

Music Genre Recognition from Audio Tracks and Metadata

Can Kocagil, Berkay Erkan, Efe Eren Ceyani, Mehmet Çalışkan, Hammad Khan Musakhel

Bilkent University

Department of EEE & CS

{berkay.erkay, can.kocagil, mehmet.caliskan, eren.ceyani, hammad.musakhel}@ug.bilkent.edu.tr

Abstract

Automatic recognition of audio tracks from music waves has been a challenging task in Music Information Retrieval (MIR). The mathematical representation of natural music, composed of rhythm, tones, intervals, patterns, and harmonies, is a complex subject for human specialists as there is a geometric interpretation in the humming of the strings and inherent subjective nature. The discrimination of genres from purely audio tracks is stochastic by nature and deeply controversial as the genres share common statistical knowledge. For instance, the musical composition emerges from the intersection of multi-cultural genres, known as fusion genres, e.g., country rock is a fusion of country music and rock music. Contextually, the difficulty level of autonomous discrimination of genres is changing over time. Music trends are changing, making the music genre classification task temporally dynamic and introducing a higher degree of conceptual complexity. For the scope of the project, we designed a computational system of end-to-end learning mechanisms based on classical machine, ensemble and deep learning algorithms for automatic recognition of music genres from Creative Commons licensed audio files, with hand-crafted track-level features.

Keywords- Audio Processing, Music Information Retrieval, Music Genre Recognition, Multi-class Time Series Classification

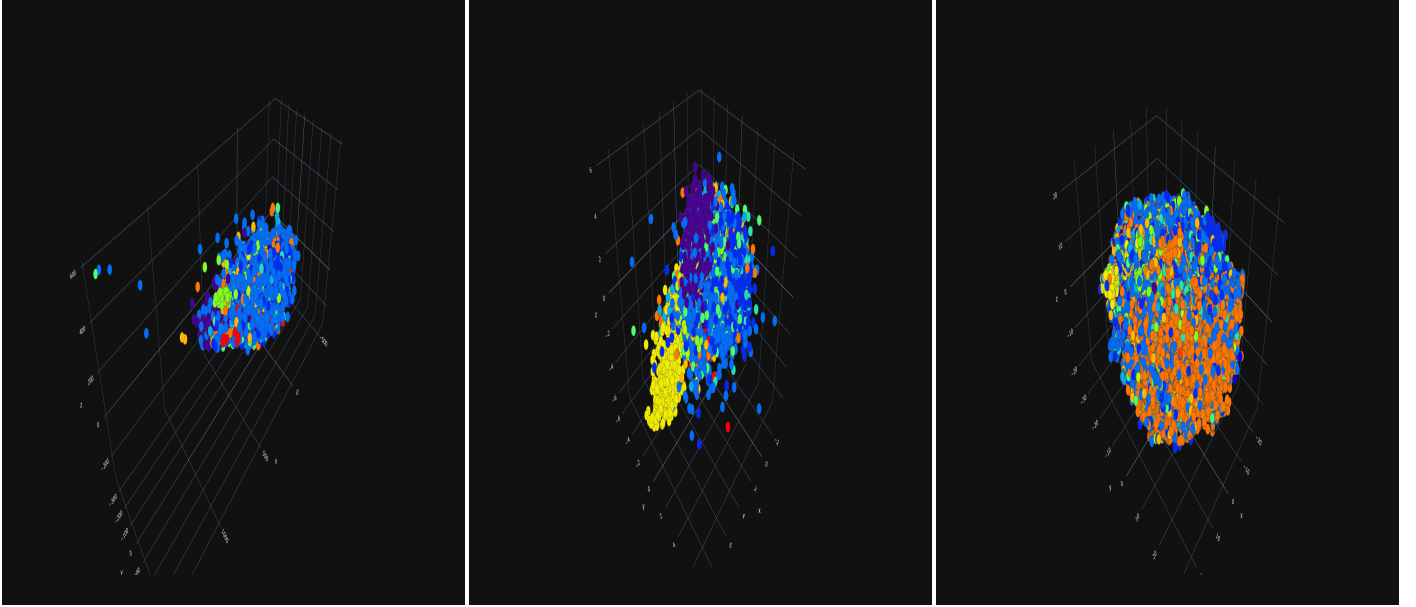


Figure 1: Left to right: Three dimensional Principal Component Analysis, Linear Discriminant Analysis, t-distributed Stochastic Neighbor Embedding (t-SNE) on metadata

1. Introduction

As the dynamics of common genres are temporally varied, representable and scalable datasets are required to perform cognitive tasks in Music Information Retrieval (MIR). The Free Music Archive (FMA) [1] is a dataset for everyday MIR analysis tasks and is supervised for browsing, searching, and organizing extensive music collections. The FMA dataset is a MIR benchmark dataset that consists of 106,574 tracks from 16,341 artists and 14,854 albums, arranged in a hierarchical taxonomy of 161 genres [1]. Its large composition is structured by 30 seconds of high-quality audio, pre-computed features, together with track- and user-level metadata, tags, and free-form text such as biographies, for scaling other music domain processing tasks [1]. FMA is a metadata-rich, future-proof, and permissively licensed dataset that offers quality audio, pre-computed music-level features, reproducibility, and long-term sustainability. Considering the computational comfort zone, the FMA dataset provides four versions, divided according to the sizes: small (30 seconds long, 8,000 tracks, eight balanced genres), medium (30 seconds long, 25,000 tracks, 16 unbalanced genres), large (30 seconds long, 106,574 tracks, 161 unbalanced genres) and full (106,574 untrimmed tracks, 161 unbalanced genres). In light of our computational resources being restricted to open-source computing, we will utilize the small or medium versions of the datasets. However, in the case of any computing device availability, we will try to process a full FMA dataset of 879 GiB of raw information.

We implemented trend-aware classifier-based learning algorithms, specialized in capturing the complex interplay of cultures to recognize genres autonomously. Our classes is music genres supervised and annotated by the artists, from contemporary to jazz. By acknowledging the socio-cultural context, we designed cognitive machine learning algorithms to capture the temporal dynamics of audio, from conventional ML algorithms like SVM to more advanced meta-learning algorithms. Various ML-based techniques are utilized to accomplish the project's scope; to start with, naive and straightforward algorithms such as Decision Trees and SVM are introduced to create our baseline model to benchmark. Then more advanced ensemble classification algorithms such as Random Forests and custom ones are proposed to compare our multi-class results. Initially, we will compute a small-sized dataset that is class-balanced; as such, our primary goal is to predict the single top genre on the balanced small subset. Then, in the case of large or full-size datasets, our main objective is to maximize multiple top genres' accuracy. Our task is a multi-class classification of signals and metadata that requires supervised performance evaluation metrics from the computational perspective. In small or medium dataset sizes, single-top class accuracy is the primary metric to compare. Hence, our primary milestone was preparing the data pipeline and creating conventional ML-based classifiers to construct a baseline.

Musical genre recognition classifies the musical collections based on the audio-level, user-level, and track-level features. There are expected challenges to overcome, such as poor labeling and non-triviality extraction of features. In the former one, musical genres are loosely defined, as people often argue over the genre of a song. In our case, the genres are supervised by the artists themselves, so it is smoother to perform data-centric computational tasks. In the latter one, automatic extraction of semantic features is not differentiable and conceptualized, so there is no standard way to extract features in music information retrieval. The Mel Spectrogram-based representation algorithms are proposed to overcome these challenges by converting audio signals to visuals and representing how the frequency spectrum varies over time [5]. The latter challenge is tried to be overcome with spectrum-based representations before feeding the ML model. From a more musical perspective, the notations of composers and audios created by their artists consist of musical representations: counting, rhythm, scales, intervals, patterns, symbols, harmonies, time signatures, overtones, tone, and pitch [7], which are highly unstructured and complex to interpret mathematically. Hence, ad-hoc statistical and structural pattern recognition frameworks is fit to capture music's conceptual and contextual sides for the extraction of genres.

