

PUTTING THE USER IN THE CENTER OF MUSIC INFORMATION RETRIEVAL

Markus Schedl

Department of Computational Perception,
Johannes Kepler University, Linz, Austria

Arthur Flexer

Austrian Research Institute for
Artificial Intelligence, Vienna, Austria

ABSTRACT

Personalized and context-aware music retrieval and recommendation algorithms ideally provide music that perfectly fits the individual listener in each imaginable situation and for each of her information or entertainment need. Although first steps towards such systems have recently been presented at ISMIR and similar venues, this vision is still far away from being a reality. In this paper, we investigate and discuss literature on the topic of user-centric music retrieval and reflect on why the breakthrough in this field has not been achieved yet. Given the different expertises of the authors, we shed light on why this topic is a particularly challenging one, taking a psychological and a computer science view. Whereas the psychological point of view is mainly concerned with proper experimental design, the computer science aspect centers on modeling and machine learning problems. We further present our ideas on aspects vital to consider when elaborating user-aware music retrieval systems, and we also describe promising evaluation methodologies, since accurately evaluating personalized systems is a notably challenging task.

1. WHY CARE ABOUT THE USER?

In our discussion of the importance and the challenges of development and evaluation in Music Information Retrieval (MIR) we distinguish between systems-based and user-centric MIR. We define systems-based MIR as all research concerned with experiments existing solely in a computer, e.g. evaluation of algorithms on digital databases. In contrast, user-centered MIR always involves human subjects and their interaction with MIR systems.

Systems-based MIR has traditionally focused on computational models to describe universal aspects of human music perception, for instance, via elaborating musical feature extractors or similarity measures. Doing so, the existence of an objective “ground truth” is assumed, against which corresponding music retrieval algorithms (e.g., playlist generators or music recommendation systems) are evaluated. To give a common example, music retrieval ap-

proaches have been evaluated via genre classification experiments for years. Although it was shown already in 2003 that musical genre is an ill-defined concept [1], genre information still serves as a proxy to assess music similarity and retrieval approaches in systems-based MIR.

On the way towards user-centered MIR, the coarse and ambiguous concept of genre should either be treated in a personalized way or replaced by the concept of similarity. When humans are asked to judge the similarity between two pieces of music, however, certain other challenges need to be faced. Common evaluation strategies typically do not take into account the musical expertise and taste of the users. A clear definition of “similarity” is often missing too. It might hence easily occur that two users apply a very different, individual notion of similarity when assessing the output of music retrieval systems. While a first person may experience two songs as rather dissimilar due to very different lyrics, a second one may feel a much higher resemblance of the very same songs because of a similar instrumentation. Similarly, a fan of Heavy Metal music might perceive a Viking Metal track as dissimilar to a Death Metal piece, while for the majority of people the two will sound alike.

The above examples illustrate that there are many aspects that influence what a human perceives as similar in a musical context. These aspects can be grouped into three different categories according to [29]: *music content*, *music context*, and *user context*. Examples for each category are given in Figure 1. It is exactly this multifaceted and individual way of music perception that has largely been neglected so far when elaborating and evaluating music retrieval approaches, but should be given more attention, in particular considering the trend towards personalized and context-aware systems.

A *personalized system* is one that incorporates information about the user into its data processing part (e.g., a particular user taste for a movie genre). A *context-aware system*, in contrast, takes into account dynamic aspects of the user context when processing the data (e.g., location and time where/when a user issues a query). Although the border between personalization and context-awareness may appear fuzzy from this definition, in summary, personalization usually refers to the incorporation of more static, general user preferences, whereas context-awareness refers to the fact that frequently changing aspects of the user’s environmental, psychological, and physiological context are considered.

Permission to make digital or hard copies of all or part 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 bear this notice and the full citation on the first page.

© 2012 International Society for Music Information Retrieval.

