Re Re	equirement already satisfied: pyparsing>=2.3.1 in c:\users\baki akgun\new folder\lib\site-packages (from matplotlib->wordcloud) (3.0.9) equirement already satisfied: python-dateutil>=2.7 in c:\users\baki akgun\appdata\roaming\python\python311\site-packages (from matplotlib->wordcloud) (2.9.0.post0) equirement already satisfied: six>=1.5 in c:\users\baki akgun\appdata\roaming\python\python311\site-packages (from python-dateutil>=2.7->matplotlib->wordcloud) (1.16.0) import pandas as pd import nltk from IPython.display import display, HTML import matplotlib.pyplot as plt import seaborn as sns from nltk.sentiment import SentimentIntensityAnalyzer
[3]: [r [r	import numpy as np from sklearn.feature_extraction.text import TfidfVectorizer nltk.download('vader_lexicon') nltk_data] Downloading package vader_lexicon to C:\Users\Baki nltk_data] Akgun\AppData\Roaming\nltk_data nltk_data] Package vader_lexicon is already up-to-date!
	True data = pd.read_csv("Music Informations and Lyrics_ from Spotify and Musixmatch.csv") df = data.copy() def getting_primary_info(df): print("
	<pre>print("Veri seti değişken tipleri:\n", df.dtypes) print("</pre>
Ve Ve Ar Ar	getting_primary_info(df) Peri setinin şekli (904, 47) Peri seti değişken tipleri: Track Name object Petist object Petist object Petist Popularity int64
Tr Re Tr Da Er Ke Lo	oudness float64 ode int64
Ac Ir Li Va Te Ti Ly	seechiness float64 cousticness float64 coustic
Av Av Va Va Va	rerage2 float64 rerage3 float64 rerage4 float64 rerage5 float64 rerience1 float64 rerience2 float64 rerience3 float64 rerience3 float64 rerience4 float64 rerience5 float64
S1 S1 S1 Mi Mi Mi	randard Deviation1 float64 randard Deviation2 float64 randard Deviation3 float64 randard Deviation3 float64 randard Deviation4 float64 randard Deviation5 float64 randard Deviation6 float64 randard Deviation8 fl
Ma Ma Ma Ma Tr Ar	Inimum5 float64 aximum1 float64 aximum2 float64 aximum3 float64 aximum4 float64 aximum5 float6
Ve	Track Artist Artist Track Release Track Danceability Energy Key Loudness Mode Speechiness Acousticness Instrumentalness Liveness Valence Tempo Time_signature Ly Popularity Genres Duration(ms) Date Popularity Center Popularity Genres Duration(ms) Date Popularity Date Pop
	just a feed to the state of the
0	Nanana Gou house Na-na-na-na-na-na-na-na-na-na-na-na-na-na
	na-na-na na\nNa na-na na-na na\n\n\n explain, a fee that I ju can't er just a fee
	t won't\nM le ber because somet th on\r
co m	Artist Popularity Mark Track Duration(ms) Popularity Duration(ms) Popularity Duration(ms) Durat
	min 22.000000 91554.000000 0.000000 0.0122000 0.006200 0.000000 -18.439000 0.000000 0.022700 0.000013 0.000000 0.026400 0.037100 65.043000 1.000000 4.545511 25% 68.750000 178358.000000 68.000000 0.545000 0.520750 1.000000 -8.114750 0.000000 0.034400 0.033775 0.000000 0.094300 0.297750 98.033000 4.000000 31.499464 50% 78.000000 205803.000000 77.000000 0.648500 0.640500 5.000000 -6.255500 1.000000 0.048350 0.136000 0.000002 0.120000 0.465000 120.002000 4.000000 4.000000 4.000000 4.000000 4.000000 53.188298 4.000000 4.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000
[7]:	<pre>import nltk from nltk.stem.porter import PorterStemmer stemmer = PorterStemmer() from sklearn.feature_extraction.text import TfidfVectorizer from sklearn.metrics.pairwise import cosine_similarity df['Lyrics'] = df['Lyrics'].str.lower().replace(r'^\w\s', ' ').replace(r'\n', ' ', regex = True)</pre>
[]: [10]: [<pre>def tokenization(txt): tokens = nltk.word_tokenize(txt) stemming = [stemmer.stem(w) for w in tokens] return " ".join(stemming) df['Lyrics'] = df['Lyrics'].apply(lambda x: tokenization(x)) from wordcloud import WordCloud import matplotlib.nyplot as plt</pre>
	<pre>import matplotlib.pyplot as plt lyrics_text = ' '.join(df['Lyrics']) wordcloud = WordCloud(width=800, height=400, background_color='white').generate(lyrics_text) plt.figure(figsize=(10, 5)) plt.imshow(wordcloud, interpolation='bilinear') plt.axis('off') plt.show()</pre>
