# Spotify Algorithm

By Eun Lee, Joa Nguyen, Gyujin
Hong, Suebeen Noh

# Table of contents

# Introduction

❏ Create cluster models using song data from spotify to find songs similar to one another

❏ Music listeners all have different tastes and preferences when it comes to their music, and being able to find similar songs will help aid users in finding their next favorite song

❏ Our group plans to create cluster models (K Means, Mixture of Gaussian) to group similar songs and see which songs may be recommended to a Spotify user after seeing their playlists/favorite songs

# Motivation

- Song varies depending on the day's mood and using the unsupervised learning to recommend similar taste of music is more convenient than having to find each individual songs

# Related works

- K-means (unsupervised)
- PCA (unsupervised)
- Mixture of Gaussians (unsupervised)

# Related Works - Literature

PCA with Recommending systems

- Dimension reduction

Yin, C.-X., & Peng, Q.-K. (2012). A careful assessment of recommendation algorithms related to dimension reduction techniques. *Knowledge-Based Systems*, *27*, 407–423. https://doi.org/10.1016/j.knosys.2011.11.022

Mixture of Gaussians

- Clustering within multi dimensional model
- Convergence to local minimum
- Elliptical boundaries with covariance matrices

Agarwal, N., Sikka, G., & Awasthi, L. K. (2022). Evaluation of web service clustering using Dirichlet Multinomial Mixture model based approach for Dimensionality Reduction in service representation. *Information Processing & Management, 57,* 102238. https://www.sciencedirect.com/science/article/pii/S0306457320300492
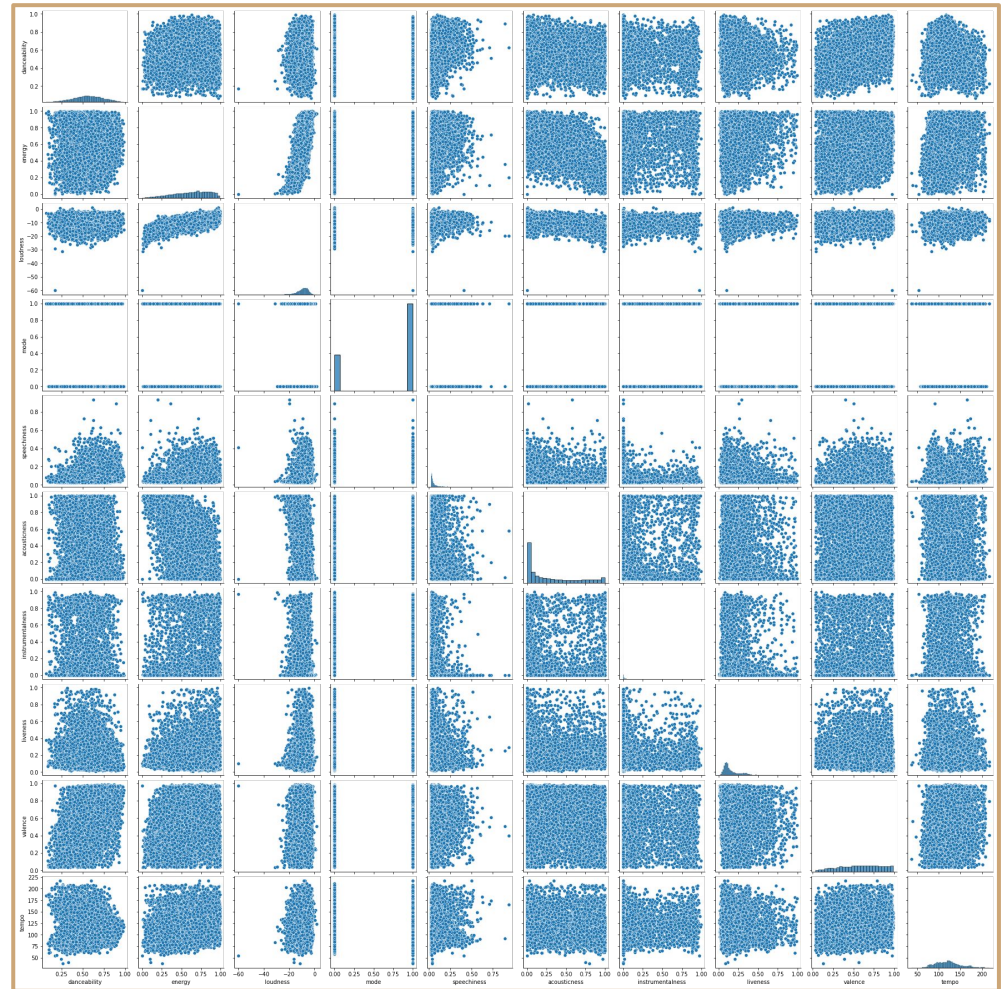
# Data

