

# Automated Generation of Music Playlists: Survey and Experiments

GEOFFRAY BONNIN and DIETMAR JANNACH, Technische Universität Dortmund

26

Most of the time when we listen to music on the radio or on our portable devices, the order in which the tracks are played is governed by so-called playlists. These playlists are basically sequences of tracks that traditionally are designed manually and whose organization is based on some underlying logic or theme. With the digitalization of music and the availability of various types of additional track-related information on the Web, new opportunities have emerged on how to automate the playlist creation process. Correspondingly, a number of proposals for automated playlist generation have been made in the literature during the past decade. These approaches vary both with respect to which kind of data they rely on and which types of algorithms they use. In this article, we review the literature on automated playlist generation and categorize the existing approaches. Furthermore, we discuss the evaluation designs that are used today in research to assess the quality of the generated playlists. Finally, we report the results of a comparative evaluation of typical playlist generation schemes based on historical data. Our results show that track and artist popularity can play a dominant role and that additional measures are required to better characterize and compare the quality of automatically generated playlists.

Categories and Subject Descriptors: H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval; H.5.5 [Information Interfaces and Presentation]: Sound and Music Computing

General Terms: Playlist Generation, Music Recommendation

Additional Key Words and Phrases: Music, playlist, algorithm, evaluation

## ACM Reference Format:

Geoffray Bonnin and Dietmar Jannach. 2014. Automated generation of music playlists: Survey and experiments. *ACM Comput. Surv.* 47, 2, Article 26 (November 2014), 35 pages.

DOI: <http://dx.doi.org/10.1145/2652481>

## 1. INTRODUCTION

Much of the music that we listen to is organized in some form of a playlist. Such playlists are traditionally created manually (e.g., by a radio disc jockey (DJ) or music producer) and usually have some internal organization logic (e.g., their tracks belong to the same musical genre or are performed by related artists). Beside playlists that are, for example, broadcast over the radio, the digitalization of music has also made it easier for users to create their own *personal* playlists. People use such personal playlists for a number of purposes, such as to create a collection of tracks used as background music for some event or activity or to share their musical preferences with others. Since the manual creation of playlists can be time intensive or can require special skills by the playlist creator, a number of proposals have been made in the literature to automate the playlist creation process. Today, popular music platforms such as last.fm or Pandora

---

Authors' address: Department of Computer Science, TU Dortmund, 44227 Dortmund, Germany; email: [geoffray.bonnin@tu-dortmund.de](mailto:geoffray.bonnin@tu-dortmund.de) and [dietmar.jannach@tu-dortmund.de](mailto:dietmar.jannach@tu-dortmund.de).

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© 2014 ACM 0360-0300/2014/11-ART26 \$15.00

DOI: <http://dx.doi.org/10.1145/2652481>

and personalized Web radio services such as Musicoverly rely strongly on automatically created playlists.

In this article, we review and categorize the various approaches to automated playlist generation. We will look at the different types of data that are used in the playlist construction process as well as the various underlying algorithms. Furthermore, we will review and discuss common approaches to evaluate the quality of such automatically generated playlists. Based on a number of experiments, we will then discuss possible limitations of existing offline evaluation schemes and current research practice. In addition, we show that simple playlist generation schemes based on popularity and artist information can—depending on the evaluation scheme—represent a comparably hard baseline. Finally, we will discuss perspectives for future research.

## 2. THE PLAYLIST GENERATION PROBLEM

In this section, we will first define the concepts of playlist and playlist generation, sketch typical challenges in the domain, and discuss the relation of playlist generation with the fields of music information retrieval (MIR) and recommender systems (RS).

### 2.1. Definitions

In this work, we define the term *playlist* as follows.

*Definition 2.1 (Playlist).* A *playlist* is a sequence of tracks (audio recordings).

This definition corresponds to what is most commonly used in the literature (e.g., see Baccigalupo and Plaza [2006], Knees et al. [2006], Baccigalupo [2009], McFee and Lanckriet [2011], and Hariri et al. [2012]). Some authors make a distinction between playlists whose tracks are arranged in a certain order and playlists in which the tracks are shuffled. However, the terminology in the literature is not always consistent. Logan [2004], for example, talks about “song sets” when the tracks are shuffled and about “playlists” otherwise. Cunningham et al. [2006], in contrast, use “playlist” for shuffled tracks and “mix” for sequentially ordered tracks.

Our definition is thus very general and covers several types of typical playlists:

- Broadcasting radio playlists* are made by DJs in radio stations. Such playlists often contain popular tracks in a relatively homogeneous style suited for a certain target audience [Ekelund et al. 2000].
- Personalized radio playlists* are generated by Web music services such as last.fm<sup>1</sup> or Pandora.<sup>2</sup> These playlists can be very different from broadcasting radio playlists because the listeners directly interact with the service and playlists can be adapted to the preferences of the listeners. Often, one of the goals of these playlists is discovery—that is, to help the listeners find tracks or artists that are unknown to them.
- Amateur playlists* are playlists made by nonprofessional music enthusiasts. Such playlists can have different purposes, and people, for example, prepare them for certain future occasions like a party or a travel [Cunningham et al. 2006] for sharing their music with others or for self-expression.
- Club playlists* are made by DJs in clubs. Such playlists usually contain “danceable” tracks [Cliff 2006].
- Album tracklists* are made by artists and labels and are usually a sequence of new tracks of one artist or a group of artists. To which extent labels allow artists to influence the selection and ordering or even recording of the tracks typically depends on the popularity of the artists [Baskerville and Baskerville 2010].

<sup>1</sup><http://last.fm>.

<sup>2</sup><http://pandora.com>.

—*Compilation tracklists* are also created by artists and labels and, for example, consist of popular tracks of one artist or are grouped around a theme such as “Christmas songs.”

In this work, we focus on techniques that support or automate the process of creating playlists. Such approaches are applicable for the first four types of playlists. In our view, the content and structure of the last two types of playlists can be governed by factors that can be hardly automated, such as specific artist preferences or marketing considerations of the producers and labels.

All considered types of playlists have in common that their construction is determined by (1) the set of available tracks and (2) the target characteristics of the playlist. When building a playlist, a human creator will pick and arrange the tracks in a way that they fulfill some characteristics and will rely on some additional information. This information can include the genre, mood, or tempo of the tracks as well as the general or current popularity of the artist within the target audience.

To automate this task, a playlist generation tool first has to include a mechanism to somehow acquire or assess the intended goals and thus what is more important—the corresponding desired target characteristics of the playlist to be generated. In addition, a mechanism or heuristic is required that determines the suitability of a given track or track ordering given these target characteristics. The playlist generation problem can thus be defined as follows.

*Definition 2.2 (Playlist generation).* Given (1) a pool of tracks, (2) a background knowledge database, and (3) some target characteristics of the playlist, create a sequence of tracks fulfilling the target characteristics in the best possible way.

Our definition is again quite broad to cover a wide range of application scenarios. Considering, for example, an interactive scenario as on last.fm, the target characteristics of the playlist can be expressed by the user by providing an artist name or a genre and then refined implicitly using the “skip” and “ban” buttons. The most recent listening history could also be considered in the target characteristics—that is, a desired characteristic could be that the continuation is coherent with the most recently listened to tracks. The background knowledge used by last.fm can include both genre information about the available tracks and the popularity of individual tracks. In a one-time generation scenario, the goal could be to create a playlist for a certain activity or context, such as doing sports in the gym with a friend [Masahiro et al. 2008]. The required knowledge could be the tempo or mood of the tracks. Obviously, the listening context and the goal of “suitability for sports” first has to be translated into target characteristics, for example, by defining a minimum tempo or by requiring that the tracks are annotated with certain tags before the selection and ordering of the tracks.

## 2.2. Challenges of Playlist Generation

A number of aspects make automated playlist generation challenging. First, even if we talk of an automated process, the large number of available tracks that can potentially be included in the playlists—already, personal music libraries can be huge [Dias and Fonseca 2010]—and the corresponding enormous amount of additional information that has to be processed can lead to high computational complexity. Although some of the computations can be done in an offline model-learning process, the problem exists that many application scenarios require the “on-demand” generation of playlists or the immediate reaction on user feedback.

Another challenge lies in the acquisition of information about the tracks. Extracting characteristics from the audio signal for a large number of tracks can require significant computation times. On the other hand, if we rely on user-provided annotations,

obtaining a set of high-quality annotations for each track can represent a major problem. Often, only a part of the tracks is annotated [Celma 2010], and the annotation quality can be low and subjective.

Finally, the taste and context of the listeners are two factors that can have a major influence on the *perceived* suitability or quality of a playlist, and two users may actually perceive the same playlist differently even if listened to in the same context. Therefore, the challenge is not only to find tracks that match the specifics of the current user taste and context but also to identify the user's taste and context correctly in the first place. Questions of how the quality of playlists can be at least partially automatically and objectively assessed will be discussed in Section 5.

### 2.3. Relation to Music Information Retrieval and Recommender Systems

Research on automated playlist generation is mainly done in two research fields: MIR and RS.

*2.3.1. Music Information Retrieval.* MIR is an interdisciplinary research area that covers a wide range of tasks based on the automatic analysis of music data, among them playlist generation [Downie 2003; Grachten et al. 2009]. The scope of MIR thus goes far beyond the application of information retrieval (IR) to music data, whose scope is often limited to the retrieval of resources matching a user query [Crestani et al. 1998].

In general, playlist generation could be considered as a retrieval task if its definition is limited to the selection of tracks satisfying a user query. In traditional IR scenarios, the optimization goal is usually to find and rank all items that are relevant for a query. The quality of a playlist, however, cannot be assessed by only looking at its individual items, as there might be desirable characteristics that concern the list as a whole, such as that it is diverse and contains items that are unknown to the user or that track transitions are smooth. Furthermore, explicit queries are only one possible user input.

Other major topics in MIR include the automated extraction of features, the categorization of tracks, or similarity analysis. These types of automatically identified information can obviously be a basis for building a playlist generation system (e.g., see Dopler et al. [2008] and Lee et al. [2011]).

*2.3.2. Recommender Systems.* RS are applications that provide personalized recommendations to users [Burke 2002]. Conceptually, the playlist generation problem can be considered as a special case of the recommendation problem—that is, the recommendation of collections [Hansen and Golbeck 2009]. Furthermore, RS can be used to recommend various musical resources, such as albums, artists, or concerts, and playlist generation can also be considered as a special case of music (track) recommendation. Again, however, the particularity of playlist generation is that the properties of the list itself and track orderings can be relevant.

Unlike in many RS research scenarios, the *immediate consumption* of the recommendations is common in playlist generation scenarios [Lamere 2012]. A playlist generator should thus be able to dynamically adapt its recommendations to the user's current situation, for example, by incorporating implicit feedback like “skip” actions or other types of information from the user [Reddy and Mascia 2006].

Finally, differently from typical e-commerce scenarios evaluated in RS research [Jannach et al. 2012], the repeated consumption of items is very common in the music domain, and listening logs like the ones available from last.fm show that many users listen to their favorite tracks over and over.

## 3. BACKGROUND KNOWLEDGE AND TARGET CHARACTERISTICS

According to our definition, three inputs are required to generate playlists: a pool of tracks, a background knowledge database, and some target characteristics. Next,

we will review existing playlist generation approaches with respect to the latter two inputs—that is, (1) the used background knowledge (2) in which form the target characteristics and user inputs are obtained.

### 3.1. Background Knowledge

The background knowledge corresponds to any information that can be used to determine whether a playlist and its tracks satisfy the target characteristics. We classify this information in the following categories: the audio signal, metadata, Social Web data, and usage data.

*3.1.1. Musical Features from the Audio Signal.* An obvious type of data that can be used in the playlist generation process is the audio content. The goal of many MIR researchers is to automatically extract musical *features* from the audio signal [Pohle et al. 2005], including, for example, pitch, loudness [Blum et al. 1999], rhythm, chord changes [Tzanetakis 2002], or simply compact representations of the frequency spectrum, particularly spectral clusters [Logan and Salomon 2001] and Gaussian mixture models [Aucouturier and Pachet 2002b]. Feature information can also be obtained through automated classification algorithms [Wold et al. 1996].

The advantage of relying on the audio signal over, say, expert annotations, is that the process is objective and can be automated to a large extent. However, extracting the features can be computationally costly [Schmädecke and Blume 2013]. Another limitation is that there might be features like the release date, the “freshness,” or popularity of a track, which can be relevant in the playlist generation process but that cannot be extracted from the audio signal [Celma and Serra 2008].

*3.1.2. Metadata and Expert Annotations.* Metadata shall be defined here as any information describing the tracks that is not derived from the audio signal, for example, track attributes such as the year of release [Van Gulik and Vignoli 2005], the label [Pachet et al. 2000], artist and genre information [Aizenberg et al. 2012; Bogdanov and Herrera 2011], and lyrics [Coelho et al. 2013]. This information is often the result of a manual annotation process; in some cases, the information can be obtained from online music platforms or music databases such as the Million Song Dataset; see the Appendix for a list of public sources.

When used in an automated process, data completeness and consistency are crucial. Some parts of the metadata, such as artist and track names, can be easy to obtain [McFee and Lanckriet 2011], whereas the year of release might be more difficult to obtain if it cannot be taken from a track database. Another potential problem is that not all types of metadata are objective, and annotations regarding, for example, the mood or the genre of a track can be imprecise or inconsistent [Celma 2010; Lee and Cho 2011].

Generally, metadata from expert annotations can be combined with automatically derived or cultural features [McKay 2010; Coelho et al. 2013]. Deriving the tempo of a track automatically can be challenging and inaccurate [Chen et al. 2009; Peeters and Flocon-Cholet 2012]. Therefore, in case of a low confidence, more reliable expert annotations could be used but induce higher costs. Relying only on expert annotations may be too costly for the often large pools of tracks.

*3.1.3. Social Web Data.* With the spread of the World Wide Web and particularly the Social Web, new sources of information about musical tracks have become available.

*Tags.* Community-provided annotations (tags) and comments made by users are common on many social platforms and can be obtained at comparably low cost through a crawling process or with predefined APIs. Although such annotations can be rich and diverse, the perception of music is again subjective and can even be influenced by the

perception of other people [McDermott 2012]. Furthermore, the tag quality can be low and less reliable than expert annotations, thus requiring data improvement or noise removal protocols. Finally, a sufficient number of tags might only be available for the most popular tracks but not for tracks from the long tail [Celma 2010].

