


Topic Modeling

EDP 618 Week 12

Dr. Abhik Roy

Setting Up



1. You can retrieve the *employee sample reviews* survey response data set and both installation and walkthrough  *scripts* by clicking on the icon below



2. Open up RStudio
3. Open up `Topic Modeling Install.R`
4. Open `Topic Modeling Script.R`

Take a look at the various types of files that can be imported in the tidyverse



Getting Prepped



In `Topic Modeling Script.R`, run the following commands

1. Setting the working directory as source

```
setwd(dirname(rstudioapi::getActiveDocumentContext())$path))
```

2. Loading the needed packages for this walkthrough

```
library(tidyverse)
library(tidytext)
library(tm)
library(textclean)
library(topicmodels)
library(ldatuning)
library(stopwords)
library(textstem)
library(broom)
```

Alternatively if you have the `pacman` package, run `pacman::p_install("tidyverse", "tidytext", "tm", "textclean", "topicmodels", "ldatuning", "stopwords", "textstem", "broom")`



3. Bringing in the survey data

```
employee_responses <-  
  read_csv("employee sample reviews.csv")
```

4. Bringing in the common names set (as a character string)

```
common_names <-  
  read_csv("most common names.csv") %>%  
  simplify_all() %>%  
  .[[1]]
```

5. Retrieving stopwords

```
data("stop_words")
```

Too Many Files?



You may have noticed that there are a lot of files. Some of them we'll use while the others are included for completeness. Below you will see a description of each

Data sets we are using

`employee_sample_reviews.csv` is a 10% random sampling of the original data set that was created to save computing time for the example given in the walkthrough

`most_common_names.csv` is a list of approximately 97,000+ common U.S. names derived from both U.S. Census and Social Security Administration data

Data sets we are not using

`employee_reviews.csv` is the original data set and may be found on [Kaggle] (<https://www.kaggle.com/datasets/fiodarryzhikau/employee-review>)

R Scripts we are using

`Topic Modeling Install.R` is the installation file needed for the walkthrough

`Topic Modeling Script.R` is a static copy of the commands used in the walkthrough

R scripts we are not using

`Get Common Names.R` provides a list of **dplyr** commands that uses the **lexicon**, **babynames**, and **genderdata** packages to create the most common names data set

`Random Sampling Rows.R` provides a list of **dplyr** commands used to create the sample data set

Before We Begin

This is the process we'll cover lightly. There is a lot more going on under the hood and you may not be able to recognize all of the terms, but if you can get a basic understanding of the process, the rest can be filled in by conducting a topic model!



If you didn't know, computers can't understand human languages...not directly anyway. Enter this idea below of using a medium to communicate with one (or multiple)

NLP

(Natural Language Processing)

- Provides computers with an ability to understand and process open text and spoken word
- Combines computational linguistics with statistical models

Abhik Roy

Here are a few things we won't be covering in this session so please read over the areas you lack familiarity with. Given that, it is absolutely fine if you cannot fully understand all of these ideas right now - they will hopefully become apparent as we progress



DOCUMENT

- a distinct piece of text
- could be an article, book, chapter, paragraph, sentence, etc.
- dependent on needs

Abhik Roy

CORPUS

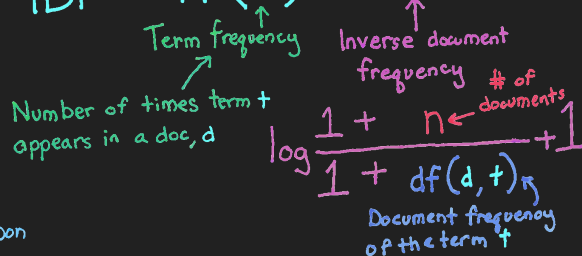
- Latin for body
- Collection of documents

Abhik Roy

TF-IDF

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$TF-IDF = TF(t, d) \times IDF(t)$$



Chris Albon

LDA

(Latent Dirichlet Allocation)

- each document is made of topics
- each topic is made of words
- discovers topics into a collection of documents
- tags each document with topics

Abhik Roy

Here are some basic terms you should try to keep while going through the walkthrough. Again it is completely fine if you do not understand what these mean in context right now!



