# Music Through Ages and Genres

A project by Prab Jaswal

# Problem, Scenario, Context

- Explore how the characteristics of music change over time and by genre

- Can we classify whether a post comes from either the r/DJs or r/musicians subreddit?

- Models: Logistic Regression vs Multinomial Naive Bayes

- Primary evaluation metric: ACCURACY score (consideration given to F1 score).

Dataset by SAURABH SHAHANE
Kaggle

"Music dataset 1950 - 2019"

# Methodology

**1. Cleaning, EDA**
- Processing lyrics / feature engineering
- spaCy
- Vectorise and Analyse

**2. Classifying genre from lyrics**

- Logistic Regression
- Gridsearch

**3. Classifying decade from lyrics**

- Logistic Regression
- Grisearch

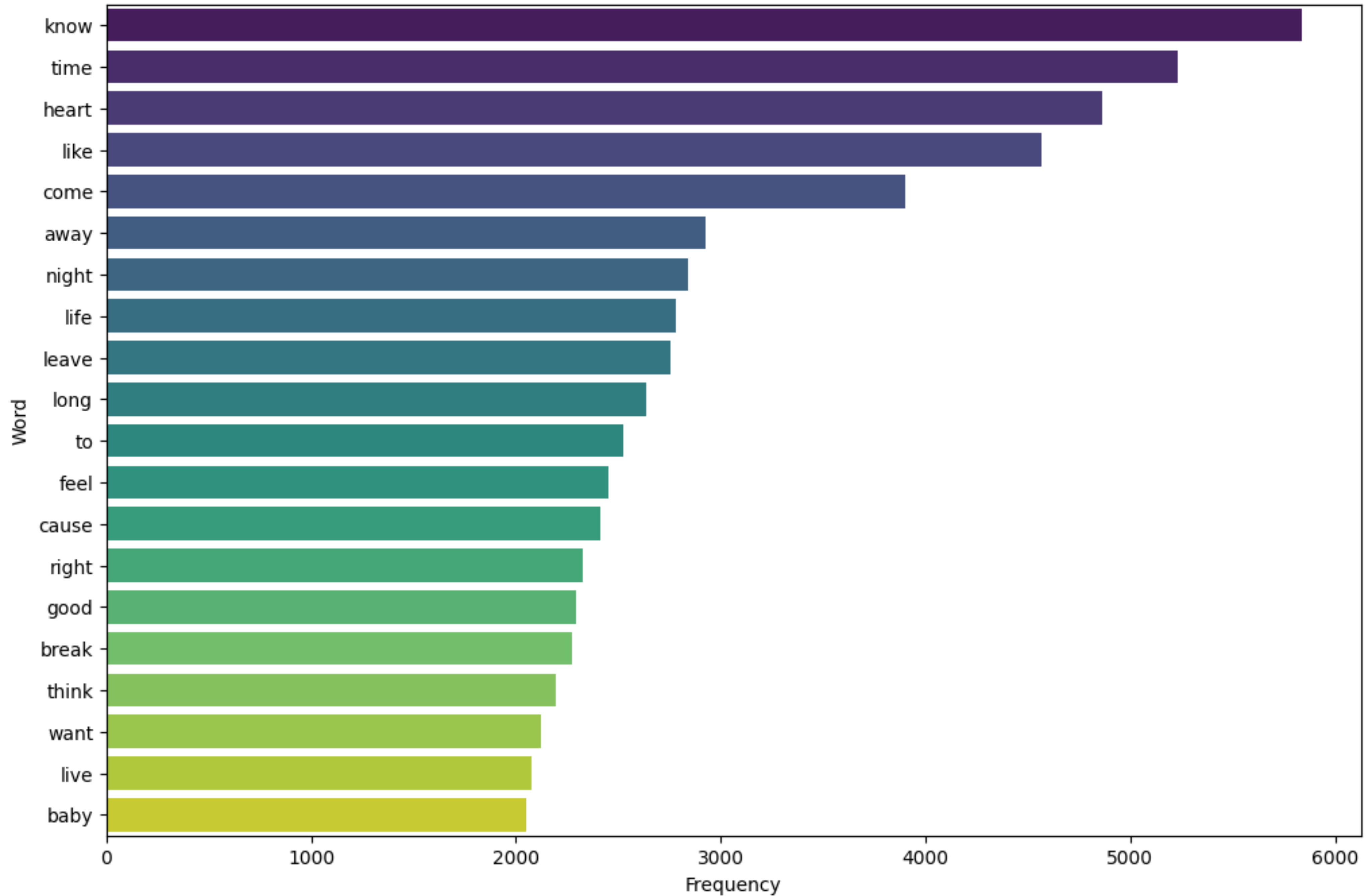**4. Classifying genres from characteristics**

