

CS4248: Natural Language Processing

Lecture 6 — Word Embeddings

Announcements

- Project
 - Progress report Deadline: Mar 7, 11.59 pm
 - Template available here: <u>https://bit.ly/cs4248-2320-iu-template</u> (you don't have to use the template, but please use a 16:9 aspect ratio for your slides)
- Assignment 2
 - Reminder Deadline: Mar 14, 11.59 pm
 - Goal: practice manual feature engineering
 - To be fair, only certain technologies already covered are allowed
- Midterm Feedback Survey
 - Deadline: Mar 14, 11.59 pm

Outline

• Motivation

- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics

Motivation

- Recall from Lecture 4: Most NLP algorithms require
 - Numerical input

- → Most common representation: vectors (a.k.a embeddings)
- Standardized input
- So far: Vector Space Model (VSM)
 - Vector representation of <u>documents</u>
 - Document vector for document d_i
 = column term-document matrix (typically using weighting schemes, e.g., tf-idf)

How to represent words as vectors?

Term-Document Matrix

	d ₁	d ₂	d ₃	d ₄	d ₅
car	0	0	0.4	0	0.4
cat	0.22	0.29	0	0	0.22
chase	0.22	0.22	0.22	0	0
dog	0.29	0	0	0.29	0.22
sit	0	0	0	0	0.7
tv	0	0	0.4	0.4	0
watch	0	0	0	0.7	0

Representing Words — **Traditional NLP**

- Words as discrete symbols: One-Hot Encoding
 - Length of vector = size of vocabulary V (number of unique words)
 - Vector values: 1 if dimension reflects word, 0 otherwise

Note: VSM document vector = aggregation over word vectors (with some weighting, e.g., sum, tf-idf)

• Toy example

• $V = \{ \text{dog}, \text{ cat}, \text{ lion}, \text{ bear}, \text{ cobra}, \text{ cow}, \text{ frog}, ... \}$

	w ₁	w ₂	w ₃	w ₄	w ₅	w ₆	w ₇	w ₈	w ₉	 w _{IVI}	
dog	1	0	0	0	0	0	0	0	0	 0	
cat	0	1	0	0	0	0	0	0	0	 0	-
lion	0	0	1	0	0	0	0	0	0	 0	
bear	0	0	0	1	0	0	0	0	0	 0	

Symbolic Representation of Words — Limitation



Symbolic Representation of Words — Limitation

- Problem: No notion of similarity
 - Words are just labels without meaning
 - Different words (syntax) → orthogonal word vectors (even for words with the same/similar meaning)



- Goal: Similar words (meaning) → similar word vectors
 - Word vectors no longer just labels but also encode "some" meaning
 - Improve basically all NLP tasks!

→ What are good embeddings, and how can we find them?

Distributional Hypothesis

"The meaning of a word is its use in the language."

(Wittgenstein, 1953)

"If A and B have almost identical environments [...], we say they are synonyms" (Harris, 1954)

"You shall know a word by the company it keeps."

(Firth, 1957)

Quick Quiz

What do you think "Fasulye" is?

I don't think **Fasulye** is already available on Blu-ray. The best part about **Fasulye** was definitely the cast. We're planning to see **Fasulye** on the next weekend. The director of **Fasulye** clearly knew what he was doing.



Outline

- Motivation
- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics

A Last Look at the Document-Term Matrix (DTM)

- Word vectors derived from DTM
 - Assumption: context of word w
 set of documents containing w
 - In principle, valid word vectors

	d ₁	d ₂	d ₃	d ₄	d ₅
car	0	0	0.4	0	0.4
cat	0.22	0.29	0	0	0.22
chase	0.22	0.22	0.22	0	0
dog	0.29	0	0	0.29	0.22
sit	0	0	0	0	0.7
tv	0	0	0.4	0.4	0
watch	0	0	0	0.7	0

Quick Quiz: Why is this arguably not a good choice for a word vector?

word vector of "cat"

Co-Occurrence Vectors

- Basic idea
 - Context of a word *w* = (small) window of words surrounding w
 - Count how often a word *w* occurs with another (w.r.t. the context of *w*)



Term-Context Matrix — Toy Example

...has shown that the movie rating reflects to overall quality...

...the cast of the show turned in a great performance and...

...is to get nlp data for ai algorithms on a large scale...

