

Faculty of Mathematics and Computer Sciences



**FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA**

Master's Thesis

Music Similarity Analysis using Big Data approaches

presented by:

Johannes Schoder

born:

03. September 1994

ID:

169197

course of studies:

M.Sc. Informatik

supervisors:

Ralf Seidler, Prof. Dr. Martin Bucker

Contents

Abbreviations	iii
List of Figures	iv
List of Tables	v
1 Introduction	1
1.1 Why and how?	1
1.2 Overview	2
2 MIR Toolkits and Data Aggregation	3
2.1 Music Information Retrieval	3
2.1.1 Audio Features and Music Similarity Measurements	3
2.1.2 MIR Toolkits	8
2.2 Datasets	11
2.2.1 Free Music Archive	11
2.2.2 Musicnet and 1517 artists	12
2.2.3 MedleyDB	12
2.2.4 Private music collection	12
2.2.5 Overview and other sources	13
2.3 Spotipy	14
2.4 Apache Hadoop and Spark	15
3 Timbre Similarity	16
3.1 Single Gaussian Model	16
3.2 Other Similarity Metrics	17
3.3 Validation	18
3.3.1 Construction Noise	18
3.3.2 Different recordings and cover versions	19
3.4 Conclusions	20
4 Rhythmic Similarity	21
4.1 Cross Correlation	21
4.2 fluctuation patterns	21

4.3	Tempo Estimation	21
4.4	Beat and Tatum Estimation	21
5	Pitch estimation/ Chroma Features	22
5.1	Pitch Curve	22
5.2	MIDI representation	22
5.3	Graph representation	22
5.4	Chroma Features	22
5.4.1	Bandpass filter	22
5.4.2	Key invariance	22
5.4.3	Beat alignment	22
5.4.4	Validation	22
6	Other Similarity metrics	24
6.1	Definition of similarity	24
6.1.1	Rhythm, Timbre, Melodic	24
6.1.2	Genre and metadata	24
6.1.3	Genre specific features	24
6.1.4	Combinations and variable model	24
6.2	Collaborative Filtering	24
7	Big Data Framework Spark	25
7.1	Data aggregation	25
7.2	MapReduce tasks	25
8	Results	26
8.1	Cover song identification	26
8.2	Genre similarity	26
	References	27
9	Appendix	31
9.1	Spotipy Crawler	31

Abbreviations

MIR	music information retrieval
MFCC	mel frequency cepstral coefficients
FFT	fast fourier transformation
DCT	discrete cosine transformation
MSD	million song dataset
MIDI	musical instrument digital interface
KL divergence	Kullback-Leibler divergence
MP	mutual proximity

List of Figures

2.1	Frequency Space	4
2.2	Features	5
2.3	Features	5
2.4	Sweep Sinal	6
2.5	MFCCs	6
2.6	Original Scores	9
2.7	Aubio	9
2.8	Melodia	10
2.9	Transcription	10
2.10	Spotify API	11
2.11	Number of songs in the free music archive dataset per genre	11
2.12	Number of songs in the 1517 artists dataset per genre	12
2.13	Number of songs in the MedleyDB dataset per genre	12
2.14	Number of songs in the private music collection per genre	13
2.15	Number of songs in the million song dataset per genre	14
3.1	Construction Noise	19
5.1	Sia Features	23
5.2	Pvris Features	23

List of Tables

2.1	Music Datasets used	13
-----	-------------------------------	----

Abstract

Calculating music similarity metrics using map reduce algorithms.

This paper is about the comparison of construction noise and modern day music. The field of music information retrieval (MIR) in computer science is mostly a data driven and purely mathematical topic. The goal of this paper is to merge the fields of computer science with music theoretical knowledge and to find potential weak spots of current music similarity algorithms.

1. Introduction

The idea originated from Dr. T. Bosse from the chair for advanced computing at the Friedrich-Schiller-University (FSU) in Jena. When proposing the idea for a master thesis with the topic of "Music similarity measurement using genre specific features" using different guitar play styles in modern day metal music, he jokingly said that he would also like to know how metal music compares to construction building noise. The idea is actually not so groundless, considering that most people would agree on the fact that metal music is often described as noise by people not used to listening to genres like death and black metal. This thesis is meant to evaluate how music similarity algorithms compare construction noise to common musical genres and how to possibly improve these algorithms.

1.1 Why and how?

Why is music similarity research even necessary?

First of all, there is no fixed definition of music similarity so far. This is one of the first problems, dealing with music similarity. This topic offers multifaceted approaches. Merging multiple approaches with different weights can offer a more diverse music recommendation system. To do this, a lot of different data is required.

Content (music features) and context (listener behaviour) data can be fed into a big data framework to speed up operations. Collecting this data for large amounts of songs results in big datasets that need to be explored efficiently.

What to improve?

"[...] Spotify Radio, iTunes Radio, Google Play Access All Areas and Xbox Music. Recommendations are typically made using (undisclosed) content-based retrieval techniques, collaborative filtering data or a combination thereof." [1, p. 9] The goal of this work is, to propose a transparent music similarity retrieval method based on various weighted contextual and content-based data. Applying different weights to different features allows similarity retrieval methods to search for different kind of similarity. E.g. weighing the tempo and beat of a song more then melodic similarity allows the creation of playlists for workout and sport, whilst melodic/ timbre etc. similarities allows to search for similar songs from musical subgenres. The user would get to decide what kind of playlist he wants to create. Adding

contextual data and feeding it to an algorithm could add more or less popular music to the playlist with the goal to discover new upcoming artists or get other popular and most listened music. For this thesis however, the focus lies on content-based data/ audio features.

How to approach this?

First of all, a lot of data is required. In the first part, different scientific datasets are evaluated. Secondly, the available features are shown and explained. In the third part, different metrics are explained using the previously explored features. Lastly, a big data approach to efficiently use the gathered features is proposed and evaluated. A way to evaluate the results is also proposed.

1.2 Overview

Structure: Basically ->

PART I: Data Aggregation

PART II: Feature Extraction

PART III: Similarity Estimation

2. MIR Toolkits and Data Aggregation

2.1 Music Information Retrieval

The field of music information retrieval is a large research area combining studies in computer science like signal processing and machine learning as well as psychology and academic music study. The algorithm used in this paper is described in [2, pp. 17ff] and in section 2.1.1. The framework used to determine the song similarity in chapter 3 is the MUSLY toolkit [3] To get started, a brief overview is given in the next section providing the most important information about publicly available datasets, MIR toolkits and approaches to music similarity.

2.1.1 Audio Features and Music Similarity Measurements

This section provides an overview about different music similarity measurements and metrics. More in-depth information about the different metrics is given in chapter 6

Frequency based

As done by: [2, pp. 17ff]

Most of the algorithms start with switching from the time domain to the frequency domain, by performing a Discrete Fourier transform as described in equation 2.1 and then compute the power spectrum (equation 2.2)

$$X_m = \sum_{k=0}^{K-1} x_k \cdot e^{-\frac{K}{2\pi \cdot i} \cdot k \cdot m} \quad (2.1)$$

$$|X_m| = \sqrt{Re(X_m)^2 + Im(X_m)^2} \quad (2.2)$$

Figure 2.1a shows the spectrogram of the first bars of the song Layla by Eric Clapton. The sound sample was recorded on an electric guitar to avoid copyright infringements. Due to the fact that the human ear perceives sound in a non-linear matter, a logarithmic or Mel-scale is better to represent different pitches.

For example the note A4 is perceived at a frequency of 440Hz, the A note of next octave (A5) is at 880Hz and the next one is at 1600Hz and so on. The Mel-scale [1, pp. 53f] was introduced to resemble the human perception of frequency (equation 2.3)

$$m = 1127 \cdot \ln\left(1 + \frac{f}{700}\right) \quad (2.3)$$

The following plots were created with the librosa library [4].

