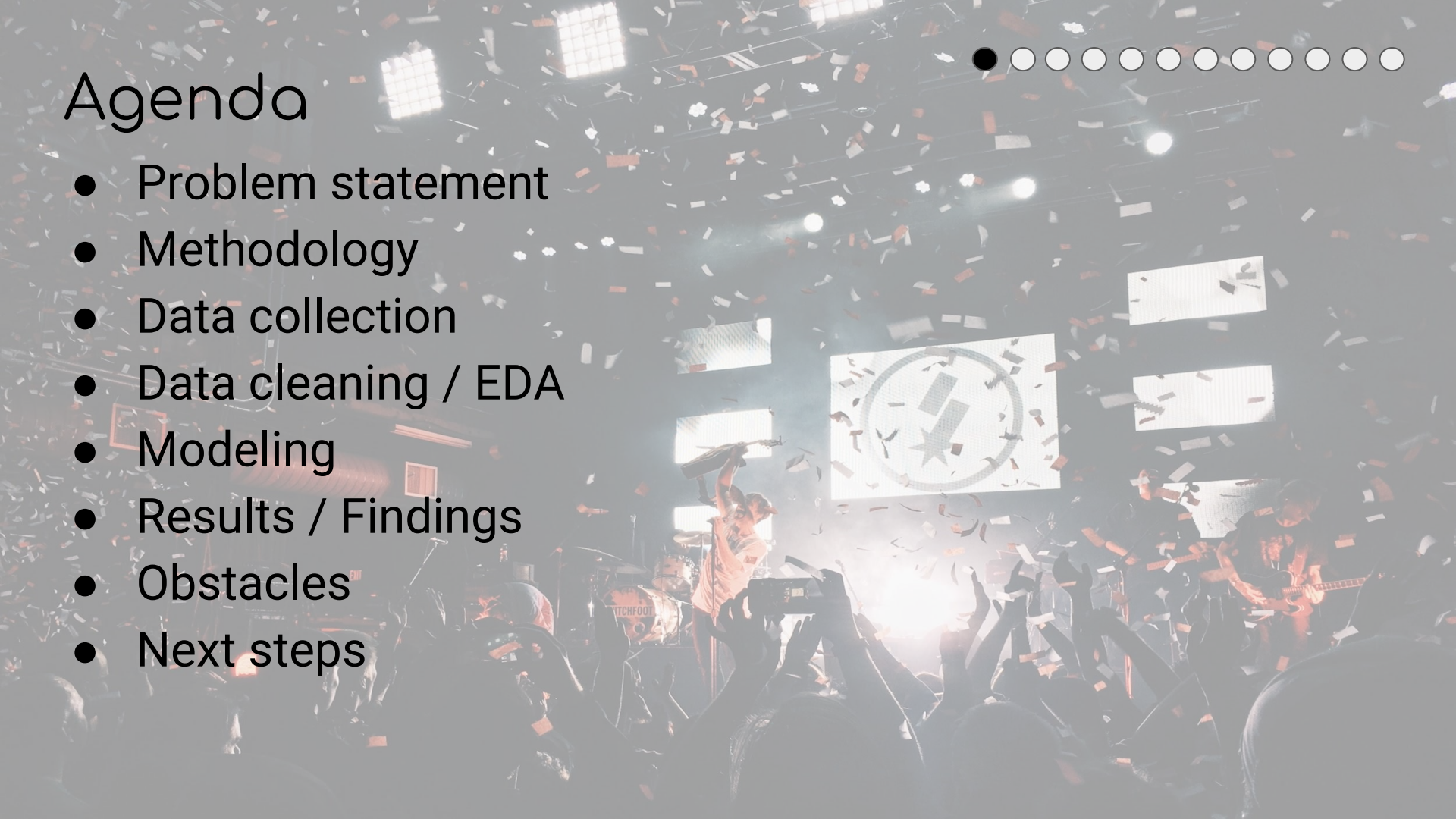# An Analysis of Hit Songs

BY DAN KIM

# Agenda

- Problem statement
- Methodology
- Data collection
- Data cleaning / EDA
- Modeling
- Results / Findings
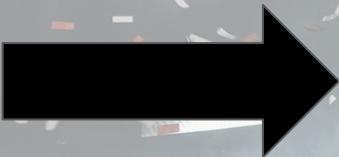- Obstacles
- Next steps

# Problem Statement

Which audio features have the most influence
on making a hit song?