*Ratings.* Explicit ratings are also typical on the Social Web and include 1-to-5 rating scales like on iTunes and “like” or “dislike” statements as on last.fm, as well as other forms of expressing an evaluation for a track, such as through a wall posting on Facebook [Germain and Chakareski 2013] or a tweet on Twitter [Hauger et al. 2013]. Ratings play a major role in RS research and are also particularly valuable for playlist generation, as they can be used to create playlists that match the taste of the individual user. When relying on rating information, some challenges exist, including the problem of data sparsity, especially for the tracks from the long tail, or a positivity bias, such as the phenomenon that most of the ratings are highly positive and negative feedback is rare [Celma 2010].

*The social graph.* On Social Web platforms, people are connected to each other through “friend” relationships or similar concepts, forming a social graph. When we assume that people who are connected on such platforms share similar tastes, we can exploit such relations and create playlists containing tracks that are liked by the friends of a given user. A special type of relationship in that context are artist networks and the connections between individual artists, such as those available on MySpace or Facebook [Germain and Chakareski 2013]. As we will see later on, artists can play an important role in the playlist generation process. Knowing which artists are connected and related in some form can thus represent a valuable information.

*3.1.4. Usage Data.* Web music platforms often provide information about the listening behavior of their users, including listening logs for individual tracks or playlists or about the general popularity of individual tracks. Such information about patterns and regularities can serve as a valuable input to the playlist generation process [Zheleva et al. 2010].

*Track popularity.* The popularity of a track can be quite easily obtained based, for example, on chart positions, playcounts, or occurrences in user playlists. However, relying only on popularity to generate playlists may not be sufficient. In particular, playlists that only contain popular tracks can be of limited value when the goal is the discovery of new tracks. Still, when limited information is available about a new user, starting with relatively popular tracks can represent a reasonable initial strategy. Moreover, depending on the available data and the chosen evaluation measure, it has been shown that very simple popularity-based algorithms can outperform sophisticated algorithms in more general music recommendation scenarios [McFee et al. 2012].

*Listening logs.* Listening logs, which are available from some music platforms, can tell us much about the taste and preferences of a user and thus represent a valuable information source for personalized playlist generation. The logs will reveal how often certain tracks are listened to, and in addition, if “skip” actions are recorded, further inferences about the user taste or the suitability of a track within a playlist or listening context can be made [Hu and Ogiwara 2011]. Listening logs have to be interpreted with care however, as users may play some tracks in their full length simply because they were not listening to them [Bosteels et al. 2009].

*Manually created playlists.* As playlists shared by users are designed in a manual process for a certain purpose, to some extent we can assume that many of these lists obey certain quality criteria, for example, in terms of homogeneity, with respect to

their overall mood or at least the genre. Such handcrafted playlists additionally might be based on relationships between tracks that cannot be captured using the audio signal or metadata [Maillet et al. 2009; McFee and Lanckriet 2011]. Once detected, these patterns and relationships can be used to guide the automated generation of playlists.

### 3.2. Specifying Target Characteristics for Playlists

Automating the playlist generation process requires that the desired target characteristics are specified in a machine-interpretable form. These characteristics correspond to the overall purpose or intended goals of the playlist. We can identify three types of such information that can be fed into the automated creation process: explicit preferences and constraints, past user preferences, and contextual and sensor information, which we discuss in the next sections.

*3.2.1. Explicit Preferences and Constraints.* We identified four typical ways for users to specify their short-term preferences and constraints in the literature—seed tracks; free-form keywords; explicit, pre-defined constraints; and real-time feedback about tracks:

- (1) The provision of seed tracks by the user can be done in different ways. The typical options are (a) the selection of a single seed track [Logan 2002], (b) the specification of a start and end track [Flexer et al. 2008], (c) a full list of tracks to be sorted or mixed with other tracks [Logan 2004], or (d) a list of tracks already contained in the playlist or recently played (playlist history) [Baur et al. 2010; Hariri et al. 2012].
- (2) Free-form inputs are common on Web music services such as last.fm. The users can enter keywords which are then matched with the tracks. Typical inputs could be artist names, track names, or genres. The keywords may also represent the user's mood or information about the user's activity or environment [Reddy and Mascia 2006; Jensen et al. 2007; Meyers 2007; Coelho et al. 2013].
- (3) Explicit predefined constraints are typically provided with standard user interface elements such as check boxes, sliders, or selection boxes [Pachet et al. 2000; Reddy and Mascia 2006; Sandvold et al. 2006]. The constraints can concern various aspects including the desired genre, tempo, year of release, and so forth.
- (4) Finally, users can express their preferences by immediately rating (like/dislike) the tracks in the created playlist [Pauws and Eggen 2002] or by skipping individual tracks broadcasted by a personalized online radio station [Pampalk et al. 2005].

*3.2.2. Past User Preferences.* The user's past preferences can represent another input guiding the playlist creation process. Such preferences can be explicitly stated (ratings) but also estimated by the system. Typical sources for estimates are the user's personal track library on the device, the listening history, or playlists that the user has created in the past.

The *library of tracks* owned by the user can be seen as an indicator for the long-term taste of the user. If the library is huge, it is unclear if the user has ever heard or even knows all of the tracks [Lamere 2012]. Thus, this sort of information can be unreliable and, in addition, does not contain “negative” feedback. The user's *listening logs* are therefore often more informative, especially when skip actions are recorded. Finally, when the user has created playlists in the past, it is intuitive to assume that the user knows and likes the tracks in the playlists. Again, however, the playlists do not provide us with information about what the user does not like.

*3.2.3. Contextual and Sensor Information.* The suitability of a certain track can finally be determined by incorporating short-term contextual information about the environment and the listeners. With the increasing capabilities of mobile devices and the immediate

Table I. Examples of Additional Information Used in the Playlist Generation Process

Information	Sensor	Reference
Physiological state	ECG, Accelerometer	Oliver et al. [Oliver and Flores-Mangas 2006; Oliver and Kreger-Stickles 2006] propose to determine the user's physiological state with the help of sensors and use it to enhance their exercise performance. Similarly, Biehl et al. [2006] select tracks based on the user's heartbeats per minute and pace.
Environment	Cameras, trackers	Cameras and trackers can provide information such as the number of clients on a dance floor in the context of an automated club DJ [Cliff 2006].
Sound, noise	Microphones	Microphones were used to obtain information such as the sound of walking [Moens et al. 2010; Jylhä et al. 2012] or surrounding noise [Reddy and Mascia 2006].
Location	GPS	The location of the user can be provided by a GPS [Reddy and Mascia 2006].
Time	Clock	Time information can easily be accessed from the computer or mobile device on which the system is deployed [Liu et al. 2009; Hu and Ogihara 2011].
Other context information	Web	Information such as the temperature of a city, the traffic level, and so forth, can be found on the Web [Reddy and Mascia 2006].

availability of various types of information on the Web, the opportunities to include additional contextual and sensor information continuously increase. In the literature, a number of such approaches have been proposed, including the types of sensors and information shown in Table I.

#### 4. PLAYLIST GENERATION ALGORITHMS

Next, we will review strategies for playlist generation. We organize these strategies into seven categories: similarity-based algorithms, collaborative filtering, frequent pattern mining, statistical models, case-based reasoning, discrete optimization, and hybrids.

##### 4.1. Similarity-Based Algorithms

The coherence of the tracks is a typical quality criterion for playlists [Logan 2004; Knees et al. 2006; Fields et al. 2008]. Therefore, selecting and ordering tracks based on their similarities is an obvious strategy to generate playlists. The core of any similarity-based approach is its *distance function*, which characterizes the closeness of two tracks. How the distance function is actually designed depends on the available data, which could include the raw audio signal along with the features that can be derived from it, but also metadata, such as the artists, the genres, playcounts, or ratings [Slaney and White 2007]. In many cases, a signature or model of each track is determined first, in which the distance function is then applied. Typical examples for such functions applied on more abstract models of a track's features are the earth-mover's distance [Logan 2004], the Kullback-Leibler (KL) divergence [Vignoli and Pauws 2005], or the Euclidean distance [Knees et al. 2006].

Once these similarities are computed, different strategies of determining the tracks in the playlist are possible. In Lehtiniemi and Seppänen [2007], for example, tracks are selected based on their distance to the seed track. Another approach is to repeatedly select a track according to the similarity with the previously selected one [Chen et al. 2012]. Quite differently, one could also simply select those tracks that have the highest similarity to the tracks that the user liked [Pampalk et al. 2005; Gartner et al. 2007], which corresponds to a typical content-based recommendation approach. Clustering

techniques also rely on similarity functions (e.g., see Pauws and Eggen [2002] and Dopler et al. [2008]) and exploit similarities based on artist information and other metadata or audio features. Ragno et al. [2005], Pohle et al. [2007], and Hartono and Yoshitake [2013] finally represent examples of a very specific way of using the track similarities, where the goal is to optimize the playlist in a way that the sum of the similarities is maximized.

One possible drawback of pure similarity-based approaches is that the similarity (homogeneity) is often considered as the main, and possibly only, quality criterion. Therefore, the danger exists that such playlists are too homogeneous or contain too many tracks from the same artist and are thus not well suited for the discovery of new tracks. Depending on the application scenario and the specific technique, similarity-based approaches can also be computationally expensive [Budhrajya et al. 2012].

One way to reduce the computational complexity is to rely only on artist similarities, as done in Hauver and French [2001], Dopler et al. [2008], and Fields et al. [2008]. The similarity of artists can be estimated in different ways, such as based on artist ratings, the similarity of artist descriptions, or based on the relationships between artists on music platforms like MySpace. Still, the resulting playlists can in some cases be not very coherent, as there are artists who perform music of different genres or styles.

## 4.2. Collaborative Filtering

Collaborative filtering is the predominant approach in the field of RS [Jannach et al. 2012] and is solely based on community-provided ratings. These ratings can be either provided explicitly or automatically collected (e.g., by recording track listening logs). In principle, any of the various collaborative filtering techniques can be used for playlist generation, provided that there exists a reasonably large dataset with ratings. Still, collaborative filtering approaches were not designed for the specific challenges of playlist generation, and aspects such as list homogeneity or smooth transitions between the tracks have to be addressed separately. In addition, in contrast to the large rating datasets in general RS research, not many datasets exist in the music domain and research is often based on very small samples [Liu et al. 2009].

Since rating information can be sparse, collaborative filtering methods are often combined with other techniques. Chedrawy and Abidi [2009], for example, propose a hybrid system that is based on a semantic similarity function and multiattribute ratings from users on playlists with respect to lyrics, rhythm, tunes, performance, and overall likability. In their approach, they first identify additional existing playlists that are relevant for the user and then construct a new *personalized* playlist using tracks from these playlists.

A different way of applying collaborative filtering techniques for playlist generation is to consider *playlists as users* and tracks as items in the sense of a binary user-item rating matrix. The idea is that when we are given a playlist that currently is being constructed, we consider the set of already-selected tracks as the positive ratings. Then, we can for example apply a nearest-neighbor approach and look for similar playlists that contain some of these tracks. Tracks appearing in similar playlists can then be considered as candidates for inclusion in the new playlist. This last approach was for instance applied and compared to the more elaborate Bayesian personalized ranking (BPR) approach of Rendle et al. [2009] in Hariri et al. [2012].

One of the main limitations of such approaches is what is called the *new user problem*—that is, the question of what to additionally include in the playlist when not much information about the user is available (or the set of start tracks is small). This problem is even more important when considering playlists as users, as each playlist to generate is equivalent to a new user.

### 4.3. Frequent Pattern Mining

The neighborhood-based method for playlist generation mentioned previously relies on the co-occurrence of items in playlists. Frequent pattern mining approaches are based on a similar principle. However, instead of only looking on local neighbors for a given playlist, they try to identify global patterns in the data. There are two types of such patterns: association rules [Agrawal et al. 1993] and sequential patterns [Agrawal and Srikant 1995]. Association rules are often applied for shopping basket analysis problems and have the form  $A \Rightarrow B$ , where both  $A$  and  $B$  are sets of items and the rule expresses “whenever  $A$  was bought, also the items in  $B$  were bought with some probability” [Han and Kamber 2006]. Sequential patterns are similar, except that the order of the elements in  $A$  and  $B$  is relevant. Whether or not sequential patterns are more accurate and preferable over association rules depends on the data characteristics.

Frequent patterns for playlist generation can be mined in samples of handcrafted playlists provided by the user community. When there are tracks that frequently appear together in such playlists, the assumption can be made that the tracks have something in common and that they both fulfill the target characteristics that the creator of the playlist had in mind. The mined patterns can be used for playlist generation as follows. Given a start track or a recent listening history, search for patterns whose left-hand side elements are contained in the history and consider the elements on the right-hand side as candidates to extend the playlist.

To the best of our knowledge, no “pure” application of frequent patterns for playlist generation has been proposed so far. The most similar approach are  $n$ -gram models, such as, Markov models of order  $n - 1$  in which the transition probabilities are derived from the track co-occurrences, which can be considered to be similar to contiguous sequential patterns of fixed size  $n$ . Crampes et al. [2007] propose to train a trigram model ( $n$ -grams with  $n = 3$ ) on playlists made by DJs. Chen et al. [2012] include a bigram model smoothed with the Witten-Bell discounting [Jurafsky and Martin 2009] in a comparative evaluation. However, mining for frequently co-occurring tracks is only one option, as one could also mine patterns of frequently co-occurring artists or genres in playlists. Hariri et al. [2012], for example, propose to mine sequential patterns of latent topics based on the tags attached to the tracks.

One of the main advantages of frequent patterns is that they allow us to implicitly reproduce the observed characteristics of manually defined playlists without looking at the “content” of the tracks. A possible drawback is that the quality of the generated playlists depends on the number and quality of the playlists used for rule extraction.

### 4.4. Statistical Models

Playlists can also be generated using statistical models.<sup>3</sup> One of the simplest approaches is to assume that each track selection only depends on the previous track, which corresponds to Markov models. The transition probabilities for such models can be estimated based on various strategies, using, for example, track co-occurrences in playlists (thus leading to the aforementioned bigram model) or the similarity of the audio signal or the metadata.

