

```
In [1]: import os
import numpy as np
import pandas as pd

import seaborn as sns
import plotly.express as px
import matplotlib.pyplot as plt
%matplotlib inline

from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.manifold import TSNE
from sklearn.decomposition import PCA
from sklearn.metrics import euclidean_distances
from scipy.spatial.distance import cdist

import warnings
warnings.filterwarnings("ignore")
```

```
In [2]: data = pd.read_csv("data/data.csv")
genre_data = pd.read_csv('data/data_by_genres.csv')
year_data = pd.read_csv('data/data_by_year.csv')
```

```
In [3]: print(data.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 170653 entries, 0 to 170652
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype  
---  -
0   valence                170653 non-null float64
1   year                  170653 non-null int64  
2   acousticness          170653 non-null float64
3   artists               170653 non-null object 
4   danceability           170653 non-null float64
5   duration_ms           170653 non-null int64  
6   energy                170653 non-null float64
7   explicit              170653 non-null int64  
8   id                    170653 non-null object 
9   instrumentalness       170653 non-null float64
10  key                    170653 non-null int64  
11  liveness              170653 non-null float64
12  loudness              170653 non-null float64
13  mode                  170653 non-null int64  
14  name                  170653 non-null object 
15  popularity             170653 non-null int64  
16  release_date           170653 non-null object 
17  speechiness           170653 non-null float64
18  tempo                 170653 non-null float64
dtypes: float64(9), int64(6), object(4)
memory usage: 24.7+ MB
None
```

```
In [4]: print(genre_data.info())
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2973 entries, 0 to 2972
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   mode                   2973 non-null  int64
1   genres                 2973 non-null  object
2   acousticness           2973 non-null  float64
3   danceability           2973 non-null  float64
4   duration_ms            2973 non-null  float64
5   energy                 2973 non-null  float64
6   instrumentalness        2973 non-null  float64
7   liveness               2973 non-null  float64
8   loudness               2973 non-null  float64
9   speechiness            2973 non-null  float64
10  tempo                  2973 non-null  float64
11  valence                 2973 non-null  float64
12  popularity              2973 non-null  float64
13  key                    2973 non-null  int64
dtypes: float64(11), int64(2), object(1)
memory usage: 325.3+ KB
None

```

In [5]: `print(year_data.info())`

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   mode                   100 non-null  int64
1   year                  100 non-null  int64
2   acousticness           100 non-null  float64
3   danceability           100 non-null  float64
4   duration_ms            100 non-null  float64
5   energy                 100 non-null  float64
6   instrumentalness        100 non-null  float64
7   liveness               100 non-null  float64
8   loudness               100 non-null  float64
9   speechiness            100 non-null  float64
10  tempo                  100 non-null  float64
11  valence                 100 non-null  float64
12  popularity              100 non-null  float64
13  key                    100 non-null  int64
dtypes: float64(11), int64(3)
memory usage: 11.1 KB
None

```

In [6]: `from yellowbrick.target import FeatureCorrelation`

```

feature_names = ['acousticness', 'danceability', 'energy', 'instrumentalness',
                 'liveness', 'loudness', 'speechiness', 'tempo', 'valence', 'duration_ms', 'explicit', 'key', 'mode', 'year']

X, y = data[feature_names], data['popularity']