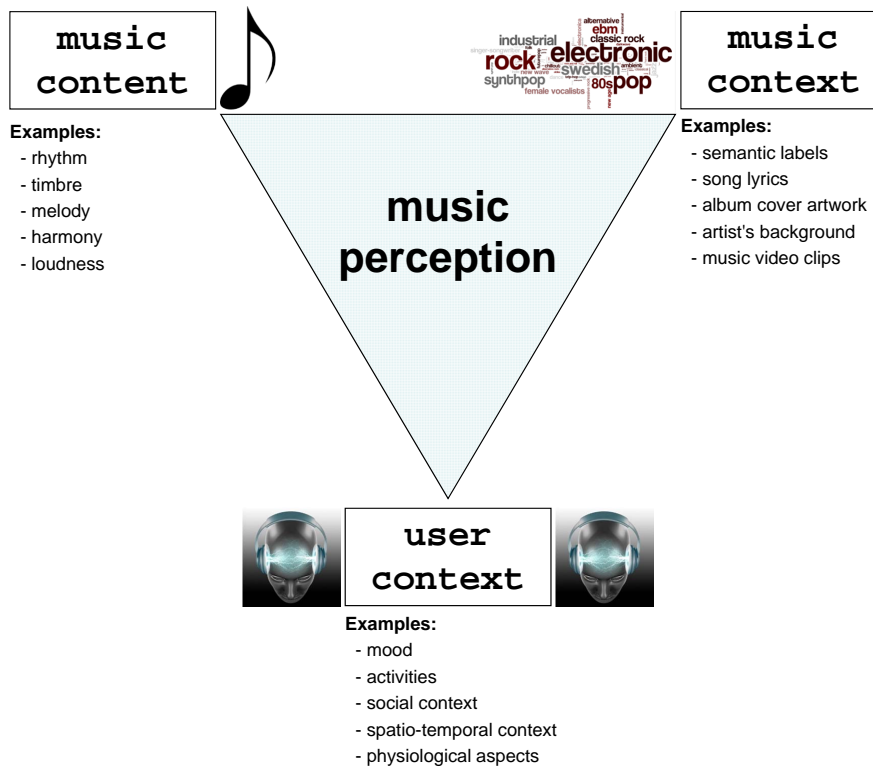


Figure 1. Factors that influence human music perception.

The remainder of this paper is organized as follows. Section 2 reviews approaches that, in one way or the other, take the user into account when building music retrieval systems. Evaluation strategies for investigating user-centric MIR are discussed in Section 3. In Section 4, we eventually summarize important factors when creating and evaluating user-aware music retrieval systems.

2. HOW TO MODEL THE USER?

Existing user-aware systems typically model the user in a very simplistic way. For instance, it is common in *collaborative filtering* approaches [22, 28] to build user profiles only from information about a user u expressing an interest in item i . As an indicator of interest may serve, for example, a click on a particular item, a purchasing transaction, or in MIR the act of listening to a certain music piece. Such indications, in their simplest form, are stored in a binary matrix where element $r(u, i)$ denotes the presence or absence of a connection between user u and item i . In common recommendation systems, a more fine-grained scale for modeling the user interest in an item is typically employed – users frequently rate items according to a Likert-type scale, e.g., by assigning one to five stars to it. Matrix factorization techniques are subsequently applied to recommend novel items [19].

Taking a closer look at literature about context-aware

retrieval and recommendation in the music domain, we can see that approaches differ considerably in terms of how the user context is defined, gathered, and incorporated. The majority of approaches rely solely on one or a few aspects (temporal features in [7], listening history and weather conditions in [21], for instance), whereas comprehensive user models are rare in MIR. One of the few exceptions is Cunningham et al.'s study [8] that investigates if and how various factors relate to music taste (e.g., human movement, emotional status, and external factors such as temperature and lightning conditions). Based on the findings, the authors present a fuzzy logic model to create playlists.

There further exists some work that assumes a mobile music consumption scenario. The corresponding systems frequently aim at matching music with the current pace of a walker or jogger, e.g., [3, 24]. Such systems typically try to match the user's heartbeat with the music played [23]. However, almost all proposed systems require additional hardware for context logging, e.g., [8, 9, 11].

