



Analyzing Lyrical Data from Historically Popular Songs

by: David Ash, Brian Blakely, and Jackson Conrad
even though we won't admit to it later



Goals

- I. Web scrape song data from Billboards
 - Analyze lyrical data from the top 100 songs from 1959 - 2019
3. Use natural language processing (NLP) to help with analysis
 - Make pretty graphs for stuff
- D. Design a good presentation for when 2.5/4 of those things fail horribly



Process

- 0) Modify a web scraper to get the data set
 1. Wait 6 hours
 2. Now that you finally have your data set
 3. Start analysis
 4. Work for 3 hours processing the data
 5. Realize your data set has full books instead of lyrics in some cases
 6. Try to modify your data set to work correctly and not have books in it
 7. Run a NLP to reprocess the data
 8. Wait 2.5 hours
 9. Watch in despair as it throws an error on the last song
10. Give up and start a new project
11. Wait 2.5 hours
12. Get locked outside
13. Give up on doing the new project and go back to the old project which works now



Summary: Progress made during the first 10 hours:

[this slide intentionally left blank]

After that

The very first data we successfully processed:

Data is from approximately

bottom 50 of 2000s - 2019




Idk why there are so many single quotes either



cont.

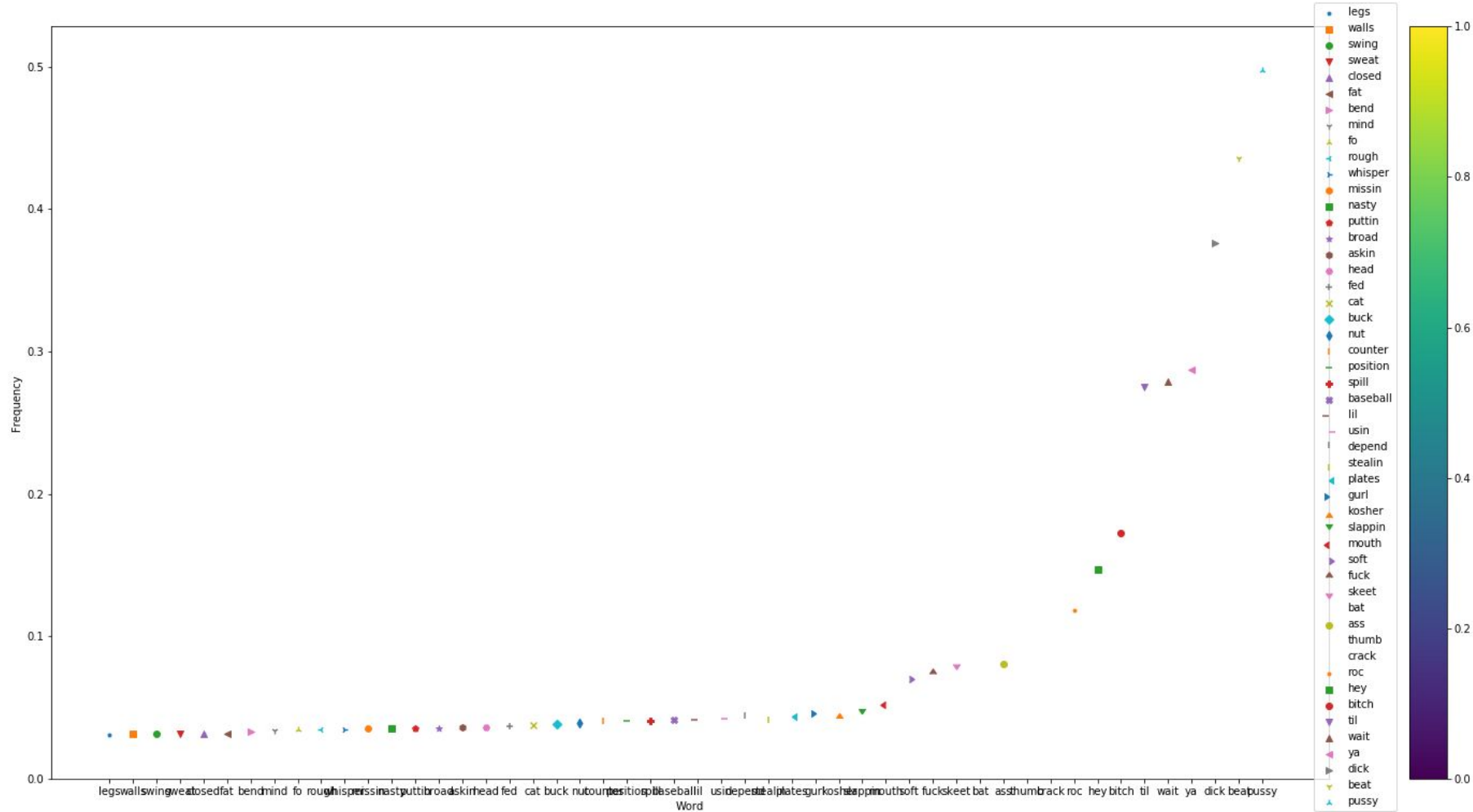
14. Realize it only works because the frequency measure gives arbitrary values for our purposes
15. Skew the data by only looking at the graphs that look nice
16. Make this presentation
17. Give the presentation
19. ???
20. Win 1st place
21. Skip class tomorrow because I'm tired and everything hurts



