# Do News Consumers Want Explanations for Personalized News Rankings?

**Maartje ter Hoeve**[*]
University of Amsterdam
Amsterdam, The Netherlands
maartje.terhoeve@student.uva.nl

**Mathieu Heruer**
Blendle
Utrecht, The Netherlands
mathieu@blendle.com

**Daan Odijk**
Blendle
Utrecht, The Netherlands
daan@blendle.com

**Anne Schuth**
Blendle
Utrecht, The Netherlands
anneschuth@blendle.com

**Martijn Spitters**
Blendle
Utrecht, The Netherlands
martijn@blendle.com

**Maarten de Rijke**
University of Amsterdam
Amsterdam, The Netherlands
derijke@uva.nl

## ABSTRACT

To gain more insight in the question whether personalized news recommender systems should be responsible for their recommendations and transparent about their decisions, we study whether news consumers want explanations of why these news articles are recommended to them and what they find the best way to explain this. We survey users of Blendle's news recommendation system, and from 120 respondents we learn that news consumers do want explanations, yet do not have a very strong preference for how explanations should be shown to them. Moreover, we perform an A/B test that shows that the open rate per user does not change if users are provided with reasons for the articles recommended for them. Most likely this is because users did not pay attention to the reasons.

## CCS CONCEPTS

•**Information systems** →**Personalization;** *Recommender systems; Relevance assessment;* **Presentation of retrieval results;**

## KEYWORDS

News recommendation; Transparency; Explainable models

## 1 INTRODUCTION

The European Union has approved the General Data Protection Regulation (GDPR) on April 14, 2016. The GDPR will be enforced on May 25, 2018, and states, amongst others, that one needs to be able to explain algorithmic decisions. At the time of writing (mid 2017), the broader implications of this regulation are not clear, but there does seem to be a broadly accepted view that citizens in a transparent society are entitled to explanations of technology-driven processes, especially as algorithmic decisions increasingly influence our daily life. To which degree do citizens actually care about this? That is, are people who base their decisions and lives on the outcomes of algorithmic decisions, interested in receiving information on why a decision was made for them?

One area in which transparency and explainability are particularly important is *news*, both concerning news content and concerning the technology used to expose citizens to news (e.g. [2, 5, 11]). We focus on one aspect of technology that helps to expose citizens to news: news search and recommendation. Increasingly, news consumers use personalized services to consume news, often based on algorithmic or mixed algorithmic/editorial selections (e.g. [4, 6]). These personalized services determine to a large extent what news items their consumers read. It is tempting to state that these services should take their responsibility and be transparent about their choices by explaining their decisions to their users. However, do consumers of personalized news services care about explanations of the way in which their personalized selections were determined? We study this question in the setting of Blendle,[1] a Dutch start-up backed by amongst others The New York Times. Every day, Blendle users receive a personalized selection of news articles, selected based on a number of features that capture their reading behavior and topical interests. On top of this, Blendle users also receive a number of *must reads* every day; these articles are selected by Blendle's editorial staff and are the same for everyone. This is one of the ways to prevent users ending up in their own filter bubble. Blendle allows users to purchase a single news article instead of having to buy an entire newspaper (using micropayments) or to pre-pay via a subscription for their personal selection (called Blendle Premium). Users have the possibility to receive a refund for an article if they are not satisfied with it.

We have three research questions. Firstly, we investigate whether users would like to see explanations about why they see the articles selected for them. Secondly, we study what users find the best way to receive these explanations. Thirdly, we would like to know whether users open more articles if they are provided with explanations. In answering these research questions, our findings contribute to our understanding of the urge that news consumers feel to read articles from a transparent news recommender system, and because of this, to what extent news recommender systems should be accountable for their decisions. More broadly, our findings contribute to our understanding of how explainability can be operationalized.