...only with enough data can ai find reliable patterns to be effective...



	aardvark	rating	story	data	cast	result	
movie	0	2	4	0	1	0	
show	0	6	3	0	2	1	
nlp	0	0	1	3	0	4	
ai	0	1	0	5	0	2	

movie ≈ show

Term-Context Matrix

- Problems with raw counts: Often very skewed
 - e.g., "the" and "of" are very frequent, but typically not very discriminative

→ Alternative: Pointwise Mutual Information (PMI)

• Do words w_i and w_j co-occur more than if they were independent?

$$PMI(w_i, w_j) = \log_2 \frac{P(w_i, w_j)}{P(w_i)P(w_j)} \longleftarrow$$

PMI can be < 0, but no good intuition for negative values for word vectors

→ Positive Pointwise Mutual Information (PPMI)

$$PPMI(w_i, w_j) = max\left(\log_2 \frac{P(w_i, w_j)}{P(w_i)P(w_j)}, 0\right)$$

PPMI Matrix — Toy Example

Assume this is the complete term-context matrix

	rating	story	data	cast	result
movie	2	4	0	1	0
show	6	3	0	2	1
nlp	0	1	3	0	4
ai	1	0	5	0	2

$$P(w = movie, c = cast) = 1/35 = 0.03$$

 $P(w = movie) = 7/35 = 0.2$
 $P(c = cast) = 3/35 = 0.09$

	rating	story	data	cast	result
movie	0.15	1.32	0	0.74	0
show	0.96	0.13	0	0.96	0
nlp	0	0	0.71	0	1.32
ai	0	0	1.45	0	0.32

$$PPMI(w = movie, c = cast) = \log_2 \frac{0.03}{0.09 \cdot 0.2} = 0.74$$

	rating	story	data	cast	result	P(w)
movie	0.06	0.11	0	0.03	0	0.20
show	0.17	0.09	0	0.06	0.03	0.34
nlp	0	0.03	0.09	0	0.11	0.23
ai	0.03	0	0.14	0	0.06	0.23
P(context)	0.26	0.23	0.23	0.09	0.20	

PPMI Word Vectors — Discussion

- Various refinements to handle (very) rare words
 - Raise context probabilities
 - Use Add-1 Smoothing

similar effects

- Consideration: Sparsity
 - Matrix is of size $|V| \times |V|$ (|V| typically between 20k and 50k)
 - PPMI word vectors are very sparse (most vector entries are 0)

Outline

- Motivation
- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics

Why Dense Word Vectors?

- Important practical benefits of dense vectors
 - More convenient features: less weights to tune, lower risk of overfitting
 - Tend to generalize better than features derived from counts
 - Tend to better capture synonymy than sparse vectors

Example sparse vector: [0 0 0 0 0.8 0 0 ... 0 0 0 0 1.2 0 ... 0 0] (e.g., for the word "actor") "movie" "film"

Each word represents a distinct dimension; fails to capture similarity between words

• Dense vector in practice

- Common dimensions: 100 to 1,000 entries
- Most to all vector elements are non-zero

Dense Word Vectors — Intuition

- Toy example: custom encoding with 2 dimensions
 - Each dimension represent a property shared between words

	furry	dangerous
og	0.90	0.15
at	0.85	0.10
on	0.80	0.95
ear	0.85	0.90
obra	0.0	0.80
ow	0.75	0.10
rog	0.05	0.05

Xion

bear

Xdog

×cat

furry

Dense Word Vectors — Intuition



Using suitable similarity metric

• $sim(w_{lion}, w_{bear}) = 1.54$

•
$$sim(w_{lion}, w_{cow}) = 0.70$$

This notion of similarity between words is what we are after!

- Problems with custom encoding
 - How to decide on the dimensions?
 - How to decide on the values?