<sup>3</sup>Notice that frequent pattern mining approaches can also be considered as statistical estimation methods. In this article, however, we dedicate a whole section to them because they represent one of the most straightforward solutions to playlist generation, a number of algorithms for efficient pattern mining already exist, and our own recent work shows that they work comparably well [Bonnin and Jannach 2013a]. Interestingly, these methods have not often been used in the literature. Similarly, collaborative filtering can use statistical models. However, considering that it is not the only possible strategy and since collaborative filtering is the predominant approach of building RS, we dedicate an entire section to this approach.

A comparison of different basic strategies can be found in McFee and Lanckriet [2011], who compare Markov models based on a uniform distribution, artist popularity, tags, and the audio signal. The experiments show that artist popularity leads to the best results in terms of the average log-likelihood measure (cf. Section 5) and that an optimized mixture of all models further enhances the results. In Chen et al. [2012], a more sophisticated Markov model is proposed—latent Markov embedding (LME). The authors represent tracks as points in the Euclidean space and derive the transition probabilities based on the Euclidean distance. The coordinates of the track are learned using a likelihood maximization heuristic. In Moore et al. [2012], the same authors extend their model by including tag information. Another statistical approach was presented in McFee and Lanckriet [2012], which is based on random walks on a hypergraph—that is, an undirected graph whose edges can connect any number of vertices. Based on the edge weights, their playlist generator first selects an edge containing the previously selected track and then chooses one of the remaining tracks of that edge in a uniform manner. The weights of the edges are learned using the L-BFGS-B algorithm [Byrd et al. 1995].

A major drawback of some of these approaches is assuming that the Markov property may actually be too strong for the domain. Researchers have therefore proposed models in which more complex dependencies can be expressed using graphical models. In Zheleva et al. [2010], two models—the taste model and the session model—are proposed. Both models select tracks based on distributions over tracks using latent clusters that were obtained by applying latent Dirichlet allocation (LDA) [Blei et al. 2003] on the usage data. The difference between the two models is that the first one only uses one latent cluster per listening session and can thus be considered as reflecting the user’s taste. The second model can use a different cluster for each track, which the authors consider as reflecting the user’s current “mood” (contextual taste). Later on, a similar approach was proposed by Hariri et al. [2013]. In their approach, the latent clusters can also incorporate tag information, which allows the use of queries for specifying preferences.

In general, one main advantage of using statistical models is that many algorithms exist to optimize the values of the parameters, as illustrated by Chen et al. [2012] and McFee and Lanckriet [2012]. Moreover, different statistical models can be easily combined through linear interpolations, which can help to further enhance the results [McFee and Lanckriet 2011]. However, one possible drawback is that the learning process can be highly time consuming, as will be discussed in Section 6.

#### 4.5. Case-Based Reasoning

The general idea of case-based reasoning techniques is to exploit information about problem settings (cases) encountered in the past to solve new problems [Althoff 2001]. Case-based reasoning approaches therefore first store a set of representative cases together with their solutions. Given a new problem instance, the case repository is searched for past similar cases whose solution is then adapted according to the specifics of the new case.

A case-based reasoning technique has been proposed for playlist generation in Baccigalupo and Plaza [2006]. The case base in their scenario consists of a set of playlists created by a user community. Within these playlists, frequent subsequences (patterns) are identified. In contrast to typical case-based reasoning approaches, the goal of their work is not to find the most similar playlists given some seed track, but to find those considered to be the most “useful,” for example, in terms of diversity. The elements of the retrieved playlists are then combined to generate a new playlist for a given seed track.

Later on, the same authors proposed to apply case-based reasoning techniques in the specific context of a broadcasting radio where the listeners can give feedback on the selected tracks [Baccigalupo and Plaza 2007]. In this setting, the system uses a case representation of the user preferences to adapt the selection of tracks for a given group of listeners.

One possible advantage of typical case-based reasoning techniques is that their computational complexity can be comparably low when only a limited number of cases is used. The specific techniques proposed by Baccigalupo et al., however, may be limited in their scalability, as they scan the entire database upon new playlist generation requests.

#### 4.6. Discrete Optimization

A quite different approach is to consider playlist generation as a discrete optimization problem. Given a set of tracks, their characteristics, and a set of explicitly specified constraints capturing the desired characteristics, the goal is to create one arbitrary or an optimal sequence of tracks that satisfies the constraints. Examples for constraints could be a given start or end track, of which the tracks in the playlist are taken from at least  $n$  different genres, or constraints related to the transitions between the tracks.

Technically, the problem can be modeled as an integer linear program as done in Alghoniemy and Tewfik [2000, 2001]. In an alternative formulation, Pachet et al. [2000] proposed to model the task as a constraint satisfaction problem (CSP) and apply the existing search algorithms to select sequences of tracks that satisfy various constraints, in that case from a comparably small catalogue.

When we adopt a standard CSP formulation, no solution can be found if the problem is overconstrained. In such situations, however, “soft constraints” approaches can be applied. Aucouturier and Pachet [2002a, 2002c] propose building playlists by iteratively enhancing a randomly chosen playlist through a local search procedure. The quality and suitability of the resulting playlist can be assessed by counting the number of satisfied constraints. From a technical perspective, different local search procedures are possible, including simulated annealing [Pauws et al. 2006, 2008] or genetic algorithms [Hsu and Chung 2011].

The advantage of these approaches is that if the used background knowledge is accurate, only playlists will be generated that satisfy all or most of the target characteristics. A potential drawback is the high computational complexity, which makes the problem intractable for larger music collections [Pauws et al. 2008].

#### 4.7. Hybrid Techniques

In the field of RS, hybridization is often used to combine the advantages of different techniques and at the same time avoid the drawbacks of individual techniques. Burke [2002] proposes a classification in seven categories for hybridization of RS, which is reproduced in Table II.

Hybrid approaches can also be found in the literature for playlist generation. McFee and Lanckriet [2011] use a weighted hybridization approach—a mixture model—that combines algorithms that rely on track similarity, track familiarity, and the uniform distribution. The same hybridization method is used by Chedrawy and Abidi [2009], who combine a multicriteria  $k$ NN-based collaborative filtering algorithm with an ontology-based item similarity approach. Hariri et al. [2012] also use a weighting strategy and combine a frequent pattern mining approach with the  $k$ NN-approach. This last work additionally uses a feature augmentation strategy, as they first extract latent topics from tags and then extract frequent patterns from them. Meyers [2007] offers an example of feature combination. Features extracted from the audio file and lyrics are combined to generate playlists matching a certain mood. Generally, approaches

Table II. Hybridization Methods

Hybridization Method	Description
Weighted	The scores from several recommendation techniques are combined.
Switching	The system switches between recommendation techniques depending on the current situation.
Mixed	Recommendations from several different recommenders are presented at the same time.
Feature combination	Features from different recommendation data sources are thrown together into a single recommendation algorithm.
Feature augmentation	The output of one technique is used as an input feature to another one.
Cascade	One recommender refines the recommendations given by another.
Metalevel	The model learned by one recommender is used as input to another.

in which features are inferred from the audio signal and then combined with manually annotated metadata can be considered as a form of feature combination (e.g., see [Aucouturier and Pachet 2002c; Vignoli and Pauws 2005]). For more information about hybrid recommendations outside of the specific scope of playlist generation, see Burke [2002], Lampropoulos et al. [2012], and Hornung et al. 2013].

## 5. EVALUATING PLAYLIST GENERATION SYSTEMS

In general, *user satisfaction* could be considered as the ultimate criterion to assess the quality or usefulness of the generated playlists. However, whether a user is satisfied with a playlist can depend on various criteria. The satisfaction can be influenced by the extent to which the list matches its intended purpose; fulfills the desired target characteristics; or is in line with the user's expectations, including aspects of taste, context, or mood [Fields 2011]. In the following, we will first discuss general quality criteria for playlists from the literature and then review approaches to evaluate the quality of playlist generation systems.

### 5.1. Quality Criteria for Playlists

One way to determine the general principles for a good playlist design is to analyze the characteristics of playlists that are created and shared by users. Such an analysis can be found in Slaney and White [2006], where the role of the *diversity* of playlists was in the focus. The chosen research approach was to use an objective measure that captures the diversity of playlists and then to analyze if users tend to create rather diverse or more homogeneous playlists. More precisely, tracks were represented as points in a genre space, and the diversity was defined based on the volume of the enclosing ellipsoid of the tracks of the playlist. This metric was applied to 887 user playlists containing about 30,000 tracks, and the output shows that diversity actually matters for playlists.<sup>4</sup> Sarroff and Casey [2012] investigated which track features are the most important for track transitions. Their results highlighted the importance of fade durations and the mean timbre of track endings and beginnings.

A different approach to determine the guiding rules of playlist creators was chosen by Cunningham et al. [2006], who analyzed the messages posted on the Art of the Mix forum. Their analysis revealed that the most typical principles for track selection mentioned in these messages are the artist (e.g., the best tracks of a given artist), the genre, the style, the event or activity, the semantics (e.g., "romance"), and the mood. Additional criteria mentioned in these messages are the maximum number of consecutive tracks of the same artist or genre or the smoothness of the transitions.

<sup>4</sup>To which extent the chosen objective diversity measure corresponds to the diversity *perceived* by the users was not analyzed in this work.

Some playlist creators also organize their playlists in a way that the most popular or most “important” tracks appear at the end of the playlist as a sort of “grand finale.”

Another possible approach is to perform user studies aimed at determining quality criteria for playlists. An example of such a study was presented in Hansen and Golbeck [2009]. Their study involved 52 participants, and one of the main results was that the choice of the first and the last tracks was considered to be particularly important. In a smaller study involving 24 participants, Andric and Haus [2005] showed that for playlists that contained only tracks that the users preferred in general, the order of the tracks and even the selection of tracks had a low impact on the perceived quality. In Reynolds et al. [2008], the results of an online music survey with 750 participants were presented, which showed that the user’s context or environment (location, activity, temperature, lighting, or weather) can have a strong influence both on the mood of the listeners and the track selection. Another smaller study involving 8 participants was done in Lee et al. [2011]. The identified key criteria included the variety of the list as well as metadata features including lyrics, theme, and time. Furthermore, having a mix of familiar tracks and new music was considered to be a factor for satisfaction. Finally, Kamalzadeh et al. [2012] presented the results of a larger study involving 222 participants, which again showed the importance of criteria related to the mood, the genres, and the artists of the tracks.

Overall, such analyses give us important insights about what users find important and which criteria can be used to judge the quality of playlists. However, in our view, more research in that respect is required, and one possible direction for future work could be a systematic approach for quality criteria identification as was done in Wang and Strong [1996] for the information quality field.

## 5.2. Evaluation Approaches for Automatically Generated Playlists

Various approaches have been proposed in the literature to evaluate the quality of playlists. McFee and Lanckriet [2011], for example, listed three categories of evaluation strategies: human evaluation, semantic cohesion, and sequence prediction. In this work, we propose to organize the possible evaluation approaches in four more general categories: user studies, log analysis, objective measures, and comparison with handcrafted playlists.

*5.2.1. User Studies.* User studies are a form of evaluation that helps us to determine the *perceived* quality of playlists with good confidence. One of the major drawbacks of such studies, however, is that they are time consuming and expensive. The particular problem in the domain is that in many experimental designs, the participants are asked to listen to all tracks of the playlist. This can take a long time, even if only excerpts are played. In settings where the participants evaluate the playlists based only on meta-information or evaluate playlists consisting of tracks they already know, the obtained results might be biased and depend on the general popularity of the individual tracks.

A typical problem in the context of user studies reported in the literature is their size. Many studies are based only on 10 to 20 participants who evaluate playlists with a questionnaire or give feedback using an application that was developed for the study. The latter form of evaluation often scales better. Barrington et al. [2009], for example, present the results of a study involving 185 subjects. Unfortunately, the reproducibility of such designs is limited, as the experiment is based on a specific software application. Finally, the studies in the literature are often based on quite different collections of tracks, and the generalizability of the findings is not always clear.

*5.2.2. Log Analysis.* Instead of asking the users about their subjective evaluations, one can look at how the generated playlists are accepted by them by analyzing their

behavior using listening and interaction logs. Such logs contain information about how often each user listened to each playlist. In addition, explicit “like” or “dislike” statements or user actions such as skipping to the next track can be recorded.

Log analysis is something that would usually be done in *A/B* tests performed by music platform providers. In such a setting, the listeners are split into different groups, and each group would receive playlists that are generated with different techniques. After the test has ended (e.g., after 4 weeks), the logs could be analyzed to determine which playlist generator worked best, for example, based on the like statements or play counts. Obviously, researchers usually cannot perform *A/B* tests with real users. Still, there are at least platforms from which information about the listening behavior can be obtained. In particular, last.fm gives access to the listening logs of its users through their public API. So far, however, only limited research has been based on such logs, the work of Bosteels et al. [2009] being one of the few examples, who used a specific evaluation protocol. They retrieved sequences of tracks of length 22 such that the last two elements were not both accepted or rejected (one was accepted and the other rejected). Different playlist generation heuristics were then compared by counting how often they could rank the accepted track before the rejected track.

**5.2.3. Objective Measures.** With the term *objective measures*, we refer to measures that can be automatically computed for a given playlist and that try to *approximate* a user-perceived *subjective* quality [Cremonesi et al. 2011]. A typical example for such a subjective measure is the diversity or homogeneity of a playlist. One basic way of assessing the diversity of a given playlist could be to simply count the number of different artists or genres appearing in the playlist. Several works base their evaluations on variants of such a strategy and determine the genre distribution only based on metadata [Pohle et al. 2005; Knees et al. 2006; Flexer et al. 2008; Dopler et al. 2008] or on a combination of metadata and the audio signal [Cai et al. 2007]. Some authors simply consider tracks to be similar if they are from the same artist [Slaney and White 2007] or the same genre [Balkema and van der Heijden 2010]. However, this might be too simplistic given that there are artists who cover a broader spectrum of genres and that the tracks of some genres, such as jazz, can be quite heterogeneous.

Considering a single quality criterion might not be sufficient to assess the quality of a playlist. Determining, for example, the homogeneity alone might be misleading, and some works indicate that diversity can be as important as homogeneity [Slaney and White 2006, 2007; Lee et al. 2011; Kamalzadeh et al. 2012]. Therefore, more comprehensive evaluation approaches are required that use multiple measures at the same time and in addition analyze typical trade-offs such as homogeneity versus diversity. Beside diversity and homogeneity, other quality criteria can be approximated through a numeric measure, such as the novelty, the freshness of the tracks, their familiarity with respect to a certain user, the consistency of the mood, and even the smoothness of the transitions (e.g., based on the tempo and instrumentation of the tracks).