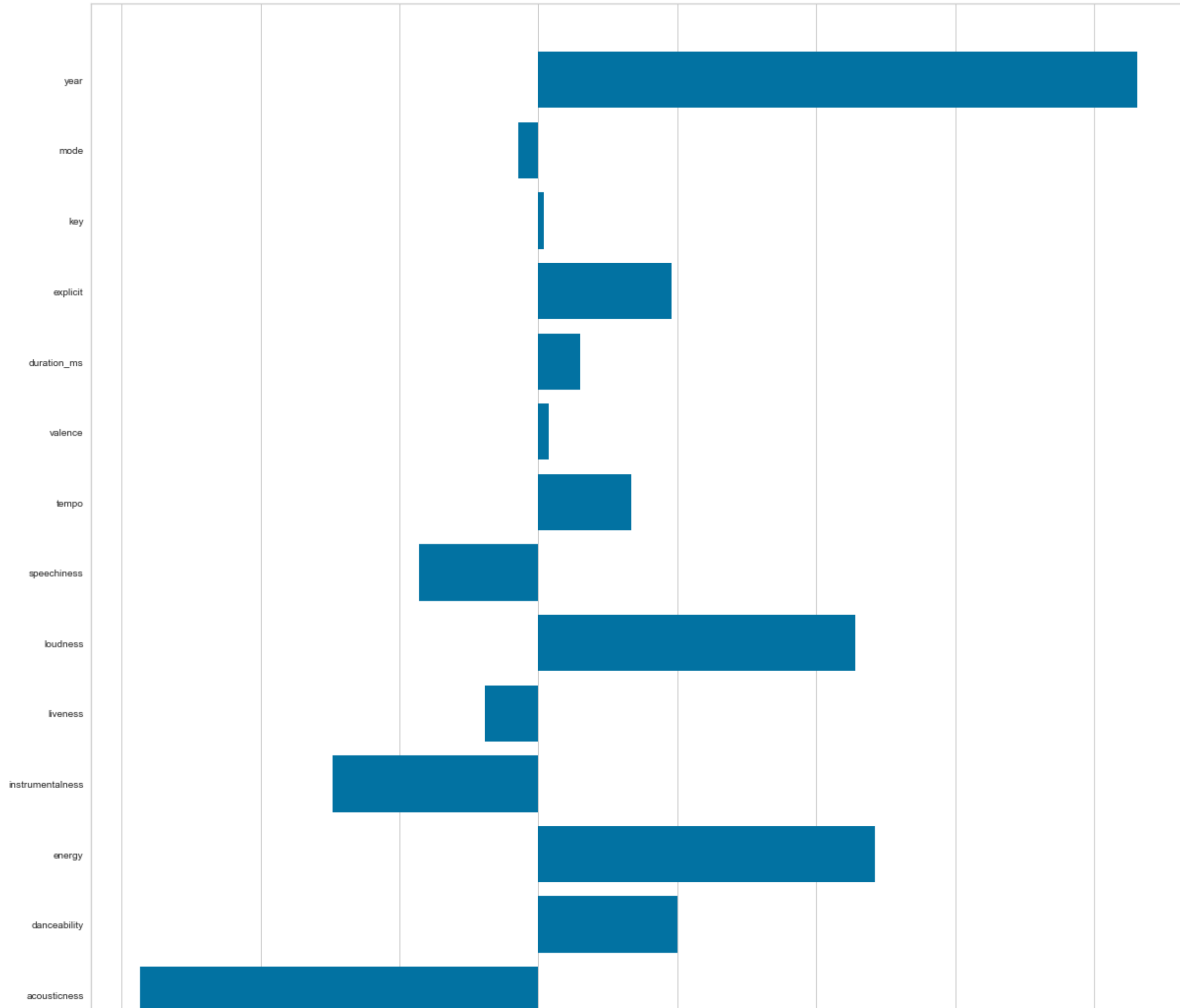
# Create a list of the feature names
features = np.array(feature_names)

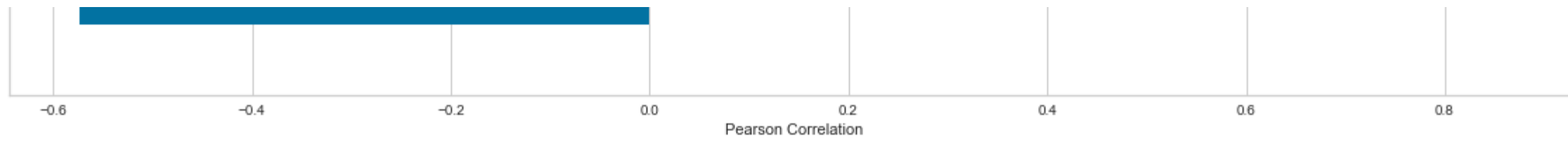
# Instantiate the visualizer
visualizer = FeatureCorrelation(labels=features)