---

---

[1]http://www.blendle.com

## 2 RELATED WORK

Tintarev and Masthoff [12] list seven possible aims when explaining the outcomes of an algorithm to users: *transparency*, *scrutability*, *trust*, *effectiveness*, *persuasiveness*, *efficiency* and *satisfaction*. Vig et al. [13] describe two explanation styles: *justifications* and *descriptions*. Justifications are focused on providing conceptual explanations that do not necessarily expose the underlying structure of the algorithm, whereas descriptions are meant to do exactly that. Several studies have investigated the explainability of recommender systems and the effects of adding explanations to the system (e.g. [1, 3, 8–10]). A number of these studies use *collaborative filtering* as recommendation technique [1, 3]. Collaborative filtering has been proven to be difficult to use for news recommendations due to what is known as the *cold start* or *first rater problem* [7, 14]. I.e., a news article needs to be recommended right after its release. At that moment the article has not been read yet and for this reason no information that can be used for collaborative filtering is known yet. In particular, Herlocker et al. [3] investigate the addition of explanations to the recommender system of *MovieLens*, that uses collaborative filtering as its recommendation technique. Users of *MovieLens* answer positively to the question whether they would like to see explanations added to the recommender system. This study differs from our study in its domain (i.e. news recommendations as opposed to movie recommendations), the underlying recommender system and because of that, the explanations that can be used (the aforementioned collaborative filtering) and it dates from the year 2000, whereas the recommender system research field has not been static since then. Several studies show that users are sensitive to the way explanations are shown [1, 9]. E.g., Bilgic and Mooney [1] find that users are more accurately able to decide which items are relevant for them based on "key-word style" explanations (a content based approach: which other items they interacted with before contain similar words) than on "neighbourhood style" explanations (how similar people rated this particular item).

## 3 RESEARCH QUESTIONS AND DESIGN

We address the following research questions: (RQ1) Do users want to receive explanations why particular news items are recommended to them? (RQ2) What way of showing news recommendations do users prefer? (RQ3) Do users open more articles if we provide explanations of why users see these articles? To answer these research questions, we design two experiments: a user study to answer RQ1 and RQ2 and an A/B test to answer RQ3. Both are detailed below.

### 3.1 User study

Our user study investigates whether users find it helpful to receive explanations about why particular news articles are selected for them and how they would like to see these explanations.

We designed five different types of reasons to explain our recommendations, to be judged by participants in the study. Table 1 summarizes all five reason types. Visible reasons are reasons that can be found on the card (e.g., the topic or the length of the article), invisible reasons are reasons that cannot be found on the card itself (e.g., the author). Figure 1 shows examples of items that were shown to participants.

**Table 1: Reason types used in the user study.**

| Reason type | Example |
|---|---|
| 1. Single reason, visible | Because you like politics |
| 2. Single reason, invisible | Because you like this author |
| 3. Multiple reasons, visible | Because you like politics and long articles |
| 4. Multiple reasons, combined | Because you like *De Tijd* and this author |
| 5. Bar chart | See Figure 1e |



**Figure 2: Example interface of the questionnaire, for a single question. Judgment at the top (Q4, see Table 2).**

We sent out an email questionnaire to a selection of Blendle users, 541 in total. Approximately two third of these users had a Blendle Premium subscription at the time of sending. The rest of these users used the micropayment system, but had used Blendle Premium at least once, for example via a free trial that lasted for one week.

Participants were shown three different types of explanations ("reason types") and subsequently asked to answer five questions per type. To limit the length of the survey, participants are asked to judge three types of explanations, out of the five described above. Figure 2 shows an example of the interface of the questionnaire. To make sure the results are not biased by the type or content of an article, three different articles were considered: 179 users were sent the first article, 180 users were sent the second, and 182 users were sent the third article.

