



Delft University of Technology

DaisyRec 2.0

Benchmarking Recommendation for Rigorous Evaluation

Sun, Zhu; Fang, Hui; Yang, Jie; Qu, Xinghua; Liu, Hongyang; Yu, Di; Ong, Yew Soon; Zhang, Jie

DOI

[10.1109/TPAMI.2022.3231891](https://doi.org/10.1109/TPAMI.2022.3231891)

Publication date

2023

Document Version

Final published version

Published in

IEEE Transactions on Pattern Analysis and Machine Intelligence

Citation (APA)

Sun, Z., Fang, H., Yang, J., Qu, X., Liu, H., Yu, D., Ong, Y. S., & Zhang, J. (2023). DaisyRec 2.0: Benchmarking Recommendation for Rigorous Evaluation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(7), 8206-8226. <https://doi.org/10.1109/TPAMI.2022.3231891>

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.

Green Open Access added to TU Delft Institutional Repository

'You share, we take care!' - Taverne project

<https://www.openaccess.nl/en/you-share-we-take-care>

Otherwise as indicated in the copyright section: the publisher is the copyright holder of this work and the author uses the Dutch legislation to make this work public.

DaisyRec 2.0: Benchmarking Recommendation for Rigorous Evaluation

Zhu Sun , Hui Fang , Jie Yang , Xinghua Qu, Hongyang Liu , Di Yu, Yew-Soon Ong , *Fellow, IEEE*, and Jie Zhang

Abstract—Recently, one critical issue looms large in the field of recommender systems – there are no effective benchmarks for rigorous evaluation – which consequently leads to unreproducible evaluation and unfair comparison. We, therefore, conduct studies from the perspectives of practical theory and experiments, aiming at benchmarking recommendation for rigorous evaluation. Regarding the theoretical study, a series of hyper-factors affecting recommendation performance throughout the whole evaluation chain are systematically summarized and analyzed via an exhaustive review on 141 papers published at eight top-tier conferences within 2017-2020. We then classify them into model-independent and model-dependent hyper-factors, and different modes of rigorous evaluation are defined and discussed in-depth accordingly. For the experimental study, we release DaisyRec 2.0 library by integrating these hyper-factors to perform rigorous evaluation, whereby a holistic empirical study is conducted to unveil the impacts of different hyper-factors on recommendation performance. Supported by the theoretical and experimental studies, we finally create benchmarks for rigorous evaluation by proposing standardized procedures and providing performance of ten state-of-the-arts across six evaluation metrics on six datasets as a reference for later study. Overall, our work sheds light on the issues in recommendation

evaluation, provides potential solutions for rigorous evaluation, and lays foundation for further investigation.

Index Terms—Benchmarks, fair comparison, recommender systems, reproducible evaluation, standardized procedures.

I. INTRODUCTION

WITH the advent of the Big Data era, we are flooded by the exponentially increased information on the Internet. To ease the severe information overload problem [1], recommender systems have been extensively studied in academia and widely applied in industry across different domains, such as e-commerce (e.g., Amazon, Tmall), location-based social networks (e.g., Foursquare, Yelp), multi-media (e.g., Netflix, Spotify), and so forth. With a massive amount of recommendation approaches being proposed, one critical issue has attracted much attention from researchers in the field of recommender systems: there are few effective benchmarks for evaluation [2], [3], [4], which, consequently, leads to unreproducible evaluation and unfair comparison. As indicated by the recent study [5], results for baselines that have been used in numerous publications over the past five years are suboptimal; with a careful setup, the baselines even outperform the reported results of any newly proposed method. This is in alignment with another latest study [6], which discovers that the recent proposed deep learning models (DLMs) can be defeated by comparably simple baselines, such as MostPop and ItemKNN [7] with fine-tuned parameters. These findings initiate an extremely heated discussion on the evaluation of recommendation methods and inspire us to deeply consider the underlying barriers that hinder the rigorous evaluation in recommendation.

As a matter of fact, there are a number of *hyper-factors* which may affect the recommendation performance throughout the whole evaluation chain, and their best settings are unknown. They can be broadly classified into two types as depicted in Fig. 1, namely model-independent and model-dependent hyper-factors. The former refers to the hyper-factors that are isolated from the model design and optimization process (e.g., dataset and comparison baseline selection); whilst the latter indicates the ones involved in the model development and parameter optimization procedure (e.g., loss function design and regularization terms). According to this categorization, three main aspects may inherently lead to such non-rigorous evaluation.

- *Diverse settings on model-independent hyper-factors.* With being prominent in different platforms, there are diverse

Manuscript received 6 June 2022; revised 3 October 2022; accepted 21 December 2022. Date of publication 26 December 2022; date of current version 5 June 2023. This work was supported in part by the National Natural Science Foundation of China under Grant 72192832, in part by the Natural Science Foundation of Shanghai under Grant 21ZR1421900, Delft Design@Scale AI Lab. This work was also supported in part by the A*Star Center for Frontier Artificial Intelligence Research and in part by the Data Science and Artificial Intelligence Research Centre, School of Computer Science and Engineering at the Nanyang Technological University (NTU), Singapore. The work of Zhang Jie was supported in part by the MOE AcRF Tier 1 funding under Grant RG90/20. This work was also supported by Shanghai Rising-Star Program 23QA1403100. Recommended for acceptance by Y. Sun. (*Corresponding author: Hui Fang.*)

Zhu Sun is with the Institute of High Performance Computing and Centre for Frontier AI Research, A*STAR, Singapore 138632 (e-mail: sun-zhuntu@gmail.com).

Hui Fang is with the Shanghai University of Finance and Economics, Shanghai 200437, China (e-mail: fang.hui@mail.shufe.edu.cn).

Jie Yang is with the Delft University of Technology, 2628, CD Delft, The Netherlands (e-mail: j.yang-3@tudelft.nl).

Xinghua Qu is with the ByteDance AI Lab, Singapore 048583 (e-mail: quxinghua17@gmail.com).

Hongyang Liu is with the Yanshan University, Qinhuangdao, Hebei 066104, China (e-mail: hylu767289@gmail.com).

Di Yu is with the Singapore Management University, Singapore 188065 (e-mail: yudi201909@gmail.com).

Yew-Soon Ong is with the Nanyang Technological University, Singapore 639798, and also with A*STAR Centre for Frontier AI Research, Singapore 138632 (e-mail: asysong@ntu.edu.sg).

Jie Zhang is with the Nanyang Technological University, Singapore 639798 (e-mail: ZhangJ@ntu.edu.sg).

This article has supplementary downloadable material available at <https://doi.org/10.1109/TPAMI.2022.3231891>, provided by the authors.