plt.rcParams['figure.figsize']=(20,20)
visualizer.fit(X, y)      # Fit the data to the visualizer
visualizer.show()

```

Features correlation with dependent variable





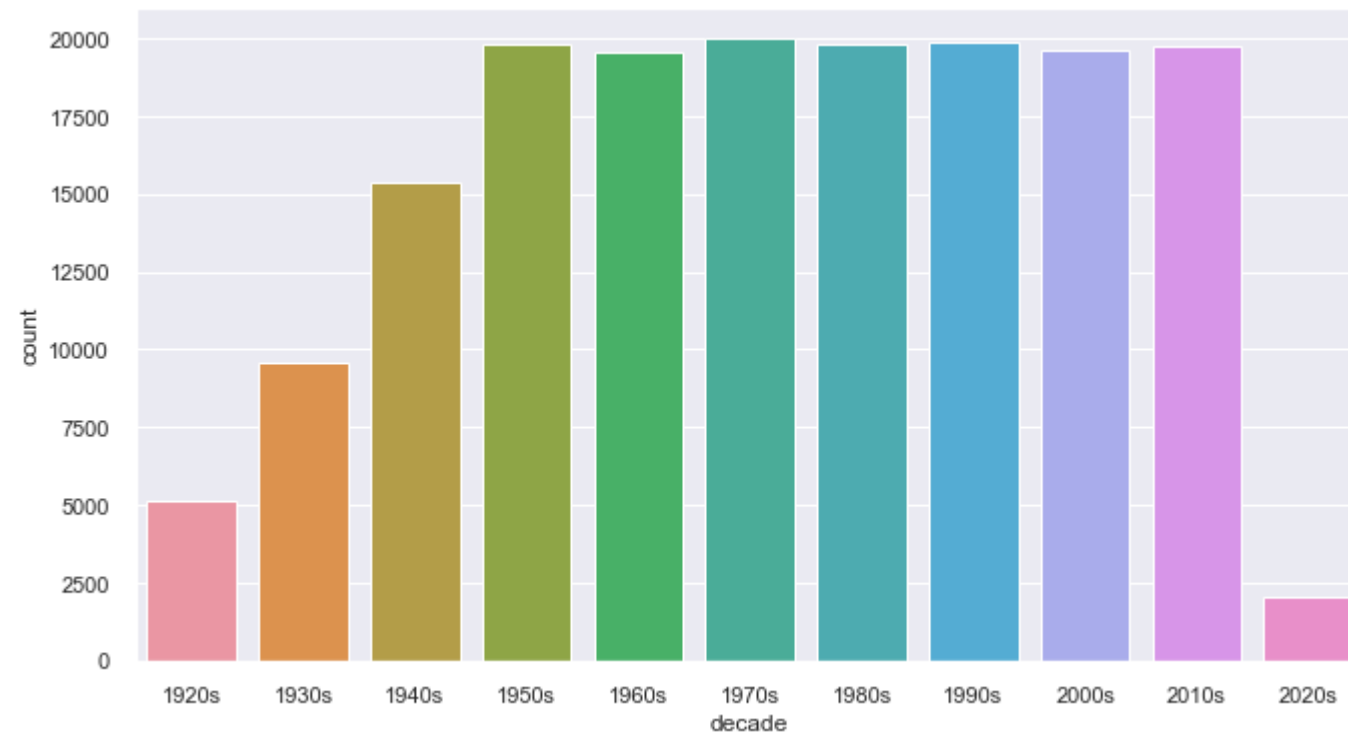
Out[6]: <AxesSubplot:title={'center':'Features correlation with dependent variable'}, xlabel='Pearson Correlation'>

```
In [7]: def get_decade(year):
    period_start = int(year/10) * 10
    decade = '{}s'.format(period_start)
    return decade

