

SANT'ANNA

DOCTORAL THESIS

Another music recommender system

Author:
Khoi Hoang NGUYEN

Supervisor:
Prof Monreale

*A thesis submitted in fulfillment of the requirements
for the degree of Doctor of Philosophy
in the*

Research Group Name
Management Department

August 31, 2017

Declaration of Authorship

I, Khoi Hoang NGUYEN, declare that this thesis titled, “Another music recommender system” and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

“Thanks to Google, Wiki, Stack Overflow for helping me with this monster”

Khoi Hoang Nguyen

Sant'Anna

Abstract

Faculty Name
Management Department

Doctor of Philosophy

Another music recommender system

by Khoi Hoang NGUYEN

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...

Acknowledgements

The acknowledgments and the people to thank go here, don't forget to include your project advisor...

Contents

Declaration of Authorship	iii
Abstract	vii
Acknowledgements	ix
1 Introduction	1
1.1 A statement of the problem	1
2 Background information	3
2.1 Overview of recommender systems	3
2.1.1 Collaborative recommendation	3
2.1.1.1 User-based nearest neighbor recommendation	4
2.1.1.2 Item-based nearest neighbor recommendation	5
2.1.1.3 Matrix factorization/ latent factor model	6
2.1.2 Content-based recommendation	8
2.1.2.1 Feature extraction	9
2.1.3 Hybrid recommendation	11
2.2 Overview of music recommender systems	11
2.2.1 Content-Based Music Recommendation	12
2.2.1.1 Metadata content	12
2.2.1.2 Audio content	13
2.2.2 Contextual Music Recommendation	14
2.2.2.1 Environment-Related Context	15
2.2.2.2 User-Related Context	15
A Appendix A	17
Bibliography	19

List of Figures

2.1	Projection of user and item on a two-dimensional space	7
2.2	Monolithic hybridization design	11
2.3	Parallelized hybridization design	11
2.4	Pipelined hybridization design	11

List of Tables

2.1 Rating for SVD-based recommendation	7
---	---

List of Abbreviations

LAH List Abbreviations **Here**
WSF What (it) Stands For

List of Symbols

a	distance	m
P	power	W (J s ⁻¹)
ω	angular frequency	rad

For/Dedicated to/To my...

Chapter 1

Introduction

1.1 A statement of the problem

Recommender system has been an interesting subject of research for long. The idea of recommender system began in early 1990s - with the popularization of the internet, to utilize the critique of millions of people online to help us acquire more useful and interesting content. The PARC Tapestry system [24] first introduced the novel idea of collaborative filtering technique in 1992, which instantaneously became a subject of interest for many other research groups and was employed in many systems, including the GroupLens system [59], the Ringo system at MIT [69], and the Bellcore Video Recommender [27].

Chapter 2

Background information

In this chapter, I will firstly make a survey of different types of recommender systems; then I will narrow down the topic to music recommender systems. As music possess different features comparing to other kind of information such as movie and news, its recommender system also needs to be tailor to bring satisfaction to its listeners.

2.1 Overview of recommender systems

There are currently three main basic approach to recommendation, namely collaborative recommendation, content-based recommendation, and knowledge-based recommendation [30]. There are also an approach, called hybrid approaches, that tries to combine different recommendation together, in order to augment the strength and limit the drawback of each separate techniques. The description, as well as the advantages and disadvantages, of each recommendation techniques will be discussed as follow.

2.1.1 Collaborative recommendation

The main idea of collaborative recommendation approaches is to predict potential items that a user would like using past behavior information of other users. Pure collaborative algorithm take only user-item ratings as input and generate a prediction suggesting to what degree the user will like a certain item or a list of n recommended items as output.

There are two kind of ratings that can be used, namely implicit and explicit ratings. Explicit rating information can be collected by explicitly asking users to rate the item on a specific scale. Different scales are applied to different domain, as the quality of recommendation is different between these scales [17]. Ratings are then convert internally to numeric values in order for recommender systems to calculate the similarities. Implicit ratings, on the other hand, are knowledge collected based on the interaction between users and the systems. They can be in various forms with distinctive characteristics, such as information about item buying, book reading, music listening, or even user browsing behavior. As implicit ratings are observed behaviors, recommenders have to interpret whether the behaviors have positive or negative impacts toward the users. Even though the interpretation might be incorrect in some cases (e.g., a user might not like all the items that she bought), a massive amount of feedback would exclude these particular cases. In fact, Shafer et al [63] report that in some domains, user model using implicit information can outperform the one with explicit ratings.

There are many algorithms develop to exploit the rating matrix; However, in this review, I will just go into detail some main approaches that have been studied

carefully in the past and have been applied widely in the industry. These approaches include user-based approach, item-based approach, and an approach using matrix factorization/ latent factor models. Some other approaches, such as probabilistic approach, or slope one predictors, will not be mentioned here because of the scope of the thesis.

2.1.1.1 User-based nearest neighbor recommendation

User-based nearest neighbor recommendation is one of the earliest algorithm for this approach. The main idea of the algorithm is to identify other users who have similar preferences to a user; then for any item i that is unknown to the current user, a prediction is given based on the similar of the ratings of other users on item i .

2.1.1.1.1 Pearson's correlation coefficient

One common method to calculate the similar between users is Pearson's correlation coefficient. The general formula for the coefficient ρ [8] is:

$$\rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y}$$

with $\text{cov}(X,Y)$ is the covariance between X and Y , and σ_X and σ_Y represent the standard deviation of X and Y .

Apply the formula to the case of collaborative filtering: let $U = \{u_1, \dots, u_n\}$ denote the set of users, $P = \{p_1, \dots, p_m\}$ for the set of items. The 2 sets form a $n \times m$ matrix of rating $r_{i,j}$ with $i \in 1 \dots n, j \in 1 \dots m$, with $r_{a,p}$ is the rating of user a on item p , and \bar{r}_a denotes the average rating of user a . The similarity of user a and b $\text{sim}(a,b)$ is defined as follow:

$$\text{sim}(a,b) = \frac{\sum_{p \in P} (r_{a,p} - \bar{r}_a)(r_{b,p} - \bar{r}_b)}{\sqrt{\sum_{p \in P} (r_{a,p} - \bar{r}_a)^2} \sqrt{\sum_{p \in P} (r_{b,p} - \bar{r}_b)^2}}$$

The Pearson correlation coefficient takes value range from -1 to +1, indicating respectively from a strong negative correlation to a strong positive correlation.

