# Evaluating Music Discovery Tools on Spotify: The Role of User Preference Characteristics

**Muh-Chyun Tang[1], Mang-Yuan Yang[2]**

## Abstract

An experimental study was conducted to assess the effectiveness of the four music discovery tools available on Spotify, a popular music streaming service, namely: radio recommendation, regional charts, genres and moods, as well as following Facebook friends. Both subjective judgment of user experience and objective measures of search effectiveness were used as the performance criteria. Other than comparison of these four tools, we also compared how consistent are these performance measures. The results show that user experience criteria were not necessarily corresponded to search effectiveness. Furthermore, three user preference characteristics: preference diversity, preference insight, and openness to novelty were introduced as mediating variables, with an aim to investigating how these attributes might interact with these four music discovery tools on performance. The results suggest that users' preference characteristics did have an impact on the performance of these music discovery tools.

Keywords: Music Retrieval; Online Music Services; Spotify; Evaluation; User Experience

## 1. Introduction

With the growing popularity of digital music, music streaming services such as Last.FM, Pandora, Spotify, and more recently Apple Music, have played an increasingly prominent role in our access to music. These services provide novel ways of music recommendation and navigational features that have the potential to greatly expand users' opportunities to come across previously unknown music.

Music, like other imaginary or creative works, is particularly suitable for content or socially based recommendation tools for two reasons: firstly, with creative works like music, what is sought after is the emotive experience it evokes, which is more difficult to express than topical knowledge in traditional information retrieval system. Secondly, unlike in the problem-solving information retrieval context where search behaviors are driven by pre-existing information need, in creative works like music, the desire to consume a product often takes place after, rather than before users' first encounter with the information. Indeed, it is more likely, in the realm of creative work consumption, that the desire to acquire or consume a work is aroused by its first

[1,2]Department and Graduate Institute of Library and Information Science, National Taiwan University, Taipei, Taiwan

* Corresponding Author: Muh-Chyun Tang, E-mail: mctang@ntu.edu.tw

being mentioned or sampled by the user. These reasons might explain why passive encounter, rather than actively seeking, becomes the major way for us to discover new music (Cunningham, Bainbridge, & McKay, 2007) and readings for leisure (Ross, 1999). Recommender systems and other navigational tools therefore play an important role in our access to creative works for their ability to expose users to items they are previously unaware of. However, little has been done in regard to how to evaluate the performance of these tools in real time, and under what circumstances might these tools be more effective than the others. One of the main purposes of this study is to explore different performance criteria for the evaluation of these tools. Besides evaluation, we are also interested in finding out whether users' preference characteristics such as preference diversity, openness, insight might influence the performance of different tools. An experimental study was conducted, using the popular music streaming service Spotify as the test site, to address these questions.

## 2. Literature Review

### 2.1 Evaluation of music discovery systems

Previous studies on real life music information behaviors with a view to supporting the design of music information retrieval (MIR) have been done (For example, Cunningham et al., 2007; Cunningham, Reeves, & Britland,

2003; Dougan, 2012. For a thorough review, see Kamalzadeh, Baur, & Möller, 2012; Laplante, 2010; Laplante & Downie, 2006; Lee, 2010; Lee & Downie, 2004; Weigl & Guastavino, 2011), the majorities of which focus on active information seeking carried out by the users for a felt information need. Relative little has been studied with regards to music discovery by passive encountering, especially by means of recommendation or browsing. Traditionally, the development of recommender systems has been driven by accuracy oriented measures such as MAE. Lately, however, it has been recognized that accuracy along cannot fully account for the success of recommender systems (Herlocker, Konstan, Terveen, & Riedl, 2004; Konstan & Riedl, 2012). Herlocker et al. (2004) argued the need to take into consideration of non-obviousness criteria such as novelty and serendipity to assess the effectiveness of recommender systems. Konstan and Riedl (2012) further argued that a broader set of evaluation criteria is needed, especially those reflect user experience, to determine the true value of recommender systems to the user. Indeed, it was found in Tang, Sie, and Ting (2014) that not all evaluation criteria agreed with each other. Their findings suggested that the value of these discovery tools are multi-faceted and should be evaluated as such. Therefore, in the present study, other than the result quality based criteria, different dimensions of user experiences,

such as satisfaction, interesting to use, future use intention, and indispensability were also introduced as evaluation criteria.