- Random Forest
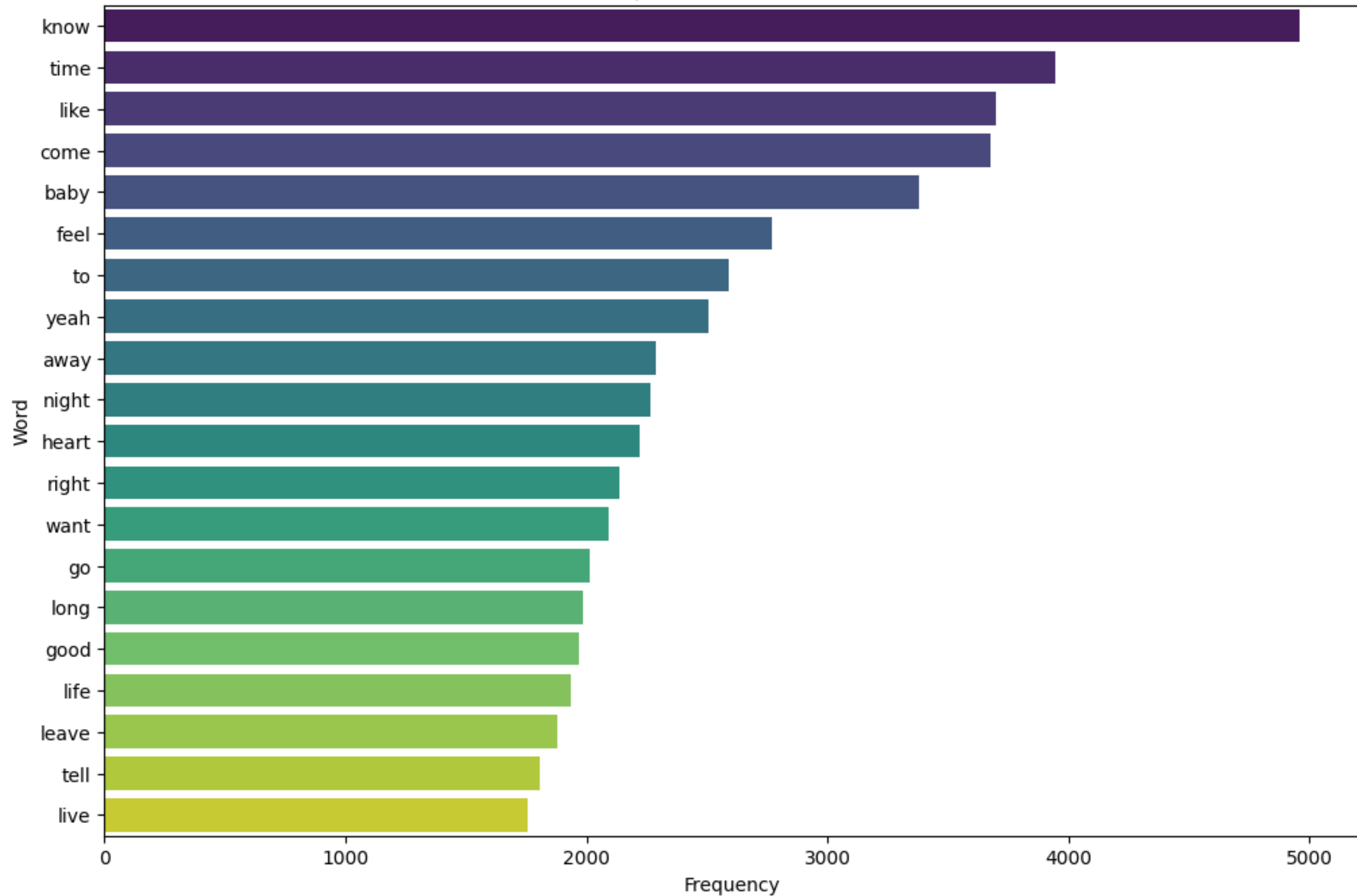- Gridearch

# Exploring the data: Over Time

- Songs have gotten more electronic (less acoustic)