Note that users were not sent the entire article, but only the introduction card to the article. This article card contains a picture, a brief introduction to the article, the name of the newspaper or the magazine, a topic, the approximate reading time of the article, how many people liked the article and the reason type. The card functions to give the news consumer a brief introduction to the article
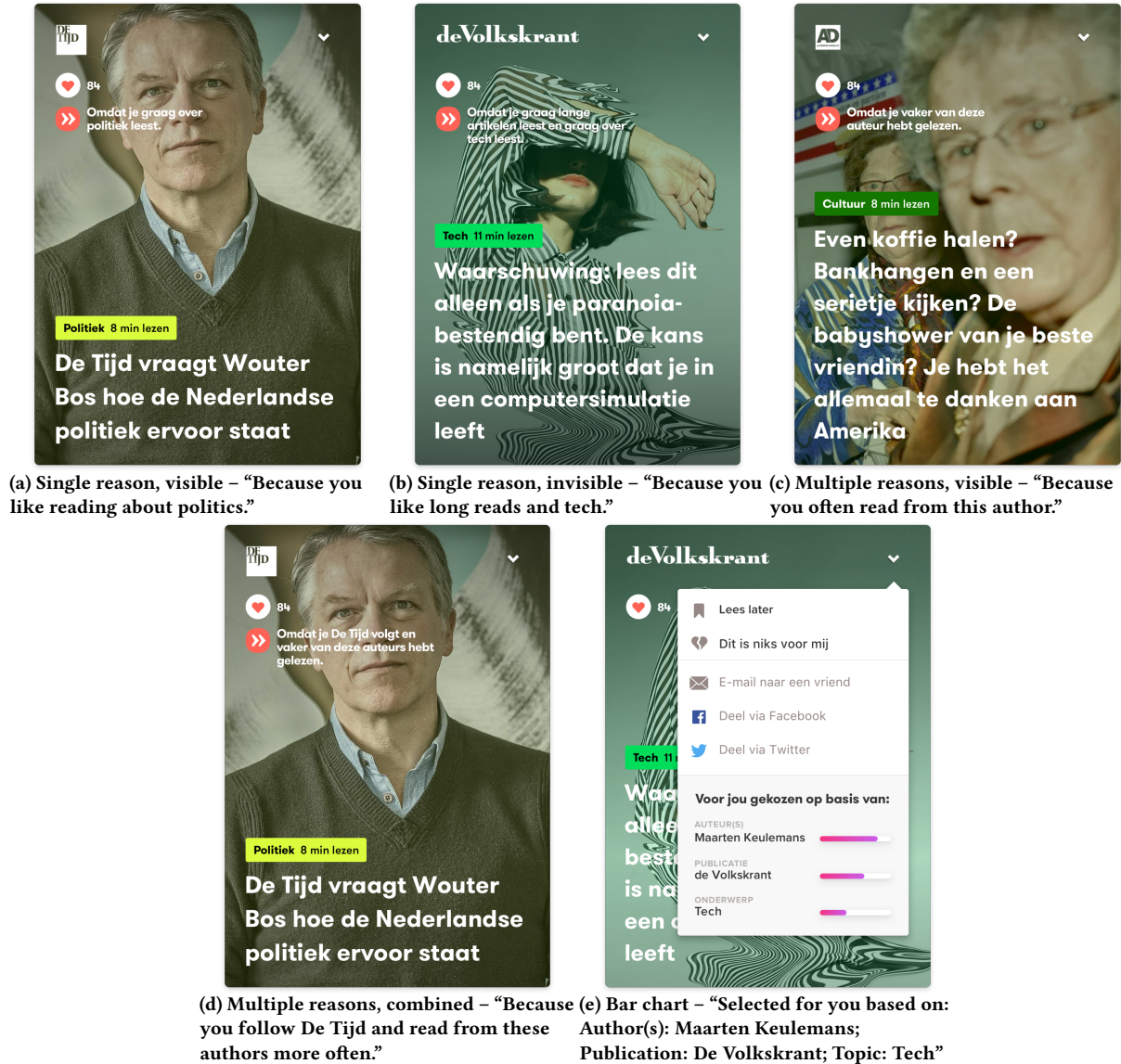
(a) Single reason, visible – "Because you like reading about politics."

(b) Single reason, invisible – "Because you like long reads and tech."

(c) Multiple reasons, visible – "Because you often read from this author."



(d) Multiple reasons, combined – "Because you follow De Tijd and read from these authors more often."

(e) Bar chart – "Selected for you based on: Author(s): Maarten Keulemans; Publication: De Volkskrant; Topic: Tech"

**Figure 1: Examples of reason types as shown to users in our user study. Textual reasons are in the lines that start with "Omdat"** (*because*)**. For the bar chart layout the reasons starts with "Voor jou gekozen" (***selected for you***). Translations are given below each article.**

Table 2: Questions used in the questionnaire as part of our user study.