A limitation of such approaches is that it is not always clear if the objective measure truly approximates the subjective quality experience. For instance, a playlist covering several genres may actually not be perceived to be diverse, because the genre annotations of tracks can be subjective and some genres might be similar to one other.

**5.2.4. Comparison with Handcrafted Playlists.** The final category of evaluation approaches discussed here is based on the estimation of the ability of the algorithms to reproduce handcrafted playlists—that is, playlists created by music enthusiasts. The assumption is that such playlists respect various of the described quality criteria. An automatically created list should therefore in some form be similar to manually created ones, and a playlist generation algorithm correspondingly should rely on patterns that are similar

to those of real users. In the literature, two measures and protocols can be found that are based on such a strategy.

*Hit rates.* The hit rate is a measure commonly used in the domain of IR. The principle is to count how many of the actually relevant items were contained in a top- $N$  result list for a given query.<sup>5</sup> The measure is also applied in the field of RS [Jannach et al. 2012]. Here, the input is rather a user profile—typically information about items that the user liked in the past—than a query. In both cases, to apply this measure, the ground truth, or whether or not an item is considered to be relevant, has to be known or approximated.

A similar principle can be applied to evaluate the ability of the algorithms to generate good playlists. The idea is to take a handcrafted playlist and “hide” some of the tracks. The task of the playlist generation system is then to guess the hidden items. Each correctly guessed hidden item corresponds to a hit, and the more often a playlist generator “agrees” with the handcrafted list, the better.

One limitation of such a strategy is that it is based on the assumption that the hidden tracks are the only relevant ones. Clearly, however, many other tracks may also be considered relevant. If the number of these tracks could be known, then recommendation lists of exactly this size could be computed. Such lists should then at least contain the test track. Although it is impossible to determine the exact number of relevant tracks, the length of the recommendation list can be considered an assumption about the average number of truly relevant tracks. Therefore, it makes sense to look at how the results change when the size of the recommendation lists varies, and to use values well beyond what users might conceivably browse through. Notice that this strategy should not be considered as a means to estimate the number of true positives, but mainly as a way to observe the results under different assumptions.

In general, any subset of the playlist elements could be hidden. However, when we remove more than one element, the playlist generation algorithm might not be able to fully take the role of the transitions between the tracks into account. One option is therefore to hide only one track at a time and consider the preceding tracks in the playlist as an input. Such a protocol was used in Hariri et al. [2012], who hide the last element of each tested playlist.<sup>6</sup>

Technically, to measure the hit rate, we split the available set of playlists into a training set *Train* and a test set *Test*. Then, each playlist in *Test* is split into the current set of given tracks and the remaining set of tracks to be recommended. If we limit ourselves to the recommendation of the last track  $t$  given the history of the  $i$  first tracks  $h = (t_1, \dots, t_i)$ , the hit rate can be measured as follows:

$$\text{HitRate}(\text{Train}, \text{Test}) = \frac{1}{\|\text{Test}\|} \sum_{(h,t) \in \text{Test}} 1_{R(h)}(t), \quad (1)$$

where  $R(h)$  is a recommendation list of a defined length computed by an algorithm based on *Train*,  $h$  is the playlist beginning, and  $1_{R(h)}(t)$  is an indicator function that indicates the membership of  $t$  in  $R(h)$ .<sup>7</sup>

*Average log-likelihood.* The average log-likelihood  $\mathcal{ALL}$  is a measure that is commonly used in the domain of natural language processing to assess the quality of the

<sup>5</sup>Alternative terms for *hit rate* are *true positive rate*, *recall*, and *sensitivity*.

<sup>6</sup>Given that the very last element can play a special role for playlist creators, this strategy might not be optimal. Still, in our evaluations in Section 6, we will rely on this scheme to make our work comparable to existing research.

<sup>7</sup>As an alternative to the use of hit rates to evaluate playlists, researchers sometimes use measures such as the “half-life of user interest” [Platt et al. 2001; Xiao et al. 2009].

modeling techniques [Rosenfeld 2000]. The measure can be used to assess how likely a system is to generate the tracks of a given set of assumedly good playlists and was first used in McFee and Lanckriet [2011]. Conceptually, the generation process is assumed to be based on a weighted random process, where the generation algorithm assigns a probability of being selected to each of the tracks. Thus, the application of this measure requires that the algorithms assign probabilities to the candidate tracks.

Technically, given a test set of playlists, the average log-likelihood  $\mathcal{ALL}$  is given by

$$\mathcal{ALL}(Train, Test) = \frac{1}{\|Test\|} \sum_{(h,t) \in Test} \log(P_{Train}(t | h)), \quad (2)$$

where  $P_{Train}(t | h)$  corresponds to the probability of observing  $t$  given  $h$  according to a model learned based on *Train*. Research works that use this measure for playlists include McFee and Lanckriet [2011], Chen et al. [2012], and Moore et al. [2012].

In contrast to the hit rate, the  $\mathcal{ALL}$  measure is not interpretable on an absolute scale: the possible values range from  $-\infty$  (at least one test track has a probability of 0) and 0 (all probabilities for all test tracks are 1). This measure therefore only allows us to compare the results of different approaches without knowing if the best one is actually good. It can thus be considered as a complementary measure to the hit rate.

Note that when an algorithm assigns a probability of zero to one single test track, we would obtain an  $\mathcal{ALL}$  value of  $-\infty$ . Such zero probabilities must therefore be avoided. This can be achieved through a smoothing step, such as by combining the probabilities of a given model, which were learned on a training dataset, with the probabilities from the uniform distribution using a mixture model (see Section 4). The corresponding weights can then be learned on yet another (optimization) dataset using some optimization algorithm [McFee and Lanckriet 2011]. However, due to the long-tailed distribution of the tracks with respect to their popularity, a large part of the tracks of the optimization dataset do not appear in the training dataset. The larger the number of these tracks, the larger the weights assigned by the optimization algorithm to the uniform distribution. The weights from the uniform distribution might thus be too large and bias the algorithms' estimates too much toward the niche items.

### 5.3. Other Evaluation Criteria for Playlist Generation Systems

So far, we have focused on measures that are related to the quality of the playlists themselves. However, there are additional criteria that can be considered when comparing playlist generation algorithms and the corresponding applications. We will discuss two of them briefly here: scalability and user effort. Other criteria can be found in Vatulkin et al. [2011].

**5.3.1. Scalability and Computational Complexity.** In many application domains, being able to dynamically compute contextualized playlists or playlist continuations is a crucial requirement. Online services such as last.fm should more or less be able to immediately create a playlist for a given set of keywords, and a personalized radio station should be able to react to the skip actions of a user in a reasonable time frame.

In contrast to the domain of RS, where this is a niche topic, scalability and computational complexity are comparably often discussed in the context of playlist generation (e.g., see Pohle et al. [2005], Knees et al. [2006], and Cai et al. [2007]). One possible explanation is that one of the most prominent approaches for playlist generation is discrete optimization, which involves the exploration of large search spaces or the maximization of the number of satisfied constraints. Finally, although learning a statistical model can be done offline, the required running times can also be problematic, particularly when the pool of available tracks is huge.

**5.3.2. User Effort.** Whether or not users are satisfied with a system can largely depend on its user interface. Design, ergonomics, and particularly how much effort is required by the user to specify the target characteristics of a playlist can play a major role. An example of research in that direction is found in Hilliges et al. [2006], where the goal is to make the playlist creation process more convenient for users by using an audio similarity map with colors for mood constrains. Another example is the Rush system [Baur et al. 2010], in which users receive repeated recommendations and can thereby construct the playlist in an incremental manner. In the experiments in Pauws and van de Wijdeven [2005], finally, the authors include measures such as the time spent creating a playlist or the number of required user actions and relate them to the resulting playlist quality.

Many popular systems, including last.fm and the iTunes Genius, provide relatively simple mechanisms for the user to specify the desired playlist characteristics. These mechanisms typically include the selection of a seed track or free text input that can correspond to the preferred genre or artist. On last.fm, users can furthermore decide to listen to a number of predefined “radios,” whose playlists are determined by the system depending on the listening history of the current user or his or her friends. However, the user cannot directly influence what is played on these radios.

Generally, in our view, more research is required with respect to what can be expected from the user. One question here is how much effort users are willing to spend on the input specification and playlist refinement process. Another question is how to enable the users to specify their preferences in an intuitive and convenient way.

## 6. A COMPARATIVE EVALUATION OF PLAYLIST GENERATION TECHNIQUES

In the following section, we report the results of a comparative evaluation of typical playlist generation approaches.<sup>8</sup> Although user studies are a reliable means to determine the quality of playlists, we argue that interesting new insights can also be obtained based on historical data. We thus adopt an offline experimental design and base our experiments on large pools of handcrafted playlists shared by music enthusiasts, as described in Section 5.2.4.

As evaluation measures of playlist quality, we use hit rates and the average log likelihood  $\mathcal{ALL}$ . Although these measures do not explicitly address aspects like diversity, homogeneity, track freshness, and others, as mentioned in Section 5.2.4, they implicitly take them all into account and reflect general user preferences. Moreover, one of the goals of our experiments is to identify possible limitations of these measures.

Through our experiments, we aim to analyze the following research questions:

- RQ1:** Are algorithms that take the order of the tracks in their models into account favorable over those that only look at track co-occurrences?
- RQ2:** How important are artist-related information and popularity aspects when creating playlists?
- RQ3:** Do the hit rate and the  $\mathcal{ALL}$  measure agree with respect to the ranking of the algorithms, and are there differences when different datasets are used?
- RQ4:** What is the computational complexity of the different techniques?

### 6.1. Experimental Procedure

Our experiments are based on pools of handcrafted playlists that we randomly split into training (90%) and test set (10%) in a 10-fold cross-validation procedure. The task of the playlist generation algorithms is to compute scores (which can be probabilities) for the possible next tracks given the first tracks of each playlist in the test set.

<sup>8</sup>See also our previous works (Bonnin and Jannach [2013a, 2013b]).

Table III. Details of the Algorithms Used in the Evaluation

Algorithm	Details
$k$ NN	<p>As for every <math>k</math>NN method, a similarity function is required to determine the neighborhood. We calculate this using a binary cosine similarity measure between the playlist history <math>h</math> and each handcrafted playlist <math>p</math>:</p> $\text{sim}(h, p) = \frac{\ h \cap p\ }{\sqrt{\ h\  \ p\ }}.$ <p>Given <math>N_h</math> nearest neighbors for a playlist history <math>h</math>, the score of a track <math>t</math> is defined as</p> $\text{score}_{k\text{NN}}(t, h) = \sum_{n \in N_h} \text{sim}_p(h, n) \cdot \mathbf{1}_n(t), \quad (3)$ <p>where <math>\mathbf{1}_n(t) = 1</math> if <math>n</math> contains <math>t</math> and 0 otherwise.</p>
Association Rules (AR) and Sequential Patterns (SP)	<p>We use the Apriori algorithm for the extraction of association rules and a simple variant thereof for the sequential patterns. For both techniques, we can vary the minimum values for support and confidence as well as the maximum size of the patterns and the maximum size of the sliding window.</p> <p>Given a set of extracted patterns, the score of a track <math>t</math> given a playlist history <math>h</math> is computed as follows:</p> $\text{score}_{\text{pattern}}(t, h) = \sum_{\omega \in \Omega} \text{confidence}(\omega \rightarrow t), \quad (4)$ <p>where <math>\Omega</math> is the set of all possible rule antecedents that can be considered in <math>h</math> and <math>\text{confidence}(\omega \rightarrow t)</math> the confidence of the rule <math>\omega \rightarrow t</math>—that is, the conditional probability of <math>t</math> given <math>\omega</math>.</p>
Same Artists—Greatest Hits (SAGH)	<p>Given a playlist history <math>h</math>, the algorithm determines the set of artists appearing in <math>h</math> and then assigns scores to their tracks depending on how often these tracks appeared in playlists in the training set. Given a history <math>h</math>, the score for a track <math>t</math> with artist <math>a</math> is computed as follows:</p> $\text{score}_{\text{SAGH}}(t, h) = \text{counts}(t) \cdot \mathbf{1}_h(a), \quad (5)$ <p>where <math>\mathbf{1}_h(a) = 1</math> if <math>a</math> is also the artist of one of the track in <math>h</math> and 0 otherwise, and <math>\text{counts}(t)</math> is the number of occurrences of <math>t</math> in the dataset, which corresponds to the greatest hits of the dataset.</p>
Collocated Artists—Greatest Hits (CAGH)	<p>We propose to select the greatest hits of artists that already appear in <math>h</math> or are similar to the artists of <math>h</math>. We compute the similarity between two artists <math>a</math> and <math>b</math> as follows:</p> $\text{sim}_a(a, b) = \frac{\sum_p (\mathbf{1}_p(a) \cdot \mathbf{1}_p(b))}{\sqrt{\sum_p \mathbf{1}_p(a) \cdot \sum_p \mathbf{1}_p(b)}},$ <p>where <math>\mathbf{1}_p(a) = 1</math> if playlist <math>p</math> contains <math>a</math> and 0 otherwise. The similarity thus depends on the collocations of artists within playlists, which can be computed offline. Our proposed formula for the computation of the score of a next track <math>t</math> with artist <math>a</math> given a playlist beginning <math>h</math> is as follows:</p> $\text{score}_{\text{CAGH}}(t, h) = \sum_{b \in A_h} \text{sim}_a(a, b) \cdot \text{counts}(t), \quad (6)$ <p>where <math>A_h</math> is the set of artist names of the tracks in <math>h</math> and <math>\text{counts}(t)</math> is the number of occurrences of <math>t</math> in the dataset. Similar to the frequent patterns, we can vary the minimum values of artist co-occurrences.</p>

**6.1.1. Algorithms.** We used six algorithms with different characteristics to answer our research questions. The specific details about the individual techniques are given in Table III:

- (1)  $k$ NN: This algorithm takes the playlist beginning  $h$  and looks for other playlists in the training data that contain the same tracks. The main assumption is that if

there are additional tracks in a very similar playlist, chances are good that these tracks suit the playlist beginning  $h$ . Recent works have shown that this technique represents a strong baseline [Hariri et al. 2012].<sup>9</sup>

- (2) *Association Rules (AR)*: This technique (see Section 4) generates association rules from the training set. Given a history  $h$ , the playlist generation algorithm looks for rules whose left-hand side matches some tracks in  $h$ . The confidence of the tracks on the right-hand side is then used to compute the scores of the tracks.
- (3) *Sequential Patterns (SP)*: The way we implemented this algorithm is equivalent to the previous one with the only difference that the order of the elements in the rules is taken into account.<sup>10</sup>
- (4) *Same Artist—Greatest Hits (SAGH)*: This method is based on a proposal from McFee et al. [2012], and the idea is simply to only consider tracks of artists that already appear in  $h$ . The score of the tracks is then based on their popularity. In our case, popularity is measured in terms of how often a track appears in the playlists in the training set. We include this method in our evaluation to assess the importance of artist information and popularity.
- (5) *Collocated Artists—Greatest Hits (CAGH)*: We propose this method as an extension to the previous one. Instead of only considering artists that appear in  $h$ , we also include the greatest hits of “similar” artists. As a measure of artist similarity, we simply count how often two artists appear together in the playlists of the training set—that is, how often they are collocated.
- (6) *PopRank*: This algorithm simply recommends the tracks by decreasing order of their overall occurrence in the playlists.

