

Plan:

1. Work through a text analysis of song lyrics

# Text Analysis: Example

Shannon E. Ellis, Ph.D  
UC San Diego



Department of Cognitive Science  
[sellis@ucsd.edu](mailto:sellis@ucsd.edu)

Today's example question: How has pop music changed in the last four years?

Goal: Understand the basics of sentiment analysis and TF-IDF

What data would we need to answer this question?

How has pop music changed in the last four years?

Data: Lyrics to the most popular songs from each year

# The data : Top songs from Feb music charts 2017-2020

2017: 152 songs

2018: 139 songs

2019: 127 songs

2020: 137 songs

Song data from **Spotify.**  
Lyrics from **genius.com**



# Questions we can ask...

1. Does the total number of words change over time?
2. Does uniqueness change over time?
3. Does the diversity or density change?
4. What words are most common?
5. What words are most unique to each year?
6. What sentiment do songs convey most frequently?
7. Has sentiment changed over time?
8. What are the sentiment of the #1 songs?
9. What words contribute to the sentiment of these #1 songs?
10. ...what about bigrams? N-grams?

# Questions we can ask...

1. Does the total number of words change over time?
2. Does uniqueness change over time?
3. Does the diversity or density change?

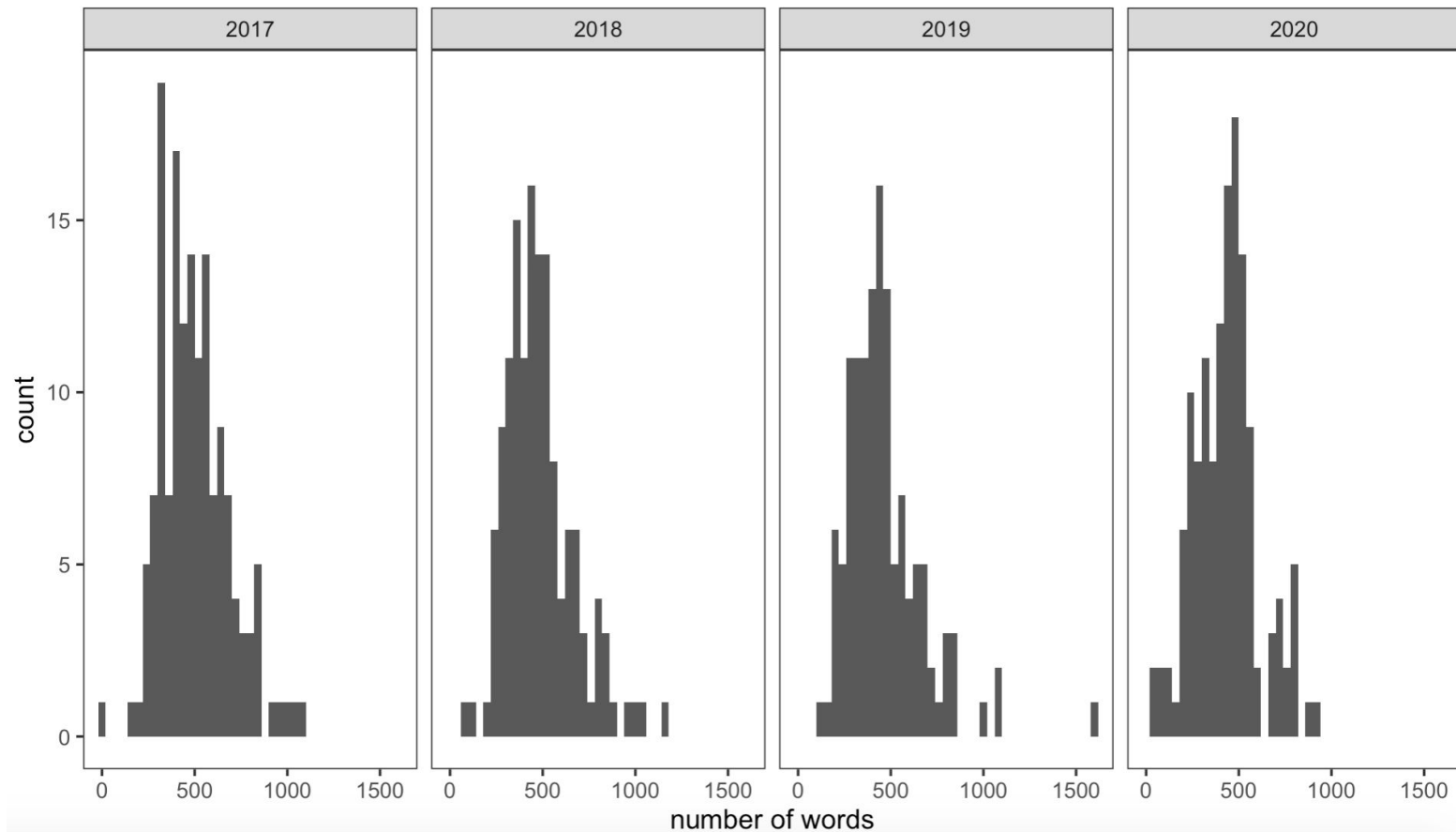
EDA

4. What words are most common?
5. What words are most unique to each year?

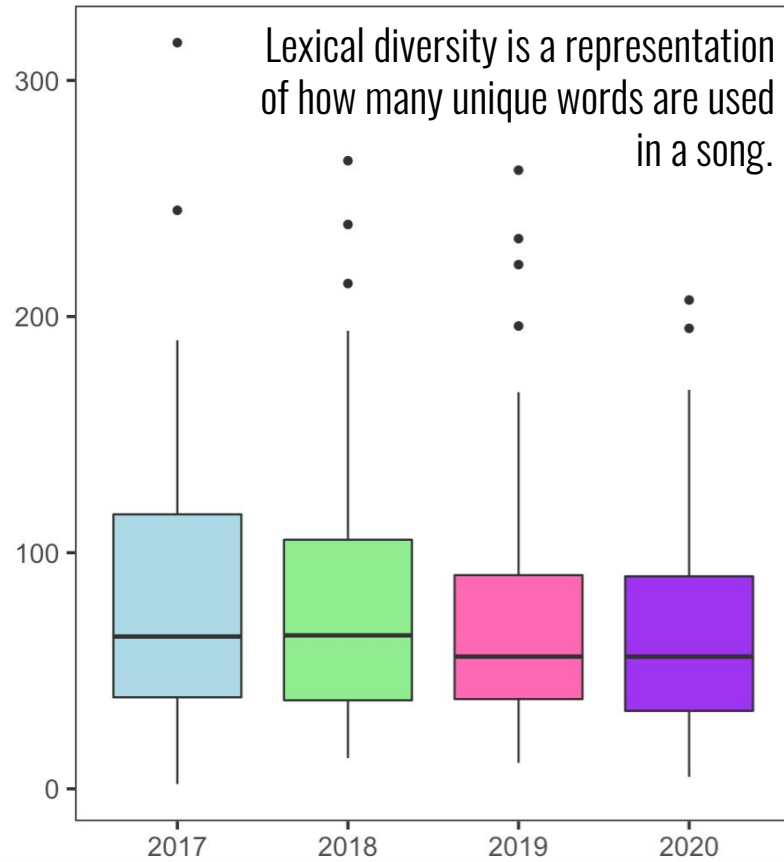
TF-IDF

6. What sentiment do songs convey most frequently?
7. Has sentiment changed over time?
8. What are the sentiment of the #1 songs?
9. What words contribute to the sentiment of these #1 songs?
10. ...what about bigrams? N-grams?