### 2.2 Users' music preference characteristics

Other than evaluation methods and criteria, we are also interested in investigating whether the effectiveness of different music discovery tools varies with users' preference characteristics. It has been shown that individuals with different preference characteristics might influence individuals' responses to recommendations of creative works such as movies (Kwon, Cho, & Park, 2009; Shen & Ball, 2011), and leisure readings (Tang et al., 2014). For example, it has been found, in the context of leisure reading seeking, that readers with higher preference insight, that is, more knowledgeable about

their preference, performed better when using author-based browsing (Tang et al., 2014), one wonders whether the impact of users' preference characteristics on tool performance can also be observed in music discovery. By "preference diversity" we mean to represent how narrow or wide one's music interests are. We suspect that individuals with diverse music interests might appreciate more tools that exposed them to music in a wide variety of styles and genres. On the other hand, "openness to novelty" represents an individual's intrinsic need to seek stimulation through novelty, i.e. previously unfamiliar genres or artists. Individuals with high openness might welcome more novel or serendipitous finds. It would be interesting, therefore, to investigate whether the mediating role of users' preference characteristics can also be observed in the realm of
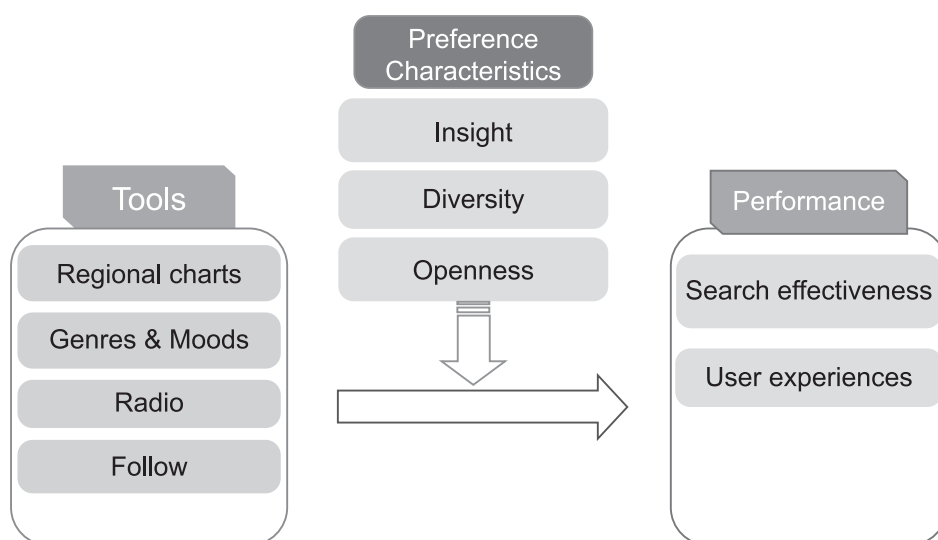


**Figure 1. Research Model**

3

music discovery. Figure 1 represents graphically our research model.

Thus, we summarize our research statement as to test the feasibility of the proposed methodology and performance criteria on the evaluation of music discovery tools, which are substantially different from traditional information retrieval context that merits the research of new evaluation methodologies. Besides evaluation methodology, the study also wish to explore the impact of users' preference characteristics on system performance.

## 3. Methodology

### 3.1 Tools and participants

Music stream service Spotify was chosen as the test site as it enjoys wide popularity and offers a plethora of music finding tools. A convenient sampling was used where the participants were recruited mainly through social media and online music forum. To take part in the study, the participants need to have at least 10 songs saved in her/his Spotify account and a minimal of 10 Facebook friends who are also using Spotify. The requirement is to make sure that all the music discovery tools can function effectively. All participants were offered 200 NTD (equivalent of 6 USD) for their time and efforts. Four music finding tools were chosen as they represent very different music finding approaches: "Charts," "Genres & Moods," "Radio," and "Follow." See

the screenshots of the four music discovery tools below. For Charts function, the user is able to browse most played music in different countries or regions (See Figure 2). Genres & Moods, as its name suggests, allows users to choose music by different genres and moods. The mood category is particularly interesting as listeners often seek music that is conducive to certain emotions (See Figure 3). The Radio tool makes recommendations of songs based on a user's profile, in other words, it makes novel recommendation based on the attributes of songs or artists previously saved or listened to by the user (See Figure 4). The Follow tool has the strongest social dimension as it allows users to follow friends, artists, or other taste-makers activities on the Facebook (See Figure 5).