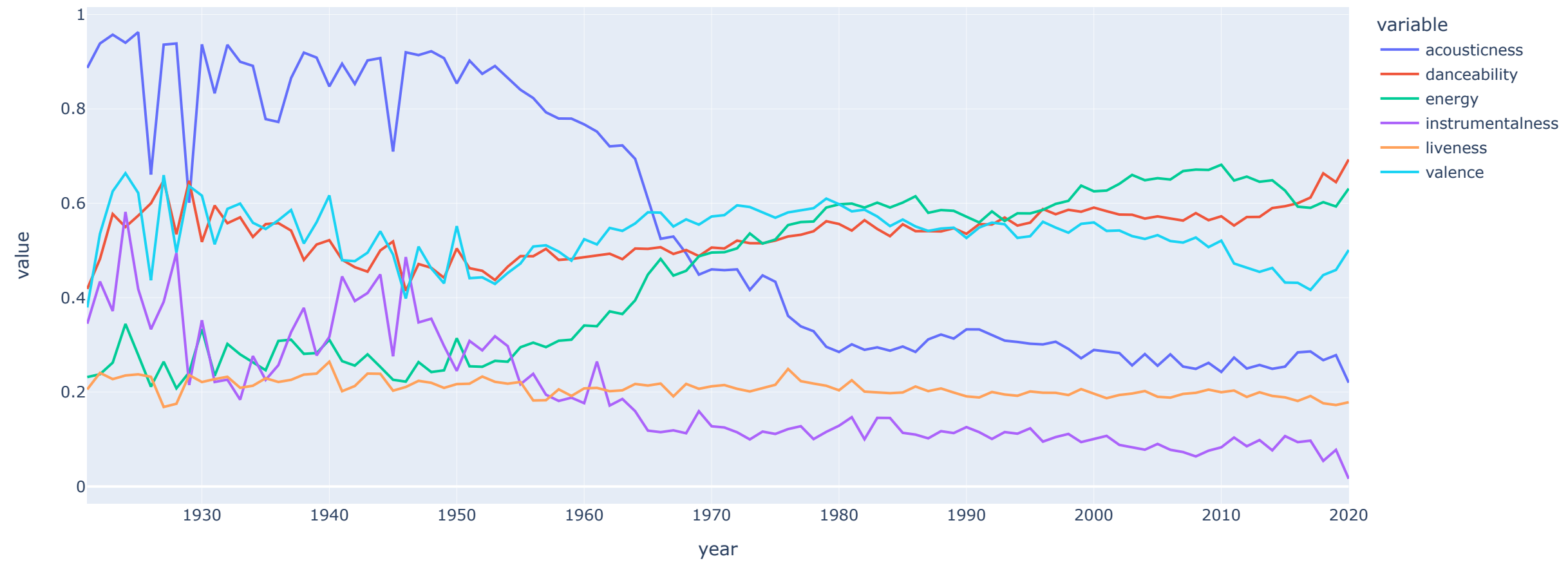
data['decade'] = data['year'].apply(get_decade)