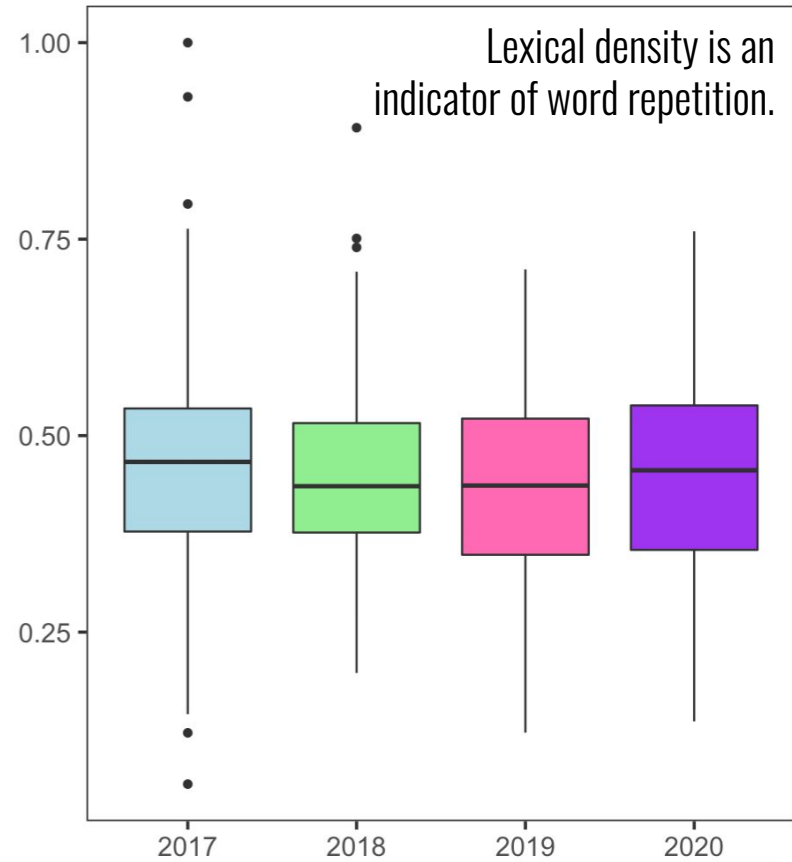
Sentiment  
Analysis



### Lexical Diversity



### Lexical Density

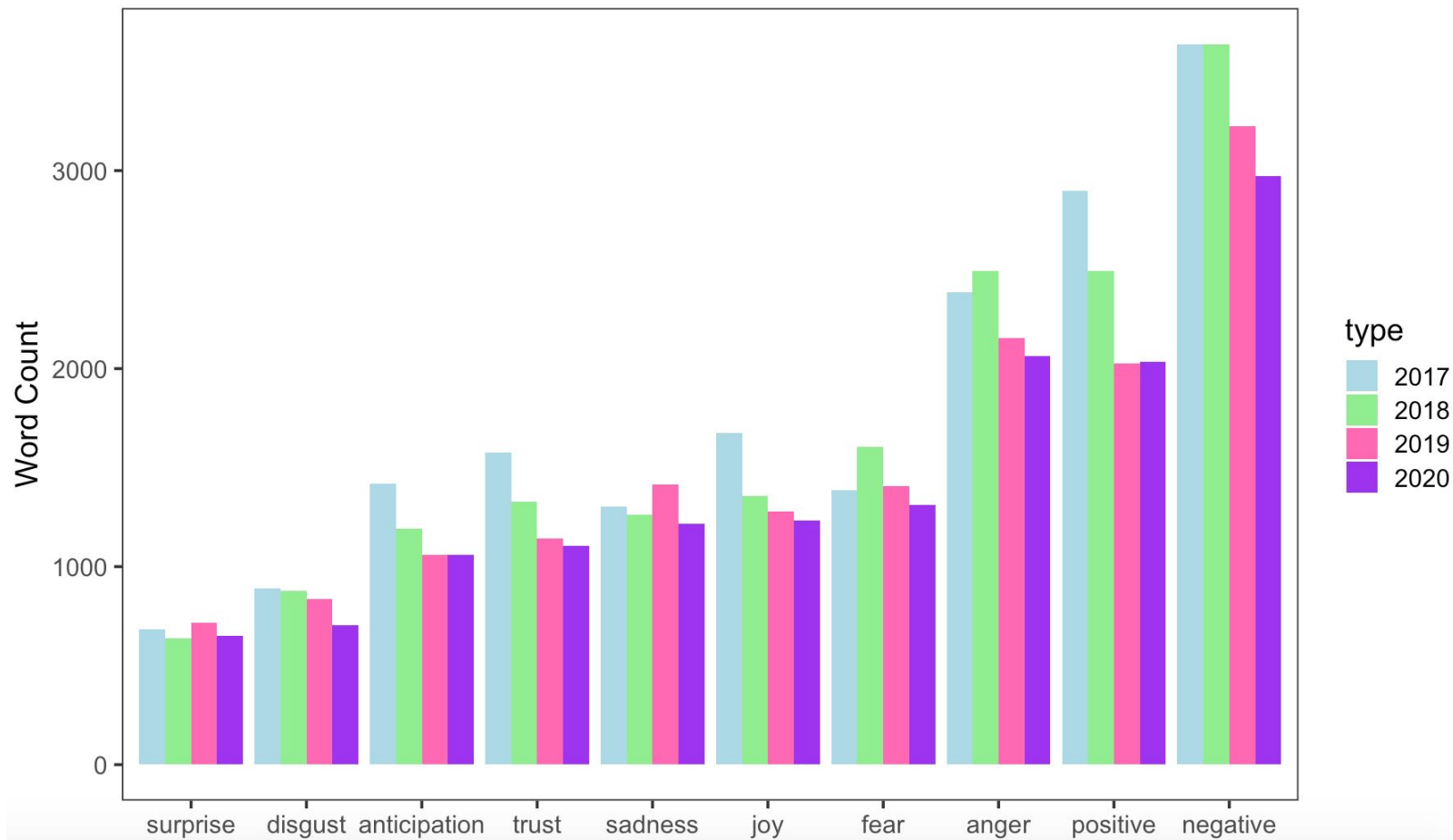




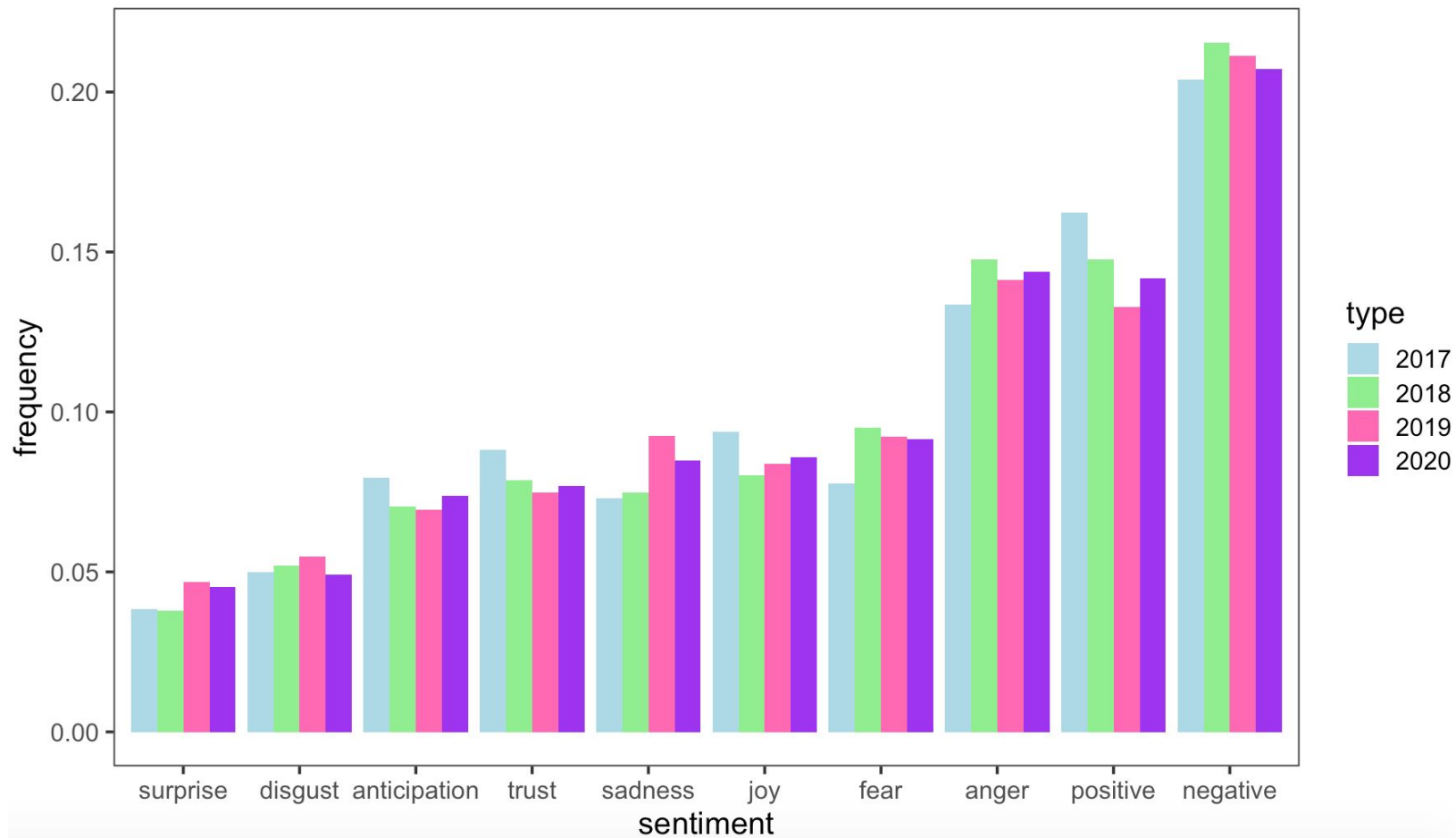