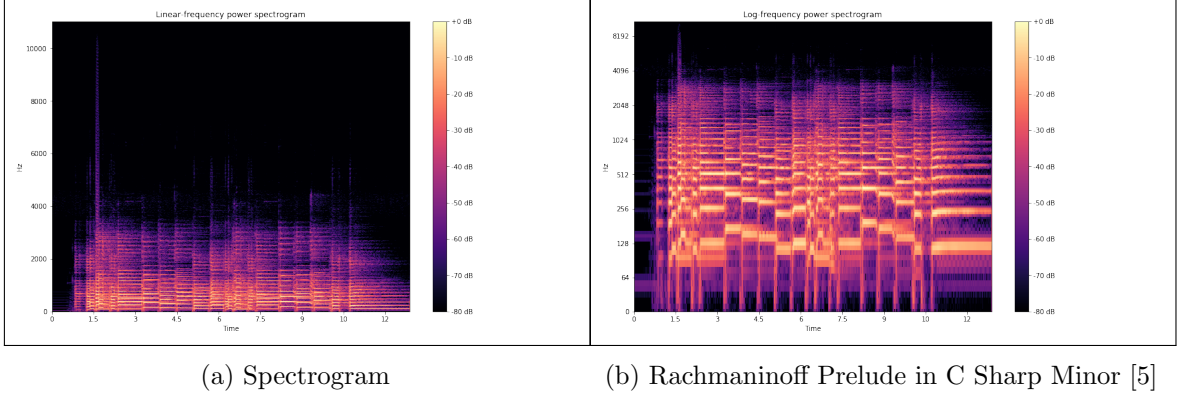


Figure 2.1: Frequency Space

The high dimensionality of the data is a problem for machine learning applications and music similarity tasks, as computation based on a vector with such a high dimensionality of the data would take too long. Given a sample rate of $f_s = 44,1kHz$ (usual CD sample-rate) and a length of a song of about $t = 180s$, the time domain contains 7938000 data points usually with 16-bit resolution.

$$K = f_s \cdot t \quad (2.4)$$

Calculating a FFT with a window size of 1024 samples and a hop size of 512 samples (resulting in the factor 1.5 in equation 2.5)[1], the full resulting spectrogram would contain 11627 frames with 1024 frequency values per frame. (eq: 2.5)

$$N_{fv} = 1.5 \cdot \left(\frac{44100 \text{ samples/s}}{1024 \text{ samples/frame}} \right) \cdot t \quad (2.5)$$

To reduce the dimensionality of the feature vector, a typical approach in MIR would be to calculate the so called Mel Frequency Cepstral Coefficients (MFCCs). They are described in more detail in the next section.

As another, better comprehensible feature set, the chroma features represent the tonal properties of a song. The chroma plot (Figure 2.3b) shows the distribution of the different pitches mapped to the various keys in one octave. The values are normalized to one by the strongest dimension. So if all values are close to one, it is most likely that there is only noise at that frame in the recording, as depicted in the first frames of figure 2.3b. Figure 2.3c figures the pitch curve of the recording. None but the most dominant frequencies are shown. Pitches below a certain threshold are filtered out. These Pitch curves can be used to estimate and transcribe musical notes from audio data as presented in 2.1.2. The Plot in figure 2.3d shows the onsets (beginning of musical notes) and estimated beats.

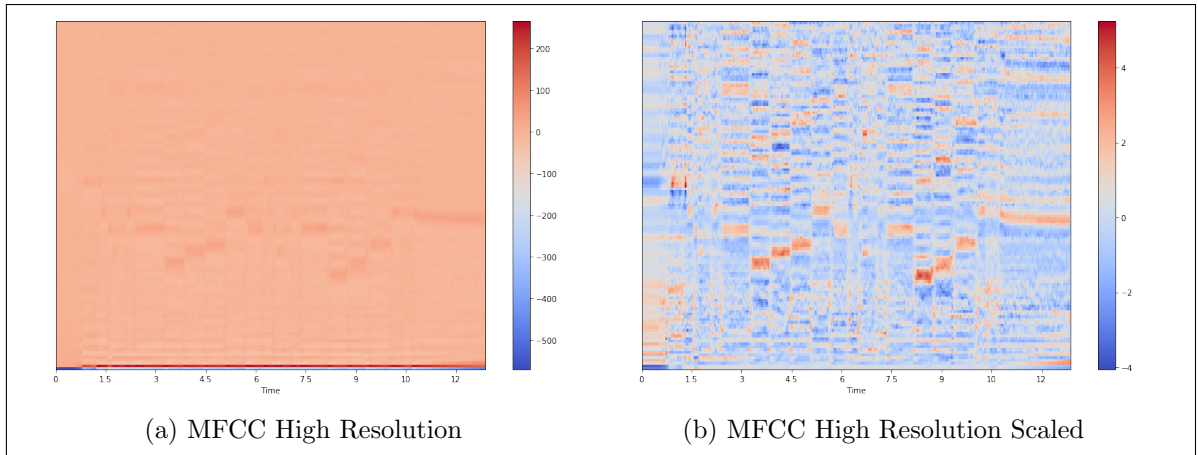


Figure 2.2: Features

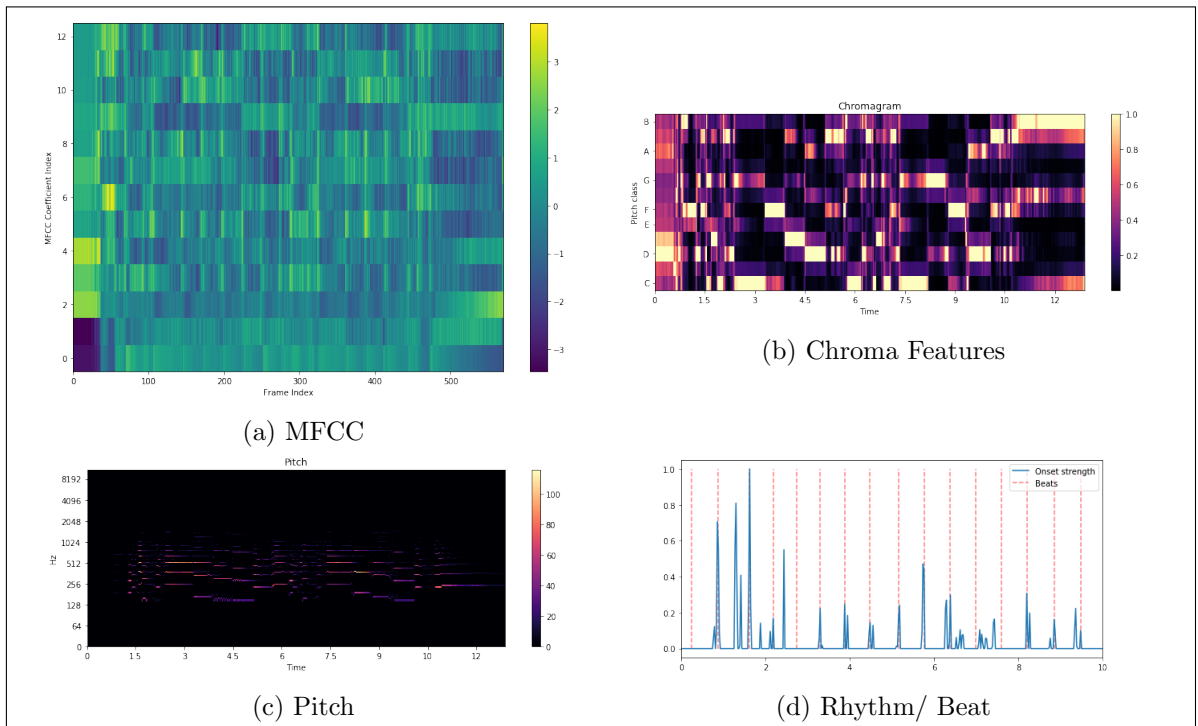


Figure 2.3: Features

MFCC

This section gives a brief overview over the computation of the MFCC as stated in [1, pp. 55ff]. Figure 2.4 shows the magnitude spectrum of a frequency sweep signal.

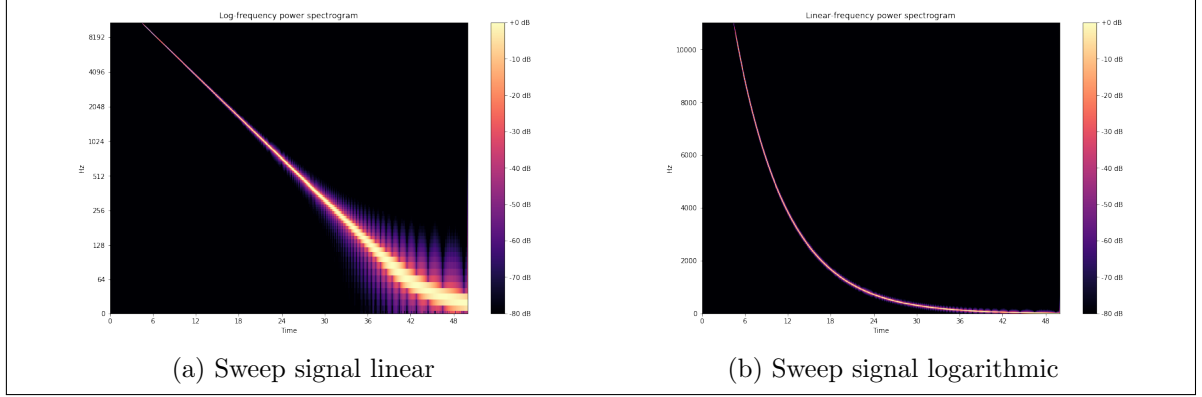


Figure 2.4: Sweep Sinal

First of all the magnitude spectrum is transformed to the Mel-scale by assigning each frequency value to a Mel- band. Doing this a dimension reduction can be done, by assigning multiple frequency values to one of typically 12 to 40 Mel-bands. The resulting vectors are then fed into a discrete cosine transformation (DCT) resulting in the MFCCs for each frame.

$$X_k = \sum_{n=0}^{N-1} x_n \cos \left[\frac{\pi}{N} \left(n + \frac{1}{2} \right) k \right] \quad (2.6)$$

Figure 2.5a shows the resulting MFCCs with a high resolution of 1024 Mel bands. This is not what in a usual application would be done, because this is nearly as high dimensional as the original spectrogram. Figure 2.5b shows the MFCC reduced to 12 Mel Bands.