Digital Object Identifier 10.1109/TPAMI.2022.3231891

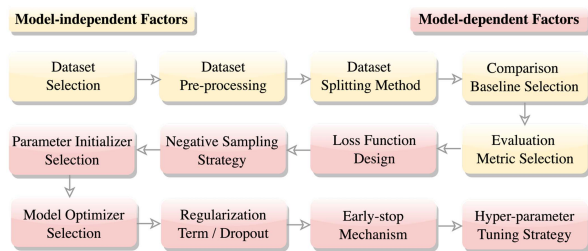


Fig. 1. Hyper-factors within the whole recommendation evaluation chain.

recommendation datasets (i.e., dataset selection) in various application domains shown in Table 8 (Appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2022.3231891>). Taking the movie domain as an example, the datasets vary from MovieLens (ML), Netflix to Amazon-Movie, etc. Even for the same dataset, it may have different versions with different sizes covering different durations, such as, ML-100 K/1 M/10 M/20 M/25 M/Latest. Different researchers may choose different datasets across different domains based on their requirements, and only report results on their selected datasets; meanwhile different settings on other model-independent hyper-factors (e.g., dataset pre-processing and splitting strategies) may generate entirely different recommendation performance.

- *Diverse settings on model-dependent hyper-factors.* There are different choices for the model-dependent hyper-factors. For instance, the loss function could be either point-wise (square error loss [149] and cross-entropy loss [134]) or pair-wise (log loss [150], hinge loss [54] and top-1 loss [151]); different types of model optimizers are also available, ranging from stochastic gradient descent (SGD) to adaptive moment estimation (Adam). The recommendation results may vary a lot with different settings on these model-dependent factors even with fixed settings on model-independent ones.
- *Missing setting details.* Most importantly, a majority of papers do not report details on the settings of either model-independent and model-dependent hyper-factors, such as data processing and parameter settings, which increases the difficulty on evaluation and leads to inconsistent results in reproduction by different researchers, thus heavily aggravating the unreproducible evaluation and unfair comparison issues.

With the above issues in mind, we are seeking at benchmarking recommendation for rigorous (i.e., reproducible and fair) evaluation, thus helping achieve a healthy and sustainable growth of research in this area. Considering the diverse recommendation tasks (e.g., temporal, sequential, diversity, explanation, location, group and cross-domain aware recommendation), we first mainly focus on the **general top- N** recommendation task, which is one of the hottest and most prominent topics in recommendation. To this end, we conduct extensive studies from the perspectives of both practical theory and experiments.

- Regarding the theoretical study, we conduct an exhaustive review on 141 papers related to top- N recommendation published in the recent four years

(2017-2020) on eight prestigious conferences as representatives, including KDD, SIGIR, WWW, IJCAI, AAAI, RecSys, WSDM and CIKM.¹ In doing so, we systematically extract and summarize hyper-factors throughout the whole evaluation chain, and classify them into model-independent and model-dependent factors in Fig. 1. Accordingly, different modes (e.g., relax, strict and mixed) of rigorous evaluation are defined and discussed in-depth, acting as valuable guidance for later study.

- For the experimental study, we release a Python-based recommendation testbed – DaisyRec 2.0 to integrate the hyper-factors throughout the evaluation chain². Our testbed advances existing libraries (e.g., LibRec [152], DeepRec [153] and RecBole [154]), which mainly aim to implement various state-of-the-art recommenders, in the light of performing rigorous evaluation in recommendation. Based on DaisyRec 2.0, a holistic empirical study has been performed to comprehensively analyze the impact of different hyper-factors on recommendation performance.

Supported by both theoretical and experimental study, we finally create benchmarks by proposing the standardized procedures to help enhance the reproducibility and fairness of evaluation. Meanwhile, the performance of ten well-tuned state-of-the-arts on six widely-used datasets across six metrics is provided as a reference for fair evaluation. Additionally, a number of interesting findings are noted throughout our study, for example, (1) the recommendation performance does not necessarily improve with denser datasets; (2) some non-deep learning based baselines, e.g., PureSVD [155] can achieve a better balance between recommendation accuracy and complexity; (3) the best hyper-parameter settings for one specific metric does not necessarily guarantee optimums w.r.t. other metrics; (4) although the objective function with pair-wise log loss generally outperforms others, different methods may have their best fit objective functions; (5) uniform sampler, though simple, performs better than the popularity based sampler; and (6) different parameter initializers and model optimizers can extensively affect the final recommendation accuracy. To sum up, our work sheds lights on the issues in evaluation for recommendation, provides potential solutions for rigorous evaluation, and paves the way for further investigation³.

¹They are most important venues (full names see Table 9 in Appendix, available in the online supplemental material) to accept high-quality recommendation papers and other related conferences and journals will be considered in our future study.

²[Online]. Available: <https://github.com/recsys-benchmark/DaisyRec-v2.0>

³A preliminary report of our work was published at RecSys'20 as a reproducibility paper [4]. In this study, we have extended it from two aspects: (1) with regards to theoretical study, we conduct a more in-depth analysis on the hyper-factors throughout the whole evaluation chain by reviewing more latest literature in 2020, whereby several new hyper-factors (e.g., regularization terms, parameter initializers and model optimizers) are further considered; we systematically classify these hyper-factors into model-independent and -dependent ones, whereby different modes of rigorous evaluation are well defined and discussed in-depth; and (2) for the experimental study, we release DaisyRec 2.0 by further fusing these new hyper-factors and extending existing ones (e.g., more types of loss function designs, negative sampling strategies, data splitting methods and deep learning based baselines). To be more user-friendly, we design a user interface tool for automatic command generation. Thereby, more holistic experiments are conducted based on DaisyRec 2.0 to unveil the impacts of different hyper-factors on recommendation performance, where more interesting and insightful observations are gained.

II. RELATED WORK

While long been recognized as a key feature of scientific discoveries, reproducibility has been increasingly characterized as a crisis recently [156], [157], [158]. It is becoming a primary concern in computer and information science, evidenced by the recently developed ACM policy on Artifact Review and Badging⁴ and emerging efforts including seminars [160], workshops [161], reproducibility checklist⁵ [162], and focused tracks at major conferences, such as ECIR [163], ACM MM [164], SIGIR⁶, and ISWC [165]. Specific to recommender systems research, besides the reproducibility track starting from 2020 on the premier conference for recommender systems – RecSys [166], the discussions have been concentrated on the fairness of comparison between newly proposed and baseline methods [5], [6]. In very recent work, Dacrema et al. [6] found neural models hardly outperform fine-tuned memory- and latent factor-based methods, and a similar finding was also discovered in [5].

Despite the importance, improving reproducibility in recommender systems research is highly challenging due to the many influential evaluation factors for recommendation performance. Said et al. [2] found large differences in the effectiveness of recommendation methods across different implementation frameworks as well as across evaluation datasets and metrics. A companion toolkit RiVal [3] was released to allow for the control of data splitting and evaluation metrics, while Elliot [167] further improves it by implementing more baselines and incorporating statistic significance tests. Beel et al. [168] found a similar phenomenon in news and research paper recommendation and identified influential factors such as user characteristics and time of recommendation. Valcarce et al. [83] specifically studied the properties of evaluation metrics for item ranking, marking precision as the most robust and NDCG presenting the highest discriminative power. More recently, Rendle et al. [5] demonstrated the importance of hyperparameter search in baseline methods, e.g., matrix factorization, and stressed the need for standardized benchmarks where methods should be extensively tuned for fair comparison. Sachdeva et al. [169] particularly examined the impact of dataset sampling strategies on model performance, and indicated that sampling methods, including the random ways do matter with regard to final performance of recommendation algorithms.