2.1.1.1.2 Other weighting metrics

Apart from Pearson's correlation coefficient, other metrics, such as adjusted cosine similarity, Spearman's rank correlation coefficient, or mean squared difference measure are also proposed to calculate user-based similarity. However, empirical study made by Herlocker et al [26] shows that for user-based recommender system, the Pearson coefficient outperforms other measures.

Still, the "pure" Pearson measure alone is not ideal as there are cases that the measure cannot handle. Consider a real life problem that there are some items that are favored by everyone, Pearson's measure would not consider that an agreement by two users on a controversial item has more weight than an agreement on a universally like item. Herlocker et al [26] also showed that applying the measure to user who has rated very few items also lead to bad predictions. Therefore, many attempts, such as significance weighting proposed by Herlocker et al [26], or case amplification suggested by Breese et al [12] have been made to fill the gap and improve the accuracy. However, the question of whether these weighting schemes are helpful in real-world settings is still opened.

2.1.1.1.3 Challenges

Although user-based approaches have been deployed successfully, they faces serious challenges when are applied to large e-commerce sites which possess millions of users and items. Specifically, the cost for scanning a vast number of potential neighbors makes it impossible for the system to predict in real time. Therefore, large scale e-commerce sites often opt for other techniques, one of them is the item-based nearest neighbor approach.

2.1.1.2 Item-based nearest neighbor recommendation

The main idea of this approach is to calculate the similarity between items instead of one between users. The advantage of this approach over user-based approach is that we can preprocess an item similarity matrix that characterize the degree of similarity between items. At run time, a prediction for product p and user u is made by detecting the most relevant items using the item similar matrix and by calculating the weighted sum of u 's rating for these items. As the number of relevant items is commonly limited, the computation of the prediction can be done within a short time frame, suitable for such online applications. A similar matrix is, theoretically, also possible with user-based approaches; however, in real time scenarios, the number of overlapping ratings for two random users is relatively small, making it unstable as a few more ratings may significantly affect the similarity between users.

For item-based approaches, the cosine similarity is found to be the standard metric [30]. The cosine similarity is defined as follows:

$$sim(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| * |\vec{b}|}$$

where \vec{a} is an item vector, \cdot denotes the dot product, and $|\vec{a}|$ is the Euclidian length of the vector, which is defined as the square root of the dot product of the vector with itself.

One drawback of cosine measure is that it does not take into account the fact that different users have different rating schemes, i.e., some users rate items highly in general, while some others give lower ratings. This drawback is solved using adjusted cosine measure, which subtracts the user average from the rating. Let U be the set of users that rate both item a and b . The adjusted cosine measure is as follows:

$$sim(a, b) = \frac{\sum_{u \in U} (r_{u,a} - \bar{r}_u)(r_{u,b} - \bar{r}_u)}{\sqrt{\sum_{u \in U} (r_{u,a} - \bar{r}_u)^2} \sqrt{\sum_{u \in U} (r_{u,b} - \bar{r}_u)^2}}$$

with \bar{r}_u is the average rating for item u .

The prediction function, which is computed in real time, is defined as follows:

$$pred(u, p) = \frac{\sum_{i \in ratedItem(u)} sim(i, p) * r_{u,i}}{\sum_{i \in ratedItem(a)} sim(i, p)}$$

Compare to user-based recommendation, item-based approaches prove to be often more scalable as for most of the case, the number of items falls behind the number of user, thus it requires less time and space to compute the similarity matrix (if necessary). Also, item-based approaches are more justifiable by users, for they can easily grasp the explanation of the recommendation, modify the list of neighbors and alter the weights. User-based methods, on the other hand, are less able to justify, as recommendations come from other users are hard to explain. However,

as item-based approaches are based on ratings on similar items, the recommender tend to suggest items that might be already familiar to that user. While this behavior leads to safe recommendations, it does not help users explore other novel items that they might like as well.

In general, nearest neighbor approaches work well for popular items. Nonetheless, there are two important drawbacks with these approaches when deal with unpopular item:

- Limited coverage: both approaches define neighbor as having ratings in common. This assumption is limiting, as users with very few common items can still have similar preferences. Moreover, coverage of such approaches can be limited, as only items rated by neighbor can be recommended.
- Sensitive to sparse data: For most system, users only rated a small portions of available items. This results in a cold star problem, when some items have very few or no ratings at all, which affects the prediction of these items. Another problem is when there are only few ratings, the weight of each item has a significant impact over the similarity between vectors, reducing the accuracy of the recommender.

To solve these problem, many small tunes for the neighborhood approaches have been made, such as the use of Significance Weighting [] for the weighting problem or ... for the Cold Start problem []. Besides that, other advance algorithms have also been developed to tackle these topics. One of the most popular approach that has gained attraction recently is the use of matrix factorization, as it was exploited to significantly improve the accuracy of the recommender system in the Netflix Prize competition in 2009.

2.1.1.3 Matrix factorization/ latent factor model

Matrix factorization methods is used to derive a set of salient patterns from user-ratings. For example, let "Gone with the wind" and "Me before you" be the set of liked item of user A, and "Romeo and Juliet" and "The fault in our star" be the set of liked item of user B, while nearest neighbor approaches would consider these books separately, matrix factorization could see that all these books belong to romantic genre and therefore recommend them to the other user. The factors, however, is not always obvious. In some cases, they can be uninterpretable.

The idea of this method is to factorize the original sparse matrix into a product of matrices. Each decomposed matrix is much denser than the original one. The technique can be used for both similarity matrix and rating matrix. There are many matrix factorization techniques with increasing complexity and accuracy. For the scope of this thesis, I will describe the Singular value decomposition (SVD) model, a basic yet effective decomposition technique. The example is an adaptation from the one from Grigorik [72].

Consider the table 2.1:

SVD takes a m -by- n matrix M and decomposes it into three factors: two unitary matrices U and V , which represent the user matrix and the item matrix accordingly, and a nonnegative diagonal matrix Σ . The main point of this decomposition is that after receiving the product matrix, we can retain only the most important values in the diagonal matrix to build back the approximation of the original matrix. The decomposition is as follows:

	User 1	User 2	User 3	User 4
Item 1	3	4	3	1
Item 2	1	3	2	6
Item 3	2	4	1	5
Item 4	3	3	5	2

TABLE 2.1: Rating for SVD-based recommendation

$$M = U\Sigma V^T$$

where V^T is the transpose matrix of V .