sns.set(rc={'figure.figsize':(11 ,6)})
sns.countplot(data['decade'])
```

Out[7]: <AxesSubplot:xlabel='decade', ylabel='count'>

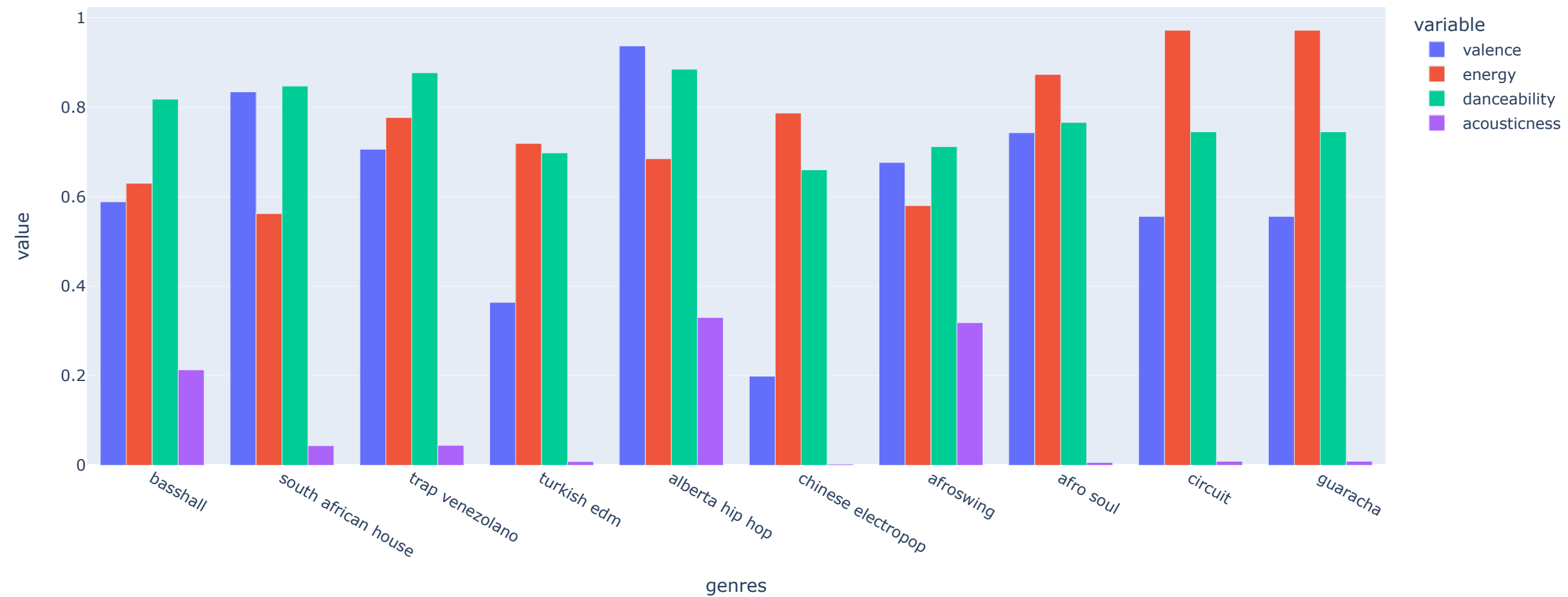


```
In [8]: sound_features = ['acousticness', 'danceability', 'energy', 'instrumentalness', 'liveness', 'valence']
fig = px.line(year_data, x='year', y=sound_features)
fig.show()
```



```
In [9]: top10_genres = genre_data.nlargest(10, 'popularity')

fig = px.bar(top10_genres, x='genres', y=['valence', 'energy', 'danceability', 'acousticness'], barmode='group')
fig.show()
```



```
In [10]: from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline

cluster_pipeline = Pipeline([('scaler', StandardScaler()), ('kmeans', KMeans(n_clusters=10))])
X = genre_data.select_dtypes(np.number)
cluster_pipeline.fit(X)
genre_data['cluster'] = cluster_pipeline.predict(X)
```

```
In [11]: # Visualizing the Clusters with t-SNE

from sklearn.manifold import TSNE

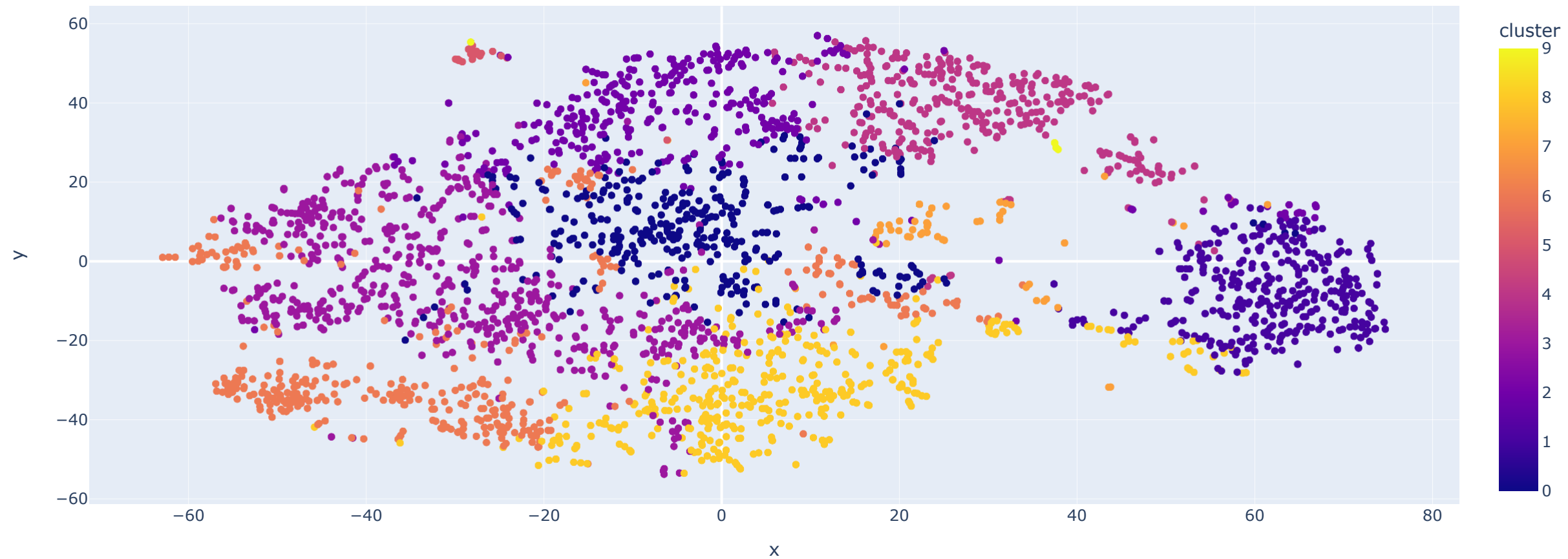
tsne_pipeline = Pipeline([('scaler', StandardScaler()), ('tsne', TSNE(n_components=2, verbose=1))])
genre_embedding = tsne_pipeline.fit_transform(X)
projection = pd.DataFrame(columns=['x', 'y'], data=genre_embedding)
projection['genres'] = genre_data['genres']
projection['cluster'] = genre_data['cluster']