Actual data (not a joke this time) (last time
wasn't a joke either fyi)

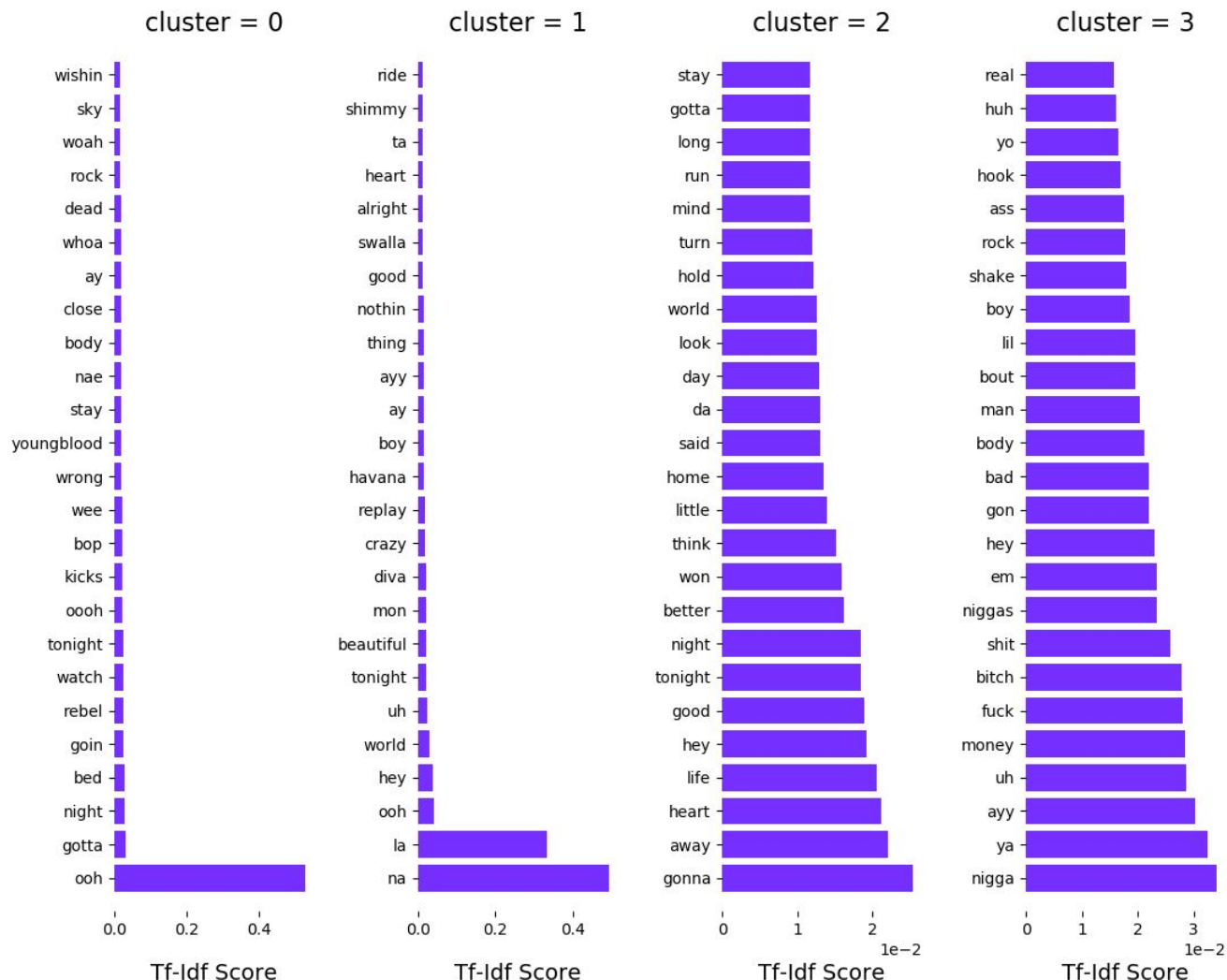
please excuse the profanity

we didn't want to skew the data

