# Spotify Tracks from 1922-2021

Corinne Smith

# Spotify Dataset

Combined 2 DataFrames: tracks.csv and artists.csv

Categorical Variables:

- mode
- explicit

Variables to Manipulate:

- key (range of 0-11, turn into dummies)
- genre (associated with artist, use first genre)
- timesignature (turn numeric into dummies)
- release_date (just want year)

Numeric:

- acousticness
- danceability
- energy
- duration_ms
- instrumentalness
- valence
- track_popularity
- tempo
- liveness
- loudness
- speechiness
- artist_followers
- artist_popularity

# Spotify Dataset

- Massive dataset: 432,000 rows!

- Limitations:
  - Genre for track is the first genre associated with the artist
    - Does not account for artists who create in different genres
  - Spotify is relatively new, so its user base is new
  - Popularity measures likely reflect ONLY demographics who use Spotify
  - Aggregated popularity loses subgroups of listeners

# Three Main Questions

1.  How can we predict the popularity of a track and what are the most important predictors?

2.  By what characteristics can we group/cluster tracks? What subcategories of tracks exist?

3.  Do different genres see varying levels of popularity, and does this change over time?

## How can we predict the popularity of a track and what are the most important predictors?
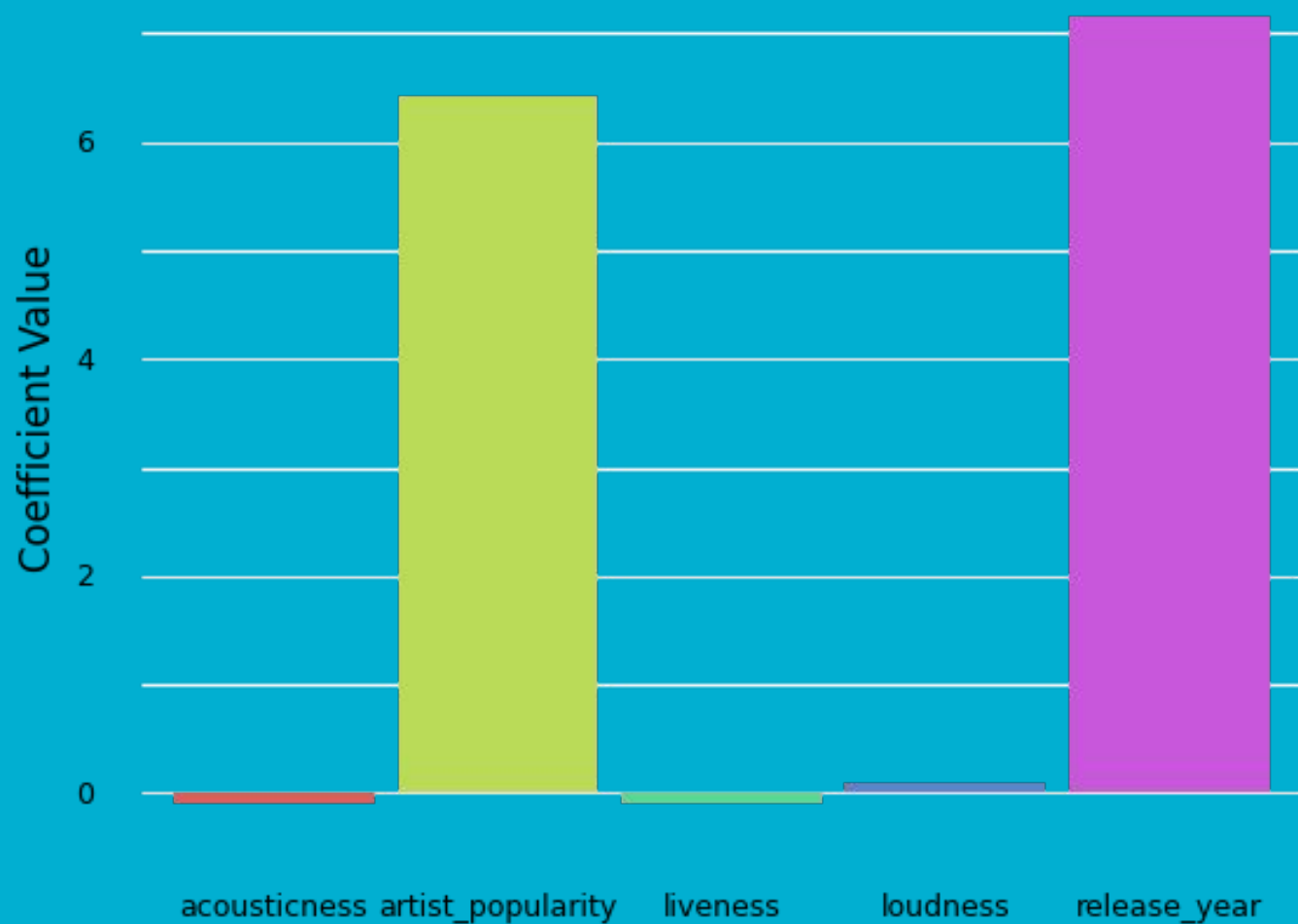
- Performed linear regression with LASSO and Ridge regularization

- R-squared values around 0.50, not ideal performance

- Odd behavior with Mean Squared Error

- Possible underfit/non-linear relationship

```
LASSO Train R2:   0.4781380768789373
LASSO Test R2:    0.47955972756038734
LASSO Train MSE:  152.41229866954444
LASSO Test MSE:   152.2105841381936

Ridge Train R2:   0.5079226693558112
Ridge Test R2:    0.508103677841667
Ridge Train MSE:  143.7135643813889
Ridge Test MSE:   143.86247662998912
```