In [15] a system that matches tags describing a particular place with tags describing music is presented. Employing text-based similarity measures between the multimodal sets of tags, Kaminskas and Ricci propose their system for location-based music recommendation. Baltrunas et al. [2] suggest a context-aware music recommender system for music consumption while driving. Although the authors take into account eight different contextual factors

(e.g., driving style, mood, road type, weather, traffic conditions), their application scenario is quite restricted and their system relies on explicit human feedback, which is burdensome.

Zhang et al. present *CompositeMap* [34], a model that takes into account similarity aspects derived from music content as well as social factors. The authors propose a multimodal music similarity measure and show its applicability to the task of music retrieval. They also allow a simple kind of personalization of this model by letting the user weight the individual music dimensions on which similarity is estimated. However, they do neither take the user context into consideration, nor do they try to learn a user's preferences.

In [26] Pohle et al. present preliminary steps towards a simple personalized music retrieval system. Based on a clustering of community-based tags extracted from *last.fm*, a small number of musical concepts are derived using *Non-Negative Matrix Factorization* (NMF) [20,32]. Each music artist or band is then described by a "concept vector". A user interface allows for adjusting the weights of the individual concepts, based on which artists that match the resulting distribution of the concepts best are recommended to the user. Zhang et al. propose in [34] a very similar kind of personalization strategy via user-adjusted weights.

Knees and Widmer present in [17] an approach that incorporates *relevance feedback* [27] into a text-based music search engine [16] to adapt the retrieval process to user preferences. The search engine proposed by Knees et al. builds a model from music content features (MFCCs) and music context features (term vector representations of artist-related Web pages). To this end, a weight is computed for each (term, music item)-pair, based on the term vectors. These weights are then smoothed, taking into account the closest neighbors according to the content-based similarity measure (Kullback-Leibler divergence on Gaussian Mixture Models of the MFCCs). To retrieve music via natural language queries, each textual query issued to the system is expanded via a *Google* search, resulting again in a term weight vector. This query vector is subsequently compared to the smoothed weight vectors describing the music pieces, and those with smallest distance to the query vector are returned.

Nürnberg and Detyniecki present in [25] a variant of the *Self-Organizing Map* (SOM) [18] that is based on a model that adapts to *user feedback*. To this end, the user can move data items on the SOM. This information is fed back into the SOM's codebook, and the mapping is adapted accordingly.

In [33] Xue et al. present a *collaborative personalized search model* that alleviates the problems of *data sparseness* and *cold-start for new users* by combining information on different levels (individuals, interest groups, and global). Although not explicitly targeted at music retrieval, the idea of integrating data about the user, his peer group, and global data to build a social retrieval model might be worth considering for MIR purposes.

The problem with the vast majority of approaches pre-

sented so far is that evaluation is still carried out without sufficient user involvement. For instance, [7, 25, 26] seemingly do not perform any kind of evaluation involving real users, or at least do not report it. Some approaches are evaluated on user-generated data, but do not request feedback from real users during the evaluation experiments. For example, [16] makes use of collaborative tags stored in a database to evaluate the proposed music search engine. Similarly, [21] relies on data sets of listening histories and weather conditions, and [33] uses a corpus of Web search data. Even if real users are questioned during evaluation, their individual properties (such as taste, expertise, or familiarity with the music items under investigation) are regularly neglected in evaluation experiments. In these cases, evaluation is typically performed to answer a very narrow question in a restricted setting. To give an example, the work on automatically selecting music while doing sports, e.g. [3, 23, 24], is evaluated on the very question of whether pace or heartbeat of the user does synchronize with the tempo of the music. Likewise Kaminskis and Ricci's work on matching music with places of interest [15], even though it is evaluated by involving real users, comprises only the single question whether the music suggested by their algorithm is suited for particular places of interest. Different dimensions of the relation between images and music are not addressed. Although this is perfectly fine for the intended use cases, such highly specific evaluation settings are not able to provide answers to more general questions of music retrieval and recommendation, foremost because these settings fail at offering explanations for the (un)suitability of the musical items under investigation.

