# *AUDIO GENRE CLASSIFICATION*

*A Project under the guidance of*
*Prof. Desire Sidibe*

PROJECT TEAM

VAMSHI KODIPAKA

BHARGAV SHAH

PARMAR HARDIKSINH
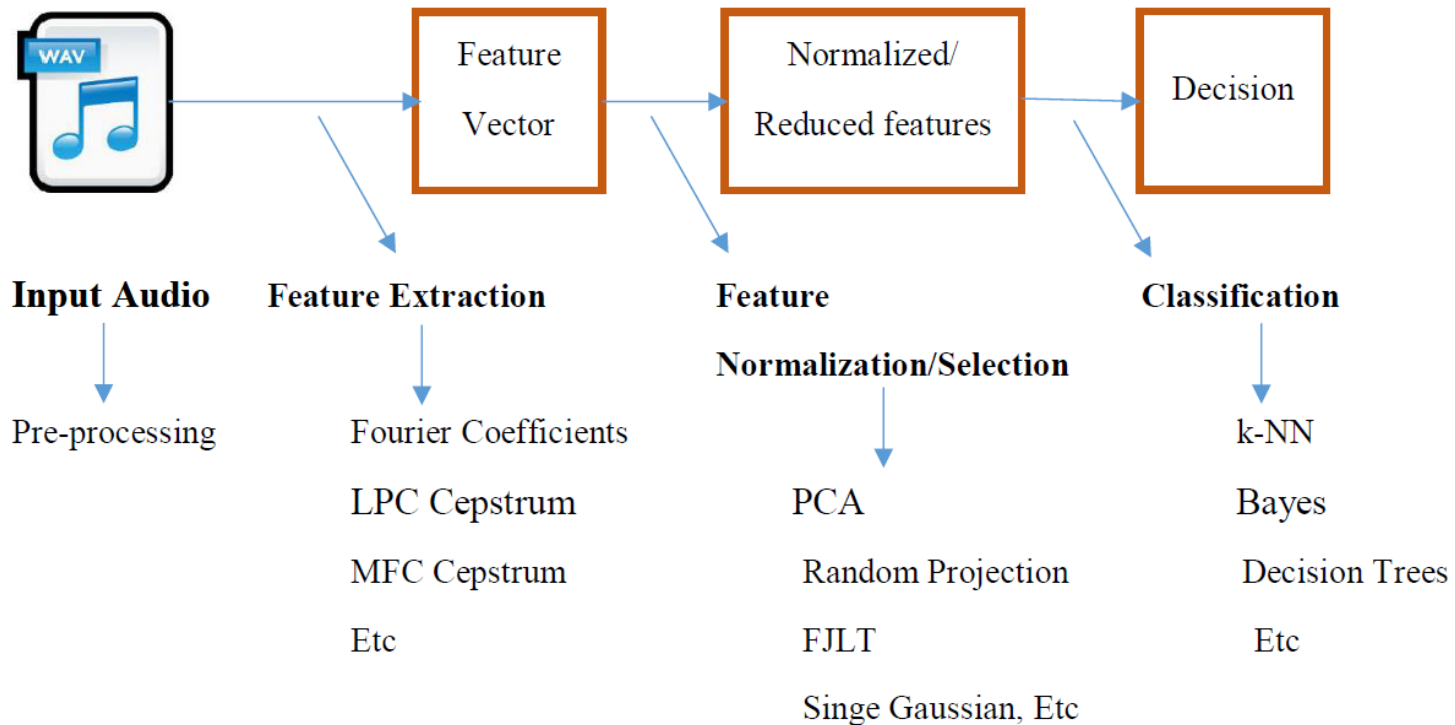
# Overview

@Research Project for LAAS CNRS - 2018

- Introduction
    - AUDIO FILES
    - MUSIC DATA SET  @ ISMIR -2004 (729 TRACKS)

- Methodology

    - Feature Extraction - MFCC
    - Feature Selection - PCA
    - Feature Classification – KNN

- Problems of Dimension Reductionality

- GTZAN Dataset -1000 au files (Quick implementation)

- Output/Results Achieved

- Future Work

# Audio Understanding : Pipeline



**Input Audio** → **Feature Extraction** → **Feature Normalization/Selection** → **Classification**

| Input Audio | Feature Extraction | Feature Normalization/Selection | Classification |
|---|---|---|---|
| Pre-processing | Fourier Coefficients | PCA | k-NN |
| | LPC Cepstrum | Random Projection | Bayes |
| | MFC Cepstrum | FJLT | Decision Trees |
| | Etc | Singe Gaussian, Etc | Etc |