BAG OF WORDS

Converts text to a matrix where every row is an observation and every feature is a unique word. The value of each element in the matrix is either a binary indicator marking the presence of that word or an integer of the number of times that word appears.

ChrisAlbon

STANDARDIZATION

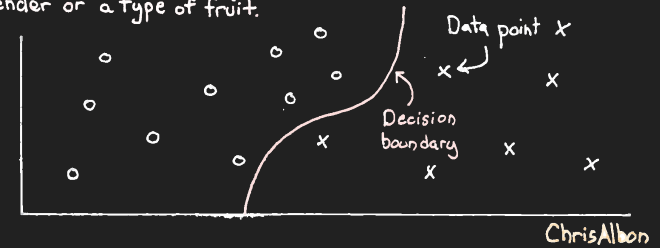
$$X'_i = \frac{\text{Value of the } i\text{th observation } X_i - \text{Mean of the feature vector } \bar{X}}{\text{Standard deviation of the feature vector } \sigma}$$

Standardization is a common scaling method. X'_i represents the number of standard deviations each value is from the the mean value. It rescales a feature to have a mean of 0 and unit variance.

ChrisAlbon

CLASSIFICATION

Classification problems are when we are training a model to predict qualitative targets. For example: gender or a type of fruit.



TOKENIZING TEXT

Splitting up text into individual units like paragraphs, sentences, or words.

Example:

"I like birds" → "I", "like", "birds"

ChrisAlbon

Topic Modeling

A type of probabilistic statistical model for

- (a) discovering the abstract "topics"
- or *hidden semantic structures* -
that occur in a collection of documents
- (b) dimensionality reduction

The Most Annoying Thing About Data



The 80/20 Rule¹: *Most data scientists spend only 20 percent of their time on actual data analysis and 80 percent of their time finding, cleaning, and reorganizing huge amounts of data*

¹ Loosely based on an idea called **Pareto's Principle** which states that *roughly 80% of outcomes come from 20% of causes*

Step 1: Assessing Data



1. Take a look at the data set and think about categorizing terms that may skew how terms are assessed

the names of the people are not important so we could replace all of them simply with the word **people**

employees are prevalent in the data so we could replace the word **people** altogether

2. Open up an empty text document and try going through on your own to consider terms that could be collapsed

Step 2: Preprocessing



```
employee_responses %>%
  select(feedback) %>% # Select the column with open ended responses
  mutate(feedback = textclean::replace_non_ascii(feedback)) %>% # Convert to a standard form
  mutate(feedback = str_to_lower(feedback)) %>% # Convert all words to lower case
  mutate(feedback = str_remove_all(feedback, "'s")) %>% # Remove all cases of 's
  mutate(feedback = str_remove_all(feedback, "[[:digit:]]")) %>% # Remove all numbers
  mutate(feedback = str_remove_all(feedback, "[[:punct:]]")) %>% # Remove all punctuation
# Remove all instances from a separate list
mutate(feedback = str_remove_all(feedback, paste0("\\b", common_names, "\\b", collapse = "|")));
# Remove any additional terms manually
mutate(feedback = str_remove_all(feedback, "mcknight|cook|cunningham|hahn|vargas")) %>%
mutate(feedback = lemmatize_strings(feedback)) %>% # Lemmatize terms
mutate(feedback = str_squish(feedback)) %>% # Remove whitespace
mutate(feedback = na_if(feedback, "")) %>% # Replace blanks with NA
drop_na() # Drop all columns with NA
```

```
## # A tibble: 88 × 1
##   feedback
##   <chr>
## 1 require additional scopee willingness to self start outperform task motivate to exceed
## 2 basically have high potential his work experience but recent day he be do performance his core i th
## 3 be a fairly average worker put out satisfactory work and show potential i would to a bite much from
## 4 i be write a review for while he be a person unfortunately that doesnt show up his performance at w
## 5 perform at a high level this past year be one of much reliable and solid performer on team have a f
## 6 be a follower rather a leader while do complete task a timely manner do without be inventive team s
## 7 be a very capable worker and always produce excellent work always finish work and diligently pursue
## 8 be a very capable worker and always produce excellent work always finish work and diligently pursue
## 9 be a very capable worker and always produce excellent work always finish work and diligently pursue
## 10 while i have find some issue with performance past specifically a lack of attention to detail i fin
## # ... with 78 more rows
```