An evaluation approach that tries to alleviate this shortcoming is presented in [4], where subjective listening tests to assess music recommendation algorithms are conducted using a multifaceted questionnaire. Besides investigating the enjoyment a user feels when listening to the recommended track ("liking"), the authors also ask for the user's "listening intention", whether or not the user knows artist and song ("familiarity"), and whether he or she would like to request more similar music ("give-me-more"). A similar evaluation scheme is suggested in [12]. However, Firan et al. only investigate liking and novelty.

In summary, almost all approaches reported are still more systems-based than user-centric.

3. HOW TO EVALUATE USER-CENTERED MIR?

In what follows we will argue that whereas evaluation of systems-based MIR has quite matured, evaluation of user-centered MIR is still in its infancy. Let us start by reviewing what the nature of experiments is in the context of MIR. The basic structure of MIR experiments is the same as in any other experimental situation: the question is whether there are effects of the variation of the independent variables (also called factors) on the dependent variables. In the case of systems-based MIR, independent variables are e.g. type and certain parameter characteristics of the algorithms used or type and characteristics of the data

set in question. Typical dependent variables are various performance measures like accuracy, precision, root mean squared error or training time. A standard computer experiment is genre classification where the independent variable is the type of classification algorithm, say algorithm A and B, and the dependent variable is the achieved accuracy. Statistical testing is used to ensure that the observed effects on the dependent variables are caused by the varied independent variables and not by mere chance, i.e. to ascertain that the observed differences are too large to attribute them to random influences only. Besides using the proper statistical instruments to establish statistical significance of results it is equally important to make sure to control all important factors in the experimental design. Any factor that is able to influence the dependent variables has to be part of the experimental design. E.g. if algorithm A, compared to algorithm B, works better for electronic dance music than for rock music then any experimental design not containing dance music will obscure differences between A and B. The important thing to note is that for systems-based MIR which uses only computer experiments it is comparably easy to control all important factors which could have an influence on the dependent variables. This is because the number of factors is both manageable and controllable since the experiments are being conducted on computers and not in the real world.

Already early on in the history of MIR research, gaps concerning the evaluation of MIR systems have been identified. Futrelle and Downie [14], in their review of the first three years of the ISMIR conference published in 2003, identify two major problems: (i) no commonly accepted means of comparing retrieval techniques, (ii) few if any attempts to study potential users of MIR systems. The first problem concerns evaluation of computer experiments and the second problem the barely existing inclusion of users in MIR studies. Flexer [13], in his review of the 2004 ISMIR conference [5], argues for the necessity of statistical evaluation of MIR experiments. He presents minimum requirements concerning statistical evaluation by applying fundamental notions of statistical hypotheses testing to MIR research. His discussion is concerned with systems-based MIR, the example used throughout the paper is that of automatic genre classification based on audio content analysis. The MIR community is criticized for the lack of statistical evaluation it uses, e.g. only two papers in the ISMIR 2004 proceedings [5] employed a statistical test to prove significance of their results. These ongoing discussions about evaluation of MIR experiments have led to a first evaluation benchmark taking place at the ISMIR conference 2004 [6] and further on to the establishment of the annual evaluation campaign for MIR algorithms (Music Information Retrieval Evaluation eXchange, MIREX) [10]. In 2011, MIREX consisted of 16 tasks ranging from audio classification, cover song identification, audio key detection to structural segmentation and audio tempo estimation. All but two tasks are concerned with systems-based MIR and a purely computer-based evaluation of algorithms. The two exceptions using human evaluations in

a more real-world setting are *Audio Music Similarity and Retrieval* and *Symbolic Melodic Similarity*. Starting with the MIREX 2006 evaluation [10] statistical tests are being used to analyze results.