**AIM:**

- Using Machine Learning Pipeline on Music Dataset we need to classify the dataset
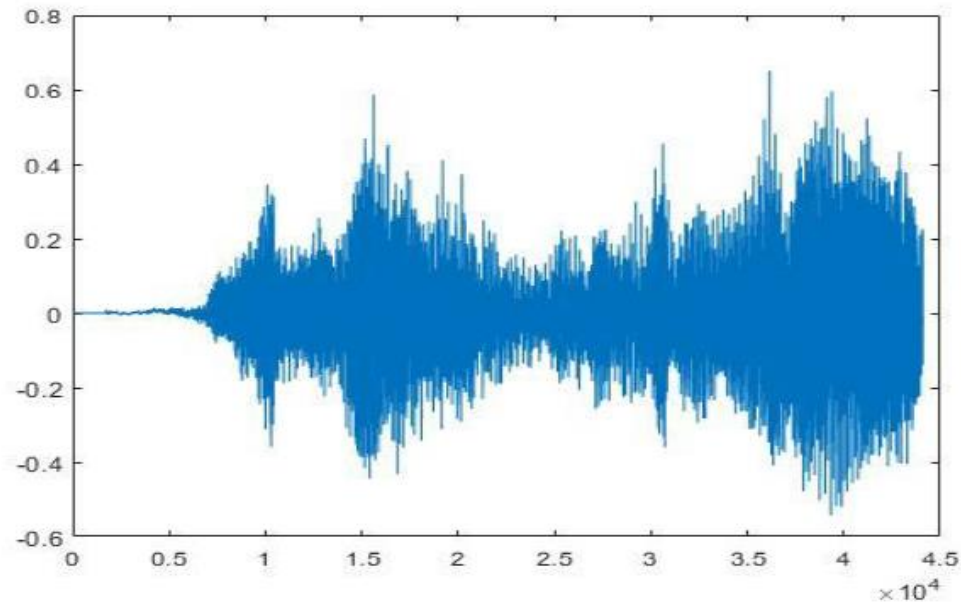
# Audio Access in MATLAB

**BASICS OF AUDIO PROCESSING:**

1. <u>To load .mp3 file:</u>    [y,Fs]=audioread('artist_1_album_1_track_1.mp3')

   Here, y gives amplitude and Fs gives frequency.

2. <u>Plot .mp3 with amplitude (vs) frequency:</u>

   plot(y(1:44100,1))



Amplitude vs Frequency Plot

# Audio Access in MATLAB

4. Underline: To play .mp3 using sound command:

    sound(y(1:441000,1),Fs)

5. Underline: To play .mp3 using audioplayer object in MATLAB:

    p =audioplayer(y,Fs)

    play (p) :: plays loaded music fil

    pause(p) :: pauses music file

    stop(p) :: stops music file

    start(p) :: starts from point of stop

    clear(p) :: clears 'p' object's music file.

```
Command Window
>> p=audioplayer(y,Fs)

p =

    audioplayer with properties:

            SampleRate: 44100
          BitsPerSample: 16
       NumberOfChannels: 2
               DeviceID: -1
          CurrentSample: 1
           TotalSamples: 1686000
                Running: 'off'
               StartFcn: []
                StopFcn: []
               TimerFcn: []
            TimerPeriod: 0.0500
                    Tag: ''
               UserData: []
                   Type: 'audioplayer'

>> play (p)
>> pause(p)
>> stop(p)
>>
```

```
Workspace
Name ▲     Value
Fs         44100
y          1686000x2 double
```

```
>> info=audioinfo('artist_1_album_1_track_1.mp3')

info =

    struct with fields:

                Filename: 'D:\AudioGenreClassifier-master\AudioGenreClassifier-master\artist_1_album_1_track_1.mp3'
       CompressionMethod: 'MP3'
             NumChannels: 2
              SampleRate: 44100
            TotalSamples: 1687360
                Duration: 38.2621
                   Title: []
                 Comment: []
                  Artist: []
                 BitRate: 128
```

Mel refers to 'Melody'.