LASSO Predictor Coefficient Values

#1

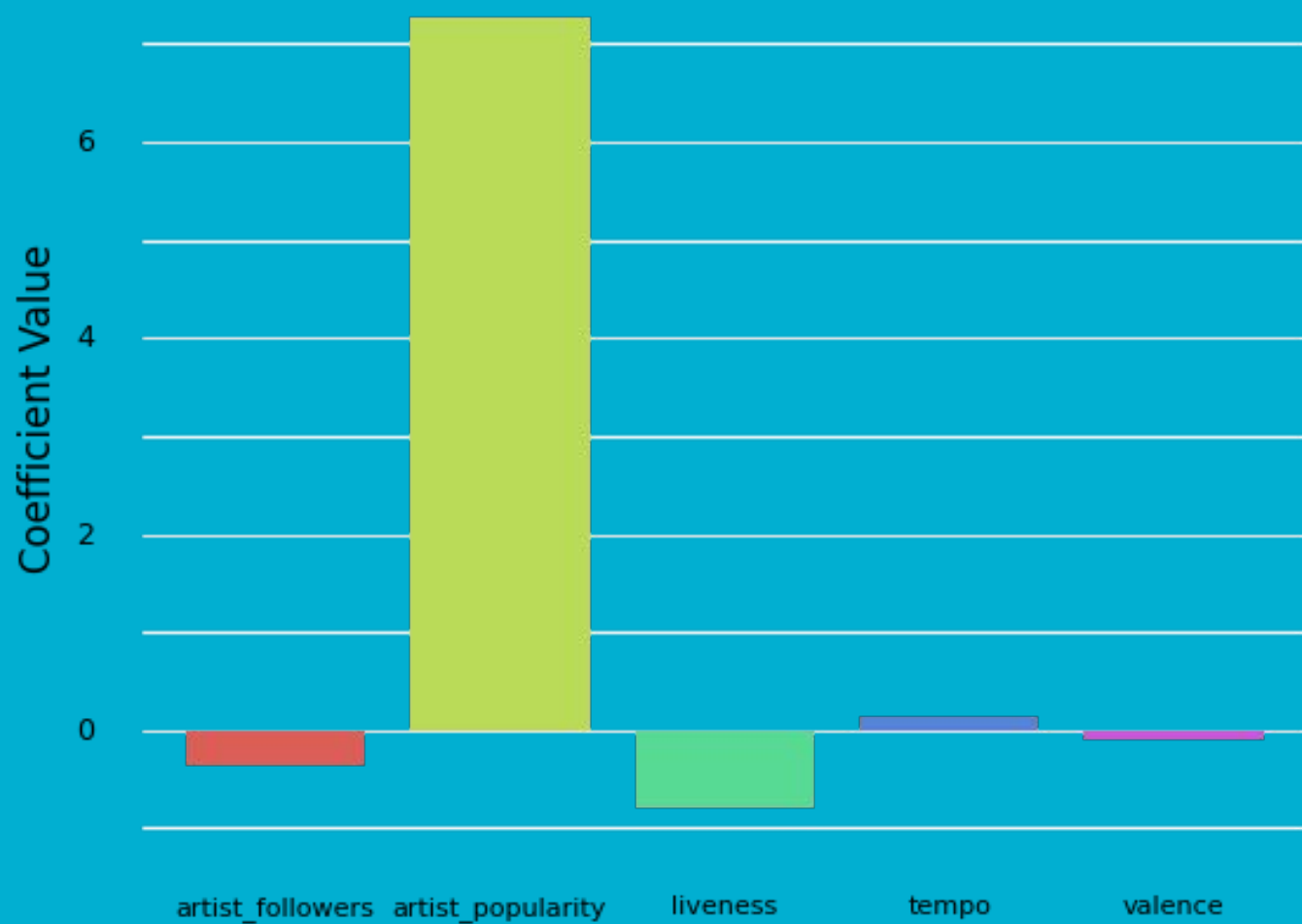Ridge Track Characteristic Coefficients

#1

#1
Ridge Track Characteristic Coefficients

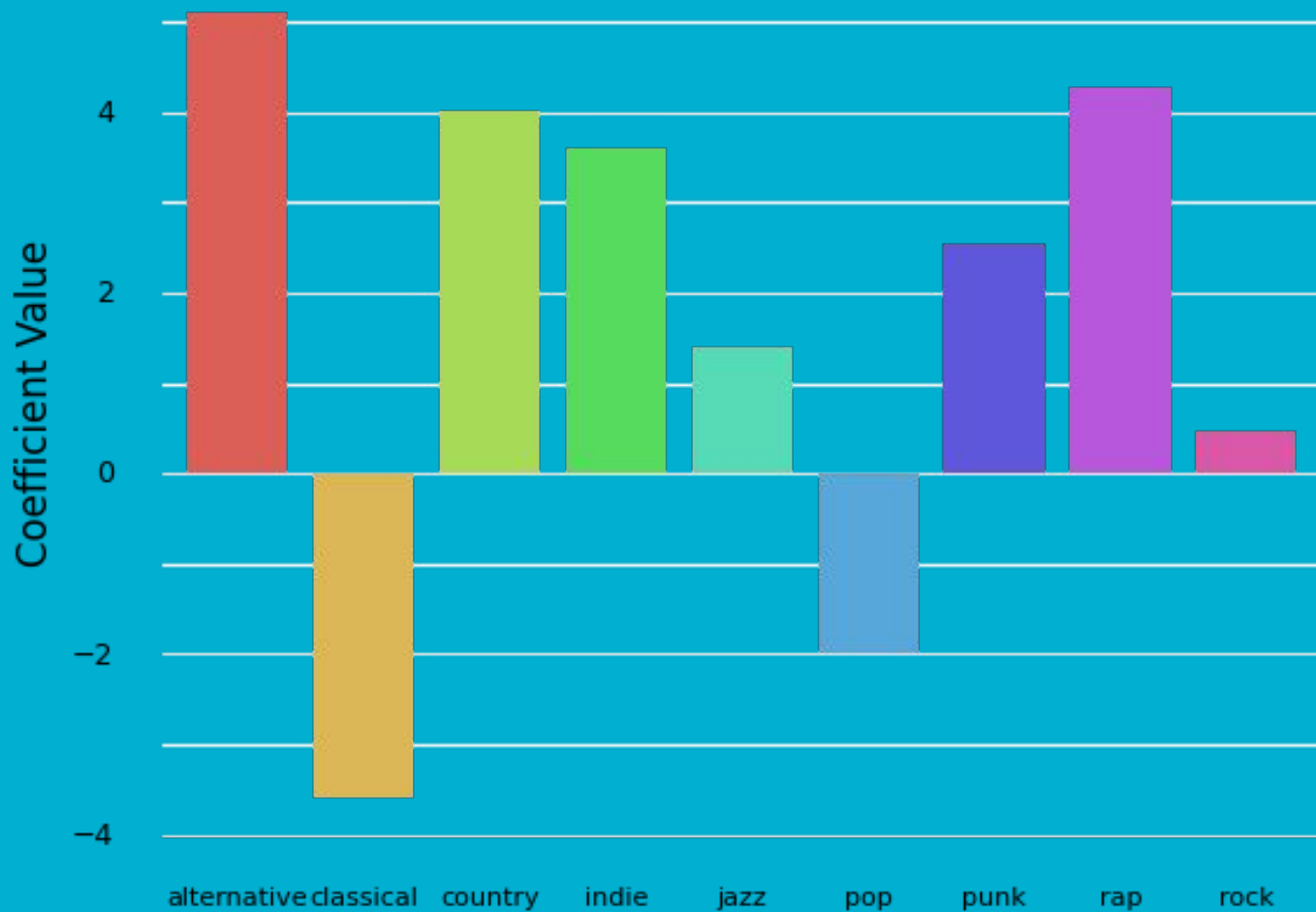Ridge Track Characteristic Coefficients

#1

Ridge Coefficients for Key