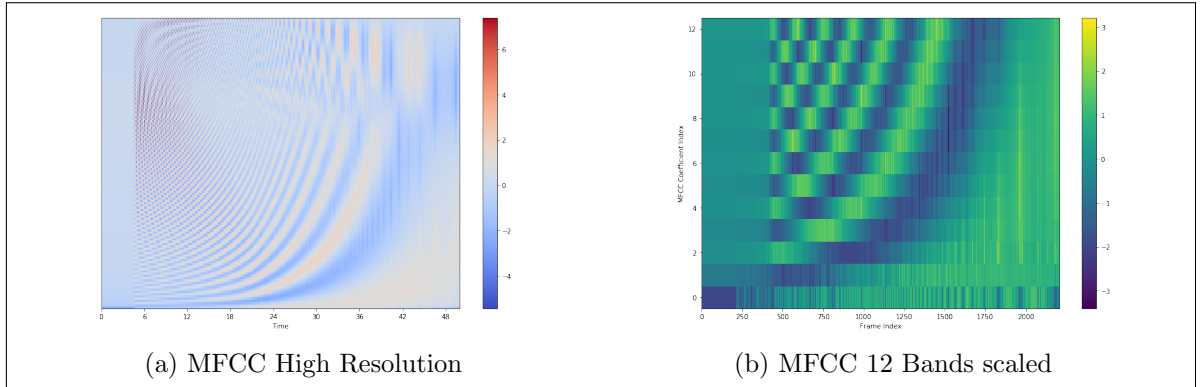


Figure 2.5: MFCCs

Rhythm based

As done by: [6]

Rhythm based music similarity algorithms use timing information of various events. The so

called rhythm spectrum can be used as a feature vector. For example the onset and beat data from the plot in Figure 2.3d could be used as a starting point for rhythmic similarity retrieval. An in depth overview is given in chapter 4

Pitch based

As done by: [10]

The proposed approach by [10] is, to take mid-level features like chroma or pitch instead of high-level features like sheet music or low-level features like gaussian mixture models of MFCCs. A more detailed analysis of this topic is given in chapter 5

Metadata based/ Collaborative Filtering

As done by: [7] and [8] using the Million Song Dataset (MSD) [9]

In 2012 the MSD Challenge was brought to the MIR community. The researches were challenged to give a list of song recommendations based on a large set of user data.

So if user X listened a lot to artist A and artist B and user Y listens mostly to artist A and artist C, then maybe user X would like artist C as well. These kind of metadata based recommendations are pretty common in large music streaming services, although not necessarily representing direct musical similarity. These kind of recommendation systems tend to propose commonly well known artists rather than not so well known ones biasing the result. On the other hand these kind of similarity algorithms can work very fast and efficient.

Note based

As done by: [11]

For comparing musical pieces by their symbolic representation (notes, tabulatures etc.) different text retrieval methods could be used. The MIDI datatype could be used, as it is a form of digital representation of music information. [11] uses a variation of the Levinshtein distance measurement. The problem with notation based algorithms is, that there are not many datasets available containing audio and MIDI information. As shown in 2.1.2 the automatic transcription of notes from raw audio does not work flawlessly. There is ongoing research to automatically annotate musical notes with the help of neural networks.[12]

Genre specific features

As done by: [13] for indian art music, by using 560 different combinations of different features. They state that: "We evaluate all possible combinations of the choices made at each step of the melodic similarity computation discussed in Section 2. We consider 5 different sampling rates of the melody representation, 8 different normalization scenarios, 2 possibilities of uniform time-scaling and 7 variants of the distance measures. In total, we evaluate 560 different variants" [13, p. 3]. This evaluation showed, that the choice of features and parameters for music similarity measurement is a critical point.

In Rock, Pop and Metal music, extraction of different guitar playstyles would be imaginable.

Guitar Tab Extraction [14] Toolkits could be used to extract information if the guitar in a song is mostly plucked or strummed for instance. Or if there are Hammer-on/ Pull-off/ side bending or tapping techniques used. In classical music, the play style of the string section of an orchestra could be taken into consideration.

2.1.2 MIR Toolkits

Music Similarity

The easiest way to test state of the art music similarity algorithms is to use the open source toolkit Musly [3]

It is based on gaussian mixture models of MFCC features and calculates the distances between songs very fast, supporting OpenMP acceleration. It offers the classical mandel-ellis similarity method [15] and a timbre based improved version of the mandel-ellis algorithm [16]

The MIR Toolkit [17] is a toolbox for Matlab [18]. A port to GNU Octave [19] is also available [20]

The short code snippet below is all it takes to compute a similarity matrix based on MFCC features, but the calculation takes quite some time.

```
mydata = cell(1, numfiles);

for k = 1:numfiles
    myfilename = sprintf( '%d.wav', k);
    mydata{k} = mirmfcc(myfilename);
    close all force
endfor

simmat = zeros(numfiles, numfiles);

for k = 1:numfiles
    for l = 1:numfiles
        simmat(k, l) = mirgetdata( ...
            mirdist(mydata{k}, ...
                mydata{l}));
    endfor
endfor
```

Audio features

To extract audio features, The YAAFE toolkit [21] is able to extract a lot of different audio features like energy, mfcc or loudness directly into the hadoop file format h5 making it ideal for big data frameworks to use. It can be used with C++, Python or Matlab.

The Essentia toolkit [22] is pretty similar to YAAFE, extending it by the calculation of the rhythm descriptors, bpm etc. It can also be used in C++ and Python

The Librosa Toolkit provides similar functionality [23] as Essentia. It is user-friendly and can be called from a Jupyter-Notebook [24], allowing rapid prototyping and testing of different algorithms. The plots from section 2.1.1 were created using librosa

Melody/ pitch extraction

To test the various pitch extraction toolkits, a piece by Rachmaninoff and by Beethoven was used. The first three bars of Rachmaninoffs Prelude can be found in figure 2.6a Figure 2.6b shows the first five bars of Beethovens Bagatelle in A Minor (Für Elise).

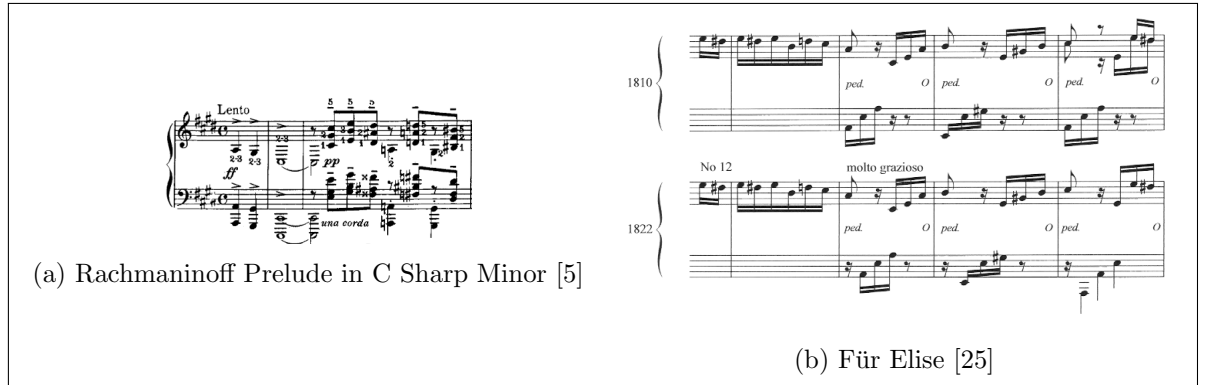


Figure 2.6: Original Scores

The first toolkit tested is Aubio [26]. The result can be seen in figure 2.7a and figure 2.7b. The upper subplot shows the waveform of the first few seconds of each piece. The second plot figures the estimated pitch with green dots. If the pitch is zero, then no pitch could be estimated, most likely because the associated frame contains silence. The blue dots resemble the estimated pitches, where the confidence (shown as the blue graph in the third subplot) is above a certain threshold (the orange line).