Applying the decomposition, we obtain $\Sigma = 12.2215, 4.9282, 2.0638, 0.2977$ and the two matrices

U				V			
-0.4312	0.4932	-0.5508	-0.5172	-0.3593	0.3677	-0.2961	0.8050
-0.5327	-0.5305	0.4197	-0.5085	-0.5675	0.0880	-0.6285	-0.5246
-0.5237	-0.4052	-0.4873	0.5693	-0.4429	0.5686	0.6590	-0.2150
-0.5059	0.5578	0.5321	0.3871	-0.5939	-0.7306	0.2882	0.1746

In this case, we can eliminate the two dimensions with the lowest values and only keep the two dimensions with value $\Sigma = 12.2215, 4.9282$. Figure 2.1 is the projection of the two matrices in a two-dimensional space.

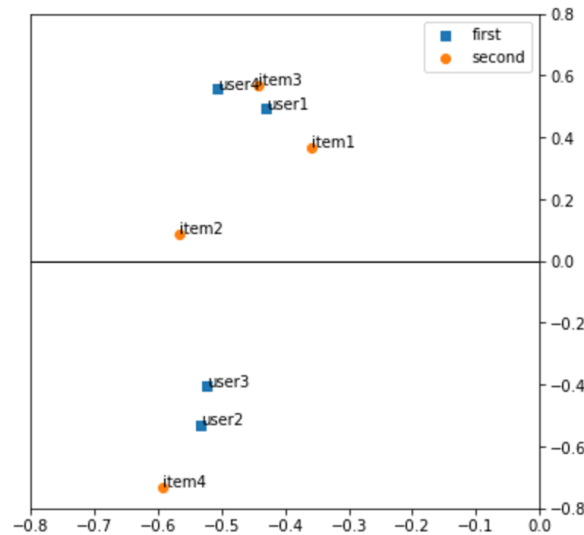


FIGURE 2.1: Projection of user and item on a two-dimensional space

After the product matrices are constructed, a new user profile can be added by multiply user's rating with the user matrix U and the diagonal matrix Σ . Once the user profile is constructed, many strategies can be used for the recommendation of item to that user. One possible approach is to again apply cosine similarity measure

to find neighbor user. Another approach is to approximate user rating by exploiting the interaction between user and item in the latent factor space [34].

2.1.2 Content-based recommendation

Although collaborative filtering approaches work well, it still have some limitations. Apart from the cold star problem and the popularity bias, collaborative approaches require lots of user rating for the system to be stable. Besides, these approach cannot take advantage of the semantic content of the item. For example, it would be obvious to recommend "Gone with the wind" to a user, if we know that (a) this book belongs to romantic catalog and (b) the user has an affair for romantic novel. Therefore, as data is becoming more abundant, new approaches have been studied to exploit these data for the recommendation process. These approaches are commonly called content-based recommendation.

The basic idea of content-based approaches is to classify similar items using a list of features of each item. For example, a book recommender could use the meta-data of the book, such as the author's name and the book's genre as a reference for similarity; or it could calculate the similarity of two books based on the overlap of keywords, using the Dice coefficient [61] as follows:

$$sim(b_i, b_j) = \frac{2 * |keywords(b_i) \cap keywords(b_j)|}{|keywords(b_i)| + |keywords(b_j)|}$$

with $|keywords(b_i)|$ is the number of keyword in book b_i .

However, metadata also need the clarification of expert. The standard approach in content-based recommendation is, not to maintain a list of metadata, but to seek for the similarity in the content of the item itself. In the case of book recommender, one common approach is to transform the content of a book into a vector in a multi-dimensional Euclidian space using *term frequency-inverse document frequency (TF-IDF)* technique. For the scope of this thesis, the detail of the technique is not described here.

Music recommender systems, however, take a different approach compare to other text recommender sytems, as the differences in characteristics between the two items. Metadata of a music track can include general tags generated by expert such as genre, artist, album, etc..., as well as other social tags by users; and the content of a track is the acoustic features that are analyzed from the signal of that track, of which the most representative ones are timbre and rhythm [14].

Basically, a content-based recommender has three basic components [60]

- Content analyzer: the mission of this component is to extract relevant information from the content of the item. In this phase, feature extraction techniques are used to exploit information and transform them to the representation that the recommender needs. The representation is the input to the *profile learner* and *filtering component* parts.
- Profile learner: in this phase, the module tries to construct user profiles using the representation of the item. Often, this phase is achieved through the use of various kind of machine learning techniques [47], depending on the nature of input data. User profiles are then passed into filtering component for generating list of recommendations.
- Filtering component: This phase generates a list of recommendation for the users by comparing the similarity between the representation of user profiles

and that of items to be recommended. After recommendation is suggested, the system might get feedback of the users for the profile learner phase to improve the user profiles.

A survey of the techniques used in content-based recommender will be discussed later. For the moment, some music features as well as the techniques that are used to extract them will be described. The profile learner and filtering component parts will not be detailed here, as they are beyond the scope of the thesis. However, the specific algorithms used for these part will be described later in chapter 4.

change number of chapter if necessary

2.1.2.1 Feature extraction

Music, in its original form, is a record of analog audio signal (i.e, electrical voltage) [28]. To store music digitally, a analog-to-digital conversion is needed. The process has two phases: sampling and quantization [32]. In the sampling part, a signal is sampled by evaluating its amplitude at a particular time, with the number of samples taken per second is called the sampling rate. The amplitude measurements are then mapped as 8 or 16-bit integers, whose process is called quantization.

Because digital music is represented as a series of integers, the raw information they contain is trivial for human at the perceptual level. Therefore, the first step of a music recommender is to extract useful features from the raw representation. In music domain, different taxonomies have been created to capture audio features in certain perspectives. Weihs et al. [78] divided audio features into four subcategories, including short term features, long term features, semantic features, and compositional features. Another taxonomy was proposed by Scaringella [62], which separates audio into three different components: timbral to denote features related to spectral content (i.e. shape) of the signal; temporal features such as loudness, tempo, onset rate; and tonal component such as harmonic, pitch, key, scale, and chords distribution.