fig = px.scatter(
    projection, x='x', y='y', color='cluster', hover_data=['x', 'y', 'genres'])
fig.show()
```

```

[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 2973 samples in 0.024s...
[t-SNE] Computed neighbors for 2973 samples in 0.701s...
[t-SNE] Computed conditional probabilities for sample 1000 / 2973
[t-SNE] Computed conditional probabilities for sample 2000 / 2973
[t-SNE] Computed conditional probabilities for sample 2973 / 2973
[t-SNE] Mean sigma: 0.777516
[t-SNE] KL divergence after 250 iterations with early exaggeration: 76.111481
[t-SNE] KL divergence after 1000 iterations: 1.389024

```



```

In [12]: song_cluster_pipeline = Pipeline([('scaler', StandardScaler()),
                                           ('kmeans', KMeans(n_clusters=20,
                                                             verbose=False)),
                                           ], verbose=False)

X = data.select_dtypes(np.number)
number_cols = list(X.columns)
song_cluster_pipeline.fit(X)
song_cluster_labels = song_cluster_pipeline.predict(X)
data['cluster_label'] = song_cluster_labels

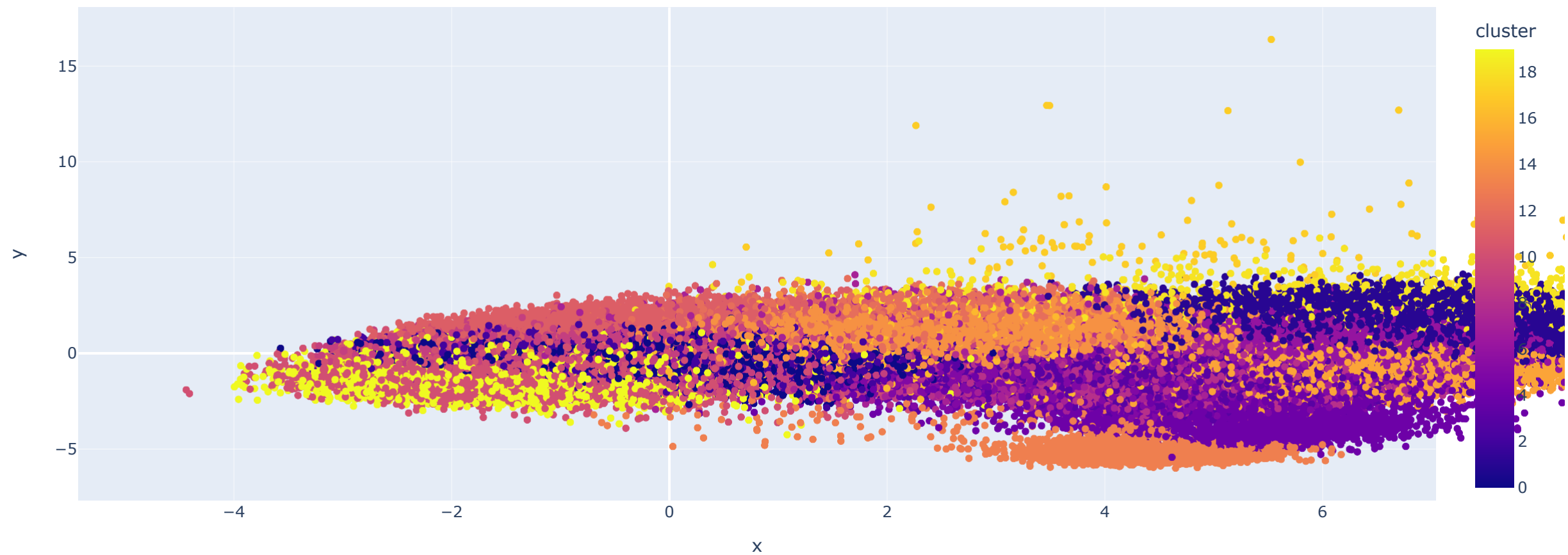
In [13]: # Visualizing the Clusters with PCA

