TRIBHUVAN UNIVERSITY

INSTITUE OF ENGINEERING

CENTRAL CAMPUS, PULCHOWK

**MUSIC CLASSIFICATION BASED ON GENRE AND MOOD**

**Submitted By:**

Ayush Shakya [069/BCT/505]

Bijay Gurung [069/BCT/512]

Mahendra Singh Thapa [069/BCT/519]

Mehang Rai [069/BCT/524]

A PROJECT WAS SUBMITTED TO THE DEPARTMENT OF ELECTRONICS AND
COMPUTER ENGINEERING IN PARTIAL FULLFILLMENT OF THE
REQUIREMENT FOR THE BACHELORS DEGREE IN COMPUTER ENGINEERING

(August, 2016)

TRIBHUVAN UNIVERSITY

INSTITUTE OF ENGINEERING

PULCHOWK CAMPUS

DEPARTMENT OF ELECTRONICS AND COMPUTER ENGINEERING

The undersigned certify that they have read, and recommended to the Institute of Engineering for acceptance, a project report entitled "Music Classification Based on Genre and Mood" submitted by Ayush Shakya, Bijay Gurung, Mahendra Singh Thapa and Mehang Rai in partial fulfilment of the requirements for the Bachelors degree in Computer Engineering.

------------------------------

Supervisor, Dr. Basanta Joshi

Department of Electronics and Computer Engineering

------------------------------          ------------------------------

Internal Examiner                        External Examiner

Name of Internal Examiner                Name of External Examiner

Title, Affiliation                       Title, Affiliation

------------------------------          ------------------------------

Dr. Nanda Bikram Adhikari                Dr. Dibakar Raj Pant

Deputy Head                              Head

Department of Electronics & Computer     Deparment of Electronics & Computer

Engineering,                             Engineering,

Institute of Engineering, Pulchowk Campus,   Institute of Engineering,Pulchowk Campus,

Tribhuvan University, Nepal              Tribhuvan University, Nepal

**Date of Approval:**

# COPYRIGHTS

The author has agreed that the Library, Department of Electronics and Computer Engineering, Pulchowk Campus, Institute of Engineering may make this report freely available for inspection. Moreover, the author has agreed that permission for extensive copying of this project report for scholarly purpose may be granted by the supervisors who supervised the project work recorded herein or, in their absence, by the Head of the Department wherein the project report was done. It is understood that the recognition will be given to the author of this report and to the Department of Electronics and Computer Engineering, Pulchowk Campus, Institute of Engineering in any use of the material of this project report. Copying or publication or the other use of this report for financial gain without approval of to the Department of Electronics and Computer Engineering, Pulchowk Campus, Institute of Engineering and authors written permission is prohibited.

Request for permission to copy or to make any other use of the material in this report in whole or in part should be addressed to:

Head

Department of Electronics and Computer Engineering

Pulchowk Campus, Institute of Engineering

Lalitpur, Kathmandu

Nepal

# ACKNOWLEDGMENT

It has been a great pleasure to work with different individuals whose perpective and ideas has directly or indirectly assisted or motivated us. We would like to gratitude to everyone whose ideas lead to this completion of the project. It is a matter of fact a great privilege for us to acknowledge their assistance and contributions to our project.

First and foremost, we would like to express our sincere gratitude towards Dr. Basanta Joshi, our supervisor under whose supervision there was successful completion of our project. Without his invaluable guidance and suggestions, it would have been a difficult journey for us. His useful suggestions for this whole work and cooperative behavior are sincerely acknowledged.

We would like to thank the Department of Electronics and Computer Engineering for adding this major project as part of final year curriculum and hence giving us this opportunity to undertake the project. The great need of research, time and sheer coding has allowed us to harness our skills, experience and knowledge.

We are also grateful to Dr. Nanda Bikram Adhikari for letting us carry out this project and co-operating with us to carry out the project smoothly. We would like to thank and express our gratitude to all out respective subject teachers for sharing their precious knowledge, constant support and guidance.

Last but not the least, we would like to thank our friends for motivating us and providing numerous assistance throughout the project development duration.

# ABSTRACT

This report describes and documents all the aspects and working functionality of our final year project titled Music Classification System Based on Genre and Mood. The project is part for the curriculum for the subject Major Project under the course of final year of B.E. in Computer Engineering. As the title itself describes the overall aim of the project is to develop a system capable of classifying music based on Genre and mood, with the availability of large number of digital media and the disorder introduced being the primary motivation.

The methodology used is that of a modular system consisting of two main stages. The first stage involves the preprocessing of the raw audio data resulting in the extraction of a number of features pertaining to music signal: Intensity, MFCC, rhythm, pitch. Each feature extractor reduces the information content in the raw data to a vector in a small number of dimensions. Or in other words we can say that feature extractor analyses the music signal and extracts its respective features compatible for further processing. It requires intensive knowledge of digital signal analysis and processing, signal sampling,etc. The second stage comprises of all the machine learning portion. In it, the set of feature vectors are classified(indexed) into certain clusters by the use of certain algorithms: K-means, Support Vector Machines and Artificial Neural Networks. This technically requires knowledge of all those respective algorithms.

This report also documents our approach towards the system development following the various aspects of Software Engineering. UML diagrams have been used to model the entire system and ERD diagrams have been used to show the relationship between the various entities in our system and iterative development method was chosen for the development of our system. Java language along with spring framework was used to build our whole system along with the GUI.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

| | |
|---|---|
| ADC | Analog to Digital Conversion |
| AMGC | Automatic Music Genre Classificaiton |
| ANN | Artificial neural network |
| API | Application Programming Interface |
| BPM | Beats Per Minute |
| DAG | Directed Acyclic Graph |
| DCT | Discrete Cosine Transform |
| DSP | Digital Signal Processing |
| DWCH | Daubechies Wavelet Coefficient Histogram |
| EM | Expectation Maximization |
| FMA | Free Music Archive |
| FFT | Fast Fourier Transform |
| FT | Fourier Transform |
| GMM | Gaussian Mixture Model |
| GTZAN | George Tzanetakis |
| GUI | Graphic User Interface |
| ILR | Intermediate Level Requirement |
| ISMIR | International Symposium on Music Information Retrieval |
| JIT | Just In Time |
| JVM | Java Virtual Machine |
| KL | Kullback-Leibler |
| KNN | K-Nearest Neighbour |
| MER | Music Emotion Recognition |
| MFCC | Mel-frequency cepstral coefficients |
| MIDI | Musical Instrument Digital Interface |
| MIR | Music Information Retrieval |

| | |
|---|---|
| MPEG | Moving Picture Experts Group |
| MP3 | MPEG Audio Layer 3 |
| RMS | Root Mean Square |
| SMILE | Statistical Machine Intelligence and Learning Engine |
| SONE | Specific Loudness sensation coefficients |
| STFT | Short Time Fourier Transform |
| SVM | Support vector machine |
| WAV | Waveform Audio File Format |

## LIST OF SYMBOLS

$\eta$       Learning rate

$\theta$       Bias

$\Delta$       Difference

$\delta$       Artifical neural network node

$\tau$       Lag at time period t

# 1   INTRODUCTION

## 1.1   Background

Music can be literally defined as the combination of soothing sounds. A more complex definition of music can be, a complex amalgam of melody, harmony, rhythm, timbre and silence in a particular structure. Music if an art form and cultural activity whose medium is sound and silence. It's a form of entertainment that puts sounds together in a way that people like or find interesting. To form a music, requirement of musical instruments are not necessary, for example a cappella, barbershop, choral, scat, plainsong, isicathamiya,etc. A group a person can simply sing in rhythm and form a music. Sometimes musician may use their voice to make noises similar to a musical intrument. Music gives the feeling of relaxation. For some people it momentarily stops the flow of time and for some it is a means of passage of time. We can find many music lovers all around the world. Its seems a bit abnormal if a person has no taste in music at all. Some might say a guitar is necessary to form a song or music, some might say a piano is a must, but some musicians may find music find in the chirping of bird, running water of river or even whistling of train.

The common elements of music are pitch (which governs melody and harmony), rhythm (and its associated concepts tempo, meter, and articulation), dynamics (loudness and softness), and the sonic qualities of timbre and texture (which are sometimes termed the "color" of a musical sound). Different styles or types of music may emphasize, de-emphasize or omit some of these elements. Music is performed with a vast range of instruments and with vocal techniques ranging from singing to rapping, and there are solely intrumental pieces, solely vocal pieces. The creation, performance, significance, and even the definition of music vary according to culture and social context.

We may date the advent of music to be centuries year old. We can point it out due to the fact of presence of tribal music which has been passed on from generation to generation like the ancient African bushman tribal song, Nepalese traditional/cultural song from each race,etc. But however music can be complex. Though some may present a pattern in themselves like

a chorus or say rift, some may have an uneven flow. Its perplex nature can not only due to its origination but also due to the evolution of music in different technological era. We have seen the evolution of music from legendary classical Beethoven symphonies to modern day hiphop which is widely popular among the youths nowadays. We have seen the rise of different genres like classical, rock, pop, metal,etc. and yet we may even don't know many of others at all and more may be yet too come like lately we have seen the rise of techno music.

It is not likely for a single person to listen to each and every genre present out there let alone all those songs. Every person may acquire a different taste in music. Some may like clasical music while come may like rock music, it's based on their choice. So a person may probably only distinguish a particular unknown song if those song were to belong to his/her genre of choice and the same goes for the mood. So realising this problem, there has been a increasing amount of research and work done in the music sector for the automatic classification of song based on genre. Though classification based on genre has been a popular one, classification based on mood catching the sight of many people lately.

With the advent of networks and internet, the number of songs are increasing exponentially throughout the internet. Website like soundcloud, youtube, facebook,etc. has given a platform for people to pursue their interest in music by forming like groups, composing and releasing songs, sharing songs,etc. Internet has made it possible for worldwide connection of whole world and it has harnessed the music industry. Because of internet only the popularity of music artists has been increasing all around the world. Their music have now been able to reach each and every corner of the world. This freedom of music throughout the internet has lead to increasing amount of songs and their databases. Due to these rapid development in the music industry, there has been an increasing amount of work in the area of automatic genre classification of music in audio format. A serious factor behind this automation can be considered as the increasing number of millions of records by different artist every year. A simple automation in classification would be much suited than a hand-to-hand task by human and its applicability can be huge. Moreover, there might be conflict regarding genre and mood issue based on the perception of a human being. So regarding these issues MIR(Music Information Retrieval) has primarily focused on automation of classification of

such music based on signal analysis. Such systems can be used as a way to evaluate features describing music content as well as a way to structure large collections of music.

## 1.2 Overview

Throughout the evolution of music, the music industry took a different path and the difference in nature, flow of music, it's tempo, etc. is huge and quite complicated. It lead to the evolution of different numerous genre. The presence of numerous genre is a source of confusion and more often than not people are overwhelmed with the sheer vastness of music available. We humans can most of the times easily categorize simple song based on genre or mood by simply listening and analysing few sample of similar song based on similarity but we are never truly able to understand its nature or features distinction. So we can sometimes never able to recognize them correctly in case of genre and mood. There are songs out there for example Bohemian Rhapsody, which we can never really point it out to a distinct genre and mood.

Moreover the advent of internet has escalated the popularity of music and various artists. Nowadays everyone wants to be a singer. They want to become famous. So, various sites like soundcloud, youtube, facebook,etc. has provided them the perfect platform for sharing their songs. Not only for some novice singers but also for whole popular artist and whole music industry it has provided the perfect platform for sharing the music and growing itself. This has lead to release of millions of songs and increase in database of the system. So given the today world in computerized technological era, automation is nowadays seen as a popular subject in every field. There is being development of automation in every field like riding cars, manufacturing factories, etc. This popularity has affected the music industry too. Realising the potential of its applicability, there has been number of research in this sector/field. Numbers of research paper are being published regarding the automation of classification of music with research paper [29] published by George Tzanetakis and Perry Cook being one of the first in this field with the primary motivation to make it easier for people to classify music (based on genre and/or mood) so that they can find songs suited to their own tastes. It can also lay the foundation for figuring out ways to represent similarity between two musical

pieces and in the making of a good recommendation system.

Given the perplexing nature of music, music classification requires specialized representations, abstraction and processing techniques for effective analysis, evaluation and classification that are fundamentally different form those used for other mediums and tasks. So focusing on these issues we created a music classification system which is web based application used for classifying music. We did not limit ourselves only to genre which is the burning issue in the music industry but we made our effort for the music classification based on mood too. In music industry there is a vast number of different genre. Most of the previous work were limited to four different genres. So, to challenge ourselves we took five different genres for our classification system, namely:-

- Classical

- Jazz

- Rock

- Pop

- Hiphop

Our application took a song as an input from the user computer and classified to its genre based on feature extracted and learning of the system.

Similar procedure was taken for classification of song based on mood. Until now not much research were done on music classification based on mood. So we made our classification system to classify that same song based on mood which was truly based on signal analysis and not lyrical features. For classification based on mood, we mapped the song among two dimensions:

- Energy

- Stress

So based upon the energy and stress level, our song is classified as:

- High Energy, High Stress = Anxious/Frantic

- High Energy, Low Stress = Exuberance

- Low Energy, High Stress = Depression

- Low Energy, Low Stress = Contentment

So, we can say that our system first extracted the required features based on the signal analysis and it's manipulation, and then used those features to classify it among one of the combination of five different genres and 4 different mood using the machine learning algorithm which is already trained on dataset.

## 1.3   Problem Statement

The evolution of music and its origination has presented us with many different genres. The advent and popularity of internet and networking has escalated the market and rise of music industry. Given the popularity of music industry, thousands of new artist are emerging every year. People are releasing song everytime as their hooby or part-time career. So we can see there are millions of songs out there world wide and is continuously increasing every year. Internet has huge contribution for it's rise. With that much of released songs, the size of database is also increasing every year. Since the subject of genre and mood depends on people's perception, it has really been a tedious job to create a quite standard one.

So we built a music classification system based on genre and mood. The choice of these genres is based on their being sufficiently distinguishable from each other. Choosing some genre thats very unique and abnormal might have made them more distinguishable and easier to classify but it would have been harder to find quality data/works for those genre. So, we chose these genre with availability of musical pieces in mind too. We chose to work on classification based on mood too because not many work had been done in the past regarding this field. But we can see this field has a wide scope of applicability. It can be used as a song recommendation system based on genre and that typical mood which the user is listening too as it is certain that the user will possibly like similar song with the similar melody. For now we are currently trying to tackle the issue of music classification based on genre and mood and not abiding to its applications.

## 1.4   Motivation

The presence of numerous number of different genres has presented tedious job for music industry. It has become a source of confusion and more often than not people are overwhelmed with the sheer vastness of music available. So, the primary motivation is to make it easier for people to classify music based on genre. Not only genre but classification based on mood has also now intrigued many people. Combined these two will provide or make a solid foundation for figuring out ways to represent similarity between two musical pieces and build a good recommendation system for music lovers who are passionate about their music and also their choices.

It can futher tackle the issue of automated music database management with large number. It can especially be useful in those cases with unknown label-genre and mood. Music player developers can then be able to make a smart playlist based on the genre and mood of some samples of song the user was currently or recently listening to. This would save a lot of time of user who had to otherwise manually maintain his/her playlist everytime based on his current mood and genre of choice.

## 1.5   Aims and Objectives

- To study and implement different preprocessing steps involved in extracting features from audio data.

- To implement suitable classification algorithm for various features of the song.

- To cross validate the result and analyze the efficiency of the algorithms used.

- To extend the compatibility of the system with different types of music formats like wav,au,etc. along with mp3 format.

- To create a web based application for music classification based on genre and mood.

## 1.6 Scope of Project

(i) The project will work on classifying music based on genre and mood. More specifcally, the classifcation will be done on western music only as the data is more easily available and lots of works have been done in the past for it. Also, only five genres will be used for genre classifcation:

- Rock

- Pop

- Classical

- Jazz

- Hiphop

(ii) The mood based classifcation will use the Thayer model, a two dimensional model based on Energy and Stress:

- High Energy, High Stress = Anxious/Frantic

- High Energy, Low Stress = Exuberance

- Low Energy, High Stress = Depression

- Low Energy, Low Stress = Contentment

(iii) Also, it is entirely possible for a song or a piece of music to fall into multiple genre or moods. The characteristics that defne the genre and the mood may change within the song itself with one part showing seeming to belong to one class while other parts may seem to belong to an entirely diferent class. The project will not cover such issues. In other words, multiple-tagging will not be done.

(iv) The classification will work on various music file formats like mp3, au, wav,etc.

## 1.7 Organization of Report

This report describes and details the design and methodology of building a music classification system based on genre and method. As this report consists documentations relating to

different field during development of a standard software product, hence the whole report is effectiely broke down to 9 chapters.

Chapter one is intended to introduce the project by simply presenting a brief background of the project field which is music and music industry, the motivation which drove us to pursue the field, the overview of the problem statement and objective of the project and at last the scope of our project.

Chapter two presents the literature review. It provides us the collective effort that has been done in the past in our project field. Since our project is music classification based on genre and mood, so at first we start by brief history of Music Information Retrieval(MIR) and music classification. We give a general review of past activities and research on music classification based on both genre and mood. We describe the procedures involved in and the quality of the datasets that we have acquired for the project/system. We analyze the different features involved in the classification. We try to distinguish and analyze the most prominent ones which have been mostly used throughout the time period until now in all research. ALong with the features we also try to we analyze different types of classification algorithms involved in it.

Chapter three describes the theoretical background. In it, we explain about the different selected features involved in our system. We also explain about the working details about the various classification algorithms involved in our system. We also describe about the testing procedures and validation mechanism involved in our system. So to be exact we explain about the cross-validation procedure and all the measures of performance done like precision, recall, fmeasure,etc.

Chapter four is all about the system analysis done at perspective of software engineering. It describe about the requirement specification which is high level requirement, functional requirement and non functional requirement. It also involves feasibility assessment which contains operational feasibility, technical feasibility and economic feasibility.

Chapter five involves system design. First there is overview of whole system design, then we

describe about the system and its various components. It is then followed by a series of use case diagram, component diagram, activity diagram and sequence diagram.

