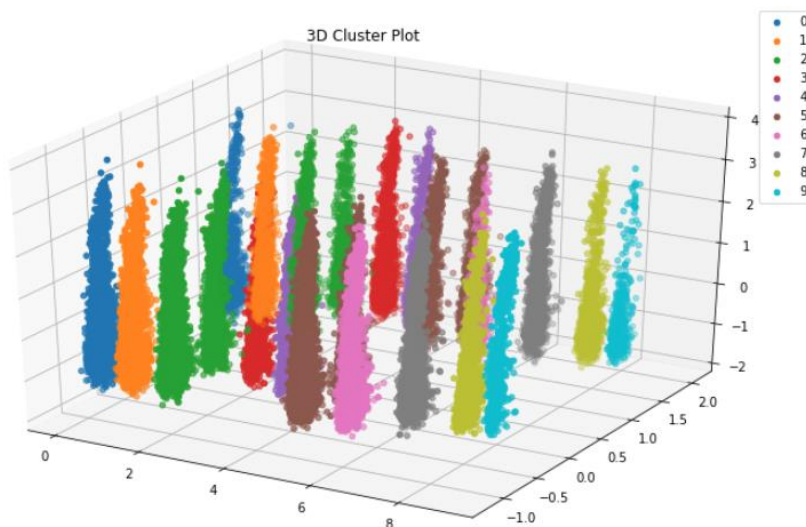


I did a train/test split so that the testing set consists of 500 randomly picked songs for one genre, and 5000 songs in total (since there are 10 genres). The training set has the rest 45000 songs.

For data cleansing, I found that there are 5 NaN values for every column. Therefore, I dropped these NaN so that the dataset becomes exactly 50000 rows. There are also missing values denoted as “?” in column “tempo” of object data type. Take column “tempo”, for example. I checked that the valid values in the “tempo” are all positive. Thus, I decided to replace all missing values with string “0” first and cast the data type from object to float. I then replaced all 0’s in “tempo” with the mean column value (since missing values are replaced as 0, this would give us the exact mean of the column) and replaced any missing values in other numerical columns with the column mean respectively. I also dropped columns “instance\_id”, “artist\_name”, “track\_name”, “obtained\_data” for the time being.

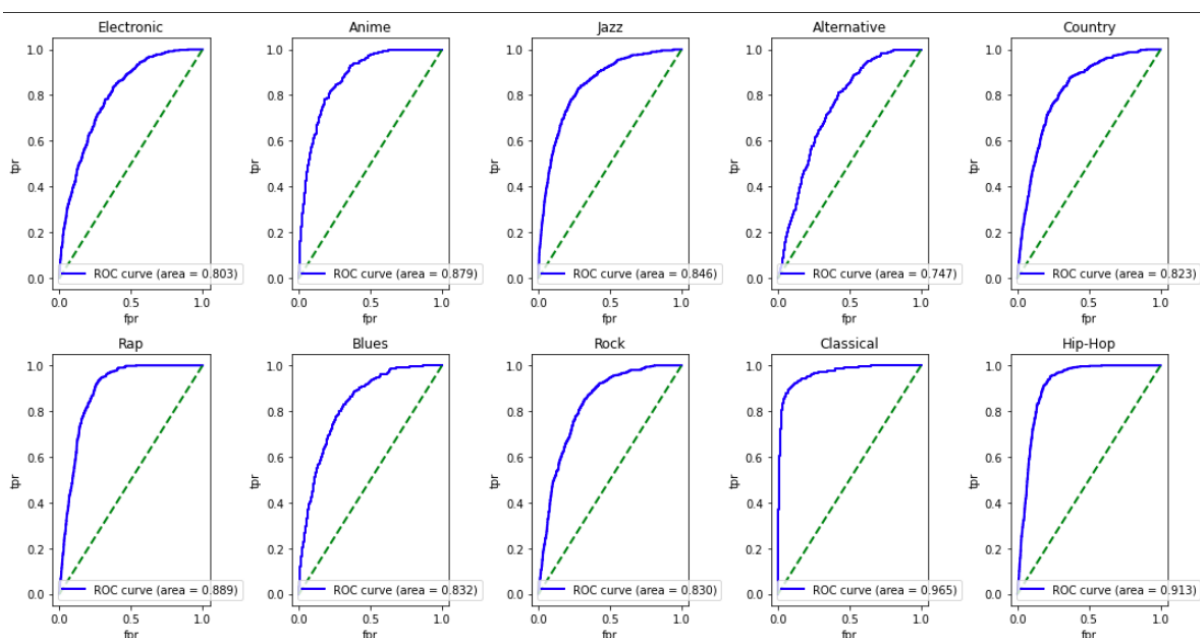
In terms of dimensional reduction, I one-hot encoded the categorical variables like “mode”, “key” using get\_dummies function. I used StandardScaler to compute z-scores for features in the training data and mapped those to the testing data for consistency. Since categorical variables have values either 0 or 1, StandardScaler would keep them intact. I applied



