Natural Language Processing Lexical semantics

Sofia Serrano sofias6@cs.washington.edu

Credit to Yulia Tsvetkov and Noah Smith for slides

Announcements

- Quiz 3 will be released on Canvas today at 2:20pm
 - Available through Thursday 2:20pm
 - 5 questions, 10 minutes
 - Will cover material from Wednesday, Friday, and Monday (so, language modeling and the first part of lexical semantics)
- Midterm course eval form (online) is out please let us know how we're doing!

Two common solutions for word weighting

tf-idf: tf-idf value for word **t** in document **d**:

$$w_{t,d} = \mathrm{tf}_{t,d} \times \mathrm{idf}_t$$

Words like "the" or "it" have very low idf

PMI: Pointwise mutual information

$$\mathsf{PMI}(w_1, w_2) = log \frac{p(w_1, w_2)}{p(w_1)p(w_2)}$$

See if words like "good" appear more often with "great" than we would expect by chance

TF-IDF

• What to do with words that are evenly distributed across many documents?

$$\mathrm{tf}_{t,d} = \log_{10}(\mathrm{count}(t,d)+1)$$



Words like "the" or "good" have very low idf

$$w_{t,d} = \mathrm{tf}_{t,d} \times \mathrm{idf}_t$$

Positive Pointwise Mutual Information (PPMI)

- In word--context matrix
- Do words w and c co-occur more than if they were independent?

$$PMI(w,c) = \log_2 \frac{P(w,c)}{P(w)P(c)}$$

$$PPMI(w,c) = \max(\log_2 \frac{P(w,c)}{P(w)P(c)}, 0)$$

- PMI is biased toward infrequent events
 - Very rare words have very high PMI values
 - Give rare words slightly higher probabilities α =0.75

$$PPMI_{\alpha}(w,c) = \max(\log_2 \frac{P(w,c)}{P(w)P_{\alpha}(c)}, 0) \qquad \qquad P_{\alpha}(c) = \frac{count(c)^{\alpha}}{\sum_c count(c)^{\alpha}}$$