Fu et al. [22] developed another taxonomy under human understanding of music perspective. According to the hierarchy, there are two level of audio features: low-level and mid-level features. Low-level features are features that can be obtained directly from the audio with proper signal processing techniques like Fourier transform, spectral analysis, autoregressive modeling, etc. Timbral and temporal are two main classes in low-level features. Low-level features have been exploited massively in music classification, due to the simple procedure to obtain and their good performance. However, they are not closely related to the nature of music that human perceive. Mid-level features, on the other hand, delimitate music using rhythm, pitch, and harmony, concepts that are more familiar to normal listeners.

rephrase

2.1.2.1.1 Low-level features

Timbral is perhaps the most exploited characteristic in the set of low-level features. It is described as the quality of sound, with different timbres belong to different types of sound sources, e.g. different instruments. Table 2.2 lists some major timbre features. Despite of the large variety in number of features, the timbre extraction phase of all these features are closely related to each other and follow some standard procedures. As this thesis does not deal with low level signal processing, the detail of these procedures is not described here.

Timbral features have been successfully used in the past for genre classification. However, it still has several disadvantages. One critical problem is that with the

Class	Feature Type	Used in
Timbre	Zero Crossing Rate	[74], [37], [9]
	Spectral Centroid	[74], [37], [9], [48]
	Mel-frequency Cepstrum Coefficient	[74], [9], [43]
	Fourier Cepstrum Coefficient	[9], [40]
	Stereo Panning Spectrum Features	[75], [76]
Temporal	Statistical Moments	[74], [37]
	Amplitude Modulation	[53], [51]
	Auto-Regressive Modeling	[70], [46]

technological advance in the recording process (e.g., tape editing, equalization, and compression), the final signals of the post production stage of the same instrument in two different tracks would be different. As timbral features techniques rely heavily on raw signal, they are heavily affected by the recording process, of which phenomenon called album effect [whitman2001artist].

Temporal features are another low-level features that are used to capture the temporal evolution of the signal. The difference between the two features is that, while each timbral is extracted in a local window of raw signal with 10- 100 ms duration, temporal extraction is performed on a series of timbre features in larger frames, allowing for the derivation of a richer set of features, such as fluctuation pattern [53], rhythmic pattern [38], rhythmic coefficient (west2009novel, etc. Morchen [48] has generated a set of operation that can be employed on top of timbre to produce new features at a coarser scale. Many of these features belong to the amplitude modulation family, as they are generated by analyzing the modulation of the amplitude spectrum.

2.1.2.1.2 Mid-level features

Low-level features were dominantly applied for genre classification; however, they do not capture the intrinsic properties of music that humans perceive. Mid-level features, such as rhythm, pitch, or harmony, are more familiar to human. Therefore, they play an important role in some specific domains, such as query by singing [29] or detect cover versions of popular songs with similar melodies [73]

Rhythm is the most widely used mid-level feature in audio-based music classification. It describes the recurrence of tension and release in music. From rhythm, "danceability" of a track can be derived by . Rhythm can also be used for mood classification [20] [79], since sad songs usually have a slow rhythm, while exciting songs usually possess a fast rhythm.

Pitch is another important mid-level feature. It is defined as the most fundamental frequency of the sound. Often, a pitch histogram is constructed and is combined with low-level feature for genre and mood classification [74] [36]. Pitch can also be used to develop pitch class profile and harmonic pitch class profile, which in turn are useful for detecting similar melodies and transcription [44] [56].

Harmony is a succession of musical chords, which are three or more notes, typically sounded simultaneously. The chord is detected by compare pitch histogram with chord template to identify the possible chords. Chord features are often used as a complimentary to pitch features in detecting melody similarity and cover song [19] [7]

To summarize, low-level audio features such as timbre are sufficient for genre classification but fail to achieve good result in song similarity detection. Mid-level

features, on the other hand, are successful in detecting similar song using pitch and harmonic features. Rhythm features have also been exploited for mood classification. Research in using acoustic contents for recommendation problems will be describe later in the overview of music recommender systems part.

2.1.3 Hybrid recommendation

Hybrid recommendation is a method that combine several kind of recommenders together. The motivation for hybrid recommender is to try to take advantage of the strength of different algorithms with fewer drawbacks than any individual one. Burke's well-known taxonomy [13] differentiate between seven kind of hybridization strategies; However, the seven approaches can be abstracted into three base design: monolithic, parallelized, and pipelined hybrids [30]. Monolithic design incorporates several recommendation strategies into one algorithm, while parallelized and pipelined designs require at least two separate recommenders. The outlines of the three designs are depicted in figure 2.2, 2.3 and 2.4

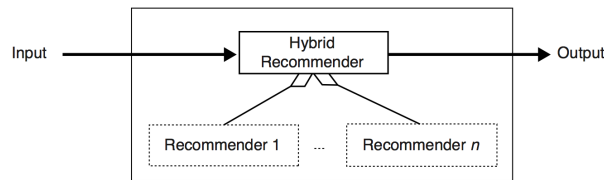


FIGURE 2.2: Monolithic hybridization design

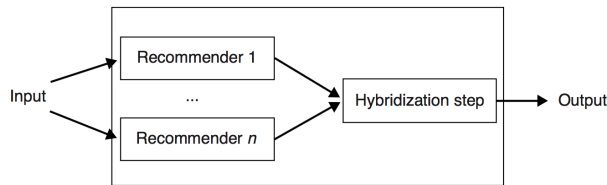


FIGURE 2.3: Parallelized hybridization design

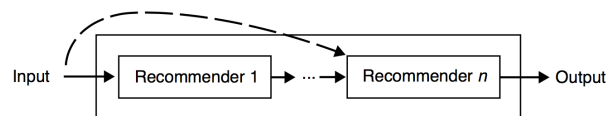


FIGURE 2.4: Pipelined hybridization design

2.2 Overview of music recommender systems

Music, differs from other content domains such as books or movies, has its own unique characteristics regarding consumption. For example, the consumption time of books and movies are quite lengthy, ranging on the average of few hours (for movies) to several days (for books); songs, on the other hand, takes a much shorter time for listener to consume, only a few minutes. Consequently, this lead to the ephemeral and disposable nature of music. Another example is the number of repetition: a single song can be consumed repeatedly (even multiple times in a row),

while books are movies are consumed a few times at most. This implies that the user might appreciate recommendations of items that they already heard in the past.

In the past, many approaches, based on the observation of the nature of music and listener's behaviors, had been tried to build an efficient music recommender system. Till now, there are three main approaches to such system [60], which will be described in the following sections.

