



TEXT MINING

EXPLORATORY DATA ANALYSIS
TO MACHINE LEARNING



TIDY TEXT

HELLO

I'm Julia Silge

Data Scientist at Stack Overflow

@juliasilge

<https://juliasilge.com/>

TIDY TEXT



TEXT DATA IS INCREASINGLY IMPORTANT

TIDY TEXT



TEXT DATA IS INCREASINGLY IMPORTANT



NLP TRAINING IS SCARCE ON THE GROUND

TIDY TEXT

TIDY DATA PRINCIPLES
+
COUNT-BASED METHODS

=



tidytext: Text mining using dplyr, ggplot2, and other tidy tools

Authors: [Julia Silge](#), [David Robinson](#)

License: [MIT](#)

[build](#) passing [build](#) passing CRAN 0.2.2 coverage 81% DOI [10.5281/zenodo.3355453](#)
JOSS [10.21105/joss.00037](#) downloads 27K/month downloads 683K



Using [tidy data principles](#) can make many text mining tasks easier, more effective, and consistent with tools already in wide use. Much of the infrastructure needed for text mining with tidy data frames already exists in packages like [dplyr](#), [broom](#), [tidyr](#) and [ggplot2](#). In this package, we provide functions and supporting data sets to allow conversion of text to and from tidy formats, and to switch seamlessly between tidy tools and existing text mining packages. Check out [our book](#) to learn more about text mining using tidy data principles.

<https://github.com/juliasilge/tidytext>

tidytext: Text mining using dplyr, ggplot2, and other tidy tools

Authors: [Julia Silge](#), [David Robinson](#)

License: [MIT](#)

[build](#) passing [build](#) passing CRAN 0.2.2 coverage 81% DOI [10.5281/zenodo.3355453](#)
JOSS [10.21105/joss.00037](#) downloads 27K/month downloads 683K



Using [tidy data principles](#) can make many text mining tasks easier, more effective, and consistent with tools already in wide use. Much of the infrastructure needed for text mining with tidy data frames already exists in packages like [dplyr](#), [broom](#), [tidyr](#) and [ggplot2](#). In this package, we provide functions and supporting data sets to allow conversion of text to and from tidy formats, and to switch seamlessly between tidy tools and existing text mining packages. Check out [our book](#) to learn more about text mining using tidy data principles.



<https://github.com/juliasilge/tidytext>

O'REILLY®

Text Mining with R

A TIDY APPROACH



Julia Silge & David Robinson

<https://tidytextmining.com/>

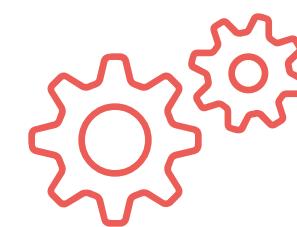
TIDY TEXT



EXPLORATORY DATA ANALYSIS



N-GRAMS AND MORE WORDS



MACHINE LEARNING

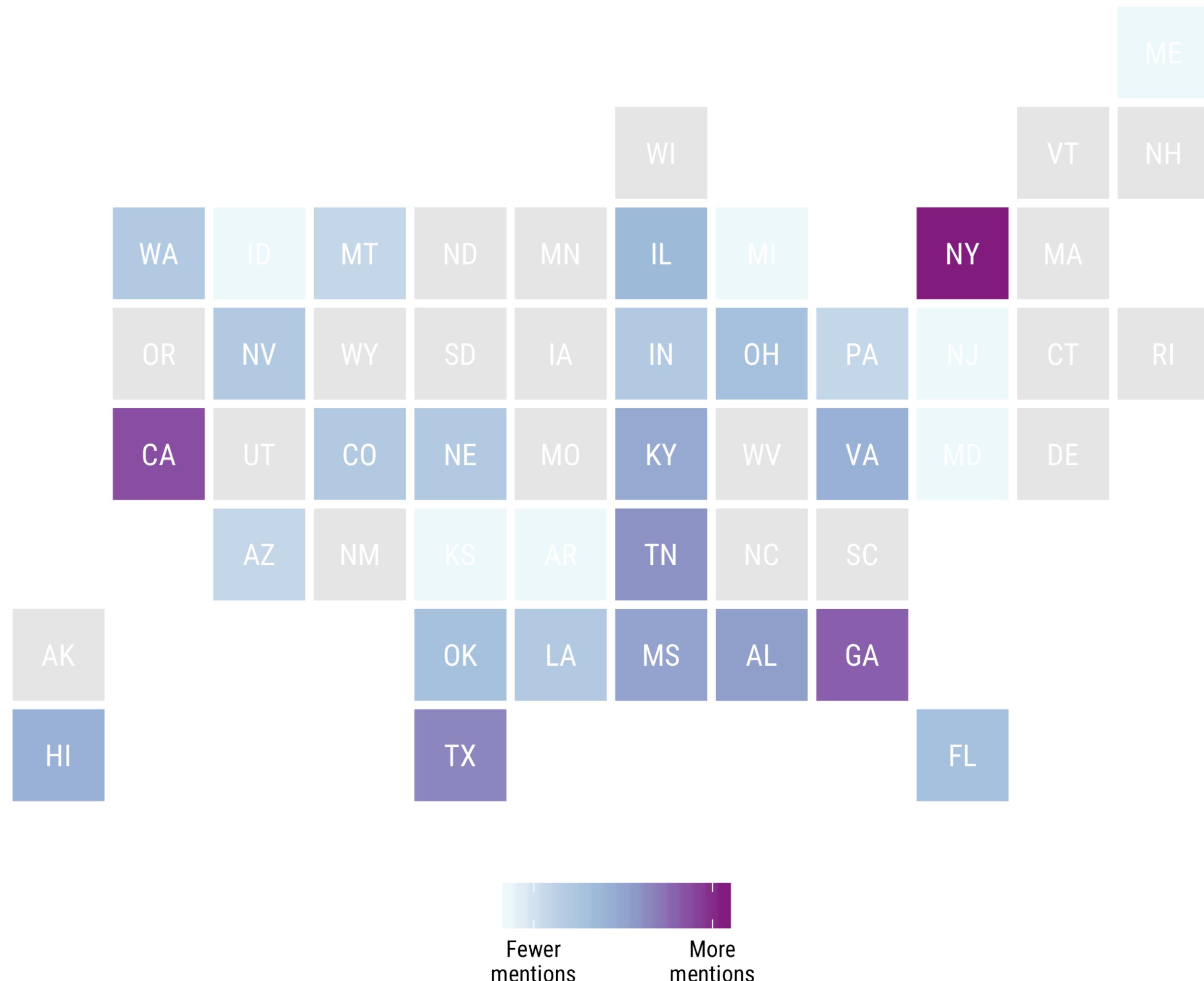
TIDY TEXT

EXPLORATORY DATA ANALYSIS



Which States Are Mentioned Most in Song Lyrics?

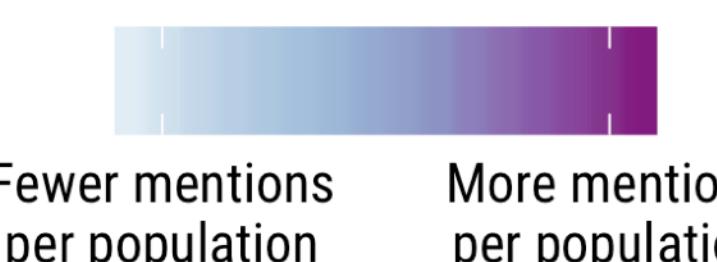
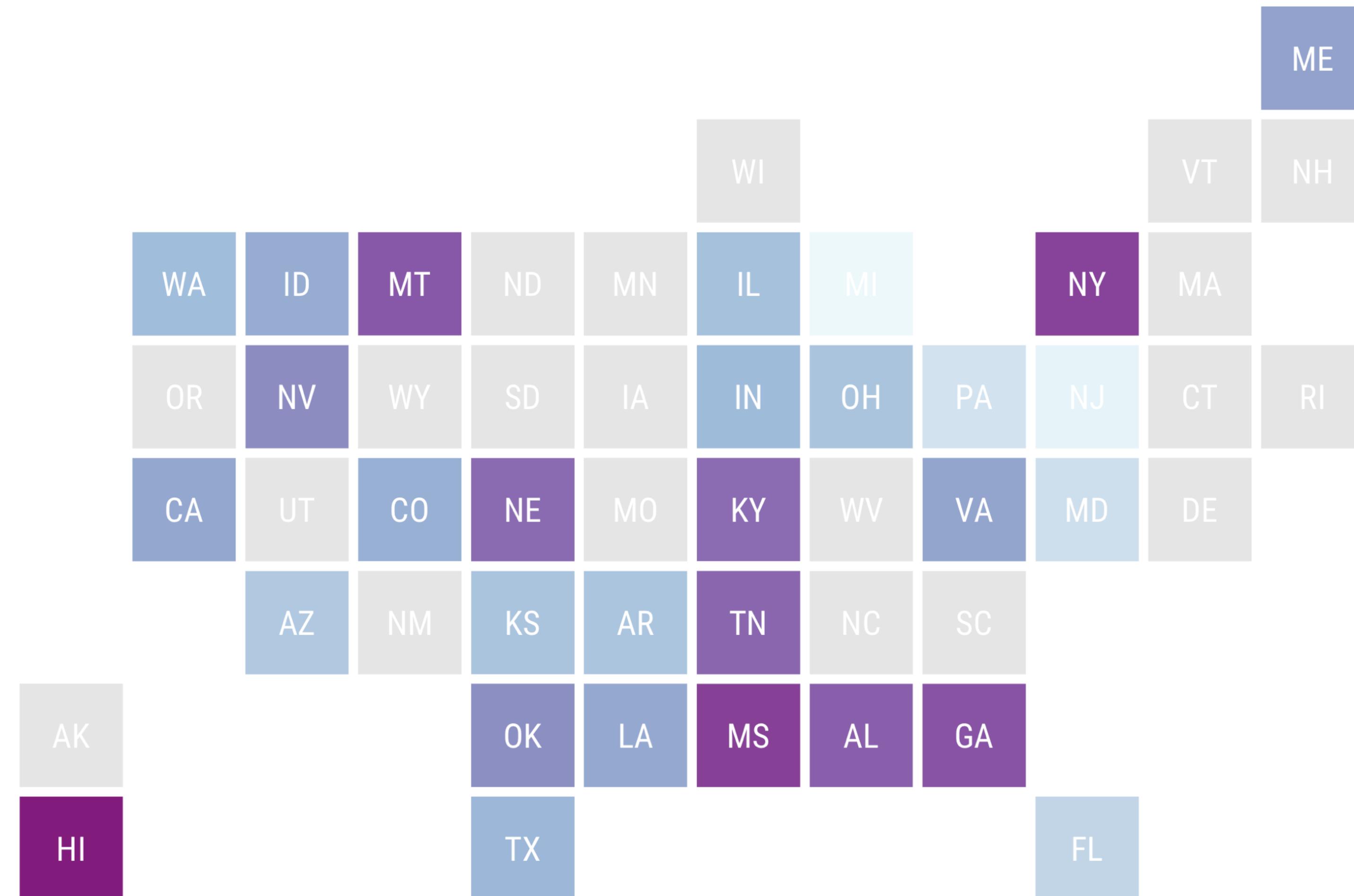
States like California are mentioned most often



from the Washington Post's Wonkblog

Which States Are Mentioned Most in Song Lyrics?

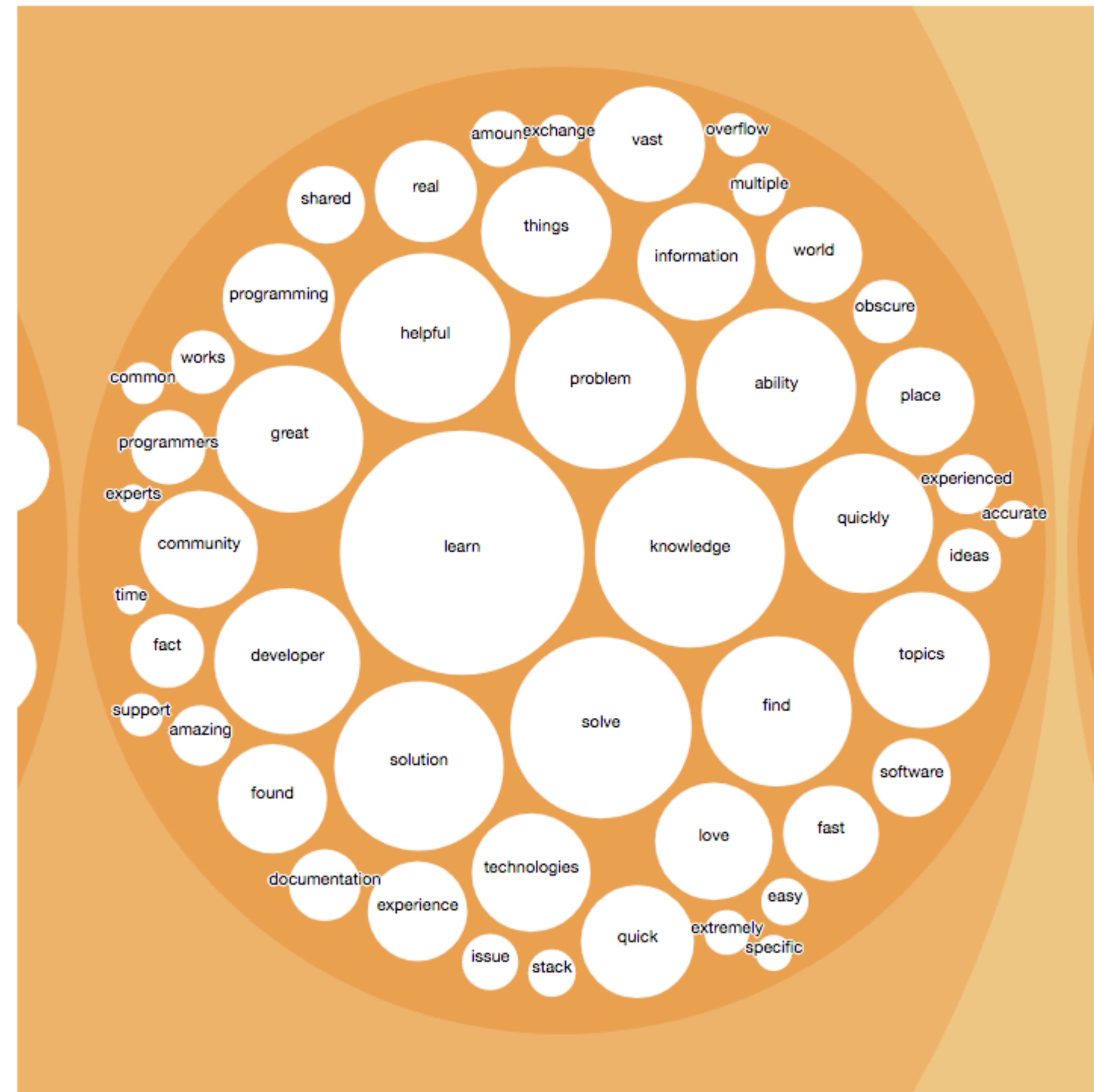
States like Hawaii and Montana are mentioned more often relative to their population



What do users say about Stack Overflow?

On the [2018 Stack Overflow Developer Survey](#), we asked users what the best, worst, most exciting, and most annoying things about Stack Overflow are.

Here are their answers.



D3 visualization [on Glitch](#)

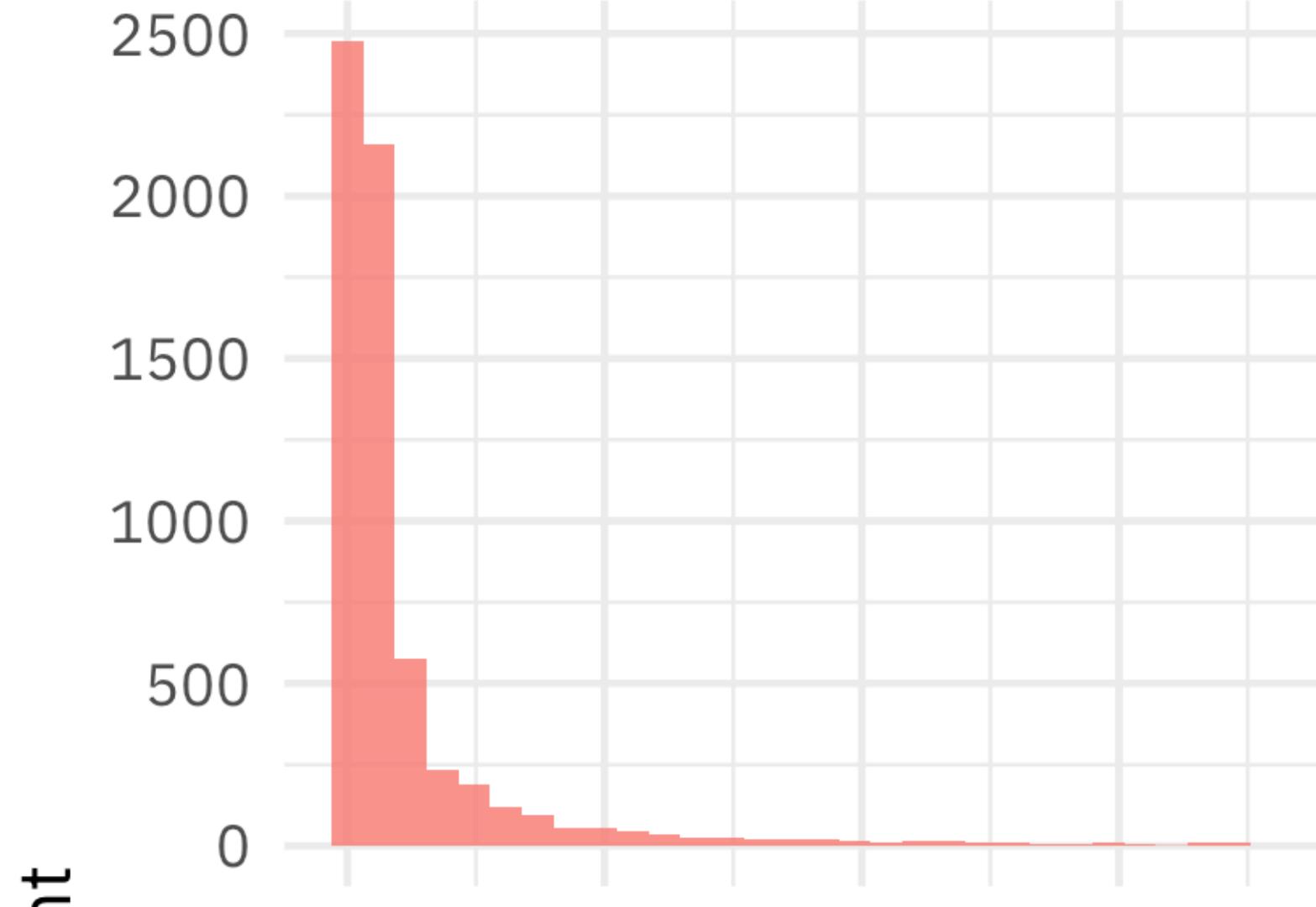
WHAT IS A DOCUMENT ABOUT?