The situation concerning evaluation of user-centric MIR research is far less well developed. In a recent comprehensive review [31] of user studies in the MIR literature by Weigl and Guastavino, papers from the first decade of ISMIR conferences and related MIR publications were analyzed. A central result is that MIR research has a mostly systems-centric focus. Only twenty papers fell under the broad category of “user studies” which is an alarmingly small number given that 719 articles have been published in the ISMIR conference series alone. To make things worse, these user studies are “predominantly qualitative in nature” and of “largely exploratory nature” [31]. The explored topics range from e.g. user requirements and information needs, insights into social and demographic factors to user-generated meta-information and ground truth. This all points to the conclusion that evaluation of user-centered MIR is at its beginning and that especially a more rigorous quantitative treatment is still missing.

In discussing the challenges of quantitative evaluation of user-centered MIR we like to turn to an illustrative example: the recent 2011 *Audio Music Similarity and Retrieval* task¹ within the annual MIREX [10] evaluation campaign. Each of 18 competing algorithms was given 7000 songs (30 second audio clips) for which they computed similarity rankings. The data consisted of 10 equally sized genre classes ranging from classic music to rock to hip-hop. From the 7000 songs, “100 songs were randomly selected from the 10 genre groups (10 per genre) as queries and the first 5 most highly ranked songs out of the 7000 were extracted for each query (after filtering out the query itself, returned results from the same artist were also omitted). Then, for each query, the returned results (candidates) from all participating algorithms were grouped and were evaluated by human graders”¹. For each individual query/candidate pair, a single human grader provided both a FINE score (from 0 (failure) to 100 (perfection)) and a BROAD score (not similar NS, somewhat similar SS, very similar VS) indicating how similar the songs are in their opinion. The independent variable here is the type of algorithm used to compute the similarity rankings. The dependent variables are the subjects’ broad and fine appraisal of the perceived similarity. But since this is a real-world experiment involving human subjects there is a whole range of factors that have not been assessed. E.g. there are social and demographic factors that might clearly influence the user’s judgment of music similarity: their age, gender, cultural background and especially their musical history, experience and knowledge. But also factors concerning their momentary situation during the actual listening experiment might have an influence: time of day, mood, physical condition. Not to forget more straightforward variables like type of speakers or headphones used for the test. As al-

¹The 2011 results and details can be found at: http://www.music-ir.org/mirex/wiki/2011:Audio_Music_Similarity_and_Retrieval_Results

ready mentioned in section 1, even the choice of dependent variable is debatable. After all, what does “similar” really mean in the context of music? Timbre, mood, harmony, melody, tempo, etc might all be valid answers for different people. This points to a certain lack of rigor concerning the instruction of subjects during the experiment. This enumeration of potential problems is not intended to badmouth this MIREX task which still is a valuable contribution and an applaudable exception to the rule of computer-only evaluation. But it is meant as a warning and to highlight the explosion of independent variables and factors that might add to the variance of observed results and might obscure significant differences. In principle, all such factors have to be recorded and made independent variables in the overall experimental design.

If MIR is to succeed in maturing from purely systems-based to user-centered research we will have to leave the nice and clean world of our computers and face the often bewilderingly complex real world of real human users and all the challenges this entails for proper design and evaluation of experiments. To make this happen it will be necessary that our community with a predominantly engineering background opens up to the so-called “soft sciences” of e.g. psychology and sociology which have developed instruments and methods to deal with the complexity of human subjects.

4. DISCUSSION AND CONCLUSIONS

Incorporating real users in both the development and assessment of music retrieval systems is of course an expensive and arduous task. However, recent trends in music distribution, in particular the emergence of music streaming services that make available millions of tracks to their users, call for intelligent personalized and context-aware systems to deal with this abundance. Concerning the development of such systems, we believe that the following two reasons have prevented major breakthroughs so far: (i) a general lack of research on user-centered systems, (ii) a lack of awareness concerning the complexity of evaluation of user-centered systems. In designing such systems, the user should already be taken into account at an early stage during the development process. We need to better understand what the user’s individual requirements are and address these requirements in our implementations. Otherwise it is unlikely that even the spiffiest personalized systems will succeed (without frustrating the user). We hence identify the following four key requirements for elaborating user-centric music retrieval systems:

Personalization aspects have to be taken into account. In this context, it is important to note the highly subjective, cognitive component in the understanding of music and judgement of its personal appeal. Therefore, designing user-aware music applications requires intelligent machine learning techniques, in particular, preference learning approaches that relate the user context to concise, situation-dependent music preferences.