# Sentiment Analysis

---

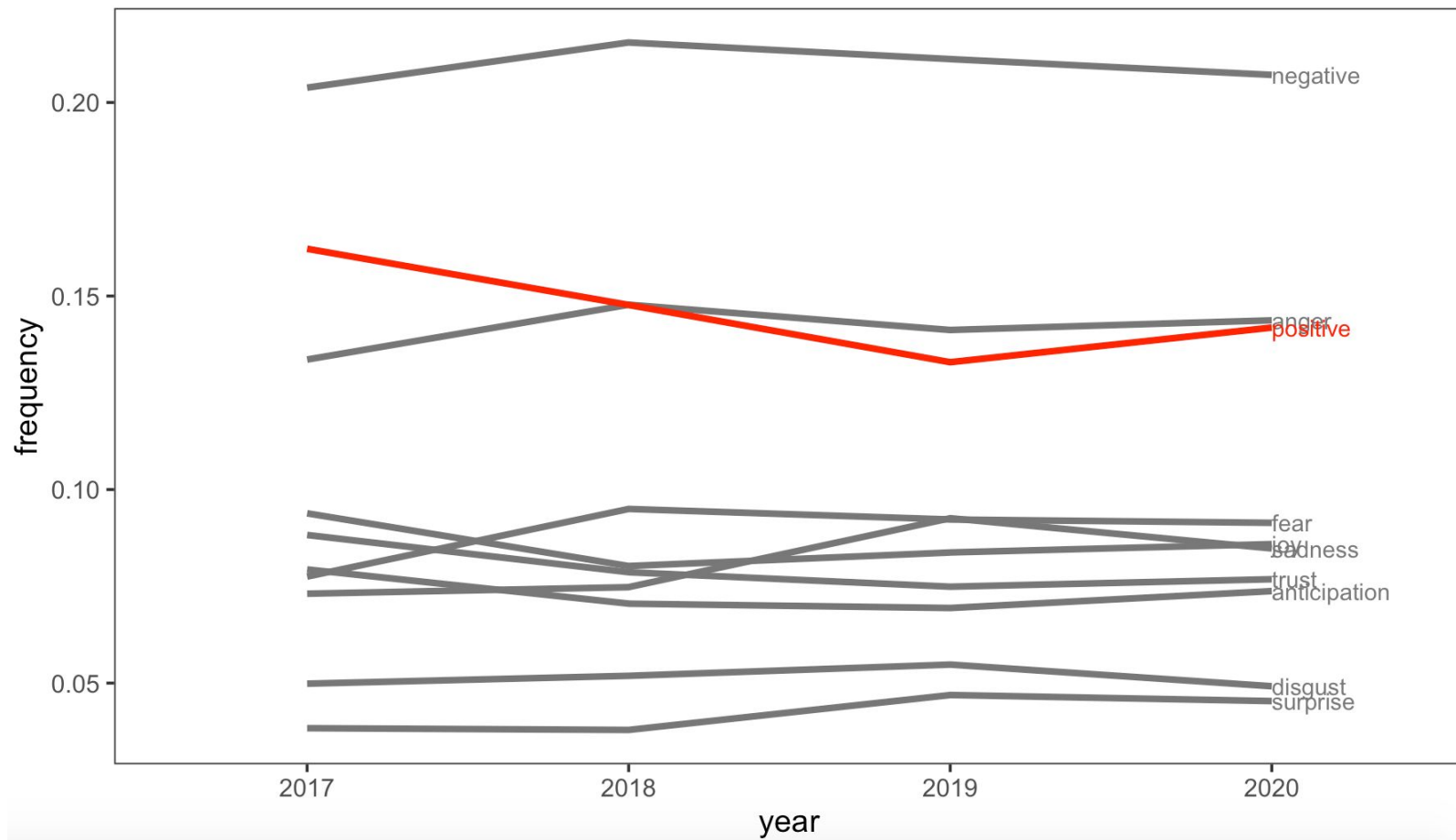
# Top Songs Sentiment



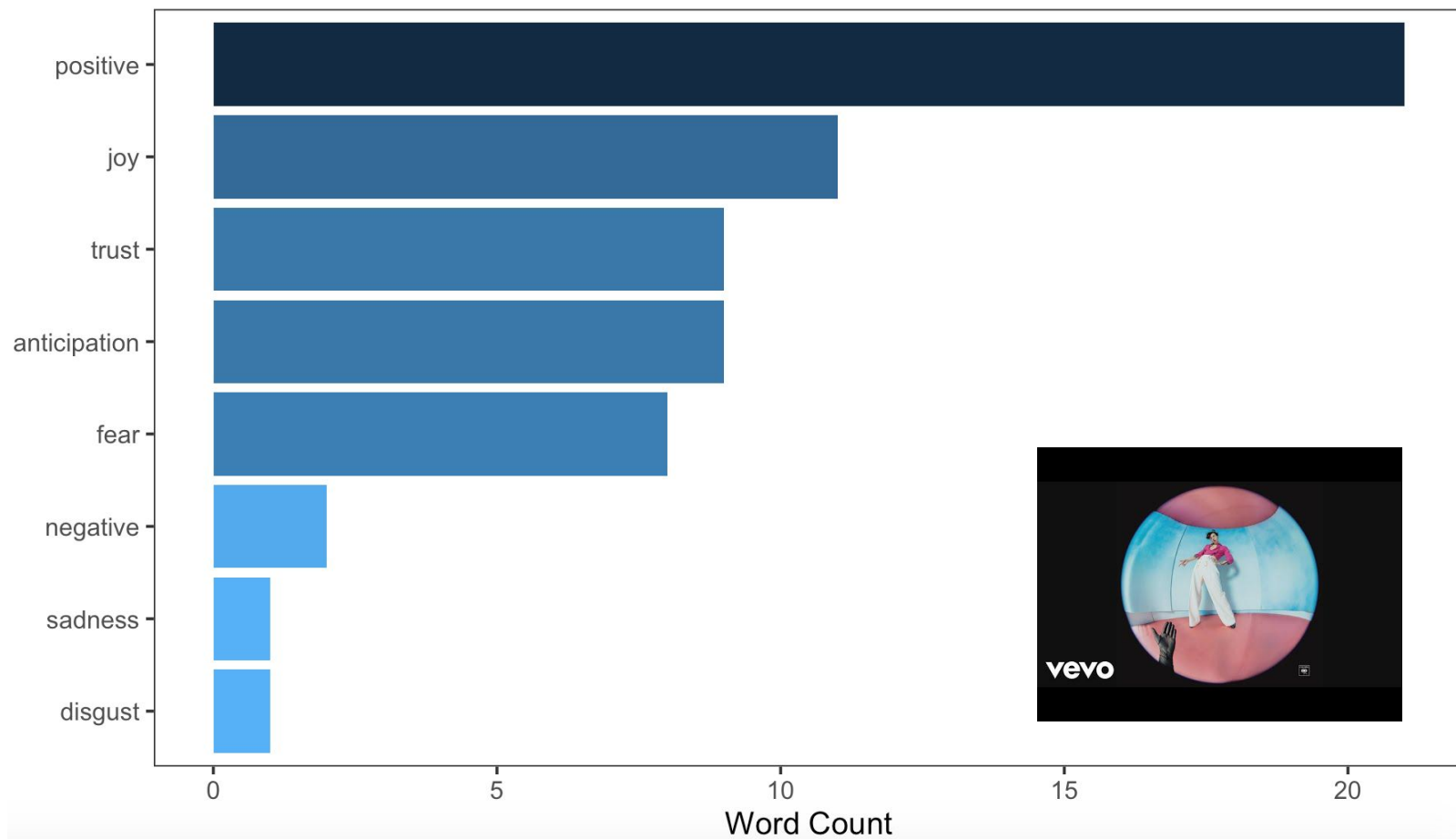
Sentiment by Year



