

Natural Language Visualization with Scattertext

Jason S. Kessler*

Global AI Conference

April 27, 2018

Code for all visualizations is available at:

<https://github.com/JasonKessler/GlobalAI2018>



*No, not *that* Jason Kessler

@jasonkessler

Lexicon speculation

	Proposed word lists	Accuracy	Ties
Human 1	positive: <i>dazzling, brilliant, phenomenal, excellent, fantastic</i> negative: <i>suck, terrible, awful, unwatchable, hideous</i>	58%	75%
Human 2	positive: <i>gripping, mesmerizing, riveting, spectacular, cool, awesome, thrilling, badass, excellent, moving, exciting</i> negative: <i>bad, cliched, sucks, boring, stupid, slow</i>	64%	39%

Figure 1: Baseline results for human word lists. Data: 700 positive and 700 negative reviews.

Lexicon mining ≈ lexicon speculation

	Proposed word lists	Accuracy	Ties
Human 1	positive: <i>dazzling, brilliant, phenomenal, excellent, fantastic</i> negative: <i>suck, terrible, awful, unwatchable, hideous</i>	58%	75%
Human 2	positive: <i>gripping, mesmerizing, riveting, spectacular, cool, awesome, thrilling, badass, excellent, moving, exciting</i> negative: <i>bad, cliched, sucks, boring, stupid, slow</i>	64%	39%

Figure 1: Baseline results for human word lists. Data: 700 positive and 700 negative reviews.

	Proposed word lists	Accuracy	Ties
Human 3 + stats	positive: <i>love, wonderful, best, great, superb, still, beautiful</i> negative: <i>bad, worst, stupid, waste, boring, ?, !</i>	69%	16%

Figure 2: Results for baseline using introspection and simple statistics of the data (including *test* data).

Bo Pang, Lillian Lee and Shivakumar Vaithyanathan. Thumbs up? Sentiment classification using machine learning techniques. EMNLP. 2002.

@jasonkessler

Language and Demographics

hobos

almond
butter

Our source:
526,000 OkCupid
profiles

My Self-Summary
In my free time I train **hobos** for the circus.
<http://www.youtube.com/watch?v=3v7MBUJcws>

What I'm doing with my life
study languages, drink gin and pickle brine **martini**, bake, **almond butter** with a spoon, draw things with a pen, write things with a different pen.

applying to graduate programs, fighting the urge to say fuck it and go to **pastry school**.

I'm really good at
selling food,
judging.

The first things people usually notice about me
if you're **japanese**: it's that I look like hermione granger.

My Self-Summary
I'm a Colorado native who loves to ski but gave up the mountains of Colorado for the city five years ago. I feel like a **Park City girl** living a Big City life - not really firmly belonging in either place. I love it here though, especially when I discover something I didn't know about before, recapturing that sense of adventure even though this is the place I live. **I like adventures!**

I am smart, sensitive, and sweet.

What I'm doing with my life
I don't know why I find this question so hard to answer. I'm pursuing an **acting career** but have had lots of (varied) jobs in between gigs. I temp sometimes. I wait tables sometimes. I've left the city to do shows. I teach kids drama and I babysit. I'm having a lot of fun in NYC.

I'm really good at
Bikram yoga. I'm an excellent cook but not a great baker.

My Self-Summary
I'm a California girl, currently in a New York state of mind. I moved to NYC sat for **grad school** and stayed on for work. I'm enjoying my time here exploring Manhattan, but ultimately planning on moving back to CA. I'm pretty **easy-going** and I like to try new things.

I am composed, thoughtful, and sarcastic.

My favorite books, movies, music, and food
Books: **Anything by Thomas Lehrn, 100 Years of Solitude, The Great Gatsby**

Movies: **Jurassic Park, Coming to America, Women on the Verge of a Nervous Breakdown**

Music: **Tupac, E-40, Jay-Z, Mariah Carey, Timbaland, Justin Timberlake, Drake**

Food: Japanese, Thai, New American, California Cuisine

100 Years of
Solitude

Bikram yoga
@jasonkessler

OKCupid: Words and phrases that distinguish white men.



tom clancy van halen **golfing**
harley davidson **ghostbusters** phish
the big lebowski soundgarden **brew** boating **nofx**
groundhog day **hockey** jeep **blazing saddles** the red sox
the dropkick murphys megadeth **grilling** ccr
robert heinlein boats **skiing** zappa **nascar** motorcycles
software dark tower **the hitchhiker's guide to the galaxy** breaking bad
band of brothers burn notice **coen brothers** michael crichton **bad religion** tenacious d
mostly rock i'm a country boy **building things** queens of the stone age **mountain biking**
i can fix anything **the offspring** a few beers **apocalypse now** lock, stock, and two smoking barrels
hunting and fishing most sports **world war z** guitar

In general, I won't comment too much on these lists, because the whole point of this piece is to let the groups speak for themselves, but I have to say that the mind of the white man is the world's greatest sausagefest. Unless you're counting **Queens of the Stone Age**, there is not even one vaguely feminine thing on his list, and as far as broad categories go we have: sweaty guitar rock, bro-on-bro comedies, things with engines, and dystopias.

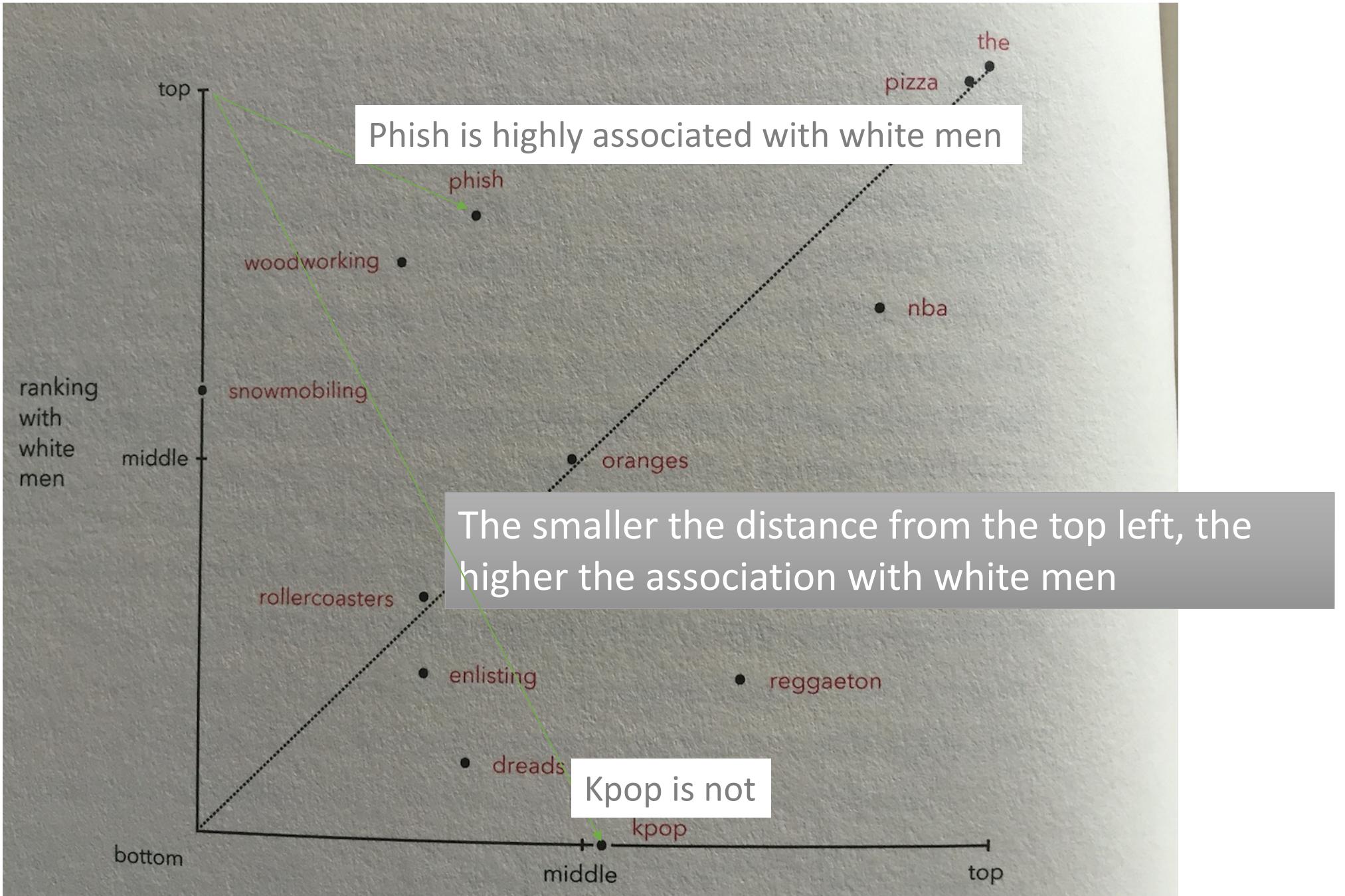
OKCupid: Words and phrases that distinguish Latin men.

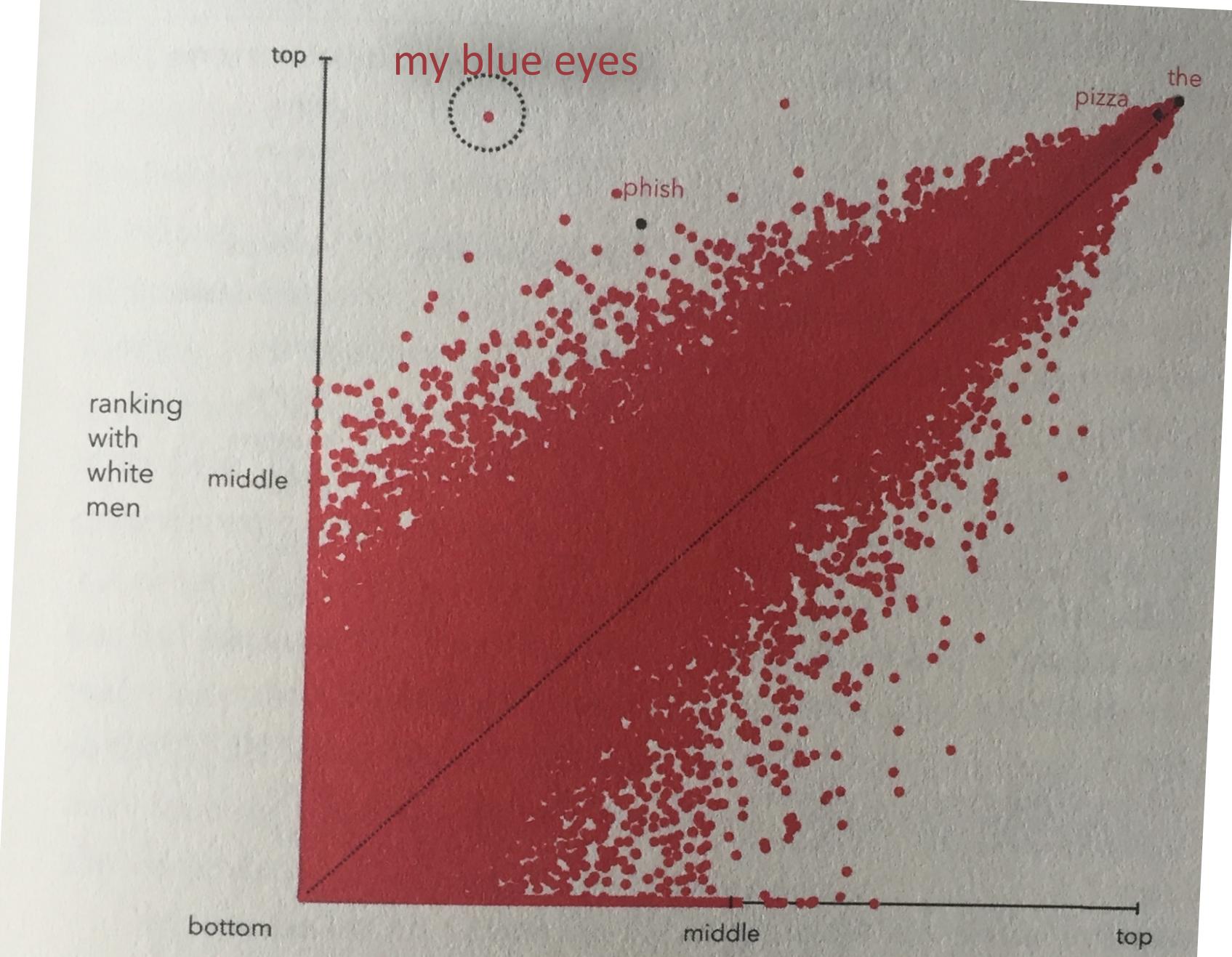


merengue **bachata** **colombian**
hispanic **latino** **dominican** **stationed** **peruvian**
reggaeton **familia** **cuban** **musica** **salsa** **soccer** **amigos** **peru**
boxing **automotive** **baseball** **hola** **marines** **mma** **hip hop** **ufc**
i'm a funny guy **respectful** **mars volta** **some drinks** **what u** **sports** **u wanna** **xbox 360**
of mice and men **chill** **comedy** **art of war** **very funny** **saving private ryan** **i'm a simple guy** **hip hop**
full metal jacket **down to earth guy** **world war z** **law enforcement** **outgoing and funny** **bars** **attending college**
forrest gump **the strokes** **all sports**