Our main study for the scope of the project as follows.

- Metadata-level descriptive data analysis is applied.
- Then, manifold learning and dimensionality reduction methods are performed on the metadata-level music data.
- We built an end-to-end discovery machine learning pipelines to classify the music genres based on the metadata.
- Model explainability operations are applied to capture which features have greater importance.
- Audio-level discovery signal analysis methods are applied.
- To convert audios to images, Mel Spectrum computation is performed.
- 2-D Convolutional Neural Networks and Vision Transformer are employed to model music genres by digesting the state-of-the-art.

2. Metadata-level Data Analysis

To gain insight into the FMA dataset, we performed a brief metadata-level data analysis. We analyzed how different audio features vary by genres. Those audio features include acousticness, danceability, speechiness and tempo, which are already given in the dataset. We also analyzed listening counts of each genre. The following figures are self-explained and best understood from the figures.

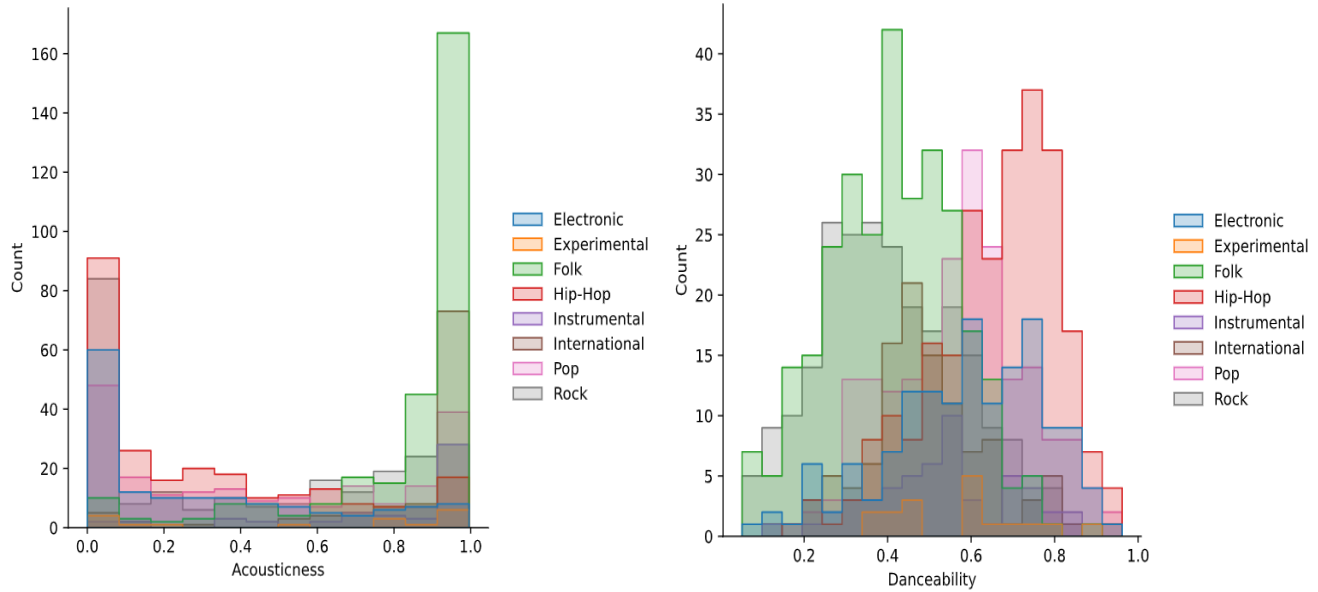


Figure 2: Left to right: Analysis of how acousticness changes for different genres, Analysis of how danceability changes for different genres

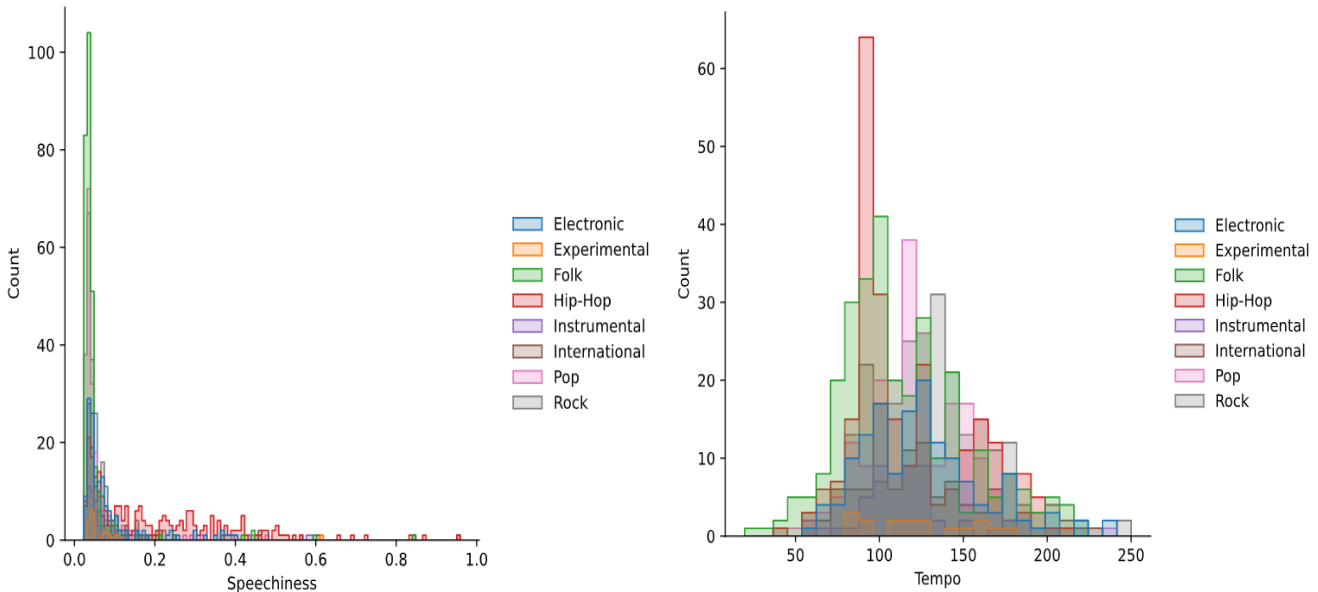


Figure 3: Left to right: Analysis of how speechiness changes for different genres, Analysis of how tempo changes for different genres

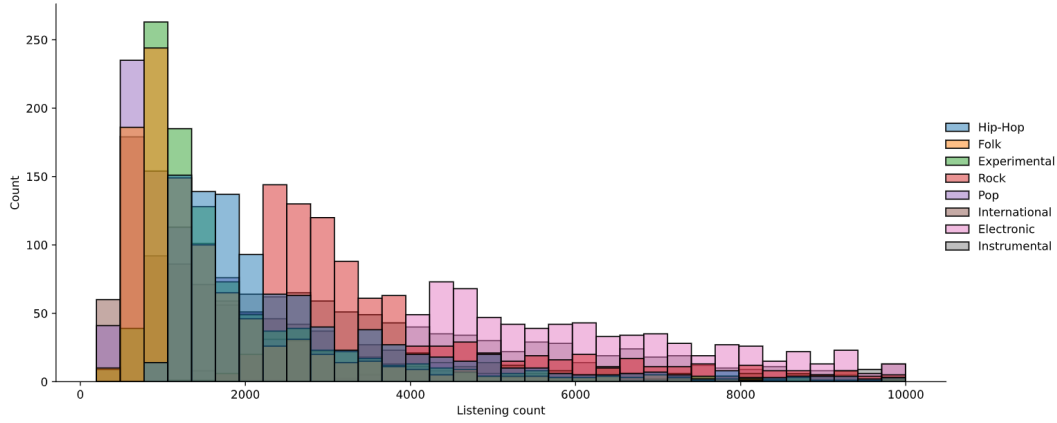


Figure 4: Analysis of how listening count changes for different genres

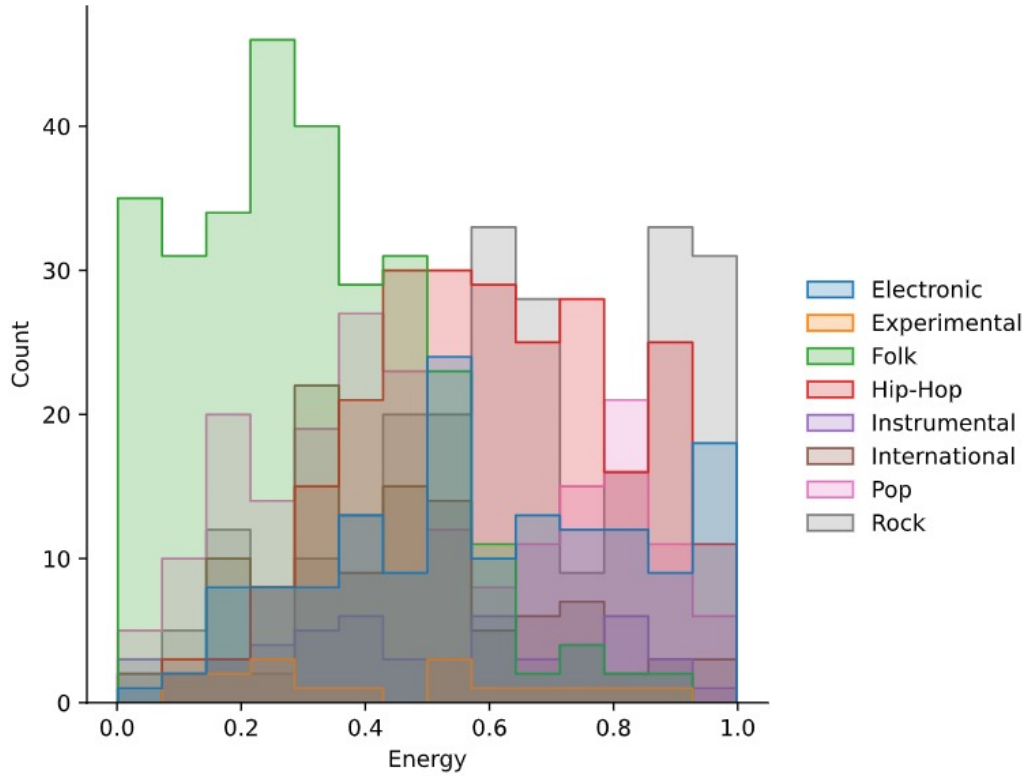


Figure 5: Analysis of how energy count changes for different genres

3. Dimension Reduction to Manifold Learning

We performed PCA and LDA as dimension reduction algorithms. Besides, t-SNE is performed to generate lower-dimensional manifolds of the metadata space. Further, we visualized all manifolds in 3-D dimensional space in figure 1.

3.1. Dimension Reduction: PCA

PCA is a linear unsupervised dimension reduction algorithm, and it computes principal vectors to change the basis of the representation [6]. PCA is a used algorithm in a broad range of topics from image compression to decorrelation of texts. Here, we performed PCA on metadata space and visualized it in figure 1 at the left corner.

3.2. Dimension Reduction: LDA

LDA is a supervised dimensionality reduction algorithm, and it is a generalization of Fisher’s linear discriminant, which aims to find linear subspace that characterizes the original data space. Since it is supervised, it is a powerful paradigm in representation learning. Here, we performed LDA on metadata space and visualized in figure 1 at medium. From the figures, we can see that LDA outperforms other methods by uniquely separating geodesic distances in the manifolds.

3.3. Manifold Learning: t-SNE

T-SNE is iterative statistical approach for producing non-linear embedding of the original data space by preserving small pairwise distances or localized similarities. It minimizes the Kullback-Leibler divergence between the joint probabilities of the low-dimensional embedding and the high-dimensional data. Here, we performed t-SNE on metadata space, and visualized in figure 7 at right corner.

4. Machine and Meta Learning

We applied more than ten genre recognition algorithms; we listed our methods and their brief descriptions.

4.1. SVC with Gaussian and Linear Kernel

Support Vector with the linear kernel is an efficient algorithm for linearly separable data spaces by fitting a hyper-plane that categorizes the metadata in our context [3]. We also implemented SVM with radial basis kernel and accumulated to mean of 60% accuracy with a standard deviation of 0.02 by 5-fold cross-validation.

4.2. SGD Classifier

SGD classifier is a linear classifier that uses stochastic gradient descent with Hinge loss to separate data spaces. The SGD classifier is implemented with l_1 regularization [3].

4.3. MLP

MLP is a simple neural network architecture that consists of a bunch of linear layers with ReLU non-linearity and is optimized by the Stochastic Gradient Descent rule [3, 4]. Multi-layer Perceptron Classifier accumulated a mean of 56% accuracy with a standard deviation of 0.02 by 5-fold cross-validation.

4.4. Logistic Regression

Logistic regression is a classical machine learning algorithm; it applies linear transformation on the data followed by sigmoidal activation to generate real-valued probabilities [3].

Given the input feature $\{X\}_{i=1}^n$, weight vector w , and the bias term c , the positive class probability can be computed as follows.

$$p(X_i) = \frac{1}{1 + e^{-(c + X_i^T w)}} \text{ for } i \in \{1, \dots, n\} \quad (1)$$

As an optimization problem, binary class ℓ_2 penalized logistic regression minimizes the following cost function:

$$\min_{w, c} \frac{1}{2} w^T w + C \sum_{i=1}^n \log(\exp(-y_i(X_i^T w + c)) + 1). \quad (2)$$

Ridge classifier is an L2 penalized version of logistic regression that helps robust decoding with lower degrees of freedom [3].

4.5. Random Forest Classifier

Random Forest classifier is a meta estimator algorithm that fits parallel decision tree classifiers with bootstrapped samples of the original dataset to improve the predictive accuracy and model reliability [3].

4.6. Gradient Boosting Classifier

Gradient boosting classifier is boosting algorithm that ensembles weak learnings, generally decision trees. It constructs an additive model in a forward stage-wise fashion that enables the model to optimize the parameters on any differentiable loss functions. We utilized 100 weak learners to construct the end classifier [3].

4.7. Quadratic Discriminant Classifier

The quadratic Discriminant classifier is a class conditional density algorithm with quadratic decision boundaries to fit separate Gaussian to each class. It is a powerful method when we have priory knowledge that individual classes exhibit distinct covariance [3].

4.8. AdaBoost Classifier

AdaBoost classifier is a type of ensemble learning algorithm that fits weak classifier on the dataset then fits additional copies of the classifier with adjusted parameters based on the incorrectly classified objects [3]. The basic concept behind Adaboost is to set the weights of classifiers and train the data sample in each iteration such that it ensures the accurate predictions of unusual observations [3].

4.9. Extra Trees Classifier

Extra Trees classifier is another ensemble learning method based on randomized decision trees on the sub-sampled version of the datasets to improve predictive quality [3]. Extra Trees Classifier is an ensemble learning algorithm similar to Random Forest Classifier except for a few key differences. When Random Forest Classifier splits a node, it first finds the optimum split among all splits and chooses that optimum split as the cut point. Extra Trees Classifier chooses this cut point randomly; hence it runs faster than Random Forest Classifier, adding more randomization.

4.10. K-Neighbors Classifier

K-Nearest Neighbor is a simple distance-based supervised ML algorithm. K-Nearest Neighbor works by finding the distances between a query and all the examples in the data, selecting the specified number examples, say K, closest to the query, then votes for the most frequent label in the case of classification or averages the labels in the case of regression [3]. As our task is classification, we take the most frequent label in K nearest neighbors.

We utilized three cartesian space distance metrics as a distance metric: Euclidean, Manhattan, and Cosine. Euclidean distance is the most common one, as it computed the second-order distance between two points.

$$d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \cdots + (x_i - y_i)^2 + \cdots + (x_n - y_n)^2}. \quad (3)$$

In the generic form, it is Minkowski distance with $p = 2$. This distance metric is valid when real-valued vector spaces and the following conditions are satisfied.

- Non-negativity: $d(x, y) \geq 0$
- Identity: $d(x, y) = 0$ if and only if $x = y$
- Symmetry: $d(x, y) = d(y, x)$
- Triangle Inequality: $d(x, y) + d(y, z) \geq d(x, z)$

Then, the Minkowski distance, in general, can be computed as follows.

$$d(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}. \quad (4)$$

If we set $p = 1$, then the $d(x, y)$ function become Manhattan distance, that is first-order distance metric. It computes the absolute value of the points, whereas Euclidean distance punishes the large distance points in a quadratic manner. Finally, we utilized the cosine distance, which can be computed by $1 - S_C(A, B)$.

$$\text{cosine similarity} = S_C(A, B) := \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}, \quad (5)$$

4.11. Naive Bayes

Naive Bayes methods are a set of supervised learning algorithms based on applying Bayes' theorem with the "naive" assumption of conditional independence between every pair of features given the value of the class variable [3]. Bayes' theorem states the following relationship, given class variable y and dependent feature vector x_1 through x_n [3]

$$P(y | x_1, \dots, x_n) = \frac{P(y)P(x_1, \dots, x_n | y)}{P(x_1, \dots, x_n)} \quad (6)$$

Using the naive conditional independence assumption that

$$P(x_i | y, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = P(x_i | y) \quad (7)$$

for all, this relationship is simplified to

$$P(y | x_1, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i | y)}{P(x_1, \dots, x_n)} \quad (8)$$

Since $P(x_1, \dots, x_n)$ is constant given the input, we can use the following classification rule [3]:

$$\begin{aligned} P(y | x_1, \dots, x_n) &\propto P(y) \prod_{i=1}^n P(x_i | y) \\ &\Downarrow \\ \hat{y} &= \arg \max_y P(y) \prod_{i=1}^n P(x_i | y) \end{aligned} \quad (9)$$

4.11.1. Bernoulli Naive Bayes Classifier

Bernoulli Naive Bayes classifier applies Bayesian rule to a dataset with the prior assumption that the data comes from multivariate Bernoulli distribution [3]. Bernoulli Naive Bayes algorithm is a variant of Naive Bayes algorithm which assumes attributes are independent and do not affect each other and gives all features equal weight [3]. A most important distinction of Bernoulli Naive Bayes is that it uses binary value features like true/false or 1/0. This algorithm assumes the Bernoulli distribution's prior distribution and uses Bayes' Rule to maximize the posterior distribution.

$$p(x) = P[X = x] = \begin{cases} q = 1 - p & x = 0 \\ p & x = 1 \end{cases} \quad (10)$$

The decision rule for Bernoulli naive Bayes is based on

$$P(x_i | y) = P(i | y)x_i + (1 - P(i | y))(1 - x_i) \quad (11)$$

4.11.2. Gaussian Naive Bayes Classifier

Gaussian Naive Bayes classifier applies Bayesian rule to a dataset with the prior assumption that the data comes from multivariate Gaussian distribution [3].

4.12. Nearest Centroid Classifier

In the Nearest Centroid classifier, each class is represented by with class centroid. New samples are classified based on the distance to the class centroid, so it is very similar to the supervised version of K-means [3].

4.13. Bagging Classifier

Bagging classifier is an ensemble learning method that fits base classifiers to random subsets of the original dataset. Then, the predictions are aggregated either by voting or averaging to produce end class-wise predictions [3].

5. Deep Learning for Mel Spectrum of Audio Waves

We further performed deep learning algorithms from CNNs to vision transformers. All models are optimized with Cross Entropy loss, Adam and SGD optimizer with starting learning rate as $1e-3$. Scheduled learning rate scaler is performed to decrease the learning rate by % 5 at every epoch. Batch size is varied from 16 to 64. One single Tesla K80 GPU is used for all experiments. Architecture details and their visuals of both CNNs and vision transformer are depicted in the presentation and the following sections.

5.1. Audio Level Signal Analysis

We begin with audio-level discovery signal analysis in both time and frequency domain. We first plot arbitrary audio samples from different classes to perceive any varying difference between them. Then, we visualized the Fourier Transform of the audios as follows.

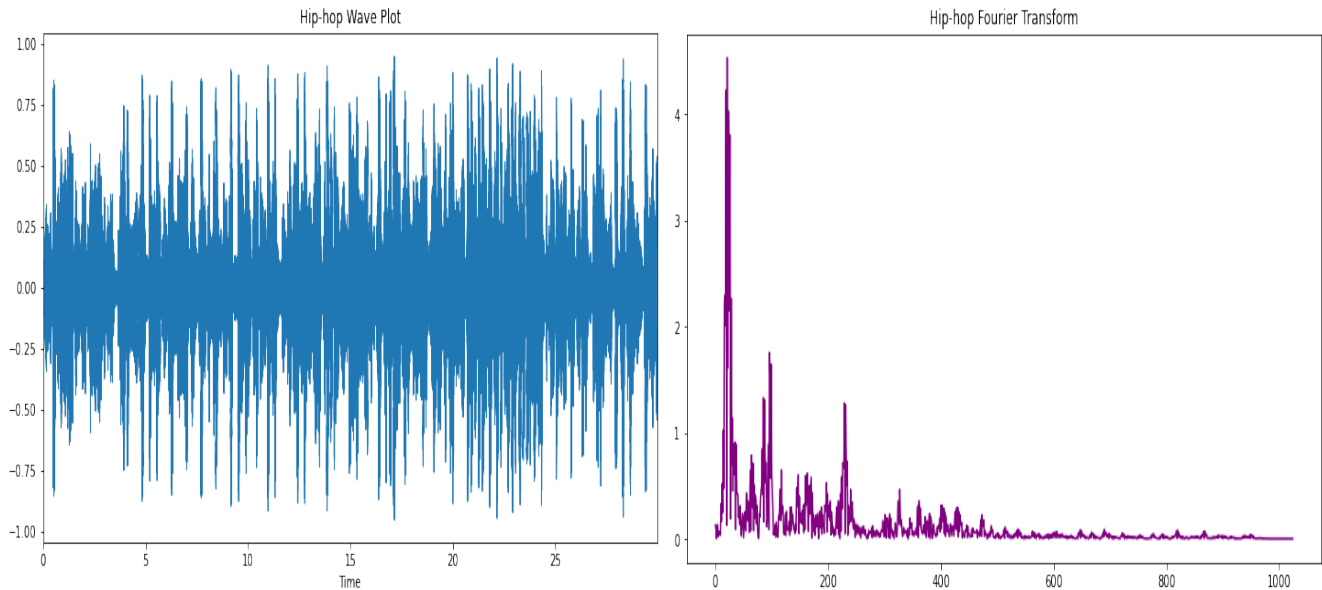


Figure 6: Left to right: Wave and Fourier Transform Visualizations of Hip-hop Audio

5.2. Mel Spectrum Computing

Mel spectrum computing consist of two essential mathematical notions - Mel Scale and Spectrum. Spectrum's are computed via Fourier Transform that is a function that gets a signal in the time domain as input, and outputs its decomposition into frequencies.

On the other hand, in mathematical terms, the Mel Scale is the outcome of a non-linear transformation of the frequency scale. This Mel Scale is set up in such a way that sounds that are equal distance apart on the Mel Scale "sound" the same to humans. Unlike the Hz scale, where the difference between 500 and 1000 Hz is clear, the gap between 7500 and 8000 Hz is barely discernible. Hence, it partitions the Hz scale into bins, and transforms each bin into a corresponding bin in the Mel Scale, using overlapping triangular filters.

Moreover, the Mel scale is a set of tones that human hearing perceives as being equally spaced apart. The interval in hertz between Mel scale values (or simply Mel's) grows as frequency rises. The name Mel comes from the word melody, and it denotes that the scale is based on pitch comparisons. The Mel spectrum converts hertz data to Mel scale values.

The linear audio spectrum is best for applications in which all frequencies are equally important, whereas Mel spectrum's are preferable for applications in which human hearing perception must be modeled. Audio classification programs can also benefit from Mel spectrum data as in our case.

We visualized the Mel Spectrum of the audio signals that are belongs to different classes to perceive any discriminating difference between them. Also, we tried to visualize the pipeline of computing Mel spectrum in figure 8.

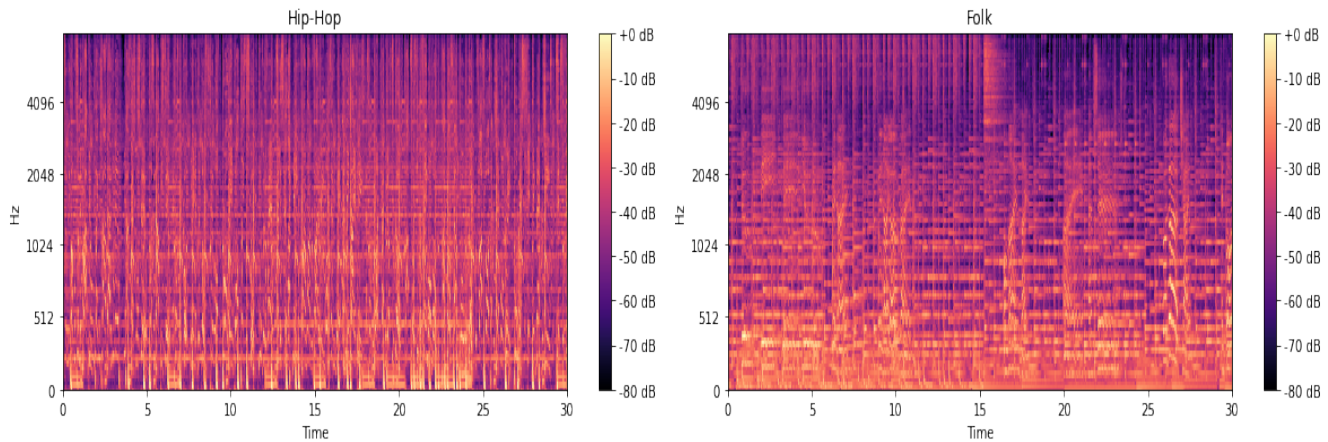


Figure 7: Mel Spectrum of Hip-hop and Folk Audios

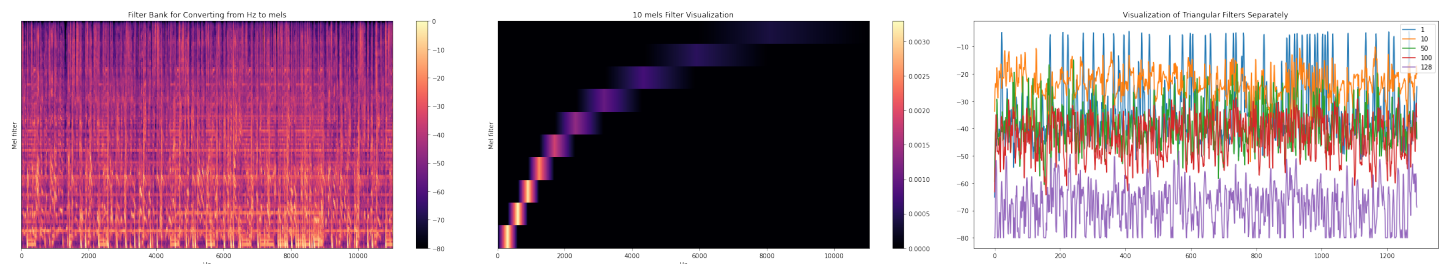


Figure 8: Mel Spectrum Computing Explanation

Finally, we visualized the overall computing of spectrums in the following figure.

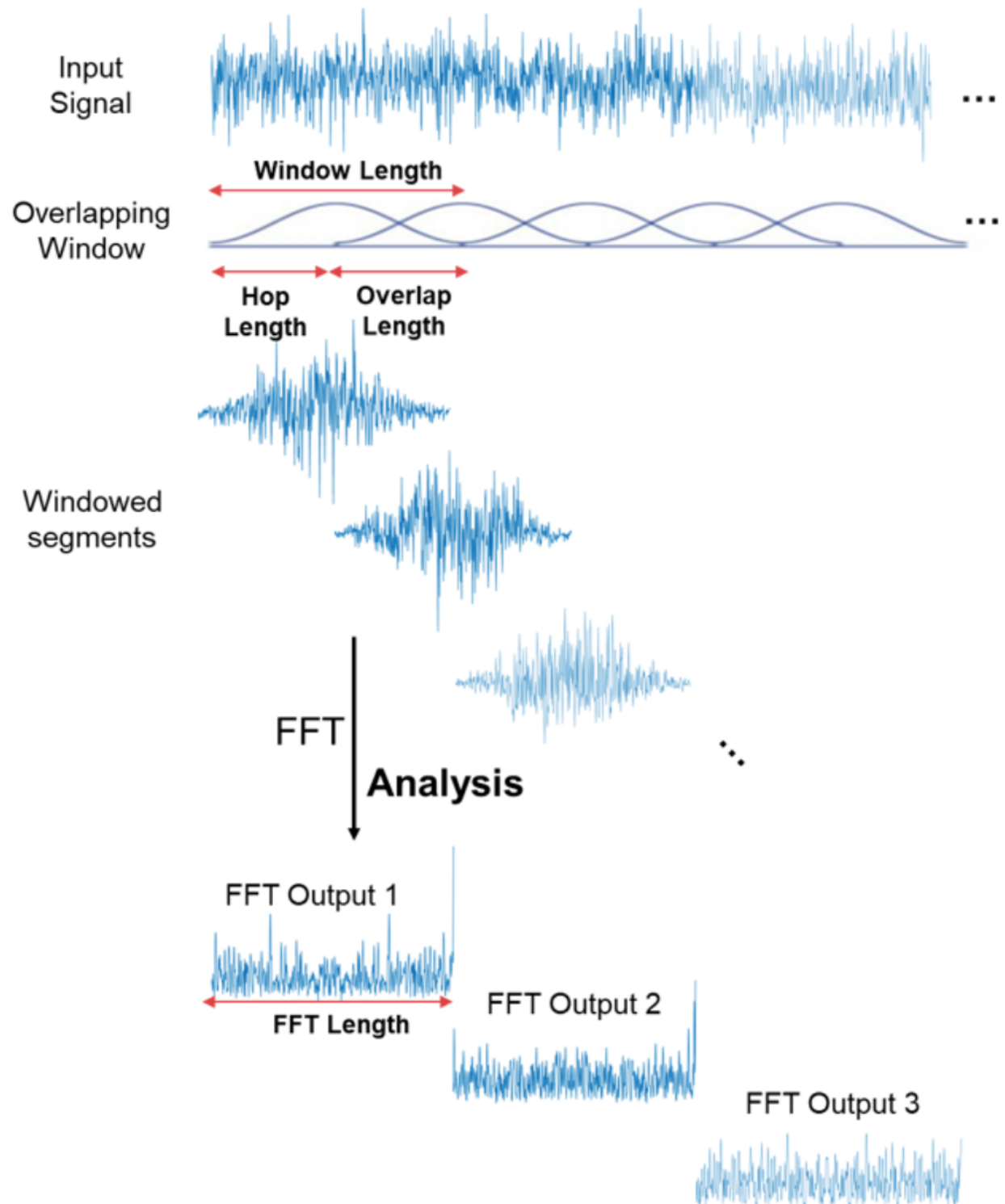


Figure 9: Mel Spectrum Computing Visualization

5.3. 2-D CNN

We developed 2-D CNN architecture to decode spectrum samples. In the architecture, there are 7 special convolutional blocks, each block consists of a sequel of Convolution, ReLU, and Max Pooling layer. To classify, the representative feature vector is propagated to linear blocks. There are two linear blocks, each block consists of a sequel of linear layer, and ReLU layer. We depicted our CNN model in the following figure, note that deeper architecture is utilized but for the visualization purposes we keep it simple. More details are provided in the appendix section.



Figure 10: ConvNet Architecture Visualization

5.4. Vision Transformer

Vision transformers are recently proposed paradigm that replaced with CNN blocks with multi-head self attention mechanism. As vision transformer, we performed experiments in ViT that is state-of-the-art vision transformer with its unique attention mechanism by taking the advantage of interactions between different streams of visual spectrum representations [2].

The Vision Transformer, or ViT, is an image classification model that uses a Transformer-like design to classify image patches. An image is divided into fixed-size patches, which are then linearly embedded, position embeddings added, and the resulting vector sequence given to a conventional Transformer encoder. The traditional strategy of adding an extra learnable "classification token" to the sequence is employed to perform classification [2].

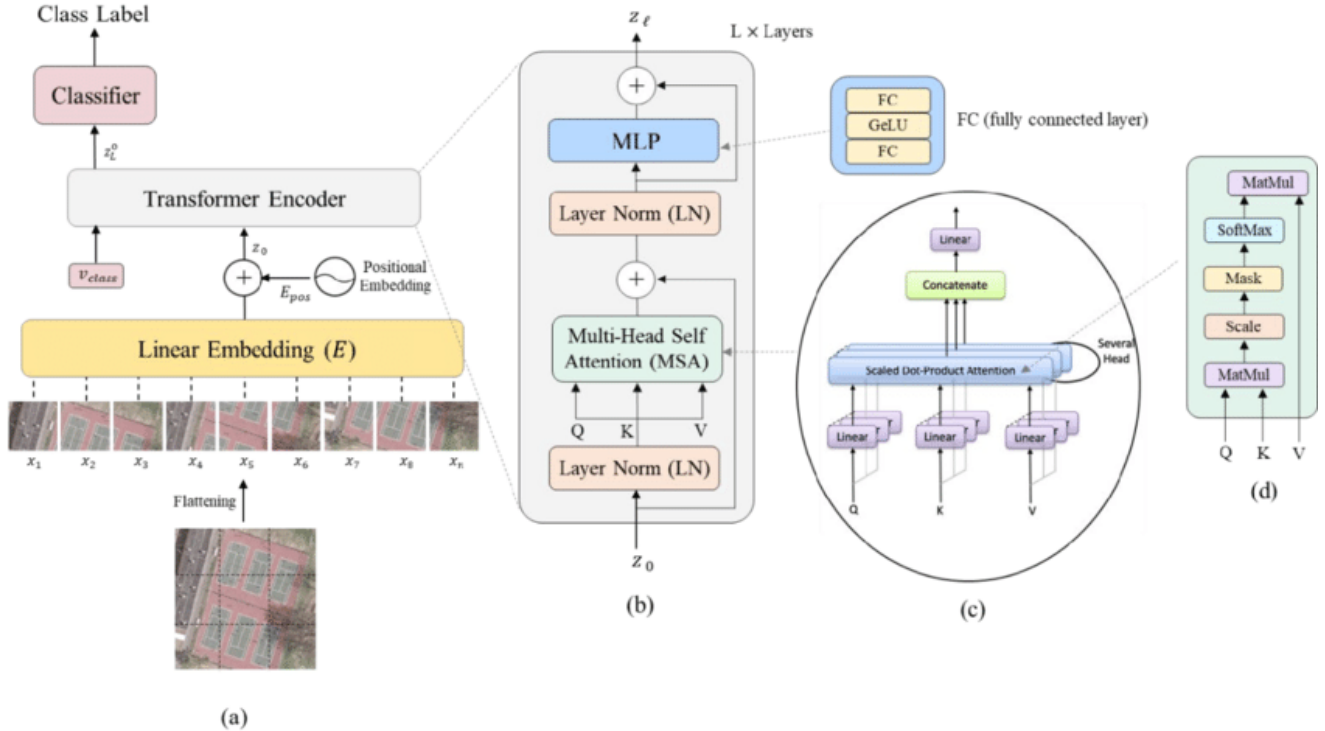


Figure 11: Vision Transformer Architecture [2]

6. Results

In this section, we provide corresponding results. All tables and figures are self-explained.

	Accuracy	Precision	Recall	F1 Score	Fitted Time (s)	Top-3 Accuracy
MLPClassifier	0.58	0.58	0.58	0.58	506.96	88
SVC	0.58	0.58	0.58	0.58	318.92	88
ExtraTreesClassifier	0.56	0.56	0.56	0.56	32.91	85
RandomForestClassifier	0.56	0.56	0.56	0.56	142.81	85
LinearSVC	0.55	0.55	0.55	0.55	253.10	84
LogisticRegression	0.54	0.54	0.54	0.54	9.08	83
KNeighborsClassifier	0.52	0.52	0.52	0.52	0.08	85
CNN	0.51	0.51	0.50	0.51	21200.71	79
BaggingClassifier	0.51	0.51	0.51	0.51	192.71	78
RidgeClassifier	0.51	0.51	0.51	0.51	0.60	77
SGDClassifier	0.49	0.49	0.49	0.49	9.37	79
Vision Transformer	0.41	0.41	0.41	0.41	90000.37	77
DecisionTreeClassifier	0.39	0.39	0.39	0.39	23.61	75
AdaBoostClassifier	0.39	0.39	0.39	0.39	103.01	76
NearestCentroid	0.38	0.38	0.38	0.38	0.13	75
Perceptron	0.37	0.37	0.37	0.37	5.93	72
GaussianNB	0.36	0.36	0.36	0.36	0.17	67

Table 1: Machine Learning Models and Performance Table. Macro averaging is used for precision, recall, and f1-score metrics. Highest metrics are bolded.

	Accuracy	Precision	Recall	F-1 Score	Fitted Time (s)	Top-3 Accuracy
Voting Ensemble of MLP, RFC, and SVC	0.58	0.58	0.58	0.58	368.62	88
Voting Ensemble of MLP and SVC	0.57	0.57	0.57	0.57	370.14	87
Voting Ensemble of k-NN, GaussianNB and SVC	0.53	0.53	0.53	0.53	310.07	81

Table 2: Meta Learning Models and Performance Table. Macro averaging is used for precision, recall, and f1-score metrics. Highest metrics are bolded.

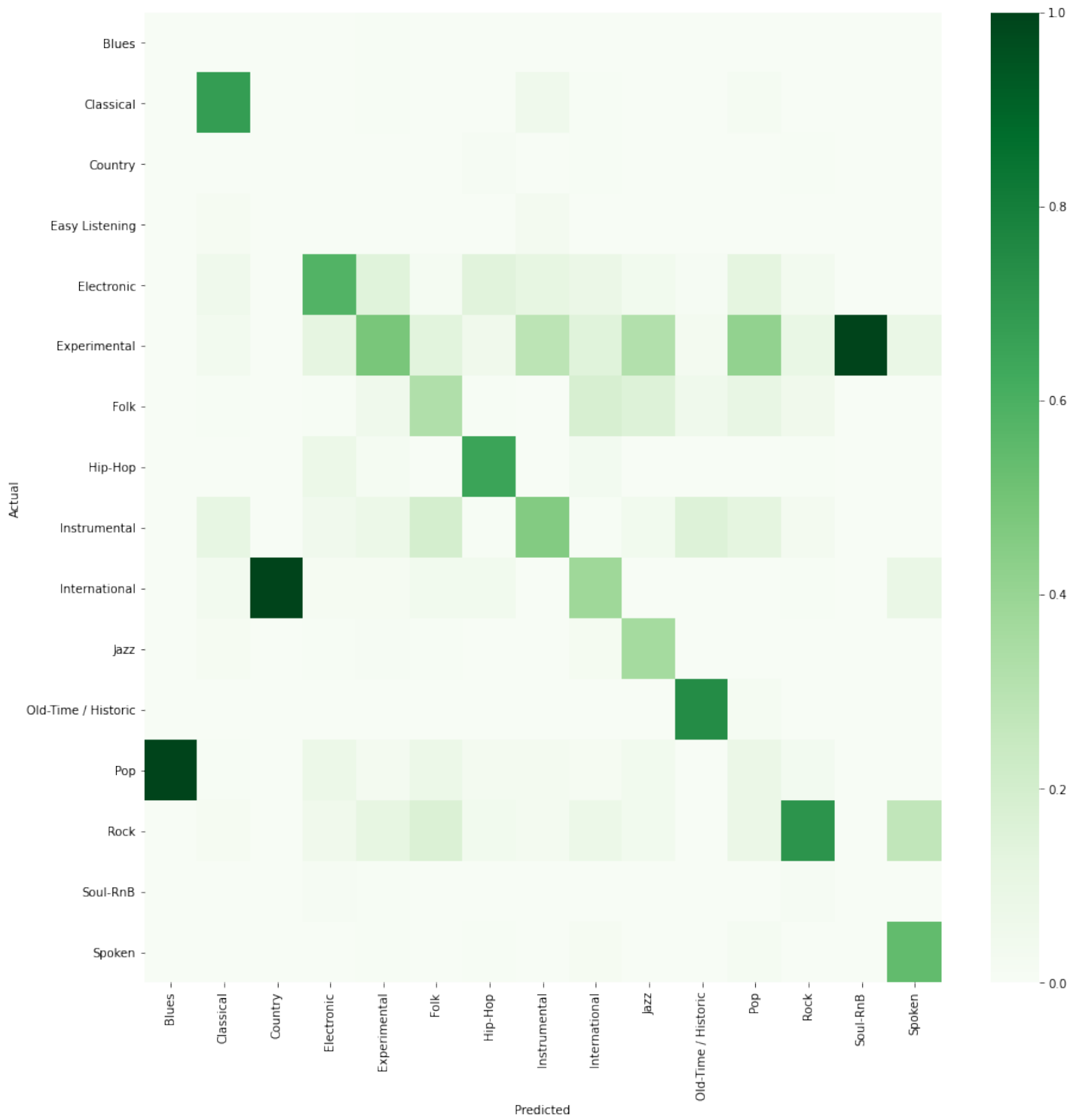


Figure 12: The heat-map of MLP Confusion Matrix

Predicted	Blues	Classical	Country	Electronic	Experimental	Folk	Hip-Hop	Instrumental	Jazz	Old-Time / Historic	Pop	Rock	Soul-RnB	Spoken
Actual														
Blues	0	0	0	0	9	0	0	0	0	0	0	4	0	0
Classical	0	72	0	0	10	1	0	2	1	0	0	1	0	0
Country	0	0	0	1	0	1	5	0	1	0	0	0	10	0
Easy Listening	0	2	0	0	3	0	0	1	0	0	0	0	0	0
Electronic	0	6	0	484	214	6	47	4	11	1	1	6	59	0
Experimental	0	4	0	98	709	37	17	10	18	8	2	20	159	1
Folk	0	1	0	4	84	88	2	0	23	4	4	5	84	0
Hip-Hop	0	0	0	65	30	1	216	0	4	0	0	0	7	0
Instrumental	0	12	0	41	116	53	3	16	1	1	11	6	49	0
International	0	4	1	14	25	11	14	0	47	0	0	0	11	1
Jazz	0	2	0	5	20	3	0	0	3	9	0	0	5	0
Old-Time / Historic	0	0	0	0	2	0	0	0	0	0	52	1	0	0
Pop	1	1	0	64	42	22	10	1	2	1	0	4	56	0
Rock	0	2	0	44	170	46	15	1	10	1	0	4	1168	3
Soul-RnB	0	0	0	10	11	0	1	0	0	0	0	0	21	0
Spoken	0	0	0	0	16	1	2	0	2	0	0	1	3	6

Table 3: MLP Confusion Matrix

7. Model Explainability

The Random Forest model is explained by the mean decrease impurity or Gini importance. The mean and standard deviation of impurity decrease accumulation inside each tree were used to determine Gini significance. The decrease in impurity is measured for features in each split.

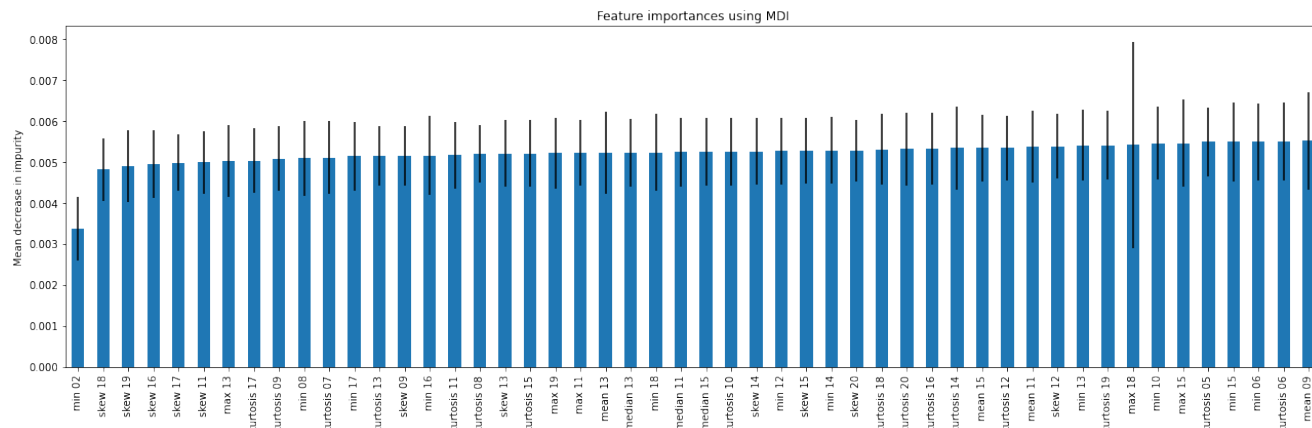


Figure 13: The error bars of Random Forest Mean Decrease Impurity feature importance



Figure 14: Sorted Random Forest Mean Decrease Impurity feature importance

8. Discussion & Conclusion

We designed machine and deep learning-oriented recognition systems for computing audio-level signals. As of the start, audio-level data analysis is made to discover and explore the dataset. Then, audio files are processed by frequency-domain and Mel spectrum processes to extract semantic features. Then, conventional machine learning-based algorithms are developed to fit audio features. As a final step, we designed learning systems by digesting the state-of-the-art neural network architectures for processing audio files to surpass our baseline accuracies.

