D, low if P(D=0|w,c) = 1 - P(D=1|w,c):  $\overline{D}$ the probability that (w,c) is from

# Word2vec Formulation (cont.)

Modeling probability as sigmoid of dot product score s(w,c):  $P(D = 1|w,c) = \frac{1}{1 + e^{-s(w,c)}}$ 

<u>Learning goal</u>: find vectors w, c for all words and contexts that maximize log-likelihood of data  $D \cup \overline{D}$ :  $\mathcal{L}(\Theta; D, \overline{D}) = \sum_{(w,c)\in D} \log P(D = 1|w, c) + \sum_{(w,c)\in \overline{D}} \log P(D = 0|w, c)$ 

Negative samples generated by original frequency, or smoothed - deemphasizing frequent words, better in practice:

$$\frac{\#(w)}{\sum_{w'} \#(w')} \quad \text{or} \quad \frac{\#(w)^{0.75}}{\sum_{w'} \#(w')^{0.75}}$$

3

# Word2vec Formulation (cont.)

The *Skip-gram* model assumes independence between the context elements.

Denote  $c_{1:k}$  a context of k elements:

$$\begin{split} P(D = 1 | w, c_i) &= \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{c}_i}} \\ P(D = 1 | w, c_{1:k}) &= \prod_{i=1}^k P(D = 1 | w, c_i) = \prod_{1=i}^k \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{c}_i}} \\ \log P(D = 1 | w, c_{1:k}) &= \log \sum_{i=1}^k \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{c}_i}} \end{split}$$

Very effective in practice, commonly used.

The training procedure results in:

- $w \cdot c$  for **good** word-context pairs is **high**.
- $w \cdot c$  for **bad** word-context pairs is **low**.
- w · c for ok-ish word-context pairs is neither high nor low.

As a result:

- Words that share many contexts get close to each other.
- Contexts that share many words get close to each other.

▲□▶▲□▶▲≡▶▲≣▶ ≣ のへで

At the end, word2vec throws away *C* and returns *W*.

Imagine we didn't throw away C. Consider the product  $WC^{\top}$ 

#### Imagine we didn't throw away C. Consider the product $WC^{\top}$



The result is a matrix *M* in which:

- Each row corresponds to a word.
- Each column corresponds to a context.
- Each cell correspond to w · c, an association measure between a word and a context.



Does this remind you of something?





Does this remind you of something?

Very similar to SVD over distributional representation:

◆□▶ ◆□▶ ◆ 三▶ ◆ 三 ● つへぐ



- A  $V_W \times V_C$  matrix
- Each cell describes the relation between a specific word-context pair



$$\vec{w} \cdot \vec{c} = ?$$

• We **prove** that for large enough *d* and enough iterations



- We **prove** that for large enough d and enough iterations
- We get the word-context PMI matrix



 $V_W$ 

- We **prove** that for large enough *d* and enough iterations
- We get the word-context PMI matrix, shifted by a global constant

$$\frac{d}{W} \approx \frac{V_c}{C} = \frac{V_c}{M^{PMI}} - \log k$$

$$Opt(\vec{w}\cdot\vec{c}) = PMI(w,c) - \log k$$

- SGNS is doing something very similar to the older approaches
- SGNS is factorizing the traditional word-context PMI matrix
- So does SVD!
- Do they capture the same similarity function?
## SGNS vs SVD

Target Word	SGNS	SVD
	dog	dog
	rabbit	rabbit
cat	cats	pet
	poodle	monkey
	pig	pig

## SGNS vs SVD

Target Word	SGNS	SVD
	wines	wines
	grape	grape
wine	grapes	grapes
	winemaking	varietal
	tasting	vintages

## SGNS vs SVD

Target Word	SGNS	SVD
	October	October
	December	December
November	April	April
	January	June
	July	March

## But word2vec is still better, isn't it?

- Plenty of evidence that word2vec outperforms traditional methods
  - In particular: "Don't count, predict!" (Baroni et al., 2014)
- How does this fit with our story?

## The Big Impact of "Small" Hyperparameters

## Hyperparameters

- word2vec is more than just an algorithm...
- Introduces many engineering tweaks and hyperparameter settings
  - May seem minor, but make a big difference in practice
  - Their impact is often more significant than the embedding algorithm's
- These modifications can be ported to distributional methods!