- ❏ Imported songs using Spotify API with the help if a python library called spotipy
- ❏ Imported data from a 10,000 song playlist called "Top 10,000 songs"
  - ❏ Has enough variety in song genre/styles to be a usable dataset
- ❏ Dropped columns we determined were not needed for song similarity

| | danceability | energy | loudness | mode | speechiness | acousticness | instrumentalness | liveness | valence | tempo |
|---|---|---|---|---|---|---|---|---|---|---|
| **9837** | 0.737 | 0.5580 | -9.814 | 1 | 0.3260 | 0.394000 | 0.021000 | 0.1110 | 0.525 | 98.546 |
| **9838** | 0.389 | 0.7360 | -5.945 | 0 | 0.0356 | 0.000004 | 0.443000 | 0.2950 | 0.224 | 74.992 |
| **9839** | 0.787 | 0.9400 | -9.707 | 1 | 0.1880 | 0.434000 | 0.000000 | 0.1510 | 0.705 | 106.892 |
| **9840** | 0.243 | 0.0969 | -16.510 | 1 | 0.0408 | 0.971000 | 0.013700 | 0.2060 | 0.385 | 47.692 |
| **9841** | 0.675 | 0.3660 | -11.566 | 1 | 0.0310 | 0.724000 | 0.002330 | 0.1860 | 0.791 | 119.122 |
| **9842** | 0.457 | 0.4800 | -15.776 | 1 | 0.0300 | 0.380000 | 0.016000 | 0.1420 | 0.683 | 157.967 |
| **9843** | 0.389 | 0.5320 | -5.632 | 1 | 0.0266 | 0.001580 | 0.000532 | 0.3270 | 0.212 | 169.918 |
| **9844** | 0.724 | 0.4840 | -5.986 | 1 | 0.0482 | 0.720000 | 0.000136 | 0.0808 | 0.933 | 77.124 |
| **9845** | 0.379 | 0.5410 | -8.121 | 1 | 0.0304 | 0.043800 | 0.015800 | 0.3380 | 0.429 | 91.434 |
| **9846** | 0.679 | 0.7090 | -5.851 | 0 | 0.1330 | 0.676000 | 0.000000 | 0.2920 | 0.673 | 84.543 |

