

Word embeddings (II)

CS 490A, Fall 2020

Applications of Natural Language Processing

https://people.cs.umass.edu/~brenocon/cs490a_f20/

Brendan O'Connor

College of Information and Computer Sciences

University of Massachusetts Amherst

including slides from Eisenstein (2019) and Jurafsky & Martin 3rd Ed.

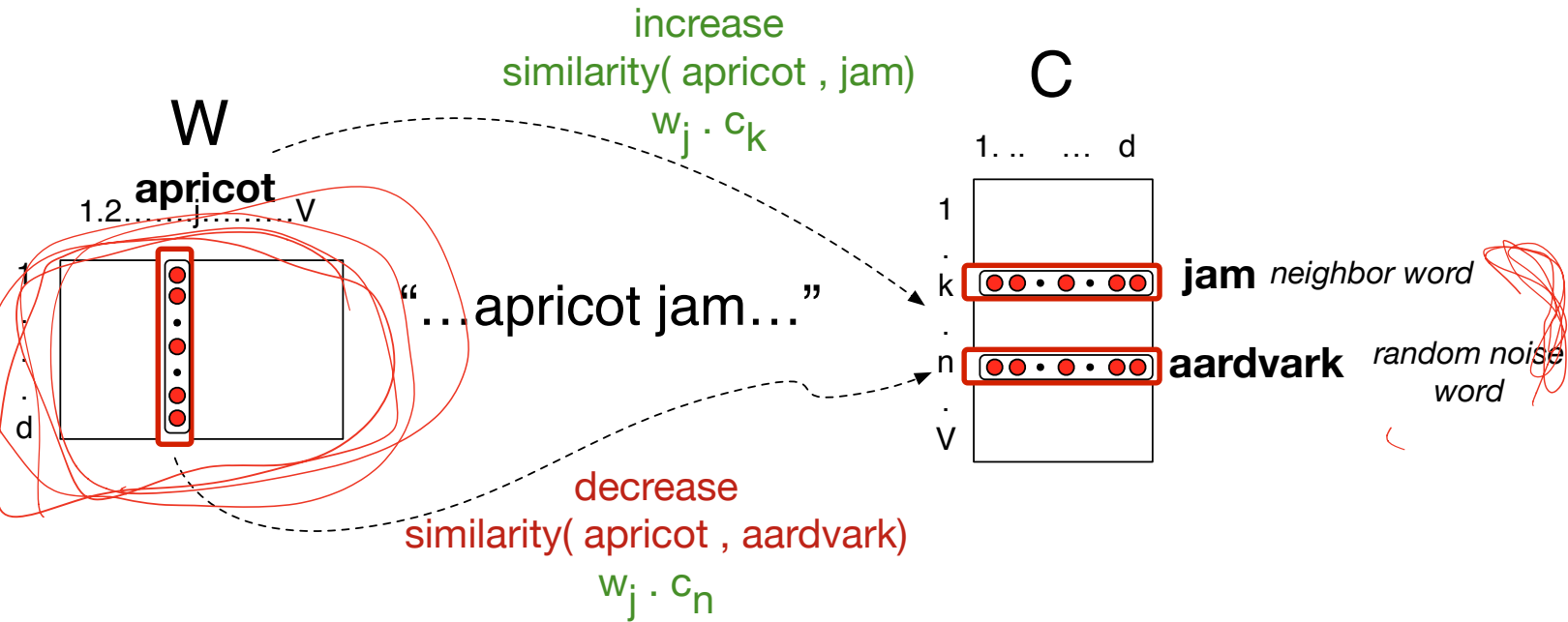
- Back to demo - what's in these embeddings?