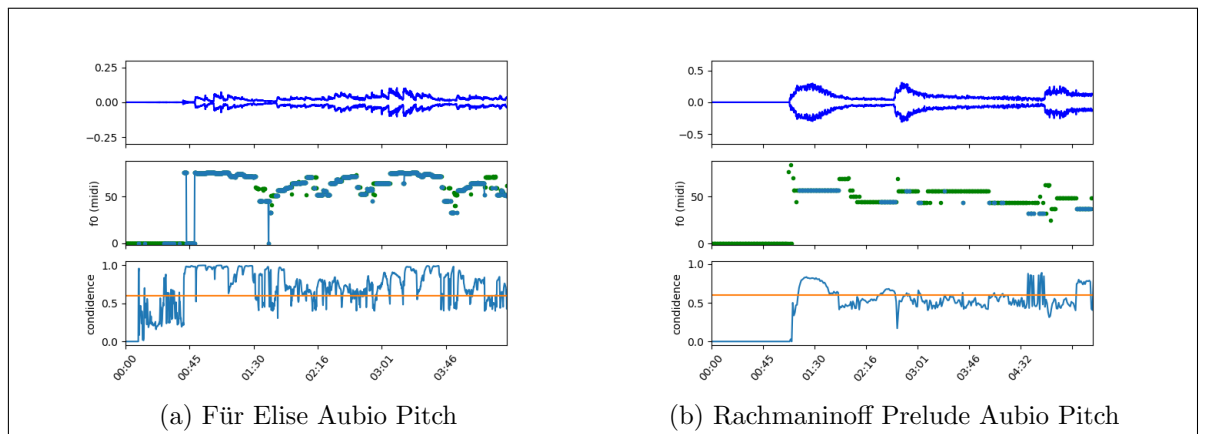


Figure 2.7: Aubio

The other melody extraction tool is Melodia[27], which is available as a VAMP plugin and can be used together with the Sonic Visualizer[28]

The results are shown in figure 2.8a and 2.8b. The purple line is the estimated pitch, however there are large jumps between different octaves of the harmonics.

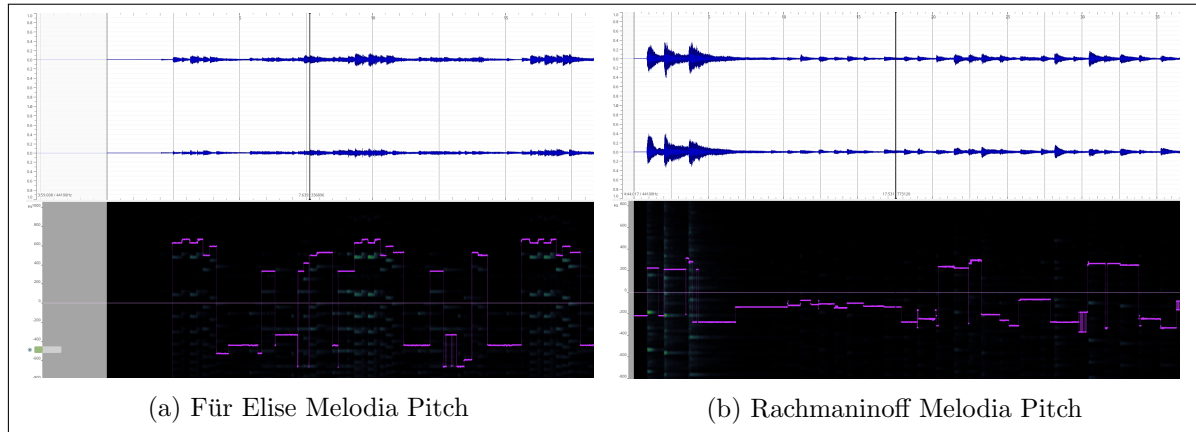


Figure 2.8: Melodia

Sadly the conversion to MIDI does not work flawlessly.

Figure 2.9a shows the output of a python script using the Melodia VAMP plugin to calculate a MIDI file containing the main melody line.

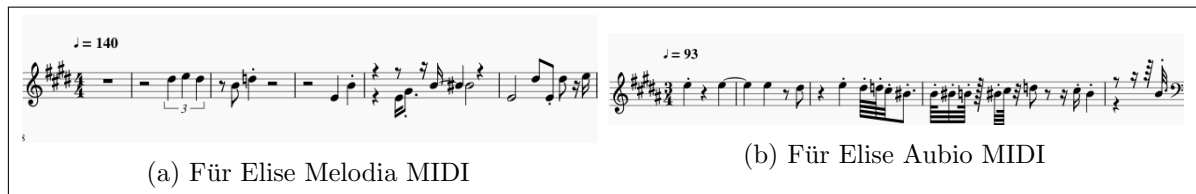


Figure 2.9: Transcription

Work in progress: A script to convert Aubio Pitch to MIDI and the output shown in figure 2.9b

Spotify API/ Echonest

Another way of getting music information, audio analysis and metadata is by using the Spotify API[29] Part of the available audio features comes from the Echo Nest[30] For each song available in Spotify, there are audiofeatures like loudness, speechiness, liveness, instrumentality, energy, danceability, acousticness, pitch, rhythm bars etc. With a small Python library named Spotipy, the available information can very easily be used and accessed. [31]

In figure 2.10a the pitch of the piano piece Für Elise by Beethoven is shown and figure 2.10b shows the beginning of the piece in more detail, including green dots, that resemble estimated bar markings. The blue dots represent the note values of one octave. That means they can resemble a value between zero and eleven with zero representing the key C and 11 is representing a B. The Spotify API actually returns a value for every single one of the keys per segment, while one segment is a section of samples that are relatively uniform in timbre and harmony. In the plots only the most dominant key per segment is shown. The pitch

values in one segment are normalized to one by their strongest dimension, so the more evenly distributed the pitch values are and the closer they all get to the value one, the more likely it is, that the sample contains only noise. The Downside using the Spotify API is, that there is no packed and ready to use test dataset containing the relevant features. So for scientific purposes, a test dataset would have to be created first.

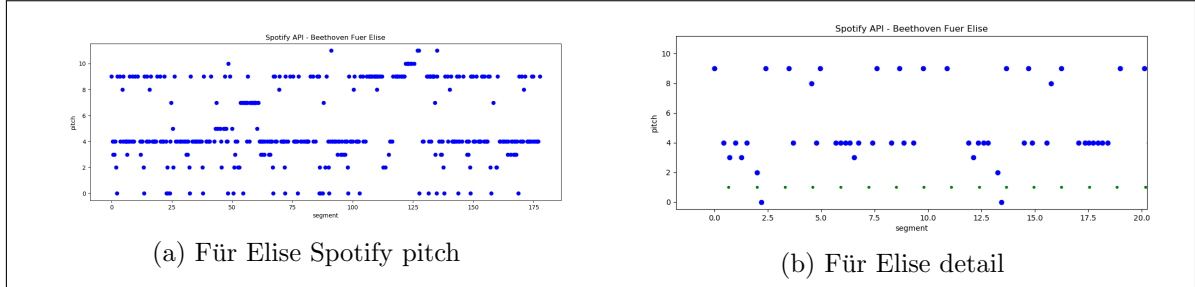


Figure 2.10: Spotify API

The pitch graph has one significant downside, because it is reduced to one octave and thus can not represent the melody of a song to its full extend.

Cover song identification

The ability to detect cover songs or different versions/ recording of a musical piece could be a good measurement for the efficiency of music similarity algorithms. In the next section one test case is presented, showing, that a MFCC based music similarity algorithm isn't able to detect different recordings of the same piano piece as most similar to each other.

2.2 Datasets

To evaluate the music similarity algorithms and metrics a lot of music data is needed. For this paper different sources of music were used.

2.2.1 Free Music Archive

The largest dataset is the Free Music Archive- dataset (fma) consisting of 106733 different songs totalling an amount of nearly one terabyte of music data from all kinds of different music genres.[32] There is also a lot of metadata information available for most of the songs.

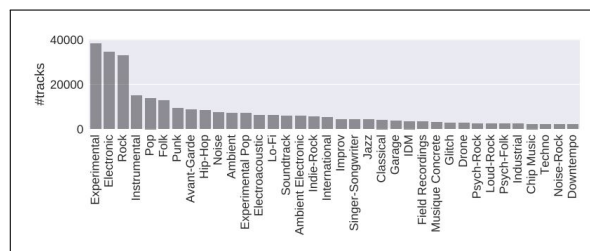


Figure 2.11: Number of songs in the free music archive dataset per genre

2.2.2 Musicnet and 1517 artists

Another source of music is the Musicnet dataset.[33] It includes 330 pieces of classical music with musical notes as annotations. Other sources of musical information would be the 1517-Artists dataset containing 3180 songs of multiple genres. [34]

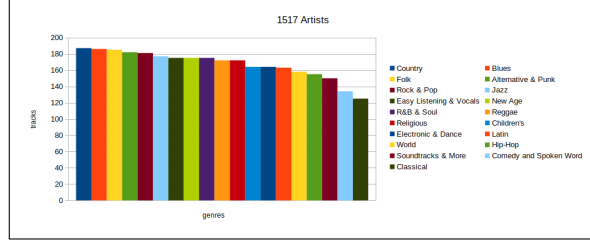


Figure 2.12: Number of songs in the 1517 artists dataset per genre

2.2.3 MedleyDB

For a melody/ pitch based similarity analysis, multitrack datasets could provide useful data, due to the fact that the pitch estimation can be done instrument by instrument. Datasets available are the MedleyDB[35] and MedleyDB2[36] datasets as well as the Open Multitrack Dataset[37] currently consisting of 593 multi-tracks in which the MedleyDB dataset is already included, leaving 481 other tracks for analysis.

