

The textmineR Package

Thomas W. Jones

16 May, 2016

1 Introduction

Introduction goes here

2 Mathematical Structures in Text Mining

Two core mathematical structures in text mining are the document term matrix (DTM) and the term co-occurrence matrix (TCM). A DTM is a matrix whose rows index documents and whose columns index linguistic features of these documents. These linguistic features called “terms”, though they may be single words, groups of words called “n-grams”, stems, lemmas, or other tokens. The (i, j) entries of a document term matrix are a frequency measure of term j within document i , for example the number of times term j is used within document i . A TCM is a square, but not necessarily symmetric, matrix whose rows and columns both index terms. The (i, j) entries of a TCM represent a relationship between term i and term j , for example the number of times term j appears within n places of term i .

Creating a DTM or TCM is the first step of a standard text mining pipeline, below. This pipeline should look familiar to many readers. It is similar to the set of steps for many traditional statistical analyses. Here, the data matrix (columns as variables, rows as observations) is the DTM/TCM. A model is fit and analysis performed with an eye towards generalizability.

1. Construct a DTM or TCM (the data matrix)
2. Fit a model or run an algorithm on the DTM or TCM
3. Analyze the model/algorithm’s results
4. Apply the model/algorithm to new documents (observations)

The core of many text mining tasks comes from the decision to make a DTM or a TCM and the definitions of rows, columns, and frequency measures of the DTM or TCM. Choices include: What is a term (or column)? This could be single words (unigrams), n ordered words (n-grams), words of only single parts of speech (for example, verbs), etc. What is a document (or row)? Raw documents may be parsed into smaller chunks, books to chapters, pages to paragraphs, etc. depending on the analyst’s goals. Finally, what is the quantitative relationship between rows and columns? For a DTM this is a question of using raw word counts, a frequency reweighting, or some other index. For a TCM, it is a question of the range of co-occurrence. Words may co-occur across a whole document or within a small window of n words.

Some choices for common tasks are below.

Probabilistic topic modeling

- DTM
- Rows are whole documents.
- Columns are unigrams or unigrams and bigrams
- Frequency measure is a raw integer count
- A probabilistic topic model, such as Latent Dirichlet Allocation (LDA), is applied to the DTM

Latent semantic analysis (LSA)

- DTM
- Rows are whole documents
- Columns are unigrams or unigrams and bigrams
- Frequency measure is a TF-IDF¹ reweighting
- The DTM is factored by single-value decomposition (SVD)

Calculating word embeddings

- TCM
- Rows and columns are unigrams
- Frequency measure is the number of times term j appears within 3 places of term i
- The TCM is de-composed by a method such as GloVe² or Word2Vec³

Document summarization

- DTM
- Each row is a sentence from a single document
- Columns are stems⁴ or lemmas⁵
- Frequency measure is TF-IDF reweighting
- Sentences (rows) are clustered, with cluster exemplars being chosen as a document summary.

3 Scalability Issues in Text Mining

Linguistic data is large and sparse by its nature. Because of the availability of digital data, many corpora comprise a large number of documents. Such size strains RAM and increases computation time. Both of these issues can make text mining challenging in practice.

Storing a DTM/TCM as a standard dense matrix may easily exceed available RAM on most systems. It is common for corpora to have many more unique terms than documents. A simple formula can be used to estimate the size a matrix of integers occupies in RAM in gigabytes. $\frac{8nk}{1,000,000,000}$

In the formula above, n is the number of rows, k is the number of columns, and 8 represents 8 bytes for each entry of the matrix. The numerator is in bytes; dividing by one billion converts it to gigabytes.

Consider the following examples of DTMs and TCMs stored as standard dense matrices in R. A 10,000 by 20,000 DTM, which represents a moderately-sized corpus with many terms removed, is about 1.6 GB. The same corpus with fewer terms removed, having 60,000 columns is about 4.8 GB. A DTM with 20,000 rows and 100,000 terms, a fairly common corpus size, is about 16 GB. A TCM of the same corpus is about 80 GB.

Sparse matrices find ways to not use RAM for entries that equal zero. The natural sparsity of linguistic data leads to significant space savings when representing a DTM/TCM in sparse form. A simple type of sparse matrix is a “triplet” matrix, so called because it has three columns. The first column, i indexes rows of the original matrix. The second column, j , indexes columns of the original matrix. The third column, v , is the value at the i, j entry of original matrix. Triplet matrices save space because i, j entries that equal zero are not given a row in the triplet matrix.

The scale of linguistic data makes computation time an issue too. Fortunately, many text mining computations are “embarrassingly parallel”. For example, removing a word from each document in a corpus may be done in parallel, as documents either do or do not contain the word, independent of other documents. Therefore, parallel computation along with sparse matrix representation eases text mining’s inherent scalability issues.

¹Term frequency inverse document frequency (TFIDF) is a reweighting where common words are penalized and rare words are promoted.

²cite

³cite

⁴Footnote or citation

⁵footnote or citation



Thomas Jones @thos_jones · 7 May 2015
Oops.