TERM FREQUENCY
INVERSE DOCUMENT FREQUENCY

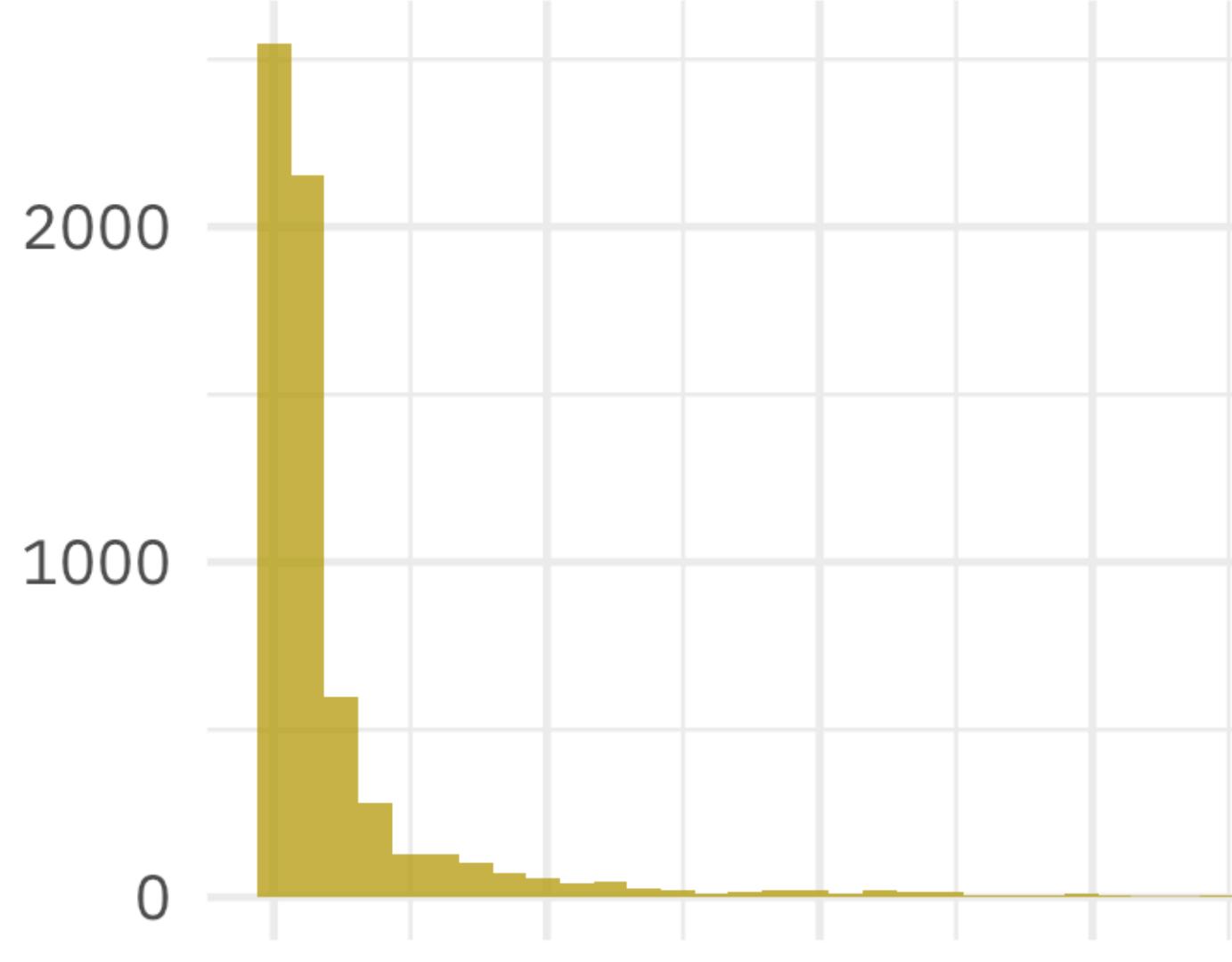
$$idf(\text{term}) = \ln \left(\frac{n_{\text{documents}}}{n_{\text{documents containing term}}} \right)$$