Word embeddings

- 1. Input: Large textual corpus
- 2. Language Model: predict words from nearby words
 - GLOVE, SVD: factorize the word-context cooccurrence matrix
 - word2vec: model is viewed as predicting one word's surrounding context words
- 3. Take out the vectors the model was forced to learn; use in downstream applications

- Wikipedia
- News
- Books
- Social Media

in practice, we learn two different sets of embeddings (W for *target* words, C for context words), but throw away C



Defining contexts

Window size C affects the nature of the similarity
something like...

syntax \leftrightarrow basic meaning \leftrightarrow topical meaning

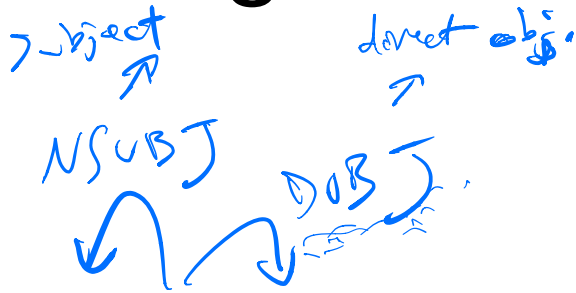
$C = \pm 2$ The nearest words to *Hogwarts*:

- *Sunnydale*
- *Evernight*

$C = \pm 5$ The nearest words to *Hogwarts*:

- *Dumbledore*
- *Malfoy*
- *halfblood*

Defining contexts



The moment one learns English, complications set in (Alfau, 1999)

Brown Clusters

{one}

WORD2VEC, $h = 2$

{moment, one, English, complications}

Structured WORD2VEC, $h = 2$

{(moment, -2), (one, -1), (English, +1), (complications, +2)}

Dependency contexts,

{(one, NSUBJ), (English, DOBJ), (moment, ACL⁻¹)}



Alternate/mis- spellings

- Distributional methods are really good at this
- Brown clusters on Twitter: [http://
www.cs.cmu.edu/~ark/TweetNLP/
cluster_viewer.html](http://www.cs.cmu.edu/~ark/TweetNLP/cluster_viewer.html)

