

Impact of Consuming Suggested Items on the Assessment of Recommendations in User Studies on Recommender Systems *

Benedikt Loepp, Tim Donkers, Timm Kleemann and Jürgen Ziegler

University of Duisburg-Essen

{benedikt.loepp, tim.donkers, timm.kleemann, juergen.ziegler}@uni-due.de

Abstract

User studies are increasingly considered important in research on recommender systems. Although participants typically cannot consume any of the recommended items, they are often asked to assess the quality of recommendations and of other aspects related to user experience by means of questionnaires. Not being able to listen to recommended songs or to watch suggested movies, might however limit the validity of the obtained results. Consequently, we have investigated the effect of consuming suggested items. In two user studies conducted in different domains, we showed that consumption may lead to differences in the assessment of recommendations and in questionnaire answers. Apparently, adequately measuring user experience is in some cases not possible without allowing users to consume items. On the other hand, participants sometimes seem to approximate the actual value of recommendations reasonably well depending on domain and provided information.

1 Introduction

Recommender Systems (RS) are well-established and widely used means for helping users in finding items (e.g. commercial goods, hotels, movies, music) that best match their preferences. For a long time, these systems and the quality of their recommendations have been investigated predominantly with respect to objective accuracy of underlying algorithms [Konstan and Riedl, 2012]. Exclusively relying on offline experiments and accuracy metrics—as usual in machine learning—has however been considered insufficient in the area of RS since long ago [McNee *et al.*, 2006], in particular as user satisfaction does not necessarily go in hand with system accuracy [Konstan and Riedl, 2012]. Nevertheless, more user-centric evaluation has gained increasing attention in RS research only in recent years [Knijnenburg and Willemsen, 2015]. In this context, conducting user studies for assessing actual user experience plays an important role, especially in academia. In these studies, it is common practice

that participants are first asked to use a RS and subsequently to fill in a questionnaire [Gunawardana and Shani, 2015; Knijnenburg and Willemsen, 2015]. Based on questionnaire items such as “I liked the products recommended by the system” or “the recommendations contained a lot of variety” [Knijnenburg *et al.*, 2011], researchers then draw inferences about various aspects of RS such as perceived recommendation quality or diversity. As known from e.g. online shopping or hotel booking, the recommended products are usually represented by textual descriptions, pictures and metadata. Only in rare cases, participants can actually consume them during the studies. Thus, they have to judge corresponding recommendations based on limited knowledge. In real-world scenarios, e.g. when people want to rate a product on *Amazon* or review a hotel on *Booking.com*, it is in contrast often required to have bought a product or visited a hotel before being able to do so. This led us to the following questions:

- What is the impact of item consumption on the assessment of recommendations in RS user studies?
- Are there domain-specific differences that determine whether users can adequately assess the value of recommendations without experiencing the items?

We addressed these questions in our paper for the *12th ACM Conference on Recommender Systems* [Loepp *et al.*, 2018]. We presented two user studies which we conducted in different domains, music and movies, to investigate possible differences in pre- and post-consumption assessment of recommendation quality and aspects related to user experience of RS. For this, recommendations were presented as usual, i.e. with descriptive data, but we also enabled participants to consume items, i.e. listen to songs or watch movies (see Fig. 1). The paper at hand represents an abridged version.

2 Related Work

While the specific influence of item consumption on the assessment of recommendations has never been investigated before, related work can be found in the stream of research that aims at explaining to users why certain items are recommended [Tintarev and Masthoff, 2015]. For instance, it has been found that recommended items get over- or underestimated depending on type and quality of explanations [Tintarev and Masthoff, 2008]. However, it was in general not possible for participants of such studies to consume products. In some cases, this was at least approximated, e.g. by

*This paper is an abridged version of our paper “Impact of Item Consumption on Assessment of Recommendations in User Studies” presented at the *12th ACM Conference on Recommender Systems*.