# Data

- ❏ PCA to decrease the
  dimensionality to 2 principle
  components (for data
  visualization) and scaled the
  data using StandardScaler()
- ❏ Pairplot to see any
  relationship between the
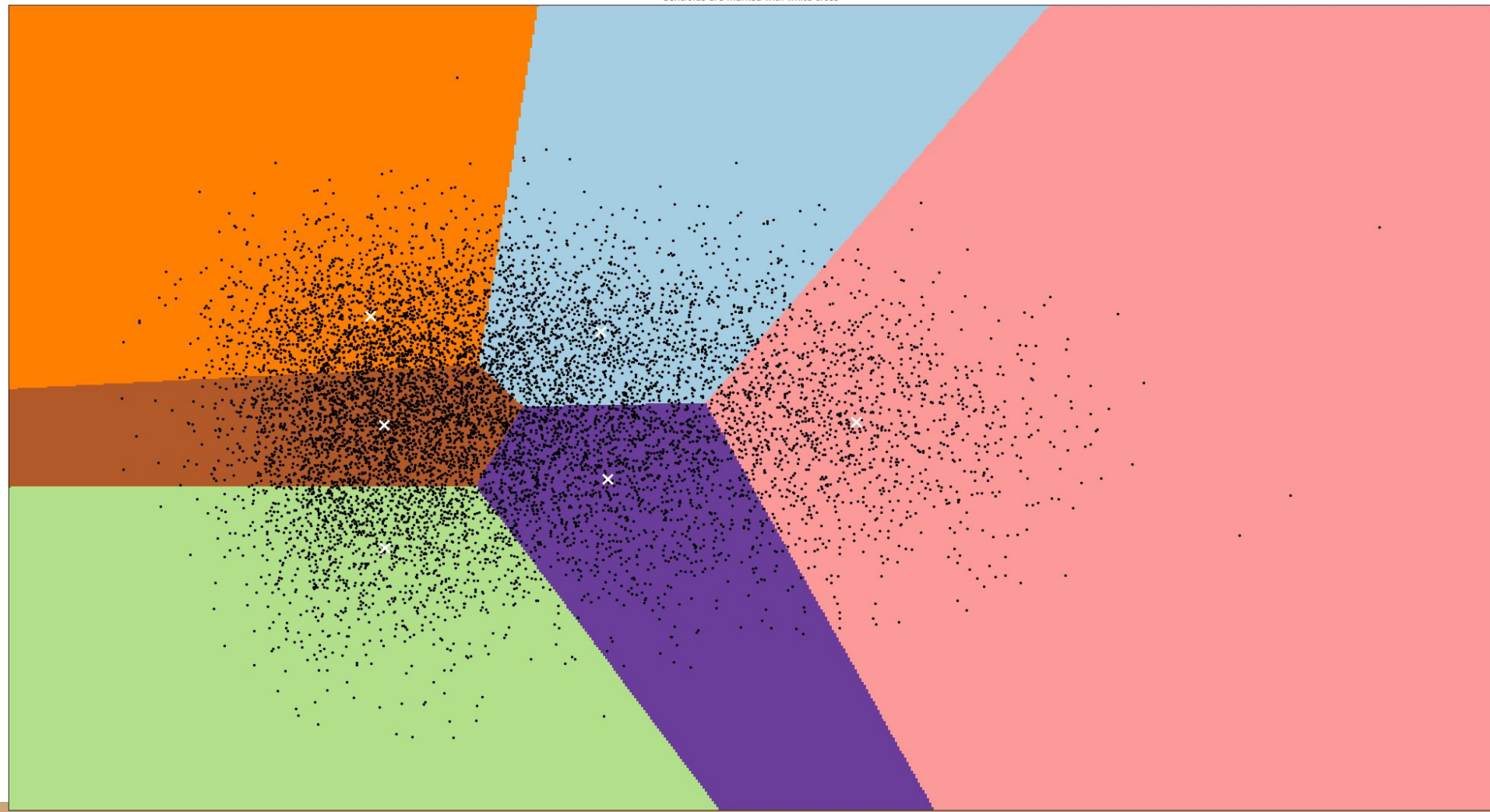  variables