- Manual assignment simply impractical/impossible!
- → Need for automated methods

Outline

- Motivation
- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics

Basic Approaches

- Popular alternatives (but not covered here)
 - Singular Value Decomposition (SVD; matrix factorization)
 - Brown Clustering
- Neural Network-based Methods
 - Inspired by (Neural) Language Models
 - Learn embeddings as part of the process of word prediction
 - Typically fast & easy to train
 - In the following: Word2Vec

Word2Vec encompasses 2 network architectures: **CBOW & Skip-gram**

Word2Vec: CBOW & Skip-Gram

Continuous Bag of Words (CBOW)



Skip-gram

CBOW — Predicting a Word from Context



CBOW — Predicting a Word from Context



Skip-Gram — Predicting a Word from Context



Skip-Gram — Predicting a Word from Context



Outline

- Motivation
- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics

Word2Vec — Basic Setup (CBOW & Skip-gram)

• Define two matrices

- $\mathcal{V} \in \mathbb{R}^{|V| \times d}$ input embedding matrix
- $\mathcal{U} \in \mathbb{R}^{d imes |V|}$ output embedding matrix
- Given a word w_i , let $v_i \in \mathcal{V}$ and $u_i \in \mathcal{U}$ be the input and output embedding of w_i





Word2Vec — Basic Setup (CBOW & Skip-gram)

- Prediction task: 1 input word w_i , 1 output word w_o (both as one-hot vectors)
 - $w_i^T \cdot V \Rightarrow v_i$ (note: one-hot vector multiplied with a matrix is just a row "lookup")
 - softmax $(v_i^T \cdot U) \rightarrow$ Probability $P(w|w_i)$ for all $w \in V$



Word2Vec — CBOW (window size m=2)





Word2Vec — Skip-Gram (window size m=2)

Training Objective — Loss Function



Training Objective — Intuition

• Main objective for Skip-gram (for CBOW it's just "mirrored")

$$P(u_{c+j} \mid v_c) = \frac{\exp(u_{c+j}^T \cdot v_c)}{\sum_{j=1}^{|V|} \exp(u_j^T \cdot v_c)}$$

$$P(u_{c+j}|v_c)$$
 larger
 \square
Dot product $u_{c+j}^T \cdot v_c$ is higher
 \square
Vectors u_{c+j} and v_c are more similar

- Goal of training
 - Make vectors of center words close to vectors of their context words
 - → Vectors of words with similar contexts will be close
 Intermediate goal

Getting the Word Embeddings

- Learning ${\mathcal U}$ and ${\mathcal V}$
 - Minimize loss using Gradient Descent (or similar optimization technique)
 - \blacksquare All trainable / learnable parameters are in ${\mathcal U}$ and ${\mathcal V}$
- Which are the final embeddings? (recall, both matrices contain embeddings for each word)
 - Use only \mathcal{U}
 - $\bullet \quad \text{Use only } \mathcal{V}$
 - \blacksquare Use average of $\, \mathcal{U} \, \, \mbox{and} \, \, \mathcal{V} \,$

Outline

- Motivation
- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics
Word2Vec — Real-World Example (but on a very small scale)

• Setup & training

- 50k movie reviews from IMDB (Source: <u>https://ai.stanford.edu/~amaas/data/sentiment/</u>)
- Dataset preparation (window size *m*=2)

Now-word removal, lowercase, lemmatization, consider only 20k most frequent words

Treat all whole dataset as a single string (i.e., context windows cross sentence boundaries)

"watching funny movie on netflix"

Χ	У			
watch funny on netflix	movie			
→ 1 CBOW samp	ole			

x	У
movie	watch
movie	funny
movie	on
movie	netflix

PyTorch implementation of CBOW and Skip-gram

5	class CE	BOW(nn.Module):
6		and the second
7	def	<pre>init(self, vocab_size, embed_dim):</pre>
8		<pre>super(CBOW, self)init()</pre>
9		<pre>self.embeddings = nn.Embedding(vocab_size, embed_dim)</pre>
10		<pre>self.linear = nn.Linear(embed_dim, vocab_size)</pre>
11		And a second
12	def	forward(self, contexts):
13		<pre>x = self.embeddings(contexts)</pre>
14		x = x.mean(axis=1)
15		<pre>x = self.linear(x)</pre>
16		<pre>x = F.log_softmax(x, dim=1)</pre>
17		return x

6		
7	def	<pre>init(self, vocab_size, embed_dim):</pre>
8		<pre>super(Skipgram, self). init ()</pre>
9		<pre>self.embeddings = nn.Embedding(vocab_size, embed_dim)</pre>
LO		<pre>self.linear = nn.Linear(embed dim, vocab size)</pre>
11		
12	def	forward(self, inputs):
13		<pre>x = self.embeddings(inputs)</pre>
14		<pre>x = self.linear(x)</pre>
15		<pre>x = F.log_softmax(x, dim=1)</pre>
16		return x

