# Building Music Recommendation System using Spotify Dataset:

```python
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")
```

```python
In [2]:  data = pd.read_csv('data.csv')
         genre_data = pd.read_csv('data_by_genres.csv')
         year_data = pd.read_csv('data_by_year.csv')
         artist_data = pd.read_csv('data_by_artist.csv')
         wgenre_data = pd.read_csv('data_w_genres.csv')
```

```python
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
```

### Feature Correlation by considering a few features by using yellowbrick:

In [6]:
```python
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']

# Creating a list of the feature names
features = np.array(feature_names)
```
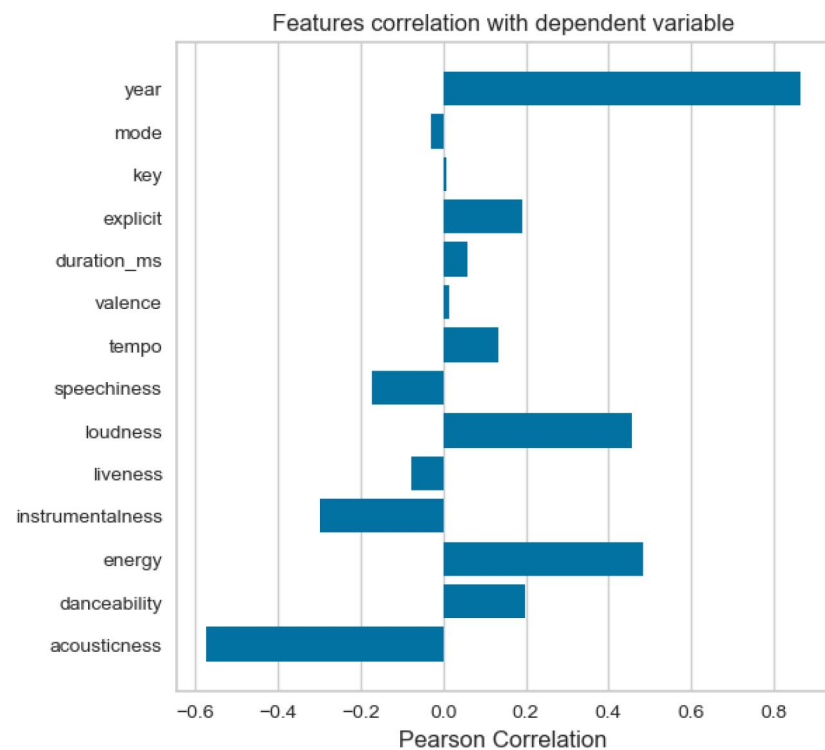
```python
# Instantiate the visualizer
visualizer = FeatureCorrelation(labels=features)

plt.rcParams['figure.figsize']=(6,6)
visualizer.fit(X, y)        # Fit the data to the visualizer
visualizer.show()
```