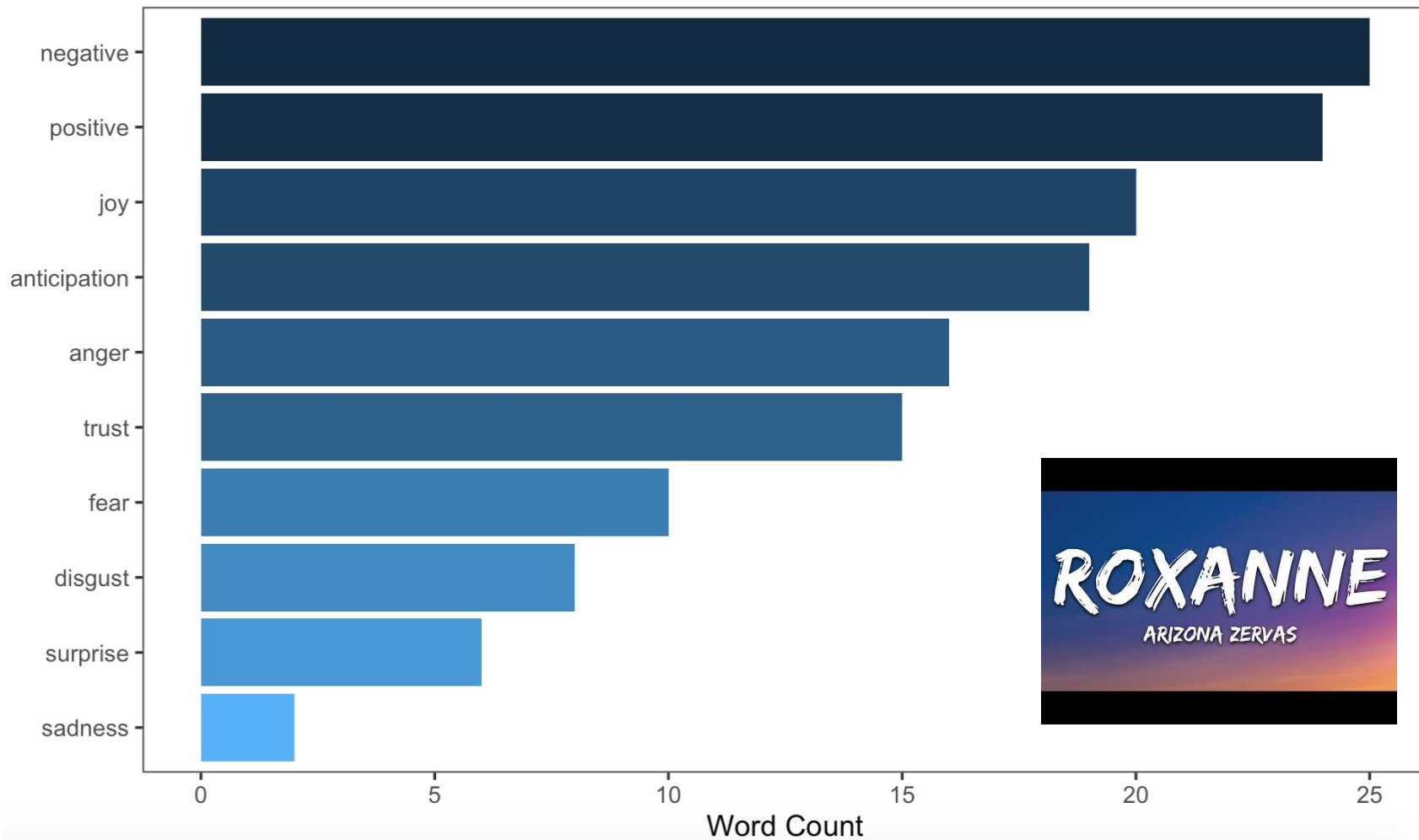
# Change in Sentiment over Time



## Sentiment: Adore You



## Sentiment: ROXANNE



**TF-IDF**

**Term Frequency - Inverse Document Frequency**

---