Chapter six describes about the system development which means all the methodology involved like data pre-processing and work-flow. We describe about different tools and environment involved. We also list out all the problem faced during the entire system development and the way to tackle them.

Chapter seven involves the result and analysis process. Since our music classification is based on genre and mood, so we analyze the accuracy involved in each with each feature involved and also with different classifiers involved and we present our perception based on the result. After that there is description of final product which is the finalized features and models involved and user interface created.

Finally, in chapter eight we present our conclusion. We present our view based on the result and analysis and give our insights on future enhancement of the system.

Along with all these there are list of references and bibliography relating to project which is included at last. There is also appendix provided which gives all the analysis and design diagrams which have been developed during the project.

# 2  LITERATURE REVIEW

## 2.1  Human Audio Perception

The human ear is an exceedingly complex organ. To make matters even more difficult, the information from two ears is combined in a perplexing neural network, the human brain. [1]

*Figure 1* illustrates the major structures and processes that comprise the human ear. The outer ear is composed of two parts, the visible flap of skin and cartilage attached to the side of the head, and the ear canal, a tubeabout 0.5 cm in diameter extending about three cm into the head. These structures direct environmental sounds to the sensitive middle and inner ear organs located safely inside of the skull bones. Stretched across the end of the ear canal is a thin sheet of tissue called the tympanic membrane or eardrum. Sound waves striking the tympanic membrane cause it to vibrate. The middle ear is a set of small bones that transfer this vibration to the cochlea (inner ear) where it is converted to neural impulses. The cochlea is a liquid filled tube roughly two mm in diameter and three cm in length.

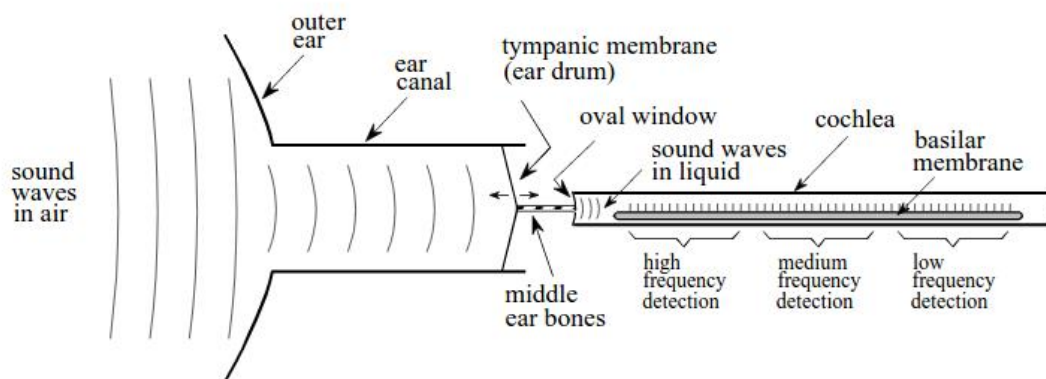Figure 2.1: Functional Diagram of Human Ear

Music can be defined as organised sound comprising the following structural elements: pitch, timbre, key, harmony, loudness (or amplitude), rhythm, meter, and tempo. Processing these elements involves almost every region of the brain and nearly every neural subsystem.

Sound does not exist outside of the brain; it is simply air molecules moving. Sound is

produced by vibrating air molecules connecting with the eardrum at varying frequencies (pitch) and velocities (amplitude). The process starts with the brains primary auditory cortex receiving a signal from the eardrum/inner ear which immediately activates our primitive brain, the cerebellum. The cerebellum is the oldest part of the brain in evolutionarily terms and plays an important part in motor control. It contributes to coordination, precision, and accurate timing of movements. The ear and the primitive brain are known collectively as the low-level processing units. They perform the main feature extraction which allows the brain to start analysing the sounds, breaking down the sensory stimulus into pitch, timbre, spatial location, amplitude, reverberant environment, tone durations, and onset times of different notes.

This data is conducted through neurons in the brain; cells specialized in transmitting information, and the basic building blocks of the nervous system. The output of these neurons connects to the high-level processing units located in the frontal lobe of the brain. It is important to note that this process is not linear. The different regions of the brain constantly update each other with new information.

## 2.2    History of MIR and Music Classification

The field of Music Information Retrieval (MIR) can be traced back to the 60s with reference to the works done by Kassler in [2]. Even Automatic Transcription of Music was attempted as early as the 70s [3]. However, there were two limiting factors that prevented progress in the field at the time. Firstly, the high computational requirements of the problem domain was simply not available. And secondly, other related fields of study such as Digital Signal Processing, Speech Processing, and Machine Learning were also not advanced enough. So, the field stalled for the next few decades.

In the 1990s, the field regained prominence as computational resources improved greatly and the rise of the internet resulted in massive online music collection. So, there was both an opportunity and demand for MIR systems. The organization of the first International Symposium on Music Information Retrieval (ISMIR 1) in 2000 highlights this resurgence of

interest in the field. 280 people from 25 different countries participated in ISMIR Conference Malaga 2015.

As for the methodologies used, MIR in the 90s was influenced by the field of Text Information Retrieval (IR), techniques for searching and retrieving text documents based on user queries. So, most of the algorithms were developed based on symbolic representations such as MIDI files [4]. One such method is described in [5].

However, as mentioned in [6], identifying approximate units of meaning in MIR, as done by the majority of text-IR methods (words serve as such units) was not easy.

Instead, statistical non-transcriptive approaches for non-speech audio signals started being adopted in the second half on the 90s [4]. This was probably influenced by progress of such methods in other fields of speech processing. For example, in [7], the authors reported 98% accuracy in distinguishing music from speech in commercial radio broadcasts. This was based on the statistics of the energy contour and the zero-crossing rate.

In [8], the authors introduced similar statistical methods for retrieval and classification of isolated sounds. Similarly, in [9], an algorithm for music-speech classification based on spectral feature was introduced. It was trained using supervised learning.

And so, starting in the 2000s, instead of methods attempting note-level transcriptions, researchers focused on direct extraction of information of audio signals using Signal Processing and Machine Learning techniques.

Currently, three basic strategies are being applied in MIR: [10]

- **Based on Conceptual Metadata** - Suited for low-specificity queries.

- **Using High-level Descriptions** - Suited for mid-specificity queries.

- **Using Low-level Signal-based Properties** - Used for all specificities.

But still most of the MIR techniques being employed at present use low-level signal features instead of high-level descriptors [11]. Thus, there exists a semantic gap between human perception of music and how MIR systems work.

## 2.3    Audio Processing

General Audio signal processing is an engineering field that focuses on the computational methods for intentionally altering sounds, methods that are used in many musical applications.

Particularly speaking, music signal processing may appear to be the junior relation of the large and mature field of speech signal processing, not least because many techniques and representations originally developed for speech have been applied to music, often with good results. However, music signals possess specific acoustic and structural characteristics that distinguish them from spoken language or other nonmusical signals. [20]

In music the most important qualities of sound are: pitch, duration, loudness, and timbre. Duration and loudness are unidimensional, while pitch and timbre are complex and multidimensional. [21]

- **Loudness** - Intensity of a tone is the physical correlate that underlies the perception of loudness. Loudness variations play an important role in music, but are less important than pitch variations.

- **Duration** - A composer or performer can alter the pace of a piece so that its apparent (virtual) time is slower or faster than clock time.

- **Timbre** - Timbre is the subjective code of the sound source or of its meaning. According to the American Standards Association, "Timbre is that attribute of auditory senstation of which a listener can judge that two steady-state tones having the same pitch and loudness are dissimilar."

- **Pitch** - Pitch is related to the frequency of a pure tone and to the fundamental frequency of a complex tone. In its musical sense, pitch has a range of about 20 to 5000 Hz. Some

five to Zseven harmonics of a complex tone can be heard out individually by paying close attention. There is a dominance region for pitch perception, roughly from 500 to 2000 or 3000 Hz. Harmonics falling in the dominance region are most influential with regard to pitch.

Again, these types of low dimensional features extracted from the acoustical signals are more popular than higher dimensional representations such as Spectrograms for Classification purposes. [22]

## 2.4 Theoretical Background

### 2.4.1 Features

2.4.1.1 Compactness. It is the measure of the noisiness of a signal. It is found by comparing the components of a windows magnitude spectrum with the magnitude spectrum of its neighbouring windows.

If M[n], M[n-1] and M[n+1] > 0, then

$$compactness = \sum_{n=2}^{N-1} ((|20*log(M[n])) - 20*(log(M[n-1]) + log(M[n]) + log(M[]n+1))/3|)$$
(1)

otherwise,

$$compactness = 0 \tag{2}$$

where M[n] is the Magnitude Spectrum at internal n.

2.4.1.2 Mel-Frequency Cepstral Coefficients. It 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.

**Algorithm**

(i) **Framing:** The process of segmenting the speech samples obtained from analog to digital conversion(ADC) into a small frame with the length within the range of 20 to

40 msec. The voice signal is divided into frames of Nsamples. Adjacent frames are being separated by M(M < N).

(ii) **Hamming Window:** Hamming window is used as window shape by considering the next block in feature extraction processing chain and integrates all the closest frequency lines. The Hamming window equation is given as:

If the window is defines as W(n), $0 \leq n \leq$ N-1 where

N = number of samples in each frame

Y[n] = Output signal X(n) = Input signal W(n) = Hamming window,

then the result of windowing signal is shown below:

$$Y(n) = X(n) \times W(n) \tag{3}$$

$$W(n) = 0.54 - 0.46cos\left(\frac{2\pi n}{N-1}\right), \qquad 0 \leq n \leq N-1 \tag{4}$$

(iii) **Fast Fourier Transform:** To convert each frame of N samples from time domain into frequency domain. The Fourier Transform is to convert the convolution of the glottal pulse U[n] and the vocal tract impulse response H[n] in the time domain. This statement supports the equation below:

$$Y(w) = FFT[h(t) * X(t)] = H(w) * X(w) \tag{5}$$

If X(w), H(w) and Y(w) are the Fourier Transform of X(t), H(t) and Y(t) respectively.

(iv) **Mel Filter Bank Processing:** The frequencies range in FFT spectrum is very wide and voice signal does not follow the linear scale as shown in figure 2.2 is then performed. This figure shows a set of triangular filters that are used to compute a weights sum of filter spectral components so that the output of process approximates to a Mel scale. Each filters magnitude frequency re- sponse is triangular in shape and equal to unity at the centre frequency and decrease linearly to zero at centre frequency of two adjacent filters [7,8]. Then, each filter output is the sum of its filtered spectral components. After that the following equation is used to compute the Mel for given frequency f in Hz: This figure shows a set of triangular filters that are used to compute a weights sum of filter spectral components so that the output of process approximates to a Mel scale.
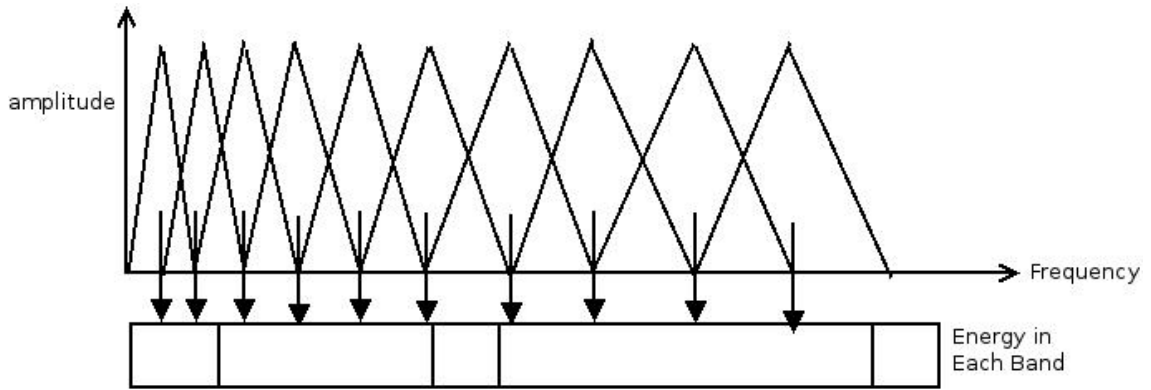
Figure 2.2: Mel scale filter bank

Each filter's magnitude frequency response is triangular in shape and equal to unity at the centre frequency and decrease linearly to zero at centre frequency of two adjacent filters [7,8]. Then, each filter output is the sum of its filtered spectral components. After that the following equation is used to compute the Mel for given frequency f in Hz:

$$M(f) = 1125 ln\left(1 + \frac{f}{700}\right) \tag{6}$$

(v) **Discrete Cosine Transform:** This is the process to convert the log Mel spectrum into time domain using Discrete Cosine Transform(DCT). The result of the conversion is called Mel Frequency Cepstrum Coefficient. The set of coefficient is called acoustic vectors. Therefore, each input utterance is transformed into a sequence of acoustic vector.

2.4.1.3   Pitch.    It is a perceptual property of sounds that allows their ordering on a frequencyrelated scale, or more commonly, pitch is the quality that makes it possible to judge sounds as "higher" and "lower" in the sense associated with musical melodies.

It is a subjective psychoacoustical attribute of sound, and hence is approximately quantified as fundamental frequency.Pitch is an auditory sensation in which a listener assigns musical tones to relative positions on a musical scale based primarily on their perception of the frequency of vibration.Pitch is closely related to frequency, but the two are not equivalent. Frequency is an objective, scientific attribute that can be measured. Pitch is each person's subjective perception of a sound wave, which cannot be directly measured. However, this

does not necessarily mean that most people won't agree on which notes are higher and lower.

**Algorithm**

(i) Model the signal $x_t$ as a periodic function with period T, by definition invariant for a time shift of T

$$x_t - x_{t+T} = 0, \qquad \forall t \tag{7}$$

The same is true after taking the square and averaging over a window

$$\sum_{d=t+1}^{t+W} (x(j) - x(j+\tau))^2 = 0 \tag{8}$$

Conversely, an unknown period may be found by forming a difference function and searching for the values of $\tau$ for which the function is zero.

(ii) The cumulative mean normalized difference function is calculated by dividing each value of the old by its average over shorter-lag values. It differs from difference function in the first step in that it starts at 1 rather than 0, tends to remain large at low lags, and drops below 1 only where the first difference function falls below average.

(iii) Set an absolute threshold and choose the smallest value of $\tau$ that gives a minimum of the difference function obtained in the second step deeper than that threshold. If none is found, the global minimum is chosen instead. With a threshold of 0.1, the error rate drops significantly.

(iv) Each local minimum of the second difference function and its immediate neighbors is fit by a parabola, and the ordinate of the interpolated minimum is used in the dip-selection process. The abscissa of the selected minimum then serves as a period estimate. Actually, one finds that the estimate obtained in this way is slightly biased. To avoid this bias, the abscissa of the corresponding minimum of the raw difference function(the first) is used instead.

## 2.4.1.4 Tempo.

The beat is the regularly occurring pattern of rhythmic stresses in music. When we count, tap or clap along with music we are experiencing the beat. Try tapping your finger along with different types of music and see what happens.

Tempo is the speed of the beat, usually expressed in Beats Per Minute(BPM). For example, at 120 BPM there will be 120 beats in one minute. Tempo can also be expressed verbally with different music terms, such as Slowly, Fast, Allegro, or Largo.

**Algorithm**

(i) Parse an audio file into samples, s[n], with a corresponding sampling rate, SR.

(ii) Break the audio file into N windows of 1024 samples.

(iii) Calculate the FFT of each window.

(iv) Calculate the power spectrum(P) of each window form the corresponding FFT results.

(v) Calculate the spectral flux(F) from the power spectrum(P) of each window, i:

$$F_i = (P_i - P_{i-1})^2 \qquad (9)$$

(vi) Find the mean flux, $F_{av}$, across all windows.

(vii) Set $F_i$, the flux for each window, to 0 if it is not at least 1.5 times $F_{av}$ (this value was experimentally determined). This gives a generous estimation of note onsets.

(viii) Use autocorrelation to find the histogram of lag frequencies, L:

$$L[lag] = \sum_i^N F[i]F[i+lag] \qquad (10)$$

(ix) Calculate the effective sampling rate, $S_{eff}$, found in F and L.

$$SR_{eff} = \frac{SR}{1024} \qquad (11)$$

(x) Convert the lag histogram L[lag] into a tempo histogram L[BPM] with bins of beats per minute by reversing the order of L and converting the bin lag indices, lag, to BPM:

$$L[BPM] = L\left(\frac{60 * SR_{eff}}{lag}\right) \qquad (12)$$

(xi) The result of step (x), L[BPM], is a tempo histogram with bin labels corresponding to beats per minute and bin frequencies correspoding to frequencies of inter-peak intervals.

### 2.4.1.5 Root Mean Square.

The root mean square(abbreviated RMS of rms) is defined as the square root of mean square(the arithmetic mean of the squares of a set of numbers).

$$x_{rms} = \sqrt{\frac{1}{n}(x_1^2 + x_1^2 + ... + x_n^2)} \tag{13}$$

where there is set of n values $\{x_1, x_2, ......, x_n\}$.

### 2.4.1.6 Spectral Centroid.

The spectral centroid is defined as the center of gravity of the magnitude spectrum of the STFT.

$$C_t = \frac{\sum\limits_{n=1}^{N} M_t[n] * n}{\sum\limits_{n=1}^{N} M_t[n]} \tag{14}$$

where $M_t[n]$ is the magnitude of the Fourier transform at frame t and frequency bin n. The centroid is a measure of spectral shape and higher centroid values correspond to "brighter" textures with more high frequencies.

### 2.4.1.7 Spectral Flux.

The spectral flux is defined as the squared difference between the normalized magnitudes of successive spectral distributions.

$$F_t = \sum_{n=1}^{N} (N_t[n] - N_{t-1}[n])^2 \tag{15}$$

where $N_t[n]$ and $N_{t-1}[n]$ are the normalized magnitude of the Fourier transform at the current frame t, and the previous frame t-1, respectively.

The spectral flux is a measure of the amount of local spectral change.

### 2.4.1.8 Spectral Roll-off Point.

The spectral roll-off is defined as the frequency $R_t$ below which 85 per cent of the magnitude distribution is concentrated.

$$\sum_{n=1}^{R_t} M_t[n] = 0.85 * \sum_{n=1}^{N} M_t[n] \tag{16}$$

The roll-off is another measure of spectral shape.