2.2.1 Content-Based Music Recommendation

Content-based recommendation exploits information describing music as material for recommender systems. There are two main approaches for content-based system: one uses metadata, such as annotations or social tags, while the other analyses audio content using machine learning techniques.

2.2.1.1 Metadata content

Metadata comes in several forms. One is the manual annotations constructed by music experts or voluntary community. This kind of annotation often obeys a strict structured taxonomy build by experts. Another form of annotations is social tags, which is built by asking casual users to provide unstructured text annotations for the item. The last kind of annotations is information collected on web pages, blogs and RSS feeds related to music items.

2.2.1.1.1 Annotation

Manual annotation contain information such as musical genre, record label, year, knowledge about tracks and artists, and albums. Some musical properties, particularly tempo, mood, and instrumentation can also be added. Many online database have been built using editorial metadata, following by many recommender systems trying to exploit them.

Bogdanov et al. [10] build an artist recommender using exclusively metadata from *Discogs*, a free and community-built database containing information about artists, records labels, and their releases. For each artist in the database, a tag weight vector is created using genre, style, label, country, and year information of the releases related to the artist. The role of the artist in each release (e.g. main artist, track artist, or extra artist) and the relations between artist, such as aliases and membership relations, are also taken into account. A sparse tag matrix is then formed from the artist vectors, and latent semantic analysis [18] is applied to reduce the dimension of the matrix. Afterwards, the authors use Pearson correlation distance [23] to measure the similarity between artists.

Apart from community-built database, some other database are developed by experts for commercial use. *Pandora*, for example, is a personalized radio that built its recommender using annotations done by experts [31]. *AllMusic* is another example that also provides mood annotations besides general editorial metadata. However, not much research has been done using these database, as they are proprietary, and there is no public data sets of this kind are available for researcher. Constructing such a data set would be costly, and they are difficult to scale to large collections.

2.2.1.1.2 Social tags

Social tags are tags provided by user using the services. Tags are personalized, arbitrary, and do not follow any particular structure, ranging from genre like "blue" or

"jazz" to event-related attributes (e.g. "live") and assertion (e.g. "my favorite song"). They also vary in scope, from broad one such as "classical" to niche terminology like "Malcolm Arnold" or "renaissance". Social tags, therefore, need to be preprocess for the data to be useful. A popular method, which is used by *Last.fm*, a social music website, for structuring social tags is to transform them into a folksonomy [71]. The tag weight vectors technique is then applied to compute the similarity [25], with the enhancement by using latent semantic analysis techniques to overcome vector sparsity problem [35]

2.2.1.1.3 Annotations by web crawling

Apart from tags made by experts and social tags, some recommender systems are built using information crawled from web pages. These recommenders often apply artist similarity metric, generated by using text mining techniques [66], as the main principle for the recommendation. Green et al. [25] compute artist-to-artist similarity using keyword extracted from *Wikipedia* entries and social tags from *Last.fm*. Similarly, McFee and Lanckriet [45] predict artist similarity based on social tags and keyword extracted from artist biographies on *Last.fm*. A deviant approach is the one made by Lim et al., as they compute song-level similarity through bag-of-words representations of lyrics found on *musiXmatch.com* [39]

2.2.1.1.2 Audio content

Audio content analysis is promoted by MIR researchers as an alternative to metadata and collaborative filtering method [6]. Content analysis is expected to solve "long tail" problem, where unpopular music items are not suggested because the lack of available user ratings, tags, and other types of metadata [16]. Music content is separated into two broad categories: acoustic features taken directly from the audio, and semantic annotations derived from acoustic features using machine learning techniques.

2.2.1.2.1 Acoustic features

Acoustic features are properties of a sound that can be recorded and analyzed using signal processing techniques. As mentioned in the content-based recommendation part, there are three main features that are often used by recommender system: timbral features, temporal and time-domain features, and tonal features.

Timbral similarity method converts timbre information to a standard representation and applies a number of methods to approximate the likelihood between two songs [4] [42]. For example, Logan [41] builds a recommender that compares Mel-frequency cepstrum coefficient (MFCC) based distance of user's music set with target play list. The approach, however, is insufficient as evaluation shows a nominal number of customer satisfaction [11].

Pampalk et al. [51] [52] study an algorithm that use, in addition to spectral similarity, loudness fluctuations and two derived descriptors from sound wave to improve the accuracy of music similarity and genre classification of a song. The algorithm is used to recommend user playlists, in which songs that are similar to the ones that are skipped by user are eliminate, and only tracks that are similar to the ones that the user wholly listens to remains.

Celma and Herrera [15] take another approach, calculating Euclidean distance using timbre, dynamics, tempo, meter, tonal strength, key, and mode information. This method is compared to an item-based collaborative filtering and a hybrid method

on a large scale evaluation. The result shows that all three algorithms work fine recommending familiar items. For unfamiliar items, collaborative filtering reveals a poor discovery ratio, while pure content-based method sometimes goes off direction, as the recommender confuse between similar sounds, such as between the sound of classical guitar and the one of harpsichord). The hybrid method performs the best, as it limit the number of possible tracks whose artists are related with the original artist, therefore reducing mistakes performing by pure content-based method.

2.2.1.2.2 Semantic annotation

Pure acoustic signal, unfortunately, cannot directly capture semantic meaning of a track. As a result, mere acoustic feature recommenders do a poor job in song suggestion, as an "energetic" song can be recommended next to a "nostalgic" track because of the similarity of the instrument. A sudden change in genre might frustrate customer, as one would not expect a mourning song in a middle of a party. Therefore, many approaches have been made in order to bridge this semantic gap by using machine learning techniques to predict annotations from audio content.

However, extracting genre or mood information from acoustic content is perplexing, as mapping between human annotation and acoustic cannot be clearly defined [5]. To solve this, Barrington et al. [6] propose a method to measure semantic similarity: they train Gaussian mixture models of MFCCs for semantic concepts such as genres, moods and instruments. Therefore, for a song, a distribution of tags is generated, which is then compared to another in order to estimate similarity.

2.2.2 Contextual Music Recommendation

The vast majority of existing recommender system approaches focus on information about users and items but not context information, such as time and place of the event. Only until recently, the topic of context-awareness starts to gain attraction in recommender system research [2]. Context, in the fields that are directly related to recommender system, is defined as "information describing where you are, whom you are with, and what resources are nearby" [68]. Therefore, context in the domain of music can be derived as a collection of factors that affects user's appreciation of music, such as time, mood, and current activities.