Explanation

Music and dancing—**merengue**, **bachata**, **reggaeton**, **salsa**—are obviously very important to Latinos of both genders. The men have two other fascinating things going on: an interest in telling you about their sense of humor (**i'm a funny guy**, **very funny**, **outgoing and funny**, etc.) and an interest in industrial strength ass-kicking (**mma**, **ufc**, **boxing**, **marines**, etc.) Basically, if a Latin dude tells you a joke, you should laugh.

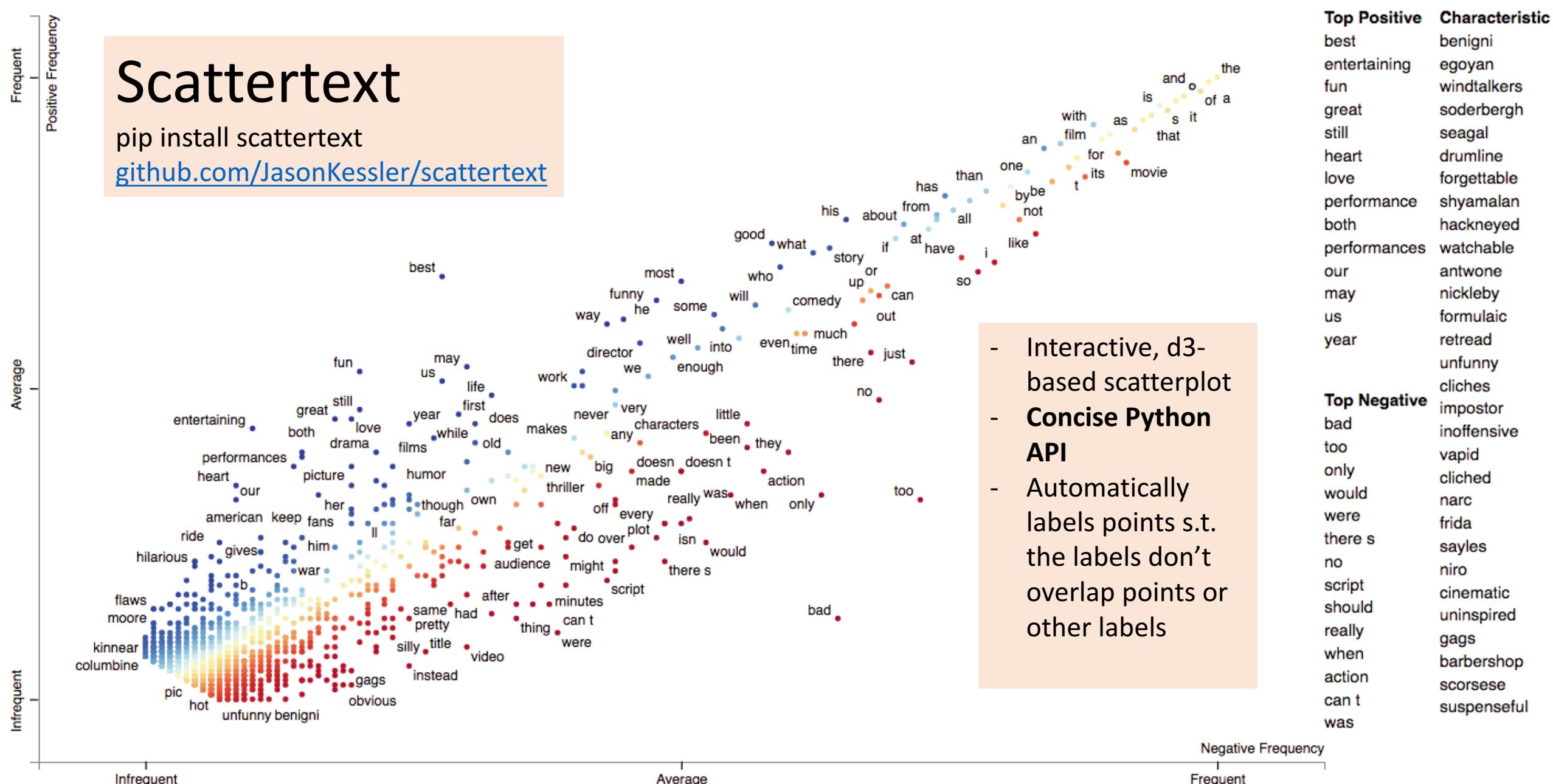




Source: Christian Rudder. Dataclysm. 2014.

ranking with everyone else

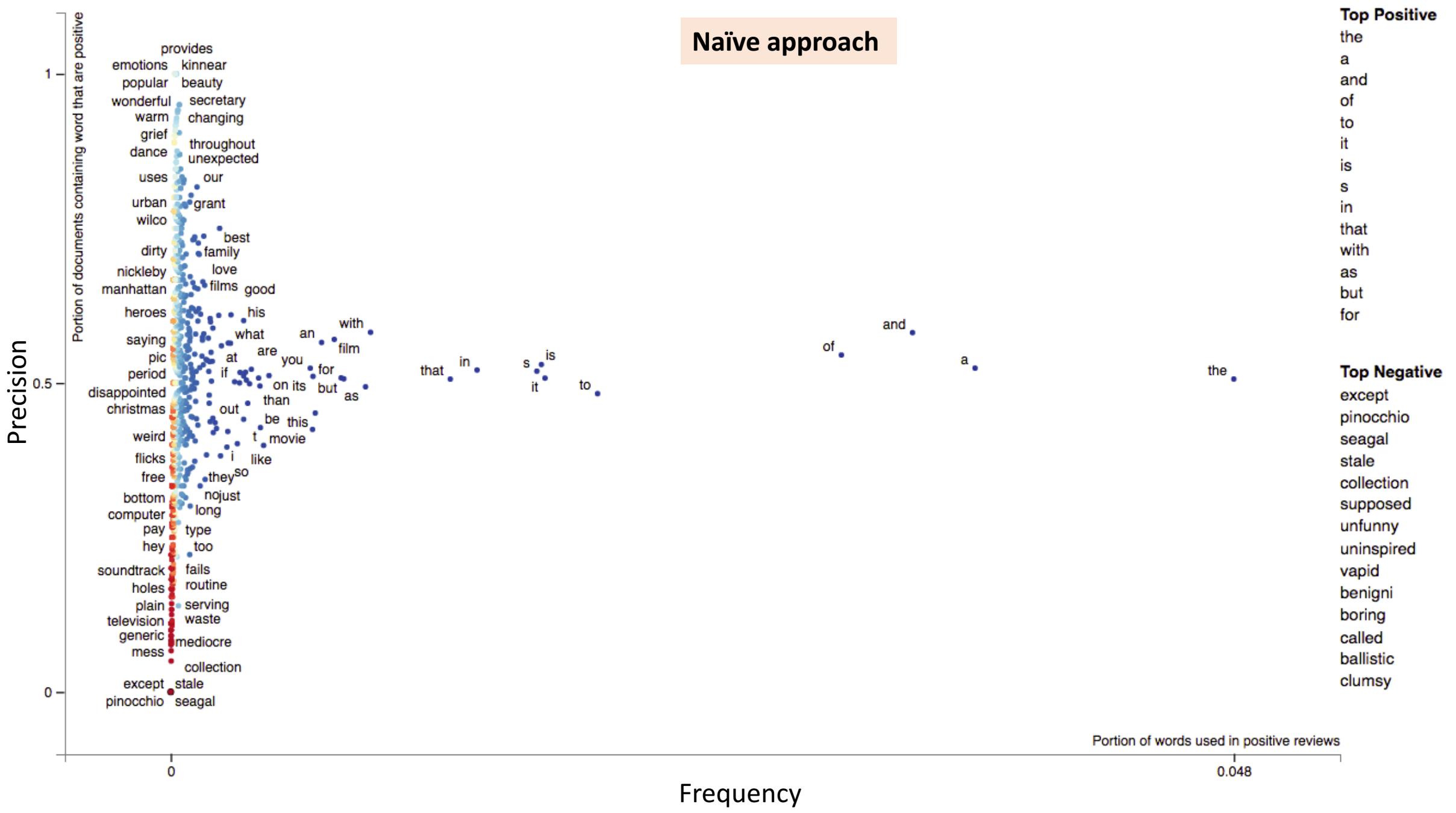
@jasonkessler



Scaled F-Score

- Term-class* associations:
 - “Good” is associated with the “positive” class
 - “Bad” with the class “negative”
- Core intuition: association relies on two necessary factors
 - **Frequency:** How often a term occurs in a class
 - **Precision:** $P(\text{class} \mid \text{document contains term})$
- F-Score:
 - Information retrieval evaluation metric
 - Harmonic mean between precision and recall
 - Requires **both** metrics to be high
- *Term: is defined to be a word, phrase or other discrete linguistic element