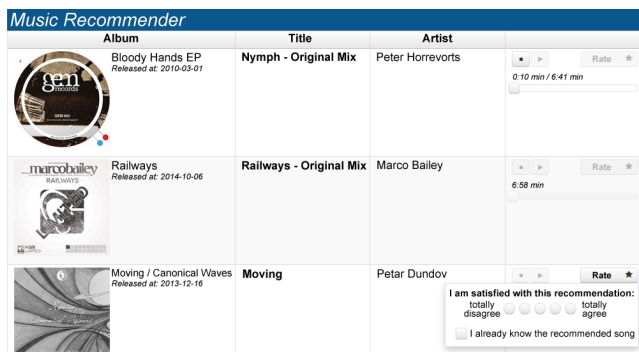


Figure 1: Screenshot of our recommender: As customary, recommendations were presented by images and metadata. However, we also enabled participants to consume items (rightmost column).

allowing to read *Amazon* detail pages of recommended books [Bilgic and Mooney, 2005]. In the experiments described by Tintarev and Masthoff [2012], watching movies was only possible in one case, while user reviews of these movies were shown otherwise. The authors compared before- and after-ratings, but focused on the effectiveness of explanations provided in addition to a very limited presentation of recommendations. Another exception is the study by Sharma and Cosley on different ways to explain music recommendations [2013], where listening to recommended songs was allowed.

Beyond that, other works in RS research have shown, for instance, that the point in time user preferences were elicited plays an important role as users provide lower ratings the longer ago they experienced an item [Bollen *et al.*, 2012]. Adomavicius *et al.* investigated anchoring effects in rating behavior [2013], but found no differences when ratings predicted by the system were shown as anchors before or after watching a TV show. As participants were only asked for their opinion after consumption, the actual influence of experiencing items remained unclear. Differences in user behavior and consistency of ratings can, however, have a considerable effect on RS performance [Said and Bellogín, 2018]. Also, the presentation of recommended items is long known for its impact [Cosley *et al.*, 2003]. For instance, Nanou *et al.* or Graus and Willemsen accompanied movie recommendations by trailer videos [2010; 2016], but did not analyze whether availability of such options affects user experience, and thus resulting questionnaire responses.

In summary, it therefore seemed to be of particular interest to examine possible differences in RS user studies between an assessment of recommendations and of related aspects before item consumption and an assessment afterwards.

3 Experiments

We hypothesized that actually listening to recommended songs or watching suggested movies makes a difference in how users assess subjective aspects of a RS, e.g. perceived quality of recommendations and their diversity, and aspects related to user experience of a RS, e.g. difficulty to settle on an item from a list of recommendations and satisfaction with this choice [Knijnenburg *et al.*, 2011; Knijnenburg and

Willemsen, 2015]. We assumed there would be differences between assessments done before or after consumption.

3.1 User Study 1: Songs

For the study on song recommendations, we set up one condition with questionnaires before and after consumption. To account for carryover effects and to control the possible influence of a pre-consumption assessment, we additionally set up a condition with a questionnaire only afterwards. However, we in this abridged paper focus on results regarding intra-individual differences, i.e. the first condition.

Method. In our controlled experiment, we had 40 participants (22 female), average age of 26.00 ($SD=8.69$). We assigned them to conditions in counter-balanced order in a between-subject design ($N=21$ for the first condition). Participants reported liking music a lot ($M=4.05$, $SD=1.04$). Yet, 28% did not know any of the recommended songs, the rest only a few ($M=1.42$, $SD=1.30$). For recommending and playing songs, we implemented a web application using the Spotify API (see Fig. 1). First, participants had to select 3 out of 110 Spotify genres. Next, they were presented with a list of 5 recommendations (generated using the API with selected genres as seed data). Song titles, artists, album titles and covers were displayed in the first condition. Participants were required to rate their satisfaction with each recommendation and fill in the questionnaire (t_1). Then, the recommendation list was shown again. Participants were asked to listen to each song for at least 30 sec with the possibility to stop, pause and forward. Finally, they again had to rate the recommendations and fill in the questionnaire (t_2).

Questionnaire. For composing the questionnaire, we relied on established instruments for user-centric evaluation of RS. We used constructs from Knijnenburg *et al.* [2011] and Pu *et al.* [2011] that have been shown to operationalize system aspects and user experience reasonably well with a limited number of questions. We generated items ourselves to ask whether participants were in doubt when selecting recommendations and which criteria they found most influential. In addition, we asked how likely they would change their ratings when they could listen to songs, and after consumption, which reasons they had to change them (open-ended question). All items were rated on a 1–5 Likert response scale.

Results and Discussion

We fitted linear mixed-effect models for each dependent variable, measured by one or more questionnaire items, with condition and point in time as a fixed factor, specified point in time as a repeated measurement (before, after), and conducted custom hypothesis tests for fixed effect parameters.