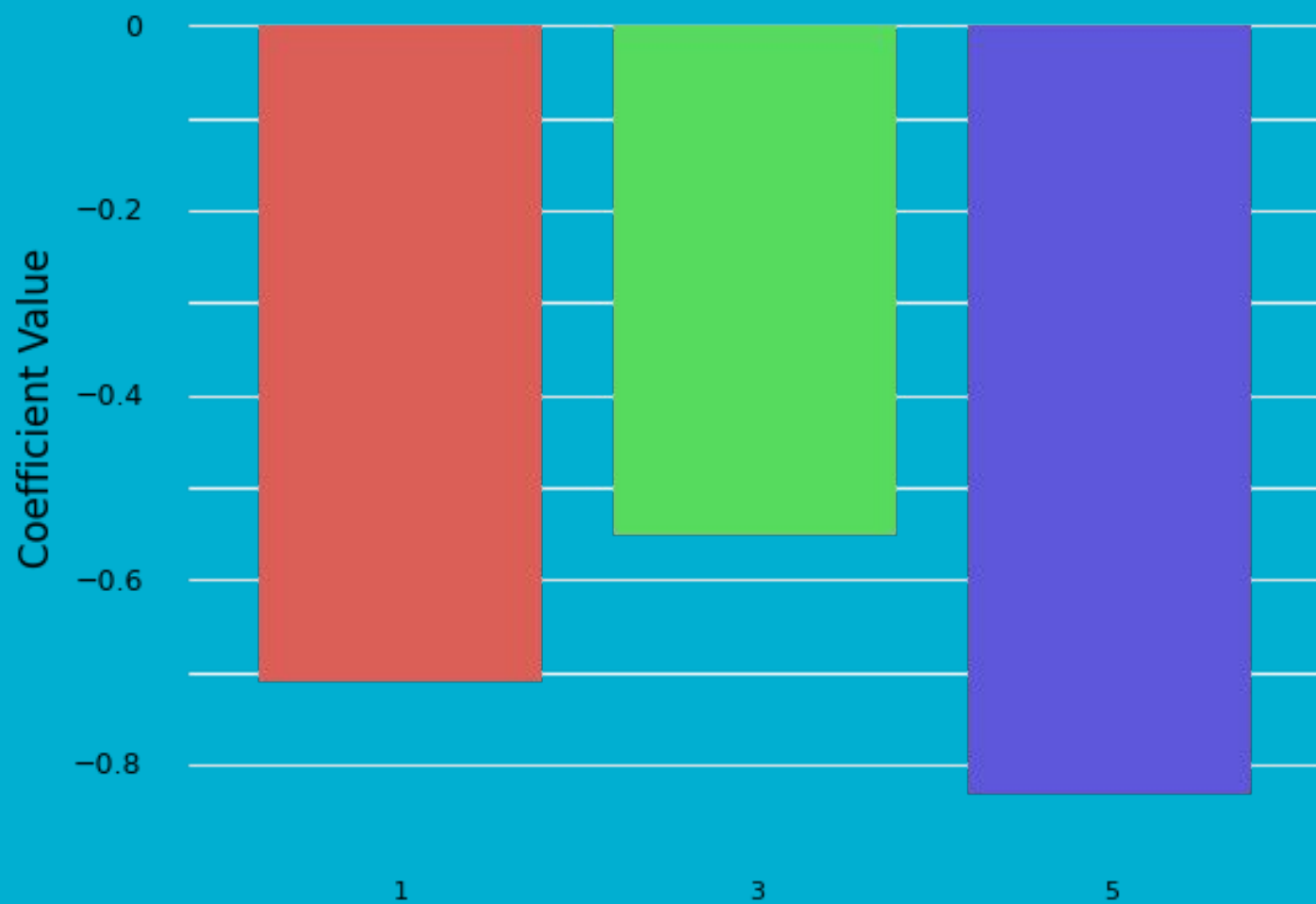
Ridge Coefficients for Genre
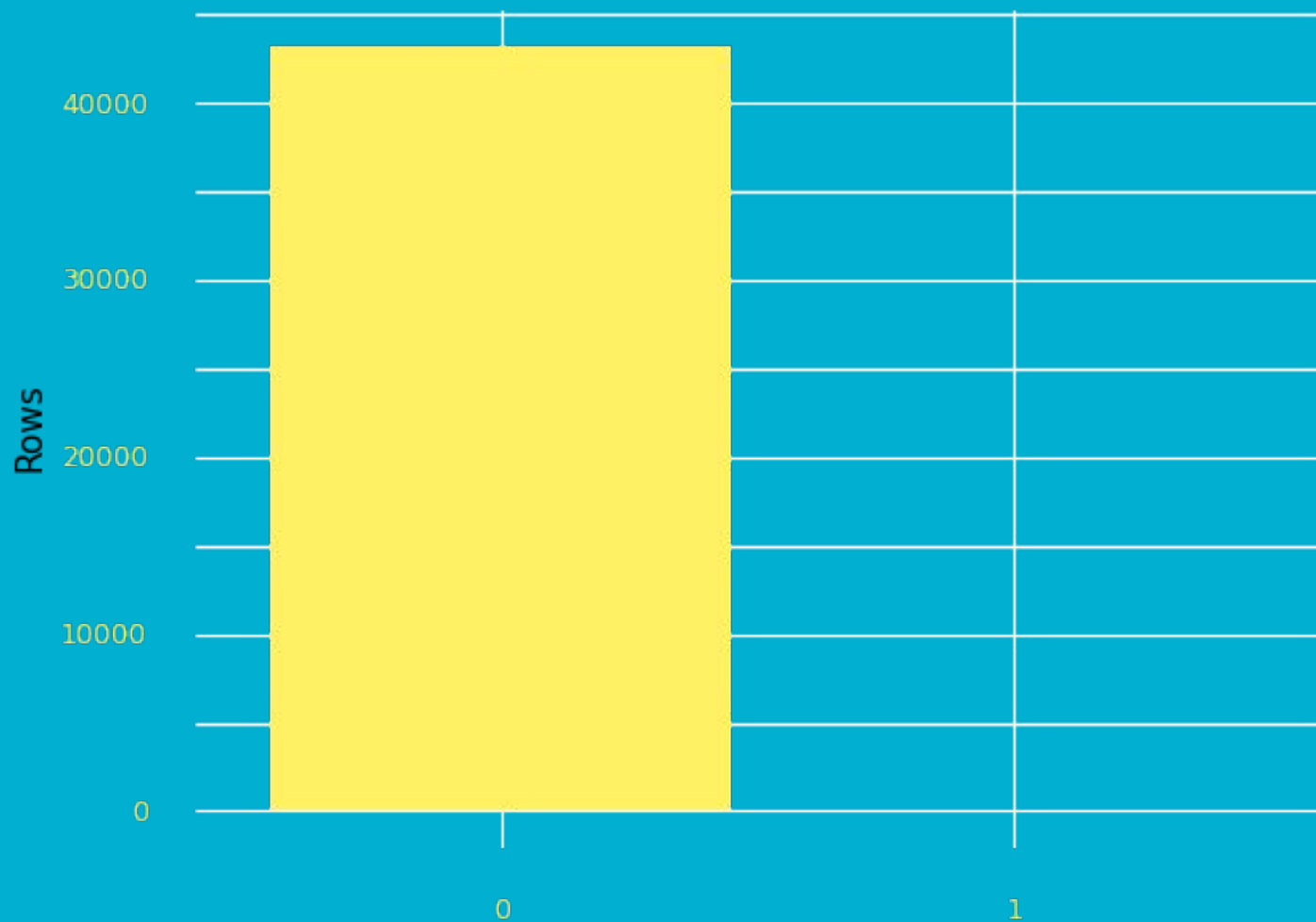
Ridge Coefficients for Time Signature Compared to 4/4

By what characteristics can we group/cluster tracks? What subcategories of tracks exist?

