



# Using Linear Regression to Predict the Streams of a Song on Spotify

Prepared by: Haneen Alhomoud



# Problem Description



Spotify

- digital music, podcast, and video service that gives you access to millions of songs and other content.
- 365 million active users




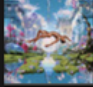
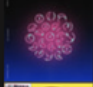





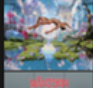



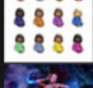


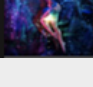
## Scope

- 2016-12-31 - 2021-09-25
- 6 years
- 1200 data point



# Spotify Charts Dataset



TRACK		STREAMS ?
	1 – <b>STAY (with Justin Bieber)</b> by The Kid LAROI	8,702,859
	2 – <b>INDUSTRY BABY (feat. Jack Harlow)</b> by Lil Nas X	7,359,778
	3 – <b>My Universe</b> by Coldplay, BTS	4,897,444
	4 ▲ <b>Pepas</b> by Farruko	4,870,261
	5 ▼ <b>Heat Waves</b> by Glass Animals	4,595,016
	6 – <b>Bad Habits</b> by Ed Sheeran	4,573,991
	7 – <b>Woman</b> by Doja Cat	4,325,652
	8 ▲ <b>Shivers</b> by Ed Sheeran	3,908,106
	9 ▼ <b>THATS WHAT I WANT</b> by Lil Nas X	3,902,738
	10 – <b>Beggin'</b> by Måneskin	3,861,687
	11 – <b>MONTERO (Call Me By Your Name)</b> by Lil Nas X	3,657,042
	12 – <b>good 4 u</b> by Olivia Rodrigo	3,586,187
	13 – <b>Way 2 Sexy (with Future &amp; Young Thug)</b> by Drake	3,235,023
	14 – <b>Need To Know</b> by Doja Cat	3,110,387
	15 – <b>Cold Heart - PNAU Remix</b> by Elton John, Dua Lipa	3,009,321
	16 ▲ <b>Kiss Me More (feat. SZA)</b> by Doja Cat	2,809,990

- 1200 Song Scraped

- Features

- Song name

- Artist name

- Popularity

- Target

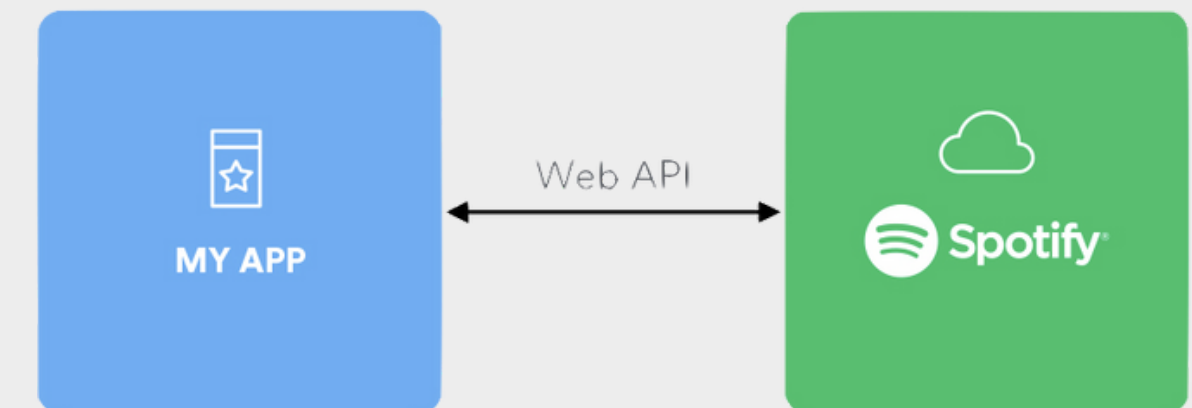
- Streams

# Audio Features Using Spotify API Dataset



## • Features

- **Danceability:** how suitable a song is for dancing
- **Energy:** how energetic tracks feel fast, loud, and noisy.
- **Loudness:** overall loudness of a track in decibels (dB)
- **Speechiness:** Detects the presence of spoken words in a track
- **Acoustics:** tells whether the track is acoustic or not
- **Instrumentals:** Predicts whether a track contains vocal
- **Liveness:** Presence of an audience in the recording
- **Valence:** Positiveness conveyed by a track
- **Tempo:** Estimated tempo of a track in beats per minute (BPM)
- **duration:** music duration time



