# FINAL REPORT

Music Genre Classification

## *Group 96*

| | |
|---|---|
| So Hee Yoon | 1004327604 |
| Ho Seok (David) Lee | 1004112177 |
| Myeong Hun (David) Song | 1004815961 |
| Ting Wei (Sherman) Lin | 1004835413 |

Word count: 2251

# 1. Introduction

Classifying music by ear is a fairly trivial human task; you can show a child one of Chopin's études and they will be able to recognize it as a classical piece. Formulating explicit musical criteria that defines a classical piece, however, is very difficult and requires detailed descriptions of musical properties such as instrumentation, tempo, and arrangement. Even with these descriptions, identification of a genre without listening to the music can become ambiguous. Despite these shortcomings, we observed that these explicit musical criteria reflect a pattern-like consistency within genres that machines can possibly learn to detect.

An image recognition algorithm can detect a dog by recognizing features that are distinctive of a canine. These machine-learned features, however, cannot easily be translatable for humans to understand. In the same manner, algorithms can learn some mechanized version of those explicit musical criteria.

An algorithm that can classify music is extremely helpful for music platforms such as Spotify or Apple Music. Users are more inclined to search for and explore through genres instead of groupings like "songs with guitars", since these are classifications that can not only be easily understood by humans, but are often reflective of one's music preference. Considering all of these factors, we believe that a machine that can recognize genres is both viable and valuable.
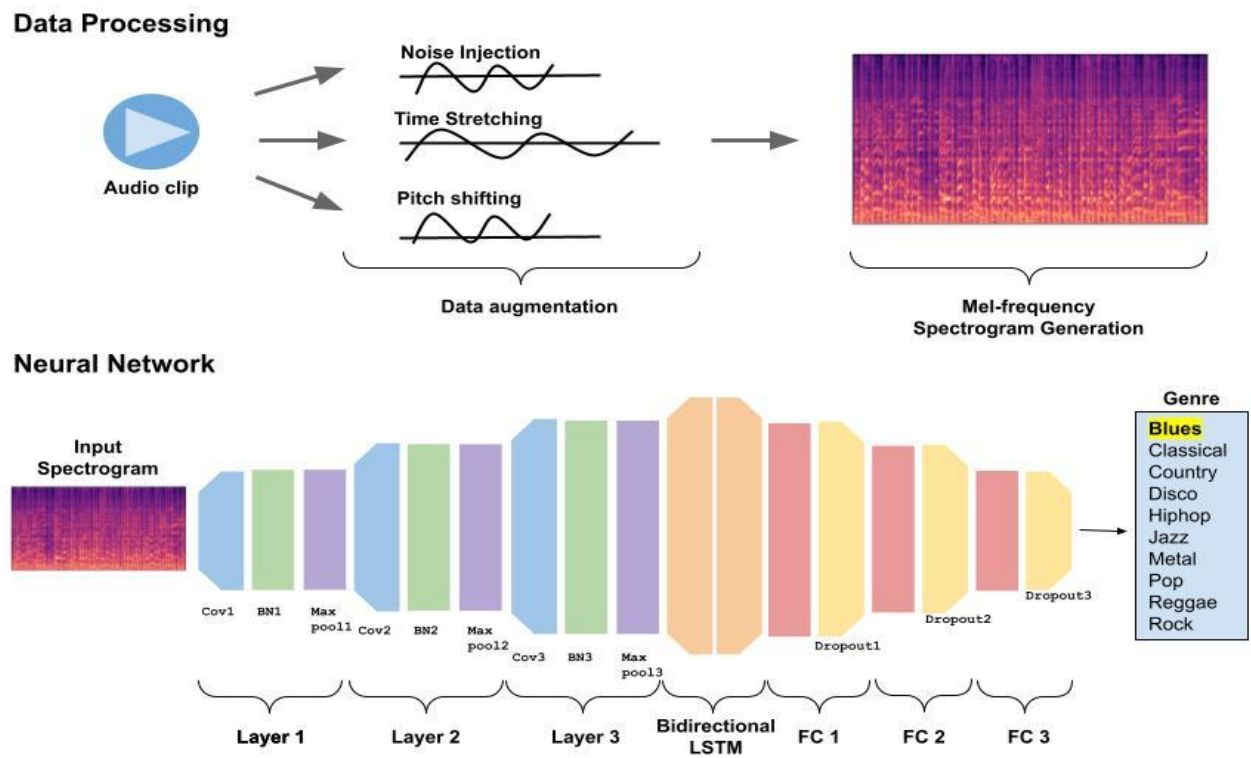
# 2. Illustration / Figure



Figure 1: Illustration showing the core idea and architecture of the project

# 3. Background & Related Work

## Analyzing lyrics and chords to classify music genre *by Timothy Greer and Shrikanth Narayanan [1]*

This paper uses natural language processing to classify music from the Billboard charts and represent the relationship between chord progression and lyric sequence in predicting genres.