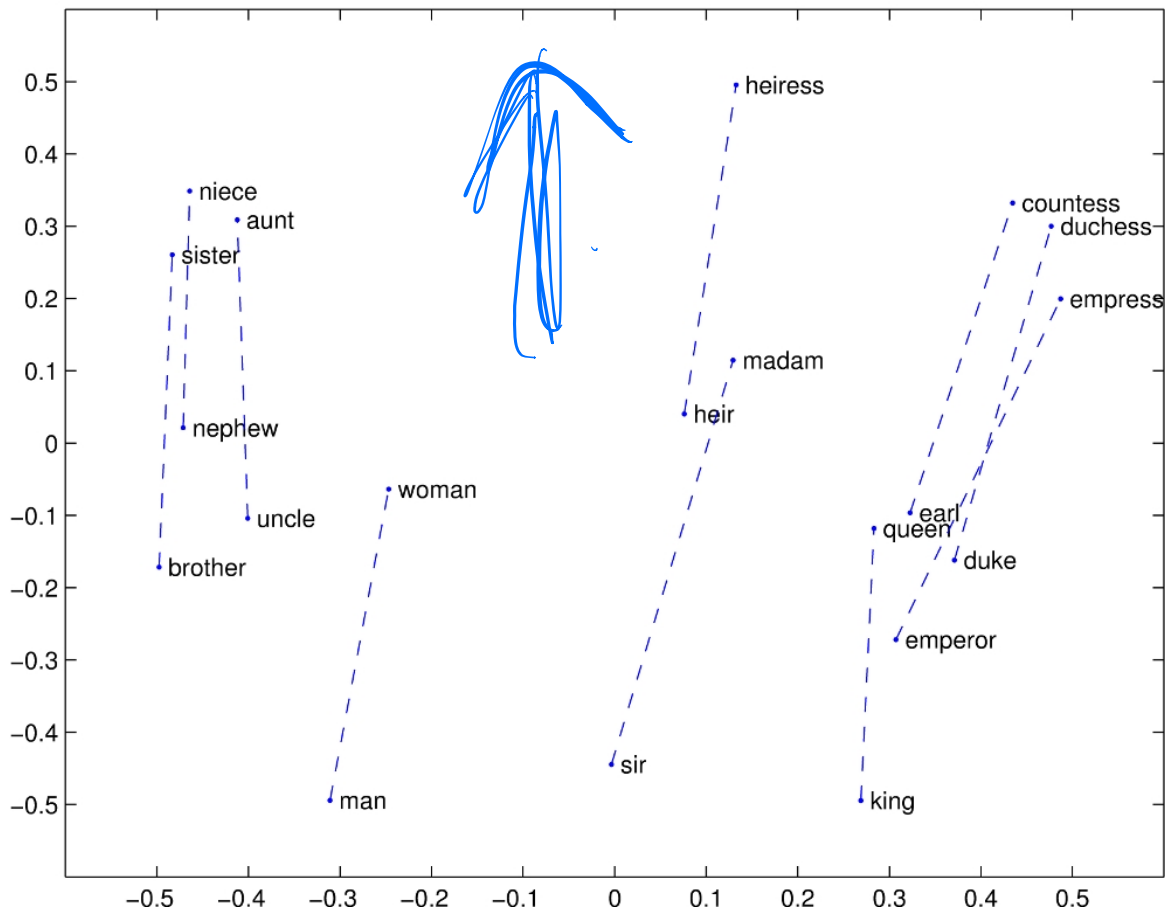
Pre-trained embeddings

- Widely useful. But make sure you know what you're getting!
 - Examples: GLOVE, fasttext, word2vec, etc.
 - Is the corpus similar to what you care about?
 - Should you care about the *data*?

