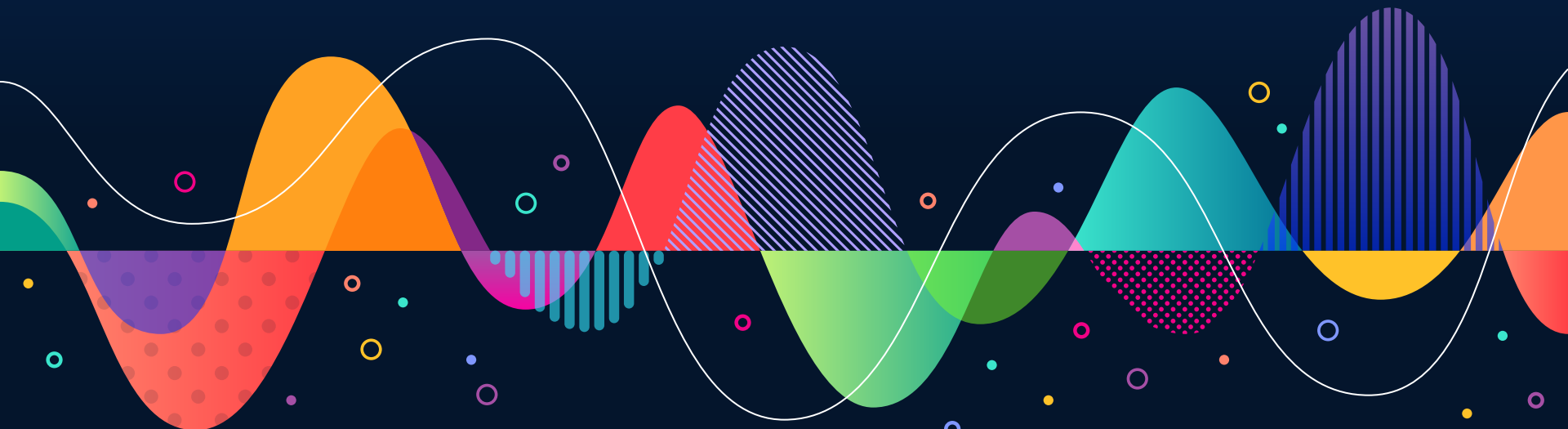# CAP5610: Machine Learning Final Project Presentation

*Would Your Favorite Song Make it to Spotify's "Top 50 - Global" Playlist?*

April 19, 2022

Bailey LaRea        Nadine Rose        Suma Marri

# Agenda

Initial Project Questions
Description of the Data
Data Preparation
Visualizations for EDA
Revised Project Question
Baseline Model
Preliminary ML models
Next Steps

# Initial Project Questions

Idea: Can we use Machine Learning techniques to help **predict** what songs will make it into Spotify's "Top 50 - Global" playlist?

Can specific track attributes be used to predict overall song popularity?

# Description of the Data

# There are 13 Track Attributes

TrackName

ArtistName

Genre

Beats Per Minute

Energy

Danceability

Loudness (dB)

Liveness

Valence

Length

Acousticness

Speechiness

Popularity

# The Datasets

## Top 50 Songs in 2019

▷ This dataset contains information for the top 50 songs in 2019.

▷ 50 rows of data

## 2019 Songs

▷ This dataset contains songs that were released in 2019.