### 2.4.1.9   Spectral Variability.   Spectral variability is the standard deviation of the magnitude spectrum. This gives the measure of the variance of a signal's magnitude spectrum.

$$Spectral variability = \sqrt{\frac{1}{N-1} - (\sum_{n=1}^{N} M[n] - mean)^2} \tag{17}$$

where M[n] = Magnitude spectrum of the signal at interval n, Mean = mean of the magnitude spectrum.

### 2.4.1.10   Zero Crossing.   The zero crossing is defined as the number of times the waveform changed sign.

$$Z_t = \frac{1}{2} \sum_{n=1}^{N} |sign(x[n]) - sign(x[n-1])| \tag{18}$$

where the sign fucntion is 1 for positive arguments and 0 for negative arguments and x[n] is the time domain signal for frame t.

Time domain zero crossings provide a measure of the noisiness of the signal.

## 2.4.2   Classifier

In machine learning and statistics, classification is the problem of identifying to which of a set of categories(sub-populations) a new observation belongs, on the basis of a training set of data containing observations(or instances) whose category membership is known. In the terminology of machine learning, classification is considered an instance of supervised learning where a training set of correctly identified observations is available. The corresponding unsupervised procedure is known as clustering, and involves grouping data into categories based on some measure of inherent similarity or distance. An algorithm that implements classification,especially in a concrete implementation, is known as a classifier.

A variety of methods have been used for music classification. Some of the popular ones are SVM, K-means, K-nearest neighbours and variants of Neural Networks. The results are also widely different. In [23] 61 per cent accuracy has been achieved using a Multilayer Perceptron based approach while in [16], the authors have managed 95 per cent (for Back Propagation Neural Network) and 83 per cent (for SVM). In [24], the authors have achieved 71 per cent accuracy using an additional rejection and verification stage.

In [25], simpler and more naive approaches (k-NN and k-Means); and more sophisticated neural networks and SVMs have been compared. The author found the latter gave better performance.

However, lots of unique methods either completely novel or a variation of a standard method have been put into use too. In [26], the authors propose a method that uses Chord labeling (ones and zeros) in conjunction with a k-windowSubsequenceMatching algorithm used to find subsequence in music sequence and a Decision tree for the actual genre classifi- cation.

It is also noted that high-level and contextual concepts can be as important as low-level content descriptors[19].

After these all research, we decided to go with K-means, artificial neural network based on backpropagation and support vector machine. Artificial neural network and support vector machine were chosen as they appeared to be famous in the field. Our choice for K-means was based on the fact that it was simple yet powerful unsupervised procedure.

2.4.2.1   K-means Clustering.   K-means is one of the most popular algorithm used for clustering of a given data sets. It is a method of vector quantization, originally from signal processing, that is popular for cluster analysis in data mining. It aims at partition n obsevations into k clusters in which each observation belongs to the cluster with the nearest mean, serving as a prototype of the cluster. This results in a partitioning of the data space into Voronoi cells. In other words, it clusters all the data points which resemble close to each other or we can that it cluster all those data points together which have the lowest cost among themselves based on the distance metric. The problem is computationally difficult(NP-hard); however, there are efficient heuristic algorithms that are commonly employed and converge quickly to a local optimum. It has loose relationship to the k-nearest neighbor classifier.

Regarding computational complexity, finding the optimal solution to the k-means cluster- ing problem for obsevations in d dimensions is NP-hard in genral Euclidean space d even for 2 clusters and NP-hard for a general number of clusters k even in the plane. It k and d

(the dimension) are fixed, the problem can be solved in time $O(n^{dk+1}logn)$, where n is the number of entities to be clustered. The choice for this algorithm was based on our research and some had already implemented it [25]. Most of the research papers has enlisted this algorithm like in [20], [39], [25]. So, our reason for implementation of this algorithm was:

(i) Most of the research papers has shown its use in music classification.

(ii) It is a pretty straightforward and simple algorithm.

(iii) We wanted to have full control over all aspects related to the implementation: the initialization method, distance metric, etc.

The implementation however is basic in the sense that no modification has been done on the algorithm to better suit the problem domain. Some finer points of the implementation are discussed below after quickly going over the algorithm itself.

**Algorithm**

Let X = $x_1, x_2, ...., x_n$ be the set of data points and V = $v_1, v_2, ..., v_c$ be the set of centers.

(i) Randomly select 'c' cluster centers.

(ii) Calculate the distance between each data point and cluster centers.

(iii) Assign the data point to the cluster center whose distance from the cluster center is minimum of all the cluster centers.

(iv) Recalculate the new cluster center:

$$v_i = \frac{1}{c_i} \sum_{j=1}^{c_i} x_i \tag{19}$$

(v) Recalculate the distance between each data point and new obtained cluster centers.

(vi) If no data point was reassigned then stop, otherwise repeat from step 3).

**Initialization method**

Commonly used initialization methods are Forgy and Random Partition. The Forgy method randomly chooses k observations from the data set and uses these as the initial means. The

Random Partition method first randomly assigns a cluster to each observation and then proceeds to the update step, thus computing the initial mean to be the centroid of the clusters randomly assigned points. The Forgy method tends to spread the initial means out, while Random Partition places all of them close to the center of the data set. According to [16], the Random Partition method is generally preferable for algorithms such as the k-harmonic means and fuzzy k-means. For expectation maximization and standard k-means algorithms, the Forgy method of initialization is preferable.

And with these facts in mind, we went for the random initialization method. This method has been used in [20] too although they have also added the constraint that the centroids be separated by at least a threshold KL-divergence distance. As the choice of initial centroids have a drastic effect on the cluster formed, we have also considered other methods of initialization such as the breakup method which uses the actual data points and the scrambled midpoints method which uses synthetic data points as suggested in [22]. So after our research on this, we found that scrambled midpoints to be much more efficient and lead to better clustering. So based on [37], scrambled midpoints technique was chosen for initialization.

In scrambled midpoints initialization method, at first our whole data range or in our case each feature range was broken down into k-paritions which were all equal. Then midpoints were taken of those equal partitions for each feature partition range. Provided there are n-features with k number of clusters needed, then we have n*k different possibilities. Now all these midpoints of each range of n-features were scrambled among each other to form k different initialization points/vectors representing all those n-features. It should be clear that there is no repetition of the midpoint of same range of same feature during scrambling of midpoints to form initialization points/vectors.

**Distance Metric**

Apart from the initialization method, K-means is also highly sensitive to the distance metric used. There are many distance metric but most popular ones are:

- Manhattan distance

- Euclidean distance

- Minkowski distance

Most of the times, K-Means is implicitly based on pairwise Euclidean distances between points, because the sum of squared deviations from centroid (that it tries to minimize) is equal to the sum of pairwise squared Euclidean distances divided by the number of points [23].

As such, we decided to use Euclidean distance with weights for the three different features added to give us a way to control the metric. Thus, the distance between two songs S1 and S2 is given by:

$$Distance(d) = \sqrt[2]{w_I(I_1 - I_2)^2 + w_M \sum_i (M_{1i} - M_{2i})^2 + w_R \sum_i (R_{1j} - R_{2j})^2} \qquad (20)$$

Where: $w_I$ , $w_M$ and $w_R$ are the weights for Intensity, MFCC and Rhythm respectively.

2.4.2.2  Artificial Neural Network.  In machine learning and cognitive science, an artificial neural networks(ANN) is a network inspired by biological neural networks(the central nervous systems of animals, in particular the brain) which are used to estimate or approxiamte functions that can depend on a large numbers of inputs that are generally unknown. Our research shows whether the research is related to our field or other, most of the time for machine learning researchers used artificial neural network. So our choice was simple as it was based on supervised learning and comparatively simple than other complicated unsupervised learning. Research papers [24], [23], ,[19], [25] and [16] shows the implementation of artificial neural network in the field of music classification.

Artificial neural networks are typically specified using three things:

- **Architecture:** It specifies what variables are involved in the network and their topological relationships- for example the variable involved in a neural network might be the weights of the connections between the neurons, along with activities of the neurons. In an artificial neural network, there are one or more hidden layers in between input and output layers with all the neurons connecting to each other.

- **Activity Rule:** Most neural network models have short tiem-scale dynamic: local rules define how the activities of the neurons change in reponse to each other. Typically the

activity rule depends on the weights(the parameters) in the network.Here in our case, the set of input neuronse is activated by the feature values like MFCC, pitch,etc. There activation function present to trigger the respective neuron. Most of the time activation function are non linear, differential mathematical functions like sigmoid, hyperbolic tangent, etc.

- **Learning Rule:** The learning rule specifies the way in which the neural network's weights change with time as the learning takes progress. This learning is usually viewed as taking place on a longer time place than time scale of the dynamicx under the activity rule. Usually the learning rule will depend on the activites of the neurons. In our case the learning depends on the values of the target values supplied by a teacher and on the current value of the weights.
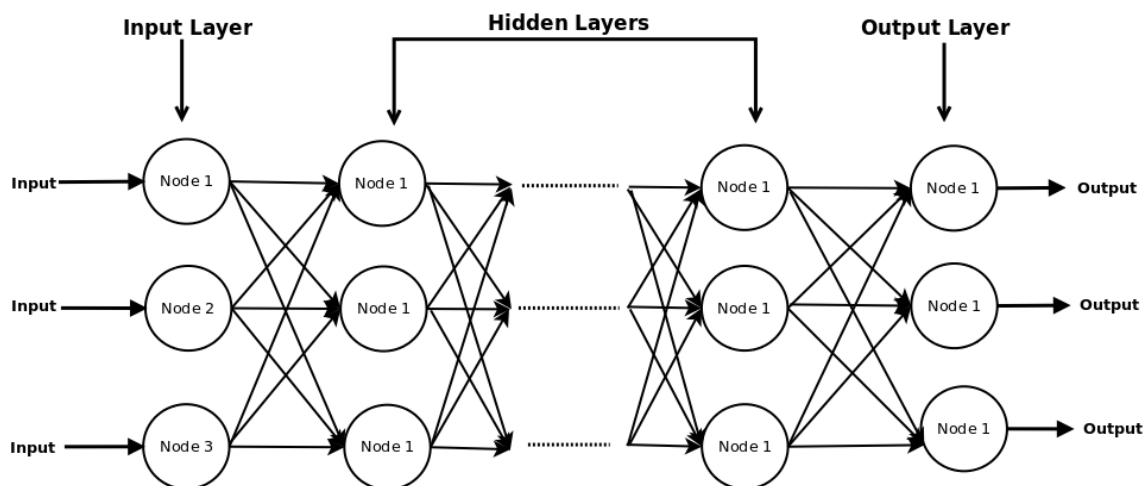


Figure 2.3: General structure of artificial neural network

For a system, generally the architecture is created as per the need. Generally for that there is lot of trial and error methodology involved to determine the correct number of hidden layers and neurons needed. In most cases, one hidden layer is suffiecient for the system but if the nature of system or data is perplexing then in accordance to expected result and current behavior, number of hidden layers of nodes can be increased.

The concept of cost function is an important in context of artificial neural network. It measures how far away a particular solution is from an optimal solution. We can also say that the cost of the optimal solution is minimum. So, the target of our artifical neural network

is to try to meet the cost of optimal solution. While it is possible to define some arbitrary ad hoc cost fucntion, frequently a particular cost will be used, either because it has desirable properties (such as convexity) or because it arises naturally from a particular formulation of problem. Ultimately, the cost fucntion will depend on the desired task. One of the mostly used cost function is squared error measure between the output value O and the target value t.

$$E = (t - y)^2 \qquad (21)$$

where E is the discrepancy or error.

Using this cost, the neural network tries to adjust it weights in order to minimize the cost function. This exact process is called learning. Supervised learning, unsupervised learning and reinforcement learning are the three major paradigms of learning. We are opting for the supervised learning. The reason for this choice is that unsupervised learning and reinforcement learning are comparatively more complex and also supervised learning have been doing great job in this music classification field based on research papers [23], [25], [16]. Supervised learning requires the correct class label to be provided along with the training dataset for the neural network so that it can adjust it's weight based on the cost function of incorrect/correct class prediction. Backpropagation algorithm is the most popular learning algorithm for neural network out there. The reason for its popularity might be its simplicity in terms of concept and wide applicability. It's also based on supervised learning. When one tries to minimize cost function using gradient descent for the class of neural networks called multilayer perceptrons(MLP), one obtains the common and well-known backpropagation algorithm for training neural networks.

Backpropagation is a common method of training artificial neural networks used in conjunction with an optimization method such as gradient descent. The method calculates the gradient of a loss function with respect to all the weights in the network. The gradient is fed to the optimization method wihich in turn uses it to update the weights, in an attempt to minimize the loss function. It is a generalization of delta rule to multi-layered feedforward networks, made possible by using the chain rule to iteratively compute gradients for each layer. Requirements of backpropagation mehod are:

- It requires a known, desired output for each input value in order to calculate the loss function gradient.

- It requires the activation function used by the artificial neurons/nodes to be differentiable.

**Algorithm**

(i) Run the network forward with your input data to get the network output.

(ii) For each output node compute

$$\delta_k = O_k(1 - O_k)(O_k - t_k) \tag{22}$$

(iii) For each hidden node calculate

$$\delta_j = O_k(1 - O_k) \sum_{k \in K} \delta_k W_{jk} \tag{23}$$

(iv) Update the weights and biases as follows:
Given

$$W = -\eta \delta_l O_{l-1} \tag{24}$$

$$\Delta\theta = -\eta \delta_l \tag{25}$$

apply

$$W + \Delta W \rightarrow W \tag{26}$$

$$\theta + \Delta\theta \rightarrow \theta \tag{27}$$

where i, j and k represents the input layer, hidden layer and ouput layer respectively,

O is the output of a neuron/node,

t is the target value,

K is the total number of output neurons/nodes,

$\Delta$ represent the difference,

l denotes every layer.

$\eta$ is the learning rate which is defined as the ratio(percentage) that influences the speed and quality of learning.The greater the ratio, the faster the neuron trains; the lower the ratio,

the more accurate the training is. $\theta$ is the bias term which is involved in adjusting the shifting of activation function. The sign of the gradient of a weight indicates where the error is increasing, this is why the weight must be updated in the opposite direction. The algorithm above is repeated until performance of the network is satisfactory.

### 2.4.2.3 Support Vector Machine.

In machine learning, support vector machines (SVM) are supervised learning models with associated learning algorithms that analyze data used for classification analysis.They are new statistical learning technique that can be seen as a new method for training classifiers based on polynomial functions, radial basis functions, neural networks, splines or other functions. The popularity of Support Vector machine is huge as lot of researcher papers [24], [19], [25], [16] shows its implementation. Not only in the field of music classification but also on various artificial intelligence field like handwriting recognition, biological and other sciences. It was also able to overthrow general neural network until the advent of deep learning.

The basic support vector machine (SVM) is a binary linear classifier which chooses the hyperplane that represents the largest separation, or margin, between the two classes. So Support vector machines use a hyperplane to create a classifier. If such a hyperplane exists, it is known as the maximummargin hyperplane and the linear classifier it defines is known as a maximum margin classifier. If there exists no hyperplane that can perfectly split the positive and negative instances, the soft margin method will choose a hyperplane that splits the instances as cleanly as possible, while still maximizing the distance to the nearest cleanly split instances. For problems that cannot be linearly separated in the input space, this machine offers a possibility to find a solution by making a non-linear transformation of the original input space into a high dimension feature space where an optimal separating hyper plane can be found. Those separating planes are optimal, which means that a maximal margin classifier with respect to the training data set can be obtained. It is developed by Vladimir Vapnik and co-workers at AT&T Bell Laboratories in 1995.