Word2Vec — Real-World Example





Word2Vec — Real-World Example



Outline

- Motivation
- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics

Word2Vec — Negative Sampling

- Observation regarding training performance
 - Basic training objective includes a Softmax
 - Normalization over entire(!) vocabulary (to ensure a valid probability distribution of outputs)
 - Each sample potentially tweaks all(!) weights (all elements in embedding matrices V and U)

$$P(u_{c+j} \mid v_c) = \frac{\exp(u_{c+j}^T \cdot v_c)}{\sum_{j=1}^{|V|} \exp(u_j^T \cdot v_c)}$$

- → Negative Sampling (in the following: SGNS Skip-Gram with Negative Sampling)
 - Convert training from a word prediction task to a binary classification task

Word2Vec — Negative Sampling

• Negative sampling — illustration

"watching funny movie on netflix"

■ Window size *m*=2



Data samples for (basic) Skip-gram



INS	
у	<i>k</i>
1	
1	2m positive samples
1	
1	
0	
0	
0	
0	2m la nogativo camplos
0	
0	
0	
0	
	y 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0

Does this look familiar?

In-Lecture Activity (+10 min break)

- Question: Which negative samples are arguably more useful?
 - Post your solution to the Canvas Discussion

(movie, with)	0
(movie, nlp)	0
(movie, on)	0
(movie, barely)	0
(movie, cluster)	0
(movie, morpheme)	0
(movie, traffic)	0
(movie, the)	0

Word2Vec — Negative Sampling

- Selection of negative samples
 - Essentially at random (error of picking a "wrong" negative sample is negligible)
 - To increase probability of rare words: Sampling using (*Q*-)weighted unigram frequency



SGNS — Training Objective

$$\left| egin{array}{l} P(+|c,m) = rac{1}{1+\exp{\left(-u_m^T v_c
ight)}}
ight.$$

Let's assume this a given mini batch *B*

(movie, watch)	1	
(movie, funny)	1	
(movie, on)	1	$\int D_{pos}$
(movie, netflix)	1	J
(movie, cake)	0	
(movie, nlp)	0	
(movie, soccer)	0	
(movie, barely)	0	
(movie, cluster)	0	$\int D_{neg}$
(movie, morpheme)	0	
(movie, traffic)	0	
(movie, nimble)	0	J

$$L = -\log\left[\prod_{(c,m)\in B_{pos}}P(+|c,m)\cdot\prod_{(c,m)\in B_{neg}}P(-|c,m)
ight]$$

$$=-\log\left[\prod_{(c,m)\in B_{pos}}P(+|c,m)\cdot\prod_{(c,m)\in B_{neg}}(1-P(+|c,m))
ight]$$

$$= - \left[\sum_{(c,m) \in B_{pos}} \log P(+|c,m) + \sum_{(c,m) \in B_{neg}} \log \left(1 - P(+|c,m)
ight)
ight]$$

SGNS — Training Objective

1	1	$1 + e^{-a}$	1	e^{-a}	1
1 -	$1 + e^{-a}$	$=\frac{1}{1+e^{-a}}$	$\frac{1}{1+e^{-a}}$	$=\frac{1}{1+e^{-a}}=$	$\overline{1+e^a}$

Let's assume this a given mini batch *B*

(movie, watch)	1	
(movie, funny)	1	
(movie, on)	1	$\int D_{pos}$
(movie, netflix)	1	J
(movie, cake)	0	
(movie, nlp)	0	
(movie, soccer)	0	
(movie, barely)	0	
(movie, cluster)	0	$\int B_{neg}$
(movie, morpheme)	0	
(movie, traffic)	0	
(movie, nimble)	0	J

$$L = - \left[\sum_{(c,m) \in B_{pos}} \log rac{1}{1 + \exp{(-u_m^T v_c)}} + \sum_{(c,m) \in B_{neg}} \log{(1 - rac{1}{1 + \exp{(-u_m^T v_c)}})}
ight]$$

$$=-\left[\sum_{(c,m)\in B_{pos}}\lograc{1}{1+\exp\left(-u_m^Tv_c
ight)}+\sum_{(c,m)\in B_{neg}}\lograc{1}{1+\exp\left(u_m^Tv_c
ight)}
ight]$$

$$= - \left[\sum_{(c,m) \in B_{pos}} \log \sigma(u_m^T v_c) + \sum_{(c,m) \in B_{neg}} \log \sigma(-u_m^T v_c)
ight]$$

SGNS — Parameters

- Sampling method to generate negative samples
 - e.g., subsampling to ignore very frequent words
- Number k of negative samples (per positive sample)
 - $2 \le k \le 5$ for large text
 - $5 \le k \le 20$ for small text.