### 3.2 Research design and procedures

Upon their arrival at the lab, the participants were asked to sign up the consent form, followed by a background questionnaire that elicits data about their music listening behaviors, and most importantly, a set of questions regarding three dimensions of the preference characteristics: preference diversity, openness to novelty, and preference insight. They were then asked to perform music exploratory task (White & Roth, 2009) in which they were to save whatever songs that they found desirable. With each tool, the participants were to find and save 1 to 20 songs that they enjoyed within 8 minutes, but they could
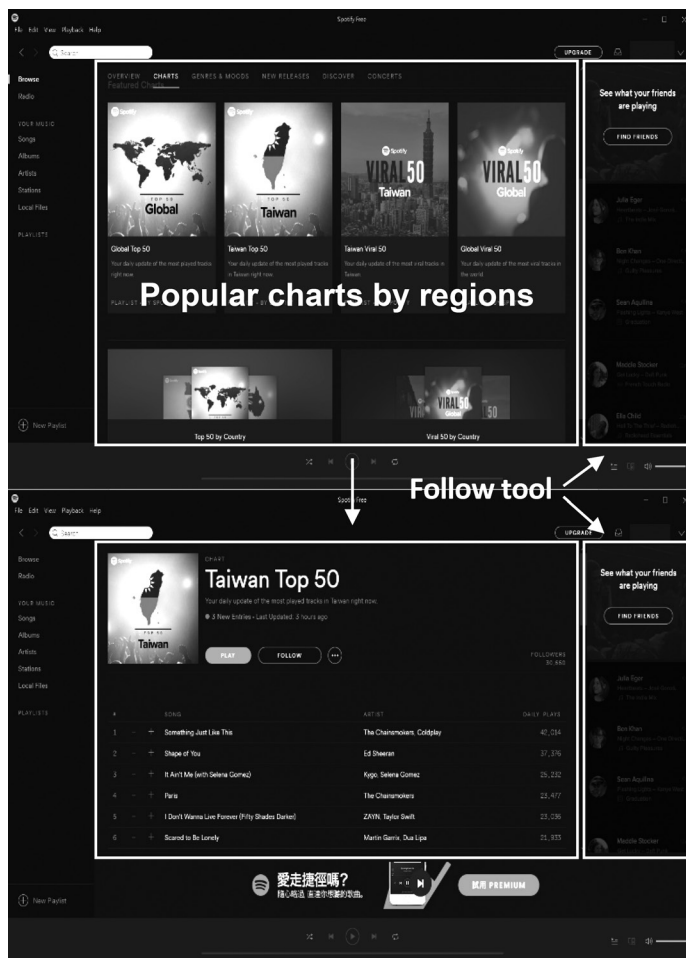
**Figure 2.   Popular Charts by Regions**

stop whenever they felt they were unable to find more interesting songs. Before performing the music finding tasks, the participants were taught by the researcher how to use the tools to be tested. A within-subject design was adopted where they were asked to perform music finding tasks alternately on all the above-mentioned four music discovery tools. To avoid order effect, the order of tools tested were alternated. All the activities performed by the participants were captured by screen recording software for further analysis. With screen capture, we were able to tally the total number of songs that has appeared on the screen when a particular tool was used. All the songs the user was exposed to when using a tool would then constitute the "awareness set," namely,
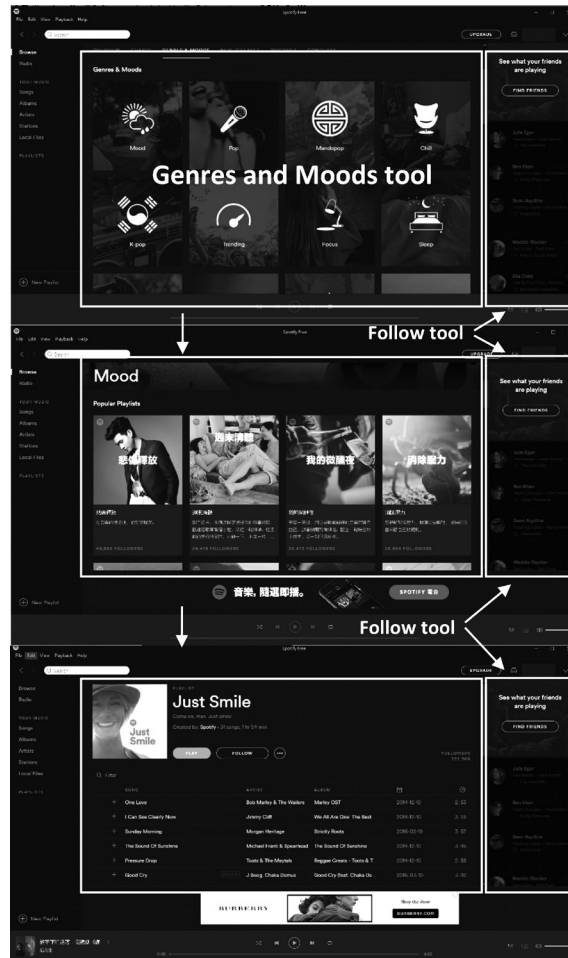
**Figure 3. Genres and Moods Tool**

