

# Computational Social Science

## Word embeddings I

Dr. Thomas Davidson

Rutgers University

March 7, 2022

# Plan

1. Course updates
2. The vector-space model review
3. Latent semantic analysis
4. Language models

# Course updates

- ▶ Project proposal assignment on Canvas, due Friday at 5pm
  - ▶ A short description of the planned project
  - ▶ Data collection
  - ▶ Data cleaning
  - ▶ Data analysis
  - ▶ Data visualization
  - ▶ Team

# The vector-space model review

## Vector representations

- ▶ Last week we looked at how we can represent texts as numeric vectors
  - ▶ Documents as vectors of words
  - ▶ Words as vectors of documents
- ▶ A document-term matrix ( $DTM$ ) is a matrix where documents are represented as rows and tokens as columns

# The vector-space model review

## Weighting schemes

- ▶ We can use different schemes to weight these vectors
  - ▶ Binary (Does word  $w_i$  occur in document  $d_j$ ?)
  - ▶ Counts (How many times does word  $w_i$  occur in document  $d_j$ ?)
  - ▶ TF-IDF (How many times does word  $w_i$  occur in document  $d_j$ , accounting for how often  $w_i$  occurs across all documents  $d \in D$ ?)
    - ▶ Recall *Zipf's Law*: a handful of words account for most words used; such words do little to help us to distinguish between documents

# The vector-space model review

## Cosine similarity

$$\cos(\theta) = \frac{\vec{u} \cdot \vec{v}}{\|\vec{u}\| \|\vec{v}\|} = \frac{\sum_i \vec{u}_i \vec{v}_i}{\sqrt{\sum_i \vec{u}_i^2} \sqrt{\sum_i \vec{v}_i^2}}$$

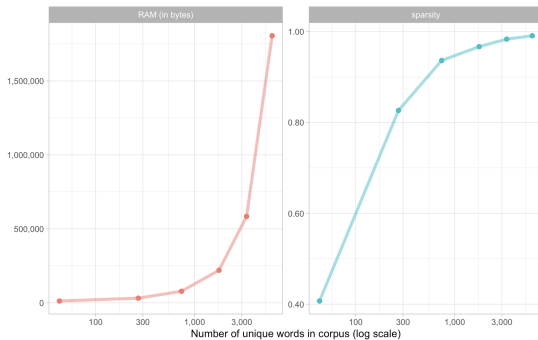
# The vector-space model review

## Limitations

- ▶ These methods produce *high-dimensional, sparse* vector representations
  - ▶ Given a vocabulary of unique tokens  $N$  the length of each vector  $|V| = N$ .
  - ▶ Many values will be zero since most documents only contain a small subset of the vocabulary.

# The vector-space model review

## Limitations



Source: <https://smiltar.com/embeddings.html>



# Latent semantic analysis

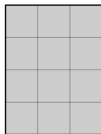
## Latent Semantic Analysis

- ▶ One approach to reduce dimensionality and better capture semantics is called **Latent Semantic Analysis (LSA)**
  - ▶ We can use a process called *singular value decomposition* to find a *low-rank approximation* of a DTM.
- ▶ This provides *low-dimensional, dense* vector representations
  - ▶ Low-dimensional, since  $|V| \ll N$
  - ▶ Dense, since vectors contain real values, with few zeros
- ▶ In short, we can “squash” a big matrix into a much smaller matrix while retaining important information.

$$DTM = U\Sigma V^T$$

# Latent semantic analysis

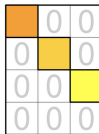
## Singular Value Decomposition



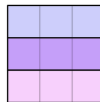
$$\mathbf{M}_{m \times n}$$



$$\mathbf{U}_{m \times m}$$



$$\mathbf{\Sigma}_{m \times n}$$



$$\mathbf{V}^*_{n \times n}$$

$$\mathbf{M} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*$$

See the Wikipedia page for video of the latent dimensions in a sparse TDM.