|     | Type           | Question asked (English translations of the Dutch questions)                                                                                                                                                                                                                                                                    |
| --- | -------------- | ------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------ |
| Q1. | Wants reasons? | On the figure below you can see what an article currently looks like on Blendle Premium. The articles that you see are chosen based on your personal preferences and what you like to read. Imagine we would give you more information about why we chose a certain article for you. Would you find that useful? |
| Q2. | Transparency   | I understand the way that is used to explain why I see this article.                                                                                                                                                                                                                                                          |
| Q3. | Sufficiency    | I get enough information to decide whether I would like to read this article.                                                                                                                                                                                                                                                  |
| Q4. | Trust          | The reason why I see this article, makes me trust the algorithm that selected this article for me.                                                                                                                                                                                                                            |
| Q5. | Satisfaction   | I am satisfied with the way in which this article is shown to me.                                                                                                                                                                                                                                                             |

to decide whether he or she would like to read it. Figures 1a, 1b, 1c show the three different types of article cards that are used. Note that users are randomly divided over all three article types and over reason types. That is, no personalization was used here. We did not, however, completely randomize the order in which participants answer questions. First, users are either shown reason type 1 or 2, then 3 or 4. All users are shown reason type 5, as reason type 5 is very different from the other reason types. In three final questions participants are asked to fill in their age and gender and whether they would like to add some final remarks (if any).

The questions that were asked for each participant are detailed in Table 2. First, we ask participants whether they would find explanations useful and we ask them to choose between *yes*, *somewhat*, *no* or *I don't know* as possible answers. We then show several examples of explanations and ask participants to judge the examples on four Tintarev and Masthoff [12]'s dimensions: *transparency*, *sufficiency*, *trust* and *satisfaction*, all on a five point scale. We decided to omit questions on Tintarev and Masthoff [12]'s *scrutability*, *efficiency* and *effectiveness* as metrics at this stage of our research, as participants are not confronted with their own personal selection of news. For this reason, they will not be able to reliably judge whether they would open this article. Note that if participants have selected *no* or *I don't know* as a reply to whether they would like explanations, we tell them we would still like to show them some possible ways of explaining their articles and ask for their judgment.
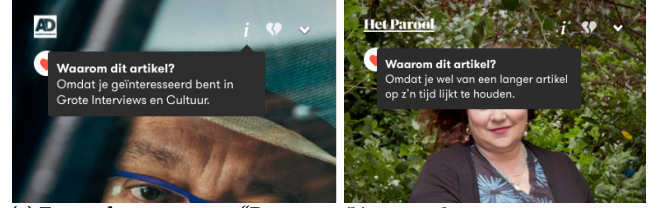
## 3.2 A/B test

In order to investigate whether users open more articles when they are provided with reasons of why they see these articles, we perform an A/B-test with two groups of Blendle users. Users are randomly assigned to a group. One of these groups is provided with explanations for the articles they see. The other group does not receive any explanations. Both groups are real Blendle users, i.e., we do not use an artificial experimental setting, but use the every day Blendle environment. The reasons shown to users in the "treatment group" are created heuristically. That is, we use a *justification* instead of an actual *description* in the sense of [13]. In our experiment, we use textual justifications. Two examples are given in Figure 3; the justifications are provided at the top of the article card, in the black boxes that pop up once a user has hovered over the "*i*" icon. This is different from the reasons tested in the user study, as we decided to launch a change in design that was as small as possible. All reasons are given in Table 4.

We run the A/B test for 24 days on 100% of our users.[2] As our objective, we measure the open rate, per day in each group.

In this study we define *open* rate as the *number of reads* over the *number of users*. We define the *number of reads* as the number articles that are opened by a user, without them asking for a refund. If users open an article multiple times (over any number of days), we only count the first time. The *number of users* is defined as the number of unique users that viewed their selection.

We test for differences in open rate between the two groups using a two-tailed paired t-test with $\alpha = 0.05$. Samples from both groups on one day form a pair. We discretize by days as news consumption

---

[2] For competitiveness reasons we cannot reveal the size of the control and treatment groups.



(a) Example reasons 1 – "Because you are interested in long interviews and Culture".

(b) Example reasons 2 – "Because you seem to like longer articles".

**Figure 3: Example reason types used during the A/B test.**

**Table 3: Participant answers to Q1: Would you like to see more information on why articles are selected for you?**

| User wants reasons | Times answered |
|---|---|
| Yes | 65 |
| Somewhat | 24 |
| No | 26 |
| I don't know | 5 |

varies over time. For the "reason group" we also count whether users have actively seen reasons, that is, hovered over the "*i*" icon. Moreover, we track whether users have seen reasons within two minutes before opening the article and if so, which reason that was.

## 4 RESULTS AND DISCUSSION

Here we answer our research questions. The first two questions are answered in Sections 4.1 and 4.2 by analyzing the results of our user study. In Section 4.3 we use the results of our A/B test to answer the last question.