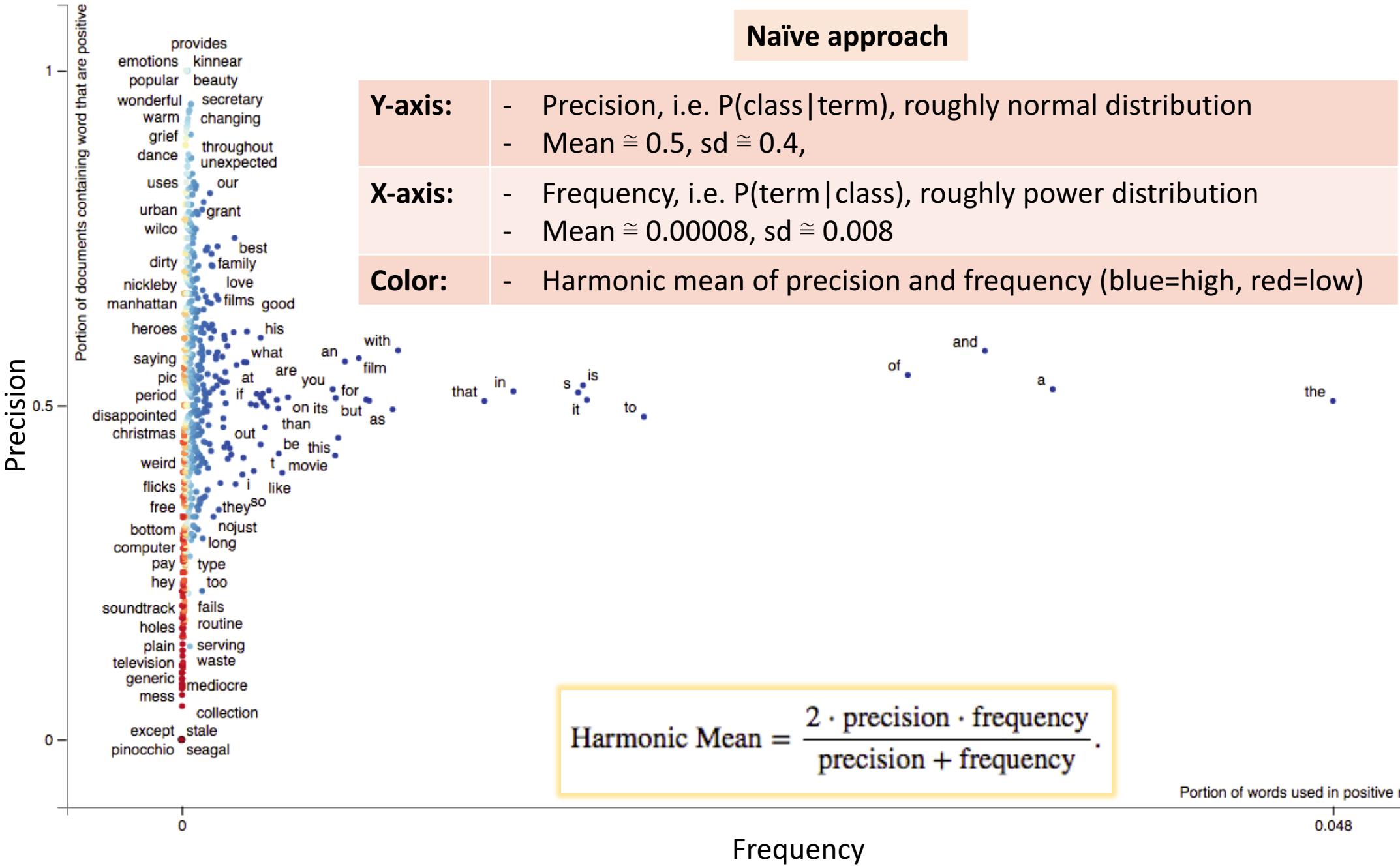
Naïve approach



Top Positive
 the
 a
 and
 of
 to
 it
 is
 s
 in
 that
 with
 as
 but
 for

Top Negative
 except
 pinocchio
 seagal
 stale
 collection
 supposed
 unfunny
 uninspired
 vapid
 benigni
 boring
 called
 ballistic
 clumsy

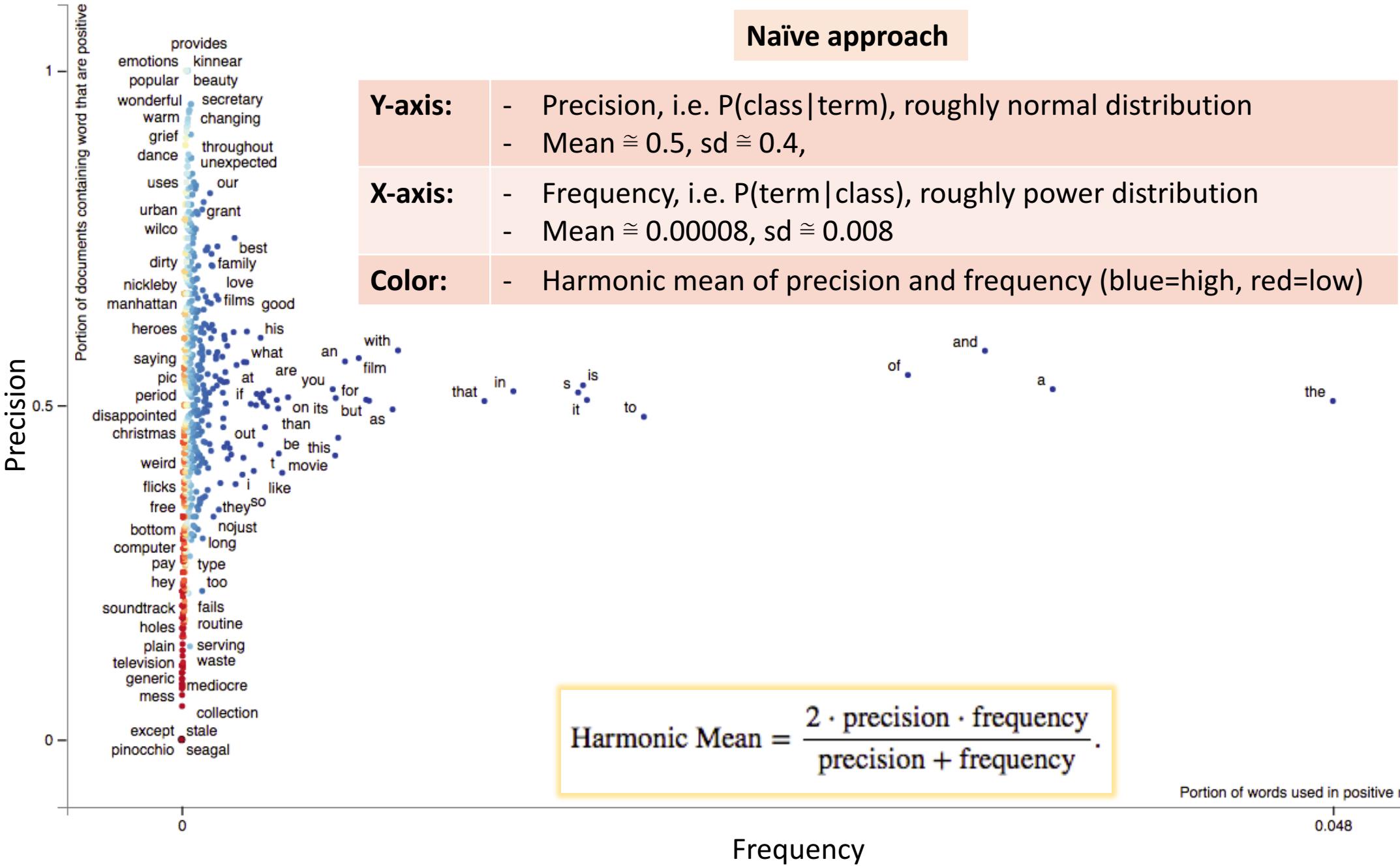
Naïve approach

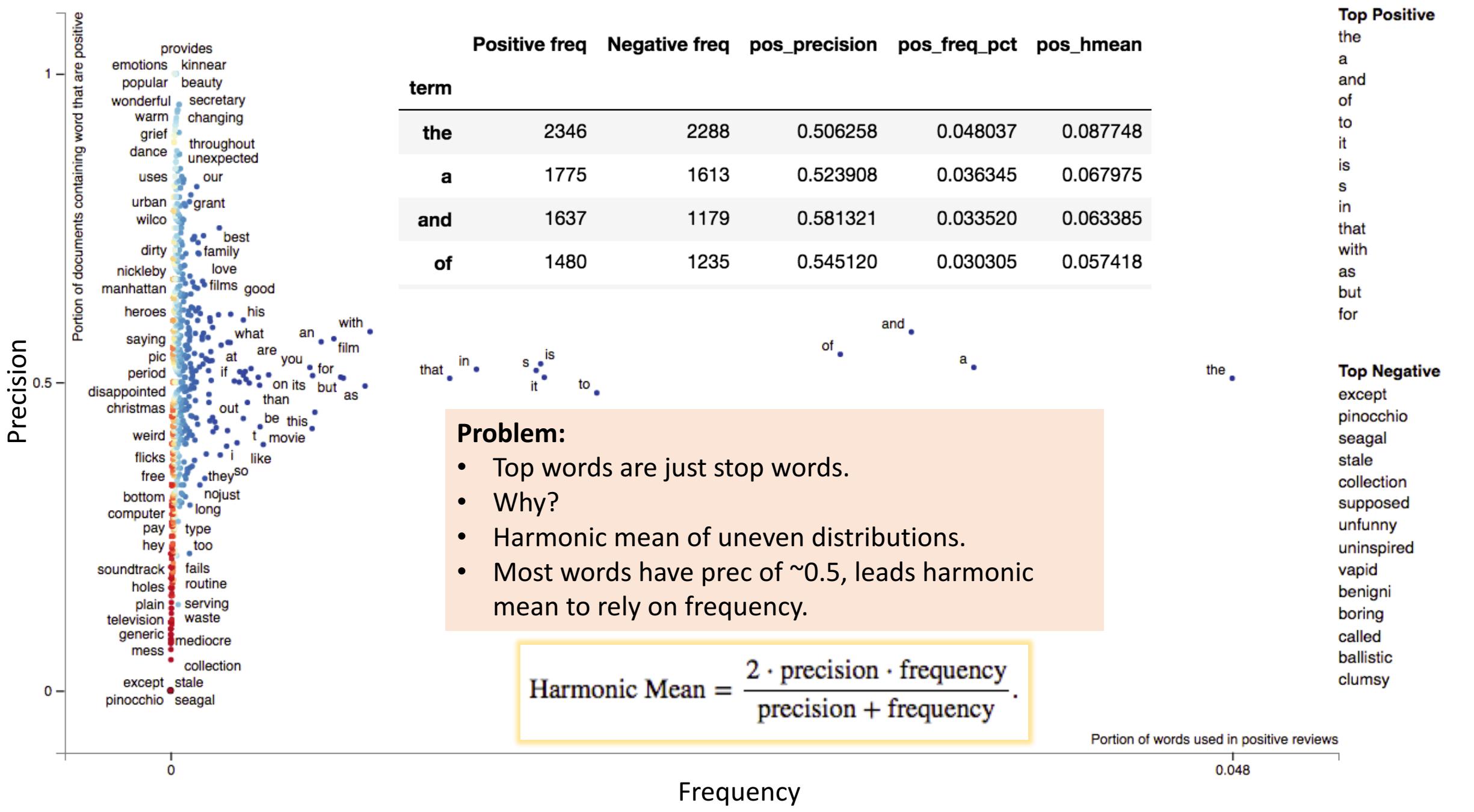


Top Positive
 the
 a
 and
 of
 to
 it
 is
 s
 in
 that
 with
 as
 but
 for

Top Negative
 except
 pinocchio
 seagal
 stale
 collection
 supposed
 unfunny
 uninspired
 vapid
 benigni
 boring
 called
 ballistic
 clumsy