# Latent semantic analysis

## Example: Shakespeare's writings

X is a TF-IDF weighted Document-Term Matrix of Shakespeare's writings from Project Gutenberg. There are 11,666 unique tokens (each of which occurs 10 or more times in the corpus) and 66 documents.

```
library(tidyverse)
library(ggplot2)
X <- as.matrix(read.table("data/shakespeare.txt"))
X <- X[, which(colSums(X) != 0)] # Drop zero columns
X <- X[, -which(colnames(X) %in% c("footnote", "sidenote"))]

dim(X)

## [1]    66 11664
```

# Latent semantic analysis

## Creating a lookup dictionary

We can construct a list to allow us to easily find the index of a particular token.

```
lookup.index.from.token <- list()

for (i in 1:length(colnames(X))) {
  lookup.index.from.token[colnames(X)[i]] <- i
}
```

# Latent semantic analysis

## Using the lookup dictionary

This easily allows us to find the vector representation of a particular word. Note how most values are zero since the character Hamlet is only mentioned in a handful of documents.

```
lookup.index.from.token["hamlet"]
```

```
## $hamlet
```

```
## [1] 8231
```

```
round(as.numeric(X[,unlist(lookup.index.from.token["hamlet"])]),3)
```

```
## [1] 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.0
```

```
## [13] 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.0
```

```
## [25] 0.046 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.0
```

```
## [37] 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.0
```

```
## [49] 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.010 0.000 0.000 0.0
```

```
## [61] 0.014 0.000 0.000 0.002 0.000 0.000
```

# Latent semantic analysis

## Calculating similarities

The following code normalizes each *column* (rather than row normalization seen last lecture) and constructs a word-word cosine-similarity matrix.

```
normalize <- function(X) {  
  for (i in 1:dim(X)[2]) {  
    X[,i] <- (X[,i]/sqrt(sum(X[,i]^2)))  
  }  
  return(X)  
}  
  
X.n <- normalize(X)  
  
sims <- t(X.n) %*% X.n  
dim(sims)  
## [1] 11664 11664
```

# Latent semantic analysis

## Most similar function

For a given token, this function allows us to find the  $n$  most similar tokens in the similarity matrix, where  $n$  defaults to 10.

```
get.top.n <- function(token, sims, n=10) {  
  top <- sort(sims[unlist(lookup.index.from.token[token]),],  
              decreasing=T)[1:n]  
  return(top)  
}
```

# Latent semantic analysis

## Finding similar words

```
get.top.n("summer",sims)
```

```
##      summer      glass      shade      minutes      winter      ages      poet
## 1.0000000 0.9449493 0.9447479 0.9424049 0.9365960 0.9217019 0.919979
##      cures      lays
## 0.9108892 0.9105933
```

```
get.top.n("fight", sims)
```

```
##      fight      retreat      sword      alarum      field      marching      strengt
## 1.0000000 0.8477011 0.8346863 0.8130756 0.8023591 0.8010200 0.783610
##      blood      soldiers
## 0.7585464 0.7579380
```

```
get.top.n("romeo", sims)
```

```
##      romeo      mercutio      tybalt      tybalts      capulet      benvolio      montague
## 1.0000000 0.9999330 0.9999318 0.9998235 0.9995140 0.9992019 0.998216
##      sampson      juliet
## 0.9923716 0.9906526
```



# Latent semantic analysis

## Singular value decomposition

The `svd` function allows us to decompose the DTM. We can then easily reconstruct it using the formula shown above.

```
# Computing the singular value decomposition
```

```
lsa <- svd(X)
```

```
# We can easily recover the original matrix from this representation
```

```
X.2 <- lsa$u %*% diag(lsa$d) %*% t(lsa$v) #  $X = U \Sigma V^T$ 
```

```
# Verifying that values are the same, example of first column
```

```
sum(round(X-X.2,5))
```

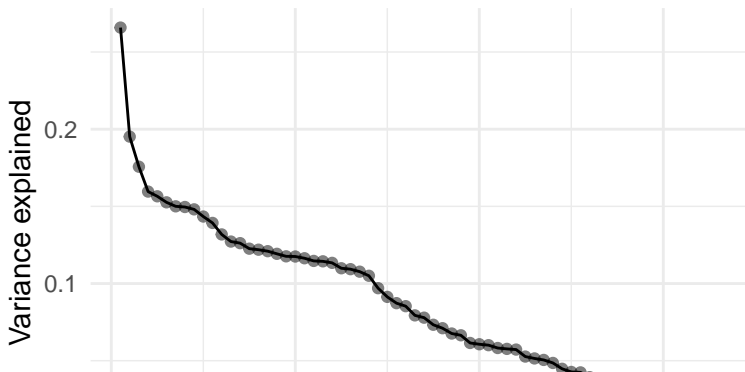
```
## [1] 0
```

# Latent semantic analysis

## Singular value decomposition

This plot shows the magnitude of the singular values (the diagonal entries of  $\Sigma$ ). The magnitude of the singular value corresponds to the amount of variance explained in the original matrix.

