```
In [29]:   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 [30]:   data = pd.read_csv("data.csv")
           genre_data = pd.read_csv('data_by_genres.csv')
           year_data = pd.read_csv('data_by_year.csv')
```

```
In [31]:   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 [32]:  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 [33]:  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 [34]:    from sklearn.pipeline import make_pipeline
            from sklearn.cluster import KMeans
            from sklearn.preprocessing import StandardScaler

            # Assuming genre_data is a DataFrame and np has been imported as numpy

            # Create a pipeline and fit it in one go
            cluster_pipeline = make_pipeline(StandardScaler(), KMeans(n_clusters=10))
            genre_data['cluster'] = cluster_pipeline.fit_predict(genre_data.select_dtypes(include=[np.number]))
```

```
In [35]:    from sklearn.manifold import TSNE
            from sklearn.manifold import TSNE
            from sklearn.pipeline import make_pipeline
            from sklearn.preprocessing import StandardScaler
            import pandas as pd
            import plotly.express as px

            # Create a pipeline and apply t-SNE transformation in one step
            tsne_pipeline = make_pipeline(StandardScaler(), TSNE(n_components=2, verbose=1, random_state=42))
            projection = pd.DataFrame(tsne_pipeline.fit_transform(genre_data.select_dtypes(include=[np.number])),
                                      columns=['x', 'y'])
            projection['genres'] = genre_data['genres']
            projection['cluster'] = genre_data['cluster']

            # Create the scatter plot
            fig = px.scatter(projection, x='x', y='y', color='cluster', hover_data=['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 0.356s...
[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: 0.808325
[t-SNE] KL divergence after 250 iterations with early exaggeration: 70.496536
[t-SNE] KL divergence after 1000 iterations: 1.187491
```

```
In [36]:    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 [37]:    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()
```