A total of 120 users filled out our survey, of which 41 answered questions about the first article type, 36 about the second and 43 about the third article type. Of these 120 users, 103 users had a Blendle Premium subscription, while 17 users used the micropayment system at the time of sending out the survey. As there are not enough responses of non-premium users to put them in a separate group, we perform our analysis on all respondents together.

## 4.1 Do users want recommendation reasons?

Table 3 shows the results of what users answered to the question whether they would like to see better explained why they see articles in their selection. The significant majority answered *yes* or *somewhat* to this question, if compared to the number of people that answered *no* or *I don't know* ($\chi^2 = 14.55, p < 0.001$).

## 4.2 Do users want a particular type of recommendation reasons?

Table 6 shows the total average and standard deviation on all three articles combined, as well as the mean and standard deviation per question per article. Table 7 shows whether the differences in scores for the different types of questions are statistically significant or not. As the answers are independent, yet not necessarily sampled from the normal distribution, we use the two-sided Mann-Whitney U test, with $\alpha = 0.05$ as significance level. The sample sizes can be found in Table 5. From these results a few points stand out. First of all, although users do want more information about why they
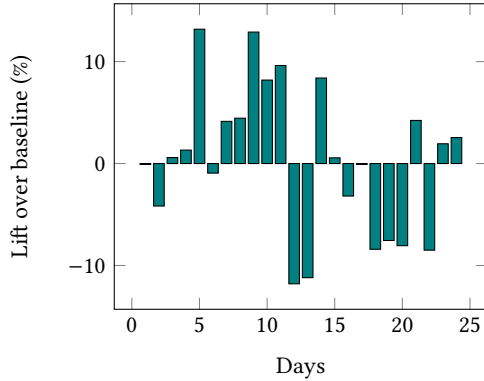
Figure 4: Lift in open rate for the group with recommendation reasons over the baseline without reasons.

see a certain article, the results do not show a clear preference as to which type of explanation users prefer. Only a few differences were significant (shown in boldface in Table 7). However, when we correct for the number of comparisons that we make, and take $\alpha = 0.001$ as significance level (using the Bonferroni correction and dividing our original $\alpha$ by 50, the number of comparisons that we make), none of the reason types scores significantly higher than another reason type. Another interesting point to make is that the standard deviations of the scores on the fifth reasoning type are, on average, bigger than the standard deviations of the scores on the other reasoning types, i.e., users either seem to like this way of showing reasons, or they do not.

## 4.3 Do users open more articles when provided with explanations?

In our A/B test, after 24 days, we see that users that were shown the recommendation reasons (the "reason group") have a lift in open rate of 0.33%. This difference is plotted in Figure 4 and is not significant ($t = -0.29, p = 0.77$).

Of all individual users in the reason group, 9.8% has seen at least one recommendation reason. Of all users who opened an article, 1.08% had seen the recommendation reason within two minutes before they opened that particular article. These users saw 1.27 reasons on average, with a standard deviation of 0.73. That is, not many users saw the reasons, which explains why we do not observe a difference in open rate per user between the two groups. Different, more prominent designs, may yield different results.

Figure 5 shows how often users saw each particular reason, in comparison to the total number of times users saw a reason. Reason type 6 is seen most often. This is the explanation that is given for the *must-reads*, i.e., not based on any form of personalization. These must-reads are on top of the user's page, which can bias these results.

## 5 CONCLUSION

In this study we investigated whether news consumers would like to receive explanations about why articles were selected for their personalized selections of news articles. We also investigated how they would prefer to receive these explanations. Moreover, we studied whether news consumers open more articles, if they are provided with reasons.
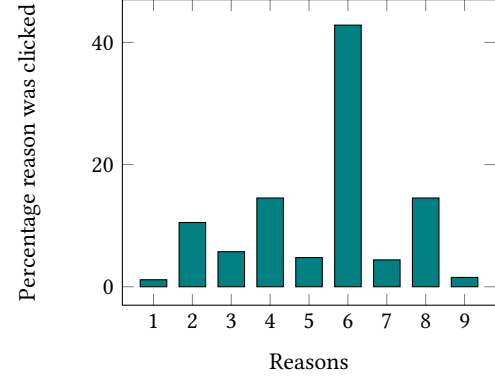


Figure 5: Reasons clicked before opening the article, see Table 4 for mapping.