Musical genre recognition classifies musical collections into categories based on audio, user, and track features. Poor labeling and the non-triviality of feature extraction are seen as major roadblocks. In the first, musical genres are loosely defined, with people frequently arguing the genre of a song. In our scenario, executing data-centric computational tasks is easier because the genres are supervised by the artists. As automatic extraction of semantic features is not differentiable and conceptualized by translating audio signals to visual signals and representing how the spectrum of frequencies evolves, there is no standard way to extract features in music information retrieval in the latter. The final hurdle is attempted to be overcome utilizing spectrum-based representations before feeding the ML model.

Trend-aware classifier-based learning algorithms were designed with the intention of capturing the dynamic interrelationships of cultures in order to detect genres autonomously. Our classes will cover a variety of musical genres, from modern to jazz, and is led and annotated by artists. We built machine learning methods to capture stochastic dynamics of music metadata and will build a neural networks for temporal dynamics of audio, from traditional ML algorithms like SVM to more advanced meta learning algorithms, while taking into account the socio-cultural environment. A multitude of machine learning-based techniques are employed to obtain the project's scope. To begin, we line with the basic and naive algorithms like Decision Trees and SVM to construct a baseline model to compare, and then we propose more complicated ensemble classification techniques like Random Forests and custom ones to compare our findings. The MLP and SVM classifiers were the winners, with a approximated accuracy of 60% in 16 classes.

About what we have learned, we have strengthened our intuitions by applying various classifiers. There are various points we have comprehended from the project. To begin with picking the most ideal classifier for an application requires an immense intuition, if there are any computational limitations. For instance, the KNN classifier is ready-to-go almost instantaneously, however, its performance is not the best. One must pick classifier carefully according to their objectives. Another example might be given in the context of CNNs. The most important property of CNNs is the translational invariance. In such applications where the input is shifted in a dimension and the output stays the same, CNNs might be considered as one of the promising classifiers as it considers the spatial information in the data.

In future directions, musical genre recognition system is a stepping stone for more complex music related machine learning tasks, such as music recommendation systems. In such complex applications, usually the objective is to perform as good as humans contextual actions. Our project might be used prior to complicated applications in order to build a framework for the application.

A. Appendix

A.1. Who did What?

A clear description of the contribution of each person are depicted in the following figure.

TASK NAME	START DATE	DAY OF MONTH*	END DATE	DURATION* (WORK DAYS)	DAYS COMPLETE*	DAYS REMAINING*	TEAM MEMBER	PERCENT COMPLETE
Project Progress								
Metadata-level Data Analysis	11/1	1	11/5	4	4	0	Mehmet Çalışkan	100%
Dimension Reduction: PCA & LDA	11/5	5	11/7	2	2	0	Berkay Erkan, Can Kocagil	100%
Monifold Learning	11/7	7	11/11	4	4	0	Can Kocagil	100%
Model Development for Classical ML algorithms	11/11	11	11/13	2	2	0	Can Kocagil, Berkay Erkan, Efe Eren Ceyani, Mehmet Çalışkan, Hammad Khan Musakhel	100%
Cross-validation for Classical ML algorithms	11/13	13	11/15				Can Kocagil, Berkay Erkan, Efe Eren Ceyani, Mehmet Çalışkan, Hammad Khan Musakhel	100%
Ensemble Learning Implementation	11/15	15	11/19	4	4	0	Can Kocagil	100%
Feature Importance Implementation	11/19	19	11/22	3	3	0	Can Kocagil	100%
Final Demo								
Audio-level Data Analysis	11/29	29	12/5	6	6	0	Can Kocagil	100%
Mel Spectram Computation on Audios	12/5	5	12/10	5	5	0	Can Kocagil	100%
Designing ML pipelines for audios							Can Kocagil, Berkay Erkan, Efe Eren Ceyani, Mehmet Çalışkan, Hammad Khan Musakhel	100%
Designing neural network architectures for audios	12/10	10	12/18	8	8	0	Can Kocagil, Berkay Erkan, Efe Eren Ceyani, Mehmet Çalışkan, Hammad Khan Musakhel	100%

Figure 15: The Visualization of the Contributions of Each Team Member

A.2. CNN Details

Layer (type:depth-idx)	Output Shape	Param #
ConvNet	--	--
Sequential: 1-1	[4, 256, 2, 20]	--
Sequential: 2-1	[4, 4, 128, 1292]	--
Conv2d: 3-1	[4, 4, 128, 1292]	112
ReLU: 3-2	[4, 4, 128, 1292]	--
Sequential: 2-2	[4, 8, 64, 646]	--
Conv2d: 3-3	[4, 8, 128, 1292]	296
ReLU: 3-4	[4, 8, 128, 1292]	--
MaxPool2d: 3-5	[4, 8, 64, 646]	--
Sequential: 2-3	[4, 16, 32, 323]	--
Conv2d: 3-6	[4, 16, 64, 646]	1,168
ReLU: 3-7	[4, 16, 64, 646]	--
MaxPool2d: 3-8	[4, 16, 32, 323]	--
Sequential: 2-4	[4, 32, 16, 161]	--
Conv2d: 3-9	[4, 32, 32, 323]	4,640
ReLU: 3-10	[4, 32, 32, 323]	--
MaxPool2d: 3-11	[4, 32, 16, 161]	--
Sequential: 2-5	[4, 64, 8, 80]	--
Conv2d: 3-12	[4, 64, 16, 161]	18,496
ReLU: 3-13	[4, 64, 16, 161]	--
MaxPool2d: 3-14	[4, 64, 8, 80]	--
Sequential: 2-6	[4, 128, 4, 40]	--
Conv2d: 3-15	[4, 128, 8, 80]	73,856
ReLU: 3-16	[4, 128, 8, 80]	--
MaxPool2d: 3-17	[4, 128, 4, 40]	--
Sequential: 2-7	[4, 256, 2, 20]	--
Conv2d: 3-18	[4, 256, 4, 40]	295,168
ReLU: 3-19	[4, 256, 4, 40]	--
MaxPool2d: 3-20	[4, 256, 2, 20]	--
Flatten: 1-2	[4, 10240]	--
Sequential: 1-3	[4, 8]	--
Linear: 2-8	[4, 5120]	52,433,920
ReLU: 2-9	[4, 5120]	--
Linear: 2-10	[4, 8]	40,968
Total params: 52,868,624		
Trainable params: 52,868,624		
Non-trainable params: 0		
Total mult-adds (G): 1.43		
Input size (MB): 7.94		
Forward/backward pass size (MB): 104.63		
Params size (MB): 211.47		
Estimated Total Size (MB): 324.04		

Figure 16: 2-D CNN Model Architecture

References

- [1] K. Benzi, M. Defferrard, P. Vandergheynst, and X. Bresson. FMA: A dataset for music analysis. *CoRR*, abs/1612.01840, 2016.
- [2] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021.
- [3] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [4] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- [5] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. J. Skerry-Ryan, R. A. Saurous, Y. Agiomyrgiannakis, and Y. Wu. Natural TTS synthesis by conditioning wavenet on mel spectrogram predictions. *CoRR*, abs/1712.05884, 2017.
- [6] S. Wold, K. Esbensen, and P. Geladi. Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3):37–52, 1987.
- [7] D. Wright. *Mathematics and music*, volume 28. American Mathematical Soc., 2009.