

Executive Summary – SPOT-POP

Objective:

The main objective of this project is to develop a predictive model that can classify the popularity of Spotify tracks based on their audio features. By analyzing a dataset containing various attributes of Spotify tracks, we aim to identify key factors that contribute to a track's popularity and create a reliable predictive system.

Data:

The data for this project is sourced from Kaggle, which can be found [here](#). The dataset contains over 100,000 song entries and 20 different features for each entry, making it a suitable dataset for data exploration, cleaning, and training.

Primary Variables and Features:

After extensive data exploration and cleaning we chose the following features as our input (*Loudness, explicit, danceability, time_signature, energy, tempo, key, liveness, valence, mode, speechiness, duration, acousticness, instrumentalness*) and the (*popularity decision*) as the output of the model.

Modeling Approach:

We are leveraging a classification modeling approach to predict whether a song entry will be “popular” or “not popular”. To prepare the data for this approach we categorized the popularity data to being popular and not popular. We are using sklearn library and we are splitting our data with the stratify option to maintain a similar value count between our splits in different categories with the test size being 30%. Logistic Regression, KNN and Random Forest classification approaches are used.

Model Performances

	LogR	KNN (K=30)	RF																																													
Accuracy	73.33%	73.49%	74.94%																																													
Conf Matrix	<table border="1"> <tr> <td>True label \ Predicted label</td> <td>Not Popular</td> <td>Popular</td> </tr> <tr> <td>Not Popular</td> <td>17120</td> <td>368</td> </tr> <tr> <td>Popular</td> <td>6124</td> <td>383</td> </tr> </table>	True label \ Predicted label	Not Popular	Popular	Not Popular	17120	368	Popular	6124	383	<table border="1"> <tr> <td>True label \ Predicted label</td> <td>Not Popular</td> <td>Popular</td> </tr> <tr> <td>Not Popular</td> <td>16549</td> <td>739</td> </tr> <tr> <td>Popular</td> <td>5015</td> <td>792</td> </tr> </table>	True label \ Predicted label	Not Popular	Popular	Not Popular	16549	739	Popular	5015	792	<table border="1"> <tr> <td>True label \ Predicted label</td> <td>Not Popular</td> <td>Popular</td> </tr> <tr> <td>Not Popular</td> <td>16719</td> <td>569</td> </tr> <tr> <td>Popular</td> <td>5343</td> <td>964</td> </tr> </table>	True label \ Predicted label	Not Popular	Popular	Not Popular	16719	569	Popular	5343	964																		
True label \ Predicted label	Not Popular	Popular																																														
Not Popular	17120	368																																														
Popular	6124	383																																														
True label \ Predicted label	Not Popular	Popular																																														
Not Popular	16549	739																																														
Popular	5015	792																																														
True label \ Predicted label	Not Popular	Popular																																														
Not Popular	16719	569																																														
Popular	5343	964																																														
Metrics	<table border="1"> <tr> <td></td> <td>precision</td> <td>recall</td> <td>f1-score</td> <td>support</td> </tr> <tr> <td>Not Popular</td> <td>0.74</td> <td>0.99</td> <td>0.84</td> <td>17288</td> </tr> <tr> <td>Popular</td> <td>0.52</td> <td>0.03</td> <td>0.05</td> <td>6307</td> </tr> </table>		precision	recall	f1-score	support	Not Popular	0.74	0.99	0.84	17288	Popular	0.52	0.03	0.05	6307	<table border="1"> <tr> <td></td> <td>precision</td> <td>recall</td> <td>f1-score</td> <td>support</td> </tr> <tr> <td>Not Popular</td> <td>0.75</td> <td>0.96</td> <td>0.84</td> <td>17288</td> </tr> <tr> <td>Popular</td> <td>0.52</td> <td>0.13</td> <td>0.20</td> <td>6307</td> </tr> </table>		precision	recall	f1-score	support	Not Popular	0.75	0.96	0.84	17288	Popular	0.52	0.13	0.20	6307	<table border="1"> <tr> <td></td> <td>precision</td> <td>recall</td> <td>f1-score</td> <td>support</td> </tr> <tr> <td>Not Popular</td> <td>0.76</td> <td>0.97</td> <td>0.85</td> <td>17288</td> </tr> <tr> <td>Popular</td> <td>0.63</td> <td>0.15</td> <td>0.25</td> <td>6307</td> </tr> </table>		precision	recall	f1-score	support	Not Popular	0.76	0.97	0.85	17288	Popular	0.63	0.15	0.25	6307
	precision	recall	f1-score	support																																												
Not Popular	0.74	0.99	0.84	17288																																												
Popular	0.52	0.03	0.05	6307																																												
	precision	recall	f1-score	support																																												
Not Popular	0.75	0.96	0.84	17288																																												
Popular	0.52	0.13	0.20	6307																																												
	precision	recall	f1-score	support																																												
Not Popular	0.76	0.97	0.85	17288																																												
Popular	0.63	0.15	0.25	6307																																												

While the models have performed well on determining “Not Popular” songs, the model is under performing on predicting “Popular” songs, which can be attributed to the smaller portion of popular data.

Next Steps

Future work includes better feature selection schemes and working with GridSearchCV to find the best fit for each model.