Existing benchmarks are, however, either restricted to pre-neural methods [2], a single evaluation factor [83], or rating prediction [5] which has been discouraged as a way to formulate the recommendation problem [170]; besides, most of the existing benchmarks consider two or three datasets (including [6]), ignoring the richness of available datasets often chosen by newly published work. The two most recent work, RecBole [154] and Elliot [167], has partially alleviated the aforementioned issues by implementing more baselines (neural ones included), considering varied datasets and recommendation scenarios

⁴[Online]. Available: www.acm.org/publications/policies/artifact-review-badging; see also SIGIR’s implementation of the policy [159].

⁵[Online]. Available: aaai.org/Conferences/AAAI-22/reproducibility-checklist/.

⁶[Online]. Available: sigir.org/sigir2022/call-for-reproducibility-track-papers/.

(e.g., temporal and context-aware ones), and incorporating hyper-parameter optimization strategies. However, similar to other recommender system libraries (e.g., Librec [152], MyMediaLite [171], and Surprise [172]), they strive to provide a unified framework for developing and reproducing algorithms for different scenarios in terms of different evaluation metrics.

Instead, aimed for a full treatment of evaluation issues, our work takes a bottom-up approach analyzing an extensive amount of literature to search for and summarize important evaluation factors, denoted as *hyper-factors* (categorized as model-dependent and model-independent ones), which might influence model performance in model evaluation, towards the goal of performing rigorous evaluation. We further present a benchmark supported by an empirical study at a bigger-than-ever scale with the hope of laying a strong foundation for future research.

III. PRACTICAL THEORY STUDY

A. Hyper-Factor Extraction

As we seek to benchmark recommendation for rigorous evaluation, we first conduct study from the angle of practical theory by an exhaustive literature review, so as to extract and summarize hyper-factors affecting recommendation performance throughout the whole evaluation chain. In particular, we review papers published in the recent four years (2017-2020) on eight top tier conferences, namely, AAAI, CIKM, IJCAI, KDD, RecSys, SIGIR WSDM and WWW. As a starting point, we mainly focus on recommendation methods for implicit feedback based top- N recommendation, which is one of the hottest topics in recommendation. Other tasks (e.g., sequential recommendation) are remained for future exploration. Specifically, we first search the accepted full paper lists ($8 * 4 = 32$) for the eight conferences in the four years. Given our interest and the 32 lists, we only consider papers with titles containing keywords ‘recommend*’ or ‘collaborative filtering’. After that, we manually select papers towards top- N recommendation adopting ranking metrics (e.g., Precision, Recall) to evaluate the *accuracy* of recommendation. In the end, we obtain a collection of 141 relevant papers as listed in Table I.

By delicately reviewing the collected papers in Table I, we find that there are a bunch of hyper-factors that may affect the recommendation performance along with the entire evaluation chain. Typically, they can be classified into two types: (1) *model-independent* factors that are isolated from the model design and optimization process (e.g., dataset and baseline selection); and (2) *model-dependent* ones involved in the model development and parameter optimization process (e.g., loss function and regularization terms). Fig. 1 illustrates the two types of hyper-factors along with the whole evaluation chain, starting with the dataset selection and ending with hyper-parameter tuning strategy. Fig. 3 shows the tree diagram of these hyper-factors in recommendation evaluation. Next, we will analyze them one by one.

B. Analysis on Model-Independent Hyper-Factors

1) *Dataset Selection*: We find two major issues on the utilized datasets by analyzing the collected papers: (1) domain

TABLE I
SUMMARY OF THE COLLECTED PAPERS

Venue	No.	Reference			
		2017	2018	2019	2020
AAAI	22	[8], [9]	[10], [11], [12]	[13], [14], [15], [16], [17], [18], [19], [20]	[21], [22], [23], [24], [25], [26], [27], [28], [29]
CIKM	19	[30], [31], [32]	[33], [34], [35]	[36], [37], [38]	[39], [40], [41], [42], [43], [44], [45], [46], [47], [48]
IJCAI	20	[49], [50], [51]	[52], [53], [54], [55], [56], [57]	[58], [59], [60], [61], [62], [63], [64]	[65], [66], [67], [68]
KDD	12	[69]	[70], [71], [72]	[73], [74], [75], [76], [77]	[78], [79], [80]
RecSys	14	[81], [82]	[83], [84], [85]	[86], [87], [88], [89], [90], [91]	[92], [93], [94]
SIGIR	22	[95]	[96], [97], [98], [99], [100]	[101], [102], [103], [104]	[105], [106], [107], [108], [109], [110], [111], [112], [113], [114], [115], [116]
WSDM	16	[117]	[118], [119], [120]	[121], [122], [123], [124]	[125], [126], [127], [128], [129], [130], [131], [132]
WWW	16	[133], [134]	[135], [136]	[137], [138], [139], [140], [141], [142]	[143], [144], [145], [146], [147], [148]
Total	141	15	28	43	55

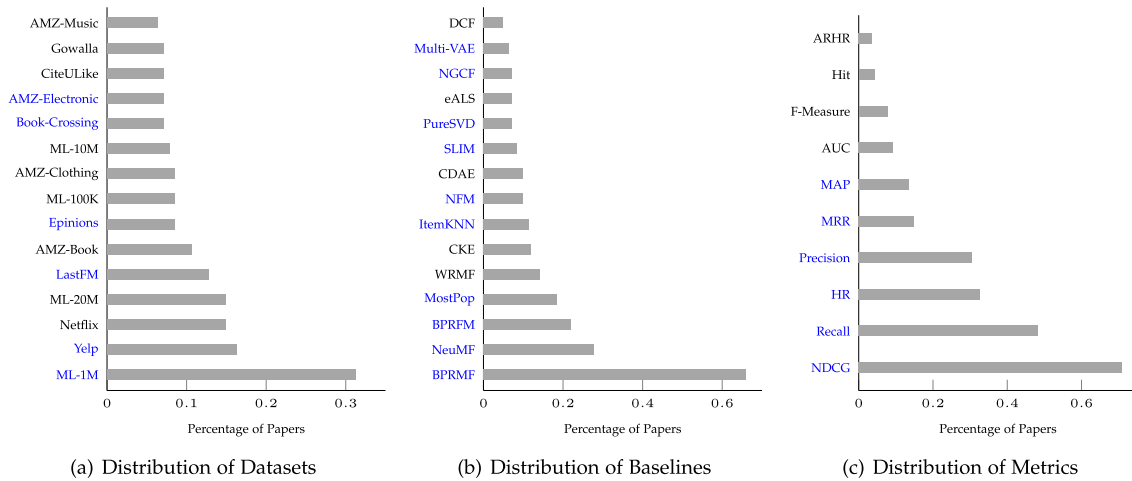


Fig. 2. (a) popularity of the top-15 datasets, where ‘ML, AMZ’ denote MovieLens and Amazon, respectively; (b) popularity of the top-15 baselines; and (c) popularity of the top-10 evaluation metrics. Note that the selected datasets, baselines and metrics in our study are highlighted in blue.

diversity, i.e., there are massive different datasets within and across various domains, as shown in Table 8; and (2) version diversity, i.e., many datasets, though with the same names, may have different versions. For example, we find more than three versions for Yelp, as it has been updated for different rounds of the challenge. By treating different versions as a same dataset, there are 84 different datasets used in the 141 papers. Fig. 2(a) shows the dataset popularity, i.e., percentage of papers for the top-15 used datasets. Around 90% of the 141 papers adopt at least one of the 15 datasets.