the number of songs made selectable to the user when using a tool. The awareness set, along with "consideration set," comprising the set of songs sampled by the user, and "choice set," the set of songs the user eventually saved, would constitute the basic elements for the calculation of the accuracy measures. Notice that for a song to be counted as being considered by the user, it has to be listened to more than 10 seconds. The threshold was set to avoid overestimating the consideration set, especially when the "Radio" tool was used, where songs were played continuously without users' active selection. Once the participants found songs they enjoy, they were to save them, which constitutes the choice set. Both the consider/awareness and choice/awareness ratio were
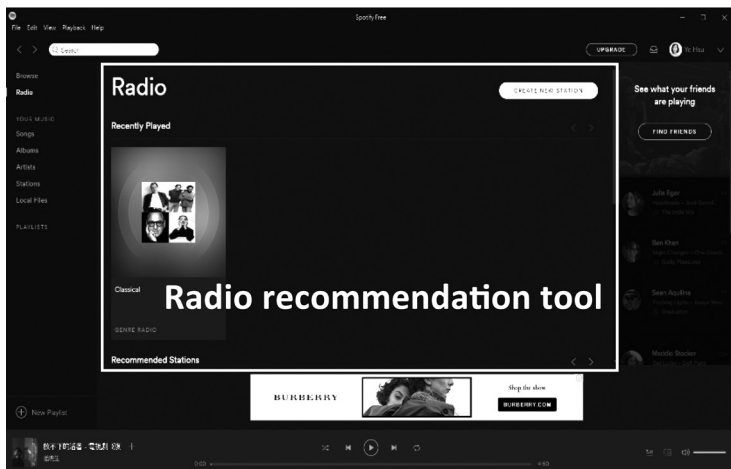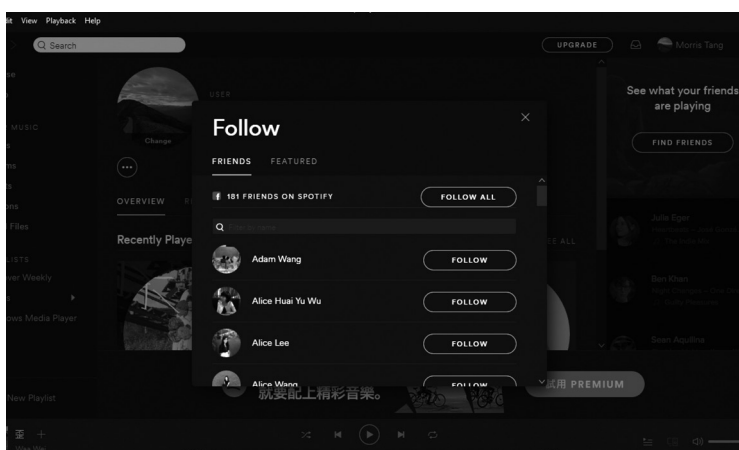
**Figure 4.   Radio Recommendation Tool**



**Figure 5.   Follow Facebook Friend/Fan Pages Tool**

used as search effectiveness measures similar to traditional precision measure in IR evaluation (See Figure 6).

After finishing each tool, the participants were to fill up a post-task questionnaire to elicit their opinions about it, which included items about different dimensions of user experiences (See Table 1 for the data collection instruments). A short interview was also conducted by the researcher to elicit participants' perception of these tools to help us better interpret the results of our quantitative analysis.
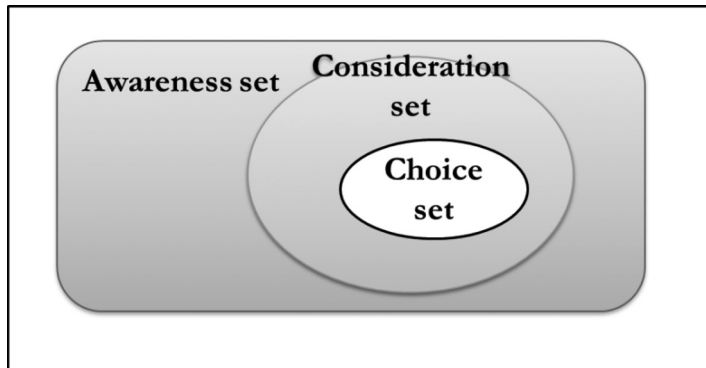
**Figure 6. Choice Set Model**

**Table 1. Data Collection Instruments**

| Data collected | Variables | Data collection instruments |
|---|---|---|
| Preference characteristics | | Pre-search questionnaire |
| User experiences | 1. Preference matching<br>2. Helpfulness<br>3. Interesting to use<br>4. Future use intention<br>5. Indispensability | Post-task questionnaire after each tool |
| Search effectiveness | 1. Consideration set ratio<br>2. Choice set ratio | Screen capture, play history, list of songs saved |

# 4. Results

A total of 26 regular music listeners took part in the study, one third of which were new to Spotify. Table 2 give the basic composition of the participants.

