Search

Home

Library

**Capstone:**

# Spotify Track Recommender

Kyle Ness

# Problem

- Overabundance and aimlessness - search for new songs by genre? By artists?
- Waste of time, less entertainment, less discovery
- Can a recommender help?
  - How would it work?
  - Can it be accurate?

# Table of Contents

# 01

## Data

Foundation for recommendation.

# Sourcing

- An accurate recommender requires a lot of data
- Although Spotipy offers an easy to use wrapper, limited time led me to Kaggle
- Close to 600,000 tracks w/ audio features, metadata available for download

# Features

- <u>Meta:</u> song name, artists, release date (yr), genre
- <u>Audio:</u> explicit, popularity, danceability, energy, key, loudness, mode, speechiness, acousticness, instrumentalness, liveness, valence, tempo, duration_ms, time signature

# Cleaning

- Plenty of data -> dropped nulls and bad entries, fixed data types
- Converted strings to their literals in various columns
- Sorted by popularity (for search purposes), dropped duplicates
- Reduction of ~100,000 tracks to 490,000.

# **Preprocessing**

- Scaled numeric features by z-score for more accurate distance calculations
- Working with cleaned tracks data and scaled tracks data going forward
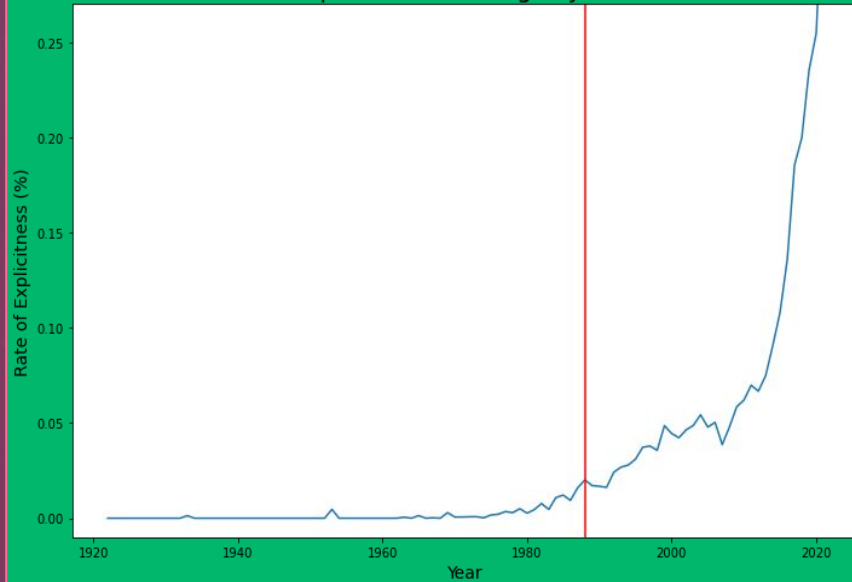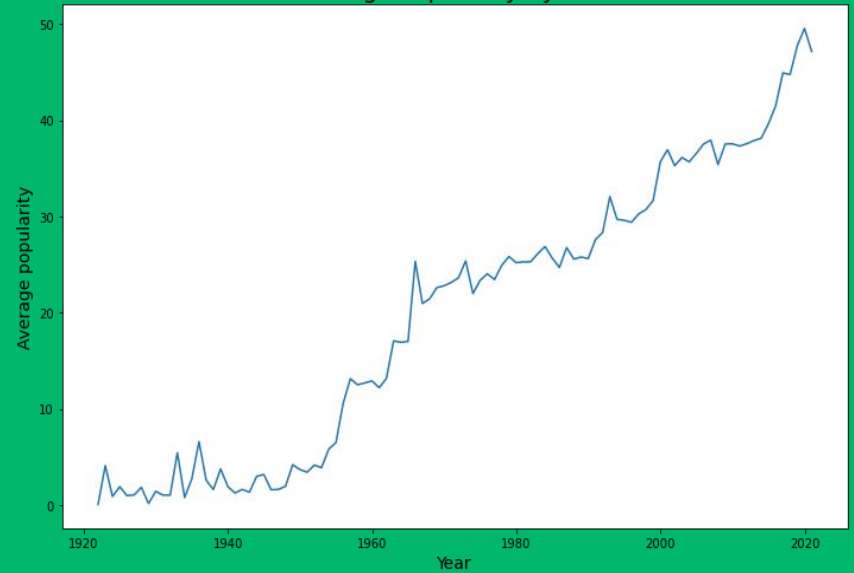- For plotting fun, scaled data was also reduced to 2 principal components with PCA

**02**

EDA

Exploring music the data scientist way.