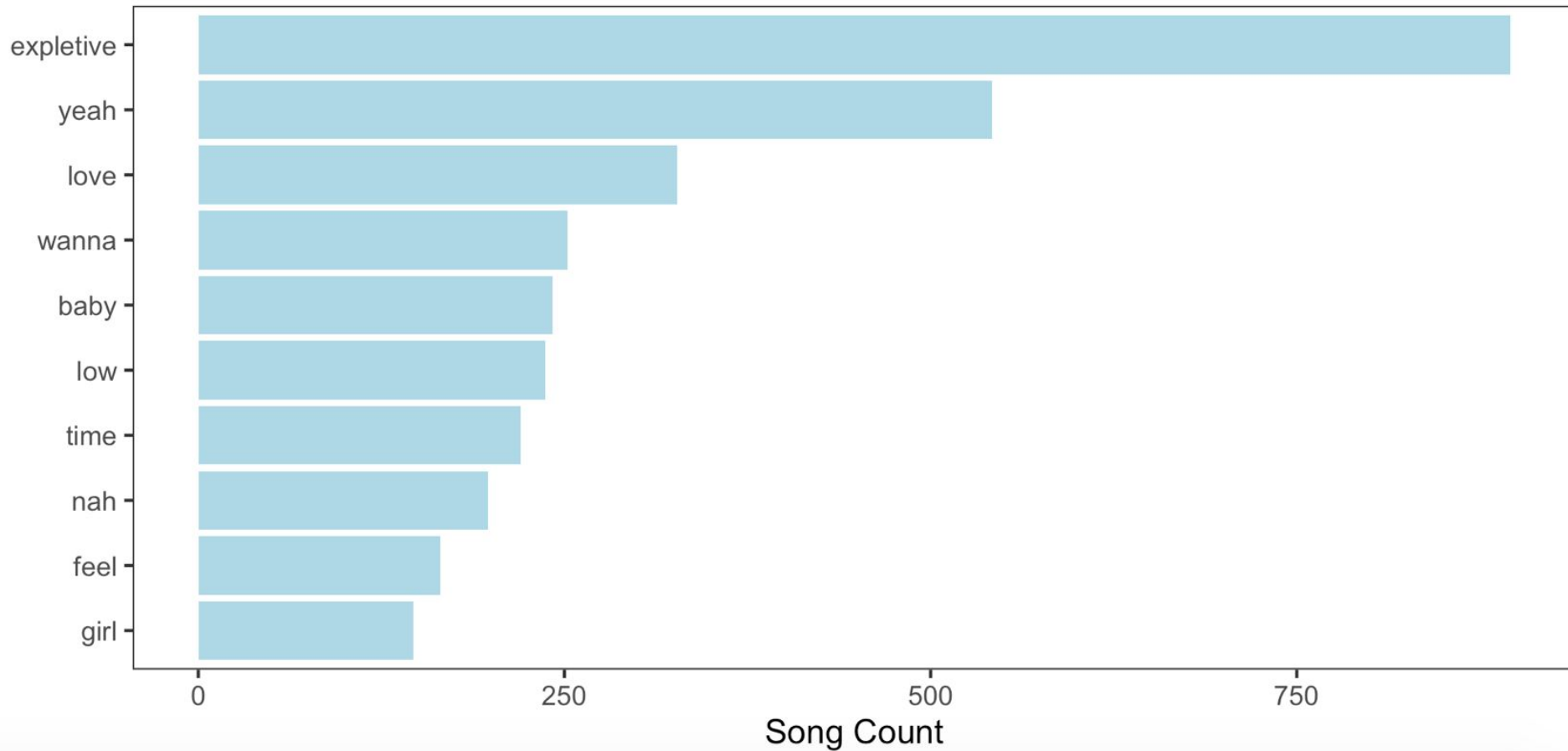
# What words are the most unique to the lyrics of each year's top hits?

Goal: to use TF-IDF to *find the important words* for the content of each document by decreasing the weight for commonly used words and increasing the weight for words that are not used very much in a collection or corpus of documents