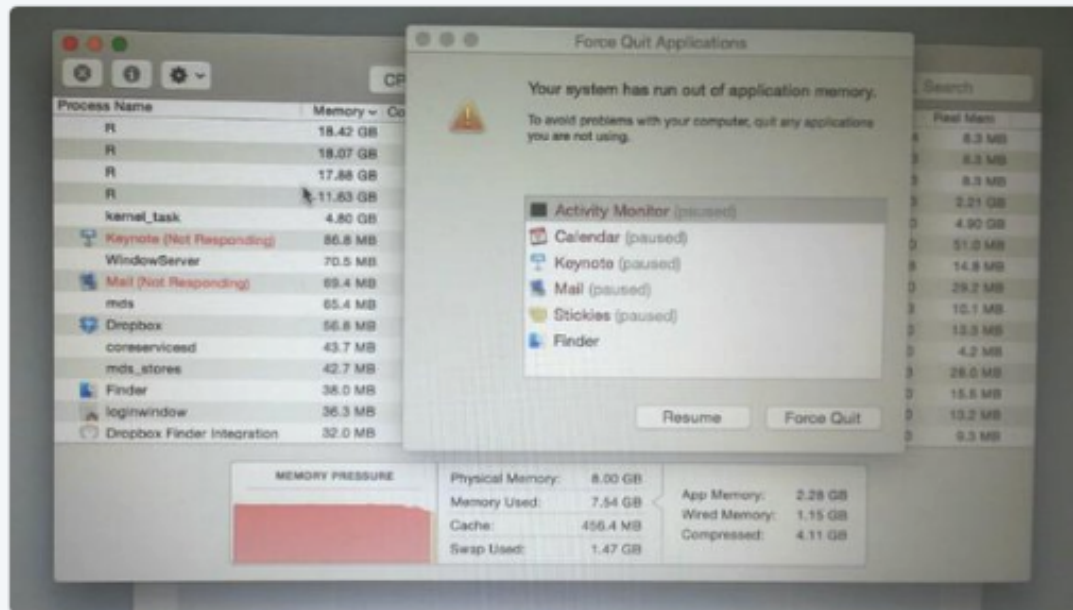


Figure 1: oops

4 An Ideal Framework for Textmining in R

4.1 The Ideal Framework

An ideal framework for text mining in R follows the below principles.

1. The DTM/TCM are the central data objects.
2. Models should, wherever possible, accept a DTM/TCM as input, not raw documents.
3. Classes, functions, and methods should be maximally interoperable within R's ecosystem.
4. Syntax should be idiomatic to R.
5. Objects and computation should scale to make working with large corpora easy.

An ideal framework holds the DTM/TCM as the central data objects for analysis. This framework has functions for creating and manipulating these matrices and its downstream modeling functions take the DTM/TCM as input. Unfortunately, many text mining programs (in and outside of R) are closed systems. They take raw documents, with some parameters given as arguments, and return the result of a model. The construction of the DTM/TCM is internal and analysts are constrained by choices made by the program's creators. Imagine a program for linear regression that took in raw survey responses, for example, and returned regression coefficients with little opportunity to transform your data before fitting a model! New regression programs would have to be built for every application area, as the raw input data may be significantly different. This is unfortunately common practice in text mining.⁶

In terms of syntax and interoperability, the ideal framework follows “the law of least astonishment”⁷. This framework uses classes that are part of R's base types except where absolutely necessary. When using non-base types (such as a sparse matrix) the class used by the ideal framework is widely adopted and has methods and functions allowing the non-base types to be manipulated as though they are base types.

4.2 `textmineR` is an Attempt at the Ideal

4.3 How `textmineR` Improves Upon Existing Frameworks in R

5 Overview of `textmineR`'s Main Functions

6 Two Example Analyses

7 Conclusion

OLD STUFF BELOW

4 An Ideal Approach to Natural Language Processing in R

An ideal NLP framework for R should have the following characteristics: It should be easy to use for an experienced R programmer. It should be maximally interoperable with other R packages, instead of being its own ecosystem. Its syntax should be intuitive for experienced R users. Finally, it should scale for large NLP tasks. `textmineR` is a best effort at these characteristics. Used in conjunction with the `text2vec` library, the framework espoused by `textmineR` scales for corpora much larger than available RAM.

⁶cite some examples, `mallet`, `word2vec`, others?

⁷cite The Tao of Programming <http://canonical.org/~kragen/tao-of-programming.html#book4>

4.1 The standard NLP pipeline

4.2 Interoperability within R’s ecosystem

To aid interoperability within R’s ecosystem, `textmineR` makes use of data types available in the `base` package wherever possible. In the case of DTMs and TCMs, it employs the `dgCMatrix` from the `Matrix` package. `dgCMatrix` objects are widely-adopted and have many methods whose syntax parallels methods available for `matrix` objects in the `base` package.