# Methods

- ❏ Pipeline includes…
  - ❏ Standardization method:
    - ❏ StandardScaler()
  - ❏ PCA
  - ❏ Cluster Model Methods
    - ❏ K-Means: KMeans()
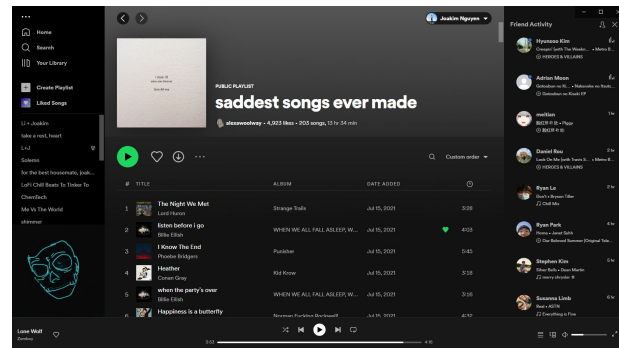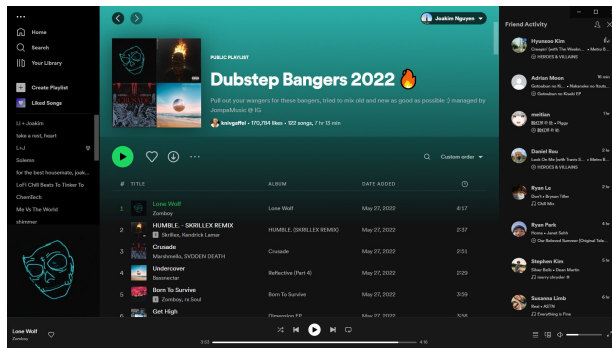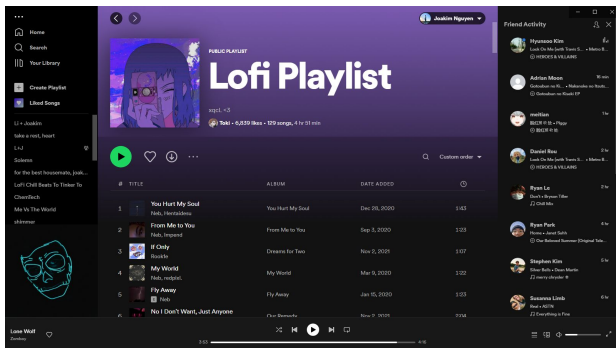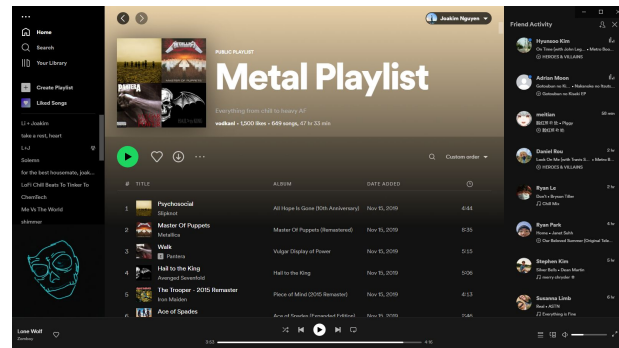    - ❏ Mixture of Gaussians: GaussianMixture()
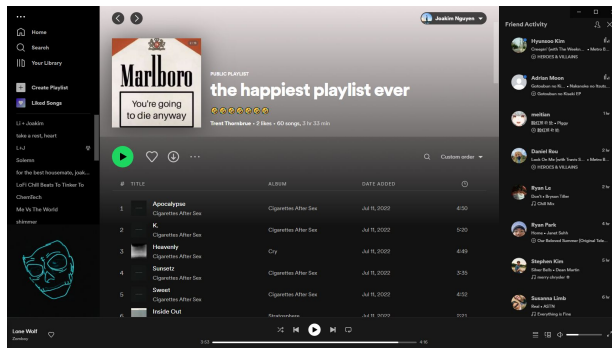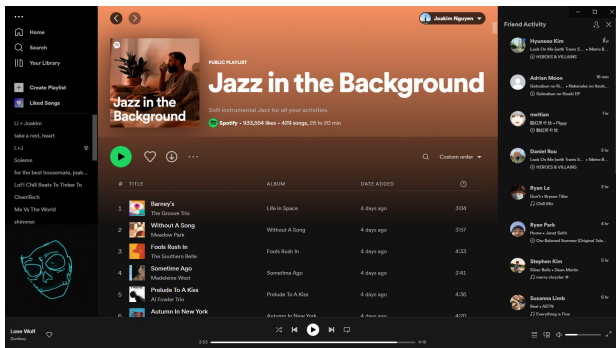
```python
pipe = Pipeline([('scaler', StandardScaler()),
                 ('pca',    PCA(n_components=2)),
                 ('kmeans', KMeans())])
```

