General guidelines when writing chart hits for the 21st Century

By: Darrin Cates, Kimberly Erin, and Rachel Pontow
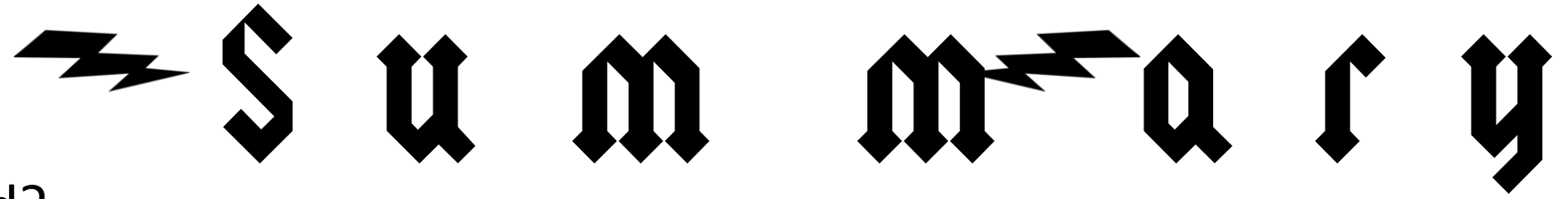
⚡ Is there a formula for creating a hit song?

⚡ Are there common components that chart-topping songs have?

⚡ Has this changed over time? If so how?

⚡ How does the positivity/negativity of lyrics affect popularity?

⚡ How does word repetition affect popularity?

⚡ How does the beats per minute affect popularity?

⚡ How does acousticness, duration, energy, instrumentalness, loudness, and valence affect popularity?

⚡ How have these changed over time?

# Summary

⚡ What did we find?

    ⚡ Acoustic songs don't tend to be as popular

    ⚡ Negative lyrics, while around 25% in popularity, are becoming more popular

    ⚡ Top songs are increasingly using repetition

    ⚡ Song length is at about 2.5 minutes to 5 minutes with few exception after 2010

    ⚡ Songs that have high energy and you can dance to have will do better than songs that you can't dance to

    ⚡ Instrumentals don't do very well on the chart toppers-best to stay away

    ⚡ Keep you music LOUD

    ⚡ Chart toppers can be happy or sad sounding. The sound seems to have no impact on chart toppers these days

    ⚡ Average 120 BPM

⚡ Has this changed over time? If so how?

# Data

- ⚡ Song Source:
  - ⚡ Billboard 1964-2015 Songs + Lyrics csv from Kaggle.com
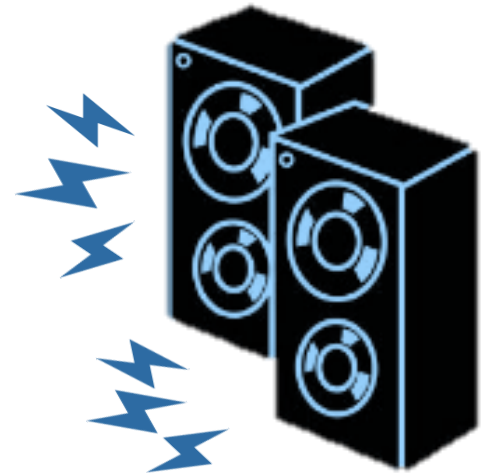
- ⚡ Sentiment Analysis:
  - ⚡ VADER Sentiment Analysis

- ⚡ Beats Per Minute:
  - ⚡ Get Song Bpm API Beats Per Minute:

- ⚡ Spotify Attributes (Spotify API):
  - ⚡ Acousticness, Danceability, Duration, Energy, Instrumentalness, Loudness, Valence

# Data

- ⚡ **The Process**
  - ⚡ Asked ourselves what we wanted to learn
  - ⚡ Identified the data sources
  - ⚡ Worked through a retrieval plan
  - ⚡ Cleaned the data
  - ⚡ Analyzed trends

- ⚡ **The Struggles**
  - ⚡ Complicated API's (Spotify)
  - ⚡ Unclear Documentation (Get Song BPM)
  - ⚡ Slowness of API calls

- ⚡ **Insights**
  - ⚡ Use time delays to make API calls. 5 seconds in a delay, makes a world of difference

# Sentiment

```python
# Loop through lyrics
for lyric in lyric_noblanks["Lyrics"]:

    word_counter_list = []

    # Run Vader Analysis
    results = analyzer.polarity_scores(lyric)

    # Store Vader Analysis Results
    compound = results["compound"]
    pos = results["pos"]
    neu = results["neu"]
    neg = results["neg"]

    # Append Vader Analysis Results
    compound_list.append(compound)
    positive_list.append(pos)
    negative_list.append(neg)
    neutral_list.append(neu)
```
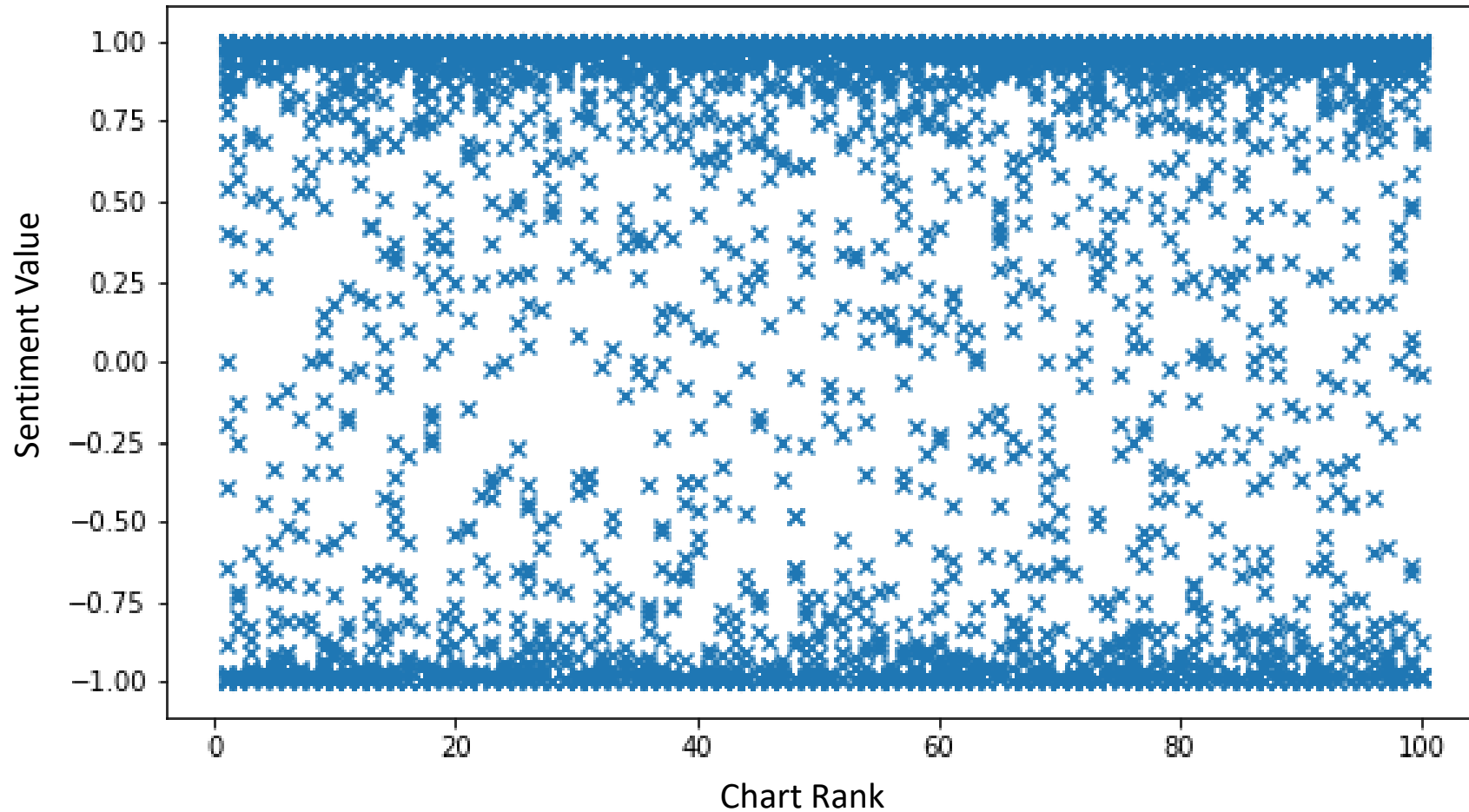
# Sentiment

