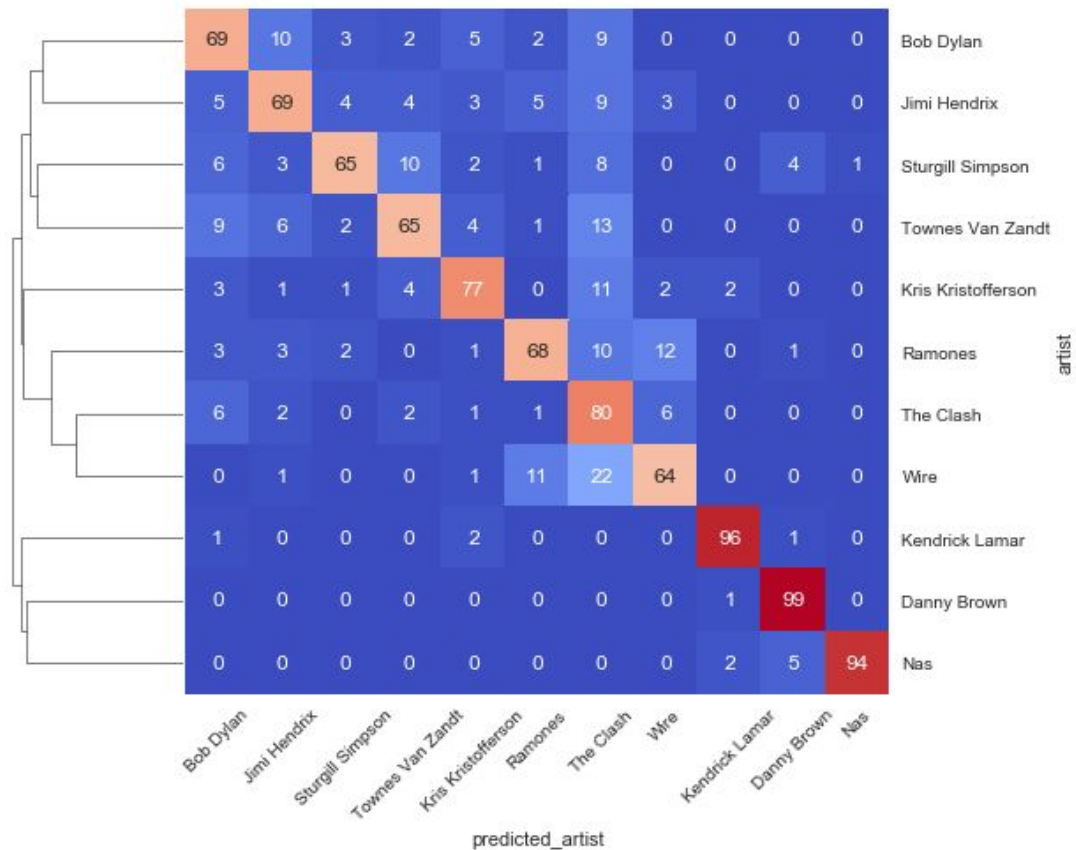# Natural Language Processing

Song Lyrics

# The Artists

- Bob Dylan
- Jimi Hendrix
- Sturgill Simpson
- Townes Van Zandt
- Kris Kristofferson
- Ramones
- The Clash
- Wire
- Kendrick Lamar
- Danny Brown
- Nas

# Data

- All data was taken from the Genius Lyrics API using R
- Regex processing to get rid of data quirks
- Combined into one line per song
- Expanded contractions
- Identified covers and created dictionary to filter

# Machine Learning Approach

- Created a pipeline using:
    - TfidfTransformer from gensim
    - CountVectorizer and SGDClassifier from sklearn
- Ran 100 random iterations of train_test_split and averaged results
- Got an average F-score of .79

# Results

- Clear clustering around Country, Rap, Punk, and Rock (whatever that means)
- Rappers are the most distinctive, likely due to the volume of words available for their genre
- Lots of artists were mistaken for The Clash
  - Wide range of topics for their songs, wide range of emotions
- Danny Brown is the most unique artist
  - The only mix-up features a song with Kendrick Lamar and thus got picked up
- Interesting cross-genre confusion

# Next Steps

- Dig in to each artist's word profile
- Answer questions like:
  - What words make an artist distinctive?
  - Why are artists confused for each other?
  - Why is that one random song predicting consistently for the wrong artist?