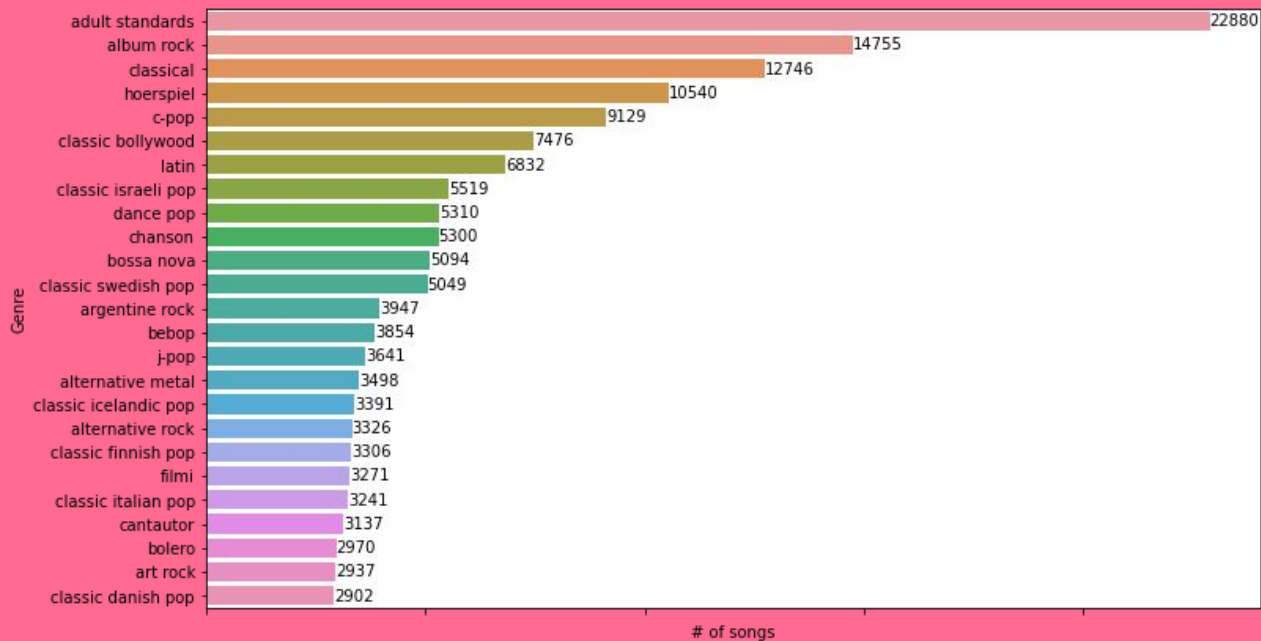# Music Over Time

# Music Over Time

# By Genre

## 25 Most Represented in Dataset



| Genre | # of songs |
|---|---|
| adult standards | 22880 |
| album rock | 14755 |
| classical | 12746 |
| hoerspiel | 10540 |
| c-pop | 9129 |
| classic bollywood | 7476 |
| latin | 6832 |
| classic israeli pop | 5519 |
| dance pop | 5310 |
| chanson | 5300 |
| bossa nova | 5094 |
| classic swedish pop | 5049 |
| argentine rock | 3947 |
| bebop | 3854 |
| j-pop | 3641 |
| alternative metal | 3498 |
| classic icelandic pop | 3391 |
| alternative rock | 3326 |
| classic finnish pop | 3306 |
| filmi | 3271 |
| classic italian pop | 3241 |
| cantautor | 3137 |
| bolero | 2970 |
| art rock | 2937 |
| classic danish pop | 2902 |

# By Genre

## Top 5: Popularity, Energy, Loudness

**Popularity**
- trap queen
- new york drill
- noise pop
- san diego rap
- png pop

**Energy**
- trance mexicano
- metallic hardcore
- garage punk
- swedish death metal
- dark black metal

**Loudness**
- singaporean metal
- aikatsu
- garage punk
- speedcore
- mainland se asia metal

# Cluster Visualization

03

Building

How it's made.

# Thought Process

- Given a list of tracks:
  - Search for each in local scaled data. If not available, pull and scale
  - With list of scaled vectors, calculate mean
  - Calculate cosine distances from mean vector to all tracks
  - Sort and grab 10 closest to the mean
- The 10 adjacent tracks should be similar in composition, and thus probably enjoyable

Search

Home

Library

04

Demo

Let's see if this works!

# Conclusion

- Recommendations can be quite accurate, but may also miss.
- Potential shortcomings:
    - Data (amount, representation)
    - Using mean vector, cosine distances
    - Only able to recommend from local
- Very fun to work with and explore, and personalization can be a lucrative avenue.

Search

Home

Library

# Thanks!

## Do you have any questions?