- Songs are getting shorter in length

- Songs are getting louder

- Songs are (as percieved at the point of data collection) getting more energetic as nightlife / club culture evolves.
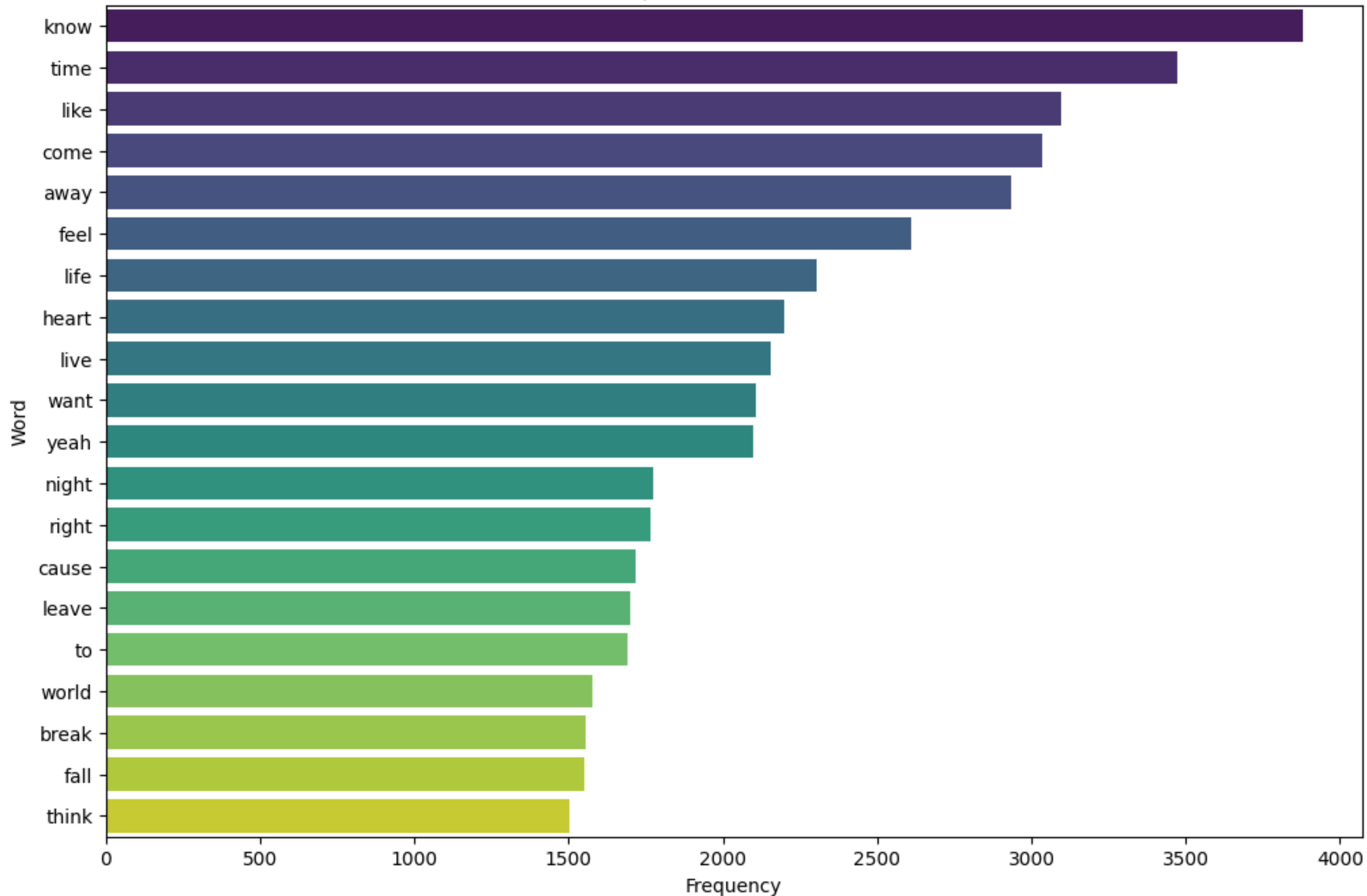
# Exploring the data: by genre
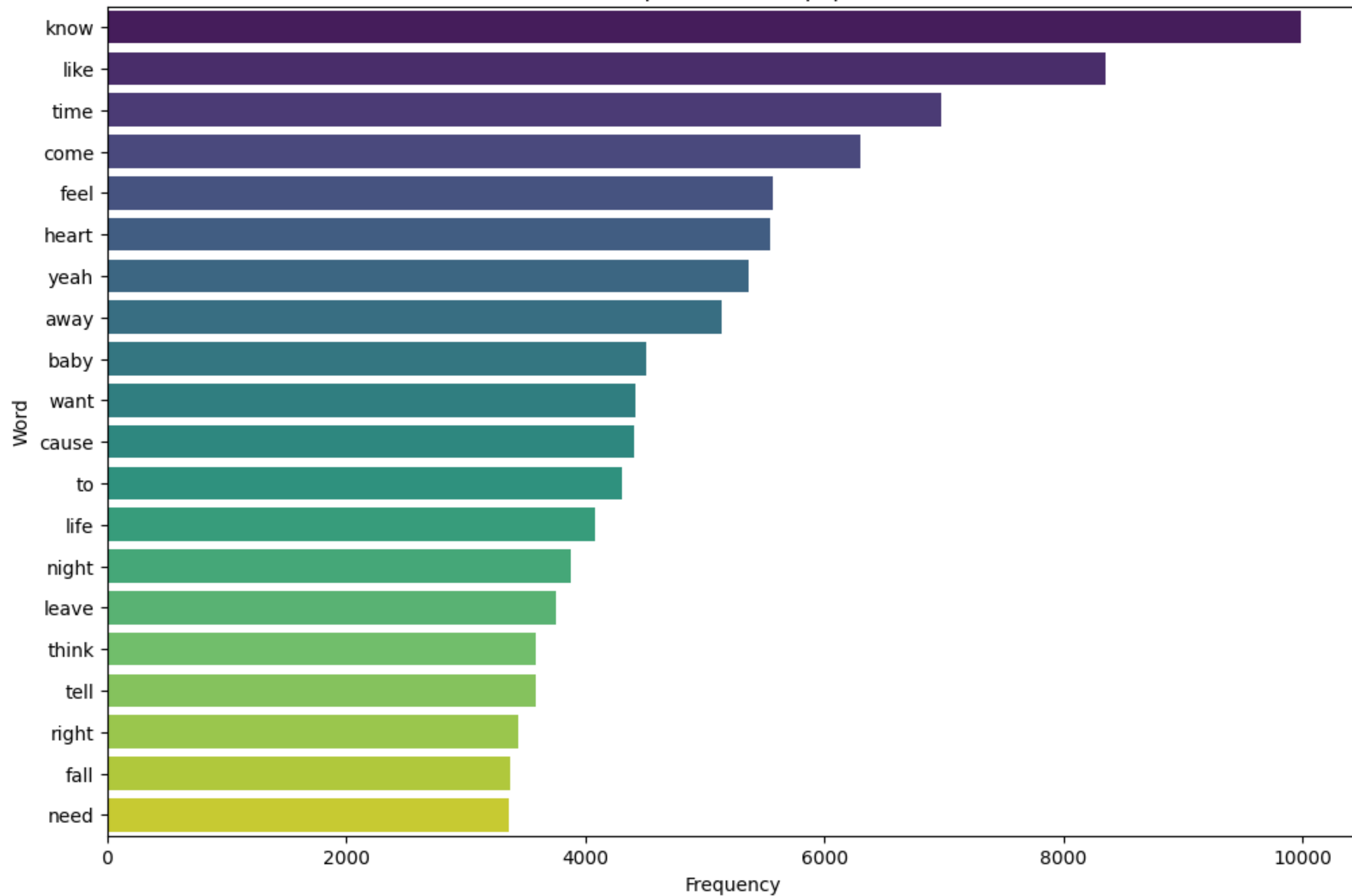
Top 20 Words in country
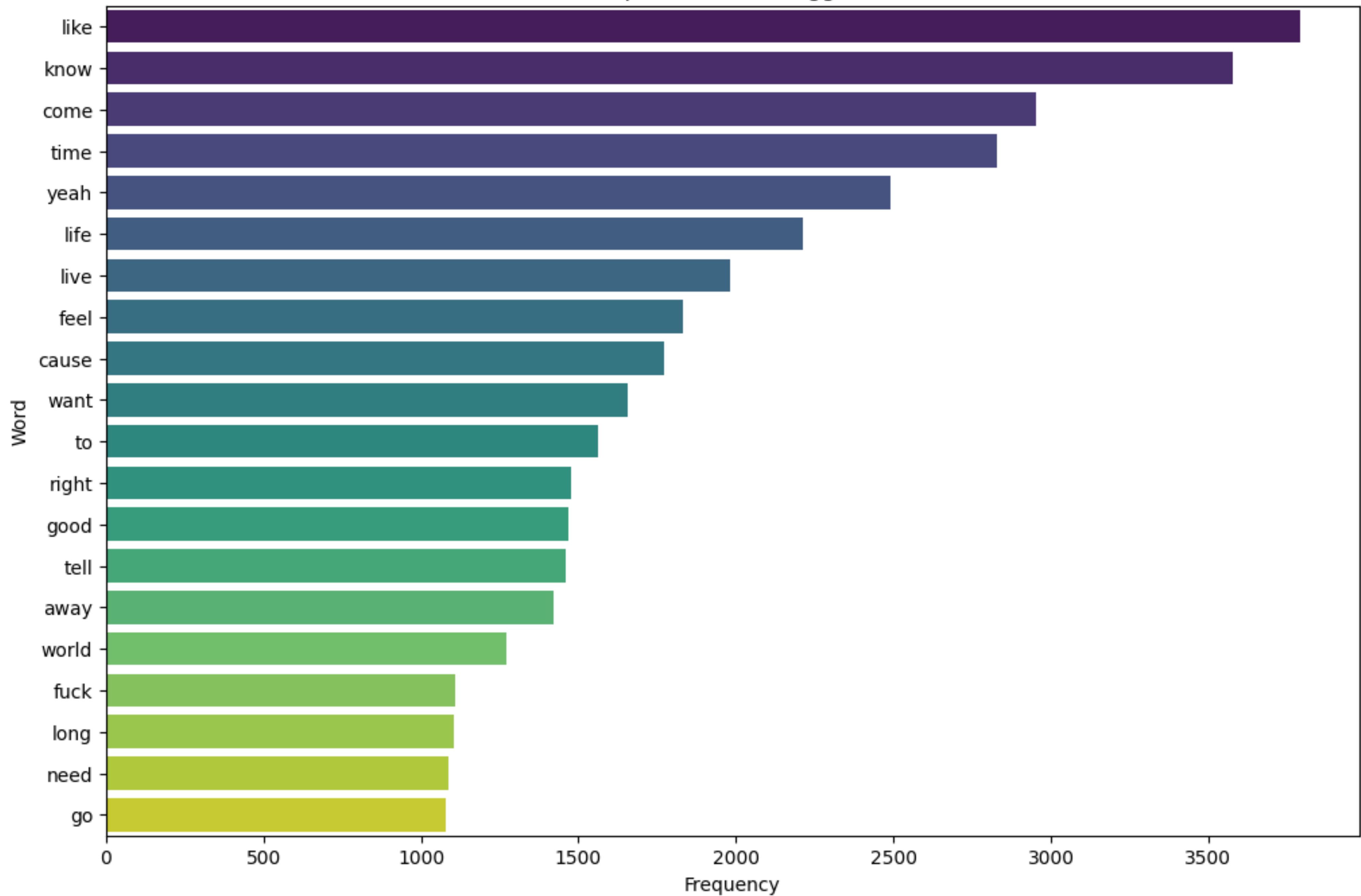