**6.1.2. Datasets.** We used three sets of playlists that we obtained from different music platforms.

- (1) *last.fm*: We retrieved this dataset, which comprises 50,000 playlists, from last.fm through their public API.
- (2) *Art of the Mix (AotM)*: The dataset was made available by McFee and Lanckriet, and the subset we finally used contained 50,000 playlists.
- (3) *8tracks*: This dataset was provided to us directly by 8tracks. The dataset used in the experiments comprises about 100,000 playlists.

The last.fm dataset contains playlists that were manually created and shared by the users of the platform. When compared to the other datasets, one difference is that creating and sharing playlists is not the main feature of last.fm but merely an additional feature. AotM is a comparably old Web site (from 1997) for creating and sharing playlists. Although it does not allow users to listen to the created playlists, it is extensively used by a very active community of music enthusiasts who exchange ideas and try to maximize the quality of their playlists. 8tracks is a more recent Web site (from 2006) dedicated to sharing playlists and has the advantage that its users can listen to the tracks of the playlists online. However, this feature comes at the cost of some limitations (e.g., only a certain number of tracks of one artist can be contained in one single playlist).

For the last.fm and Art of the Mix datasets, we removed playlists of size 1 and random playlists that contain tracks appearing only once (i.e., only in the considered playlist) until reaching a size of 50,000 playlists. To enhance the quality of the data, we

<sup>9</sup>In fact, the more sophisticated hybrid approach from Hariri et al. [2012] only slightly outperforms the  $k$ NN method for low values of  $k$ . They could also show that the  $k$ NN method is better than the recent learning-to-rank technique BPR [Rendle et al. 2009] and a content-based approach for their dataset.

<sup>10</sup>For a detailed description of similar versions of these approaches, see Nakagawa and Mobasher [2003].

Table IV. Properties of the Datasets

	last.fm	AotM	8tracks
Playlists	50,000	50,000	99,542
Users	47,603	11,655	51,743
Tracks	69,022	453,270	179,779
Avg. tracks/playlist	7.6	19.4	9.7
Avg. track usage count	5.5	2.1	5.3
Artists	11,788	119,878	29,352
Avg. artists/playlist	4.5	16.9	8.9
Avg. artist usage count	32.2	8.1	32.7
Artist reuse rate	31.1%	21.3%	13.8%

used the API of last.fm to correct artist and track misspellings and used the existing MusicBrainz IDs of the resulting tracks as a basis to identify tracks.

More details about the characteristics of the datasets are given in Table IV. The datasets are partially quite different, for example, with respect to the average length of the playlists or how many artists appear on average in a playlist. The row “artist reuse rate” in the table is the percentage of cases when the artist of the last track of a playlist already appeared in the same playlist before. The “usage count” statistics express how often a track or artist on average appeared in the playlists.

Using three datasets with quite different characteristics should allow us to analyze how the different algorithms perform in different situations. At least two characteristics are of major importance: the size of the datasets and the sparsity of track usage. Regarding these characteristics, our datasets have complementary properties: one dataset has a limited size and sparsity (last.fm), one dataset has a limited size and a higher sparsity (AotM), and one dataset has a larger size and a smaller sparsity. Overall, generating playlists using the AotM data is expected to be the most difficult task. However, other factors may play a major role as well, particularly that 8tracks users are only allowed to add two tracks of the same artist in one playlist.

## 6.2. Evaluation Results

### *RQ1: Relevance of Track Order in Models*

Our first set of experiments is concerned with the question, “Does the order matter?” To that purpose, we made a comparison between the association rules and the sequential patterns method. In this experiment, we used hit rates as the only measure to assess the playlist quality and manually optimized the support and confidence threshold values for each dataset.<sup>11</sup> To make our results comparable, we applied the experimental procedure that was used in Hariri et al. [2012].

Figure 1 shows the hit rates for the three datasets. When comparing the absolute values for the different datasets, we can observe that they can be quite different across the datasets. This is not surprising given the different characteristics of the datasets. As expected, the highest values are observed for the last.fm dataset, which has a low sparsity, and the lowest values were obtained for AotM.

Regarding RQ1 and the role of the track orderings, both for the AotM and the last.fm datasets, the sequential patterns lead to statistically significant higher hit rates than

<sup>11</sup>We do not present the results in terms of the  $ACL$  because of the limitations of this measure, as will be presented in RQ3.

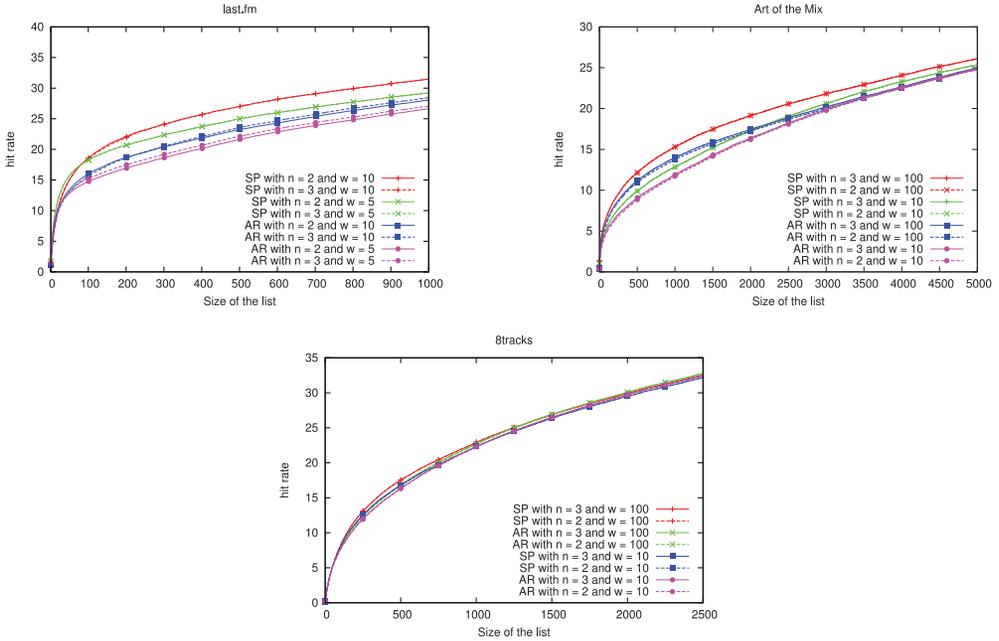


Fig. 1. Hit rates of the different frequent pattern mining techniques for the three datasets. AR, association rules; SP, sequential patterns; n, max. size of patterns; w, window size.

the association rules according to a one-way ANOVA analysis and a Tukey post hoc test (the largest  $p$ -value equals 0.03 on last.fm and  $10^{-5}$  on AotM, and the smallest Cohen's  $d$ -value equals 1.44 on last.fm and 1.21 on AotM). This indicates that integrating information about track orderings can actually help to increase the quality of the generated playlists. The effect, however, seems to depend on the dataset, and no significant difference could be found for the 8tracks dataset.

A side observation in this context is that larger window sizes seem to be preferable for the sequential patterns approach. This means that it is better not to look only at the very recent playlist history. In the subsequent experiments, we will therefore use the sequential patterns approach with larger window sizes.

### RQ2: The Role of Artist-Related Information and Track Popularity

To determine to which extent artist and popularity information are relevant for the playlist generation process, we compared the hit rates of the artist-agnostic techniques  $k$ NN, Sequential Patterns,<sup>12</sup> and PopRank to the SAGH and CAGH algorithms. We ran the  $k$ NN method with three comparably large neighborhood sizes to see how it influences the results. We also manually optimized a minimum threshold regarding artist co-occurrences for the computation of the artist similarities. Again, we used hit rates as the measure to assess the playlist quality and hid only the last track.

Figure 2 shows the hit rates obtained by the different approaches for the three datasets. Overall, we can see that the CAGH approach outperforms the other approaches. These differences were again statistically significant for the last.fm and the AotM datasets (the largest  $p$ -value equals  $10^{-4}$  on last.fm and  $10^{-8}$  on AotM, and the smallest Cohen's  $d$ -value equals 2.73 on last.fm and 5.11 on AotM). With respect to

<sup>12</sup>While sequential patterns do not exploit information about artists, the popularity of tracks is implicitly taken into account, as the support and confidence of the rules are based on the occurrences of their elements in the playlists.

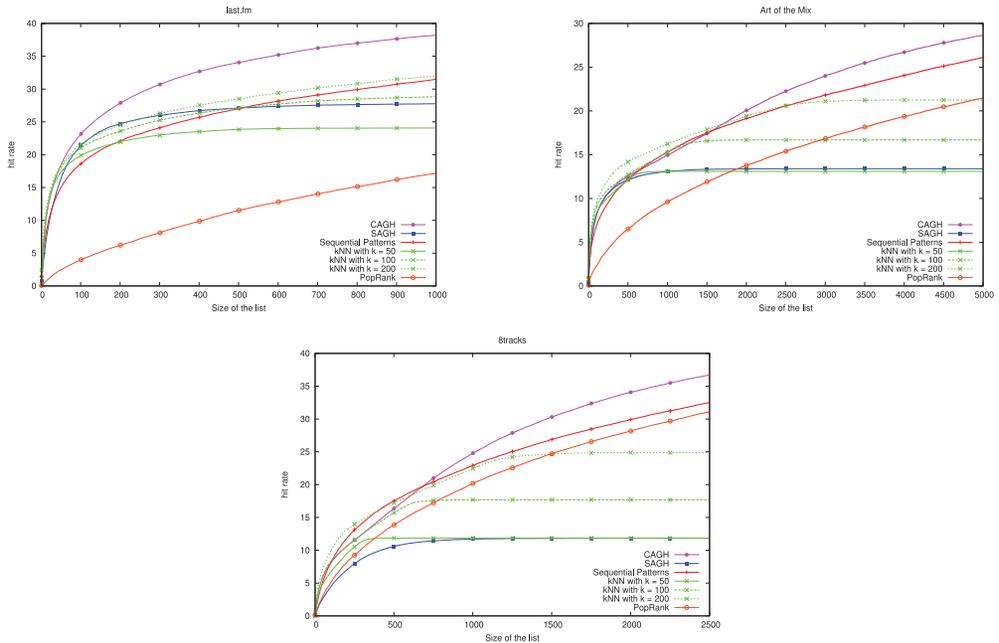


Fig. 2. Hit rates of the different algorithms.

RQ2, this means that including tracks from similar artists can be a good strategy to improve the quality of the generated playlists. In fact, even a very simple strategy like CAGH seems to be able to capture interesting characteristics of what users consider to be a good playlist.

It is worth noting that playlist creators on 8tracks are not allowed to have the same artist more than twice in one playlist. Still, the CAGH algorithm performs well and would perhaps be even better if this constraint were not there. Another interesting particularity of the results of this dataset is the surprisingly high hit rates of the popularity-based approach for lists longer than 1,500 tracks. In other words, playlist creators on 8tracks include one of the 1,500 most popular tracks (out of about 180,000) as the last track of the playlist about 25% of the time, which means that the head of the long-tailed distribution of track usage is particularly large. The sequential patterns approach generally works well on all datasets. Considering the track ordering and to some extent the popularity of the tracks thus appears to be a good strategy in general.

These results also give an insight to how the sparsity of track usage affects the performance of the different algorithms. As for the previous results, we can see that the lowest values were obtained for the AotM dataset, which has the highest sparsity. However, this is only true for some of the algorithms. In particular, the SAGH algorithm has higher hit rates on AotM than on 8tracks, which could be caused by the aforementioned constraint for the playlists of 8tracks. However, the same phenomenon can be observed for the  $k$ NN algorithm with 50 neighbors, on which the constraint of the 8tracks dataset should not have such strong effects. This can in turn be explained by the aforementioned large size of the head of long-tailed distribution on this dataset (the higher the number of neighbors, the more popular tracks can be taken into account).

Beside the experiments whose results are shown in Figure 2, we also experimented with models based on the Markov property, among them the simple bigram model and the LME model of Chen et al. [2012]. Despite the long time that can be required to train these models—for example, several weeks for the LME model—these methods led to

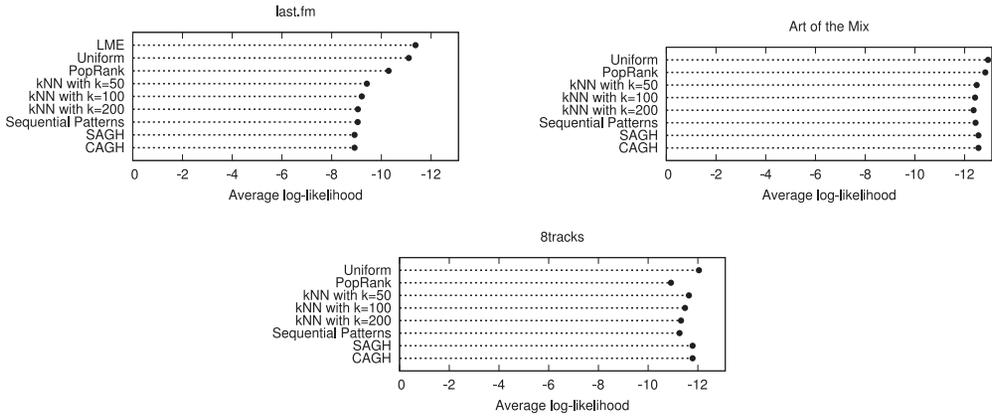


Fig. 3. Average log-likelihood of the evaluated approaches on the different datasets.

particularly low hit rates values that consistently were below 10% on the three datasets. We therefore omit these results in this article. In general, given these differences, assuming the Markov property might be too strong for this problem setting.