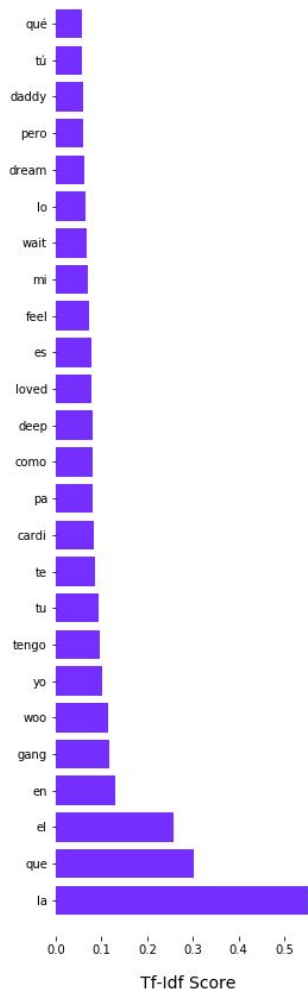


2005 - 2019

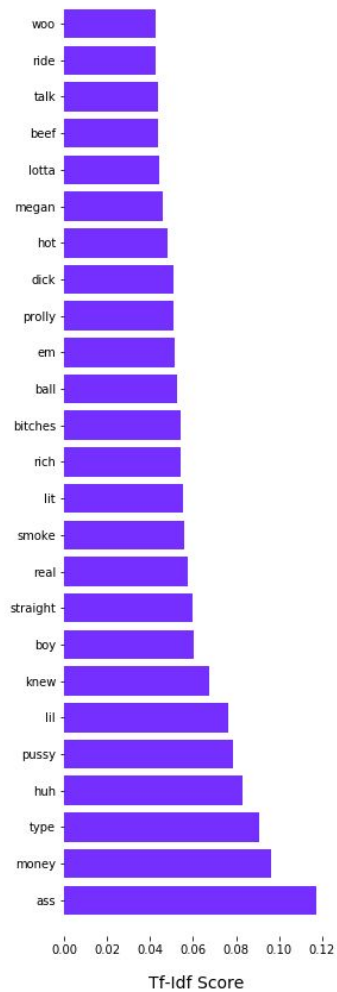
(clustering algorithm
grouped bad words
with each other)