▷ ≈ 11,000 rows of data

*Both datasets are sourced from Kaggle*

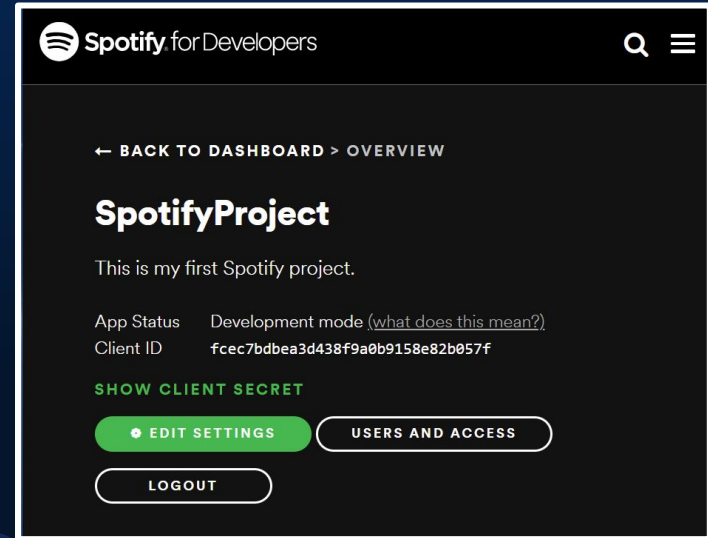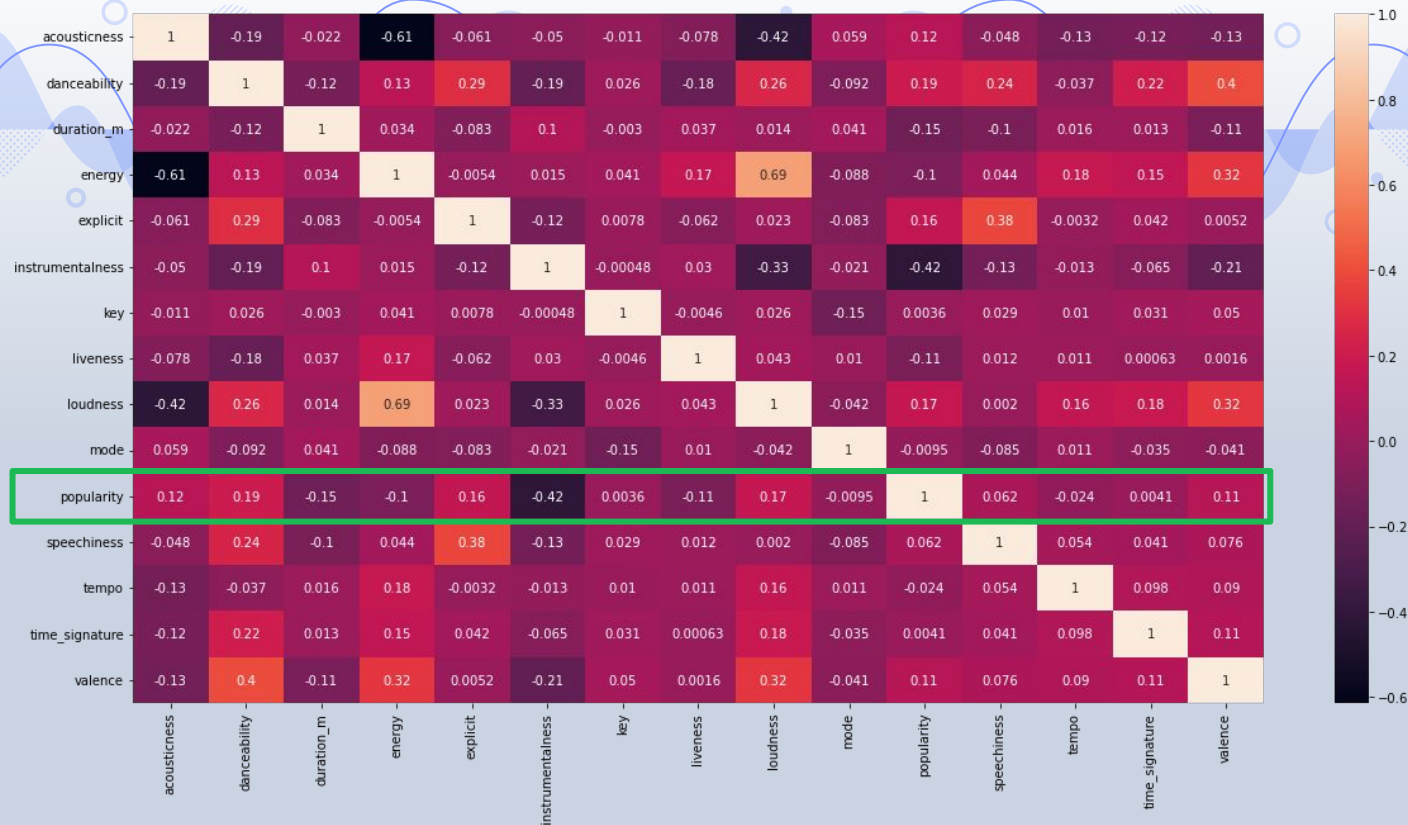# Data Preparation

Kaggle

- We used this for our initial dataset and for the top 50 songs.