Term Frequency Distribution in Jane Austen's Novels

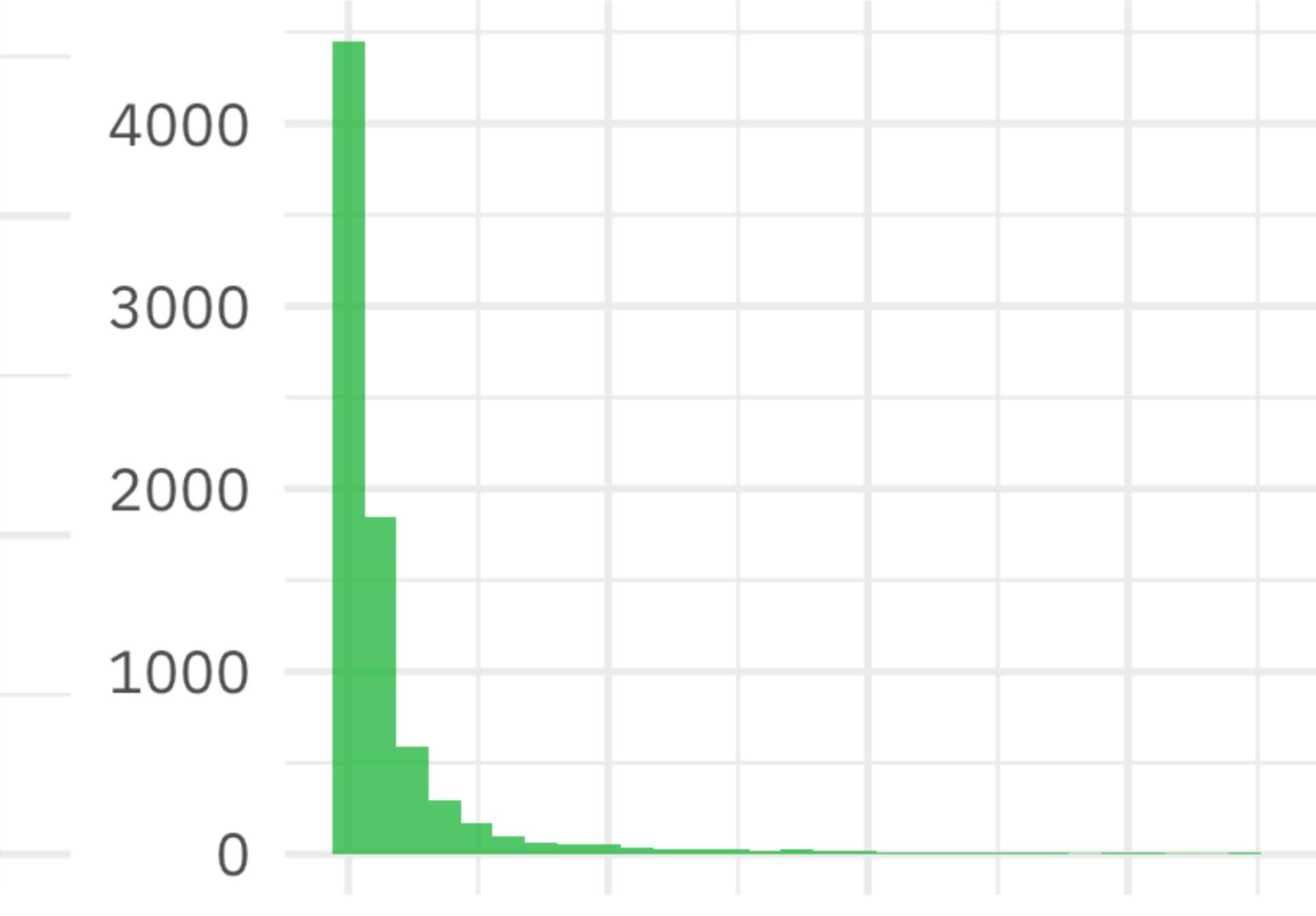
Sense & Sensibility



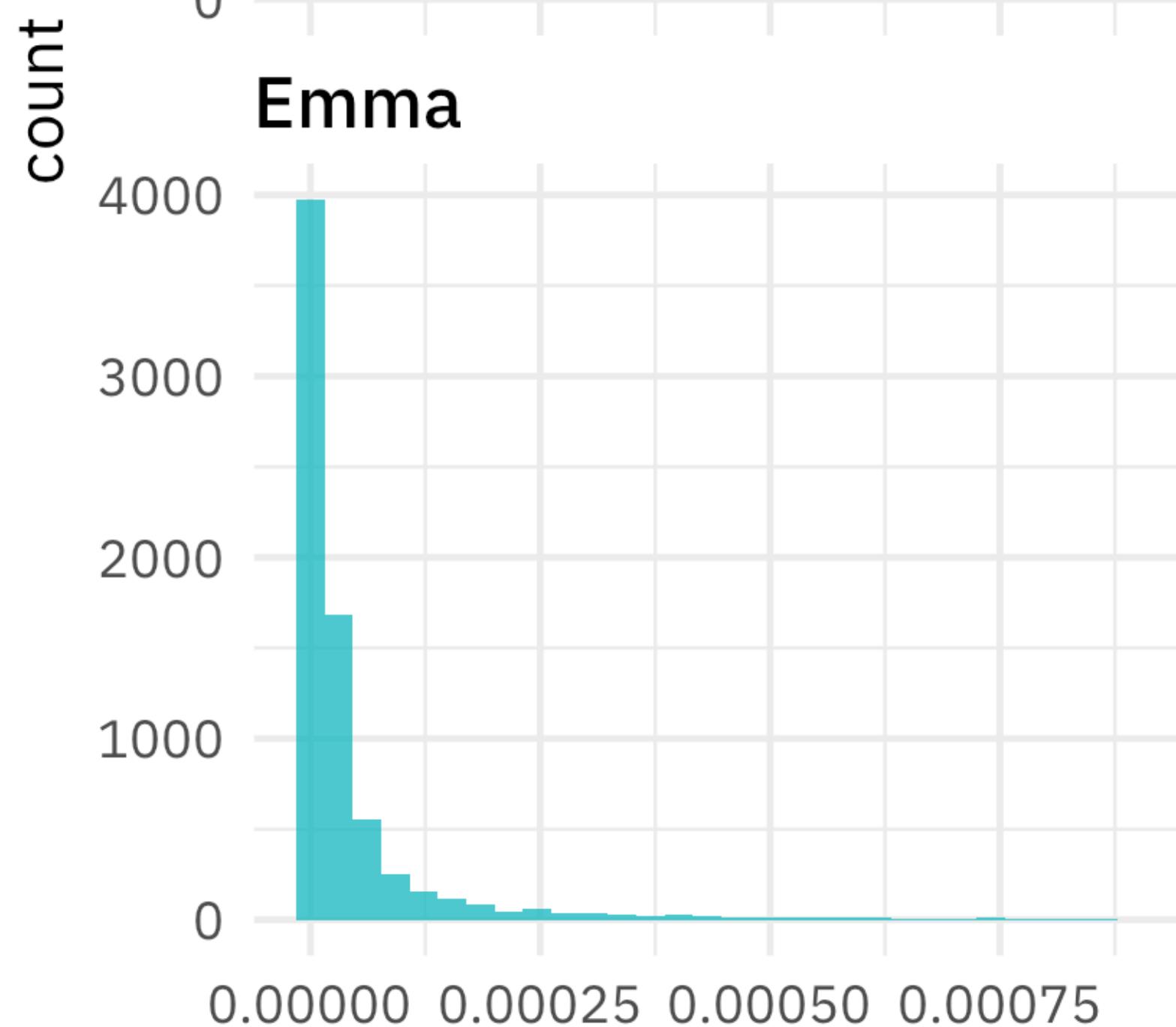
Pride & Prejudice



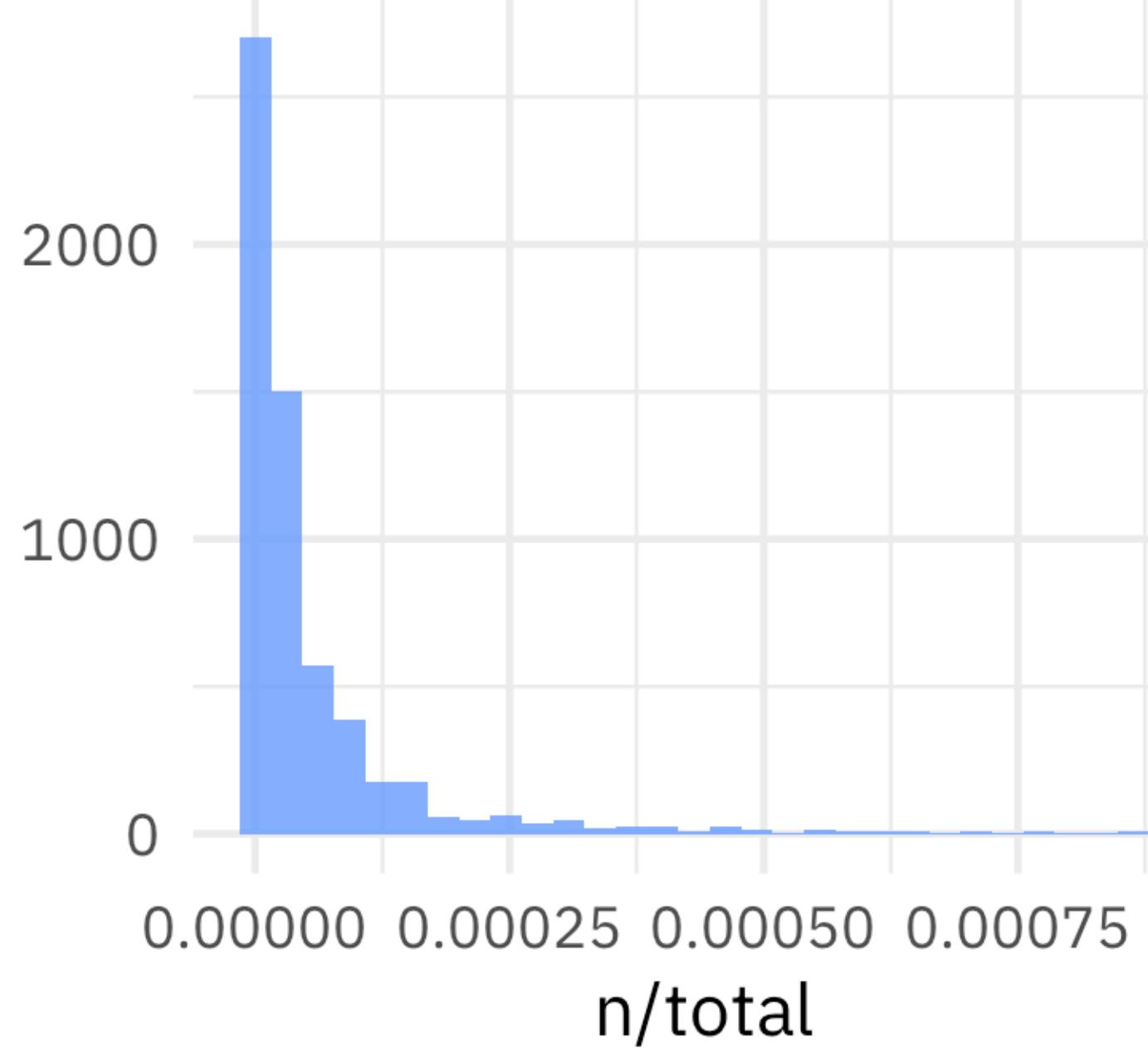
Mansfield Park



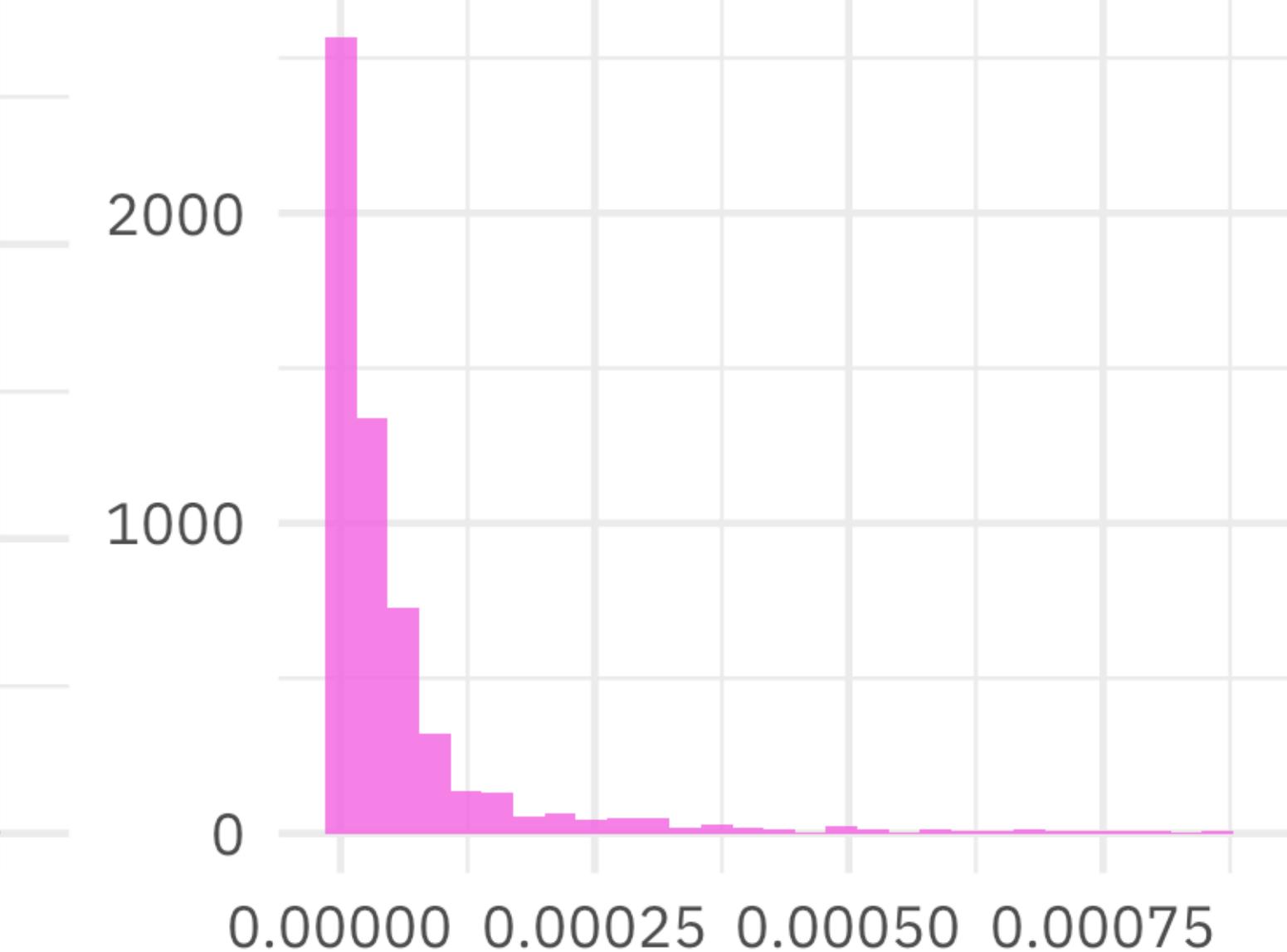
Emma



Northanger Abbey

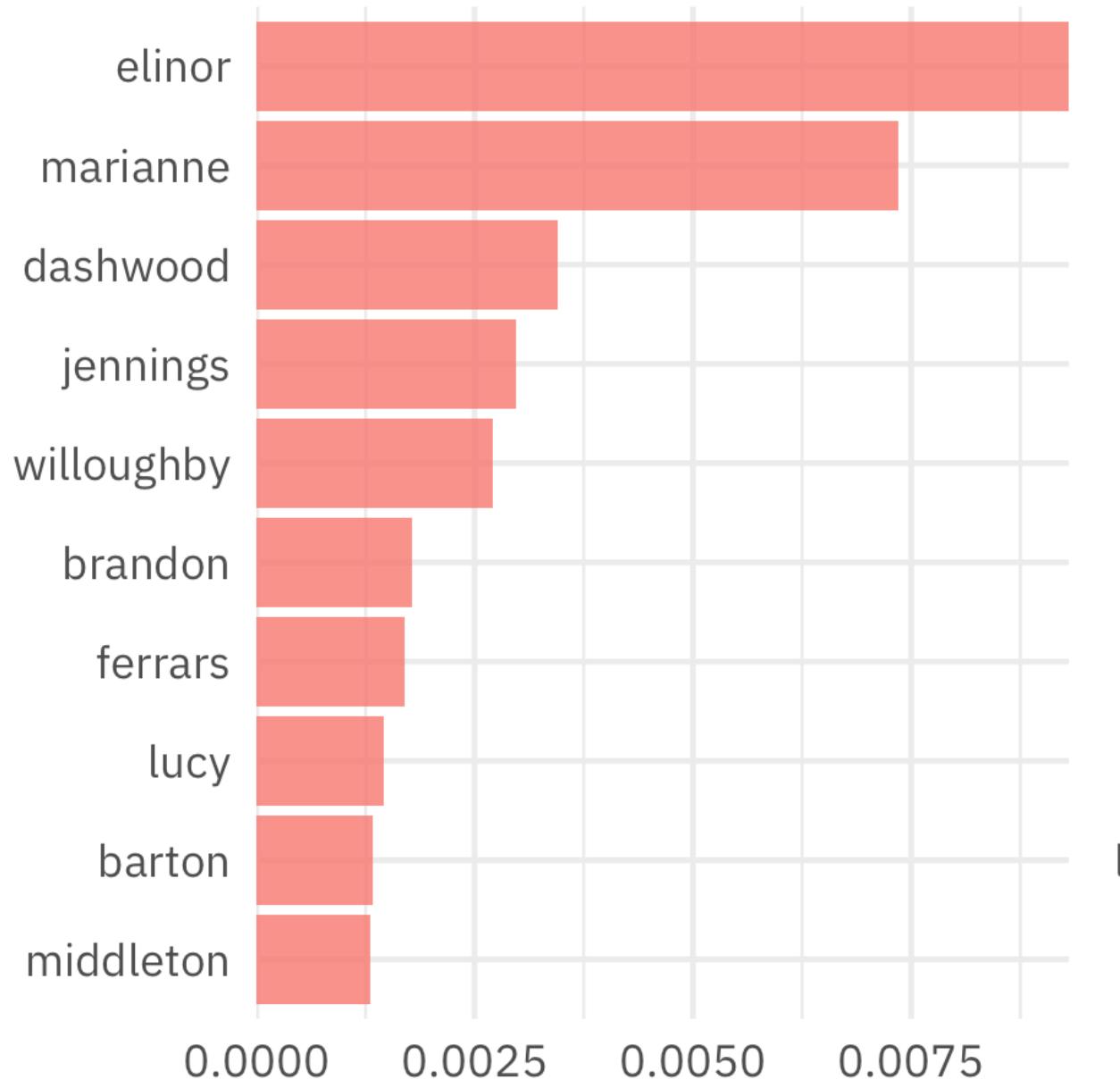


Persuasion

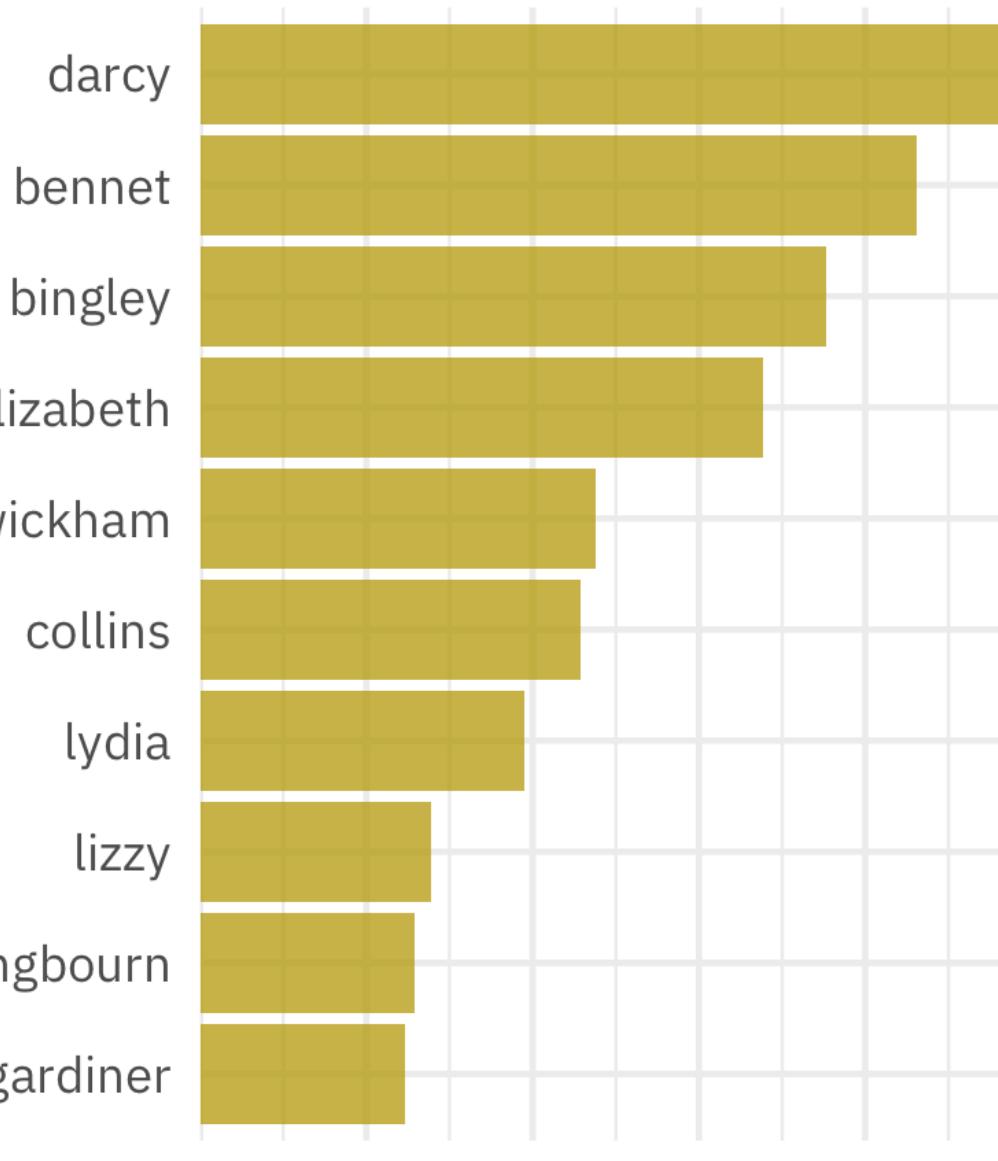


Highest tf-idf words in Jane Austen's Novels

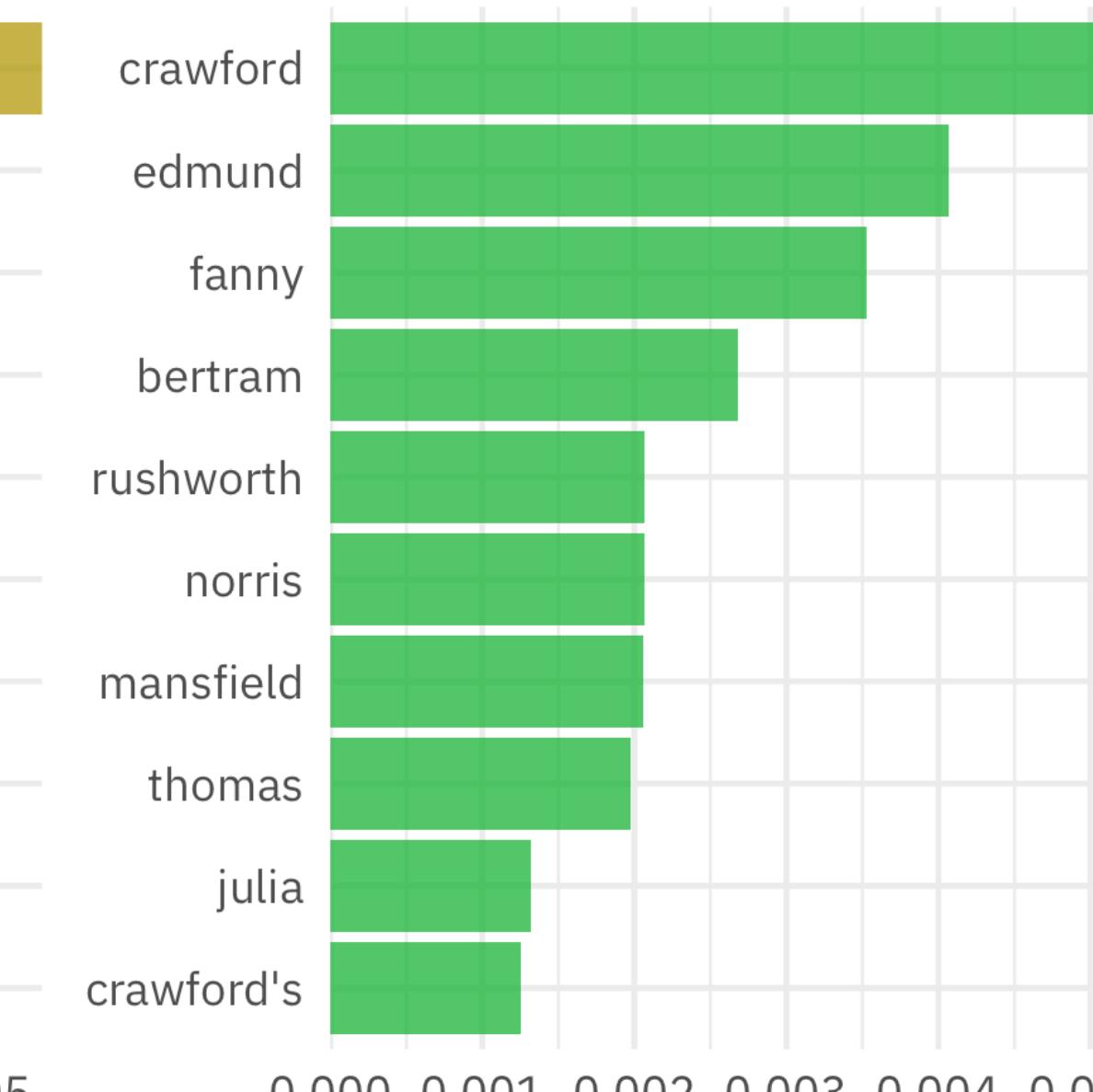
Sense & Sensibility



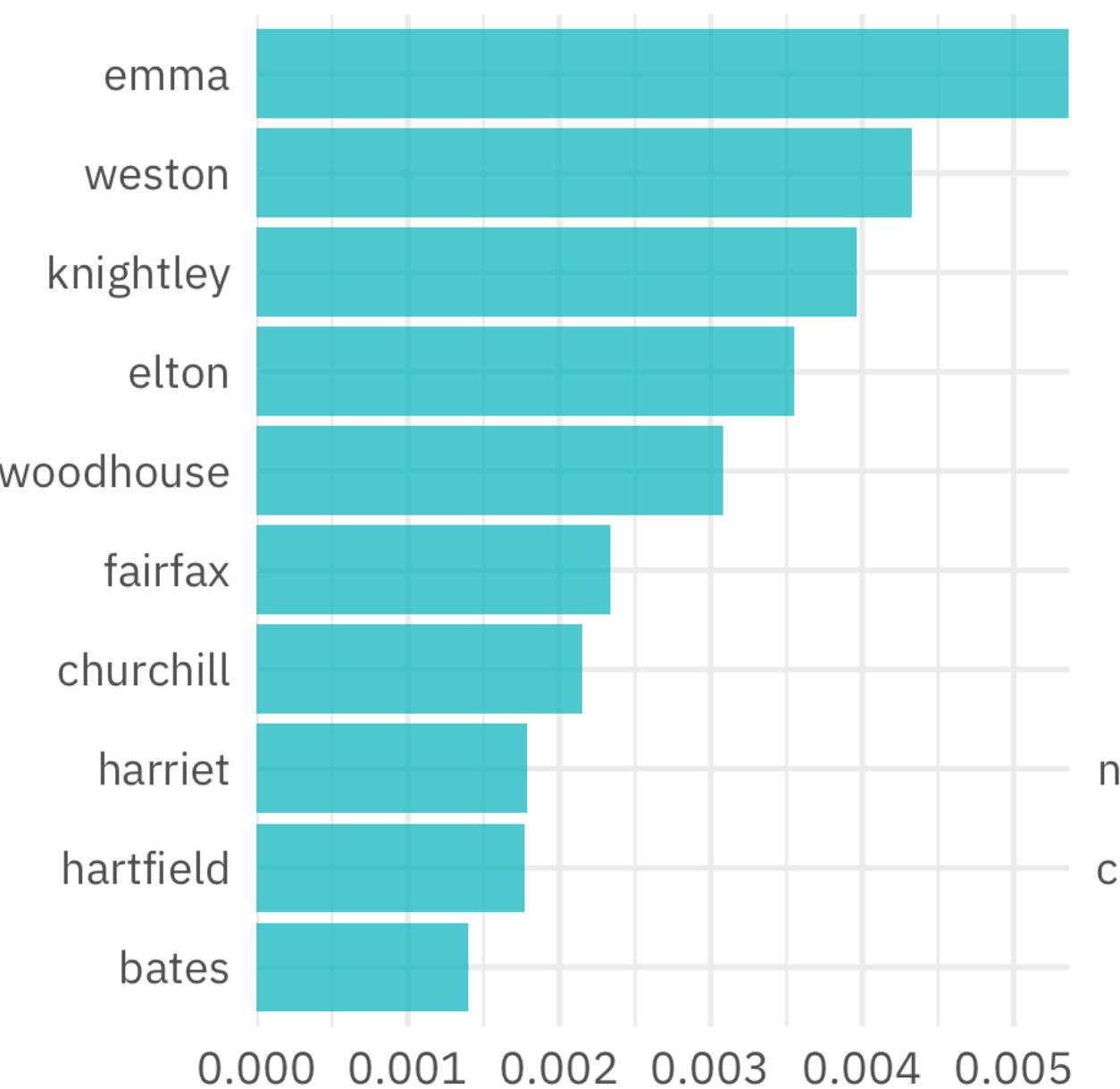
Pride & Prejudice



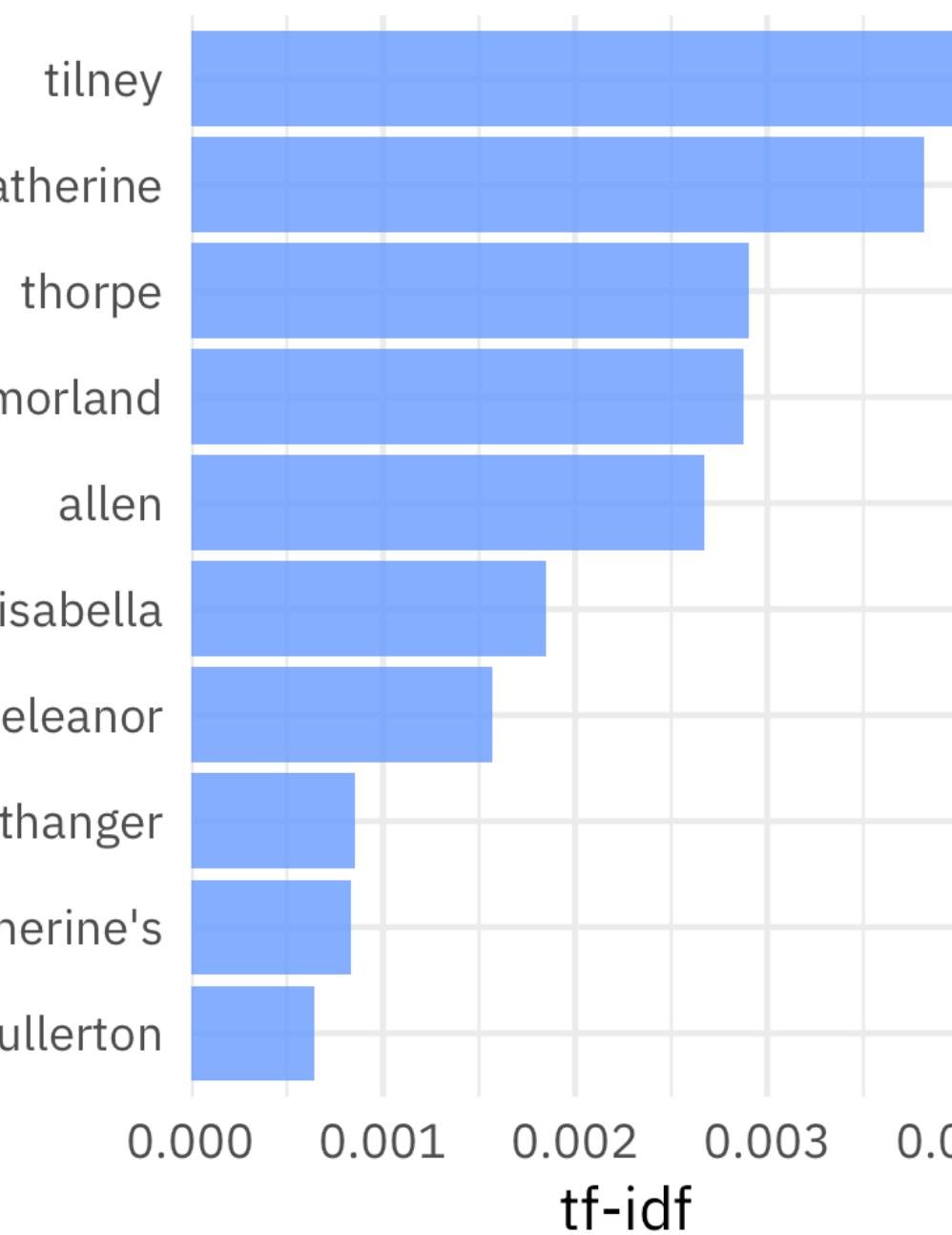
Mansfield Park



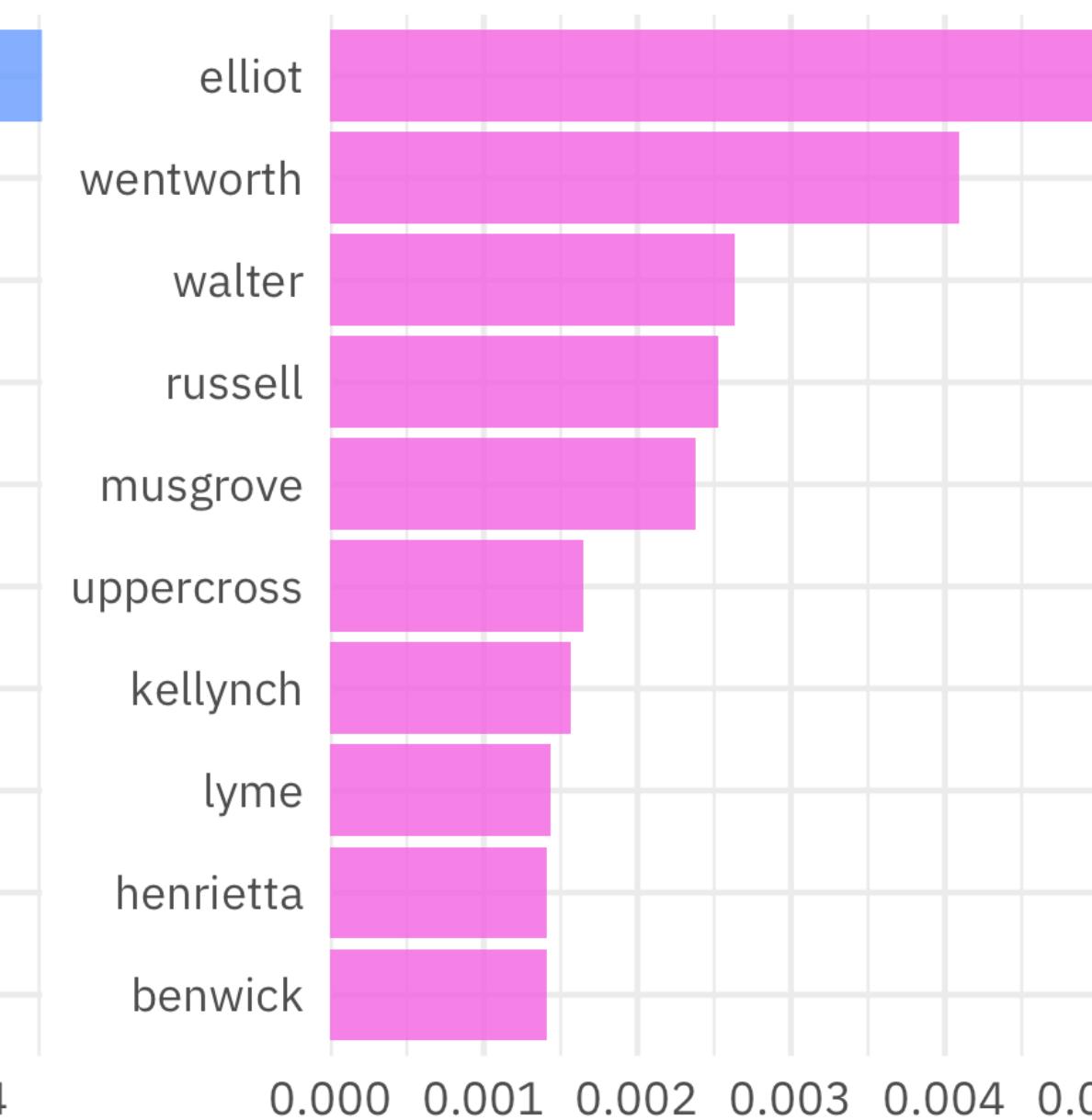
Emma



Northanger Abbey



Persuasion

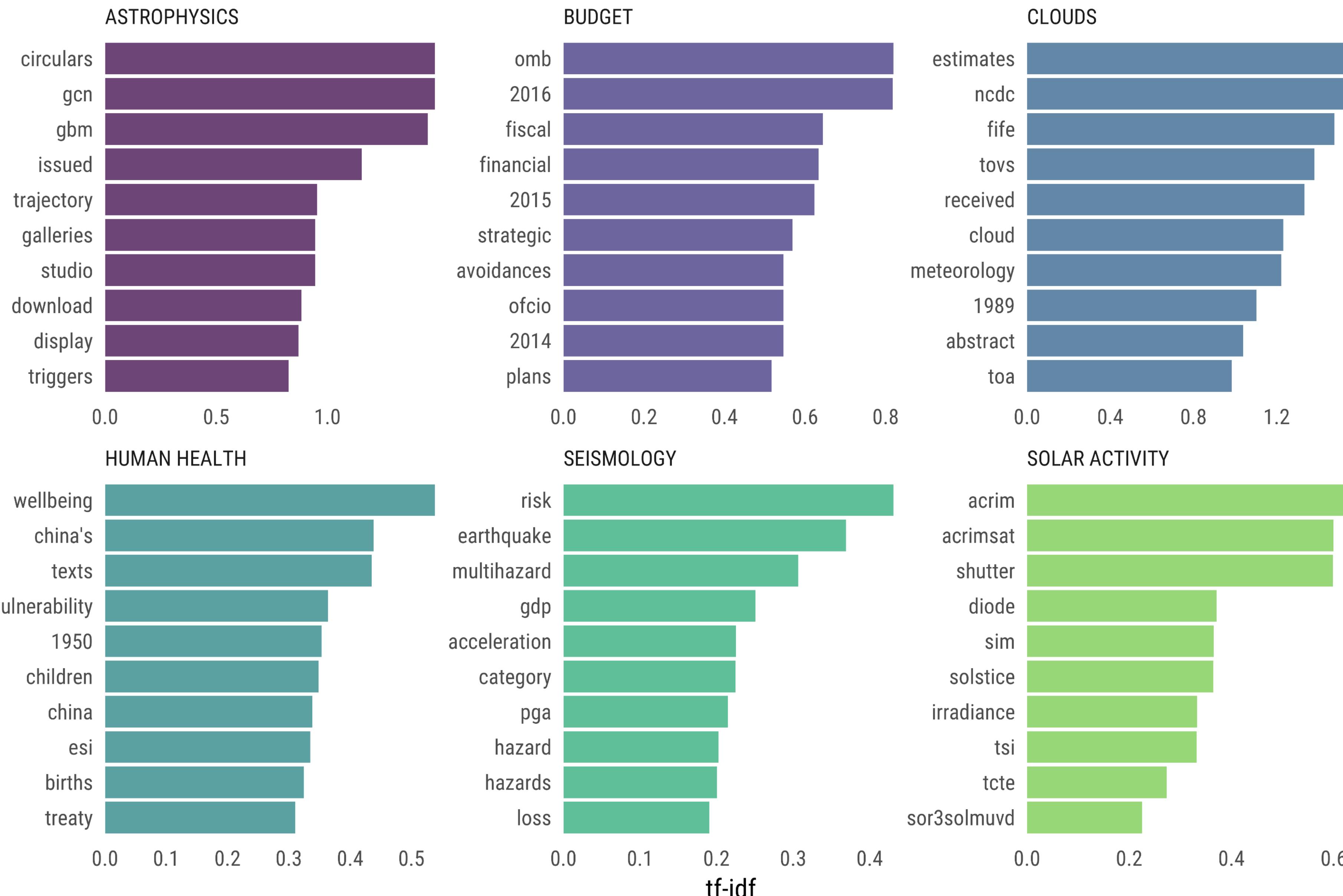


- As part of the NASA Daternauts program, I worked on a project to understand NASA datasets
- Metadata includes title, description, keywords, etc



