

Spotify Genres - Data Mining

Business objective

Maintaining radio stations or constructing playlists of high quality can be an expensive and daunting task. People enjoying these forms of listening to music may expect that songs offered by any particular source should be as similar to each other as possible - combining somber and lively pieces steer the audience towards the other providers. Therefore, it is crucial to find a number of sets/clusters of music tracks that could be viably combined together into a single playlist or radio station. The number of presented options should not be too high (maintaining many stations/playlists is not profitable) and tracks/songs should obviously be relatively similar to others composing the same cluster.

Dataset

The dataset we are using for this task can be found on kaggle (<https://www.kaggle.com/yamaerenay/spotify-dataset-19212020-160k-tracks>) and comprises data obtained using Spotify API. Although a variety of data can be found in this dataset, we will mostly limit ourselves to one of the files: `data_o.csv` (data regarding particular tracks), with every song being described with parameters portraying musical properties like:

- acousticness (0-1)
- danceability (0-1)
- energy (0-1)
- duration (usually between 200k and 300k in miliseconds)
- instrumentalness (0-1)
- valence (0-1)
- tempo (50-150)
- liveness (0-1)
- loudness (-60-0)
- speechiness (0-1)
- mode (0=minor, 1=major)
- key (categorical, 0-11)

																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																					</
--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	----

The dataset contains over 170k tracks; however, most of them are obscure and unknown. In order to provide a meaningful set of static playlists, it makes sense to limit ourself to a subset of the most popular songs.

3	0.607000	0.875175	0.9300	0.342880	0.741905	1.90067		0.309617	1.181892	0.679862	0.978999	0.713693
4	0.691763	0.6558124	0.810465	0.741040	0.816461	1.90067		0.309617	1.468977	0.608432	0.697381	0.713693
29995	1.010694	0.312578	0.355918	0.280777	0.123775	0.52613		0.309525	1.060786	0.549018	0.600168	1.401162
29996	0.248165	2.370524	0.438563	0.860227	1.064012	0.52613		0.208684	0.906807	0.295340	0.949114	0.713693
29997	0.531129	1.310955	0.683046	0.915093	0.961034	0.52613		0.309617	0.498616	0.613399	0.685938	0.713693
29998	0.770911	0.681769	1.358969	1.448862	0.028454	1.90067		0.309617	1.622595	0.771396	0.684025	0.713693
29999	0.335306	0.503984	1.515398	0.139958	0.114821	0.52613		0.309617	0.906807	0.375449	0.674391	0.713693

30000 rows × 13 columns

Quantiles of your data:

```
poplar_stand.quantile([0, .05, .25, .5, .75, .95, 1])
```

	valence	acousticness	danceability	duration_ms	energy	explicit	instrumentalness	key	liveness	loudness	mode
0.00	0.203376	0.941483	0.340198	2.587014	2.774258	0.52613	0.309617	1.468977	1.229943	10.264512	1.401162
0.05	1.602158	0.938995	1.781043	1.155939	1.882361	0.52613	0.309617	1.468977	0.844005	1.725365	1.401162
0.25	0.786897	0.839628	0.629917	0.951972	0.643144	0.52613	0.309617	0.906807	0.600421	0.340053	1.401162

Preprocessing

There are no missing values in our dataset regarding Spotify tracks, hence we do not need to impute any values. However, that does not mean there is nothing we should take care of before applying clusterisation algorithm.

Feature selection

The data being clustered should be free of any IDs or any attributes that would be able to distract the algorithm from "real" distinguishing factors. Therefore, the first action we take is removing the id/name/artist attributes.