# Who can this benefit?

Record Labels, independent artists

# Data Collection

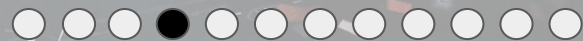**Steps:**

1. Web Scrape:
   a. BillBoard Year-End Hot 100 Songs (2008, 2013, 2018)
   b. SongFacts (2008, 2013, 2018)
2. Interface with API:
   a. Genius - song lyrics
   b. Spotify - audio features

# What are audio features?

- Danceability
- Energy
- Loudness
- Mode
- Speechiness
- Acousticness
- Instrumentalness
- Liveness
- Valence
- Duration (milliseconds)

- Type
- ID
- URI
- Track Reference URL
- Track Analysis URL

# Data Cleaning / Preprocessing

Audio Features
- Create labels
- Alleviate imbalanced classes
- Create song titles column
- Convert Duration column to seconds
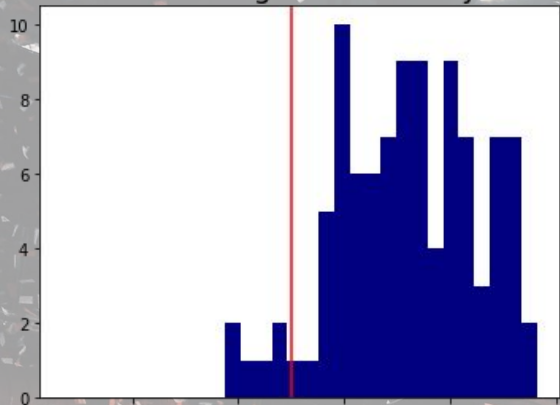- Drop duplicates and unnecessary columns

Lyrics
- Tokenize and Stem
- Create labels
- Alleviate imbalanced classes
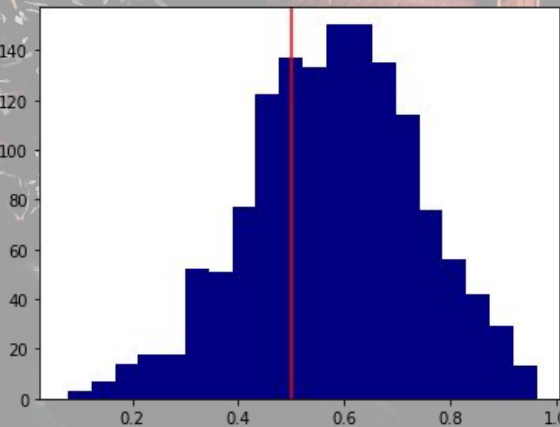
# EDA - danceability