## Sentiment Vs Chart Ranking



⚡ Positive: > 0.75
⚡ Negative: <-0.75

# Sentiment

```python
# Analysis for which is more popular, positive or negative songs
# Positive is >0.75 and Negative is <-0.75
bins = [0, 1970, 1980, 1990, 2000, 2010, 2016]
group_names = ['<1970', '1970-1979', '1980-1989', '1990-1999', '2000-2009', '2010+']

sentiment_analysis_decade = lyric_filter.copy()
sentiment_analysis_decade["Decade"] = pd.cut(sentiment_analysis_decade["Year"], bins, labels=group_names)

pos_songs = []
neg_songs = []
tot_songs = []

for x in group_names:
    decade_count = sentiment_analysis_decade.loc[sentiment_analysis_decade["Decade"] == x]

    count_pos = decade_count.loc[decade_count["Sentiment"] > 0.75]
    pos_songs.append(len(count_pos["Sentiment"]))

    count_neg = decade_count.loc[decade_count["Sentiment"] < -0.75]
    neg_songs.append(len(count_neg["Sentiment"]))

    tot_songs.append(len(decade_count["Sentiment"]))

percent_pos = [a/b*100 for a,b in zip(pos_songs,tot_songs)]
percent_neg = [a/b*100 for a,b in zip(neg_songs,tot_songs)]
percent_neu = [100-a-b for a,b in zip(percent_pos,percent_neg)]

pos_analysis_df = pd.DataFrame({"Decade": group_names, "Percent Positive": percent_pos, "Percent Negative"
: percent_neg, "Percent Neutral": percent_neu})
```
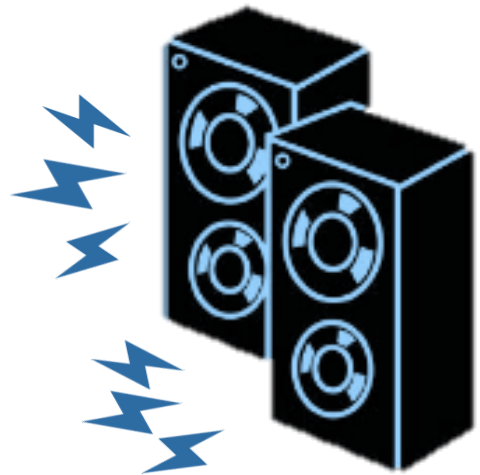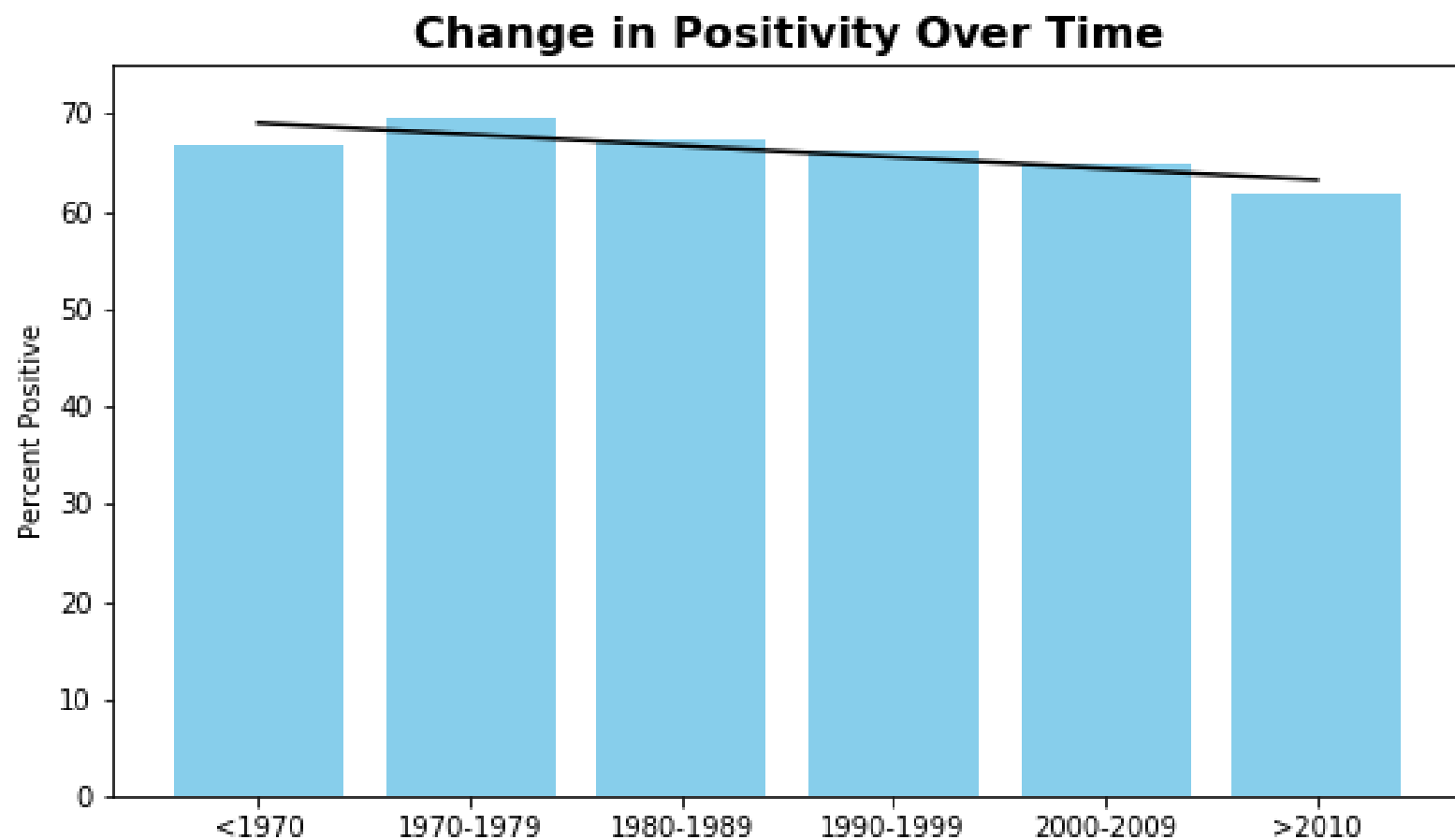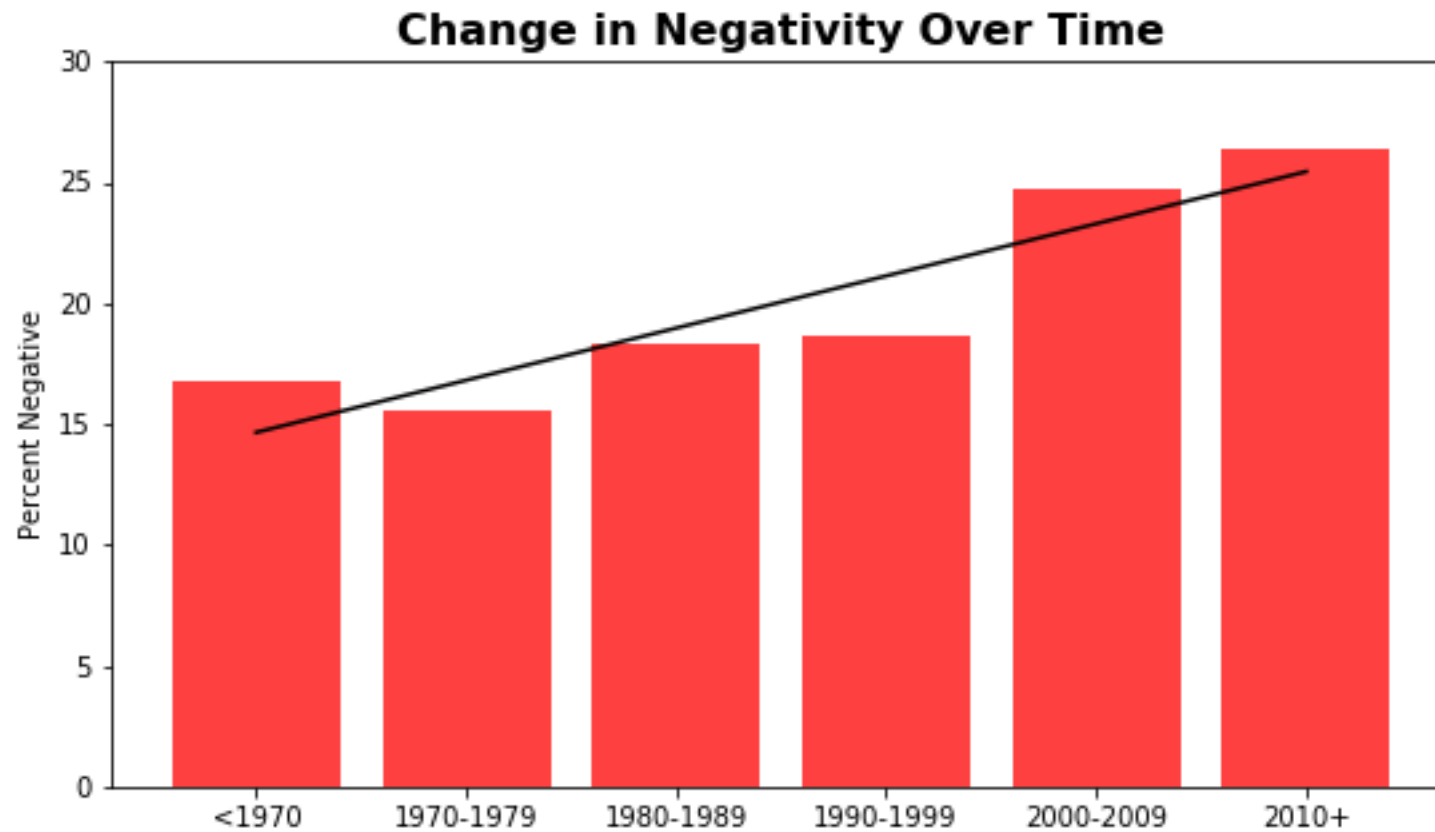
# Positivity