In [6]:	<pre>tracks_popular = tracks_df.iloc[:30000] tracks_stand = tracks_df.drop(['id', 'name', 'release_date', 'year', 'artists', 'popularity'], axis=1)</pre>
---------	---

Standardisation

Since both outlier removal and clustering algorithm used by us are based on distance metrics, we need to standardise the values of numerical attributes in order to ensure equal impact of all attributes beforehand and after any operation that may "de-standardise" the data.

```
scaler = preprocessing.StandardScaler()

popular_stand = pd.DataFrame(
    scaler.fit_transform(popular_stand),
    index=popular_stand.index,
    columns=popular_stand.columns
)

Data following Outlier Removal and standardisation:

popular_stand.head(5)
```

	valence	acousticness	danceability	duration_ms	energy	explicit	instrumentalness	key	liveness	loudness	mode	speechiness
0	-1.454292	0.434836	0.774489	-0.361418	-0.213412	1.900193	-0.304748	-0.344686	-0.474724	-0.486443	-1.401677	-0.000000
1	0.986103	-0.182634	0.589854	-1.351256	0.457074	1.900193	-0.305009	0.499441	0.591935	0.980925	-1.401677	-0.000000
2	0.939096	-0.902330	0.863828	-0.453957	0.650570	-0.526262	-0.305009	0.218065	-0.604870	0.788616	-1.401677	-0.000000
3	-0.614046	-0.874201	1.989502	-0.630464	-0.748901	1.900193	-0.305009	-1.188813	-0.680006	0.0891128	0.713431	-0.000000
4	0.689334	0.664673	0.810224	-0.863742	0.817067	1.900193	-0.305009	-1.470188	-0.608225	0.707134	0.713431	-0.000000

30000 rows x 13 columns

Quantiles of standardised data:

```
pca = PCA().fit(popular_stand)
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.title('Explained variance vs number of PCA components')
plt.xlabel('Number of components')
plt.ylabel('Cumulative explained variance')
plt.show()
```

Outlier removal

Outliers can lead to abnormal behaviour of ML and DM algorithms; therefore, we use Local Outlier Factor-based outlier detection algorithm to get rid of them.

Local Outlier Factor measures the local deviation of density of a given sample with respect to its neighbours and in this way is able to determine how isolated any given sample is.

```

pca_fit_transform[popular_stand],
index=popular_stand.index,
columns=["PCA_"+i]).format(i for i in range(8))
)

tracks_pca.head(5)

```

	PCA_0	PCA_1	PCA_2	PCA_3	PCA_4	PCA_5	PCA_6	PCA_7
0	0.278460	1.855064	0.719560	0.409059	-0.077108	0.866782	0.297073	-1.336122
1	-1.694874	1.287312	-0.248167	0.946499	0.895249	-0.876597	0.225151	-1.259284
2	-1.662148	-0.046207	-0.715231	1.265405	0.154577	0.134244	0.955327	-0.556029
3	-1.458801	3.418890	0.320703	-1.819108	-0.542795	1.064744	0.248062	0.601677
4	-1.320683	0.709881	-0.837004	-1.757726	-0.568912	0.443914	0.625821	-0.063106

Data after mapping it using 8 PCA components after standardisation:

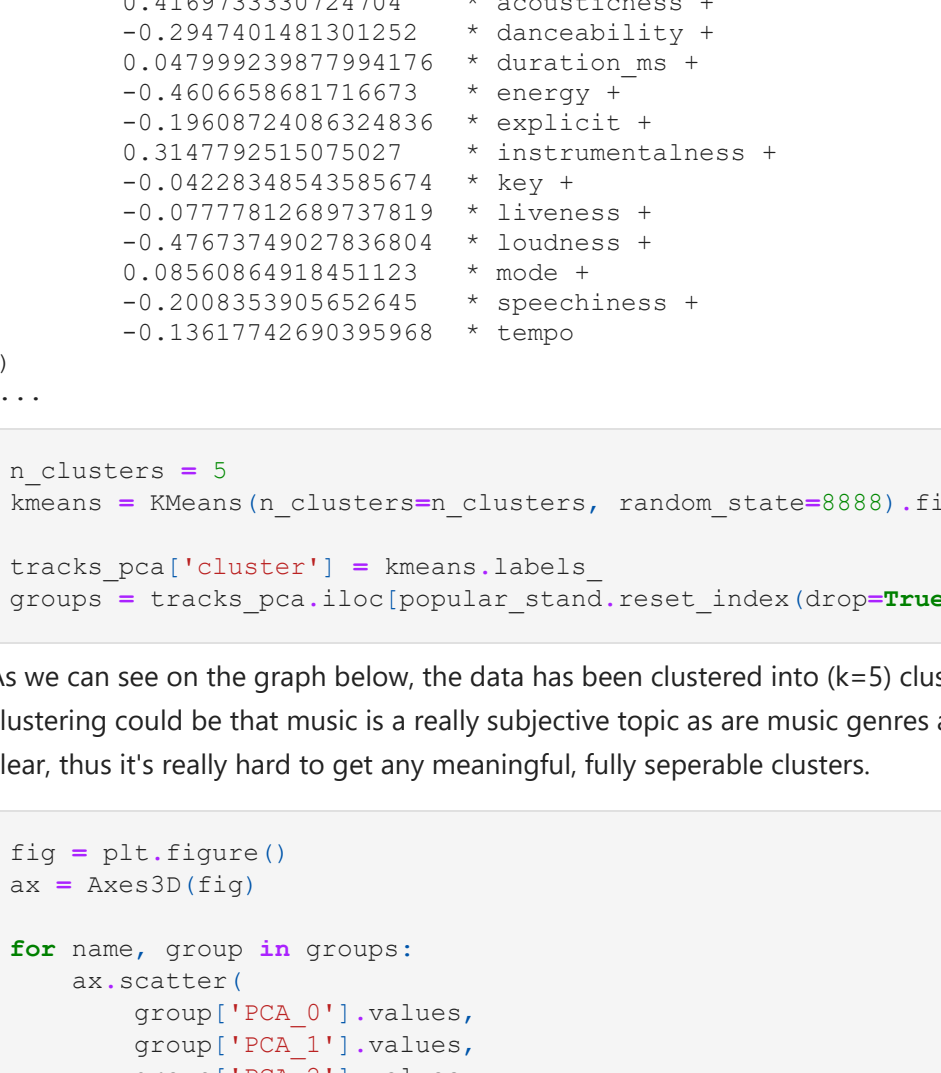
```

tracks_pca = pd.DataFrame(
    scaler.fit_transform(tracks_pca),

```

Dimensionality reduction

Distance-based methods may encountered negative impact of phenomena called the "curse of dimensionality" that occurs when a distance between two objects is calculated in a multi-dimensional environment is calculated. Dimensionality reduction algorithms, such as PCA, can be therefore used to reduce the number of dimensions and ensure correct representation of distances.

In [15]:	<pre>pca = PCA(n_fit(popular_stand)) plt.plot(np.cumsum(pca.explained_variance_ratio_)) plt.title("Explained variance vs number of PCA components") plt.xlabel("Number of components") plt.ylabel("Cumulative explained variance") plt.show()</pre>
Out [15]:	

```

    marker="o", label=name
)

ax.legend()
ax.set_xlabel("PCA_0")
ax.set_ylabel("PCA_1")
ax.set_zlabel("PCA_2")
plt.title("Clusters obtained after PCA preprocessing")
plt.show()

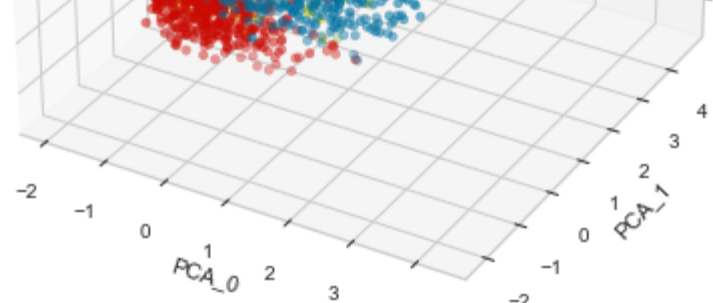
```