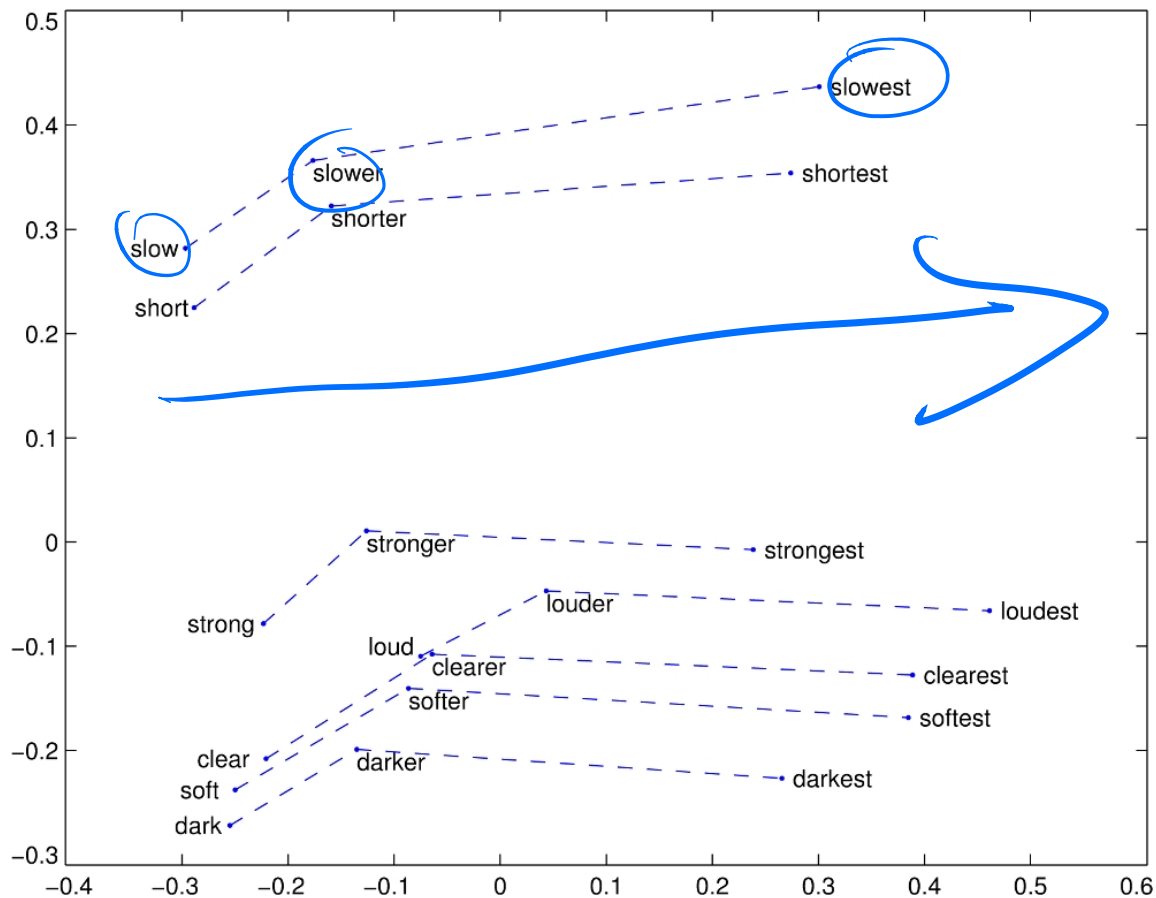
Evaluating embeddings

- Anecdotal inspection (not a real evaluation, but better than nothing; but seen next slides)
- Intrinsic evaluations
 - Compare embeddings' word pair similarities to human judgments
 - TOEFL: "*Levied* is closest to *imposed*, *believed*, *requested*, *correlated*"
 - Numerical similarity judgments (e.g. Wordsim-353)
 - There some other attempts at this (word analogies) but IMO not trustworthy (e.g. Linzen 2016)
 - Extrinsic evaluation: use embeddings in some task

PCA dim. reduction of selected embeddings

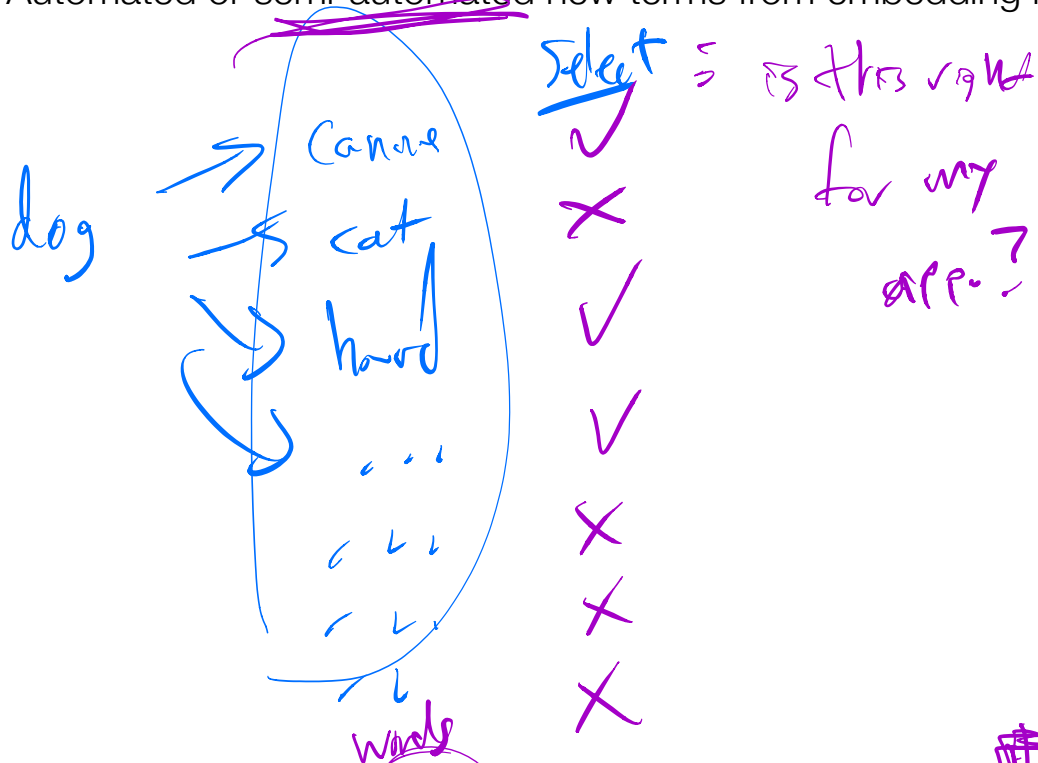


PCA dim. reduction of selected embeddings



Application: keyword expansion

- I have a few keywords for my task. Are there any I missed?
- Automated or semi-automated new terms from embedding neighbors



- Other non-embedding lexical resources can do this too (e.g. WordNet), but word embeddings typically cover a *lot* of diverse vocabulary

why better??