Highest tf-idf words in NASA Metadata Description Fields

Distribution of tf-idf for words from datasets labeled with select keywords

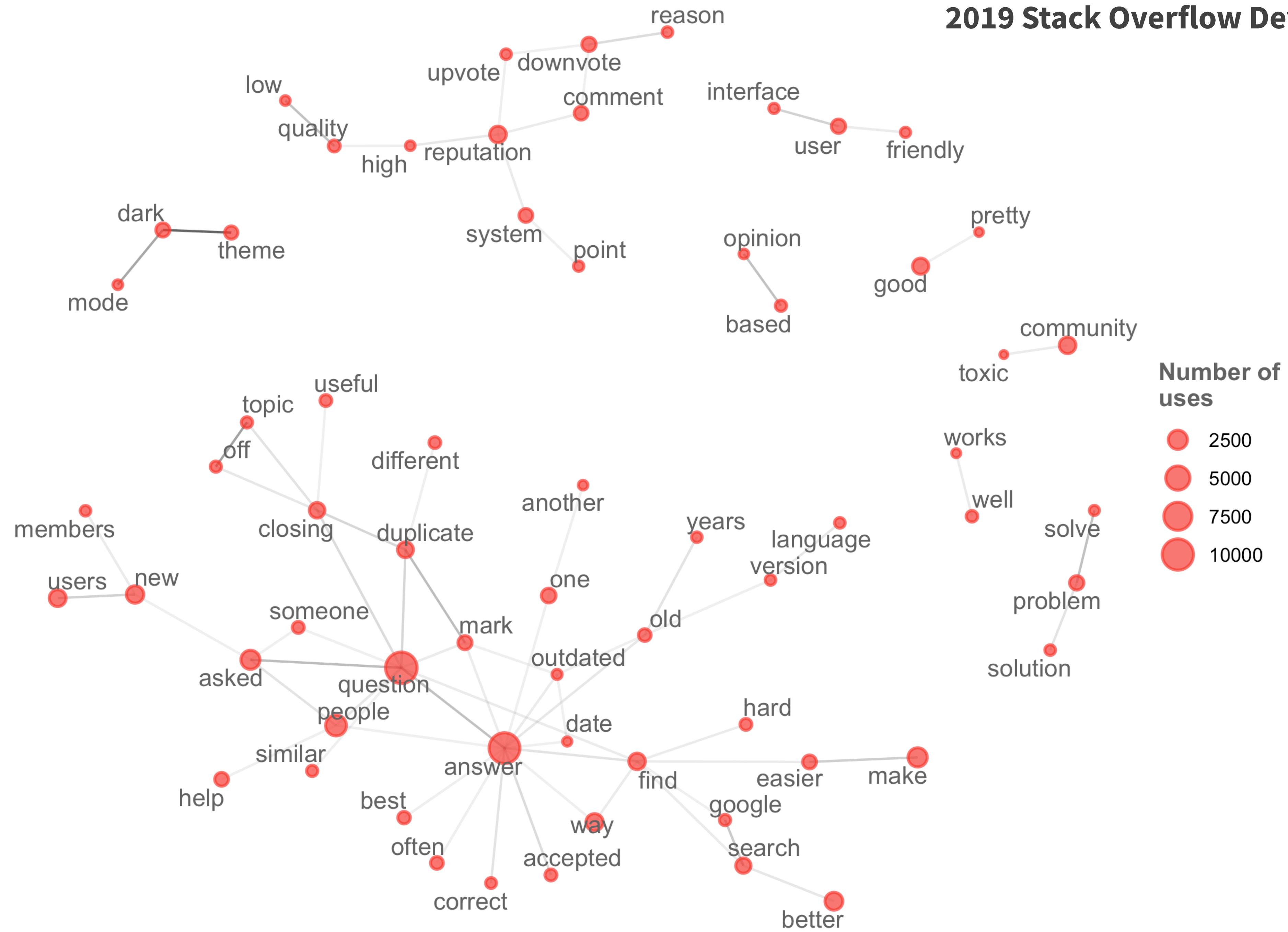




TAKING TIDY TEXT TO
THE NEXT LEVEL

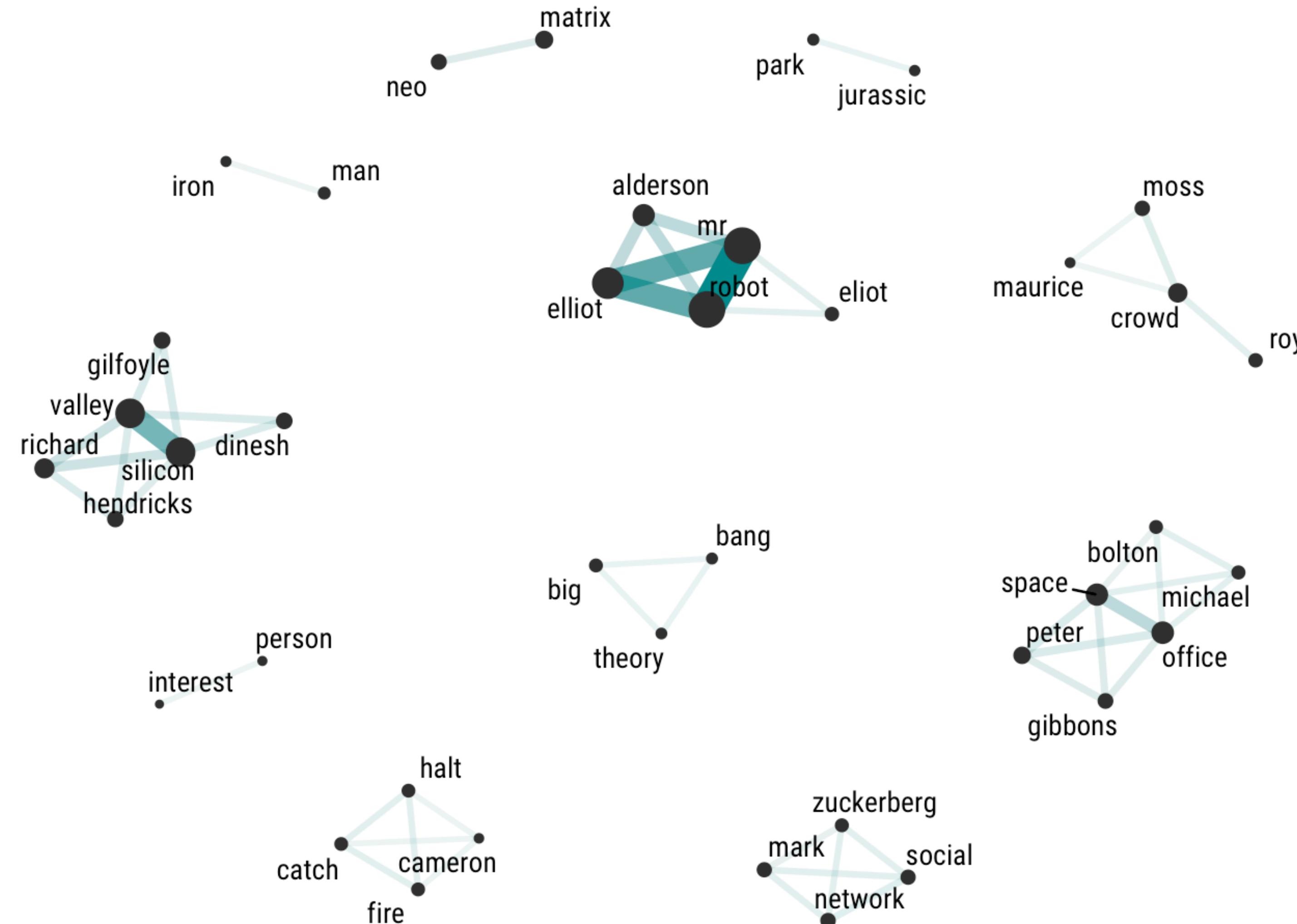
N-GRAMS,
NETWORKS, &
NEGATION

2019 Stack Overflow Developer Survey



Most Realistic Fictional Characters

As identified by respondents in the 2017 Stack Overflow Survey

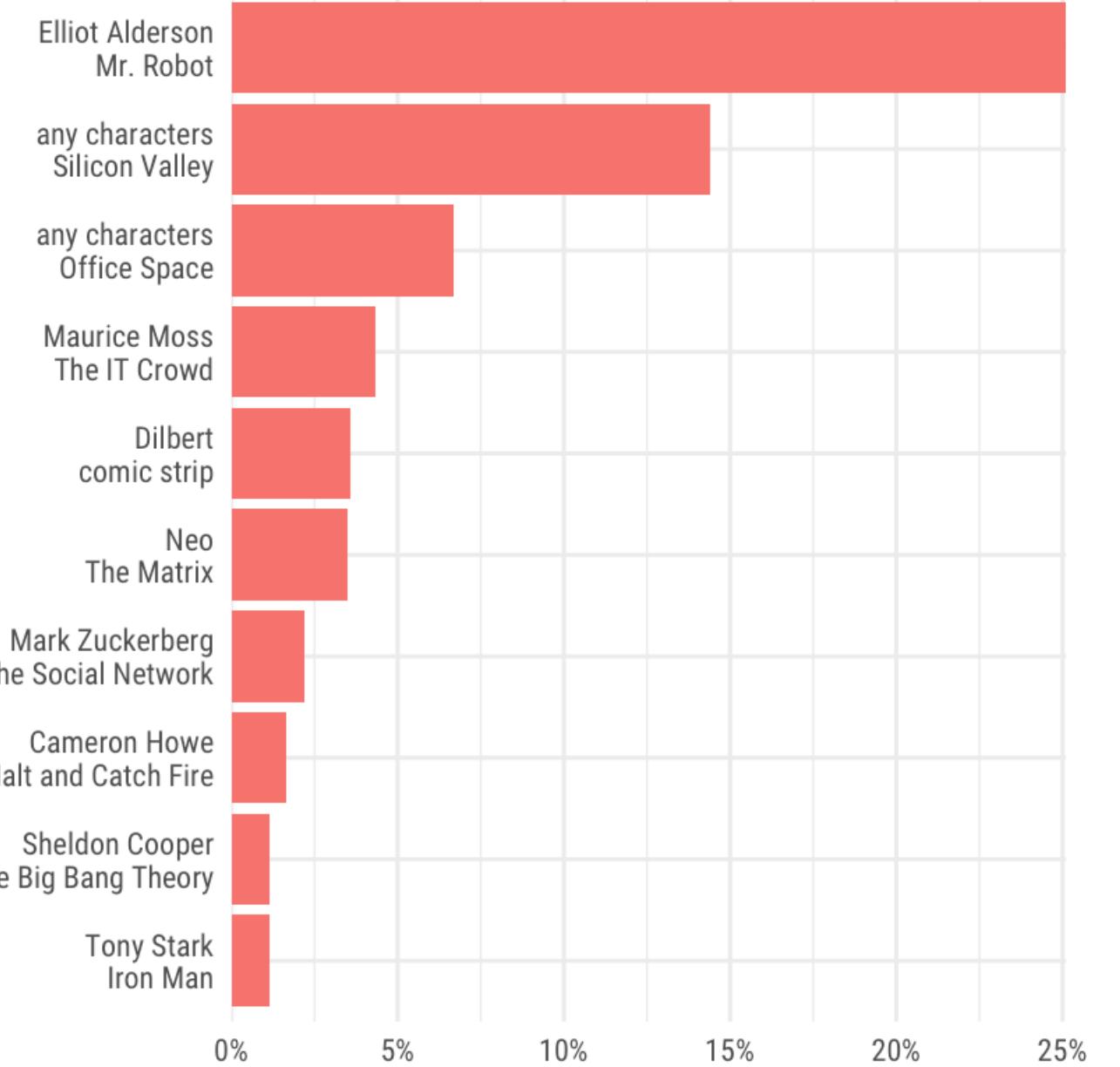


Words with larger points are used more often, and heavier connections indicate words are used more often together

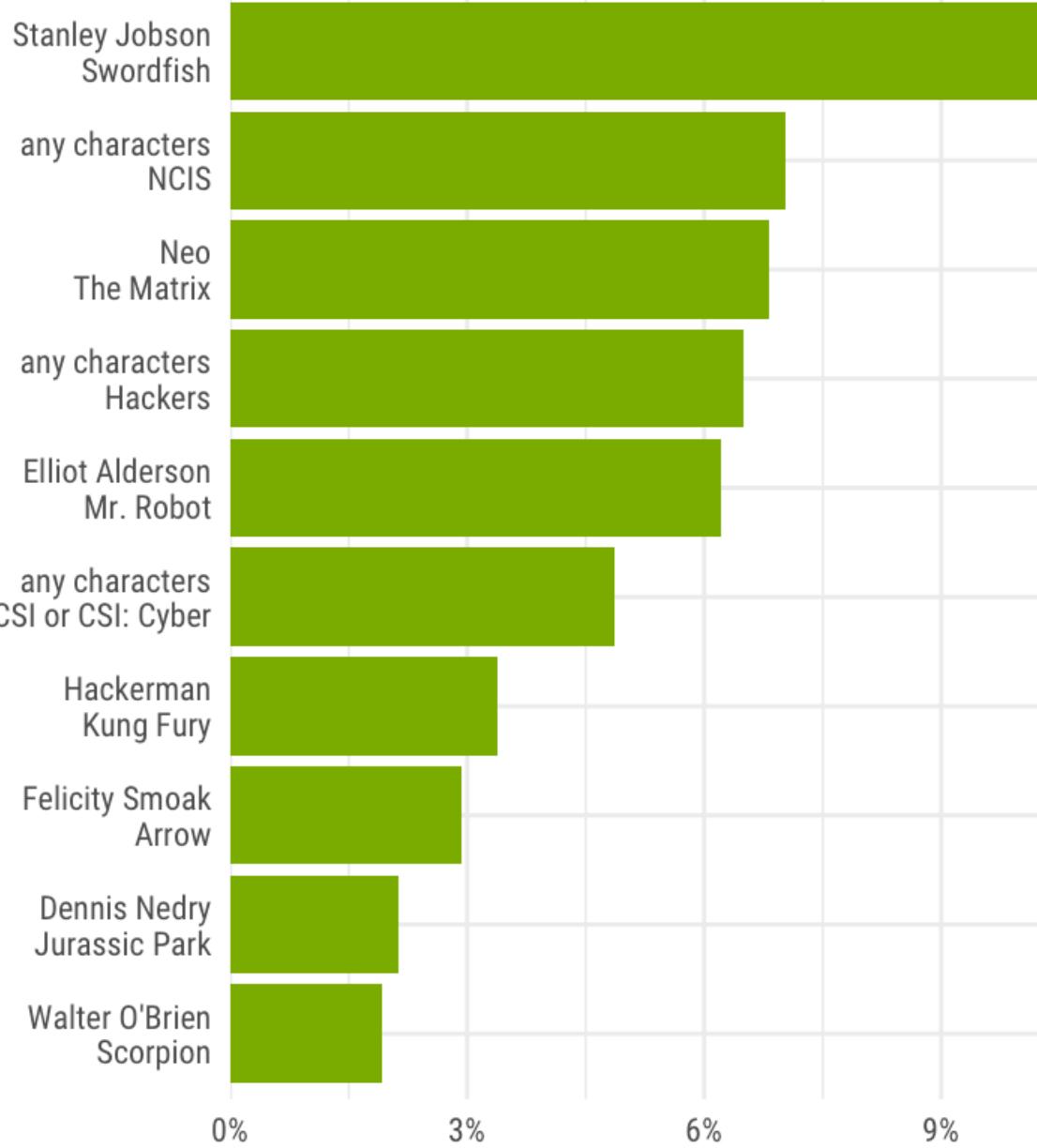
Hackers, Programmers, and IT Professionals in Fiction