Clusters obtained after PCA preprocessing

Legend for Clusters:

- 0 (Blue)
- 1 (Green)
- 2 (Red)
- 3 (Purple)
- 4 (Yellow)

Data after mapping it using 8 PCA components after standardisation:



Clustering algorithm

For this clustering task, we have decided to use kMeans algorithm, which works by creating and recalculating centroids as long as any sample shifts its cluster. It needs to be initialised for a specific number of clusters (parameter k) that must be found experimentally.


Parameter tuning (number of clusters)

In order to establish a good number of clusters for our data, we use heuristic method called "elbow method" that involves finding an elbow of the curve of a measurement (in our case - sum of squared distances of sample to their cluster's centroid) as the number of clusters grows.

Created columns by PCA are really hard to interpret as each column is a linear combinations of all attributes. So for example the first column `PCA_0` (presented above) consists of combinations.

In [19]:	<pre>for ind, comp in enumerate(pca.components_): text = f"PCA_{ind} = " for i, val in enumerate(comp): if i != len(comp) - 1: text += f"{val}<2> + " else: text += f"{val}<2> + {popular_stand.columns[i]} + \n" print(text) break print("...")</pre>
Out [19]:	<pre>PCA_0 = -0.2930403707982217 * valence + 0.4169733330724704 * acousticness + -0.2947401481301252 * danceability + 0.047992923977994716 * duration_ms + -0.460658681716673 * energy + -0.19608724086324836 * explicit + 0.314792515075027 * instrumentalness + -0.04228348543958574 * key + -0.0777812689737819 * liveness + -0.07673749027836804 * loudness + 0.08560864918451123 * mode + -0.2008339305632645 * speechiness + 0.1361774260393948 * tempo ...</pre>

In [20]:	<pre>n_clusters = 5 kmeans = KMeans(n_clusters=n_clusters, random_state=8888).fit(tracks_pca) tracks_pca['cluster'] = kmeans.labels_ groups = tracks_pca.iloc[popular_stand.reset_index(drop=True).iloc[:20000].index.values].groupby('cluster')</pre>
Out [20]:	<pre>As we can see on the graph below, the data has been clustered into k=5 clusters but they are densely packed. One reason for this clustering could be that music is a really subjective topic as are music genres and the boundaries between each song may or may not be so clear, thus it's really hard to get any meaningful, fully separable clusters.</pre>

In [21]:	<pre>fig = plt.figure() ax = Axes3D(fig) for name, group in groups: ax.scatter(group['PCA_0'].values, group['PCA_1'].values, group['PCA_2'].values, marker="o", label=name) ax.legend() ax.set_xlabel("PCA_0") ax.set_ylabel("PCA_1") ax.set_zlabel("PCA_2") plt.title("Clusters obtained after PCA preprocessing") plt.show()</pre>
Out [21]:	