Tab. 1 shows the interaction terms as well as the differences between the two points in time t_1 and t_2 , i.e. the within comparison of the first condition (see the original paper for between-subject effects). The hypothesis tests in case of significant interaction terms confirmed, among others, that participants gave higher ratings to recommendations after listening to recommended songs (cf. Fig. 2). Questionnaire results were in line: Prior to consumption, perceived recommendation quality showed a significant correlation with mean recommendation rating ($r=.603$, $p=.004$). Correlations of other

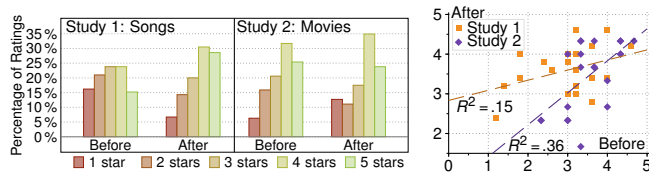


Figure 2: Distribution of ratings for recommendations (left) and scatter plot of mean recommendation ratings (right).

constructs such as overall satisfaction ($r = .575, p = .006$) confirmed this result. The difference of mean recommendation rating between the two assessments was larger, the lower perceived recommendation quality was a priori ($r = -.709, p < .001$). Listening seemed to have more influence when recommendations were initially perceived less appropriate, while participants saw little reason to change ratings otherwise. Also, ratings were normally distributed with variance of 0.77 prior to consumption, while the distribution was bounded with variance of 0.33 afterwards (cf. Fig. 2, left). Apparently, participants had difficulties to form a strong opinion before consumption [Tintarev and Masthoff, 2012], but became more certain through listening.

Choice and overall satisfaction were higher after consumption (Tab. 1). Since participants already needed to settle for an item a priori, choice difficulty was lower when they were confronted with the same recommendation list again. Mean recommendation rating tended to correlate between the two points in time ($r = .390, p = .080$, see Fig. 2, right), which is consistent with prior work on re-rating behavior [Hill *et al.*, 1995; Cosley *et al.*, 2003]. Overall, Fig. 2 underlines that scores were on average higher after consumption. We found similar correlations for questionnaire results, e.g. choice ($r = .510, p = .018$) and overall ($r = .600, p = .004$) satisfaction.

Listening to songs also had a significant effect on the perceived sufficiency of the information that came with the recommendations (Tab. 1). The difference between assessments was higher, the fewer items were known ($r = -.492, p = .023$), i.e. in typical RS scenarios where novelty of recommended items is a widely pursued goal, consumption seems especially important to support users. Artist information was the most influential criterion before consumption, followed by song title and album cover. 81% stated it would be very useful to listen to songs or at least extracts. Accordingly, listening was considered most influential afterwards, underlining the need for consuming items from a subjective perspective.

Qualitative comments supported this: Participants reported that “artist or album title are not meaningful, while it is important to like how a song sounds” and that “listening allowed imagining how well a song fits the own taste”. Others “had bad expectations when reading the artist’s name” or even found “the name of an artist very misleading”, but eventually wrote: “I was positively surprised when I heard the song” or that “as soon as I played the song, it seemed quite okay”. Another reason for changing their opinion was that they “knew one song, but could only remember and rate it after listening”.

Overall, participants were more satisfied when they found

information sufficient a priori ($r = .745, p < .001$). When sufficiency was low, they had more doubts ($r = .463, p = .034$). The more doubts, the more they reported they would change their ratings due to listening ($r = -.682, p = .001$). One participant summarized: “The information was not sufficient to build an opinion on the quality of the recommendations. Without having listened to at least an extract it was impossible to rate the songs: Even when the cover is bad or the artist has a strange name, a song might still be good”.

In conclusion, it seemed that participants had difficulties to correctly assess several aspects of the music RS, especially those related to user experience. This shows that the typical design of user studies may contribute to an inaccurate picture compared to when users can actually experience the items.

3.2 User Study 2: Movies

The second study was designed similar to study 1. We again had two conditions, but focus on the within comparison in this abridged version.

Method. We again had 40 participants (30 female), average age of 21.78 ($SD = 3.77$). We assigned them to conditions as in study 1 ($N = 21$ in the first condition). They reported liking movies ($M = 3.63, SD = 0.98$). As item data, we used 13 short movies available at *YouTube* recommended in an online article of the German newspaper *Zeit* (<http://bit.ly/zeit-movies>). Only 2 participants knew one of the movies before. First, participants had to provide demographics. Then, they were presented with a list of 3 pseudo movie recommendations (chosen randomly from a movie category they initially selected). In the first condition, we displayed movie titles, posters, genres, further metadata and (subjective) description texts by the article’s author. Participants were required to rate their satisfaction with each recommendation and fill in the questionnaire. Next, they had to select a movie they would like to watch and answer questions regarding this choice (t_1). Then, they had to watch this movie and were subsequently asked to re-rate their satisfaction with this recommendation and answer corresponding questions (t_2). Afterwards, they had to watch and assess the two remaining movies. Eventually, they again had to choose one movie (independent of their previous choice) and answer questions regarding this choice (t_3).