# Artist Rank Dataset



- Feature
  - Artist Rank

Pos	Artist	Total Streams
1	Drake	21,277,495,267
2	Bad Bunny	17,628,242,359
3	J Balvin	17,534,317,610
4	Justin Bieber	15,197,896,200
5	Post Malone	14,295,061,426
6	Ozuna	12,467,100,127
7	Ed Sheeran	11,952,846,730
8	The Weeknd	10,728,857,095
9	Ariana Grande	10,634,032,845
10	Khalid	9,889,653,507
11	Billie Eilish	9,535,416,511
12	Juice WRLD	9,375,460,193
13	Dua Lipa	9,296,883,619
14	Daddy Yankee	9,290,520,482
15	Travis Scott	8,885,765,334
16	Maluma	8,462,789,141
17	Anuel Aa	7,367,301,612
18	XXXTENTACION	7,297,900,502
19	Cardi B	7,183,374,176
20	Nicky Jam	6,891,721,152
21	Farruko	6,328,794,127

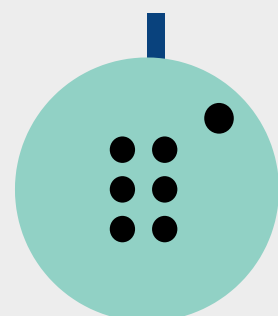


## 12 Features

- popularity
- artist\_rank
- danceability
- energy
- loudness
- speechiness
- acoustics
- instrumentals
- liveness
- valence
- tempo
- duration