Outline

- Motivation
- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics

Word Embeddings — (Desired) Properties

● Vector differences yield semantic relationships → linear substructures



Word Embeddings — (Desired) Properties

• Other meaningful linear substructures



Note: Getting these semantic relationships prohibit the use of stemming to lemmatization!

Outline

- Motivation
- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics

• Data preprocessing steps

- Choice of tokenizer
- Case-folding (yes/no)
- Stemming/lemmatization (yes/no)
- Stopword removal (yes/no)
- Cross-sentence contexts (yes/no)

• Parameters

- Window size *m*
- Number of negative samples (e.g., 2mk for Skip-gram)

- Unable to represent phrases
 - "New York", "snow cat", "ice cream", "land mine", "hot dog", "disc drive", etc.
- Unable to handle polysemy and part of speech
 - Polysemy: multiple meanings for the same word
 - Part of speech: the same word used as noun, verb, or adjective

```
1 word2vec_wikipedia.wv.most_similar("light", topn=10)
[('lights', 0.5668156743049622),
('illumination', 0.5530915260314941),
('glow', 0.5415263175964355),
('sunlight', 0.5396571159362793),
('lamp', 0.5024341344833374),
('flame', 0.48772770166397095),
('lamps', 0.47849947214126587),
('dark', 0.4764614701271057),
('luminous', 0.4740492105484009),
('lighting', 0.47177615761756897)]
```

- Distributional representation does not capture all semantics
 - Common case: words with opposite polarity (sentiment) → Why?



• Embeddings dependent on application / dataset

Dataset: Wikipedia

```
1 word2vec_wikipedia.wv.most_similar("house")
[('mansion', 0.7079392075538635),
('cottage', 0.6541333198547363),
('farmhouse', 0.6259987950325012),
('barn', 0.5747625827789307),
('bungalow', 0.5724436044692993),
('townhouse', 0.567018449306488),
('houses', 0.5506472587585449),
('parsonage', 0.5426527857780457),
('tavern', 0.5370140671730042),
('summerhouse', 0.5307810306549072)]
```

Dataset: Google News

1	<pre>word2vec_googlenews.most_similar("house")</pre>
[('	houses', 0.7072390913963318),
('	bungalow', 0.6878559589385986),
('	apartment', 0.6628996729850769),
('	bedroom', 0.6496936678886414),
('	townhouse', 0.6384080052375793),
('	residence', 0.6198420524597168),
('	mansion', 0.6058192253112793),
('	farmhouse', 0.5857570171356201),
('	duplex', 0.5757936239242554),
('	appartment', 0.5690325498580933)]

Word2Vec in Practice (credits to Google <u>https://code.google.com/archive/p/word2vec/</u>)

- Architecture:
 - Skip-gram: slower, better for infrequent words
 - CBOW (fast)
- Training:
 - Hierarchical softmax: better for infrequent words
 - Negative sampling: better for frequent words, better with low dimensional vectors
- Sub-sampling of frequent words
 - can improve both accuracy and speed for large data sets (useful values are in range 1e-3 to 1e-5)
- Dimensionality of the word vectors
 - usually more is better, but not always
- Context (window) size
 - For skip-gram usually around 10, for CBOW around 5

Outline

- Motivation
- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics

NLP Ethics — Why this Topic?

- Real-world NLP applications have real-world impacts
 - Wide range of very common and popular services based on NLP we all use (online search / information retrieval, machine translation, chatbots, text summarization, etc.)
 - Many NLP applications making decisions affecting people's lives (e.g., what content we see — or don't see — on social media → Think about it: What is needed for that?)
- Language does not exist in isolation
 - Natural language → humans gave, give, and will give meaning to written and spoken word
 - Humans have different knowledge, beliefs, biases, preconceptions, etc.

"The common misconception is that language has to do with words and what they mean. It doesn't. It has to do with people and what they mean."