4.2.1 Corpus and metadata management

R has two core data types well-suited for storing documents and metadata: data frames and lists. A corpus with relatively simple metadata may be stored as a data frame, with one variable containing the text of the documents. An example is the sample NIH dataset included with `textmineR`⁸. This dataset, callable by `data(nih_sample)` is stored as a data frame. There are five columns containing unstructured or semi-structured textual data: `ABSTRACT_TEXT`, `NIH_SPENDING_CATS`, `PROJECT_TERMS`, `PHR`, and `PROJECT_TITLE`. These columns are described in the table below. The `nih_sample` dataset contains 39 additional columns of metadata. Metadata include administrative information, geographic information, temporal information, and more. More complex metadata may be stored in a list. There is no need for any special corpus object with its own functions and methods. R’s `base` library has this covered.

Variable Name	Variable Description
<code>ABSTRACT_TEXT</code>	The free-text abstract of an NIH-funded research project
<code>NIH_SPENDING_CATS</code>	Several key words tagged by NIH to the corresponding project corresponding to NIH’s reports of spending to Congress
<code>PROJECT_TERMS</code>	Additional key words tagged to the corresponding project
<code>PHR</code>	The free-text field describing the public-health relevance of the corresponding research project
<code>PROJECT_TITLE</code>	The free-text title of the research project

4.2.2 Sparse matrices for DTMs and TCMs

The `dgCMatrix`

`textmineR` need not be *the* framework for NLP in R. By building packages that operate on `dgCMatrix` objects, any number of packages can develop NLP techniques. `textmineR` and `text2vec` are already two libraries using `dcCMatrix` sparse matrices. In fact, `textmineR`’s functions `CreateDtm` and `CreateTcm` are wrappers simplifying the syntax (at the expense of some flexibility) of `text2vec`.

⁸cite.

4.3 textmineR’s philosophy on syntax

4.4 Scaling textmineR

5 Why textmineR Improves on tm and other NLP Frameworks

The CRAN Task View on Natural Language Processing promotes the `tm` package as a standard framework for NLP in R. Specifically it states

In recent years, we have elaborated a framework to be used in packages dealing with the processing of written material: the package `tm`. Extension packages in this area are highly recommended to interface with `tm`’s basic routines and useRs are cordially invited to join in the discussion on further developmens of this framework package.

`tm` is thoroughly programmed and well-established, having been around since at least 2008⁹. It has 29 reverse imports¹⁰ and depends as of this writing. `tm` has thorough vignettes demonstrating `tm`’s applicability to many use cases in NLP and instructions for writing extensions to `tm`. The `tm` package shines in its extensions for platforms such as LexisNexis¹¹, Factiva¹² and more.

Yet the `tm` package has two less-than-desireable properties: it is excessively object-oriented and it uses a somewhat esoteric data structure—with esoteric methods—as its sparse matrix data type. Unfortunately, this leads to a syntax that does not feel like an “R way” of doing things. The learning curve can be steep, even if the user is an experienced R programmer.

Consider the case of importing documents into `tm`’s framework as an example. Users cannot create a DTM directly from a character vector or other core R data type. Instead, they must first create a `Corpus` object. This `Corpus` object is designed to hold the text and metadata of documents. (This report argues that a `data.frame` or `list` is sufficient. More on this in the next section.) As a prerequisite for creating a `Corpus`, users must first create a `Source` object. There are different methods and functions for creating a `Source` object, depending on how the documents are stored internally or externally to R. Readers of this report may type `help(Source, package = "tm")` or `help(Corpus, package = "tm")` to see the complexity of these objects. Typing `help(meta, package = "tm")` gives insight to `tm`’s approach to metadata management.

The below code shows the differences in creating a DTM with `tm` and `textmineR`. The syntax of `textmineR` is much simpler, with transformations being passed as arguments to a single function. See section 6 for a deeper discussion of `textmineR`’s `CreateDtm` function and on extending the framework employed by `textmineR` through the `text2vec` library.

```
data(nih_sample)

docs <- nih_sample$ABSTRACT_TEXT

names(docs) <- nih_sample$APPLICATION_ID

stopwords <- c(tm::stopwords("english"), tm::stopwords("SMART"))

### Creating a DTM with the tm package from a character vector -----

# Create a corpus object from a vector source
corp <- tm::Corpus(tm::VectorSource(docs))
```