The nonlinear SVMs are created by applying the kernel trick to maximummargin hyperplanes. The resulting algorithm is formally similar, except that every dot product is replaced by a nonlinear kernel function.
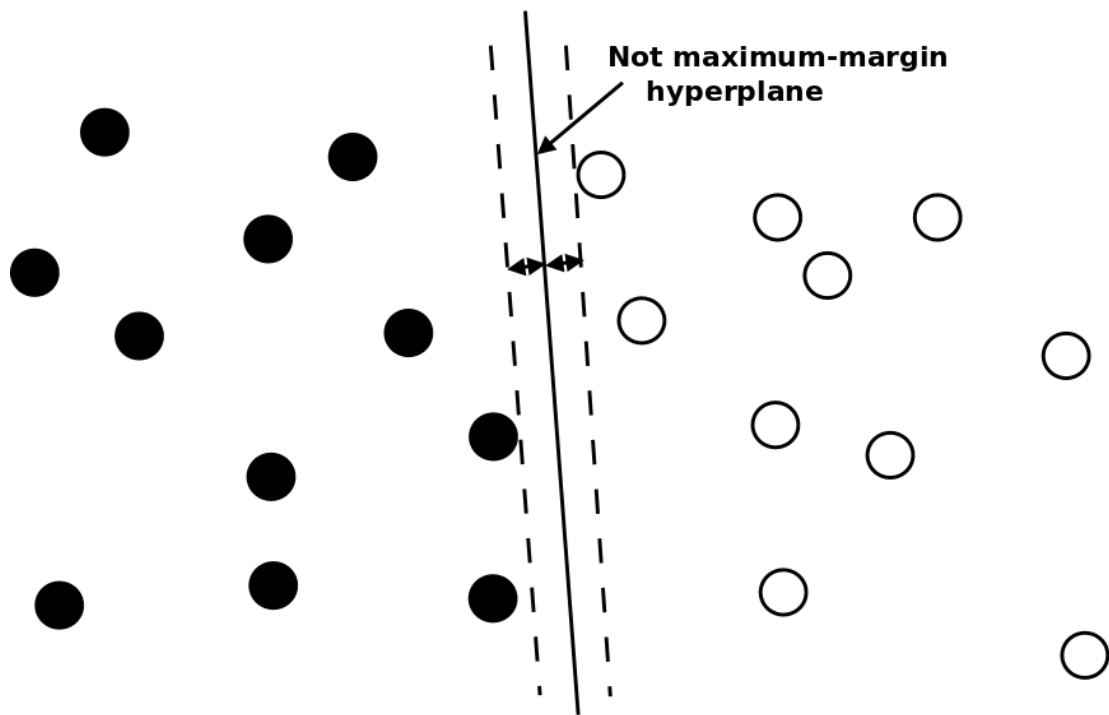
Figure 2.4: Misleading hyperplane

**Kernel Function**

The simplest way to divide two groups is with a straight line, flat plane or an N-dimensional hyper plane. But what if the points are separated by a non-linear region. In such case we would need a non-linear dividing line. Rather than fitting non-linear curves to the data, support vector machine handles this by using a kernel function to map the data into a different space where a hyperplane can be used to do the separation. [19] shows the use of second order polynomial kernel in the support vector machine. So, kernel function allows the algorithm to fit the maximummargin hyperplane in a transformed feature space. The transformation may be nonlinear and the transformed space be high dimensional. For example, the feature space corresponding to Gaussian kernel is a Hilbert space of infinite dimension. Thus though the classifier is a hyperplane in the high dimensional feature space, it may be nonlinear in the original input space. Maximum margin classifiers are well regularized, so the infinite dimension does not spoil the result as the separation will be performed even with very complex boundaries.
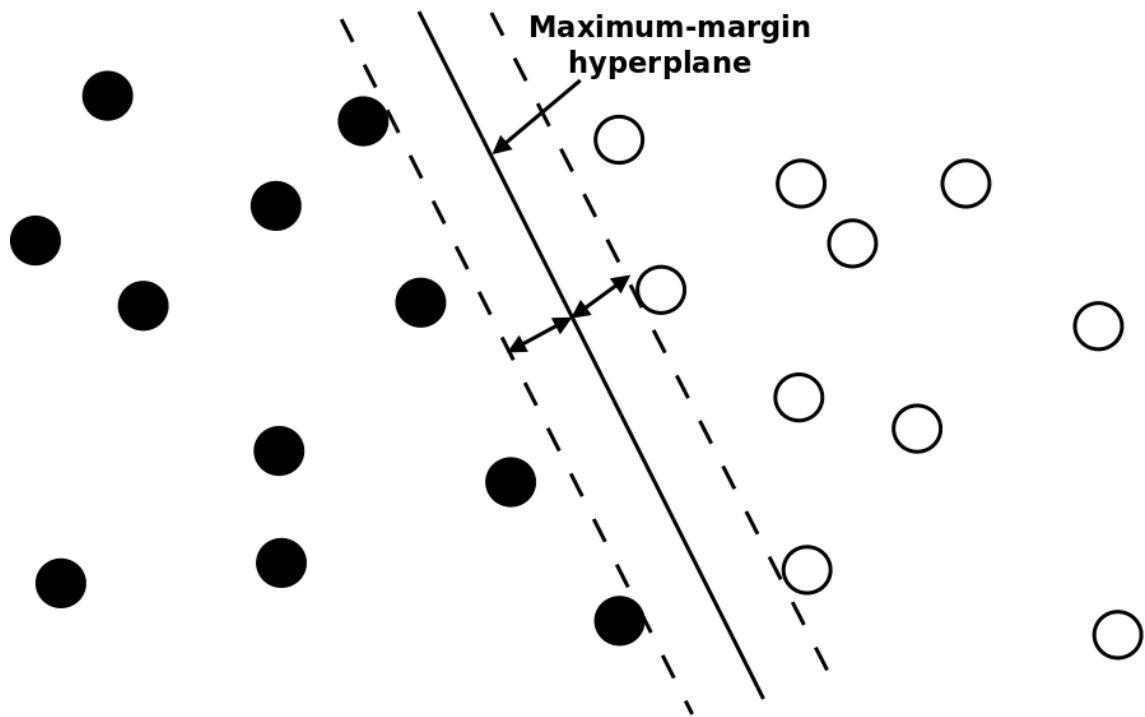
Figure 2.5: Support vector machine with maximum-margin hyperplane

The effectiveness of SVM depends on the selection of kernel, the kernels parameters, and soft margin parameter c. Given a kernel, best combination of c and kernels parameters is often selected by a gridsearch with cross validation. The dominant approach for creating multiclass SVMs is to reduce multiclass problem into multiple binary classification problems. Common methods for such reduction is to build binary classifiers which distinguish between (i) one of the labels to the rest (oneversusall) or (ii) between every pair of classes (oneversusone). Classification of new instances for oneversusall case is done by a winnertakesall strategy, in which the classifier with the highest output function assigns the class. For the one versusone approach, classification is done by a maxwins voting strategy, in which every classifier assigns the instance to one of the two classes, then the vote for the assigned class is increased by one vote, and finally the class with most votes determines the instance classification. To tackle the same multiclass problem [25] has the utilization of DAG(Directed Acyclic Graph) SVMs in which a directed acyclic graph(DAG) of two-class SVMs is trained on each pair of classes in the data set.

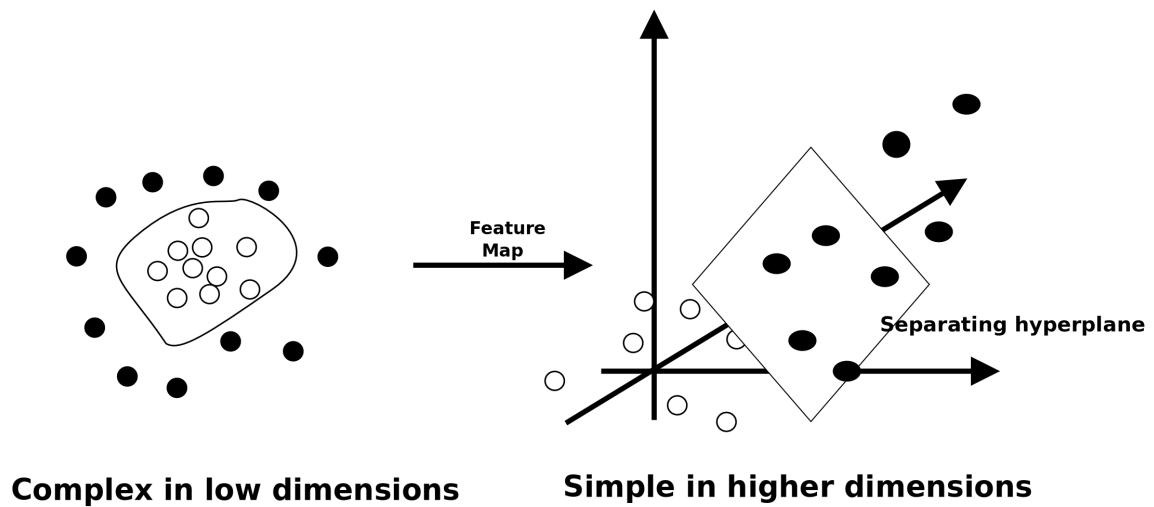**Complex in low dimensions**    **Simple in higher dimensions**

Figure 2.6: Support vector machine with kernel function

## 2.5 Testing and Validation

While building a software product, the concept of testing and validation plays an important role throughout its development. Along with the development of a software product, numerous testing and validation are need to be performed regulary. In fact, most of the development period is given to testing and validation. It's important as it conforms to whether we are building product right and whether we are building the right product. After series of testing and validation only we can point that the system conforms to the desired specifications and need.

Since our software project is also research oriented, so we need series of rigorous testing. Not only this, being a predictive system/model a type of validation and measure for performance is a must, so that it can analyze the result and decide how to progress further accounting the fact. For our system, we chose cross-validation after we confronted it's popularity for a predictive model and ease of use. Similarly for measure of performance we chose recall, presicion and F-measure.

### 2.5.1   Cross-validation

Cross-validation is a technique for estimating the performance of a predictive model. Cross-validation, sometimes called rotation estimation is a model validation technique for assessing how the results of a statistical analysis will generalize to an independent data set. In a predictive problem like our project, a model is usually given a dataset of known data(training dataset) on which the model is trained, and a dataset of unknown data (or first seen data) against which the model is tested (testing dataset). The goal of cross validation is to define a dataset to "test" the model in the training phase(validation dataset), in order to limit problems like overfitting, give an insight on how the model will generalize to an independent dataset(and unknown dataset, for instance from a real problem),etc.

One of the main reasons for using cross-validation instead of using the conventional validation is that there is not enough data available to partition it into separate training and test sets without losing significant modelling or testing capability. There are several cross-validation method out there like leave-p-out cross-validation, leave-one-out cross validation,2-fold cross validaiton, repeated random sub-sampling validation, etc. But we decide to got with k-fold cross validation. The main reason for this selection is that k-fold cross-validation estimator can prove to have a lower variance than a single hold-out set estimator if the amount of data available is limited. Our main goal with k-fold cross validation is to estimate the expected level of fit of a model to a dataset that is independent of the data that were used to train the model.

In k-fold cross-validation, the original sample is randomly partitioned into k equal sized subsamples. Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k  1 subsamples are used as training data. The cross-validation process is then repeated k times (the folds), with each of the k subsamples used exactly once as the validation data. The k results from the folds can then be averaged to produce a single estimation. The advantage of this method over repeated random sub-sampling (see below) is that all observations are used for both training and validation, and each observation is used for validation exactly once. ten-fold cross-validation is currently applied in our system. Our music data set is first divided into ten equal samples of songs, and on every

run for ten times, a different sample of songs are used for testing and the rest for training.

This k-fold validation validation process is applied with each performance metrics and the appropriate features are taken to classify in our system.

## 2.5.2 Measure of Performance

In a predictive problem, we always need a measure of performance for that predictive model along with cross validation. Based on the value of different attributes falling under the measure of performance, we can then have the knowledge of how well the system is performing. It is important because as a matter of fact it provides much deeper insight of the performance of the predictive model, not only accuracy.

Our measure of performance include precision, recall and F-measure. We can say our measure of performance are based on the relevant and non-relevent documents. In this case, relevant documents are those one which simply belong to the relevant category. Relevant category is also called the true positives. True positive is the number of items correctly labeled as belong to the positive class or say which are predicted correctly. The not relevant documents are those one which simply are not relevant to the category or say they are false negative. False negative is the number of negative items which are predicted to fall under the positive class. False Positive is the number of negative items which are predicted to be in true class. Similarly, true negative is the number of negative items which are correctly predicted to fall under negative class.

- **Precision**

  The precision of a measurement system, related to reproducibility and repeatability, is the degree to which repeated measurements under unchanged conditions show the same results. In a classification task, the precision for a class is the number of true positives divided by the total number of elements labeled as belonging to the positive class. A perfect precision score of 1.0 means that every result retrieved by a search was relevant which means falling to the specified class. But precision does not say

whether all items of that class were correctly predicted.

$$Precision = \frac{TruePositive}{TruePositive + FalsePositive}$$ (28)

- **Recall**

  Recall literally is how many of the true positives were recalled (found), i.e. how many of the correct hits were also found. Recall is also sometimes called sensitivity or probability of detection. It is the true positive rate as it measure the proportion of positives that are correctly identified as such. Recall is the number of true positives divided by the total number of elements that actually belong to the positive class. A perfect recall score of 1.0 means that all relevant documents were retrieved by search. We can also say that perfect recall shows that all items of a class were labeled fall under same class. But recall does not say that how many other items were misclassified to fall under same class.

$$Recall = \frac{TruePositive}{TruePositive + FalseNegative}$$ (29)

- **F measure** A measure that combines precision and recall is the harmonic mean of precision and recall. F measure is the weighted harmonic mean of its the precision and recall of a system. It is sometimes called as balanced F-score. This measure is approximately the average of the two when they are close, and is more generally the harmonic mean which for the case of two numbers coincides with the square of the geometric mean divided by the arithmetic mean. There are several reasons that the F-score can be criticized in particular circumstances due to its bias as an evaluation metric. This is also known as the $F_1$ measure, because recall and precision are evenly weighted.

$$F = 2 * \frac{Precision * Recall}{Precision + Recall}$$ (30)

All three performance metrics have been used to determine the best features to be used for classification in cases of genre, arousal and valence.

As depicted in the result and analysis section, we can see that not all the features are suitable for the classification. Some features such as intensity or pitch do not contribute the classification but some features such as Mel Frequency Cepstral Coefficients(MFCC) contribute greatly to the classifiers performance.

## 2.6  Related Works

### 2.6.1  Genre Based Classification

2.6.1.1  <u>Overview.</u>   Automatic Music Genre Classification (AMGC) is one of the tasks focused by MIR. However, it is not a straightforward one.

In [13], Scaringella et al. discuss how and why musical genres are a poorly defined concept making the task of automatic classification non-trivial. Still, although the boundaries between genres are fuzzy and there are no well-defined definitions, it is still one of the widely used method of classification of music. If we look at human capability in genre classification, Perrot et al [14] found that people classified songs–in a ten-way classification setup–with an accuracy of 70% after listening to 3s excerpts.

2.6.1.2  <u>Features.</u>   The features used for genre based classification have been heavily influenced by the related field of speech recognition. For instance, Mel-frequency Cepstral Coefficients (MFCC), a set of perceptually motivated features that is widely used in music classification, was first used in speech recognition.

The seminal paper on musical genre classification by Tzanetakis et al. [12] presented three feature sets for representing timbral texture, rhythmic content and pitch content. With the proposed feature set, they achieved a classification accuracy of 61% for ten musical genre.

Timbral features are usually calculated for every short-time frame of sound based on the Short Time Fourier Transform (STFT). So, these are low-level features. Typical examples are Spectral Centroid, Spectral Rolloff, Spectral Flux, Energy, Zero Crossings, and the aforementioned Mel-Frequency Cepstral Coefficients (MFCCs). Among these, MFCC is the most widely preferred feature [15][16]. Logan [17] investigated the applicability of MFCCs to music modeling and found it to be "at least not harmful".

Rhythmic features capture the recurring pattern of tension and release in music while pitch is the perceived fundamental frequency of the sound. These are usually termed as mid-level

features.

Apart from these, many non-standard features have been proposed in the literature.

Li et al.[18] proposed a new set of features based on Daubechies Wavelet Coefficient Histograms (DWCH), and also presented a comparative study with the features included in the MARSYAS framework. They showed that it significantly increased the accuracy of the classifier.

Anglade, Amlie, et al.[19] propose the use of Harmony as a high-level descriptor of music, focusing on the structure, progression, and relation of chords.

2.6.1.3  Classifer.    A variety of methods have been used for music classification. Some of the popular ones are SVM, K Nearest Neighbours and variants of Neural Networks. The results are also widely different.  In [23], 61 per cent accuracy has been achieved using a Multilayer Perceptron based approach. While in [24], the authors have achieved 71 per cent accuracy through the use of an additional rejection and verification stage. Haggblade et al. [25], compared simpler and more naive approaches (k-NN and k-Means) with more sophisticated neural networks and SVMs. They found that the latter gave better results.

Standard statistical pattern recognition classifiers are also used for AMGC. They may be simple Gaussian Classifiers or Gaussian mixture model (GMM) classifier, where each class pdf is assumed to consist of a mixture of a specific number of multidimensional Gaussian distributions. In such an approach, the parameters of each Gaussian component and the mixture weights are estimated using the iterative EM algorithm.

However, lots of unique methods – either completely novel or a variation of a standard method – have been put into use too. In [26], the authors propose a method that uses Chord labeling (ones and zeros) in conjunction with a k-windowSubsequenceMatching algorithm used to find subsequence in music sequence and a Decision tree for the actual genre classifi-

cation.

It is also noted that high-level and contextual concepts can be as important as low-level content descriptors. [19]

### 2.6.2 Mood Based Classification

<u>2.6.2.1 Overview.</u>  As mood is a very human thing, Mood Based Classification, also known as Mood Emotion Recognition (MER), requires knowledge of both technical aspects as well as the human emotional system. So, the conceptualization of emotion and understanding of the associated emotion taxonomy is vital. However, it is a difficult thing to do, because

(i) It is subjective and

(ii) We cannot agree on a model to depict emotional states.

Usually, two approaces to emotion conceptualization are taken:

- **Categorical Conceptualization** - This approach to MER categorizes emotions into a number of distinct classes. It requires the belief of base emotions (happiness, anger, sadness, etc) from which all other secondary emotion classes can be derived.[30] However, the major drawback of the categorical approach is that the number of primary emotion classes is too small in comparison with the richness of music emotion perceived by humans.

- **Dimensional Conceptualization** - It defines Musical Values as numerical values over a number of emotion dimensions. So, the focus is on distinguishing emotions based on their position on a predefined space. Most of these conceptualizations map to three axes of emotions: valence (pleasantness), arousal (activation) and potency (dominance). By placing emotions on a continuum instead of trying to label them as discrete, this approach can encompass a wide variety of general emotions.

### 2.6.2.2 Circumplex and Thayer Mood Model.

One of the Dimensional conceptualization was proposed by Russell (1980) [31]. As shown in *Figure 2*, the model consists of a two-dimensional structure involving the dimensions of valence and arousal. General emotions are placed within thic circular framework.
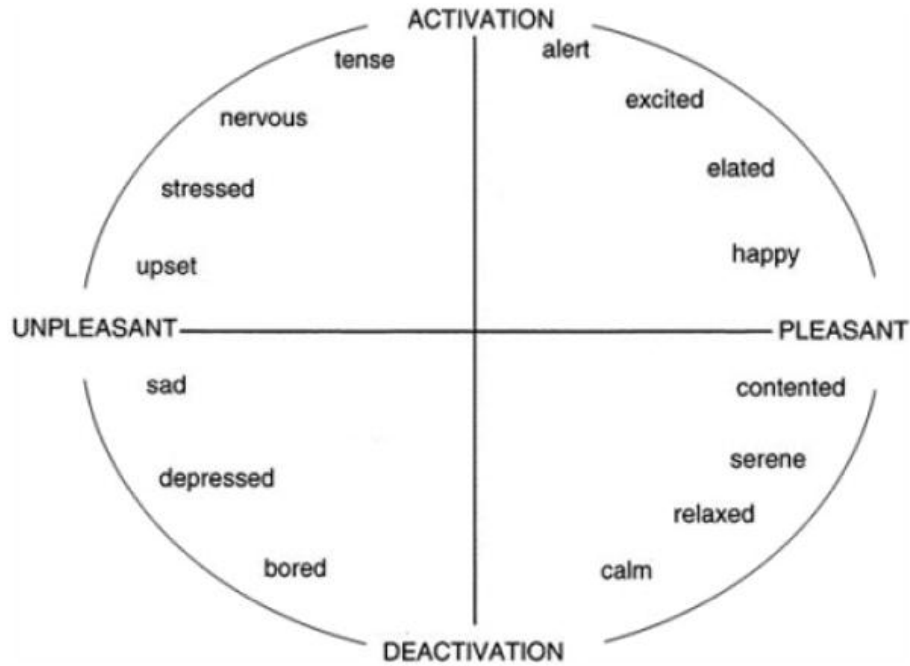


Figure 2.7: A graphical representation of the circumplex model of affect with the horizontal axis representing the valence dimension and the vertical axis representing the arousal or activation dimension.