## Target

- streams



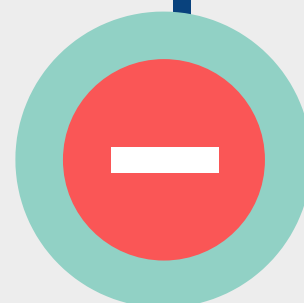
## Outliers

Target: Streams



## Duplicated rows

The same song is on top 200 list more than once



## Null Values

None

# Modeling - Linear Regression



# Baseline Model



**Training Score: 0.517**



**Validation Mean Score: 0.524**

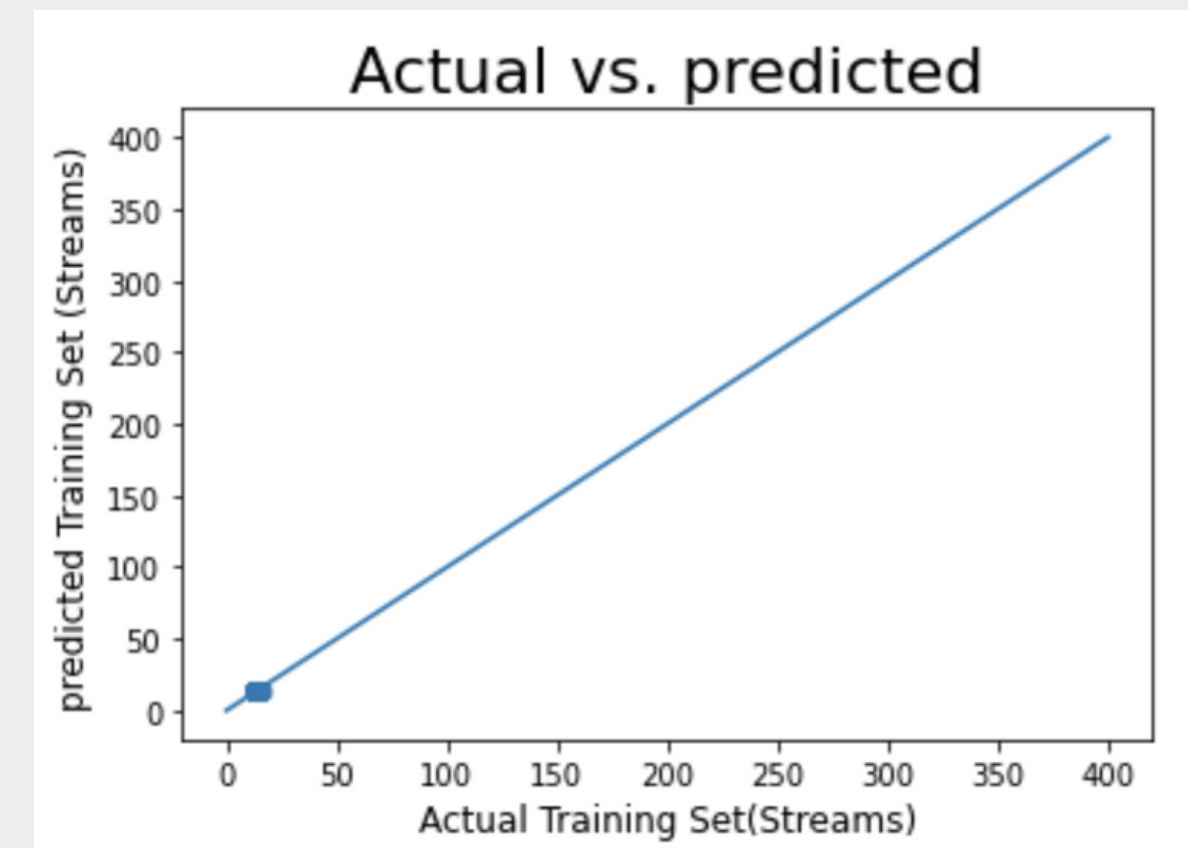
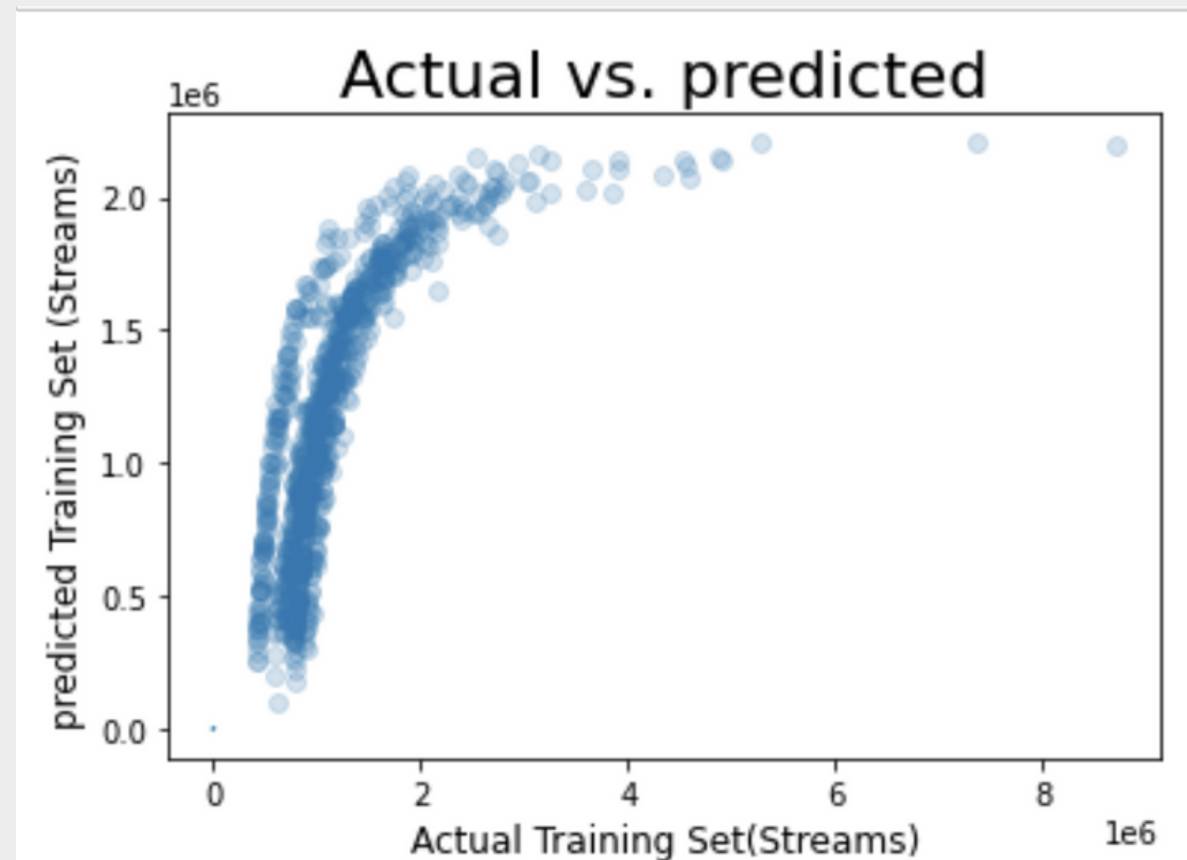
# Experiment#1 ( log Experiment to Handle Skewness of the Streams)