# EDA - loudness



2018 song's loudness    2013 song's loudness    2008 song's loudness

2018 song's instrumentalness

2013 song's instrumentalness

2008 song's instrumentalness

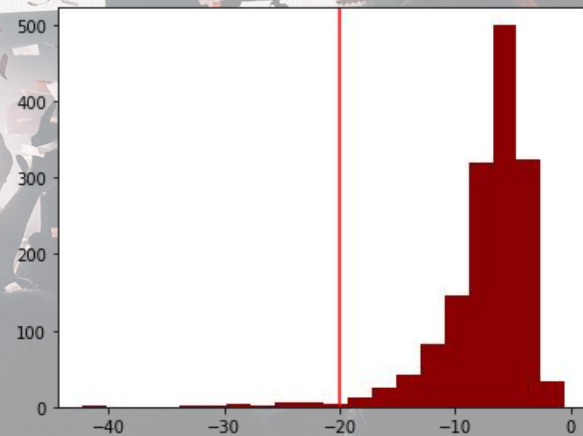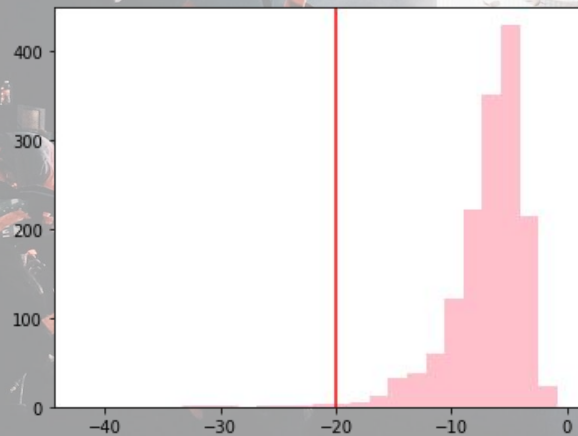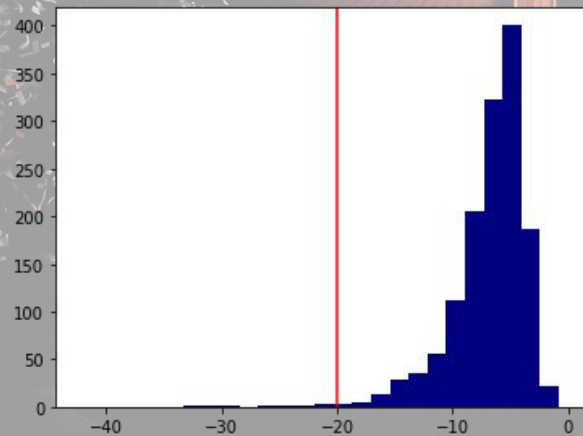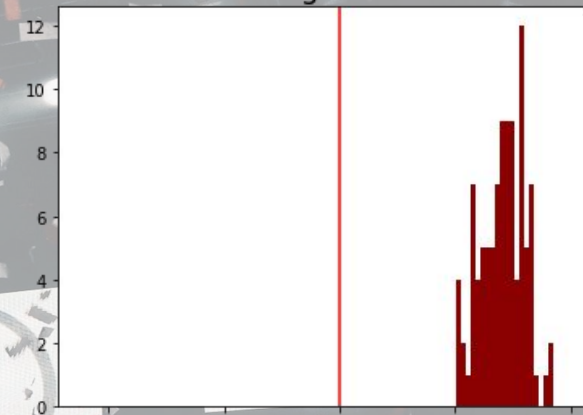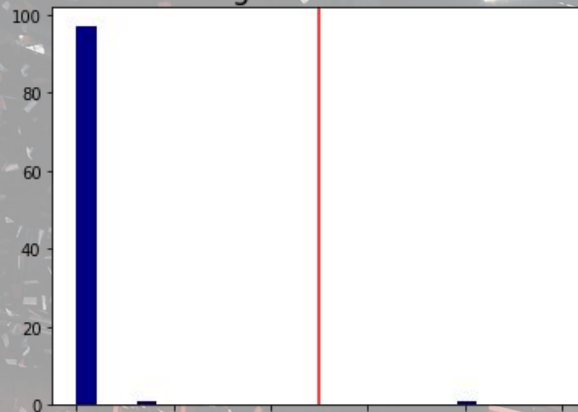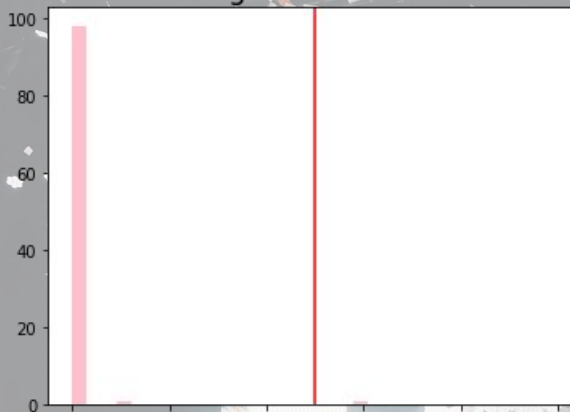