CS378: Natural Language Processing Lecture 8: Bias in Embeddings, Multilingual Embeddings





Assignment 2 due in one week

- Survey on Instapoll

Bias in embeddings response due next Tuesday (submit on Canvas)



Recap





Playing around with embeddings



Using Word Embeddings

- Approach 1: learn embeddings as parameters from your data
 - Often works pretty well
- Approach 2: initialize using GloVe, keep fixed
 - Faster because no need to update these parameters
- Approach 3: initialize using GloVe, fine-tune
 - Works best for some tasks

Beyond Word Embeddings



Same as SGNS, but break words down into n-grams with n = 3 to 6

where:

- 3-grams: <wh, whe, her, ere, re>
- 4-grams: <whe, wher, here, ere>,
- 5-grams: <wher, where, here>,
- 6-grams: <where, where>
- Replace $w \cdot c$ in skip-gram computation with $\left(\sum_{g \in ngram} w_g \cdot c\right)$

fastText: Sub-word Embeddings

Bojanowski et al. (2017)





What if we want embedding representations for whole sentences?

- sentence level (more later)
- Summing?

Sentence Embeddings

Skip-thought vectors (Kiros et al., 2015), similar to skip-gram generalized to a

Is there a way we can compose vectors to make sentence representations?

Will return to this in a few weeks as we move on to syntax and semantics





How to handle different word senses? One vector for bats



- in the sentence, use its internal representations as word vectors
- Context-sensitive word embeddings: depend on rest of the sentence
- Huge improvements across nearly all NLP tasks over GloVe

Preview: Context-dependent Embeddings



ELMo: train a neural language model to predict the next word given previous words

Peters et al. (2018)







Bias in Word Embeddings



- What's wrong with learning a word's "meaning" from its usage? Maybe some words are used in ways we don't want to replicate?
- What data are we learning from?
- What are we going to learn from this data?

What can go wrong with word embeddings?



Answer the following in <=3 sentences each. Consider learning word embeddings from a corpus of news articles.

in the data (e.g., give an example of how it could appear in a news story)

a task for which this biased association might lead to bias in the system?

Now consider learning word embeddings from a corpus of social media data comments (think about reddit + Twitter).

(b) why you think this is present in social media data

Bias Exercise

- 1. Think about a similarity association a model might learn that you believe constitutes **bias.** For this association, list (a) what the word pair is; (b) why you think this is present
- 2. Embeddings are often used at the input layer of a neural network. Can you think of
- 3. Do you think you're likely to see the bad association from above? Why or why not?
- 4. Come up with a new biased similarity association; list (a) what the word pair is;





Compare distance (using cosine similarity) of many occupations to the vectors for he and she

$$\cos(u, v) = \frac{u \cdot v}{\|u\| \|v\|}^{7.10}$$

- These regularities are not restricted to gendered pronouns. receptionist is closer to softball than football
- This work focuses on binary gender stereotypes, but it can be extended

What do we mean by bias?

Extreme	she	occupations
---------	-----	-------------

- 1. homemaker
- 4. librarian
- 7. nanny
- 10. housekeeper
- 5. socialite 8. bookkeeper

2. nurse

- 11. interior designer
- 3. receptionist
- 6. hairdresser
- 9. stylist
- 12. guidance counselor

Extreme *he* occupations

- 1. maestro
- 4. philosopher
 - financier
 - 0. magician
- 2. skipper 5. captain
- 8. warrior
- 11. figher pilot
- 3. protege
- 6. architect
- 9. broadcaster
- 12. boss

Bolukbasi et al. (2016)







What do we mean by bias?

Extreme *she* occupations

- 1. homemaker
- 4. librarian
- 7. nanny
- 10. housekeeper
- 1. maestro
- 4. philosopher
- 7. financier
- 10. magician

- 2. nurse
- 5. socialite
- 8. bookkeeper
- 11. interior designer
- 6. hairdresser

3. receptionist

- 9. stylist
- 12. guidance counselor

Extreme *he* occupations

- 2. skipper
- 5. captain
- 8. warrior
- 11. figher pilot
- 3. protege
- 6. architect
- 9. broadcaster
- 12. boss

Bolukbasi et al. (2016)

Racial Analogies								
$black \rightarrow homeless$	caucasian \rightarrow servicemen							
caucasian \rightarrow hillbilly	asian \rightarrow suburban							
asian \rightarrow laborer	$black \rightarrow landowner$							
Religious Analogies								
$jew \rightarrow greedy$	$muslim \rightarrow powerless$							
$christian \rightarrow familial$	$muslim \rightarrow warzone$							
muslim \rightarrow uneducated	christian \rightarrow intellectually							

Manzini et al. (2019)

Nearest neighbor of (b - a + c)







- Identify gender subspace with gendered words (avg "male" - avg "female" word)
- Project words onto this subspace
- Subtract those projections from the original word









- Not that effective...and the male and female words are still clustered together
- Bias pervades the word embedding space and isn't just a local property of a few words

Gonen and Goldberg (2019)

Hardness of Debiasing





"Toxic degeneration": neural models that generate toxic stuff



[Trump supporters]....|

training data

Toxicity



System trained on a big chunk of the Internet: conditioning on "SJW", "black" gives the system a chance of recalling bad stuff from its

https://toxicdegeneration.allenai.org/





Multilingual Word Embeddings



- Input: a large corpus of text in some language (English)
- Output: embedding for each word
- What if we have multiple corpora of text in different languages?
- Learning embeddings on each language individually: these embeddings aren't expected to have any relation

Recall: Training Embeddings



Multilingual Embeddings

Input: corpora in many languages. Output: embeddings where similar words in different languages have similar embeddings

I have an apple 47 24 18 427

J' ai des oranges 47 24 89 1981

MultiCluster: use bilingual dictionaries to form clusters of words that are translations of one another, replace corpora with cluster IDs, train "monolingual" embeddings over all these corpora

Works okay but not all that well



Ammar et al. (2016)





- What if you already have embeddings in two languages and you just want to align them?
- Given: dictionary of pairs (x_i, z_i) , where x are word embeddings in a source lang (English) and z are word embeddings in a target lang (French)
- Learn a matrix W to minimize the following:



(Looks like a loss function! Can learn with SGD on the pairs)

Aligning existing embeddings

Mikolov et al. (2013)







Aligning existing embeddings



Rotation learns to align these word embedding spaces! Does this cartoon match reality?

Conneau et al. (2017)





Table 2: Accuracy of the word translation methods using the WMT11 datasets. The Edit Distance uses morphological structure of words to find the translation. The Word Co-occurrence technique based on counts uses similarity of contexts in which words appear, which is related to our proposed technique that uses continuous representations of words and a Translation Matrix between two languages.

Translation	Edit Distance		Word Co-occurrence		Translation Matrix		ED + TM		Coverage
	P@1	P@5	P@1	P@5	P@1	P@5	P@1	P@5	
$En \rightarrow Sp$	13%	24%	19%	30%	33%	51%	43%	60%	92.9%
$Sp \rightarrow En$	18%	27%	20%	30%	35%	52%	44%	62%	92.9%
$En \rightarrow Cz$	5%	9%	9%	17%	27%	47%	29%	50%	90.5%
$Cz \rightarrow En$	7%	11%	11%	20%	23%	42%	25%	45%	90.5%

Aligning existing embeddings

Mikolov et al. (2013)