## Music genre classification *by Derek A. Huang [2]*

Convolutional Neural Networks are used to predict genres of the music, out of 10 distinct music genres. The model inputs two different data: raw amplitude data and their transformed mel-spectrograms to determine which input resulted in a higher accuracy.

# 4. Data Processing

## *Source of data*

We decided to use the GTZAN genre dataset [3], which feature clips of music (of file type .wav) of 30 seconds in length of categorized genres: 10 genres with 100 audio clips in each. Data augmentation tripled the dataset by generating three new clips from each through noise injection, pitch shift, and time stretching.

## *Mel-Spectrogram*

1. Load audio file using (librosa.load)
2. Create mono audio signal (librosa.to_mono)
3. Generate mel spectrogram (librosa.feature.melspectrogram)
4. Convert power spectrogram into decibel units (librosa.power_to_db)
5. Save figures as image file in the genre folder
6. Resize the 1200*800 pixel image to 240*160 pixel



*Figure 2: Example of generated spectrogram from audio file in blues genre.*



*Figure 3: Spectrograms generated from an audio clip through data augmentation.*

# Mel-Frequency Cepstrum Coefficient

1) Load audio file using (librosa.load)
2) Extract mfcc features (librosa.feature.mfcc)
3) Save mfcc values into data.csv file



*Figure 4: Example of generated MFCC(mel-frequency cepstrum coefficient) logged in excel sheet.*

# Statistics



*Figure 5: Cleaned data distribution for each genre.*

Total Number of dataset
● 400 audio clips in each genre
● Total of 4000 audio clips

Dataset split ratio
● 70% training data (2800 audio clips)
● 15% validation data (600 audio clips)
● 15% test data (600 audio clips)

# 5. Architecture

Between the option of training with mel-spectrogram images or MFCC numbers, we believe that mel-spectrograms offer more possible information for the machine to learn from. Building upon this, we decided to use a convolutional neural network with a fully-connected classification neural network with ReLU activation functions.

## Convolutional Neural Network

There are three increasing convolutional neural networks layers channel sizes of 32, 64, and 128 respectively, each with a kernel size of 3 and stride and padding of 1. After each convolutional layer, we decided to perform batch normalization to improve training efficiency and maximum pooling layers to help consolidate learned information.

## Fully-Connected Neural Network

The classifier is a simple fully-connected neural network with hidden units of 128, 64, and 10 respectively, with dropout layers of 50% to avoid overfitting to the training data. The last output layer represents the classifications for each of the genres, and must be applied cross entropy loss.

## Long Short-Term Memory

After noticing that mel-spectrograms not only contained information about the frequencies but also the order in which they appear, or the arrangement of the song, we decided to incorporate a bidirectional long short-term memory module between the CNN and the classification ANN. The feature maps extracted through the CNN are fed in sequence to the LSTM, which are converted into the embedding for the classifier. We decided to use 128 units in the hidden layer, as making the model too complex resulted in overfitting or slower training. This allows the overall model to learn as much about the mel-spectrograms as possible.

# 6. Baseline Model

To compare the effectiveness of our model, we will compare the accuracy with a few rudimentary supervised baseline models. The two baseline models that we will be using are the $k$-nearest neighbours algorithm ($k$-NN) and a support-vector machine (SVM). Both of these models will be implemented in the scikit-learn Python library [4][5].

## k-Nearest Neighbours Algorithm

The $k$-nearest neighbours algorithm is a classification algorithm that processes labeled data [6]. Given an unlabeled data point, the algorithm finds the closest labeled data point within the parameter vector space.

Each mel-spectrogram from the labeled training data will be converted into 20 mel-frequency cepstral coefficients (MFCC) to populate the vector space. For each unobserved datum, using the Minkowski distance

$$dist(x,\ z)\ =\ \left(\sum_{r=1}^{d} \left|x_r - z_r\right|^p\right)^{1/p},$$

the algorithm will find 5 nearest neighbours (i.e. closest sounding clips). The predicted label will be the most commonly occurring label from the 5 neighbours, weighted by their distance.