Contextual information can be classified using various kind of classifications. Adomavicius et al. [2] suggest categorizing context into three distinct classes: fully observable, partially observable, and unobservable context. Dey and Abowd [1], in attempt to construct another classification, propose to classify context into primary context, which is four most important factors that describe user situation: location, identity, activity, and time; and secondary context, which is data derived from primary information. Applying the classification of Dey and Abowd, M. Schedl et al [67] divide context information into two generic classes: environment-related context and user-related context. Environment-related context refers to information that can be obtained by user's computer or mobile phone, such as location, time, weather, etc., while user-related context indicates information that can be derived from the environment-related one, such as user's activity, emotion state, and social environment.

2.2.2.1 Environment-Related Context

Surrounding environment has been proven to have an influence on user's preference of music. Adrian C. North and David J. Hargreaves [49] find a correlation between musical descriptors, such as arousal, sensuality, spirituality, and listening situation, such as activity, spirituality, and social constraint. Pettijohn et al. [55] find that different seasons, such as winter and summer, also affect musical preferences.

Many attempts have been made to build recommender using environment-related information. Reddy and Mascia [58] used space information capturing using GPS coordinates, internal time data, kinetic information derived from difference in GPS signals, and even meteorological info to build a recommender for a mobile music player called Lifetrak. Songs have to be tagged manually by users using system predefined tagging system, and are played in appropriate situation based on users' preferences. For example, the app may play rock music when a user is in a gym, and classical music when the user tries to study in a cafeteria.

Ankolekar and Sandholm from HP labs [3] propose a mobile audio application, Foxtrot, that exploits crowd-sourced geo-tagged audio information to provide a stream of location-aware audio content to the users. The recommender, however, generated poor user experiences, as an environment generates different meaning to different people, leading to diverse music preferences. Indeed, a research made by Okada et al. [50] shows that not only the algorithm, but also the architectural design and usability of the application that matter, as user feedback suggests the need for explanations of the recommendations and more control over the playlist.

2.2.2.2 User-Related Context

One's music preference is not only affected by geological and activity component, but also by factors such as emotions and social background. Schäfer and Sedlmeier [64] discovered that one's music preference is also linked with one's sociocultural and physiological functions. In other words, people use music preference as a mean to express their identities and personal values. User-related context can be divided into the following groups:

- Activity information: information implies user's actions (e.g., walking, driving, working) or user's state (e.g., walking pace or heart rate). Foley [21] showed that people with different occupation have different favored music tempo. For example, those who work with power machines like a slow allegro, while typists prefer faster tempo like presto.
- Emotional information: current mood of user has a direct impact on the choice of music. For example, a user may want to listening energetic music while he is happy, and calm music vice versa. Schäfer and Sedlmeier [64] found that music has a function to moderate listener's mood by energizing him or making he feel better.
- Social context information: music preference can be affected by the presence of other people. People may choose music taken into account the event they participate in. Many researchers have address the issue of group recommender system. For example, Popescu and Pu [57] proposed using probabilistic weighted sum as the algorithm to recommend group playlist.
- Cultural context information: information about user's culture characteristics. Koenigstein et al. [33] exploited file sharing information on a Peer to

Peer network in US to predict the success of a song on Billboard Hot 100 Chart. Schedl [65] built a location-aware recommender system by retrieving geo-tagged Twitter tweets to detect listening trend at a particular place.

Compare to environment-related context, user-related context is more difficult to infer using electronic devices. Many attempts have been made to predict user's emotion or daily activities [54] [77] by extracting environment-related feature such as the time of day, temperature, weather, etc... Emotion-based music recommender has gained attention recently, due to advances in automatic music emotion recognition [80]

Appendix A

Appendix A

Bibliography

- [1] Gregory Abowd et al. "Towards a better understanding of context and context-awareness". In: *Handheld and ubiquitous computing*. Springer. 1999, pp. 304–307.
- [2] Gediminas Adomavicius and Alexander Tuzhilin. "Context-aware recommender systems". In: *Recommender systems handbook*. Springer, 2011, pp. 217–253.
- [3] Anupriya Ankolekar and Thomas Sandholm. "Foxtrot: a soundtrack for where you are". In: *Proceedings of Interacting with Sound Workshop: Exploring Context-Aware, Local and Social Audio Applications*. ACM. 2011, pp. 26–31.
- [4] J-J Aucouturier, François Pachet, and Mark Sandler. "'The way it Sounds': timbre models for analysis and retrieval of music signals". In: *IEEE Transactions on Multimedia* 7.6 (2005), pp. 1028–1035.
- [5] Jean-Julien Aucouturier. "Sounds like teen spirit: Computational insights into the grounding of everyday musical terms". In: *Language, evolution and the brain* (2009), pp. 35–64.
- [6] Luke Barrington, Reid Oda, and Gert RG Lanckriet. "Smarter than Genius? Human Evaluation of Music Recommender Systems." In: *ISMIR*. Vol. 9. 2009, pp. 357–362.
- [7] Juan Pablo Bello. "Audio-Based Cover Song Retrieval Using Approximate Chord Sequences: Testing Shifts, Gaps, Swaps and Beats." In: *ISMIR*. Vol. 7. 2007, pp. 239–244.
- [8] Jacob Benesty et al. "Pearson correlation coefficient". In: *Noise reduction in speech processing*. Springer, 2009, pp. 1–4.
- [9] James Bergstra et al. "Aggregate features and AdaBoost for music classification". In: *Machine learning* 65.2-3 (2006), pp. 473–484.
- [10] Dmitry Bogdanov and Perfecto Herrera. "Taking advantage of editorial meta-data to recommend music". In: *9th International Symposium on Computer Music Modeling and Retrieval (CMMR)*. 2012, pp. 618–632.
- [11] Dmitry Bogdanov et al. "Semantic audio content-based music recommendation and visualization based on user preference examples". In: *Information Processing & Management* 49.1 (2013), pp. 13–33.
- [12] John S Breese, David Heckerman, and Carl Kadie. "Empirical analysis of predictive algorithms for collaborative filtering". In: *Proceedings of the Fourteenth conference on Uncertainty in artificial intelligence*. Morgan Kaufmann Publishers Inc. 1998, pp. 43–52.
- [13] Robin Burke. "Hybrid recommender systems: Survey and experiments". In: *User modeling and user-adapted interaction* 12.4 (2002), pp. 331–370.
- [14] Pedro Cano, Markus Koppenberger, and Nicolas Wack. "An industrial-strength content-based music recommendation system". In: *Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM. 2005, pp. 673–673.