Spotify for Developers  - API

- The Spotify for Developers allowed us to extract and collect the data from Spotify
- We used the Client Keys to access the data on Spotify.
- After extracting the data, we had deleted and duplicate values and any records that had null values.

**Visualizations - Heatmap**

# Visualizations - Scatterplots

## Instrumentalness vs. Popularity

We can see that there is a **negative correlation** between both variables.



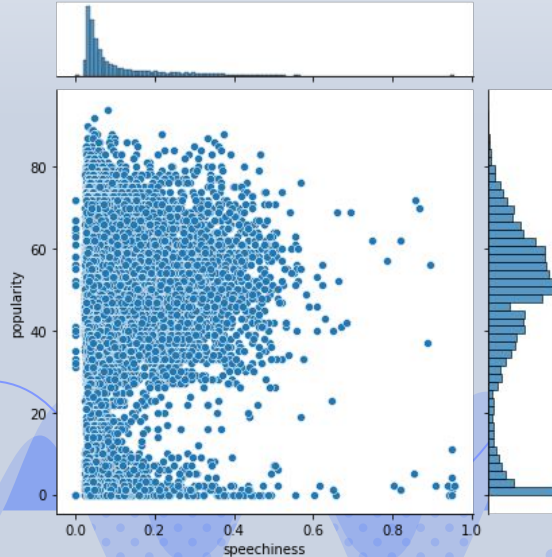## Speechiness vs. Popularity

We can see that there is **no correlation** between both variables.



## Loudness vs. Popularity

We can see that there is a **positive correlation** between both variables.

# Visualizations - Histogram

Here are the following distributions of Speechiness and Loudness from another dataset. They follow distributions that align with the prior slides.

# Revised Project Question

Forecasting → <u>Time Series</u> for Baseline Model → ❌ Data following sequential order with respect to time
❌ Time as the independent variable

**Classification** → ✅ Allows us to predict if a song will be classified as a song within the 2019 Top 50 - Global Playlist.

# Baseline Model: Linear Regression

```
from sklearn.model_selection import cross_val_score
from sklearn import datasets, linear_model

scores = cross_val_score(model,X,Y, cv=5)
print("Print all scores: ", scores)
print("Mean Accuracy: ", scores.mean())
```

```
Print all scores:  [0.21003084 0.16980323 0.08999494 0.15596551 0.07952351]
Mean Accuracy:  0.14106360537917992
```

```
R2 Score:  0.2059066468330093
MAE:  13.970961045076779
MSE:  342.95193847341307
RMSE: 18.518961592740915
```

With a mean accuracy score of only ~14%, some
changes need to be made!

# Preliminary Version of ML Models

**Decision Tree - Plot of RMSE**



| | feature | importance |
|---|---|---|
| 4 | instrumentalness | 0.408603 |
| 2 | duration_ms | 0.128069 |
| 0 | acousticness | 0.120920 |
| 7 | loudness | 0.065383 |
| 1 | danceability | 0.060481 |
| 3 | energy | 0.051191 |
| 9 | valence | 0.047880 |
| 8 | tempo | 0.046519 |
| 6 | speechiness | 0.043091 |
| 5 | liveness | 0.027864 |

There is a little bit of improvement.... However, there must be a **better model** out there?

Random Forest & KNN are also used!