Table 4: Reason mapping for reasons used in A/B test.

| Number | Reason |
|--------|--------|
| Reason 1 | Because you often read about TOPIC |
| Reason 2 | Because you are interested in TOPIC |
| Reason 3 | Because we think NEWSPAPER could be interesting for you |
| Reason 4 | The editors really liked this piece |
| Reason 5 | Because you follow NEWSPAPER |
| Reason 6 | According to the editors, this is one of the best stories of the day. No matter your preferences |
| Reason 7 | Because you often read from NEWSPAPER |
| Reason 8 | Because you seem to like a long read every now and then |
| Reason 9 | Because you often read from AUTHOR |

Table 5: Sample sizes per reason type

| Type 1 | Type 2 | Type 3 | Type 4 | Type 5 |
|--------|--------|--------|--------|--------|
| 66 | 56 | 63 | 55 | 120 |

Our questionnaire showed that a large majority of the respondents would like to receive these explanations, yet they do not show a clear preference as to how they would like to see these. Our A/B test shows that the open rate per user does not increase by adding explanations. In fact, in many cases, users do not read the the explanations.

More broadly, our research shows that users nowadays still attach importance to explanations of algorithmic decisions broader than the domain described in [3] and it motivates us to strive for transparent, responsible and accountable recommender systems.

Even though we tested several designs for explanations in our questionnaire, the number of options that we were able to expose to our participants was limited. It could very well be that alternative designs would be preferred by news consumers.

Hence, as future work, we recommend that A/B tests with additional designs are conducted. They may either result in a clearer preference for a particular way of explaining recommendations or further strengthen our conclusions. We especially recommend conducting A/B tests with reasons clearly visible, that is, not behind an icon as in the work reported here. More research in different domains, with different user groups, should lead to insights into the generalizability of our findings.

**Table 6: Mean and standard deviations of the scores on different types of judgments in the user study. The "reason types" refer back to the types of reason listed in Table 1.**

| Reason type | 1 | | 2 | | 3 | | 4 | | 5 | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Question** | **Mean** | **Std** | **Mean** | **Std** | **Mean** | **Std** | **Mean** | **Std** | **Mean** | **Std** |
| Transparency | 3.697 | 1.141 | 3.786 | 1.129 | 3.587 | 1.107 | **3.873** | 1.096 | 3.650 | **1.339** |
| Sufficiency | 3.530 | 1.076 | 3.625 | 1.028 | 3.333 | 1.098 | **3.764** | 0.953 | 3.408 | **1.275** |
| Trust | 3.000 | 1.115 | 3.250 | 1.122 | 3.032 | 1.023 | 3.400 | 0.984 | **3.500** | **1.258** |
| Satisfaction | 3.606 | 0.919 | **3.661** | 0.969 | 3.317 | 1.096 | 3.582 | 1.073 | 3.233 | **1.327** |
| *Average* | 3.458 | 0.798 | 3.580 | 0.836 | 3.317 | 0.916 | **3.655** | 0.798 | 3.448 | **1.154** |