| # name | formula | referen |
|--------------------------------|---|------------------------------|
| 1. Joint probability | p(xy) | (Giuliano, 19 |
| 2. Conditional probability | p(y x) | (Gregory et al., 199 |
| 3. Reverse cond. probability | p(x y) | (Gregory et al., 199 |
| 4. Pointwise mutual inf. (MI) | $\log \frac{p(xy)}{p(x+)p(xy)}$ | (Church and Hanks, 199 |
| 5. Mutual dependency (MD) | $\log \frac{p(xy)^2}{p(x+)p(+y)}$ | (Thanopoulos et al., 200 |
| 6. Log frequency biased MD | $\log \frac{p(xy)^2}{p(x+)p(xy)} + \log p(xy)$ | (Thanopoulos et al., 200 |
| 7. Normalized expectation | $\frac{2f(xy)}{f(x+)+f(+y)}$ | (Smadja and McKeown, 199 |
| 8. Mutual expectation | $\frac{2f(xy)}{f(x+)+f(xy)} \cdot p(xy)$ | (Dias et al., 200 |
| 9. Salience | $\log \frac{p(xy)^2}{p(x+)p(xy)} \cdot \log f(xy)$ | (Kilgarriff and Tugwell, 200 |
| 10. Pearson's χ^2 test | $\sum_{i,j} \frac{(f_{ij} - \hat{f}_{ij})^2}{f_{ij}}$ | (Manning and Schütze, 199 |
| 11. Fisher's exact test | $\frac{f(x*)!f(\bar{x}*)!f(*y)!f(*\bar{y})!}{N!f(xy)!f(x\bar{y})!f(\bar{x}y)!f(\bar{x}\bar{y})!}$ | (Pedersen, 199 |
| 12. t test | $\frac{f(xy) - \hat{f}(xy)}{\sqrt{f(xy)(1 - (f(xy)/N))}}$ | (Church and Hanks, 199 |
| 13. z score | $\frac{f(xy) - \hat{f}(xy)}{\sqrt{\hat{f}(xy) + \hat{f}(xy)}}$ | (Berry-Rogghe, 192 |
| 14. Poisson significance | $\frac{f(xy)(1-(f(xy)/N))}{f(xy)-f(xy)\log f(xy)+\log f(xy)!}$ | (Quasthoff and Wolff, 200 |
| 15. Log likelihood ratio | $-2\sum_{i,j} f_{ij} \log \frac{f_{ij}}{f_{ij}}$ | (Dunning, 199 |
| 16. Squared log likelihood rat | io $-2\sum_{i,j} \frac{\log r_{ij}^2}{r_{ij}}$ | (Inkpen and Hirst, 200 |
| 17. Russel-Rao | a a+b+c+d | (Russel and Rao, 194 |
| 18. Sokal-Michiner | a+d a+b+c+d | (Sokal and Michener, 193 |
| 19. Rogers-Tanimoto | $\frac{a+d}{a+2b+2c+d}$ | (Rogers and Tanimoto, 190 |
| 20. Hamann | $\frac{(a+d)-(b+c)}{a+b+c+d}$ | (Hamann, 190 |
| 21. Third Sokal-Sneath | b+c a+d | (Sokal and Sneath, 196 |
| 22. Jaccard | a a+b+c | (Jaccard, 19) |
| 23. First Kulczynsky | a b+c | (Kulczynski, 192 |
| 24. Second Sokal-Sneath | a a+2(b+c) | (Sokal and Sneath, 190 |
| 25. Second Kulczynski | $\frac{1}{2}\left(\frac{a}{a+b}+\frac{a}{a+c}\right)$ | (Kulczynski, 192 |
| 26. Fourth Sokal-Sneath | $\frac{1}{4}\left(\frac{a}{a+b} + \frac{a}{a+c} + \frac{d}{d+b} + \frac{d}{d+c}\right)$ | (Kulczynski, 192 |
| 27. Odds ratio | ad bc | (Tan et al., 200 |
| 28. Yulle's ω | Vad-Vbc | (Tan et al., 200 |
| 29. Yulle's Q | ad-bc ad+bc | (Tan et al., 200 |
| 30. Driver-Kroeber | a /(a+b)(a+c) | (Driver and Kroeber, 193 |

| reference | # name |
|-------------|------------------------|
| no, 1964) | 31. Fifth Sokal-Sneath |
| al., 1999) | 32. Pearson |
| al., 1999) | 33. Baroni-Urbani |
| uks, 1990) | 34 Braun-Blanquet |
| al., 2002) | 25 Simpson |
| al., 2002) | 35. Michael |
| wn, 1990) | 36. Michael |
| al., 2000) | 37. Mountford |
| ell, 2001) | 38. Fager |
| tze, 1999) | 39. Unigram subtuples |
| en, 1996) | 40. U cost |
| uks, 1990) | 41. S cost |
| ;he, 1973) | 42. R cost |
| olff, 2002) | 43. T combined cost |
| ng, 1993) | 44. Phi |
| rst, 2002) | 45. Kappa |
| lao, 1940) | 46. J measure |
| ner, 1958) | |
| oto, 1960) | 47. Gini index |
| nn, 1961) | |
| ath, 1963) | |
| rd, 1912) | |
| ski, 1927) | 48. Confidence |
| ath, 1963) | 49. Laplace |
| ski, 1927) | 50. Conviction |
| ski, 1927) | 51. Piatersky-Shapiro |
| al., 2002) | 52. Certainity factor |
| al., 2002) | 53. Added value (AV) |
| al., 2002) | 54 Collective strength |
| ber, 1932) | 55 Klosgen |
| | b. Riosgen |