⁹cite JSS paper

¹⁰Including `textmineR`

¹¹See package `tm.plugin.lexisnexis` on CRAN

¹²See package `tm.plugin.factiva` on CRAN

```

# Perform document curation, lowering, removing non-alpha character,
# removal of stopwords, etc.
corp <- tm::tm_map(corp, tm::content_transformer(tolower))

corp <- tm::tm_map(x=corp, tm::removeWords, stopwords)

corp <- tm::tm_map(corp, tm::removePunctuation)

corp <- tm::tm_map(corp, tm::removeNumbers)

corp <- tm::tm_map(corp, tm::stripWhitespace)

# Create a final DTM
dtm <- tm::DocumentTermMatrix(corp)

### Creating a DTM with the textmineR package from a character vector -----

# note lowering, removal of non-alpha characters is the default behavior
dtm <- CreateDtm(doc_vec = docs,
                 ngram_window = c(1, 1),
                 stopword_vec = stopwords,
                 lower = TRUE,
                 remove_punctuation = TRUE,
                 remove_numbers = TRUE)

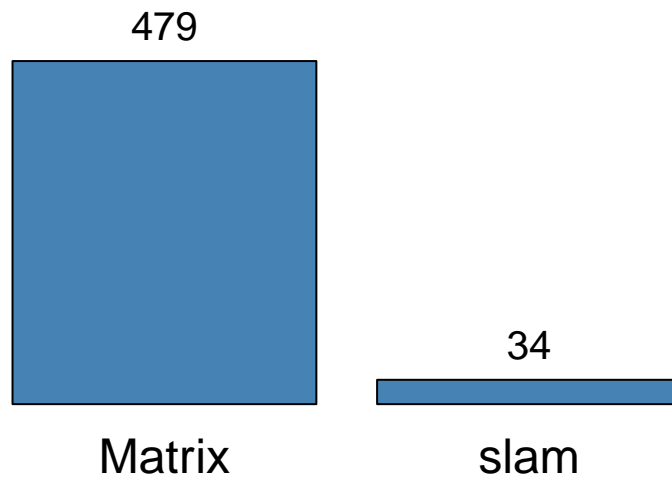
```

It is clear from the two types of `Corpus` object, that the makers of `tm` are concerned with scalability issues. `Corpus` objects of type `VCorpus` are designed to be held in RAM, whereas `Corpus` types of `PCorpus` are designed to be held on disk, facilitating analyses of massive corpora. For most users, this much data is an exception rather than the rule. (cite HH analyzing the analyzers) It compels the question of whether or not the additional complexity and overhead are necessary for a corner case.

The `tm` package stores DTMs as a `simple_triplet_matrix` object from the `slam` package. Simple triplet matrices are a form of sparse matrix where only non-zero entries are stored. The “triplet” comes from the three standard columns: `i` is a row index, `j` is a column index, and `v` for the value at the (i, j) position. (i, j) entries that are not indexed in the `i` and `j` columns are assumed to be zero. There are two limitations to `slam`’s `simple_triplet_matrix` class, however. Methods and functions for matrix manipulation and matrix math differ significantly from standard R dense matrices. This, again, makes for a steeper learning curve even when users are experienced R programmers. Second, `simple_triplet_matrix` objects are less-commonly used than other sparse matrix classes, notably from the `Matrix` package. (More on this in the next section.) This limits the availability of out-of-the-box statistical methods that can be performed on a DTM created by the `tm` package.

The figure below illustrates

Number of Reverse Dependencies, Imports, Suggests, etc.

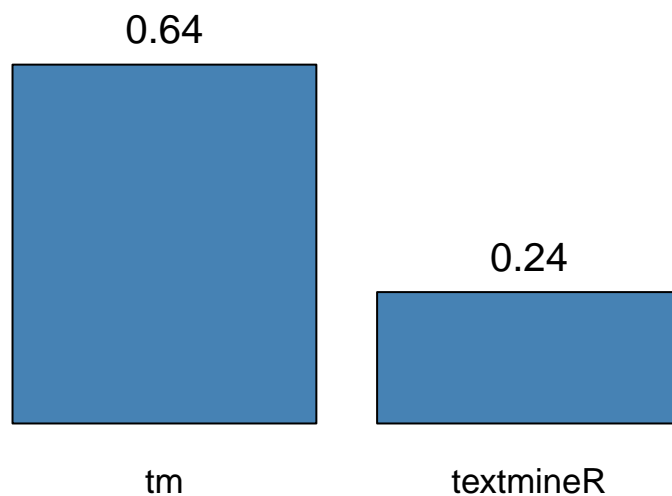


6 An Overview of textmineR's Functions

6.1 textmineR's core functions: CreateDtm and CreateTcm

`textmineR` has two core functions. `CreateDtm` creates document term matrices from a character vector. `CreateTcm` creates term co-occurrence matrices from a character vector. Both functions have arguments allowing users to remove stop words, remove numbers, remove punctuation, tokenize unigrams and n-grams on spaces, convert terms to lowercase, and pass custom stemming and lemmatization functions. In addition, both automatically make use of parallelism on Windows and Unix-based operating systems. This parallelism is executed at the C++ level through `text2vec` and at the R level through `textmineR`'s `TmParallelApply` function.