User models that encompass different social scopes are needed. They may aggregate an individual model, an in-

terest group model, a cultural model, and a global model. Furthermore, the user should be modeled as comprehensively as possible, in a fine-grained and multifaceted manner. With today’s sensor-packed smartphones and other intelligent devices it has become easy to perform extensive context logging. Of course, privacy issues must also be taken seriously.

Multifaceted similarity measures that combine different feature categories (music content, music context, and user context) are required. The corresponding representation models should then not only allow to derive similarity between music via content-related aspects, such as beat strength or instruments playing, or via music context-related properties, such as the geographic origin of the performer or a song’s lyrics, but also to describe users and user groups in order to compute a listener-based similarity score.

Evaluation of user-centric music retrieval approaches has to include all independent variables that are able to influence the dependent variables into the experimental design. In particular, such factors may well relate to individual properties of the human assessors. Furthermore, it is advisable to make use of recent approaches that minimize the amount of labor required by the human assessors, while at the same time maintaining the significance of the experiments. This can be achieved, for instance, by employing the concept of “Minimal Test Collections” in the evaluation of music retrieval systems [30].

By paying attention to these advices, we are sure that the exciting field of user-centric music information retrieval will continue to grow and eventually provide us with algorithms and systems that offer personalized and context-aware access to music in an unintrusive way.

5. ACKNOWLEDGMENTS

This research is supported by the Austrian Science Fund (FWF): P22856-N23 and P24095 as well as by the European Commission, FP7 (Seventh Framework Programme), ICT-2011.1.5 Networked Media and Search Systems, grant agreement no. 287711.

6. REFERENCES

- [1] J.-J. Aucouturier and F. Pachet. Representing Musical Genre: A State of the Art. *Journal of New Music Research*, 32(1):83–93, 2003.
- [2] L. Baltrunas, M. Kaminskas, B. Ludwig, O. Moling, F. Ricci, K.-H. Lüke, and R. Schwaiger. InCarMusic: Context-Aware Music Recommendations in a Car. In *Proc. EC-Web*, 2011.
- [3] J. T. Biehl, P. D. Adamczyk, and B. P. Bailey. DJogger: A Mobile Dynamic Music Device. In *CHI 2006: Extended Abstracts*, 2006.
- [4] D. Bogdanov and P. Herrera. How Much Metadata Do We Need in Music Recommendation? A Subjective Evaluation Using Preference Sets. In *Proc. ISMIR*, 2011.