Questionnaire. The questionnaire was similar to study 1. Due to slightly different design, questions regarding the chosen item were now asked separately (at t_1, t_2 and t_3).

Results and Discussion

We fitted mixed models as in the first study. Tab. 1 shows the within comparison, the original paper also addresses between-subject effects.

In contrast to study 1, only few interaction terms were significant, with custom hypothesis tests showing no differences between points in time in the first condition. Still, mean ratings for individual recommendations were in line with questionnaire results: Before consumption, we found high correlations for perceived recommendation quality ($r = .515, p = .017$) and overall satisfaction ($r = .634, p = .002$).

In comparison to study 1, the correlation between points in time with respect to mean recommendation rating was higher ($r = .600, p = .004$, see Fig. 2, right). At the same time, Fig.

	Study 1: Songs				Study 2: Movies			
	Interact.	Before vs. after consum.		Interact.	Before vs. after consum.			
	Sig.	Est. Diff.	Std. Err.	Sig.	Sig.	Est. Diff.	Std. Err.	Sig.
Perc. Rec. Quality ¹	.390	0.38	0.28	.183	.467	-0.14	0.17	.411
Mean Rec. Rating	.009*	0.59	0.18	.004*	.771	-0.08	0.14	.578
Choice Satisfaction ¹	.000*	0.71	0.21	.003*	.020*	-0.19	0.25	.450
Choice Difficulty ¹	.001*	1.14	0.29	.001*	.968	0.05	0.31	.877
Effort ¹	.415	0.21	0.16	.196	.012*	-0.07	0.08	.383
Effectiveness ¹	.000*	0.81	0.19	.000*	.479	-0.14	0.22	.520
Diversity ¹	.056	-0.38	0.26	.151	.117	0.24	0.19	.224
Novelty ²	.288	-0.19	0.13	.144	.218	0.14	0.09	.106
Info. Sufficiency ²	.000*	1.48	0.38	.000*	.041*	-0.33	0.23	.149
Transparency ²	.104	0.48	0.22	.051	.763	-0.14	0.21	.499
Confidence & Trust ²	.017*	0.54	0.20	.014*	.787	0.04	0.16	.826
Doubts	.000*	2.19	0.33	.000*	.680	-0.14	0.27	.605
Overall Satisfact. ²	.005*	0.62	0.20	.005*	.442	-0.14	0.22	.525

Table 1: Results of our mixed models: Positive differences indicate better results after consumption (*Choice Diff.*, *Effort* and *Doubts* are reversed accordingly). Items marked with 1 are from [Knijnenburg *et al.*, 2011], items marked with 2 from [Pu *et al.*, 2011].

2 is aligned with questionnaire results: There were no differences in ratings (only a slight decrease with more 1-star ratings afterwards, but also more 4-star ratings). The richer information seemed to make it easier for participants to form a strong opinion already prior to consumption, leading to no significant differences with respect to questionnaire responses and a bounded distribution of ratings both before and afterwards (cf. [Tintarev and Masthoff, 2012]). Overall, this underlined that participants indeed can assess recommendations consistently—depending on domain and means to approximate their value, e.g. subjective descriptions as taken from the newspaper. Questionnaire results supported this, e.g. perceived recommendation quality correlated as well between assessments ($r = .660, p = .001$). Yet, the initial score did not seem to affect the difference found with respect to mean recommendation rating this time ($r = .157, p = .496$).

The most influential information before consumption was the newspaper description, followed by poster, genre and title. Some participants reported that “information was too basic, directors not helping (all unknown) and casts not allowing to conclude about movie quality” so that they “relied entirely on description and genre”. Although quantitative results did not differ, participants settled on a different movie when this was possible at t_3 in 62 % of all cases. Indeed, this might be influenced by participants who assumed that they would have to watch the movie again, and wanted to circumvent this. Still, selection was altered more often when perceived recommendation quality was lower prior to consumption ($r = -.470, p = .031$) and when satisfaction with their initial choice was lower after watching the respective movie (at $t_2, r = -.548, p = .010$).