From 10,983 responses on the 2017 Developer Survey

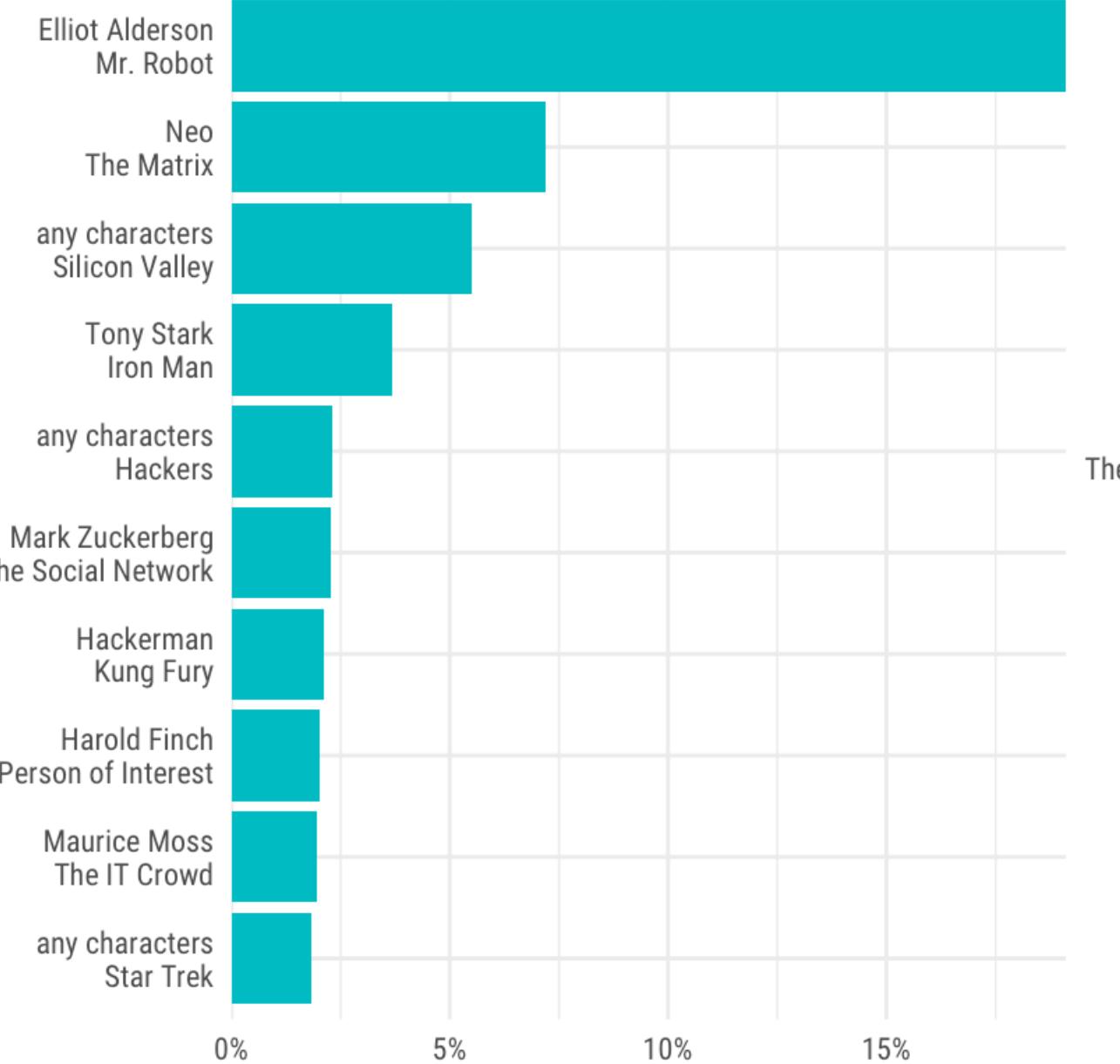
Most Realistic



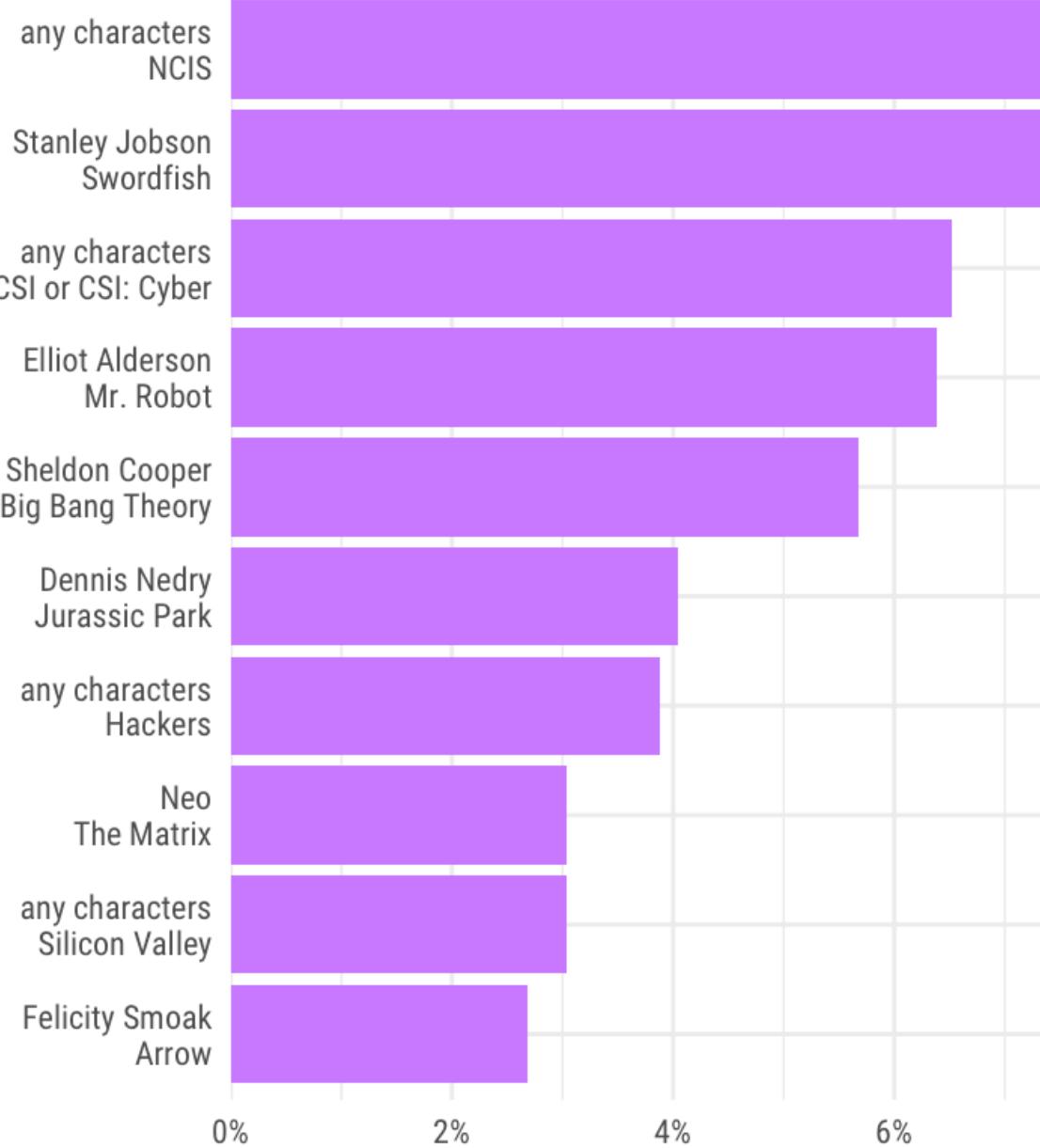
Least Realistic



Inspiring



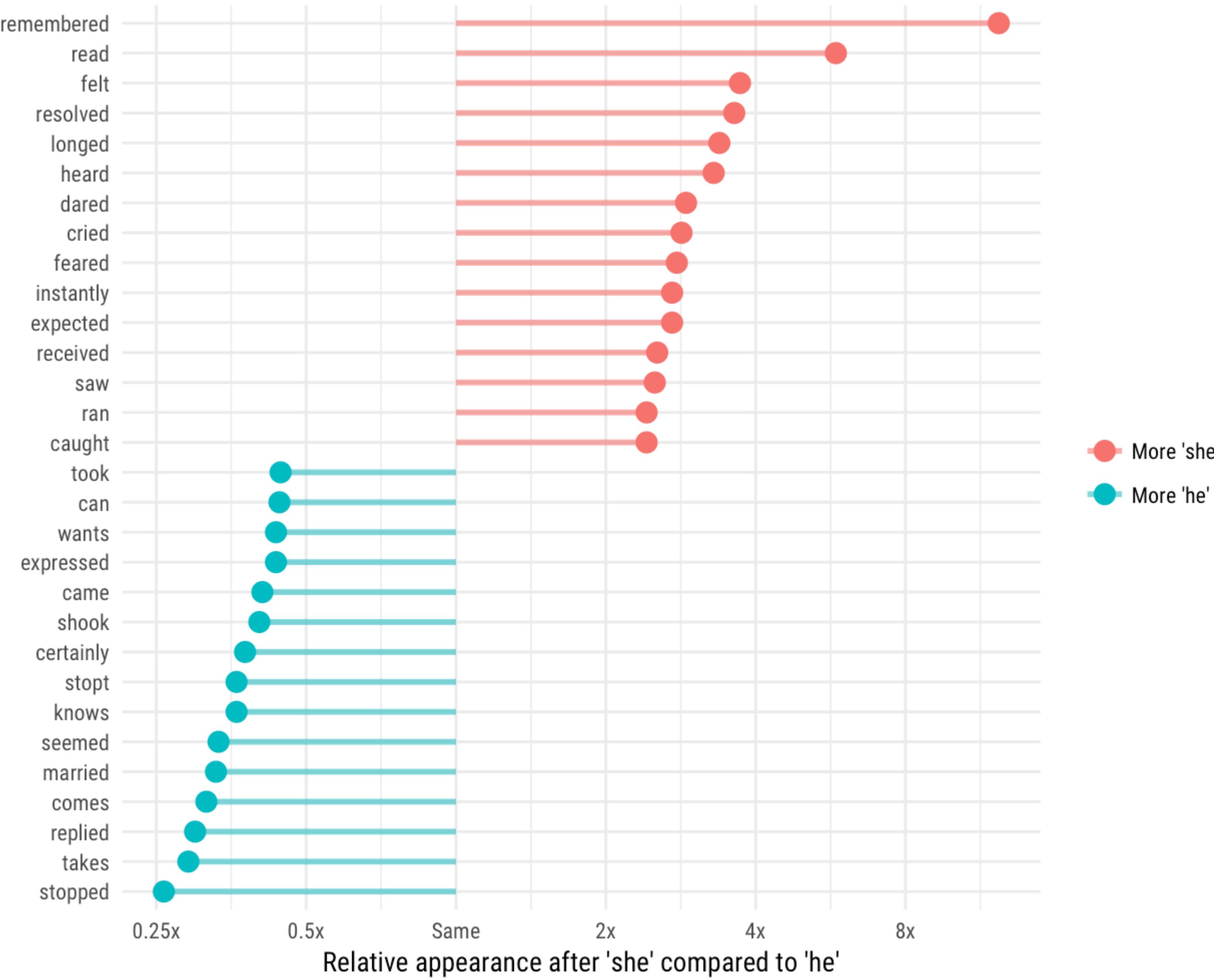
Annoying



% who responded with each character

Words paired with 'he' and 'she' in Jane Austen's novels

Women remember, read, and feel while men stop, take, and reply



She Giggles, He Gallops

Analyzing gender tropes in film with screen direction
from 2,000 scripts.

By Julia Silge

+

Russell Goldenberg Amber Thomas Hanah Anderson

The top 800 words paired with “she” or “he”

Underlined words contain examples of their usage in screen direction.

EVEN

MORE "SHE"

MORE "HE"

snuggles giggles squeals sobs weeps blushes clings rocks shrieks hugs shrinks gasps responds trembles pets flinches arches skips utters shudders startles buries swats murmurs resists
hovers caresses awakens shivers screams dances beats absently flees cleans stirs straddles cries moans bites realises mouths accepts wore smiles laughs wrote serves scoots liked
arranges scampers storms **twirls** softens ignores softly faints wonders fades sags hesitates casts applies hisses fiddles kisses sings awkwardly smokes stretches sips unbuttons stiffens
hurriedly hurries dries looks
rakes relents reluctantly fr
yelps ends allows flashes
thrashes becomes types a
returns forces closes fixes
floats left whacks blows i
squeezes begins kneels tu
She twirls, looking at herself in the mirror...
She twirls toward the door, grabbing her purse.
She twirls off. He chases her, beer and entries in hand.
She twirls a wooden katana as part of a floor routine.
She twirls a seductive finger around his tie. They kiss.
She twirls her finger around Milo's palm.
She twirls the string of her basketball shorts...
early refuses tiptoes lingers beams pivots curls glides strokes meant abruptly retrieves bursts
trails frowns retreats gonna licks touches reacts nearly sighs backs embraces squirms panics
shakes instinctively replies told freezes resumes creeps calms gives rushes sails tentatively
s blinks dabs meets rests regards tilts attacks darts eyes brushes descends gently nervously
s halts wakes bolts slaps fumbles quickly wears clasps faces feeds barely shrugs believes
s slices runs leans sounds washes swallows cranes observes accidentally marches rifles
grimaces frees put calls glares tucks like plops scurries whispers tries remains actually
continues pinches tells lets yawns disappears heads looks chokes discovers plugs springs watches cautiously opens clutches studies dropped massages obeys suddenly loves scowls
crosses packs scribbles spits sneaks puts likes just lashes topples hangs lies angles starts claps adds flicks slowly cups angrily really stops pours jabs traces unzips crawls died
grasps slugs steadies breathes glances pushes directs inches hands reads comes unfolds winds attempts rubs snorts walks drifts sinks hides goes swivels feels keeps sprays sways
means ducks races repeats used steps sends finishes talking trips locks waits gets snatches chooses obviously takes decides plucks slides moves peers holds claws already stomps
swipes asks admires met said stands buttons owns says sets strips must wanted focuses fell unwraps edges unlocks remembers shines reaches jumps always places climbs stumbles
handles leafs needed passed surreptitiously concentrates helps interrupts reels scrambles steals blocks clenches floors gags splashes rips notices knocks enters finally rises listens
quietly jerks bumps wants lays kicks flies makes picks throws casually scans winges might dangles hefts flails waves dresses realizes passes catches plants thinks knows rings empties
hears sweeps may signs wrenches swims shares started recovers hops props tightens indicates hear dashes got peels will silently probably grows whips finds spins pulls switches
lights assumes ties digs hopes wheels lines performs settles tenses sniffs stabs wanders downs pounds notes slams twists shows weaves bends bounces curves expects hastily pokes
can sees acts scrolls brings arrives needs follows get zips manages proceeds deposits hardly strikes exhales still smells came tears yanks lifts knew shifts presses grabs excuses
straightens hustles speeds recognizes carries pays also falls stays tastes drags never speaks dials turned slows relaxes brought dumps sticks changes know come lowers flips bought
sleeps greets registers succeeds now even withdraws cracks writhe saw nudges scratches hobbles steals pauses collapses offers puffs grinds braces thought expected levels gave hits
drops spots lurches clambers ever eventually snares writes searches gazes approaches selects eases dips cocks groans lives works winks swerves grips fills rummages loved joins
regains wiggles talks beckons gulps maneuvers zooms stalks seizes vanishes points thrusts hauls leaps hit adjusts heard shoots refers deals honks rams releases clears nears pries
burps mutters extends curls trudges found removes accelerates orders produces inspects tumbles cuts calmly presents scratches none checks drives cries pushes double