| formula | reference |
|--|--------------------|
| $\frac{ad}{\sqrt{(a+b)(a+c)(d+b)(d+c)}}$ (Sokal a | and Sneath, 1963) |
| $\frac{ad-bc}{\sqrt{(a+b)(a+c)(d+b)(d+c)}}$ | (Pearson,1950) |
| <u>a+vad</u> (Baroni-Urbani | and Buser, 1976) |
| a+b+c+vad (Brau | n-Blanquet 1932) |
| a (Druce | (Simpson 1943) |
| $\frac{\min\{a+b,a+c\}}{4(ad-bc)}$ | (Michael 1020) |
| $(a+d)^2 + (b+c)^2$ | (Michael, 1920) |
| 2bc+ab+ac (Kaufman and | Kousseeuw, 1990) |
| $\frac{a}{\sqrt{(a+b)(a+c)}} - \frac{1}{2} \max(b,c)$ (Kaufman and) | Rousseeuw, 1990) |
| $\log \frac{ad}{bc} - 3.29\sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$ (Blaheta and | nd Johnson, 2001) |
| $log(1 + \frac{min(b,c)+a}{max(b,c)+a})$ | (Tulloss, 1997) |
| $\log(1 + \frac{\min(b,c)}{a+1})^{-\frac{1}{2}}$ | (Tulloss, 1997) |
| $\log(1 + \frac{a}{a+b}) \cdot \log(1 + \frac{a}{a+b})$ | (Tulloss, 1997) |
| $\sqrt{U \times S \times R}$ | (Tulloss, 1997) |
| p(xy)-p(x+)p(+y) | (Tan et al., 2002) |
| $\sqrt{p(x*)p(*y)(1-p(x*))(1-p(*y))}$ $p(xy)+p(\tilde{x}\tilde{y})-p(x*)p(xy)-p(\tilde{x}*)p(*\tilde{y})$ | (, |
| $\frac{1 - p(x *)p(*y) - p(\hat{x} *)p(*\hat{y})}{p(x)}$ | (Ian et al., 2002) |
| $\max[p(xy)\log\frac{p(y x)}{p(*y)} + p(x\bar{y})\log\frac{p(y x)}{p(*\bar{y})},$ | (Tan et al., 2002) |
| $p(xy)\log\frac{p(x y)}{p(x+)} + p(\bar{x}y)\log\frac{p(x y)}{p(\bar{x}+)}]$ | |
| $\max[p(x*)(p(y x)^2 + p(\bar{y} x)^2) - p(*y)^2]$ | (Tan et al., 2002) |
| $+p(\bar{x*})(p(y \bar{x})^2 + p(\bar{y} \bar{x})^2) - p(*\bar{y})^2,$ | |
| $p(*y)(p(x y)^2 + p(\bar{x} y)^2) - p(x*)^2$ | |
| $+p(*\bar{y})(p(x \bar{y})^2+p(\bar{x} \bar{y})^2)-p(\bar{x}*)^2]$ | |
| $\max[p(y x), p(x y)]$ | (Tan et al., 2002) |
| $\max[\frac{Np(xy)+1}{Np(xy)+2}, \frac{Np(xy)+1}{Np(xy)+2}]$ | (Tan et al., 2002) |
| $\max[\frac{p(x*)p(*y)}{p(xy)}, \frac{p(x*)p(*y)}{p(xy)}]$ | (Tan et al., 2002) |
| p(xy) - p(x*)p(xy) | (Tan et al., 2002) |
| $\max[\frac{p(y x)-p(xy)}{1-p(xy)}, \frac{p(x y)-p(x+)}{1-p(x+)}]$ | (Tan et al., 2002) |
| $\max[p(y x) - p(xy), p(x y) - p(x*)]$ | (Tan et al., 2002) |
| $\frac{p(xy) + p(xy)}{1 - p(x*)p(*y) - p(x*)p(*y)}$ | (Tan et al., 2002) |
| $p(x+)p(y)+p(x+)p(*y) = 1-p(xy)-p(xg)$ $\sqrt{p(xy)} = AV$ | (Tan et al. 2002) |
| V P(Ag) . AV | (lan et al., 2002) |



Dense vectors (part 1)

Term-document matrix from Monday

| | As You Like It | Twelfth Night | Julius Caesar | Henry V |
|---------|-------------------|------------------|------------------|---------|
| battle | 1 | 0 | 7 | 17 |
| soldier | 2 | 80 | 62 | 89 |
| fool | 36 | 58 | 1 | 4 |
| clown | 20 | 15 | 2 | 3 |

These word vectors are still the length of our number of documents! Hmmm...