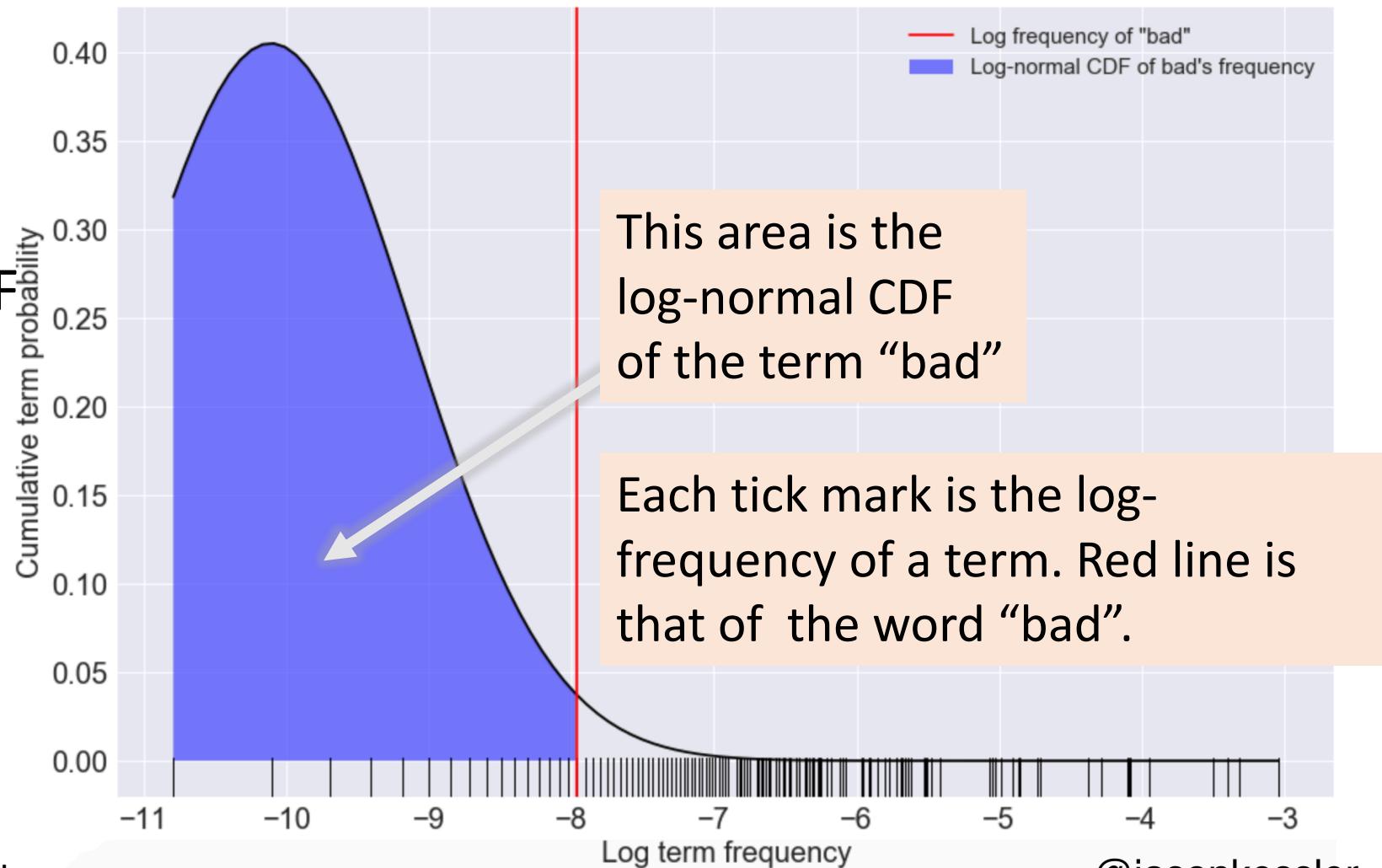
Naïve approach





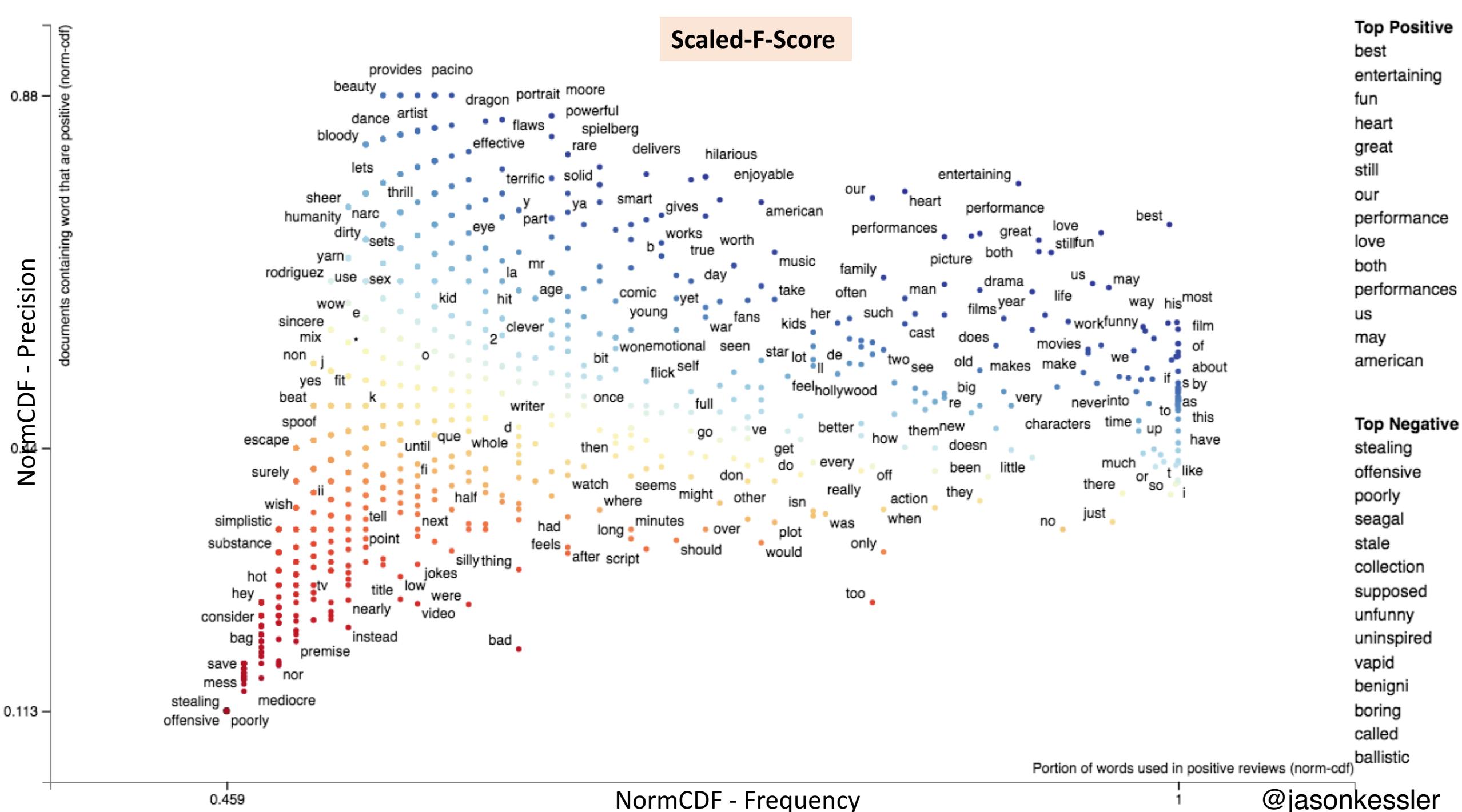
Fix: Normalize Precision and Frequency

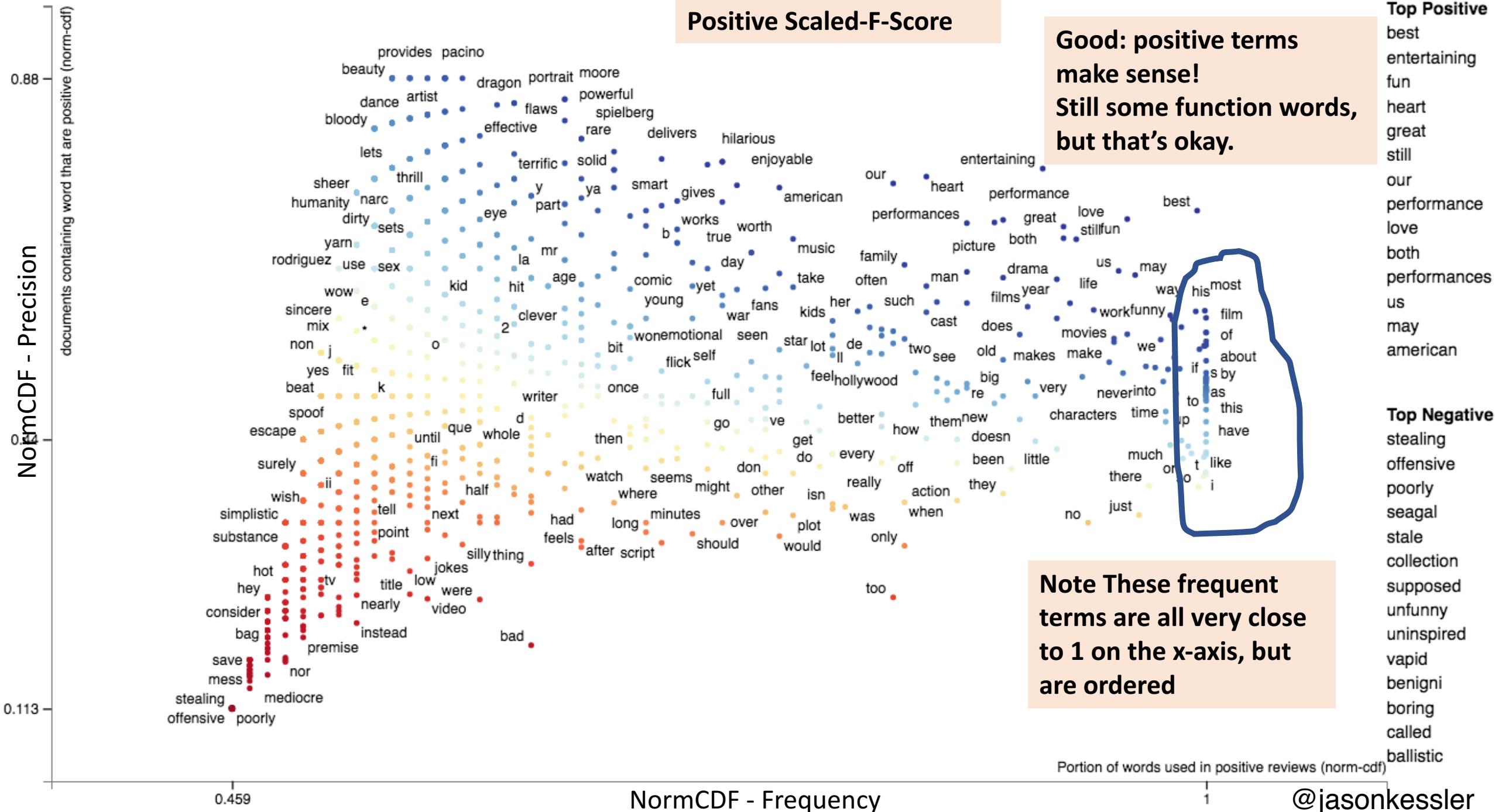
- Task: make precision and frequency similarly distributed
- How: take normal CDF of each term's precision and frequency
- Mean and std. computed from data
- Right: log normal CDF

