

# Artificial Intelligence 21F Final Project

Tina Buzanis<sup>1</sup>

December 18, 2021

## About the dataset

This dataset contains 13 audio features, along with the genre, for each of 11,741 songs. These songs were collected using Spotify's API, and organized in an sqlite database, which is included with this file. There was, however, a unique challenge in that Spotify does not tag their songs with a genre. Rather than manually assigning genres to each of the 11,741 songs, I instead chose the songs based on genre: for each of the genres I chose to include (chill, indie/alt, edm/dance, jazz, rock, metal, pop, and hip-hop), I downloaded the top playlists available within that genre on Spotify. This gave me a roughly equal distribution of samples from each genre. From there, I was able to proceed to the exploratory data analysis.

## I. EDA

Each of the 11,741 songs in this dataset belong to one or more genres, though the distribution is roughly equal:

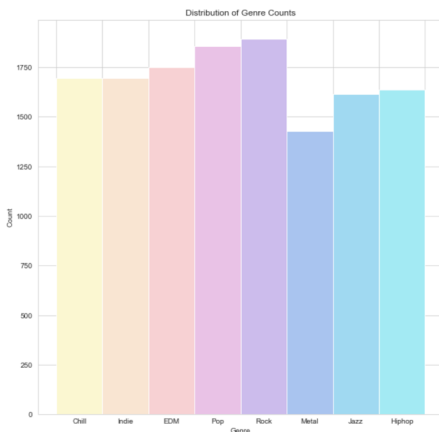


Figure 1: Genre Distribution

However, despite the distribution being roughly equal, there was a unique challenge in this data set in that there is not much correlation between the features:

Through continued trial, error, and inspection of

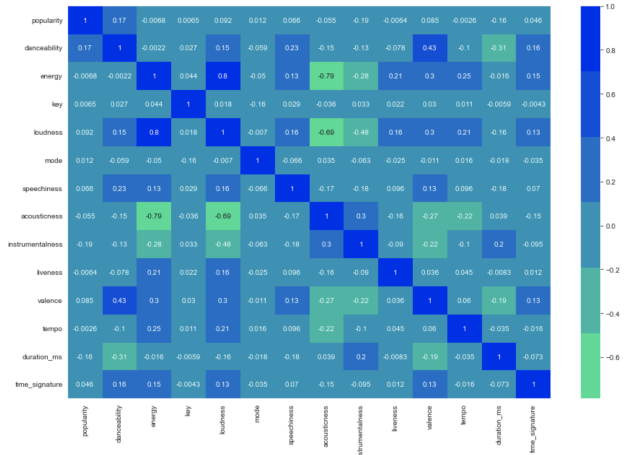


Figure 2: Heatmap of Genre Features

the features' relationships, I narrowed the list of features to be used down to seven: popularity, danceability, energy, key, valence, tempo, acousticness, and instrumentalness. Spotify delivers the audio data pre-scaled. That is, the features directly relating to the audio are of a continuous nature, scaled already to be between zero and one. Features that do not directly correspond to the song audio, such as popularity, key, and tempo, are not scaled. To fix that, I applied both a StandardScaler and a MinMaxScaler to all features that were not standardized. This brought all of the data features onto the same scale, and around a similar distribution.

## II. Naive Bayes - Gaussian

The first model I used was a Gaussian Naive Bayes model. As one may infer from the name, Gaussian Naive Bayes models perform most accurately on data which follows a Gaussian distribution. It was for this reason that I applied the scalers during the EDA. The confusion matrix and F1 score for both the training and test sets is shown below.

Surprisingly, running the model on unscaled data did not provide much of a difference - about one percent improvement for both the training and test cases. This may be because even scaling will not make cor-

```
=====TRAINING PREDS=====
[[676  89  46 153 103  11  62 108]
 [ 90 566  58  44  25 105  28 142]
 [ 39  92 431  24   9  14  65  35]
 [253 135  72 160  21  49  62 295]
 [390 103  65 117 216  19   3  71]
 [ 29  22   6  10   5 713   0 170]
 [155 142 136  60   2  20 145 221]
 [ 72  62  14  61  12 261  31 471]]
0.41136398485452785
```

Figure 3: Training

```
=====TRAINING PREDS=====
[[598  71  36 138 230   3 106  49]
 [ 70 557  54  44  43  77  73 106]
 [ 36  44 521  35   7   8  92  31]
 [237 127  85 184  50  62 139 169]
 [229 104  46  89 432  10  31  21]
 [ 31  47   0  14  18 681   2 143]
 [130 137 188  81  14  12 165 132]
 [ 72  88  12  91   4 254  85 348]]
0.43841268744943396
```