Application: document embedding

- Instead of bag-of-words, can we derive a latent embedding of a document/sentence?

$w_1 = P$

$$\Rightarrow V(w_1) = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

$w_2 = \text{Sam}$

$$\Rightarrow V(w_2) = \begin{bmatrix} 2 & 1 & 1 & 1 & 1 \end{bmatrix}$$

$w_3 = \text{the}$

⋮

$w_4 = \text{more}$

⋮

⋮

⋮

⋮



Avg

$$\text{Doc Embed} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

①

Avg them!
"Bag of Embeddings"

$$\vec{x} = \frac{1}{N_{\text{tok}}} \sum_{i=1}^{N_{\text{tok}}} V(w_i)$$

② use \vec{x} for Logistic

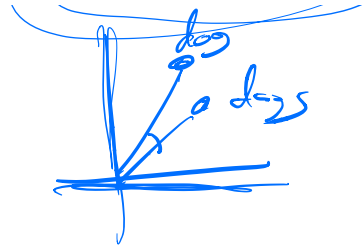
$$p(y=1/x) = \frac{1}{1 + e^{-\beta \vec{x}}}$$

See: Arora et al. 2017

$$\text{dog} = [5, 3, -2]$$

$$\text{dog} = [5.0, 2.0, -1.2]$$

$$B = [2, 2, 2]$$



inv freq weighting

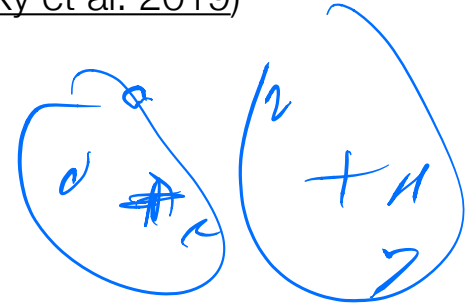
$$\vec{x} = \frac{1}{\sum_i \frac{a}{a+p(w_i)}} V(w_i)$$

$2 =$ \rightarrow $\underbrace{\hspace{10em}}_{\text{freq. in corpus}}$

10^{-3}

Exploratory usage

- Example: tweets about mass shootings (Demszky et al. 2019)
 1. Average word embeddings => tweet embeddings
 2. Cluster tweets (kmeans)
 3. Interpret clusters' words (closest to centroid)



Topic	10 Nearest Stems
news (19%)	break, custodi, #breakingnew, #updat, confirm, fatal, multipl, updat, unconfirm, sever
investigation (9%)	suspect, arrest, alleg, apprehend, custodi, charg, accus, prosecutor, #break, ap
shooter's identity & ideology (11%)	extremist, radic, racist, ideolog, label, rhetor, wing, blm, islamist, christian
victims & location (4%)	bar, thousand, california, calif, among, los, southern, veteran, angel, via
laws & policy (14%)	sensibl, regul, requir, access, abid, #gunreformnow, legisl, argument, allow, #guncontolnow
solidarity (13%)	affect, senseless, ach, heart, heartbroken, sadden, faculti, pray, #prayer, deepest
remembrance (6%)	honor, memori, tuesday, candlelight, flown, vigil, gather, observ, honour, capitol
other (23%)	dude, yeah, eat, huh, gonna, ain, shit, ass, damn, guess

Table 1: Our eight topics (with their average proportions across events) and nearest-neighbor stem embeddings to the cluster centroids. Topic names were manually assigned based on inspecting the tweets.