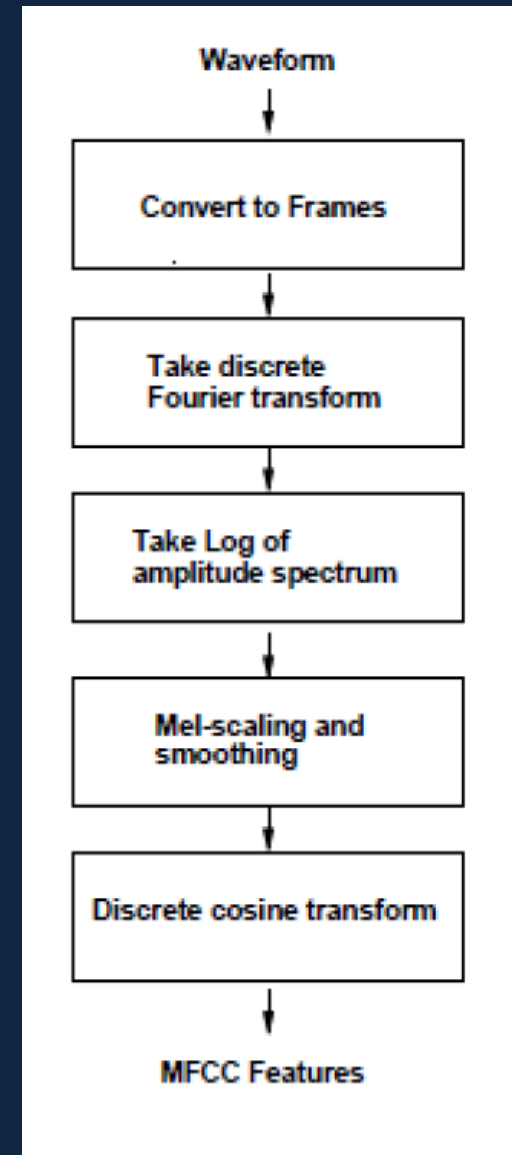Cepstrum means the IFT of the logarithm of the estimated spectrum of a signal.

**Definition:**

In sound processing, the **Mel-Frequency Cepstrum** (**MFC**) is a representation of the short-term power spectrum of a sound, based on a linear cosine transform of a log power spectrum on a nonlinear mel scale of frequency

To retrieve info of music file:

Audioinfo("filename")

# FEATURE EXTRACTION : MFCC

1. Take the Fourier transform of (a windowed excerpt of) a signal.

2. Map the powers of the spectrum obtained above onto the mel scale, using triangular overlapping windows.

3. Take the logs of the powers at each of the mel frequencies.

4. Take the discrete cosine transform of the list of mel log powers, as if it were a signal.

5. The MFCCs are the amplitudes of the resulting spectrum.

Waveform

↓

Convert to Frames

↓

Take discrete Fourier transform

↓

Take Log of amplitude spectrum

↓

Mel-scaling and smoothing

↓

Discrete cosine transform

↓

MFCC Features

# FEATURE EXTRACTION : TYPES

**From conventional spectral analysis:**

1. IFT

   a. Positive Cepstrum
   b. Negative Cepstrum

**Linear Predictive Coding Cepstrum(LPC Cpestrum):**

The LPC vector is defined by $[a_0, a_1, a_2, \ldots a_p]$ and the CC vector is defined by $[c_0 c_1 c_2 \ldots c_p \ldots c_{n-1}]$

| LPC Cepstrum $(c_m)$ | |
|---|---|
| $c_0 = \log G^2$ | |
| $c_m = a_m + \sum_{k=1}^{m-1} \left(\frac{k}{m}\right) c_k a_{m-k}, \quad 1 \le m \le p$ | $G = e^{c_0/2}$ |
| $c_m = \sum_{k=1}^{m-1} \left(\frac{k}{m}\right) c_k a_{m-k}, \quad m > p$ | $a_m = c_m - \sum_{k=1}^{m-1} \left(\frac{k}{m}\right) c_k a_{m-k}, \quad 1 \le m \le p$ |

# MEL-FREQUENCY

**Note:** Mel-Scale is approximately linear for low-frequency (f<500Hz) and logarithmic for high frequencies

1. Noise Sensitivity

2. Use of MFCC

3. Pre-processing

4. Calculation:

$$M(mel\_freq) = 1127 * \log(1 + f/700)$$

where f is frequency in linear scale
and M is frequency is mel scale.

This modulation of log acts as a weight vector like in $(y = w^T X - t)$ in a classical regression model

Let x[n] is framed through the entire signal ,

DFT

has a window size fixed          x[n]  →   X[k]  -- |X[k]|

Sample plot ::
amplitude X[k] (vs) discrete frequency 'k'

k= N/2;   fmax > Fs/2 ;   Fs = 8KHz;     fs becomes 4KHz

If triangular train filter = 100Hz  → Fs=4000/100 = 40points

Now suppose, this triangular filter is applied to 8KHz =80points will be generated.