- Randomly sampled 10% of data

- Used Hierarchical Agglomerative Clustering

- Attempted to find genres/subgenres

- Silhouette Score: 0.8772

- 2 clusters

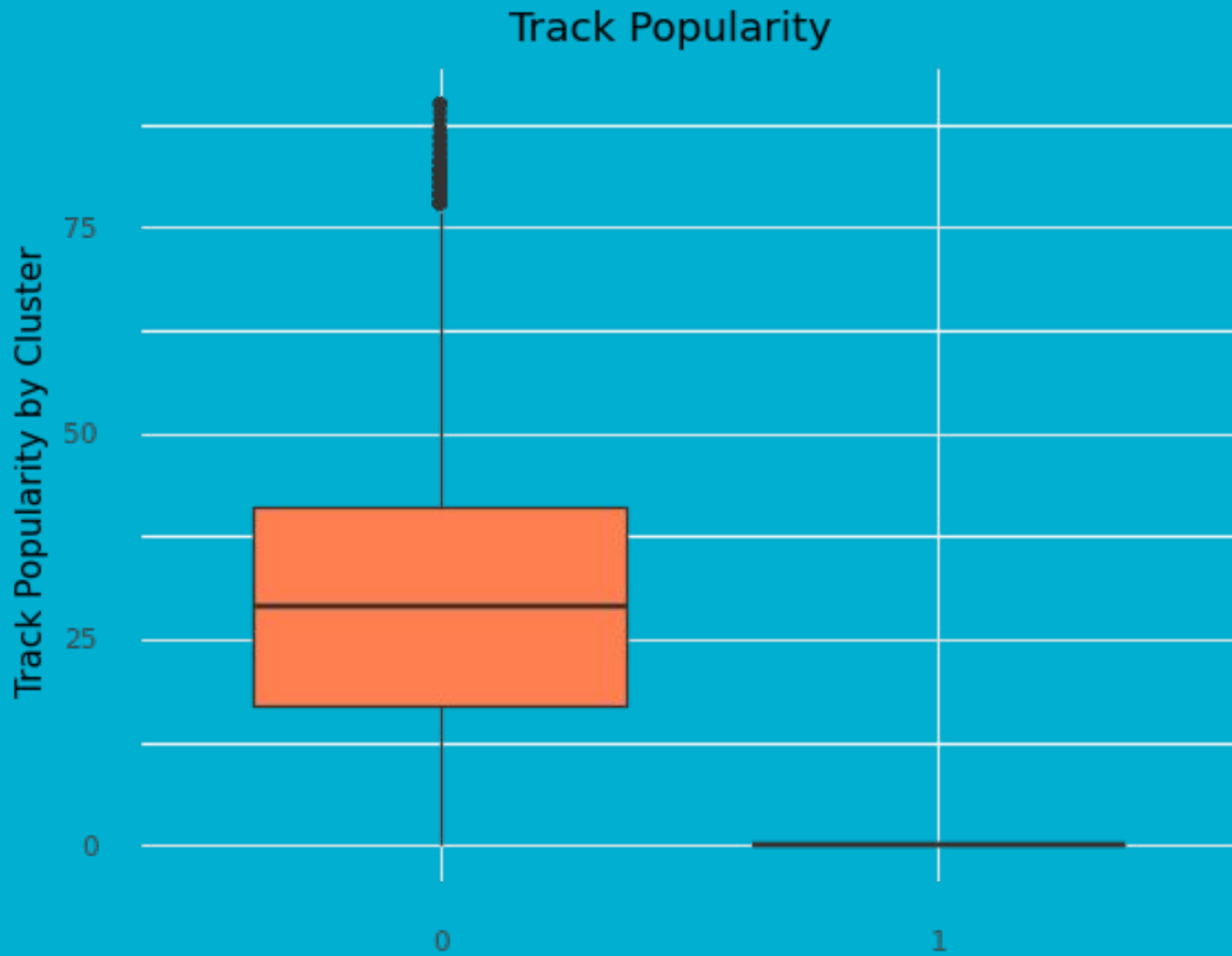- 43,199 members for Cluster 0
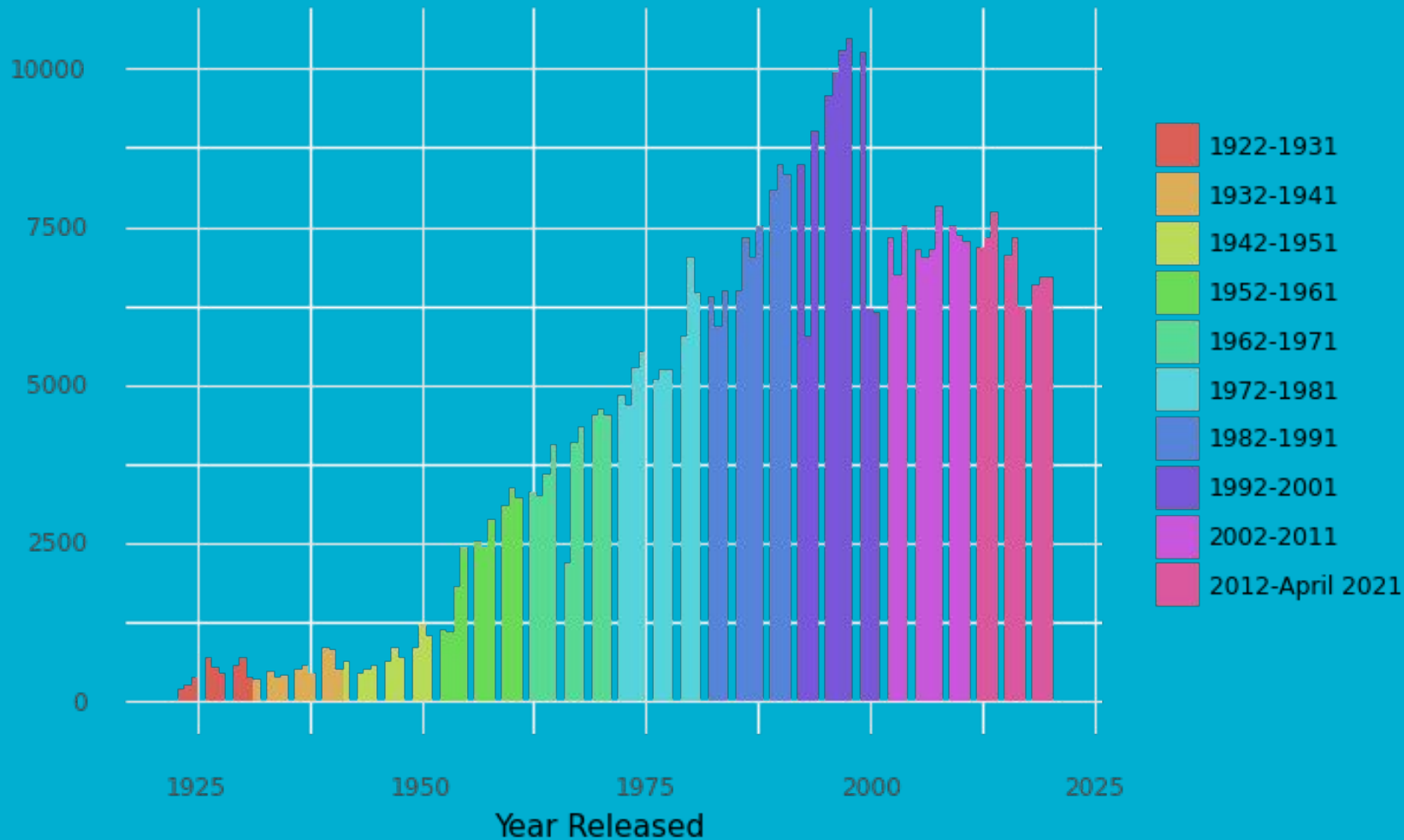
- 1 member for Cluster 1