Mean Seconds to Create a DTM of 500 NIH Grant Abstracts



6.2 Topic models available in textmineR

6.2.1 Latent Semantic Analysis

```
data(nih_sample_dtm)

lsa <- FitLsaModel(dtm = nih_sample_dtm, k = 10)
```

6.2.2 Latent Dirichlet Allocation

```
data(nih_sample_dtm)

# LDA Model with Gibbs Sampling
lda <- FitLdaModel(dtm = nih_sample_dtm, k = 10, alpha = 0.1, beta=0.05,
                  method = "gibbs")
```

```
data(nih_sample_dtm)

# LDA Model with Variational EM
lda <- FitLdaModel(dtm = nih_sample_dtm, k = 10, alpha = 0.1, beta = 0.05,
                  method = "vem")
```

6.2.3 Correlated Topic Models

```
data(nih_sample_dtm)

ctm <- FitCtmModel(dtm = nih_sample_dtm, k = 10)
```

6.2.4 Document Clustering as a Topic Model

```
data(nih_sample)
data(nih_sample_dtm)

# Take the IC_NAME column as a clustering of abstracts
model <- Cluster2TopicModel(dtm = nih_sample_dtm,
                           clustering = nih_sample$IC_NAME)
```

6.3 Topic model utility functions in textmineR

6.3.1 CalcTopicModelR2

```
data(nih_sample_dtm)
data(nih_sample_topic_model)
```

```
# Get the R-squared of the model
r2 <- CalcTopicModelR2(dtm = nih_sample_dtm,
                      phi = nih_sample_topic_model$phi,
                      theta = nih_sample_topic_model$theta)
```

6.3.2 CalcLikelihood

```
data(nih_sample_dtm)
data(nih_sample_topic_model)

# Get the likelihood of the data given the fitted model parameters
ll <- CalcLikelihood(dtm = nih_sample_dtm,
                    phi = nih_sample_topic_model$phi,
                    theta = nih_sample_topic_model$theta)
```

6.3.3 CalcProbCoherence

```
data(nih_sample_topic_model)
data(nih_sample_dtm)

CalcProbCoherence(phi = nih_sample_topic_model$phi,
                  dtm = nih_sample_dtm,
                  M = 5)
```

6.3.4 CalcPhiPrime

```
data(nih_sample_topic_model)
data(nih_sample_dtm)

# Make a phi_prime matrix, P(topic/words)
phi_prime <- GetPhiPrime(phi = nih_sample_topic_model$phi,
                        theta = nih_sample_topic_model$theta,
                        p_docs = rowSums(nih_sample_dtm))
```

6.3.5 GetTopTerms

```
data(nih_sample_dtm)
data(nih_sample_topic_model)

top_terms <- GetTopTerms(phi = nih_sample_topic_model$phi,
                        M = 5)
```

6.3.6 LabelTopics

```

data(nih_sample)
data(nih_sample_topic_model)

# Create a new DTM with n-grams only
dtm_ngram <- CreateDtm(nih_sample$ABSTRACT_TEXT,
                      doc_names = nih_sample$APPLICATION_ID,
                      ngram_window = c(2, 2))

# Create a matrix of assignments that give a hard in/out topic assignment per document
assignments <- t(apply(nih_sample_topic_model$theta, 1, function(x){
  x[ x < 0.05 ] <- 0
  x / sum(x)
})))

# Label topics with our new n-gram DTM
labels <- LabelTopics(assignments = assignments, dtm = dtm_ngram, M = 2)

```

6.4 Broadly-applicable functions available in textmineR

6.4.1 CalcHellingerDist

```

# Generate some random vectors
x <- rchisq(n = 100, df = 8)

y <- x ^ 2

myamat <- rbind(x, y)

# Get the Hellinger distance between them
CalcHellingerDist(x = myamat)

```

6.4.2 CalcJSDivergence

```

# Generate some random vectors
x <- rchisq(n = 100, df = 8)

y <- x ^ 2

myamat <- rbind(x, y)

# Get the Jensen-Shannon Divergence between them
CalcJSDivergence(x = myamat)

```

6.4.3 Dtm2Docs

```

data(nih_sample)
data(nih_sample_dtm)

```

```

# see the original documents
nih_sample$ABSTRACT_TEXT[ 1:3 ]

# see the new documents re-structured from the DTM
new_docs <- Dtm2Docs(dtm = nih_sample_dtm)

new_docs[ 1:3 ]

```

6.4.4 TmParallelApply

6.4.5 RecursiveRbind

6.5 Using textmineR with other NLP frameworks

7 Examples

7.1 Document Clustering

Describe document clustering

describe dtm creation and vocab curation

```

data(nih_sample)

# Create a DTM of abstracts with unigrams and bigrams, stopwords removed
dtm <- CreateDtm(doc_vec = nih_sample$ABSTRACT_TEXT,
                 doc_names = nih_sample$APPLICATION_ID,
                 ngram_window = c(1, 2))

```

describe tf-idf frequency reweighting

```

# TF-IDF Frequency re-weighting
idf <- log(nrow(dtm) / colSums(dtm > 0))

tfidf <- t(dtm) * idf

tfidf <- t(tfidf)

```

describe cosine similarity describe how dist object now works with any R clustering function

```

# Calculate document-to-document cosine similarity
csim <- tfidf / sqrt(rowSums(tfidf * tfidf))

csim <- csim %*% t(csim)

# Create a dist object of cosine distances
cdist <- as.dist(as.matrix(csim))

# Create an hclust object
doc_hclust <- hclust(cdist, method = "ward.D")