cluster = 0



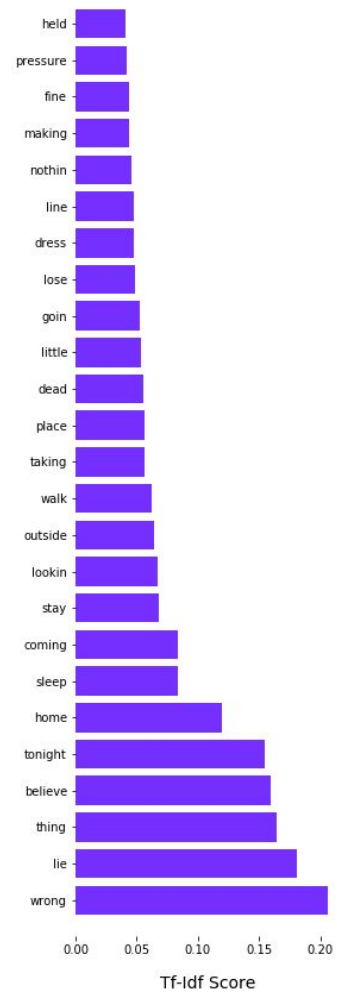
cluster = 1



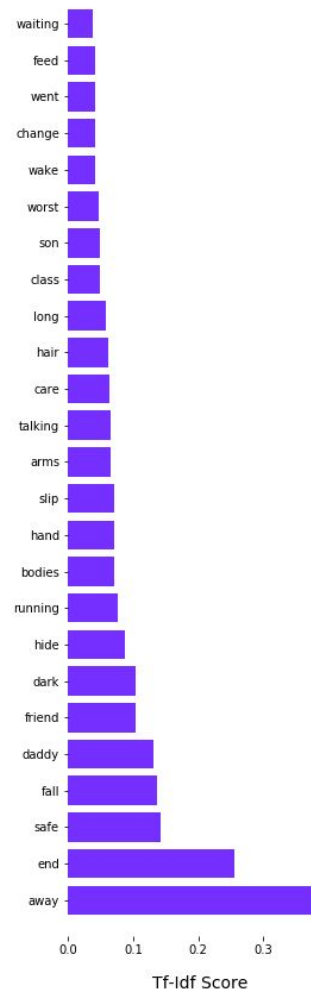
cluster = 2



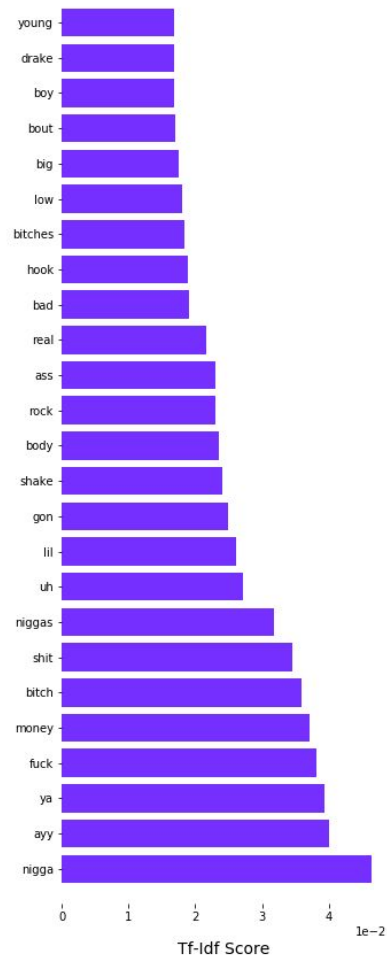
cluster = 3



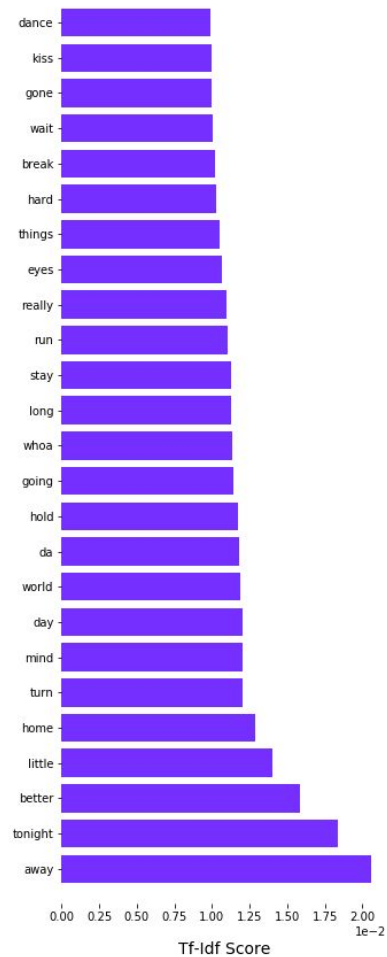
cluster = 4



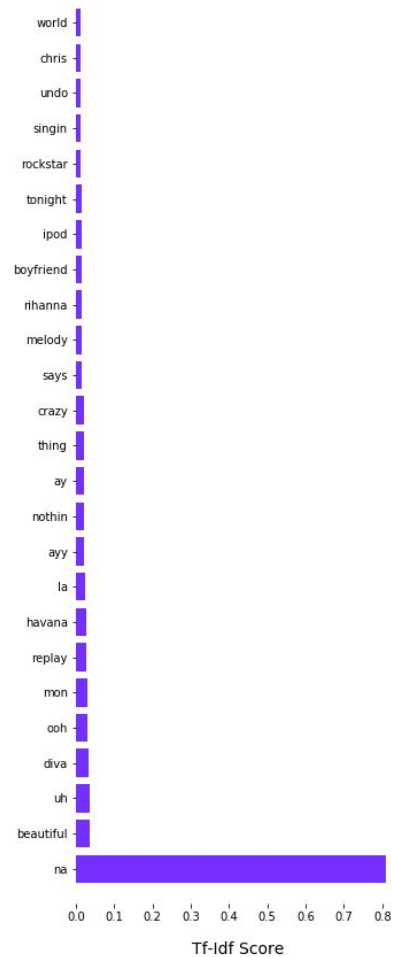
cluster = 0



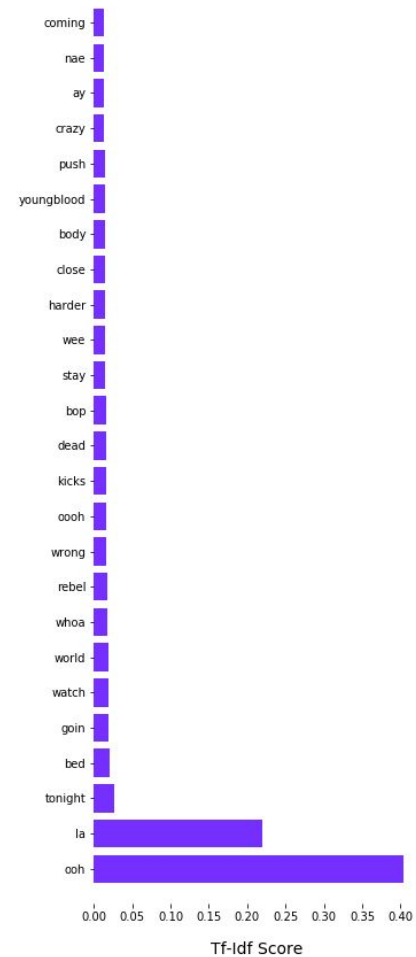
cluster = 1



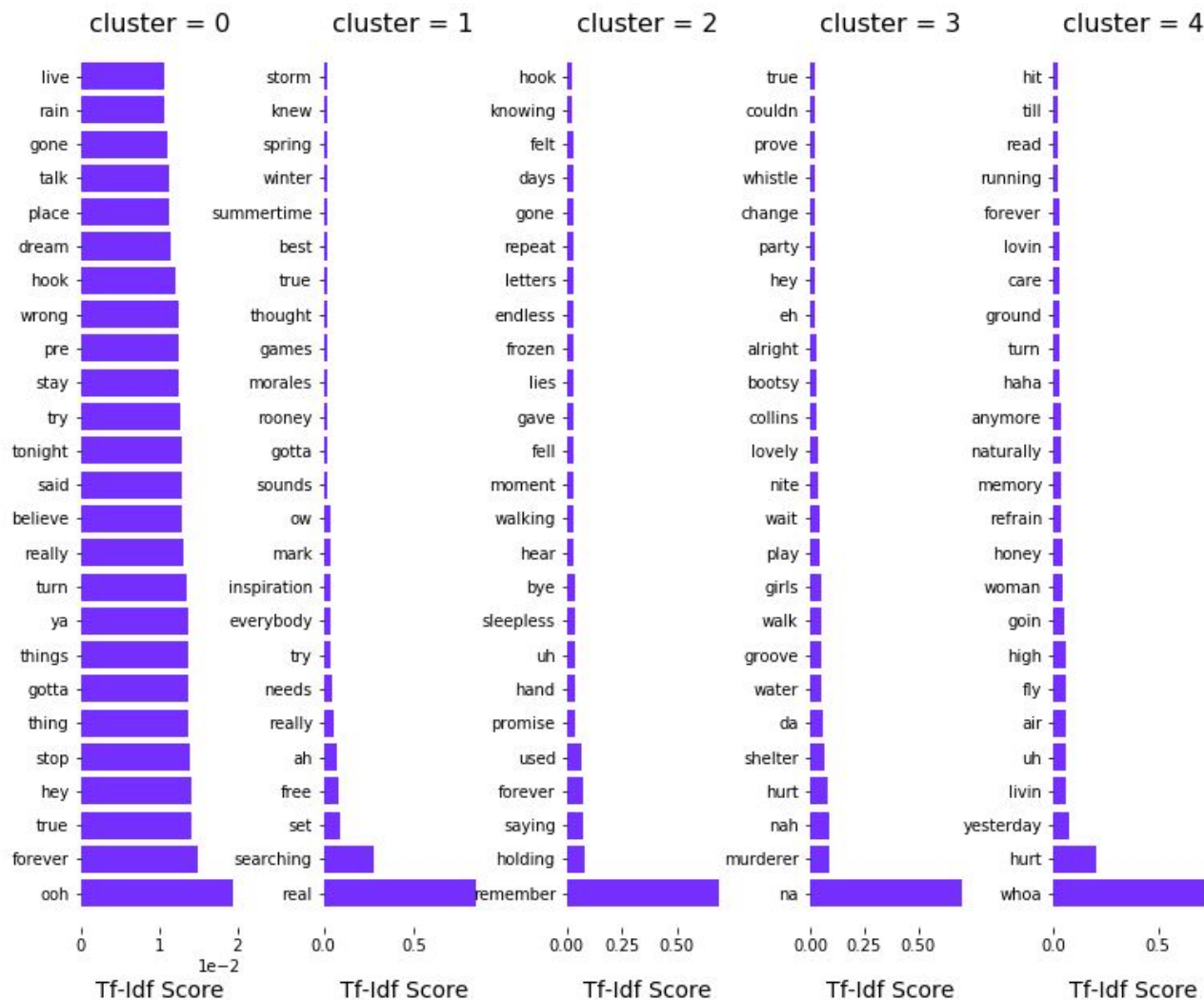
cluster = 2



cluster = 3



Mid 1980s



EVERLASTING' UNDERSTAND' FINISHED' HEART' CHARMING'
HOPED' WAIT' FILLED' FAR' FILES' KISS' MAN'
COLLISION' ROMANTIC' EARTH' JULIET' FIRED' HUG'
FIGURES' PRINCE' FINNA' MEAN' SWEAR' FOREVER' PROMISE'
TENDERNESS' CINDERELLA' PRECIOUS' FIRES' KISSED' DREAM' FELT' OPEN' FIRESIDE' FINALLY'
EYES' CHANCE' WAITING' WORLD' HOLDING' ROMEO' READ' COMES' MOVIES' TRUE' RESCUE'
FORGET' HEAVEN' DAY'

1990s - 2000s



1980s - 1990s



What you could do with this project

- ★ Collect data about the frequency of...
 - Parts of speech (nouns, verbs, etc.)
 - Different kinds of nouns (People, nationalities, products, etc.)
 - Bad words (I'm not typing them)
- ★ Divide the data by...
 - Time (specific years, decades, etc.)
 - Genre
 - Artist
 - Lyrical content (length, unique words, etc.)
- ★ Create cool graphs with the data collected
 - Stuff we wanted to do but couldn't do in time includes...
 - Collect information on the usage of product and company names over time
 - Display information about use of profanity over the decades
 - Display most popular verb by decade
 - How closely words follow Zipf's Law