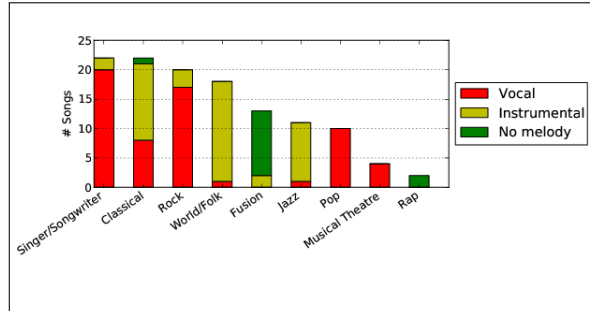


Figure 2.13: Number of songs in the MedleyDB dataset per genre

2.2.4 Private music collection

The private music collection used in this work consist mainly of metal music. The music was legally purchased, all rights belong to the respective owners. The distribution of different songs per genre for this dataset is visualized in figure 2.14.

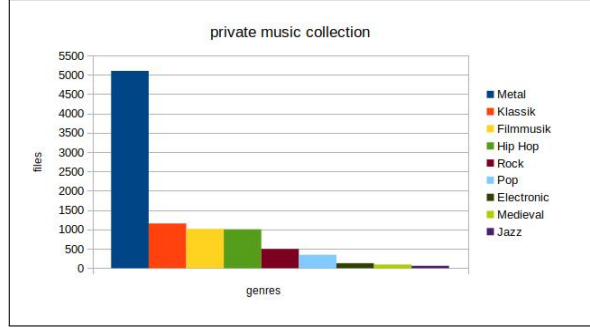


Figure 2.14: Number of songs in the private music collection per genre

Additionally a private recording dataset was used, consisting of ambient recordings and self produced music. Most of this music can be downloaded from soundcloud[38] The music sources and amounts of songs used for the task at hand is listed in table 2.1.

2.2.5 Overview and other sources

Table 2.1: Music Datasets used

fma	106.733 Songs
private	8484 Songs
1517 artists	3180 Songs
Maestro	1184 Songs (piano) + MIDI
musicnet	330 Songs (classical) + note annotation
Open Multitrack Testbed	593(481) Songs/ Multitracks
MedleyDB	122 Songs/ Multitracks
MedleyDB2	74 Songs/ Multitracks

A special and very large dataset is available with the Million Song Dataset (MSD)[9]. It contains a large set of metadata per track as well as a lot of supplementary datasets, like the Tagtraum genre annotation (figure 2.15)[39], the last.fm dataset[40] and the Echo Nest API dataset[41]. Although the MSD does not contain any music files in the first place, 30 second samples can be gathered through simple scripts from 7digital.com. On top of that the Echo Nest API data already contains a lot of audio features like pitch, loudness, energy and danceability to name just a few. Another addition is the secondhand dataset, containing a list of cover songs in the million song dataset[42]

Due to the fact that the Spotify API[29] also works with audio features from the Echo Nest[30], the MSD could be used in a big data environment to simulate the work with Spotify data, without manually mining the actual data. The MSD was actually already used in Big Data frameworks for music similarity retrieval based on metadata and user information[8]

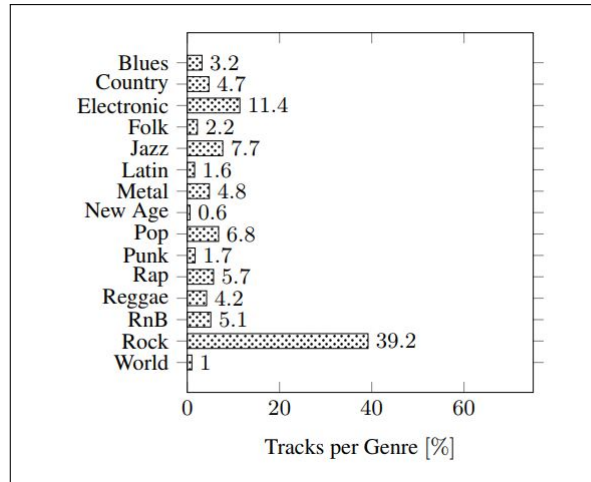


Figure 2.15: Number of songs in the million song dataset per genre

2.3 Spotipy

For the purpose of this thesis, the option of creating an own dataset using the spotify API and spotipy (section 2.1.2) was considered. Ten very small test playlists of different genres were created using the Spotify Playlist Miner [43]. Appendix 9.1 lists a small script, that downloads all audio features and analysis data from all of the songs of a playlist, that contains a preview URL with a 30 second audio snippet.

The audio features and analysis data is saved as a JSON file containing information about (beyond others)

- acousticness
- danceability
- instrumentalness
- liveness
- loudness
- speechiness
- valence
- predicted key
- tempo

as well as pitch and timbre information, beats and bars.

Together with the 30 second audio sample from which more data like MFCCs could be extracted, these crawler could provide all the information needed to build a large dataset for MIR. However the terms and conditions explicitly prohibits crawling the Spotify service. As stated by the Spotify Terms and Conditions of Use, section 9 (User guidelines):

"The following is not permitted for any reason whatsoever: [...]"

12. “crawling” the Spotify Service or otherwise using any automated means (including bots, scrapers, and spiders) to view, access, or collect information from Spotify or the Spotify Service;” [44]

Therefore a larger user created dataset can not be used without the risk of legal infringements. However one could argue, that there is a difference between data mining and data crawling and for small datasets with the purpose of creating Spotify playlists, these restrictions may not apply.

In the sense of the Spotify Developer Terms of Service [45] there may be no legal infringements by creating a non-commercial playlist creation tool. [46] states, that by creating algorithmically-generated playlists similar to the “Discover Weekly” Playlists one may run into challenges if using such features commercially. However it does not prohibit it for non-commercial cases.

2.4 Apache Hadoop and Spark

3. Timbre Similarity

MFCCs explained in chapter 1, here: how many melbands, visualize GMM

Proposed methods according to [1, pp.51ff]

Mel Frequency Cepstral Coefficients have already been introduced in chapter 2.1.1. This chapter focuses on the different similarity metrics.

To reduce the dimensionality of the data even further, a statistical summarization of the MFCC feature can be calculated [1, p. 58]. For each of the Mel- Bands (12 in this case) the mean and standard deviation over all frames is calculated, resulting in a vector of 12 mean values, a 12 by 12 co-variance matrix (actually $\frac{12*11}{2}$ values, because of the triangular shape - the upper triangle contains the co-variances and the main diagonal contains the variances) and 12 variances. These vectors are therefore not dependent on the length of the actual song. Using such a model, the distance between two songs can be calculated as in equation 3.1, where x and y are the n -dimensional feature vectors of two different musical pieces:

$$d(x, y) = ||x - y||_p = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}} \quad (3.1)$$

also known as the L_p distance. Most of the times, the Euclidean (L_2) or the Manhattan (L_1) distance would be used in real world scenarios. This very basic approach has been refined and improved over the past years.

3.1 Single Gaussian Model

This approach was first proposed by Mandel and Ellis [15] in 2005 and is briefly summarized in [1, pp. 65f].

After computing the mean value of each MFCC (μ_P and μ_Q) and the covariance matrix of the different MFCC vectors (Σ_P and Σ_Q) of two musical pieces P and Q , the Kullback-Leibler divergence (KL divergence) can be calculated as follows, with $Tr(\cdot)$ being the trace (i.e. the sum of the diagonal of a matrix), d being the dimensionality (number of MFCCs) and $|\Sigma_P|$ being the determinant of Σ_P

$$KL_{(P||Q)} = \frac{1}{2} \left[\log \frac{|\Sigma_P|}{|\Sigma_Q|} + Tr(\Sigma_P^{-1} \Sigma_Q) + (\mu_P - \mu_Q)^T \Sigma_P^{-1} (\mu_Q - \mu_P) - d \right] \quad (3.2)$$

As a second step the result has to be symmetrized.

$$d_{KL}(P, Q) = \frac{1}{2}(KL_{(P||Q)} + KL_{(Q||P)}) \quad (3.3)$$

This approach is one of the two available similarity metrics of the musly [3] toolkit introduced in section 2.1.2

3.2 Other Similarity Metrics

different similarity metrics, we use easiest one

The second available metric in the musly toolkit by Schnitzer is using the Jensen-Shannon Divergence (in an slightly adapted way). "The Jensen-Shannon (JS) divergence is another symmetric divergence derived from the Kullback-Leibler divergence. To compute it, a mixture X_m of the two distributions is defined" [2, p. 43]. "To use the Jensen-Shannon divergence [...] to estimate similarities between Gaussians, an approximation of X_m as a single multivariate Gaussian can be used [...] This approximation of X_m is exactly the same as the left-type Kullback-Leibler centroid of the two Gaussian distributions [...]" [2, p. 44]

$$\mu_m = \frac{1}{2}\mu_P + \frac{1}{2}\mu_Q \quad (3.4)$$

$$\Sigma_m = \frac{1}{2}(\Sigma_P + \mu_P\mu_P^T) + \frac{1}{2}(\Sigma_Q + \mu_Q\mu_Q^T) - \mu_m\mu_m^T \quad (3.5)$$

$$JS(P, Q) = \frac{1}{2}\log|\Sigma_m| - \frac{1}{4}\log|\Sigma_P| - \frac{1}{4}\log|\Sigma_Q| \quad (3.6)$$

After calculating a similarity matrix for all songs, it normalizes the similarities with mutual proximity (MP). [16] This method wants to reduce the effect of a phenomenon called hubness that appears as a general problem of machine learning in high-dimensional data spaces. "Hubs are data points which keep appearing unwontedly often as nearest neighbors of a large number of other data points." [2, p. 66].