To convert into k to f, we have:    $f = 2\pi k/N$;

# MEL-FREQUENCY



$$\text{Pitch } (mels) = 3322 \log_{10}(1 + f/1000)$$

Alternatively, we can approximate curve as:
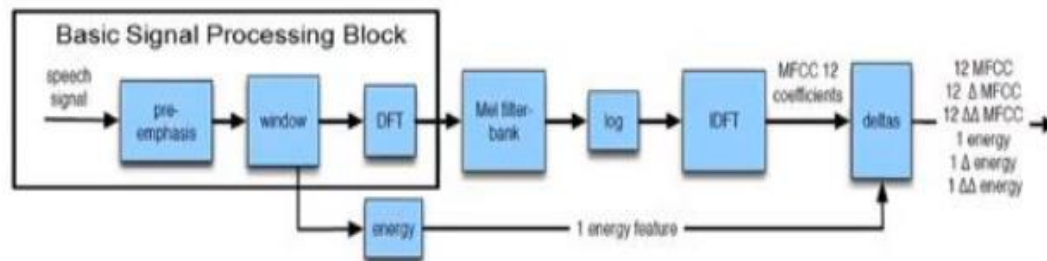
$$\text{Pitch } (mels) = 1127 \log_e(1 + f/700)$$
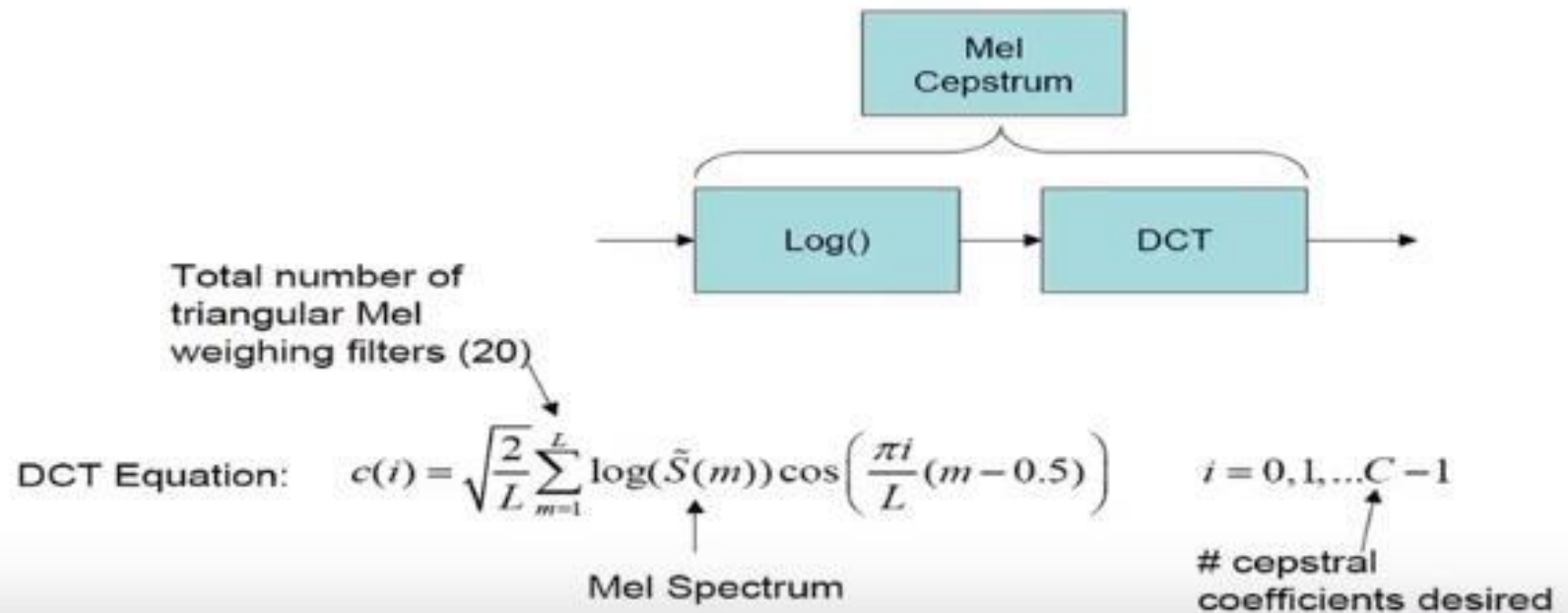
Log- Representation of MEL-FREQ