## 4.1 Search effectiveness and user experience

A repeated-measured ANOVA was performed with the tools as the factor and the number of songs in the consideration/awareness set ratio as dependent variable. The result was significant, $F (3, 100) = 33.82, p < .001$. Post-hoc tests showed that, the consideration/awareness ratio ($M = 55.89\%$) of Follow was significantly higher than Chart ($M = 18.82\%$), Genres & Moods (15.27%), and Radio (14.52%). Similar patterns were found in the ANOVA results when using choice/awareness set ratio as the dependent variable, $F (3, 100) = 27.58, p < .001$. The choice set ratio of Follow was found to significantly higher than the other three tools. A closer examination of the awareness set by four tools

**Table 2.  Participant Backgrounds**

|  |  | Count | Percentage |
|---|---|---|---|
| Gender | Male | 10 | 38.5 |
|  | Female | 16 | 61.5 |
| Use experience | New users | 9 | 34.6 |
|  | Less than 6 months | 5 | 19.2 |
|  | 6 to 12 months | 6 | 23.1 |
|  | More than 12 months | 6 | 23.1 |
| Use frequency | Daily | 5 | 19.2 |
|  | 1-3 times/week | 5 | 19.2 |
|  | 1-3 times/month | 6 | 23.1 |
|  | Less than once/month | 10 | 38.5 |

revealed that the apparently high precision of the Follow tool was mainly due to its lack of choices. On average, the Follow tool produced a very small awareness set, only about 50 songs, which was significantly fewer than other tools, with Genres & Moods having the highest awareness set of over 220, followed by Radio, and Chart (See Figure 7).

Next we compared different aspects of user experiences, again, using ANOVA. After using each tool, the participants were asked to rate, on a 0-5 scale, on how accurate it matched one's preference (Match), how helpful it was to find songs (Helpfulness), and how interesting it was to use (Interesting), how willing one was to use it in the future (Future use intention), and how indispensable it was as a discovery tool (Indispensability). Significant differences were found in Match, $F (3, 100) = 9.59, p <$ .001; Helpfulness, $F (3, 100) = 4.74, p = .004;$ Interesting, $F (3, 100) = 4.63, p = .005;$ Future use intention, $F (3, 100) = 3.57, p = .017,$ but not in Indispensability. From Figure 8, it can be observed that, in general, the Radio and Genres & Moods provided better user experiences than the Charts and Follow tools, with the Follow tool doing the worst.

### 4.2 Preference characteristics and tools

A questionnaire was administered to elicit data about three aspects of user preference characteristics: "preference diversity," "openness to novelty," and "preference insight." The scale purification resulted in 9 questions, out of 15 questions asked. The rotated Varimax solution yielded three interpretable preference factors, "preference diversity," "openness to novelty," and
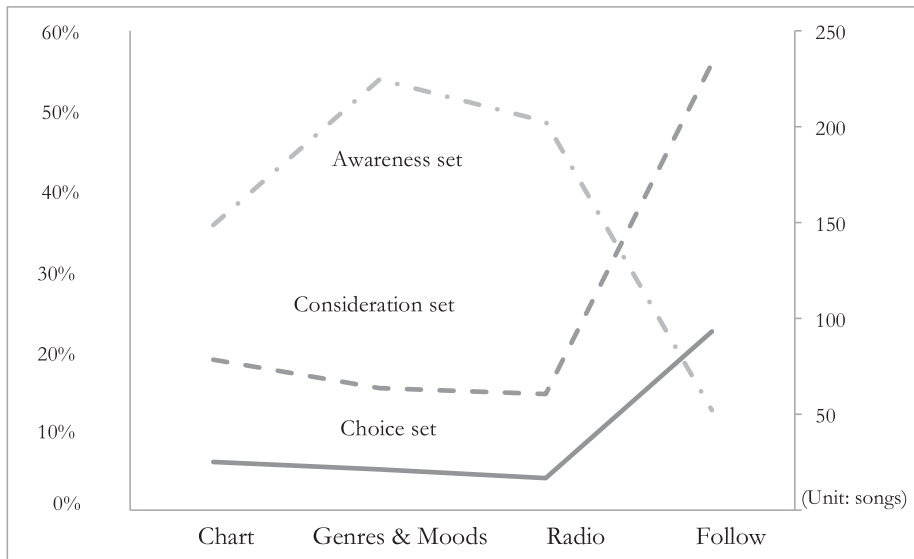
9