from sklearn.decomposition import PCA

pca_pipeline = Pipeline([('scaler', StandardScaler()), ('PCA', PCA(n_components=2))])
song_embedding = pca_pipeline.fit_transform(X)
projection = pd.DataFrame(columns=['x', 'y'], data=song_embedding)
projection['title'] = data['name']
projection['cluster'] = data['cluster_label']

```

```
fig = px.scatter(
    projection, x='x', y='y', color='cluster', hover_data=['x', 'y', 'title'])
fig.show()
```



In [14]: !pip install spotipy

```
Requirement already satisfied: spotipy in c:\users\lenovo\appdata\local\programs\python\python310\lib\site-packages (2.19.0)
Requirement already satisfied: requests>=2.25.0 in c:\users\lenovo\appdata\local\programs\python\python310\lib\site-packages (from spotipy) (2.27.1)
Requirement already satisfied: six>=1.15.0 in c:\users\lenovo\appdata\local\programs\python\python310\lib\site-packages (from spotipy) (1.16.0)
Requirement already satisfied: urllib3>=1.26.0 in c:\users\lenovo\appdata\local\programs\python\python310\lib\site-packages (from spotipy) (1.26.9)
Requirement already satisfied: charset-normalizer~=2.0.0 in c:\users\lenovo\appdata\local\programs\python\python310\lib\site-packages (from requests>=2.25.0->spotipy) (2.0.12)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\lenovo\appdata\local\programs\python\python310\lib\site-packages (from requests>=2.25.0->spotipy) (2021.10.8)
Requirement already satisfied: idna<4,>=2.5 in c:\users\lenovo\appdata\local\programs\python\python310\lib\site-packages (from requests>=2.25.0->spotipy) (3.3)
WARNING: You are using pip version 22.0.4; however, version 23.1.2 is available.
You should consider upgrading via the 'C:\Users\lenovo\AppData\Local\Programs\Python\Python310\python.exe -m pip install --upgrade pip' command.
```

```
In [15]: import spotipy
from spotipy.oauth2 import SpotifyClientCredentials
from collections import defaultdict

sp = spotipy.Spotify(auth_manager=SpotifyClientCredentials(client_id='f0e48736cb9840069532c86d12d4740d',
                                                         client_secret='4bffb4bc44f74e2d8144b89db6b94d08'))

def find_song(name, year):
    song_data = defaultdict()
    results = sp.search(q= 'track: {} year: {}'.format(name,year), limit=1)
    if results['tracks']['items'] == []:
        return None

    results = results['tracks']['items'][0]
```



```

track_id = results['id']
audio_features = sp.audio_features(track_id)[0]