## Change in Positivity Over Time

# Negativity



Change in Negativity Over Time

# Repetition

```python
# Word Repetition
total_number_words = len(word_split)
unique_words = set(word_split)
unique_filtered = [word for word in unique_words if word not in stop_words]
unique_words = list(unique_words)
word_counter_list = []

for x in unique_filtered:
    word_counter = 0
    for y in word_split:
        if x==y:
            word_counter += 1
    word_counter_list.append(word_counter)

mode = max(word_counter_list)
mode_index = word_counter_list.index(max(word_counter_list))
mode_word = unique_filtered[mode_index]
word_count = len(unique_filtered)

word_repetition.append(mode)
word_repeated.append(mode_word)
word_counts.append(unique_filtered)
unique_word_count.append(word_count)
total_word_count.append(total_number_words)
```
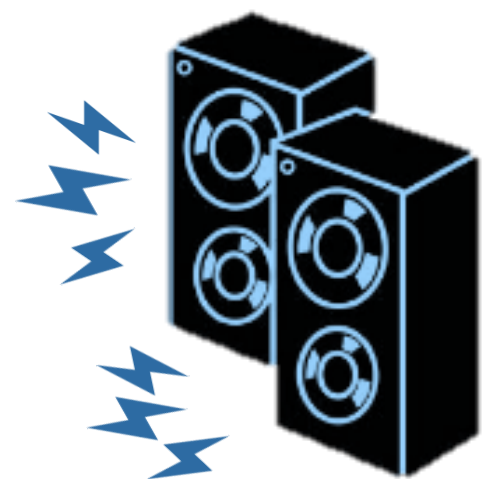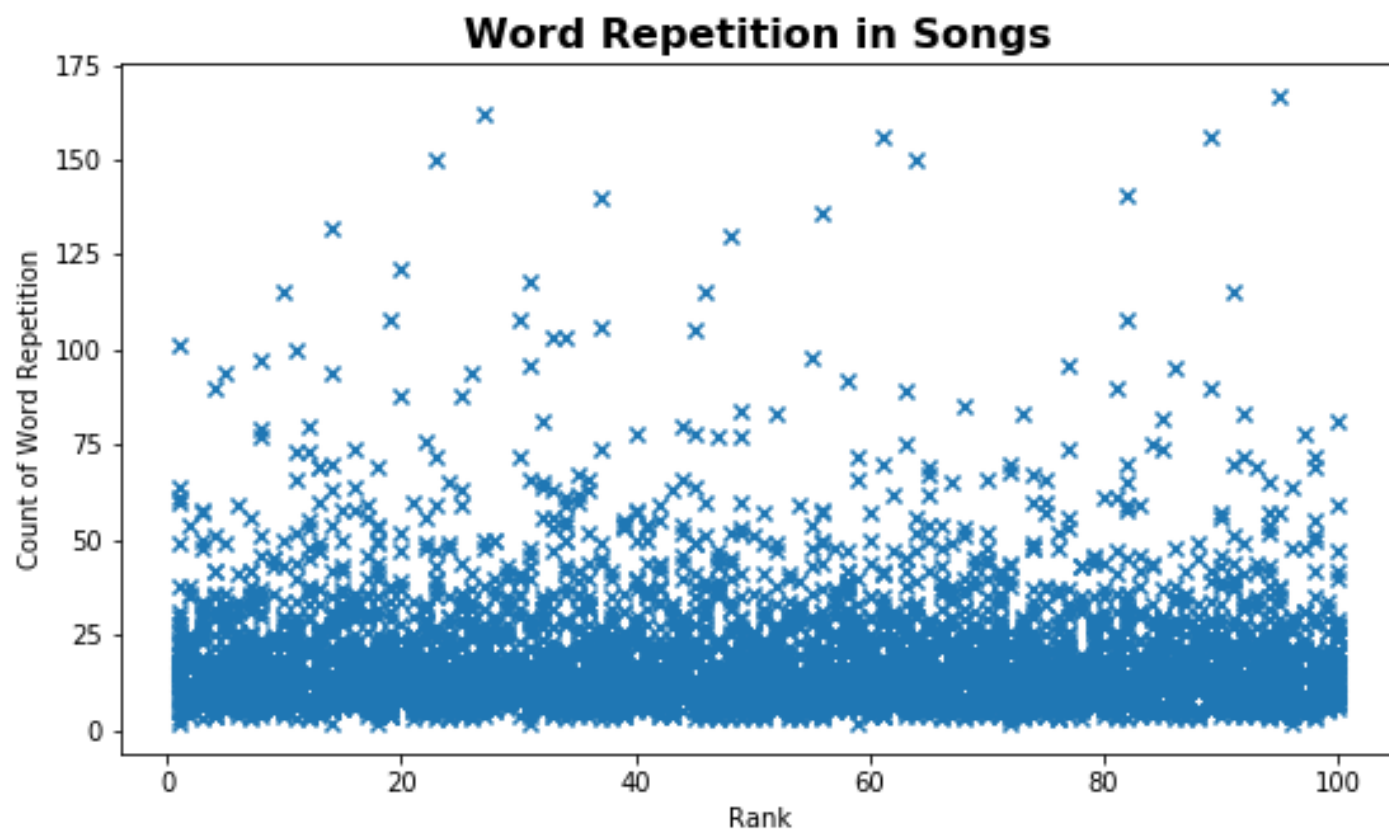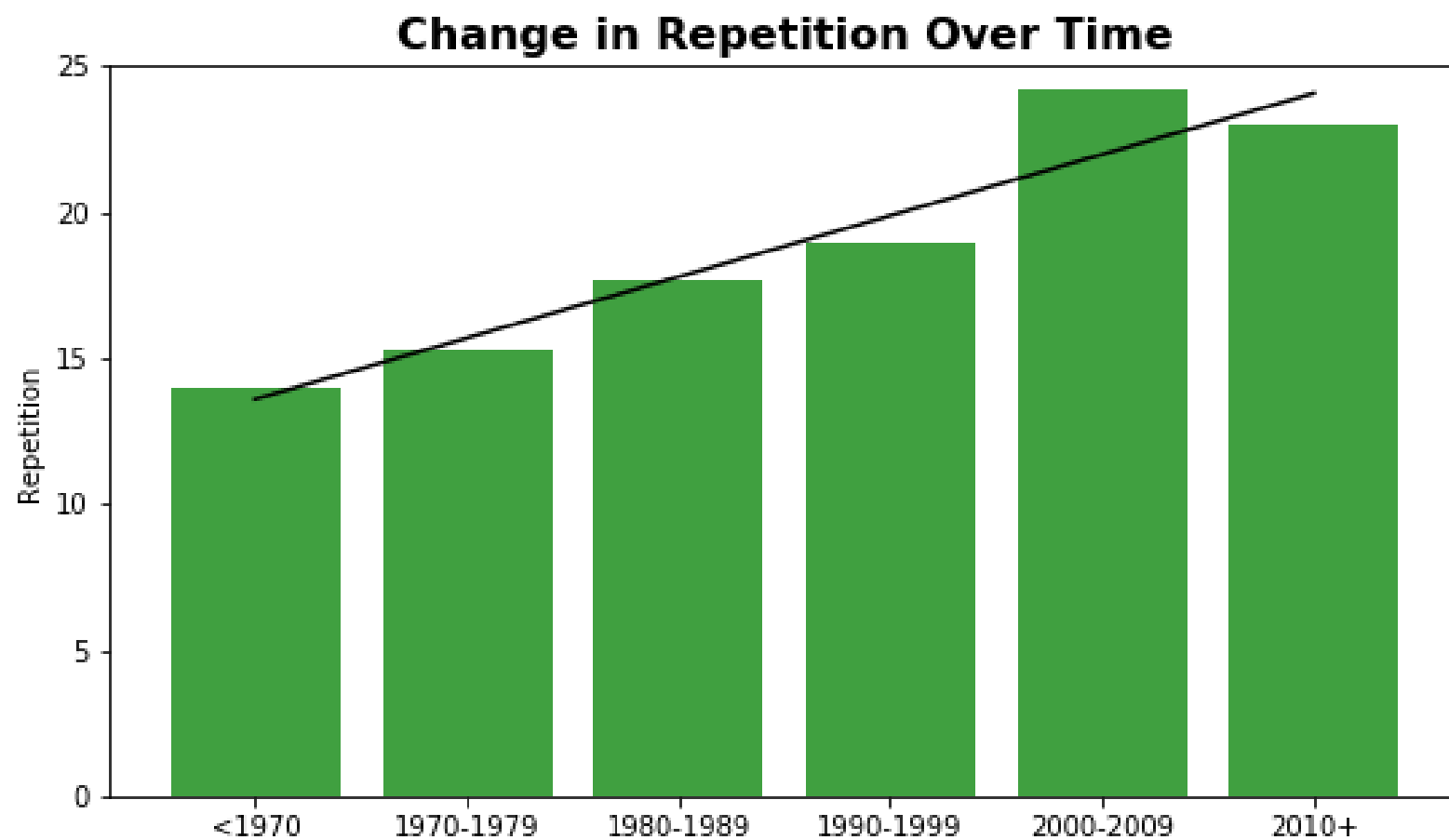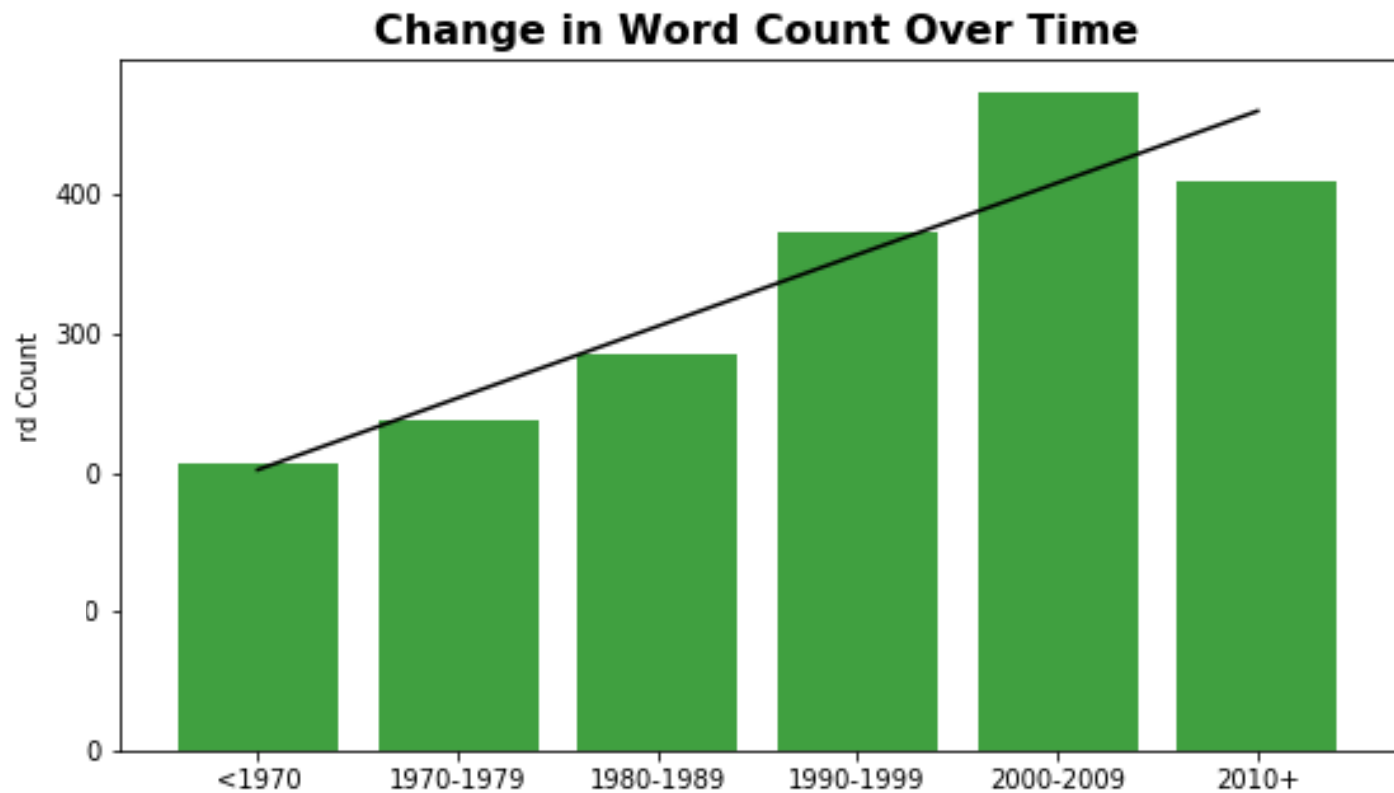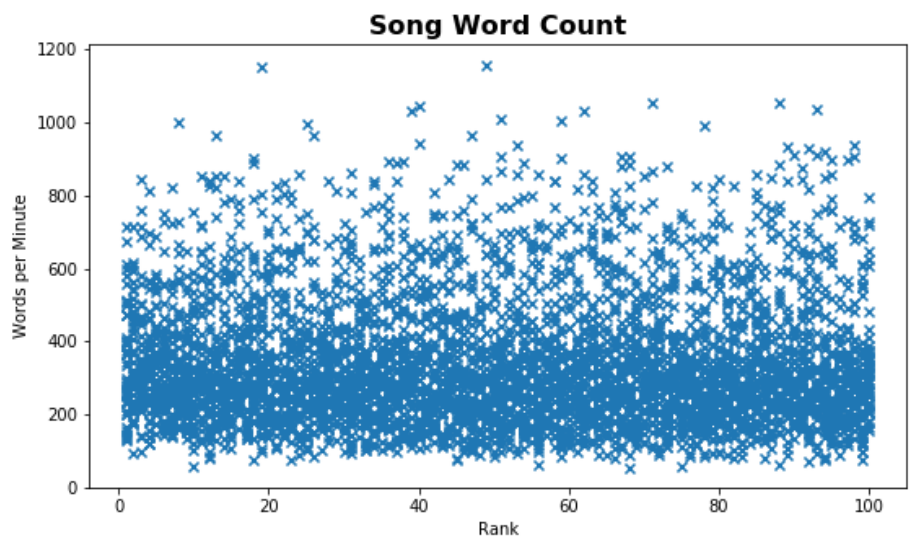
# Repetition