Variance explained by singular values



# Latent semantic analysis

## Truncated singular value decomposition

In the example above retained the original matrix dimensions. The point of latent semantic analysis is to compute a *truncated* SVD such that we have a new matrix in a sub-space of  $X$ . In this case we only want to retain the first  $k$  dimensions of the matrix.

```
k <- 20 # Dimensions in truncated matrix

# We can take the SVD of X but only retain the first k singular values
lsa.2 <- svd(X, nu=k, nv=k)

# In this case we reconstruct X just using the first k singular values
X.trunc <- lsa.2$u %*% diag(lsa.2$d[1:k]) %*% t(lsa.2$v)

# But the values will be slightly different since it is an approximation
# Some information is lost due to the compression
sum(round(X-X.trunc,2))

## [1] 7.75
```

# Latent semantic analysis

## Recalculating similarities using the LSA matrix

```
words.lsa <- t(lsa.2$v)
colnames(words.lsa) <- colnames(X)

round(as.numeric(words.lsa[,unlist(lookup.index.from.token["hamlet"])]))
## [1] 0.00 0.01 0.00 -0.03 0.11 0.03 -0.03 0.01 -0.04 0.01 -0.
## [13] -0.01 0.00 0.00 -0.04 0.00 0.00 -0.01 -0.02
```

# Latent semantic analysis

## Recalculating similarities using the LSA matrix

```
words.lsa.n <- normalize(words.lsa)
sims.lsa <- t(words.lsa.n) %*% words.lsa.n
```

# Latent semantic analysis

## Comparing similarities

```
get.top.n("summer",sims)
```

```
##      summer      glass      shade  minutes    winter      ages      poet
## 1.0000000 0.9449493 0.9447479 0.9424049 0.9365960 0.9217019 0.919979
##      cures      lays
## 0.9108892 0.9105933
```

```
get.top.n("summer",sims.lsa)
```

```
##      summer      dance      eyes    deface    breath    flesh    flam
## 1.0000000 0.9159871 0.9061238 0.8969820 0.8934762 0.8930085 0.891249
##      half      birds
## 0.8886182 0.8863785
```

```
bind_cols(names(get.top.n("summer",sims)), names(get.top.n("summer",sims.lsa)))
```

```
## # A tibble: 10 x 2
```

```
##   ...1    ...2
```

```
##   <chr>  <chr>
```

```
## 1 summer summer
```

```
## 2 glass  dance
```

# Latent semantic analysis

## Comparing similarities

```
get.bottom.n <- function(token, sims, n=10) {  
  bottom <- sort(sims[unlist(lookup.index.from.token[token]),],  
                 decreasing=F)[1:n]  
  return(bottom)  
}
```

```
get.bottom.n("summer", sims)
```

```
##   abness   aegeon syracuse eglamour   silvia   silvius elsinore ophel  
##       0         0           0         0         0         0         0  
##   osric reynaldo  
##       0         0
```

# Latent semantic analysis

## Comparing similarities

```
get.top.n("fight",sims)
```

```
##      fight  retreat      sword      alarum      field  marching  strengt
## 1.0000000 0.8477011 0.8346863 0.8130756 0.8023591 0.8010200 0.783610
##      blood  soldiers
## 0.7585464 0.7579380
```

```
get.top.n("fight",sims.lsa)
```

```
##      fight      sword      blood      seat  children      chance      thron
## 1.0000000 0.9375718 0.9265156 0.9222373 0.9192967 0.9182052 0.912325
##      field      whos
## 0.9015076 0.8990767
```

```
bind_cols(names(get.top.n("fight",sims)), names(get.top.n("fight",sims.
```

```
## # A tibble: 10 x 2
```

```
##   ...1   ...2
```

```
##   <chr>  <chr>
```

```
## 1 fight  fight
```

```
## 2 retreat retreat
```