As shown in *Figure 3*, Thayer [32] proposed a similar two-dimensional approach that adopts the theory that mood is entailed from two factors: -Stress (happy/anxious) -Energy (calm/ energetic). This divides music mood into four clusters: Contentment, Depression, Exuberance and Anxious/Frantic.

Although, the two-dimensional approach has been criticized as deficient (leading to a proposal of the third dimension of potency), it seems to offer the right balance between sufficient "verbosity" and low complexity [33].

### 2.6.2.3 Features.

Some of the commonly used features in MER are:

- **Energy**: Energy related features such as audio power, loudness, specific loudness sensation coefficients (SONE), are correlated to the perception of arousal. Lu et al. [34] used it to classify arousal.

Figure 2.8: Thayer's two-dimensional model of mood

- **Rhythm**: Flowing/fluent rhythm is associated with positive valence while firm rhythms with negative valence.

- **Melody**: These include features such as Pitch (perceived fundamental frequency), chromogram centroid, etc.

- **Timbre**: As with the AMGC problem, MFCC is widely used in MER too. Apart from MFCC, octave-based spectral contrast as well as DWCH (Daubechies wavelets coefficient histogram) are also proposed in literature.

So, we see that the features used in MER are almost the same as those in AMGC. However, Fu et al. note in their extensive survey on Audio-based Music Classification [35] that although their effectiveness is debatable, mid-level features such as Rhythm seem to be more popular in MER.

<u>2.6.2.4   Classifiers.</u>    The algorithms used in AMGC are also popular in MER. So, support vector machines, Gaussian mixture models, neural networks, and k-nearest neighbor are the ones regularly used.

# 3   SYSTEM ANALYSIS

## 3.1   Requirement Specification

### 3.1.1   High Level Requirements

Our music classification system will perform classification of audio files on the basis of genre and mood. Genres include:

- Hip-hop

- Rock

- Jazz

- Classical

- Pop

Another classification that has been accomplished is the mood based classification. Under mood based classification, the audio files can be classified along the lines of:

- Depressive

- Frantic

- Exuberant

- Contentment

The classification of the audio files based on either mood or genre can be then be used in the creation of auto generated playlists.After each high level requirements are identified, corresponding intermediate level requirements (ILRs) are also identified. These list the feature metadata used to classify the audio.

### 3.1.2 Functional Requirements

- The classification of music based on genre and mood.

- The classification will work on various music file format like mp3, wav, etc.

### 3.1.3  Non-Functional Requirements

Nonfunctional requirement is a requirement that specifies criteria that can be used to judge the operation of a system, rather than specific behaviors.

The nonfunctional requirement in our project are:

- Code Documentation

- Project Documentation

- Code Quality

- Performance

- Fault Tolerance

- Log Maintenance

- Scalability

- Testability

- Maintainability

## 3.2  Feasibility Assessment

A feasibility assessment or feasibility analysis is a preliminary study undertaken before the real work of a project starts to ascertain the likelihood fo the project's success. It is an analysis of possible alternative solutions to a problem and a recommendation on the best alternative. It, for example, can decide whether an order processing be carried out by a new system more efficiently than the previous one. It can be thought as an assessment of the practicality of a proposed project.

A feasibility study aims to objectively and raitonally uncover the strengths and weakness of an existing project paradigm, opportunities and threats present in the environment, the resources required to carry through, and ultimately the prospects of success. In simplest terms,

the two criteria to judge feasibility are cost required and value to be attained. A well-designed feasibility study should provide a historical background of the project, a description of the project and details of operations and technicality needed.Generally feasibility studies always precede technical development and project implementation. A feasibility study evaluates the project's potential for success. It must there be conducted with an objective, unbiased approach to provide information upon which decisions can be based.

The acronym TELOS refer to the five areas of feasibility:

- Technical feasibility

- Economic feasibility

- Legal feasibility

- Operational feasibility

- Scheduling feasibility

### 3.2.1 Technical Feasibility

This assessment is based on an outline design of system requirements, to determine whether the company has the technical expertise to handle completion of the project. After our thorough research we decided to go with the features mentioned above. We also chose simple but powerful algorithm to do the classification. Since our first priority is to do all the coding from the scratch, hence we might face somewhat difficulty as we don't have that much experience of that of library developer. Moreover, doing so we might face difficulty in validating the feature extraction process. But based on our experience we can do most of the most by coding ourselves or at least try and then analyze the result. Based on our novice nature in software development, the project might also be somewhat not highy optimized

### 3.2.2 Economic feasibility

The purpose of the economic feasibility assessment is to determine the positive economic benefits to the organization that the proposed system will provide. This assessment typically

involves a cost/ benefits analysis.

Since our project is just the beginning portion which could be a huge leading factor in future and bring revolution in music industry. Our current product might not be able to provide that much of montary benefit but sure it is contributing to a huge research ahead in future. Then, talking about the cost, as no hardware is required in our project so we can say there is almost no loss involved in monetary terms but we are surely obtaining huge amount of knowledge and insight.

### 3.2.3 Operational Feasibility

Operational feasibility is a measure of how well a proposed system solves the problems, and takes advantage of the opportunities identified during scope definition and how it satisfies the requirements identified in the requirements analysis phase of system development. The operational feasibility assessment focuses on the degree to which the proposed development projects fits in with the existing business environment. Since our project is one of the most talked about project in the music industry, hence we can say that it will surely contribute to the music industry. The possibility of automated song management, smart playlist,etc. might be new thing in future. Our research has already pointed to the prominent feature to be included like MFCC, pitch,etc. The design of the data flow diagram shows feature extraction not to so hard and same goes for the classfication. Moreover as we are developing in the LINUX system so we are contributing towards it's environment.

### 3.2.4 Schedule feasibility

A project will fail if it takes too long to be completed before it is useful. Typically this means estimating how long the system will take to develop, and if it can be completed in a given time period using some methods like payback period. The gantt chart in our proposal is pretty much convincing to move along with but however somewhat delay might be obtained if any problem like lack of accuracy and misleading feature value arises.

### 3.2.5 Legal feasibility

Legal feasibility determines whether the proposed system conflicts with legal requirements. After our research, we found that our conflicting factor may only be on the use of the database provided provided we misuse it or violate some conditions.
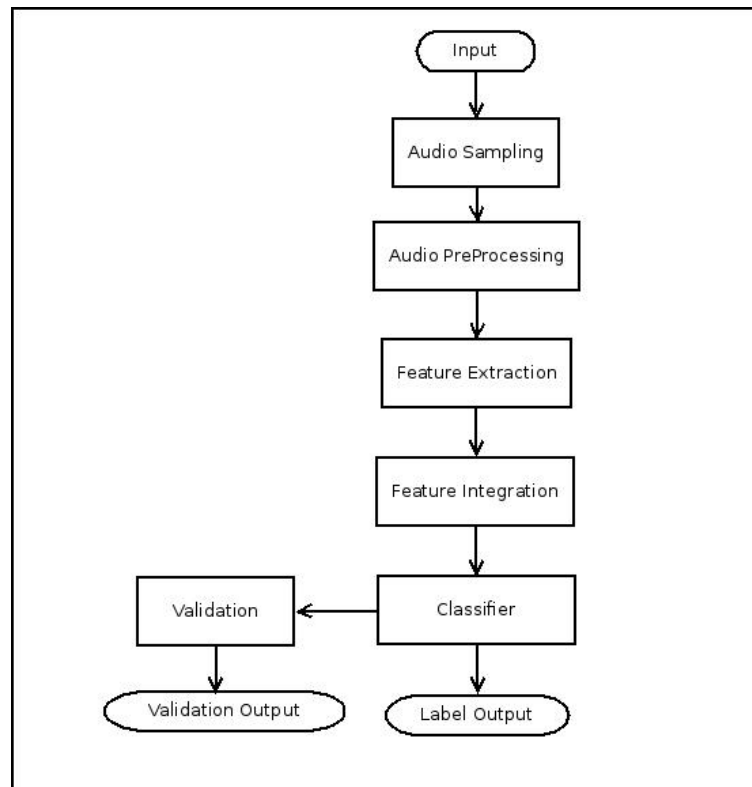
# 4 METHODOLOGY

## 4.1 System Block Diagram



Figure 4.1: General workflow

## 4.2 Audio Sampling

First of all to extract any features from a audio, samples are first taken out of it. Theses samples are then fine tuned to the correct audio format. This process involves converting to mono signals if they are stereo and adjusting bit depths if it is deemed to be necessary. Sample extraction and audio format is correction is achieved through the java sound API. Audio Sampling is then followed by audio framing and then a hamming window is used. Windows are necessary because whenever we do a finite Fourier transform, it is implicitly being applied to an infinitely repeating signal. So, for instance, if the start and end of a finite sample doesnt match then that will look just like a discontinuity in the signal, and show up

as lots of high-frequency nonsense in the Fourier transform, which is harmful. If the sample happens to be a perfect sinusoid but with an integer number of periods then it doesn't fit exactly into the finite sample and the FT will show appreciable energy in all sorts of places nowhere near the real frequency. Windowing the data makes sure that the ends match up while keeping everything reasonably smooth.

Pre processing the audio signal involves the following sub steps:

- Sampling the audio and fixing the audio format (Bit depth, number of channels etc)

- Audio framing to carry out feature extraction in each frame individually.

- Finally the use a hamming window to minimize the signal side lobe (unwanted radiation). Thus improving the quality or harmonics of the sound

Feature extraction involves the calculation of the following features:Root Mean Square(RMS), Compactness, MFCC, Rhythm, Pitch, Zero Crossing, Spectral variability, Spectral roll off point, Spectral centroid and Spectral flux. These features are extracted using the algorithms described in the previous chapter.

Feature integration involves the use of various measures of performance such as precision, recall and F measure to determine which feature is the best to calculate the genre/arousal/valence of an audio signal. So, using the accuracy produced by each feature when used in the SVM classifier, we select the best features, or a combination of features for the particular classification. We then integrate the features together by simply appending them.

Finally in the classification part of the system, we utilize two different types of classifiers, Artificial Neural networks and Support Vector machines to classify the audio signals and find the proper label for each based on genre/arousal/valence.

Validation is done when building the model. We use k-fold cross-validation with k=10 to validate the classifier. Obviously, this step is not done in the workflow of the final product.