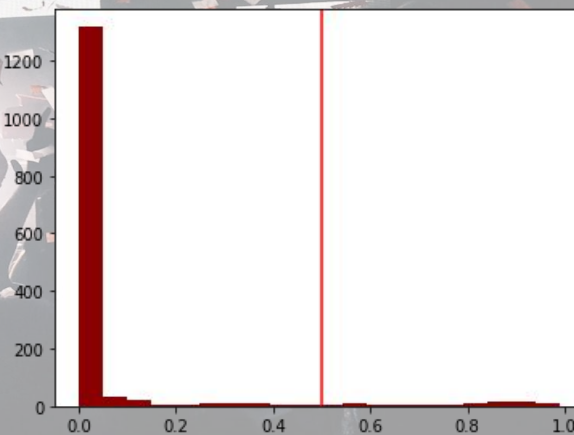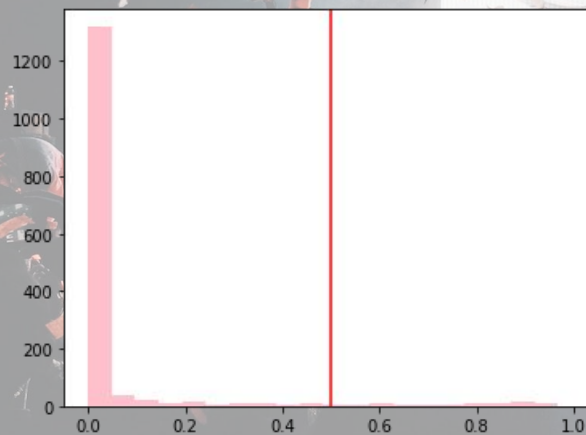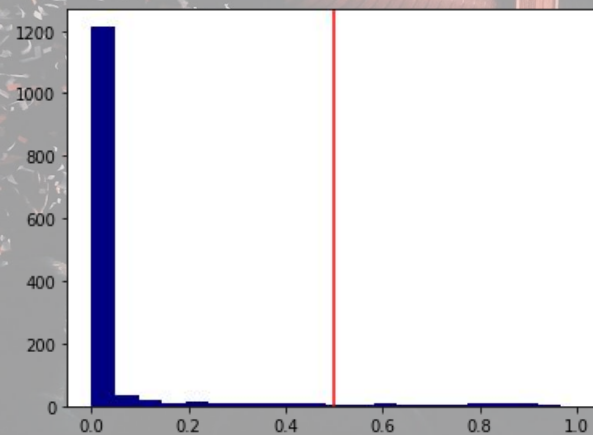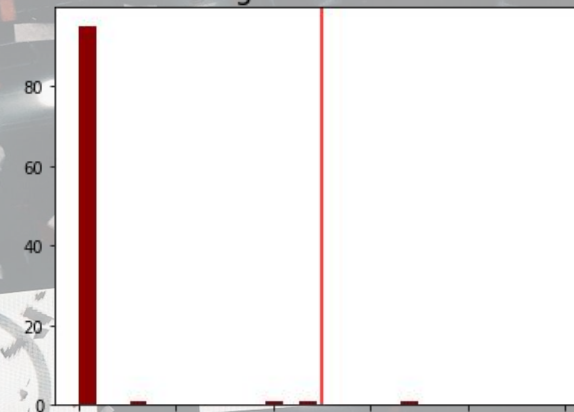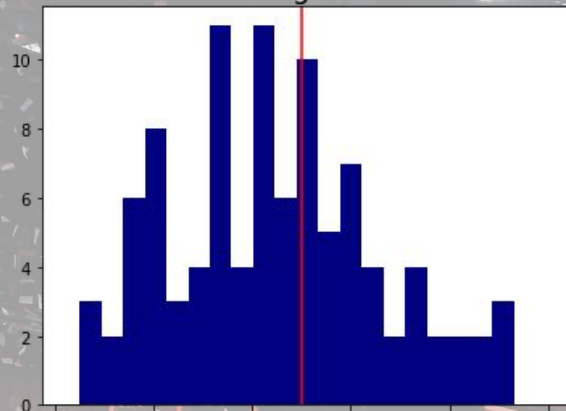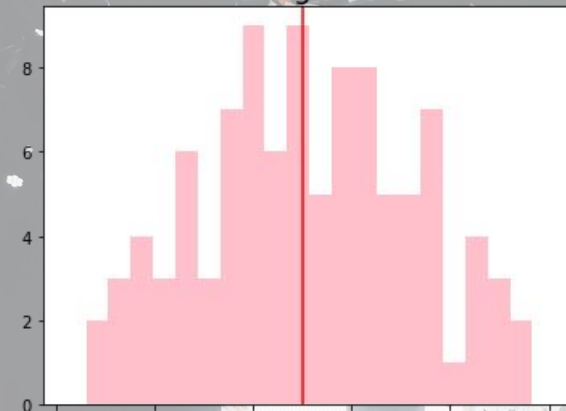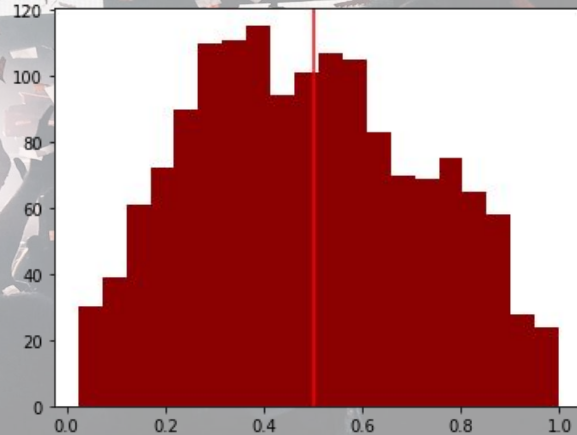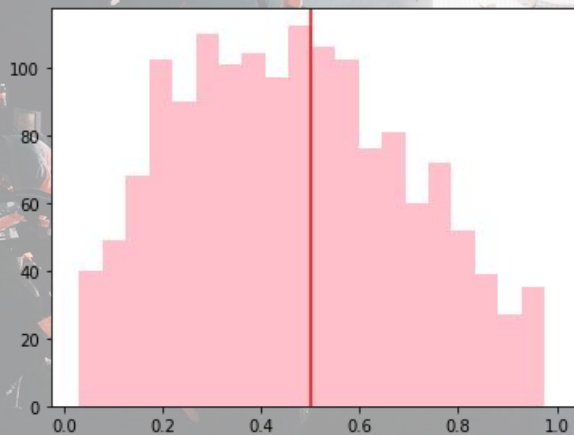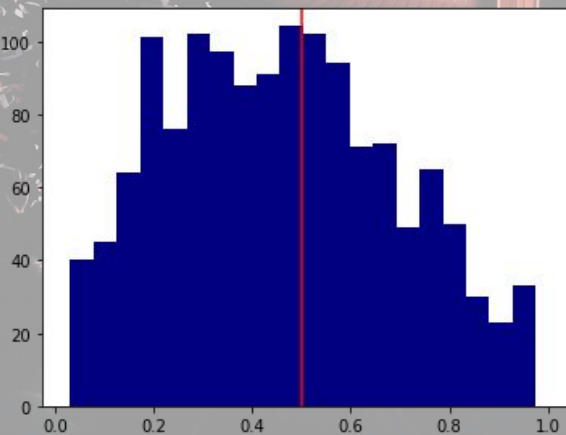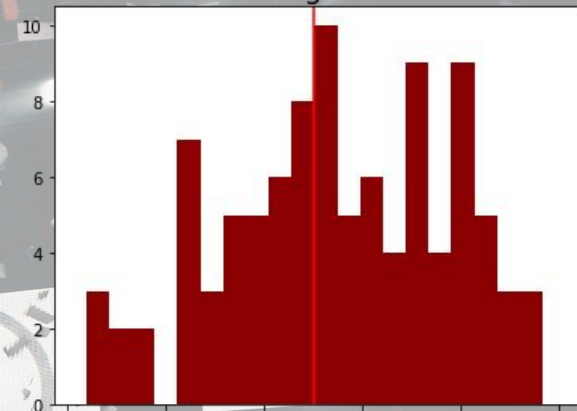# EDA - valence