**Table 7: Statistical differences between reason types, between different questions.**

| | Type 2 | Type 3 | Type 4 | Type 5 |
|---|---|---|---|---|
| **Type 1** | | | | |
| Transparency | $U = 1811.0, p > 0.05$ | $U = 2059.0, p > 0.05$ | $U = 1597.0, p > 0.05$ | $U = 3858.0, p > 0.05$ |
| Sufficiency | $U = 1868.0, p > 0.05$ | $U = 2147.5, p > 0.05$ | $U = 1571.5, p > 0.05$ | $U = 4005.0, p > 0.05$ |
| Trust | $U = 1860.5, p > 0.05$ | $U = 2016.5, p > 0.05$ | $U = 1512.0, p > 0.05$ | $U = 3001.0, p < 0.05$ |
| Satisfaction | $U = 1748.0, p > 0.05$ | $U = 2347.0, p > 0.05$ | $U = 1740.0, p > 0.05$ | $U = 4304.0, p > 0.05$ |
| Average | $U = 1684.0, p > 0.05$ | $U = 2257.5, p > 0.05$ | $U = 1591.0, p > 0.05$ | $U = 3848.0, p > 0.05$ |
| **Type 2** | | | | |
| Transparency | | $U = 1838.5, p > 0.05$ | $U = 1422.5, p > 0.05$ | $U = 3404.0, p > 0.05$ |
| Sufficiency | | $U = 1899.5, p > 0.05$ | $U = 1397.0, p > 0.05$ | $U = 3529.0, p > 0.05$ |
| Trust | | $U = 1938.0, p > 0.05$ | $U = 1472.0, p > 0.05$ | $U = 2900.0, p > 0.05$ |
| Satisfaction | | $U = 2038.5, p > 0.05$ | $U = 1523.0, p > 0.05$ | $U = 3734.5, p > 0.05$ |
| Average | | $U = 2041.0, p > 0.05$ | $U = 1493.0, p > 0.05$ | $U = 3505.0, p > 0.05$ |
| **Type 3** | | | | |
| Transparency | | | $U = 1417.5, p > 0.05$ | $U = 3476.0, p > 0.05$ |
| Sufficiency | | | $U = 1324.0, p < 0.05$ | $U = 3472.0, p > 0.05$ |
| Trust | | | $U = 1411.5, p > 0.05$ | $U = 2847.5, p < 0.05$ |
| Satisfaction | | | $U = 1441.5, p > 0.05$ | $U = 3657.5, p > 0.05$ |
| Average | | | $U = 1369.5, p < 0.05$ | $U = 3416.5, p > 0.05$ |
| **Type 4** | | | | |
| Transparency | | | | $U = 3469.5, p > 0.05$ |
| Sufficiency | | | | $U = 3676.0, p > 0.05$ |
| Trust | | | | $U = 2992.0, p > 0.05$ |
| Satisfaction | | | | $U = 3586.0, p > 0.05$ |
| Average | | | | $U = 3575.0, p > 0.05$ |

## REFERENCES

[1] M. Bilgic and R. J. Mooney. Explaining recommendations: Satisfaction vs. promotion. In *Beyond Personalization Workshop, IUI*, volume 5, page 153, 2005.
[2] N. Diakopoulos. Algorithmic accountability: Journalistic investigation of computational power structures. *Digital Journalism*, 3(3):398–415, 2015.
[3] J. L. Herlocker, J. A. Konstan, and J. Riedl. Explaining collaborative filtering recommendations. In *Proceedings of the 2000 ACM conference on Computer supported cooperative work*, pages 241–250. ACM, 2000.
[4] I. Ilievski and S. Roy. Personalized news recommendation based on implicit feedback. In *Proceedings of the 2013 International News Recommender Systems Workshop and Challenge*, pages 10–15. ACM, 2013.
[5] M. Karlsson. The immediacy of online news, the visibility of journalistic processes and a restructuring of journalistic authority. *Journalism*, 12(3):279–295, 2011.
[6] J. Liu, P. Dolan, and E. R. Pedersen. Personalized news recommendation based on click behavior. In *Proceedings of the 15th international conference on Intelligent user interfaces*, pages 31–40. ACM, 2010.
[7] P. Melville, R. J. Mooney, and R. Nagarajan. Content-boosted collaborative filtering for improved recommendations. In *Aaai/iaai*, pages 187–192, 2002.
[8] C. Musto, F. Narducci, P. Lops, M. De Gemmis, and G. Semeraro. Explod: A framework for explaining recommendations based on the linked open data cloud. In *Proceedings of the 10th ACM Conference on Recommender Systems*, pages 151–154. ACM, 2016.
[9] P. Pu and L. Chen. Trust-inspiring explanation interfaces for recommender systems. *Knowledge-Based Systems*, 20(6):542–556, 2007.
[10] M. T. Ribeiro, S. Singh, and C. Guestrin. Why should I trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144. ACM, 2016.
[11] J. B. Singer. Contested autonomy: Professional and popular claims on journalistic norms. *Journalism studies*, 8(1):79–95, 2007.
[12] N. Tintarev and J. Masthoff. A survey of explanations in recommender systems. In *Data Engineering Workshop, 2007 IEEE 23rd International Conference on*, pages 801–810. IEEE, 2007.
[13] J. Vig, S. Sen, and J. Riedl. Tagsplanations: explaining recommendations using tags. In *Proceedings of the 14th international conference on Intelligent user interfaces*, pages 47–56. ACM, 2009.
[14] E. Vozalis and K. G. Margaritis. Analysis of recommender systems algorithms. In *The 6th Hellenic European Conference on Computer Mathematics & its Applications*, pages 732–745, 2003.