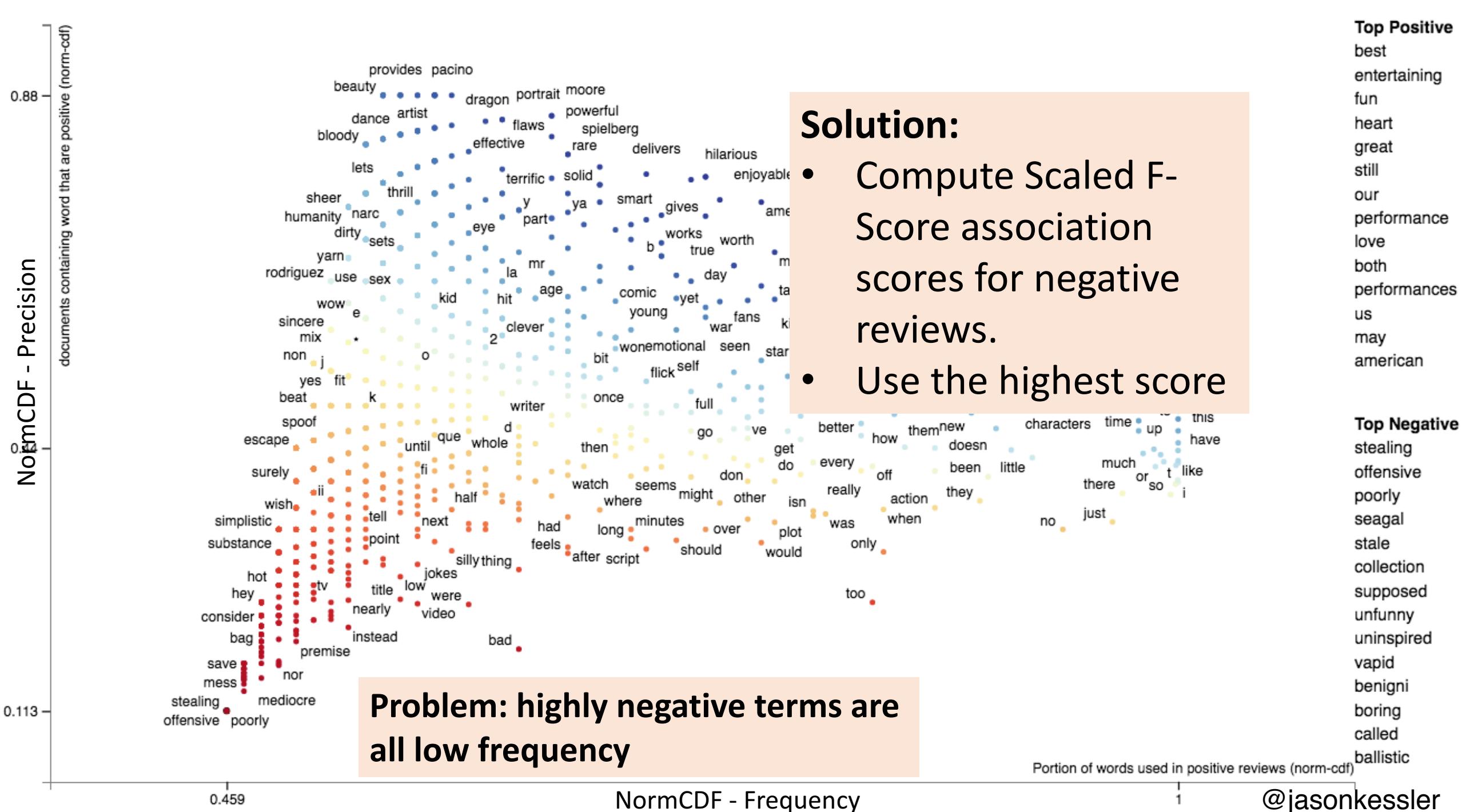


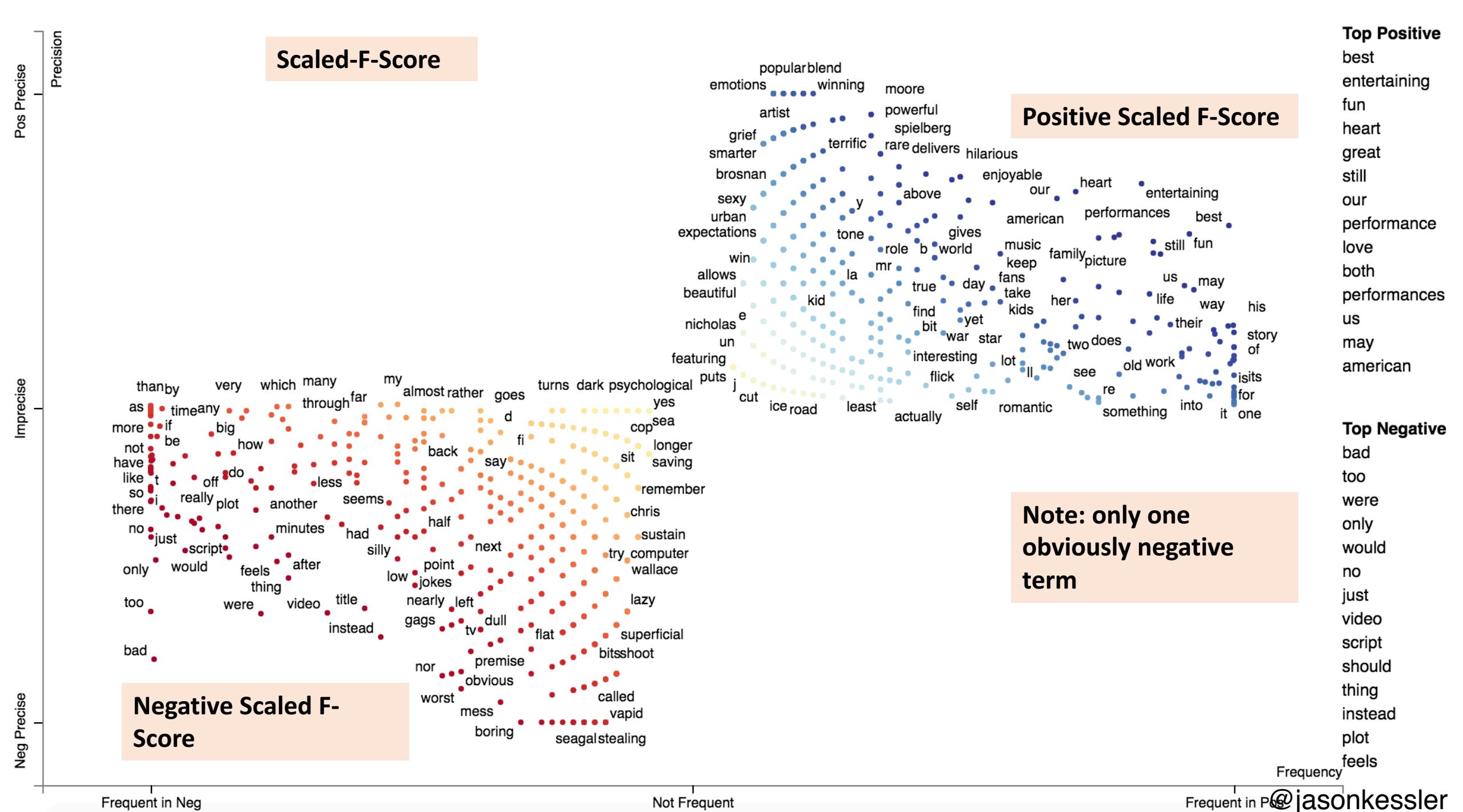
*log-normal CDF isn't used in these charts