Overall, given the results of this and the previous measurements, we can conclude that the choice of the “best” playlist generation strategy can to some extent depend on the dataset characteristics. At the same time, our measurements suggest that considering different information sources, including artist information or general track popularity, can help to improve the results of individual techniques. In our view, future research could thus (1) aim to develop appropriate hybridization designs and (2) explore the use of algorithms that take the particularities of the datasets into account.<sup>13</sup>

### RQ3: Comparing Hit Rate and $ACL$ Results

In a next step, we evaluated the different algorithms using the  $ACL$  measure to compare it with our findings regarding the hit rates. The results are shown in Figure 3.

As described in Section 5, this evaluation approach requires a smoothing step to avoid zero probabilities. To make the measurement, we therefore developed variants of the described algorithms in which the returned scores correspond to probability values and combined them with a uniform distribution in a mixture model. These values were obtained using a linear normalization of the scores of the tracks:

$$p_M(t|h) = \frac{\text{score}_M(t, h)}{\sum_i \text{score}_M(i, h)}.$$

Notice that using a linear normalization is not the only possible strategy. Each training set of the cross-validation process was additionally split into a learning set to obtain the model probabilities and a validation set to compute optimized weights using the expectation-maximization algorithm. We also provide the results of the aforementioned LME model with standard parameters for the last.fm dataset. On both other datasets, the training process lasted more than a month.

The figure shows that the  $ACL$  results are different from the hit rate results and that there is, for example, no consistent “winner” across the datasets. The CAGH method only outperforms the others very slightly on the last.fm dataset. On the AotM dataset, generally no strong differences between the algorithms were observed. On the 8tracks dataset, finally, the PopRank method outperformed the others. Generally, the

<sup>13</sup>In the context of frequent pattern mining, the usage of techniques that are particularly suited for long-tail distributions and that are capable of handling “rare” itemsets [Weiss 2004] could be explored.

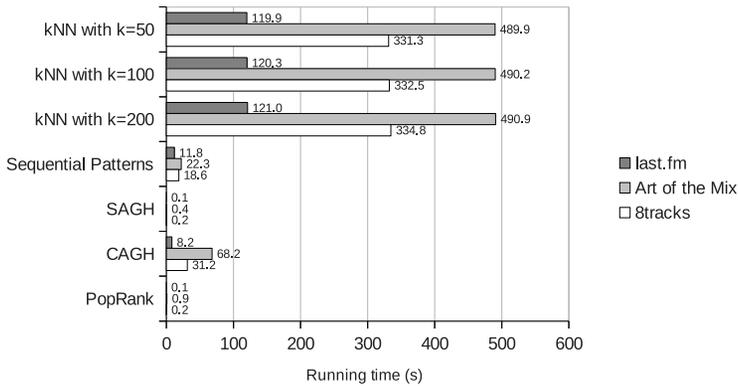


Fig. 4. Time required for 1,000 predictions by each algorithm.

differences between the methods are very small. The explanation for this phenomenon lies in the quite extreme long-tail distributions regarding the popularity of the tracks. When the probabilities of the playlist generation methods are combined with the uniform distribution, too much weight is put on the long-tail tracks. As a result, the differences between the techniques are dampened.

To verify our explanation about the similar results, we also evaluated all mixture models of this section in terms of the hit rate and obtained very similar outputs for all algorithms, which we do not include for space reasons. In particular, the hit rate of the CAGH model is lowered to an extent that approaches CAGH and SAGH both lead to similar low hit rates.

Regarding RQ3, we conclude that the ranking of the algorithms can depend on the used measure and that  $ALC$  measurements do not necessarily follow the trend of the hit rate measure. At the same time, using the  $ALC$  measure exhibits major limitations. First, it can only be used to compare algorithms, and it does not tell us if the generated playlists are actually good. Second, the measure requires model smoothing to avoid zero probabilities. This in turn can result in generally lower values of the estimated quality, which can finally lead to very small differences between different algorithms, at least when a linear normalization is used.

#### RQ4: Computational Complexity

Finally, we draw our attention to questions of computational complexity and the running times (in seconds) required to generate playlists. Figure 4 shows the running times that the different algorithms required to compute 1,000 predictions on the three datasets. For each dataset, we averaged the running time of each algorithm on each test set of the 10-fold cross-validation. To minimize the effect of possible perturbations, each of these running times corresponds the minimum of 10 runs. The experiments were run on a Linux computer with an AMD Opteron Processor 6272 CPU and required around 500MB of RAM in the worst case.

As expected, the running time of the  $k$ NN algorithm is much higher than that of the other algorithms. The results of this algorithm strongly depend on the number and size of playlists that have to be scanned to determine the neighbors for each prediction. Therefore, the running times on the 8tracks dataset are much higher than on the last.fm dataset, and much lower than the results on the AotM dataset. Notice that in the case of playlist generation, neighborhoods cannot be precomputed, as the playlist history to be continued is not known in advance. The measurements also show that the value of  $k$  only had a small influence on the running times. Another expected result is the very short running time of the SAGH and PopRank algorithms. Finally,

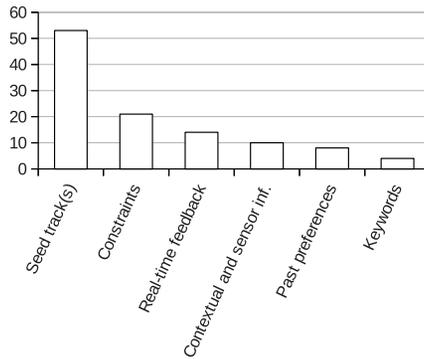


Fig. 5. Methods used in the literature to specify target characteristics.

the results show that when using the CAGH algorithm or the sequential patterns, the recommendations can be generated quite fast.

Since the observed running times depend on a number of factors, including the effectiveness of the implementation, we made a theoretical analysis of the properties of the different techniques compared in this work. The observed running times reported here correspond to the insights of the theoretical analysis. The details of this analysis can be found in the Appendix.

## 7. RESEARCH PERSPECTIVES AND CONCLUSIONS

In this work, we have reviewed the recent literature on automated playlist generation and categorized the different proposals with respect to the types of background data they use, the types of inputs they can process, and the algorithms on which they base the generation process. In the course of the literature review, we collected a number of statistics about which kinds of approaches are popular in the research community. Based on these statistics and our observations from the empirical evaluation reported in Section 6, we derived a number of limitations of current research and a subjective list of potential directions for future research.

### 7.1. Data Perspective

We first draw our attention to the different ways of specifying the target characteristics, as discussed in Section 3.2. Figure 5 shows a statistic on how often each method was used in the literature reviewed for this article.

In the research literature, approaches that rely on a given seed track or on explicit constraints, such as genre or mood, are the most common ones. Quite interestingly, even though most online music services allow their user to enter free-text input (e.g., in the form of keywords), this type of user interaction is barely addressed in the research community. Interpreting user inputs based on individual keywords can be challenging. Typing in, for example, the term *eclipse* on last.fm’s music search field results in dozens, if not hundreds, of artists with the keyword appearing in the band name, hundreds of track and album titles, and even two dozen different labels [Lamere 2012]. Another type of input that is not in the focus of the research community is the use of the long-term user preferences model, which can be based on past user interactions, listening logs, or explicit preference statements in the user profile.

Future research could therefore aim to explore the challenges and opportunities when different types of user inputs should be supported. The interpretation of individual keywords could be based on information from the current user’s profile. At the same time, such profile information could also be the basis for personalizing the generated playlists. Another interesting perspective is that of contextual playlisting, as only a few

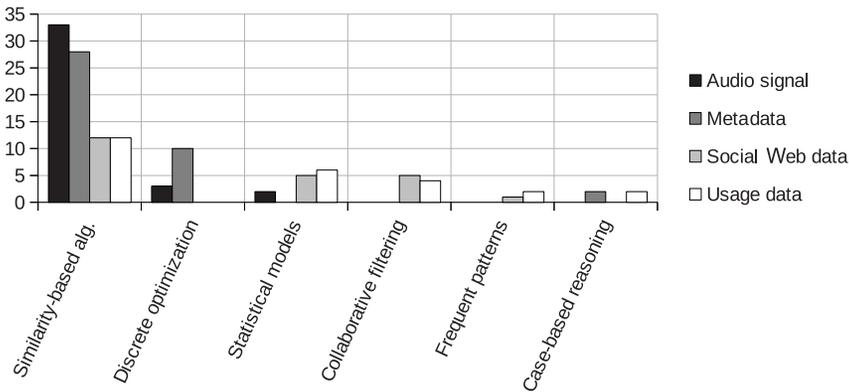


Fig. 6. Types of background knowledge per type of algorithm in the literature.

works focus on it (e.g., Reddy and Mascia [2006] and Chiarandini et al. [2011]). Some inspiration could be taken from the more general domain of music recommendation, where this is an emerging topic [Kaminskas and Ricci 2011].

Figure 6 gives an overview of which types of background knowledge are used in the literature and at the same time gives statistics about the corresponding algorithmic approaches.<sup>14</sup> The figure shows that most existing research is based on features extracted from the audio signal or some additionally available information, such as that about genres or artists. This observation generally is not particularly surprising, as many of the reviewed papers appeared in publication outlets of the MIR community. However, the figure clearly highlights that potentially valuable sources of information that can be obtained from Social Web sources or from listening logs are not yet fully explored and thus represent a promising area for future research. Although the number of papers that we reviewed is not very large, it seems that more and more papers deal with such type of background knowledge. However, the majority of recent research that we reviewed still uses the audio signal and metadata.

## 7.2. Algorithms Perspective

Figure 6 also reveals that similarity-based algorithms are the most common approach to generate playlists, which is an obvious approach when the goal is to maximize the homogeneity of the playlists. Correspondingly, however, the diversity and serendipity of the tracks can be low. Interestingly, collaborative filtering techniques, which dominate the RS research field, are not very frequently used in the literature even though platforms such as last.fm at least partially base their radios on such techniques.<sup>15</sup>

Looking at the dates of publication, we could see that most of the research based on discrete optimization dates back to around the year 2000. Since then, only a few papers have dealt with this type of algorithm, and approaches based on similarity analysis represent the majority even in the most recent research.

In future works, the goal could thus be to better exploit the additionally available information about the tastes of users and the community to build not only more diverse playlists but, in particular, generate playlists that are *personalized* according to user preferences. In that context, hybrid systems that combine different types of information sources such as Hariri et al. [2012] appear to be particularly promising, as the inner

<sup>14</sup>The numbers were obtained by classifying only the papers that were reviewed in the context of this work. Therefore, they should be interpreted as being rough indicators.

<sup>15</sup><http://www.last.fm/help/faq>.

consistency of a playlist can be relevant, which perhaps can only be achieved if metadata about the tracks is considered.

Since smooth transitions of tracks are often mentioned as an important quality criterion in the literature, the value of better incorporating such aspects into playlists generation approaches could be explored in more depth in future research. Finally, no substantial amount of research has been done on automatic playlisting for social groups. However, this is an interesting topic, as music is often shared and used for social interaction (parties, motivation for team sports, romantic evening, etc.).

### 7.3. Evaluation Perspective

We finally looked at what represents the “common practice” in research to evaluate automatically generated playlists. Given the subjectivity of how playlists are perceived by users, user studies surprisingly are not an important means used by researchers to assess the quality of playlists. Overall, a little more than 30% of the papers that we reviewed used some form of user studies to validate their hypotheses; about the same amount of papers relied on historical data, offline experimental designs, and objective quality measures. About 10% of the papers (additionally) used individual examples to discuss the quality of the generated playlists. Another 10%, however, did not contain any proper evaluation at all. The remaining papers evaluated other aspects of the generation system as a whole, such as the computational complexity.

When the evaluation was based on historical data (e.g., on handcrafted playlists) or other experimental designs not involving users, IR measures like the hit rate are common. The *ALL* measure only plays a minor role. As another “objective” measure, the homogeneity of the generated playlists is very often in the focus of researchers.

Overall, although user studies are probably the best means for evaluation, these studies often have their limitations, particularly with respect to their size and the selection of participants. We therefore argue that the available (log) data—including listening logs—should be better exploited in future research, as there are a number of aspects that could be assessed in offline experimental designs. In particular, future research could try to assess multiple criteria at the same time and explore existing trade-offs such as homogeneity versus diversity or familiarity versus novelty [Jannach et al. 2013]. The result of such research could be a guideline as to which playlist generation strategy is best suited for a given playlist goal. If, for example, the goal is the discovery of new songs, collaborative filtering methods could be favorable. If the playlist tracks should be on a certain theme, other techniques could be better.

Furthermore, additional objective measures beyond homogeneity could be designed to assess different quality criteria, including novelty, serendipity, or the smoothness of the track transitions. In any case, however, experimental studies should be made to analyze the correspondence with the user-perceived quality.

Finally, although a number of papers exist that explicitly address questions of computational complexity and scalability, in our view, this topic needs more attention, especially considering that our experiments have shown that some advanced playlist generation methods do not scale well and require weeks of training already for medium-sized problem settings.

### 7.4. Conclusions

In this article, we have reviewed different approaches to automatically generating music playlists. After having characterized the playlist generation problem and its challenges, we have looked at both what kinds of data are used by such systems and which technical approach they adopt to create the playlists. Our literature review furthermore included a discussion of common approaches on how to evaluate the quality of the generated playlists.

In the second part of the article, we compared several of the existing techniques from the literature on different datasets using different measures of quality and complexity. Among other aspects, our evaluation showed the importance of taking into account the order of the tracks, artist information, and track popularity in the generation process. The results indicate that some of the measures applied in the literature have their limitations and that the outcomes can strongly depend on the characteristics of the underlying datasets. Last, the results of the experiments on the computational complexity confirmed the importance of this dimension when designing algorithms for playlist generation.

Based on our literature review and the experiments, we have sketched a number of possible directions for future research. These directions include the incorporation of additional knowledge sources and the use of additional measures that can help us to better characterize the quality of playlists.