(Clark & Schober, 1982)

Dual Use and Adversarial NLP (a.k.a. Malicious NLP)

(Arguably) Useful

(Potentially) Harmful

Authorship attribution

(NLP/AI meets Linguistic Forensics)

Text Generation

User content analysis

Censorship

Historical documents, ransom notes, plagiarism

Authors of political dissent

Fake content detection

Personalized content and recommendations

Fake content generation

Privacy intrusion, "over-personalized" content (e.g., echo chambers / filter bubbles)

Censorship evasion

More robust censorship

(Unintentionally) Harmful NLP

- Biased humans/society → biased data → biased model
 - NLP models are very likely to pick up such biases
 - Example: machine translation

												0
	Trans	late							Turn	off ins	tant tran	slation
	Bengali	English	Hungarian	Detect language	•	€.	English	Spanish	Hungarian	*	Trans	slate
Hungarian is an example of a gender-neutral language	ő egy ő egy ő egy ő egy ő egy ő egy ő egy	ápoló. tudós. mérnöl pék. tanár. esküvő vezérig	k. ji szervezd jazgatója.	5.		×	she's he is a she's he is a She is he's a he's a	a nurse a scienti an engir a baker a teache s a wede CEO.	st. neer. er. ding orgar	nizer.		
	•)	· •			110)/5000						

Biased NLP Technologies

- Biases identified in many NLP tasks/technologies
 - Bias in Word Embeddings (Bolukbasi et al. 2017; Caliskan et al. 2017; Garg et al. 2018)
 - Bias in Language identification (Blodgett & O'Connor. 2017; Jurgens et al. 2017)
 - Bias in Visual Semantic Role Labeling (Zhao et al. 2017)
 - Bias in Natural Language Inference (Rudinger et al. 2017)
 - Bias in Coreference Resolution (Rudinger et al. 2018; Zhao et al. 2018)
 - Bias in Automated Essay Scoring (Amorim et al. 2018)
 - Bias in Machine Translation (Prates et al. 2019)

Bias in Word Embeddings

• Recall: Desired properties of word embeddings

> $v(king) - v(man) + v(woman) \approx v(queen)$ $v(France) - v(Paris) + v(Berlin) \approx v(Germany)$

 $v(programmer) - v(man) + v(woman) \approx v(homemaker)$



Bias in Sentiment Analysis

- Simple sentiment analysis
 - Sentiment lexicon + word embeddings (replace pos/neg words with their pretrained embedding)
 - Train a model to predict word sentiments (input: word vectors; target: sentiment label)

```
Looks alright
text_to_sentiment("this example is pretty cool")
3.889968926086298
text_to_sentiment("this example is okay")
2.7997773492425186
text_to_sentiment("meh, this example sucks")
-1.1774475917460698
```

Yeah, well...

```
text_to_sentiment("Let's go get Italian food")
```

```
2.0429166109408983
```

```
text_to_sentiment("Let's go get Chinese food")
```

```
1,4094033658140972
```

text_to_sentiment("Let's go get Mexican food")

```
0.38801985560121732
```

Bias in Sentiment Analysis

- Here: Looking at common first names
 - more specifically: word vectors of names

	sentiment	group
mohammed	0.834974	Arab/Muslim
alya	3.916803	Arab/Muslim
terryl	-2.858010	Black
josé	0.432956	Hispanic
luciana	1.086073	Hispanic
hank	0.391858	White
megan	2.158679	White



In-Lecture Activity (5 min)

- Question: Why would you go about avoiding or removing biases in word embeddings?
 - Just brainstorm possible approaches but also try to briefly justify them
 - Post your idea(s) to the Canvas Discussion

```
text_to_sentiment("Let's go get Italian food")
2.0429166109408983
text_to_sentiment("Let's go get Chinese food")
1.4094033658140972
text_to_sentiment("Let's go get Mexican food")
0.38801985560121732
```

Towards Debiasing — Measuring Bias

- How to check if your word embeddings contain biases?
 - Example: gender biases
 - Approach: find nearest occupations to "he" and "she"
 (e.g. the word vector for "homemaker" is very close to the word vector of "she")

common female occupations

VS.

common male occupations

Extreme he
1. maestro
2. skipper
3. protege
4. philosopher
5. captain
6. architect
7. financier
8. warrior
9. broadcaster
10. magician