Please note that removing all instances from a separate list may take up to a minute to complete

Assigning a Variable



Let's save the entire cleaning process

```
responses_cleaned <-  
  employee_responses %>%  
  select(feedback) %>%  
  mutate(feedback = textclean::replace_non_ascii(feedback)) %>%  
  mutate(feedback = str_to_lower(feedback)) %>%  
  mutate(feedback = str_remove_all(feedback, "'s")) %>%  
  mutate(feedback = str_remove_all(feedback, "[[:digit:]]")) %>%  
  mutate(feedback = str_remove_all(feedback, "[[:punct:]]")) %>%  
  mutate(feedback = str_remove_all(feedback, paste0("\\b", common_names, "\\b", collapse = "|"))) %>%  
  mutate(feedback = str_remove_all(feedback, "mcknight|cook|cunningham|hahn|vargas")) %>%  
  mutate(feedback = lemmatize_strings(feedback)) %>%  
  mutate(feedback = str_squish(feedback)) %>%  
  mutate(feedback = na_if(feedback, "")) %>%  
  drop_na()
```


What Just Happened?



Let's try doing something similar but with shorter and simpler text taken from the very funny skit [Sharknado Pitch Meeting](#)

```
example_text <-  
  c("Excerpt from Sharknado Pitch Meeting.  
    Creator: Ryan George.  
  
    (1) It's peer reviewed.  
    (2) Multiple scientists looked over that and approved of it?  
    (3) No some drunk guy on the pier checked it out. He loved it!  
    (4) That is technically peer reviewed. I think we're good.  
  
    --The End--  
  ")
```



```
example_text %>%  
  read_lines() %>% # Parse text into individual lines  
  as_tibble_col("text") %>% # Create a single tidy column  
  slice(4:n()) %>% # Remove unnecessary text  
  mutate(text = textclean::replace_non_ascii(text)) %>% # Convert to a standard format  
  mutate(text = str_to_lower(text)) %>% # Convert all words to lower case  
  mutate(text = str_remove_all(text, "'s")) %>% # Remove all cases of `s  
  mutate(text = str_remove_all(text, "[[:digit:]]")) %>% # Remove all numbers  
  mutate(text = str_remove_all(text, "[[:punct:]]")) %>% # Remove all punctuation  
  mutate(text = str_remove_all(text, "the end")) %>% # Remove term  
  mutate(text = str_replace_all(text, "multiple scientists", "scientists")) %>% # Replace term  
  mutate(text = str_replace_all(text, "it", "paper")) %>% # Replace term  
  mutate(text = str_replace_all(text, "that", "paper")) %>% # Replace term  
  mutate(text = lemmatize_strings(text)) %>% # Lemmatize term  
  mutate(text = str_remove_all(text, c("paper"))) %>% # Remove term  
  mutate(text = str_squish(text)) %>% # Remove whitespace  
  mutate(text = na_if(text, "")) %>% # Replace blanks with NA  
  drop_na() # Drop all columns with NA
```

```
## # A tibble: 4 × 1  
##   text  
##   <chr>  
## 1 peer review  
## 2 scientist look over and approve of  
## 3 no some drink guy on the pier check out he love  
## 4 be technically peer review i think be good
```

Normalization of Remaining Wording

is used to reduce word randomness which allows some level of standardization to help to reduce the amount of different information that a computer has to process therefore improving efficiency

the overall goal is to reduce inflectional forms and sometimes derivationally related forms of a word to a common base form

two popular normalization techniques are *lemmatization* and *stemming*

Lemmatization vs. Stemming

Lemmatization

the process of reducing words to their base word
(takes more time)