## APPENDIX

### A.1. Computational Complexity Analysis

Algorithm	Analysis
$k$ NN	<p>The major drawback of this algorithm is its complexity. Because each tested playlist history corresponds to a new playlist, no similarity can be calculated in advance in an offline process. This means that all the playlists in the pool of available playlists <math>P</math> have to be scanned for each track selection. The time complexity for the computation of each similarity value between <math>h</math> and each playlist <math>p</math> can be reduced to <math>\mathcal{O}( h  +  p )</math> using a direct matching algorithm. Assuming that the playlists have a maximum size of <math>M</math>, the complexity for the computation of the neighborhood of <math>h</math> is <math>\mathcal{O}( P  \cdot ( h  + M))</math>. If all the tracks in the set of tracks <math>T</math> of the database are considered, the overall complexity is <math>\mathcal{O}( P  \cdot ( h  + M) + k *  T )</math>. In practice, the running time can be reduced by considering only the tracks contained in the resulting neighborhood of <math>h</math>.</p> <p>The space complexity depends on the space required for storing the pool of playlists. Speaking in terms of number of tracks, if the maximum size of playlists is <math>M</math> and the number of playlists <math> P </math>, the space complexity is <math>\mathcal{O}(M \cdot  P )</math>.</p>
Frequent patterns	<p>By extracting a limited set of relevant patterns in an offline process, this approach theoretically allows a lower time and space complexity than that of the <math>k</math>NN approach. Both association rules and sequential patterns consider sub-histories of variable length inside a sliding window of size <math>w</math>. The number of sub-histories considered is thus:</p> $\sum_{n=1}^{w-1} \binom{w-1}{n} = 2^{w-1} + 1$ <p>Determining the rules that match each sub-history in the model can be done in both cases in <math>\mathcal{O}(n)</math> with <math>n</math> the size of the patterns extracted [Nakagawa and Mobasher 2003]. The overall time complexity is thus:</p> $\sum_{n=1}^{w-1} \mathcal{O}(n) \cdot \binom{w-1}{n} \leq \mathcal{O}(w \cdot 2^w)$ <p>The space complexity of sequential patterns corresponds to the number of permutations with repetition of <math>n</math> elements from <math> T </math>. The space complexity is thus <math>\mathcal{O}( T ^n)</math>. The space complexity of association rules corresponds to the number of combinations with repetition of <math>n</math> elements from <math> T </math>, which is less than <math> T ^n</math>. In practice, long patterns are not frequent enough to be of any value. In our experiments, the maximum sizes that led to the best results were 2 and 3. The time complexity is then reduced to <math>\mathcal{O}(w)</math> and <math>\mathcal{O}(w^2)</math> and the space complexities to <math>\mathcal{O}( T ^2)</math> and <math>\mathcal{O}( T ^3)</math>, respectively. Finally, the selection of the most relevant patterns in the offline mining process in practice means that the number of elements to scan is much lower than the theoretical complexity bounds.</p>

Continued

Algorithm	Analysis
SAGH, CAGH	<p>The time complexity of the SAGH algorithm is particularly low. Given a playlist history <math>h</math>, the algorithm only has to retrieve the corresponding artists. The complexity of this task is <math>\mathcal{O}( h )</math>. Retrieving the greatest hits of a given artist can be done in <math>\mathcal{O}(1)</math>. The overall time complexity for the computation of the score of all tracks is thus <math>\mathcal{O}( h )</math>, which is very low, especially when considering that playlists typically have average sizes between 7 and 20 in our datasets.</p> <p>The time complexity of the CAGH is a little higher because the algorithm also considers similarities between artists. As these similarities can be computed offline and stored, the algorithm only has to retrieve the similar artists for each artist in <math>A_h</math> from the database, and then to multiply the greatest hit counts with the calculated similarity values. The time complexity is thus <math>\mathcal{O}( h  \cdot  A  \cdot  T )</math>, where <math>A</math> is the set of artists.</p> <p>The space complexity of the SAGH algorithm corresponds to the set of greatest hits stored, i.e., <math>\mathcal{O}( A  \cdot \gamma)</math> where <math>A</math> is the set of artists and <math>\gamma</math> the maximum number of greatest hits stored per artist.</p> <p>The space complexity of the CAGH algorithm is higher as the similarities between the artists have to be stored. The space complexity is thus: <math>\mathcal{O}( A  \cdot \gamma +  A ^2)</math>.</p>

## A.2. Examples of Publicly Available Music Datasets

### *Datasets Provided by Companies*

It is likely that the most prominent source of music data provided by a commercial service is the database of last.fm, which contains information about 45 million tracks [Grabham 2011] and 40 million users. This includes track attributes, tags, number of listeners per track, play counts, artist information, album information, comments, top fans, user playlists, and so forth.

A similar source for music data is MusicBrainz (musicbrainz.org), which is thought of as an encyclopedia where users can collaboratively collect information on music. Submitted information is subject to a review and voting process by other users, leading to a highly reliable database. It contains track attributes, tags, artist information, album information, ratings, and the like but is smaller (15 million tracks) and less rich than the database of last.fm.

Finally, another publicly available music database is provided by The Echo Nest (the.echonest.com), which contains information about 30 million tracks [Zax 2011]. One particularity of this database is that it contains various types of information that is inferred from the audio signal, such as the danceability, the energy, and the loudness.

### *Datasets Provided by Research Groups*

Some smaller databases were released by researchers to allow others to experiment and compare results. One of the most prominent ones is the Million Song Dataset [Bertin-Mahieux et al. 2011]. The dataset consists of one million tracks and corresponding metadata. Each track is described by 55 attributes, some of which taken from The Echo Nest and MusicBrainz. Other smaller similar datasets also exist, such as the dataset of the CUIDADO project [Pachet 2001; Vinet et al. 2002].

A dataset of handcrafted playlists was made available by McFee and Lanckriet [2012]. The dataset contains all playlists of the Web site Art of the Mix (<http://www.artofthemix.org>) that were created before June 2011. It contains about 100,000 playlists created by about 16,000 users. Each playlist contains artist names, track names, creator name, a categorical label (e.g., “Reggae”), and the identifiers of the tracks that could be found in the Million Song Dataset.

## REFERENCES

- Rakesh Agrawal, Tomasz Imieliński, and Arun Swami. 1993. Mining association rules between sets of items in large databases. In *Proceedings of SIGMOD*. 207–216.
- Rakesh Agrawal and Ramakrishnan Srikant. 1995. Mining sequential patterns. In *Proceedings of ICDE*. 3–14.
- Natalie Aizenberg, Yehuda Koren, and Oren Somekh. 2012. Build your own music recommender by modeling internet radio streams. In *Proceedings of WWW*. 1–10.
- Masoud Alghoniemy and Ahmed H. Tewfik. 2000. User-defined music sequence retrieval. In *Proceedings of Multimedia*. 356–358.
- Masoud Alghoniemy and Ahmed H Tewfik. 2001. A network flow model for playlist generation. In *Proceedings of ICME*. 445–448.
- Klaus-Dieter Althoff. 2001. Case-based reasoning. *Handbook on Software Engineering and Knowledge Engineering*, Vol. 1. World Scientific, Singapore, 549–588.
- Andreja Andric and Goffredo Haus. 2005. Estimating quality of playlists by sight. In *Proceedings of AXMEDIS*. 68–74.
- Jean-Julien Aucouturier and François Pachet. 2002a. Finding songs that sound the same. In *Proceedings of MPCA*. 1–8.
- Jean-Julien Aucouturier and François Pachet. 2002b. Music similarity measures: What’s the use? In *Proceedings of ISMIR*. 157–163.
- Jean-Julien Aucouturier and François Pachet. 2002c. Scaling up music playlist generation. In *Proceedings of ICME*, Vol. 1. 105–108.
- Claudio Baccigalupo and Enric Plaza. 2006. Case-based sequential ordering of songs for playlist recommendation. In *Proceedings of ECCBR*. 286–300.
- Claudio Baccigalupo and Enric Plaza. 2007. A case-based song scheduler for group customised radio. In *Case-Based Reasoning Research and Development*. Lecture Notes in Computer Science, Vol. 4626. Springer, 433–448.
- Claudio Giovanni Baccigalupo. 2009. *Poolcasting: An Intelligent Technique to Customise Musical Programmes for Their Audience*. Ph.D. Dissertation. Universitat Autònoma de Barcelona, Barcelona, Spain.
- Wietse Balkema and Ferdi van der Heijden. 2010. Music playlist generation by assimilating GMMs into SOMs. *Pattern Recognition Letters* 31, 11, 1396–1402.
- Luke Barrington, Reid Oda, and Gert R. G. Lanckriet. 2009. Smarter than genius? Human evaluation of music recommender systems. In *Proceedings of ISMIR*. 357–362.
- David Baskerville and Tim Baskerville. 2010. *Music Business Handbook and Career Guide*. Sage.
- Dominikus Baur, Sebastian Boring, and Andreas Butz. 2010. Rush: Repeated recommendations on mobile devices. In *Proceedings of IUI*. 91–100.
- Thierry Bertin-Mahieux, Daniel P. W. Ellis, Brian Whitman, and Paul Lamere. 2011. The million song dataset. In *Proceedings of ISMIR*. 591–596.
- Jacob T. Biehl, Piotr D. Adamczyk, and Brian P. Bailey. 2006. DJogger: A mobile dynamic music device. In *Proceedings of CHI*. 556–561.
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent Dirichlet allocation. *Journal of Machine Learning Research* 3, 993–1022.
- Thomas L. Blum, Douglas F. Keislar, James A. Wheaton, and Erling H. Wold. 1999. Method and Article of Manufacture for Content-Based Analysis, Storage, Retrieval, and Segmentation of Audio Information. U.S. Patent 5,918,223.
- Dmitry Bogdanov and Perfecto Herrera. 2011. How much metadata do we need in music recommendation? A subjective evaluation using preference sets. In *Proceedings of ISMIR*. 97–102.
- Geoffray Bonnin and Dietmar Jannach. 2013a. A comparison of playlist generation strategies for music recommendation and a new baseline scheme. In *Proceedings of ITWP*. 16–23.
- Geoffray Bonnin and Dietmar Jannach. 2013b. Evaluating the quality of playlists based on hand-crafted samples. In *Proceedings of ISMIR*. 263–268.
- Klaas Bosteels, Elias Pampalk, and Etienne E. Kerre. 2009. Evaluating and analysing dynamic playlist generation heuristics using radio logs and fuzzy set theory. In *Proceedings of ISMIR*. 351–356.
- Karan Kumar Budhreja, Ashutosh Singh, Gaurav Dubey, and Arun Khosla. 2012. Probability based playlist generation based on music similarity and user customization. In *Proceedings of NCCCS*. 1–5.
- Robin Burke. 2002. Hybrid recommender systems: Survey and experiments. *User Modeling and User-Adapted Interaction* 12, 4, 331–370.

- Richard H. Byrd, Peihuang Lu, Jorge Nocedal, and Ciyou Zhu. 1995. A limited memory algorithm for bound constrained optimization. *SIAM Journal on Scientific Computing* 16, 5, 1190–1208.
- Rui Cai, Chao Zhang, Lei Zhang, and Wei-Ying Ma. 2007. Scalable music recommendation by search. In *Proceedings of Multimedia*. 1065–1074.
- Òscar Celma. 2010. *Music Recommendation and Discovery: The Long Tail, Long Fail, and Long Play in the Digital Music Space*. Springer.
- Òscar Celma and Xavier Serra. 2008. FOAFing the music: Bridging the semantic gap in music recommendation. *Web Semantics: Science, Services and Agents on the World Wide Web* 6, 4, 250–256.
- Zeina Chedrawy and Syed Sibte Raza Abidi. 2009. A Web recommender system for recommending, predicting and personalizing music playlists. In *Proceedings of WISE*. 335–342.
- Ching-Wei Chen, Kyogu Lee, and Ho-Hsiang Wu. 2009. Towards a class-based representation of perceptual tempo for music retrieval. In *Proceedings of ICMLA*. 602–607.
- Shuo Chen, Josh L. Moore, Douglas Turnbull, and Thorsten Joachims. 2012. Playlist prediction via metric embedding. In *Proceedings of KDD*. 714–722.
- Luca Chiarandini, Massimiliano Zanoni, and Augusto Sarti. 2011. A system for dynamic playlist generation driven by multimodal control signals and descriptors. In *Proceedings of MMSP*. 1–6.
- Dave Cliff. 2006. hpDJ: An automated DJ with floorshow feedback. In *Consuming Music Together*. Computer Supported Cooperative Work, Vol. 35. Springer, 241–264.
- Filipe Coelho, José Devezas, and Cristina Ribeiro. 2013. Large-scale crossmedia retrieval for playlist generation and song discovery. In *Proceedings of OAIR*. 61–64.
- Michel Crampes, Jean Villerd, Andrew Emery, and Sylvie Ranwez. 2007. Automatic playlist composition in a dynamic music landscape. In *Proceedings of SADPI*. 15–20.
- Paolo Cremonesi, Franca Garzotto, Sara Negro, Alessandro Vittorio Papadopoulos, and Roberto Turrin. 2011. Looking for “good” recommendations: A comparative evaluation of recommender systems. In *Proceedings of INTERACT*. 152–168.
- Fabio Crestani, M. Lalmas, and Cornelius J. van Rijsbergen. 1998. *Information Retrieval: Uncertainty and Logics*. Advanced Models for the Representation and Retrieval of Information: The Information Retrieval Series, Vol. 4. Springer.
- Sally Jo Cunningham, David Bainbridge, and Annette Falconer. 2006. ‘More of an art than a science’: Supporting the creation of playlists and mixes. In *Proceedings of ISMIR*. 240–245.
- Ricardo Dias and Manuel J. Fonseca. 2010. MuVis: An application for interactive exploration of large music collections. In *Proceedings of MM*. 1043–1046.
- Markus Dopler, Markus Schedl, Tim Pohle, and Peter Knees. 2008. Accessing music collections via representative cluster prototypes in a hierarchical organization scheme. In *Proceedings of ISMIR*. 179–184.
- J. Stephen Downie. 2003. Music information retrieval. *Annual Review of Information Science and Technology* 37, 1, 295–340.
- Robert B. Ekelund, George S. Ford, and Thomas Koutsky. 2000. Market power in radio markets: An empirical analysis of local and national concentration. *Journal of Law and Economics* 43, 1, 157–184.
- Benjamin Fields. 2011. *Contextualize Your Listening: The Playlist as Recommendation Engine*. Ph.D. Dissertation. Department of Computing Goldsmiths, University of London, London, UK.
- Benjamin Fields, Christophe Rhodes, Michael Casey, and Kurt Jacobson. 2008. Social playlists and bottleneck measurements: Exploiting musician social graphs using content-based dissimilarity and pairwise maximum flow values. In *Proceedings of ISMIR*. 559–564.
- Arthur Flexer, Dominik Schnitzer, Martin Gasser, and Gerhard Widmer. 2008. Playlist generation using start and end songs. In *Proceedings of ISMIR*. 173–178.
- Don Gartner, Florian Kraft, and Thomas Schaaf. 2007. An adaptive distance measure for similarity based playlist generation. In *Proceedings of ICASSP*. 229–232.
- Arthur Germain and Jacob Chakareski. 2013. Spotify me: Facebook-assisted automatic playlist generation. In *Proceedings of MMSP*. 25–28.
- Dan Grabham. 2011. Last.fm: ‘Free On-Demand Not Successful for Us.’ *TechRadar*. Retrieved August 19, 2014, from <http://www.techradar.com/us/news/internet/audio/last-fm-free-on-demand-not-successful-for-us-946276>.
- Maarten Grachten, Markus Schedl, Tim Pohle, and Gerhard Widmer. 2009. The ISMIR cloud: A decade of ISMIR conferences at your fingertips. In *Proceedings of ISMIR*. 63–68.
- Jiawei Han and Micheline Kamber. 2006. *Data Mining: Concepts and Techniques* (2nd ed.). Morgan Kaufmann.