Features correlation with dependent variable



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

## Data Visualization:

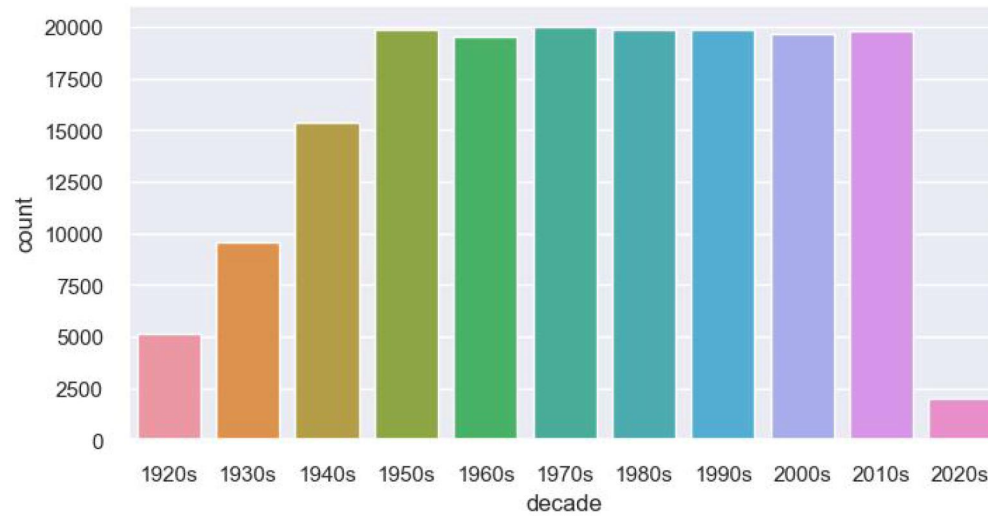### Grouping the music over the years:

```python
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':(8 ,4)})
         sns.countplot(x=data['decade'])
```
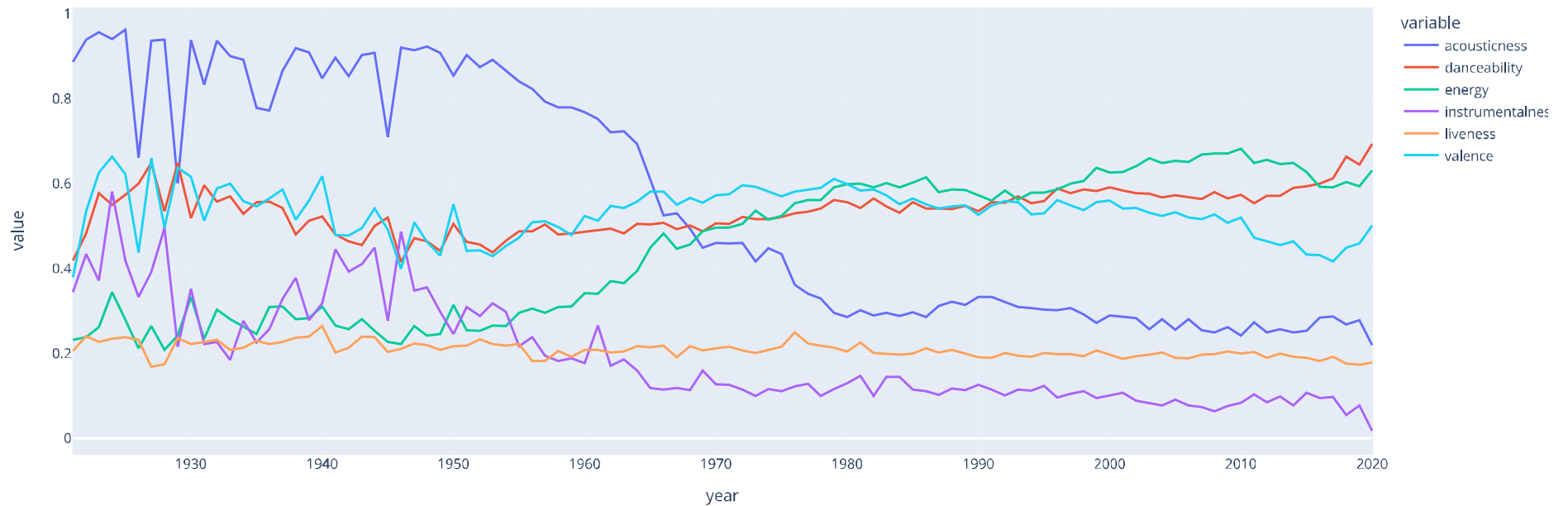
```
Out[7]:   <Axes: xlabel='decade', ylabel='count'>
```
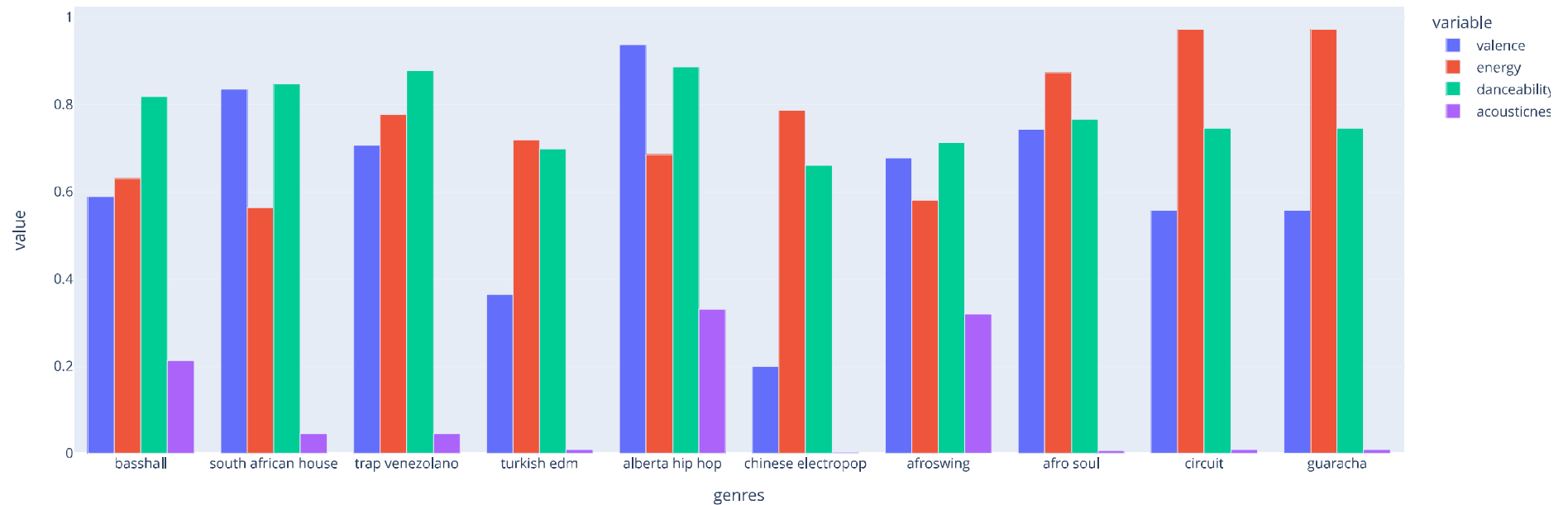
```
In [8]: sound_features = ['acousticness', 'danceability', 'energy', 'instrumentalness', 'liveness', 'valence']
        fig = px.line(year_data, x='year', y=sound_features)
        fig.show()
```