Word Repetition in Songs

# Repetition

## Change in Repetition Over Time

# Word

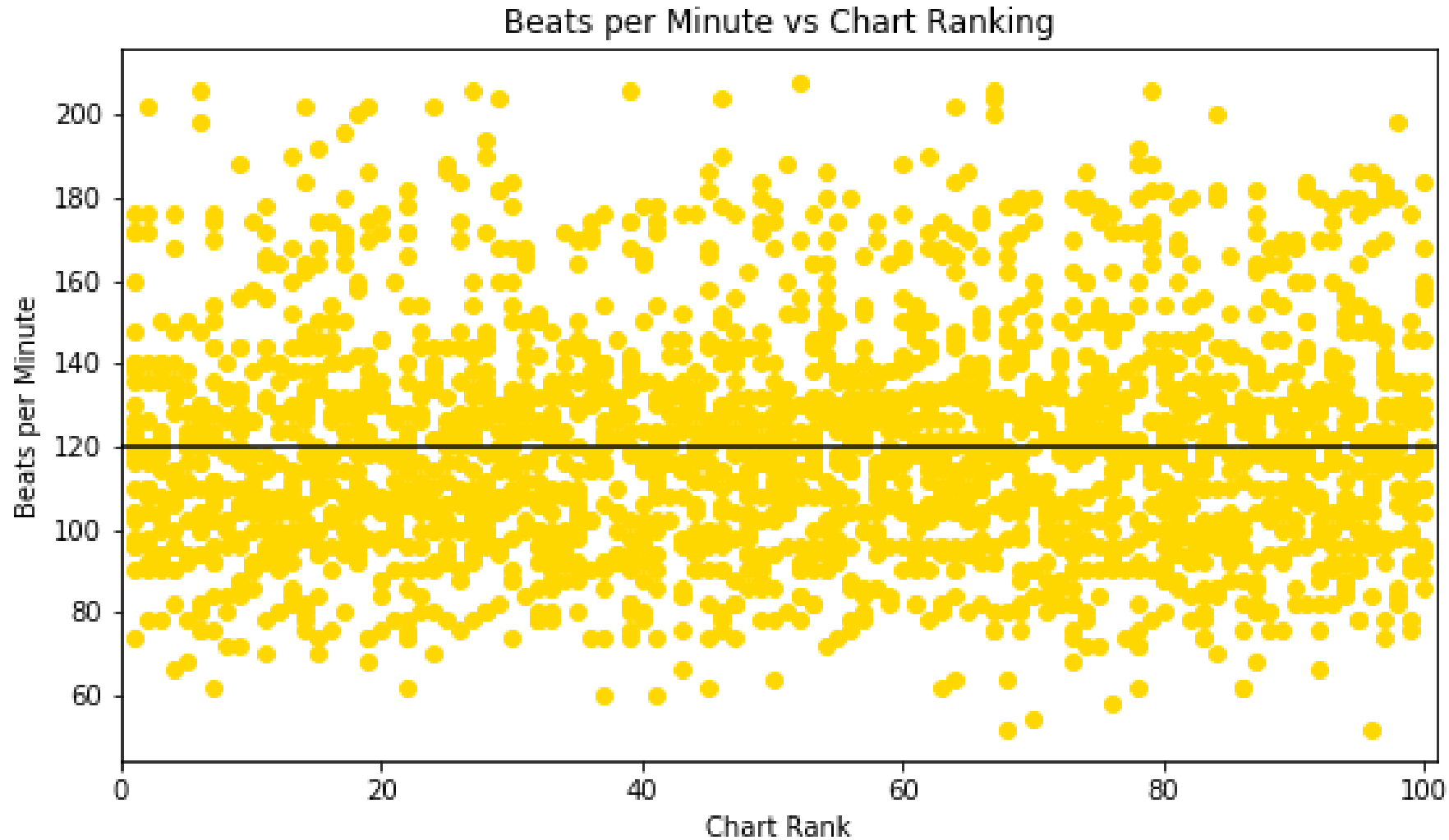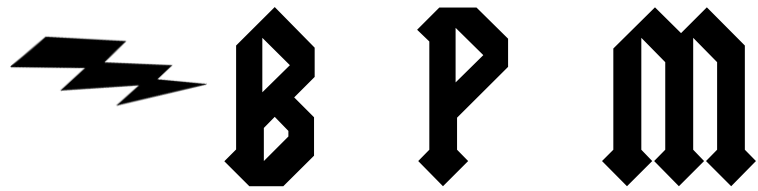## Change in Word Count Over Time