Example

the term *better* has the lemma *good*

Stemming

the process of reducing words to their word stem or root form
by removing word endings or other affixes
(takes less time)

Example

the term *flooding* has the stem *flood*

Sentence

My friend had a beautiful singing voice

Lemmatization

my friend have a beautiful singing voice

Stemming

my friend had a beauti sing voic

Tokenizing Handled Data



A process of distinguishing and classifying sections of a string of input characters

What you should take from this is that **unnesting** data successfully is a requirement to be able to **tokenize**. While the next set of commands should look familiar, please consider taking a bit of time to really see what occurs in each step



```
responses_cleaned %>%  
  unnest_tokens(word, feedback) %>%  
  anti_join(stop_words) %>%  
  count(word, sort = TRUE) %>%  
  add_column(document = 1)
```

```
## Joining, by = "word"  
  
## # A tibble: 579 × 3  
##   word          n document  
##   <chr>      <int>   <dbl>  
## 1 time         40         1  
## 2 improve     38         1  
## 3 team        31         1  
## 4 performance 27         1  
## 5 potential   26         1  
## 6 company     23         1  
## 7 task        23         1  
## 8 skill       18         1  
## 9 worker      17         1  
## 10 complete   15         1  
## # ... with 569 more rows
```

Assign a Variable



Let's save the tokenized data frame

```
responses_tokens <-  
  responses_cleaned %>%  
  unnest_tokens(word, feedback) %>%  
  anti_join(stop_words) %>%  
  count(word, sort = TRUE) %>%  
  add_column(document = 1)  
  
## Joining, by = "word"
```


Step 3: Statistical Classification and Modeling

```
responses_tokens %>%
```

```
cast_dtm(document, word, n)
```

```
## <<DocumentTermMatrix (documents: 1, terms: 16)>>  
## Non-/sparse entries: 579/0  
## Sparsity           : 0%  
## Maximal term length: 16  
## Weighting          : term frequency (tf)
```



Assigning a Variable

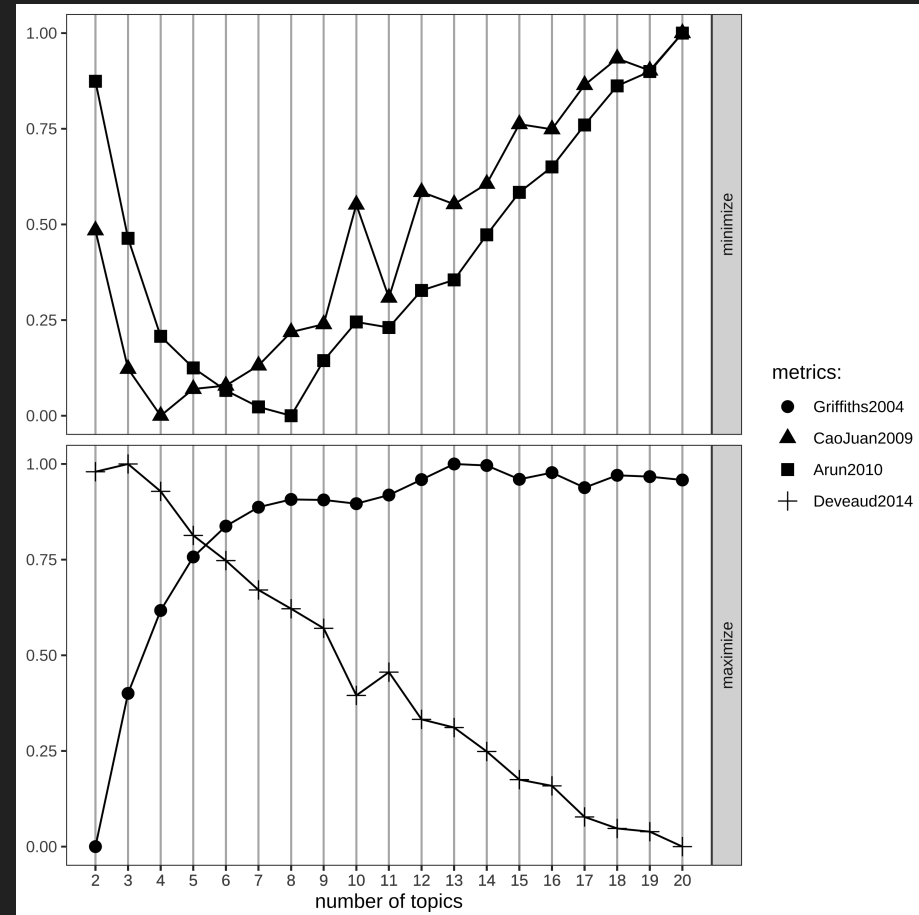
```
responses_dtm <-  
  responses_tokens %>%  
  cast_dtm(document, word, n)
```



```
FindTopicsNumber(
  responses_dtm,
  topics = seq(from = 2, to = 20, by = 1),
  metrics = c("Griffiths2004",
              "CaoJuan2009",
              "Arun2010",
              "Deveaud2014"),
  method = "Gibbs",
  control = list(seed = 77),
  mc.cores = 2L,
  verbose = TRUE
) %>%
```