Similar to study 1, participants reported that “only after watching, it became clear which movies were of most interest, while descriptions were not sufficient to decide”, that “from the movie initially chosen, more was expected after reading its description” and that their “initial impression seemed not right anymore, because the short summary was not able to convey the atmosphere”. However, this time, numerous participants stated that “the watching experience met the expectations raised by the information provided”, “ratings remained constant as descriptions allowed to get a pretty good

impression” and “summaries helped to quickly grasp what to expect, making eager to watch the movies”.

Overall, the possibility to watch recommended movies seemed to have no considerable effect, which was clearly in contrast to study 1. Apparently, domain as well as type and amount of provided information had an influence on whether participants obtained an adequate picture of user experience. Probably, it was naturally easier in the movie domain to comprehend why certain items were recommended, even without consumption. In contrast, rather abstract emotional content such as music needed to be experienced first. For example, participants wrote that “when a song evokes good mood, one automatically rates it better” and that “listening may bring up good memories, leading to higher ratings”. Furthermore, information was richer and more subjective (newspaper texts including the author’s opinion vs. song metadata). Comments supported that “descriptions corresponded well to movies” and were “written subjective and emotionally”. One participant explicitly stated that “the description revealed so much, there was no reason to change the movie’s rating after watching it”. In conclusion, it seems participants were able to estimate well whether they will like recommended movies.

4 Conclusions and Outlook

The results presented in our original paper show that it seems necessary to take questionnaire results of user studies with a grain of salt. Depending on domain as well as type and amount of information provided alongside recommendations, participants in some cases cannot adequately assess all aspects of a RS, while in others, the experience can sufficiently be substituted. For instance, we found that participants in the music domain tended to underrate songs and were less satisfied with their choice when only receiving descriptive information. Consumption in turn had a positive influence on satisfaction, leading to significantly higher scores especially for questionnaire items related to user experience. Stability of these results is however subject of future work since it has been shown, for instance, that people tend to overestimate the impact of past events [Wilson *et al.*, 2003]. On the other hand, subjective system aspects such as perceived recommendation quality were rated equally independent of consumption. For movies, this seems true in more aspects, especially if high-quality textual descriptions are available, which is more likely the case for movies than for abstract emotional content.

As a consequence, we suggest to avoid comparisons across different settings and to pay attention in user experiments when participants may not have sufficient knowledge to properly assess the system’s outcome. Only when adequate information is available, questionnaire responses may be reliable: In this case, consumption may not be needed as participants seem able to form a mental model in which they also take their individual preferences into account, allowing to judge results the same way as if the corresponding content really had been experienced. In this context, domain knowledge might play a role. For example, the fact that songs in study 1 were more often known than movies in study 2 (which were nearly completely unknown), could have introduced bias. Thus, the potential impact of item familiarity

should be further investigated. Nevertheless, user studies appear to provide at least lower bound results, which is particularly relieving for domains where it is not feasible to let participants consume recommended items (e.g. full movies, books, hotels). Beyond that, while we generally promote dissemination of user-centric evaluation methods—not only for RS—the possible impact of extensively using questionnaires on participants should not remain unconsidered as thinking consciously about decisions was found not always beneficial [Dijksterhuis and van Olden, 2006].