## Song Word Count

# Beats

```python
#Loop through and get beats per minute
from bpm_config import api_key
import requests
bpm = []
#tracks_per_year = lyric_noblanks.loc[lyric_noblanks["Year"]==2015]
for index, row in lyric_noblanks.iterrows():
    base_url = "https://api.getsongbpm.com/search/?"
    track = row["Song"]
    artist = row["Artist"]
    info = requests.get(f"{base_url}api_key={api_key}&type=both&lookup=song:{track} artist:{artist}").json
()
    bpm.append(info)
print(bpm)
```
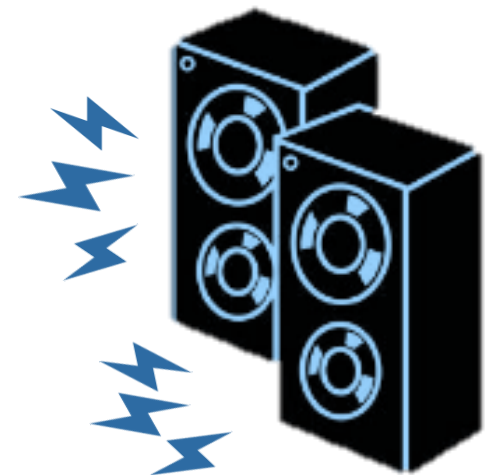
```python
bpm_list = []
for track in bpm:
    try:
        bpm_list.append(track["search"][0]["tempo"])
    except (KeyError):
        bpm_list.append("N/A")
```

```python
#Remove blank rows
clean_bpm_df = bpm_df.replace(" ","NaN")
clean_bpm_df = bpm_df.dropna(subset=["BPM"])
clean_bpm_df = clean_bpm_df.reset_index(drop=True)
clean_bpm_df = clean_bpm_df[["Rank","Song","Artist","Year","Lyrics", "Source", "BPM"]]
clean_bpm_df.head()
#clean_bpm_df.count()
```
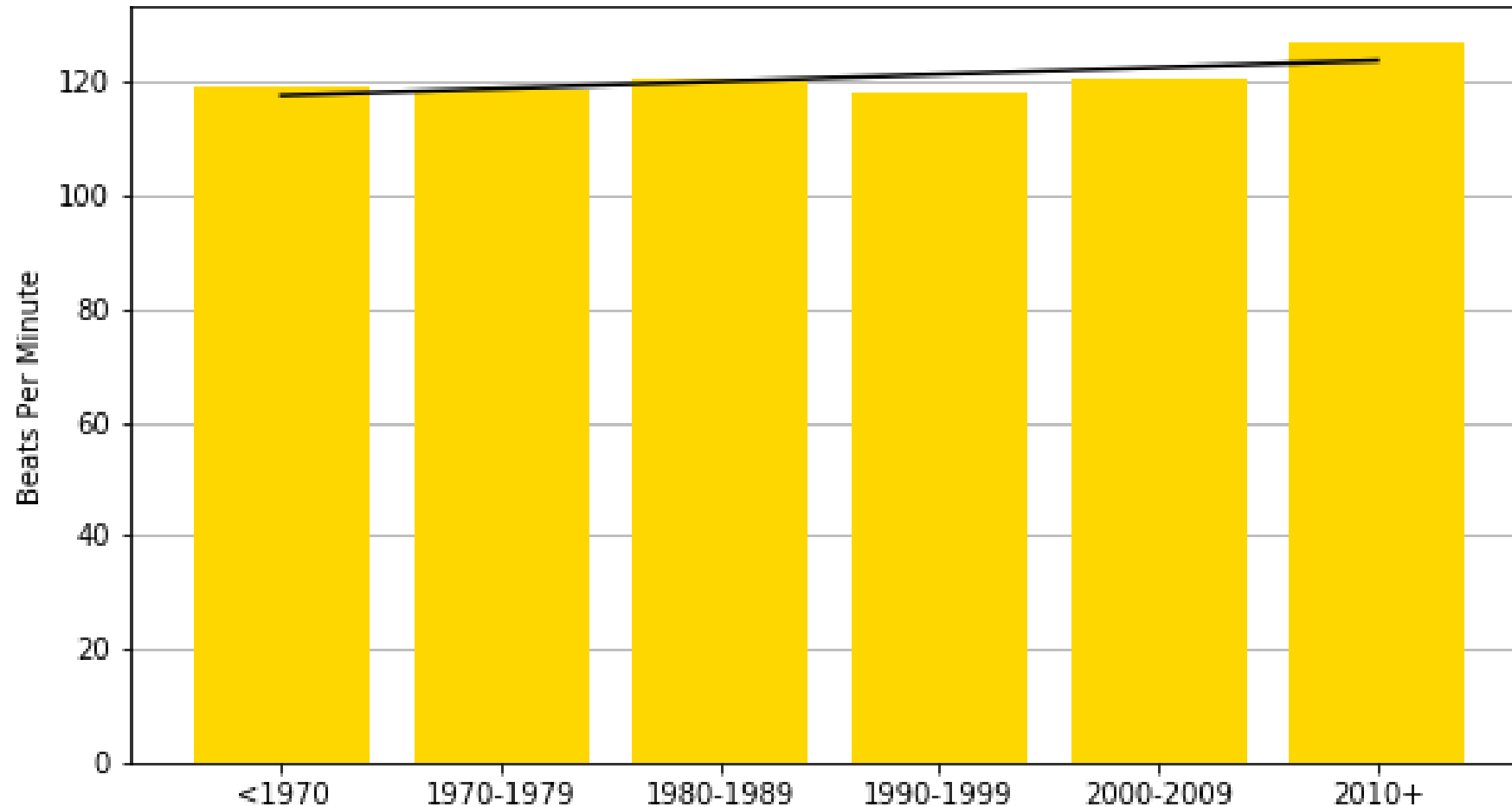
# BPM

## Change in Beats Per Minute by Decade



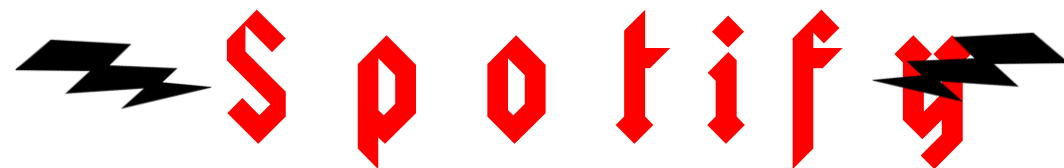| Decade | Average BPM |
|---|---|
| <1970 | 119.042373 |
| 1970-1979 | 118.607375 |
| 1980-1989 | 120.415771 |
| 1990-1999 | 118.209446 |
| 2000-2009 | 120.700219 |
| 2010+ | 126.850000 |

# Spotify

The quest for song ID's

```
1   #commenting out, so it's not run again
2   #looking up spotify id's
3
4   #creating lists to store data
5   #song_name = []
6   #artist = []
7   #spotify_id = []
8
9   #looping through and pulling all songs that match the song name in the data set
10  #for id_lookup in song_list:
11      #try:
12  #       time.sleep(10)
13
14  #       name = id_lookup
15
16  #       spotify = spotipy.Spotify(client_credentials_manager=client_credentials_manager)
17  #       results1 = spotify.search(q='track:' + name, type='track')
18  #       results1
19  #       count = 0
20
21  #       for x in np.arange(len(results1["tracks"]["items"])):
22
23  #           spotify = spotipy.Spotify(client_credentials_manager=client_credentials_manager)
24  #           results = spotify.search(q='track:' + name, type='track')
25
26  #           song_name.append(results1["tracks"]["items"][count]["name"])
27
28  #           artist.append(results1["tracks"]["items"][count]["album"]["artists"][0]["name"])
29
30  #           spotify_id.append(results1["tracks"]["items"][count]["id"])
31  #           count += 1
32      #except IndexError:
33       #   next
34  #creating new dataframe
35      #df['song name']  = song_name
36      #df['artist']  = artist
37      #df['id']  = spotify_id
38  #print(f'song {song_name} artist {artist} id {spotify_id}')
```

⚡ Was unable to call Spotify to look up only 1 specific ID with Artist and song

⚡ Made API call with the song name and pulled all track ID's for future lookups

⚡ Took all songs and pulled out only the songs that matched both the song name and artist name from original dataset and appended the Spotify ID
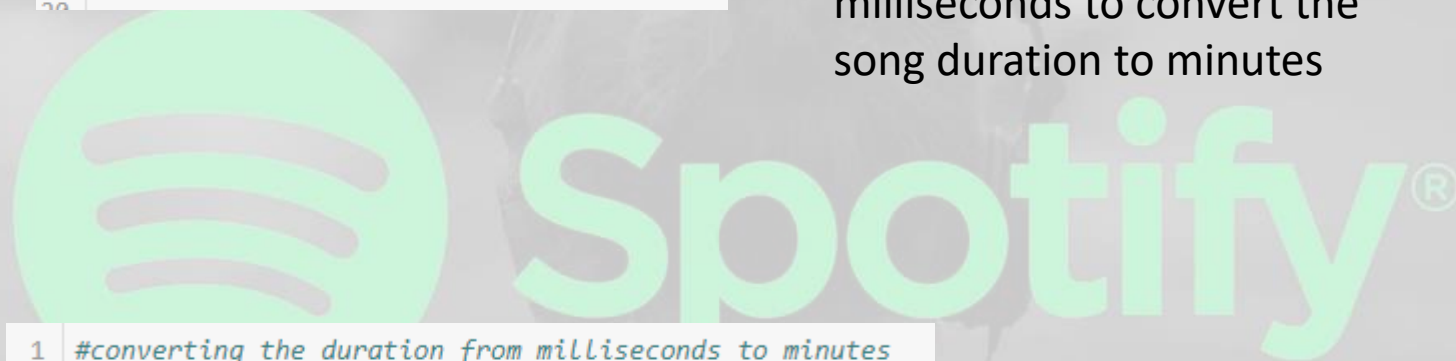
```python
1   #credentials needed for api call
2   from spotify.oauth2 import SpotifyClientCredentials
3   client_credentials_manager = SpotifyClientCredentials(client_id=token, client_secret=secret)
4
5   #pulling only first 100 records to keep the api call from getting to many requests at once
6   x = 0
7   y = 99
8   id10 = spotify_ids[x:y]
9   idsong = songnames[x:y]
10  idartist = artistnames[x:y]
11  idrank = rank[x:y]
12  idyear = year[x:y]
13
14  #creating lists to store data retrieval
15  sng_nam = []
16  art_nam = []
17  rnk = []
18  yr = []
19  acousticness = []
20  danceability = []
21  duration_ms = []
22  energy = []
23  instrumentalness = []
24  key = []
25  liveness = []
26  loudness = []
27  mode = []
28  speechiness = []
29  tempo = []
30  time_signature = []
31  valence = []
32
33  #creating counter for looping through records
34  cnt = 0
35  for id_lookup in id10:
36      #credentials for api call
37      spotify = spotipy.Spotify(client_credentials_manager=client_credentials_manager)
38      results = spotify.audio_features(id10)
39
40      #appending results into the lists created
41      acousticness.append(results[cnt]["acousticness"])
42      danceability.append(results[cnt]["danceability"])
43      duration_ms.append(results[cnt]["duration_ms"])
44      energy.append(results[cnt]["energy"])
45      instrumentalness.append(results[cnt]["instrumentalness"])
46      key.append(results[cnt]["key"])
47      liveness.append(results[cnt]["liveness"])
48      loudness.append(results[cnt]["loudness"])
49      mode.append(results[cnt]["mode"])
50      speechiness.append(results[cnt]["speechiness"])
51      tempo.append(results[cnt]["tempo"])
52      time_signature.append(results[cnt]["time_signature"])
53      valence.append(results[cnt]["valence"])
54      sng_nam.append(idsong[cnt])
55      art_nam.append(idartist[cnt])
56      rnk.append(idrank[cnt])
57      yr.append(idyear[cnt])
58
59      #adding to the counter
60      cnt += 1
```

```python
1   #CREATING DATA FRAME FROM SPOTIFY ID'S
2   df = pd.DataFrame()
3   df['song_name'] = sng_nam
4   df['rank'] = rnk
5   df['year'] = yr
6   df['artist_name'] = art_nam
7   df['acousticness']   = acousticness
8   df['danceability']   = danceability
9   df['duration_ms']    = duration_ms
10  df['energy']    = energy
11  df['instrumentalness']   = instrumentalness
12  df['key']   = key
13  df['liveness']   = liveness
14  df['loudness']   = loudness
15  df['mode']   = mode
16  df['speechiness']   = speechiness
17  df['tempo']   = tempo
18  df['time_signature']    = time_signature
19  df['valence']   = valence
20
```