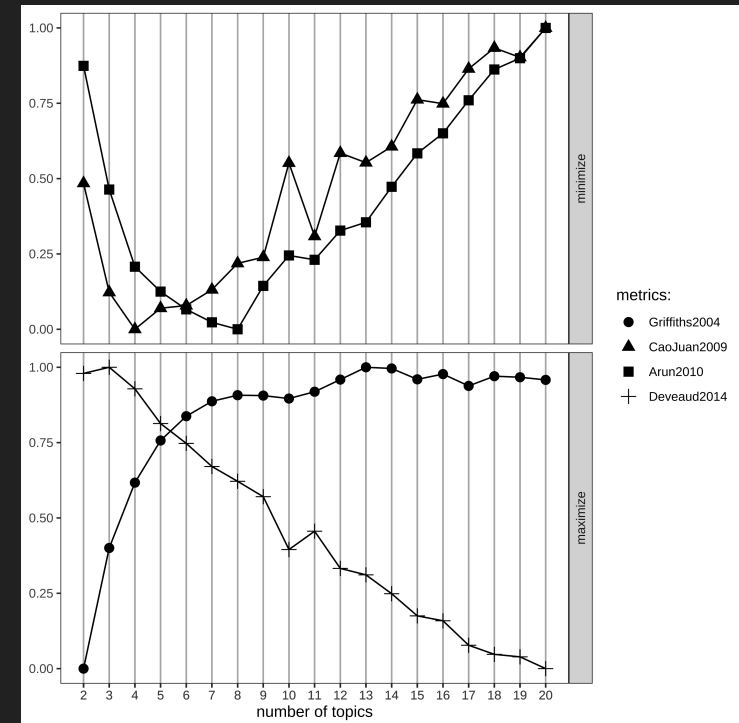
```
FindTopicsNumber_plot()
```

```
## fit models... done.
## calculate metrics:
## Griffiths2004... done.
## CaoJuan2009... done.
## Arun2010... done.
## Deveaud2014... done.
```