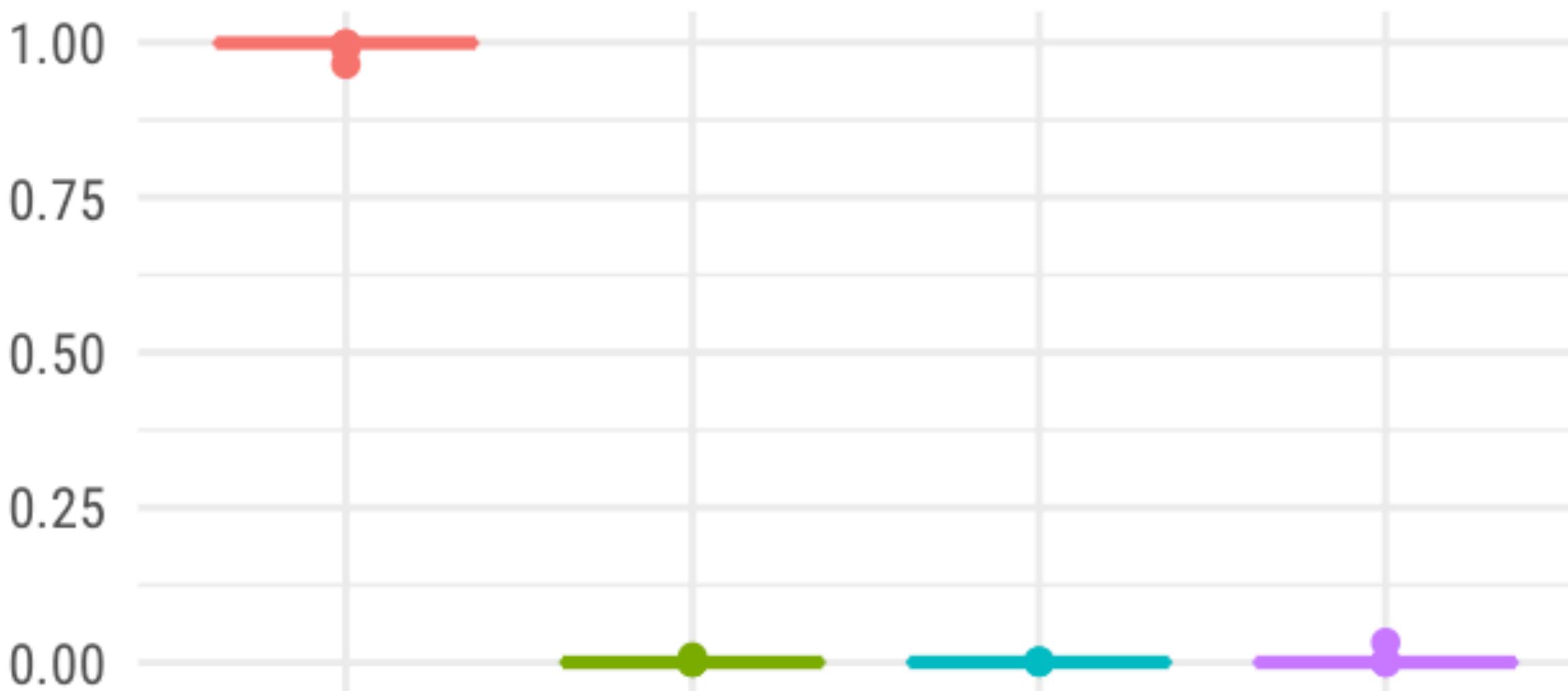
TAKING TIDY TEXT TO
THE NEXT LEVEL

TOPIC MODELING

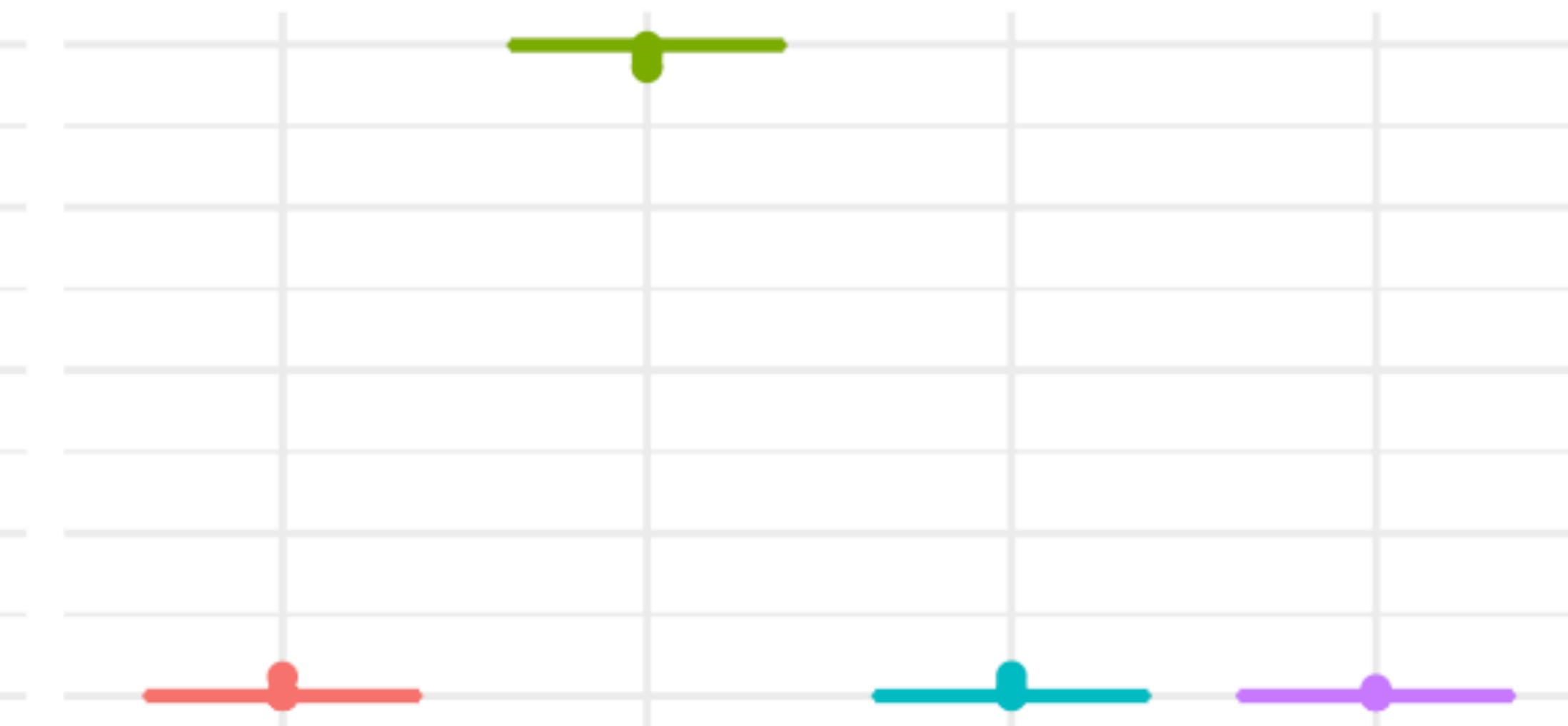
TOPIC MODELING

- Each DOCUMENT = mixture of topics
- Each TOPIC = mixture of words

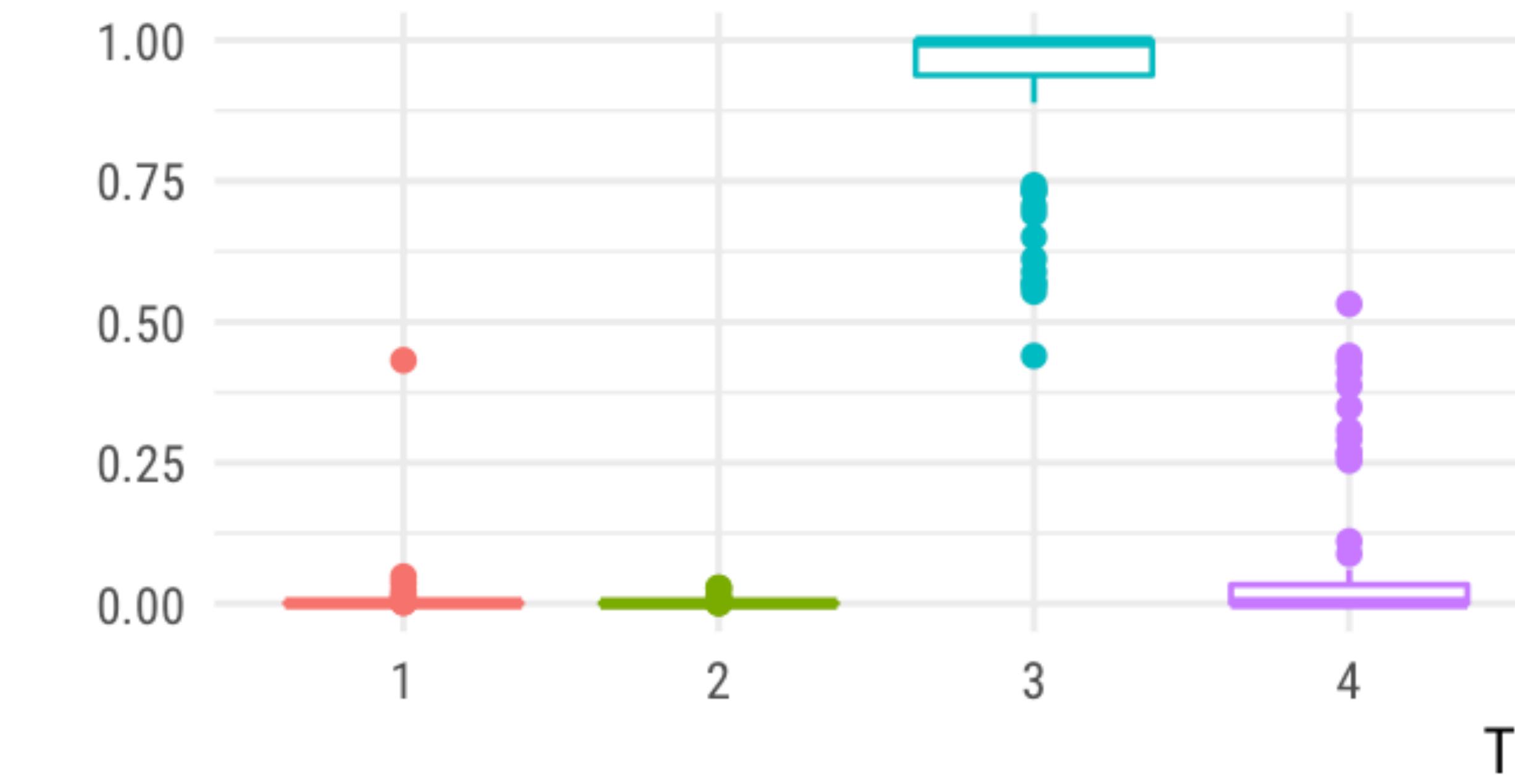
Pride and Prejudice



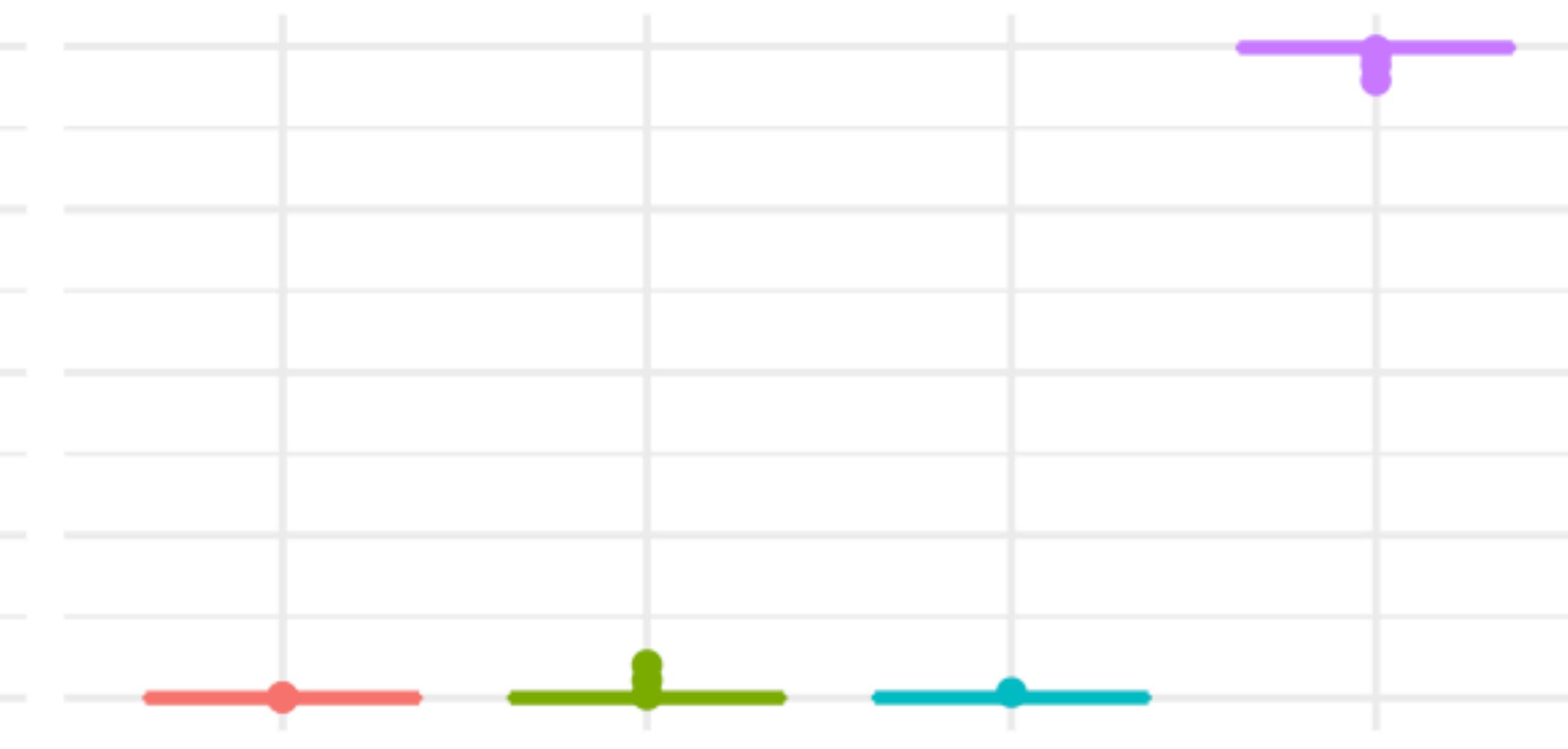
Twenty Thousand Leagues under the Sea



Great Expectations



The War of the Worlds

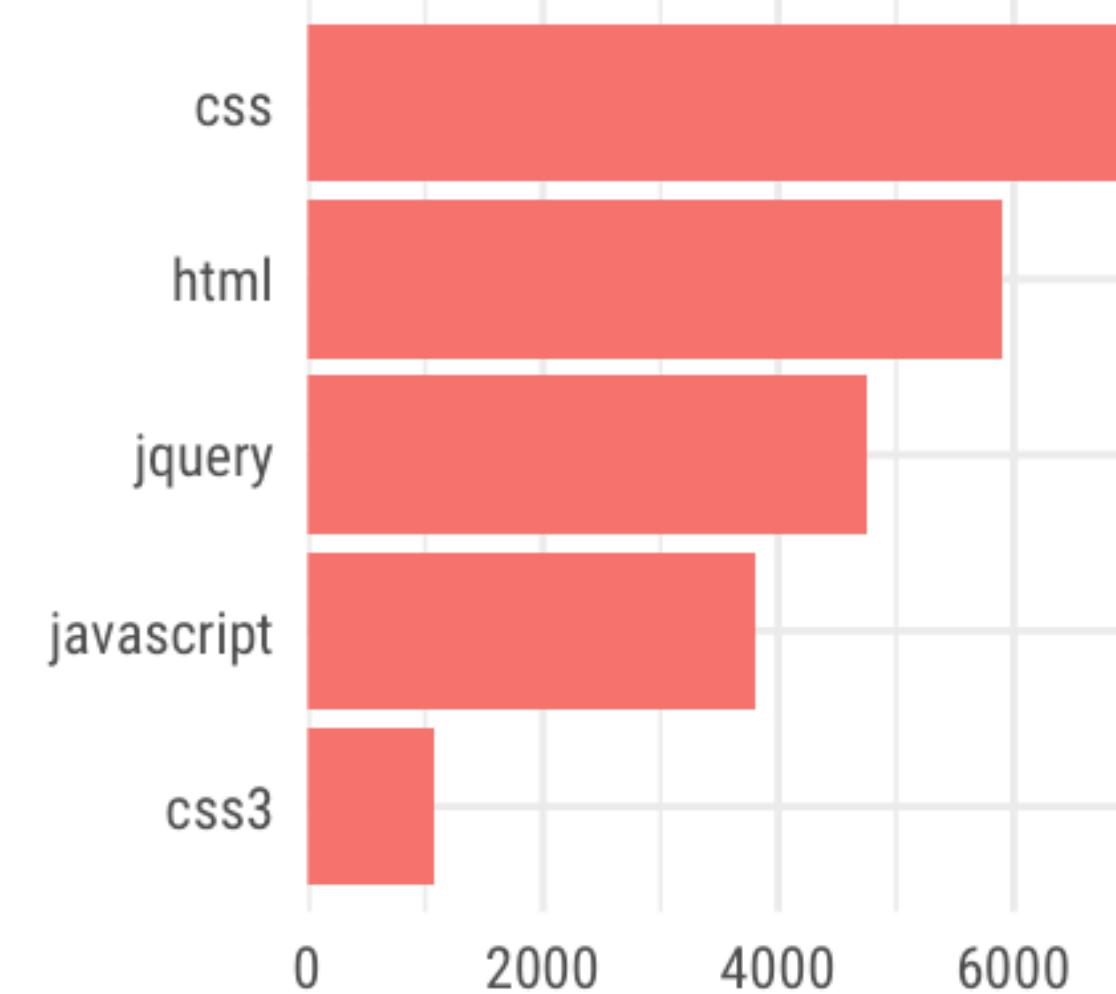


Topic

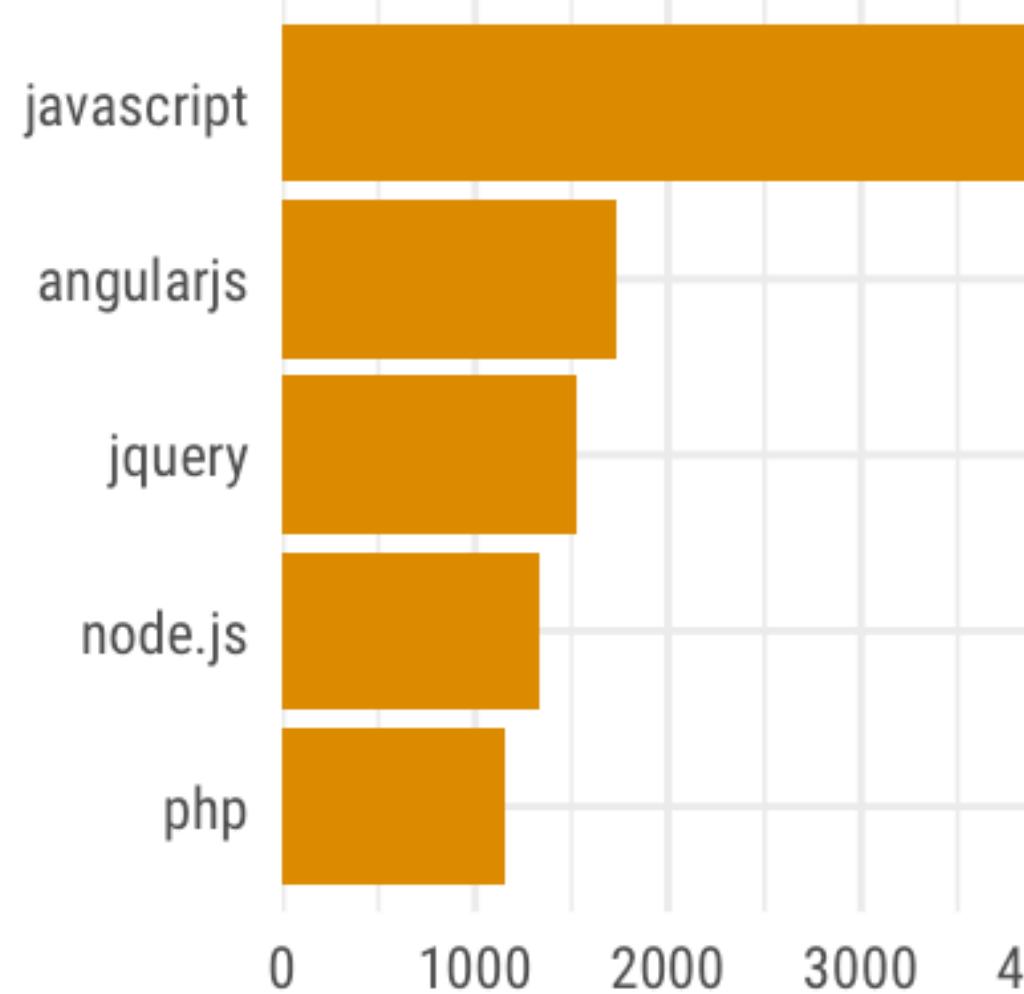
Top tags for each LDA topic

For questions with >80% probability for that topic

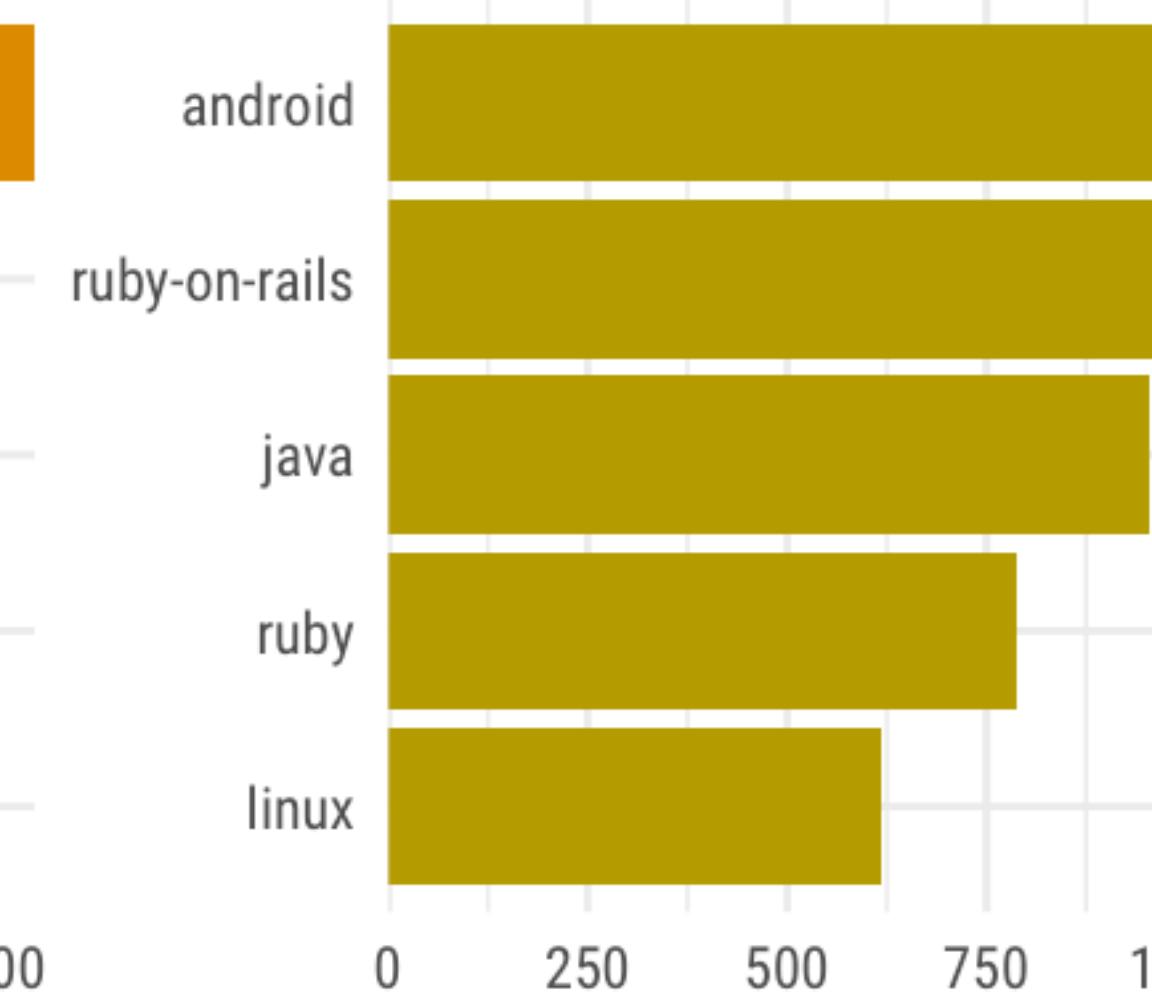
topic 1



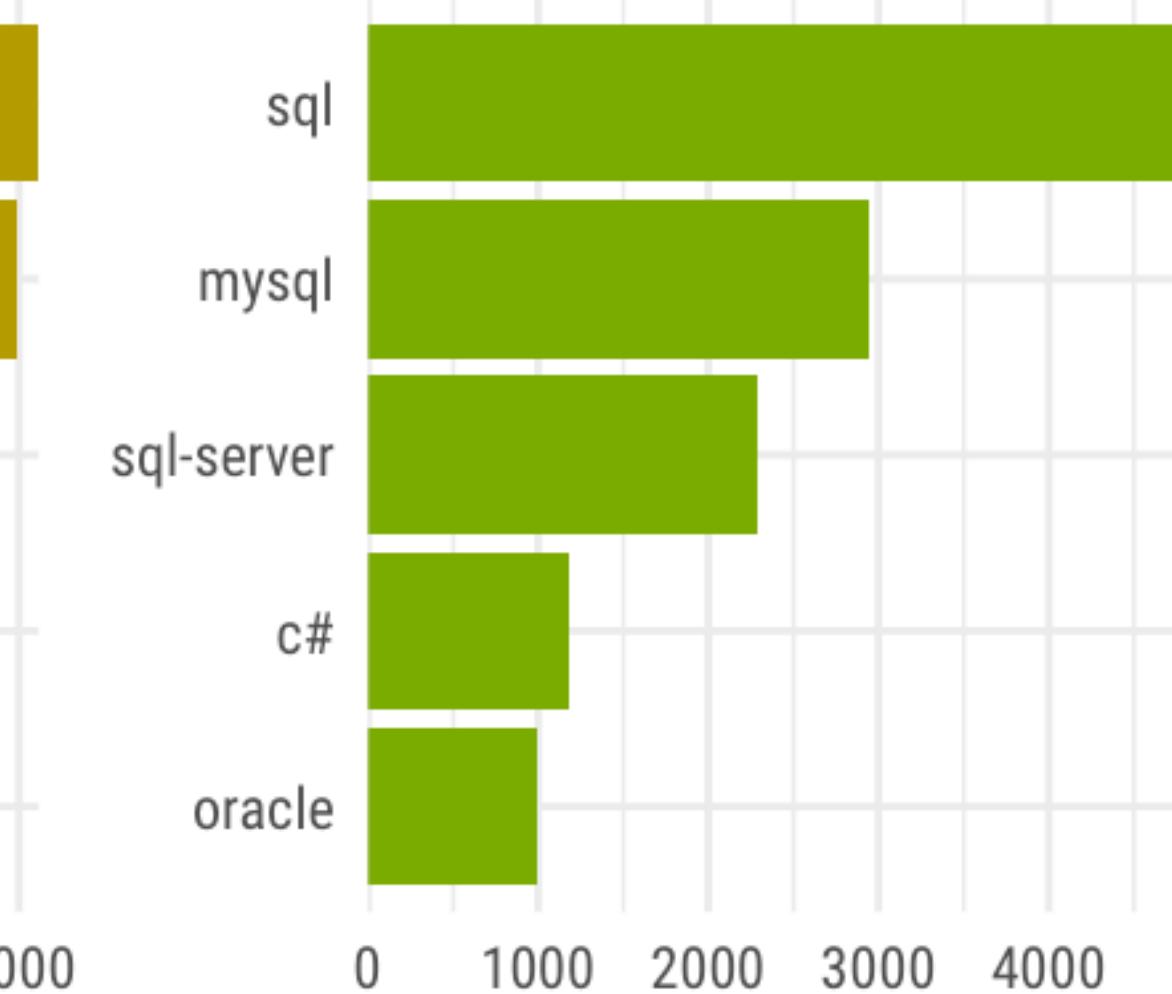
topic 2



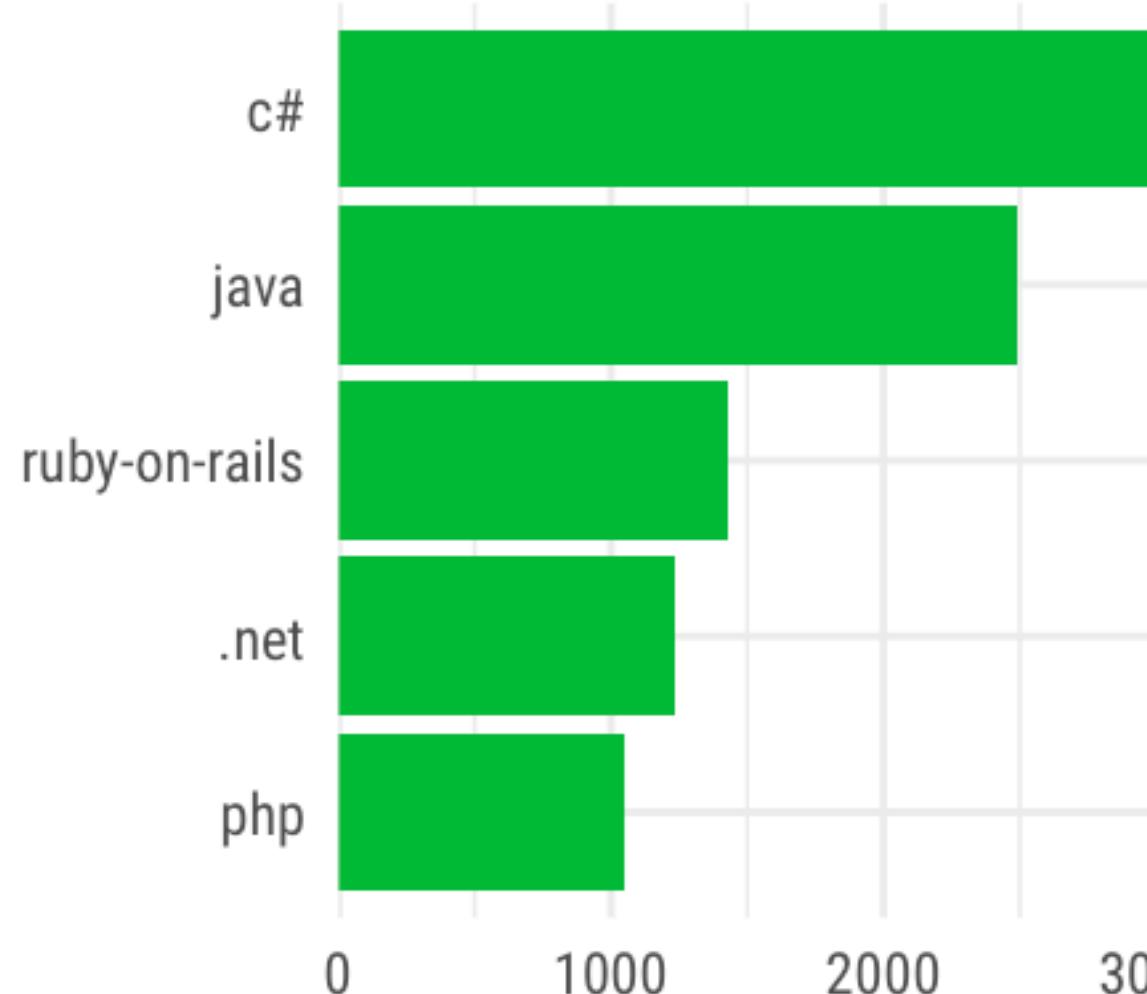