# Methods – K Means

- ❏ GridSearchCV() method to determine the number of means to use for the model

    - ❏ Found that data was so clumped it was difficult to find optimal k-means using gridsearch

    - ❏ After adding additional playlists (sad, happy, jazz, etc.) we found 6 fit the model the best

- ❏ Matplotlib to create a meshgrid to visualize model

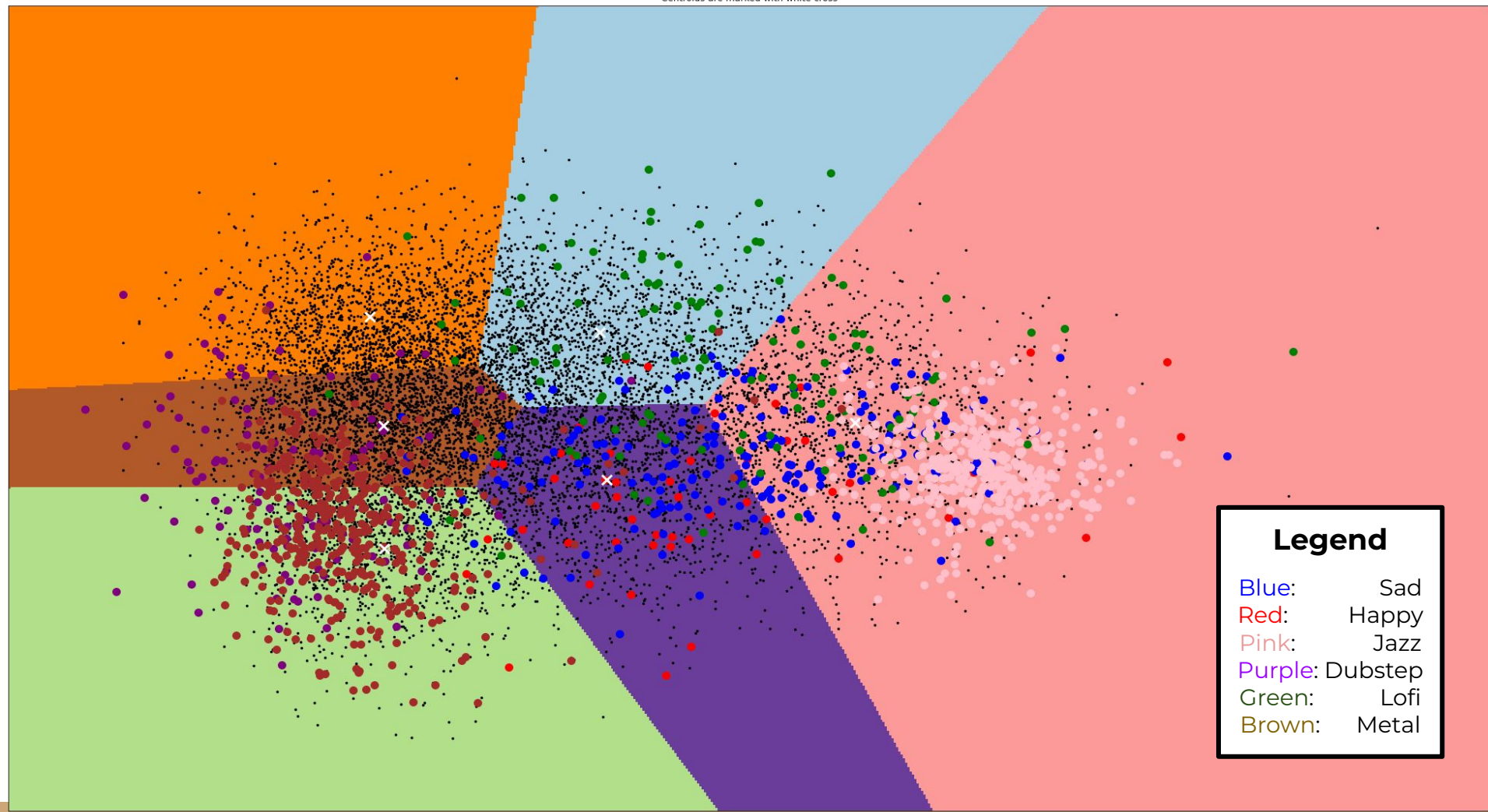- ❏ Added the 6 new playlists onto the model fitted by the original data