Assigning a Variable

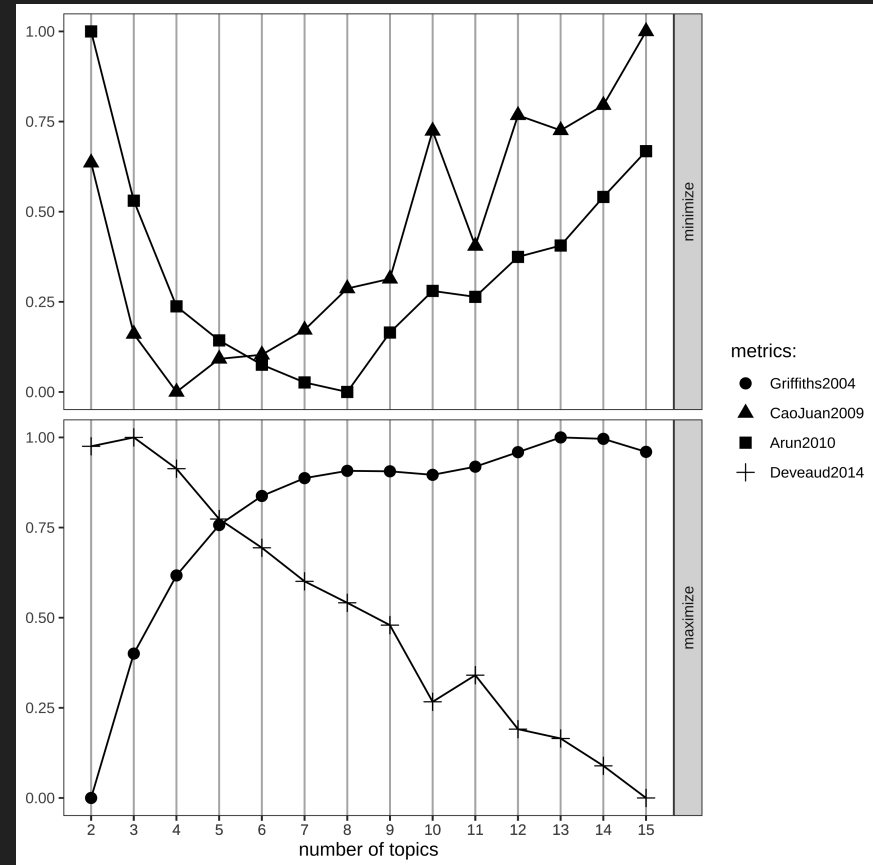
```
responses_topic_est <-  
  FindTopicsNumber(  
    responses_dtm,  
    topics = seq(from = 2, to = 20, by = 1), # amend these  
    metrics = c("Griffiths2004",  
               "CaoJuan2009",  
               "Arun2010",  
               "Deveaud2014"),  
    method = "Gibbs",  
    control = list(seed = 77),  
    mc.cores = 2L,  
    verbose = TRUE  
  ) %>%  
  FindTopicsNumber_plot()
```



The estimate for a total number of topics can be a lowest single value or range of values. We do this by observing where the metric curves tend to plateau and get as close to each other as possible along the horizontal axis. This is known as a limit

From the plot, the metric symbolized by + is diverging away from the rest. While they may head back towards the horizontal axis in the future, the metrics symbolized by \square , \circ , and \triangle look to be the closest between 5 and 6

We want to use the lowest possible count so let's start by modeling 5 topics!



```
LDA(responses_dtm,  
     k = 5, # Number of topics  
     control = list(seed = 1234)) %>%  
tidy(matrix = "beta")
```

```
## # A tibble: 2,895 × 3  
##   topic term      beta  
##   <int> <chr>    <dbl>  
## 1     1   time  0.0189  
## 2     2   time  0.0241  
## 3     3   time  0.0401  
## 4     4   time  0.0237  
## 5     5   time  0.0361  
## 6     1 improve 0.0292  
## 7     2 improve 0.00549  
## 8     3 improve 0.0190  
## 9     4 improve 0.0323  
## 10    5 improve 0.0495  
## # ... with 2,885 more rows
```