```

describe selection by silhouette

```
# Choose number of clusters with silhouette
silh <- parallel::mclapply(2:99, function(k){

  clust <- cutree(doc_hclust, k= k)

  s <- cluster::silhouette(clust, dist = cdist)

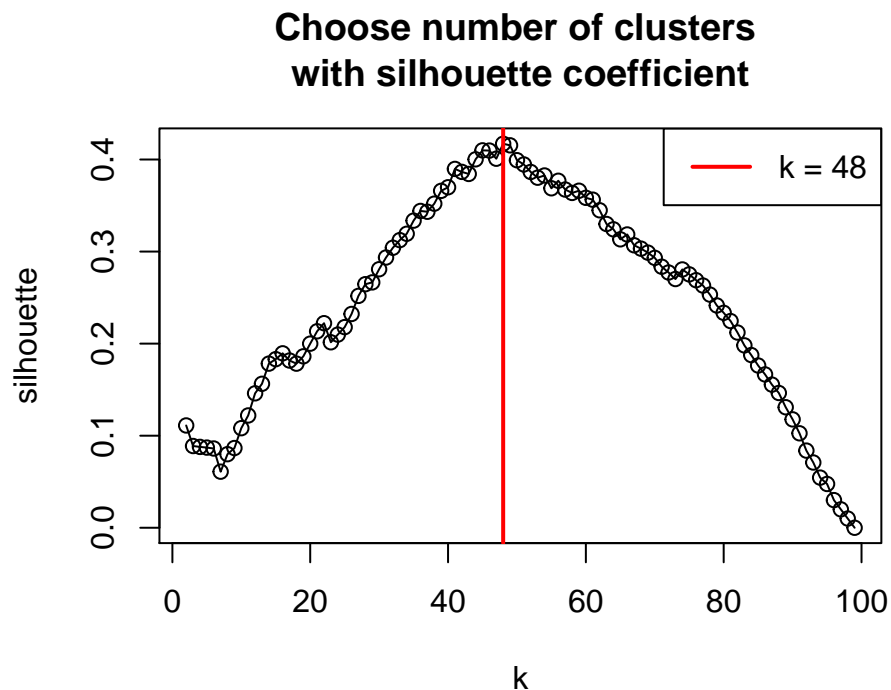
  s <- summary(s)[ "avg.width" ]

  data.frame(k = k, silhouette = as.numeric(s), stringsAsFactors = F)

}, mc.cores = parallel::detectCores())

silh <- do.call(rbind, silh)
```

describe selection by silhouette



final k

```
k <- silh$k[ silh$silhouette == max(silh$silhouette) ]

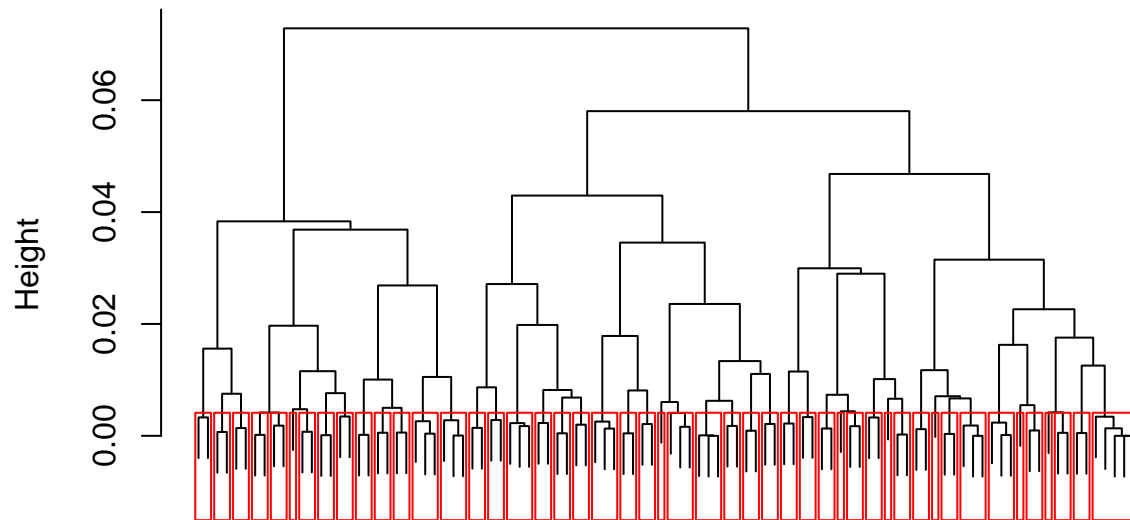
doc_hclust$clustering <- cutree(doc_hclust, k = k)
```

describe final clustering

```
# Plot the document clustering
plot(doc_hclust, labels=rep("", nrow(dtm)))

rect.hclust(doc_hclust, k = k)
```

Cluster Dendrogram



cdist
hclust (*, "ward.D")