- [15] Òscar Celma and Perfecto Herrera. "A new approach to evaluating novel recommendations". In: *Proceedings of the 2008 ACM conference on Recommender systems*. ACM. 2008, pp. 179–186.
- [16] Òscar Celma Herrada et al. *Music recommendation and discovery in the long tail*. Universitat Pompeu Fabra, 2009.
- [17] Dan Cosley et al. "Is seeing believing?: how recommender system interfaces affect users' opinions". In: *Proceedings of the SIGCHI conference on Human factors in computing systems*. ACM. 2003, pp. 585–592.
- [18] Scott Deerwester et al. "Indexing by latent semantic analysis". In: *Journal of the American society for information science* 41.6 (1990), p. 391.
- [19] Daniel PW Ellis and Graham E Poliner. "Identifying cover songs' with chroma features and dynamic programming beat tracking". In: *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on*. Vol. 4. IEEE. 2007, pp. IV–1429.
- [20] Yazhong Feng, Yueting Zhuang, and Yunhe Pan. "Music information retrieval by detecting mood via computational media aesthetics". In: *Web Intelligence, 2003. WI 2003. Proceedings. IEEE/WIC International Conference on*. IEEE. 2003, pp. 235–241.
- [21] John P Foley Jr. "The Occupational Conditioning of Preferential Auditory Tempo: A Contribution toward an Empirical Theory of Aesthetics". In: *The Journal of Social Psychology* 12.1 (1940), pp. 121–129.
- [22] Zhouyu Fu et al. "A survey of audio-based music classification and annotation". In: *IEEE transactions on multimedia* 13.2 (2011), pp. 303–319.
- [23] Jean Dickinson Gibbons and Subhabrata Chakraborti. "Nonparametric statistical inference". In: *International encyclopedia of statistical science*. Springer, 2011, pp. 977–979.
- [24] David Goldberg et al. "Using collaborative filtering to weave an information tapestry". In: *Communications of the ACM* 35.12 (1992), pp. 61–70.
- [25] Stephen J Green et al. "Generating transparent, steerable recommendations from textual descriptions of items". In: *Proceedings of the third ACM conference on Recommender systems*. ACM. 2009, pp. 281–284.
- [26] Jonathan L Herlocker et al. "An algorithmic framework for performing collaborative filtering". In: *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*. ACM. 1999, pp. 230–237.
- [27] Will Hill et al. "Recommending and evaluating choices in a virtual community of use". In: *Proceedings of the SIGCHI conference on Human factors in computing systems*. ACM Press/Addison-Wesley Publishing Co. 1995, pp. 194–201.
- [28] Jay Hodgson. *Understanding Records: A Field Guide to Recording Practice*. Bloomsbury Publishing, 2010.
- [29] Jyh-Shing Roger Jang and Hong-Ru Lee. "A general framework of progressive filtering and its application to query by singing/humming". In: *IEEE Transactions on Audio, Speech, and Language Processing* 16.2 (2008), pp. 350–358.
- [30] Dietmar Jannach et al. *Recommender systems: an introduction*. Cambridge University Press, 2010.

- [31] Nicolas Jones and Pearl Pu. "User technology adoption issues in recommender systems". In: *Proceedings of the 2007 Networking and Electronic Commerce Research Conference*. HCI-CONF-2008-001. 2007, pp. 379–394.
- [32] Daniel Jurafsky and James H Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*.
- [33] Noam Koenigstein, Yuval Shavitt, and Noa Zilberman. "Predicting billboard success using data-mining in p2p networks". In: *Multimedia, 2009. ISM'09. 11th IEEE International Symposium on*. IEEE. 2009, pp. 465–470.
- [34] Yehuda Koren, Robert Bell, and Chris Volinsky. "Matrix factorization techniques for recommender systems". In: *Computer* 42.8 (2009).
- [35] Mark Levy and Mark Sandler. "Learning latent semantic models for music from social tags". In: *Journal of New Music Research* 37.2 (2008), pp. 137–150.
- [36] Tao Li and Mitsunori Ogihara. "Detecting emotion in music". In: (2003).
- [37] Tao Li, Mitsunori Ogihara, and Qi Li. "A comparative study on content-based music genre classification". In: *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*. ACM. 2003, pp. 282–289.
- [38] Thomas Lidy et al. "Improving Genre Classification by Combination of Audio and Symbolic Descriptors Using a Transcription Systems." In: *ISMIR*. 2007, pp. 61–66.
- [39] Daryl Lim, Brian McFee, and Gert R Lanckriet. "Robust structural metric learning". In: *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*. 2013, pp. 615–623.
- [40] Chien-Chang Lin et al. "Audio classification and categorization based on wavelets and support vector machine". In: *IEEE Transactions on Speech and Audio Processing* 13.5 (2005), pp. 644–651.
- [41] Beth Logan. "Music Recommendation from Song Sets." In: *ISMIR*. 2004, pp. 425–428.
- [42] Beth Logan and Ariel Salomon. "A Music Similarity Function Based on Signal Analysis." In: *ICME*. 2001, pp. 22–25.
- [43] M Mandel and D Ellis. "Song-level features and SVM for music classification". In: *In Proceedings of the 6th International Conference on Music Information Retrieval, ISMIR*. Vol. 5. 2006.
- [44] Matija Marolt. "A Mid-level Melody-based Representation for Calculating Audio Similarity." In: *ISMIR*. 2006, pp. 280–285.
- [45] Brian McFee and Gert Lanckriet. "Learning multi-modal similarity". In: *Journal of machine learning research* 12.Feb (2011), pp. 491–523.
- [46] Anders Meng et al. "Temporal feature integration for music genre classification". In: *IEEE Transactions on Audio, Speech, and Language Processing* 15.5 (2007), pp. 1654–1664.
- [47] Tom M Mitchell. "Machine learning. 1997". In: *Burr Ridge, IL: McGraw Hill* 45.37 (1997), pp. 870–877.
- [48] F Morchen et al. "Modeling timbre distance with temporal statistics from polyphonic music". In: *IEEE Transactions on Audio, Speech, and Language Processing* 14.1 (2006), pp. 81–90.