- ⚡ Made another Spotify API Call to collect the attributes for each of the songs in the list
- ⚡ This was ran multiple times to make all of the API calls as only 100 songs could be looked up at once
- ⚡ Created a new dataframe with the new data
- ⚡ Ran a conversion on milliseconds to convert the song duration to minutes

```python
1   #converting the duration from milliseconds to minutes
2   df['duration_minutes']   = df['duration_ms']/60000
```
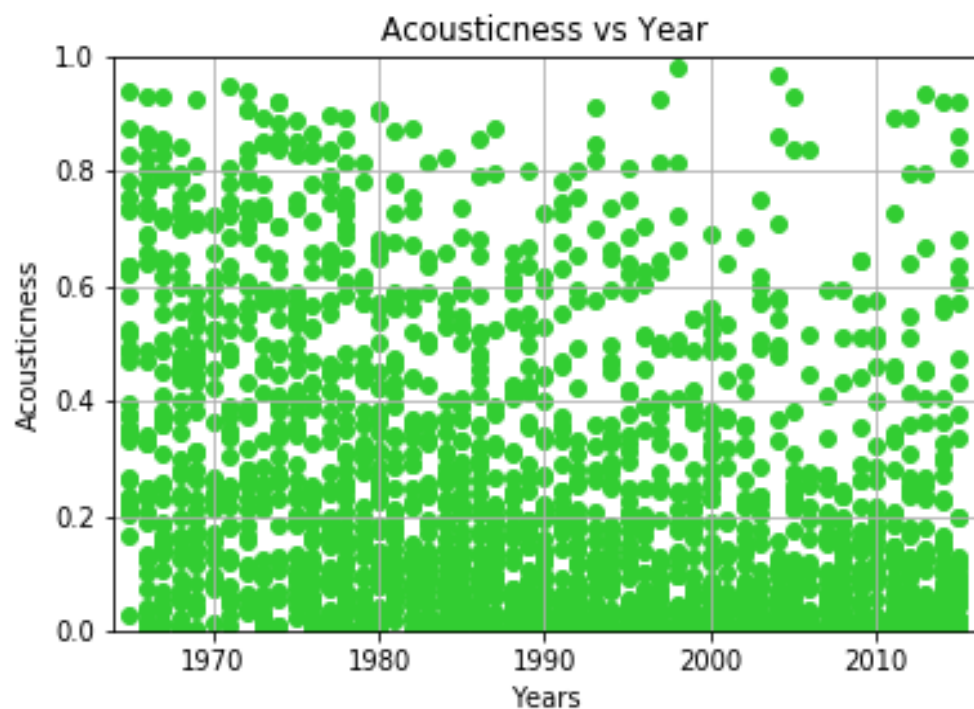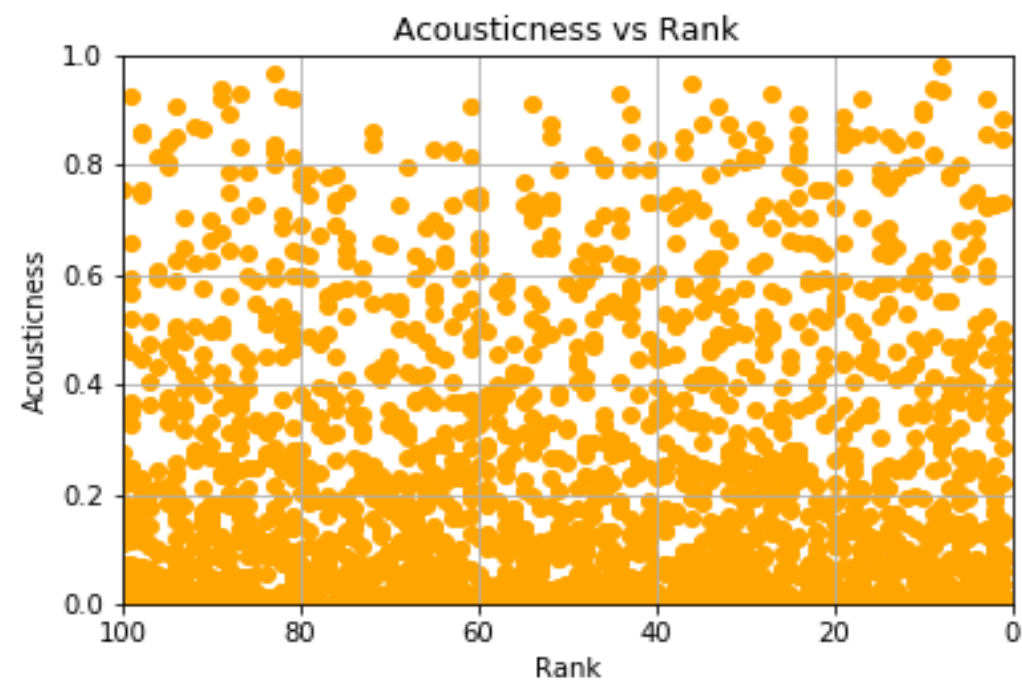
# Acoustitnes

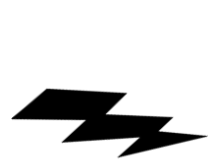**Definition:** A confidence measure from 0.0 to 1.0 of whether the track is acoustic. 1.0 represents high confidence the track is acoustic
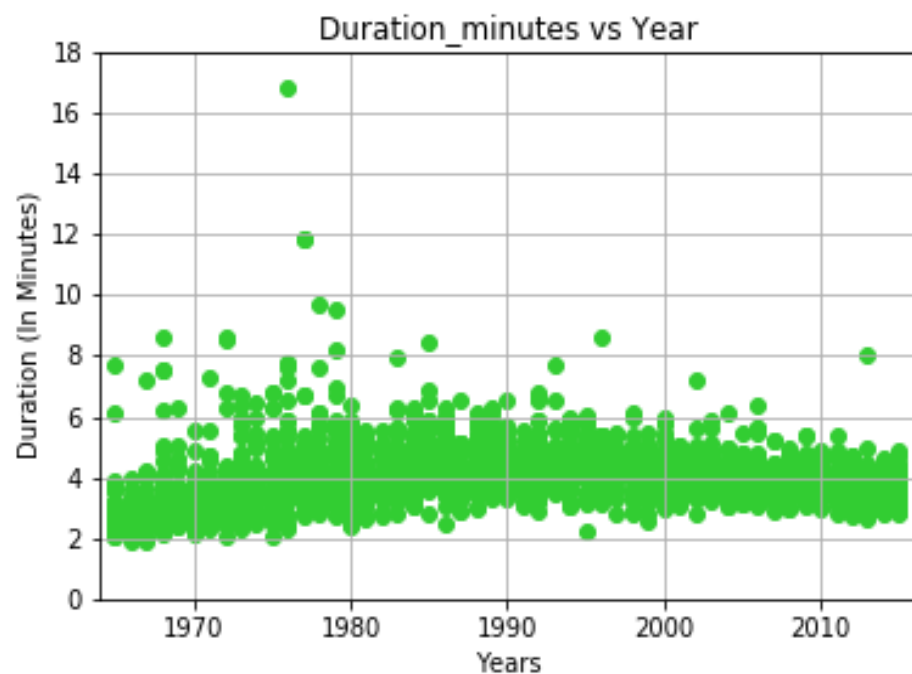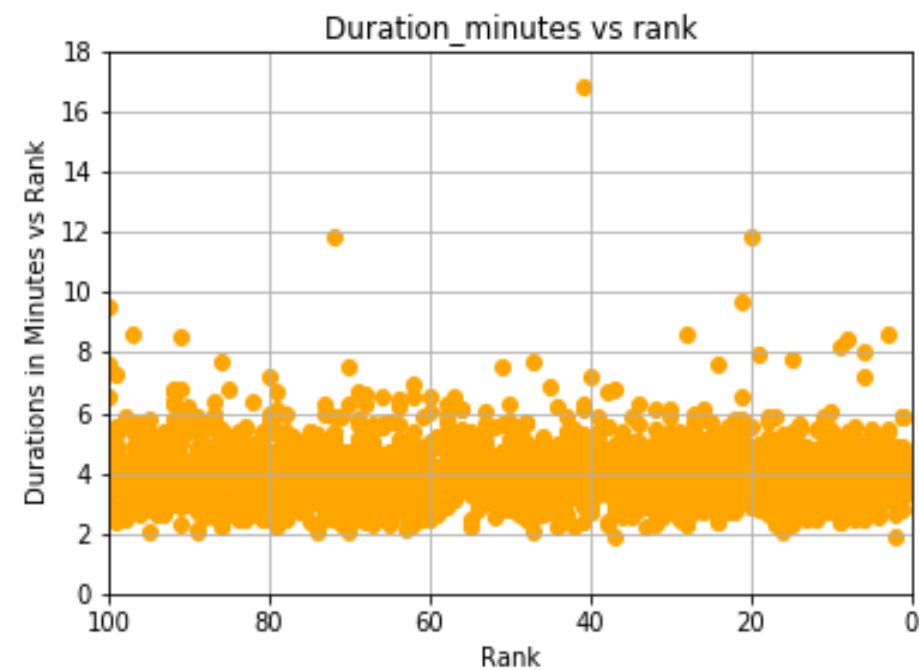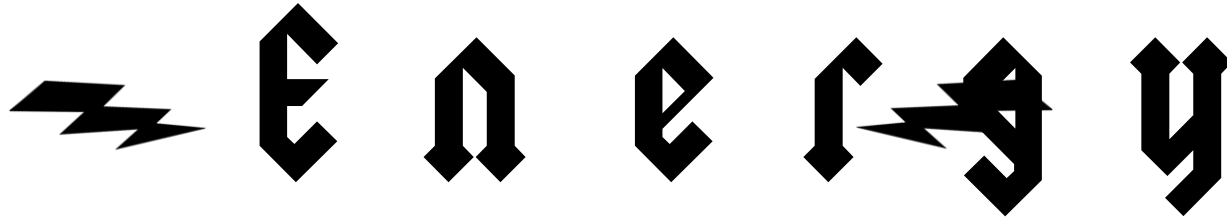
## Year