Dimensionality Reduction

- Wikipedia: ~29 million English documents. Vocab: ~1M words.
 - High dimensionality of word--document matrix
 - Sparsity
 - The order of rows and columns doesn't matter
- Goal:
 - good similarity measure for words or documents
 - dense representation
- Sparse vs Dense vectors
 - Short vectors may be easier to use as features in machine learning (less weights to tune)
 - Dense vectors may generalize better than storing explicit counts
 - They may do better at capturing synonymy
 - In practice, they work better





Solution idea

- Find a projection into a low-dimensional space (~300 dim)...
- ... that, up to a certain vector-length budget, preserves the most important information

We turn to Singular Value Decomposition (SVD)

Any matrix can be decomposed into



Orthonormal, unitary (Rectangular) diagonal

Orthonormal, unitary

Any matrix can be decomposed into



Let's trim away the zero scaling factors





Truncated SVD

We can approximate the full matrix by only considering the leftmost k terms in the diagonal matrix (the k largest singular values) dense document vectors



 $A_{m \times n} \approx U_{m \times k} \Sigma_{k \times k} V_{k \times n}^{\top}$ $k \ll m, n$

Latent Semantic Analysis

| #0 | #1 | #2 | #3 | #4 | #5 |
|-------|------------|-----------|----------|-------------|------|
| we | music | company | how | program | 10 |
| said | film | mr | what | project | 30 |
| have | theater | its | about | russian | 11 |
| they | mr | inc | their | space | 12 |
| not | this | stock | or | russia | 15 |
| but | who | companies | this | center | 13 |
| be | movie | sales | are | programs | 14 |
| do | which | shares | history | clark | 20 |
| he | show | said | be | aircraft | sept |
| this | about | business | social | ballet | 16 |
| there | dance | share | these | its | 25 |
| you | its | chief | other | projects | 17 |
| are | disney | executive | research | orchestra | 18 |
| what | play | president | writes | development | 19 |
| if | production | group | language | work | 21 |

How do we tell whether a set of word embeddings is any good?

Evaluation

- Intrinsic
- Extrinsic
- Qualitative

_

| WORD | d1 | d2 | d3 | d4 | d5 | | d50 |
|----------|-----------|------|------|-----------|------|------|------|
| summer | 0.12 | 0.21 | 0.07 | 0.25 | 0.33 | | 0.51 |
| spring | 0.19 | 0.57 | 0.99 | 0.30 | 0.02 | | 0.73 |
| fall | 0.53 | 0.77 | 0.43 | 0.20 | 0.29 | •••• | 0.85 |
| light | 0.00 | 0.68 | 0.84 | 0.45 | 0.11 | | 0.03 |
| clear | 0.27 | 0.50 | 0.21 | 0.56 | 0.25 | | 0.32 |
| blizzard | 0.15 | 0.05 | 0.64 | 0.17 | 0.99 | | 0.23 |

Extrinsic Evaluation

- Chunking
- POS tagging
- Parsing
- MT
- SRL
- Topic categorization
- Sentiment analysis
- Metaphor detection

19

• etc.

Intrinsic Evaluation



20

• MEN-3k (<u>Bruni et al. '12</u>)

SimLex-999 dataset (<u>Hill et al., 2015</u>)

Computing word similarity

The dot product between two vectors is a scalar:

dot product(
$$\mathbf{v}, \mathbf{w}$$
) = $\mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^{N} v_i w_i = v_1 w_1 + v_2 w_2 + \dots + v_N w_N$

- The dot product tends to be high when the two vectors have large values in the same dimensions
- Dot product can thus be a useful similarity metric between vectors

Problem with raw dot-product

- Dot product favors long vectors
 - Dot product is higher if a vector is longer (has higher values in many dimension) Vector length:

$$|\mathbf{v}| = \sqrt{\sum_{i=1}^{N} v_i^2}$$

- Frequent words (of, the, you) have long vectors (since they occur many times with other words).
 - So dot product overly favors frequent words

Alternative: cosine for computing word similarity

$$\operatorname{cosine}(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}| |\vec{w}|} = \frac{\sum_{i=1}^{N} v_i w_i}{\sqrt{\sum_{i=1}^{N} v_i^2} \sqrt{\sum_{i=1}^{N} w_i^2}}$$