# MFCC FEATURES



Mel Filter bank



Block diagram of Extracting a sequence of 39-dimensional MFCC feature vectors

# MFCC FEATURES



Instead of IFT we can take DCT: MFCC
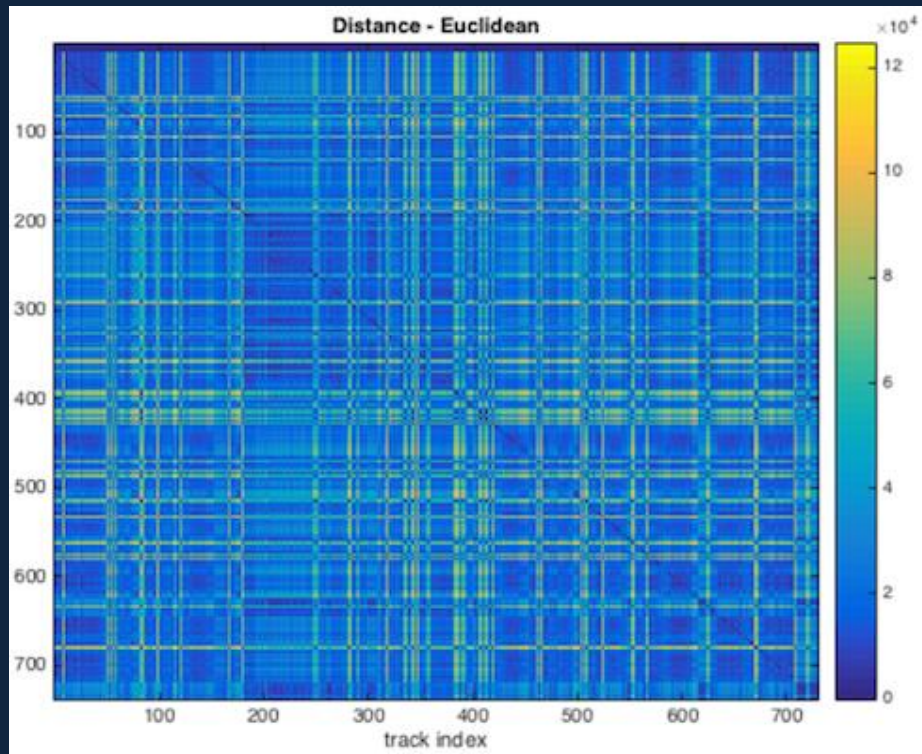
# MFCC FEATURE EXTRACTION: OUTPUTS



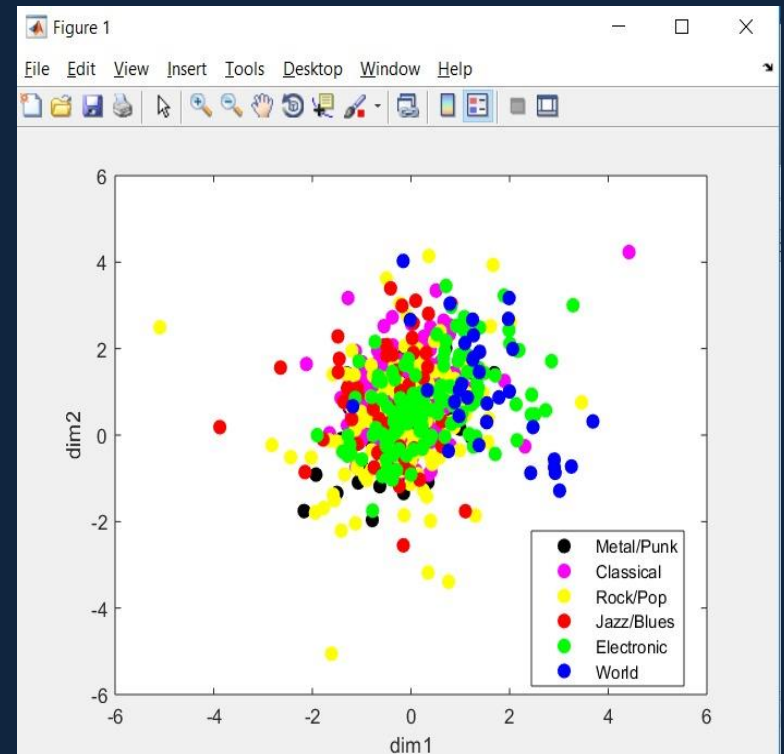| Name ▲ | Value |
|---|---|
| distance | 743x729 double |
| euclideanDistance | 1x79 double |
| features | 729x79 double |
| filename | 1x729 cell |
| files | 729x1 struct |
| frobeniusNorm | 0 |
| GMModel | 1x1 gmdistributio... |
| groupNames | 729x1 cell |
| groups | 729x1 double |
| i | 745 |
| j | 745 |
| k | 6 |
| matrix1 | 79x79 double |
| matrix2 | 79x79 double |
| mfcc_artist_100_album_1_track_1 | 79x79 double |
| mfcc_artist_100_album_1_track_2 | 79x79 double |
| mfcc_artist_100_album_1_track_3 | 79x79 double |
| mfcc_artist_100_album_1_track_4 | 79x79 double |
| mfcc_artist_100_album_2_track_1 | 79x79 double |
| mfcc_artist_100_album_2_track_2 | 79x79 double |
| mfcc_artist_100_album_2_track_3 | 79x79 double |
| mfcc_artist_100_album_3_track_1 | 79x79 double |
| mfcc_artist_100_album_3_track_2 | 79x79 double |
| mfcc_artist_100_album_3_track_3 | 79x79 double |
| mfcc_artist_100_album_4_track_1 | 79x79 double |
| mfcc_artist_100_album_4_track_2 | 79x79 double |
| mfcc_artist_100_album_4_track_3 | 79x79 double |
| mfcc_artist_100_album_4_track_4 | 79x79 double |
| mfcc_artist_100_album_5_track_1 | 79x79 double |
| mfcc_artist_100_album_5_track_2 | 79x79 double |
| mfcc_artist_100_album_5_track_3 | 79x79 double |
| mfcc_artist_101_album_1_track_1 | 79x79 double |
| mfcc_artist_101_album_1_track_2 | 79x79 double |
| mfcc_artist_101_album_1_track_3 | 79x79 double |
| mfcc_artist_102_album_1_track_1 | 79x79 double |
| mfcc_artist_102_album_1_track_2 | 79x79 double |