Acousticness vs Year

## Rank

Acousticness vs Rank

# Duration

**Definition:** Length of Song

## Year



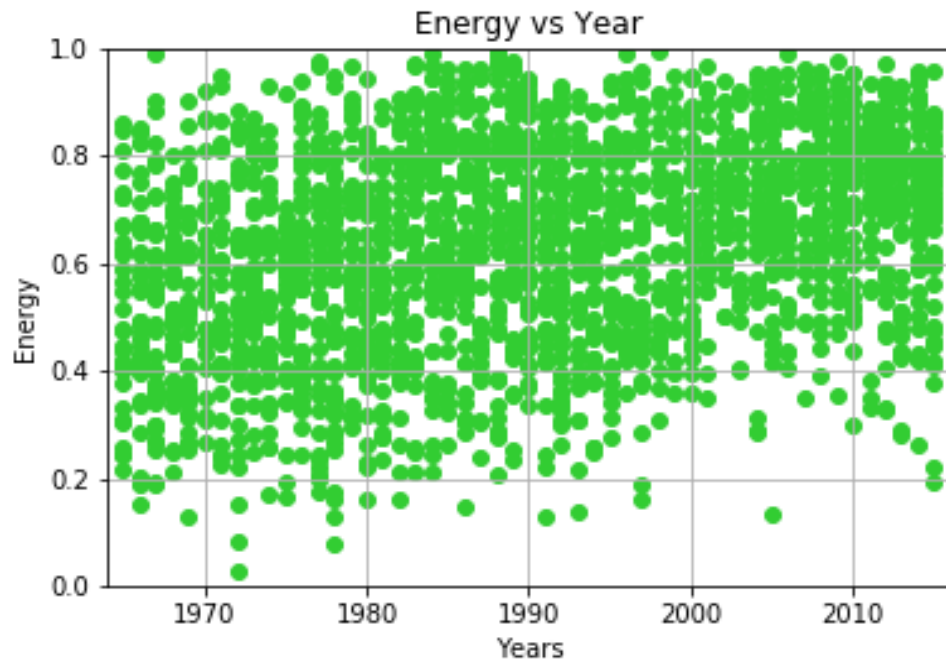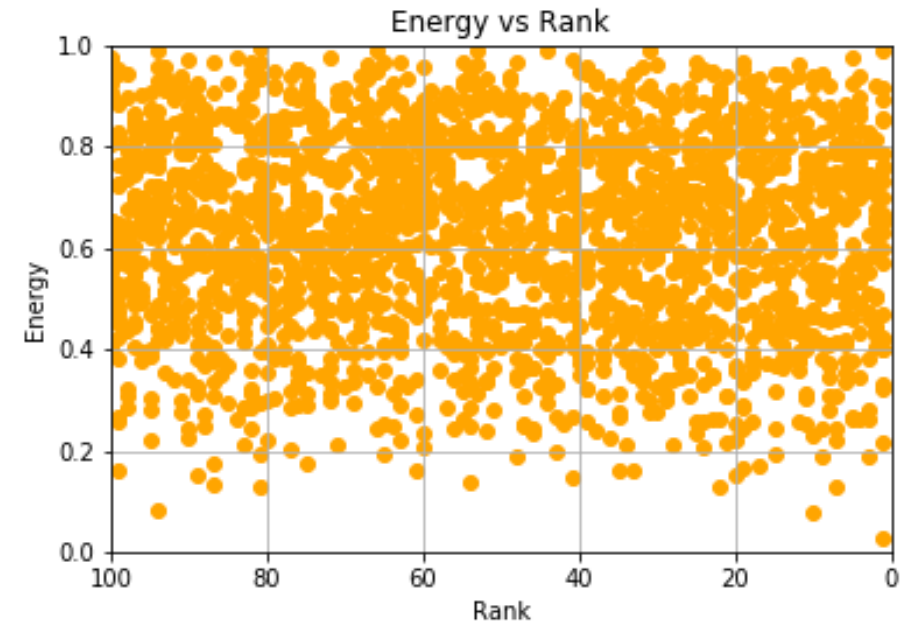Duration_minutes vs Year

## Rank



Duration_minutes vs rank

# Energy

**Definition:** Energy is a measure from 0.0 to 1.0 and represents a perceptual measure of intensity and activity. Typically, energetic tracks feel fast, loud, and noisy. For example, death metal has high energy, while a Bach prelude scores low on the scale. Perceptual features contributing to this attribute include dynamic range, perceived loudness, timbre, onset rate, and general entropy
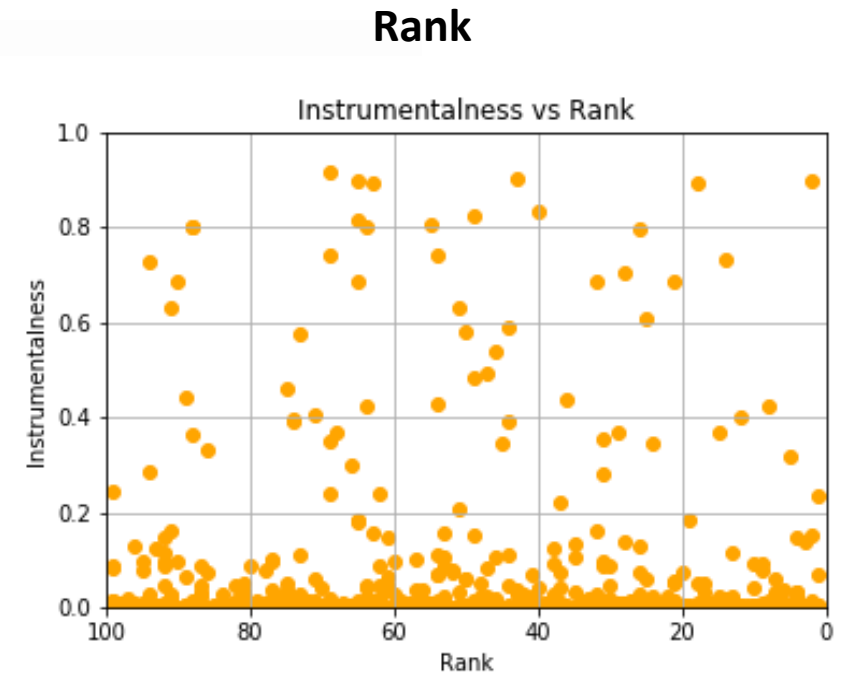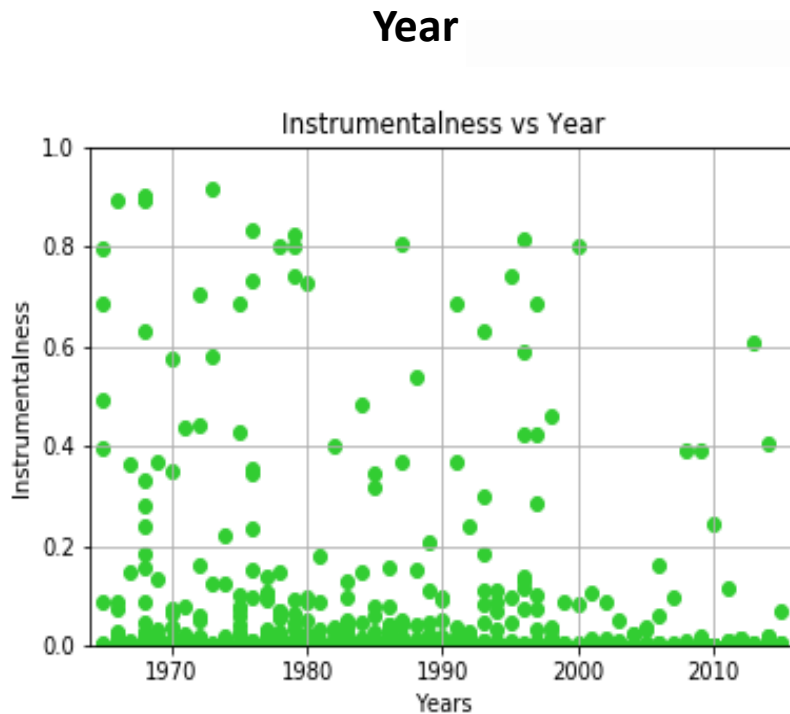
## Year