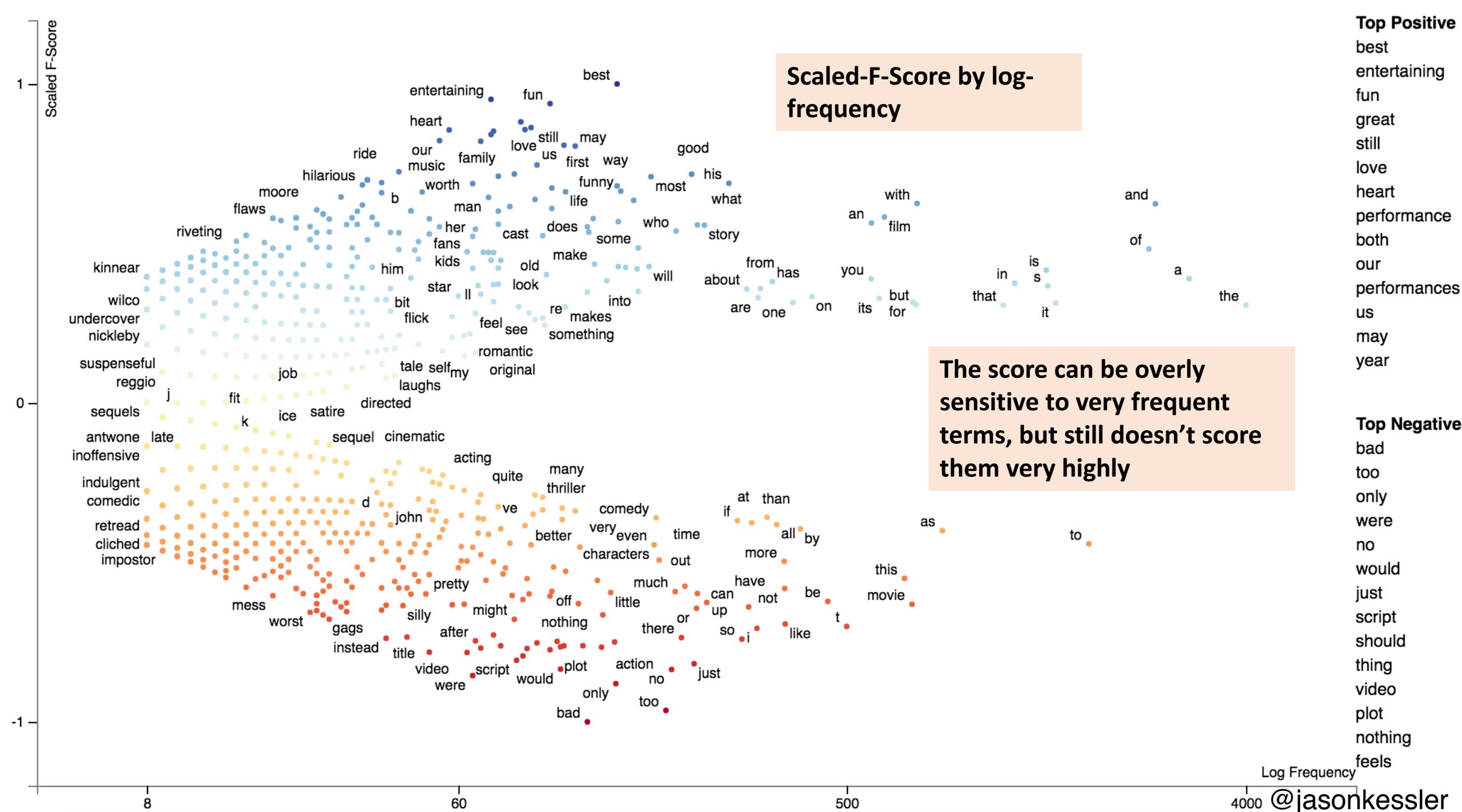
@jasonkessler



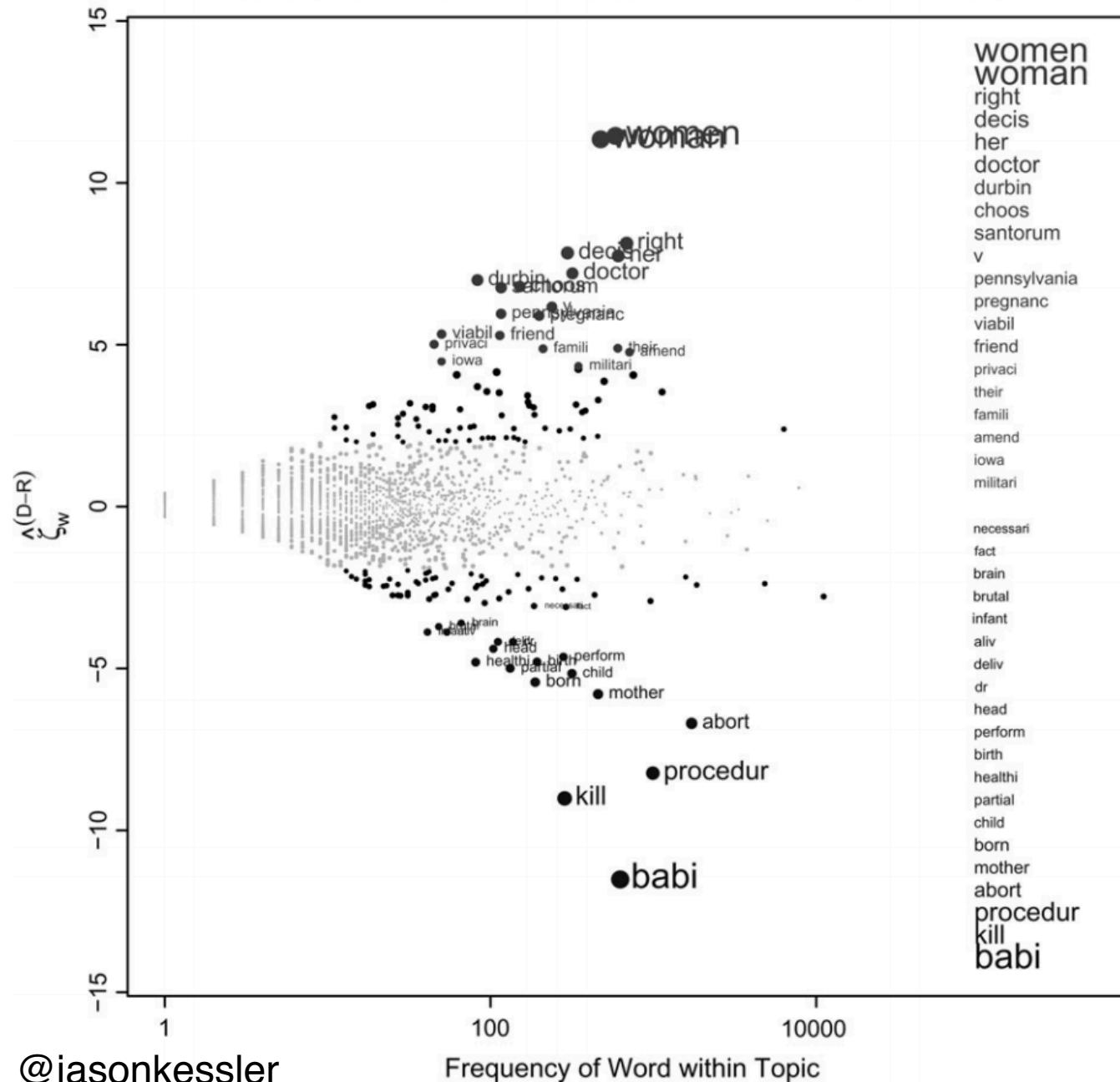








Partisan Words, 106th Congress, Abortion
(Weighted Log-Odds-Ratio, Informative Dirichlet Prior)



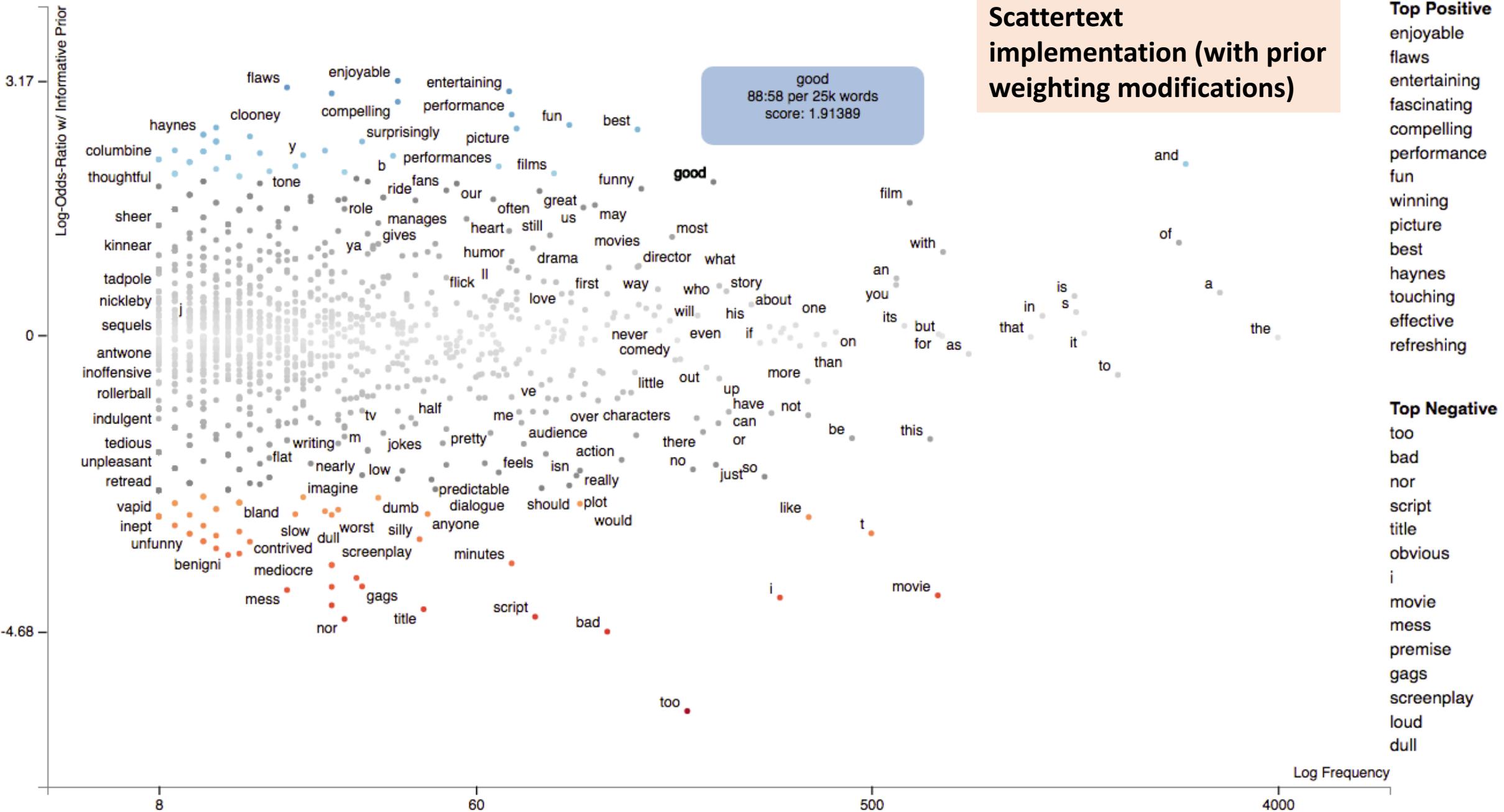
Monroe et. al (2009) approach

- Bayesian approach to term-association
- **Likelihood:** Z-score of log-odds-ratio
- **Prior:** Term frequency in a background corpus
- **Posterior:** Z-score of log-odds-ratio with background counts as smoothing values

Popular, but much more tweaking to get to work than Scaled F Score.

Burt Monroe, Michael Colaresi and Kevin Quinn. Fightin' words: Lexical feature selection and evaluation for identifying the content of political conflict. Political Analysis. 2008.

Scattertext implementation (with prior weighting modifications)



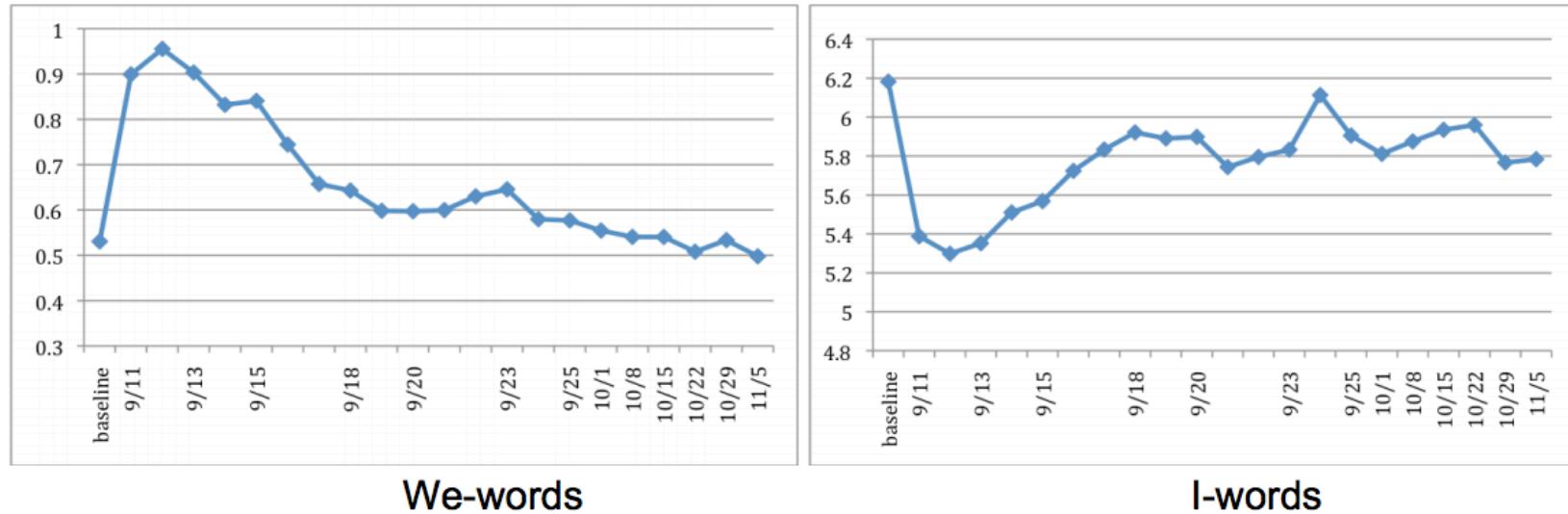
Scattertext reimplementations of Monroe et al. See

<http://nbviewer.jupyter.org/github/JasonKessler/PuPPyTalk/blob/master/notebooks/Class-Association-Scores.ipynb> for code.

@jasonkessler

In defense of stop words

Figure 1. Pronoun Use by Bloggers Before and After September 11, 2001



Note. Graphs reflect percentage of we-words (left) and I-words (right) within daily blog entries of 1,084 bloggers in the two months surrounding September 11, 2001.



Cindy K. Chung and James W. Pennebaker. Counting Little Words in Big Data: The Psychology of Communities, Culture, and History. EASP. 2012

In times of shared crisis, “we” use increases, while “I” use decreases.

I/we: age, social integration

I: lying, social rank

Function words and gender

LIWC Dimension Bold: entirely stop words	Effect Size (Cohen's d) (>0 F, <0 M) MANOVA p<.001
All Pronouns (esp. 3rd person)	0.36
Present tense verbs (walk, is, be)	0.18
Feeling (touch, hold, feel)	0.17
Certainty (always, never)	0.14
Word count	NS
Numbers	-0.15
Prepositions	-0.17
Words >6 letters	-0.24
Swear words	-0.22
Articles	-0.24

- Performed on a variety of language categories, including speech.
- Other studies have found that function words are the best predictors of gender.

Newman, ML; Groom, CJ; Handelman LD, Pennebaker, JW. Gender Differences in Language Use: An Analysis of 14,000 Text Samples. 2008.

Function word usage is counter-intuitive

- Pennebaker et al:
 - Testosterone levels (in two therapeutic settings) predict:
 - Modest but significant decreases in:
 - Pronouns referring to others (ex: we, she, they)
 - Communication verbs (ex: hear, say)
 - Modest but significant decreases in:
 - Optimism words (ex: energy, upbeat)
 - Negative (but statistically insignificant) correlation with subject's beliefs about testosterone and category usage!
 - Does not entirely explain gender differences.
- Herring et al:
 - Subjects tasked with impersonating opposite gender in online game
 - Discussed stereotypical topics (cars, shopping) but **didn't change stylistic cues**

- James W. Pennebaker, Carla J. Groom, Daniel Loew, James M. Dabbs. Testosterone as a Social Inhibitor: Two Case Studies of the Effect of Testosterone Treatment on Language. 2004.
- Susan C. Herring, Anna Martinson. Assessing Gender Authenticity in Computer-Mediated Language Use: Evidence From an Identity Game. Journal of Language and Social Psychology. 2004.

Clickbait: what works?