K-means clustering on the digits dataset (PCA-reduced data)
Centroids are marked with white cross
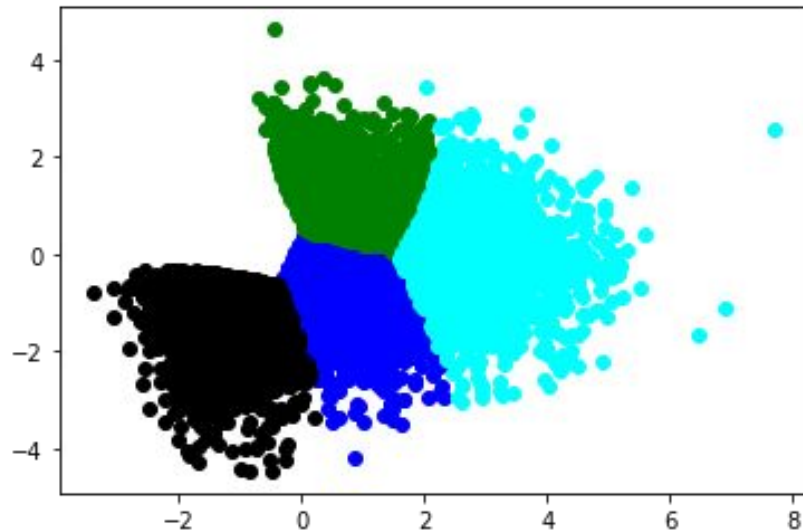
**6 unique playlists with their own genres**

K-means clustering on the digits dataset (PCA-reduced data)
Centroids are marked with white cross

Legend

Blue:      Sad
Red:       Happy
Pink:      Jazz
Purple: Dubstep
Green:    Lofi
Brown:   Metal

# Methods – Mixture of Gaussians

❏ Added ('MoG', GaussianMixture()) into

    the pipeline

❏ Visualized Data using sklearn.cluster

❏ Ultimately found that since the song

    data is not randomly distributed,

    clusters were not gaussians

# Result - K Means

❏ After visualizing the data, we compared the resulting graphs ranging from five to twelve cluster points (k) and concluded that six clusters (k=6) fitted our data best.

❏ After the best k value was determined, we found six different genres of playlists with distinctively strong features of song data.

❏ After overlaying the six playlists' clusters identified by different colors on top of our k-means model on the 10,000 random songs, we tried to identify which k-cluster certain playlists overlapped with the best

# Discussion

- Deeper understanding of:

    - Implementing spotipy (Spotify API wrapper), pipeline to find a recommender algorithm

    - How K-means and PCA's dimensionality reduction worked well to reduce noise

- Could have increase the principal components for a more accurate model

- Mixture of Gaussian isn't the best model since the data is not random

- Our playlist was only selected with sample of 10,000 songs from a single playlist (which could potentially be biased)

# References

- https://www.sciencedirect.com/science/article/pii/S0167865509002323?casa_token=HXAcvkia2XMAAAAA:DbZtC

  3vPao5Js7T53qhcVI8svx_OiHZ9NEzfBwH_o2qlfEOSyJqsy48H6JjP_IAmeNvOI4iPrN0

- https://www.analyticsvidhya.com/blog/2019/10/gaussian-mixture-models-clustering/

- https://developer.spotify.com/documentation/web-api/

- https://spotipy.readthedocs.io/en/2.21.0/