Towards Debiasing — Measuring Bias

- Identifying biases using analogies
 - Approach: Find pairs of words (x, w) that minimize:

```
to ensure that x and y are semantically similar "enough"
```

cosine(v(she) - v(he), v(x) - v(y)) if $||v(x) - v(y)|| \le \delta$

Gender	stereotype	she-he	analogies	
		10.000		20.00

sewing-carpentry	registered nurse-physician	housewife-shopkeeper
nurse-surgeon	interior designer-architect	softball-baseball
blond-burly	feminism-conservatism	cosmetics-pharmaceuticals
giggle-chuckle	vocalist-guitarist	petite-lanky
sassy-snappy	diva-superstar	charming-affable
volleyball-football	cupcakes-pizzas	lovely-brilliant

Gender appropriate she-he analogies

queen-king	sister-brother	mother-father
waitress-waiter	ovarian cancer-prostate cancer	convent-monastery

Towards Debiasing — Identify Gender Subspace

- Which "direction" of the 300-dim embedding space encodes gender?
 - Approach: Pick top pairs of gender words → interpret difference as direction(s) of gender
 - Problem: Language is "messy" → differences point not exactly in the same direction(s)



Use top PCs (principal components, vectors in the embedding space) as gender subspace

Towards Debiasing — Identify Gender-Neutral Words

- Split vocabulary into gender-neutral words (*N*) and gender-specific words (*S*)
 - Manually identify a set of gender-specific words → training data (we are interesting in *N*, but there are much more gender-neutral words, so that's easier)
 - Train a binary classifier to predict if a word (vector) is gender-neutral or gender-specific
 - Generate *N* and *S* by predicting the class for each word (vector)

Towards Debiasing — Embedding Transformation



- Goal: Transform embeddings to remove gender bias
 - $\hfill\blacksquare$ Idea: Find a transformation matrix $\,T$ the transforms the embedding matrix $\,W$
 - Approach: Find T by minimizing

 $m_T^{in} \| (TW)^T (TW) - (W^TW) \|_F^2 + \lambda \| (TN)^T (TB) \|_F^2$ inner products of inner products of transformed transformed embeddings after embeddings before gender-neutral gender transformation transformation word vectors subspace Keep difference (here: Frobenius norm) small ---Minimal if gender subspace removed from vectors of gender-neutral words preserve the original embeddings as much as possible!

Outline

- Motivation
- Sparse Word Embeddings
 - Co-occurrence Vectors
 - Discussion & Limitations

• Dense Word Embeddings

- Basic Idea
- Word2Vec (CBOW & Skipgram)
- Negative Sampling
- Basic Properties
- Practical Considerations & Limitations

• NLP Ethics

Summary

• (Dense) word embeddings

- Core component of many to most NLP applications (particularly applications based on neural network solutions)
- Dense vectors automatically learned from data
- Support a intuitive notion of similarity between words (words with similar meanings → words have similar word vectors)

• NLP ethics

- Like with all technologies: use and misuse (accidentally or maliciously)
- Focus here: biases in word embeddings (due to biases in the data, due to biases in society)
Pre-Lecture Activity for Next Week

- Assigned Task
 - Post a 1-2 paragraph answer to the following question in the [Pre-Lecture] discussion
 - Watch the panel discussion of the recent Generative AI, Free Speech, & Public Discourse (Available on YouTube: <u>https://www.youtube.com/live/BBhewsinQwU?si=ODkIpYjqOCLZh8xD&t=8659</u>)



Relate an opinion by one of the panelists to the lecture materials presented today. Why did you pick this opinion to highlight and what is your own opinion on it?

Side notes:

- We will talk about this in the next lecture
- You can just copy-&-paste others' answers or use AI Tools, but please consider your original stance and opinion too

Solutions to Quick Quizzes

- Slide 9: A (or maybe D)
 - It should be clear from the example sentences (i.e., the surrounding words)
 - Arguably best indicators: *cast*, *Blu-ray*, *director*
- Slide 11:
 - These word vectors do not reflect the Distributional Hypothesis
 - Here: 2 words a similar if they appear in many shared documents
- Slide 38
 - For example, the words *"scary"* and *"funny"* very similar → bad for sentiment analysis
 - scary" and "funny" are similar because they are often used in the same context (movie description)
 - Distributional Hypothesis does not capture all aspects of a word (e.g., polarity)