Schedl and Knees state: "To apply MP to a distance matrix, it is assumed that the distances $D_{x,i=1..N}$ from an object x to all other objects in the data set follow a certain probability distribution; thus, any distance $D_{x,y}$ can be reinterpreted as the probability of y being the nearest neighbor of x , given the distance $D_{x,y}$ and the probability distribution $P(x)$ [...] MP is then defined as the probability that y is the nearest neighbor of x given $P(x)$ and x is the nearest neighbor of y given $P(y)$ " [1, p. 80]

Resulting in:

$$P(X > D_{x,y}) = 1 - P(X \leq D_{x,y}) = 1 - \mathcal{F}(D_{x,y}) \quad (3.7)$$

$$MP(D_{x,y}) = P(X > D_{x,y} \cap Y > D_{x,y}) \quad (3.8)$$

according to [1, p. 80]

Another, more compute heavy distance measurement would make use of Gaussian Mixture

Models of MFCCs. As Knees and Schedl state "Other work on music audio feature modeling for similarity has shown that aggregating the MFCC vectors of each song via a single Gaussian may work almost as well as using a GMM [...] Doing so decreases computational complexity by several magnitudes, in comparison to GMM-based similarity computations." [1, p. 65] Therefore the usage of GMMs is not further considered in this thesis.

The last method mentioned in this thesis for timbral similarity is to use block-level features as proposed by Seyerlehner [47] and described in short by Knees and Schedl [1, p. 67]. Instead of using single frames and summarizing them into statistical or probabilistic models, block-level features use larger, e.g. multiple second long, audio frames and features like fluctuation patterns are computed for these frames.

3.3 Validation

For this thesis the approach by Mandel and Ellis is used to compute the similarity matrix, because it is . There is always a trade-off of complexity and functionality. Using the musly toolkit, a first evaluation is presented in the next section. The feature extraction and distance calculation can also be done in python using the librosa library.

A possible measure for the efficiency of a timbre similarity algorithm is the ability to find songs of the same genre.

3.3.1 Construction Noise

Comparing a construction noise sound sample with the private music collection containing mostly metal, rock, pop, classical and hip hop music, the following six best results could be achieved in descending order:

- Ziegenmühlen Session - Down On The Corner (Folk Musik)
- While She Sleeps - The Divide (Metalcore)
- Delain - Mother Machine (Live) (Symphonic Metal)
- Within Temptation - Sanctuary (Intro Live) (Symphonic Metal)
- Without A Martyr - Medusa's Gaze (Death Metal)
- 100 Meisterwerke der Klassik - Orpheus In The Underworld (Orphée aux enfers) - Can-Can (Live At Grosser Saal, Musikverein) (Klassik)

Figure 3.1a and 3.1b show the distribution of the genres of 100 most similar songs compared to the construction noise sample.

Using the full dataset consisting of the private music collection, private field recordings, the full fma-dataset and the musicnet data, the following results could be achieved:

- Born Pilot - Birds Fell (FMA, Electronic, Noise)
- mrandmrsBrian - sun is boring (FMA, Avant-Garde, field recordings)
- steps in snow (private field recording)

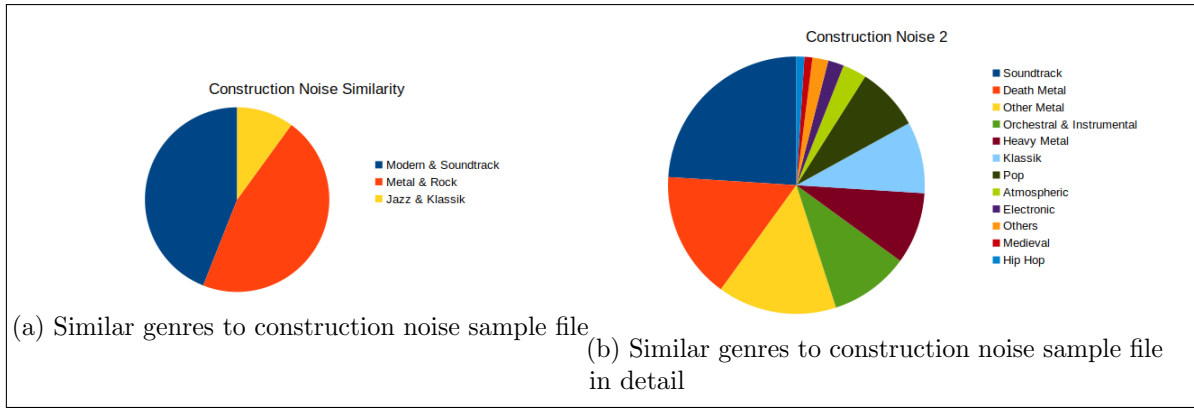


Figure 3.1: Construction Noise

- Sawako - Paris Children (FMA, field recordings)
- Jeremy Gluck and Michael Dent - Olivier (FMA, Ambient Electronic)

3.3.2 Different recordings and cover versions

Another experiment was, to get the most similar songs to the famous 'Rondo alla Turca' by Mozart. The recording used as a starting point is from the CD "100 Meisterwerke der Klassik" and has a length of 3:33 minutes. This piece by Mozart appears overall four times in the dataset and is recorded by different pianists. Every recording has a different length as listed in the following overview of the recordings by CD

- 100 Meisterwerke der Klassik (3:33)
- Piano Perlen (3:30)
- The Piano Collection - Disk 18 (3:28)
- Mozart Premium Edition - Disk 31 (4:29)

The top ten most similar songs to the 3 minutes and 33 seconds version are listed below:

- Mozart - Concert No. 10 for 2 Pianos and Orchestra in E Flat Major, KV 365 - 2. Andante
- Schubert - Sonata in B Flat, D. 960 - III. Scherzo (Allegro vivace con delicatezza)
- Albeniz - Iberia, Book I - Evocación
- Mozart Sonate Nr. 11 in A-Dur, K. 33 - Mozart - Alla Turca Allegretto (3:28)
- Beethoven - Bagatellen Op 119 -Allemande in D major
- Mozart - Rondo No. 1 in D Major, K. 485
- Mozart - Sonata For Piano No. 8 KV 310 A Minor - Allegro Maestoso
- Sonata For Piano No. 16 KV 545 C Major - Rondo: Allegretto
- Mozart Sonate Nr. 11 in A-Dur, K. 33 - III. Tuerkischer Marsch (3:30)

- Mozart - Piano Sonata No. 13 in B flat major, K. 333 (K. 315c): Allegretto grazioso

The interesting conclusion is that only 2 out of the 3 other versions were considered as most similar songs. The slower recording wasn't even in the top 30 list of the most similar songs. The same was observable for other cover versions of songs in the dataset like Serj Tankians song "Lie Lie Lie" from the CD "Harakiri" and an orchestral recording of the same piece. This is probably due to the usage of GMMs of MFCCs representing and valuing the timbre of the music predominantly instead of the pitches.

3.4 Conclusions

Why using a big data framework would help: music similarity is not well defined. It is a rather subjective value that differs from listener to listener. Two tracks could be considered as "similar" when they are equal in tempo, loudness, melody, instrumentation, key, rhythm mood, lyrics or a combination of more than a few of these features. The usage of a big data framework allows to create a variable/ fuzzy metric definition. Various parameters could easily be taken into consideration when calculating the musical distance of two different pieces. Using a Big Data framework, the problem of the fuzzy definition of music similarity could be avoided, if a metric can be found, that takes multiple of the accounted features of this paper into consideration. Available information includes metadata, user data, audio features, sheet music and more.

The idea of using genre specific features could be evaluated any further

Another important question is how to measure similarity algorithms. There are a few possibilities like genre, composer/ interpret or cover song identification. Or actual user data. MSD Challenge Dataset usable?[48]

TODO:

Comparison of Echo Nest Pitch Features vs. Full MIDI notation?