For a practical study, we further delicately select six among them by considering popularity and domain coverage, thus resolving the domain diversity issue. They are **ML-1M** (Movie), **Yelp** (LBSNs), **LastFM** (Music), **Epinions** (SNs), **Book-X** (Book) and **AMZe** (Consumable), covering 62% papers of the collection. To ease the version diversity issue, we conduct a careful selection by considering the authority and information richness of data sources, which could benefit the study on diverse recommendation models. Specifically, we use MovieLens-1 M (ML-1 M) released by GroupLens⁷; Yelp was created by Kaggle in 2018⁸; LastFM was released by

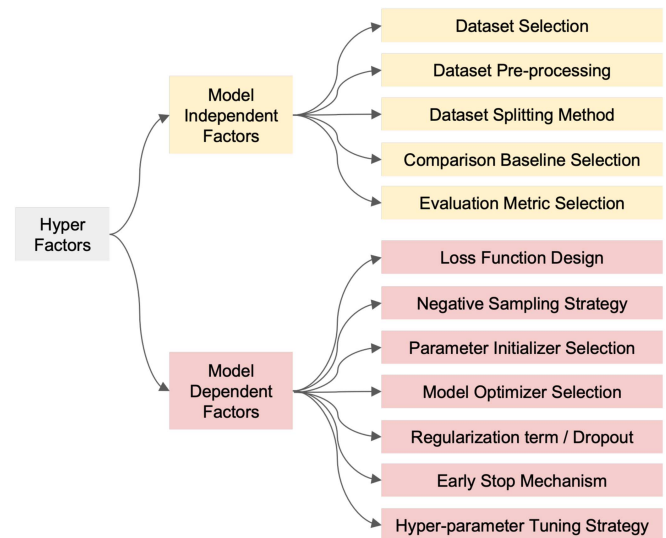


Fig. 3. The tree diagram of hyper-factors in recommendation evaluation.

the 2nd international workshop HetRec 2011 [173]; Epinions was crawled by [174] containing timestamp and item category information; Book-Crossing (Book-X) [175] was collected

⁷[Online]. Available: grouplens.org/datasets/movielens/

⁸[Online]. Available: www.kaggle.com/yelp-dataset/yelp-dataset

TABLE II
STATISTICS OF THE SIX SELECTED DATASETS

Dataset	ML-1M	Yelp	LastFM	Epinions	Book-X	AMZe	
origin	#User	6,038	1,326,101	1,892	22,164	105,283	4,201,696
	#Item	3,533	174,567	17,632	296,277	340,556	476,002
	#Record	575,281	5,261,669	92,834	922,267	1,149,780	7,824,482
	Density	2.697e-2	2.273e-5	2.783e-3	1.404e-4	3.207e-5	3.912e-6
5-filter	#User	6,034	227,109	1,874	21,995	22,072	253,994
	#Item	3,125	123,985	2,828	31,678	43,748	145,199
	#Record	574,376	3,419,587	71,411	550,117	623,405	2,109,869
	Density	3.046e-2	1.214e-4	1.348e-2	7.895e-4	6.456e-4	5.721e-5
10-filter	#User	5,950	96,168	1,867	21,111	12,720	63,161
	#Item	2,811	80,351	1,530	14,030	18,318	85,930
	#Record	571,549	2,458,153	62,984	434,162	443,196	949,416
	Density	3.412e-2	3.181e-4	2.205e-2	1.466e-3	1.902e-3	1.749e-4
Timestamp	✓	✓	×	✓	×	✓	

by Cai-Nicolas Ziegler from the Book-Crossing community⁹; Amazon Electronic (AMZe) was released by Julian McAuley¹⁰. The statistics of all datasets are listed in Table II and all links for the datasets are available at the homepage of DaisyRec 2.0.

2) *Dataset Pre-Processing*: There are typically two core steps for the dataset pre-processing, namely binarization and filtering.

Binarization. As our current study focuses on the implicit feedback, all the datasets with explicit feedback (e.g., ratings or counts) should be binarized into implicit data. Let $u \in \mathcal{U}$, $i \in \mathcal{I}$ denote user u and item i ; \mathcal{U}, \mathcal{I} are user and item sets; and $r_{ui} \in \mathcal{R}$ is the binary feedback of user u over item i . For each user u , we transform all her explicit feedback with no less than a threshold (denoted by r) into positive feedback ($r_{ui} = 1$); otherwise, negative feedback ($r_{ui} = 0$). Different papers may have different settings for r (e.g., $r = 1/2/3/4$). By following the majority studies [74], [76], [103], we recommend to set $r = 4$ for ML-1 M, and $r = 1$ for the rest datasets.

Filtering. The original datasets are generally quite sparse, where some users (items) only interact with few items (users), e.g., less than 5. To ensure the quality, the filtering strategy is usually adopted to help remove the inactive users and items. By analyzing the paper collection, we have found out that around 57% of papers adopt filtering strategies; while 22% of papers utilize the original datasets; and 21% of papers do not report details on data filtering. Among the papers adopting pre-processing strategies, more than 58% of them utilize 5- or 10-filter/core setting [38], [73], [102] on the datasets, which filter out users and items with less than 5 or 10 interactions, respectively. While, others adopt, such as 1-, 2-, 3-, 4-, 20- or 30-filter/core settings. Therefore, to check the performance and robustness of different methods w.r.t. various data sparsity levels, besides original datasets, we also take the two most common settings (i.e., **5-** and **10-filter**) on all selected datasets, the statistics of which are summarized in Table II. Note that F -filter is different from F -core: the former means that users and items are only filtered with less than F interactions in one pass; by contrast, the latter indicates a recursive filtering until all users and items have at least F interactions.