## Characteristics of Different Genres

```
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()
```



### Applying K-means clustering algorithm based on geners:

```
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, n_init="auto"))])
         X = genre_data.select_dtypes(np.number)
         cluster_pipeline.fit(X)
         genre_data['cluster'] = cluster_pipeline.predict(X)
```
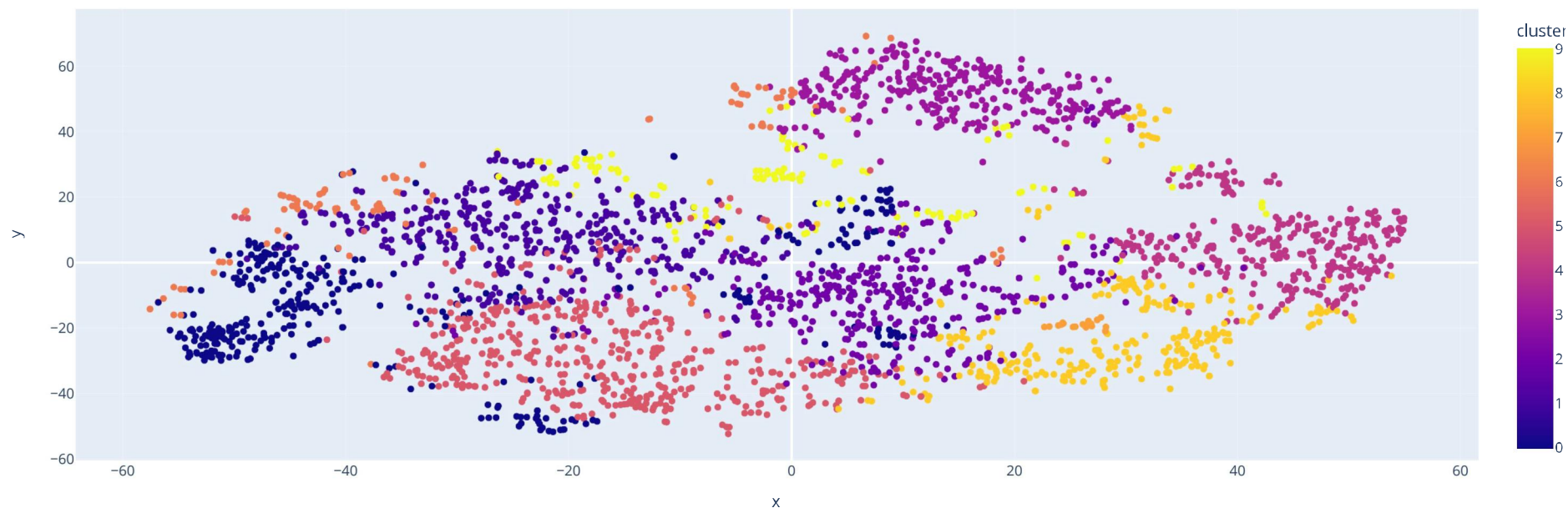
```
In [11]: # Visualizing the Clusters with t-SNE

         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 0.008s...
[t-SNE] Computed neighbors for 2973 samples in 0.310s...
[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: 76.106255
[t-SNE] KL divergence after 1000 iterations: 1.394141
```