0 (++	left gonna thought turn of the start alone name dream na na without fuck think ooh ooh light of turn of the start alone name dream na na without fuck think ooh ooh light of turn of the start alone name dream na na without fuck think ooh ooh light of turn of the start alone name dream na na without fuck think ooh ooh light of turn of the start alone name dream na na without fuck think ooh ooh light of turn of the start alone name dream na na without fuck think ooh ooh light of turn of the start alone name dream na na without fuck think ooh ooh light of turn of the start alone name dream na na without fuck think ooh ooh light of turn of tur
	que city stop oh nigga will break show nigga will hold no break show never heart hear ah ah ah told no break show hold no break
11]:	
	<pre>def recommendation(song_df): print(song_df) idx = df[df['Track Name'] == song_df].index[0] distances = sorted(list(enumerate(similarity[idx])),reverse=True,key=lambda x:x[1]) songs = [] for m_id in distances: song_name = df.iloc[m_id[0]]["Track Name"]</pre>
	<pre>song_name = df.iloc[m_id[0]]["Track Name"] cosine_similarity = distances[m_id[0]][1] songs.append((song_name, cosine_similarity)) return songs max_popularity_index = df['Track Popularity'].idxmax() song_with_max_popularity = df.loc[max_popularity_index, 'Track Name']</pre>
Er 14]:	print("En popüler şarkı:", song_with_max_popularity) n popüler şarkı: Cruel Summer recommendations = recommendation("Cruel Summer") recommendations_df = pd.DataFrame(recommendations, columns=["Track Name", "Cosine Similarity"]) df["Cosine Similarity"] = df["Track Name"].apply(lambda x: recommendations_df[recommendations_df["Track Name"] == x]["Cosine Similarity"].values[0] if x in recommendations_df["Track Name"].values else np.nan)
	Track Name Artist Popularity Genres Duration(ms) Date Popularity Danceability Energy Key Minimum4 Minimum5 Maximum1 Maximum2 Maximum3 Maximum4 Maximum4 Maximum5 Popularity Out (It Goes
	Like) Nanana - Edit
	Month Harmony ohio Contemporary country, country pop, mo Chayce 61 Modern 228144 2021- 82 0.626 0.583 0.152565 0.583 0.152565 0.377560 185 25354 1338 29640 768 20624 578 84380 364 90048 NaN Na
16]:	<pre>forows × 48 columns sid = SentimentIntensityAnalyzer() sentiments = [] for review in df['Lyrics']:</pre>
17]:	<pre>sentiment_score = sid.polarity_scores(review) sentiments.append(sentiment_score) sentiments_df = pd.DataFrame(sentiments) df = pd.concat([df, sentiments_df], axis=1) df.drop(['Track Popularity Output', 'Artist Popularity Output'], axis=1, inplace=True) len(df["Artist"].unique())</pre>
	<pre>def ms_to_min(milliseconds): minutes = milliseconds / (1000 * 60) return minutes df['Track Duration(dk)'] = df['Track Duration(ms)'].apply(lambda x: ms_to_min(x))</pre>
21]:	<pre>df.drop(columns = "Track Duration(ms)",inplace = True) # import seaborn as sns # import matplotlib.pyplot as plt # numeric_columns = df.select_dtypes(include=['number']) # # Cross-correlation matrisini hesaplayalım # correlation_matrix = numeric_columns.corr() # correlation_matrix</pre>
22]:	<pre># correlation_matrix def custom_mean(group): return group.mean() rolling_avg = df.groupby('Artist')['Energy'].rolling(window=2).apply(custom_mean, raw=False) # Hesaplanan kayan ortalamayı yeni bir sütun olarak DataFrame'e ekleyelim df['rolling_avg_energy'] = rolling_avg_reset_index(level=0, drop=True).values</pre>
23]: [24]: [[]: [<pre>df['rolling_avg_energy'] = df['rolling_avg_energy'].fillna(df['Energy']) df['Year'] = 2024 - df['Release Date'].str[:4].astype(int) df.drop(columns = "Release Date" , inplace = True)</pre>
25]:	<pre>general_genres = ["Pop", "Rock", "Hip Hop", "R&B", "Jazz", "Electronic", "Country", "Reggae",</pre>
	"Folk", "Classical", "Blues", "Metal", "Indie", "Punk", "Funk", "Disco", "Alternative",