# Latent semantic analysis

## Comparing similarities

```
get.top.n("romeo", sims)
```

```
##      romeo  mercutio   tybalt   tybalts   capulet  benvolio montague  
## 1.0000000 0.9999330 0.9999318 0.9998235 0.9995140 0.9992019 0.998216  
## sampson   juliet  
## 0.9923716 0.9906526
```

```
get.top.n("romeo", sims.lsa)
```

```
##      romeo montagues   tybalts   tybalt  mercutio mercutios   capule  
## 1.0000000 0.9999960 0.9999946 0.9999940 0.9999925 0.9999911 0.999988  
## capulets   romeos  
## 0.9999838 0.9999668
```

# Latent semantic analysis

## Comparing similarities

```
get.top.n("hamlet", sims)
```

```
##      hamlet  horatio marcellus  ophelia  polonius  barnardo  laerte  
## 1.0000000  0.9829677  0.9824643  0.9607921  0.9600178  0.9591448  0.958729  
## voltemand  lucianus  
## 0.9465517  0.9308989
```

```
get.top.n("hamlet", sims.lsa)
```

```
##      hamlet fortinbras  laertes  hamlets  horatio  ophelia  
## 1.0000000  0.9997027  0.9992680  0.9992437  0.9988308  0.9987964  0  
## gertrude  ophelias  danish  
## 0.9986344  0.9985602  0.9985404
```

# Latent semantic analysis

## Exercise

Re-run the code above with a different value of  $k$  (try lower or higher). Compare some terms in the original similarity matrix and the new matrix. How does changing  $k$  affect the results?

```
get.top.n("", sims)
```

```
## [1] NA NA NA NA NA NA NA NA NA NA
```

```
get.top.n("", sims.lsa)
```

```
## [1] NA NA NA NA NA NA NA NA NA NA
```

# Latent semantic analysis

## Inspecting the latent dimensions

We can analyze the meaning of the latent dimensions by looking at the terms with the highest weights in each row. In this case I use the raw LSA matrix without normalizing it. What do you notice about the dimensions?

```
for (i in 1:dim(words.lsa)[1]) {  
  top.words <- sort(words.lsa[i,], decreasing=T)[1:5]  
  print(paste(c("Dimension: ",i), collapse=" "))  
  print(top.words)  
}
```

```
## [1] "Dimension: 1"
```

```
##          amy          bened          bero          botes          cas  
## -1.183045e-06 -1.183045e-06 -1.183045e-06 -1.183045e-06 -1.183045e-06
```

```
## [1] "Dimension: 2"
```

```
##      iago  othello  cassio  desdemona  emilia  
## 0.6565950 0.4270426 0.4018033 0.3006999 0.1657854
```

```
## [1] "Dimension: 3"
```

```
## benedick  leonato  beatrice  pedro  claudio
```

# Latent semantic analysis

## Limitations of Latent Semantic Analysis

- ▶ Bag-of-words assumptions and document-level word associations
  - ▶ We still treat words as belonging to documents and lack finer context about their relationships
    - ▶ Although we could theoretically treat smaller units like sentences as documents
- ▶ Matrix computations become intractable with large corpora
- ▶ A neat linear algebra trick, but no underlying *language model*

# Language models

## Intuition

- ▶ A language model is a probabilistic model of language use
- ▶ Given some string of tokens, what is the most likely token?
  - ▶ Examples
    - ▶ Auto-complete
    - ▶ Google search completion

# Language models

## Bigram models

- ▶  $P(w_i|w_{i-1})$  = What is the probability of word  $w_i$  given the last word,  $w_{i-1}$ ?
  - ▶  $P(\textit{Jersey}|\textit{New})$
  - ▶  $P(\textit{Brunswick}|\textit{New})$
  - ▶  $P(\textit{York}|\textit{New})$
  - ▶  $P(\textit{Sociology}|\textit{New})$