## Hyperparameters

- Preprocessing
- Association Metric
- Postprocessing

## Hyperparameters

- Preprocessing
- Association Metric
- Postprocessing

## Association Metric Hyperparameters

- Since SGNS and PMI are strongly related, we can import 2 of SGNS's hyperparameters to traditional PMI:
- 1. Shifted PMI
- 2. Negative Sampling Smoothing
- Both stem from the negative sampling procedure

## Negative Sampling Smoothing

- Recall that SGNS picks  $w' \sim P$  to form negative (w', c) examples
- Our analysis assumes P is the unigram distribution

$$P(w) = \frac{\#w}{\sum_{w' \in V_W} \#w'}$$

## Negative Sampling Smoothing

- Recall that SGNS picks  $w' \sim P$  to form negative (w', c) examples
- Our analysis assumes P is the unigram distribution
- In practice, it's a **smoothed** unigram distribution

$$P^{0.75}(w) = \frac{(\#w)^{0.75}}{\sum_{w' \in V_W} (\#w')^{0.75}}$$

• This little change makes a big difference

## Negative Sampling Smoothing

- This smoothing has an analogue in PMI
- Replace P(w) with  $P^{0.75}(w)$ :

$$PMI^{0.75}(w,c) = \log \frac{P(w,c)}{P^{0.75}(w)P(c)}$$

• Yields a dramatic improvement with every method on every task

## Experiments & Results

- We compared several methods, while controlling for hyperparameters
  - PPMI, SVD(PPMI), SGNS, GloVe
- Methods perform on-par in most tasks
  - Slight advantage to SVD in word similarity
  - SGNS is better at syntactic analogies
  - SGNS is robust in general
- Negative sampling smoothing accounts for much of the differences observed in "Don't count, predict!"

## Other Hyperparameters

- There are many other hyperparameters that can be investigated
- Perhaps the most interesting one is **the type of context**

## What's in a Context?

## What's in a Context?

- Importing ideas from embeddings improves distributional methods
- Can distributional ideas also improve embeddings?
- Idea: change SGNS's default BoW contexts into dependency contexts



#### Australian scientist discovers star with telescope

## Target Word

#### Australian scientist discovers star with telescope

## Bag of Words (BoW) Context

#### Australian scientist discovers star with telescope

## Bag of Words (BoW) Context

#### Australian scientist discovers star with telescope

## Bag of Words (BoW) Context

#### Australian scientist discovers star with telescope

## Syntactic Dependency Context

#### Australian scientist discovers star with telescope

## Syntactic Dependency Context



## Syntactic Dependency Context



# Embedding Similarity with Different Contexts

Target Word	Bag of Words (k=5)	Dependencies
	Dumbledore	Sunnydale
	hallows	Collinwood
Hogwarts	half-blood	Calarts
(Harry Potter's school)	Malfoy	Greendale
	Snape	Millfield
	Related to Harry Potter	Schools

# Embedding Similarity with Different Contexts

Target Word	Bag of Words (k=5)	Dependencies
	nondeterministic	Pauling
	non-deterministic	Hotelling
Turing	computability	Heting
(computer scientist)	deterministic	Lessing
	finite-state	Hamming
	Related to computability	Scientists

# Embedding Similarity with Different Contexts

Target Word	Bag of Words (k=5)	Dependencies
	singing	singing
	dance	rapping
dancing	dances	breakdancing
(dance gerund)	dancers	miming
	tap-dancing	busking
	Related to dance	Gerunds

# What is the effect of different context types?

- Thoroughly studied in distributional methods
  - Lin (1998), Padó and Lapata (2007), and many others...

#### **General Conclusion:**

- Bag-of-words contexts induce topical similarities
- Dependency contexts induce *functional* similarities
  - Share the same semantic type
  - Cohyponyms
- Holds for embeddings as well

- Same algorithm, different inputs -- very different kinds of similarity.
- Inputs matter much more than algorithm.
- Think about your inputs.

