Distributional Semantics Crash Course

September 11, 2018 CSCI 2952C: Computational Semantics Instructor: Ellie Pavlick HTA: Arun Drelich UTA: Jonathan Chang

Agenda

- Quick Background
- Your Discussion Questions
- Step through of VMS/word2vec
- Announcements

Agenda



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

— J. R. Firth

Firth, Harris







1926: Conditioned Reflexes



Behaviorism

"Behaviorism was developed with the mandate that only observations that satisfied the criteria of the scientific method, namely that they must be repeatable at different times and by independent observers, were to be admissible as evidence. This effectively dismissed introspection, the main technique of psychologists following Wilhelm Wundt's experimental psychology, the dominant paradigm in psychology in the early twentieth century. Thus, behaviorism can be seen as a form of materialism, denying any independent significance to processes of the mind."

Firth (1957)

- Language is a learned behavior, no different than other learned behaviors
- Restricted languages and registers
- Collocations: word types -> meaning
- Colligations: word categories -> syntax

Look-ahead: Frege's Sense and Reference

(for this Thursday)











Discussion! Firth

- different contexts for same word "meaning"
- non-linguistic context, including collocation vs. context, augmented datasets (e.g. tagging)
- emphasis/speech patterns
- language vs. dialect
- slips of the tongue—semantic or prosodic?
- Alice in Wonderland...what else is lost in translation?
- learning "online" without first enumerating all the collocations

Discussion! VMSs

- This paper is from 2010—have there been any fundamental advances since?
- Matrix: multiple levels of context (words, subwords, phrases)? how are patterns chosen? do they make sense out of context? how does context size effect meaning captured? can we model longer phrases and/or morphological roots on the rows? can we put ngrams on the columns?
- Frequencies: how should frequent vs. rare events factor into meaning? should/shouldn't we care more about rare events? what happens with unknown words in the test set?
- Linear Algebraic Assumptions: what to make of the assumptions about vector spaces, e.g. inverses/associativity? is it fair to say that dimensionality reduction -> "higher order features"? why can't we represent arbitrary FOL statements?
- **Applications**: plagiarism detection? text processing (tokenization/normalization)?
- Evaluation/Similarity Metrics: should we model relational similarity directly (pair-pattern) or implicitly, via vector arithmetic? could we reduce attributional similarity to relational similarity/when would this help? do these models only work well on "passive" tasks, or can they work in generation tasks which require knowledge/state?
- **Bias/Ethics**: how do we prevent these models from encoding biases in the data/evaluations? what are the ethical implications e.g. "gaming the system" on resume cites, mining personal information?

Discussion! word2vec

- Matrix: word ordering, size of context
- Frequency: effect of low frequency words, both on rows and columns +
- **Representations**: what differs between parts of speech? what do polysemous words look like? can these capture different senses and more fine-grained "meanings" (e.g. speaker-dependent, context-dependent)? generalizing to new languages?
- Vector Arithmetic: what to make of it? why does France Paris != capitol? can this structure be used to build e.g. ontologies? is the a + b - c order-sensitive, or are they hiding some limitations by focusing on this one type of operation?
- Evaluation/Similarity: can these spaces capture different notions of similarity? why does syntax appear to be easier than semantics? why is it "not surprising" that the NN LM does better than the RNN LM? why is skipgram better than CBOW at semantics? does it have to do with averaging?
- Loss Functions: would more complex loss functions help to learn e.g. transitive verbs? can analogical reasoning relationships be trained directly/incorporated into loss? can multiple loss functions be combined?
- Efficiency: does computational complexity matter that much? is the point moot as machines get faster?

You shall know a word by the company it keeps! Words that occur in similar contexts tend to have similar meanings. If words have similar row vectors in a wordcontext matrix, then they tend to have similar meanings.



markets below levinson olsen remorse schuyler rodents scrambled likely minnesota

Term-Document

Documents as bags of words?





SL

σ

0 0 9 ad

markets below levinson olsen remorse schuyler rodents scrambled likely minnesota administrative berths backup roam operative oalaiologos chrissie supernova andowner bS

Word-Context

Turney and Pantel note that VSMs aren't by limited to text context





peace/region enjoyable/block of/surprise duties/received to/morakot 1942/field returns/goldeh g/overtaken space/second infiltrated/hong \succ the X was X has Y X has X Y has X Y is not X Y or X Y or X Y'S X X and Y

Pair-Pattern

Relationship to Firth's ideas of word classes/ abstraction? Colligation?

	chrissie	supern ova	berths	landow ner	backup	roam	ps	palaiolo gos	operativ e	adminis trative
markets	1000	40	500	700	400	3	80	100	15	6



markets



https://towardsdatascience.com/word2vec-skip-gram-model-part-1-intuition-78614e4d6e0b



https://towardsdatascience.com/word2vec-skip-gram-model-part-1-intuition-78614e4d6e0b

Clarifications/ Procrastinations

- (Neural) Language Modeling:
 - The quick brown fox ____?
- Stochastic gradient descent ("SGD")
- Back-propagation ("Backprop")



SkipGram



Similarity Metrics

- Cosine cares about angle but not length
- Dice/Jaccard for sets/sparse vectors
- Metrics with high vs. low frequency biases What would Firth say?
- Use as features in ML models ("pretraining")

Optimizations/ Approximations

Optimizations/ Approximations

- How much should things like efficiency/scalability matter in a theory of linguistic representation?
- What about computing exactly vs. approximately vs. heuristically? Word embeddings vs. "representation learning"?

Linguistic Preprocessing

- Types vs. tokens
- Tokenization/Phrasal Collocations what should we consider to be the "basic units" of the language?
- Punctuation "okay..." vs. "okay!"
- Normalization "Trump" vs. "trump"
- Stop words "pb and jelly" vs. "pb or jelly"
- Tagging "fish fish fish fish fish"

Mathematical Preprocessing

- Counts: one-hot, frequency, tf-idf/PMI
- Limiting vocab size problems?
- Subsampling in Skipgram: drop words relative to their frequency—what would Firth say about this?
- Dimensionality/sparsity does a "bottle neck" lead to better representations?

Loss Functions

- Softmax: is the predicted distribution (over all words in the vocabulary) the right one?
- Hierarchical Softmax: represent loss function using binary tree, so compute loss for log(V) nodes per word, rather than V words per word.
- NCE/Negative Sampling: can you distinguish the real word from a randomly drawn word (or actually, k randomly drawn words)

If it isn't 11:40 or later, then the fact that I am on this slide means you didn't interrupt enough.

If it is 11:40 or later: well done, team!

Announcements

- Reading for Thursday...there is less of it
- Welcome Jonathan! Office hours TBD (?)
- Arun's office hours: 5pm Wednesday
- My office hours: 5pm this Friday (or some other time?), Monday thereafter 4pm

Assignment 1 is up!

- Quick overview (Arun)
- Due September 25 (in 2 weeks)