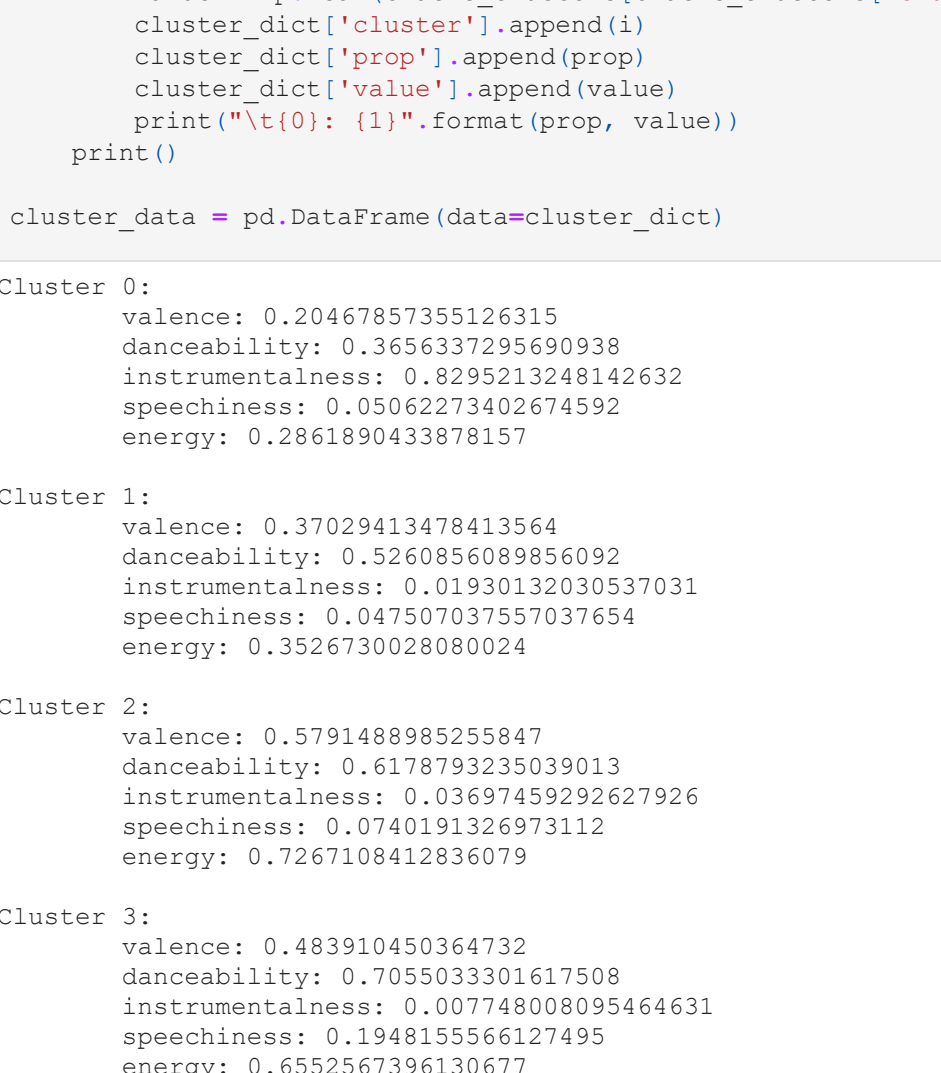
Clustering algorithm

For this clusterisation task, we have decided to use kMeans algorithm, which works by creating and recalculating clusters as long as any sample shifts its cluster. It needs to be initialised for a specific number of clusters (parameter k) that must be found experimentally.

Parameter tuning (number of clusters)