**Applying K-means clustering algorithm based on songs:**

```
In [12]:  song_cluster_pipeline = Pipeline([('scaler', StandardScaler()), ('kmeans', KMeans(n_clusters=20, verbose=False, n_init=4))], 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
```
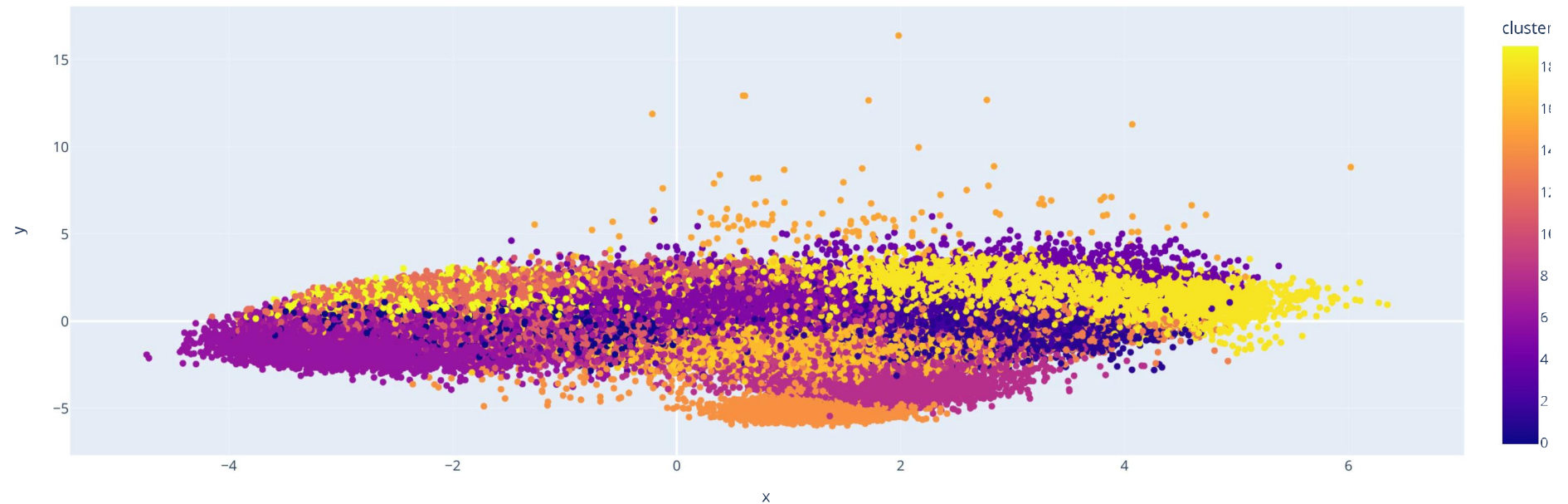
```
In [13]:  # Visualizing the Clusters with PCA

          from sklearn.decomposition import PCA
```

```python
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()
```



## Model Building of Music Recommendation System:

```python
import spotipy
import os
from spotipy.oauth2 import SpotifyClientCredentials
from collections import defaultdict

sp = spotipy.Spotify(auth_manager=SpotifyClientCredentials(client_id="SPOTIFY_CLIENT_ID",client_secret="SPOTIFY_CLIENT_SECRET"))

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]
```

```python
        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 [15]:
```python
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',


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)
```

```python
        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 [16]:
```python
recommend_songs([{'name': 'Come As You Are', 'year':1991},
                 {'name': 'Smells Like Teen Spirit', 'year': 1991},
                 {'name': 'Lithium', 'year': 1992},
                 {'name': 'All Apologies', 'year': 1993}], data)
```

Out[16]:
```
[{'name': 'Hanging By A Moment', 'year': 2000, 'artists': "['Lifehouse']"},
 {'name': 'Kiss Me', 'year': 1997, 'artists': "['Sixpence None The Richer']"},
 {'name': "Breakfast At Tiffany's",
  'year': 1995,
  'artists': "['Deep Blue Something']"},
 {'name': 'Otherside', 'year': 1999, 'artists': "['Red Hot Chili Peppers']"},
 {'name': "It's Not Living (If It's Not With You)",
  'year': 2018,
  'artists': "['The 1975']"},
 {'name': 'No Excuses', 'year': 1994, 'artists': "['Alice In Chains']"},
 {'name': 'Wherever You Will Go', 'year': 2001, 'artists': "['The Calling']"},
 {'name': 'Ballbreaker', 'year': 1995, 'artists': "['AC/DC']"},
 {'name': 'Runaway (U & I)', 'year': 2015, 'artists': "['Galantis']"},
 {'name': "Club Can't Handle Me (feat. David Guetta)",
  'year': 2010,
  'artists': "['Flo Rida', 'David Guetta']"}]
```