- [49] Adrian C North and David J Hargreaves. "Situational influences on reported musical preference." In: *Psychomusicology: A Journal of Research in Music Cognition* 15.1-2 (1996), p. 30.
- [50] Karla Okada et al. "ContextPlayer: Learning contextual music preferences for situational recommendations". In: *SIGGRAPH Asia 2013 Symposium on Mobile Graphics and Interactive Applications*. ACM. 2013, p. 6.
- [51] Elias Pampalk, Arthur Flexer, Gerhard Widmer, et al. "Improvements of Audio-Based Music Similarity and Genre Classification." In: *ISMIR*. Vol. 5. London, UK. 2005, pp. 634–637.
- [52] Elias Pampalk, Tim Pohle, and Gerhard Widmer. "Dynamic Playlist Generation Based on Skipping Behavior." In: *ISMIR*. Vol. 5. 2005, pp. 634–637.
- [53] Elias Pampalk, Andreas Rauber, and Dieter Merkl. "Content-based organization and visualization of music archives". In: *Proceedings of the tenth ACM international conference on Multimedia*. ACM. 2002, pp. 570–579.
- [54] Han-Saem Park, Ji-Oh Yoo, and Sung-Bae Cho. "A context-aware music recommendation system using fuzzy bayesian networks with utility theory". In: *International Conference on Fuzzy Systems and Knowledge Discovery*. Springer. 2006, pp. 970–979.
- [55] Terry F Pettijohn, Greg M Williams, and Tiffany C Carter. "Music for the seasons: seasonal music preferences in college students". In: *Current psychology* 29.4 (2010), pp. 328–345.
- [56] Graham E Poliner et al. "Melody transcription from music audio: Approaches and evaluation". In: *IEEE Transactions on Audio, Speech, and Language Processing* 15.4 (2007), pp. 1247–1256.
- [57] George Popescu and Pearl Pu. "Probabilistic game theoretic algorithms for group recommender systems". In: *Colocated with ACM RecSys 2011 Chicago, IL, USA October 23, 2011* (2011), p. 30.
- [58] Sasank Reddy and Jeff Mascia. "Lifetrak: music in tune with your life". In: *Proceedings of the 1st ACM international workshop on Human-centered multimedia*. ACM. 2006, pp. 25–34.
- [59] Paul Resnick. "An Open Architecture for Collaborative Filtering of Netnews". In: *Proc. of CSCW'94*. 1994, pp. 175–186.
- [60] Francesco Ricci, Lior Rokach, and Bracha Shapira. "Introduction to recommender systems handbook". In: *Recommender systems handbook*. Springer, 2011, pp. 1–35.
- [61] Gerard Salton. "Syntactic approaches to automatic book indexing". In: *Proceedings of the 26th annual meeting on Association for Computational Linguistics*. Association for Computational Linguistics. 1988, pp. 204–210.
- [62] Nicolas Scaringella, Giorgio Zoia, and Daniel Mlynek. "Automatic genre classification of music content: a survey". In: *IEEE Signal Processing Magazine* 23.2 (2006), pp. 133–141.
- [63] J Ben Schafer, Joseph A Konstan, and John T Riedl. "Recommender Systems for the Web". In: *Visualizing the Semantic Web*, Springer (2006), pp. 102–123.
- [64] Thomas Schäfer and Peter Sedlmeier. "From the functions of music to music preference". In: *Psychology of Music* 37.3 (2009), pp. 279–300.

- [65] Markus Schedl. "Leveraging Microblogs for Spatiotemporal Music Information Retrieval." In: *ECIR*. Springer. 2013, pp. 796–799.
- [66] Markus Schedl et al. "Exploring the music similarity space on the web". In: *ACM Transactions on Information Systems (TOIS)* 29.3 (2011), p. 14.
- [67] Markus Schedl et al. "Music recommender systems". In: *Recommender Systems Handbook*. Springer, 2015, pp. 453–492.
- [68] Bill Schilit, Norman Adams, and Roy Want. "Context-aware computing applications". In: *Mobile Computing Systems and Applications, 1994. WMCSA 1994. First Workshop on*. IEEE. 1994, pp. 85–90.
- [69] Upendra Shardanand and Pattie Maes. "Social information filtering: algorithms for automating "word of mouth"". In: *Proceedings of the SIGCHI conference on Human factors in computing systems*. ACM Press/Addison-Wesley Publishing Co. 1995, pp. 210–217.
- [70] JS Shawe-Taylor and Anders Meng. "An investigation of feature models for music genre classification using the support vector classifier". In: (2005).
- [71] Mohamed Sordo et al. "The quest for musical genres: Do the experts and the wisdom of crowds agree?" In: *ISMIR*. 2008, pp. 255–260.
- [72] *SVD recommendation system in ruby*. <http://www.igvita.com/2007/01/15/svd-recommendation-system-in-ruby/>. Accessed: 2017-08-15.
- [73] Wei-Ho Tsai, Hung-Ming Yu, Hsin-Min Wang, et al. "Query-By-Example Technique for Retrieving Cover Versions of Popular Songs with Similar Melodies." In: *ISMIR*. Vol. 5. 2005, pp. 183–190.
- [74] George Tzanetakis and Perry Cook. "Musical genre classification of audio signals". In: *IEEE Transactions on speech and audio processing* 10.5 (2002), pp. 293–302.
- [75] George Tzanetakis, Randy Jones, and Kirk McNally. "Stereo Panning Features for Classifying Recording Production Style." In: *ISMIR*. 2007, pp. 441–444.
- [76] George Tzanetakis et al. "Stereo panning information for music information retrieval tasks". In: *Journal of the Audio Engineering Society* 58.5 (2010), pp. 409–417.
- [77] Xinxi Wang, David Rosenblum, and Ye Wang. "Context-aware mobile music recommendation for daily activities". In: *Proceedings of the 20th ACM international conference on Multimedia*. ACM. 2012, pp. 99–108.
- [78] Claus Weihs et al. "Classification in music research". In: *Advances in Data Analysis and Classification* 1.3 (2007), pp. 255–291.
- [79] Dan Yang and Won-Sook Lee. "Disambiguating Music Emotion Using Software Agents." In: *ISMIR*. Vol. 4. 2004, pp. 218–223.
- [80] Yi-Hsuan Yang and Homer H Chen. "Machine recognition of music emotion: A review". In: *ACM Transactions on Intelligent Systems and Technology (TIST)* 3.3 (2012), p. 40.