Assigning a Variable



Let's save the topics list

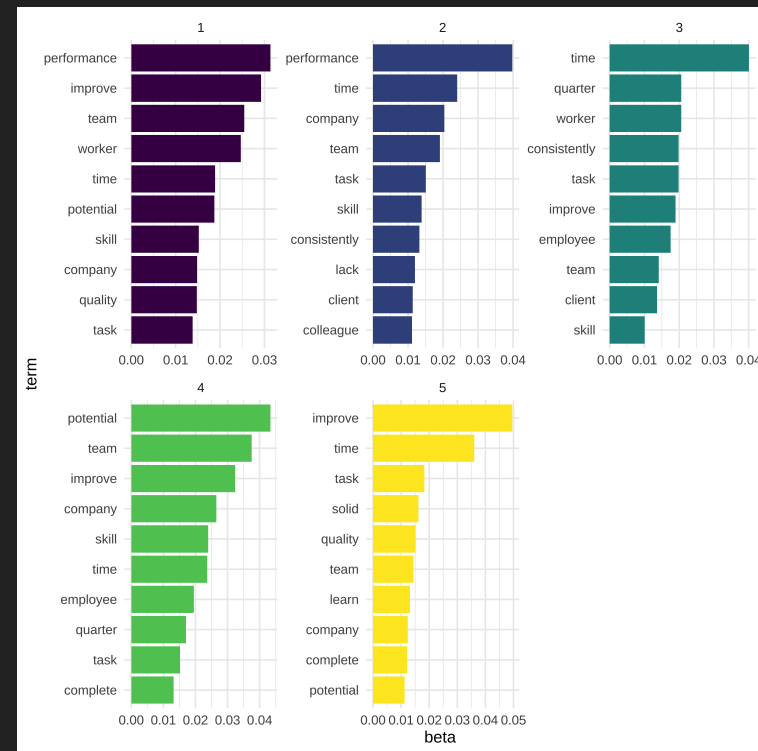
```
responses_topics <-  
  LDA(responses_dtm,  
      k = 5, # Amend this to test a certain number of topics  
      control = list(seed = 1234)) %>%  
  tidy(matrix = "beta")
```


Step 4: Visualization and Interpretation

```

responses_topics %>%
group_by(topic) %>%
slice_max(beta, n = 10) %>%
ungroup() %>%
arrange(topic, -beta) %>%
mutate(term = reorder_within(term, beta, topic)) %>%
ggplot(aes(beta, term, fill = factor(topic))) +
geom_col(show.legend = FALSE) +
scale_fill_viridis_d() +
facet_wrap(~ topic, scales = "free") +
scale_y_reordered() +
theme_minimal()

