

Music Genre Classification

CS 4824

Group: G5-CS

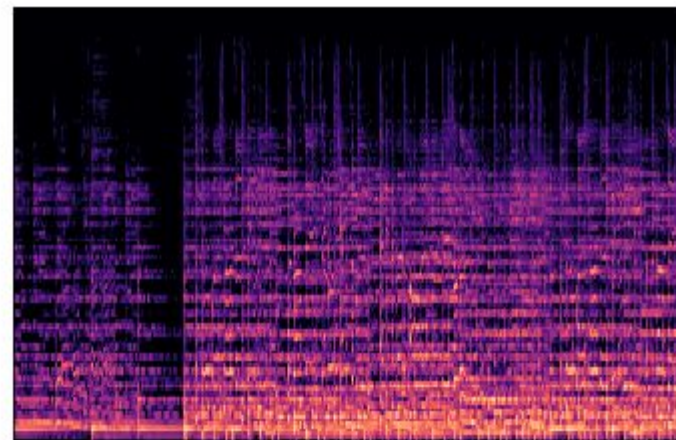
Dataset

- GTZAN dataset.
- It has 9990 data points, categorized into 10 different genres.
- It holds audio files, image files, and a CSV file of the features derived from the said audio files.

Dataset

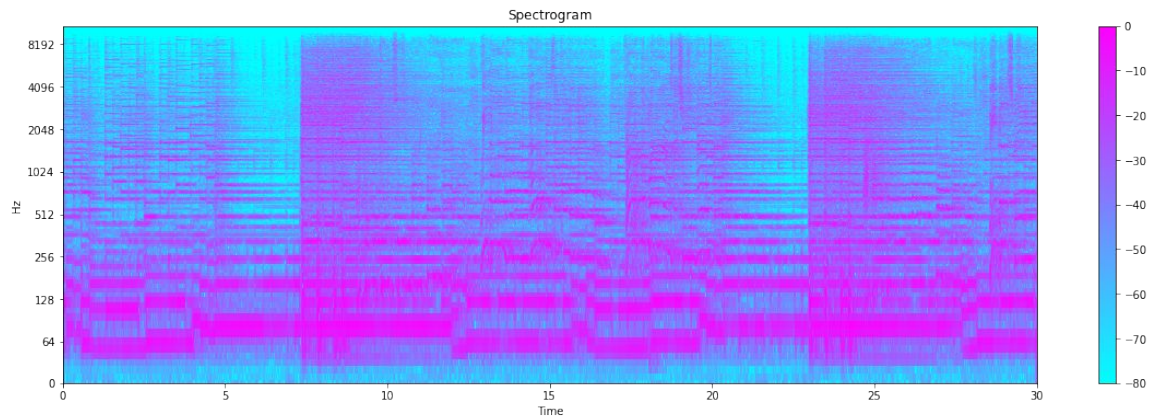
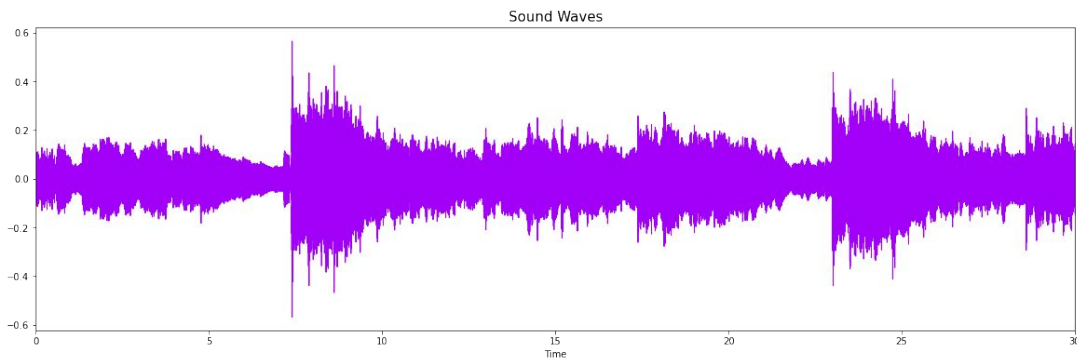
filename	length	chroma_stft_mean	chroma_stft_var	rms_mean	rms_var	spectral_centroid...	spectral_centroid...	spe
9990 unique values								
blues_00000_0.wav	66149	0.3354863639104865	0.89184829281568527	0.1304858236948384	0.8033218842268587215	1773.8658319944662	167541.6388866573	1972.
blues_00000_1.wav	66149	0.3438653512477875	0.8861465268386467	0.1126992478966713	0.8014436854273682366	1816.4937765281948	98525.69886581325	2810.
blues_00000_2.wav	66149	0.34681475162586184	0.89224288918627365	0.13208338184833527	0.804626399326886844	1788.539718722745	111407.43761296556	2884.
blues_00000_3.wav	66149	0.3636397884616852	0.88685615658768071	0.1325647234916687	0.8024475634563714266	1655.289445486881	111952.2845173748	1968.
blues_00000_4.wav	66149	0.33557942589651184	0.88812854439828157	0.14328888688881818	0.801788886175967753	1638.4561993518775	79667.26765440499	1948.
blues_00000_5.wav	66149	0.376697347164154	0.89578218781227188	0.1326178814278412	0.8035825634848675864	1994.9152192785897	211788.6195693862	2152.
blues_00000_6.wav	66149	0.3799887485284773	0.8888273119264526	0.1383347786281706	0.8031662711487889264	1962.1588958726888	177443.87884548853	2146.
blues_00000_7.wav	66149	0.33187994360923767	0.89211885184849686	0.140688323677863	0.802545942886726427	1701.8989235781698	35678.13861562547	1979.
blues_00000_8.wav	66149	0.347877832321167	0.89428917183158782	0.13312998963458815	0.802538164146244526	1746.4739823276876	138873.93124392346	1887.
blues_00000_9.wav	66149	0.3588612548245856	0.88295789639787674	0.11531286228388412	0.8018468289924858749	1763.9489423859536	61493.42312137413	1874.
blues_00001_0.wav	66149	0.4824888818528172	0.89839742588891782	0.09382367269992828	0.808387555478585823	1279.1822416189314	486513.8167436446	1921.
blues_00001_1.wav	66149	0.34588726413728807	0.891838487238884	0.09465641528387996	0.8014944711963458849	1513.7639693588	214748.8044331843	2891.
blues_00001_2.wav	66149	0.33811864256858826	0.8838818367428131	0.09777681888389555	0.80138564663939178	1388.8697277374274	154289.12244842653	1588.
blues_00001_3.wav	66149	0.33878858499758519	0.89393848999758519	0.08536521345376968	0.8026413983186884	1479.538437863734	64888.3731518855	1986.

60 features

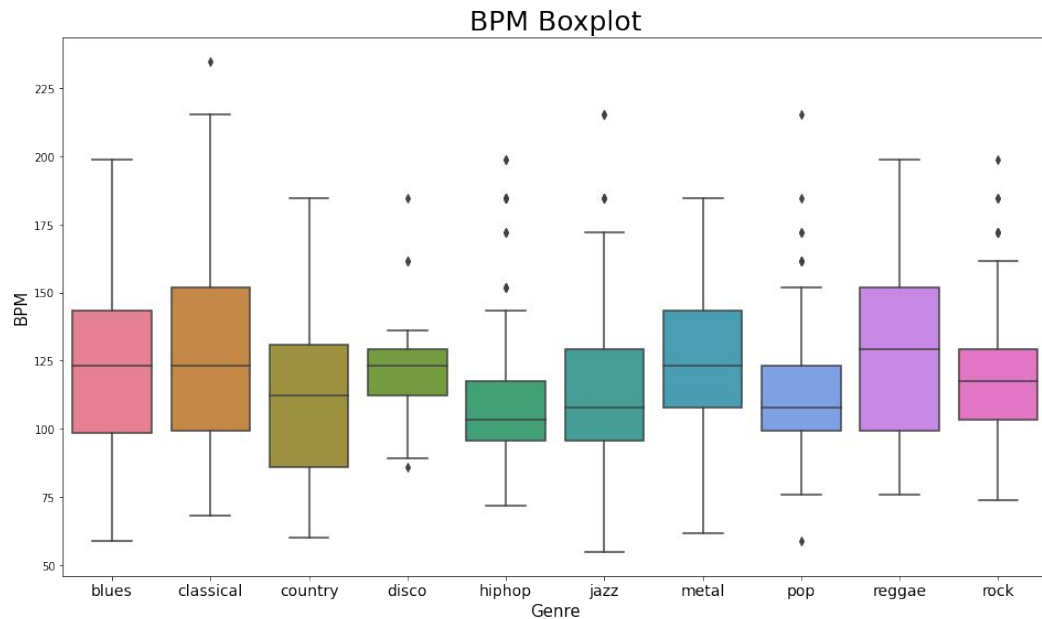
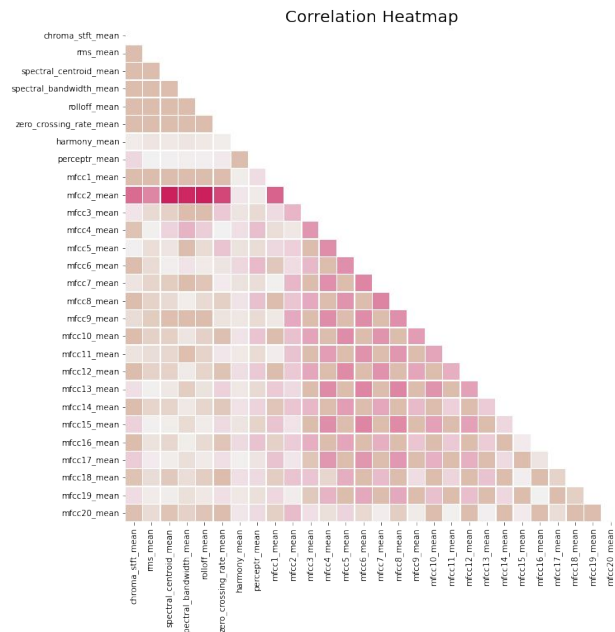


Mel Spectrogram

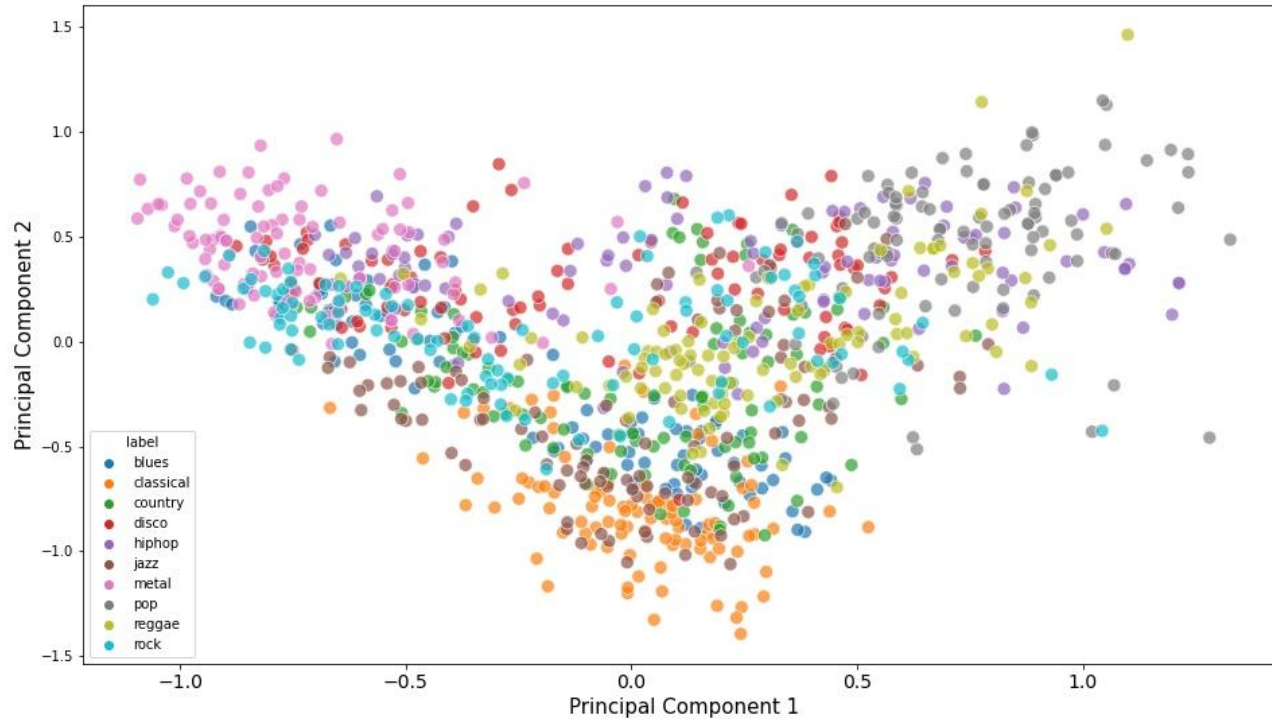
Data Exploration: Raw data



Data Exploration: Features



PCA on Genres



Model

- Raw Audio processing
 - SVM
 - Deep Neural Network
- Spectrogram Images
 - Convolutional NN
- Audio features dataset
 - Naive Bayes
 - Random Forest
 - K-NN
 - Support Vector Machine
 - Neural Networks

Raw Audio Processing

Short-time fourier transform

- Split the 30 second audio into 10 slices of 3 sec each
- Apply the STFT on the 3 sec files
- Extract the principal components
- 9990 Samples with 10 principal components
- Split ratio 70 - 15 - 15

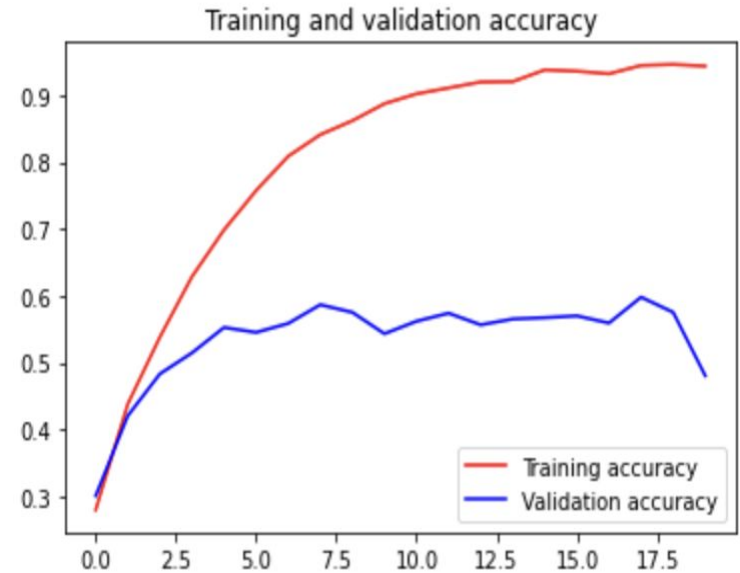
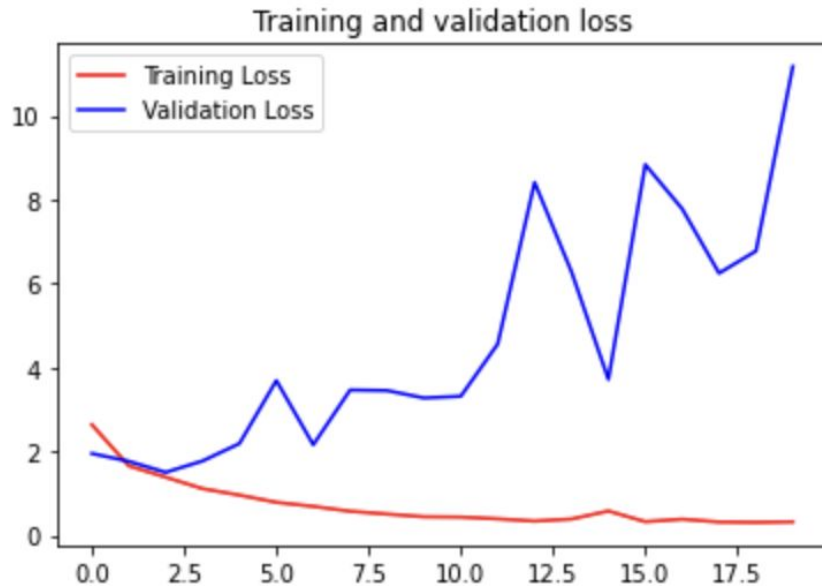
Raw Audio Processing

Short-time fourier transform

- Neural Network with 6 hidden layers
 - 2570 -> 1024 -> 512 -> 256 -> 128 -> 64 -> 10
 - Testing accuracy 55.2%
- SVM with RBF kernel
 - Testing accuracy 61.1%

Raw Audio Processing

Short-time fourier transform



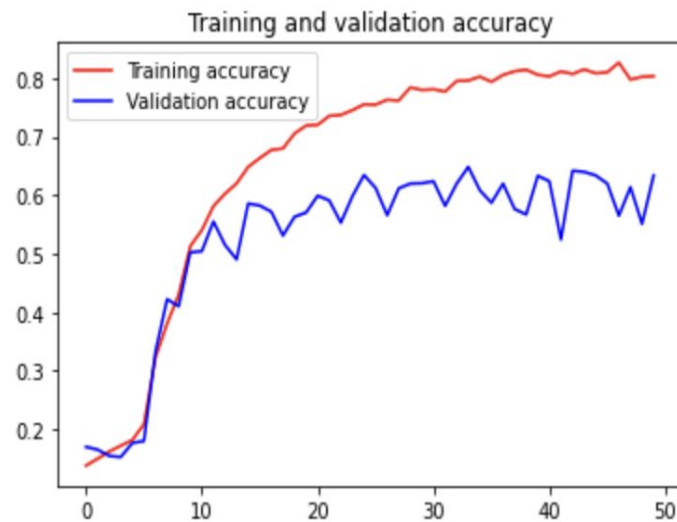
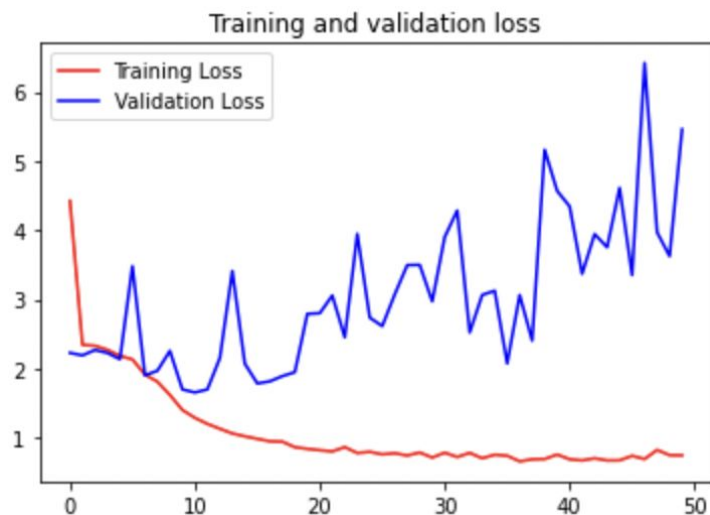
Raw Audio Processing

MFCC features

- Neural Network with 6 hidden layers
 - ReLU activation for hidden layers
 - 2570 -> 1024 -> 512 -> 256 -> 128 -> 64 -> 10
 - Testing accuracy 62.1%
- SVM with RBF kernel
 - Testing accuracy 56.7%

