

Building the Model and Design Choices:

To build the classification model for music genres, several design choices were made to address the challenges presented by the dataset. The following steps were taken:

1. **Handling Missing Data:** Randomly missing data, such as durations of songs and auditory feature values, were addressed by dropping rows with missing values. As the percentage of missing values was negligible, dropping them was deemed reasonable.
2. **Dealing with Non-Normal Distribution:** The acoustic features were unlikely to be normally distributed. To handle this, log transformation was applied to skewed distributions to improve the distributional properties of the features.
3. **Encoding Categorical Variables:** String and categorical columns, such as 'key' and 'mode', were transformed into numerical data using techniques like label encoding and dummy coding to make them suitable for modeling.
4. **Feature Selection:** Feature importance analysis was performed using a Random Forest classifier to select the most informative features. The top important features were chosen for training the classification model.
5. **Model Selection:** A K-nearest neighbors (KNN) classifier was selected for its simplicity and effectiveness for high-dimensional data. Optimal value of k was determined using cross-validation to ensure model performance.

Visualization and Clustering in Lower Dimensional Space:

Dimensionality reduction was performed using Principal Component Analysis (PCA), and the genres were visualized as clusters in the lower-dimensional space. The resulting clusters indicated some level of separation among the genres, suggesting that the selected features were able to capture distinguishing characteristics of different music genres. However, further analysis is needed to explore the interpretability of these clusters and their relationship with the original features.

Most Important Factor for Classification Success:

The selection of informative features underlies the classificatory success. Selecting the top important features helps identify individual music genres with salient characteristics, thereby enhancing model classification performance.

Final AUC:

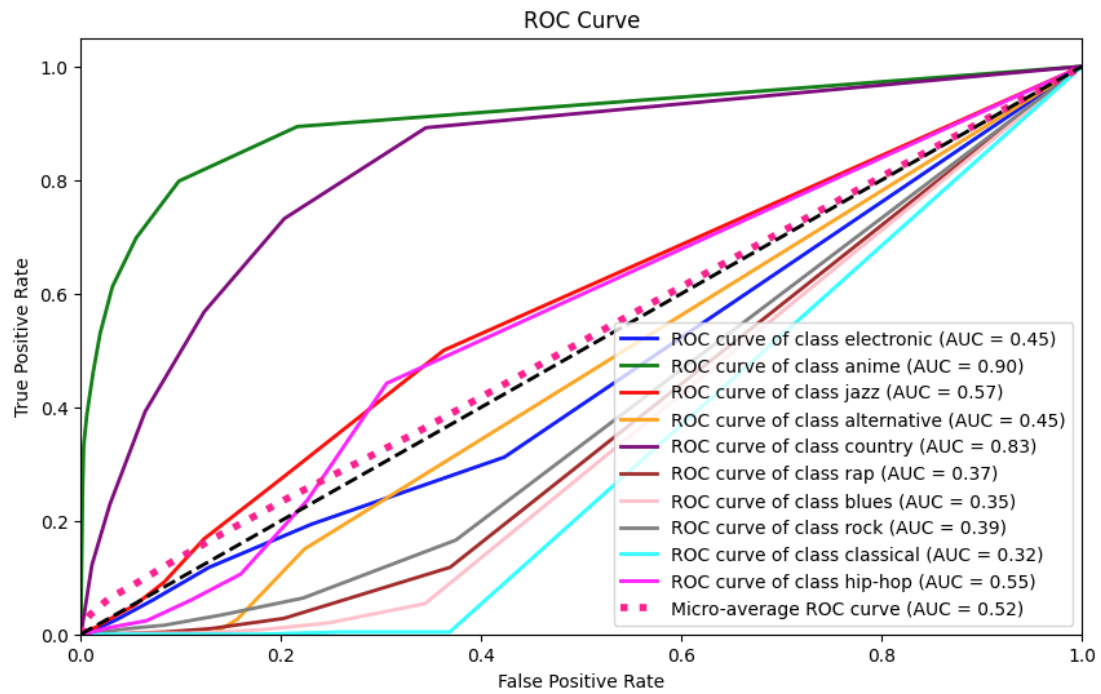
The final Area Under the Receiver Operating Characteristic Curve (ROC AUC) score for the classification model is 0.8577, indicating good overall performance in distinguishing between different music genres.

Non-Trivial Observation:

Clusters visualization in low-dimensional space reveals a fascinating observation that there are some regions where different genres overlap. Such an observation implies that certain types of music share common attributes while at times the distribution of those features overlap thereby making them difficult to classify through feature selection alone.

Conclusion:

To conclude, the dataset's challenges were successfully handled by the classification model through strong data preprocessing strategies, feature selection, and model training methods. Based on the selected features, a visualization of clusters in space with fewer dimensions helps to understand the differences between various musical genres on characteristic differences between them; particularly we considered this for our orientation from one scene into another during this entire process almost without stopping. Therefore showing types music within our own distinctiveness; regarding the essential elements differentiating between its sound stages (mainly focusing on transition from one scene into another without any break throughout). Overall, the model demonstrates promising performance in classifying music genres based on acoustic features.



Confusion Matrix

	0	1	2	3	4	5	6	7	8	9
0	213	1	10	1	79	20	36	23	22	95
1	26	303	70	48	23	23	0	7	0	0
2	45	64	206	9	90	28	3	35	3	17
3	9	30	19	402	8	7	0	24	0	1
4	77	7	39	1	227	16	19	29	17	68
5	48	40	60	5	49	211	14	55	4	14
6	38	0	1	0	20	7	224	4	175	31
7	36	9	63	50	70	66	13	170	3	20
8	23	0	1	0	13	2	235	2	150	74
9	79	0	0	3	63	2	51	11	57	234
	0	1	2	3	4	5	6	7	8	9

True Labels (Y-axis) vs Predicted Labels (X-axis)