| Name ▲ | Value |
|---|---|
| mfcc_artist_96_album_1_track_2 | 79x79 double |
| mfcc_artist_97_album_1_track_1 | 79x79 double |
| mfcc_artist_97_album_1_track_2 | 79x79 double |
| mfcc_artist_97_album_1_track_3 | 79x79 double |
| mfcc_artist_98_album_1_track_1 | 79x79 double |
| mfcc_artist_98_album_1_track_2 | 79x79 double |
| mfcc_artist_98_album_1_track_3 | 79x79 double |
| mfcc_artist_99_album_1_track_1 | 79x79 double |
| mfcc_artist_99_album_1_track_2 | 79x79 double |
| mfcc_artist_99_album_1_track_3 | 79x79 double |
| mfcc_artist_99_album_1_track_4 | 79x79 double |
| mfcc_artist_9_album_1_track_1 | 79x79 double |
| mfcc_artist_9_album_1_track_2 | 79x79 double |
| mfcc_artist_9_album_1_track_3 | 79x79 double |
| mfcc_artist_9_album_1_track_4 | 79x79 double |
| mfcc_artist_9_album_1_track_5 | 79x79 double |
| mfcc_artist_9_album_1_track_6 | 79x79 double |
| mfcc_artist_9_album_2_track_1 | 79x79 double |
| mfcc_artist_9_album_2_track_10 | 79x79 double |
| mfcc_artist_9_album_2_track_11 | 79x79 double |
| mfcc_artist_9_album_2_track_2 | 79x79 double |
| mfcc_artist_9_album_2_track_3 | 79x79 double |
| mfcc_artist_9_album_2_track_4 | 79x79 double |
| mfcc_artist_9_album_2_track_5 | 79x79 double |
| mfcc_artist_9_album_2_track_6 | 79x79 double |
| mfcc_artist_9_album_2_track_7 | 79x79 double |
| mfcc_artist_9_album_2_track_8 | 79x79 double |
| mfcc_artist_9_album_2_track_9 | 79x79 double |
| mfcc_artist_9_album_3_track_1 | 79x79 double |
| mfcc_artist_9_album_3_track_2 | 79x79 double |
| mfcc_artist_9_album_3_track_3 | 79x79 double |
| mfcc_artist_9_album_3_track_4 | 79x79 double |
| mfcc_artist_9_album_3_track_5 | 79x79 double |
| mfcc_artist_9_album_3_track_6 | 79x79 double |
| values | 1x729 double |
| variables | 745x1 cell |

