

In [50]:

```
import numpy as np
import pandas as pd

import seaborn as sns
import plotly.express as px
import matplotlib.pyplot as plt
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")
```

In [51]:

```
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
```

In [52]:

```
data = pd.read_csv("data.csv")
genre_data = pd.read_csv('data_by_genre.csv')
year_data = pd.read_csv('data_by_year.csv')
```

In [53]:

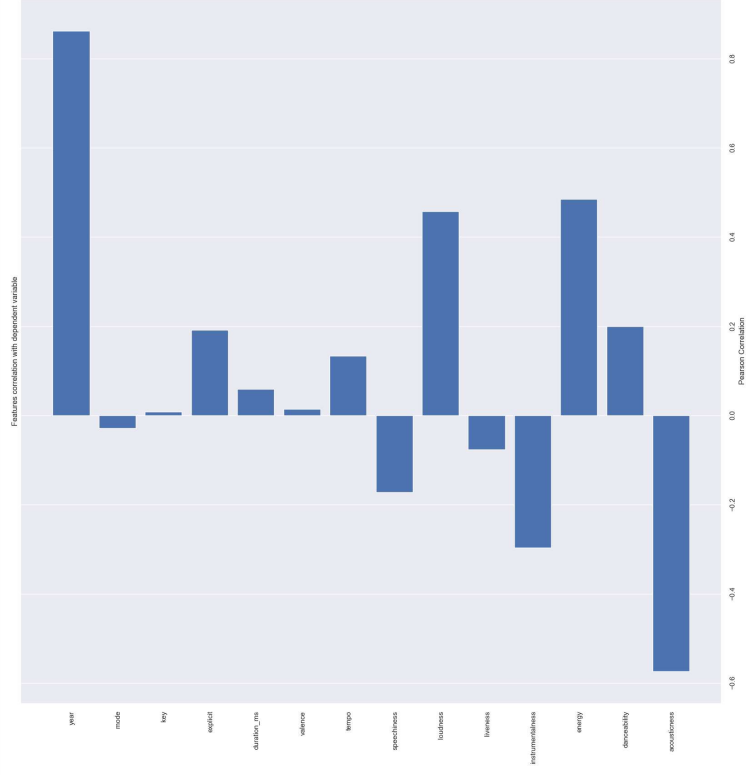
```
#!pip install yellowbrick
from yellowbrick.target import FeatureCorrelation

feature_names = ['acousticness', 'danceability', 'energy', 'instrumentalness',
                 'liveness', 'loudness', 'speechiness', 'tempo', 'valence', 'duration_ms']
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()
```



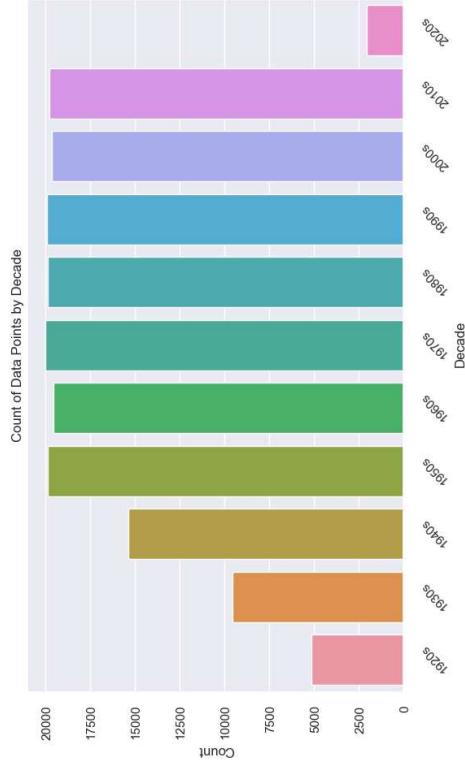
Out[53]: <Axes: title={'center': 'Features correlation with dependent variable'}, xlabel='Pearson Correlation'>

In [54]:

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

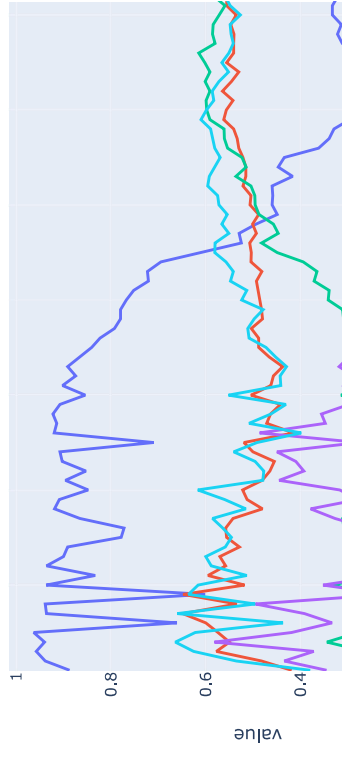
data['decade'] = data['year'].apply(get_decade)
# Specify the order of x-axis ticks
decade_order = sorted(data['decade'].unique()) # Sort the decades

sns.set(rc={'figure.figsize': (11, 6)})
sns.countplot(data=data, x='decade', order=decade_order) # Use 'x' to specify
plt.xticks(rotation=45) # Rotate x-axis labels for better readability
plt.xlabel('Decade')
plt.ylabel('Count')
plt.title('Count of Data Points by Decade')
plt.show()
```



In [55]:

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



In [56]:

```
top10_genres = genre_data.nlargest(10, 'popularity')
fig = px.bar(top10_genres, x='genres', y=['valence', 'energy', 'popularity', 'danceability'],
fig.show()
```

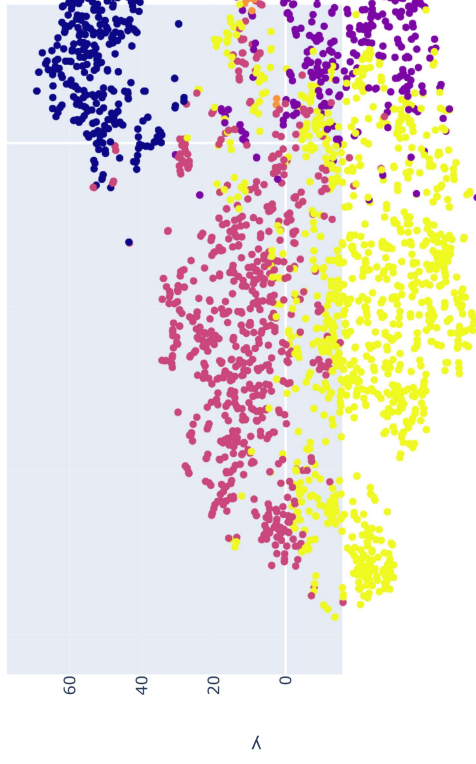


In [58]:

```
from sklearn.manifold import TSNE
tsne_pipeline = Pipeline([('scaler', StandardScaler()), ('tsne', TSNE(n_compon
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 0.0055...
[t-SNE] Computed neighbors for 2973 samples in 0.1565...
[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: 76.106102
[t-SNE] KL divergence after 1000 iterations: 1.392876
```



In [57]:

```
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
cluster_pipeline = Pipeline([('scaler', StandardScaler()), ('kmeans', KMeans(n
X = genre_data.select_dtypes(np.number)
cluster_pipeline.fit(X)
genre_data['cluster'] = cluster_pipeline.predict(X)
```

In [59]:

```
song_cluster_pipeline = Pipeline([('scaler', StandardScaler()),
                                   ('kmeans', KMeans(n_clusters=5,
                                                       verbose=False))
                                   ], verbose=False)

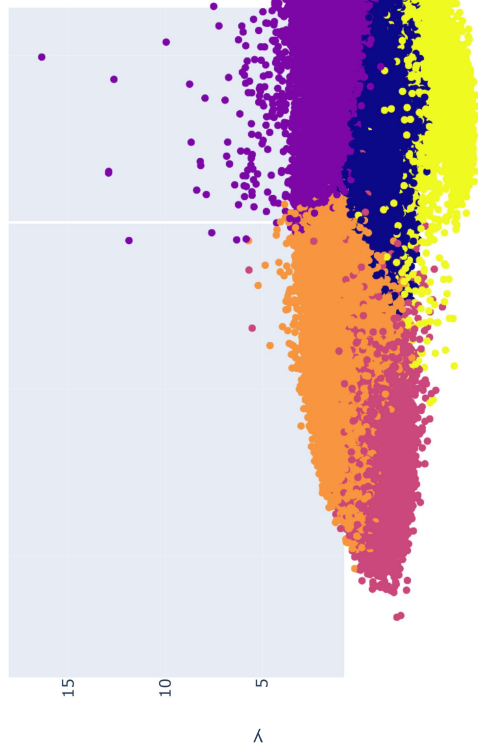
X = data.select_dtypes(np.number)
number_cols = list(X.columns)
song_cluster_pipeline.fit(X)
song_cluster_labels = song_cluster_pipeline.predict(X)
data['cluster_label'] = song_cluster_labels
```

In [60]:

```
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()
```

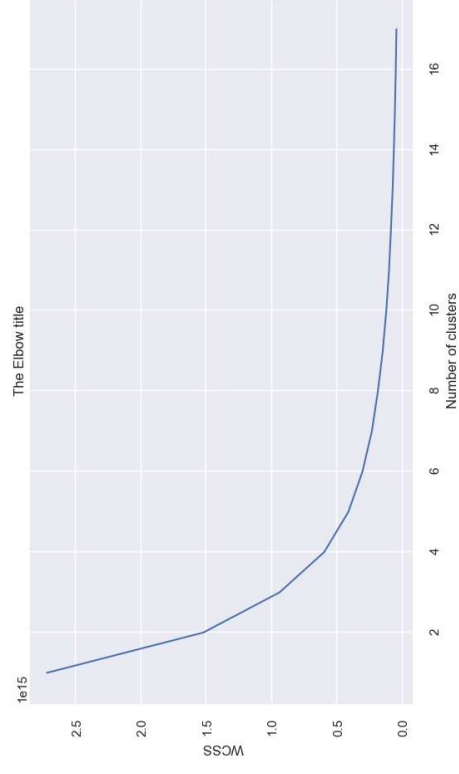


In [61]:

```
wcss=[]
for i in range(1,18):
    kmeans = KMeans(i)
    kmeans.fit(X)
    wcss_iter = kmeans.inertia_
    wcss.append(wcss_iter)

number_clusters = range(1,18)
plt.plot(number_clusters,wcss)
plt.title('The Elbow title')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
```

Out[61]: Text(0, 0.5, 'WCSS')



In []: