



Exploring the Landscape of Recommender Systems Evaluation: Practices and Perspectives

CHRISTINE BAUER, Paris Lodron University Salzburg, Austria

EVA ZANGERLE, University of Innsbruck, Austria

ALAN SAID, University of Gothenburg, Sweden

Recommender systems research and practice are fast-developing topics with growing adoption in a wide variety of information access scenarios. In this article, we present an overview of research specifically focused on the evaluation of recommender systems. We perform a systematic literature review, in which we analyze 57 papers spanning six years (2017–2022). Focusing on the processes surrounding evaluation, we dial in on the methods applied, the datasets utilized, and the metrics used. Our study shows that the predominant experiment type in research on the evaluation of recommender systems is offline experimentation and that online evaluations are primarily used in combination with other experimentation methods, e.g., an offline experiment. Furthermore, we find that only a few datasets (MovieLens, Amazon review dataset) are widely used, while many datasets are used in only a few papers each. We observe a similar scenario when analyzing the employed performance metrics—a few metrics are widely used (precision, normalized Discounted Cumulative Gain, and Recall), while many others are used in only a few papers. Overall, our review indicates that beyond-accuracy qualities are rarely assessed. Our analysis shows that the research community working on evaluation has focused on the development of evaluation in a rather narrow scope, with the majority of experiments focusing on a few metrics, datasets, and methods.

CCS Concepts: • **Information systems** → **Recommender systems**; **Evaluation of retrieval results**; • **Human-centered computing** → **HCI design and evaluation methods**;

Additional Key Words and Phrases: Evaluation, survey, systematic literature review, recommender systems

ACM Reference format:

Christine Bauer, Eva Zangerle, and Alan Said. 2024. Exploring the Landscape of Recommender Systems Evaluation: Practices and Perspectives. *ACM Trans. Recomm. Syst.* 2, 1, Article 11 (March 2024), 31 pages. <https://doi.org/10.1145/3629170>

1 INTRODUCTION

Recommender systems aim to alleviate choice overload by providing personalized item recommendations to users. In the development and maintenance of these systems, evaluating their

This research was funded in whole, or in part, by the Austrian Science Fund (FWF): P33526. This research was funded in whole or in part by Vinnova.

Authors' addresses: C. Bauer, Paris Lodron University Salzburg, Salzburg, 5020, Jakob-Haringer-Strasse 1, Austria; e-mail: christine.bauer@plus.ac.at; E. Zangerle, University of Innsbruck, Technikerstr. 21A, Innsbruck, 6020, Austria; e-mail: eva.zangerle@uibk.ac.at; A. Said, University of Gothenburg, Sweden; e-mail: alansaid@acm.org.



This work is licensed under a [Creative Commons Attribution-NonCommercial-NoDerivs International 4.0 License](https://creativecommons.org/licenses/by-nc-nd/4.0/).

© 2024 Copyright held by the owner/author(s).

2770-6699/2024/03-ART11

<https://doi.org/10.1145/3629170>

performance is crucial. This work provides an overview of research specifically focused on the *evaluation* of recommender systems from 2017 to 2022. While evaluation is a significant aspect of the recommender systems field, our systematic literature review concentrates on research that specifically addresses the evaluation of recommender systems, covering papers that delve into methodological evaluation issues. This includes, for instance, papers describing research on new evaluation methods or metrics, papers analyzing how the design and implementation of the evaluation can impact the outcome of an analysis, research highlighting flaws in evaluation—or how evaluation can be improved. On the contrary, works that, for instance, propose a new recommendation model and validate it through evaluation or in other ways use evaluation to gauge the performance of a recommender system, thus, fall outside of the scope of this literature review.

The evaluation of recommender systems has been explored in previous works, but no systematic literature review has comprehensively examined datasets, metrics, or experiment types, and performed a quantitative analysis of the reviewed literature. One notable study by Herlocker et al. [50] focuses on collaborative filtering systems and proposes various recommendation tasks, such as identifying good items or recommending in sequence. The work also discusses the suitability of datasets and metrics for evaluating recommendation-specific tasks prevalent during that era of recommender systems research. More recently, Gunawardana et al. [45] provide an extensive overview of the evaluation processes involved in assessing recommender systems. The study examines a wide range of properties that impact user experience and explores methods for measuring these properties, encompassing the entire evaluation pipeline from research hypotheses and experimental design to metrics for quantification. Taking a specialized approach, Pu et al. [78] presents a survey on recommender system evaluation from the users' perspective. The research particularly focuses on the initial preference elicitation process, preference refinement, and the final presentation of recommendations. From the survey results, Pu et al. [78] distills a set of usability and user interface design guidelines for user-centered evaluation of recommender systems. Beel et al. [14, 15] surveyed evaluation approaches in the field of research paper recommender systems and found that 69% of the papers featured an offline evaluation while 21% do not provide an evaluation. A survey conducted by Ihemelandu and Ekstrand [51] examines the use of statistical inference in recommender systems research and reveals that 59% of the surveyed papers did not perform significance testing. The authors argue for the inclusion of statistical inference tests in recommender systems evaluation while also acknowledging the associated challenges. More recently, Zangerle and Bauer [96] present a survey on the evaluation of recommender systems, introducing the Framework for EValuating Recommender systems (FEVR). This framework conceptualizes the evaluation space of recommender systems, providing a systematic overview of essential evaluation aspects and their application. The proposed FEVR framework encompasses a wide variety of facets required for evaluating recommender systems, accommodating comprehensive evaluations that address the multi-faceted dimensions found in this domain.

In addition to survey papers, several works offer critical retrospectives and analyses of evaluation procedures and setups. For example, Ferrari Dacrema et al. [40, 41] critically analyze the performance of neural recommendation approaches published from 2015 and 2018. They compare these approaches against well-tuned, non-neural baseline methods, such as nearest-neighbor or content-based approaches, and find that the simpler methods outperform 11 of the 12 analyzed approaches. These findings suggest that limited progress has been made due to weak baselines and insufficient optimization of their parameters. Similarly, Rendle et al. [79] analyze the use of baselines in research, focusing on the MovieLens 10M and the Netflix Prize datasets. They compare the reported results of baselines with the results obtained through a re-run of the baselines, revealing substantial divergences, particularly for the MovieLens 10M dataset. They then

introduce stronger and well-tuned baselines, which outperform the proposed methods. Following the same line of investigation, Ludewig et al. [66] perform a similar analysis of evaluation for session-based recommendation approaches. They compare neural sequential recommendation approaches from 2016 to 2019 with well-tuned baseline approaches, such as nearest neighbor. Like previous works, they conclude that the claimed progress is mostly illusory, attributing it to weak baselines that are insufficiently or not at all tuned. Ludewig et al. [66] argue that this limitation is a critical drawback in current evaluation practices.

The goal of our study is to provide a quantitative snapshot of the landscape of research on the evaluation of recommender systems over the past six years. Through a systematic literature review [57] of major conferences and journals from 2017 to 2022, we analyze the evaluation methods, datasets, and metrics employed in the recommender systems community. Initially screening 339 papers, we apply defined inclusion and exclusion criteria to narrow down our review to a final sample of 57 papers. Our focus lies on three key aspects of recommender systems evaluation: (1) experiment type (offline experiments, user study, online experiment), (2) datasets, and (3) evaluation metrics.

This article is structured as follows: In Section 2, we detail the stepwise procedure for the systematic literature review. In Section 3, we present the results of our analysis with a focus on experiment type, datasets, and evaluation metrics. Finally (Section 4), we discuss the findings of this review and provide an outlook on future work.

2 MATERIAL AND METHODS

Our approach to identifying papers that are concerned with the evaluation of recommender systems relies on a systematic literature review [57]. A systematic literature review represents a systematic search for papers on a predefined topic and the analysis of the respective paper landscape. In this section, we outline the stepwise procedure for searching, filtering, categorizing, and analyzing the papers, which is visualized in Figure 1 and described in detail in the following subsections.

2.1 Literature Search

For data collection, we rely on the systematic literature review procedure as outlined in the guidelines by Kitchenham et al. [57]. To develop and pursue an effective search strategy, we performed a so-called scoping review on relevant published literature. In this scoping review, we, for instance, identified that the keyword *recommendation systems* is used interchangeably with *recommender systems*, with the latter being more common in the research community centered around the ACM Conference on Recommender Systems (RecSys), while both alternatives are used broadly in other research outlets. Moreover, as our article aims to cover research that revolves around methodological issues of evaluation, we identified that a search with the keywords *reproducible* or *reproducibility* has strong overlaps with a search for the keyword *evaluation* but also yields additional hits. Similarly, using the keywords *method* or *methodology* has proven useful to identify additional works. Further, we identified that some papers were miscategorized (e.g., as a short paper instead of research paper), necessitating the use of a broader query followed by manual filtering.

The search strategy to identify eligible papers to be included in our sample consisted of several consecutive stages. As the ACM Digital Library¹ does not only contain papers published by ACM but also by other publishers, we could use this library to search for papers in the main established conferences and journals where research on recommender system evaluation is published. Besides

¹<https://dl.acm.org>

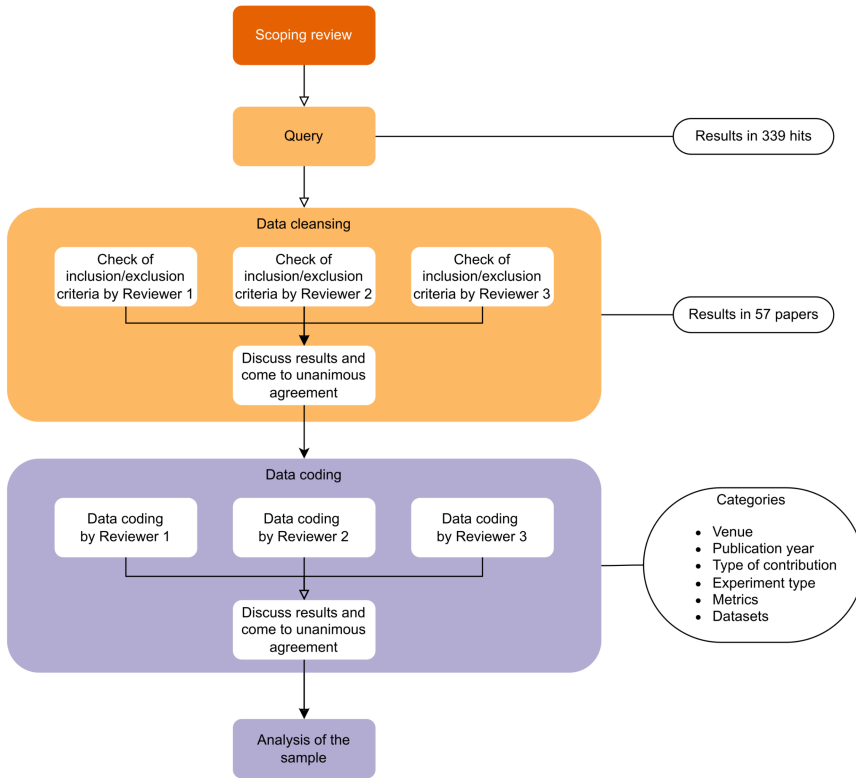


Fig. 1. Stepwise procedure for searching, filtering, categorizing, and analyzing the surveyed papers.

the main conference on recommender systems (RecSys), this embraces conferences such as SIGIR, CIKM, UMAP, and WSDM. Journals include, for instance, TOIS, UMUAI, and CSUR.

Accordingly, we sampled papers that we found in the ACM Digital Library (The ACM Guide to Computing Literature), which describes as “the most comprehensive bibliographic database in existence today focused exclusively on the field of computing.”² For reasons of reproducibility, we consider papers in an encapsulated time frame of six years, for which we can assume that the employed databases and search engines have already completed indexing the papers from conferences and journals (2017–2022). As our literature review is concerned with research on the evaluation of recommender systems, we searched for papers that were indexed with the keywords *recommend** (to cover both, *recommender systems* and *recommendation systems*), and either *evalua** (to cover *evaluation* and *evaluability*) or *reproducib** (to cover *reproducible* and *reproducibility*) or *method* or *methodology*. For papers appearing in the ACM Conference on Recommender Systems, we presume that the keywords *recommender systems* or *recommendation systems* are not necessarily used; hence, for papers appearing in RecSys, we relied solely on the keywords *evalua** or *reproducib** or *method* or *methodology*. Altogether, this resulted in the following query:³

²<https://libraries.acm.org/digital-library/acm-guide-to-computing-literature>

³https://dl.acm.org/action/doSearch?fillQuickSearch=false&target=advanced&expand=all&AfterMonth=1&AfterYear=2017&BeforeMonth=12&BeforeYear=2022&AllField=Keyword%3A%28recommend*%29+AND+Keyword%3A%28reproducib*+OR+method+OR+methodology+OR+evalua*%29+OR+ContentGroupTitle%3A%28%22ACM+Conference+on+Recommender+Systems%22%29+AND+Keyword%3A%28reproducib*+OR+method+OR+methodology+OR+evalua*%29

```

"query": {
  Keyword:(recommend*)
  AND
  Keyword:(reproducib* OR method OR methodology OR evalua*)
  OR
  ContentGroupTitle:("ACM Conference on Recommender Systems")
  AND
  Keyword:(reproducib* OR method OR methodology OR evalua*)
}
"filter": { E-Publication Date: (01/01/2017 TO 12/31/2022) }

```

This query returns a total of 339 hits (as of June 10, 2023).

We note that the query did not return any papers from the conferences CHI, CSCW, and IUI. To validate this result, for each conference separately, we searched for papers with the respective keywords without time restriction. The latest papers on the evaluation of recommender systems at CSCW and IUI were published in 2013 and at CHI in 2016.

2.2 Data Cleansing and Selection of Papers for the Sample

We retrieved the 339 papers and reviewed them against the ex-ante-defined inclusion and exclusion criteria described below.

A paper was included if it fulfilled *each and every* of the following criteria (ex-ante inclusion criteria):

- (A) The paper revolves around methodological issues of the evaluation of recommender systems.
- (B) The paper is a full research paper.
- (C) The paper is published within the time range from 01/01/2017 until and including 12/31/2022.

A paper was excluded if *any* of the following criteria were met (ex-ante exclusion criteria):

- (a) The paper is not a research paper.
- (b) The paper is a short paper, an abstract, a demo paper, a tutorial paper, or a workshop paper.⁴
- (c) The paper is not concerned with recommender systems.
- (d) The paper does not make a contribution regarding the evaluation of recommender systems.

Next, three reviewers independently screened the retrieved 339 papers against these inclusion and exclusion criteria by examining titles and abstracts, as well as the results and methodology sections. Any disagreement on paper selection was resolved by discussions to reach unanimous consensus among the three reviewers. These discussions resulted in the formulation of more specific inclusion criteria, further specifying the ex-ante inclusion criterion (A) that a paper is included if it “revolves around methodological issues of the evaluation of recommender systems.” Hence, the ex-ante inclusion criterion (A) was considered fulfilled if *any* of the following criteria was fulfilled (ex-post inclusion criteria):

- (A.1) The paper provides a literature survey on the evaluation of recommender systems.
- (A.2) The paper introduces one or more novel metrics of evaluation.

⁴We note that we did not consider the search criterion *research paper* in the query, because essential full papers were not returned by the query due to miscategorization as a short paper in the database (e.g., Reference [10]).

- (A.3) The paper analyzes metrics of evaluation.
- (A.4) The paper contributes an extensive critical evaluation across a set of approaches.
- (A.5) The paper contributes a conceptual framework for evaluation.
- (A.6) The paper contributes a framework for evaluation in the form of a toolkit.
- (A.7) The paper contributes a novel evaluation model; e.g., related to off-policy learning.
- (A.8) The paper proposes a novel sampling approach for (offline) evaluation.
- (A.9) The paper contributes to evaluation by analyzing sampling approaches.
- (A.10) The paper demonstrates or discusses how the results inform the evaluation of recommender systems.

Further, the ex-ante inclusion criterion (A) was *not* considered fulfilled if *any* of the following criteria was fulfilled (ex-post exclusion criteria):⁵

- (A.i) The paper proposes a recommendation model with or without validating it through evaluation but does not contribute to methodological issues of evaluation.
- (A.ii) The paper presents an exploratory evaluation of a recommender system but does not contribute to methodological issues of evaluation.
- (A.iii) The paper presents an experiment but does not contribute to methodological issues of evaluation.
- (A.iv) The paper analyses recommendation approaches but does not contribute to methodological issues of evaluation.
- (A.v) The paper studies psychological effects influencing the design and development of recommender systems.

This data cleansing and selection procedure led to the exclusion of 282 papers (see Appendix). The remaining 57 papers make up our final sample resulting from the query. Table 1 provides an overview of all papers in the sample.

2.3 Review of the Selected Papers in Full Text (Coding)

For each paper, we obtained meta-information on the paper from the citation information, i.e., author, year, title, type of venue—conference or journal—and venue name. In addition, to address the main purpose of this paper, we extracted the following information from the full text: experiment type, used dataset(s), used metric(s), and type of contribution. To this end, three reviewers examined the full text of the papers and extracted the respective information.

Concerning datasets and metrics, the respective information was extracted directly from the full text of the papers. Concerning the experiment type, we relied on the established differentiation between offline experiment, user study, and online experiment [96]: Offline evaluation refers to a computational evaluation without human subjects being involved in the evaluation process; user studies refer to evaluations (in live or laboratory settings) with a set of human participants that carry out tasks as defined by the researcher; and online evaluations refer to field experiments where users carry out their self-selected tasks in a real-world setting. For the type of contribution, the categorization scheme was developed inductively from raw data. The categorization scheme allowed each paper to belong to exclusively one type of contribution. An overview of the types is presented in Table 2; the specified types are benchmark, framework, metrics, model, and survey, respectively. The initial inter-rater reliability was at an acceptable level (Krippendorff's $\alpha = 0.8214$). Disagreement was resolved by discussions to reach unanimous consensus (Krippendorff's $\alpha = 1$).

⁵Note, these are also a further specification of the ex-ante exclusion criterion (d).

Table 1. Surveyed Papers, Sorted by Venue (Alphabetically) and Year

Papers	Venues	Year
Saraswat et al. [84]	AIML Systems	2021
Jannach [52]	ARTR	2023
Eftimov et al. [38]	BDR	2021
Sonboli et al. [88], Zhu et al. [99]	CIKM	2021
Ekstrand [39]	CIKM	2020
Alhijawi et al. [5], Sánchez and Bellogín [83], Zangerle and Bauer [96]	CSUR	2022
Jin et al. [54]	HAI	2021
Belavadi et al. [16]	HCII	2021
Peska and Vojtas [77]	HT	2020
Ostendorff et al. [75]	ICADL	2021
Afolabi and Toivanen [2]	IJEHMC	2020
Bellogín et al. [17]	IRJ	2017
Latifi et al. [62]	ISCI	2022
Carraro and Bridge [23]	JiIS	2022
Krichene and Rendle [60], Li et al. [63], McInerney et al. [69]	KDD	2020
Dehghani Champiri et al. [36]	KIS	2019
Latifi and Jannach [61]	RecSys	2022
Dallmann et al. [35], Narita et al. [73], Parapar and Radlinski [76], Saito et al. [82]	RecSys	2021
Cañamares and Castells [22], Kouki et al. [59], Sun et al. [90], Symeonidis et al. [91]	RecSys	2020
Ferrari Dacrema et al. [41]	RecSys	2019
Yang et al. [95]	RecSys	2018
Xin et al. [94]	RecSys	2017
Ali et al. [6]	Scientometrics	2021
Diaz and Ferraro [37], Silva et al. [87]	SIGIR	2022
Anelli et al. [10], Li et al. [64], Lu et al. [65]	SIGIR	2021
Balog and Radlinski [11], Mena-Maldonado et al. [70]	SIGIR	2020
Cañamares and Castells [21]	SIGIR	2018
Cañamares and Castells [20]	SIGIR	2017
Chen et al. [25]	TheWebConf	2019
Al Jurdi et al. [4]	TKDD	2021
Guo et al. [47]	TOCHI	2022
Zhao et al. [98]	TOIS	2022
Ferrari Dacrema et al. [40], Mena-Maldonado et al. [71]	TOIS	2021
Anelli et al. [9]	UMAP	2022
Frumerman et al. [42]	UMAP	2019
Bellogín and Said [19]	UMUAI	2021
Said and Bellogín [80]	UMUAI	2018
Chin et al. [26], Kiyohara et al. [58]	WSDM	2022
Cotta et al. [31]	WSDM	2019
Gilotte et al. [44]	WSDM	2018

Table 2. The Five Types Used to Describe the Type of Contribution Made in the Reviewed Literature

Types of Contribution	Description
Benchmark	Providing an extensive critical evaluation across a (wide) set of approaches or datasets
Framework	Introducing a framework for evaluation, which may take the form of a toolkit or a conceptual framework
Metrics	Analyzing existing or introducing novel metrics of evaluation
Model	Introducing a novel recommendation or evaluation model
Survey	A literature survey

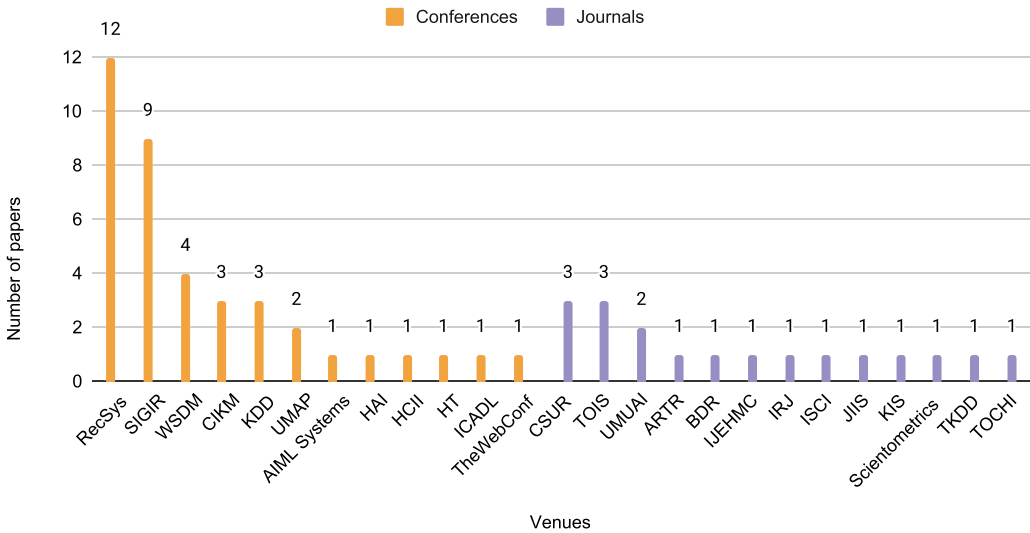


Fig. 2. Number of papers per venue, sorted by venue type (journals vs. conferences) and number of papers.

In all phases of extracting and categorizing data, all authors were engaged. Where disagreement emerged in rare cases, the authors discussed the categorization in question, drawing upon domain expertise on a case-by-case basis, until unanimous consensus was established.

3 RESULTS

In this section, we first give a general overview of papers on the evaluation of recommender systems in the analyzed time frame 2017–2022 (Section 3.1). Then, we detail the types of contributions to the discourse (Section 3.2). Further, we provide an overview of the experiment types used in the papers (Section 3.3). Section 3.4 provides an overview and discussion of the datasets used. In Section 3.5, we detail the metrics used and discussed in the papers.

3.1 General Overview

Most papers on evaluation in recommender systems are published at RecSys (the main conference concerning the research topic *recommender systems*) (12) and at SIGIR (9) (the main conference concerning the closely related research topic of *information retrieval*) (Figure 2). Notably, as can be seen from Figure 2, papers on the evaluation of recommender systems are published in a wide scale of venues (12 conference venues and 13 journal venues) where it is often only 1 paper at the respective venue in the set time frame of our review. The majority of papers on evaluation are

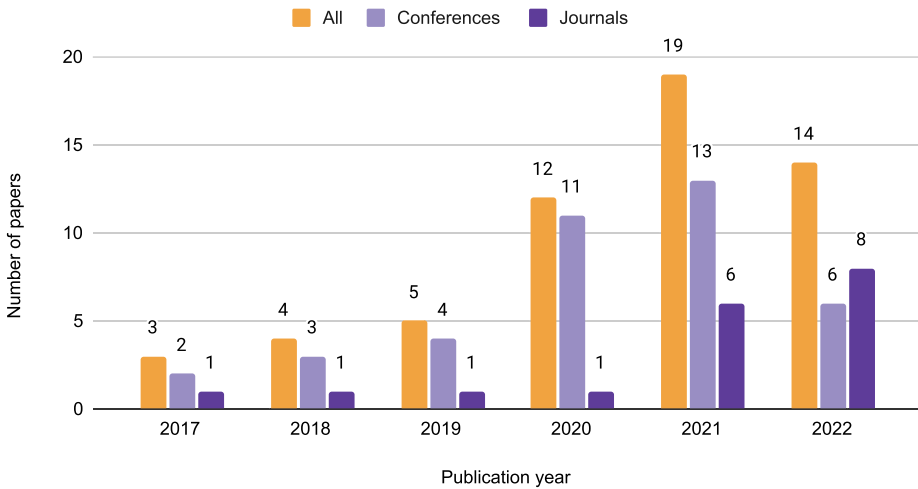


Fig. 3. Number of papers per year.

published at conferences (39 papers) compared to 18 papers published in journals. Further, from Figure 2, we see that there is a clear concentration across conference venues (RecSys and SIGIR), whereas papers on evaluation are particularly scattered across journal venues.

Concerning the temporal evolution of evaluation papers, we observe an increasing number of papers on the evaluation of recommender systems in the analyzed time frame 2017–2022 (Figure 3). Starting in 2017, there were only 3 papers on the evaluation of recommender systems published, while this number peaked in 2021 with 19 papers on that topic. While there is a continuous upward trend of papers on that topic in conference venues, there is a sharp increase of papers on that topic in journal venues (only one journal paper in the years 2017–2020, respectively; then 6 and 8 journal papers in 2021 and 2022, respectively). We note that two of the journal papers published in 2021 (Ferrari Dacrema et al. [40] and Mena-Maldonado et al. [71]) are extended versions of previously published conference papers (Ferrari Dacrema et al. [41] from 2019 and Mena-Maldonado et al. [70] 2020, respectively). Further, the increase of journal papers on evaluation in the years 2020 and 2021 aligns with the COVID-19 pandemic, during which all conferences were either canceled or held online; which points to having led researchers to focus on journal submissions instead of conferences.

3.2 Type of Contribution

This section provides a detailed overview of the types of papers included in the literature review. The types as specified in Table 2 (i.e., benchmark, framework, metrics, model, and survey) were inferred according to the description in Section 2.3.

Figure 4 provides an overview of the number of papers per type of contribution in our sample. Most of the papers in our sample contribute to models (19); these papers provide a conceptual and empirical basis for improved recommendation or evaluation models. Considerably fewer papers (13) investigate metrics. Nine papers provide a survey, another 9 papers provide benchmarks of various approaches, and 7 papers propose frameworks.

Among the model papers, the majority focus on evaluation models, specifically on issues related to off-policy learning [23, 31, 44, 58, 69, 73, 82, 95], which helps to obtain unbiased estimates for improved offline evaluation [55]. Cañamares and Castells [20] propose a probabilistic reformulation of memory-based collaborative filtering. While the core contribution of

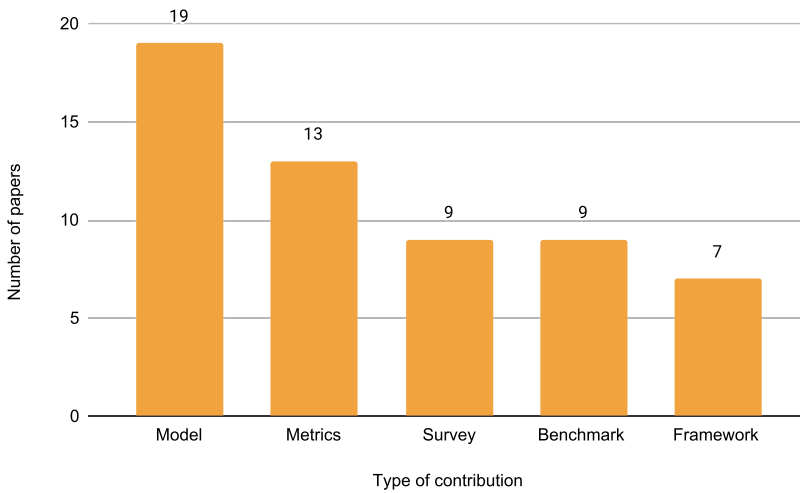


Fig. 4. Number of papers per type of contribution.

that work is a recommendation model, it also contributes to evaluation, because the experiments demonstrate that performance measurements may heavily depend on statistical properties of the input data, which the authors discuss in detail. With a probabilistic analysis, Cañamares and Castells [21] address the question of whether popularity is an effective or misleading signal in recommendation. Their work illustrates the contradictions between the accuracy that would be measured in common biased offline experimental settings and the measured with unbiased observations. Cañamares and Castells [22] demonstrate the importance of item sampling in offline experiments. Based on a thorough literature review, Carraro and Bridge [23] propose a new sampling approach to debiasing offline experiments. A second line of model papers considers user-related aspects as an important ingredient of recommender systems. For example, Frummerman et al. [42] investigate the meaning of “rejected” recommendations in a more fine-grained manner. Symeonidis et al. [91] consider short-term intentions to inform models. Jin et al. [54] rely on a psychometric modeling method to study the key qualities of conversational recommender systems. In a large-scale user study, Chen et al. [25] investigate how serendipity improves user satisfaction with recommendations; their results inform the modeling for recommendations. Ostendorff et al. [75] study users’ preferences for link-based versus text-based recommendations using qualitative evaluation methods. Lu et al. [65] investigate whether and how annotations made by external assessors (thus, not the recommender system’s users) are a viable source for preference labeling. Guo et al. [47] study order effects in recommendation sequences, which has implications for the design of recommender systems. Said and Bellogín [80] evaluate and model inconsistencies in user rating behavior to improve the performance of recommendation methods. These papers considering user-related aspects have in common that each work primarily studies phenomena to improve recommendation models and the discussion of the results also contributes to methodological issues regarding the evaluation of recommender systems.

Among papers focusing on metrics, one set of papers compares metrics (e.g., References [70, 71, 77]), whereas some papers focus their analysis on a specific type of metrics; for instance, sampling metrics (e.g., References [60, 63]) and folding metrics (e.g., Reference [94]). In a similar spirit, Bellogín et al. [17] study biases in information retrieval metrics. Another line of metrics papers aims for harmonization of metrics (e.g., References [2, 76]) or metric improvements (e.g., Reference [64]). Balog and Radlinski [11] propose how to measure the quality of explanations in

recommender systems. Saraswat et al. [84] propose combining both performance and user satisfaction metrics in offline evaluation, leading to improved correlation with desired business metrics. Finally, Diaz and Ferraro [37] makes a metrics analysis and discussion leading into the proposal of an altogether metric-free evaluation method.

Papers discussing infrastructural aspects of recommender systems can be categorized into two types of framework papers: Those that contribute with a recommendation toolkit and those proposing a conceptual framework. The presented toolkits are iRec [87], ELLIOT [10], LensKit [39], and librec-auto [88].⁶ The framework by Bellogín and Said [19] provides guidelines for reproducibility; their paper also provides an in-depth analysis to support their guidelines. Eftimov et al. [38] propose a general framework that fuses different evaluation measures and aims at helping users to rank systems. Considering users' expectations and perceptions, Belavadi et al. [16] study the relationships between several user evaluation criteria.

Several papers provide an extensive critical evaluation across a (wide) set of approaches (Table 3). Dallmann et al. [35] study sampling strategies for sequential item recommendation. They compare four methods across five datasets and find that both sampling strategies—uniform random sampling and sampling by popularity—can produce inconsistent rankings compared with the full ranking of the models. Ferrari Dacrema et al. [41] and its extended version Ferrari Dacrema et al. [40] perform a reproducibility study, critically analyzing the performance of 12 neural recommendation approaches in comparison to well-tuned, established, non-neural baseline methods. Their work identifies several methodological issues and finds that 11 of the 12 analyzed approaches are outperformed by far simpler, yet well-tuned, methods (e.g., nearest-neighbor or content-based approaches). In a similar vein, Latifi and Jannach [61] perform a reproducibility study where they benchmark Graph Neural Networks (GNN) against an effective session-based nearest neighbor method. Also, this work finds that the conceptually simpler method outperforms the GNN-based method. Anelli et al. [9] perform a reproducibility study, systematically comparing 10 collaborative filtering algorithms (including approaches based on nearest-neighbors, matrix factorization, linear models, and techniques based on deep learning). Different to Ferrari Dacrema et al. [40, 41], Anelli et al. [9] benchmark all algorithms using the very same datasets (MovieLens-1M [48], Amazon Digital Music [74], and epinions [92]) and the identical evaluation protocol. Based on their study on modest-sized datasets, they conclude—similarly to other works—that the latest models are often not the best-performing ones. Kouki et al. [59] compare 14 models (8 baseline and 6 deep learning) for session-based recommendations using 8 different popular evaluation metrics. After an offline evaluation, they selected the 5 algorithms that performed the best and ran a second round of evaluation using human experts (user study). Reference [90] provides benchmarks across several datasets, recommendation approaches, and metrics; beyond that, this work introduces the toolkit daisyRec. Zhu et al. [99] compare 24 models for click-through rate (CTR) prediction on multiple dataset settings. Their evaluation framework for CTR (including the benchmarking tools, evaluation protocols, and experimental settings) is publicly available. Latifi et al. [62] focus on sequential recommendation problems, for which they compare the Transformer-based BERT4Rec method [89] to nearest-neighbor methods, showing that the nearest-neighbor methods achieve comparable performance to BERT4Rec for the smaller datasets, whereas BERT4Rec outperforms the simple methods when the datasets are larger.

Table 4 provides an overview of survey papers on the evaluation of recommender systems. Some of the papers provide an extensive critical evaluation across a (wide) set of datasets and approaches on a specialized topic (e.g., References [26, 40, 41, 59, 61]). Others provide a (systematic) review of

⁶Note that the work by Sun et al. [90]—besides providing benchmarks across several datasets, recommendation approaches, and metrics—also proposes the toolkit daisyRec.

Table 3. Benchmark Papers

Papers	Details
Anelli et al. [9]	Reproducibility study. An in-depth, systematic, and reproducible comparison of 10 collaborative filtering algorithms (including approaches based on nearest-neighbors, matrix factorization, linear models, and techniques based on deep learning) using three datasets and the identical evaluation protocol. Provide a guide for future research with respect to baselines and systematic evaluation.
Dallmann et al. [35]	Study sampling strategies for sequential item recommendation. Compare four methods across five datasets and find that both, uniform random sampling and sampling by popularity, can produce inconsistent rankings compared with the full ranking of the models.
Ferrari Dacrema et al. [40, 41]	Reproducibility study. Critical analysis of the performance of 12 neural recommendation approaches with reproducible setup. Comparison against well-tuned, established, non-neural baseline methods. Identification of several methodological issues, including choice of baselines, propagation of weak baselines, and a lack of proper tuning of baselines.
Kouki et al. [59]	Compare 14 models (8 baseline and 6 deep learning) for session-based recommendations using 8 different popular evaluation metrics.
Latifi and Jannach [61]	Reproducibility study. Benchmark Graph Neural Networks against an effective session-based nearest neighbor method. The conceptually simpler method outperforms the GNN-based method both in terms of Hit Ratio and the MRR.
Latifi et al. [62]	Compare the Transformer-based BERT4Rec method [89] to nearest-neighbor methods for sequential recommendation problems across four datasets using exact and sampled metrics. The nearest-neighbor methods achieve comparable or better performance than BERT4Rec for the smaller datasets, whereas BERT4Rec outperforms the simple methods for the larger ones.
Sun et al. [90]	Benchmarks across several datasets, recommendation approaches, and metrics; in addition, it introduces the toolkit daisyRec.
Zhu et al. [99]	Open benchmarking for click-through rate prediction with a rigorous comparison of 24 existing models on multiple dataset settings in a reproducible manner. The evaluation framework for CTR (including the benchmarking tools, evaluation protocols, and experimental settings) are publicly available.

the literature landscape on a specialized topic (e.g., References [4–6, 36, 52, 83, 98]). The framework by Zangerle and Bauer [96] is based on a survey of previous literature on the respective topic. Similarly, Zhao et al. [98] starts with a survey of literature on aspects related to offline evaluation for top- N recommendation, which builds the basis for their systematic comparison of a selected set of 12 algorithms across eight datasets.

3.3 Experiment Types

While many types of experiments can be performed, the results presented in this section rely on the established definitions of online, offline, and user study, respectively.

As shown in Figure 5, the vast majority of the papers (38) use offline experiments. Considerably fewer papers (12) report user studies. Comparably few (6) report on online experiments. Ten papers do not report any evaluation, these are mainly survey papers [4–6, 36, 52, 83], papers on metrics [6, 60, 94], and one paper contributing with a framework [88].

While most papers (39) employ one experiment type, there are seven papers that combine two types, and one paper [59] combining all three types (Table 5). Interestingly, all papers using an

Table 4. Survey Papers on the Evaluation of Recommender Systems

Papers	Details
Al Jurdi et al. [4]	Classification of natural noise management (NNM) techniques and analysis of their strengths and weaknesses. Comparative statistical analysis of the NNM mechanisms.
Alhijawi et al. [5]	Specifically address the objectives: relevance, diversity, novelty, coverage, and serendipity. Reviews the definitions and measures associated with these objectives. Classifies over 100 articles (published from 2015 to 2020) regarding objective-oriented evaluation measures and methodologies. Collect 43 objective-oriented evaluation measures.
Ali et al. [6]	Survey on the evaluation of scholarly recommender systems. Analysis suggests that there is a focus on offline experiments, whereby either simple/trivial baselines are used or no baselines at all.
Chin et al. [26]	Compare 45 datasets used for implicit feedback based top- k recommendation based on characteristics (similarities and differences) and usage patterns across papers. For 15 datasets, they evaluate and compare the performance of five different recommendation algorithms.
Dehghani Champiri et al. [36]	Focus on context-aware scholarly recommender systems. Classification evaluation methods and metrics on usage.
Jannach [52]	Provide an overview of evaluation aspects as reported in 127 papers on conversational recommender systems. Argue for a mixed methods approach, combining objective (computational) and subjective (perception-oriented) techniques for the evaluation of conversational recommenders, because these are complex multi-component applications, consisting of multiple machine learning models and a natural language user interface.
Sánchez and Bellogín [83]	Focus on point-of-interest recommender systems. Systematic review covering 10 years of research on that topic, categorizing the algorithms and evaluation methodologies used. The common problems are that both, the algorithms and the used datasets (statistics), are described in insufficient detail.
Zangerle and Bauer [96]	Introduce “Framework for EValuating Recommender systems,” derived from the discourse on recommender systems evaluation. Categorization of the evaluation space of recommender systems evaluation. Emphasis on the required multi-facetedness of a comprehensive evaluation of a recommender system.
Zhao et al. [98]	Survey of 93 offline evaluation for top- N recommendation algorithms. Provide an overview of aspects related to evaluation metrics, dataset construction, and model optimization. In addition, this work presents a systematic comparison of 12 top- N recommendation algorithms (covering both traditional and neural-based algorithms) across eight datasets.

online experiment, combine it with another experiment type; four papers using an online experiment [44, 73, 77, 91], also carry out offline experiments, one combines online experiments with user studies [16], and one paper combines all three experiment types [59]. Further two papers [42, 80] use offline experiments and user studies.

3.4 Datasets

Table 6 provides an overview of the datasets used in the papers. In total, our analysis contains 80 datasets. We distinguish between papers that use pre-collected, established datasets (65 datasets) and papers that propose a custom dataset (15 datasets, see the last row of Table 6). In a graphical overview, Figure 6 presents the number of papers relying on each dataset. Note that in this chart, we have aggregated different versions of a dataset into a single dataset category (for instance, we

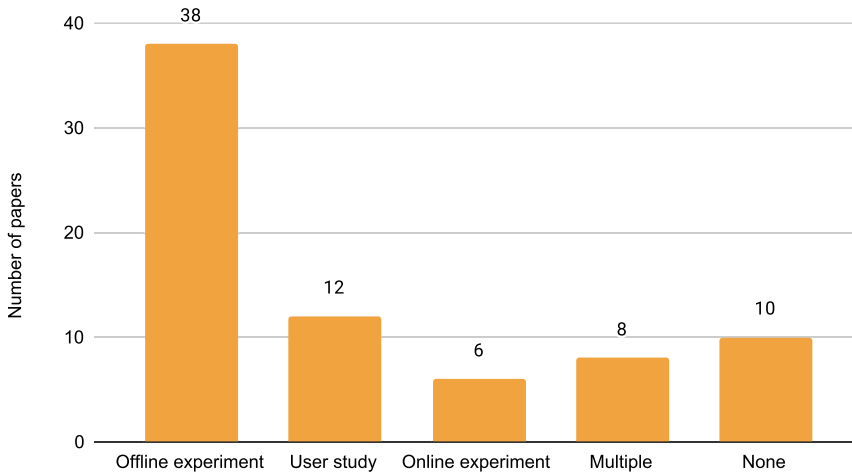


Fig. 5. Number of papers per experiment type.

Table 5. Papers Using More Than One Experiment Type

Papers	Online experiment	Offline experiment	User study
Gilotte et al. [44]	x	x	
Narita et al. [73]	x	x	
Peska and Vojtas [77]	x	x	
Symeonidis et al. [91]	x	x	
Frumerman et al. [42]		x	x
Said and Bellogin [80]		x	x
Belavadi et al. [16]	x		x
Kouki et al. [59]	x	x	x

combined the widely used MovieLens datasets MovieLens 100k, 1M, 10M, 20M, 25M, Latest, and HetRec).

Table 6 and Figure 6 show that the dataset usage distribution for established (pre-collected) datasets is dominated by the MovieLens datasets. MovieLens datasets are used 32 times in the papers investigated, with MovieLens 1M being the most popular dataset (19 usages). Furthermore, the Amazon review datasets are used in 24 papers, followed by the LastFM dataset, appearing in the evaluation of 9 papers. We also observe that 43 and hence, 66.15% of the listed datasets are only used in a single paper. Further 8 datasets are used in 2 of the papers in our study and another 14 datasets are employed in three or more papers.

Generally, the majority of papers relied on existing, pre-collected datasets: Of 146 dataset usages, 15 were custom datasets. These findings are in line with a previous analysis of datasets being used for recommender systems evaluation [13], with a focus on the use of data pruning methods for the years 2017 and 2018. Generally, the high number of datasets employed at a low rate makes a direct comparison of recommendation approaches hardly possible. Particularly, given the vastly different characteristics of these. In contrast, we also observe that established datasets like the MovieLens dataset family, are used frequently, allowing for a better comparison of approaches.

A further aspect to consider regarding the comparability of approaches is dataset pre-processing. Typical pre-processing steps include removing users, items, or sessions with a low number of

Table 6. Overview of Datasets Used in Surveyed Papers

Datasets	Papers	# Papers
Amazon Beauty [74]	[26, 35, 62]	3
Amazon Book [74]	[95]	1
Amazon Digital Music [74]	[9, 26]	2
Amazon Electronics [74]	[26, 90, 98]	3
Amazon Home & Kitchen [74]	[64]	1
Amazon Instant Video [74]	[41]	1
Amazon Kindle Store [74]	[87]	1
Amazon Movies & TV [74]	[26, 40, 98]	3
Amazon Musical Instruments [74]	[26, 40]	2
Amazon Patio, Lawn & Garden [74]	[26]	1
Amazon Sports & Outdoors [74]	[64]	1
Amazon Toys & Games [74]	[26, 98]	2
Amazon Video Games [74]	[26, 35, 98]	3
Avazu ⁷	[99]	1
BeerAdvocate [68]	[37]	1
Book crossing [100]	[90]	1
citeulike-a [93]	[40, 41, 95]	3
citeULike-t [93]	[26, 40, 64]	3
Clothing Fit [72]	[87]	1
CM100k [21]	[70, 71]	2
CoatShopping [86]	[23]	1
Criteo ⁸	[99]	1
epinions [92]	[9, 40, 64, 90]	4
Filmtrust [46]	[40]	1
Flixster ⁹	[26]	1
Frappe [12]	[40]	1
Good Books ¹⁰	[87]	1
Good Reads ¹¹	[87]	1
Gowalla [27]	[40, 61]	2
LastFM [24]	[17, 19, 20, 26, 40, 61, 87, 90, 98]	9
Library-Thing [97]	[37]	1
Million Playlist Dataset ¹²	[38]	1
Million Post Corpus [85]	[16]	1
MovieLens 100k [48]	[26, 40, 41]	3
MovieLens 1M [48]	[9, 10, 17, 19, 20, 22, 35, 37, 40, 41, 60, 62, 63, 70, 71, 80, 87, 90, 98]	19
MovieLens 10M [48]	[26, 94]	2
MovieLens 20M [48]	[26, 35, 40, 62, 76]	5
MovieLens 25M [48]	[84]	1
MovieLens Latest [48]	[65]	1
MovieLens HetRec ¹³	[40]	1
MoviePilot ¹⁴	[80]	1

(Continued)

⁷<https://www.kaggle.com/c/avazu-ctr-prediction>⁸<https://www.kaggle.com/c/criteo-display-ad-challenge>⁹<https://sites.google.com/view/mohsenjamali/home>¹⁰<https://github.com/zygmuntz/goodbooks-10k>¹¹<https://www.kaggle.com/datasets/jealousleopard/goodreadsbooks>¹²<https://research.atspotify.com/datasets/>¹³<https://grouplens.org/datasets/hetrec-2011/>¹⁴<http://www.moviepilot.de/>

Table 6. Continued

Datasets	Papers	# Papers
NetflixPrize ¹⁵	[20, 40, 41, 87, 98]	5
Open Bandit [81]	[82]	1
Pinterest [43]	[40, 41]	2
Steam [56]	[35, 62]	2
Ta Feng Grocery Dataset ¹⁶	[40]	1
Tradesy [49]	[95]	1
TREC Common Core 2017 [7] ¹⁷	[37]	1
TREC Common Core 2018 ¹⁸	[37]	1
TREC Deep Learning Document Ranking 2019 [32]	[37]	1
TREC Deep Learning Document Ranking 2020 [32]	[37]	1
TREC Deep Learning Passage Ranking 2019 [32]	[37]	1
TREC Deep Learning Passage Ranking 2020 [33]	[37]	1
TREC Robust 2004 ¹⁹	[37]	1
TREC Web 2009 [28]	[37]	1
TREC Web 2010 ²⁰	[37]	1
TREC Web 2011 ²¹	[37]	1
TREC Web 2012 [29]	[37]	1
TREC Web 2013 ²²	[37]	1
TREC Web 2014 [30]	[37]	1
Webscope R3 [67]	[23]	1
Yelp ²³	[19, 40, 80, 90, 98]	5
Yahoo R3 (Music) ²⁴	[22, 70, 71, 87]	4
Yahoo R4 ²⁵	[26]	1
Xing [1]	[42]	1
Custom	[2, 11, 21, 25, 31, 44, 47, 54, 58, 59, 69, 73, 75, 77, 91]	15

interactions or converting explicit ratings to binary relevance values. As Ferrari Dacrema et al. [40] note in their survey on the reproducibility of deep learning recommendation approaches, it is important that all pre-processing steps are clearly stated in the paper and that the removal of data is justified and motivated. Also, pre-processing should be included in the code published. Inspecting the papers of our survey, we find that eight papers mention that they convert explicit rating data to a binary relevance score or song play counts to explicit ratings [17, 23, 26, 37, 38, 62, 64, 90]. Furthermore, users, items or sessions with fewer and/or more interactions than a given threshold are removed in 12 papers [9, 22, 26, 35, 42, 61, 62, 64, 77, 90, 91, 98]. Zhao et al. [98] refer to this pre-processing step as *n*-core filtering. They perform a study on three aspects in the context of evaluating recommender systems: evaluation metrics, dataset construction, and model optimization. For dataset construction, they find that 44% of the papers in their study do not provide any information about pre-processing, and 34% of the papers apply *n*-core filtering with *n* set to 5 or 10.

¹⁵<https://www.kaggle.com/datasets/netflix-inc/netflix-prize-data>

¹⁶<https://www.kaggle.com/datasets/chiranjivdas09/ta-feng-grocery-dataset>

¹⁷<https://github.com/trec-core/2017>

¹⁸<https://github.com/trec-core/2018>

¹⁹https://trec.nist.gov/data/t13_robust.html

²⁰<https://trec.nist.gov/data/web10.html>

²¹<https://trec.nist.gov/data/web2011.html>

²²<https://github.com/trec-web/trec-web-2013>

²³<https://www.kaggle.com/datasets/yelp-dataset/yelp-dataset>

²⁴<https://webscope.sandbox.yahoo.com/catalog.php?datatype=r&did=3>

²⁵<https://webscope.sandbox.yahoo.com/catalog.php?datatype=r&did=4>

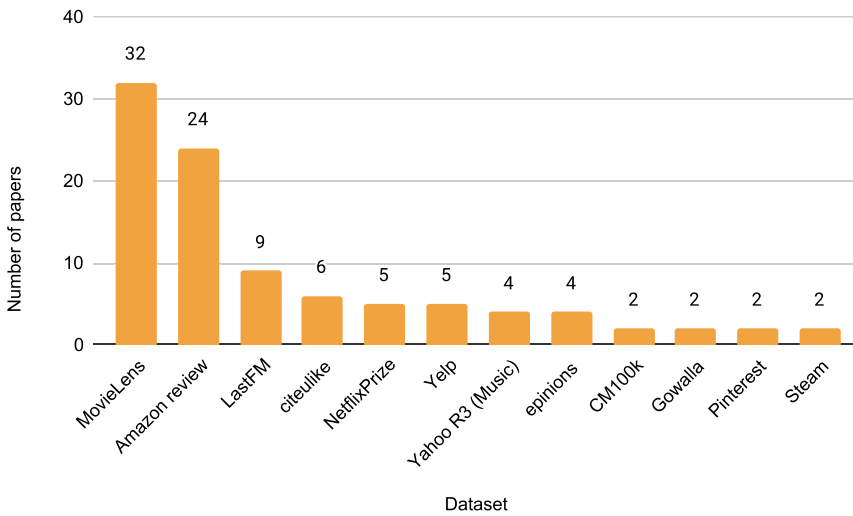


Fig. 6. Overview of datasets used in at least two papers, where different versions of a dataset are aggregated into a single dataset category for the Amazon review, MovieLens, and citeulike datasets.

Sun et al. [90] also study the impact of different thresholds for filtering users and items. Here it is important to note that, for instance, the MovieLens datasets are already pre-processed to some extent as they only include users with more than or equal to 20 interactions.

In the following, we focus our analysis on datasets that have been used at least three times in the surveyed papers. Table 7 provides an overview of these 12 datasets, where we list the domain, the feedback type (hence, whether the dataset features explicit or implicit data; in the case of explicit ratings, we also add the rating scale), the size of the dataset captured by the number of interactions, and the type of side information contained. Notably, 5 of the 12 most popular datasets stem from the movie or music domain. In terms of the type of ratings contained, the citeulike and LastFM datasets provide implicit feedback (0 or 1), while the other datasets provide explicit ratings on a scale from 0 (or 1) to 5 stars. Interestingly, when inspecting the size of the datasets, the most popular datasets appear to be relatively small, with the most popular dataset (MovieLens 1M) holding 1,000,000 interactions.

Another interesting aspect when investigating the choice of datasets for the evaluation of recommender systems is the number of different datasets used by individual papers. Evaluating a recommender system on diverse datasets is critical to gaining insights into the generalizability and robustness of the recommender system proposed. When inspecting the number of different datasets used in the experiments, we find that 26 papers (45.61% of all papers contained in the study) rely on a single dataset, 5 papers (8.77%) rely on two datasets, 7 papers (12.28%) use 3 datasets and another 10 papers (17.54%) use four or more datasets. Of these, 3 papers used more than 10 different datasets: In extensive experiments, Ferrari Dacrema et al. [41] benchmark deep learning-based recommender systems against a set of relatively simple baselines. Diaz and Ferraro [37] showcase a metric-free evaluation method for recommendation and retrieval based on a set of 16 datasets. Chin et al. [26] conduct an empirical study on the impact of datasets on the evaluation outcome and resulting conclusions. Their study shows a different distribution of dataset popularity among recommender systems evaluation than we observe in the analysis at hand. However, we conjecture that this is due to the diverse inclusion criteria of the studies. For instance, Chin et al.'s study is restricted to implicit feedback top- k recommendation tasks. Notably, our analysis also contains

Table 7. Details of Datasets Used in at Least Three Papers

Datasets	Domains	Feedback	Interactions	Side Information
Amazon Electronics, Products, Video Games [74]	Products	[1,5]	20,994,353 (E), 371,345 (B), 2,565,349 (V)	product information (e.g., description, color, product images, technical details), timestamp
citeulike-a, citeulike-t	Scientific Papers	{0,1}	204,987 (a), 134,860 (t)	tags, bag-of-words, and raw text for each article, citations between articles
epinions [92]	Products	[0,5]	922,267	explicit trust relationships among users, timestamps
LastFM [24]	Music	{0,1}	19,150,868	artist, song name, timestamp
MovieLens (100k, 1M, 20M) [48]	Movies	[0,5]	100,000 (100k)– 20,000,000 (20M)	movie metadata (e.g., title, genre), user metadata (e.g., age, gender), rating timestamp
NetflixPrize ²⁶	Movies	[1,5]	100,000,000	movie metadata (title, release year), rating date
Yelp ²⁷	Business	[0,5]	6,990,280	business metadata (address, category, etc.), user metadata (user name, user stats (no. of reviews, user votes, etc.)), rating timestamp

We list the domain of the dataset, the type of feedback, number of interactions contained, and side information provided.

9 papers (15.79%) that did not use any dataset. The reason here is that most of these papers are surveys [4–6, 36, 52, 83, 96]. Furthermore, Ekstrand [39] describes the Python LensKit software framework and Sonboli et al. [88] describe the librec-auto toolkit.

Our analysis contains 13 versions of the Amazon review datasets, seven different versions (or subsets) of the MovieLens dataset, and two versions of the citeulike dataset. Considering the usage of different versions of the same dataset, we find that five papers use different versions of the same aggregated dataset. In their survey on dataset usage, Chin et al. [26] use eight versions of the Amazon reviews dataset and three versions of the MovieLens dataset (of a total of 15 individual datasets used). In their reproducibility study, Ferrari Dacrema et al. [40] used four versions of the MovieLens datasets, both versions of the citeulike datasets, and two versions of the Amazon reviews dataset (of 17 individual datasets used). In their prior reproducibility study, Ferrari Dacrema et al. [41] used two versions of the MovieLens dataset.

We further investigate which datasets are jointly used in evaluations. For this analysis, analyze the sets of datasets co-used in the papers (note that the co-usage of individual datasets is already presented in Table 6). We employed a frequent itemset approach (i.e., the Apriori algorithm [3]) and present the results in Table 8. This table shows the set of datasets employed together and the number of papers that co-use these datasets. The most frequently combined datasets are LastFM

²⁶<https://www.kaggle.com/datasets/netflix-inc/netflix-prize-data>

²⁷<https://www.kaggle.com/datasets/yelp-dataset/yelp-dataset>

Table 8. Combinations of Datasets (Pairs and Triples) Frequently Co-occurring in Experiments

Dataset Combinations	# Papers
{LastFM, ML 1M}	7
{ML 1M, NetflixPrize}, {ML 1M, Yelp}	5
{ML 1M, Yahoo R3}, {LastFM, Yelp}, {LastFM, NetflixPrize}, {LastFM, ML 1M, NetflixPrize}, {LastFM, ML 1M, Yelp}	4
{Amazon Movies & TV, LastFM}, {Amazon Electronics, LastFM}, {Amazon Beauty, ML 20M}, {epinions, ML 1M}, {ML 100k, ML 20M}, {ML 100k, ML 1M}, {ML 1M, ML20M}	3

We list all sets of datasets that co-occur in at least three papers (ML = MovieLens).

and MovieLens 1M (appearing in seven papers). The MovieLens 1M dataset appears in pairs with the NetflixPrize and the Yelp datasets in five papers. In the list of sets of datasets that appear in four papers, we find not only pairs but also triples of datasets that are jointly used for evaluation in three papers. Unsurprisingly, the MovieLens datasets and other popular datasets are dominant. This aspect has also been raised by Chin et al. [26] and our results are in line with these previous findings.

Inspecting the papers that use custom datasets, we observe that the majority of these papers feature (or create) a custom dataset for three distinctive reasons. One reason is user surveys [2, 25] and user studies being conducted [11, 47, 54, 75], where the result of the user study itself is presented as a novel dataset. For instance, Chen et al. [25] perform a user study to get a deeper understanding of the impact of serendipity on user satisfaction on a popular mobile e-commerce platform in China. A further reason for using custom datasets is the recent trend toward counterfactual (off-policy) learning, which requires an unbiased, missing-at-random dataset [22, 31, 44, 58, 73]. Furthermore, several papers perform evaluations based on proprietary data provided by a private sector business entity [44, 59, 69, 73, 77, 91].

3.5 Metrics

The reviewed literature features an extensive range of datasets, as depicted in Section 3.4. This variety is also mirrored in the selection of evaluation metrics. We divide the metrics into two categories: conventional metrics widely utilized in the field and specific metrics proposed for the unique problem addressed within a certain paper. We refer to these as custom metrics (see the final row of Table 9). A visual representation of the most frequently used metrics—those employed in at least two papers within our surveyed literature—is provided in Figure 7.

Traditionally, recommender systems research has relied on a standard set of metrics, including Precision, Recall, and normalized Discounted Cumulative Gain (nDCG) [18, 45]. These metrics have gained significant popularity in the examined literature. However, our analysis also uncovers the existence of a diverse array of less prevalent metrics, as illustrated in Table 9. In essence, a selected group of metrics is featured prominently: Precision is employed in 22 of the 57 reviewed papers (approximately 36%), nDCG in 20 papers (around 35%), and Recall in 17 papers (nearly 30%). These findings resonate with the notion that ranking and relevance metrics align more closely with actual user preferences than a minimized rating prediction error does [34, 45]. Yet, metrics associated with rating prediction, such as RMSE, MAE, and MSE, still figure prominently in a considerable portion of the reviewed literature, appearing in a total of 7 papers (about 12%). While a vast majority of papers do not employ rating prediction metrics, the fact that more than 1 in 10 papers uses them contradicts the general consensus in the recommender systems research field, which holds that rating prediction is an inadequate surrogate for actual user preference [8].

Table 9. Overview of the Metrics Used in Surveyed Papers

Metrics	Abbr.	Papers	#
Area Under Curve	AUC	[25, 35, 38, 60, 77]	5
Average Coverage of Long Tail	ACLT	[9]	1
Average Percentage of Long Tail	APLT	[9]	1
Average Precision	AP	[37, 60, 64, 95]	4
Average Recommendation Popularity	ARP	[9]	1
Binary Preference-based measure	bpref	[17]	1
Clickthrough rate	CTR	[77, 84, 91, 99]	4
Conversion rate	CVR	[31]	1
Coverage (item)	Coverage	[38, 59, 98]	3
Coverage (user)		[87]	1
Discounted Cumulative Gain	DCG	[95]	1
Expected Free Discovery	EFD	[9]	1
Expected Popularity Complement	EPC	[9, 87]	2
Expected Profile Distance	EPD	[87]	1
F-measure	F1	[9]	1
Fallout		[71]	1
Gini		[9, 87]	2
Hit Rate	HR	[35, 38, 40, 59, 61, 62, 90, 98]	8
Hits		[87]	1
Intra-list Diversity	ILD	[87]	1
Inferred Average Precision	InfAP	[17]	1
Item Coverage	IC	[9]	1
Jaccard coefficient		[65]	1
Logistic Loss	Logloss	[99]	1
Mean Absolute Error	MAE	[95]	1
Mean Average Precision	MAP	[9, 23, 37, 40, 59, 77, 90, 98]	8
Mean Reciprocal Rank	MRR	[9, 40, 59, 61, 62, 77, 90, 98]	8
Mean Squared Error	MSE	[58, 73]	2
normalized Discounted Cumulative Gain	nDCG	[9, 17, 19–23, 26, 35, 37, 40, 41, 59, 60, 62, 64, 76, 77, 90, 98]	20
Novelty		[98]	1
Overlap		[65]	1
Pearson Correlation Coefficient	PCC	[65]	1
Popularity		[59]	1
Popularity-based Ranking-based Equal Opportunity	PREO	[9]	1
Popularity-based Ranking-based Statistical Parity	PRSP	[9]	1
Precision	P	[9, 17, 19–23, 38, 40–42, 44, 59, 64, 65, 70, 71, 77, 87, 90, 91, 98]	22
Recall	R	[9, 19, 22, 23, 26, 37, 40, 41, 59, 60, 63, 65, 77, 87, 90, 95, 98]	17
Reciprocal Rank	RR	[37, 64]	2
Root Mean Squared Error	RMSE	[65, 69, 73, 80, 94]	5
Custom		[2, 11, 25, 37–39, 54, 75, 81, 94]	12
Total number of metrics: 40			

Figure 7 portrays the disparity in popularity among various metrics. Precision, nDCG, and Recall are roughly twice as favored as any of the other top metrics. These three metrics epitomize the core characteristics of recommender and information retrieval systems, notably relevance and ranking.

Furthermore, it is worth mentioning that of the total 40 metrics employed in the reviewed papers, 23 metrics (approximately 58%) are each applied in just a single paper. Some of these uniquely applied metrics are specific to individual papers that utilize an extensive range of metrics. For example, Silva et al. [87] introduce metrics such as user-coverage, EPC, EPD, Gini, and Hits, while Anelli et al. [9] introduce various non-accuracy metrics like Average Coverage of Long Tail, Average

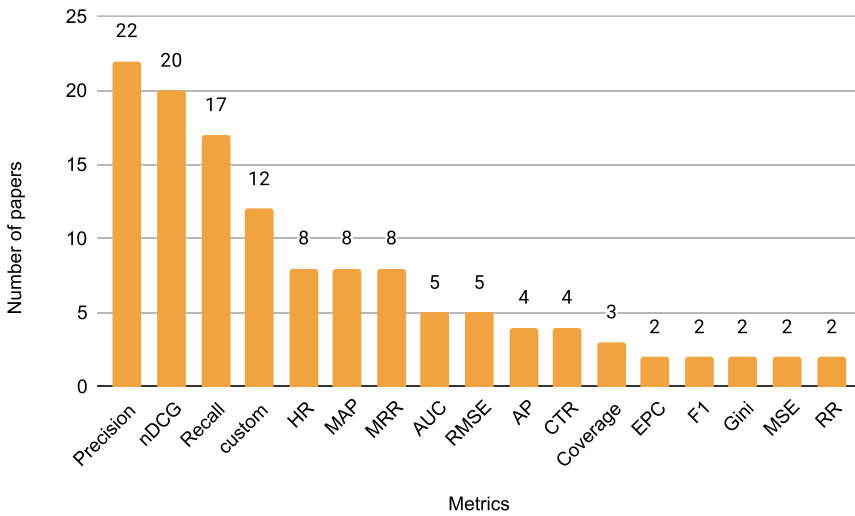


Fig. 7. Overview of metrics used in at least two papers (NB: Coverage refers to item coverage).

Table 10. The Categories of Value the Metrics Express

Categories	Metrics
Relevance	AP, AUC, F1, fallout, Hits, HR, InfAP, Logloss, MAP, P, R
Success Rate	CTR, CVR
Rating Prediction Accuracy	bpref, MAE, MSE, RMSE
Ranking	DCG, nDCG, MRR, RR
Non-accuracy	ACLT, APLT, Coverage, EFD, EPC, EPD, Gini, IC, ILD, Jaccard, Overlap, PCC, Popularity, PREO, PRSP

Percentage of Long Tail, Expected Free Discovery, and Popularity-based Ranking-based Equal Opportunity, among others. Moreover, five metrics appear in only two papers each, and a single metric is utilized in three papers. The variation in metric usage complicates the comparison and benchmarking across different papers, as emphasized in the discussion on dataset usage (see Section 3.4).

Similarly, we scrutinize the number of metrics utilized per paper. It is crucial to emphasize that the quantity of metrics employed does not necessarily reflect the quality or completeness of a paper or recommender system. Nonetheless, the use of multiple metrics can yield insights into different facets of a system. When analyzing our data, we discover that 18 papers (32%) use only a single metric, and surprisingly, and 10 papers (18%) do not use any metrics whatsoever. Although the majority of papers that abstain from using metrics are categorized as literature reviews (refer to Table 4), there are exceptions. Furthermore, 9 papers (16%) apply two metrics, while 5 papers (9%) employ three metrics. In total, 42 papers (74%) utilize three or fewer metrics. With this understanding, we now probe into the variety of metrics. In Table 10, we present a classification of evaluation metrics into overarching categories that correspond to specific recommendation tasks, like ranking, rating prediction, and relevance. Despite the absence of a universally accepted classification of metrics in the recommender systems research field, our categorization resonates with the general application scenarios of recommendations and the desired attributes of a recommender system.

Table 11. Combinations of Metrics Used Frequently in the Surveyed Papers

Metric combinations	# Papers
{nDCG, P} [*]	14
{nDCG, R} [*]	13
{P, R}	12
{nDCG, P, R} [*]	10
{nDCG, MAP} [*] , {R, MAP}, {nDCG, R, MAP} [*]	8
{nDCG, P, MAP} [*] , {P, MAP}, {nDCG, P, R, MAP} [*] , {nDCG, MRR}, {P, R, MAP}	7
{nDCG, MAP, MRR, R} [*] , {MRR, P, MAP, R} [*] , {nDCG, MRR, MAP} [*] , {MRR, MAP, R} [*] , {MRR, P, MAP} [*] , {nDCG, P, MRR, MAP} [*] , {MRR, P} [*] , {MRR, R} [*] , {MRR, MAP} [*] , {nDCG, HR} [*] , {nDCG, P, MRR, R} [*] , {MRR, HR} [*] , {MRR, P, R} [*] , {nDCG, P, MRR} [*] , {nDCG, MRR, R} [*] , {nDCG, MAP, MRR, P, R} [*]	6
{P, HR}, {nDCG, HR, MRR} [*]	5
{nDCG, P, HR, MAP} [*] , {P, HR, R, MAP}, {nDCG, HR, R} [*] , {nDCG, HR, R, MAP} [*] , {MRR, P, HR, R} [*] , {nDCG, P, HR, MRR} [*] , {nDCG, HR, MRR, R} [*] , {nDCG, P, HR, R} [*] , {MRR, MAP, HR, R} [*] , {MAP, MRR, P, HR, R} [*] , {nDCG, MRR, P, HR, MAP} [*] , {nDCG, MAP, MRR, HR, R} [*] , {nDCG, MAP, P, HR, R} [*] , {nDCG, MRR, P, HR, R} [*] , {nDCG, HR, MRR, MAP} [*] , {nDCG, MRR, P, HR, R, MAP} [*] , {MRR, P, HR, MAP} [*] , {nDCG, P, HR} [*] , {P, HR, R}, {MRR, HR, R} [*] , {MRR, P, HR} [*] , {nDCG, HR, MAP} [*] , {HR, R, MAP}, {P, HR, MAP}, {MRR, HR, MAP} [*] , {HR, R}, {HR, MAP}	4
{Coverage, HR} [*] , {P, AUC}, {AUC, R}, {nDCG, AUC} [*] , {P, Coverage, HR} [*] , {P, Coverage} [*] , {nDCG, AUC, R} [*] , {nDCG, AP} [*] , {AP, R}	3

Tuples with asterisks contain metrics from at least two of the categories in Table 10, excluding custom metrics. (NB: Coverage in refers to item coverage).

In the context of metrics, it is interesting to explore the combinations of metric types, that is, the characteristics being measured in tandem. Given that recommendations apply across diverse contexts, the extensive array of metrics used mirrors the various goals pursued by recommendation applications and the stakeholders involved. By concentrating on metrics adopted in three or more papers, we examine the employed combinations in the surveyed literature (refer to Table 11). A key observation from this table is that the majority of combinations encompass ranking and relevance metrics, while combinations incorporating other metric types are less prevalent. This observation contrasts with current discussions in the recommender systems community, with the only beyond-accuracy metric appearing in the table being item coverage. This indicates that beyond-accuracy metrics are seldom used in combination with other metrics, including other beyond-accuracy metrics such as novelty, fairness, or any of the metrics in the bottom row of Table 10. A similar comment can be made regarding the utilization of success rate metrics.

Additionally, in agreement with the discourse within the recommender systems community, particularly regarding rating prediction, it is worth mentioning that no rating prediction error metrics are present in this table. This could signal a decrease in the overall usage of these metrics. Even when acknowledging that some papers use these metrics (as noted above), they do so without merging them with the more widely accepted evaluation tools and metrics.

4 DISCUSSION

With this survey paper, we aim to provide an analysis of a snapshot of research on the evaluation of recommender systems. We gain insights into the type of experiments the community performs

when researching on evaluation aspects, the data it focuses on, and the metrics that are seen as important.

First, we find that, within research on evaluation aspects of recommender systems, there is a strong focus on offline experiments, a result that is in line with what has been shown in earlier overviews of recommender systems research in general, e.g., References [6, 53]. We observe that several papers combine two types of experiments; this is seen as contributing to getting a more comprehensive picture than when using one experiment type only (see, e.g., Zangerle and Bauer [96]). However, with 8 of 57 papers that employ such a multi-method approach, the number of papers taking this approach is low.²⁸ Interestingly, when investigating the use of online experiments, we find that online experiments are predominantly combined with another experiment—typically with an offline experiment. Overall, this indicates that the landscape of research on the evaluation of recommender systems is a narrow one, with a strong focus on offline experiments, at least in published literature. As our review concentrates on research that specifically focuses on the evaluation of recommender systems, it does not allow for drawing conclusions concerning evaluation practices of the recommender systems research at large. Still, suppose that the broader landscape of recommender systems research embraces the full spectrum of experiment types (i.e., online experiments, user studies, offline experiments), then research on the evaluation of recommender systems needs to reflect the broad spectrum, too. In case the broader landscape of recommender systems research has a strong focus on offline evaluations (as, for instance, shown in Jannach [52] and Jannach and Bauer [53]), the community is encouraged to embrace the wider spectrum in their evaluation efforts. For the specialized topic of conversational recommender systems, Jannach [52] provides a good rationale for why it is essential to involve humans in the evaluation process of such systems (thus, encouraging to use user studies and online experiments). With their FEVR framework, Zangerle and Bauer [96] provide guidance concerning the multifaceted aspects that need to be considered in a comprehensive evaluation (thus, encouraging to use the full spectrum of experiment types). In the realm of research that specifically focuses on the evaluation of recommender systems, it appears worthwhile to embrace the full spectrum and possibly demonstrate how the results of different experiment types may diverge or complement each other. In this regard, we want to point to Kouki et al. [59], which is the only work covered by our survey that embraces all three experiment types.

Second, we observe a popularity gap in the use of datasets. On the one hand, the same few (and relatively old) datasets (i.e., MovieLens, Amazon review dataset) are used extensively; on the other hand, as many as 50% of the datasets (32) are used in only one single paper each. While the use of the same (or similar) datasets across multiple papers can increase comparability and benchmarking, in many cases it is disputable whether those few datasets indeed represent the best choice. First, older datasets are typically significantly smaller than newer, or current, datasets. This, in turn, raises questions regarding generalizability and applicability in the current landscape but also points to a lack of validation concerning the scalability of the evaluated recommendation models and approaches to larger datasets. Second, we have to be aware that older datasets may not be good proxies of the user behavior and preferences of today's users. As a result, good performance results with outdated datasets may not work sufficiently well in current practice. Third, with the focus on MovieLens and Amazon reviews, it is difficult to assess whether, and how, the evaluation results generalize to other domains. Yet, while the newly created datasets may better reflect these issues, these do not allow for comparison because of their one-time use. Against this background, we encourage the community to use more recent datasets and—where feasible—demonstrate generalizability by including datasets from multiple domains.

²⁸Note that 10 papers in our sample (for instance, several survey papers) do not use any experiment type.

To facilitate reproducibility, researchers are strongly encouraged to make datasets publicly available.

Third, when analyzing the employed performance metrics, we observe a similar picture as for dataset usage: Only a few metrics are widely used, i.e., Precision, nDCG, and Recall. There are a number of metrics that are, comparatively, rarely used in experiments validating the performance of recommendation approaches. Interestingly, next to Precision, nDCG, and Recall, a large number of papers (22) introduce specific custom metrics. These custom metrics make it difficult, if not impossible, to compare recommendation quality across, and even within, papers. The observation of the (still) high popularity of error metrics (used in 8 papers, 13%) goes against the general consensus in the recommender systems research field that these are poor proxies to assess recommender performance related to actual user preferences. Further, our review indicates that beyond-accuracy metrics are rarely used in research on the evaluation of recommender systems, which is not aligned with the discourse in the recommender systems field that evaluation concerning beyond-accuracy qualities are crucial. We note that our review surveys papers that focus on the evaluation of recommender systems; thus, while the consideration of beyond-accuracy metrics is also essential for papers with a focus on evaluation, this observation does not allow to draw conclusions about the use of beyond-accuracy metrics in recommender systems research practice in general. However, other surveys that cover evaluation practice in recommender systems show a similar picture: For instance, the recent review by Alhijawi et al. [5], drawing a sample from works published from 2015 to 2020, found that the main objective of all reviewed papers was to generate relevant recommendations, whereas other objectives did not get the same attention as relevance (only 21.3% of the reviewed works considered diversity, 6.1% coverage, 3.4% serendipity, and 6.1% novelty) and, in the recent survey on offline evaluation for top- N recommendation algorithms by Zhao et al. [98], only 2 of 93 papers (2.15%) used beyond-accuracy metrics. In short, the community is encouraged to use appropriate metrics and, particularly, include beyond-accuracy metrics in their evaluation efforts, as both are essential for both, research on the evaluation of recommender systems and also for research on recommender systems at large.

Our literature review comes with certain limitations. In our search strategy, we relied on the paper keywords provided by the authors. This may have caused relevant papers contributing to evaluation being excluded from our datasets, because these were not tagged with the keywords used in our query. For example, we observe that some papers do not put the evaluation of recommender systems at the core of the investigation, but—in addition—also contribute to evaluation. For instance, the core contribution of Cañamares and Castells [20] is a recommendation model. In addition, their work demonstrates that the performance measurements may heavily depend on the statistical properties of the input data sample, which is a significant contribution to evaluation and is also discussed accordingly in the paper. Other papers with a core contribution outside the evaluation field might not use the keyword “evaluation” and our query might have missed those. However, a query using only the keywords “recommender systems” or “recommendation systems” to an enormous number of papers (1,698 hits as of July 19, 2023) for the time frame 2017–2022, which was not reasonable to process manually for this review. Moreover, we note that our review provides a snapshot of research on the evaluation of recommender systems in the limited time frame of 2017–2022. Accordingly, this review does not allow for deriving conclusions about how the evaluation practices have evolved over (longer) time. Given the observations in our snapshot—namely, that offline experiments are the dominant experiment type; that long-established but small datasets are commonly used; and that novel metrics have been shown to be of little value to assess the performance of recommender systems—, we conjecture that the advancements in these regards are limited overall. A longitudinal analysis would be a worthwhile research path to follow to gain a deeper insight into the developments made in the field of recommender systems evaluation. A

further limitation is that we restricted our literature search to the ACM Digital Library. While we searched the extended collection of this library, which includes the essential conferences and journals where recommender systems research is typically published, we may have missed relevant papers published outside the typical venues, especially those outside of the general research space related to “computing.” As the recommender systems field is increasingly viewed as an interdisciplinary research field, papers may be dispersed across a much wider scale of venues.

5 CONCLUSIONS

To gain insight into recent research focused on the evaluation of recommender systems, we conducted a systematic literature review. Our analysis covered papers published from 2017 to 2022, providing a thorough understanding of the current state of research on the evaluation of recommender systems within the research and practitioner communities. Throughout our review, we identified and discussed strengths and weaknesses in the field of recommender systems evaluation research. We observed notable strengths that demonstrate the continuous evolution and refinement of evaluation practices. These strengths are exemplified by the ongoing development of metrics, experiment types, and datasets that better accommodate the diverse use cases and requirements of recommender systems.

However, our analysis also brought to light certain weaknesses that require attention and improvement. One significant weakness is the persistent focus on recommendation problems that are deemed suboptimal proxies for user preferences, such as rating prediction. Additionally, the utilization of small and outdated datasets remains a challenge that hampers the overall advancement of recommender systems. To drive further progress and development in the realm of recommender systems, it is imperative for the research community to embrace the identified strengths and move away from outdated perspectives that contribute to the weaknesses. Achieving this objective is a collaborative effort that necessitates the collective expertise and participation of the entire recommender systems research community.

REFERENCES

- [1] Fabian Abel, Yashar Deldjoo, Mehdi Elahi, and Daniel Kohlsdorf. 2017. Recsys challenge 2017: Offline and online evaluation. In *Proceedings of the 11th ACM Conference on Recommender Systems (RecSys'17)*. Association for Computing Machinery, New York, NY, 372–373.
- [2] Adekunle Oluseyi Afolabi and Pekka Toivanen. 2020. Harmonization and categorization of metrics and criteria for evaluation of recommender systems in healthcare from dual perspectives. *Int. J. E-Health Med. Commun.* 11, 1 (2020), 69–92.
- [3] Rakesh Agrawal, Ramakrishnan Srikant, et al. 1994. Fast algorithms for mining association rules. In *Proceedings of the 20th International Conference of Very Large Data Bases (VLDB'94)*. 487–499.
- [4] Wissam Al Jurdi, Jacques Bou Abdo, Jacques Demerjian, and Abdallah Makhoul. 2021. Critique on natural noise in recommender systems. *ACM Trans. Knowl. Discov. Data* 15, 5, Article 75 (May 2021), 30 pages. <https://doi.org/10.1145/3447780>
- [5] Bushra Alhijawi, Arafat Awajan, and Salam Fraihat. 2022. Survey on the objectives of recommender systems: Measures, solutions, evaluation methodology, and new perspectives. *Comput. Surv.* 55, 5, Article 93 (Dec. 2022), 38 pages. <https://doi.org/10.1145/3527449>
- [6] Zafar Ali, Irfan Ullah, Amin Khan, Asim Ullah Jan, and Khan Muhammad. 2021. An overview and evaluation of citation recommendation models. *Scientometrics* 126, 5 (May 2021), 4083–4119. <https://doi.org/10.1007/s11192-021-03909-y>
- [7] James Allan, Donna Harman, Evangelos Kanoulas, Dan Li, Christophe Van Gysel, and Ellen M. Voorhees. 2017. TREC 2017 common core track overview. In *Proceedings of the Text Retrieval Conference (TREC'17)*. 14 pages. <https://trec.nist.gov/pubs/trec26/papers/Overview-CC.pdf>
- [8] Xavier Amatriain and Justin Basilico. 2016. Past, present, and future of recommender systems: An industry perspective. In *Proceedings of the 10th ACM Conference on Recommender Systems (RecSys'16)*. Association for Computing Machinery, New York, NY, USA, 211–214. <https://doi.org/10.1145/2959100.2959144>

- [9] Vito Walter Anelli, Alejandro Bellogin, Tommaso Di Noia, Dietmar Jannach, and Claudio Pomo. 2022. Top-N recommendation algorithms: A quest for the state-of-the-art. In *Proceedings of the 30th ACM Conference on User Modeling, Adaptation and Personalization (UMAP'22)*. Association for Computing Machinery, New York, NY, 121–131. <https://doi.org/10.1145/3503252.3531292>
- [10] Vito Walter Anelli, Alejandro Bellogin, Antonio Ferrara, Daniele Malitesta, Felice Antonio Merra, Claudio Pomo, Francesco Maria Donini, and Tommaso Di Noia. 2021. Elliot: A comprehensive and rigorous framework for reproducible recommender systems evaluation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'21)*. Association for Computing Machinery, New York, NY, 2405–2414. <https://doi.org/10.1145/3404835.3463245>
- [11] Krisztian Balog and Filip Radlinski. 2020. Measuring recommendation explanation quality: The conflicting goals of explanations. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'20)*. Association for Computing Machinery, New York, NY, 329–338. <https://doi.org/10.1145/3397271.3401032>
- [12] Linas Baltrunas, Karen Church, Alexandros Karatzoglou, and Nuria Oliver. 2015. Frappe: Understanding the usage and perception of mobile app recommendations in-the-wild. arXiv:1505.03014. Retrieved from <http://arxiv.org/abs/1505.03014>
- [13] Jöran Beel and Victor Brunel. 2019. Data pruning in recommender systems research: Best-practice or malpractice? In *Proceedings of ACM RecSys'19 Late-Breaking Results co-located with the 13th ACM Conference on Recommender Systems, RecSys'19 Late-Breaking Results (CEUR Workshop Proceedings, Vol. 2431)*, Marko Tkalcić and Sole Pera (Eds.). CEUR-WS.org, 26–30.
- [14] Joeran Beel, Bela Gipp, Stefan Langer, and Corinna Breitingner. 2015. Research paper recommender systems: A literature survey. *Int. J. Digit. Libr.* 17, 4 (2015), 305–338. <https://doi.org/10.1007/s00799-015-0156-0>
- [15] Joeran Beel, Stefan Langer, Marcel Genzmehr, Bela Gipp, Corinna Breitingner, and Andreas Nürnberger. 2013. Research paper recommender system evaluation: A quantitative literature survey. In *Proceedings of the International Workshop on Reproducibility and Replication in Recommender Systems Evaluation (RepSys'13)*. Association for Computing Machinery, New York, NY, 15–22. <https://doi.org/10.1145/2532508.2532512>
- [16] Poornima Belavadi, Laura Burbach, Stefan Ahlers, Martina Ziefle, and André Calero Valdez. 2021. Expectation, perception, and accuracy in news recommender systems: Understanding the relationships of user evaluation criteria using direct feedback. In *HCI International 2021—Late Breaking Papers: Design and User Experience*, Constantine Stephanidis, Marcelo M. Soares, Elizabeth Rosenzweig, Aaron Marcus, Sakae Yamamoto, Hirohiko Mori, Pei-Luen Patrick Rau, Gabriele Meiselwitz, Xiaowen Fang, and Abbas Moallem (Eds.). Springer International Publishing, Cham, 179–197. https://doi.org/10.1007/978-3-030-90238-4_14
- [17] Alejandro Bellogin, Pablo Castells, and Iván Cantador. 2017. Statistical biases in information retrieval metrics for recommender systems. *Inf. Retrieval* 20, 6 (2017), 606–634.
- [18] Alejandro Bellogin and Alan Said. 2018. Offline and online evaluation of recommendations. In *Collaborative Recommendations*, Shlomo Berkovsky, Iván Cantador, and Domonkos Tikk (Eds.). World Scientific, 295–328. https://doi.org/10.1142/9789813275355_0009
- [19] Alejandro Bellogin and Alan Said. 2021. Improving accountability in recommender systems research through reproducibility. *User Model. User-Adapt. Interact.* 31, 5 (2021), 941–977. <https://doi.org/10.1007/s11257-021-09302-x>
- [20] Rocío Cañamares and Pablo Castells. 2017. A probabilistic reformulation of memory-based collaborative filtering: Implications on popularity biases. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'17)*. Association for Computing Machinery, New York, NY, 215–224. <https://doi.org/10.1145/3077136.3080836>
- [21] Rocío Cañamares and Pablo Castells. 2018. Should I follow the crowd? A probabilistic analysis of the effectiveness of popularity in recommender systems. In *Proceedings of the 41st International ACM SIGIR Conference on Research & Development in Information Retrieval (SIGIR'18)*. Association for Computing Machinery, New York, NY, 415–424. <https://doi.org/10.1145/3209978.3210014>
- [22] Rocío Cañamares and Pablo Castells. 2020. On target item sampling in offline recommender system evaluation. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys'20)*. Association for Computing Machinery, New York, NY, 259–268. <https://doi.org/10.1145/3383313.3412259>
- [23] Diego Carraro and Derek Bridge. 2022. A sampling approach to debiasing the offline evaluation of recommender systems. *J. Intell. Inf. Syst.* 58, 2 (Apr. 2022), 311–336. <https://doi.org/10.1007/s10844-021-00651-y>
- [24] Óscar Celma. 2010. *Music Recommendation and Discovery in the Long Tail*. Springer, Berlin. <https://doi.org/10.1007/978-3-642-13287-2>
- [25] Li Chen, Yonghua Yang, Ningxia Wang, Keping Yang, and Quan Yuan. 2019. How serendipity improves user satisfaction with recommendations? A large-scale user evaluation. In *Proceedings of the World Wide Web Conference (TheWebConf'19)*. Association for Computing Machinery, New York, NY, 240–250. <https://doi.org/10.1145/3308558.3313469>

- [26] Jin Yao Chin, Yile Chen, and Gao Cong. 2022. The datasets dilemma: How much do we really know about recommendation datasets?. In *Proceedings of the 15th ACM International Conference on Web Search and Data Mining (WSDM'22)*. Association for Computing Machinery, New York, NY, 141–149. <https://doi.org/10.1145/3488560.3498519>
- [27] Eunjoon Cho, Seth A. Myers, and Jure Leskovec. 2011. Friendship and mobility: User movement in location-based social networks. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'11)*. Association for Computing Machinery, New York, NY, 1082–1090. <https://doi.org/10.1145/2020408.2020579>
- [28] Charles L. Clarke, Nick Craswell, and Ian Soboroff. 2009. *Overview of the Trec 2009 Web Track*. Technical Report. University of Waterloo, Ontario.
- [29] Charles L. Clarke, Nick Craswell, and Ellen M. Voorhees. 2012. *Overview of the TREC 2012 Web Track*. Technical Report. National Institute of Standards and Technology, Gaithersburg, MD.
- [30] Kevyn Collins-Thompson, Craig Macdonald, Paul Bennett, Fernando Diaz, and Ellen M. Voorhees. 2015. *TREC 2014 Web Track Overview*. Technical Report. University of Michigan, Ann Arbor.
- [31] Randell Cotta, Mingyang Hu, Dan Jiang, and Peizhou Liao. 2019. Off-policy evaluation of probabilistic identity data in lookalike modeling. In *Proceedings of the 12th ACM International Conference on Web Search and Data Mining (WSDM'19)*. Association for Computing Machinery, New York, NY, 483–491. <https://doi.org/10.1145/3289600.3291033>
- [32] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. 2020. Overview of the TREC 2019 deep learning track. arXiv:2003.07820. Retrieved from <https://arxiv.org/abs/2003.07820>.
- [33] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, Ellen M. Voorhees, and Ian Soboroff. 2021. TREC deep learning track: Reusable test collections in the large data regime. In *Proceedings of the Annual Conference of the Association for Computing Machinery Special Interest Group in Information Retrieval (SIGIR'21)*. Association for Computing Machinery, New York, NY, 2369–2375. <https://doi.org/10.1145/3404835.3463249>
- [34] Paolo Cremonesi, Yehuda Koren, and Roberto Turrin. 2010. Performance of recommender algorithms on top-n recommendation tasks. In *Proceedings of the 4th ACM Conference on Recommender Systems (RecSys'10)*. Association for Computing Machinery, New York, NY, 39–46. <https://doi.org/10.1145/1864708.1864721>
- [35] Alexander Dallmann, Daniel Zoller, and Andreas Hotho. 2021. A case study on sampling strategies for evaluating neural sequential item recommendation models. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys'21)*. Association for Computing Machinery, New York, NY, 505–514. <https://doi.org/10.1145/3460231.3475943>
- [36] Zohreh Dehghani Champiri, Adeleh Asemi, and Salim Siti Salwah Binti. 2019. Meta-analysis of evaluation methods and metrics used in context-aware scholarly recommender systems. *Knowl. Inf. Syst.* 61, 2 (Nov. 2019), 1147–1178. <https://doi.org/10.1007/s10115-018-1324-5>
- [37] Fernando Diaz and Andres Ferraro. 2022. Offline retrieval evaluation without evaluation metrics. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'22)*. Association for Computing Machinery, New York, NY, 599–609. <https://doi.org/10.1145/3477495.3532033>
- [38] Tome Eftimov, Bibek Paudel, Gorjan Popovski, and Dragi Kocev. 2021. A framework for evaluating personalized ranking systems by fusing different evaluation measures. *Big Data Res.* 25 (2021), 13 Pages. <https://doi.org/10.1016/j.bdr.2021.100211>
- [39] Michael D. Ekstrand. 2020. LensKit for python: Next-generation software for recommender systems experiments. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management (CIKM'20)*. Association for Computing Machinery, New York, NY, 2999–3006. <https://doi.org/10.1145/3340531.3412778>
- [40] Maurizio Ferrari Dacrema, Simone Boglio, Paolo Cremonesi, and Dietmar Jannach. 2021. A troubling analysis of reproducibility and progress in recommender systems research. *ACM Trans. Inf. Syst.* 39, 2, Article 20 (Jan. 2021), 49 pages. <https://doi.org/10.1145/3434185>
- [41] Maurizio Ferrari Dacrema, Paolo Cremonesi, and Dietmar Jannach. 2019. Are we really making much progress? A worrying analysis of recent neural recommendation approaches. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys'19)*. Association for Computing Machinery, New York, NY, 101–109. <https://doi.org/10.1145/3298689.3347058>
- [42] Shir Frumerman, Guy Shani, Bracha Shapira, and Oren Sar Shalom. 2019. Are all rejected recommendations equally bad? Towards analysing rejected recommendations. In *Proceedings of the 27th ACM Conference on User Modeling, Adaptation and Personalization (UMAP'19)*. Association for Computing Machinery, New York, NY, 157–165. <https://doi.org/10.1145/3320435.3320448>
- [43] Xue Geng, Hanwang Zhang, Jingwen Bian, and Tat-Seng Chua. 2015. Learning image and user features for recommendation in social networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV'15)*. IEEE Press, 4274–4282. <https://doi.org/10.1109/ICCV.2015.486>
- [44] Alexandre Gilotte, Clément Calauzènes, Thomas Nedelec, Alexandre Abraham, and Simon Dollé. 2018. Offline A/B testing for recommender systems. In *Proceedings of the 11th ACM International Conference on Web Search and Data Mining (WSDM'18)*. Association for Computing Machinery, New York, NY, 198–206. <https://doi.org/10.1145/3159652.3159687>

- [45] Asela Gunawardana, Guy Shani, and Sivan Yogev. 2022. Evaluating recommender systems. In *Recommender Systems Handbook*, Francesco Ricci, Lior Rokach, and Bracha Shapira (Eds.). Springer US, New York, NY, 547–601. https://doi.org/10.1007/978-1-0716-2197-4_15
- [46] G. Guo, J. Zhang, and N. Yorke-Smith. 2013. A novel Bayesian similarity measure for recommender systems. In *Proceedings of the 23rd International Joint Conference on Artificial Intelligence (IJCAI'13)*. AAAI Press, 2619–2625.
- [47] Xunhua Guo, Lingli Wang, Mingyue Zhang, and Guoqing Chen. 2022. First things first? Order effects in online product recommender systems. *ACM Trans. Comput.-Hum. Interact.* 30, 1 (2022), 1–35. <https://doi.org/10.1145/3557886>
- [48] F. Maxwell Harper and Joseph A. Konstan. 2015. The MovieLens datasets: History and context. *ACM Trans. Interact. Intell. Syst.* 5, 4 (Dec. 2015), 1–19.
- [49] Ruining He and Julian McAuley. 2016. VBPR: Visual Bayesian personalized ranking from implicit feedback. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence*, 144–150. <https://doi.org/10.1609/aaai.v30i1.9973>
- [50] Jonathan L. Herlocker, Joseph A. Konstan, Loren G. Terveen, and John T. Riedl. 2004. Evaluating collaborative filtering recommender systems. *ACM Trans. Inf. Syst.* 22, 1 (Jan. 2004), 5–53. <https://doi.org/10.1145/963770.963772>
- [51] Ngozi Ihemelandu and Michael D. Ekstrand. 2021. Statistical inference: The missing piece of RecSys experiment reliability discourse. In *Proceedings of the Perspectives on the Evaluation of Recommender Systems, Workshop co-located with the 15th ACM Conference on Recommender Systems (RecSys'21) (CEUR Workshop Proceedings, Vol. 2955)*, Eva Zangerle, Christine Bauer, and Alan Said (Eds.). CEUR-WS.org, Aachen, Germany, 10 pages. <https://ceur-ws.org/Vol-2955/paper9.pdf>
- [52] Dietmar Jannach. 2023. Evaluating conversational recommender systems: A landscape of research. *Artif. Intell. Rev.* 56, 3 (2023), 2365–2400. <https://doi.org/10.1007/s10462-022-10229-x>
- [53] Dietmar Jannach and Christine Bauer. 2020. Escaping the McNamara fallacy: Towards more impactful recommender systems research. *AI Mag.* 41, 4 (Dec. 2020), 79–95. <https://doi.org/10.1609/aimag.v41i4.5312>
- [54] Yucheng Jin, Li Chen, Wanling Cai, and Pearl Pu. 2021. Key qualities of conversational recommender systems: From users' perspective. In *Proceedings of the 9th International Conference on Human-Agent Interaction (HAI'21)*. Association for Computing Machinery, New York, NY, 93–102. <https://doi.org/10.1145/3472307.3484164>
- [55] Thorsten Joachims, Ben London, Yi Su, Adith Swaminathan, and Lequn Wang. 2021. Recommendations as treatments. *AI Mag.* 42, 3 (2021), 19–30. <https://doi.org/10.1609/aaai.12014>
- [56] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *Proceedings of the IEEE International Conference on Data Mining*. IEEE, 197–206. <https://doi.org/10.1109/ICDM.2018.00035>
- [57] Barbara Kitchenham, Stuart Charters, David Budgen, Pearl Brereton, Mark Turner, Steve Linkman, Magne Jørgensen, Emilia Mendes, and Giuseppe Visaggio. 2007. *Guidelines for Performing Systematic Literature Reviews in Software Engineering*. EBSE Technical Report EBSE-2007-01, version 2.3. Keele University and University of Durham.
- [58] Haruka Kiyohara, Yuta Saito, Tatsuya Matsuhira, Yusuke Narita, Nobuyuki Shimizu, and Yasuo Yamamoto. 2022. Doubly robust off-policy evaluation for ranking policies under the cascade behavior model. In *Proceedings of the 15th ACM International Conference on Web Search and Data Mining (WSDM'22)*. Association for Computing Machinery, New York, NY, 487–497. <https://doi.org/10.1145/3488560.3498380>
- [59] Pigi Kouki, Ilias Fountalis, Nikolaos Vasiloglou, Xiquan Cui, Edo Liberty, and Khalifeh Al Jadda. 2020. From the lab to production: A case study of session-based recommendations in the home-improvement domain. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys'20)*. Association for Computing Machinery, New York, NY, 140–149. <https://doi.org/10.1145/3383313.3412235>
- [60] Walid Krichene and Steffen Rendle. 2020. On sampled metrics for item recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD'20)*. Association for Computing Machinery, New York, NY, 1748–1757. <https://doi.org/10.1145/3394486.3403226>
- [61] Sara Latifi and Dietmar Jannach. 2022. Streaming session-based recommendation: When graph neural networks meet the neighborhood. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys'22)*. Association for Computing Machinery, New York, NY, 420–426. <https://doi.org/10.1145/3523227.3548485>
- [62] Sara Latifi, Dietmar Jannach, and Andrés Ferraro. 2022. Sequential recommendation: A study on transformers, nearest neighbors and sampled metrics. *Inf. Sci.* 609 (2022), 660–678. <https://doi.org/10.1016/j.ins.2022.07.079>
- [63] Dong Li, Ruoming Jin, Jing Gao, and Zhi Liu. 2020. On sampling Top-K recommendation evaluation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD'20)*. Association for Computing Machinery, New York, NY, 2114–2124. <https://doi.org/10.1145/3394486.3403262>
- [64] Roger Zhe Li, Julián Urbano, and Alan Hanjalic. 2021. New insights into metric optimization for ranking-based recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'21)*. Association for Computing Machinery, New York, NY, 932–941. <https://doi.org/10.1145/3404835.3462973>
- [65] Hongyu Lu, Weizhi Ma, Min Zhang, Maarten de Rijke, Yiqun Liu, and Shaoping Ma. 2021. Standing in your shoes: External assessments for personalized recommender systems. In *Proceedings of the 44th International ACM SIGIR*

- Conference on Research and Development in Information Retrieval (SIGIR'21)*. Association for Computing Machinery, New York, NY, 1523–1533. <https://doi.org/10.1145/3404835.3462916>
- [66] Malte Ludewig, Noemi Mauro, Sara Latifi, and Dietmar Jannach. 2019. Performance comparison of neural and non-neural approaches to session-based recommendation. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys'19)*. Association for Computing Machinery, New York, NY, 462–466. <https://doi.org/10.1145/3298689.3347041>
- [67] Benjamin M. Marlin, Richard S. Zemel, Sam Roweis, and Malcolm Slaney. 2007. Collaborative filtering and the missing at random assumption. In *Proceedings of the 23rd Conference on Uncertainty in Artificial Intelligence (UAI'07)*. AUAI Press, Arlington, VA, 267–275.
- [68] Julian McAuley, Jure Leskovec, and Dan Jurafsky. 2012. Learning attitudes and attributes from multi-aspect reviews. In *Proceedings of the IEEE 12th International Conference on Data Mining (ICDM'12)*. IEEE Press, 1020–1025.
- [69] James McInerney, Brian Brost, Praveen Chandar, Rishabh Mehrotra, and Benjamin Carterette. 2020. Counterfactual evaluation of slate recommendations with sequential reward interactions. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD'20)*. Association for Computing Machinery, New York, NY, 1779–1788. <https://doi.org/10.1145/3394486.3403229>
- [70] Elisa Mena-Maldonado, Rocío Cañamares, Pablo Castells, Yongli Ren, and Mark Sanderson. 2020. Agreement and disagreement between true and false-positive metrics in recommender systems evaluation. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'20)*. Association for Computing Machinery, New York, NY, USA, 841–850. <https://doi.org/10.1145/3397271.3401096>
- [71] Elisa Mena-Maldonado, Rocío Cañamares, Pablo Castells, Yongli Ren, and Mark Sanderson. 2021. Popularity bias in false-positive metrics for recommender systems evaluation. *ACM Trans. Inf. Syst.* 39, 3, Article 36 (May 2021), 43 pages. <https://doi.org/10.1145/3452740>
- [72] Rishabh Misra, Mengting Wan, and Julian McAuley. 2018. Decomposing fit semantics for product size recommendation in metric spaces. In *Proceedings of the 12th ACM Conference on Recommender Systems*. Association for Computing Machinery, New York, NY, 422–426.
- [73] Yusuke Narita, Shota Yasui, and Kohei Yata. 2021. Debaised off-policy evaluation for recommendation systems. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys'21)*. Association for Computing Machinery, New York, NY, 372–379. <https://doi.org/10.1145/3460231.3474231>
- [74] Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP'19)*. Association for Computational Linguistics, 188–197. <https://doi.org/10.18653/v1/D19-1018>
- [75] Malte Ostendorff, Corinna Breitingner, and Bela Gipp. 2021. A qualitative evaluation of user preference for link-based vs. text-based recommendations of wikipedia articles. In *Towards Open and Trustworthy Digital Societies: Proceedings of the 23rd International Conference on Asia-Pacific Digital Libraries (ICADL'21)*, Hao-Ren Ke, Chei Sian Lee, and Kazunari Sugiyama (Eds.). Springer International Publishing, Cham, Germany, 63–79. https://doi.org/10.1007/978-3-030-91669-5_6
- [76] Javier Parapar and Filip Radlinski. 2021. Towards unified metrics for accuracy and diversity for recommender systems. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys'21)*. Association for Computing Machinery, New York, NY, 75–84. <https://doi.org/10.1145/3460231.3474234>
- [77] Ladislav Peska and Peter Vojtas. 2020. Off-line vs. On-line evaluation of recommender systems in small E-Commerce. In *Proceedings of the 31st ACM Conference on Hypertext and Social Media (HT'20)*. Association for Computing Machinery, New York, NY, 291–300. <https://doi.org/10.1145/3372923.3404781>
- [78] Pearl Pu, Li Chen, and Rong Hu. 2012. Evaluating recommender systems from the user's perspective: Survey of the state of the art. *User Model. User-Adapt. Interact.* 22, 4 (Oct. 2012), 317–355. <https://doi.org/10.1007/s11257-011-9115-7>
- [79] Steffen Rendle, Li Zhang, and Yehuda Koren. 2019. On the difficulty of evaluating baselines: A study on recommender systems. <https://doi.org/10.48550/ARXIV.1905.01395>
- [80] Alan Said and Alejandro Bellogín. 2018. Coherence and inconsistencies in rating behavior: Estimating the magic barrier of recommender systems. *User Model. User-Adapt. Interact.* 28, 2 (2018), 97–125. <https://doi.org/10.1007/s11257-018-9202-0>
- [81] Yuta Saito, Shunsuke Aihara, Megumi Matsutani, and Yusuke Narita. 2020. Open bandit dataset and pipeline: Towards realistic and reproducible off-policy evaluation. <https://doi.org/10.48550/ARXIV.2008.07146>
- [82] Yuta Saito, Takuma Udagawa, Haruka Kiyohara, Kazuki Mogi, Yusuke Narita, and Kei Tateno. 2021. Evaluating the robustness of off-policy evaluation. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys'21)*. Association for Computing Machinery, New York, NY, 114–123. <https://doi.org/10.1145/3460231.3474245>

- [83] Pablo Sánchez and Alejandro Bellogin. 2022. Point-of-interest recommender systems based on location-based social networks: A survey from an experimental perspective. *Comput. Surv.* 54, 11s, Article 223 (Sep. 2022), 37 pages. <https://doi.org/10.1145/3510409>
- [84] Prabhath Kumar Saraswat, Samuel William, and Eswar Reddy. 2021. A hybrid approach for offline A/B evaluation for item ranking algorithms in recommendation systems. In *Proceedings of the 1st International Conference on AI-ML-Systems (AIMLSystems'21)*. Association for Computing Machinery, New York, NY, Article 21, 6 pages. <https://doi.org/10.1145/3486001.3486241>
- [85] Dietmar Schabus, Marcin Skowron, and Martin Trapp. 2017. One million posts: A data set of german online discussions. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'17)*. Association for Computing Machinery, 1241–1244. <https://doi.org/10.1145/3077136.3080711>
- [86] Tobias Schnabel, Adith Swaminathan, Ashudeep Singh, Navin Chandak, and Thorsten Joachims. 2016. Recommendations as treatments: Debiasing learning and evaluation. In *International Conference on Machine Learning*. PMLR, 1670–1679.
- [87] Thiago Silva, Nicollas Silva, Heitor Werneck, Carlos Mito, Adriano C.M. Pereira, and Leonardo Rocha. 2022. IRec: An interactive recommendation framework. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'22)*. Association for Computing Machinery, New York, NY, 3165–3175. <https://doi.org/10.1145/3477495.3531754>
- [88] Nasim Sonboli, Masoud Mansoury, Ziyue Guo, Shreyas Kadekodi, Weiwen Liu, Zijun Liu, Andrew Schwartz, and Robin Burke. 2021. Librec-auto: A tool for recommender systems experimentation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management (CIKM'21)*. Association for Computing Machinery, New York, NY, 4584–4593. <https://doi.org/10.1145/3459637.3482006>
- [89] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. BERT4Rec: Sequential recommendation with bidirectional encoder representations from transformer. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management (CIKM'19)*. Association for Computing Machinery, New York, NY, 1441–1450. <https://doi.org/10.1145/3357384.3357895>
- [90] Zhu Sun, Di Yu, Hui Fang, Jie Yang, Xinghua Qu, Jie Zhang, and Cong Geng. 2020. Are we evaluating rigorously? Benchmarking recommendation for reproducible evaluation and fair comparison. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys'20)*. Association for Computing Machinery, New York, NY, 23–32. <https://doi.org/10.1145/3383313.3412489>
- [91] Panagiotis Symeonidis, Andrea Janes, Dmitry Chaltsev, Philip Giuliani, Daniel Morandini, Andreas Unterhuber, Ludovik Coba, and Markus Zanker. 2020. Recommending the video to watch next: An offline and online evaluation at YOUTV.De. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys'20)*. Association for Computing Machinery, New York, NY, 299–308. <https://doi.org/10.1145/3383313.3412257>
- [92] Jiliang Tang, Huiji Gao, and Huan Liu. 2012. MTrust: Discerning multi-faceted trust in a connected world. In *Proceedings of the 5th ACM International Conference on Web Search and Data Mining (WSDM'12)*. Association for Computing Machinery, New York, NY, 93–102. <https://doi.org/10.1145/2124295.2124309>
- [93] Hao Wang, Naiyan Wang, and Dit-Yan Yeung. 2015. Collaborative deep learning for recommender systems. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'15)*. Association for Computing Machinery, New York, NY, 1235–1244. <https://doi.org/10.1145/2783258.2783273>
- [94] Doris Xin, Nicolas Mayoraz, Hubert Pham, Karthik Lakshmanan, and John R. Anderson. 2017. Folding: Why good models sometimes make spurious recommendations. In *Proceedings of the 11th ACM Conference on Recommender Systems (RecSys'17)*. Association for Computing Machinery, New York, NY, 201–209. <https://doi.org/10.1145/3109859.3109911>
- [95] Longqi Yang, Yin Cui, Yuan Xuan, Chenyang Wang, Serge Belongie, and Deborah Estrin. 2018. Unbiased offline recommender evaluation for missing-not-at-random implicit feedback. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys'18)*. Association for Computing Machinery, New York, NY, 279–287. <https://doi.org/10.1145/3240323.3240355>
- [96] Eva Zangerle and Christine Bauer. 2022. Evaluating recommender systems: Survey and framework. *Comput. Surv.* 55, 8, Article 170 (2022), 38 pages. <https://doi.org/10.1145/3556536>
- [97] Tong Zhao, Julian McAuley, and Irwin King. 2015. Improving latent factor models via personalized feature projection for one class recommendation. In *Proceedings of the 24th ACM International Conference on Information and Knowledge Management (CIKM'15)*. Association for Computing Machinery, New York, NY, 821–830. <https://doi.org/10.1145/2806416.2806511>
- [98] Wayne Xin Zhao, Zihan Lin, Zhichao Feng, Pengfei Wang, and Ji-Rong Wen. 2022. A revisiting study of appropriate offline evaluation for Top-N recommendation algorithms. *ACM Trans. Inf. Syst.* 41, 2, Article 32 (Dec. 2022), 41 pages. <https://doi.org/10.1145/3545796>

- [99] Jieming Zhu, Jinyang Liu, Shuai Yang, Qi Zhang, and Xiuqiang He. 2021. Open benchmarking for click-through rate prediction. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management (CIKM'21)*. Association for Computing Machinery, New York, NY, 2759–2769. <https://doi.org/10.1145/3459637.3482486>
- [100] Cai-Nicolas Ziegler, Sean M. McNee, Joseph A. Konstan, and Georg Lausen. 2005. Improving recommendation lists through topic diversification. In *Proceedings of the 14th International Conference on World Wide Web (WWW'05)*. Association for Computing Machinery, New York, NY, 22–32. <https://doi.org/10.1145/1060745.1060754>

Received 14 December 2022; revised 19 July 2023; accepted 10 September 2023

Appendix

Exploring the Landscape of Recommender Systems Evaluation: Practices and Perspectives

[CHRISTINE BAUER](#), Paris Lodron University Salzburg, Austria

[EVA ZANGERLE](#), University of Innsbruck, Austria

[ALAN SAID](#), University of Gothenburg, Sweden

Table: Excluded papers, sorted alphabetically

Papers	Not on recommender systems	No contribution to evaluation	Workshop/ Challenge/ Tutorial overview paper, or Research plan	Short paper, Workshop paper, or Demo paper	Retracted
Abbas et al. [1]	x				
Abdollahpouri et al. [3]	x			x	
Abdollahpouri et al. [2]		x		x	
Abid et al. [4]		x		x	
Ajoudanian and Abadeh [5]		x		x	
ALRossais [6]		x		x	
Alslaity and Tran [7]				x	
Alves et al. [8]		x			
Amer et al. [9]		x			
Amin et al. [10]		x			
Anelli et al. [11]				x	
Anwar et al. [12]		x			
Assami et al. [13]		x			
Bagnato [14]		x			
Barraza-Urbina and Glowacka [15]			x		
Barzegar Nozari and Koochi [16]		x			
Bauer [17]			x		
Bauer [18]			x		
Bayle et al. [19]	x				
Berendt et al. [20]				x	
Bernardis et al. [21]		x			
Brusilovsky et al. [22]			x		

Authors' addresses: [Christine Bauer](#), christine.bauer@plus.ac.at, Paris Lodron University Salzburg, Salzburg, 5020, Jakob-Haringer-Strasse 1, Austria; [Eva Zangerle](#), eva.zangerle@uibk.ac.at, University of Innsbruck, Technikerstr. 21A, Innsbruck, 6020, Austria; [Alan Said](#), University of Gothenburg, Sweden, alansaid@acm.org.

Table: Excluded papers, sorted alphabetically

Papers	Not on recommender systems	No contribution to evaluation	Workshop/ Challenge/ Tutorial overview paper, or Research plan	Short paper, Workshop paper, or Demo paper	Retracted
Brusilovsky et al. [23]			x		
Brusilovsky et al. [24]			x		
Brusilovsky et al. [25]			x		
Bukowski et al. [26]	x				
Burke et al. [27]			x		
Cañamares and Castells [28]				x	
Cai et al. [29]		x			
Caro-Álvaro et al. [30]		x			
Caselles-Dupré et al. [31]		x			
Celik et al. [32]				x	
Chamberlain et al. [33]		x			
Chatterjee [34]		x			
Chen et al. [35]			x		
Cinelli et al. [36]		x			
Clemente et al. [37]		x			
Coba et al. [38]		x			
Coba et al. [39]		x			
Colaço et al. [40]		x			
Contreras et al. [41]		x			
Coró et al. [42]		x			
Costa and Tamzalit [43]		x			
Cugny et al. [44]		x			
Cui et al. [45]	x				
Cui et al. [46]		x			
Da et al. [47]		x			
Dareddy [48]			x		
Dehghani Champiri et al. [49]		x			
Deldjoo et al. [50]		x			
Deldjoo et al. [51]				x	
Di Penta [52]		x			
Di Rocco et al. [53]		x			
Dokoupil and Peska [54]				x	

Table: Excluded papers, sorted alphabetically

Papers	Not on recommender systems	No contribution to evaluation	Workshop/ Challenge/ Tutorial overview paper, or Research plan	Short paper, Workshop paper, or Demo paper	Retracted
Doryab et al. [55]		x			
Dror et al. [56]		x			
Ekstrand et al. [57]			x		
Ekstrand and Kluver [58]		x			
El Alaoui et al. [59]		x			
Elahi et al. [60]		x			
Fails et al. [61]			x		
Fang et al. [62]		x			
Felicioni et al. [64]				x	
Felicioni et al. [63]				x	
Feng et al. [65]		x			
Ferraro et al. [66]				x	
Ferwerda et al. [67]				x	
Figueroa et al. [68]		x			
Fopa et al. [69]		x			
Forster et al. [70]		x			
Freire and de Castro [71]		x			
Gabbolini and Bridge [72]		x			
Gao and Shah [73]				x	
Garcia-Gathright et al. [74]			x		
Garcia-Gathright et al. [75]		x			
Gharahighehi et al. [76]		x			
Goyal [77]	x				
Grace et al. [78]		x			
Guinea et al. [79]		x			
Guo and Xie [80]		x			
Haggag [81]	x				
Haggag et al. [82]	x				
Han et al. [83]		x			
Hao and Zhang [84]		x			
Hendrickx et al. [85]		x			
Hiraishi et al. [86]		x			

Table: Excluded papers, sorted alphabetically

Papers	Not on recommender systems	No contribution to evaluation	Workshop/ Challenge/ Tutorial overview paper, or Research plan	Short paper, Workshop paper, or Demo paper	Retracted
Hong et al. [87]		x			
Houtti et al. [88]	x				
Hu et al. [89]		x			
Huang et al. [90]		x			
Ihemelandu [91]		x			
Iizuka et al. [92]		x			
Isinkaye et al. [93]		x			
Jadidinejad et al. [94]		x			
Jannach et al. [95]		x			
Jeunen [96]				x	
Jeunen et al. [97]			x		
Ji et al. [98]		x			
Jiang [99]		x			
Joachims et al. [100]			x		
Joachims et al. [101]			x		
Joachims et al. [102]			x		
Jugovac et al. [103]				x	
Kale et al. [104]			x		
Kang et al. [105]		x			
Kang et al. [106]		x			
Kang et al. [107]	x				
Kartoglu and Spratling [108]		x			
Keller and Munz [109]		x			
Khan et al. [110]		x			
Kim et al. [111]		x			
Kitazawa and Yui [112]		x			
Kokkalas et al. [113]		x			
Kondratova and Emond [114]		x			
Kotkov et al. [115]		x			
Kowald et al. [116]				x	
Kowald et al. [117]		x			
Kowald et al. [118]		x			

Table: Excluded papers, sorted alphabetically

Papers	Not on recommender systems	No contribution to evaluation	Workshop/ Challenge/ Tutorial overview paper, or Research plan	Short paper, Workshop paper, or Demo paper	Retracted
Kuhnle et al. [119]	x				
Kumar and Parimala [120]		x			
Kurbatova et al. [121]				x	
Lacic et al. [122]				x	
Laishram and Padmanabhan [123]		x			
Lavee et al. [124]		x			
Lee et al. [125]		x			
Leng and Yu [126]		x			
Lex and Schedl [127]			x		
Liang [128]		x			
Liang and Willemsen [129]		x			
Linda et al. [130]		x			
Liu et al. [131]		x			
Liu et al. [132]		x			
Ludewig et al. [133]				x	
Ludewig et al. [134]		x			
Luo et al. [135]		x			
Ma et al. [136]		x			
Maccatrozzo et al. [137]		x			
Maleszka [138]		x			
Mangili et al. [139]		x			
Manzoor and Jannach [140]		x			
Margaris et al. [141]		x			
Margaris et al. [142]	x				
Margaris et al. [143]		x			
Margaris et al. [144]		x			
Margaris and Vassilakis [145]		x			
McGill et al. [146]	x				
Mello et al. [147]	x				
Meng et al. [148]				x	
Mesas and Bellogin [149]				x	
Michiels et al. [150]				x	

Table: Excluded papers, sorted alphabetically

Papers	Not on recommender systems	No contribution to evaluation	Workshop/ Challenge/ Tutorial overview paper, or Research plan	Short paper, Workshop paper, or Demo paper	Retracted
Mishra et al. [151]		x			
Mo et al. [152]	x				
Moens et al. [153]				x	
Mogenet et al. [154]				x	
Mohammadian et al. [155]		x			
Morales et al. [156]		x			
Moses and Babu [157]		x			
Mueller [158]		x			
Muellner et al. [159]		x			
Nagashima and He [160]		x			
Neidhardt et al. [161]			x		
Neidhardt et al. [162]			x		
Neidhardt et al. [163]			x		
Neophytou et al. [164]			x		
Neves et al. [165]		x			
Nikolakopoulos et al. [166]		x			
Novozhilov et al. [167]				x	
Ojokoh et al. [168]		x			
Ong et al. [169]		x			
Ospanova and Shelestova [170]	x				
Ou et al. [171]	x				
Pálovics et al. [172]				x	
Pérez Maurera et al. [183]				x	
Pérez Maurera et al. [182]		x			
Pandey et al. [173]		x			
Paraschakis and Nilsson [174]		x			
Parizi et al. [175]	x				
Park et al. [176]	x				
Paryudi et al. [177]		x			
Paun [178]		x			
Pellegrini et al. [179]		x			
Peng et al. [180]	x				

Table: Excluded papers, sorted alphabetically

Papers	Not on recommender systems	No contribution to evaluation	Workshop/ Challenge/ Tutorial overview paper, or Research plan	Short paper, Workshop paper, or Demo paper	Retracted
Perez et al. [181]		x			
Peska and Balcar [184]		x			
Pinheiro et al. [185]	x				
Piscopo et al. [186]			x		
Pitt et al. [187]	x				
Poirson and Cunha [188]		x			
Polatidis et al. [189]		x			
Polatidis et al. [190]		x			
Polato and Aiolli [191]		x			
Prato et al. [192]		x			
Qiao and Wang [193]		x			
Qin [194]		x			
Qiu and Niu [195]		x			
Quadrana et al. [196]		x			
Rahmani et al. [197]		x			
Rappaz et al. [198]		x			
Raza and Ding [199]		x			
Ricci et al. [200]				x	
Sá et al. [201]		x			
Sánchez [205]		x			
Sánchez and Bellogín [206]		x			
Sachdeva and McAuley [202]		x			
Saito and Joachims [203]				x	
Saito and Joachims [204]					
Sánchez and Dietz [207]				x	
Sanz-Cruzado and Castells [208]		x			
Sanz-Cruzado and Castells [209]		x			
Sanz-Cruzado et al. [210]		x			
Sar Shalom et al. [211]			x		
Sar Shalom et al. [212]			x		
Sasaki et al. [213]	x				
Sato et al. [214]		x			

Table: Excluded papers, sorted alphabetically

Papers	Not on recommender systems	No contribution to evaluation	Workshop/ Challenge/ Tutorial overview paper, or Research plan	Short paper, Workshop paper, or Demo paper	Retracted
Sato et al. [215]		x			
Schaffer et al. [216]		x			
Schwarzer et al. [217]		x			
Sen et al. [218]	x				
Seo et al. [219]		x			
Seo et al. [220]		x			
Sesagiri Raamkumar and Foo [221]				x	
Sharma et al. [222]		x			
Shen et al. [223]	x				
Shu et al. [224]		x			
Silva et al. [225]		x			
Silva-Rodríguez et al. [226]		x			
Silveira et al. [227]				x	
Sinha and Dhanalakshmi [228]		x			
Slokom et al. [229]		x			
Smirnov and Ponomarev [230]		x			
Song et al. [231]		x			
Stavinova et al. [232]		x			
Sudhakar et al. [233]	x				
Sun and Lee [234]		x			
Sun et al. [235]		x			
Tamm et al. [236]				x	
Tan et al. [237]		x			
Tang et al. [238]		x			
Thirumuruganathan et al. [239]		x			
Umemoto [240]				x	
Ustalov et al. [241]					
Uta et al. [242]				x	
Vaquero-Patricio et al. [243]		x			
Vasile et al. [244]			x		
Velásquez et al. [245]	x				
Vij et al. [246]	x				

Table: Excluded papers, sorted alphabetically

Papers	Not on recommender systems	No contribution to evaluation	Workshop/ Challenge/ Tutorial overview paper, or Research plan	Short paper, Workshop paper, or Demo paper	Retracted
Villegas et al. [247]		x			
Vougioukas et al. [248]	x				
Vougioukas et al. [249]	x				
Wang et al. [250]		x			
Wang et al. [251]	x				
Wang et al. [252]		x			
Wang et al. [253]		x			
Wasilewski and Hurley [254]		x			
Wen et al. [255]				x	
Wen et al. [256]		x			
Wong et al. [257]				x	
Xu [258]		x			
Xu et al. [259]		x			
Xu et al. [260]		x			
Yalcin [261]		x			
Yang et al. [262]					x
Yang et al. [263]		x			
Yang et al. [264]		x			
Yoshida et al. [265]		x		x	
You et al. [266]	x				
Younes and Boukerche [267]	x				
Yu et al. [268]	x				
Zamani et al. [269]			x		
Zangerle et al. [270]			x		
Zangerle et al. [271]			x		
Zhang et al. [272]		x			
Zhang et al. [273]		x			
Zhao et al. [274]		x			
Zhao et al. [275]	x				
Zhao et al. [276]				x	
Zhao et al. [277]				x	
Zhitomirsky-Geffet and Zadok [278]		x			

Table: Excluded papers, sorted alphabetically

Papers	Not on recommender systems	No contribution to evaluation	Workshop/ Challenge/ Tutorial overview paper, or Research plan	Short paper, Workshop paper, or Demo paper	Retracted
Zhou and Tan [279]		x			
Ziarani and Ravanmehr [280]		x			
Zipori and Sarne [281]		x			
Zolaktaf et al. [282]				x	

REFERENCES

- [1] Asad Abbas, Hussein Haruna, Arturo Arrona-Palacios, Claudia Camacho-Zuñiga, Sandra Núñez-Daruich, Jose Francisco Enríquez de la O, Raquel Castaño-Gonzalez, Jose Escamilla, and Samira Hosseini. 2022. Students' Evaluations of Teachers and Recommendation Based on Course Structure or Teaching Approaches: An Empirical Study Based on the Institutional Dataset of Student Opinion Survey. *Education and Information Technologies* 27, 9 (2022), 12049–12064. <https://doi.org/10.1007/s10639-022-11119-z>
- [2] Himan Abdollahpouri, Robin Burke, and Bamshad Mobasher. 2017. Controlling Popularity Bias in Learning-to-Rank Recommendation. In *Proceedings of the Eleventh ACM Conference on Recommender Systems (RecSys '17)*. Association for Computing Machinery, New York, NY, USA, 42–46. <https://doi.org/10.1145/3109859.3109912>
- [3] Himan Abdollahpouri, Robin Burke, and Bamshad Mobasher. 2017. Recommender Systems as Multistakeholder Environments. In *Proceedings of the 25th Conference on User Modeling, Adaptation and Personalization (UMAP '17)*. Association for Computing Machinery, New York, NY, USA, 347–348. <https://doi.org/10.1145/3079628.3079657>
- [4] Shamsa Abid, Hamid Abdul Basit, and Shafay Shamail. 2022. Context-Aware Code Recommendation in IntelliJ IDEA. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2022)*. Association for Computing Machinery, New York, NY, USA, 1647–1651. <https://doi.org/10.1145/3540250.3558937>
- [5] Shohreh Ajoudanian and Maryam Nooraei Abadeh. 2019. Recommending Human Resources to Project Leaders Using a Collaborative Filtering-Based Recommender System: Case Study of GitHub. *IET Software* 13, 5 (2019), 379–385. <https://doi.org/10.1049/iet-sen.2018.5261>
- [6] Nourah A. ALRossais. 2018. Integrating Item Based Stereotypes in Recommender Systems. In *Proceedings of the 26th Conference on User Modeling, Adaptation and Personalization (UMAP '18)*. Association for Computing Machinery, New York, NY, USA, 265–268. <https://doi.org/10.1145/3209219.3213593>
- [7] Alaa Alslaity and Thomas Tran. 2021. Goal Modeling-Based Evaluation for Personalized Recommendation Systems. In *Adjunct Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '21)*. Association for Computing Machinery, New York, NY, USA, 276–283. <https://doi.org/10.1145/3450614.3464619>
- [8] Rodrigo Alves, Antoine Ledent, and Marius Kloft. 2021. Burst-induced Multi-Armed Bandit for Learning Recommendation. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys 21)*. Association for Computing Machinery, New York, NY, USA, 292–301. <https://doi.org/10.1145/3460231.3474250>
- [9] Ali A. Amer, Hassan I. Abdalla, and Loc Nguyen. 2021. Enhancing Recommendation Systems Performance Using Highly-Effective Similarity Measures. *Knowledge-Based Systems* 217, C (2021). <https://doi.org/10.1016/j.knosys.2021.106842>
- [10] Sana Abida Amin, James Philips, and Nasseh Tabrizi. 2019. Current Trends in Collaborative Filtering Recommendation Systems. In *SERVICES 2019: 15th World Congress, Held as Part of the Services Conference Federation, SCF 2019, San Diego, CA, USA, June 25–30, 2019, Proceedings (SERVICES 2019)*. Springer, 46–60. https://doi.org/10.1007/978-3-030-23381-5_4
- [11] Vito Walter Anelli, Alejandro Bellogin, Antonio Ferrara, Daniele Malatesta, Felice Antonio Merra, Claudio Pomo, Francesco Maria Donini, and Tommaso Di Noia. 2021. V-Elliot: Design, Evaluate and Tune Visual Recommender Systems. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys '21)*. Association for Computing Machinery, New York, NY, USA, 768–771. <https://doi.org/10.1145/3460231.3478881>
- [12] Khalid Anwar, Jamshed Siddiqui, and Shahab Saquib Sohail. 2020. Machine Learning-Based Book Recommender System: A Survey and New Perspectives. *International Journal of Intelligent Information and Database Systems* 13, 2–4 (2020), 231–248. <https://doi.org/10.1504/ijids.2020.109457>
- [13] Sara Assami, Najima Daoudi, and Rachida Ajhoun. 2020. Learning Actor Ontology for a Personalised Recommendation in Massive Open Online Courses. *International Journal of Technology Enhanced Learning* 12, 4 (2020), 390–410. <https://doi.org/10.1504/ijtel.2020.110048>

- [14] Domenica Bagnato. 2022. Recommendation CM/REC(2017)5 of the Council of Europe and an Analysis of EVoting Protocols. In *Proceedings of the Central and Eastern European EDem and EGov Days (CEEeGov '22)*. Association for Computing Machinery, New York, NY, USA, 169–178. <https://doi.org/10.1145/3551504.3551519>
- [15] Andrea Barraza-Urbina and Dorota Glowacka. 2020. Introduction to Bandits in Recommender Systems. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys '20)*. Association for Computing Machinery, New York, NY, USA, 748–750. <https://doi.org/10.1145/3383313.3411547>
- [16] Reza Barzegar Nozari and Hamidreza Koochi. 2022. Novel Implicit-Trust-Network-Based Recommendation Methodology. *Expert Systems with Applications* 186, C (2022). <https://doi.org/10.1016/j.eswa.2021.115709>
- [17] Christine Bauer. 2020. Multi-Method Evaluation: Leveraging Multiple Methods to Answer What You Were Looking For. In *Proceedings of the 2020 Conference on Human Information Interaction and Retrieval (CHIIR '20)*. Association for Computing Machinery, New York, NY, USA, 472–474. <https://doi.org/10.1145/3343413.3378015>
- [18] Christine Bauer. 2021. Multi-Method Evaluation of Adaptive Systems. In *Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '21)*. Association for Computing Machinery, New York, NY, USA, 323–325. <https://doi.org/10.1145/3450613.3457122>
- [19] Yann Bayle, Matthias Robine, and Pierre Hanna. 2019. SATIN: A Persistent Musical Database for Music Information Retrieval and a Supporting Deep Learning Experiment on Song Instrumental Classification. *Multimedia Tools and Applications* 78, 3 (2019), 2703–2718. <https://doi.org/10.1007/s1042-018-5797-8>
- [20] Bettina Berendt, Özgür Karadeniz, Stefan Mertens, and Leen d’Haenens. 2021. Fairness beyond “Equal”: The Diversity Searcher as a Tool to Detect and Enhance the Representation of Socio-Political Actors in News Media. In *Companion Proceedings of the Web Conference 2021 (TheWebConf '21)*. Association for Computing Machinery, New York, NY, USA, 202–212. <https://doi.org/10.1145/3442442.3452303>
- [21] Cesare Bernardis, Maurizio Ferrari Dacrema, and Paolo Cremonesi. 2019. Estimating Confidence of Individual User Predictions in Item-Based Recommender Systems. In *Proceedings of the 27th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '19)*. Association for Computing Machinery, New York, NY, USA, 149–156. <https://doi.org/10.1145/3320435.3320453>
- [22] Peter Brusilovsky, Marco de Gemmis, Alexander Felfernig, Elisabeth Lex, Pasquale Lops, Giovanni Semeraro, and Martijn C. Willemsen. 2021. Joint Workshop on Interfaces and Human Decision Making for Recommender Systems (InRS'21). In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys '21)*. Association for Computing Machinery, New York, NY, USA, 783–786. <https://doi.org/10.1145/3460231.3470927>
- [23] Peter Brusilovsky, Marco de Gemmis, Alexander Felfernig, Pasquale Lops, John O’Donovan, Giovanni Semeraro, and Martijn C. Willemsen. 2019. RecSys '19 Joint Workshop on Interfaces and Human Decision Making for Recommender Systems. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys '19)*. Association for Computing Machinery, New York, NY, USA, 560–561. <https://doi.org/10.1145/3298689.3346971>
- [24] Peter Brusilovsky, Marco de Gemmis, Alexander Felfernig, Pasquale Lops, John O’Donovan, Giovanni Semeraro, and Martijn C. Willemsen. 2020. Interfaces and Human Decision Making for Recommender Systems. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys '20)*. Association for Computing Machinery, New York, NY, USA, 613–618. <https://doi.org/10.1145/3383313.3411539>
- [25] Peter Brusilovsky, Marco de Gemmis, Alexander Felfernig, Pasquale Lops, Marco Polignano, Giovanni Semeraro, and Martijn C. Willemsen. 2022. Joint Workshop on Interfaces and Human Decision Making for Recommender Systems (InRS'22). In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 667–670. <https://doi.org/10.1145/3523227.3547413>
- [26] Mark Bukowski, Sandra Geisler, Thomas Schmitz-Rode, and Robert Farkas. 2020. Feasibility of Activity-Based Expert Profiling Using Text Mining of Scientific Publications and Patents. *Scientometrics* 123, 2 (2020), 579–620. <https://doi.org/10.1007/s11192-020-03414-8>
- [27] Robin Douglas Burke, Masoud Mansoury, and Nasim Sonboli. 2020. Experimentation with Fairness-Aware Recommendation Using LibreC-Auto: Hands-on Tutorial. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (EAT* '20)*. Association for Computing Machinery, New York, NY, USA, 700. <https://doi.org/10.1145/3351095.3375670>
- [28] Rocío Cañamares and Pablo Castells. 2018. From the PRP to the Low Prior Discovery Recall Principle for Recommender Systems. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval (SIGIR '18)*. Association for Computing Machinery, New York, NY, USA, 1081–1084. <https://doi.org/10.1145/3209978.3210076>
- [29] Xingjuan Cai, Lijie Xie, Rui Tian, and Zhihua Cui. 2022. Explicable Recommendation Based on Knowledge Graph. *Expert Systems with Applications* 200, C (2022). <https://doi.org/10.1016/j.eswa.2022.117035>
- [30] Sergio Caro-Álvaro, Eva García-López, Antonio García-Cabot, Luis de Marcos, and Adrián Domínguez-Díaz. 2021. Applying Usability Recommendations When Developing Mobile Instant Messaging Applications. *IET Software* 16, 1 (2021), 73–93. <https://doi.org/10.1049/sfw.2.12039>
- [31] Hugo Caselles-Dupré, Florian Lesaint, and Jimena Royo-Letelier. 2018. Word2vec Applied to Recommendation: Hyperparameters Matter. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys '18)*. Association for Computing Machinery, New York, NY, USA, 352–356. <https://doi.org/10.1145/3240323.3240377>
- [32] Ilknur Celik, Ilaria Torre, Frosina Kocева, Christine Bauer, Eva Zangerle, and Bart Knijnenburg. 2018. UMAP 2018 Intelligent User-Adapted Interfaces: Design and Multi-Modal Evaluation (IUadaptMe) Workshop Chairs’ Welcome & Organization. In *Adjunct Publication of the 26th Conference on User Modeling, Adaptation and Personalization (UMAP '18)*. Association for Computing Machinery, New York, NY, USA, 137–139. <https://doi.org/10.1145/3213586.3226202>
- [33] Benjamin P. Chamberlain, Emanuele Rossi, Dan Shiebler, Suvash Sedhain, and Michael M. Bronstein. 2020. Tuning Word2vec for Large Scale Recommendation Systems. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys '20)*. Association for Computing Machinery, New York, NY, USA, 732–737. <https://doi.org/10.1145/3383313.3418486>

- [34] Dr. Swagato Chatterjee. 2019. Explaining Customer Ratings and Recommendations by Combining Qualitative and Quantitative User Generated Contents. *Decision Support Systems* 119, C (2019), 14–22. <https://doi.org/10.1016/j.dss.2019.02.008>
- [35] Ching-Wei Chen, Paul Lamere, Markus Schedl, and Hamed Zamani. 2018. Recsys Challenge 2018: Automatic Music Playlist Continuation. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys '18)*. Association for Computing Machinery, New York, NY, USA, 527–528. <https://doi.org/10.1145/3240323.3240342>
- [36] Marco Cinelli, Peter Burgherr, Milosz Kadziński, and Roman Slowiński. 2022. Proper and improper uses of MCDA methods in energy systems analysis. *Decision Support Systems* 163, C (2022). <https://doi.org/10.1016/j.dss.2022.113848>
- [37] Julia Clemente, Héctor Yago, Javier de Pedro-Carracedo, and Javier Bueno. 2022. A Proposal for an Adaptive Recommender System Based on Competences and Ontologies. *Expert Systems with Applications* 208, C (2022). <https://doi.org/10.1016/j.eswa.2022.118171>
- [38] Ludovik Coba, Panagiotis Symeonidis, and Markus Zanker. 2017. Visual Analysis of Recommendation Performance. In *Proceedings of the Eleventh ACM Conference on Recommender Systems (RecSys '17)*. Association for Computing Machinery, New York, NY, USA, 362–363. <https://doi.org/10.1145/3109859.3109982>
- [39] Ludovik Coba, Panagiotis Symeonidis, and Markus Zanker. 2018. Replicating and Improving Top-N Recommendations in Open Source Packages. In *Proceedings of the 8th International Conference on Web Intelligence, Mining and Semantics (WIMS '18)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3227609.3227671>
- [40] Fábio Colaço, Márcia Barros, and Francisco M. Couto. 2020. DRecPy: A Python Framework for Developing Deep Learning-Based Recommenders. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys 20)*. Association for Computing Machinery, New York, NY, USA, 675–680. <https://doi.org/10.1145/3383313.3418483>
- [41] David Contreras, Maria Salamó, and Ludovico Boratto. 2021. Integrating Collaboration and Leadership in Conversational Group Recommender Systems. *ACM Transactions on Information Systems* 39, 4 (2021). <https://doi.org/10.1145/3462759>
- [42] Federico Coró, Gianlorenzo D'angelo, and Yllka Velaj. 2021. Link Recommendation for Social Influence Maximization. *ACM Trans. Knowl. Discov. Data* 15, 6 (2021). <https://doi.org/10.1145/3449023>
- [43] Mateus Barcellos Costa and Dalila Tamzalit. 2017. Recommendation Patterns for Business Process Imperative Modeling. In *Proceedings of the Symposium on Applied Computing (SAC '17)*. Association for Computing Machinery, New York, NY, USA, 735–742. <https://doi.org/10.1145/3019612.3019619>
- [44] Robin Cugny, Julien Aligon, Max Chevalier, Geoffrey Roman Jimenez, and Olivier Teste. 2022. AutoXAI: A Framework to Automatically Select the Most Adapted XAI Solution. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management (CIKM '22)*. Association for Computing Machinery, New York, NY, USA, 315–324. <https://doi.org/10.1145/3511808.3557247>
- [45] Jipeng Cui, Chungang Yan, and Cheng Wang. 2021. A Credible Individual Behavior Profiling Method for Online Payment Fraud Detection. In *2021 4th International Conference on Data Storage and Data Engineering (DSDE '21)*. Association for Computing Machinery, New York, NY, USA, 22–30. <https://doi.org/10.1145/3456146.3456151>
- [46] Liang-Zhong Cui, Fu-Liang Guo, and Ying-jie Liang. 2018. Research Overview of Educational Recommender Systems. In *Proceedings of the 2nd International Conference on Computer Science and Application Engineering (CSAE '18)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3207677.3278071>
- [47] Yifei Da, Hailong Sun, and Xudong Liu. 2018. A Hybrid Approach to Developer Recommendation Based on Multi-Relationship. In *Proceedings of the 10th Asia-Pacific Symposium on Internetware (Internetware '18)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3275219.3275232>
- [48] Manoj Reddy Dareddy. 2017. Recommender Systems: Research Direction. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining (WSDM '17)*. Association for Computing Machinery, New York, NY, USA, 831. <https://doi.org/10.1145/3018661.3022748>
- [49] Zohreh Dehghani Champiri, Brian Fisher, Loo Chu Kiong, and Mahmoud Danaee. 2020. How Contextual Data Influences User Experience with Scholarly Recommender Systems: An Empirical Framework. In *HCI International 2020 - Late Breaking Papers: User Experience Design and Case Studies: 22nd HCI International Conference, HCII 2020, Copenhagen, Denmark, July 19–24, 2020, Proceedings (HCI International 2020)*. Springer, 635–661. https://doi.org/10.1007/978-3-030-60114-0_42
- [50] Yashar Deldjoo, Mihai Gabriel Constantin, Hamid Eghbal-Zadeh, Bogdan Ionescu, Markus Schedl, and Paolo Cremonesi. 2018. Audio-Visual Encoding of Multimedia Content for Enhancing Movie Recommendations. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys '18)*. Association for Computing Machinery, New York, NY, USA, 455–459. <https://doi.org/10.1145/3240323.3240407>
- [51] Yashar Deldjoo, Paolo Cremonesi, Markus Schedl, and Massimo Quadrana. 2017. The Effect of Different Video Summarization Models on the Quality of Video Recommendation Based on Low-Level Visual Features. In *Proceedings of the 15th International Workshop on Content-Based Multimedia Indexing (CBMI '17)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3095713.3095734>
- [52] Massimiliano Di Penta. 2021. How Empirical Research Supports Tool Development: A Retrospective Analysis and New Horizons. In *Proceedings of the 15th ACM / IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM) (ESEM '21)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3475716.3484488>
- [53] Juri Di Rocco, Davide Di Ruscio, Claudio Di Sipio, Phuong T. Nguyen, and Riccardo Rubei. 2021. Development of Recommendation Systems for Software Engineering: The CROSSMINER Experience. *Empirical Software Engineering* 26, 4 (2021). <https://doi.org/10.1007/s10664-021-09963-7>
- [54] Patrik Dokoupil and Ladislav Peska. 2022. Robustness Against Polarity Bias in Decoupled Group Recommendations Evaluation. In *Adjunct Proceedings of the 30th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '22 Adjunct)*. Association for Computing Machinery,

- New York, NY, USA, 302–307. <https://doi.org/10.1145/3511047.3537650>
- [55] Afsaneh Doryab, Victoria Bellotti, Alaaeddine Yousfi, Shuobi Wu, John M. Carroll, and Anind K. Dey. 2017. If It's Convenient: Leveraging Context in Peer-to-Peer Variable Service Transaction Recommendations. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 1, 3 (2017). <https://doi.org/10.1145/3130913>
- [56] Rotem Dror, Amir Kantor, Michael Masin, and Segev Shlomov. 2017. Set-SMAA for Finding Preferable Multi-Objective Solutions. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion (GECCO '17)*. Association for Computing Machinery, New York, NY, USA, 149–150. <https://doi.org/10.1145/3067695.3076005>
- [57] Michael D Ekstrand, Allison Chaney, Pablo Castells, Robin Burke, David Rohde, and Manel Slokom. 2021. SimuRec: Workshop on Synthetic Data and Simulation Methods for Recommender Systems Research. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys 21)*. Association for Computing Machinery, New York, NY, USA, 803–805. <https://doi.org/10.1145/3460231.3470938>
- [58] Michael D. Ekstrand and Daniel Kluger. 2021. Exploring Author Gender in Book Rating and Recommendation. *User Modeling and User-Adapted Interaction* 31, 3 (2021), 377–420. <https://doi.org/10.1007/s11257-020-09284-2>
- [59] Driss El Alaoui, Jamal Riffi, Badraddine Aghoutane, Abdelouahed Sabri, Ali Yahyaouy, and Hamid Tairi. 2020. Collaborative Filtering: Comparative Study Between Matrix Factorization and Neural Network Method. In *Networked Systems: 8th International Conference, NETYS 2020, Marrakech, Morocco, June 3–5, 2020, Proceedings (NETYS 2020)*. Springer, 361–367. https://doi.org/10.1007/978-3-030-67087-0_24
- [60] Mehdi Elahi, Himan Abdollahpour, Masoud Mansoury, and Helma Torkamaan. 2021. Beyond Algorithmic Fairness in Recommender Systems. In *Adjunct Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '21)*. Association for Computing Machinery, New York, NY, USA, 41–46. <https://doi.org/10.1145/3450614.3461685>
- [61] Jerry Alan Fails, Maria Soledad Pera, Natalia Kucirkova, and Franca Garzotto. 2018. International and Interdisciplinary Perspectives on Children & Recommender Systems (KidRec). In *Proceedings of the 17th ACM Conference on Interaction Design and Children (IDC '18)*. Association for Computing Machinery, New York, NY, USA, 705–712. <https://doi.org/10.1145/3202185.3205866>
- [62] Hui Fang, Danning Zhang, Yiheng Shu, and Guibing Guo. 2020. Deep Learning for Sequential Recommendation: Algorithms, Influential Factors, and Evaluations. *ACM Transactions on Information Systems* 39, 1 (2020). <https://doi.org/10.1145/3426723>
- [63] Nicolò Felicioni, Maurizio Ferrari Dacrema, and Paolo Cremonesi. 2021. Measuring the User Satisfaction in a Recommendation Interface with Multiple Carousels. In *ACM International Conference on Interactive Media Experiences (IMX '21)*. Association for Computing Machinery, New York, NY, USA, 212–217. <https://doi.org/10.1145/3452918.3465493>
- [64] Nicolò Felicioni, Maurizio Ferrari Dacrema, and Paolo Cremonesi. 2021. A Methodology for the Offline Evaluation of Recommender Systems in a User Interface with Multiple Carousels. In *Adjunct Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '21)*. Association for Computing Machinery, New York, NY, USA, 10–15. <https://doi.org/10.1145/3450614.3461680>
- [65] Liang Feng, Qianchuan Zhao, and Cangqi Zhou. 2020. Improving Performances of Top-N Recommendations with Co-Clustering Method. *Expert Systems with Applications* 143, C (2020). <https://doi.org/10.1016/j.eswa.2019.113078>
- [66] Andres Ferraro, Gustavo Ferreira, Fernando Diaz, and Georgina Born. 2022. Measuring Commonality in Recommendation of Cultural Content: Recommender Systems to Enhance Cultural Citizenship. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 567–572. <https://doi.org/10.1145/3523227.3551476>
- [67] Bruce Ferwerda, Mark P. Graus, Andreu Vall, Marko Tkaleic, and Markus Schedl. 2017. How Item Discovery Enabled by Diversity Leads to Increased Recommendation List Attractiveness. In *Proceedings of the Symposium on Applied Computing (SAC '17)*. Association for Computing Machinery, New York, NY, USA, 1693–1696. <https://doi.org/10.1145/3019612.3019899>
- [68] Cristhian Figueroa, Iacopo Vagliano, Oscar Rodriguez Rocha, Marco Torchiano, Catherine Faron Zucker, Juan Carlos Corrales, and Maurizio Morisio. 2017. Allied: A Framework for Executing Linked Data-Based Recommendation Algorithms. *International Journal on Semantic Web and Information Systems* 13, 4 (2017), 134–154. <https://doi.org/10.4018/IJSWIS.2017100107>
- [69] Medjeu Fopa, Modou Gueye, Samba Ndiaye, and Hubert Naacke. 2022. A Parameter-Free KNN for Rating Prediction. *Data & Knowledge Engineering* 142, C (2022). <https://doi.org/10.1016/j.datak.2022.102095>
- [70] Yannick Forster, Sebastian Hergeth, Frederik Naujoks, and Josef F. Krems. 2018. How Usability Can Save the Day - Methodological Considerations for Making Automated Driving a Success Story. In *Proceedings of the 10th International Conference on Automotive User Interfaces and Interactive Vehicular Applications (AutomotiveUI '18)*. Association for Computing Machinery, New York, NY, USA, 278–290. <https://doi.org/10.1145/3239060.3239076>
- [71] Mauricio Noris Freire and Leandro Nunes de Castro. 2021. E-Recruitment Recommender Systems: A Systematic Review. *Knowledge-Based Systems* 63, 1 (2021), 1–20. <https://doi.org/10.1007/s10115-020-01522-8>
- [72] Giovanni Gabbolini and Derek Bridge. 2022. A User-Centered Investigation of Personal Music Tours. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 25–34. <https://doi.org/10.1145/3523227.3546776>
- [73] Ruoyuan Gao and Chirag Shah. 2020. Counteracting Bias and Increasing Fairness in Search and Recommender Systems. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys 20)*. Association for Computing Machinery, New York, NY, USA, 745–747. <https://doi.org/10.1145/3383313.3411545>
- [74] Jean Garcia-Gathright, Christine Hosey, Brian St. Thomas, Ben Carterette, and Fernando Diaz. 2018. Mixed Methods for Evaluating User Satisfaction. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys '18)*. Association for Computing Machinery, New York, NY, USA, 541–542. <https://doi.org/10.1145/3240323.3241622>

- [75] Jean Garcia-Gathright, Brian St. Thomas, Christine Hosey, Zahra Nazari, and Fernando Diaz. 2018. Understanding and Evaluating User Satisfaction with Music Discovery. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval (SIGIR '18)*. Association for Computing Machinery, New York, NY, USA, 55–64. <https://doi.org/10.1145/3209978.3210049>
- [76] Alireza Gharahighehi, Celine Vens, and Konstantinos Pliakos. 2021. An Ensemble Hypergraph Learning Framework for Recommendation. In *Discovery Science: 24th International Conference, DS 2021, Halifax, NS, Canada, October 11–13, 2021, Proceedings (DS 2021)*. Springer, 295–304. https://doi.org/10.1007/978-3-030-88942-5_23
- [77] Anjali Goyal. 2017. Effective Bug Triage for Non Reproducible Bugs. In *Proceedings of the 39th International Conference on Software Engineering Companion (ICSE-C '17)*. IEEE Press, 487–488. <https://doi.org/10.1109/ICSE-C.2017.41>
- [78] Kazjon Grace, Elanor Finch, Natalia Gulbransen-Diaz, and Hamish Henderson. 2022. Q-Chef: The Impact of Surprise-Eliciting Systems on Food-Related Decision-Making. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems (CHI '22)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3491102.3501862>
- [79] Maria Guinea, Isabel Litton, Rigel Smiroldo, Irma Nitsche, and Eric Sax. 2021. A Proactive Context-Aware Recommender System for In-Vehicle Use. In *Proceedings of the 2020 4th International Conference on Vision, Image and Signal Processing (ICVISP 2020)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3448823.3448852>
- [80] Xiaohan Guo and Dejun Xie. 2021. Optimizing the Design of Recommendation Systems. In *Proceedings of the 9th International Conference on Computer and Communications Management (ICCCM '21)*. Association for Computing Machinery, New York, NY, USA, 79–84. <https://doi.org/10.1145/3479162.3479174>
- [81] Omar Haggag. 2022. Better Identifying and Addressing Diverse Issues in MHealth and Emerging Apps Using User Reviews. In *Proceedings of the 26th International Conference on Evaluation and Assessment in Software Engineering (EASE '22)*. Association for Computing Machinery, New York, NY, USA, 329–335. <https://doi.org/10.1145/3530019.3535346>
- [82] Omar Haggag, John Grundy, Mohamed Abdelrazek, and Sherif Haggag. 2022. Better Addressing Diverse Accessibility Issues in Emerging Apps: A Case Study Using COVID-19 Apps. In *Proceedings of the 9th IEEE/ACM International Conference on Mobile Software Engineering and Systems (MOBILESoft '22)*. Association for Computing Machinery, New York, NY, USA, 50–61. <https://doi.org/10.1145/3524613.3527817>
- [83] Soyeon Caren Han, Taejun Lim, Siqu Long, Bernd Burgstaller, and Josiah Poon. 2021. GLocal-K: Global and Local Kernels for Recommender Systems. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management (CIKM '21)*. Association for Computing Machinery, New York, NY, USA, 3063–3067. <https://doi.org/10.1145/3459637.3482112>
- [84] Yaojun Hao and Fuzhi Zhang. 2018. Detecting Shilling Profiles in Collaborative Recommender Systems via Multidimensional Profile Temporal Features. *IET Information Security* 12, 4 (2018), 362–374. <https://doi.org/10.1049/iet-ifs.2017.0012>
- [85] Iris Hendrickx, Federica Cena, Erkan Basar, Luigi Di Caro, Florian Kunne, Elena Musi, Cataldo Musto, Amon Rapp, and Jelte van Waterschoot. 2021. Towards a New Generation of Personalized Intelligent Conversational Agents. In *Adjunct Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '21)*. Association for Computing Machinery, New York, NY, USA, 373–374. <https://doi.org/10.1145/3450614.3461453>
- [86] Yuuki Hiraishi, Takayoshi Kitamura, Tomoko Izumi, and Yoshio Nakatani. 2018. Experimental Verification of Sightseeing Information as a Weak Trigger to Affect Tourist Behavior. In *Social Computing and Social Media. User Experience and Behavior: 10th International Conference, SCSSM 2018, Held as Part of HCI International 2018, Las Vegas, NV, USA, July 15–20, 2018, Proceedings, Part I (SCSSM 2018)*. Springer, 303–317. https://doi.org/10.1007/978-3-319-91521-0_22
- [87] Yan Hong, Xianyi Zeng, Yuyang Wang, Pascal Bruniaux, and Yan Chen. 2018. CBCRS: An open case-based color recommendation system. *Knowledge-Based Systems* 141, C (2018), 113–128. <https://doi.org/10.1016/j.knsys.2017.11.014>
- [88] Mo Houtti, Isaac Johnson, Joel Cepeda, Soumya Khandelwal, Aviral Bhatnagar, and Loren Terveen. 2022. “We Need a Woman in Music”: Exploring Wikipedia’s Values on Article Priority. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (2022). <https://doi.org/10.1145/3555156>
- [89] Xiao Hu, Jeremy T. D. Ng, Chengrui Yang, and Samuel K. W. Chu. 2020. Personalized Book Recommendation to Young Readers: Two Online Prototypes and A Preliminary User Evaluation. In *Proceedings of the ACM/IEEE Joint Conference on Digital Libraries in 2020 (JCDL '20)*. Association for Computing Machinery, New York, NY, USA, 413–416. <https://doi.org/10.1145/3383583.3398604>
- [90] Ying Huang, Hong-Yu Zhang, and Jian-Qiang Wang. 2018. A Comprehensive Mechanism for Hotel Recommendation to Achieve Personalized Search Engine. *Journal of Intelligent & Fuzzy Systems* 35, 3 (2018), 3733–3745. <https://doi.org/10.3233/JIFS-18547>
- [91] Ngozi Ihemelandu. 2022. Best Practices for Top-N Recommendation Evaluation: Candidate Set Sampling and Statistical Inference Techniques. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management (CIKM '22)*. Association for Computing Machinery, New York, NY, USA, 5120–5123. <https://doi.org/10.1145/3511808.3557816>
- [92] Kojiro Iizuka, Takeshi Yoneda, and Yoshifumi Seki. 2019. Greedy optimized multileaving for personalization. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys 19)*. Association for Computing Machinery, New York, NY, USA, 413–417. <https://doi.org/10.1145/3298689.3347008>
- [93] Folasade O. Isinkaye, Yetunde O. Folajimi, and Adesesan B. Adeyemo. 2020. On Collaborative Filtering Model Optimised with Multi-Item Attribute Information Space for Enhanced Recommendation Accuracy. *International Journal of Intelligent Systems Technologies and Applications* 19, 3 (2020), 207–215. <https://doi.org/10.1504/ijista.2020.108054>
- [94] Amir H. Jadidinejad, Craig Macdonald, and Iadh Ounis. 2020. Using Exploration to Alleviate Closed Loop Effects in Recommender Systems. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20)*. Association for

- Computing Machinery, New York, NY, USA, 2025–2028. <https://doi.org/10.1145/3397271.3401230>
- [95] Dietmar Jannach, Iman Kamehkhosh, and Lukas Lerche. 2017. Leveraging Multi-Dimensional User Models for Personalized next-Track Music Recommendation. In *Proceedings of the Symposium on Applied Computing (SAC '17)*. Association for Computing Machinery, New York, NY, USA, 1635–1642. <https://doi.org/10.1145/3019612.3019756>
- [96] Olivier Jeunen. 2019. Revisiting offline evaluation for implicit-feedback recommender systems. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys 19)*. Association for Computing Machinery, New York, NY, USA, 596–600. <https://doi.org/10.1145/3298689.3347069>
- [97] Olivier Jeunen, Thorsten Joachims, Harrie Oosterhuis, Yuta Saito, and Flavian Vasile. 2022. CONSEQUENCES—Causality, Counterfactuals and Sequential Decision-Making for Recommender Systems. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 654–657. <https://doi.org/10.1145/3523227.3547409>
- [98] Yitong Ji, Aixin Sun, Jie Zhang, and Chenliang Li. 2020. A Re-Visit of the Popularity Baseline in Recommender Systems. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20)*. Association for Computing Machinery, New York, NY, USA, 1749–1752. <https://doi.org/10.1145/3397271.3401233>
- [99] Wenrong Jiang. 2021. Research on the Recommendation Algorithm Based on 0-1 Knapsack Problem. In *The 2nd International Conference on Computing and Data Science (CONF-CDS 2021)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3448734.3450921>
- [100] Thorsten Joachims, Maria Dimakopoulou, Adith Swaminathan, Yves Raimond, Olivier Koch, and Flavian Vasile. 2019. REVEAL 2019: Closing the Loop with the Real World: Reinforcement and Robust Estimators for Recommendation. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys '19)*. Association for Computing Machinery, New York, NY, USA, 568–569. <https://doi.org/10.1145/3298689.3346975>
- [101] Thorsten Joachims, Yves Raimond, Olivier Koch, Maria Dimakopoulou, Flavian Vasile, and Adith Swaminathan. 2020. REVEAL 2020: Bandit and Reinforcement Learning from User Interactions. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys '20)*. Association for Computing Machinery, New York, NY, USA, 628–629. <https://doi.org/10.1145/3383313.3411536>
- [102] Thorsten Joachims, Adith Swaminathan, Yves Raimond, Olivier Koch, and Flavian Vasile. 2018. REVEAL 2018: Offline Evaluation for Recommender Systems. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys '18)*. Association for Computing Machinery, New York, NY, USA, 514–515. <https://doi.org/10.1145/3240323.3240334>
- [103] Michael Jugovac, Dietmar Jannach, and Mozghan Karimi. 2018. Streamingrec: A Framework for Benchmarking Stream-Based News Recommenders. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys '18)*. Association for Computing Machinery, New York, NY, USA, 269–273. <https://doi.org/10.1145/3240323.3240384>
- [104] Ajinkya Kale, Surya Kallumadi, Tracy Holloway King, Shervin Malmasi, Maarten de Rijke, and Jacopo Tagliabue. 2022. ECom'22: The SIGIR 2022 Workshop on ECommerce. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*. Association for Computing Machinery, New York, NY, USA, 3485–3487. <https://doi.org/10.1145/3477495.3531701>
- [105] Jie Kang, Kyle Condiff, Shuo Chang, Joseph A. Konstan, Loren Terveen, and F. Maxwell Harper. 2017. Understanding How People Use Natural Language to Ask for Recommendations. In *Proceedings of the Eleventh ACM Conference on Recommender Systems (RecSys '17)*. Association for Computing Machinery, New York, NY, USA, 229–237. <https://doi.org/10.1145/3109859.3109873>
- [106] Jingda Kang, Juntao Zhang, Wei Song, and Xiandi Yang. 2021. Friend Relationships Recommendation Algorithm in Online Education Platform. In *Web Information Systems and Applications: 18th International Conference, WISA 2021, Kaifeng, China, September 24–26, 2021, Proceedings (WISA 21)*. Springer, 592–604. https://doi.org/10.1007/978-3-030-87571-8_51
- [107] Zhao Kang, Chong Peng, and Qiang Cheng. 2017. Kernel-Driven Similarity Learning. *Neurocomputing* 267, C (2017), 210–219. <https://doi.org/10.1016/j.neucom.2017.06.005>
- [108] Ismail Emre Kartoglu and Michael W. Spratling. 2018. Two Collaborative Filtering Recommender Systems Based on Sparse Dictionary Coding. *Knowledge and Information Systems* 57, 3 (2018), 709–720. <https://doi.org/10.1007/s10115-018-1157-2>
- [109] Jüri Keller and Leon Paul Mondrian Munz. 2022. Evaluating Research Dataset Recommendations In A Living Lab. In *Experimental IR Meets Multilinguality, Multimodality, and Interaction: 13th International Conference of the CLEF Association, CLEF 2022, Bologna, Italy, September 5–8, 2022, Proceedings (CLEF 2022)*. Springer, 135–148. https://doi.org/10.1007/978-3-031-13643-6_11
- [110] Nasrullah Khan, Zongmin Ma, Aman Ullah, and Kemal Polat. 2022. Categorization of Knowledge Graph Based Recommendation Methods and Benchmark Datasets from the Perspectives of Application Scenarios: A Comprehensive Survey. *Expert Systems with Applications* 206, C (2022). <https://doi.org/10.1016/j.eswa.2022.117737>
- [111] Hyun Jeong Kim, So Yeon Park, Minju Park, and Kyogu Lee. 2020. Do Channels Matter? Illuminating Interpersonal Influence on Music Recommendations. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys '20)*. Association for Computing Machinery, New York, NY, USA, 663–668. <https://doi.org/10.1145/3383313.3418489>
- [112] Takuya Kitazawa and Makoto Yui. 2018. Query-Based Simple and Scalable Recommender Systems with Apache Hivemall. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys '18)*. Association for Computing Machinery, New York, NY, USA, 502–503. <https://doi.org/10.1145/3240323.3241592>
- [113] Alexandros Kokkalas, Athanasios T. Patenidis, Evangelos A. Stathopoulos, Eirini E. Mitsopoulou, Sotiris Diplaris, Konstadinos Papadopoulos, Stefanos Vrochidis, Konstantinos Votis, Dimitrios Tzovaras, and Ioannis Kompatsiaris. 2022. E-Tracer: A Smart, Personalized And Immersive Digital Tourist Software System. In *Information Integration and Web Intelligence: 24th International Conference, IiWAS 2022, Virtual Event, November 28–30, 2022, Proceedings (IiWAS 2022)*. Springer, 581–587. https://doi.org/10.1007/978-3-031-21047-1_53

- [114] Irina Kondratova and Bruno Emond. 2020. Voice Interaction for Training: Opportunities, Challenges, and Recommendations from HCI Perspective. In *Learning and Collaboration Technologies. Human and Technology Ecosystems: 7th International Conference, LCT 2020, Held as Part of the 22nd HCI International Conference, HCII 2020, Copenhagen, Denmark, July 19–24, 2020, Proceedings, Part II (HCII 2020)*. Springer, 59–75. https://doi.org/10.1007/978-3-030-50506-6_6
- [115] Denis Kotkov, Jari Veijalainen, and Shuaiqiang Wang. 2020. How Does Serendipity Affect Diversity in Recommender Systems? A Serendipity-Oriented Greedy Algorithm. *Computing* 102, 2 (2020), 393–411. <https://doi.org/10.1007/s00607-018-0687-5>
- [116] Dominik Kowald, Simone Kopeinik, and Elisabeth Lex. 2017. The TagRec Framework as a Toolkit for the Development of Tag-Based Recommender Systems. In *Adjunct Publication of the 25th Conference on User Modeling, Adaptation and Personalization (UMAP '17)*. Association for Computing Machinery, New York, NY, USA, 23–28. <https://doi.org/10.1145/3099023.3099069>
- [117] Dominik Kowald, Markus Schedl, and Elisabeth Lex. 2020. The Unfairness of Popularity Bias in Music Recommendation: A Reproducibility Study. In *Advances in Information Retrieval: 42nd European Conference on IR Research, ECIR 2020, Lisbon, Portugal, April 14–17, 2020, Proceedings, Part II (ECIR 2020)*. Springer, 35–42. https://doi.org/10.1007/978-3-030-45442-5_5
- [118] Dominik Kowald, Paul Seitlinger, Tobias Ley, and Elisabeth Lex. 2018. The Impact of Semantic Context Cues on the User Acceptance of Tag Recommendations: An Online Study. In *Companion Proceedings of the The Web Conference 2018 (TheWebConf '18)*. International World Wide Web Conferences Steering Committee, 1–2. <https://doi.org/10.1145/3184558.3186899>
- [119] Alexander Kuhnle, Miguel Aroca-Ouellette, Anindya Basu, Murat Sensoy, John Reid, and Dell Zhang. 2021. Reinforcement Learning for Information Retrieval. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*. Association for Computing Machinery, New York, NY, USA, 2669–2672. <https://doi.org/10.1145/3404835.3462813>
- [120] Gaurav Kumar and N. Parimala. 2020. A Weighted Sum Method MCDM Approach for Recommending Product Using Sentiment Analysis. *International Journal of Business Information Systems* 35, 2 (2020), 185–203. <https://doi.org/10.1504/ijbis.2020.110172>
- [121] Zarina Kurbatova, Ivan Veselov, Yaroslav Golubev, and Timofey Bryksin. 2020. Recommendation of Move Method Refactoring Using Path-Based Representation of Code. In *Proceedings of the IEEE/ACM 42nd International Conference on Software Engineering Workshops (ICSEW '20)*. Association for Computing Machinery, New York, NY, USA, 315–322. <https://doi.org/10.1145/3387940.3392191>
- [122] Emanuel Lacic, Markus Reiter-Haas, Tomislav Duricic, Valentin Slawicek, and Elisabeth Lex. 2019. Should We Embed? A Study on the Online Performance of Utilizing Embeddings for Real-Time Job Recommendations. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys '19)*. Association for Computing Machinery, New York, NY, USA, 496–500. <https://doi.org/10.1145/3298689.3346989>
- [123] Ayangleima Laishram and Vineet Padmanabhan. 2019. Discovery of User-Item Subgroups via Genetic Algorithm for Effective Prediction of Ratings in Collaborative Filtering. *Applied Intelligence* 49, 11 (2019), 3990–4006. <https://doi.org/10.1007/s10489-019-01495-4>
- [124] Gal Lavee, Noam Koenigstein, and Oren Barkan. 2019. When Actions Speak Louder than Clicks: A Combined Model of Purchase Probability and Long-Term Customer Satisfaction. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys '19)*. Association for Computing Machinery, New York, NY, USA, 287–295. <https://doi.org/10.1145/3298689.3347044>
- [125] Chang-Shing Lee, Mei-Hui Wang, Yung-Chang Hsiao, and Bing-Heng Tsai. 2019. Ontology-Based GFML Agent for Patent Technology Requirement Evaluation and Recommendation. *Soft Computing* 23, 2 (2019), 537–556. <https://doi.org/10.1007/s00500-017-2859-1>
- [126] Youfang Leng and Li Yu. 2022. Incorporating Global and Local Social Networks for Group Recommendations. *Pattern Recognition* 127, C (2022). <https://doi.org/10.1016/j.patcog.2022.108601>
- [127] Elisabeth Lex and Markus Schedl. 2022. Psychology-Informed Recommender Systems: A Human-Centric Perspective on Recommender Systems. In *Proceedings of the 2022 Conference on Human Information Interaction and Retrieval (CHIIR '22)*. Association for Computing Machinery, New York, NY, USA, 367–368. <https://doi.org/10.1145/3498366.3505841>
- [128] Yu Liang. 2019. Recommender System for Developing New Preferences and Goals. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys '19)*. Association for Computing Machinery, New York, NY, USA, 611–615. <https://doi.org/10.1145/3298689.3347054>
- [129] Yu Liang and Martijn C. Willemsen. 2019. Personalized Recommendations for Music Genre Exploration. In *Proceedings of the 27th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '19)*. Association for Computing Machinery, New York, NY, USA, 276–284. <https://doi.org/10.1145/3320435.3320455>
- [130] Sonal Linda, Sonajharia Minz, and K.K. Bharadwaj. 2019. Fuzzy-Genetic Approach to Context-Aware Recommender Systems Based on the Hybridization of Collaborative Filtering and Reclusive Method Techniques. *AI Communications* 32, 2 (2019), 125–141. <https://doi.org/10.3233/AIC-180593>
- [131] Fang Liu, Ge Li, Zhiyi Fu, Shuai Lu, Yiyang Hao, and Zhi Jin. 2022. Learning to Recommend Method Names with Global Context. In *Proceedings of the 44th International Conference on Software Engineering (ICSE '22)*. Association for Computing Machinery, New York, NY, USA, 1294–1306. <https://doi.org/10.1145/3510003.3510154>
- [132] Junrui Liu, Zhen Yang, Tong Li, Di Wu, and Ruiyi Wang. 2022. SPR: Similarity Pairwise Ranking for Personalized Recommendation. *Knowledge-Based Systems* 239, C (2022). <https://doi.org/10.1016/j.knosys.2021.107828>
- [133] Malte Ludewig, Noemi Mauro, Sara Latifi, and Dietmar Jannach. 2019. Performance Comparison of Neural and Non-Neural Approaches to Session-Based Recommendation. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys '19)*. Association for Computing Machinery, New York, NY, USA, 462–466. <https://doi.org/10.1145/3298689.3347041>
- [134] Malte Ludewig, Noemi Mauro, Sara Latifi, and Dietmar Jannach. 2021. Empirical Analysis of Session-Based Recommendation Algorithms: A Comparison of Neural and Non-Neural Approaches. *User Modeling and User-Adapted Interaction* 31, 1 (2021), 149–181. <https://doi.org/10.1007/s11257->

- 020-09277-1
- [135] Cheng Luo, Bo Zhang, Yang Xiang, and Man Qi. 2019. Gaussian-Gamma Collaborative Filtering: A Hierarchical Bayesian Model for Recommender Systems. *Journal of Computer Science and Technology* 102, C (2019), 42–56. <https://doi.org/10.1016/j.jcss.2017.03.007>
- [136] Shunan Ma, Xunbo Shuai, and Yongcai Chai. 2020. An Intelligent Optimization Method for Information Recommendation. In *Proceedings of the 16th EAI International Conference on Mobile and Ubiquitous Systems: Computing, Networking and Services (MobiQuitous '19)*. Association for Computing Machinery, New York, NY, USA, 110–115. <https://doi.org/10.1145/3360774.3360782>
- [137] Valentina Maccatrozzo, Manon Terstall, Lora Aroyo, and Guus Schreiber. 2017. SIRUP: Serendipity In Recommendations via User Perceptions. In *Proceedings of the 22nd International Conference on Intelligent User Interfaces (IUI '17)*. Association for Computing Machinery, New York, NY, USA, 35–44. <https://doi.org/10.1145/3025171.3025185>
- [138] Bernadetta Maleszka. 2019. A Framework for Research Publication Recommendation System. In *Computational Collective Intelligence: 11th International Conference, ICCCI 2019, Hedaye, France, September 4–6, 2019, Proceedings, Part I (ICCCI 2019)*. Springer, 167–178. https://doi.org/10.1007/978-3-030-28377-3_14
- [139] Francesca Mangili, Denis Broggin, and Alessandro Antonucci. 2020. Conversational Recommender System by Bayesian Methods. In *Scalable Uncertainty Management: 14th International Conference, SUM 2020, Bozen-Bolzano, Italy, September 23–25, 2020, Proceedings (SUM 2020)*. Springer, 200–214. https://doi.org/10.1007/978-3-030-58449-8_14
- [140] Ahtsham Manzoor and Dietmar Jannach. 2021. Generation-Based vs. Retrieval-Based Conversational Recommendation: A User-Centric Comparison. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys '21)*. Association for Computing Machinery, New York, NY, USA, 515–520. <https://doi.org/10.1145/3460231.3475942>
- [141] Dionisis Margaris, Dimitris Spiliotopoulos, Gregory Karagiorgos, Costas Vassilakis, and Dionysios Vasilopoulos. 2021. On Addressing the Low Rating Prediction Coverage in Sparse Datasets Using Virtual Ratings. *SN Computer Science* 2, 4 (2021). <https://doi.org/10.1007/s42979-021-00668-8>
- [142] Dionisis Margaris, Dimitris Spiliotopoulos, Dionysios Vasilopoulos, and Costas Vassilakis. 2021. A User Interface for Personalising WS-BPEL Scenarios. In *HCI in Business, Government and Organizations: 8th International Conference, HCIBGO 2021, Held as Part of the 23rd HCI International Conference, HCII 2021, Virtual Event, July 24–29, 2021, Proceedings (HCII 2021)*. Springer, 399–416. https://doi.org/10.1007/978-3-030-77750-0_25
- [143] Dionisis Margaris, Dimitris Spiliotopoulos, and Costas Vassilakis. 2022. Anchoring Effect Mitigation for Complex Recommender System Design. In *HCI International 2022 - Late Breaking Papers. Design, User Experience and Interaction: 24th International Conference on Human-Computer Interaction, HCII 2022, Virtual Event, June 26 – July 1, 2022, Proceedings (HCII 2022)*. Springer, 424–436. https://doi.org/10.1007/978-3-031-17615-9_29
- [144] Dionisis Margaris, Dimitris Spiliotopoulos, Costas Vassilakis, and Dionysios Vasilopoulos. 2020. Improving Collaborative Filtering’s Rating Prediction Accuracy by Introducing the Experiencing Period Criterion. *Neural Computing and Applications* 35, 1 (2020), 193–210. <https://doi.org/10.1007/s00521-020-05460-y>
- [145] Dionisis Margaris and Costas Vassilakis. 2017. Exploiting Internet of Things Information to Enhance Venues’ Recommendation Accuracy. *Service Oriented Computing and Applications* 11, 4 (2017), 393–409. <https://doi.org/10.1007/s11761-017-0216-y>
- [146] Monica M McGill, Rebecca Zarch, Stacey Sexton, Julie M. Smith, Christine Ong, Melissa Raspberry, and Shelly Hollis. 2021. Evaluating Computer Science Professional Development for Teachers in the United States. In *Proceedings of the 21st Koli Calling International Conference on Computing Education Research (Koli Calling '21)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3488042.3488054>
- [147] Rafael Ferreira Mello, Rodrigues Neto, Giuseppe Fiorentino, Gabriel Alves, Verenna Arêdes, João Victor Galdino Ferreira Silva, Taciana Pontual Falcão, and Dragan Gašević. 2022. Enhancing Instructors’ Capability To Assess Open-Response Using Natural Language Processing And Learning Analytics. In *Educating for a New Future: Making Sense of Technology-Enhanced Learning Adoption: 17th European Conference on Technology Enhanced Learning, EC-TEL 2022, Toulouse, France, September 12–16, 2022, Proceedings (EC-TEL 2022)*. Springer, 102–115. https://doi.org/10.1007/978-3-031-16290-9_8
- [148] Zaiqiao Meng, Richard McCreddie, Craig Macdonald, and Iadh Ounis. 2020. Exploring Data Splitting Strategies for the Evaluation of Recommendation Models. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys '20)*. Association for Computing Machinery, New York, NY, USA, 681–686. <https://doi.org/10.1145/3383313.3418479>
- [149] Rus M. Mesas and Alejandro Bellogín. 2017. Evaluating Decision-Aware Recommender Systems. In *Proceedings of the Eleventh ACM Conference on Recommender Systems (RecSys 17)*. Association for Computing Machinery, New York, NY, USA, 74–78. <https://doi.org/10.1145/3109859.3109888>
- [150] Lien Michiels, Robin Verachtert, and Bart Goethals. 2022. RecPack: An(Other) Experimentation Toolkit for Top-N Recommendation Using Implicit Feedback Data. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 648–651. <https://doi.org/10.1145/3523227.3551472>
- [151] Nitin Mishra, Vimal Mishra, and Saumya Chaturvedi. 2017. Tools and Techniques for Solving Cold Start Recommendation. In *Proceedings of the 1st International Conference on Internet of Things and Machine Learning (IML '17)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3109761.3109772>
- [152] Dongmei Mo, Xingxing Zou, and WaiKeung Wong. 2022. Neural Stylist: Towards Online Styling Service. *Expert Systems with Applications* 203, C (2022). <https://doi.org/10.1016/j.eswa.2022.117333>
- [153] Sandy Moens, Olivier Jeunen, and Bart Goethals. 2019. Interactive Evaluation of Recommender Systems with SNIPER: An Episode Mining Approach. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys '19)*. Association for Computing Machinery, New York, NY, USA, 538–539. <https://doi.org/10.1145/3298689.3346965>

- [154] Adrien Mogenet, Tuan Anh Nguyen Pham, Masahiro Kazama, and Jialin Kong. 2019. Predicting online performance of job recommender systems with offline evaluation. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys 19)*. Association for Computing Machinery, New York, NY, USA, 477–480. <https://doi.org/10.1145/3298689.3347032>
- [155] Mojtaba Mohammadian, Yahya Forghani, and Masood Niazi Torshiz. 2021. An Initialization Method to Improve the Training Time of Matrix Factorization Algorithm for Fast Recommendation. *Soft Computing* 25, 5 (2021), 3975–3987. <https://doi.org/10.1007/s00500-020-05419-0>
- [156] Pedro Ramaciotti Morales, Lionel Tabourier, and Raphaël Fournier-S’Niehotta. 2021. Testing the Impact of Semantics and Structure on Recommendation Accuracy and Diversity. In *Proceedings of the 12th IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM ’20)*. IEEE Press, 250–257. <https://doi.org/10.1109/ASONAM49781.2020.9381334>
- [157] J. Sharon Moses and L.D. Dhinesh Babu. 2018. Evaluating Prediction Accuracy, Developmental Challenges, and Issues of Recommender Systems. *International Journal of Web Portals* 10, 2 (2018), 61–79. <https://doi.org/10.4018/IJWP.2018070105>
- [158] Juergen Mueller. 2017. Combining Aspects of Genetic Algorithms with Weighted Recommender Hybridization. In *Proceedings of the 19th International Conference on Information Integration and Web-Based Applications & Services (iiWAS ’17)*. Association for Computing Machinery, New York, NY, USA, 13–22. <https://doi.org/10.1145/3151759.3151765>
- [159] Peter Muellner, Dominik Kowald, and Elisabeth Lex. 2021. Robustness of Meta Matrix Factorization Against Strict Privacy Constraints. In *Advances in Information Retrieval: 43rd European Conference on IR Research, ECIR 2021, Virtual Event, March 28 – April 1, 2021, Proceedings, Part II (ECIR 2021)*. Springer, 107–119. https://doi.org/10.1007/978-3-030-72240-1_8
- [160] Yutaka Nagashima and Yilun He. 2018. PaMpeR: Proof Method Recommendation System for Isabelle/HOL. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering (ASE ’18)*. Association for Computing Machinery, New York, NY, USA, 362–372. <https://doi.org/10.1145/3238147.3238210>
- [161] Julia Neidhardt, Wolfgang Wörndl, Tsvi Kuflik, Dmitri Goldenberg, and Markus Zanker. 2022. Workshop on Recommenders in Tourism (RecTour). In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys ’22)*. Association for Computing Machinery, New York, NY, USA, 678–679. <https://doi.org/10.1145/3523227.3547416>
- [162] Julia Neidhardt, Wolfgang Wörndl, Tsvi Kuflik, and Markus Zanker. 2018. ACM Recsys Workshop on Recommenders in Tourism (RecTour 2018). In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys ’18)*. Association for Computing Machinery, New York, NY, USA, 525–526. <https://doi.org/10.1145/3240323.3240341>
- [163] Julia Neidhardt, Wolfgang Wörndl, Tsvi Kuflik, and Markus Zanker. 2021. Workshop on Recommenders in Tourism (RecTour). In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys ’21)*. Association for Computing Machinery, New York, NY, USA, 815–816. <https://doi.org/10.1145/3460231.3470930>
- [164] Nicola Neophytou, Bhaskar Mitra, and Catherine Stinson. 2022. Revisiting Popularity And Demographic Biases in Recommender Evaluation and Effectiveness. In *Advances in Information Retrieval: 44th European Conference on IR Research, ECIR 2022, Stavanger, Norway, April 10–14, 2022, Proceedings, Part I (ECIR 2022)*. Springer, 641–654. https://doi.org/10.1007/978-3-030-99736-6_43
- [165] Angelo B. Neves, Rodrigo G. G. de Oliveira, Luiz André P. Paes Leme, Giseli Rabello Lopes, Bernardo P. Nunes, and Marco A. Casanova. 2018. Empirical Analysis of Ranking Models for an Adaptable Dataset Search. In *The Semantic Web: 15th International Conference, ESWC 2018, Heraklion, Crete, Greece, June 3–7, 2018, Proceedings (ESWC 2018)*. Springer, 50–64. https://doi.org/10.1007/978-3-319-93417-4_4
- [166] Athanasios N. Nikolakopoulos, Vassilis Kalantzis, Efstratios Gallopoulos, and John D. Garofalakis. 2019. EigenRec: Generalizing PureSVD for Effective and Efficient Top-N Recommendations. *Knowledge-Based Systems* 58, 1 (2019), 59–81. <https://doi.org/10.1007/s10115-018-1197-7>
- [167] Evgenii Novozhilov, Ivan Veselov, Mikhail Pravilov, and Timofey Bryksin. 2019. Evaluation of Move Method Refactorings Recommendation Algorithms: Are We Doing It Right?. In *Proceedings of the 3rd International Workshop on Refactoring (IWOR ’19)*. IEEE Press, 23–26. <https://doi.org/10.1109/IWoR.2019.00012>
- [168] Bolanle Adefowoke Ojokoh, Oluwatosin Olatunbosun Aboluje, and Tobore Igbe. 2020. A Collaborative Content-Based Movie Recommender System. *International Journal of Business Intelligence and Data Mining* 17, 3 (2020), 298–320. <https://doi.org/10.1504/ijbidm.2020.109293>
- [169] Kyle Ong, Su-Cheng Haw, and Kok-Why Ng. 2020. Deep Learning Based-Recommendation System: An Overview on Models, Datasets, Evaluation Metrics, and Future Trends. In *Proceedings of the 2019 2nd International Conference on Computational Intelligence and Intelligent Systems (CIIS ’19)*. Association for Computing Machinery, New York, NY, USA, 6–11. <https://doi.org/10.1145/3372422.3372444>
- [170] A. K. Ospanova and T. Y. Shelestova. 2020. Examining Students Views on the Implementation of the Two-Way Immersion Program in the Educational Process of Kazakhstan. In *Proceedings of the 6th International Conference on Engineering & MIS 2020 (ICEMIS ’20)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3410352.3410757>
- [171] Wei Ou, Entao Luo, Zhiyuan Tan, Lihong Xiang, Qin Yi, and Chen Tian. 2019. A Multi-Attributes-Based Trust Model of Internet of Vehicle. In *Network and System Security: 13th International Conference, NSS 2019, Sapporo, Japan, December 15–18, 2019, Proceedings (NSS 2019)*. Springer, 706–713. https://doi.org/10.1007/978-3-030-36938-5_45
- [172] Róbert Pálovics, Domokos Kelen, and András A. Benczúr. 2017. Tutorial on Open Source Online Learning Recommenders. In *Proceedings of the Eleventh ACM Conference on Recommender Systems (RecSys ’17)*. Association for Computing Machinery, New York, NY, USA, 400–401. <https://doi.org/10.1145/3109859.3109937>
- [173] Mamta Pandey, Ratnesh Litoriya, and Prateek Pandey. 2019. Novel Approach for Mobile Based App Development Incorporating MAAF. *Wireless Personal Communications* 107, 4 (2019), 1687–1708. <https://doi.org/10.1007/s11277-019-06351-9>

- [174] Dimitris Paraschakis and Bengt J. Nilsson. 2020. FlowRec: Prototyping Session-Based Recommender Systems in Streaming Mode. In *Advances in Knowledge Discovery and Data Mining: 24th Pacific-Asia Conference, PAKDD 2020, Singapore, May 11–14, 2020, Proceedings, Part I (PAKDD 2020)*. Springer, 65–77. https://doi.org/10.1007/978-3-030-47426-3_6
- [175] Rafael Parizi, Marina Moreira, Igor Couto, Sabrina Marczak, and Tayana Conte. 2021. A Design Thinking Techniques Recommendation Tool: An Initial and On-Going Proposal. In *Proceedings of the XIX Brazilian Symposium on Software Quality (SBQS '20)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3439961.3439997>
- [176] Joonseok Park, Ungsoo Kim, Donggyu Yun, and Keunhyuk Yeom. 2019. C-RCE: An Approach for Constructing and Managing a Cloud Service Broker. *Journal of Grid Computing* 17, 1 (2019), 137–168. <https://doi.org/10.1007/s10723-017-9422-2>
- [177] Iman Paryudi, Ahmad Ashari, and A. Min Tjoa. 2020. Personality Estimation Using Demographic Data in a Personality-Based Recommender System: A Proposal. In *Proceedings of the 21st International Conference on Information Integration and Web-Based Applications & Services (iiWAS 2019)*. Association for Computing Machinery, New York, NY, USA, 156–160. <https://doi.org/10.1145/3366030.3366098>
- [178] Iulia Paun. 2020. Efficiency-Effectiveness Trade-Offs in Recommendation Systems. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys '20)*. Association for Computing Machinery, New York, NY, USA, 770–775. <https://doi.org/10.1145/3383313.3411452>
- [179] Roberto Pellegrini, Wenjie Zhao, and Iain Murray. 2022. Don't recommend the obvious: estimate probability ratios. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 188–197. <https://doi.org/10.1145/3523227.3546753>
- [180] Zhenhui Peng, Jeehoon Yoo, Meng Xia, Sunghun Kim, and Xiaojuan Ma. 2018. Exploring How Software Developers Work with Mention Bot in GitHub. In *Proceedings of the Sixth International Symposium of Chinese CHI (ChineseCHI '18)*. Association for Computing Machinery, New York, NY, USA, 152–155. <https://doi.org/10.1145/3202667.3202694>
- [181] Mijael R. Bueno Perez, Elmar Eisemann, and Rafael Bidarra. 2021. A Synset-Based Recommender Method For Mixed-Initiative Narrative World Creation. In *Interactive Storytelling: 14th International Conference on Interactive Digital Storytelling, ICIDS 2021, Tallinn, Estonia, December 7–10, 2021, Proceedings (ICIDS 2021)*. Springer, 13–28. https://doi.org/10.1007/978-3-030-92300-6_2
- [182] Fernando Benjamin Pérez Maurera, Maurizio Ferrari Dacrema, and Paolo Cremonesi. 2022. An Evaluation Study of Generative Adversarial Networks for Collaborative Filtering. In *Advances in Information Retrieval: 44th European Conference on IR Research, ECIR 2022, Stavanger, Norway, April 10–14, 2022, Proceedings, Part I (ECIR 2022)*. Springer, 671–685. https://doi.org/10.1007/978-3-030-99736-6_45
- [183] Fernando Benjamin Pérez Maurera, Maurizio Ferrari Dacrema, and Paolo Cremonesi. 2022. Towards the Evaluation of Recommender Systems with Impressions. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 610–615. <https://doi.org/10.1145/3523227.3551483>
- [184] Ladislav Peska and Stepan Balcar. 2022. The Effect of Feedback Granularity on Recommender Systems Performance. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 586–591. <https://doi.org/10.1145/3523227.3551479>
- [185] Mateus Pinheiro, Nayana Carneiro, and Ticiane Darin. 2017. Recommendations for the Design of Urban Mobility Applications Based on the Study of the User Experience. In *Proceedings of the XVI Brazilian Symposium on Human Factors in Computing Systems (IHC 2017)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3160504.3160517>
- [186] Alessandro Piscopo, Oana Inel, Sanne Vrijenhoek, Martijn Millemcamp, and Krisztian Balog. 2022. QUARE: 1st Workshop on Measuring the Quality of Explanations in Recommender Systems. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*. Association for Computing Machinery, New York, NY, USA, 3478–3481. <https://doi.org/10.1145/3477495.3531699>
- [187] Caroline Pitt, Adam Bell, Edgar Onofre, and Katie Davis. 2019. A Badge, Not a Barrier: Designing for-and Throughout-Digital Badge Implementation. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI '19)*. Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3290605.3300920>
- [188] Emilie Poirson and Catherine Da Cunha. 2019. A Recommender Approach Based on Customer Emotions. *Expert Systems with Applications* 122, C (2019), 281–288. <https://doi.org/10.1016/j.eswa.2018.12.035>
- [189] Nikolaos Polatidis, Stelios Kapetanakis, Elias Pimenidis, and Yannis Manolopoulos. 2022. Fast and Accurate Evaluation of Collaborative Filtering Recommendation Algorithms. In *Intelligent Information and Database Systems: 14th Asian Conference, ACIIDS 2022, Ho Chi Minh City, Vietnam, November 28–30, 2022, Proceedings, Part I (ACIIDS 2022)*. Springer, 623–634. https://doi.org/10.1007/978-3-031-21743-2_50
- [190] Nikolaos Polatidis, Antonios Papaleonidas, Elias Pimenidis, and Lazaros Iliadis. 2020. An Explanation-Based Approach for Experiment Reproducibility in Recommender Systems. *Neural Computing and Applications* 32, 16 (2020), 12259–12266. <https://doi.org/10.1007/s00521-019-04274-x>
- [191] Mirko Polato and Fabio Aiolli. 2018. Boolean Kernels for Collaborative Filtering in Top-N Item Recommendation. *Neurocomputing* 286, C (2018), 214–225. <https://doi.org/10.1016/j.neucom.2018.01.057>
- [192] Gabriele Prato, Federico Sallemi, Paolo Cremonesi, Mario Scriminaci, Stefan Gudmundsson, and Silvio Palumbo. 2020. Outfit Completion and Clothes Recommendation. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems (CHI EA '20)*. Association for Computing Machinery, New York, NY, USA, 1–7. <https://doi.org/10.1145/3334480.3383076>
- [193] Jingjing Qiao and Li Wang. 2022. Modeling User Micro-Behaviors and Original Interest via Adaptive Multi-Attention Network for Session-Based Recommendation. *Knowledge-Based Systems* 244, C (2022). <https://doi.org/10.1016/j.knsys.2022.108567>

- [194] Jing Qin. 2021. Long-Tail Recommendation Framework Using Frequent Neighbors. In *Proceedings of the 2020 8th International Conference on Information Technology: IoT and Smart City (ICIT '20)*. Association for Computing Machinery, New York, NY, USA, 6–12. <https://doi.org/10.1145/3446999.3447001>
- [195] Ping Qiu and Zhendong Niu. 2021. TCIC_FS: Total Correlation Information Coefficient-Based Feature Selection Method for High-Dimensional Data. *Knowledge-Based Systems* 231, C (2021). <https://doi.org/10.1016/j.knosys.2021.107418>
- [196] Massimo Quadrana, Paolo Cremonesi, and Dietmar Jannach. 2018. Sequence-Aware Recommender Systems. *Comput. Surveys* 51, 4 (2018). <https://doi.org/10.1145/3190616>
- [197] Hossein A. Rahmani, Yashar Deldjoo, and Tommaso di Noia. 2022. The Role of Context Fusion on Accuracy, beyond-Accuracy, and Fairness of Point-of-Interest Recommendation Systems. *Expert Systems with Applications* 205, C (2022). <https://doi.org/10.1016/j.eswa.2022.117700>
- [198] Jérémie Rappaz, Julian McAuley, and Karl Aberer. 2021. Recommendation on Live-Streaming Platforms: Dynamic Availability and Repeat Consumption. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys '21)*. Association for Computing Machinery, New York, NY, USA, 390–399. <https://doi.org/10.1145/3460231.3474267>
- [199] Shaina Raza and Chen Ding. 2022. News Recommender System: A Review of Recent Progress, Challenges, and Opportunities. *Artificial Intelligence Review* 55, 1 (2022), 749–800. <https://doi.org/10.1007/s10462-021-10043-x>
- [200] Francesco Ricci, David Massimo, and Antonella De Angeli. 2021. Challenges for Recommender Systems Evaluation. In *CHIItaly 2021: 14th Biannual Conference of the Italian SIGCHI Chapter (CHIItaly '21)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3464385.3464733>
- [201] João Sá, Vanessa Queiroz Marinho, Ana Rita Magalhães, Tiago Lacerda, and Diogo Goncalves. 2022. Diversity Vs Relevance: A Practical Multi-Objective Study in Luxury Fashion Recommendations. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*. Association for Computing Machinery, New York, NY, USA, 2405–2409. <https://doi.org/10.1145/3477495.3531866>
- [202] Naveen Sachdeva and Julian McAuley. 2020. How Useful Are Reviews for Recommendation? A Critical Review and Potential Improvements. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '20)*. Association for Computing Machinery, New York, NY, USA, 1845–1848. <https://doi.org/10.1145/3397271.3401281>
- [203] Yuta Saito and Thorsten Joachims. 2021. Counterfactual Learning and Evaluation for Recommender Systems: Foundations, Implementations, and Recent Advances. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys '21)*. Association for Computing Machinery, New York, NY, USA, 828–830. <https://doi.org/10.1145/3460231.3473320>
- [204] Yuta Saito and Thorsten Joachims. 2022. Counterfactual Evaluation and Learning for Interactive Systems: Foundations, Implementations, and Recent Advances. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '22)*. Association for Computing Machinery, New York, NY, USA, 4824–4825. <https://doi.org/10.1145/3534678.3542601>
- [205] Pablo Sánchez. 2019. Exploiting Contextual Information for Recommender Systems Oriented to Tourism. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys '19)*. Association for Computing Machinery, New York, NY, USA, 601–605. <https://doi.org/10.1145/3298689.3347062>
- [206] Pablo Sánchez and Alejandro Bellogin. 2021. On the Effects of Aggregation Strategies for Different Groups of Users in Venue Recommendation. *Information Processing & Management* 58, 5 (2021). <https://doi.org/10.1016/j.ipm.2021.102609>
- [207] Pablo Sánchez and Linus W. Dietz. 2022. Travelers vs. Locals: The Effect of Cluster Analysis in Point-of-Interest Recommendation. In *Proceedings of the 30th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '22)*. Association for Computing Machinery, New York, NY, USA, 132–142. <https://doi.org/10.1145/3503252.3531320>
- [208] Javier Sanz-Cruzado and Pablo Castells. 2018. Enhancing Structural Diversity in Social Networks by Recommending Weak Ties. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys '18)*. Association for Computing Machinery, New York, NY, USA, 233–241. <https://doi.org/10.1145/3240323.3240371>
- [209] Javier Sanz-Cruzado and Pablo Castells. 2022. RELISON: A Framework for Link Recommendation in Social Networks. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*. Association for Computing Machinery, New York, NY, USA, 2992–3002. <https://doi.org/10.1145/3477495.3531730>
- [210] Javier Sanz-Cruzado, Sofia M. Pepa, and Pablo Castells. 2018. Structural Novelty and Diversity in Link Prediction. In *Companion Proceedings of the The Web Conference 2018 (TheWebConf '18)*. International World Wide Web Conferences Steering Committee, 1347–1351. <https://doi.org/10.1145/3184558.3191576>
- [211] Oren Sar Shalom, Dietmar Jannach, and Ido Guy. 2019. First Workshop on the Impact of Recommender Systems at ACM RecSys 2019. In *Proceedings of the 13th ACM Conference on Recommender Systems (RecSys '19)*. Association for Computing Machinery, New York, NY, USA, 556–557. <https://doi.org/10.1145/3298689.3347060>
- [212] Oren Sar Shalom, Dietmar Jannach, and Joseph A. Konstan. 2020. Second Workshop on the Impact of Recommender Systems at ACM RecSys '20. In *Proceedings of the 14th ACM Conference on Recommender Systems (RecSys '20)*. Association for Computing Machinery, New York, NY, USA, 630–631. <https://doi.org/10.1145/3383313.3411471>
- [213] Jun Sasaki, Shuang Li, and Enrique Herrera-Viedma. 2019. A Classification Method of Photos in a Tourism Website by Color Analysis. In *Advances and Trends in Artificial Intelligence. From Theory to Practice: 32nd International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems, IEA/AIE 2019, Graz, Austria, July 9–11, 2019, Proceedings (IEA/AIE 2019)*. Springer, 265–278. https://doi.org/10.1007/978-3-030-22999-3_24

- [214] Masahiro Sato, Budrul Ahsan, Koki Nagatani, Takashi Sonoda, Qian Zhang, and Tomoko Ohkuma. 2018. Explaining Recommendations Using Contexts. In *23rd International Conference on Intelligent User Interfaces (IUI '18)*. Association for Computing Machinery, New York, NY, USA, 659–664. <https://doi.org/10.1145/3172944.3173012>
- [215] Masahiro Sato, Koki Nagatani, and Takuji Tahara. 2017. Exploring an Optimal Online Model for New Job Recommendation: Solution for RecSys Challenge 2017. In *Proceedings of the Recommender Systems Challenge 2017 (RecSys Challenge '17)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3124791.3124797>
- [216] James Schaffer, John O'Donovan, and Tobias Höllerer. 2018. Easy to Please: Separating User Experience from Choice Satisfaction. In *Proceedings of the 26th Conference on User Modeling, Adaptation and Personalization (UMAP '18)*. Association for Computing Machinery, New York, NY, USA, 177–185. <https://doi.org/10.1145/3209219.3209222>
- [217] Malte Schwarzer, Corinna Breiter, Moritz Schubotz, Norman Meuschke, and Bela Gipp. 2017. Citolytics: A Link-Based Recommender System for Wikipedia. In *Proceedings of the Eleventh ACM Conference on Recommender Systems (RecSys '17)*. Association for Computing Machinery, New York, NY, USA, 360–361. <https://doi.org/10.1145/3109859.3109981>
- [218] Procheta Sen, Debasis Ganguly, and Gareth Jones. 2018. Procrastination is the Thief of Time: Evaluating the Effectiveness of Proactive Search Systems. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval (SIGIR '18)*. Association for Computing Machinery, New York, NY, USA, 1157–1160. <https://doi.org/10.1145/3209978.3210114>
- [219] Young-Duk Seo, Young-Gab Kim, Euijong Lee, and Hyungjin Kim. 2021. Group Recommender System Based on Genre Preference Focusing on Reducing the Clustering Cost. *Expert Systems with Applications* 183, C (2021). <https://doi.org/10.1016/j.eswa.2021.115396>
- [220] Young-Duk Seo, Young-Gab Kim, Euijong Lee, Kwang-Soo Seol, and Doo-Kwon Baik. 2018. An Enhanced Aggregation Method Considering Deviations for a Group Recommendation. *Expert Systems with Applications* 93, C (2018), 299–312. <https://doi.org/10.1016/j.eswa.2017.10.027>
- [221] Aravind Sesagiri Raamkumar and Schubert Foo. 2018. Multi-Method Evaluation in Scientific Paper Recommender Systems. In *Adjunct Publication of the 26th Conference on User Modeling, Adaptation and Personalization (UMAP '18)*. Association for Computing Machinery, New York, NY, USA, 179–182. <https://doi.org/10.1145/3213586.3226215>
- [222] Chhavi Sharma, Punam Bedi, Sabu M. Thampi, and El-Sayed M. El-Alfy. 2017. CCFRS – Community based Collaborative Filtering Recommender System. *Journal of Intelligent & Fuzzy Systems* 32, 4 (2017), 2987–2995. <https://doi.org/10.3233/jifs-169242>
- [223] Mengmeng Shen, Shiwei Zhao, Jun Wang, and Liting Ding. 2021. A Review Expert Recommendation Method Based on Comprehensive Evaluation in Multi-Source Data. In *Proceedings of the 5th International Conference on Control Engineering and Artificial Intelligence (CCEAI '21)*. Association for Computing Machinery, New York, NY, USA, 36–40. <https://doi.org/10.1145/3448218.3448236>
- [224] Jiangbo Shu, Xiaoxuan Shen, Hai Liu, Baolin Yi, and Zhaoli Zhang. 2018. A Content-Based Recommendation Algorithm for Learning Resources. *Multimedia Systems* 24, 2 (2018), 163–173. <https://doi.org/10.1007/s00530-017-0539-8>
- [225] Nelson Silva, Tobias Schreck, Eduardo Veas, Vedran Sabol, Eva Eggeling, and Dieter W. Fellner. 2018. Leveraging Eye-Gaze and Time-Series Features to Predict User Interests and Build a Recommendation Model for Visual Analysis. In *Proceedings of the 2018 ACM Symposium on Eye Tracking Research & Applications (ETRA '18)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3204493.3204546>
- [226] Viridiana Silva-Rodríguez, Sandra Edith Nava-Muñoz, Luis A Castro, Francisco E Martínez-Pérez, Héctor G Pérez-González, and Francisco Torres-Reyes. 2020. Classifying Design-Level Requirements Using Machine Learning for a Recommender of Interaction Design Patterns. *IET Software* 14, 5 (2020), 544–552. <https://doi.org/10.1049/iet-sen.2019.0291>
- [227] Thiago Silveira, Leonardo Rocha, Fernando Mourão, and Marcos Gonçalves. 2017. A Framework for Unexpectedness Evaluation in Recommendation. In *Proceedings of the Symposium on Applied Computing (SAC '17)*. Association for Computing Machinery, New York, NY, USA, 1662–1667. <https://doi.org/10.1145/3019612.3019760>
- [228] Bam Bahadur Sinha and R. Dhanalakshmi. 2021. Building a Fuzzy Logic-Based McCulloch-Pitts Neuron Recommendation Model to Uplift Accuracy. *Journal of Supercomputing* 77, 3 (2021), 2251–2267. <https://doi.org/10.1007/s11227-020-03344-5>
- [229] Manel Slokom, Alan Hanjalic, and Martha Larson. 2021. Towards user-oriented privacy for recommender system data: A personalization-based approach to gender obfuscation for user profiles. *Information Processing & Management* 58, 6, Article 102722 (2021). <https://doi.org/10.1016/j.ipm.2021.102722>
- [230] A. V. Smirnov and A. V. Ponomarev. 2020. Multicriteria Context-Driven Recommender Systems: Model and Method. *Scientific and Technical Information Processing* 47, 5 (2020), 298–303. <https://doi.org/10.3103/S014768822005007X>
- [231] Bolin Song, Xiaoyu Wang, Peihan Li, Peng Sun, and Azzedine Boukerche. 2022. A Novel Machine Learning-Assisted Policy Recommendation Method on COVID-19 Vaccination Campaign. In *Proceedings of the 2021 IEEE/ACM 25th International Symposium on Distributed Simulation and Real Time Applications (DS-RT '21)*. IEEE Press. <https://doi.org/10.1109/DS-RT52167.2021.9576138>
- [232] Elizaveta Stavinova, Andrey Gurov, Anton Lysenko, and Petr Chunaev. 2022. Performance Ranking of Recommender Systems on Simulated Data. *Procedia Computer Science* 212, C (2022), 142–151. <https://doi.org/10.1016/j.procs.2022.10.216>
- [233] Katherapaka Sudhakar, M.A.H. Farquad, and G. Narshimha. 2019. Effective Convolution Method for Privacy Preserving in Cloud over Big Data Using Map Reduce Framework. *IET Software* 13, 3 (2019), 187–194. <https://doi.org/10.1049/iet-sen.2018.5258>
- [234] Chih-Yuan Sun and Anthony J.T. Lee. 2017. Tour Recommendations by Mining Photo Sharing Social Media. *Decision Support Systems* 101, C (2017), 28–39. <https://doi.org/10.1016/j.dss.2017.05.013>
- [235] Zhu Sun, Qing Guo, Jie Yang, Hui Fang, Guibing Guo, Jie Zhang, and Robin Burke. 2019. Research Commentary on Recommendations with Side Information: A Survey and Research Directions. *Electron. Commer. Rec. Appl.* 37, C (2019). <https://doi.org/10.1016/j.elerap.2019.100879>

- [236] Yan-Martin Tamm, Rinchin Damdinov, and Alexey Vasilev. 2021. Quality Metrics in Recommender Systems: Do We Calculate Metrics Consistently?. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys '21)*. Association for Computing Machinery, New York, NY, USA, 708–713. <https://doi.org/10.1145/3460231.3478848>
- [237] Wenan Tan, Xin Zhou, Xiao Zhang, Xiaojuan Cai, and Weinan Niu. 2021. Cold Start Recommendation Algorithm Based on Latent Factor Prediction. In *Web Information Systems and Applications: 18th International Conference, WISA 2021, Kaifeng, China, September 24–26, 2021, Proceedings (WISA 2021)*. Springer, 617–624. https://doi.org/10.1007/978-3-030-87571-8_53
- [238] Ziyang Tang, Yiheng Duan, Steven Zhu, Stephanie Zhang, and Lihong Li. 2022. Estimating Long-term Effects from Experimental Data. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 516–518. <https://doi.org/10.1145/3523227.3547398>
- [239] Saravanan Thirumuruganathan, Soon-gyo Jung, Dianne Ramirez Robillos, Joni Salminen, and Bernard J. Jansen. 2021. Forecasting the Nearly Unforecastable: Why Aren't Airline Bookings Adhering to the Prediction Algorithm? *Electronic Commerce Research* 21, 1 (2021), 73–100. <https://doi.org/10.1007/s10660-021-09457-0>
- [240] Kazutoshi Umemoto. 2022. ML-1M++: MovieLens-Compatible Additional Preferences for More Robust Offline Evaluation of Sequential Recommenders. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management (CIKM '22)*. Association for Computing Machinery, New York, NY, USA, 4540–4544. <https://doi.org/10.1145/3511808.3557643>
- [241] Dmitry Ustalov, Natalia Fedorova, and Nikita Pavlichenko. 2022. Improving Recommender Systems with Human-in-the-Loop. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 708–709. <https://doi.org/10.1145/3523227.3547373>
- [242] Mathias Uta, Alexander Felfernig, Viet-Man Le, Andrei Popescu, Thi Ngoc Trang Tran, and Denis Helic. 2021. Evaluating Recommender Systems in Feature Model Configuration. In *Proceedings of the 25th ACM International Systems and Software Product Line Conference - Volume A (SPLC '21)*. Association for Computing Machinery, New York, NY, USA, 58–63. <https://doi.org/10.1145/3461001.3471144>
- [243] Carlos Vaquero-Patricio, Nikki van Ommeren, and Santiago Gil-Begue. 2021. Recommenders in Banking: An End-to-End Personalization Pipeline within ING. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys '21)*. Association for Computing Machinery, New York, NY, USA, 587–589. <https://doi.org/10.1145/3460231.3474612>
- [244] Flavias Vasile, David Rohde, Olivier Jeunen, and Amine Benhalloum. 2020. A Gentle Introduction to Recommendation as Counterfactual Policy Learning. In *Proceedings of the 28th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '20)*. Association for Computing Machinery, New York, NY, USA, 392–393. <https://doi.org/10.1145/3340631.3398666>
- [245] Ignacio Velásquez, Angélica Caro, and Alfonso Rodríguez. 2018. Kontun: A Framework for Recommendation of Authentication Schemes and Methods. *Information and Software Technology* 96, C (2018), 27–37. <https://doi.org/10.1016/j.infsof.2017.11.004>
- [246] Ritvik Vij, Rohit Raj, Madhur Singhal, Manish Tanwar, and Srikanta Bedathur. 2022. VizAI: Selecting Accurate Visualizations of Numerical Data. In *5th Joint International Conference on Data Science & Management of Data (9th ACM IKDD CODS and 27th COMAD) (CODS-COMAD 2022)*. Association for Computing Machinery, New York, NY, USA, 28–36. <https://doi.org/10.1145/3493700.3493717>
- [247] Norha M. Villegas, Cristian Snchez, Javier Daz-Cely, and Gabriel Tamura. 2018. Characterizing Context-Aware Recommender Systems. *Knowledge-Based Systems* 140, C (2018), 173–200. <https://doi.org/10.1016/j.knosys.2017.11.003>
- [248] Michail Vougioukas, Ion Androutsopoulos, and Georgios Paliouras. 2017. A Personalized Global Filter To Predict Retweets. In *Proceedings of the 25th Conference on User Modeling, Adaptation and Personalization (UMAP '17)*. Association for Computing Machinery, New York, NY, USA, 393–394. <https://doi.org/10.1145/3079628.3079655>
- [249] Michail Vougioukas, Ion Androutsopoulos, and Georgios Paliouras. 2018. Identifying Retweetable Tweets with a Personalized Global Classifier. In *Proceedings of the 10th Hellenic Conference on Artificial Intelligence (SETN '18)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3200947.3201019>
- [250] Peng Wang, Yunsheng Jiang, Chunxu Xu, and Xiaohui Xie. 2019. Overview of Content-Based Click-Through Rate Prediction Challenge for Video Recommendation. In *Proceedings of the 27th ACM International Conference on Multimedia (MM '19)*. Association for Computing Machinery, New York, NY, USA, 2593–2596. <https://doi.org/10.1145/3343031.3356085>
- [251] Shangwen Wang, Ming Wen, Bo Lin, and Xiaoguang Mao. 2021. Lightweight Global and Local Contexts Guided Method Name Recommendation with Prior Knowledge. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2021)*. Association for Computing Machinery, New York, NY, USA, 741–753. <https://doi.org/10.1145/3468264.3468567>
- [252] Wei Wang, Shiyong Zheng, Rizwan Ali, and Jiaying Li. 2022. Relevancy or Diversity? Recommendation Strategy Based on the Degree of Multi-Context Use of News Feed Users. *Journal of Global Information Management* 30, 1 (2022), 1–24. <https://doi.org/10.4018/jgim.310929>
- [253] Xi Wang, Chuantao Yin, Xin Fan, Si Wu, and Lan Wang. 2021. An IoT Ontology Class Recommendation Method Based on Knowledge Graph. In *Knowledge Science, Engineering and Management: 14th International Conference, KSEM 2021, Tokyo, Japan, August 14–16, 2021, Proceedings, Part I (KSEM 2021)*. Springer, 666–678. https://doi.org/10.1007/978-3-030-82136-4_54
- [254] Jacek Wasilewski and Neil Hurley. 2018. Are You Reaching Your Audience? Exploring Item Exposure over Consumer Segments in Recommender Systems. In *Proceedings of the 26th Conference on User Modeling, Adaptation and Personalization (UMAP '18)*. Association for Computing Machinery, New York, NY, USA, 213–217. <https://doi.org/10.1145/3209219.3209246>
- [255] Bingbing Wen, Yunhe Feng, Yongfeng Zhang, and Chirag Shah. 2022. ExpScore: Learning Metrics for Recommendation Explanation. In *Proceedings of the ACM Web Conference 2022 (TheWebConf '22)*. Association for Computing Machinery, New York, NY, USA, 3740–3744. <https://doi.org/10.1145/3523227.3547398>

- 1145/3485447.3512269
- [256] Hongyi Wen, Longqi Yang, Michael Sobolev, and Deborah Estrin. 2018. Exploring Recommendations under User-Controlled Data Filtering. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys '18)*. Association for Computing Machinery, New York, NY, USA, 72–76. <https://doi.org/10.1145/3240323.3240399>
- [257] Chin Lin Wong, Diego De Oliveira, Farhad Zafari, Fernando Mourão, Rafael Colares, and Sabir Ribas. 2021. Offline Evaluation Standards for Recommender Systems. In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys '21)*. Association for Computing Machinery, New York, NY, USA, 567–568. <https://doi.org/10.1145/3460231.3474608>
- [258] Chonghuan Xu. 2018. A Novel Recommendation Method Based on Social Network Using Matrix Factorization Technique. *Information Processing & Management* 54, 3 (2018), 463–474. <https://doi.org/10.1016/j.ipm.2018.02.005>
- [259] Ni Xu, Yu-Hsuan Chen, Ping-Yu Hsu, Ming-Shien Cheng, and Chi-Yen Li. 2022. Recommendation Model for Tourism by Personality Type Using Mass Diffusion Method. In *Human Interface and the Management of Information: Applications in Complex Technological Environments: Thematic Area, HIMI 2022, Held as Part of the 24th HCI International Conference, HCII 2022, Virtual Event, June 26 – July 1, 2022, Proceedings, Part II (HIMI 2022)*. Springer, 80–95. https://doi.org/10.1007/978-3-031-06509-5_6
- [260] Zhichao Xu, Hansi Zeng, and Qingyao Ai. 2021. Understanding the Effectiveness of Reviews in E-Commerce Top-N Recommendation. In *Proceedings of the 2021 ACM SIGIR International Conference on Theory of Information Retrieval (ICTIR '21)*. Association for Computing Machinery, New York, NY, USA, 149–155. <https://doi.org/10.1145/3471158.3472258>
- [261] Emre Yalcin. 2022. Exploring Potential Biases towards Blockbuster Items in Ranking-Based Recommendations. *Data Mining and Knowledge Discovery* 36, 6 (2022), 2033–2073. <https://doi.org/10.1007/s10618-022-00860-1>
- [262] Jinfei Yang, Jiajia Li, and Shouqiang Liu. 2018. RETRACTED ARTICLE: A Novel Technique Applied to the Economic Investigation of Recommender System. *Multimedia Tools and Applications* 77, 4 (2018), 4237–4252. <https://doi.org/10.1007/s11042-017-4752-4>
- [263] Li Yang, Wei Huang, and Xinxin Niu. 2017. Defending Shilling Attacks in Recommender Systems Using Soft Co-Clustering. *IET Information Security* 11, 6 (2017), 319–325. <https://doi.org/10.1049/iet-ifs.2016.0345>
- [264] Zhen Yang, Ming Ding, Bin Xu, Hongxia Yang, and Jie Tang. 2022. STAM: A Spatiotemporal Aggregation Method for Graph Neural Network-Based Recommendation. In *Proceedings of the ACM Web Conference 2022 (TheWebConf '22)*. Association for Computing Machinery, New York, NY, USA, 3217–3228. <https://doi.org/10.1145/3485447.3512041>
- [265] Norihiro Yoshida, Seiya Numata, Eunjong Choi, and Katsuro Inoue. 2019. Proactive Clone Recommendation System for Extract Method Refactoring. In *Proceedings of the 3rd International Workshop on Refactoring (IWOR '19)*. IEEE Press, 67–70. <https://doi.org/10.1109/IWoR.2019.00020>
- [266] Mingyu You, Xuan Han, Yangliu Xu, and Li Li. 2020. Systematic Evaluation of Deep Face Recognition Methods. *Neurocomputing* 388, C (2020), 144–156. <https://doi.org/10.1016/j.neucom.2020.01.023>
- [267] Maram Bani Younes and Azzedine Boukerche. 2020. Traffic Efficiency Applications over Downtown Roads: A New Challenge for Intelligent Connected Vehicles. *Comput. Surveys* 53, 5 (2020). <https://doi.org/10.1145/3403952>
- [268] Chenguang Yu, Hao Ding, Houwei Cao, Yong Liu, and Can Yang. 2017. Follow Me: Personalized IPTV Channel Switching Guide. In *Proceedings of the 8th ACM on Multimedia Systems Conference (MMSys '17)*. Association for Computing Machinery, New York, NY, USA, 147–157. <https://doi.org/10.1145/3083187.3083194>
- [269] Hamed Zamani, Markus Schedl, Paul Lamere, and Ching-Wei Chen. 2019. An Analysis of Approaches Taken in the ACM RecSys Challenge 2018 for Automatic Music Playlist Continuation. *ACM Transactions on Intelligent Systems and Technology* 10, 5 (2019). <https://doi.org/10.1145/3344257>
- [270] Eva Zangerle, Christine Bauer, and Alan Said. 2021. Perspectives on the Evaluation of Recommender Systems (PERSPECTIVES). In *Proceedings of the 15th ACM Conference on Recommender Systems (RecSys '21)*. Association for Computing Machinery, New York, NY, USA, 794–795. <https://doi.org/10.1145/3460231.3470929>
- [271] Eva Zangerle, Christine Bauer, and Alan Said. 2022. Second Workshop: Perspectives on the Evaluation of Recommender Systems (PERSPECTIVES 2022). In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22)*. Association for Computing Machinery, New York, NY, USA, 652–653. <https://doi.org/10.1145/3523227.3547408>
- [272] Li Zhang, Wei Lu, Haihua Chen, Yong Huang, and Qikai Cheng. 2022. A Comparative Evaluation of Biomedical Similar Article Recommendation. *Journal of Biomedical Informatics* 131, C (2022). <https://doi.org/10.1016/j.jbi.2022.104106>
- [273] Zhao Zhang, Armelle Brun, and Anne Boyer. 2020. New Measures for Offline Evaluation of Learning Path Recommenders. In *Addressing Global Challenges and Quality Education: 15th European Conference on Technology Enhanced Learning, EC-TEL 2020, Heidelberg, Germany, September 14–18, 2020, Proceedings (EC-TEL 2020)*. Springer, 259–273. https://doi.org/10.1007/978-3-030-57717-9_19
- [274] Qian Zhao, F. Maxwell Harper, Gediminas Adomavicius, and Joseph A. Konstan. 2018. Explicit or Implicit Feedback? Engagement or Satisfaction? A Field Experiment on Machine-Learning-Based Recommender Systems. In *Proceedings of the 33rd Annual ACM Symposium on Applied Computing (SAC '18)*. Association for Computing Machinery, New York, NY, USA, 1331–1340. <https://doi.org/10.1145/3167132.3167275>
- [275] Shenglin Zhao, Irwin King, Michael R. Lyu, Jia Zeng, and Mingxuan Yuan. 2017. Mining Business Opportunities from Location-Based Social Networks. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '17)*. Association for Computing Machinery, New York, NY, USA, 1037–1040. <https://doi.org/10.1145/3077136.3080712>
- [276] Wayne Xin Zhao, Junhua Chen, Pengfei Wang, Qi Gu, and Ji-Rong Wen. 2020. Revisiting Alternative Experimental Settings for Evaluating Top-N Item Recommendation Algorithms. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management (CIKM '20)*. Association for Computing Machinery, New York, NY, USA, 2329–2332. <https://doi.org/10.1145/3340531.3412095>

- [277] Wayne Xin Zhao, Yupeng Hou, Xingyu Pan, Chen Yang, Zeyu Zhang, Zihan Lin, Jingsen Zhang, Shuqing Bian, Jiakai Tang, Wenqi Sun, Yushuo Chen, Lanling Xu, Gaowei Zhang, Zhen Tian, Changxin Tian, Shanlei Mu, Xinyan Fan, Xu Chen, and Ji-Rong Wen. 2022. RecBole 2.0: Towards a More Up-to-Date Recommendation Library. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management (CIKM '22)*. Association for Computing Machinery, New York, NY, USA, 4722–4726. <https://doi.org/10.1145/3511808.3557680>
- [278] Maayan Zhitomirsky-Geffet and Avital Zadok. 2018. Risk analysis and prediction in welfare institutions using a recommender system. *AI & Society* 33, 4 (2018), 511–525. <https://doi.org/10.1007/s00146-017-0735-2>
- [279] Xin Zhou and Wenan Tan. 2020. An Improved Collaborative Filtering Algorithm Based on Filling Missing Data. In *Human Centered Computing: 6th International Conference, HCC 2020, Virtual Event, December 14–15, 2020, Revised Selected Papers (HCC 2020)*. Springer, 220–226. https://doi.org/10.1007/978-3-030-70626-5_23
- [280] Reza Jafari Ziarani and Reza Ravanmehr. 2021. Serendipity in Recommender Systems: A Systematic Literature Review. *Journal of Computer Science and Technology* 36, 2 (2021), 375–396. <https://doi.org/10.1007/s11390-020-0135-9>
- [281] Or Zipori and David Sarne. 2021. ML-Based Arm Recommendation in Short-Horizon MABs. In *Proceedings of the 9th International Conference on Human-Agent Interaction (HAI '21)*. Association for Computing Machinery, New York, NY, USA, 377–381. <https://doi.org/10.1145/3472307.3484673>
- [282] Zainab Zolaktaf, Omar AlOmeir, and Rachel Pottinger. 2018. Bridging the Gap Between User-Centric and Offline Evaluation of Personalized Recommendation Systems. In *Adjunct Publication of the 26th Conference on User Modeling, Adaptation and Personalization (UMAP '18)*. Association for Computing Machinery, New York, NY, USA, 183–186. <https://doi.org/10.1145/3213586.3226216>