- [5] C.L. Buyoli and R. Loureiro. *Fifth International Conference on Music Information Retrieval*. Universitat Pompeu Fabra, 2004.
- [6] P. Cano, E. Gómez, F. Gouyon, P. Herrera, M. Koppenberger, B. Ong, X. Serra, S. Streich, and N. Wack. ISMIR 2004 Audio Description Contest. 2006.
- [7] T. Cebrián, M. Planagumà, P. Villegas, and X. Amatriain. Music Recommendations with Temporal Context Awareness. In *Proc. RecSys*, 2010.
- [8] S. Cunningham, S. Caulder, and V. Grout. Saturday Night or Fever? Context-Aware Music Playlists. In *Proc. Audio Mostly*, 2008.
- [9] S. Dornbush, J. English, T. Oates, Z. Segall, and A. Joshi. XPod: A Human Activity Aware Learning Mobile Music Player. In *Proc. Workshop on Ambient Intelligence, IJCAI*, 2007.
- [10] J. S. Downie. The Music Information Retrieval Evaluation eXchange (MIREX). *D-Lib Magazine*, Dec 2006.
- [11] G. T. Elliott and B. Tomlinson. PersonalSoundtrack: Context-aware Playlists that Adapt to User Pace. In *CHI 2006: Extended Abstracts*, 2006.
- [12] Claudiu S. Firan, Wolfgang Nejdl, and Raluca Paiu. The Benefit of Using Tag-Based Profiles. In *Proceedings of the 5th Latin American Web Congress (LA-WEB)*, pages 32–41, Santiago de Chile, Chile, October–November 2007.
- [13] A. Flexer. Statistical Evaluation of Music Information Retrieval Experiments. *Journal of New Music Research*, 35(2):113–120, June 2006.
- [14] J. Futrelle and J. S. Downie. Interdisciplinary Research Issues in Music Information Retrieval: ISMIR 2000–2002. *Journal of New Music Research*, 32(2):121–131, 2003.
- [15] M. Kaminskas and F. Ricci. Location-Adapted Music Recommendation Using Tags. In *Proc. UMAP*, 2011.
- [16] P. Knees, T. Pohle, M. Schedl, and G. Widmer. A Music Search Engine Built upon Audio-based and Web-based Similarity Measures. In *Proc. SIGIR*, 2007.
- [17] P. Knees and G. Widmer. Searching for Music Using Natural Language Queries and Relevance Feedback. In *Proc. AMR*, 2007.
- [18] T. Kohonen. *Self-Organizing Maps*, volume 30 of *Springer Series in Information Sciences*. Springer, Berlin, Germany, 3rd edition, 2001.
- [19] Y. Koren, R. Bell, and C. Volinsky. Matrix Factorization Techniques for Recommender Systems. *Computer*, 42, Aug 2009.
- [20] D. D. Lee and H. S. Seung. Learning the Parts of Objects by Non-negative Matrix Factorization. *Nature*, 401(6755):788–791, 1999.
- [21] J. S. Lee and J. C. Lee. Context Awareness by Case-Based Reasoning in a Music Recommendation System. In *Proc. UCS*, 2007.
- [22] G. Linden, B. Smith, and J. York. Amazon.com Recommendations: Item-to-Item Collaborative Filtering. *IEEE Internet Computing*, 4(1), 2003.
- [23] H. Liu, J. Hu, and M. Rauterberg. Music Playlist Recommendation Based on User Heartbeat and Music Preference. In *Proc. ICCTD*, 2009.
- [24] B. Moens, L. van Noorden, and M. Leman. D-Jogger: Syncing Music with Walking. In *Proc. SMC*, 2010.
- [25] A. Nürnberger and M. Detyniecki. Weighted Self-Organizing Maps: Incorporating User Feedback. In *Proc. ICANN/ICONIP*, 2003.
- [26] T. Pohle, P. Knees, M. Schedl, and G. Widmer. Building an Interactive Next-Generation Artist Recommender Based on Automatically Derived High-Level Concepts. In *Proc. CBMI*, 2007.
- [27] J. J. Rocchio. Relevance Feedback in Information Retrieval. In Gerard Salton, editor, *The SMART Retrieval System - Experiments in Automatic Document Processing*, pages 313–323. Englewood Cliffs, NJ: Prentice-Hall, 1971.
- [28] B. Sarwar, G. Karypis, J. Konstan, and J. Reidl. Item-based Collaborative Filtering Recommendation Algorithms. In *Proc. WWW*, 2001.
- [29] M. Schedl and P. Knees. Personalization in Multimodal Music Retrieval. In *Proc. AMR*, 2011.
- [30] J. Urbano and M. Schedl. Towards Minimal Test Collections for Evaluation of Audio Music Similarity and Retrieval. In *Proc. AdMIRe*, 2012.
- [31] D. Weigl and C. Guastavino. User Studies in the Music Information Retrieval Literature. In *Proc. ISMIR*, 2011.
- [32] W. Xu, X. Liu, and Y. Gong. Document Clustering Based on Non-negative Matrix Factorization. In *Proc. SIGIR*, 2003.
- [33] G.-R. Xue, J. Han, Y. Yu, and Q. Yang. User Language Model for Collaborative Personalized Search. *ACM Transactions on Information Systems*, 27(2), Feb 2009.
- [34] B. Zhang, J. Shen, Q. Xiang, and Y. Wang. CompositeMap: A Novel Framework for Music Similarity Measure. In *Proc. SIGIR*, 2009.