Using genre specific features

Use extracted MIDI vs. professionally annotated sheet music

Echo Nest Features invariant to key

4. Rhythmic Similarity

4.1 Cross Correlation

- simple approach using only onset and beat features extracted with librosa - cross correlation of 2 songs

dependent on beat extraction and onset detection algorithms

4.2 fluctuation patterns

4.3 Tempo Estimation

4.4 Beat and Tatum Estimation

5. Pitch estimation/ Chroma Features

TODO: pros and cons

bring all songs to standard key

usage (cover song identification)

text based (levenshtein distance)

different representations of melody (Chroma Features)

5.1 Pitch Curve

As presented in section 2.1.2 there are tools for pitch curve extraction of the main melody line. However in polyphonic music these kind of algorithms struggle to get reasonable results, even in popcultural music. In musical genres like Metal it gets even worse.

5.2 MIDI representation

5.3 Graph representation

5.4 Chroma Features

Chroma Features as described in section 2.1.1 are a good and low dimensional way to describe the melodic features of a song. The reduction of dimensionality however comes with a loss of information, especially what octaves the notes are played in.

5.4.1 Bandpass filter

5.4.2 Key invariance

5.4.3 Beat alignment

5.4.4 Validation

A good measure for the efficiency of a melodic similarity algorithm is the ability to find cover songs.

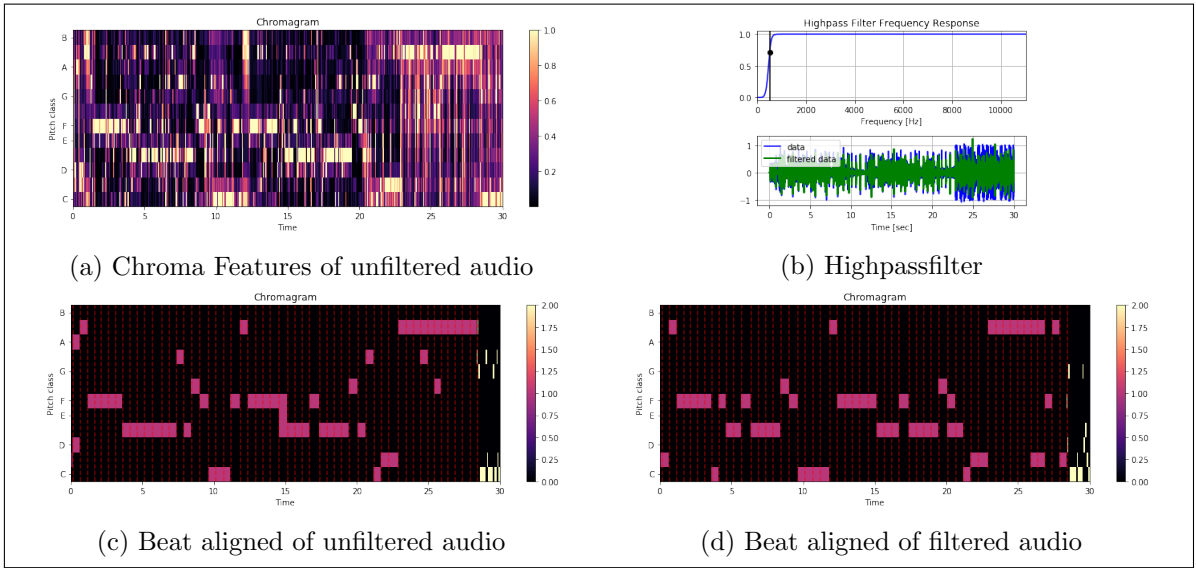


Figure 5.1: Sia Features

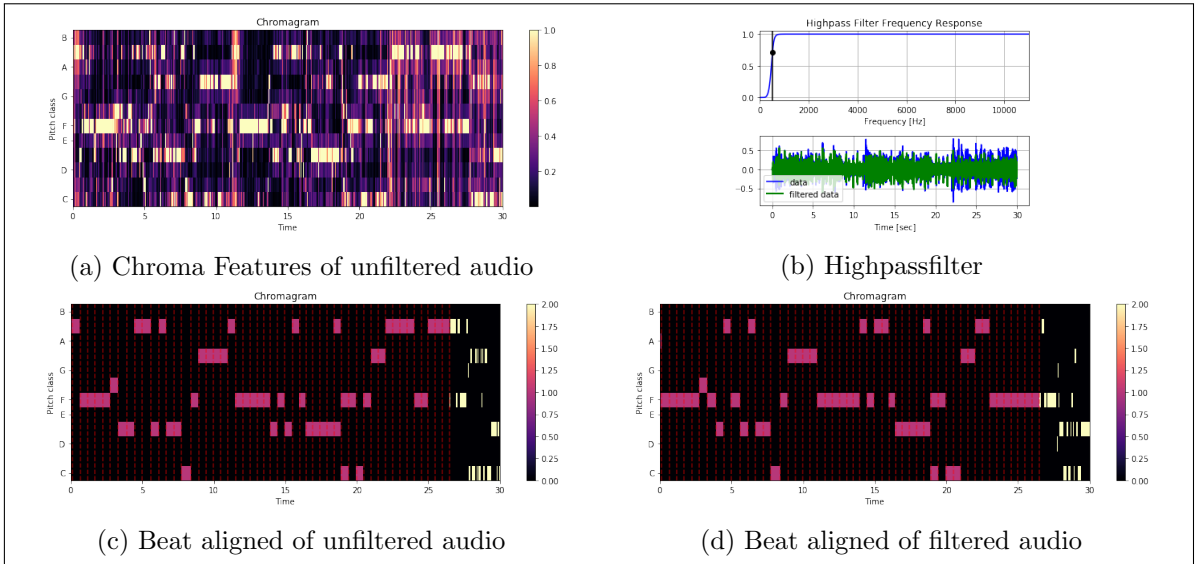


Figure 5.2: Pvriz Features

6. Other Similarity metrics

6.1 Definition of similarity

6.1.1 Rhythm, Timbre, Melodic

6.1.2 Genre and metadata

6.1.3 Genre specific features

6.1.4 Combinations and variable model

6.2 Collaborative Filtering

7. Big Data Framework Spark

7.1 Data aggregation

7.2 MapReduce tasks

8. Results

8.1 Cover song identification

8.2 Genre similarity

References

- [1] Peter Knees and Markus Schedl. Music Similarity and Retrieval. An introduction to web audio- and web-based strategies. In: Springer, 2016. ISBN: 9783662497203.
- [2] Dr. Dominik Schnitzer. Dealing with the Music of the World: Indexing Content-Based Music Similarity Models for Fast Retrieval in Massive Databases, 1st ed. In: 2012. ISBN: 9781477494158.
- [3] Dr. Dominik Schnitzer. Audio Music Similarity. In: URL: <http://www.musly.org/index.html>.
- [4] Brian McFee et al. LibROSA: Audio and Music Signal Analysis in Python. In: Proceedings of the 14th Python in Science Conference.
- [5] Prélude cis-Moll (Rachmaninow). In: URL: [https://imslp.org/wiki/Morceaux_de_fantaisie%2C_Op.3_\(Rachmaninoff%2C_Sergei\)](https://imslp.org/wiki/Morceaux_de_fantaisie%2C_Op.3_(Rachmaninoff%2C_Sergei)).
- [6] Jonathan Foote, Matthew Cooper, and Unjung Nam. Audio Retrieval by Rhythmic Similarity. In: Jan. 2002.
- [7] Zhao Yufeng and Li Xinwei. Design and Implementation of Music Recommendation System Based on Hadoop. In: 2018. DOI: [10.2991/icsnce-18.2018.36](https://doi.org/10.2991/icsnce-18.2018.36)..
- [8] B. McFee and G.R.G. Lanckriet. Large-scale music similarity search with spatial trees. In: ISMIR '11. URL: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.226.5060>.
- [9] Thierry Bertin-Mahieux et al. The Million Song Dataset. In: Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011). 2011.
- [10] Matija Marolt. A Mid-level Melody-based Representation for Calculating Audio Similarity. In: ISMIR 2006, 7th International Conference on Music Information Retrieval, Victoria, Canada, 8-12 October 2006. 2006.
- [11] Guangyu Xia et al. MidiFind: Similarity Search and Popularity Mining in Large MIDI Databases, 259-276. In: 2013. DOI: [10.1007/978-3-319-12976-1_17](https://doi.org/10.1007/978-3-319-12976-1_17)..
- [12] Jong Wook Kim et al. CREPE: A Convolutional Representation for Pitch Estimation. In: Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP), 2018.
- [13] Sankalp Gulati, Joan Serra, and Xavier Serra. An evaluation of methodologies for melodic similarity in audio recordings of Indian art music, 678-682. In: 2015. DOI:

10.1109/ICASSP.2015.7178055..