# Cluster Membership

#2

Track Popularity

Do different genres see varying levels of popularity, and does this change over time?
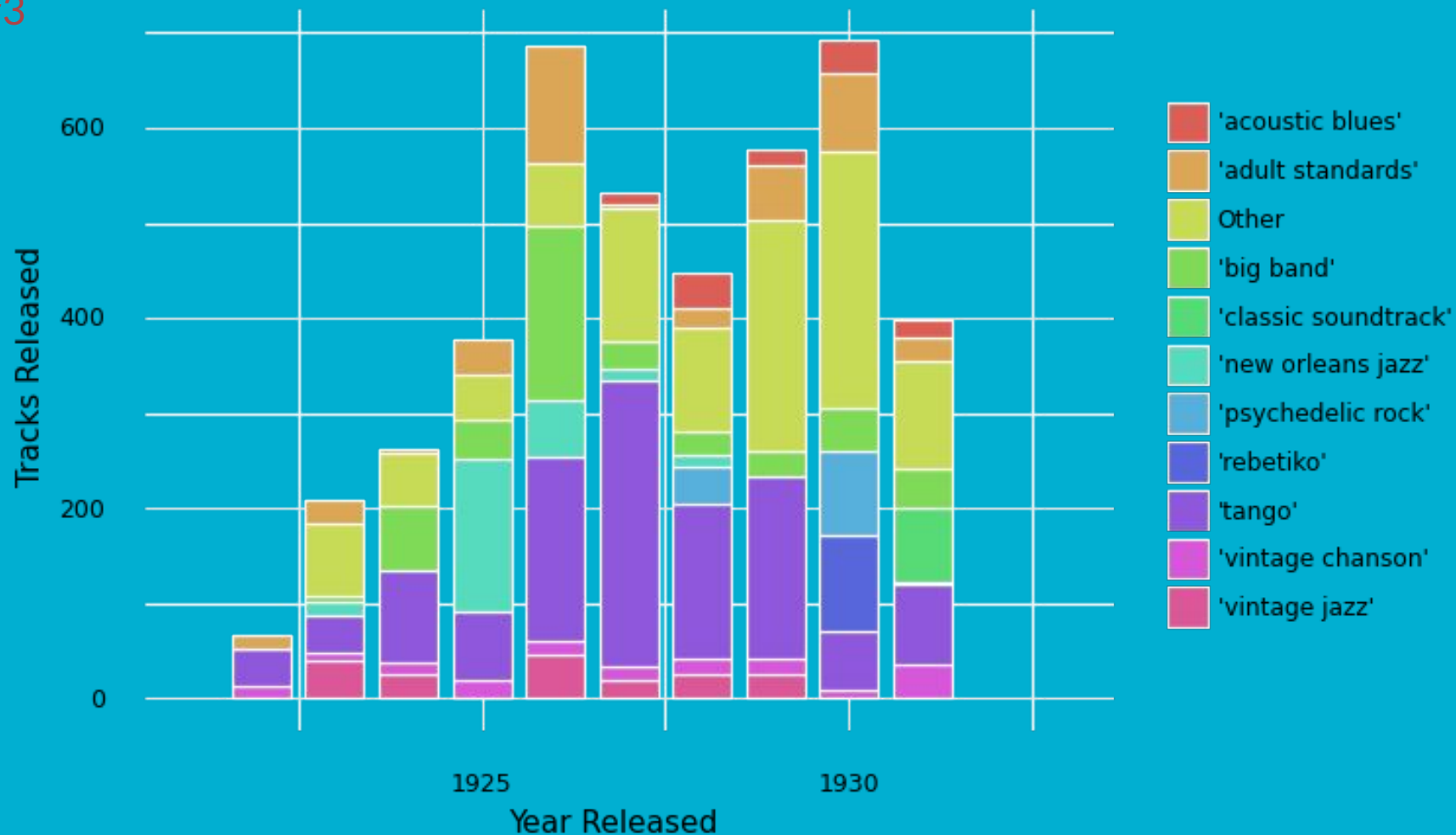
- Visualization questions:

- How many songs from each year are on Spotify?

- What are the most popular genres per decade?

- How do those change over years and decades?

- Most produced genres on Spotify?