Figure 2:

```
# Get the titles of documents in cluster 2
nih_sample$PROJECT_TITLE[ doc_hclust$clustering == 2 ]
```

```
## [1] "Single molecule imaging of HIV Env"
## [2] "Overall Administration of Rare Diseases Clinical Research Consortia (RDCRC) "
## [3] "Direct Stimulation of Bacterial Metabolism with Applied Electrical Current"
```

7.2 Latent Dirichlet Allocation

Explain LDA overview

```
# Load some data into the workspace
data(nih_sample)

# Create a document term matrix
dtm <- CreateDtm(nih_sample$ABSTRACT_TEXT[ 1:99 ],
                 doc_names = nih_sample$APPLICATION_ID[ 1:99 ],
                 ngram_window = c(1, 2))

# explore basic frequencies & curate vocabulary
tf <- TermDocFreq(dtm = dtm)
```

```
# Eliminate words appearing less than 2 times or in more than half of the
# documents
vocabulary <- tf$term[ tf$term_freq > 1 & tf$doc_freq < nrow(dtm) / 2 ]

dtm <- dtm[ , vocabulary]
```

discuss selecting the number of topics and modeling choices

```
model <- FitLdaModel(dtm = dtm, k = 10,
                    alpha = 0.1, beta = 0.05,
                    iterations = 500,
                    method = "gibbs")

names(model) # phi is P(words | topics), theta is P(topics | documents)
```

```
## [1] "theta" "phi"
```

explain R-squared

```
# Get the R-squared of this model
model$r2 <- CalcTopicModelR2(dtm = dtm, phi = model$phi, theta = model$theta)

model$r2
```

```
## [1] 0.2606902
```

Explain model top terms, probabilistic coherence, topic labeling, topic prevalence

```
# top 5 terms of the model according to phi & phi-prime
model$top_terms <- GetTopTerms(phi = model$phi, M = 5)

# probabilistic coherence counting top 5 terms in each topic
model$coherence <- CalcProbCoherence(phi = model$phi, dtm=dtm)

# give a hard in/out assignment of topics in documents
model$assignments <- model$theta

model$assignments[ model$assignments < 0.05 ] <- 0

model$assignments <- model$assignments / rowSums(model$assignments)

model$assignments[ is.na(model$assignments) ] <- 0

# Get some topic labels using n-grams from the DTM
model$labels <- LabelTopics(assignments = model$assignments,
                           dtm = dtm,
                           M = 2)

model$doc_count <- colSums(model$assignments > 0)
```

```

# Create a summary matrix to view topics
model$topic_summary <- data.frame(topic = rownames(model$phi),
                                top_terms = apply(model$top_terms, 2,
                                                  function(x) paste(x, collapse=" ", ")),
                                labels = apply(model$labels, 1,
                                              function(x) paste(x, collapse=" ", ")),
                                coherence = round(model$coherence, 3),
                                doc_count = model$doc_count,
                                stringsAsFactors=FALSE)

```

explain classifying a new document

```

# phi-prime, P(topic | words) for classifying new documents
model$phi_prime <- CalcPhiPrime(phi = model$phi, theta = model$theta, p_docs = rowSums(dtm))

# vectorize our new document
new_doc <- nih_sample$ABSTRACT_TEXT[ 100 ]

names(new_doc) <- nih_sample$APPLICATION_ID[ 100 ]

new_dtm <- CreateDtm(new_doc, ngram_window = c(1, 2))

# classify new topics with phi_prime
common_vocab <- intersect(colnames(new_dtm), colnames(model$phi_prime))

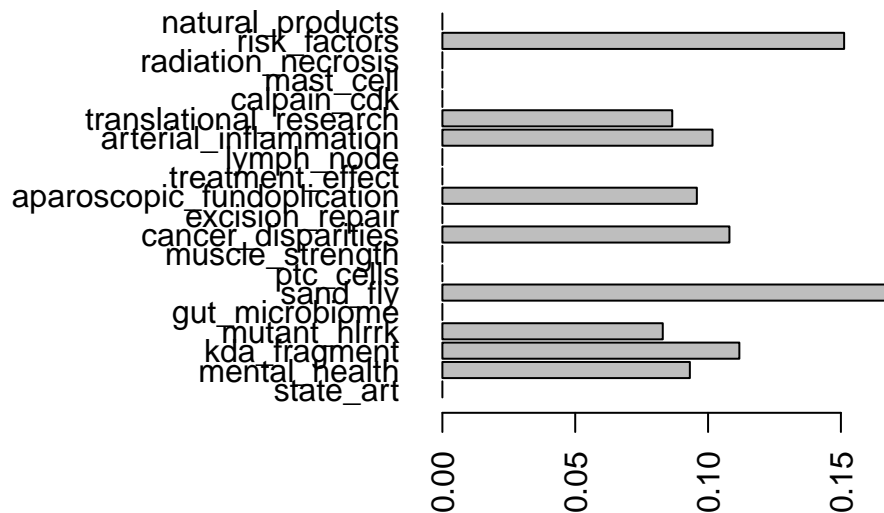
new_topics <- new_dtm[ , common_vocab ] %*% t(model$phi_prime[ , common_vocab])

# normalize rows
new_topics <- new_topics / rowSums(new_topics)

# give a hard in/out assignment and re-normalize
new_topics[ new_topics < 0.05 ] <- 0

new_topics <- new_topics / rowSums(new_topics)