Based on the definition of the dot product between two vectors a and b

$$\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}| |\mathbf{b}| \cos \theta$$
$$\frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{a}| |\mathbf{b}|} = \cos \theta$$

Cosine as a similarity metric

- -1: vectors point in opposite directions
- +1: vectors point in same directions
- **0**: vectors are orthogonal



• But since raw frequency values are non-negative, the cosine for term-term matrix vectors ranges from 0–1

Cosine examples

$$\cos(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}| |\vec{w}|} = \frac{\vec{v}}{|\vec{v}|} \cdot \frac{\vec{w}}{|\vec{w}|} = \frac{\sum_{i=1}^{N} v_i w_i}{\sqrt{\sum_{i=1}^{N} v_i^2} \sqrt{\sum_{i=1}^{N} w_i^2}}$$

| | pie | data | computer |
|-------------|-----|------|----------|
| cherry | 442 | 8 | 2 |
| digital | 114 | 80 | 62 |
| information | 36 | 58 | 1 |

 $\cos(\text{cherry}, \text{information}) =$

$$\frac{442*5+8*3982+2*3325}{\sqrt{442^2+8^2+2^2}\sqrt{5^2+3982^2+3325^2}} = .017$$

 $\cos(\text{digital}, \text{information}) =$

$$\frac{5*5+1683*3982+1670*3325}{\sqrt{5^2+1683^2+1670^2}\sqrt{5^2+3982^2+3325^2}} = .996$$

Visualizing angles



Visualisation



Figure 6.5: Monolingual (top) and multilingual (bottom; marked with apostrophe) word projections of the antonyms (shown in red) and synonyms of "beautiful".

Visualizing Data using t-SNE (van der Maaten & Hinton '08)

Dense vectors (part 2)

Distributed representations

Word Vectors

| WORD | d1 | d2 | d3 | d4 | d5 | | d50 |
|----------|------|------|------|------|------|-------|------|
| summer | 0.12 | 0.21 | 0.07 | 0.25 | 0.33 | | 0.51 |
| spring | 0.19 | 0.57 | 0.99 | 0.30 | 0.02 | ••• | 0.73 |
| fall | 0.53 | 0.77 | 0.43 | 0.20 | 0.29 | • • • | 0.85 |
| light | 0.00 | 0.68 | 0.84 | 0.45 | 0.11 | | 0.03 |
| clear | 0.27 | 0.50 | 0.21 | 0.56 | 0.25 | | 0.32 |
| blizzard | 0.15 | 0.05 | 0.64 | 0.17 | 0.99 | | 0.23 |

"One hot" vectors and dense word vectors (embeddings)



Low-dimensional word representations

- Learning representations by back-propagating errors
 - Rumelhart, Hinton & Williams, 1986
- A neural probabilistic language model
 - Bengio et al., 2003
- Natural Language Processing (almost) from scratch
 - Collobert & Weston, 2008
- Word representations: A simple and general method for semi-supervised learning
 - Turian et al., 2010
- Distributed Representations of Words and Phrases and their Compositionality
 - Word2Vec; Mikolov et al., 2013

Word2Vec

- Popular embedding method
- Very fast to train
- Code available on the web
- Idea: predict rather than count



Word2Vec

PROJECTION INPUT OUTPUT INPUT PROJECTION OUTPUT w(t-2) w(t-2) w(t-1) w(t-1) SUM w(t) w(t) w(t+1) w(t+1) w(t+2) w(t+2) Skip-gram **CBOW**

• Predict vs Count





the cat sat on the mat

• Predict vs Count



• Predict vs Count



Predict vs Count



• Predict vs Count





• Predict vs Count



• Predict vs Count





• Predict vs Count



Skip-gram



How to compute p(+|t,c)?



FastText



Typical traits of these embeddings

Automatically learn some analogies pretty well



Figure from Sutor et al. MIPR 2019



What we've learned

- The contexts in which a word typically appears (i.e., the tokens that typically appear around it) tell us a lot about that word
- We can use those contexts to automatically learn more powerful representations of words than just a one-hot encoding
- These "word embeddings" can plug in as parameters in models of your choice

Next class

Neural Networks I (Feedforward networks and LSTMs)