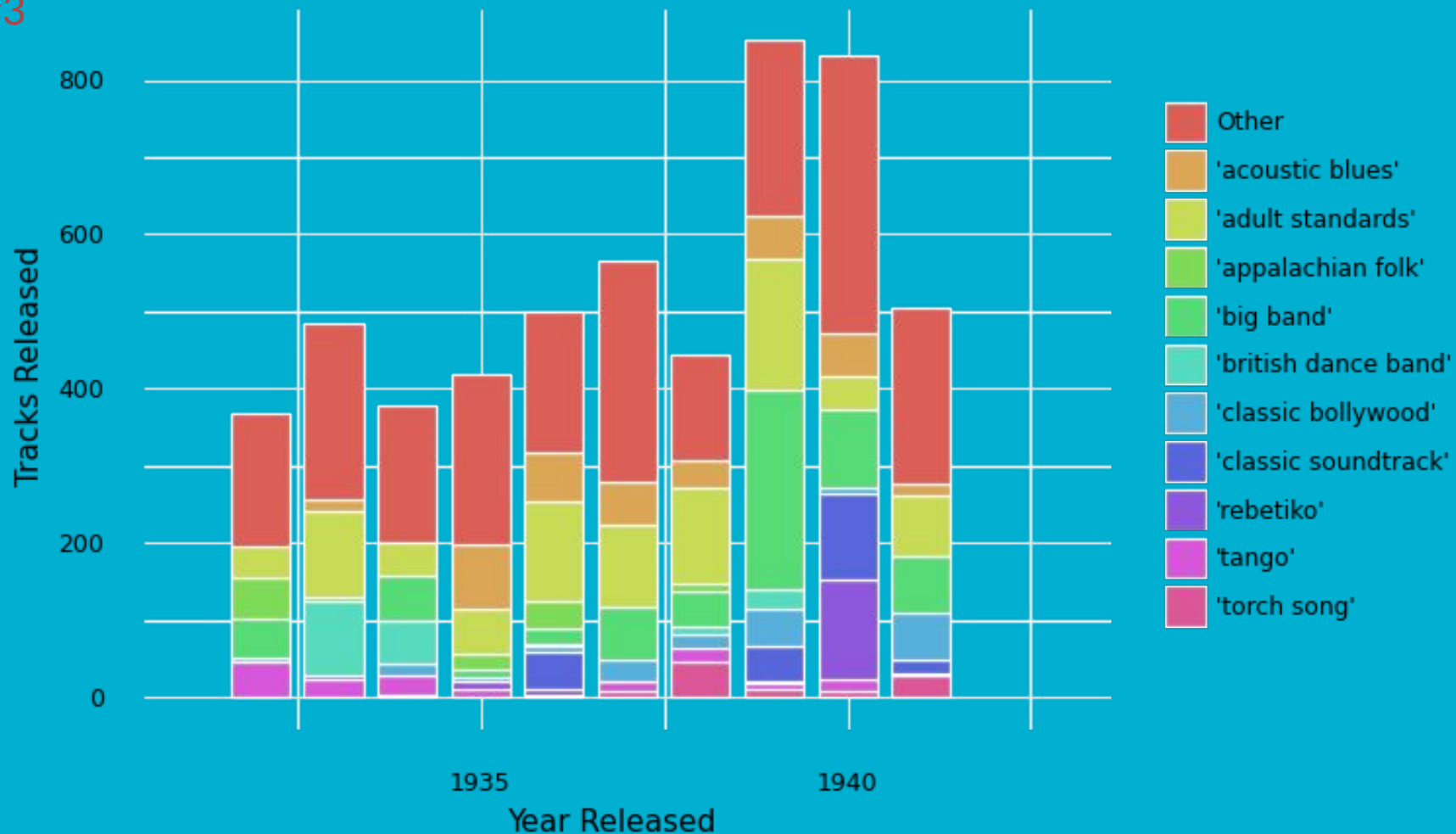- All visualized through bar charts

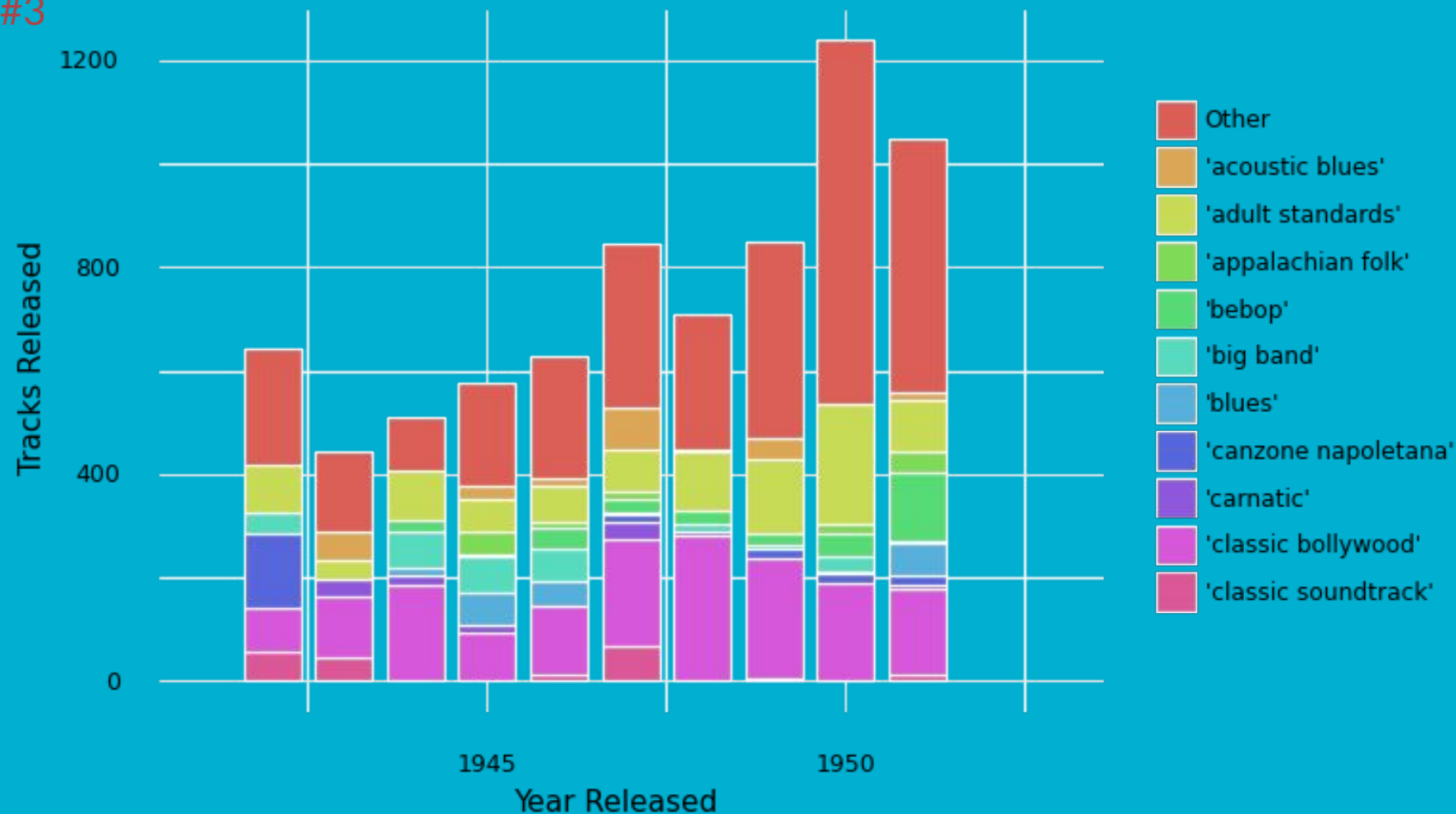Spotify Tracks by Year Released

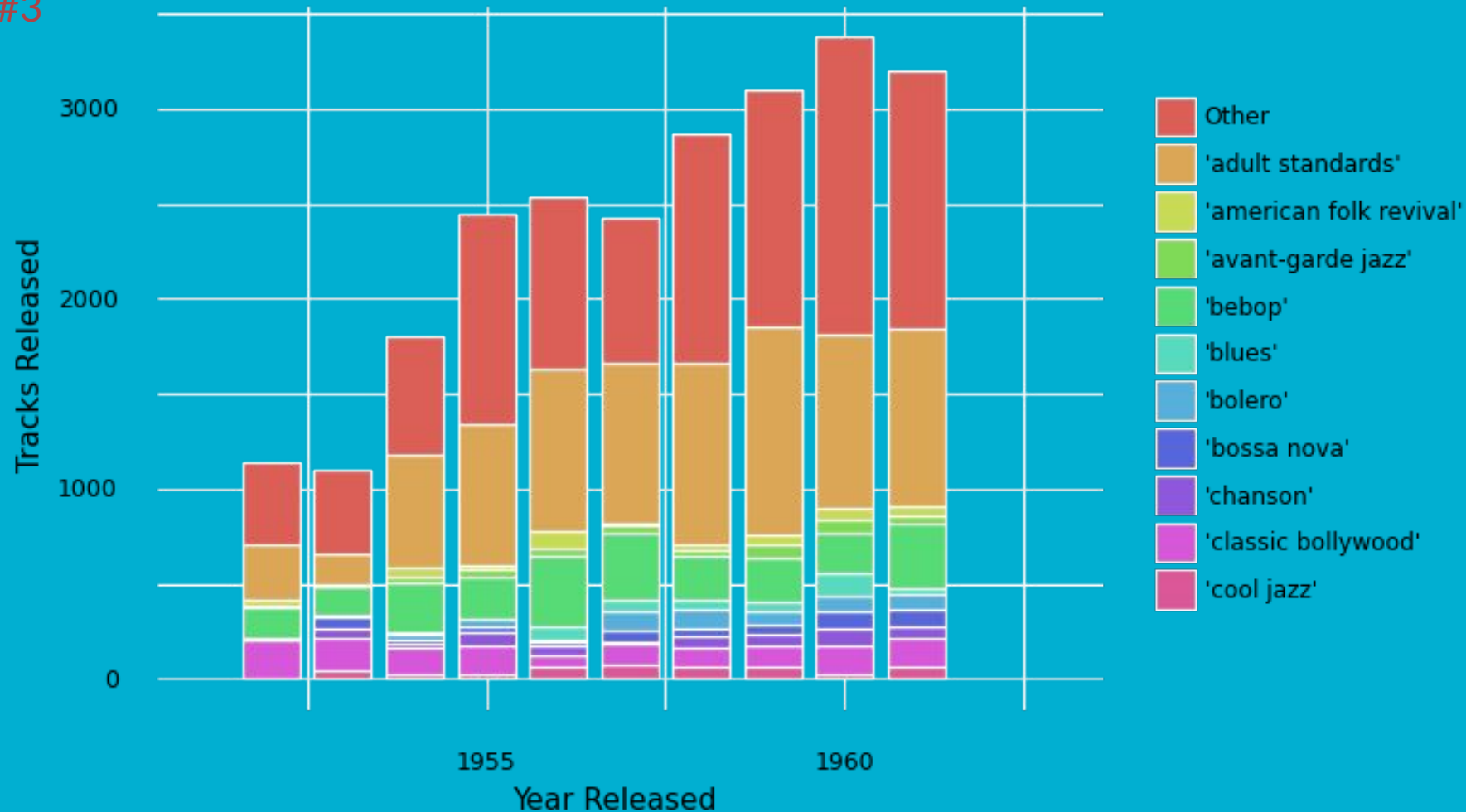Genre Distribution for 1922-1931

#3

Genre Distribution for 1932-1941

Genre Distribution for 1942-1951

#3

Legend:
- Other
- 'acoustic blues'
- 'adult standards'
- 'appalachian folk'
- 'bebop'
- 'big band'
- 'blues'
- 'canzone napoletana'
- 'carnatic'
- 'classic bollywood'
- 'classic soundtrack'

Y-axis: Tracks Released (0, 400, 800, 1200)
X-axis: Year Released (1945, 1950)

# Genre Distribution for 1952-1961

#3



Legend:
- Other
- 'adult standards'
- 'american folk revival'
- 'avant-garde jazz'
- 'bebop'
- 'blues'
- 'bolero'
- 'bossa nova'
- 'chanson'
- 'classic bollywood'
- 'cool jazz'