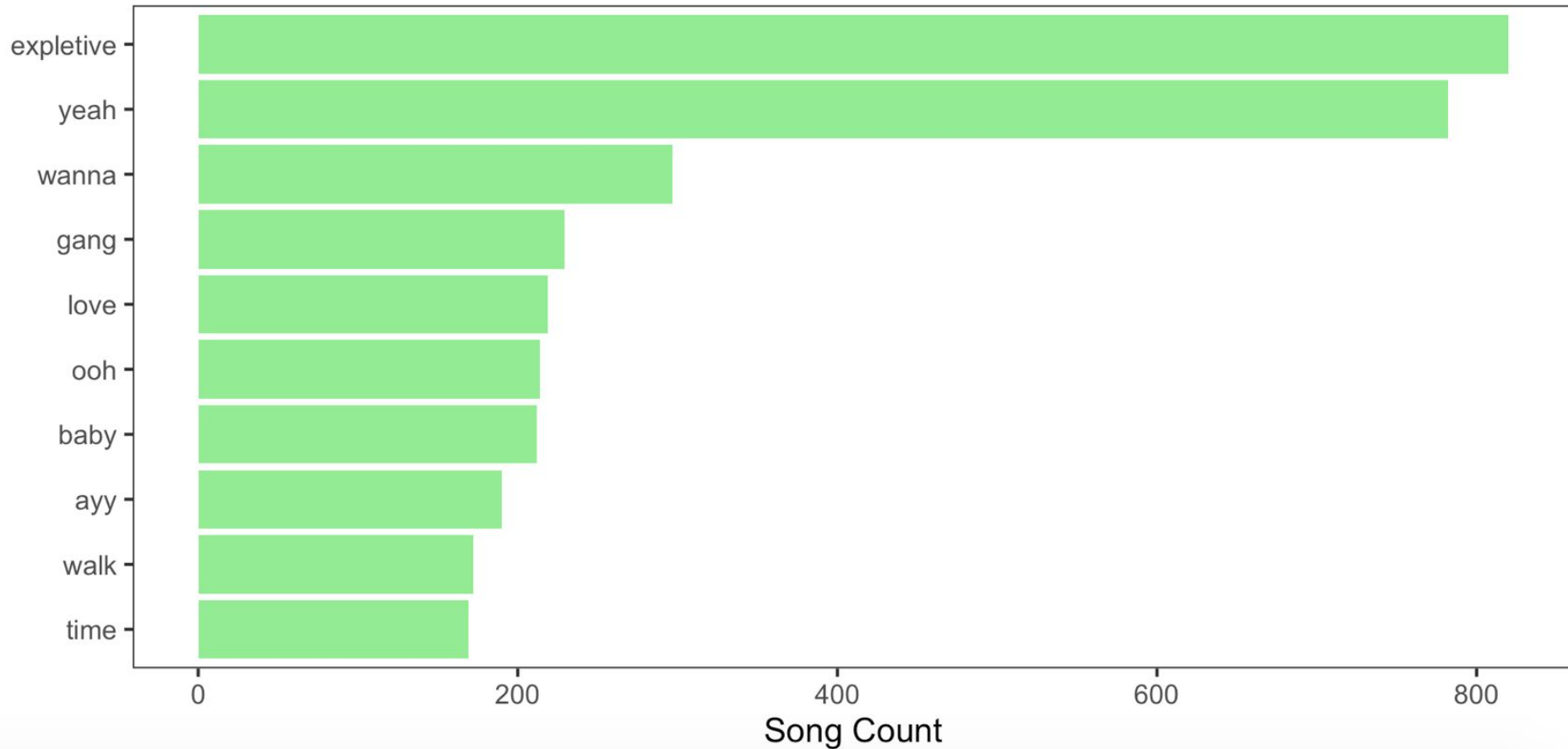
Calculating TF-IDF attempts to find the words that are important (i.e., common) in a text, but not too common



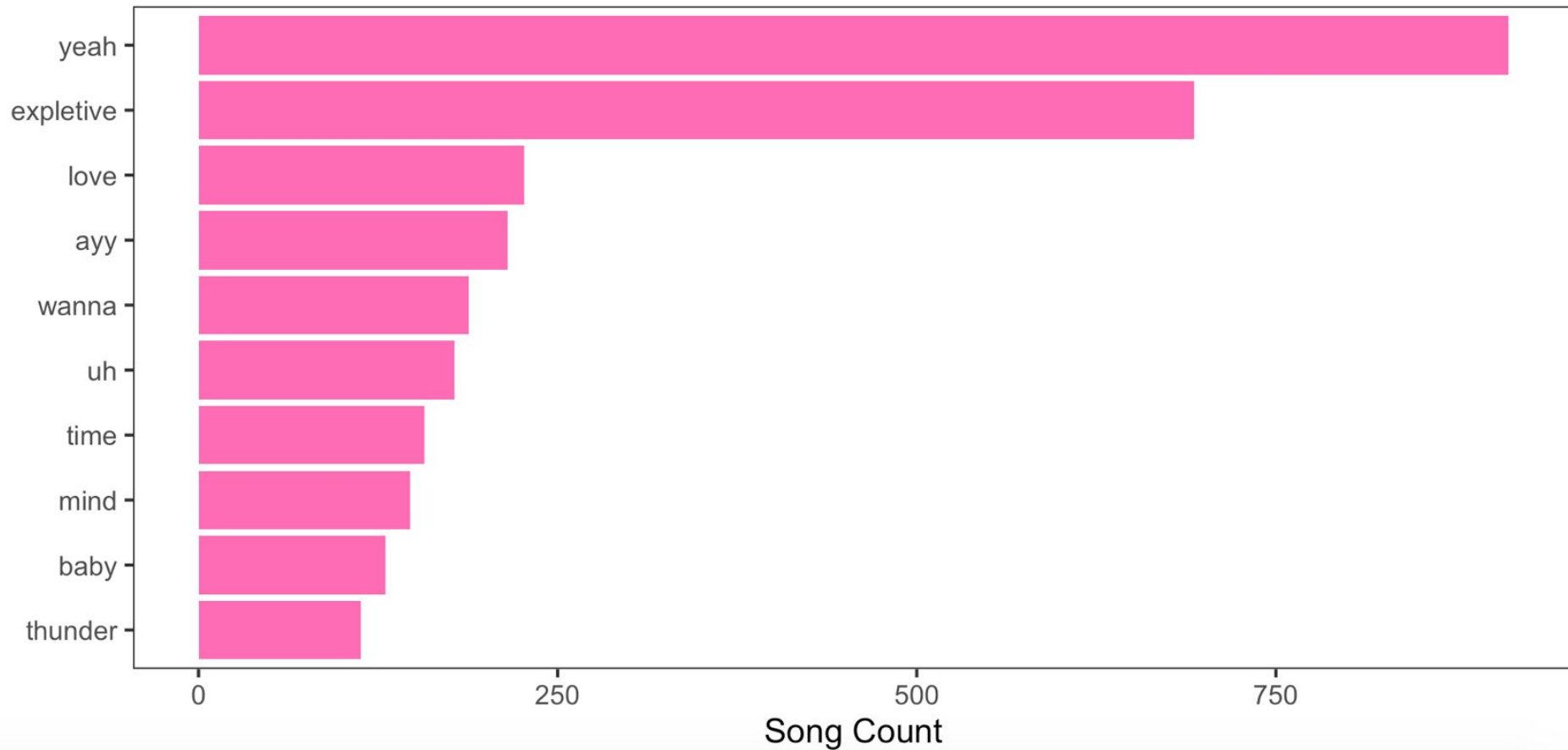
# Most Frequently Used Words in top 200 songs (2017)