# 5  SYSTEM DEVELOPMENT

## 5.1  Data Flow Diagram

### 5.1.1  Level Zero
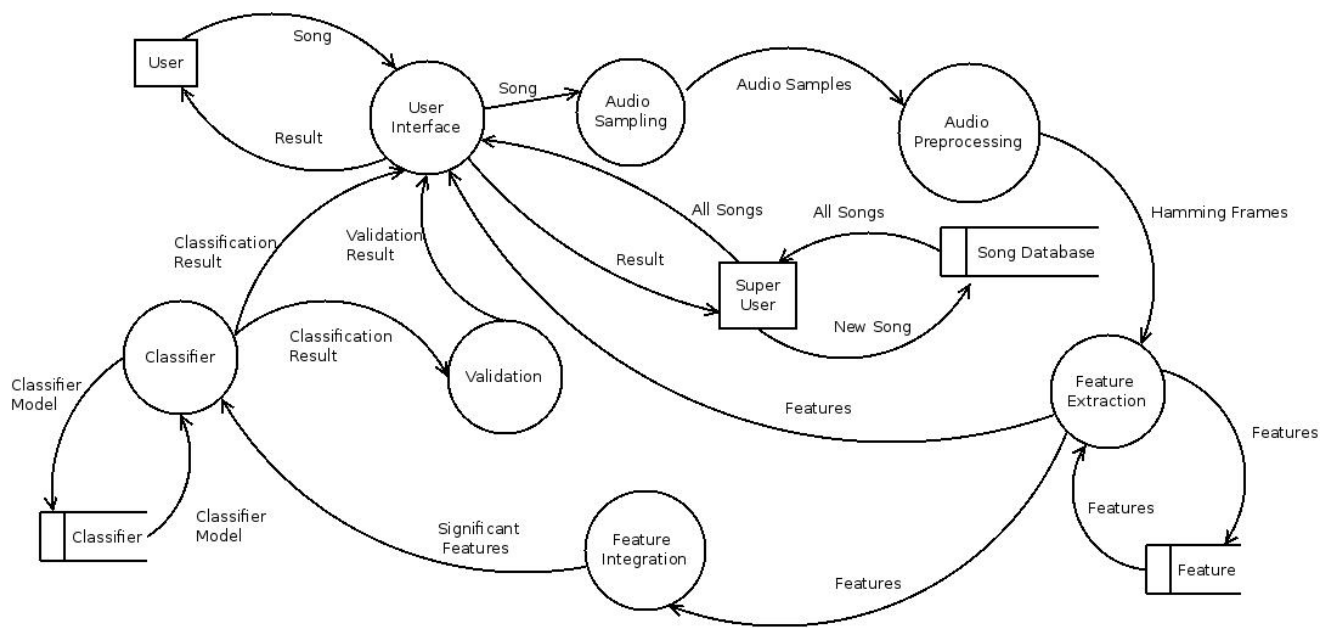


Figure 5.1: Context Level Zero

### 5.1.2  Level One

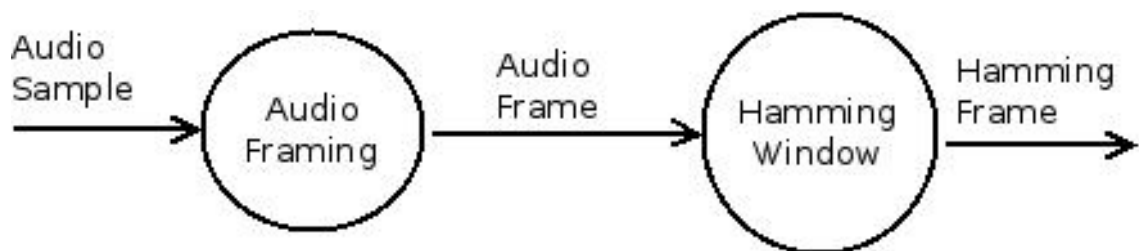Figure 5.2: Level One

## 5.1.3   Level Two
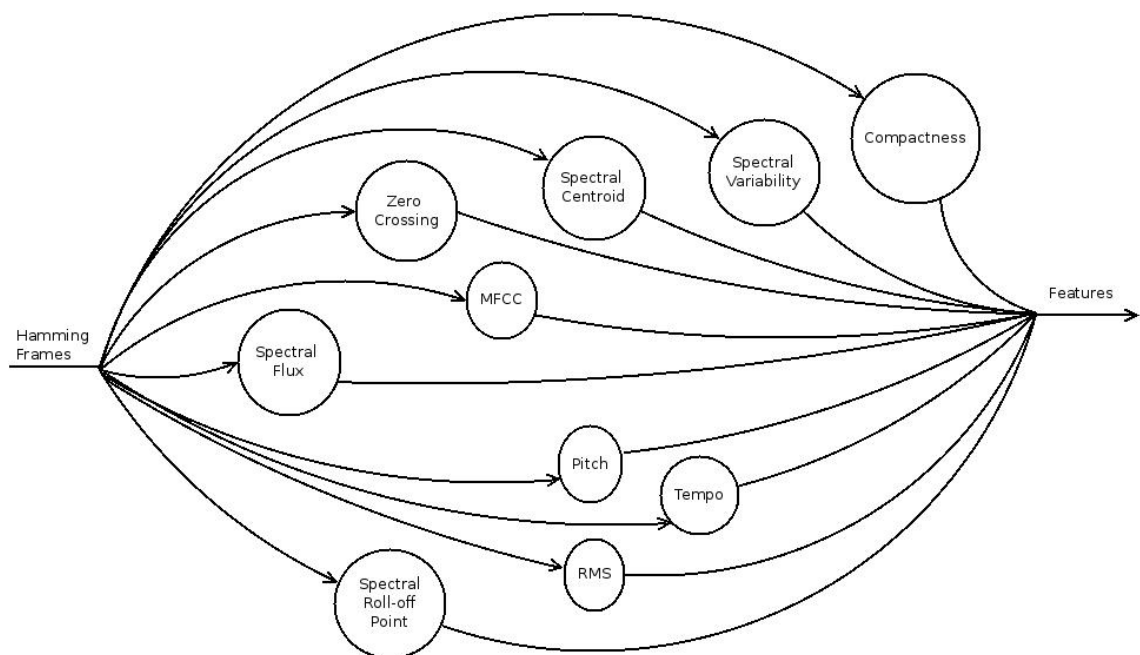


Figure 5.3: Pre-processing
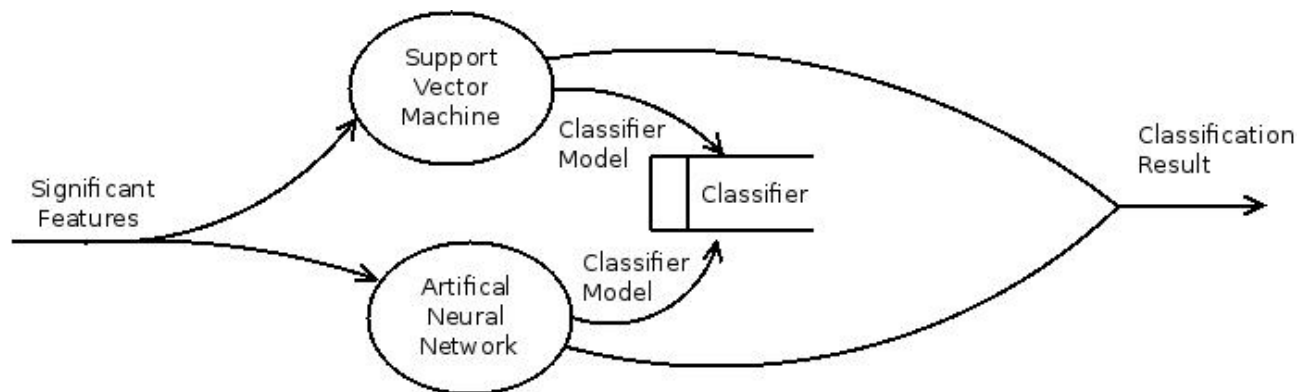
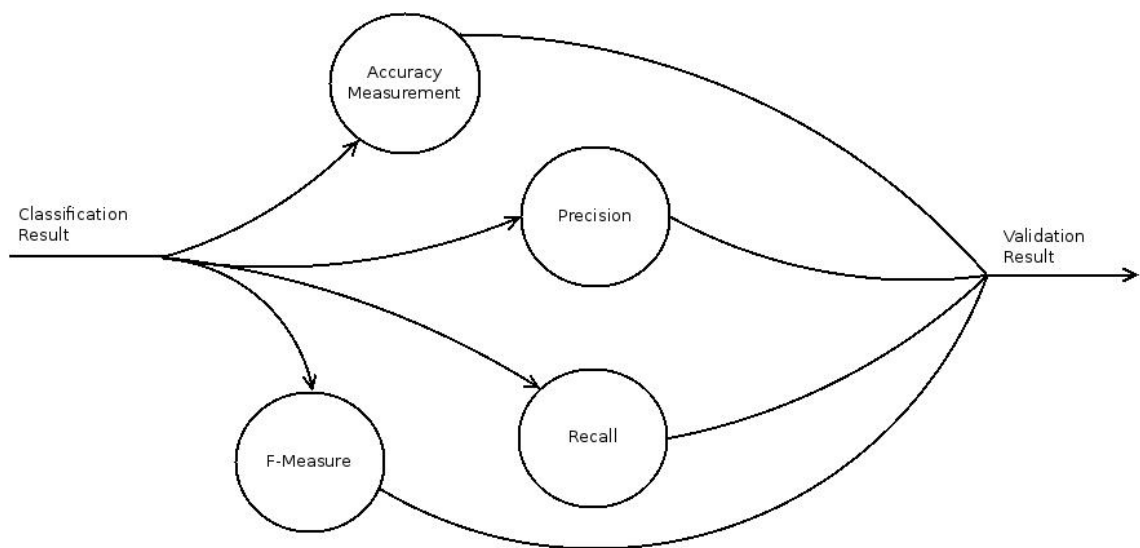Figure 5.4: Feature



Figure 5.5: Classifier

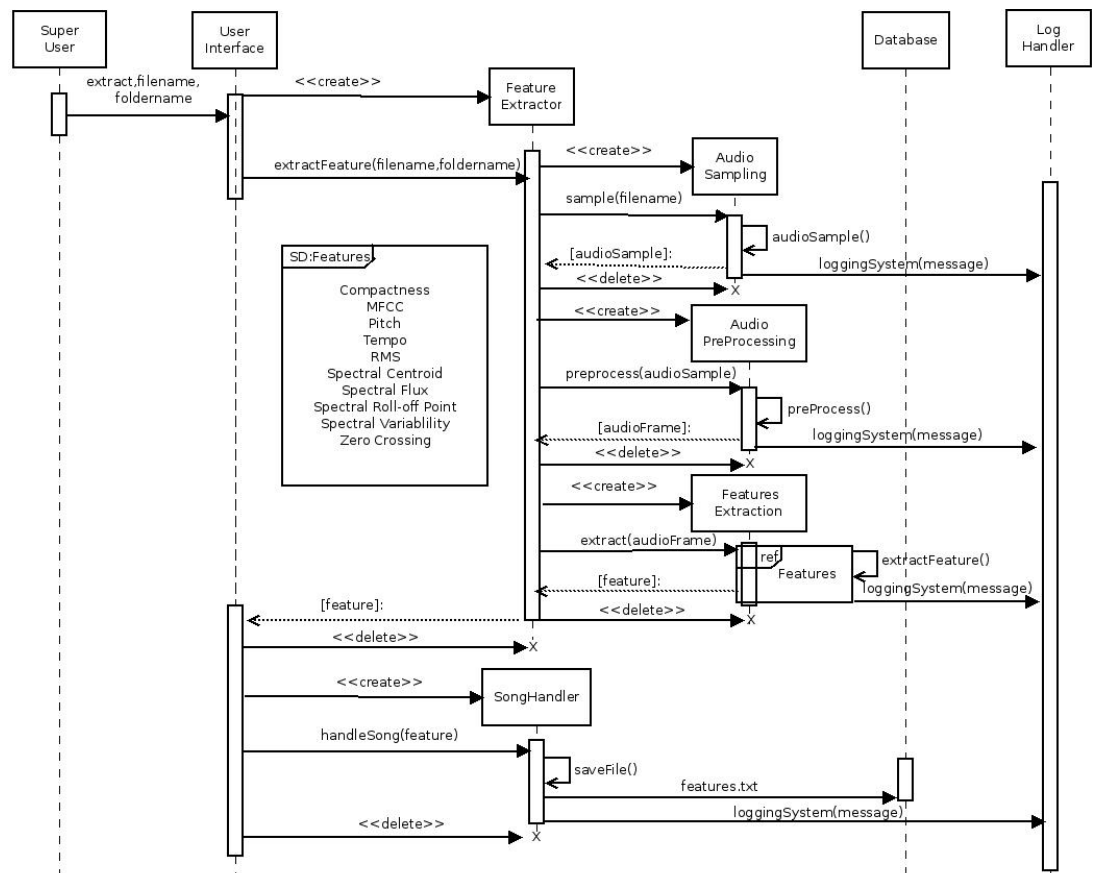Figure 5.6: Validation

## 5.2 Sequence Diagram

### 5.2.1 Feature Extraction

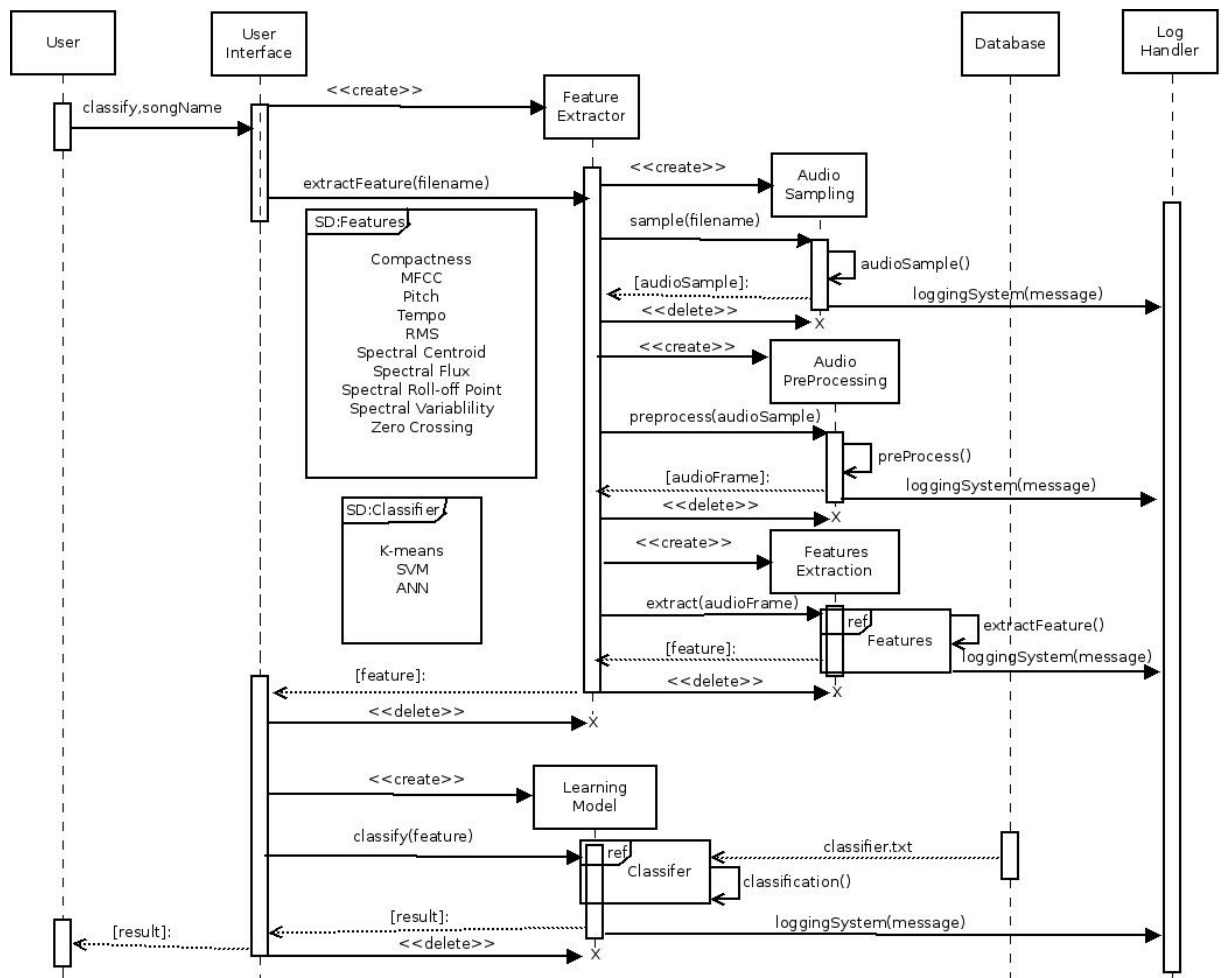Figure 5.7: Sequence diagram of extracting feature

## 5.2.2 Classification

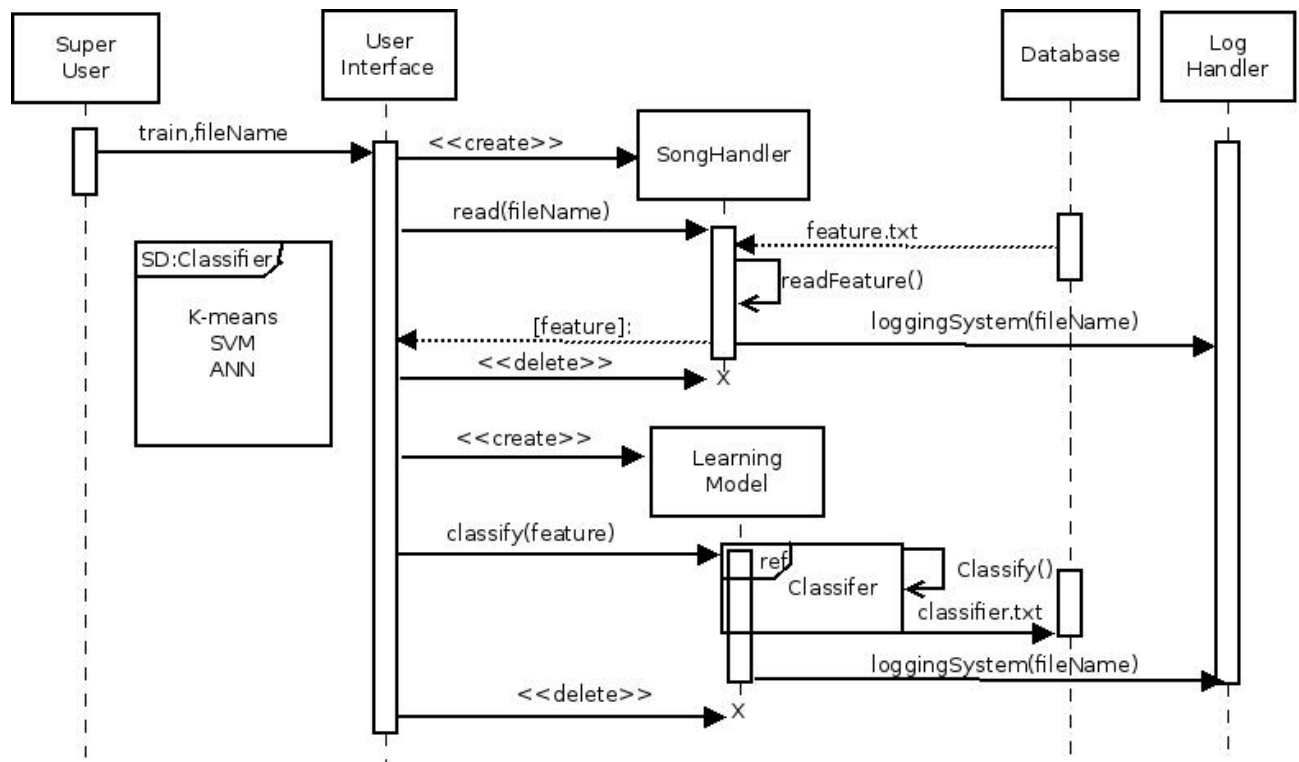Figure 5.8: Sequence diagram of classification

## 5.2.3 Training

Figure 5.9: Sequence diagram of training

## 5.2.4 Testing

## 5.3 Tools and Environment

### 5.3.1 Java Programming Language

Java is a general-purpose computer programming language that is concurrent class-based, object-oriented. We chose java due to many reasons:

- We were familiar with Java having worked on pretty decent sized Java projects.

- Although Java may be too verbose for small applications, we believe it is

- Portability - Java runs well on all the popular platforms.

- Speed - The latest JIT compilers for Suns JVM approach the speed of C/C++ code, and in some memory allocation intensive circumstances, exceed it.

- Standard APIs that encompass lots of common funcionalities.

Figure 5.10: Sequence diagram of testing

- GarbageCollection - the programmer doesn't have to worry about memory (most of the time).

- Availability of lots of Libraries, tools and frameworks in Java.

## 5.3.2 SMILE - Statistical Machine Intelligence and Learning Engine

Smile (Statistical Machine Intelligence and Learning Engine) is a fast and comprehensive machine learning system. With advanced data structures and algorithms, Smile delivers state-of-art performance. As mentioned in its documentation: "Smile covers every aspect of machine learning, including classification, regression, clustering, association rule mining, feature selection, manifold learning, multidimensional scaling, genetic algorithms, missing

value imputation, efficient nearest neighbor search, etc."

### 5.3.3 TarsosDSP

TarsosDSP is a Java library for audio processing. Its aim is to provide an easy-to-use interface to practical music processing algorithms implemented, as simply as possible, in pure Java and without any other external dependencies. The library tries to hit the sweet spot between being capable enough to get real tasks done but compact and simple enough to serve as a demonstration on how DSP algorithms works. We used TarsosDSP to process the low level song features.

### 5.3.4 Google Gson

Gson is a Java library that can be used to convert Java Objects into their JSON representation. It can also be used to convert a JSON string to an equivalent Java object. Gson can work with arbitrary Java objects including pre-existing objects that you do not have source-code of. We used Gson to store the features as JSON and read them back. We chose JSON because it is easily readable through a lot of methods, libraries and languages.

### 5.3.5 Xstream

XStream is a simple library to serialize objects to XML and back again. We used XStream to store the model as a XML file so that it can be read back.

### 5.3.6 MP3SPI

MP3SPI is a Java Service Provider Interface that adds MP3 (MPEG 1/2/2.5 Layer 1/2/3) audio format support for Java Platform. We used it to add a more reliable support for mp3 files.

### 5.3.7    Python Programming Language

Python is a high-level, general-purpose, interpreted, dynamic programming language. We used Python scripts to handle the audio files and modify them to create the mood based dataset.

# 6 RESULTS AND ANALYSIS

## 6.1 Dataset pre-processing

### 6.1.1 Genre Dataset

The GTZAN dataset for genre classification didn't require any special treatment as it is one of the most widely used dataset in Automatic Music Classification. So, the format of the songs (22 050 Hz, 16-bit, mono audio files) didn't require any special pre-processing.

### 6.1.2 Mood Dataset

We acquired a filtered version (with some redundancies removed) of the 1000 songs dataset for emotional analysis resulting in a final set of 744 songs.