- Derek L. Hansen and Jennifer Golbeck. 2009. Mixing it up: Recommending collections of items. In *Proceedings of CHI*. 1217–1226.
- Negar Hariri, Bamshad Mobasher, and Robin Burke. 2012. Context-aware music recommendation based on latent topic sequential patterns. In *Proceedings of RecSys*. 131–138.
- Negar Hariri, Bamshad Mobasher, and Robin Burke. 2013. Personalized text-based music retrieval. In *Proceedings of ITWP*. 24–30.
- Pitoyo Hartono and Ryo Yoshitake. 2013. Automatic playlist generation from self-organizing music map. *Journal of Signal Processing* 17, 1, 11–19.
- David Hauger, Johannes Kepler, Markus Schedl, Andrej Košir, and Marko Tkalcić. 2013. The million musical tweets dataset: What can we learn from microblogs. In *Proceedings of ISMIR*. 189–194.
- David B. Hauver and James C. French. 2001. Flycasting: Using collaborative filtering to generate a playlist for online radio. In *Proceedings of WEDELMUSIC*. 123–130.
- Otmar Hilliges, Phillipp Holzer, Rene Klüber, and Andreas Butz. 2006. AudioRadar: A metaphorical visualization for the navigation of large music collections. In *Proceedings of Smart Graphics*. 82–92.
- Thomas Hornung, Cai-Nicolas Ziegler, Simon Franz, Martin Przyjaciel-Zablocki, Alexander Schatzle, and Georg Lausen. 2013. Evaluating hybrid music recommender systems. In *Proceedings of WI-IAT*. 57–64.
- Jia-Lien Hsu and Shuk-Chun Chung. 2011. Constraint-based playlist generation by applying genetic algorithm. In *Proceedings of SMC*. 1417–1422.
- Yajie Hu and Mitsunori Ogihara. 2011. NextOne player: A music recommendation system based on user behavior. In *Proceedings of ISMIR*. 103–108.
- Dietmar Jannach, Lukas Lerche, Fatih Gedikli, and Geoffroy Bonnin. 2013. What recommenders recommend—an analysis of accuracy, popularity, and sales diversity effects. In *User Modeling, Adaptation, and Personalization*. 25–37.
- Dietmar Jannach, Markus Zanker, Mouzhi Ge, and Marian Gröning. 2012. Recommender systems in computer science and information systems—a landscape of research. In *Proceedings of EC-Web*. 76–87.
- Claus Aage Jensen, Ester M. Mungure, Torben Bach Pedersen, and Kenneth Sørensen. 2007. A data and query model for dynamic playlist generation. In *Proceedings of ICDE*. 65–74.
- Daniel Jurafsky and James H. Martin. 2009. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition* (2nd ed.). Prentice Hall.
- Antti Jylhä, Stefania Serafin, and Cumhur Erku. 2012. Rhythmic walking interactions with auditory feedback: An exploratory study. In *Proceedings of AM*. 68–75.
- Mohsen Kamalzadeh, Dominikus Baur, and Torsten Möller. 2012. A survey on music listening and management behaviours. In *Proceedings of ISMIR*. 373–378.
- Marius Kaminskas and Francesco Ricci. 2011. Location-adapted music recommendation using tags. In *Proceedings of UMAP*. 183–194.
- Peter Knees, Tim Pohle, Markus Schedl, and Gerhard Widmer. 2006. Combining audio-based similarity with Web-based data to accelerate automatic music playlist generation. In *Proceedings of MIR*. 147–154.
- Paul B. Lamere. 2012. I've got 10 million songs in my pocket: Now what? In *Proceedings of RecSys*. 207–208.
- Aristomenis S. Lampropoulos, Dionysios N. Sotiropoulos, and George A. Tsihrintzis. 2012. Evaluation of a cascade hybrid recommendation as a combination of one-class classification and collaborative filtering. In *Proceedings of ICTAI*. 674–681.
- Jin Ha Lee, Bobby Bare, and Gary Meek. 2011. How similar is too similar? Exploring users' perceptions of similarity in playlist evaluation. In *Proceedings of ISMIR*. 109–114.
- Kyogu Lee and Minsu Cho. 2011. Mood classification from musical audio using user group-dependent models. In *Proceedings of ICMLA*. 130–135.
- Arto Lehtiniemi and Jarno Seppänen. 2007. Evaluation of automatic mobile playlist generator. In *Proceedings of MC*. 452–459.
- Ning-Han Liu, Szu-Wei Lai, Chien-Yi Chen, and Shu-Ju Hsieh. 2009. Adaptive music recommendation based on user behavior in time slot. *Computer Science and Network Security* 9, 2, 219–227.
- Beth Logan. 2002. Content-based playlist generation: Exploratory experiments. In *Proceedings of ISMIR*. 295–296.
- Beth Logan. 2004. Music recommendation from song sets. In *Proceedings of ISMIR*. 425–428.
- Beth Logan and Ariel Salomon. 2001. A music similarity function based on signal analysis. In *Proceedings of ICME*. 190–193.
- François Maillat, Douglas Eck, Guillaume Desjardins, and Paul Lamere. 2009. Steerable playlist generation by learning song similarity from radio station playlists. In *Proceedings of ISMIR*. 345–350.

- Niitsuma Masahiro, Hiroshi Takaesu, Hazuki Demachi, Masaki Oono, and Hiroaki Saito. 2008. Development of an automatic music selection system based on runner's step frequency. In *Proceedings of ISMIR*. 193–198.
- Josh McDermott. 2012. Auditory preferences and aesthetics: Music, voices, and everyday sounds. In *Neuroscience of Performance and Choice*. Elsevier, 227–256.
- Brian McFee, Thierry Bertin-Mahieux, Daniel, and Gert Lanckriet. 2012. The million song dataset challenge. In *Proceedings of WWW'12 Companion*. 909–916.
- Brian McFee and Gert R. G. Lanckriet. 2011. The natural language of playlists. In *Proceedings of ISMIR*. 537–542.
- Brian McFee and Gert R. G. Lanckriet. 2012. Hypergraph models of playlist dialects. In *Proceedings of ISMIR*. 343–348.
- Cory McKay. 2010. *Automatic Music Classification with jMIR*. Ph.D. Dissertation. McGill University, Montreal, Quebec, Canada.
- Owen Craigie Meyers. 2007. *A Mood-Based Music Classification and Exploration System*. Master's Thesis. Massachusetts Institute of Technology, Cambridge, MA.
- Bart Moens, Leon van Noorden, and Marc Leman. 2010. D-Jogger: Syncing music with walking. In *Proceedings of SMC*. 451–456.
- Joshua L. Moore, Shuo Chen, Thorsten Joachims, and Douglas Turnbull. 2012. Learning to embed songs and tags for playlist prediction. In *Proceedings of ISMIR*. 349–354.
- Miki Nakagawa and Bamshad Mobasher. 2003. Impact of site characteristics on recommendation models based on association rules and sequential patterns. In *Proceedings of ITWP*.
- Nuria Oliver and Fernando Flores-Mangas. 2006. MPTrain: A mobile, music and physiology-based personal trainer. In *Proceedings of MobileHCI*. 21–28.
- Nuria Oliver and Lucas Kreger-Stickles. 2006. PAPA: Physiology and purpose-aware automatic playlist generation. In *Proceedings of ISMIR*. 393–396.
- François Pachet. 2001. Metadata for music and sounds: The CUIDADO project. In *Proceedings of CBMI*. 411–415.
- François Pachet, Pierre Roy, and Daniel Cazaly. 2000. A combinatorial approach to content-based music selection. *Multimedia* 7, 1, 44–51.
- Elias Pampalk, Tim Pohle, and Gerhard Widmer. 2005. Dynamic playlist generation based on skipping behavior. In *Proceedings of ISMIR*. 634–637.
- Steffen Pauws and Berry Eggen. 2002. PATS: Realization and user evaluation of an automatic playlist generator. In *Proceedings of ISMIR*. 222–231.
- Steffen Pauws and Sander van de Wijdeven. 2005. User evaluation of a new interactive playlist generation concept. In *Proceedings of ISMIR*. 638–643.
- Steffen Pauws, Wim Verhaegh, and Mark Vossen. 2006. Fast generation of optimal music playlists using local search. In *Proceedings of ISMIR*. 138–143.
- Steffen Pauws, Wim Verhaegh, and Mark Vossen. 2008. Music playlist generation by adapted simulated annealing. *Information Sciences* 178, 3, 647–662.
- Geoffroy Peeters and Joachim Flocon-Cholet. 2012. Perceptual tempo estimation using GMM-regression. In *Proceedings of MIRUM*. 45–50.
- John C. Platt, Christopher J. C. Burges, Steven Swenson, Christopher Weare, and Alice Zheng. 2001. Learning a Gaussian process prior for automatically generating music playlists. In *Proceedings of NIPS*. 1425–1432.
- Tim Pohle, Peter Knees, Markus Schedl, Elias Pampalk, and Gerhard Widmer. 2007. Reinventing the wheel: A novel approach to music player interfaces. *Multimedia* 9, 3, 567–575.
- Tim Pohle, Elias Pampalk, and Gerhard Widmer. 2005. Generating similarity-based playlists using traveling salesman algorithms. In *Proceedings of DAFx*. 220–225.
- Robert Ragno, Chris J. C. Burges, and Cormac Herley. 2005. Inferring similarity between music objects with application to playlist generation. In *Proceedings of MIR*. 73–80.
- Sasank Reddy and Jeff Mascia. 2006. Lifetrak: Music in tune with your life. In *Proceedings of HCM*. 25–34.
- Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian personalized ranking from implicit feedback. In *Proceedings of UAI*. 452–461.
- Gordon Reynolds, Dan Barry, Ted Burke, and Eugene Coyle. 2008. Interacting with large music collections: Towards the use of environmental metadata. In *Proceedings of ICME*. 989–992.
- Ronald Rosenfeld. 2000. Two decades of statistical language modeling: Where do we go from here? *Proceedings of the IEEE* 88, 8, 1270–1278.

- Vegard Sandvold, Thomas Aussenac, Oscar Celma, and Perfecto Herrera. 2006. Good vibrations: Music discovery through personal musical concepts. In *Proceedings of ISMIR*. 322–323.
- Andy M. Sarroff and Michael Casey. 2012. Modeling and predicting song adjacencies in commercial albums. In *Proceedings of SMC*.
- Ingo Schmädecke and Holger Blume. 2013. High performance hardware architectures for automated music classification. In *Algorithms from and for Nature and Life*. 539–547.
- Malcolm Slaney and William White. 2006. Measuring playlist diversity for recommendation systems. In *Proceedings of AMCMM*. 77–82.
- Malcolm Slaney and William White. 2007. Similarity based on rating data. In *Proceedings of ISMIR*. 479–484.
- George Tzanetakis. 2002. *Manipulation, Analysis and Retrieval Systems for Audio Signals*. Ph.D. Dissertation. Princeton University, Princeton, NJ.
- Rob Van Gulik and Fabio Vignoli. 2005. Visual playlist generation on the artist map. In *Proceedings of ISMIR*. 520–523.
- Igor Vatolkin, Mike Preuß, and Günter Rudolph. 2011. Multi-objective feature selection in music genre and style recognition tasks. In *Proceedings of GECCO*. 411–418.
- Fabio Vignoli and Steffen Pauws. 2005. A music retrieval system based on user driven similarity and its evaluation. In *Proceedings of ISMIR*. 272–279.
- Hugues Vinet, Perfecto Herrera, and François Pachet. 2002. The CUIDADO project. In *Proceedings of ISMIR*. 197–203.
- Richard Y. Wang and Diane M. Strong. 1996. Beyond accuracy: What data quality means to data consumers. *Journal of Management Information Systems* 12, 4, 5–33.
- Gary M. Weiss. 2004. Mining with rarity: A unifying framework. *SIGKDD Explorations Newsletter* 6, 1, 7–19.
- Erling Wold, Thom Blum, Douglas Keislar, and James Wheaton. 1996. Content-based classification, search, and retrieval of audio. *Multimedia* 3, 3, 27–36.
- Linxing Xiao, Lie Lu, Frank Seide, and Jie Zhou. 2009. Learning a music similarity measure on automatic annotations with application to playlist generation. In *Proceedings of ICASSP*. 1885–1888.
- David Zax. 2011. The Echo Nest makes Pandora look like a transistor radio. *Fast Company*. Retrieved August 19, 2014, from <http://www.fastcompany.com/1734773/echo-nest-makes-pandora-look-transistor-radio>.
- Elena Zheleva, John Guiver, Eduarda Mendes Rodrigues, and Nataša Milić-Frayling. 2010. Statistical models of music-listening sessions in social media. In *Proceedings of WWW*. 1019–1028.

Received January 2014; revised July 2014; accepted July 2014