3) *Dataset Splitting Methods*: Three types of data splitting methods are mainly used in the collected papers, including *split-by-ratio* (69% of the papers), *leave-one-out* (21% of the papers)

⁹[Online]. Available: grouplens.org/datasets/book-crossing/

¹⁰[Online]. Available: jmcauley.ucsd.edu/data/amazon/links.html

TABLE III
AVERAGE SIZE OF TRAINING & TEST SETS FOR EACH USER

Type	Setting	ML-1M	Yelp	LastFM	Epinions	Book-X	AMZe
Train	origin	86	4	39	35	10	2
	5-filter	86	13	30	23	23	7
	10-filter	86	21	27	20	28	12
Test	origin	64	2	10	30	5	2
	5-filter	64	5	8	13	7	3
	10-filter	64	8	7	10	8	5

and *split-by-time* (6% of the papers). There are also 4% of papers not reporting their data splitting methods. In particular, *split-by-ratio* means that a proportion ρ (e.g., $\rho = 80\%$) of the dataset (i.e., user-item interaction records) is treated as training set, and the rest ($1 - \rho = 20\%$) as test set; *leave-one-out* refers to that for each user, only one record is kept as test set and the remaining are for training; and *split-by-time* indicates directly dividing training and test sets by a fixed timestamp, that is, the data before the fixed timestamp is used as training set, and the rest as test set.

Although 69% of papers adopt *split-by-ratio*, they are quite different from each other due to: (1) different proportion settings, e.g., $\rho = 50\%, 60\%, 70\%, 80\%, 90\%$; (2) global- or user-level split. That is, some globally split the entire records into training and test sets regardless of different users; whilst others split training and test sets on the user basis; and (3) random- or time-aware split. Among papers exploiting *split-by-ratio*, 87% of papers merely randomly split the data, whereas 13% of papers split the data based on the timestamp, i.e., the earlier (e.g., $\rho = 80\%$) records as training and the later ones as test. In terms of *leave-one-out*, the split is generally on the user basis and timestamp is taken into account by 60% of papers. To validate the impact of different data splitting methods, in our study, we compare the recommendation performance with random-/time-aware *split-by-ratio* at *global-level* with $\rho = 80\%$ and random-/time-aware *leave-one-out* at *user-level*, and *split-by-time* as our future exploration.

Besides, to improve the test efficiency, they usually randomly sample a number of negative items (e.g., $neg_test = 99, 100, 999, 1,000$) that are not interacted by each user, and then rank each test item among the ($neg_test + 1$) items for recommendation [44], [104], [110], [138]. To speed up the test process, we randomly sample negative items for each user to ensure her test candidates to be 1,000, and then rank all test items among the 1,000 items w.r.t. both *split-by-ratio* and *leave-one-out*. Table III depicts the average number of test items for each user on the six datasets across origin, 5- and 10-filter settings w.r.t. *split-by-ratio*, where all values are smaller than 100, indicating that 1,000 test candidates are sufficient to examine the performance of recommenders.

4) *Comparison Baseline Selection*: As observed, the compared baselines vary a lot in different collected papers. We show the top-15 widely-compared baselines in these papers in Fig. 2(b), covering 98% of papers in total, that is, 98% of the papers consider at least one of the 15 baselines. They can be classified into three types, (1) memory-based methods (MMs): MostPop, ItemKNN [7]; (2) latent factor methods (LFMs): BPRMF [150], FM [176], WRMF [177], SLIM [178], PureSVD [155], eALS [179] and DCF [180]; and (3) deep

learning methods (DLMs): NeuMF [134], CKE [181], NeuMF [134], CDAE [182], NGCF [102] and Multi-VAE [136].

For a practical study, we ultimately take 10 baselines into account, as highlighted in blue in Fig. 2(b). Specifically, two MMs are considered. **MostPop** is a non-personalized method and recommends most popular items to all users; and **ItemKNN** is a K -nearest neighborhood based method recommending items based on item similarity. We adapt it for implicit feedback data by following [177], and adopt cosine similarity. In terms of LFM, **BPRMF** is selected as the representative of matrix factorization method (WRMF, eALS and DCF are remained for future exploration); **BPRFM** (factorization machine) considers the second-order feature interactions between inputs and we train it by optimizing the BPR loss [150]; **SLIM** [178] learns a sparse item similarity matrix by minimizing a constrained reconstruction square loss; and **PureSVD** directly performs conventional singular value decomposition on the user-item implicit interaction matrix, where all the unobserved entries are set as 0. Regarding DLMs, **NeuMF** [134] is a state-of-the-art neural method taking advantage of both generalized matrix factorization and multi-layer perceptron (MLP); **NFM** [134] seamlessly combines the linearity of FM in modelling second-order feature interactions and the non-linearity of neural network in modelling higher-order feature interactions, and we train it by optimizing the BPR loss; **NGCF** [102] leverages graph neural networks to capture the high-order connectivity in the user-item graph; and **Multi-VAE** [136] is a generative model with multinomial likelihood and extends variational autoencoders to collaborative filtering. CKE and CDAE are remained for future study, as CKE involves textual and visual information besides user-item interaction data; both CDAE and Multi-VAE are in the family of autoencoders, while Multi-VAE has proven to be a stronger baseline [6].

5) *Evaluation Metric Selection*: The evaluation metrics vary a lot in different papers in the collection. Fig. 2(c) depicts the popularity of the used evaluation metrics. We thus adopt the top-6 evaluation metrics covering 99% of the collected papers. That is to say, 99% of these papers adopt at least one of the six metrics. They are *Precision*, *Recall*, Mean Average Precision (**MAP**), Hit Ratio (**HR**), Mean Reciprocal Rank (**MRR**) and Normalized Discounted Cumulative Gain (**NDCG**). In particular, the first four metrics intuitively measure whether a test item is present in the top- N recommendation list, whilst the latter two accounts for the ranking positions of test items. Detailed formulas are given by Table 10 (Appendix, available in the online supplemental material), where $R(u), T(u)$ represent the recommendation set and the test set for user u , respectively; $rel_j = 1/0$ indicates whether the item at rank j is in the intersection of $R(u)$ and $T(u)$, i.e., $(R(u) \cap T(u))$; $\delta(x) = 1$ if x is true, otherwise 0; and IDCG means the maximum possible DCG through ideal ranking.

C. Analysis on Model-Dependent Hyper-Factors

1) *Loss Function Design*: Two types of objective functions are widely utilized by the collected papers: **point-wise** (55% of the collected papers) and **pair-wise** (40% of the collected

TABLE IV
OBJECTIVE FUNCTIONS OF DIFFERENT BASELINES

Method	Origin	To explore		
BPRMF	$\mathcal{L}_{poi} + f_{ll}$	$\mathcal{L}_{poi} + f_{cl}$	$\mathcal{L}_{pai} + f_{tl}$	$\mathcal{L}_{pai} + f_{hl}$
BPRFM	$\mathcal{L}_{poi} + f_{ll}$	$\mathcal{L}_{poi} + f_{cl}$	$\mathcal{L}_{pai} + f_{tl}$	$\mathcal{L}_{pai} + f_{hl}$
SLIM	$\mathcal{L}_{poi} + f_{sl}$	-	-	-
NeuMF	$\mathcal{L}_{poi} + f_{cl}$	$\mathcal{L}_{pai} + f_{ll}$	$\mathcal{L}_{pai} + f_{tl}$	$\mathcal{L}_{pai} + f_{hl}$
NFM	$\mathcal{L}_{poi} + f_{ll}$	$\mathcal{L}_{poi} + f_{cl}$	$\mathcal{L}_{pai} + f_{tl}$	$\mathcal{L}_{pai} + f_{hl}$
NGCF	$\mathcal{L}_{poi} + f_{ll}$	$\mathcal{L}_{poi} + f_{cl}$	$\mathcal{L}_{pai} + f_{tl}$	$\mathcal{L}_{pai} + f_{hl}$
Multi-VAE	$\mathcal{L}_{poi} + f_{cl}$	-	-	-

papers). The former only relies on the accuracy of the prediction of individual preferences; whilst the latter approximates ranking loss by considering the relative order of the predictions for pairs of items. Regardless of which one is deployed, it is critical to properly exploit unobserved feedback within the model, as merely considering the observed feedback fails to account for the fact that feedback is not missing at random, thus being not suitable for top- N recommenders [182]. Let \mathcal{L} denote the objective function of recommendation task. The point- and pair-wise objectives are thus given by:

$$\begin{aligned}\mathcal{L}_{poi} &= \sum_{(u,i) \in \tilde{\mathcal{O}}} f(r_{ui}, \hat{r}_{ui}) + \lambda \Omega(\Theta); \\ \mathcal{L}_{pai} &= \sum_{(u,i,j) \in \tilde{\mathcal{O}}} f(r_{uij}, \hat{r}_{uij}) + \lambda \Omega(\Theta),\end{aligned}\quad (1)$$

where $\tilde{\mathcal{O}} = \{\mathcal{O}^+ \cup \mathcal{O}^-\}$ is the augmented dataset with the unobserved user-item set $\mathcal{O}^- = \{(u, j) | r_{uj} = 0\}$ in addition to the observed user-item set $\mathcal{O}^+ = \{(u, i) | r_{ui} = 1\}$; $f(\cdot)$ is the loss function; r_{ui}, \hat{r}_{ui} are the observed and estimated feedback of user u on item i , respectively; (u, i, j) is the triple meaning that u prefers positive item i to negative item j ; $r_{uij} = r_{ui} - r_{uj}$, $\hat{r}_{uij} = \hat{r}_{ui} - \hat{r}_{uj}$; $\Omega(\Theta)$ is the regularization term, whose impact is illustrated in Section III-C5; and Θ is the set of model parameters.

W.r.t. the loss function $f(\cdot)$, point-wise objective usually adopts square loss and cross-entropy (CE) loss, whereas pair-wise objective generally employs log loss, top-1 loss and hinge loss:

$$\begin{aligned}\mathcal{L}_{poi} &= \begin{cases} \text{Square Loss} & f_{sl} = \frac{1}{2}(r_{ui} - \hat{r}_{ui})^2 \\ \text{CE Loss} & f_{cl} = -r_{ui} \log(\hat{r}_{ui}) \\ & \quad - (1 - r_{ui}) \log(1 - \hat{r}_{ui}) \end{cases} \\ \mathcal{L}_{pai} &= \begin{cases} \text{Log Loss} & f_{ll} = -\log(\sigma(\hat{r}_{uij})) \\ \text{Top-1 Loss} & f_{tl} = \sigma(-\hat{r}_{uij}) + \sigma(\hat{r}_{uj}^2) \\ \text{Hinge Loss} & f_{hl} = \max(0, 1 - \hat{r}_{uij}). \end{cases}\end{aligned}\quad (2)$$

Table IV shows the original objectives used by BPRMF, BPRFM, SLIM, NeuMF, NFM, NGCF and Multi-VAE. Besides, we vary different objectives on these baselines to further examine their respective impacts. Note that MostPop, ItemKNN and PureSVD do not have objective functions; we did not consider the square loss, as it is more suitable for rating prediction task instead of ranking problem [150]; Multi-VAE cannot be easily adapted with different objective functions; and we did not study the impacts of different objectives on SLIM due to its high complexity and low scalability, which will be discussed in Section IV-B3.

2) *Negative Sampling Strategies*: As pointed out in Section III-C1, properly exploiting the unobserved feedback (i.e., negative samples) helps learn users' relative preferences and benefits more accurate top- N recommendation. This can be further supported by the fact that around 65% of the collected papers consider the unobserved feedback when designing objective functions regardless of point- and pair-wise ones. However, it is not practical to leverage all unobserved feedback in large volume, as most users only provide feedback for a small number of items. Negative sampling is, therefore, adopted to balance the efficiency and effectiveness. It is noteworthy that we follow majority studies [73], [84], [134], [137] and directly treat the unobserved feedback as negative feedback. There may be different explanations behind the unobserved feedback [183], but we leave it for further exploration.

There are various kinds of negative sampling strategies. Specifically, **uniform sampler** [134], where all unobserved items of each user are sampled with an equal probability, has been adopted by almost all papers in the collection. To better study the impact of negative sampling, we additionally consider item popularity-based samplers, which have also been adopted in recommendation [38], [179]. *Low-popularity sampler* refers to that for each user, her unobserved items with a lower popularity are sampled with a higher probability. This is based on the assumption that a user is less likely to prefer less popular items. *High-popularity sampler* is opposite to the low-popularity sampler, where the unobserved items of each user with a higher popularity are more likely to be sampled. The rationale behind is that if a user provides no feedback for a quite popular item favored by a large number of users, it indicates that she may be not into this item. Moreover, we also compare two types of hybrid samplers by leveraging both uniform and popularity samplers, namely *uniform+low-popularity sampler* and *uniform+high-popularity sampler*. In these cases, half of the unobserved items are sampled via the uniform sampler, while the rest half are sampled via popularity samplers.

3) *Parameter Initializer Selection*: There are normally a set of learnable parameters (e.g. user/item representation matrix and the network weights) for the recommendation models, ranging from early LFM to recently emerged DLMs. A proper parameter initializer will assist in a faster model convergence and better model performance. Specifically, the core of LFM is to learn accurate user and item representations, which are generally initialized based on either a *Uniform distribution* in the range of $(0, a)$ or a *Normal/Gaussian distribution* with zero mean and a variance of σ^2 [28], [92], given by,

$$\mathbf{v}_{uf}/\mathbf{v}_{vf} \sim U(0, a); \quad \mathbf{v}_{uf}/\mathbf{v}_{vf} \sim \mathcal{N}(0, \sigma^2), \quad (3)$$

where the common settings are $a = 1$; and $\sigma = 0.01$.

Regarding DLMs, besides user and item representations, initializing with proper weights helps ensure the network to converge in a reasonable amount of time; otherwise the network loss function does not go anywhere even after hundreds of thousands of iterations. Given too small weights, the variance of the input diminishes as it passes through each layer in the network, and eventually drops to a really low value thus failing

to work. Contrarily, a too-large weight leads to exploding gradients. Xavier initialization [184] has been proven as an effective fashion, and widely adopted by DLMs in recommendation [105], [114], [116], [130]. Typically, the weights are also initialized based on either a Uniform or Normal distribution, defined as¹¹,

$$\mathbf{W}_{ij} \sim U\left(-\frac{\sqrt{6}}{\sqrt{n_{in} + n_{out}}}, \frac{\sqrt{6}}{\sqrt{n_{in} + n_{out}}}\right);$$

$$\mathbf{W}_{ij} \sim \mathcal{N}\left(0, \frac{2}{n_{in} + n_{out}}\right), \quad (4)$$

where \mathbf{W}_{ij} is the network weight; n_{in}, n_{out} are the number of input and output neurons, respectively.

By analyzing the collected papers, we find that around 59% of them do not report the parameter initializers. Among the papers mentioning parameter initializers, 36% of them are based on a Normal distribution; 10% of them use a Uniform distribution; 18% of them adopt the Xavier Initialization; 13% of them employ the pre-trained embeddings (e.g., via BPRMF) for initialization; and the rest 23% utilize other methods. The impacts of different parameter initializers are investigated in Section IV-C3.

4) *Model Optimizer Selection*: Optimizer is used to update model parameters, thus minimizing the loss function; meanwhile loss function acts as a guide to the terrain telling optimizer if it is moving in the right direction to reach the bottom of the valley, i.e., global minimum. Different optimizers also affect the recommendation performance. By looking into the collected papers, we find that 23% of them do not report their optimizers. Among the papers mentioning used optimizers, 50% of them adopt Adam [134], [136], [182]; 23% of them use SGD [150], [176]; and the rest 27% employ other optimizers (e.g., ALS [35], AdaGrad [104] and RMSProp [85]).

Here, we discuss six commonly-selected optimizers as shown in Table 11 (Appendix, available in the online supplemental material). (1) Gradient Descent (**GD**) is a first-order optimization algorithm which is dependent on the first order derivative of a loss function. The parameters are then updated in the negative gradient direction to minimize the loss; (2) Stochastic Gradient Descent (**SGD**) is a variant of GD to update the model's parameters after computation of loss on each training example, whilst parameters are changed after calculating gradient on the whole dataset by GD; (3) Mini Batch Stochastic Gradient Descent (**MB-SGD**) is an improvement on both SGD and standard GD, where the dataset is divided into various batches and after every batch, the parameters are updated; (4) Adaptive Gradient (**AdaGrad**) is an algorithm for gradient-based optimization that adapts the learning rate to the parameters, performing smaller updates (i.e. low learning rates) for frequently parameters, whilst larger updates (i.e. high learning rates) for infrequent parameters; (5) Root Mean Square Propagation (**RMSProp**) is devised to resolve AdaGrad's radically diminishing learning rates, and divides the learning rate by the average of the exponential decay of squared gradients; and (6) Adaptive Moment Estimation (**Adam**) calculates the adaptive learning rate for each parameter

¹¹[Online]. Available: <https://pytorch.org/docs/stable/nn.init.html>

from estimates of first and second moments of the gradients. In addition to the decaying average of past squared gradients like RMSprop, it also keeps a decaying average of past gradients.

5) *Strategies to Avoid Over-Fitting*: In machine learning, different strategies are exploited to combat the issue of over-fitting, which refers to the model over-fits the training data, thus achieving poor performance on the validation or test data. As a matter of fact, the most widely used regularization techniques include regularization terms, dropout and early-stop mechanism.

Regularization. It is generally integrated into the loss function, so as to help avoid over-fitting while training a recommendation model. Two types of terms are mainly adopted, namely $L1$ and $L2$ regularization (i.e. norm). $L1$ norm is also known as Manhattan Distance, which is the most natural way of measure distance between vectors. It is the sum of the magnitudes of the vectors in a space, where all the components of the vector are weighted equally. $L2$ norm is the most popular norm, also known as the euclidean norm, which is the shortest distance between two points. Different from $L1$ norm, each component of the vector is squared for $L2$ norm, indicating that the outliers have more weighting, so it can skew results. The main difference between the $L1$ and $L2$ regularization lies in (1) $L1$ regularization attempts to estimate the median of the data, whereas $L2$ regularization tries to estimate the mean of the data to avoid over-fitting; and (2) $L1$ regularization helps in feature selection by eliminating less important features, which is helpful given a large number of feature points.

Dropout. It has been widely adopted in DLMs to help avoid over-fitting [185]. The key idea is to randomly drop units (along with their connections) from the neural network during training, which prevents units from co-adapting too much. Hence, an extra hyper-parameter, i.e., the probability of retaining a unit p , is introduced, controlling the intensity of dropout. For instance, $p = 1$ implies no dropout and low values of p mean more dropout. Smaller p could lead to under-fitting; whereas large p may not produce enough dropout to prevent over-fitting. Typical values of p for hidden units are in the range of 0.5 to 0.8 [185].

Early-stop Mechanism. Early-stop is also a form of regularization used to avoid over-fitting. A major issue with training recommenders (e.g., LFM and DLM) is in the choice of the number of training epochs to use. Too many epochs can lead to over-fitting of the training dataset, whereas too few may result in an under-fit model. Early-stop is a method that allows us to specify an arbitrary large number of training epochs and stop training once the model performance stops improving on a hold out validation dataset. To be more specific, if the validation loss stops decreasing for several epochs in a row, the training stops. Through analyzing on the collected papers, only 11% of them point out early-stop strategy is adopted in their papers. In our study, we also investigate the impacts of early-stop mechanism on recommendation evaluation.

6) *Hyper-Parameter Tuning Strategies*: Hyper-parameter tuning, including validation and searching strategies, plays a vital role in the training process of recommendation approaches, and greatly influences the final recommendation performance.

Validation Strategy. Through the paper analysis, we notice that more than 33% of papers directly tune hyper-parameters

according to the performance on the test set. That is to say, they use the same data to tune model parameters and evaluate the model performance. Information may thus leak into the model and overfit the historical data. As a matter of fact, besides the training and test sets, an extra validation set should be retained to help tune the hyper-parameters, which is called *nested validation*¹². With nested validation, the optimal hyper-parameter settings are obtained when the model achieves the best performance on the validation set. By doing so, the information leak issue is well avoided in the model training and evaluation process. Therefore, in our study, we adopt the nested validation strategy. To be more specific, we further select **10%** of records from the training set as the validation set for split-by-ratio; and for leave-one-out, we retain one record from the training set as the validation set to tune hyper-parameters. Finally, the performance on the test set is reported. Due to the computational requirements of certain baselines, we were unable to search the hyper-parameter space for cross-validation in a reasonable amount of time.