- [14] Cumhur Erkut et al. Extraction of Physical and Expressive Parameters for Model-based Sound Synthesis of the Classical Guitar. In: 2002.
- [15] M. Mandel and D. Ellis. Song-level features and support vector machines for music classification. In: Proceedings of the 6th International Conference on Music Information Retrieval, ISMIR, 2005.
- [16] D. Schnitzer et al. Using mutual proximity to improve content-based audio similarity. In: Proceedings of the 12th International Society for Music Information Retrieval Conference, ISMIR, 2011.
- [17] Olivier Lartillot and Petri Toivainen. MIR in Matlab (II): A Toolbox for Musical Feature Extraction from Audio. In: Proceedings of the 8th International Conference on Music Information Retrieval, ISMIR 2007, Vienna, Austria, September 23-27, 2007.
- [18] MATLAB. in: URL: <https://de.mathworks.com/help/matlab/index.html>.
- [19] John W. Eaton et al. GNU Octave version 4.2.0 manual: a high-level interactive language for numerical computations. In: 2016. URL: <http://www.gnu.org/software/octave/doc/interpreter/>.
- [20] Martin Ariel Hartmann. A port of MIRToolbox for Octave. In: 2016. URL: <https://github.com/martinarielhartmann/mirtooloct>.
- [21] B.Mathieu et al. YAAFE, an Easy to Use and Efficient Audio Feature Extraction Software. In: Proceedings of the 11th ISMIR conference, Utrecht, Netherlands, 2010.
- [22] D. Bogdanov et al. ESSENTIA: an Audio Analysis Library for Music Information Retrieval. In: International Society for Music Information Retrieval Conference (ISMIR'13). 2013, pp. 493–498.
- [23] Michael I. Mandel and Daniel P. W. Ellis. Labrosa's audio music similarity and classification submissions. In: 2007.
- [24] Jupyter. In: URL: <https://jupyter.org/index.html>.
- [25] Für Elise. In: URL: https://upload.wikimedia.org/wikipedia/commons/a/a9/BH\116_Vergleich.png.
- [26] Paul Brossier et al. aubio/aubio: 0.4.8 (Version 0.4.8). In: 2018. URL: <http://doi.org/10.5281/zenodo.1494152>.
- [27] J. Salamon and E. Gómez. Melody Extraction from Polyphonic Music Signals using Pitch Contour Characteristics. In: IEEE Transactions on Audio, Speech and Language Processing, 20(6):1759-1770. 2012.
- [28] C. Cannam, C. Landone, and M. Sandler. Sonic Visualiser: An Open Source Application for Viewing, Analysing, and Annotating Music Audio Files. In: Proceedings of the ACM Multimedia 2010 International Conference. Firenze, Italy, 2010, pp. 1467–1468.
- [29] Spotify API. in: URL: <https://developer.spotify.com/documentation/>.

- [30] The Echo Nest. In: URL: <http://the.echonest.com/>.
- [31] Spotipy - a Python client for The Spotify Web API. in: URL: <https://github.com/plamere/spotipy>.
- [32] Mich  el Defferrard et al. FMA: A Dataset for Music Analysis. In: 18th International Society for Music Information Retrieval Conference, 2017. URL: <https://arxiv.org/abs/1612.01840>.
- [33] John Thickstun, Zaid Harchaoui, and Sham M. Kakade. Learning Features of Music from Scratch. In: International Conference on Learning Representations (ICLR). 2017. URL: <https://arxiv.org/abs/1611.09827>.
- [34] Klaus Seyerlehner. 1517-Artists Dataset. In: 2010. URL: http://www.seyerlehner.info/index.php?p=1_3_Download.
- [35] R. Bittner et al. "MedleyDB: A Multitrack Dataset for Annotation-Intensive MIR Research". In: 15th International Society for Music Information Retrieval Conference, 2014.
- [36] R. Bittner et al. MedleyDB 2.0: New Data and a System for Sustainable Data Collection. In: New York, NY, USA: International Conference on Music Information Retrieval (ISMIR-16). 2016.
- [37] B. De Man et al. The Open Multitrack Testbed. In: 137th Convention of the Audio Engineering Society, 2014.
- [38] Soundcloud bqpd. In: URL: https://soundcloud.com/bq_pd.
- [39] Hendrik Schreiber. Improving Genre Annotations for the Million Song Dataset. In: In Proceedings of the 16th International Society for Music Information Retrieval Conference (ISMIR), pages 241-247, M  laga, Spain, Oct. 2015.
- [40] Last.fm dataset, the official song tags and song similarity collection for the Million Song Dataset. In: URL: <http://labrosa.ee.columbia.edu/millionsong/lastfm>.
- [41] The Echo Nest Taste profile subset, the official user data collection for the Million Song Dataset. In: URL: <http://labrosa.ee.columbia.edu/millionsong/tasteprofile>.
- [42] SecondHandSongs dataset, the official list of cover songs within the Million Song Dataset. In: URL: <http://labrosa.ee.columbia.edu/millionsong/secondhand>.
- [43] Spotify Playlist Miner. In: URL: <https://developer.spotify.com/community/showcase/playlist-miner-spotify/>.
- [44] Spotify Terms and Conditions of Use. In: URL: <https://www.spotify.com/lt/legal/end-user-agreement/plain/#s9>.
- [45] Spotify Developer Terms of Service. In: URL: <https://developer.spotify.com/terms/>.
- [46] Spotify Commercial Restrictions. In: URL: <https://developer.spotify.com/legal/commercial-restrictions/>.

- [47] Klaus Seyerlehner and Markus Schedl. Block-Level Audio Features for Music Genre Classification. In: 2009.
- [48] B. McFee et al. The Million Song Dataset Challenge, AdMIRe '12. In:

9. Appendix

9.1 Spotipy Crawler

```
from __future__ import print_function
from spotipy.oauth2 import SpotifyClientCredentials
import json, sys, spotipy, time, os.path
import requests, urllib
import matplotlib.pyplot as plt
import h5json, scipy
import numpy as np
from scipy.spatial import distance

reload(sys)
sys.setdefaultencoding('utf8')
client_credentials_manager = SpotifyClientCredentials()
sp = spotipy.Spotify(client_credentials_manager=client_credentials_manager)
if len(sys.argv) > 1:
    uri = sys.argv[1]
else:
    uri = 'spotify:user:bqpd:playlist:5oF8D71X38WwzeRUdyvpmd'
username = uri.split(':')[2]
playlist_id = uri.split(':')[4]
playlist = sp.user_playlist(username, playlist_id)
results = sp.user_playlists(username, limit=50)
playlist_length = playlist['tracks']['total']
path = os.getcwd()
path = path + "/crawled_data"
playlist_name = playlist['name']
directory = path + "/" + playlist_name
if not os.path.exists(directory):
    os.makedirs(directory)
t_start = time.time()
f_feat = open(path + "/" + playlist_name + "/featurevector.txt", "w")
```

```

f_feat.write("Features: \n")
f_feat.close()
feat_vec = []
feat_num = []
feat_name = []

for num in range(0, playlist_length, 100):
    results = sp.user_playlist_tracks(username, playlist_id, limit=100, offset=int(num))
    tracks = results
    for i, item in enumerate(tracks['items']):
        track = item['track']
        track_id = str(track['id'])
        path = os.getcwd()
        path = path + "/crawled_data"
        artist = str(track['artists'][0]['name'])
        songtitle = str(track['name'])
        artist = artist.replace("/", "")
        songtitle = songtitle.replace("/", "")
        artist = artist.replace("$", "")
        songtitle = songtitle.replace("$", "")
        number = i + num
        name = str(number) + " - " + artist + " - " + songtitle
        directory = path + "/" + playlist_name + "/" + name
        prev_url = track['preview_url']
        if not prev_url == None:
            if not os.path.exists(directory):
                os.makedirs(directory)
            filename = directory + "/" + artist + " - " + songtitle + ".mp3"
            urllib.urlretrieve(prev_url, filename)

        tid = 'spotify:track:' + track['id']

        analysis = sp.audio_analysis(tid)
        with open(directory + "/" + songtitle + '_analysis.json', 'w') as outfile:
            json.dump(analysis, outfile)
        outfile.close()
        segments = analysis["segments"]
        bars = analysis["bars"]
        beats = analysis["beats"]

        tid = str(tid)

```

```

features = sp.audio_features(tid)
with open(directory + "/" + songtitle + '_features.json', 'w') as outfile:
    json.dump(features, outfile)
outfile.close()
acousticness = features[0]['acousticness']
danceability = features[0]['danceability']
energy = features[0]['energy']
instrumentalness = features[0]['instrumentalness']
liveness = features[0]['liveness']
loudness = features[0]['loudness']
speechiness = features[0]['speechiness']
valence = features[0]['valence']
feat_vec.append(scipy.array([acousticness, danceability, instrumentalness, liveness,
loudness, speechiness, valence]))

else:
    print("no url - entry: " + artist + " - " + songtitle)
    print(track_id + "\n")

t_delta = time.time() - t_start
print ("features retrieved in %.2f seconds" % (t_delta,))
dist = distance.euclidean(feat_vec[0], feat_vec[1])

```


Acknowledgement

Declaration of Authorship

I hereby declare that the thesis submitted is my own unaided work. All direct or indirect sources used are acknowledged as references. This paper was not previously presented to another examination board and has not been published in German, English or any other language.

Hermisdorf, June 5, 2019

Johannes Schoder