The New York Times



2:04
The Woman Who Brought
Down Bill Cosby

Video by NEETI UPADHYE. Photo by Mark Makela/Getty Images

Critic's Notebook

Cliff Huxtable Was Bill Cosby's Sickest Joke

By WESLEY MORRIS 8:58 PM ET
Mr. Cosby's signature TV character was a patient, wise father figure. In other words, our critic writes, the ideal cover for terrible behavior.

- Did the #MeToo Movement Affect the Cosby Verdict?

Cosby Has a Few Words for the Court

By LIAM STACK 10:20 PM ET
Mr. Cosby has not said much about the sexual assault allegations against him, but on Thursday he erupted with a vulgarity aimed at the prosecutor.

- Did the #MeToo Movement Affect the Cosby Verdict?

North Korea's Phony Peace Ploy

By NICHOLAS EBERSTADT

If the past is any guide, Pyongyang will offer Seoul unenforceable verbiage at this week's summit meeting.

The Supreme Court and the New Civil War

By LINDA GREENHOUSE

The battle between the White House and blue states raises questions about the limits of federal authority.

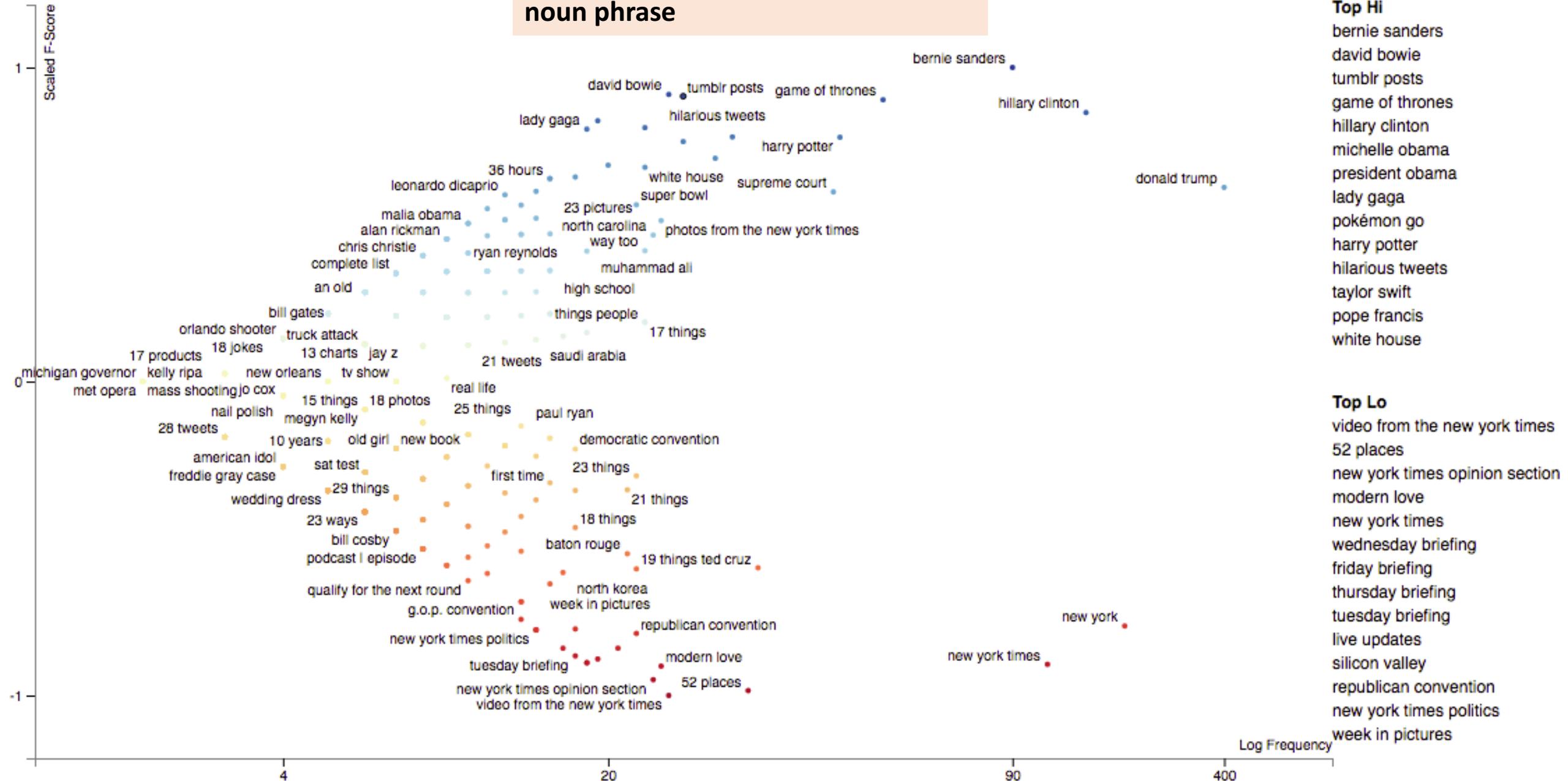
- When Misogynists Become Terrorists
- Krugman: Trump's War on the Poor
- Bruni: Ronny Jackson, From Fawning to Fiasco



Clickbait corpus

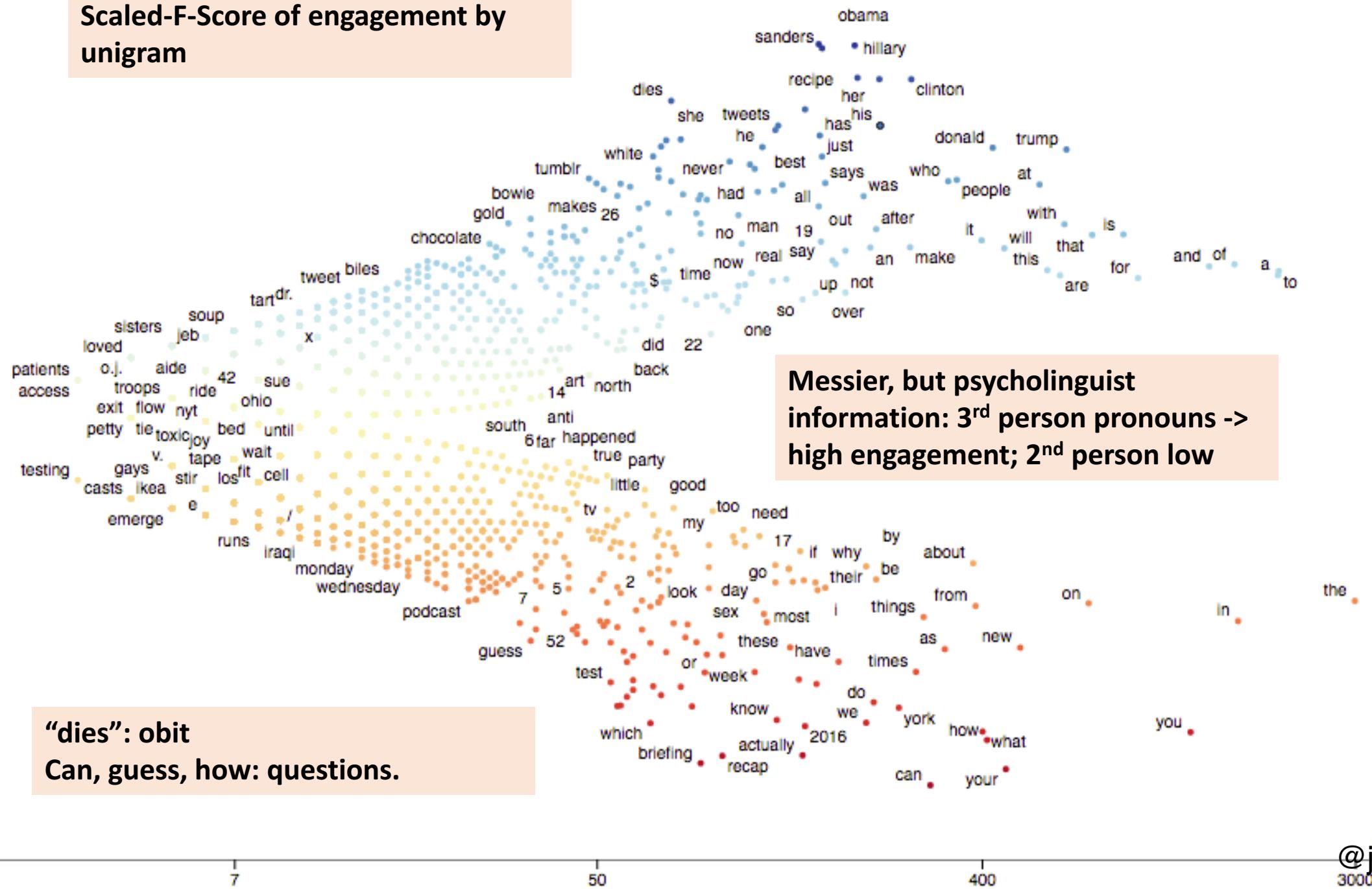
- Facebook posts from BuzzFeed, NY Times, etc/ from 2010s.
- Includes headline and the number of Facebook likes
- Scrapped by researcher Max Woolf at github.com/minimaxir/clickbait-cluster.
- We'll separate articles from 2016 into the upper third and lower third of likes.
- Identify words and phrases that predict likes.
- Begin with noun phrases identified from Phrase Machine (Handler et al. 2016)
- Filter out redundant NPs.

Scaled-F-Score of engagement by noun phrase



@jasonkessler

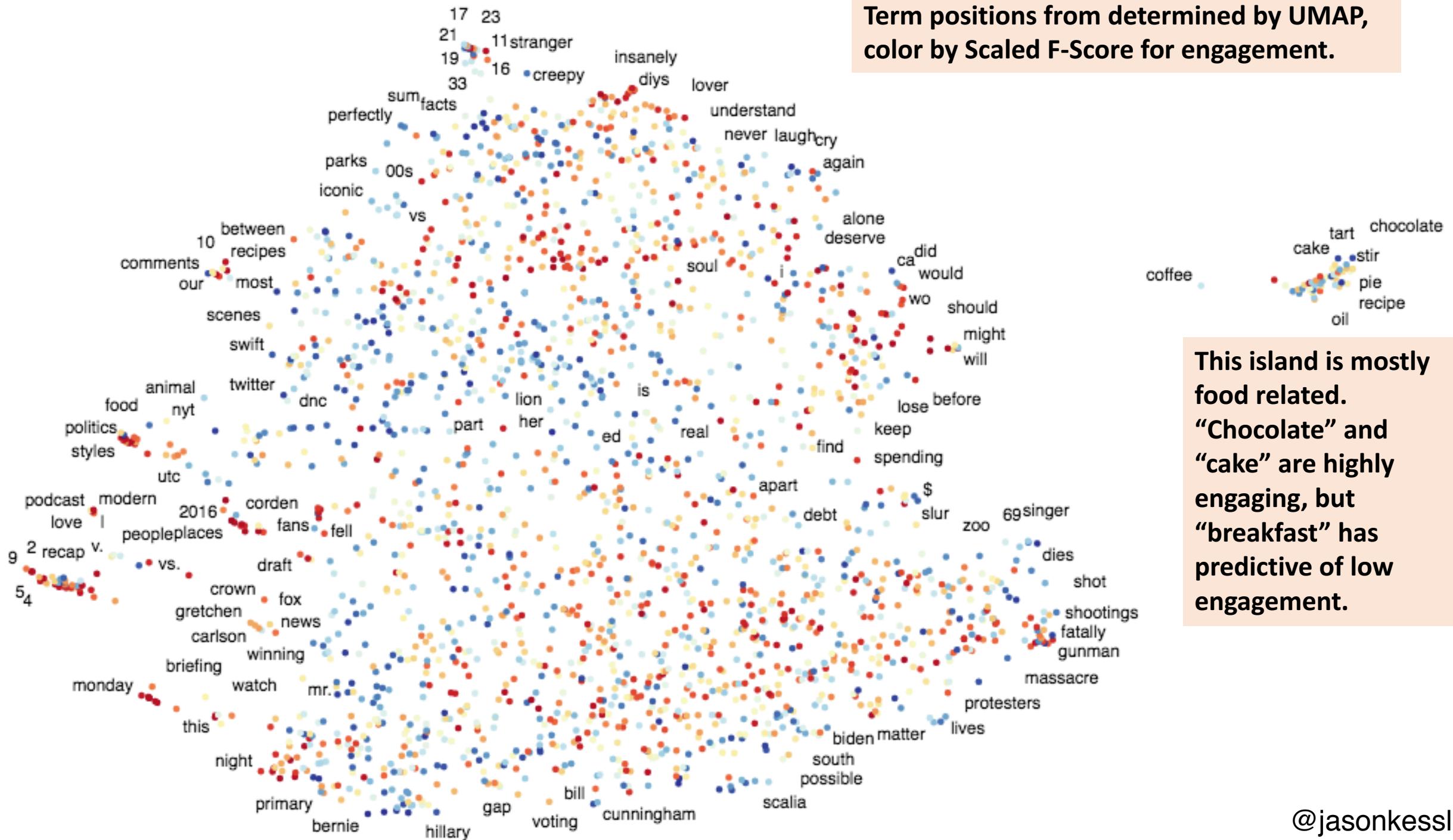
Scaled-F-Score of engagement by unigram



Clickbait corpus

- How do terms with similar meanings differ in terms of their engagement rates?
- Use Gensim (<https://radimrehurek.com/gensim/>) to find word embeddings
- Use UMAP (McInnes and Healy 2018) to project them into two dimensions, and explore them with Scattertext.
 - Locally groups words with similar embeddings together.
 - Better alternative to T-SNE; allows for cosine instead of Euclidean distance criteria