## Rank

# Instrumenta

**Definition:** Predicts whether a track contains no vocals. "Ooh" and "aah" sounds are treated as instrumental in this context. Rap or spoken word tracks are clearly "vocal". The closer the instrumentalness value is to 1.0, the greater likelihood the track contains no vocal content. Values above 0.5 are intended to represent instrumental tracks, but confidence is higher as the value approaches 1.0.
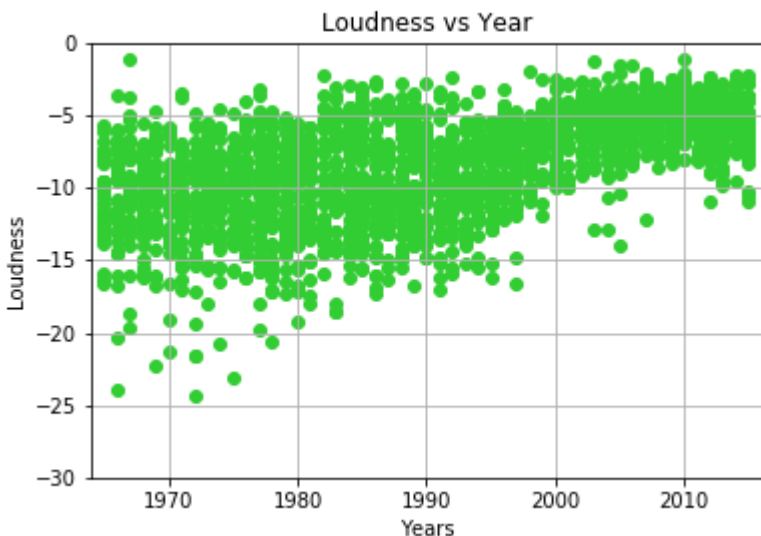
## Year



Instrumentalness vs Year

## Rank



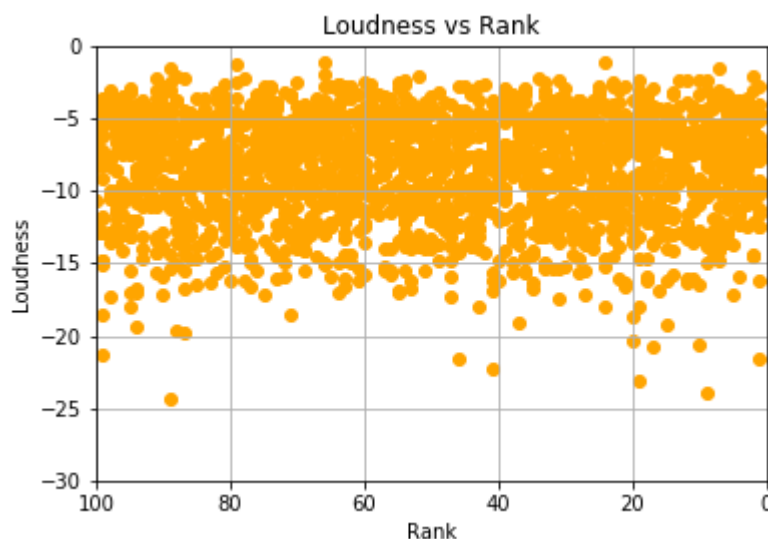Instrumentalness vs Rank

# ⚡ Loudness

**Definition:** The overall loudness of a track in decibels (dB). Loudness values are averaged across the entire track and are useful for comparing relative loudness of tracks. Loudness is the quality of a sound that is the primary psychological correlate of physical strength (amplitude). Values typical range between -60 and 0 db.
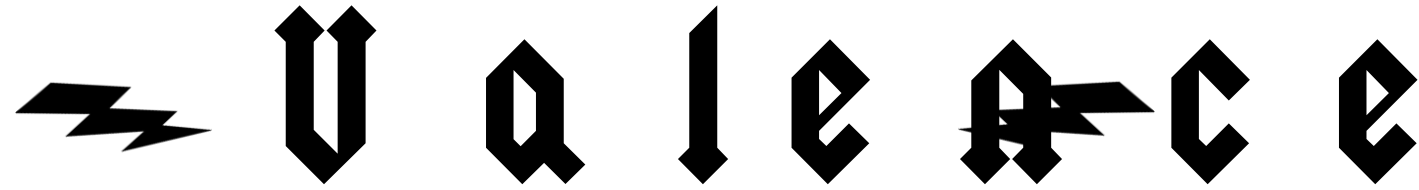
## Year



Loudness vs Year

## Rank



Loudness vs Rank

## LOUDNESS war

The loudness war (or loudness race) refers to the trend of increasing audio levels in recorded music which many critics believe reduces sound quality and listener enjoyment. Increasing loudness was first reported as early as the 1940s, with respect to mastering practices for 7" singles.
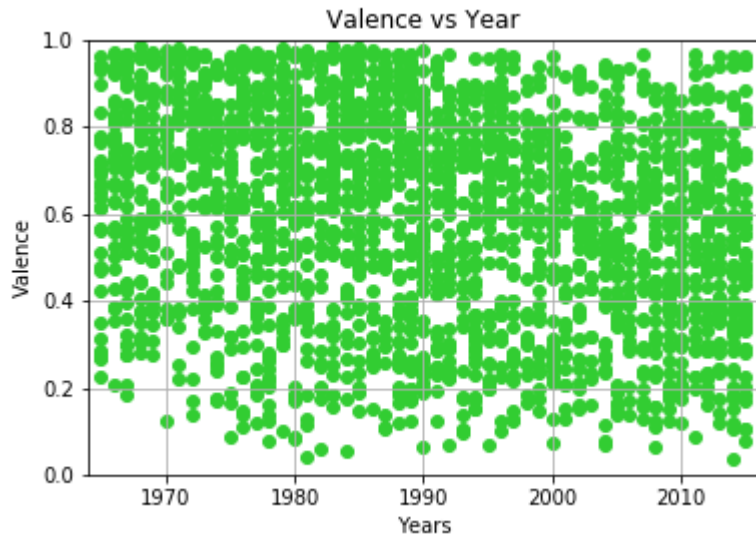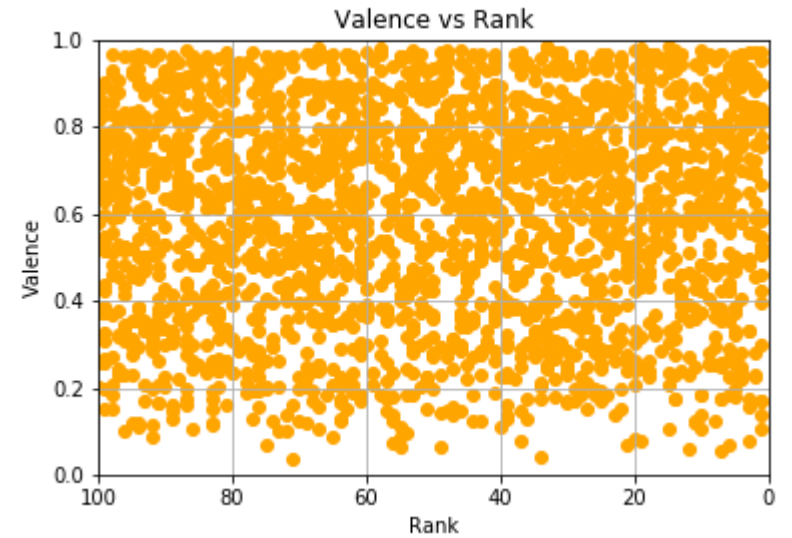https://en.wikipedia.org/wiki/Loudness_war

# Valence

**Definition:** A measure from 0.0 to 1.0 describing the musical positiveness conveyed by a track. Tracks with high valence sound more positive (e.g. happy, cheerful, euphoric), while tracks with low valence sound more negative (e.g. sad, depressed, angry).
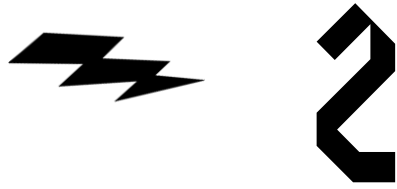
## Year



## Rank

# Conclusion

⚡Keep your song length between 2.5 – 5 minutes

⚡Don't write instrumentals

⚡More likely to have a pop hit with positive words in your song

⚡Keep the beats per minute around 120 bpm

- Break out songs based on Genre and review for trends
- Trend what the guidelines would be for the year 2050
- Complete regression analysis for the Spotify data
- Add in a review of Time Signatures and Key Signatures
- Analyze the song structures
- Address the record label, producer, artist

# Questions?