X-axis: Year Released (1955, 1960)
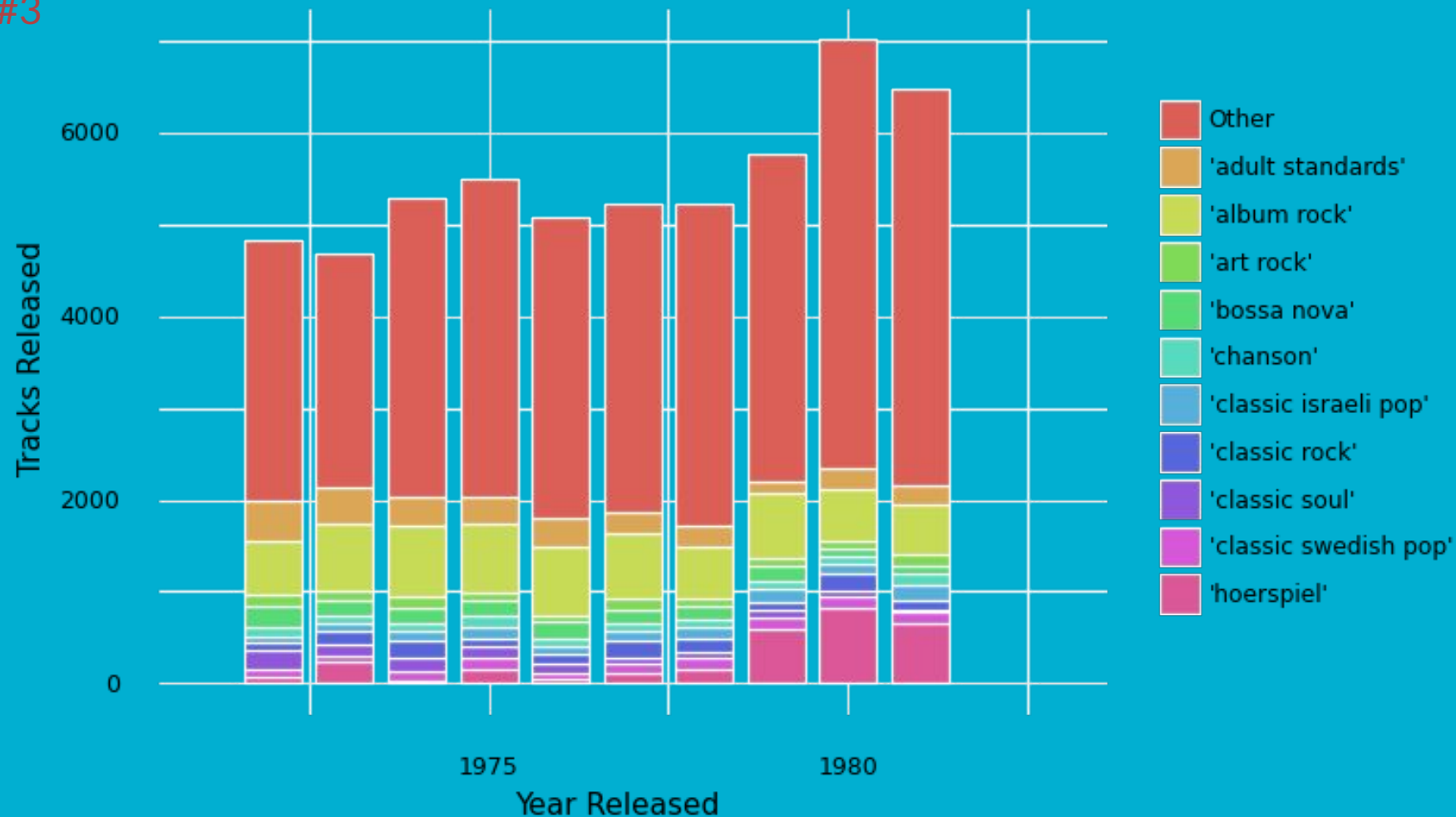
Y-axis: Tracks Released (0, 1000, 2000, 3000)