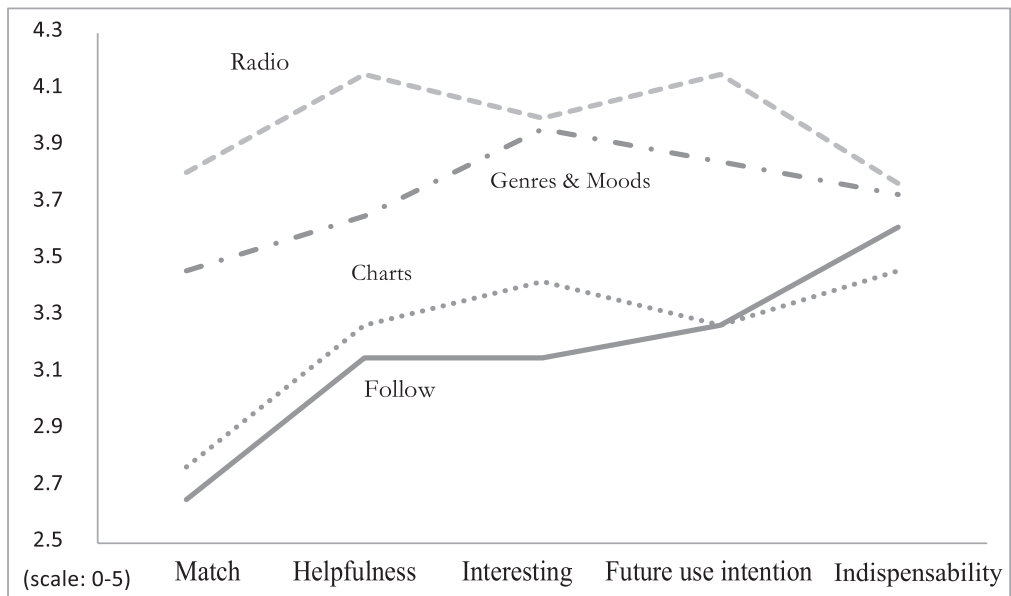**Figure 7. Performance of the Tools on Search Effectiveness Measures**



**Figure 8. Comparison of 4 Discovery Tools on 5 Aspects of User Experiences**

"preference insight," each of which accounted for 23.52%, 22.98%, and 14.26% of variance, respectively. The average scores of the item associated with each component were then used as the score on these dimensions.

Correlation analyses between users' preference characteristics and tool performance were then conducted to examine whether there is a selective compatibility between these preference characteristics and preference criteria. Significant correlations were found only when the participants used the Chart tool. As shown in Table 4, individuals with high preference diversity tended to find the Chart tool better match their preference and had a higher willingness for future use. On the other hand, a negative correlation was found between preference insight and choice set ratio with Chart tool, which suggests that those who had

better understanding of their preference were more reluctant to save songs when using the Chart tool.

## 5. Discussion and Conclusion

As digital music services grow in popularity, individuals' access to music has also undergone tremendous change. Users can now discover music by different navigational or recommendation tools offered by these services. Little has been studied, however, about how effective these music discovery tools are.

An experiment was conducted to test four music discovery tools available on the online music services Spotify using both search effectiveness and user experience measures. Search effectiveness was measured by the percentage of songs sampled or saved divided by

**Table 3.  Factor Analysis of User Reading Preference Characteristics**

| Items | Components | | |
|---|---|---|---|
| | 1 | 2 | 3 |
| 1-6 I appreciate music of a variety of styles and artists | .960 | -.084 | .050 |
| 1-7 I enjoy a wide range of music, not limited certain genres | .951 | -.080 | -.001 |
| 1-15 I actively seek out unfamiliar music genres and artists | .754 | .382 | -.028 |
| 1-5 I have trusted recommended sources which I follow almost exclusively[a] | .237 | -.677 | .074 |
| 1-14 I am receptive to music recommended by others or sources | .063 | .675 | .180 |
| 1-10 Other than Chinese music, I also enjoy music from other countries | .244 | .633 | .131 |
| 1-13 I mostly rely on previously known artists to choose music[b] | .132 | .378 | -.757 |
| 1-2 I have a good understanding regarding what music I enjoy | -.045 | .362 | .743 |
| 1-3 I know where I can find music I like | .165 | .211 | .730 |

[a, b] Reverse question.

**Table 4.  Correlations between Performance Criteria and Preference Characteristics**

|                       | Preference diversity | Openness | Preference insight |
| --------------------- | -------------------- | -------- | ------------------ |
| Match                 | .402*                | .179     | -.096              |
| Helpfulness           | .243                 | .183     | -.011              |
| Interesting           | .250                 | -.163    | -.242              |
| Future use intention  | .450*                | .105     | -.291              |
| Indispensability      | .046                 | .064     | .063               |
| Consideration set     | -.097                | -.022    | -.154              |
| Choice set            | -.098                | -.316    | -.576**            |