# Random Forest: Classification

**Random Forest**

**Baseline: LR**

RMSE: 18.518961592740915

Mean Accuracy: 0.697

| | feature | importance |
|---|---|---|
| 4 | instrumentalness | 0.202412 |
| 0 | acousticness | 0.126836 |
| 2 | duration_ms | 0.116127 |
| 7 | loudness | 0.093867 |
| 3 | energy | 0.088406 |
| 1 | danceability | 0.079084 |
| 9 | valence | 0.076297 |
| 5 | liveness | 0.073639 |
| 8 | tempo | 0.073136 |
| 6 | speechiness | 0.070195 |

There is great improvement…. However, is there a **better model** out there?

# K Nearest Neighbor (KNN)

**KNN**

Mean Accuracy Score:   0.9917773561037319

**Baseline: LR**
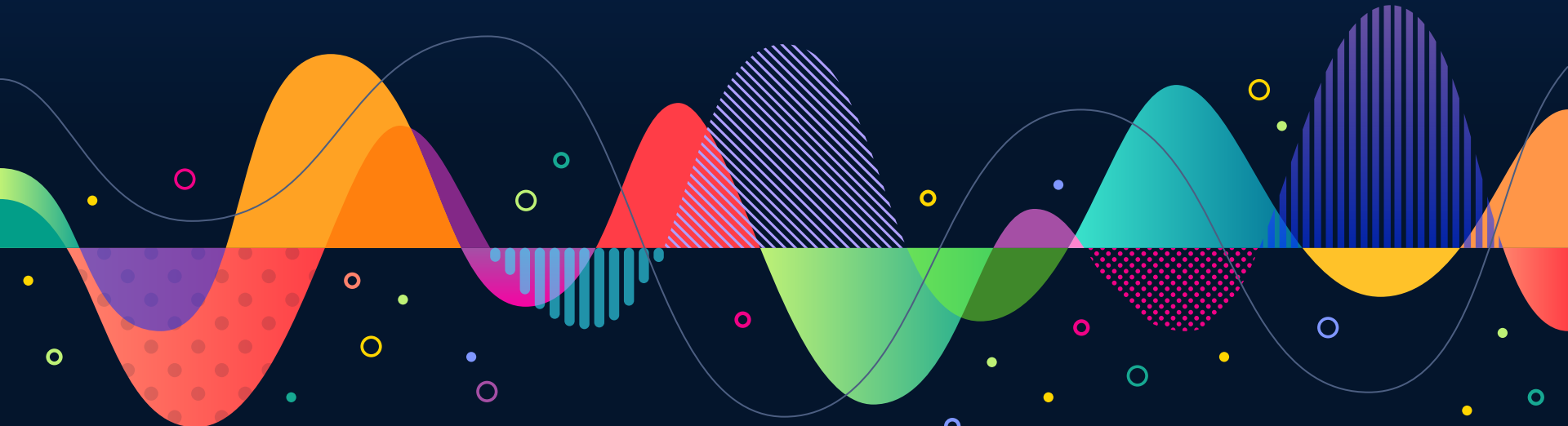
Mean Accuracy:   0.14106360537917992

There is significant improvement.... However, there
is the risk of **overfitting**!

# Next Steps

▷ Investigate potential overfitting of KNN model

▷ Use the finalized ML models to see if successful classification is achieved

▷ Data Preparation

# Questions?

# References

- https://www.kaggle.com/theoverman/the-spotify-hit-predictor-dataset?select=dataset-of-00s.csv

- https://www.kaggle.com/leonardopena/top50spotify2019

- https://towardsdatascience.com/predicting-popularity-on-spotify-when-data-needs-culture-more-than-culture-needs-data-2ed3661f75f1

# Notes from Dr. Yousefi

▷ RMSE is used for regression models

▷ Mean Accuracy is used for classification models

▷ Need to use more regression models to compare with the baseline of linear regression

▷ Set a baseline for classification models and have at least 2 other models to compare to baseline

▷ For overfitting, it is suggested to use diagnostic plots

▷ All learning models should have at least 0.5 mean accuracy, check linear model again

▷ Decision Tree is defined as regression model in our code not classification

▷ Good to give the statistics of the data in the report

▷ 94% is good and usually not overfitting

▷ Regression = predict

▷ Argue why we are choosing the features in the analysis