References

- [Adomavicius *et al.*, 2013] Gediminas Adomavicius, Jesse C. Bockstedt, Shawn P. Curley, and Jingjing Zhang. Do recommender systems manipulate consumer preferences? A study of anchoring effects. *Information Systems Research*, 24(4):956–975, 2013.
- [Bilgic and Mooney, 2005] Mustafa Bilgic and Raymond J. Mooney. Explaining recommendations: Satisfaction vs. promotion. In *Proceedings of the Beyond Personalization Workshop*, 2005.
- [Bollen *et al.*, 2012] Dirk Bollen, Mark P. Graus, and Martijn C. Willemsen. Remembering the stars? Effect of time on preference retrieval from memory. In *RecSys '12: Proceedings of the 6th ACM Conference on Recommender Systems*, pages 217–220. ACM, 2012.
- [Cosley *et al.*, 2003] Dan Cosley, Shyong K. Lam, Istvan Albert, Joseph A. Konstan, and John Riedl. Is seeing believing? How recommender interfaces affect users' opinions. In *CHI '03: Proceedings of the 21st ACM Conference on Human Factors in Computing Systems*, pages 585–592. ACM, 2003.
- [Dijksterhuis and van Olden, 2006] Ap Dijksterhuis and Zeger van Olden. On the benefits of thinking unconsciously: Unconscious thought can increase post-choice satisfaction. *Journal of Experimental Social Psychology*, 42(5):627–631, 2006.
- [Graus and Willemsen, 2016] Mark P. Graus and Martijn C. Willemsen. Can trailers help to alleviate popularity bias in choice-based preference elicitation? In *IntRS '16: Proceedings of the 3rd Joint Workshop on Interfaces and Human Decision Making for Recommender Systems*, pages 22–27, 2016.
- [Gunawardana and Shani, 2015] Asela Gunawardana and Guy Shani. *Recommender Systems Handbook*, chapter Evaluating Recommender Systems, pages 265–308. Springer US, 2015.
- [Hill *et al.*, 1995] Will Hill, Larry Stead, Mark Rosenstein, and George Furnas. Recommending and evaluating choices in a virtual community of use. In *CHI '95: Proceedings of the 13th ACM Conference on Human Factors in Computing Systems*, pages 194–201. ACM, 1995.
- [Knijnenburg and Willemsen, 2015] Bart P. Knijnenburg and Martijn C. Willemsen. *Recommender Systems Handbook*, chapter Evaluating Recommender Systems with User Experiments, pages 309–352. Springer US, 2015.
- [Knijnenburg *et al.*, 2011] Bart P. Knijnenburg, Martijn C. Willemsen, and Alfred Kobsa. A pragmatic procedure to support the user-centric evaluation of recommender systems. In *RecSys '11: Proceedings of the 5th ACM Conference on Recommender Systems*, pages 321–324. ACM, 2011.
- [Konstan and Riedl, 2012] Joseph A. Konstan and John Riedl. Recommender systems: From algorithms to user experience. *User Modeling and User-Adapted Interaction*, 22(1-2):101–123, 2012.
- [Loepp *et al.*, 2018] Benedikt Loepp, Tim Donkers, Timm Kleemann, and Jürgen Ziegler. Impact of item consumption on assessment of recommendations in user studies. In *RecSys '18: Proceedings of the 12th ACM Conference on Recommender Systems*, pages 49–53. ACM, 2018.
- [McNee *et al.*, 2006] Sean M. McNee, John Riedl, and Joseph A. Konstan. Being accurate is not enough: How accuracy metrics have hurt recommender systems. In *CHI '06: Extended Abstracts on Human Factors in Computing Systems*, pages 1097–1101. ACM, 2006.
- [Nanou *et al.*, 2010] Theodora Nanou, George Lekakos, and Konstantinos Fouskas. The effects of recommendations' presentation on persuasion and satisfaction in a movie recommender system. *Multimedia Systems*, 16(4-5):219–230, 2010.
- [Pu *et al.*, 2011] Pearl Pu, Li Chen, and Rong Hu. A user-centric evaluation framework for recommender systems. In *RecSys '11: Proceedings of the 5th ACM Conference on Recommender Systems*, pages 157–164. ACM, 2011.
- [Said and Bellogín, 2018] Alan Said and Alejandro Bellogín. Coherence and inconsistencies in rating behavior: Estimating the magic barrier of recommender systems. *User Modeling and User-Adapted Interaction*, 2018.
- [Sharma and Cosley, 2013] Amit Sharma and Dan Cosley. Do social explanations work? Studying and modeling the effects of social explanations in recommender systems. In *WWW '13: Proceedings of the 22nd International Conference on World Wide Web*, pages 1133–1144. ACM, 2013.
- [Tintarev and Masthoff, 2008] Nava Tintarev and Judith Masthoff. Over- and underestimation in different product domains. In *Proceedings of the ECAI Workshop on Recommender Systems*, pages 14–19, 2008.
- [Tintarev and Masthoff, 2012] Nava Tintarev and Judith Masthoff. Evaluating the effectiveness of explanations for recommender systems. *User Modeling and User-Adapted Interaction*, 22(4-5):399–439, 2012.
- [Tintarev and Masthoff, 2015] Nava Tintarev and Judith Masthoff. *Recommender Systems Handbook*, chapter Explaining Recommendations: Design and Evaluation, pages 353–382. Springer US, 2015.
- [Wilson *et al.*, 2003] Timothy D. Wilson, Jay Meyers, and Daniel T. Gilbert. “How happy was I, anyway?” A retrospective impact bias. *Social Cognition*, 21(6):421–446, 2003.