```

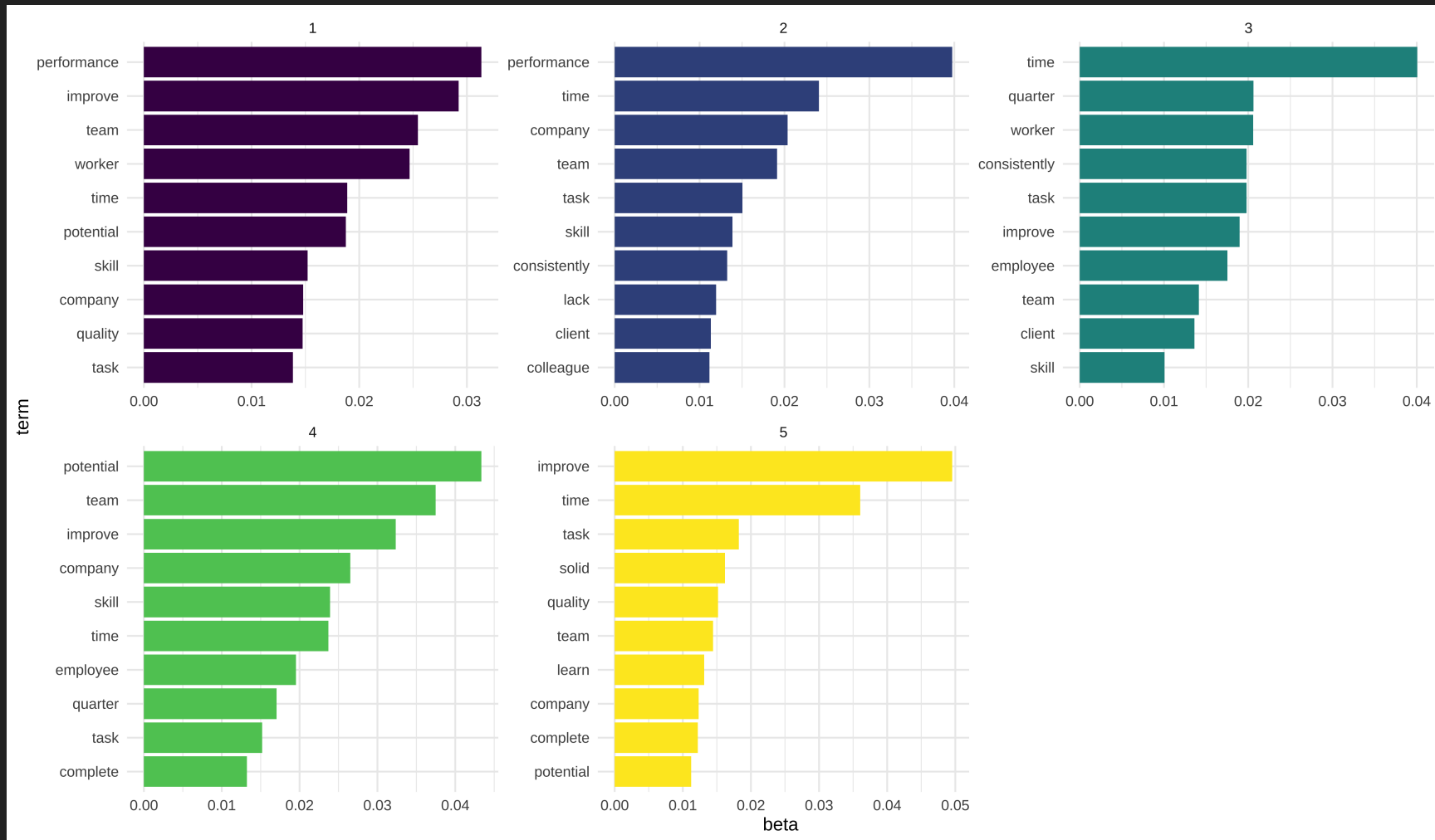


Assigning a Variable



Let's save the plot

```
responses_top_terms <-
  responses_topics %>%
  group_by(topic) %>%
  slice_max(beta, n = 10) %>%
  ungroup() %>%
  arrange(topic, -beta) %>%
  mutate(term = reorder_within(term, beta, topic)) %>%
  ggplot(aes(beta, term, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  scale_fill_viridis_d() +
  facet_wrap(~ topic, scales = "free") +
  scale_y_reordered() +
  theme_minimal()
```



Tip: you can save high (or really any) resolution visuals easily using `ggsave`

What Just Happened?



LDA is a form of (unsupervised) learning that views documents as bags-of-words (BoW) where order does not matter. Not having to track the placement of every term saves a lot of time and computational energy

LDA works by first making a key assumption: the way a document was generated was by picking a set of topics and then for each topic picking a set of words

Steps to Finding Topics



In a nutshell for each document m

1. Assume there are k topics across all of the documents
2. Create a distribution α where the k topics are symmetric or asymmetrically spread across each document m by assigning each word a topic
3. For each word w in every document m , assume its topic is associated incorrectly but every other word is assigned the correct topic
4. Probabilistically assign word w a topic based on two things:
 - what topics are in document m
 - Create a distribution β to assess how many times word w has been assigned a particular topic across all of the documents
5. Repeat this process a number of times for each document until saturation

Interpret



Much like you would assess a factor or component, the topics are unlabeled and it is up to you to figure out what they could mean. Not every topic may be directly applicable, but should still be interpreted and reported. Discarding topics means that you are removing potentially relevant information

Here is a brief assessment of some possible topics that are represented in the topic model with reference to *employee responses*



Topic	Label
1	production and gains reliance on employee abilities
2	teams' abilities to tackle client needs affect on-time completion
3	achievement varies by timeframe and workers' talent
4	possible growth is helped or hindered by worker characteristics
5	increases tied to employees' capacity and capabilities

Your assessment would likely differ to varying degrees and that is the point - in that qualitative concepts such as triangulation and saturation still play a large and impactful role in the interpretation phase. Note with a much larger text data set, this task could be significantly easier

That's It!

Any questions?



This work is licensed under a
Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License