2018 song's valence | 2013 song's valence | 2008 song's valence
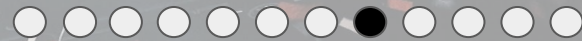
# Modeling - Song Lyrics

Model 1: Neural Network

Model 2: Logistic Regression

# Modeling - Audio Features

Model: Logistic Regression

# Results - Logistic Regression

| | Audio Features | Lyrics |
|---|---|---|
| Accuracy Score | 0.628 | 0.939 |
| Recall Score | 0.623 | 0.08 |

# Results - Confusion Matrices

## Audio Features

|  | predicted not hit | predicted hit |
|---|---|---|
| actual not hit | 449 | 265 |
| actual hit | 26 | 43 |

## Lyrics

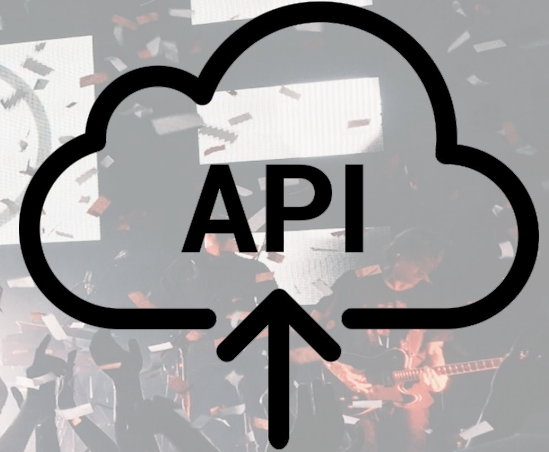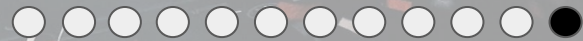|  | predicted not hit | predicted hit |
|---|---|---|
| actual not hit | 1472 | 27 |
| actual hit | 69 | 6 |

# Findings

# Recommendation

Invest in creating louder songs that
have upbeat tempos and strong
beats that are ideal for dancing!

# Obstacles - Data Collection

# Next Steps

- Combine audio features and lyrics into one model

- Implement clustering algorithm on lyrics

- Fill in the gaps of years of analysis

- Run project on AWS