song_data['name'] = [name]
song_data['year'] = [year]
song_data['explicit'] = [int(results['explicit'])]
song_data['duration_ms'] = [results['duration_ms']]
song_data['popularity'] = [results['popularity']]

for key, value in audio_features.items():
    song_data[key] = value

return pd.DataFrame(song_data)

```

```

In [16]: from collections import defaultdict
from sklearn.metrics import euclidean_distances
from scipy.spatial.distance import cdist
import difflib

number_cols = ['valence', 'year', 'acousticness', 'danceability', 'duration_ms', 'energy', 'explicit',
               'instrumentalness', 'key', 'liveness', 'loudness', 'mode', 'popularity', 'speechiness', 'tempo']

def get_song_data(song, spotify_data):
    try:
        song_data = spotify_data[(spotify_data['name'] == song['name'])
                                & (spotify_data['year'] == song['year'])].iloc[0]
        return song_data
    except IndexError:
        return find_song(song['name'], song['year'])

def get_mean_vector(song_list, spotify_data):
    song_vectors = []

    for song in song_list:
        song_data = get_song_data(song, spotify_data)
        if song_data is None:
            print('Warning: {} does not exist in Spotify or in database'.format(song['name']))
            continue
        song_vector = song_data[number_cols].values
        song_vectors.append(song_vector)

    song_matrix = np.array(list(song_vectors))
    return np.mean(song_matrix, axis=0)

def flatten_dict_list(dict_list):
    flattened_dict = defaultdict()
    for key in dict_list[0].keys():
        flattened_dict[key] = []

    for dictionary in dict_list:
        for key, value in dictionary.items():
            flattened_dict[key].append(value)

    return flattened_dict

```

```
def recommend_songs( song_list, spotify_data, n_songs=10):

    metadata_cols = ['name', 'year', 'artists']
    song_dict = flatten_dict_list(song_list)

    song_center = get_mean_vector(song_list, spotify_data)
    scaler = song_cluster_pipeline.steps[0][1]
    scaled_data = scaler.transform(spotify_data[number_cols])
    scaled_song_center = scaler.transform(song_center.reshape(1, -1))
    distances = cdist(scaled_song_center, scaled_data, 'cosine')
    index = list(np.argsort(distances)[: , :n_songs][0])

    rec_songs = spotify_data.iloc[index]
    rec_songs = rec_songs[~rec_songs['name'].isin(song_dict['name'])]
    return rec_songs[metadata_cols].to_dict(orient='records')
```

```
In [17]: recommend_songs([{'name': 'Savage', 'year':2020},
                        {'name': 'DNA.', 'year': 2017},
                        {'name': 'Freedom', 'year': 2016},
                        {'name': 'Summertime Magic', 'year': 2018},
                        {'name': 'I need U', 'year': 2015}], data)
```

```
Out[17]: [{ 'name': 'GATTI',
            'year': 2019,
            'artists': "['JACKBOYS', 'Pop Smoke', 'Travis Scott']"},
  { 'name': "King's Dead (with Kendrick Lamar, Future & James Blake)",
    'year': 2018,
    'artists': "['Jay Rock', 'Kendrick Lamar', 'Future', 'James Blake']"},
  { 'name': 'Slidin', 'year': 2020, 'artists': "['21 Savage', 'Metro Boomin']"},
  { 'name': 'Black Skinhead', 'year': 2013, 'artists': "['Kanye West']"},
  { 'name': 'THICK', 'year': 2020, 'artists': "['DJ Chose', 'Beatking']"},
  { 'name': 'All Dat (with Megan Thee Stallion)',
    'year': 2020,
    'artists': "['Moneybagg Yo', 'Megan Thee Stallion']"},
  { 'name': 'Fight Back', 'year': 2018, 'artists': "['NEFFEX']"},
  { 'name': 'Shake The Room (feat. Quavo)',
    'year': 2020,
    'artists': "['Pop Smoke', 'Quavo']"},
  { 'name': 'Double G (feat. Pop Smoke)',
    'year': 2020,
    'artists': "['French Montana', 'Pop Smoke']"},
  { 'name': 'GOOBA', 'year': 2020, 'artists': "['6ix9ine']"}]
```

```
In [ ]:
```

```
In [ ]:
```