**Set of feature Extracted : ISMIR2004**

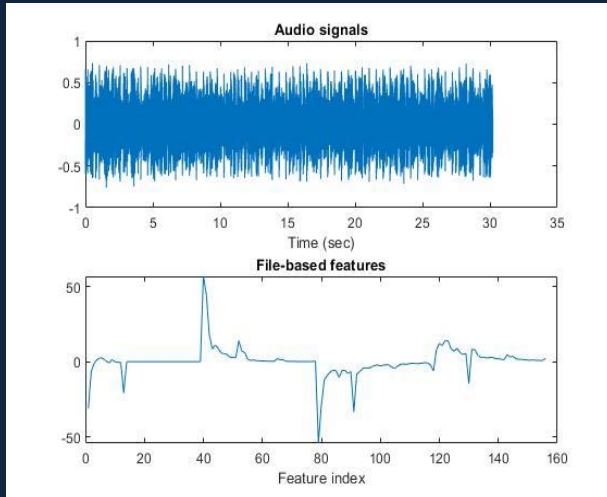# MFCC FEATURE EXTRACTION: OUTPUTS



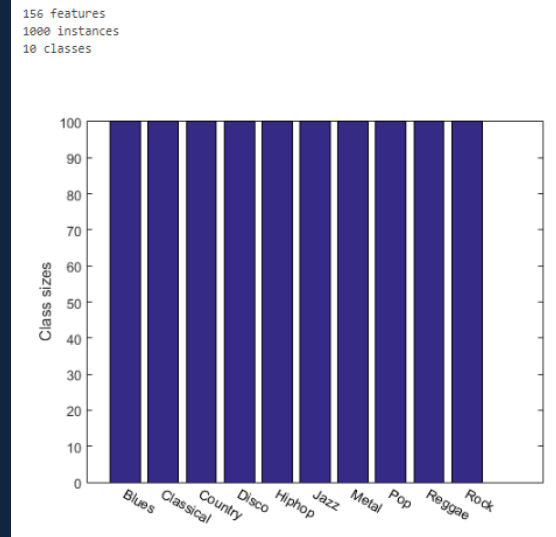Calculation of Frobenius norm of the MFCC
Euclidean Distance