topic 3



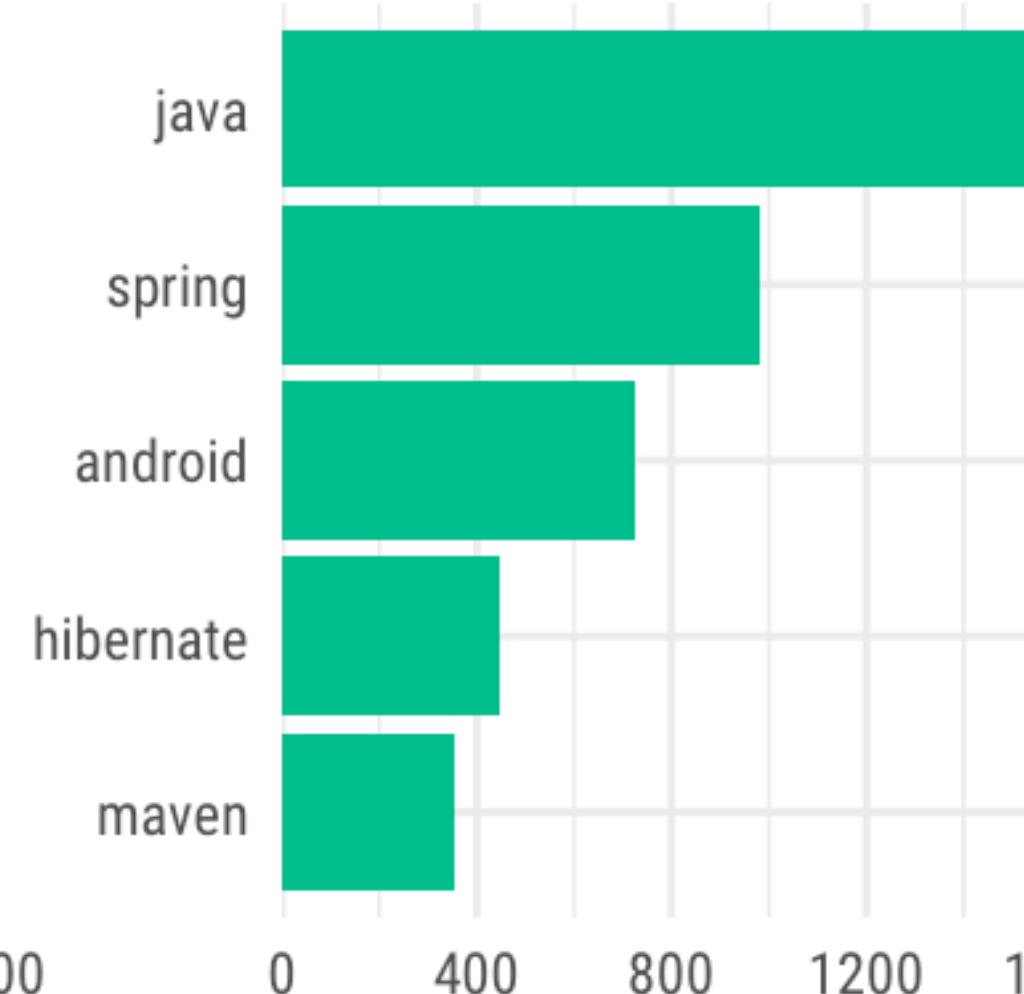
topic 4



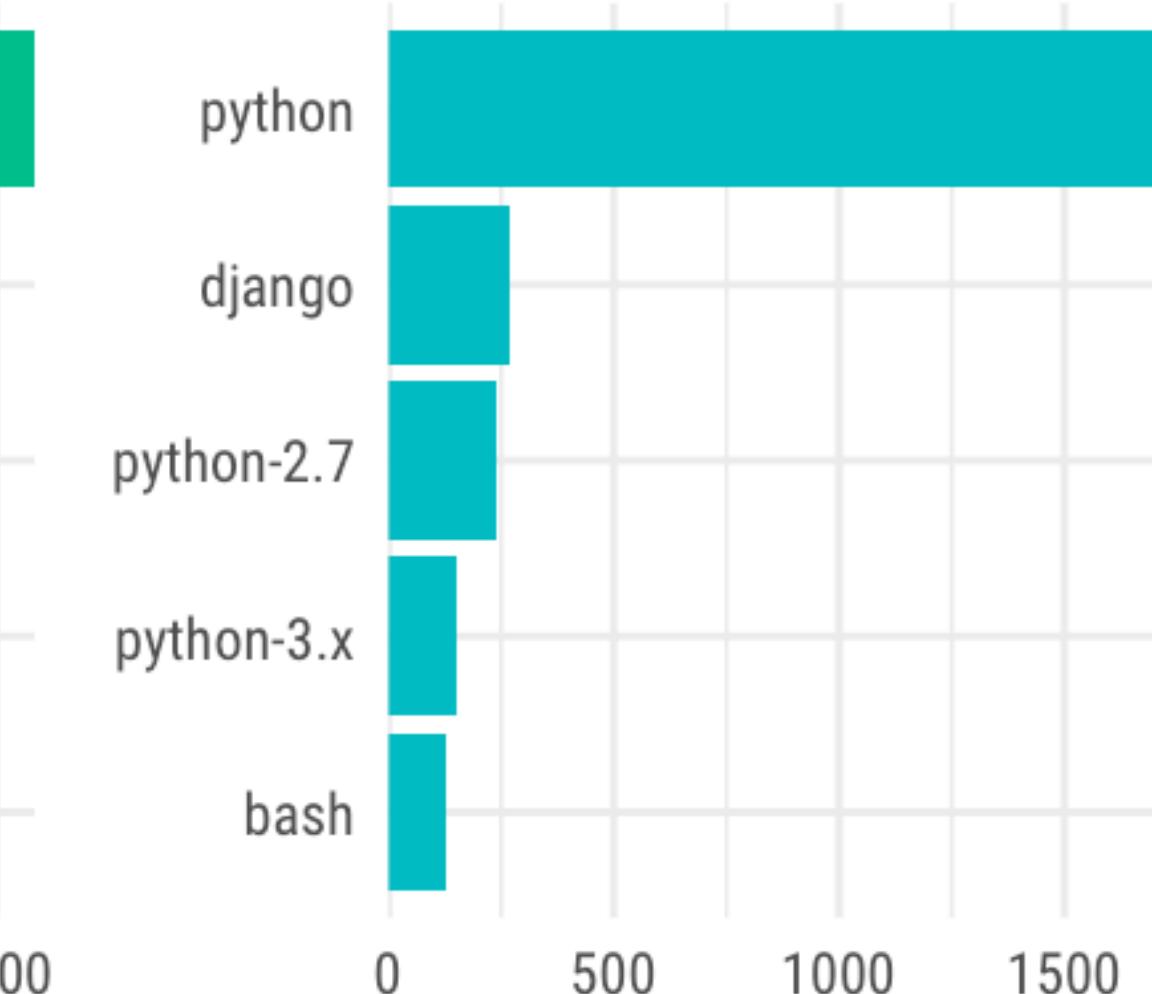
topic 5



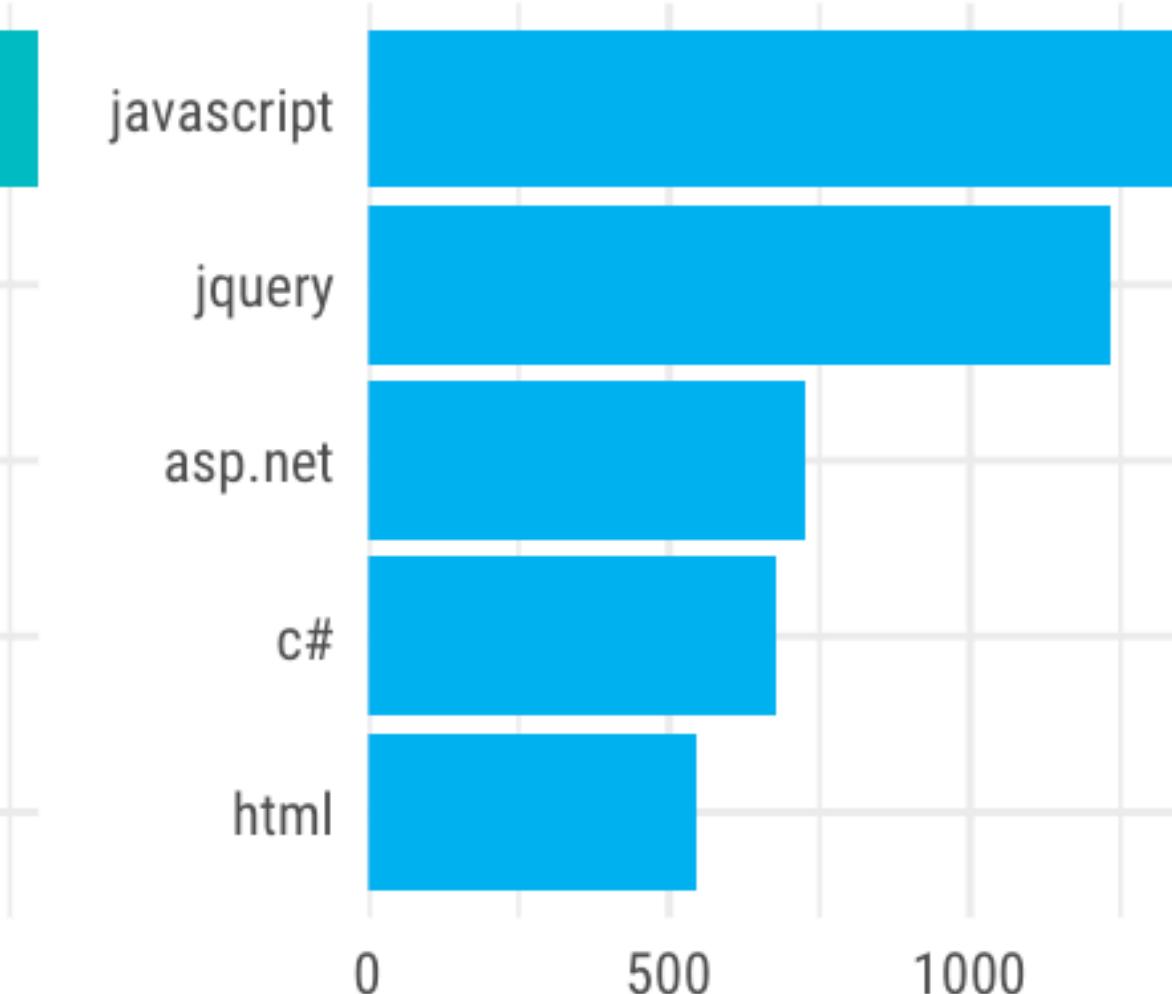
topic 6



topic 7



topic 8



topic 9



topic 10

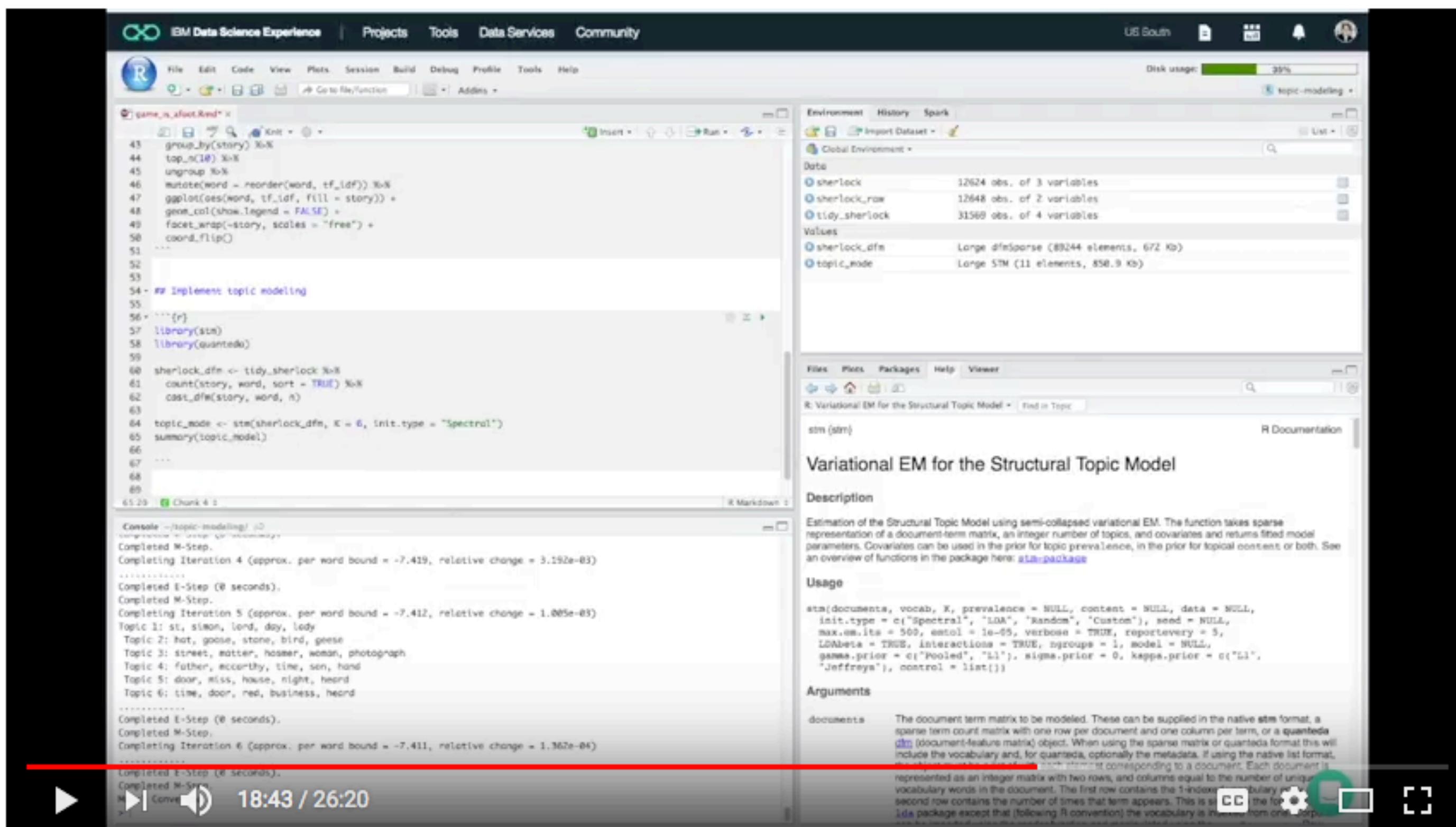


topic 11



topic 12





Topic modeling with R and tidy data principles

1,372 views

42 likes · 0 comments · SHARE · ...



Julia Silge

Published on Dec 18, 2017

SUBSCRIBE 59

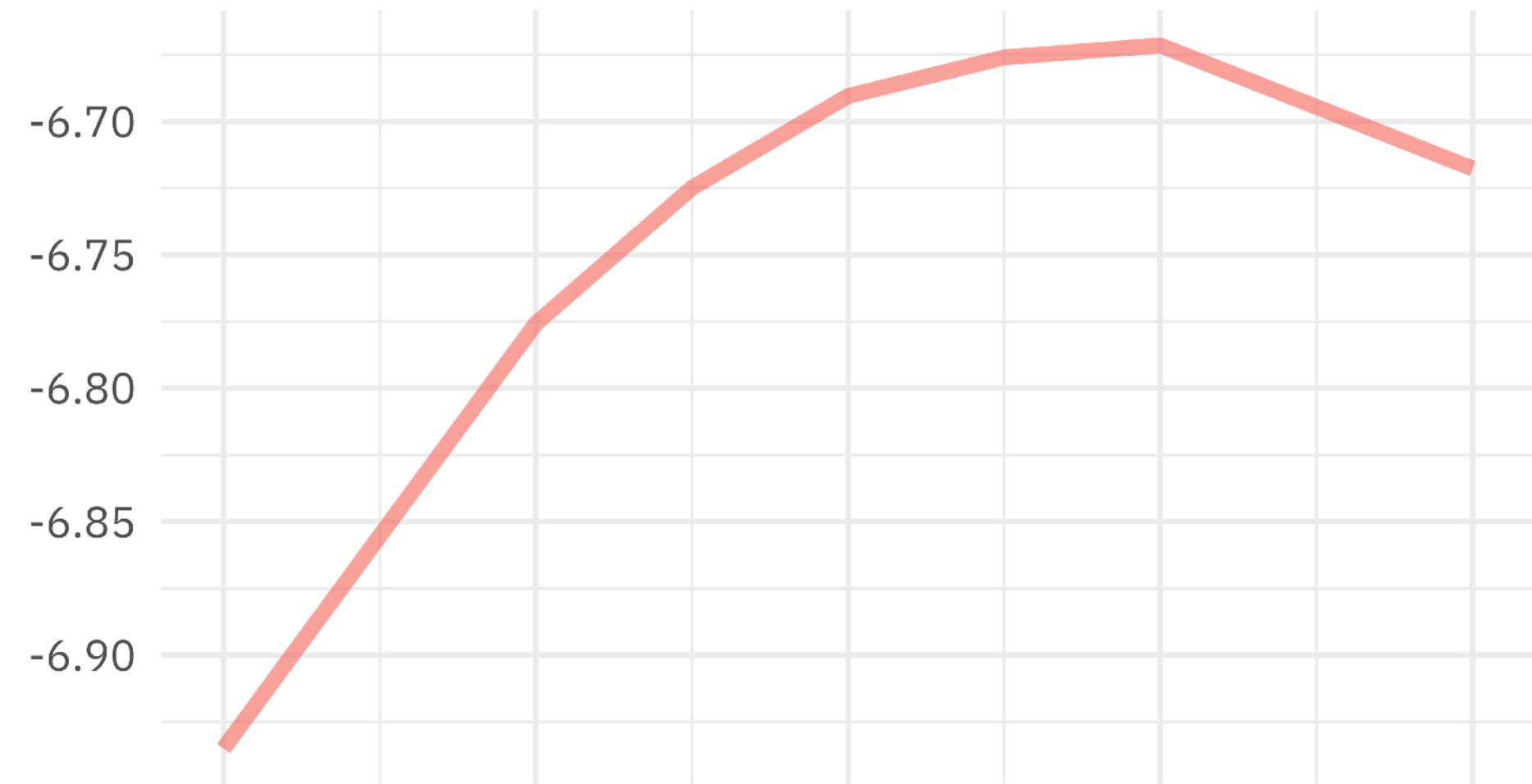
Watch along as I demonstrate how to train a topic model in R using the tidytext and stm packages on a collection of Sherlock Holmes stories. In this video, I'm working in IBM Cloud's Data Science Experience environment.