Genre Distribution for 1962-1971

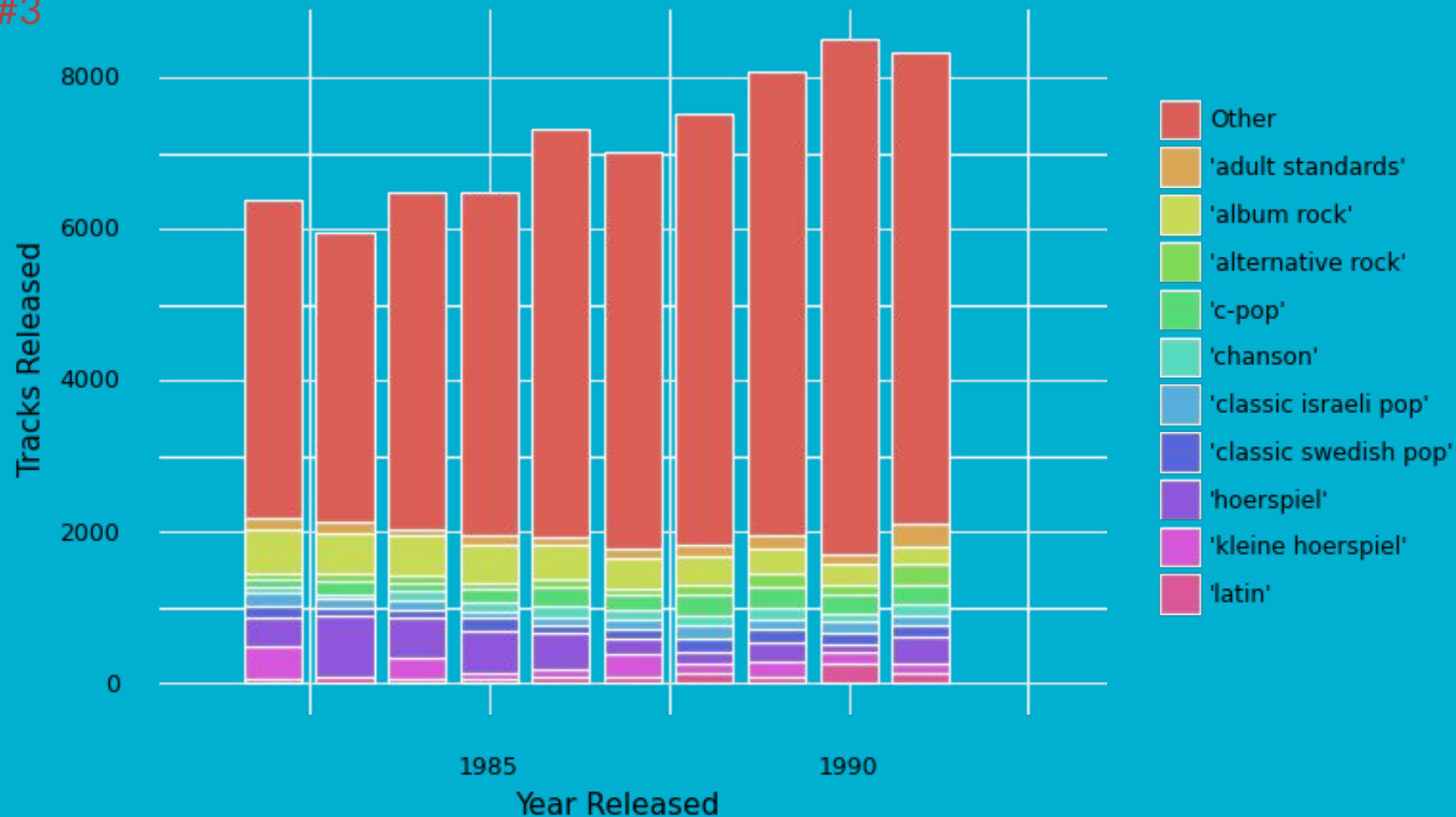Genre Distribution for 1972-1981

#3

Genre Distribution for 1982-1991

Genre Distribution for 1992-2001

#3

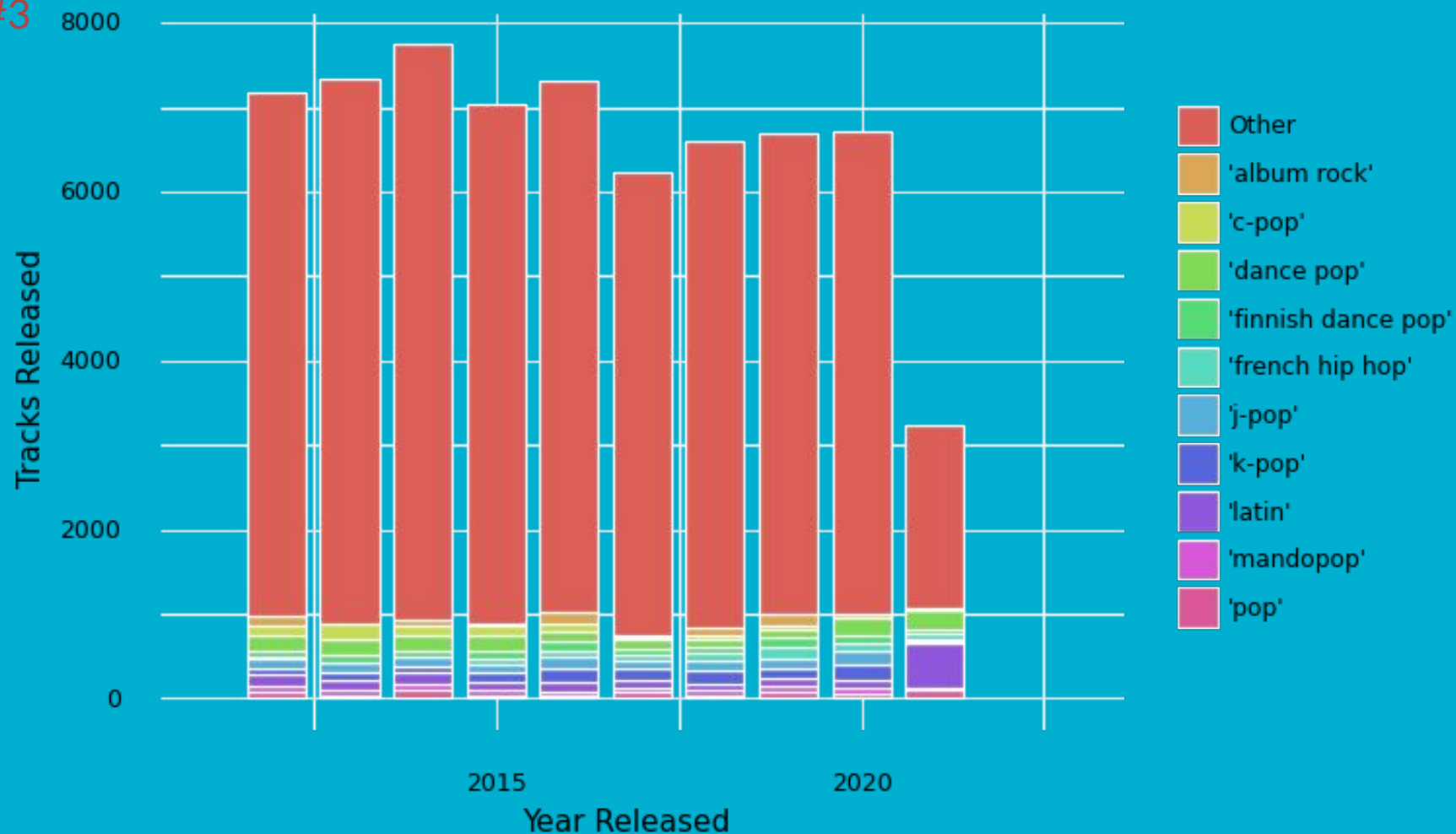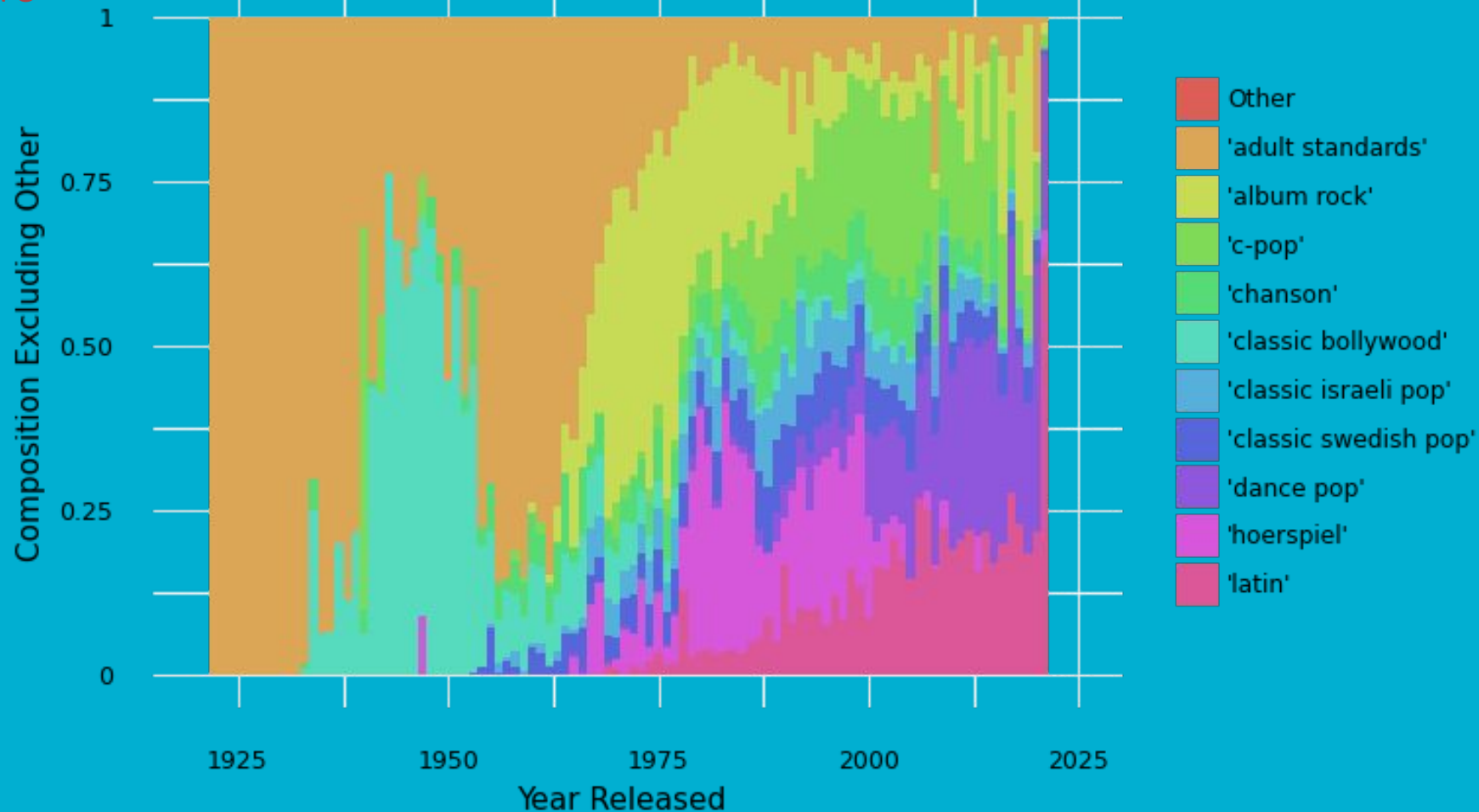Genre Distribution for 2002-2011

Genre Distribution for 2012-2021

#3

Spotify Genre Distributions for 1922-2021

#3

# Conclusion

- Linear model was okay

- HAC method did not work here

- Visualizations showed clear shifts over time

- In the future:

    - Use different predictive model, such as Decision Trees

    - Try different clustering algorithms

    - Try to get better genre classifications from data