**Term positions from determined by UMAP,
color by Scaled F-Score for engagement.**



This island is mostly food related.
“Chocolate” and “cake” are highly engaging, but “breakfast” has predictive of low engagement.

Clickbait corpus

- How do the Times and Buzzfeed differ in what they talk about, and their content engages their readers?
- Scattertext can easily create visualizations to help answer these questions.
- First, we'll look at how what engages for Buzzfeed contrasts with what engages for the Times, and vice versa

Appeals to all

her, obama, bernie, sanders, women, bernie sanders, president, dies at, has, of thrones

NY Times: ↑ Engagement

bernie sanders, bernie,
dies, sanders, dies at,
obama, women, white,
wins, president

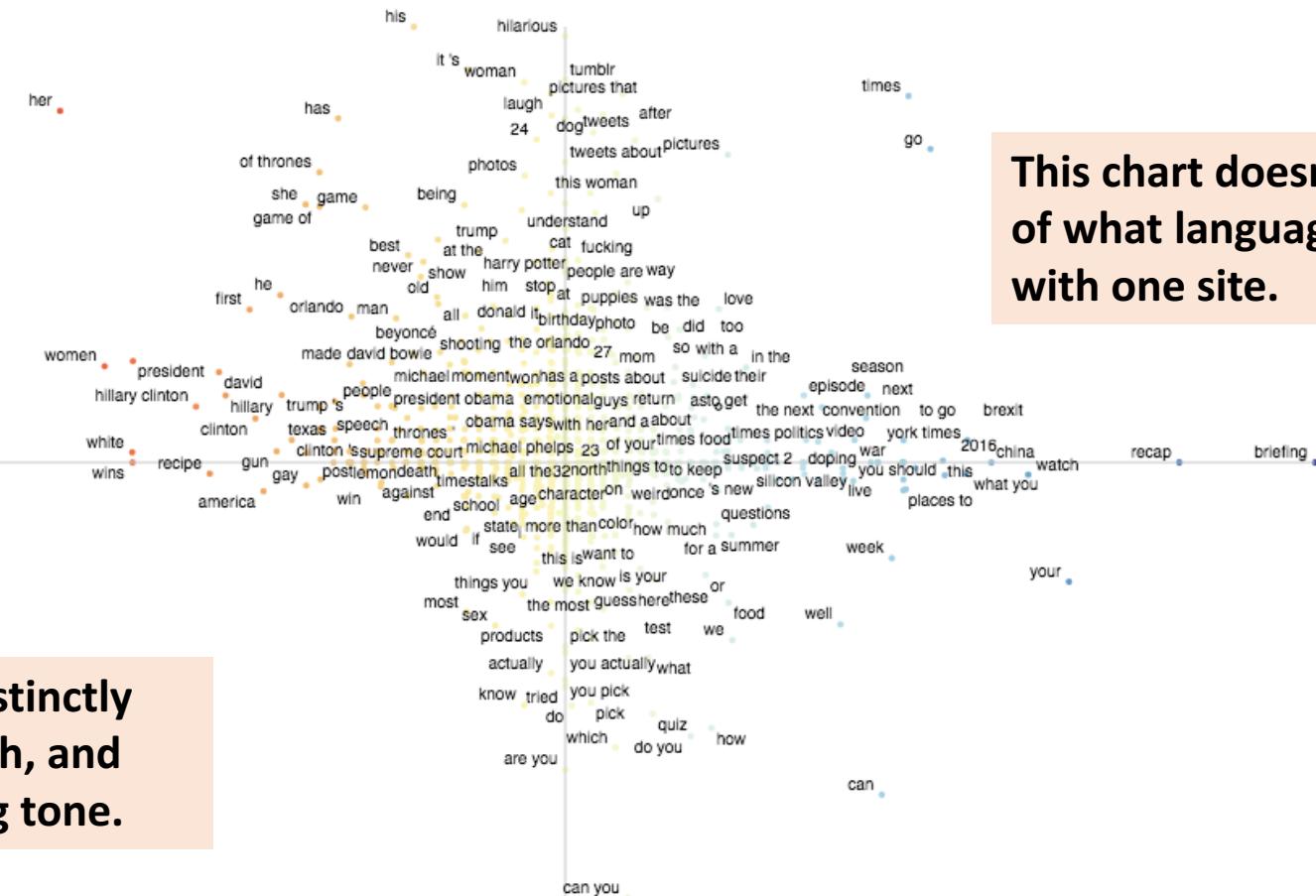
Oddly, NY Times readers distinctly like articles about sex, death, and which are written in a smug tone.

Appeals to elite, ignored by masses

know, actually, sex, most, products, the most, how well, things you, dies, we know

BuzzFeed: ↑ Engagement

hilarious, tumblr, woman, laugh, tweets, after, dog, it's, pictures that, 26



BuzzFeed: ↓ Engagement

Ignored by elite, appeals to masses

go, times, brexit, this, season, to go, pictures, next, war, the new

This chart doesn't give a good sense of what language is more associated with one site.

NY Times: ↓ Engagement

briefing, recap, watch,
china, brexit, what you,
your, 2016, this, watch this

Ignored by all
can, can you, how, do you,
your, well, which, quiz, are
you, we

Highbrow Engagement

sanders, obama, bernie,
bernie sanders, hillary
clinton, hillary, recipe,
clinton, trump 's, dies

bernie
sanders

NYTimes

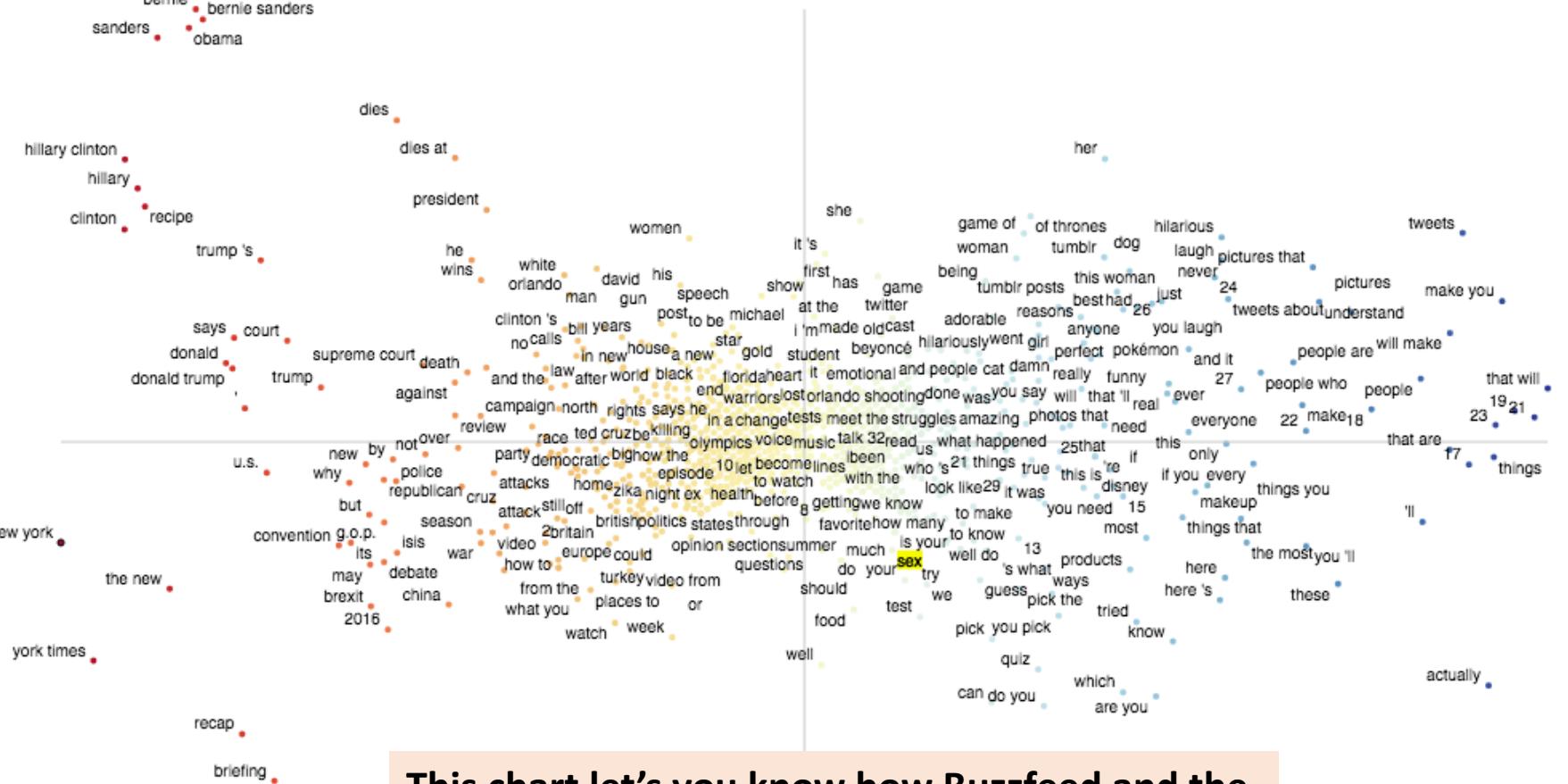
new york, ', donald, donald
trump, u.s., says, the new,
court, clinton, trump

Highbrow Ignored

briefing, recap, york times,
the new, 2016, new york,
brexit, may, china,
convention

Many Facebook Reactions

she, it 's, women, has, first, game, of thrones, game of, her, at the



This chart let's you know how Buzzfeed and the Times are distinct, while still distinguishing engaging content,

Few Facebook Reactions

can, well, do you, week, quiz, pick, food, test, which, watch

Lowbrow Engagement

tweets, make you, pictures
that, will make, pictures,
understand, hilarious,
laugh, that will, 24

BuzzFeed

21, 19, 23, things, that will,
17, that are, people, will
make, 'll

Lowbrow Ignored

can you, actually, are you,
these, which, know, here 's,
here, 'll, do you

Thank you! Questions?

Jason S. Kessler
Global AI Conference
April 27, 2018

<https://github.com/JasonKessler/GlobalAI2018>