## Support Vector Machine

A support vector machine attempts to maximize the separation between clusters of labeled vectors [7]. Since we do not know the general shape and behaviour of our data, we can blindly assume that the classes are not linearly separable. Changing the SVM kernel from a linear function to the Gaussian radial basis function (RBF)

$$rbf(x,\ z)\ =\ exp\left(-\ \gamma||x - z||^2\right),$$

allows the vector space to be mapped to an infinite dimensional feature space for the SVM to create non-linear decision boundaries. Once trained, the SVM can begin predicting the labels of unobserved data.

# 7. Quantitative Results

The KNN algorithm had a testing accuracy of 66.85%, and validation accuracy of 44.67%. The SVM had a testing accuracy of 47.57%, and validation accuracy of 42.67%, which are lower than K-nearest Neighbours algorithm. We have decided to use the KNN algorithm performance as the baseline model.



```
Final Training Accuracy: 0.9027181688125894
Final Validation Accuracy: 0.54
Total time elapsed: 1060.66 seconds
```

*Figure 6: Training and validation accuracy curves with 1000 samples, 50 epochs, learning rate $1 \times 10^{-4}$.*

The first trial run had 1000 samples in dataset, and the model had 50 epochs with Adam optimizer which had learning rate of $1 \times 10^{-4}$. We obtained training accuracy of 90%, but the validation accuracy was only 54%, indicating overfitting.



```
Final Training Accuracy: 0.9660714285714286
Final Validation Accuracy: 0.785
Total time elapsed: 3760.61 seconds
```

To help generalize, we decided to increase the dataset to 4000 samples through data augmentation. The hyperparameters of the model were not changed for the second trial, and we obtained training and validation accuracy of 96.6% and 78.5%. However, the training time was longer than an hour, and we realized that accuracies are not improving significantly after epoch 30.



```
Final Training Accuracy: 0.9475
Final Validation Accuracy: 0.7883333333333333
Total time elapsed: 1868.63 seconds
Test Accuracy:  0.7783333333333333
```

*Figure 8: Training and validation accuracy curves with 4000 samples, 30 epochs, learning rate* $1 \times 10^{-4}$.

In the third trial, we decided to use 30 epochs, which took 31 minutes for training and validation. The training and validation accuracy was 94.75% and 78.83%, which was not too different from the result we obtained in the second trial. This final model produced a test accuracy of 77.83%.

# 8. Qualitative Results

In order to check if its prediction is correct, the model has compared the actual label of the song and predictions, and increased the number of correct predictions when the actual label corresponds to predictions.



*Figure 9: Sample prediction with correct classification.*

This image shows sample prediction with correct classification. It is labelled as classical, and prediction of our model was aligned with the label. To compare it with another mel-spectrogram from classical, they look similar to each other. Both of them are having some black area at the above, followed by a long range of purple and orange, and a short range of orange below. They also showed similar patterns.



*Figure 10: Sample prediction with wrong classification.*

This image shows sample prediction with wrong classifications. It was labelled reggae, but the model predicted it as disco. Comparing it with examples from reggae and disco, all three were showing similar patterns, but reggae does not have black region at the top, while disco does. The image had a black region on top, so it appears like the music clip had some disco traits which could confuse our model.

*Figure 11: Confusion Matrix.*

According to Figure 11, our model has predicted classical and metal clips accurately, while it showed lower accuracy on hip-hop and rock. The most inaccurate prediction was predicting hip-hop as reggae (26%), followed by predicting rock as country (22%), blues as country (11%), and disco as reggae (10%). It confirms that the model could be confused by music clips which show traits of different genres.

# 9. Evaluate Model on New Data



*Figure 12. Confusion matrix from new dataset*

To test our model on a completely new dataset, we created a new dataset consisting of 3 songs each for each genre, which is 30 songs in total. In the new dataset, we obtained an accuracy of 61.725%. Even though it is slightly lower than test accuracy in [Figure 11], it is still considered better than the baseline test accuracy of 50.67% provided by K-nearest Neighbours algorithm. For this new dataset, hip-hop and rock are the two genres that were predicted most accurately, while our model doesn't perform as well when trying to identify country, metal and reggae genres.