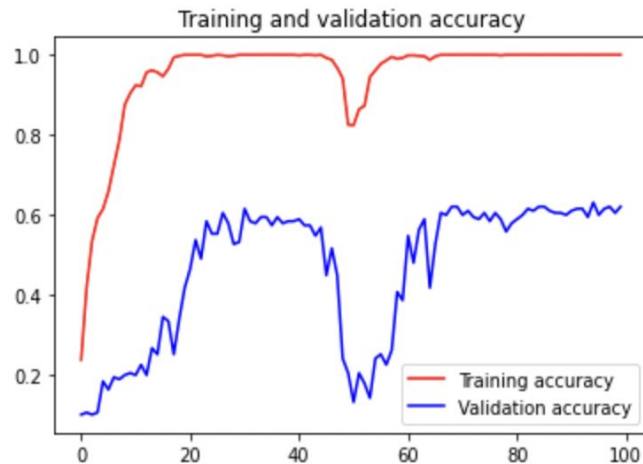
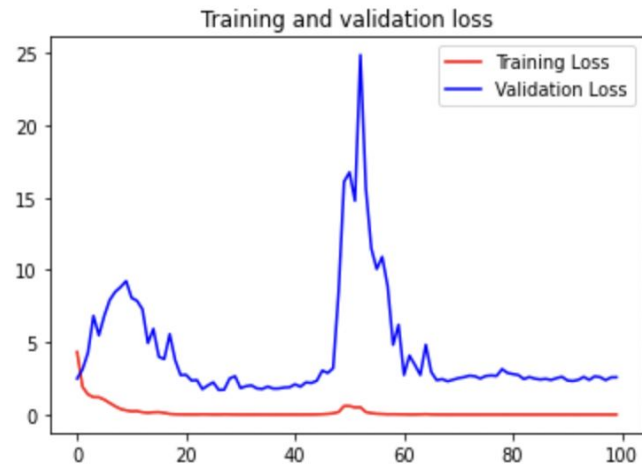
Raw Audio Processing

MFCC features



Spectrogram Images

- Convolutional NN
 - Batch normalization
 - Max pooling
 - Relu activation
- Dataset size 999 samples
- Testing Accuracy 61.9%



Audio features dataset

- We normalize all the features.
- We use a 60-20-20 train-validation-test split.
- We tune our hyperparameters (if any) using the validation set.
- We finally test on our set-out test set and report the results as follows

Results

Naive Bayes	Normalized audio features	9990	0.523
1-NN	Normalized audio features	9990	0.914
9-NN	Normalized audio features	9990	0.85
Decision tree	Normalized audio features	9990	0.640
Random Forest	Normalized audio features	9990	0.868
SVM	Normalized audio features	9990	0.747
Logistic Regression	Normalized audio features	9990	0.690
Neural Network	Normalized audio features	9990	0.672

Understanding our best models

KNN				
	precision	recall	f1-score	support
blues	0.95	0.94	0.94	208
classical	0.90	0.93	0.92	203
country	0.85	0.84	0.85	186
disco	0.90	0.93	0.91	199
hiphop	0.95	0.89	0.92	218
jazz	0.87	0.90	0.88	192
metal	0.96	0.98	0.97	204
pop	0.95	0.93	0.94	180
reggae	0.92	0.93	0.92	211
rock	0.89	0.86	0.87	197
accuracy			0.91	1998
macro avg	0.91	0.91	0.91	1998
weighted avg	0.91	0.91	0.91	1998

Random Forest				
	precision	recall	f1-score	support
blues	0.90	0.86	0.88	208
classical	0.92	0.99	0.95	203
country	0.72	0.81	0.76	186
disco	0.84	0.85	0.84	199
hiphop	0.94	0.85	0.89	218
jazz	0.85	0.90	0.87	192
metal	0.87	0.96	0.91	204
pop	0.92	0.91	0.91	180
reggae	0.89	0.86	0.87	211
rock	0.87	0.69	0.77	197
accuracy			0.87	1998
macro avg	0.87	0.87	0.87	1998
weighted avg	0.87	0.87	0.87	1998

KNN can pick up on metal, but not on country.

Random Forest is good with classical but not with rock.

Questions?