## Peeking into Skip-Gram's Black Box

- In explicit representations, we can **look** at the features and analyze
- But embeddings are a black box!
- Dimensions are latent and don't necessarily have any meaning

## Peeking into Skip-Gram's Black Box

- Skip-Gram allows a peek...
- Contexts are embedded in the same space!
- Given a word w, find the contexts c it "activates" most:

 $\arg\max_{c}(\vec{w}\cdot\vec{c})$ 

## Associated Contexts

Target Word	Dependencies
Hogwarts	students/prep_at <sup>-1</sup>
	educated/prep_at <sup>-1</sup>
	student/prep_at <sup>-1</sup>
	stay/prep_at⁻¹
	learned/prep_at <sup>-1</sup>

## Associated Contexts

Target Word	Dependencies
Turing	machine/nn <sup>-1</sup>
	test/nn <sup>-1</sup>
	theorem/poss <sup>-1</sup>
	machines/nn <sup>-1</sup>
	tests/nn <sup>-1</sup>

## Associated Contexts

Target Word	Dependencies
dancing	dancing/conj
	dancing/conj <sup>-1</sup>
	singing/conj <sup>-1</sup>
	singing/conj
	ballroom/nn

# let's take a step back

- We don't really care about the vectors.
- We care about the **similarity function** they induce.
  - (or, maybe we want to use them in an external task)
- We want similar words to have similar vectors.
- So evaluating on word-similarity tasks is great.
- But what does similar mean?

# many faces of similarity

- dog -- cat
- dog -- poodle
- dog -- animal
- dog -- bark
- dog -- leash
# many faces of similarity

- dog -- cat
- dog -- poodle
- dog -- animal
- dog -- bark
- dog -- leash

- dog -- chair
- dog -- dig
- dog -- god
- dog -- fog
- dog -- 6op

# many faces of similarity

- dog -- cat
- dog -- poodle
- dog -- animal
- dog -- bark
- dog -- leash

- dog -- chair
- dog -- dig

• dog -- god

- edit distance
- same letters

rhyme

- dog -- fog
- dog -- 6op shape

### some forms of similarity look more useful than they really are

- Almost every algorithm you come up with will be good at capturing:
  - countries
  - cities
  - months
  - person names

### some forms of similarity look more useful than they really are

- Almost every algorithm you come up with will be good at capturing:
  - countries

useful for tagging/parsing/NER

- cities
- months
- person names

### some forms of similarity look more useful than they really are

- Almost every algorithm you come up with will be good at capturing:
  - countries
  - cities
  - months
  - person names

useful for tagging/parsing/NER

but do we really want "John went to China in June" to be similar to "Carl went to Italy in February" ??

# there is no single downstream task

- Different tasks require different kinds of similarity.
- Different vector-inducing algorithms produce different similarity functions.
- No single representation for all tasks.
- If your vectors do great on task X, I don't care that they suck on task Y.

#### Choose the correct contexts for your application

- larger window sizes more topical
- dependency relations more functional

#### Choose the correct contexts for your application

< □ ▶ < □ ▶ < 三 ▶ < 三 ▶ < 三 りへぐ

- larger window sizes more topical
- dependency relations more functional
- only noun-adjective relations
- only verb-subject relations

#### Choose the correct contexts for your application

- larger window sizes more topical
- dependency relations more functional
- only noun-adjective relations
- only verb-subject relations
- context: time of the current message
- context: user who wrote the message

▲□▶▲□▶▲≡▶▲≡▶ ≡ 釣�?

#### Choose the correct contexts for your application

- larger window sizes more topical
- dependency relations more functional
- only noun-adjective relations
- only verb-subject relations
- context: time of the current message
- context: user who wrote the message

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ● ⑦ � @

- ▶ ...
- the sky is the limit

# what happens in Hebrew?

(based on my work with Oded Avraham)

# what happens when we look outside of English?

- Things don't work nearly as well.
- Known problems from English become more extreme.
- We get some new problems as well.

## a quick look at Hebrew

### word senses

#### ספר

book(N). barber(N). counted(V). tell!(V). told(V).