The accuracy from the new dataset is lower than validation or test accuracy, and it can be a consequence of having "bad" data. GTZAN is the dataset that is popular for machine learning projects involving music audio, because it is processed and labeled properly. We used GTZAN datasets for training our models and hyperparameters tuning. However, for the new dataset we obtained our data from other sources on the internet, which has risks of data not being labelled properly or data that are of low quality. If the music clips are not labelled properly, then there is a chance that the accuracy from the dataset is lower. For example, if a song is labeled as "metal" when in reality it is more similar to being classified as "rock", then this would result in a wrong prediction from our model if our model classifies it as "rock". Since music can have traits of multiple genres, labelling a song into one genre may not represent the song properly. If we chose songs which show traits of multiple genres, then it is likely that the model would not classify the song's genre properly.

# 10. Discussion

As mentioned in the introduction, we believe that genre-detection is not a trivial task. There are large diversities between genres itself, caused by differences in musicians and creative liberties taken by certain composers. Despite the patterns formed by musical properties within genres, these factors often cause ambiguity in mel-spectrograms and cannot be clearly parsed by the machine. For example, in [Figure 10], the mel spectrogram was labelled as reggae, but our model classified it as disco, since it had similar characteristics with disco's mel-spectrogram. This can justify our low testing accuracy, and may also justify the existence of a performance plateau of genre-detection machines in general.

One of the more important things we learned was the importance of having large datasets in complex problems, as it can greatly improve a model's accuracy. We were able to increase validation accuracy by 24.5% after creating a larger dataset using data augmentation. Since data augmentation was done by simply manipulating existing data to create new data, we didn't expect large improvements in performance. The new data, however, introduced enough diversity to allow the model to generalize efficiently.

# 11. Ethical Consideration

Music is a form of artistic and creative expression. There isn't a definite description of how a certain genre of music has to sound like. A song can be very dynamic, composed of elements from many genres of music. By trying to give a definite answer and classifying a song to be only one genre of music, we are restricting its artistic expression. This puts a preconceived notion on a song that may go against the ideas of its creator. For example, our model may classify a song to be Rock music, but its author may claim otherwise.

Most music streaming platforms have playlists of different music genres. If a song is classified as only one genre and is only placed in one of the many genre-specific playlists, it will lose many potential listeners who don't listen to other playlists because of its labeled genre. As a result, an artist may lose profits on music royalties or the reach of a bigger audience. Furthermore, there is a potential for our project to replace jobs in the music industry that involves music classification, like playlist curators on music streaming platforms.

# 12. Project Difficulty

Throughout the report, we have repeatedly noted how the difficulty in identifying genres stems from the ambiguous similarities between genres. For example, the two genres rock and metal are often misclassified by even human ears, a decision that usually comes down to personal opinion.

We attempted to thwart this inherent and unavoidable property by including an LSTM module to learn from musical arrangement, as humans don't actively consider arrangement when determining a genre. With a fairly deep network and the LSTM in place, we only managed to acquire a 77.83% accuracy.

At the start of hyperparameter tuning, while attempting to deepen the network, we noticed that we were not only overfitting to the training dataset but severely slowing down the training speed. To mitigate this, we worked with a smaller network and implemented dropout layers to help generalize the model, and batch normalization layers for training efficiency.

We only managed to achieve a validation accuracy of 50% with this architecture, and decided that the only way to improve generalization was to generate more data. With a limited audio file data set, we used data augmentation techniques of noise injection, pitch shifting and time stretching to increase the amount of data that we have to train our network. As a result, we were able to  drastically increase the validation performance from 50% to around 75%.

# 13. References

[1] T. Greer and S. Narayanan, "Using shared vector representations of words and chords in music for genre classification," SMM19, Workshop on Speech, Music and Mind 2019, 2019.

[2] D. A. Huang, S. A. Arianna, and P. J. Eli, "Music Genre Classification," 2018.

[3] Index of /sound/genres. [Online]. Available: http://opihi.cs.uvic.ca/sound/genres/. [Accessed: 12-Feb-2021].

[4] "1.4. Support Vector Machines¶," scikit. [Online]. Available: https://scikit-learn.org/stable/modules/svm.html. [Accessed: 12-Feb-2021].

[5] "sklearn.neighbors.NearestNeighbors¶," scikit. [Online]. Available: https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.NearestNeighbors.html#sklearn.neighbors.NearestNeighbors. [Accessed: 12-Feb-2021].

[6] "Lecture 2: k-nearest neighbors / Curse of Dimensionality," Cornell Computer Science - CS4780. [Online]. Available: https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote02_kNN.html. [Accessed: 12-Feb-2021].

[7] R. Berwick, "An Idiot's guide to Support vector machines (SVMs)," in 6.034 Artificial Intelligence - MIT.