$*p < .05. **p < .01.$

title made selectable by the tool. User experiences were measured along the five dimensions: preference matching, helpfulness, interesting to use, future use intention, and indispensability. In terms of user experiences, Radio, which has the strongest personalization character as it makes recommendation based on users' listening profile, and Genres & Moods, which also allows users to browse song based on genres and the moods conveyed by the song, consistently performed better than other tools. As for search effectiveness, the Follow tool, which allows users to look into what their Facebook friends have been listening to, has significantly higher consideration and choice set than the rest. However, the result could be misleading because it was mainly because users got to know very few suggestions by using Follow. The main reason therefore was due to the fact that only a small percentage of the participants' Facebook friends were using Spotify when the experiment was conducted. One should use caution to infer the efficacy of socially-based recommendation for music based on our finding here. Our results suggest the effectiveness of personalized recommendation to provide a satisfactory way of discovering music. Also shown to be effective was the novel Genres & Moods feature. Indeed, it is oftentimes the mood carried by the music that music listeners seek after. The findings demonstrated the feasibility of our proposed methodology in which users were to freely explore with different music discovery tools while measuring tool performance with both behavior and questionnaire based metrics. By capturing users' online browsing and saving activities, we were able to create search effective measures using choice set model without using the relevance set model in traditional IR evaluation, which is ill-equipped for search task where no objectively determined relevant set is available.

After all, it's users' preference rather than relevance of search topic that is at stake when it comes to finding music, or other types of creative works.

Another innovative aspect of our inquiry explored the relationship between the effectiveness of the tools and users' preference characteristics. It was found that the Chart tool, which allows users to browse popular hits in different part of the worlds, performed better for users who have diverse music interests. From the post-search interview, it was found that the participants attributed the appeal of the Chart tool to its allowing them to explore music from foreign countries, especially in English, which should not come as a surprise considering the wide influence of Western music. Yet for individuals who have a better insight into their preference tended to find the popularity-based Chart tool less appealing and much less likely to select songs when using this tool. It is believed that the influence of users' preference characteristics on tool performance has significant implications on recommendation strategy as it suggests that user with different music preference characteristics might be better served by different music discovery tools.

There are obvious limitations to our study. Firstly, a relative small sample size greatly limits the power of the statistical analysis. Furthermore, up to 30 percent of the participants were new to Spotify, which might introduce usability issues that might compound the results, even though a training session for each tool was given before the search task. It turned out that only relative few of the participants' Facebook friends were also using Spotify, which also made the results less generalizable. An incongruity was found between system performance and the consideration/awareness set based measure as the Follow function failed to produce comparable recommendations to the users. Cautions need to take, therefore, when applying the measure, to ensure tools tested are able to generate comparable amount of recommendations.

# References

Cunningham, S. J., Bainbridge, D., & McKay, D. (2007). Finding new music: A diary study of everyday encounter with novel songs. In S. Dixon, D. Bainbridge, & R. Typke (Eds.), *Proceedings of the 8th International Conference on Music Information Retrieval, ISMIR2007* (pp. 83-88). Vienna, Austria: Austrian Computer Society.

Cunningham, S. J., Reeves, N., & Britland, M. (2003). An ethnographic study of music information seeking: Implications for the design of a music digital library. In *Proceedings of the 3rd ACM/IEEE-CS joint conference on digital libraries* (pp. 5-17). Washington, DC: IEEE Computer Society. doi: 10.1109/JCDL.2003.1204839

Dougan, K. (2012). Information seeking behaviors of music students. *Reference*

*Services Review, 40*(4), 558-573. doi: 10.1108/00907321211277369

Herlocker, J. L., Konstan, J. A., Terveen, L. G., & Riedl, J. T. (2004). Evaluating collaborative filtering recommender systems. *ACM Transactions on Information Systems (TOIS), 22*(1), 5-53. doi: 10.1145/963770.963772

Kamalzadeh, M., Baur, D., & Möller, T. (2012). A survey on music listening and management behaviours. In F. Gouyon, P. Herrera, L. G. Martins, & M. Müller (Eds.), *Proceedings of the 13th International Society for Music Information Retrieval Conference, ISMIR 2012* (pp. 373-378). Porto, Portugal: FEUP Edições.

Konstan, J. A., & Riedl, J. (2012). Recommender systems: From algorithms to user experience. *User Modeling and User-Adapted Interaction, 22*(1/2), 101-123. doi: 10.1007/s11257-011-9112-x

Kwon, K., Cho, J., & Park, Y. (2009). Influences of customer preference development on the effectiveness of recommendation strategies. *Electronic Commerce Research and Applications, 8*(5), 263-275. doi: 10.1016/j.elerap.2009.04.004

Laplante, A. (2010). Users' relevance criteria in music retrieval in everyday life: An exploratory study. In J. S. Downie & R. C. Veltkamp (Eds.), *Proceedings of the 11th International Society for Music Information Retrieval Conference, ISMIR 2010* (pp. 601-606). Victoria, Canada: International Society for Music Information Retrieval.