Figure 5: Training

```
=====TEST PREDS=====
[[311  49  23  48  73   7  31  55]
 [ 34 261  27  23  11  48  12  61]
 [ 13  62 251  19   1   8  31  20]
 [123  98  45  68  17  44  16 135]
 [191  68  26  50 109   5   3  29]
 [ 11  21   3   4  12 352   0  78]
 [ 55  78  90  27   7   9  53  94]
 [ 33  46  11  28  11 132  10 204]]
0.3915903376789576
```

Figure 4: Test

```
=====TEST PREDS=====
[[287  30  14  56 129   4  53  32]
 [ 37 276  28  25  29  40  34  38]
 [ 21  20 224  11   0   3  33  15]
 [113  59  41  85  23  32  78  86]
 [121  50  23  42 226   7  23   5]
 [ 13  17   2  16   7 360   0  65]
 [ 73  59  84  41   8   7  98  51]
 [ 52  40  11  51   6 128  24 173]]
0.4424405968760996
```

Figure 6: Test

relations appear where there are none.

### III Naive Bayes - Categorical

The second model tried used a Categorical Naive Bayes algorithm. This classifier works best on data sets that are discrete, or categorical, in nature. This data set has a mixture of categorical and continuous features. As was done for the Gaussian classifier, all features were made categorical prior to feeding the data to the model. This model performed similarly to the Multiclass Logistic Regression model - slightly better than the Naive Bayes. This is not surprising, considering that a good part of the data from this dataset is categorical in nature. Given more time, I would go and carefully formulate how I would discretize the data, instead of doing it all at once. I would also like to explore the different effects different ways of filtering affect the results. The confusion matrices for the training and test predictions are below:

### IV Multiclass Logistic Regression

Next, a multiclass logistic regression model was tried. As Regression models favor continuous data, I once again scaled and centered the data prior to feeding it to this model. This algorithm performed about as well as the Gaussian Naive Bayes, with about a 43% accuracy.

### V. Neural Network

The last algorithm evaluated on this data set was a Neural Network. Having implemented the algorithm from scratch, I did not have access to the same metrics as I would through sklearn. However, after implementing a model with six hidden layers, and running it over 1000 epochs with a learning rate of 0.005, the results appeared to be the best out of the models so far:

Although the model had a difficult time classifying

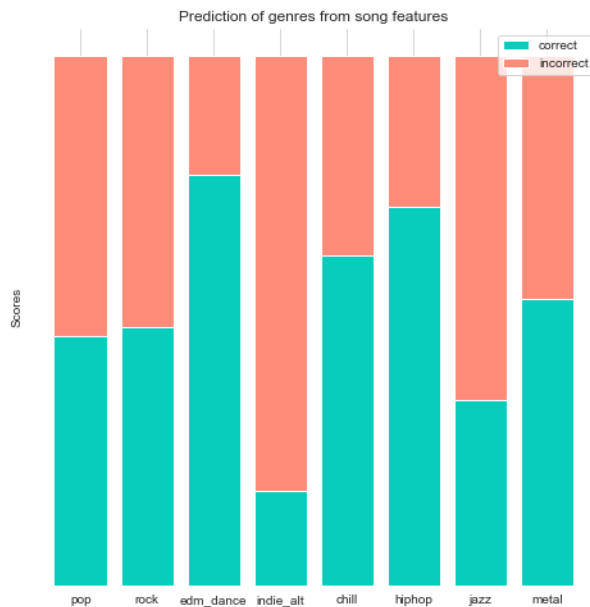


Figure 7: Neural Network

indie-alternative songs, it did well with edm-dance, chill, and hiphop. A deeper exploration of the data is necessary, as is the time to tune the hyperparameters of this model. Several iterations were more "balanced" in that the number of correct predictions were lower, but more evenly distributed across the genres. However, in seeing that this model is quite successful at classifying several genres, it is worth a deeper look into the data to see the trends the algorithm may be capitalizing on.

## VI. Conclusion

Although the performance of these models is not perfect, their performance is acceptable given the size and nature of the dataset. One difficulty that was present was that running the algorithms on the full dataset took a rather long time - about thirty minutes each run. That made the process of iterative improvement time-consuming, as I had to wait a while after making a change to see the results. From the data collected, it appears that the Neural Network performs the best, although it likely presents with a similar accuracy to the Naive Bayes and Logistic Regression models on average. Given more time, I would spend more time exploring the data, and engineering new features. I would also explore different genres, or perhaps a different number of genres. Additionally, I would compare the results from the algorithms I wrote to sklearn's implementation, and see which version per-

forms better on this dataset.