Training Score

0.517 → 0.691

Validation Mean Score

0.524 → 0.680



# Experiment#2 ( Polynomial )

Polynomial Features = 3

Training Score

0.691 → 0.852

Validation Mean Score

0.680 → -2.181

Polynomial Features = 2

Training Score

0.691 → 0.656

Validation Mean Score

0.680 → 0.432



**OVERFITTING**

# Experiment#3 ( Lasso )

**$\lambda = 100$**

**Training Score**

0.656 → 0.693

**Validation Mean Score**

0.432 → 0.673

**$\lambda = 10$**

**Training Score**

0.656 → 0.719

**Validation Mean Score**

0.432 → 0.695

**$\lambda = 1$**

**Training Score**

0.656 → 0.735

**Validation Mean Score**

0.432 → 0.704

# Experiment#4 ( Ridge )

**$\lambda = 100$**

**Training Score**

0.656 → 0.790

**Validation Mean Score**

0.432 → 0.744

**$\lambda = 10$**

**Training Score**

0.656 → 0.793

**Validation Mean Score**

0.432 → 0.742

**$\lambda = 1$**

**Training Score**

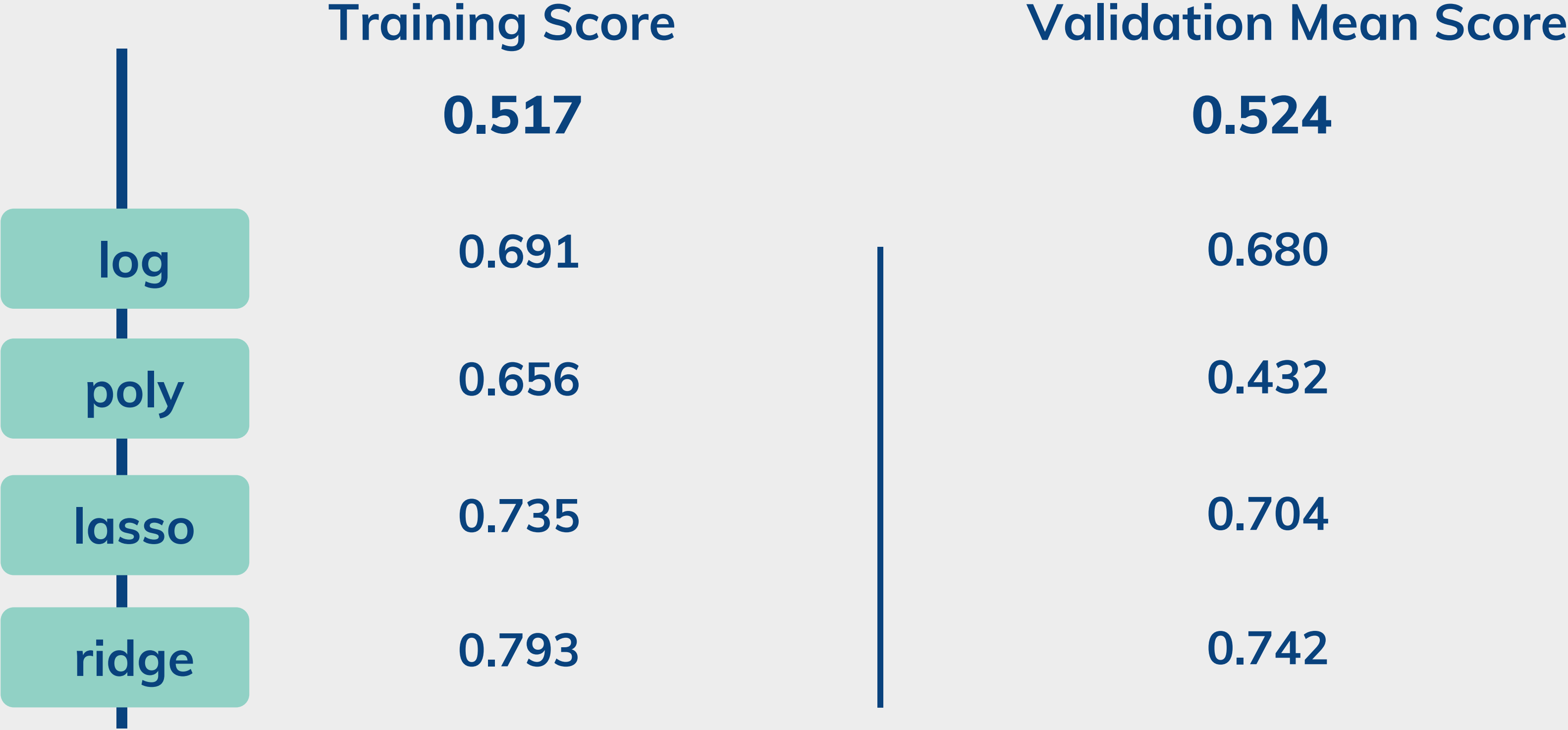
0.656 → 0.796

**Validation Mean Score**

0.432 → 0.733



# Experiments Summary





**Thank You**