```



7.3 Word embeddings with GloVe from text2vec

```
data(nih_sample)

tcm <- CreateTcm(doc_vec = nih_sample$ABSTRACT_TEXT, skipgram_window = 3)

glove_model <- text2vec::glove(tcm = tcm,
                              word_vectors_size = 10,
                              x_max = 10, learning_rate = 0.2,
                              num_iters = 50, grain_size = 1e5,
                              max_cost = 100, convergence_threshold = 0.005 )
```

```
## 2016-05-16 00:20:07 - epoch 1, expected cost 0.0401

## 2016-05-16 00:20:07 - epoch 2, expected cost 0.0297

## 2016-05-16 00:20:07 - epoch 3, expected cost 0.0258

## 2016-05-16 00:20:07 - epoch 4, expected cost 0.0231

## 2016-05-16 00:20:07 - epoch 5, expected cost 0.0212

## 2016-05-16 00:20:07 - epoch 6, expected cost 0.0196

## 2016-05-16 00:20:07 - epoch 7, expected cost 0.0183

## 2016-05-16 00:20:07 - epoch 8, expected cost 0.0172

## 2016-05-16 00:20:07 - epoch 9, expected cost 0.0163

## 2016-05-16 00:20:07 - epoch 10, expected cost 0.0155

## 2016-05-16 00:20:07 - epoch 11, expected cost 0.0149

## 2016-05-16 00:20:07 - epoch 12, expected cost 0.0143

## 2016-05-16 00:20:07 - epoch 13, expected cost 0.0137

## 2016-05-16 00:20:08 - epoch 14, expected cost 0.0133

## 2016-05-16 00:20:08 - epoch 15, expected cost 0.0128

## 2016-05-16 00:20:08 - epoch 16, expected cost 0.0124

## 2016-05-16 00:20:08 - epoch 17, expected cost 0.0121

## 2016-05-16 00:20:08 - epoch 18, expected cost 0.0117
```

2016-05-16 00:20:08 - epoch 19, expected cost 0.0115

2016-05-16 00:20:08 - epoch 20, expected cost 0.0112

2016-05-16 00:20:08 - epoch 21, expected cost 0.0109

2016-05-16 00:20:08 - epoch 22, expected cost 0.0107

2016-05-16 00:20:08 - epoch 23, expected cost 0.0105

2016-05-16 00:20:08 - epoch 24, expected cost 0.0102

2016-05-16 00:20:08 - epoch 25, expected cost 0.0101

2016-05-16 00:20:08 - epoch 26, expected cost 0.0099

2016-05-16 00:20:08 - epoch 27, expected cost 0.0097

2016-05-16 00:20:08 - epoch 28, expected cost 0.0095

2016-05-16 00:20:08 - epoch 29, expected cost 0.0094

2016-05-16 00:20:08 - epoch 30, expected cost 0.0092

2016-05-16 00:20:08 - epoch 31, expected cost 0.0091

2016-05-16 00:20:08 - epoch 32, expected cost 0.0090

2016-05-16 00:20:08 - epoch 33, expected cost 0.0088

2016-05-16 00:20:08 - epoch 34, expected cost 0.0087

2016-05-16 00:20:08 - epoch 35, expected cost 0.0086

2016-05-16 00:20:08 - epoch 36, expected cost 0.0085

2016-05-16 00:20:08 - epoch 37, expected cost 0.0084

2016-05-16 00:20:08 - epoch 38, expected cost 0.0083

2016-05-16 00:20:08 - epoch 39, expected cost 0.0082

2016-05-16 00:20:08 - epoch 40, expected cost 0.0081

2016-05-16 00:20:08 - epoch 41, expected cost 0.0080

2016-05-16 00:20:08 - epoch 42, expected cost 0.0080

```
## 2016-05-16 00:20:08 - epoch 43, expected cost 0.0079
## 2016-05-16 00:20:08 - epoch 44, expected cost 0.0078
## 2016-05-16 00:20:09 - epoch 45, expected cost 0.0077
## 2016-05-16 00:20:09 - epoch 46, expected cost 0.0077
## 2016-05-16 00:20:09 - epoch 47, expected cost 0.0076
## 2016-05-16 00:20:09 - epoch 48, expected cost 0.0075
## 2016-05-16 00:20:09 - epoch 49, expected cost 0.0075
## 2016-05-16 00:20:09 - epoch 50, expected cost 0.0074
```

```
glove_model$phi <- t(Reduce("+", glove_model$word_vectors) / 2)
colnames(glove_model$phi) <- colnames(tcm)
```

7.4 Document summarization