Each of the songs were of optimum length (45 seconds), and in suitable format too. However, we had to generate the labels for the song from the numerical values given in the dataset. Each song had two associated numerical values in the range of 1-9. We used empirically determined values close to 5 to separate the dataset based on arousal and valence using a Python script.
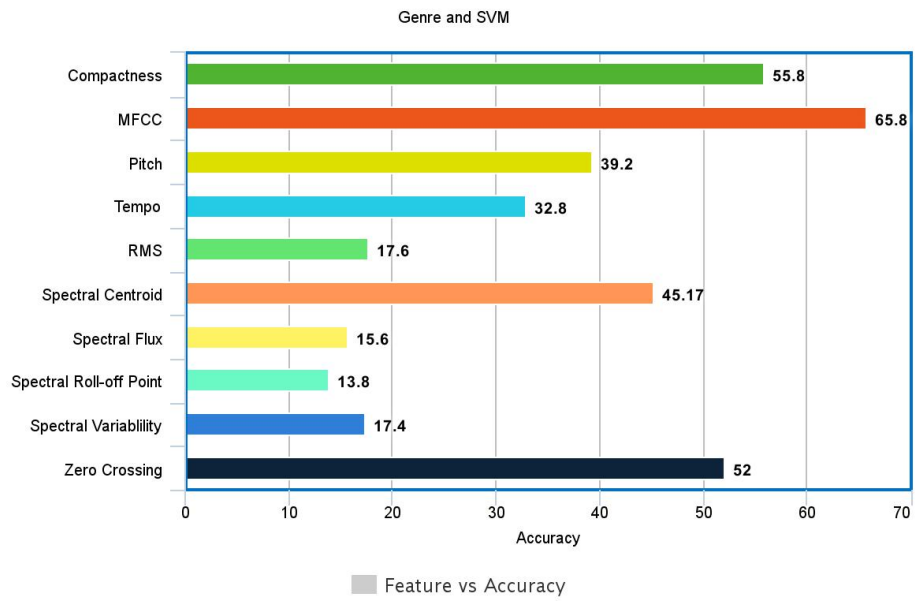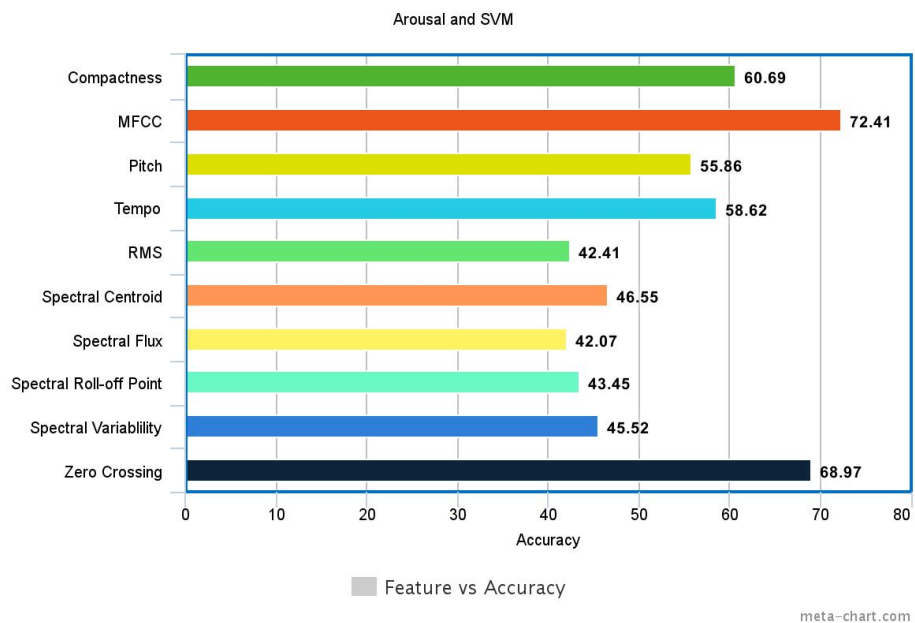
## 6.2 SVM

Figure 6.1: Genre using SVM



Figure 6.2: Arousal using SVM

We found MFCC to give the best accuracy for Genre and Arousal. For valence, tempo gave the best result.

So, for Genre, we only took the following features into consideration in our final model:
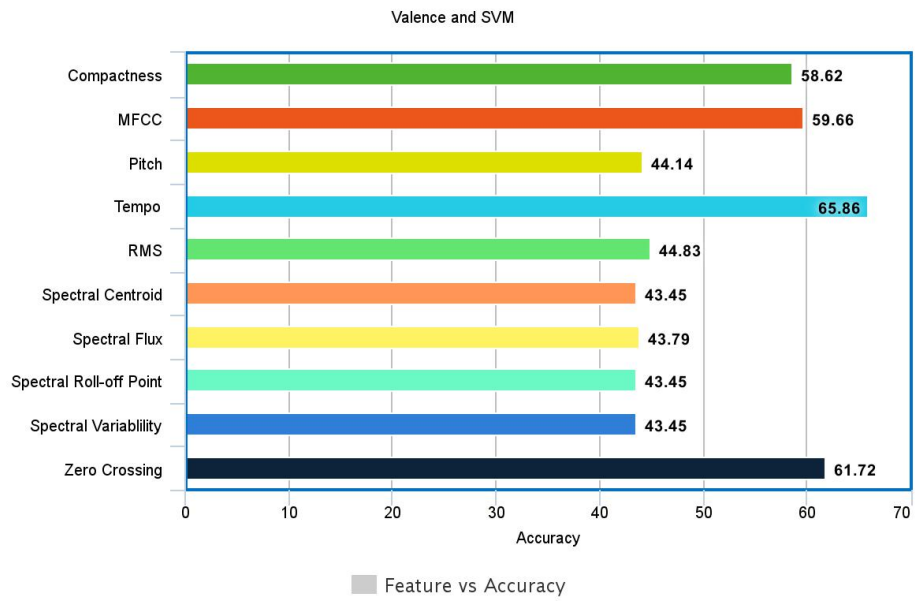
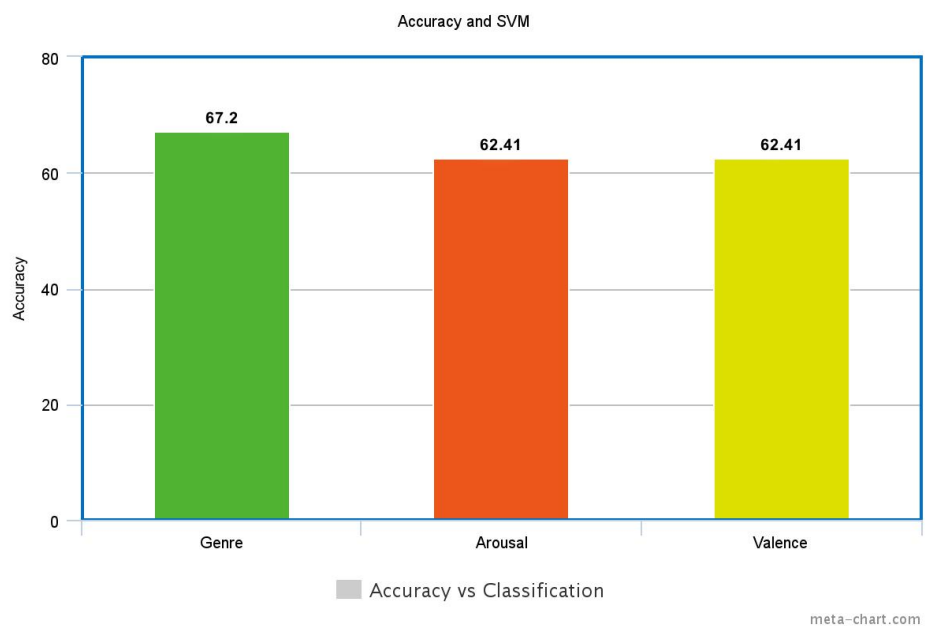- Compactness

- MFCC

Figure 6.3: Valence using SVM



Figure 6.4: Accuracy of SVM

- Zero Crossing

For Arousal, we also factored in Tempo as it also provided good results.
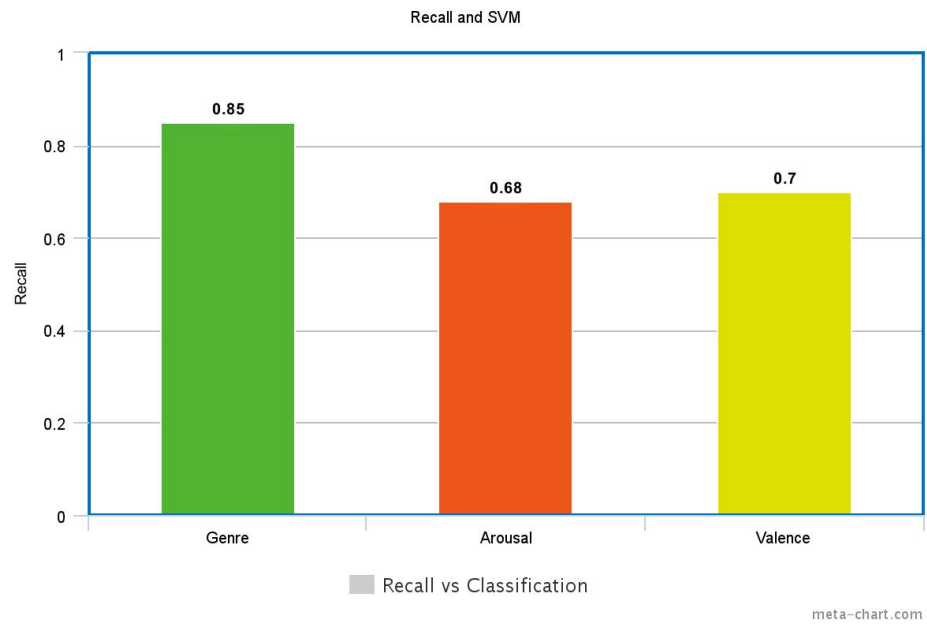
- Compactness

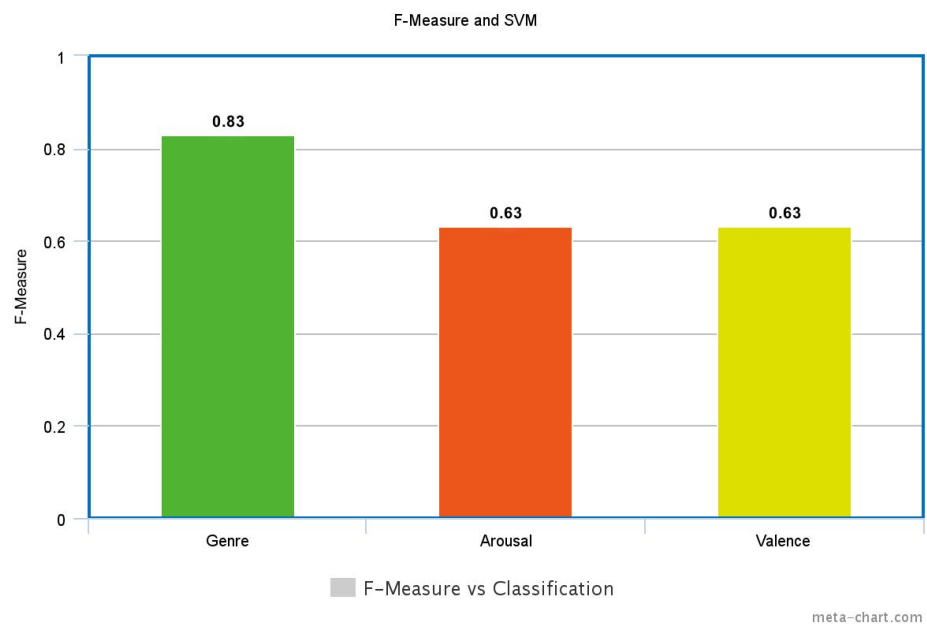- MFCC

Figure 6.5: Recall of SVM



Figure 6.6: F-measure of SVM

- Tempo

- Zero Crossing

For Valence, we left out MFCC and retained the three other good performers.
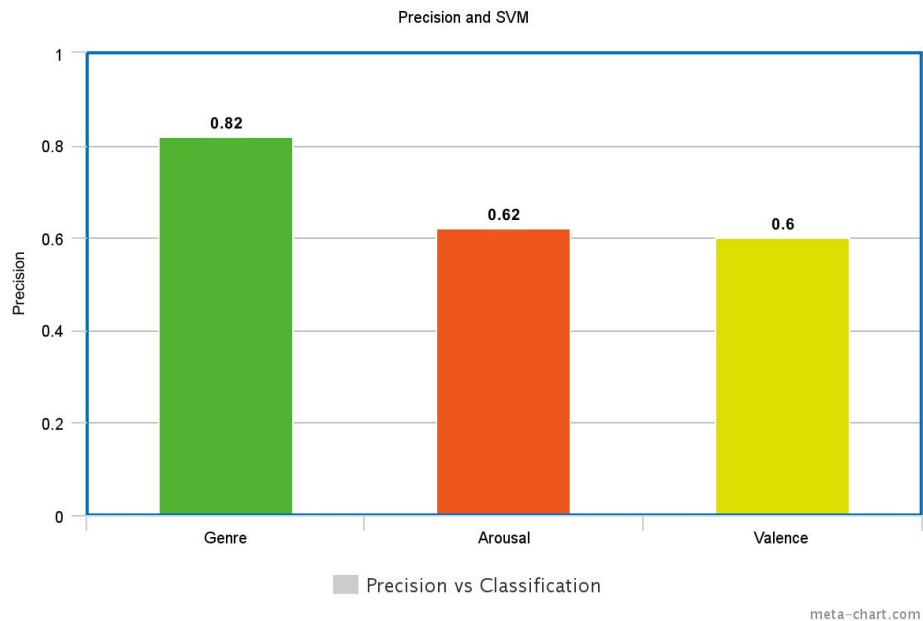
- Compactness

Figure 6.7: Precision of SVM

- Tempo

- Zero Crossing

## 6.3 ANN

MFCC performed well in Genre Classification as well as for Arousal Classification. For valence too, MFCC performed well but it is only a little better than random. This is probably because, the people annotating the dataset had the most disagreement on valence.

So, in the final classsifier, we used MFCC and RMS for all classifiers as they gave the best result.

## 6.4 Problems and Solutions

We faced the following problems during the project:

- Decoding Mp3: Mp3 files are lossy compression of audio. So, we had a hard time decoding mp3 files on our custom implementation.

- Distance Metric: We couldn't find a decent distance metric between songs.
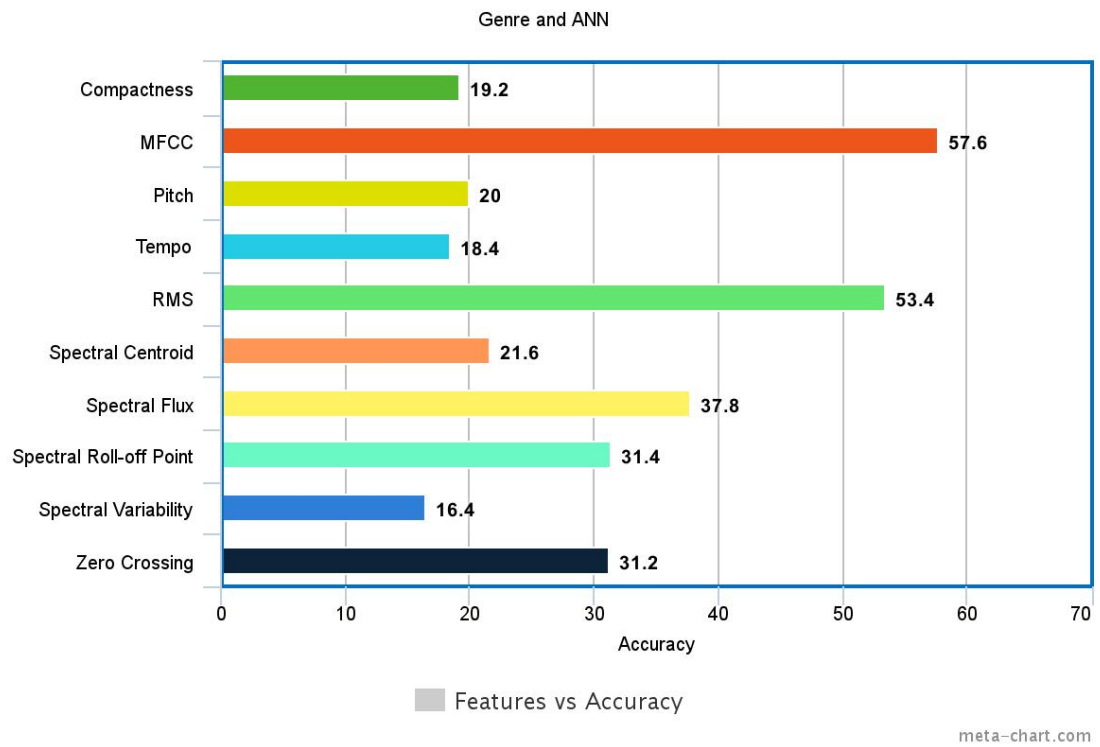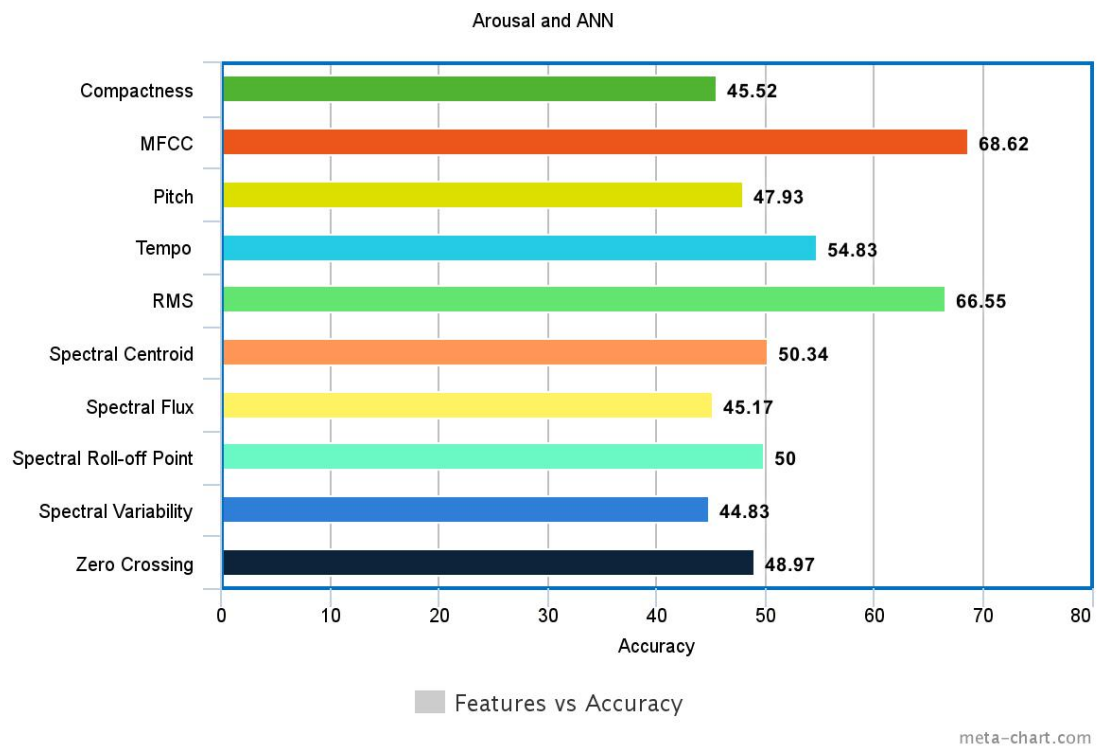
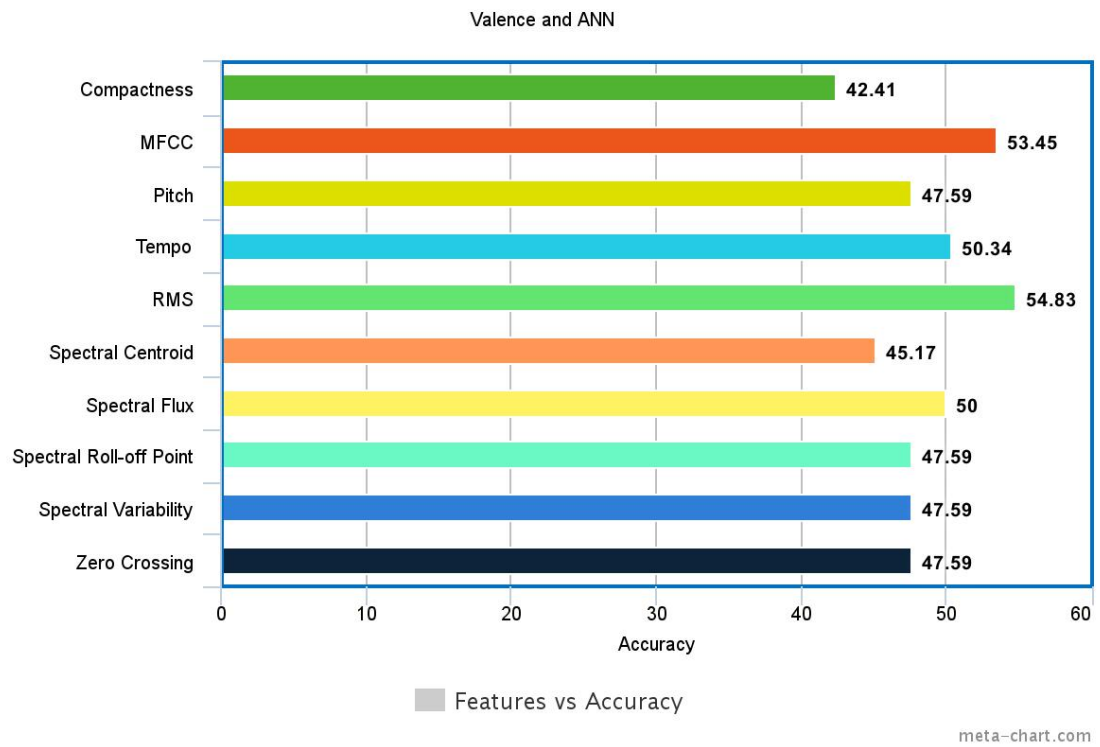Figure 6.8: Genre using ANN