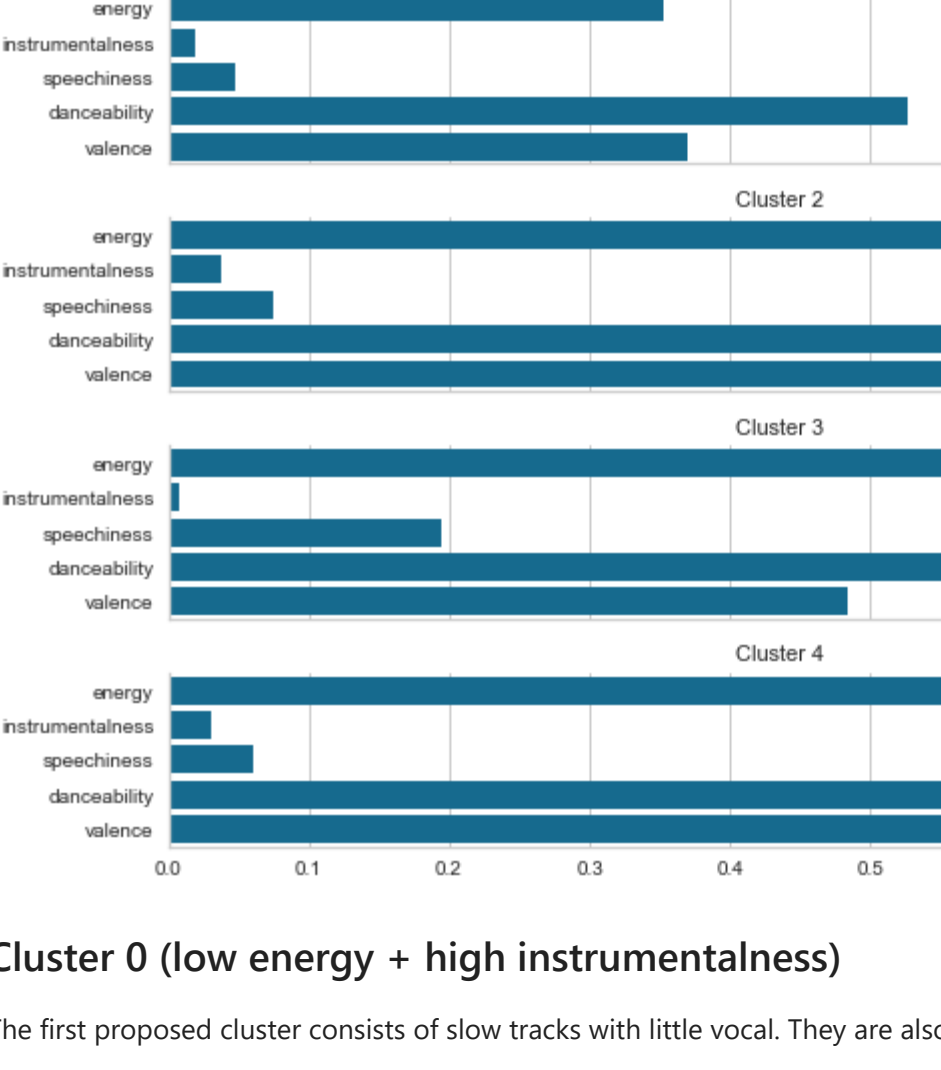
In order to establish a good number of clusters for our data, we use heuristic method called "elbow method" that involves finding an elbow of the curve of a measurement (in our case - sum of squared distances of sample to their cluster's centroid) as the number of clusters increases.

In [22]:	<pre>scores = [] for n_clusters in range(1, 20): clustering = KMeans(n_clusters = n_clusters, random_state=2500) clustering.fit(popular_stand) scores.append(clustering.inertia_)</pre>
Out [22]:	<pre>scores.append(clustering.inertia_)</pre>

In [23]:	<pre>sn.lineplot(x = range(1, 20), y = scores) plt.xlabel("Number of clusters") plt.ylabel("Sum of squared errors") plt.show()</pre>
Out [23]:	

The curve has no clearly defined elbow, which makes it difficult to select an optimal value of parameter k. Moreover, we do not know whether our data actually has any meaningful clusters. We can check it using Silhouette score - a measure of consistency of data within clusters that compares how similar an object is to its own vs other clusters. Consistently low values of that measure may suggest that the data has no natural division into clusters.

In [24]:	<pre>scores = [] for n_clusters in range(2, 10): clustering = KMeans(n_clusters = n_clusters, random_state=2500) clustering.fit(popular_stand) scores.append(silhouette_score(popular_stand, clustering.labels_))</pre>
Out [24]:	<pre>scores.append(silhouette_score(popular_stand, clustering.labels_))</pre>

In [25]:	<pre>sn.lineplot(x = range(2, 10), y = scores) plt.xlabel("Number of clusters") plt.ylabel("Mean silhouette score") plt.show()</pre>
Out [25]:	

In [26]:	<pre>clustering = KMeans(n_clusters = </pre>
----------	--