# Language models

## Bigram models

- ▶ We use a corpus of text to calculate these probabilities from word co-occurrence.
  - ▶  $P(\text{Jersey}|\text{New}) = \frac{C(\text{New Jersey})}{C(\text{New})}$ , e.g. proportion of times “New” is followed by “Jersey”, where  $C()$  is the count operator.
- ▶ More frequently occurring pairs will have a higher probability.
  - ▶ We might expect that  $P(\text{York}|\text{New}) \approx P(\text{Jersey}|\text{New}) > P(\text{Brunswick}|\text{New}) \gg P(\text{Sociology}|\text{New})$



# Language models

## Incorporating more information

- ▶ We can also model the probability of a word, given a sequence of words
- ▶  $P(x|S)$  = What is the probability of some word  $x$  given a partial sentence  $S$ ?
- ▶  $A = P(\text{Jersey} | \text{Rutgers University is in New})$
- ▶  $B = P(\text{Brunswick} | \text{Rutgers University is in New})$
- ▶  $C = P(\text{York} | \text{Rutgers University is in New})$
- ▶ In this case we have more information, so “York” is less likely to be the next word. Hence,
  - ▶  $A \approx B > C$ .

# Language models

## Estimation

We can compute the probability of an entire sequence of words by using considering the *joint conditional probabilities* of each pair of words in the sequence. For a sequence of  $n$  words, we want to know the joint probability of  $P(w_1, w_2, w_3, \dots, w_n)$ . We can simplify this using the chain rule of probability:

$$\begin{aligned} P(w_{1:n}) &= P(w_1)P(w_2|w_1)P(w_3|w_{1:2})\dots P(w_n|w_{1:n-1}) \\ &= \prod_{k=1}^n P(w_k|w_{1:k-1}) \end{aligned}$$

# Language models

## Estimation

The bigram model simplifies this by assuming it is a first-order Markov process, such that the probability  $w_k$  only depends on the previous word,  $w_{k-1}$ .

$$P(w_{1:n}) \approx \prod_{k=1}^n P(w_k | w_{k-1})$$

These probabilities can be estimated by using Maximum Likelihood Estimation on a corpus.

See <https://web.stanford.edu/~jurafsky/slp3/3.pdf> for an excellent review of language models

# Language models

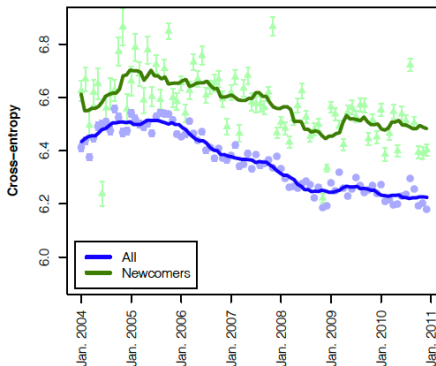
## Empirical applications

- ▶ Danescu-Niculescu-Mizil et al. 2013 construct a bigram language model for each month on *BeerAdvocate* and *RateBeer* to capture the language of the community
  - ▶ For any given comment or user, they can then use a measure called *cross-entropy* to calculate how “surprising” the text is, given the assumptions about the language model
- ▶ The theory is that new users will take time to assimilate into the linguistic norms of the community

[https://en.wikipedia.org/wiki/Cross\\_entropy](https://en.wikipedia.org/wiki/Cross_entropy)

# Language models

## Empirical applications



(a) BeerAdvocate

Danescu-Niculescu-Mizil, Cristian, Robert West, Dan Jurafsky, Jure Leskovec, and Christopher Potts. 2013. "No Country for Old Members: User Lifecycle and Linguistic Change in Online Communities." In Proceedings of the 22nd International Conference on World Wide Web, 307–18. ACM. <http://dl.acm.org/citation.cfm?id=2488416>.

# Language models

## Limitations of N-gram language models

- ▶ Language use is much more complex than N-gram language models
- ▶ Three limitations
  - ▶ Insufficient data to sufficiently model language generation
  - ▶ Complex models become intractable to compute
  - ▶ Limited information on word order
- ▶ Next lecture we will see how advances in neural network models and the availability of large text corpora have opened up new avenues for language modeling and semantic analysis

# Summary

- ▶ Limitations of sparse representations of text
  - ▶ LSA allows us to project sparse matrix into a dense, low-dimensional representation
- ▶ Probabilistic language models allow us to directly model language use
- ▶ Next lecture: How neural language models allow us to create more meaningful semantic representations of texts