#### חומה

brown (feminine, singular) wall (noun) her fever (possessed noun)

- בית שימוש
- ראש עיר •
- יושב ראש •
- שומר ראש •
- בית ספר
- עורך דין •

### multi-word units

### words vs. tokens

### וכשמהבית

and when from the house

### words vs. tokens

### וכשמהבית

#### and when from the house

בצל in shadow

> בצל onion

- nouns, pronouns and adjectives
  --> are inflected for *number* and *gender*
- verbs
  - --> are inflected for *number*, *gender*, *tense*, *person*
- syntax requires *agreement* between
  - nouns and adjectives
  - verbs and subjects

she **saw** a **brown** fox

he saw a brown fence

[fem] [masc] she saw a brown fox

### he **saw** a **brown** fence [masc] [fem]



## inflections and dist-sim

- More word forms -- more sparsity
- But more importantly: *agreement patterns affect the resulting similarities.*

## adjectives

green [m,sg]	green [f,sg]	green [m,pl]
ירוק	ירוקה	ירוקים
כחול	אפורה	<mark>אפורים</mark>
blue [m,sg]	gray [f,sg]	gray [m,pl]
כתום	כתומה	כחולים
orange [m,sg]	orange [f,sg]	blue [m,pl]
צהוב	<b>צהובה</b>	שחורים
yellow [m,sg]	yellow [f,sg]	black [m,pl]
אדום	קסומה	<b>שמימיים</b>
red [m,sg]	magical [f,g]	heavenly [m,pl]

### verbs

(he) walked)	(she) thought	(they) ate
הלך	חשבה	אכלו
הלכו	חושבת	יאכלו
(they) walked	(she) is thinking	(they) will eat
הולך	הרגישה	אוכלים
he) is walking)	(she) felt	(they) are eating
פנה	משוכנעת	אכל
he) turned)	(she) is convinved	(he) ate
התקרב	התעקשה	שתה
he) came closer)	(she) insisted	(they) drank

Doctor [m,sg]	Doctor [f, sg]
רופא	רופאה
<b>פסיכיאטר</b>	סטודנטית
psychiatrist [m,sg]	student [f, sg]
<b>פסיכולוג</b>	נזירה
psychologist [m, sg]	nun [f, sg]
נוירולוג	<b>מלצרית</b>
neurologist [m, sg]	waitress [f, sg]
מהנדס	<b>צלמת</b>
engineer [m, sg]	photographer [f, sg]

### Bias

- The word vectors capture language use.
- Language use may contain biases.
- We need to be aware of it when using the vectors.

sweater סוודר	shirt חולצה
ג׳קט jacket	suit חליפה
חלוק down	גלימה robe
overall	dress
טורבאן turban	קסדה helmet

sweater סוודר	shirt חולצה
ג׳קט jacket	suit חליפה
חלוק down	גלימה robe
overall	שמלה dress
טורבאן turban	קסדה helmet
masculine	feminine

sweater סוודר	shirt חולצה	
ג׳קט jacket	suit חליפה	
חלוק down	גלימה robe	
overall	שמלה dress	
טורבאן turban	קסדה helmet	
masculine	feminine	
completely arbitrary		

# When arbitrary gender does matter

נפל נפלה

# When arbitrary gender does matter

נפל נפלה

נהרג נכבשה

# inflections and dist-sim

- Inflections and agreement really influence the results.
- We get a mix of syntax and semantics.
- Which aspect of the similarity we care about? what does it mean to be similar?
- Need better control of the different aspects.

# inflections and dist-sim

- Work with lemmas instead of words!!
- Sure, but where do you get the lemmas?
- ...for unknown words?
- And what should you lemmatize? everything? somethings? context-dependent?
- Ongoing work in my lab -- but still much to do.

#### Software

#### word2vecf

https://bitbucket.org/yoavgo/word2vecf

- Extension of word2vec.
- Allows saving the context matrix.
- Allows using arbitraty contexts.
  - Input is a (large) file of word context pairs.

◆□▶ ◆□▶ ◆ 三▶ ◆ 三 ● のへで

#### Software

#### hyperwords

https://bitbucket.org/omerlevy/hyperwords/

- Python library for working with either sparse or dense word vectors (similarity, analogies).
- Scripts for creating dense representations using word2vecf or SVD.
- Scripts for creating sparse distributional representations.

#### Software

dissect

http://clic.cimec.unitn.it/composes/toolkit/

- Given vector representation of words...
- ... derive vector representation of phrases/sentences

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Implements various composition methods
## Summary

## **Distributional Semantics**

- Words in similar contexts have similar meanings.
- Represent a word by the contexts it appears in.
- But what is a context?

## Neural Models (word2vec)

- Represent each word as dense, low-dimensional vector.
- Same intuitions as in distributional vector-space models.
- Efficient to run, scales well, modest memory requirement.
- Dense vectors are convenient to work with.
- Still helpful to think of the context types.

## Software

Build your own word representations.