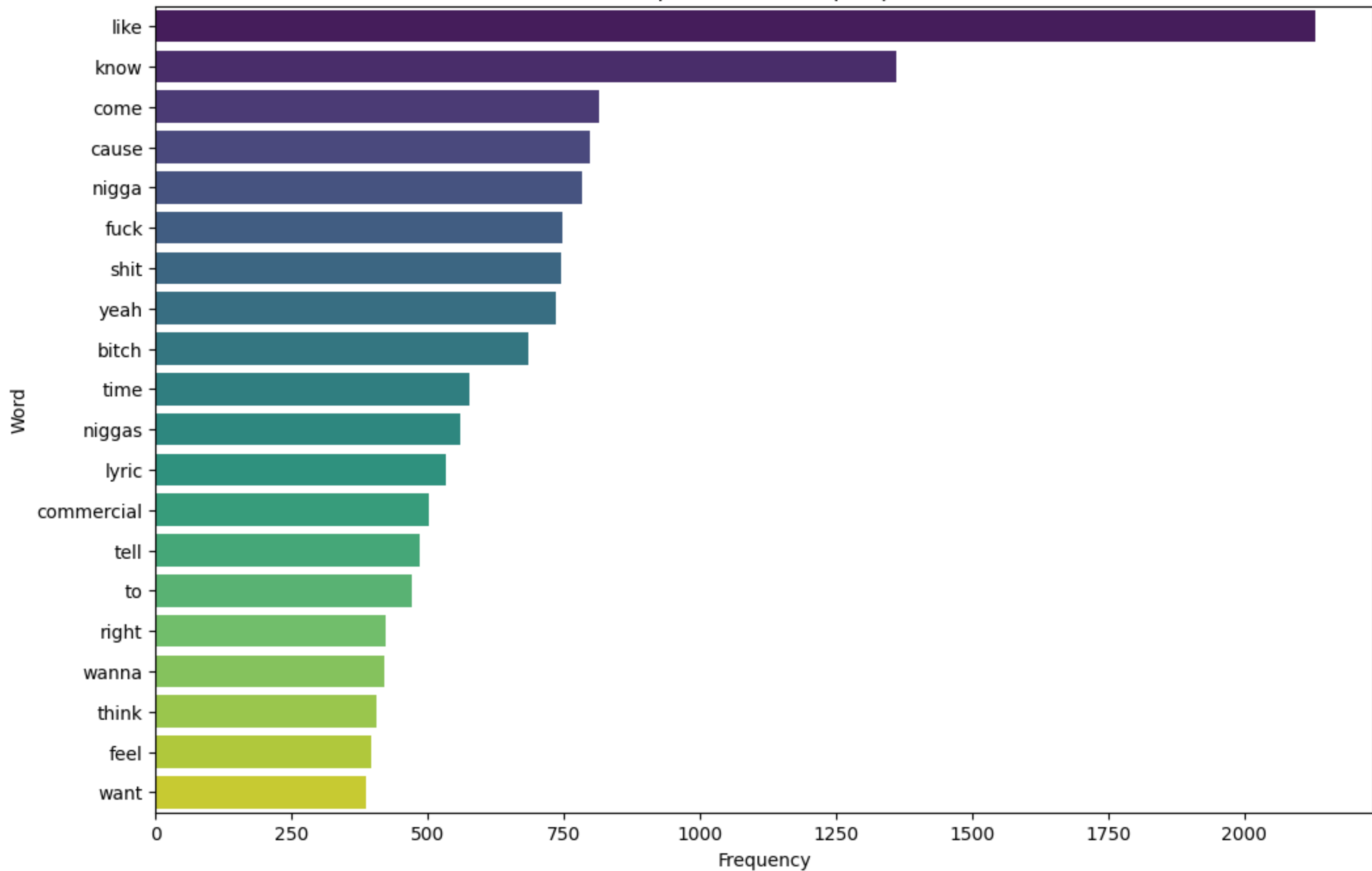Top 20 Words in blues

Top 20 Words in rock

Top 20 Words in pop

Top 20 Words in reggae

# Top 20 Words in hip hop

# Exploring the data: genres by characteristics

- Hip-hop has highest scores in characteristics of Obscenity and Danceability

- Jazz leads in acousticness and instrumentalness with hip-hop being the lowest of these

- Country leads in sadness

- Rock leads in energy and violence.

# Models - Classification

- Classifying genres from lyrics

- Classifying decade from lyrics

- Classifying genres from characteristics

- Imbalanced Classes (challenge) - F1 Score

# Next Steps

- Trying different models (e.g., Random Forest, SVM, Neural Networks).

- Engineering new features or using different text vectorization methods (e.g., Word2Vec, GloVe).

- Increasing the dataset size or improving the quality of the data through webscraping lyrics for songs.

- Given further time to undestand ROC AUC curve, I would also incorporate this as a metric for model evaluation however my current understanding at this time is too rocky for it to be a useful metric in this project.

I welcome a discussion if you choose!