SHOW MORE

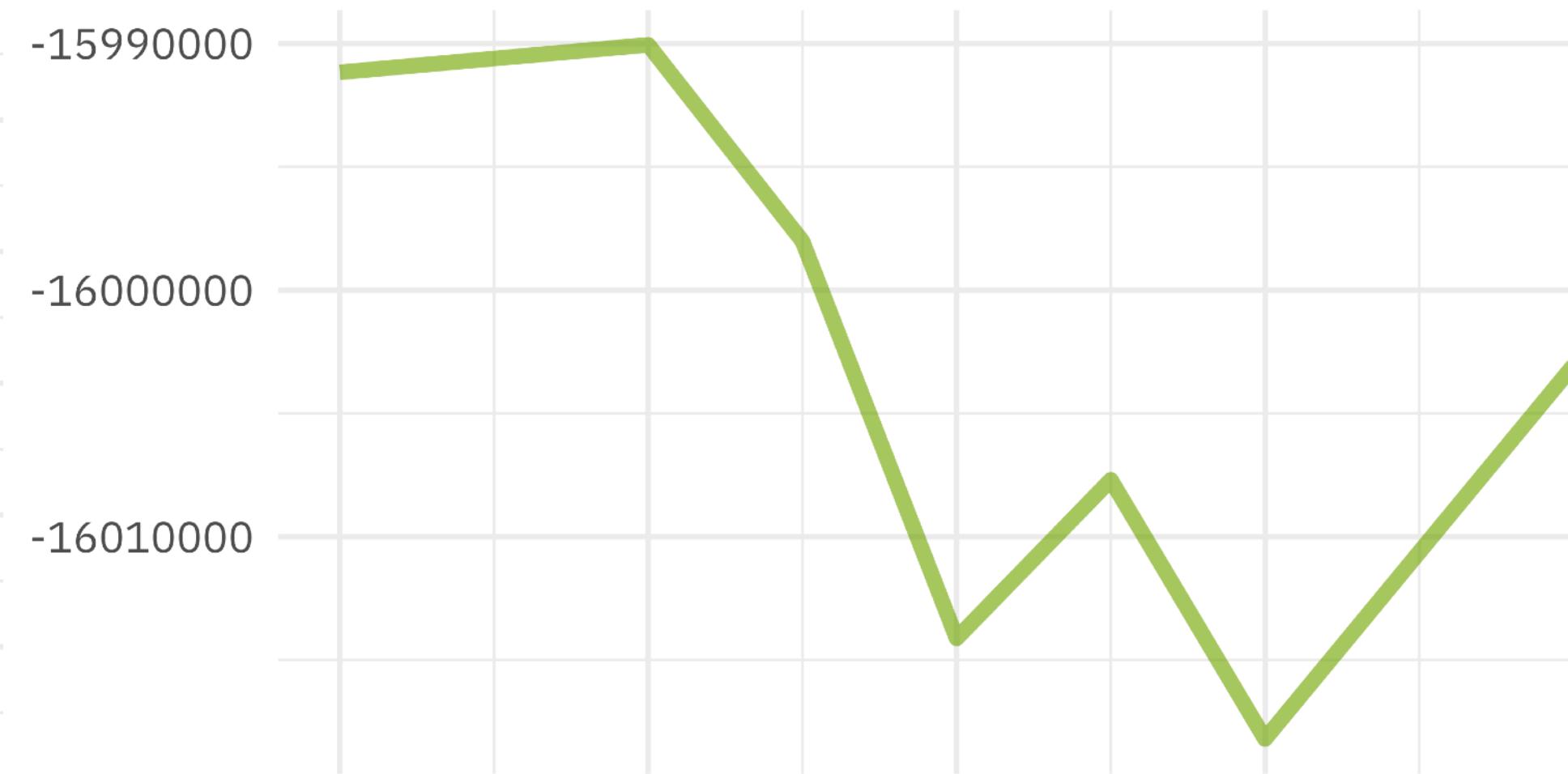
Model diagnostics by number of topics

These diagnostics indicate that a good number of topics would be around 60

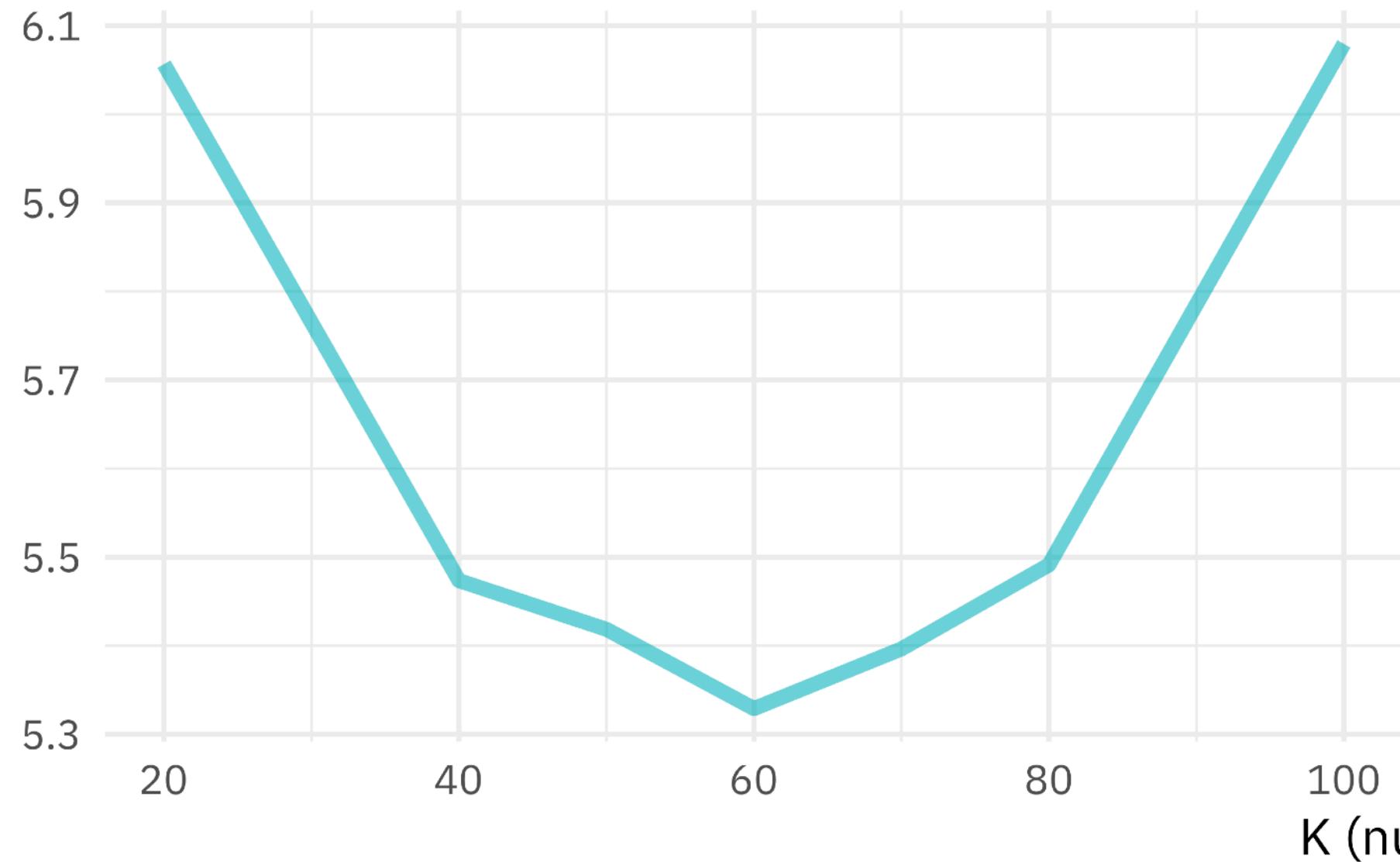
Held-out likelihood



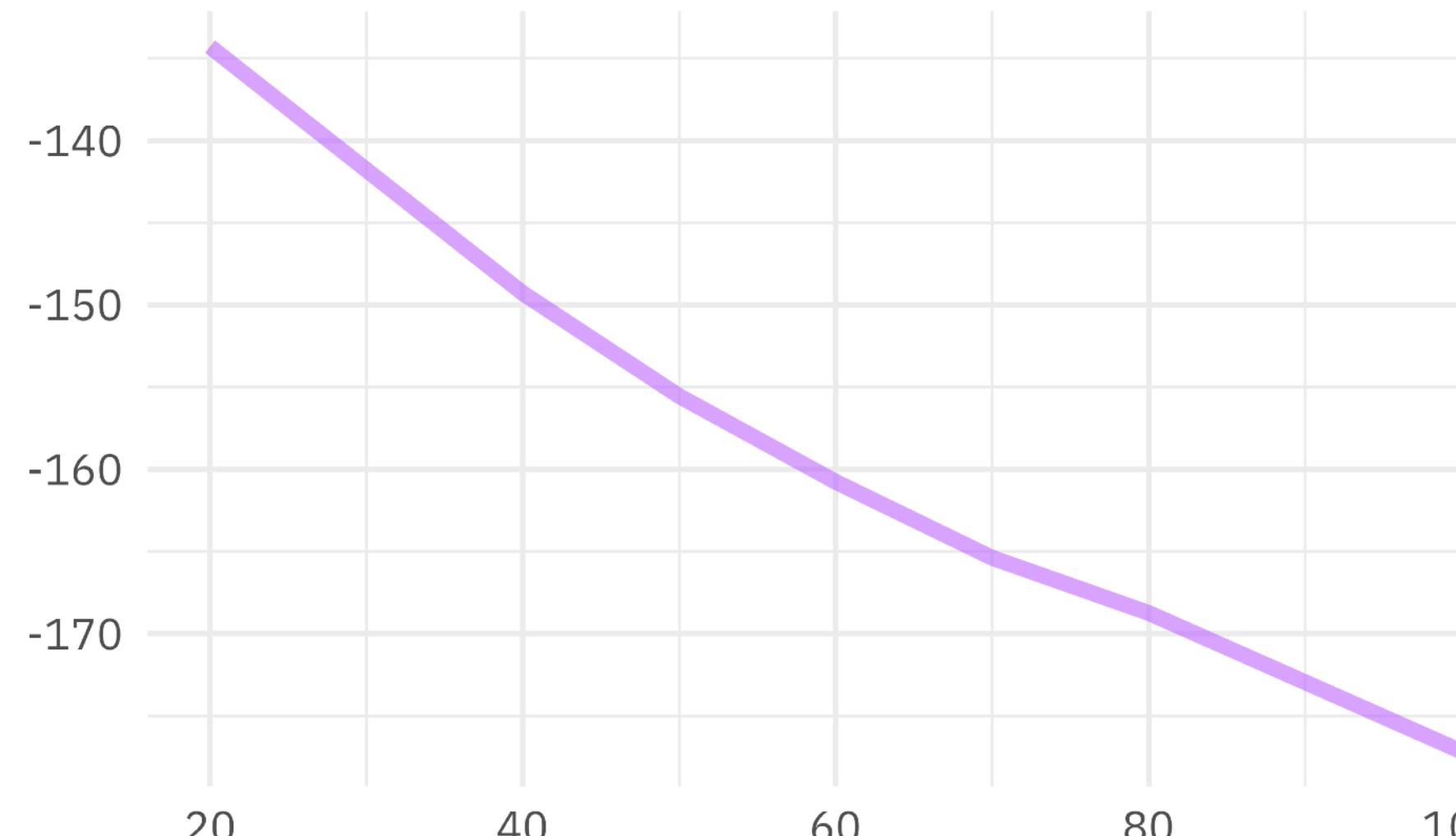
Lower bound



Residuals



Semantic coherence





TAKING TIDY TEXT TO
THE NEXT LEVEL

TEXT CLASSIFICATION

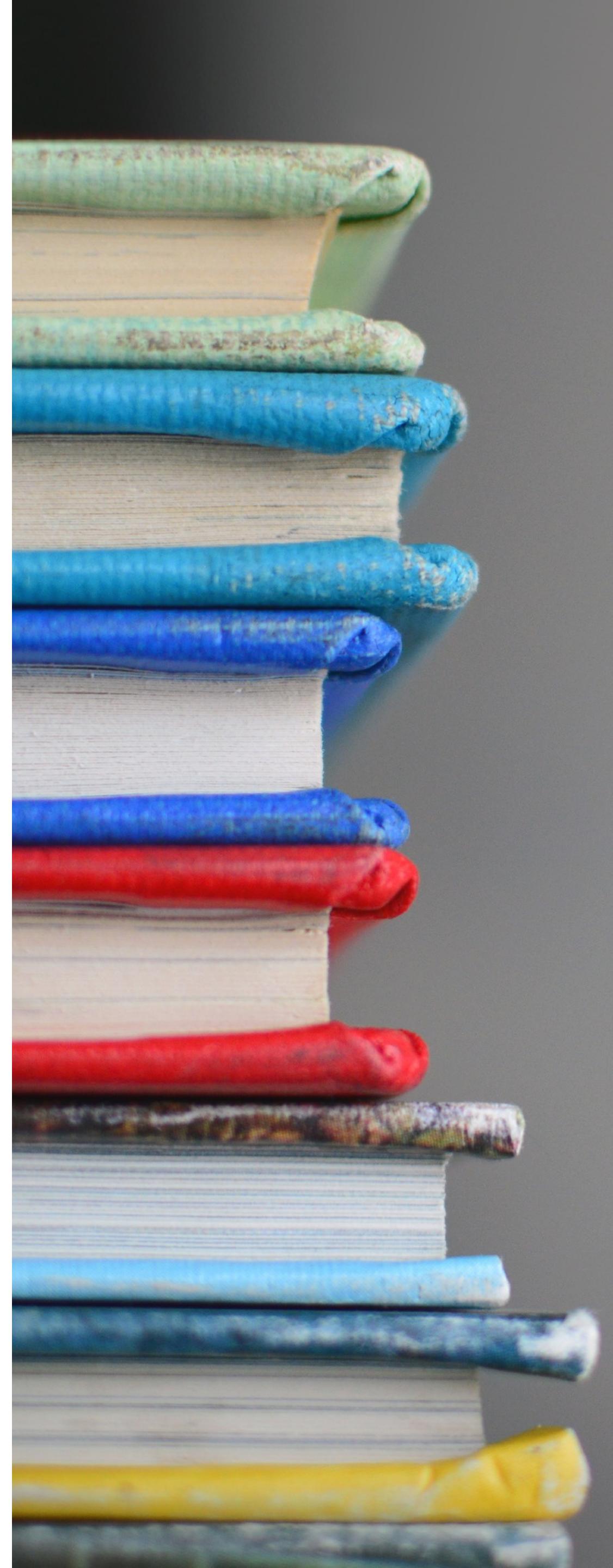
TIDY TEXT

TRAIN A GLMNET MODEL



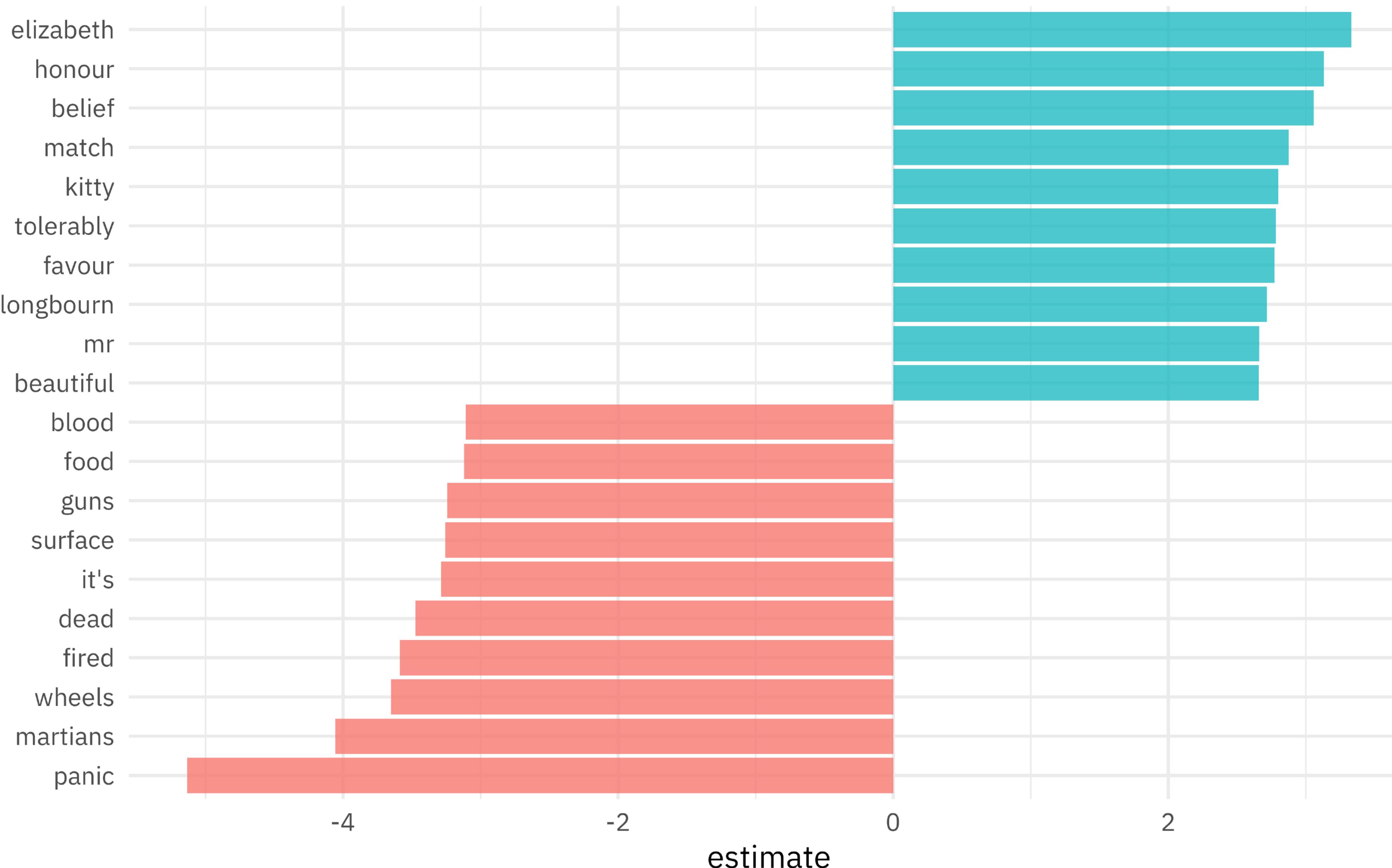
TEXT CLASSIFICATION

```
> library(glmnet)
> library(doMC)
> registerDoMC(cores = 8)
>
> is_jane <- books_joined$title == "Pride and Prejudice"
>
> model <- cv.glmnet(sparse_words, is_jane, family = "binomial",
+                      parallel = TRUE, keep = TRUE)
```



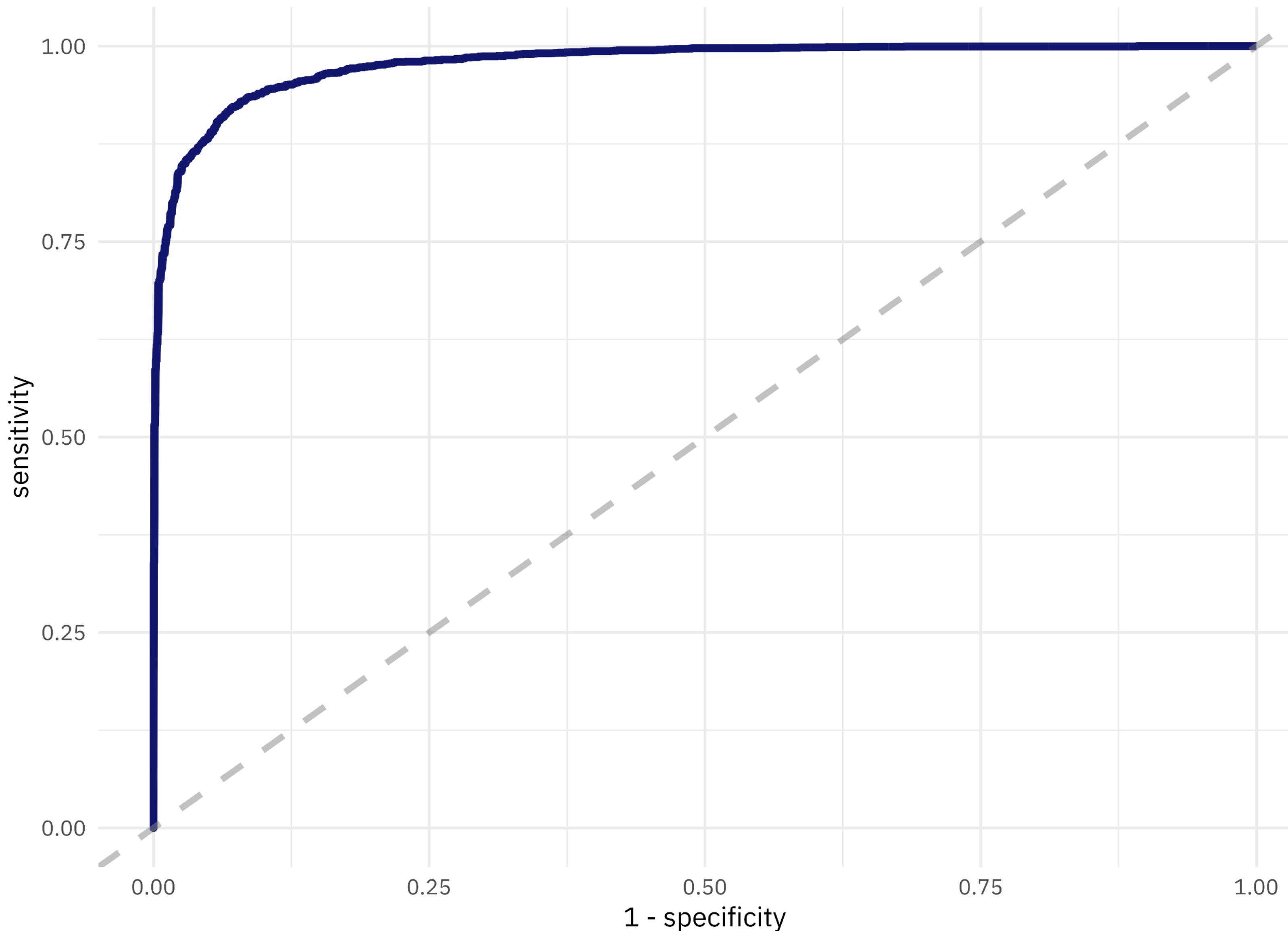
Coefficients that increase/decrease probability the most

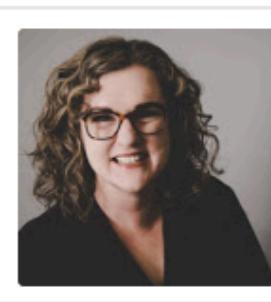
A document mentioning Martians is unlikely to be written by Jane Austen



ROC curve for text classification using regularized regression

Predicting whether text was written by Jane Austen or H.G. Wells





Julia Silge

Data Scientist at Stack Overflow

📍 Salt Lake City, UT, United States ↗ http://juliasilge.com/ 🐦 juliasilge
✉️ juliasilge

I enjoy making beautiful charts, the statistical programming language R, black coffee, red wine, and the mountains of my adopted home here in Utah. I have a PhD in astrophysics and an abiding love for Jane Austen. My work involves analyzing and modeling complex data sets while communicating about technical topics with diverse audiences.

Favorite editor: RStudio

I want to work with

- r
- ggplot2
- dplyr
- data-visualization
- shiny
- text-mining
- knitr
- machine-learning

TOP 5% r rstudio TOP 10% dplyr tidyverse TOP 30% ggplot2

Position • Dec 2016 → Current (2 years, 10 months)

Data Scientist at Stack Overflow

- r
- tidyverse
- ggplot2
- shiny

- Analyze large datasets and build models to understand developers and

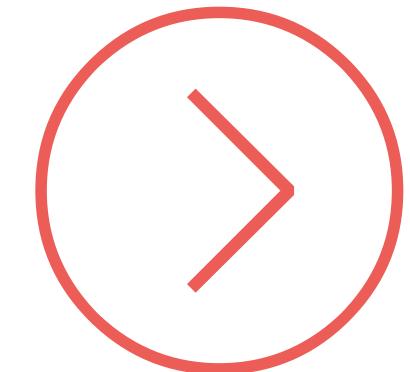
[Read more](#)

tidytext: Text mining using dplyr, ggplot2, and other tidy tools

Last commit on Sep 01, 19

351 Commits / 383,231 ++ / 314,272 --

Using tidy data principles can make many text mining tasks easier, more effective, and



search jobs search companies

Search all jobs Located anywhere Save Search

Remote Tech Compensation Perks Background More

Debug

9,735 results

FEATURED

★ **R Developer**
YouGov - No office location
Remote

r shiny python etl

✓ Select min. experience to Select max. experience

Student Junior Mid-Level Senior Lead Manager

Full-time Contract Internship

Developer

Apply filters Cancel

TIDY TEXT

THANK YOU



JULIA SILGE

@juliasilge

<https://juliasilge.com>

TIDY TEXT

THANK YOU



JULIA SILGE

@juliasilge

<https://juliasilge.com>

Author portraits from Wikimedia

Photos by Glen Noble and Kimberly Farmer on Unsplash