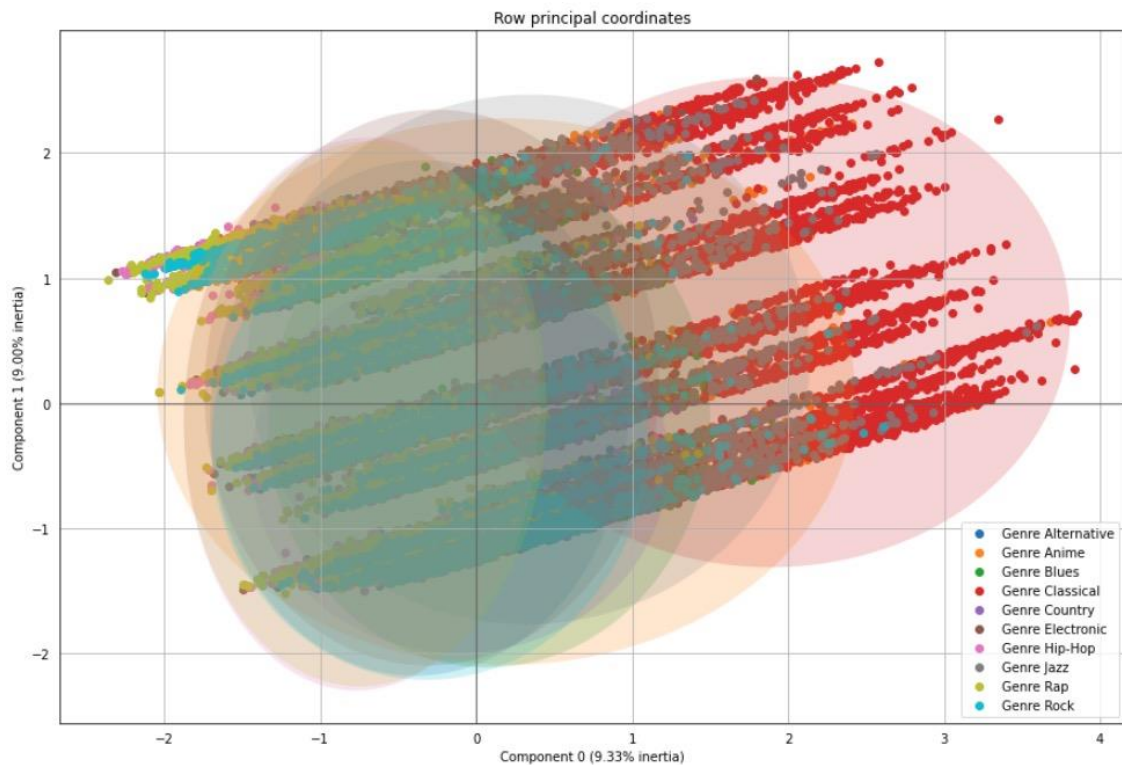
PCA on the training set and looked for eigenvalues that are greater than 1, which are factors that contain more amount of information than a single variable. I found that there are 14 eigenvalues greater than 1, so I kept the first 14 columns of PCs. A 3D visualization of the genres as clusters is given above. We can see that for each genre of music, there are also clear distinctions, as there are 2 or 4 (green and brown) sticks of the same color. This shows diversity across genres. Among jazz and country music, they each have 4 clusters for each genre.

Next, I did a clustering step and included the labels in the dataset as a new feature. It would also boost our classification performance potentially. Elbow method and Silhouette method both showed that the optimal number of clusters is 12. Thus, I included this feature for every row of data in both training and testing sets.

For the multi-class classification problem, the neural network would be able to identify whether the classification was correct or not. Since models with activation function yields better performance than those without generally, I considered cases with activation functions. I built feedforward neural network of one-hidden and two-hidden layers (with ReLU, Tanh, Sigmoid respectively) and assessed their accuracy. The input size is 15 features, and the output size is 10,







FAMD Dimensional Reduction 2D Plot

Final AUC: 0.866

Clustering technique is referenced from <https://towardsdatascience.com/cluster-then-predict-for-classification-tasks-142fdcdc87d6>.

ROC curves plotting is referenced from <https://towardsdatascience.com/multiclass-classification-evaluation-with-roc-curves-and-roc-auc-294fd4617e3a>.