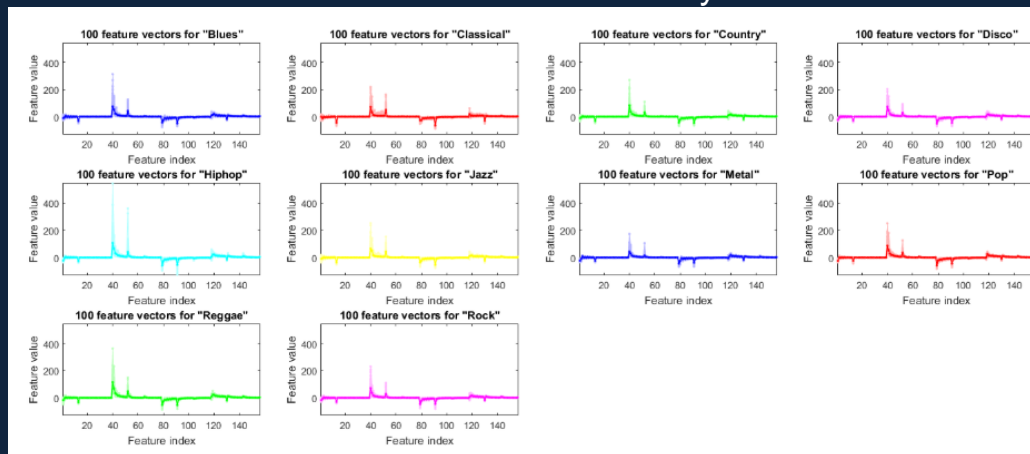Mean of MCC and GMM

# WORKING ON GTZAN DATASET



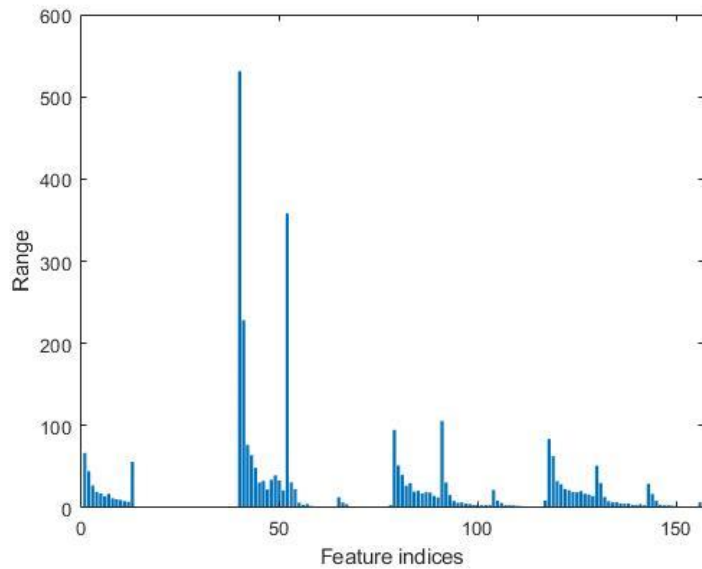**Feature extraction**



**Data Visualization**

Class-wise Features Density

# GTZAN



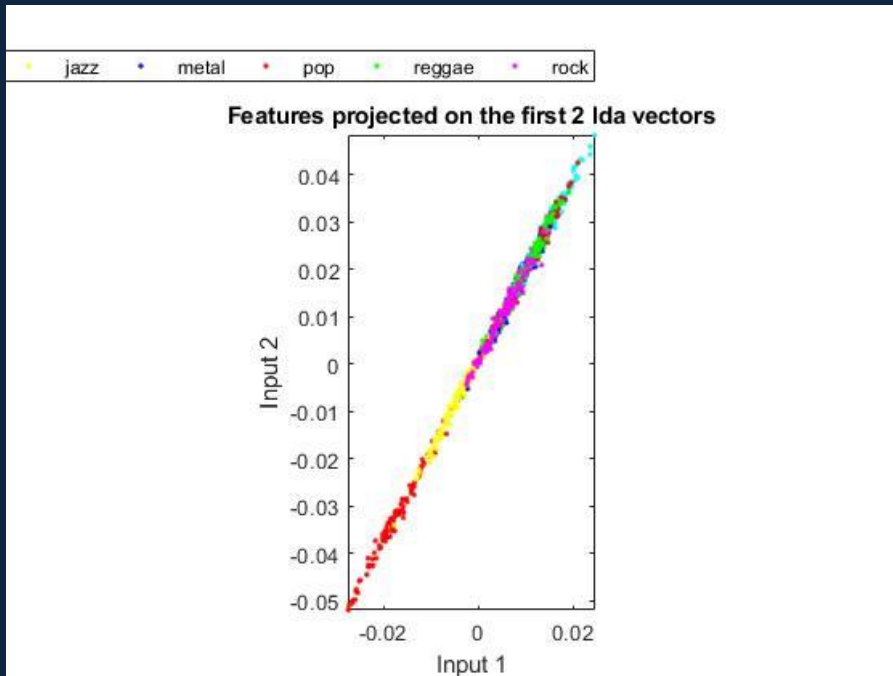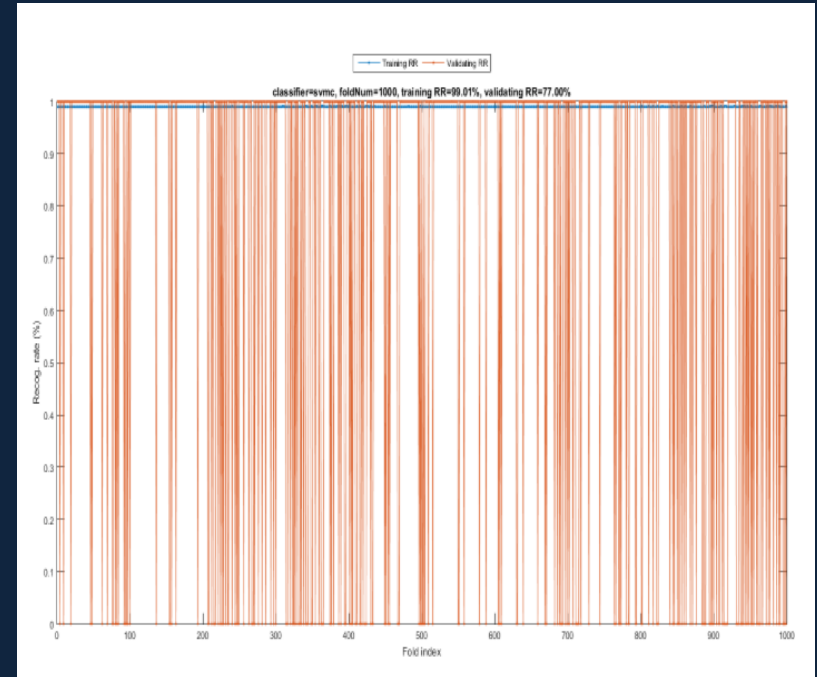Feature Density



**Dimension Reduction**

# GTZAN



LDA Reduction



SVM Classifier

# GTZAN



**Confusion Matrix**

# FUTURE WORK

- 1. Problem of Dimension Reduction – PCA

- 2. KNN Classifier

- 3. Testing and Classifying the Sample Audio

# Thank you…