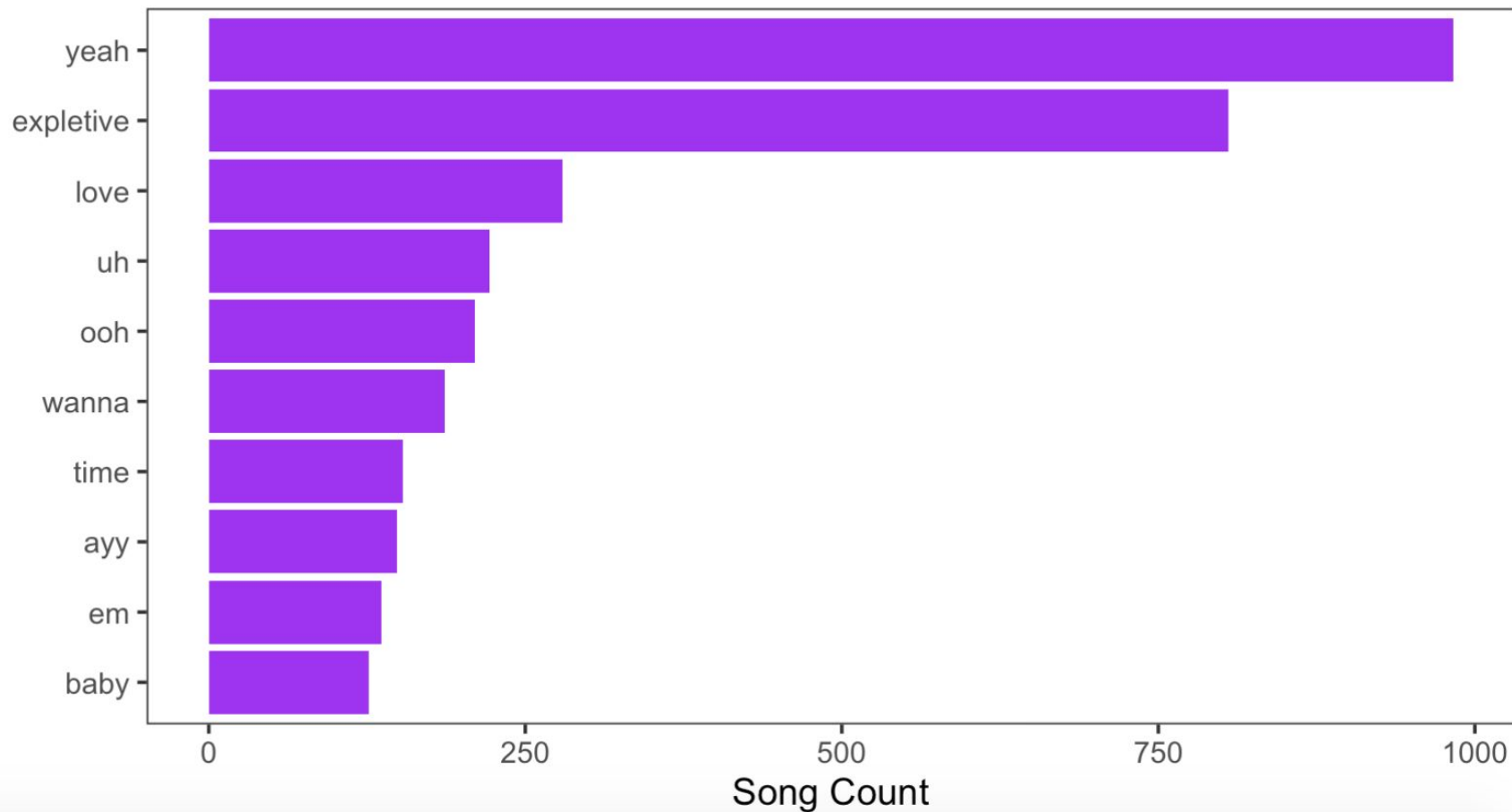
Most Frequently Used Words in top 200 songs (2018)



Most Frequently Used Words in top 200 songs (2019)



# Most Frequently Used Words in top 200 songs (2020)

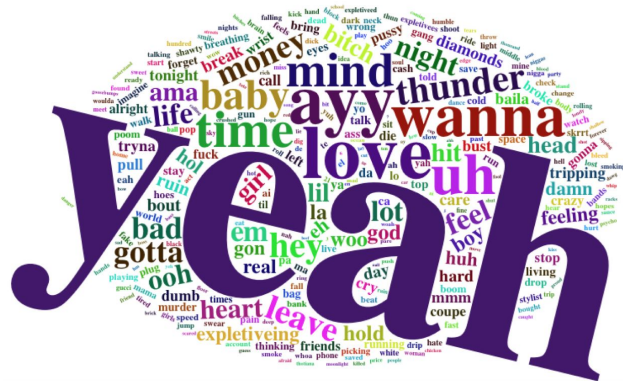




2017



2018



2019



2020

Term Frequency  
can only tell us  
so much....

# TF-IDF:

## Term Frequency - Inverse Document Frequency

the frequency of a term adjusted for how rarely it is used

$$w_{x,y} = tf_{x,y} \times \log \left( \frac{N}{df_x} \right)$$

### TF-IDF

Term  $x$  within document  $y$

$tf_{x,y}$  = frequency of  $x$  in  $y$

$df_x$  = number of documents containing  $x$

$N$  = total number of documents

## Important Words using TF-IDF by Year