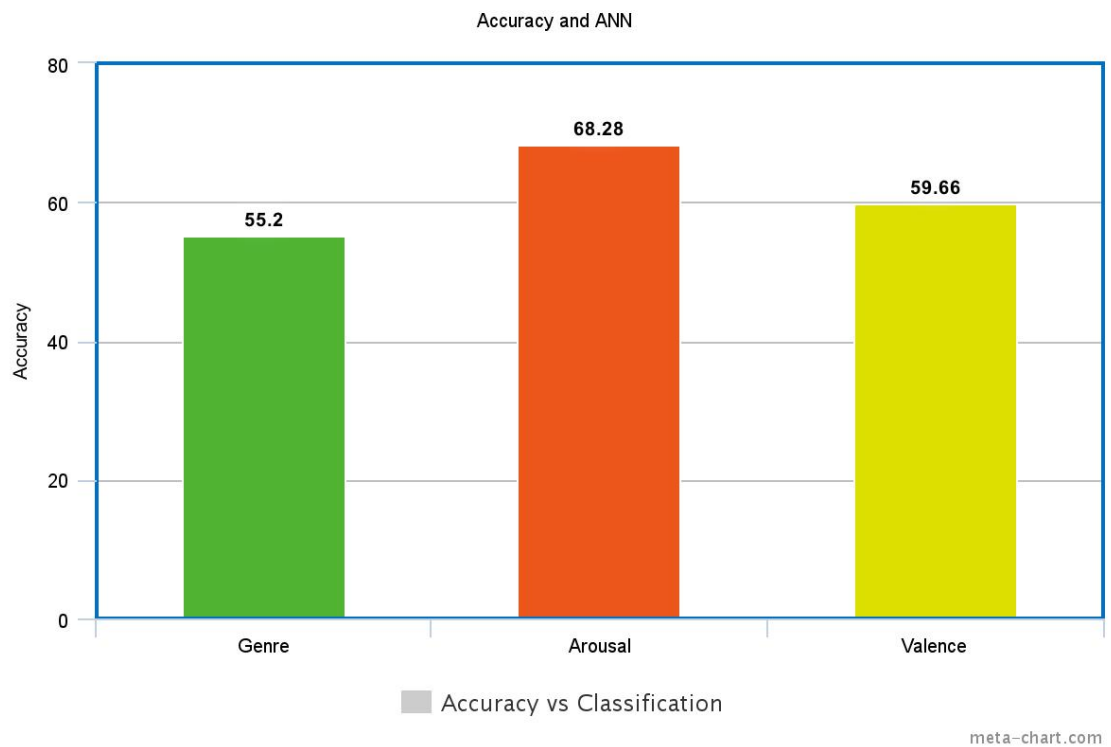Figure 6.9: Arousal Using ANN

Figure 6.10: Valence Using ANN



Figure 6.11: Accuracy of ANN

Figure 6.12: F-measure of ANN



Figure 6.13: Precision of ANN

Figure 6.14: Recall of ANN

- Data Refining: Our mood based dataset was not directly labeled. They had numerical values instead.

- Computational Intensiveness For Preprocessing: As the samples of the song needed to be kept in memory during pre-processing, for long songs, this resulted in a lot of memory being consumed.

- Computational Intensiveness For Classification: Some classifiers such as those based on GMM were computationally intensive in nature.

### 6.4.1 Solved Problems

- Decoding Mp3: We solved it by using the MP3SPI library.

- Data Refining: We wrote custom Python scripts to generate the dataset.

### 6.4.2 Unsolved Problems

- Distance Metric: We circumvented the problem by using classifiers that didn't need distance metrics. So, we didn't use K-NN and K-means.

- Computational Intensiveness For Preprocessing: We still have this problem. And so, we might get Heap error for long songs ($>$5 mins).

- Computational Intensiveness For Classification: Instead of decreasing the computational intensiveness of the implementation, we used other classifiers instead.

# 7 CONCLUSION

## 7.1 Project Conclusion

In our musical classification system we have successfully been able to produce satisfactory results in all scenarios of genre,arousal and valence based music classification As with many projects we faced many problems some of which we were able to solve and some which still remain as limitations.

So, overall we can say that we have achieved the satisfactory accuracy in various aspects and improving the performance through the continuous integration developmental practice.

## 7.2 Future Enhancements

A person can never be satisfied with one's work and so is our case. Our project is just a beginning phase which will be ultimately leading to huge revolution in future. Hence future enhancements are possible in our project are possible in our project. Mentioning few are:

(i) **GUI Improvement:** As with the current state of the GUI in our system , it is very basic. With every future iteration we hope to improve and introduce newer and attractive elements and animations that pleasantly improve the user experience and encourages user to use our system as much as possible.

(ii) **Inclusion of additional genres:** With only five basic recognizable genres we hope to include further sub genres with these five primary genres such as PopRock, PunkRock, Grunge,IndieRock for the Rock songs etc

(iii) **Instrumentals:** We also hope to add the ability to add instrumental tracks (i.e tracks without vocals) in future enhancements

(iv) **Preference changes:** Owing to the ever changing nature of music preference of an average user, we may also provide system with ability to change the various genre and mood template so that the system also changes with the change in user preference.

(v) **Improve performance:** Performance can be improved by using Multithread concept.

(vi) **Business Value:** Develop the product as the API so that it can be used in others commercial sites.

# REFERENCES

[1] Smith, Steven W. *The scientists and engineers guide to digital signal processing.* California Technical Publishing, 2013. Print.

[2] Kassler, Michael. *Toward musical information retrieval.* Perspectives of New Music (1966): 59-67.

[3] Andel, James. *On the segmentation and analysis of continuous musical sound by digital computer.* Diss. Stanford University, 1975.

[4] Tzanetakis, George, and Perry Cook. *Manipulation, analysis and retrieval systems for audio signals.* Princeton, NJ, USA: Princeton University, 2002.

[5] Alghoniemy, Masoud, and Ahmed H. Tewfik. *Rhythm and periodicity detection in polyphonic music.* Multimedia Signal Processing, 1999 IEEE 3rd Workshop on. Ieee, 1999.

[6] Byrd, Donald, and Tim Crawford. *Problems of music information retrieval in the real world.* Information processing and management 38.2 (2002): 249-272.

[7] Saunders, John. *Real-time discrimination of broadcast speech/music.* icassp. Vol. 96. 1996.

[8] Wold, Erling, et al. *Content-based classification, search, and retrieval of audio.* IEEE multimedia 3.3 (1996): 27-36.

[9] Scheirer, Eric, and Malcolm Slaney. *Construction and evaluation of a robust multifeature speech/music discriminator.* Acoustics, Speech, and Signal Processing, 1997. ICASSP-97., 1997 IEEE International Conference on. Vol. 2. IEEE, 1997.

[10] Casey, Michael A., et al. *Content-based music information retrieval: Current directions and future challenges.* Proceedings of the IEEE 96.4 (2008): 668-696.

[11] Kaminskas, Marius, and Francesco Ricci. *Contextual music information retrieval and recommendation: State of the art and challenges.* Computer Science Review 6.2 (2012): 89-119.

[12] Tzanetakis, George, and Perry Cook, Musical genre classification based on audio signals *IEEE Transactions on Speech and Audio Processing*, 10.5 (2002)

[13] Scaringella, Nicolas, Giorgio Zoia, and Daniel Mlynek. *Automatic genre classification of music content: a survey.* IEEE Signal Processing Magazine 23.2 (2006): 133-141.

[14] Perrot, David, and Robert Gjerdigen. *Scanning the dial: An exploration of factors in the identification of musical style.* Proceedings of the 1999 Society for Music Perception and Cognition. 1999.

[15] Lippens, Stefaan, Jean-Pierre Martens, and Tom De Mulder. *A comparison of human and automatic musical genre classification.* Acoustics, Speech, and Signal Processing, 2004. Proceedings.(ICASSP'04). IEEE International Conference on. Vol. 4. IEEE, 2004.

[16] Kour, Gursimran, and Neha Mehan. *Music Genre Classification using MFCC, SVM and BPNN.* International Journal of Computer Applications (09758887) Volume (2015).

[17] Logan, Beth. *Mel Frequency Cepstral Coefficients for Music Modeling.* ISMIR. 2000.

[18] Li, Tao, Mitsunori Ogihara, and Qi Li. *A comparative study on content-based music genre classification.* Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval. ACM, 2003.

[19] Anglade, Amlie, et al. *Improving music genre classification using automatically induced harmony rules.* Journal of New Music Research 39.4 (2010): 349-361.

[20] Muller, Mathias, et al. *Signal processing for music analysis.*, Selected Topics in Signal Processing, IEEE Journal of 5.6 (2011): 1088-1110.

[21] Dooling, Robert J and Stewart H Hulse. *The Comparative Psychology Of Audition.* Hillsdale, N.J.: L. Erlbaum Associates, 1989. Print.

[22] Prasad, Bhanu, and SR Mahadeva Prasanna, eds. *Speech, audio, image and biomedical signal processing using neural networks.* Vol. 83. Springer, 2007.

[23] Neumayer, Robert. *Musical genre classification.* (2004).

[24] Koerich, Alessandro L. *Improving the Reliability of Music Genre Classification using Rejection and Verification.* ISMIR. 2013.

[25] Haggblade, Michael, Yang Hong and Kenny Kao. Music genere classification. *Department of Computer Science*, Stanford University, 2011.

[26] Nasridinov, Aziz, and Young-Ho Park. *A Study on Music Genre Recognition and Classification Techniques.* International Journal of Multimedia and Ubiquitous Engineering 9.4 (2014): 31-42.

[27] Li, Tao, and George Tzanetakis. *Factors in automatic musical genre classification of audio signals.* Applications of Signal Processing to Audio and Acoustics, 2003 IEEE Workshop on.. IEEE, 2003.

[28] Downie, J. Stephen. *Toward the scientific evaluation of music information retrieval systems.* ISMIR. 2003.

[29] Tzanetakis, George, and Perry Cook, Musical genre classification *Journal of Personality*, 60(2):225251, 1992.

[30] Ekman, Paul. *Are there basic emotions?.* (1992): 550.

[31] James A. Russell. *A circumplex model of affect.* Journal of Personality and Social Psychology, 39(6): 1161-1178, 1980.

[32] Thayer, Robert E. *The biopsychology of mood and arousal.* Oxford University Press, 1990.

[33] Juslin, Patrik N., and John A. Sloboda. *Music and emotion: Theory and research.* Oxford University Press, 2001.

[34] Lu, Lie, Dan Liu, and Hong-Jiang Zhang. *Automatic mood detection and tracking of music audio signals.* IEEE Transactions on audio, speech, and language processing 14.1 (2006): 5-18.

[35] Fu, Zhouyu, et al. *A survey of audio-based music classification and annotation.* IEEE Transactions on Multimedia 13.2 (2011): 303-319.

[36] Soleymani, Mohammad, et al. *1000 songs for emotional analysis of music.* Proceedings of the 2nd ACM international workshop on Crowdsourcing for multimedia. ACM, 2013.

[37] Apon, Amy, et al. *Inital Starting Point Analysis for K-Means Clustering: A Case Study.* (2006).

[38] Jain, Anil K. *Data clustering: 50 years beyond K-means.* Pattern recognition letters 31.8 (2010): 651-666.

[39] Hamerly, Greg, and Charles Elkan. *Alternatives to the k-means algorithm that find better clusterings.* In Proceedings of the eleventh international conference on Information and knowledge management, pp. 600-607. ACM, 2002.

# APPENDIX A

## 8.3  Signal Processsing

In signal processing, a window function (also known as an apodization function or tapering function) is a mathematical function that is zero-valued outside of some chosen interval. For instance, a function that is constant inside the interval and zero elsewhere is called a rectangular window, which describes the shape of its graphical representation. When another function or waveform/data-sequence is multiplied by a window function, the product is also zero-valued outside the interval: all that is left is the part where they overlap, the "view through the window".

Applications of window functions include spectral analysis, filter design, and beam forming. In typical applications, the window functions used are non-negative smooth "bell-shaped" curves, though rectangle, triangle, and other functions can be used. A more general definition of window functions does not require them to be identically zero outside an interval, as long as the product of the window multiplied by its argument is square integral, and, more specifically, that the function goes sufficiently rapidly toward zero. One of the major applications of window functions includes the design of finite impulse response filters and the spectral analysis.

## 8.4  Spectral Analysis

The Fourier transform of the function cos t is zero, except at frequency . However, many other functions and waveforms do not have convenient closed form transforms. Alternatively, one might be interested in their spectral content only during a certain time period. In either case, the Fourier transform (or something similar) can be applied on one or more finite intervals of the waveform. In general, the transform is applied to the product of the waveform and a window function. Any window (including rectangular) affects the spectral

estimate computed by this method.

## 8.5   Windowing

Windowing of a simple waveform like cos t causes its Fourier transform to develop non- zero values (commonly called spectral leakage) at frequencies other than . The leakage tends to be worst (highest) near  and least at frequencies farthest from . If the waveform under analysis comprises two sinusoids of different frequencies, leakage can interfere with the ability to distinguish them spectrally. If their frequencies are dissimilar and one component is weaker, then leakage from the larger component can obscure the weaker ones presence.  But if the frequencies are similar, leakage can render them irresolvable even when the sinusoids are of equal strength.

The rectangular window has excellent resolution characteristics for sinusoids of comparable strength, but it is a poor choice for sinusoids of disparate amplitudes. This characteristic is sometimes described as low-dynamic-range.

At the other extreme of dynamic range are the windows with the poorest resolution. These high-dynamic-range low-resolution windows are also poorest in terms of sensitivity; this is, if the input waveform contains random noise close to the frequency of a sinusoid, the response to noise, compared to the sinusoid, will be higher than with a higher-resolution window.  In other words, the ability to find weak sinusoids amidst the noise is diminished by a high- dynamic-range window.  High-dynamic-range windows are probably most often justified in wideband applications, where the spectrum being analyzed is expected to contain many different components of various amplitudes.

In between the extremes are moderate windows, such as Hamming and Hann.  They are commonly used in narrowband applications, such as the spectrum of a telephone channel. In summary, spectral analysis involves a tradeoff between resolving comparable strength components with similar frequencies and resolving disparate strength components with dissimilar frequencies. That tradeoff occurs when the window function is chosen.

## 8.6   Filter Bank

In signal processing, a filter bank is an array of band-pass filters that separates the input signal into multiple components, each one carrying a single frequency sub-band of the original signal. One application of a filter bank is a graphic equalizer, which can attenuate the components differently and recombine them into a modified version of the original signal. The process of decomposition performed by the filter bank is called analysis (meaning analysis of the signal in terms of its components in each sub-band); the output of analysis is referred to as a sub band signal with as many sub bands as there are filters in the filter bank. The reconstruction process is called synthesis, meaning reconstitution of a complete signal resulting from the filtering process.

In digital signal processing, the term filter bank is also commonly applied to a bank of receivers. The difference is that receivers also down-convert the sub bands to a low center frequency that can be re-sampled at a reduced rate. The same result can sometimes be achieved by under sampling the band pass sub bands.

Another application of filter banks is signal compression, when some frequencies are more important than others. After decomposition, the important frequencies can be coded with a fine resolution. Small differences at these frequencies are significant and a coding scheme that preserves these differences must be used. On the other hand, less important frequencies do not have to be exact. A coarser coding scheme can be used, even though some of the finer (but less important) details will be lost in the coding.

The vocoder uses a filter bank to determine the amplitude information of the sub bands of a modulator signal (such as a voice) and uses them to control the amplitude of the sub bands of a carrier signal (such as the output of a guitar or synthesizer), thus imposing the dynamic characteristics of the modulator on the carrier.

## 8.7   GUI

Figure 8.1: Graphics user interface