Laplante, A., & Downie, J. S. (2006). Everyday life music information-seeking behaviour of young adults. In R. B. Dannenberg, K. Lemström, A. Tindale, & University of Victoria, Music Intelligence and Sound Technology Interdisciplinary Colloquium (Eds.), *Proceedings of the 7th International Conference on Music Information Retrieval, ISMIR 2006* (pp. 381-382). Victoria, Canada: University of Victoria.

Lee, J. H. (2010). Analysis of user needs and information features in natural language queries seeking music information. *Journal of the American Society for Information Science and Technology, 61*(5), 1025-1045. doi: 10.1002/asi.21302

Lee, J. H., & Downie, J. S. (2004). Survey of music information needs, uses, and seeking behaviours: Preliminary findings. In R. Loureiro & C. L. Buyoli (Eds.), *Proceedings of the 5th International Conference on Music Information Retrieval, ISMIR 2004* (pp. 441-446). Barcelona, Spain: Audiovisual Institute, Pompeu Fabra University.

Ross, C. S. (1999). Finding without seeking: The information encounter in the context of reading for pleasure. *Information Processing & Management, 35*(6), 783-799. doi: 10.1016/S0306-4573(99)00026-6

Shen, A., & Ball, A. D. (2011). Preference stability belief as a determinant of response to personalized recommendations. *Journal of Consumer Behaviour, 10*(2), 71-79. doi: 10.1002/cb.350

Tang, M. C., Sie, Y. J., & Ting, P. H. (2014). Evaluating books finding tools on social media: A case study of aNobii. *Information*

*Processing & Management, 50*(1), 54-68. doi: 10.1016/j.ipm.2013.07.005

Weigl, D., & Guastavino, C. (2011). User studies in the music information retrieval literature. In A. Klapuri & C. Leider (Eds.), *Proceedings of the 12th International Society for Music Information Retrieval Conference,* *ISMIR 2011* (pp. 335-340). Coral Gables, FL: University of Miami.

White, R. W., & Roth, R. A. (2009). Exploratory search: Beyond the query-response paradigm. *Synthesis Lectures on Information Concepts, Retrieval, and Services, 1*(1), 1-98. doi: 10.2200/S00174ED1V01Y200901ICR003

15

# 使用者偏好屬性對音樂發掘工具效能的影響
## —以Spotify音樂串流服務為例

# Evaluating Music Discovery Tools on Spotify:
## The Role of User Preference Characteristics

**唐牧群[1]　楊莽原[2]**
**Muh-Chyun Tang[1], Mang-Yuan Yang[2]**

## 摘　要

　　本研究以 Spotify 為研究平台，探討音樂社交軟體的使用者使用不同音樂發掘工具進行音樂欣賞時的主觀評價和客觀推薦成效，以及與使用者偏好結構之間的關係。本研究以實驗法為主，一共有26位參與者，採用拉丁方格的組內設計，每位參與者都使用了4種音樂發掘工具（地區排行導覽工具、情境風格導覽工具、曲目電臺推薦工具、音樂追蹤導覽工具）在限定時間內探索並存取喜好的歌曲，所有參與者和系統互動的過程都以螢幕錄製的方式記錄下來。為了能從多維度、更準確地評估音樂發掘工具之效用，我們使用了主觀評價和客觀推薦成效兩個測量項目：(1)通過實證型的小型實驗來測量受試者之主觀評價，自變項為 Spotify所提供的四種音樂發掘工具；中介變項為受試者的偏好結構（偏好洞見、偏好多樣性、偏好開放性）；依變項為實驗後問卷中收集的受試者主觀評價；(2)客觀推薦成效則由受試者在實驗中產生的曲目集合數量之比例決定，即以受測者所感興趣的曲目相較於工具所推薦的歌曲數目的比例。質化研究的部份，採用訪談法，通過實驗後對受試者進行針對性的訪談，為量化研究的結果提供檢定、補充和解釋。研究結果發現：一、不同音樂發掘工具的推薦效用的確有所差異。二、使用者面對不同音樂發掘工具時的主觀評價與客觀推薦成效並不一致。三、使用者的個人偏好結構的確會影響音樂發掘工具的推薦效用。

關鍵字：Spotify、音樂串流服務、系統評估、偏好屬性、推薦系統

[1,2]國立臺灣大學圖書資訊學系暨研究所
　Department and Graduate Institute of Library and Information Science, National Taiwan University, Taipei, Taiwan
* 通訊作者Corresponding Author: 唐牧群Muh-Chyun Tang, E-mail: mctang@ntu.edu.tw