	"Alternative", "Soul", "Gospel", "Dance", "Latin", "World", "Experimental", "Ambient", "Acoustic", "Instrumental"
26]:	"Instrumental"] df = df.dropna(subset=['Artist Genres']) for genre in general_genres: df[genre] = df['Artist Genres'].apply(lambda x: 1 if genre.lower() in x.lower() else 0) df.drop(columns = "Artist Genres",inplace = True) df.to_csv("deneme.csv" , index = False)
28]: 29]: 30]:	<pre>import pandas as pd df = pd.read_csv("Music Informations and Lyrics_ from Spotify and Musixmatch.csv") dff = pd.read_csv("deneme.csv") dff.head()</pre>
	Track Name Artist Popularity Popu
5	3 2016 Sam Hunt 68 54 0.520 0.384 6 -8.168 1 0.0362 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
	<pre>from sklearn.linear_model import LinearRegression from sklearn.metrics import r2_score from sklearn.metrics import r2_score from sklearn.model_selection import train_test_split from sklearn.ensemble import RandomForestRegressor from sklearn.ensemble import GradientBoostingRegressor X , y = dff.drop(columns=["Artist Popularity", "Track Popularity", "Artist", "Track Name", "Lyrics"]), dff["Artist Popularity"] X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=47)</pre>
	<pre># Linear Regression model = LinearRegression() model.fit(X_train, y_train) y_pred = model.predict(X_test) r2 = r2_score(y_test, y_pred) print("Linear Regression Test R^2 Score:", r2)</pre>
	<pre>y_pred_train = model.predict(X_train) r2_train = r2_score(y_train, y_pred_train) print("Linear Regression Train R^2 Score:", r2_train) print("") # Random Forest model_random_forest = RandomForestRegressor(n_estimators=20, min_samples_split=2) model_random_forest.fit(X_train, y_train)</pre>
	<pre>y_pred_rf = model_random_forest.predict(X_test) r2_rf = r2_score(y_test, y_pred_rf) print("Random Forest Test R^2 Score:", r2_rf) y_pred_rf_train = model_random_forest.predict(X_train) r2_rf_train = r2_score(y_train, y_pred_rf_train) print("Random Forest Train R^2 Score:", r2_rf_train)</pre>
	<pre>print("") # Gradient Boosting model_gradient_boosting = GradientBoostingRegressor(n_estimators=80, min_samples_split=2) model_gradient_boosting.fit(X_train, y_train) y_pred_gb = model_gradient_boosting.predict(X_test) r2_gb = r2_score(y_test, y_pred_gb) print("Gradient Boosting Test R^2 Score:", r2_gb)</pre>
Li Li	
Ra Gr Gr 32]:	Andom Forest Train R^2 Score: 0.8978323740515428 Aradient Boosting Test R^2 Score: 0.2977105848852125 Aradient Boosting Train R^2 Score: 0.7179555979966388 From sklearn.linear_model import LinearRegression From sklearn.metrics import r2_score From sklearn.model_selection import train_test_split
	<pre>from sklearn.ensemble import RandomForestRegressor from sklearn.ensemble import GradientBoostingRegressor X , y = dff.drop(columns=["Artist Popularity", "Track Popularity", "Artist", "Track Name", "Lyrics"]), dff["Track Popularity"] X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=47) # Linear Regression model = LinearRegression()</pre>
	<pre>model = Linearregression() model.fit(X_train, y_train) y_pred = model.predict(X_test) r2 = r2_score(y_test, y_pred) print("Linear Regression Test R^2 Score:", r2) y_pred_train = model.predict(X_train) r2_train = r2_score(y_train, y_pred_train) print("Linear Regression Train R^2 Score:", r2_train)</pre>
	<pre>print("Linear Regression Train R^2 Score:", r2_train) print("") # Random Forest model_random_forest = RandomForestRegressor(n_estimators=20, min_samples_split=2) model_random_forest.fit(X_train, y_train) y_pred_rf = model_random_forest.predict(X_test)</pre>
	<pre>y_pred_rT = model_random_forest.predict(X_test) r2_rf = r2_score(y_test, y_pred_rf) print("Random Forest Test R^2 Score:", r2_rf) y_pred_rf_train = model_random_forest.predict(X_train)</pre>
	r2_rf_train = r2_score(y_train, y_pred_rf_train) print("Random Forest Train R^2 Score:", r2_rf_train) print("")
	r2_rf_train = r2_score(y_train, y_pred_rf_train) print("Random Forest Train R^2 Score:", r2_rf_train)