Searching Strategy. From our observation, almost all collected papers employ the most straightforward and simple method – *grid search* [73], [134] to find out the optimal parameter settings. In particular, each hyper-parameter is provided with a set of possible values (i.e., search space) based on the prior-knowledge, and the optimal setting is then obtained by traversing the entire search space. Suppose a model has m parameters, where each parameter has an average of n possible values, the model needs to be executed for n^m times to find out optimal settings for all parameters. Hence, grid search is more suitable for models with less hyper-parameters; otherwise, it may suffer from the combination explosion issue. To improve the tuning efficiency, other strategies have been introduced. Given the search space of each parameter, *random search* [186] randomly chooses trials for a pre-defined times (e.g., 30) instead of traversing the entire search space. It is able to find models that are as good or slightly worse but within a smaller fraction of the computation time. In contrast, *Bayesian HyperOpt* [187] is not a brute force but more intelligent technique compared to grid and random search. It makes use of information from past trials to inform the next set of hyper-parameters to explore, while not compromising the quality of the results [6]. Therefore, for each baseline, we adopt Bayesian HyperOpt to perform hyper-parameter optimization on **NDCG**, which is the most popular metrics as shown in Fig. 2(c); and other metrics are expected to be simultaneously optimized with the optimal results on NDCG.

D. Categorization on Evaluation Modes

Based on the model-independent and model-dependent hyper-factors introduced in Sections III-B-III-C, we define four modes of rigorous evaluation as below.

- *Relax Mode* keeps exactly the same settings for all model-independent hyper-factors and follows the original settings as per individual approach for model-dependent hyper-factors.

¹²[Online]. Available: https://scikit-learn.org/stable/auto_examples/model_selection/plot_nested_cross_validation_iris.html

- *Hard-strict Mode* keeps exactly the same settings for both model-independent and model-dependent hyper-factors for all approaches.
- *Soft-strict Mode* keeps exactly the same settings for all model-independent hyper-factors and empirically finds out the optimal settings for per individual approach for model-dependent hyper-factors.
- *Mixed Mode* keeps exactly the same settings for all model-independent hyper-factors; while applies hard-strict mode on some model-dependent hyper-factors (e.g., the same initializer/optimizer), and relax or soft-strict modes on the others (e.g., empirically searching the optimal hyper-parameter settings for different baselines).

Through analysis, we find that most of the collected papers adopt the mixed mode for evaluation [42], [46], [116], for example, using the same model optimizer and parameter initializer for all approaches; adopting different loss functions as indicated in the original papers; and empirically finding out best parameter settings for all approaches. Regardless of different modes, it is essential to keep exactly the same settings for all model-independent hyper-factors for a rigorous evaluation. W.r.t. the model-dependent hyper-factors, different modes have their own pros and cons.

- Relax mode can extremely reduce the cost on exploring the optimal performance. However, it relies on the original settings indicated by the individual approach, which are not always available and may lead to a unfair comparison, e.g., one model defeats the others merely because it adopts a different loss function.
- Hard-strict mode ensures a fair comparison among different baselines, while it may not always be reasonable to, e.g., have the same settings for all shared hyper-parameters for all baselines, as the optimal hyper-parameter settings for different baselines may vary a lot across different datasets.
- Soft-strict mode could help find out the optimal performance for per individual approach, which, however, may be quite expensive due to the large amount of combinations of model-dependent hyper-factors.
- Mixed mode would be a better balance among complexity, performance and fairness for per individual approach. For instance, soft-strict mode can be applied regarding, e.g., optimal hyper-parameter settings, to maintain model performance; hard-strict mode can be used for, e.g., parameter initializer and model optimizer, to ensure fair comparison and less exploration complexity.

In summary, *mixed mode* could be the most practical way for achieving rigorous evaluation in recommendation.

IV. EXPERIMENTAL STUDY

A. Introduction to DaisyRec 2.0

To support the empirical study, we release a user-friendly Python toolkit named as **DaisyRec 2.0**, where Daisy is short for ‘Multi-Dimension fAIRly compARison for recommender SYstem’. Different from existing open-source libraries (e.g., LibRec [152], OpenRec [188] and DeepRec [153]), which mainly aim to reproduce various state-of-the-art recommenders,

DaisyRec 2.0 is designed with the goal of performing rigorous evaluation in recommendation by seamlessly accommodating the extracted hyper-factors in Section III. It is built upon the widely-used deep learning framework Pytorch (pytorch.org), and Fig. 4 depicts its overall structure consisting of four modules: GUI Command Generator, Loader, Recommender and Evaluator.

In particular, GUI Command Generator¹³ is used to help generate tune and test commands in a more user-friendly fashion. Taking the tune command generator as an example, users first need to select values for the basic settings (e.g., algorithm name and dataset) from a drop-down menu. Based on the selected algorithm, it then automatically displays the algorithm-specific parameters (e.g., KL regularization for Multi-VAE). Accordingly, users can select and set the search space for the algorithm-specific parameters. Lastly, it generates the corresponding tune command (shown in Fig. 5) based on all selected settings, which can be directly copied and pasted into the terminal to run.

Loader mainly aims to: (1) load and pre-process the dataset; (2) split it into training and test sets based on the selected Splitter; (3) divide validation set from training set by choosing proper Splitter according to Step 2; (4) sample negative items for training by choosing different samplers; and (5) convert the data into the specific format to fit the Recommender. Four components are included in Recommender, where ‘Algorithms’ implements the ten selected state-of-the-arts in Section III-B4 (more recommenders will be implemented); ‘LossSelector’ makes it flexible to change different objective functions for the algorithms; ‘ParameterInitializer’ allows to select different initialization methods (e.g., Xavier uniform distribution); and ‘Regularizer’ provides options for different regularization terms to avoid overfitting (e.g., $L1$ and $L2$). Evaluator is equipped with ‘Tuner,’ ‘ModelOptimizer,’ ‘Metric,’ and ‘Early-stop,’ where ‘Tuner’ helps accomplish hyper-parameter optimization; ‘ModelOptimizer’ provides options for different optimizers; ‘Metric’ implements the classic ranking metrics, e.g., Precision; and ‘Early-stop’ helps further avoid over-fitting.

To sum up, all modules in DaisyRec 2.0 are wrapped friendly for users to deploy, and new algorithms can be easily added into this extensible and adaptable framework. We keep DaisyRec 2.0 updated by adding more features.

B. Analysis on Model-Independent Hyper-Factors

1) *Impacts of Dataset Pre-Processing:* To study the impacts of pre-processing strategies (origin, 5- and 10-filter), we adopt Bayesian HyperOpt to perform hyper-parameter optimization w.r.t. NDCG@10 for each baseline under each view on each dataset for 30 trails [6]. We keep original objective functions for each baseline (see Table IV), employ the uniform sampler, and adopt time-aware split-by-ratio at global level ($\rho = 80\%$) as the data splitting method. Besides, 10% of the latest training set is held out as the validation set to tune the hyper-parameters. Once the optimal hyper-parameters are decided, we feed the whole training set to train the final model and report the performance

¹³[Online]. Available: <https://daisyrec.netlify.app/>

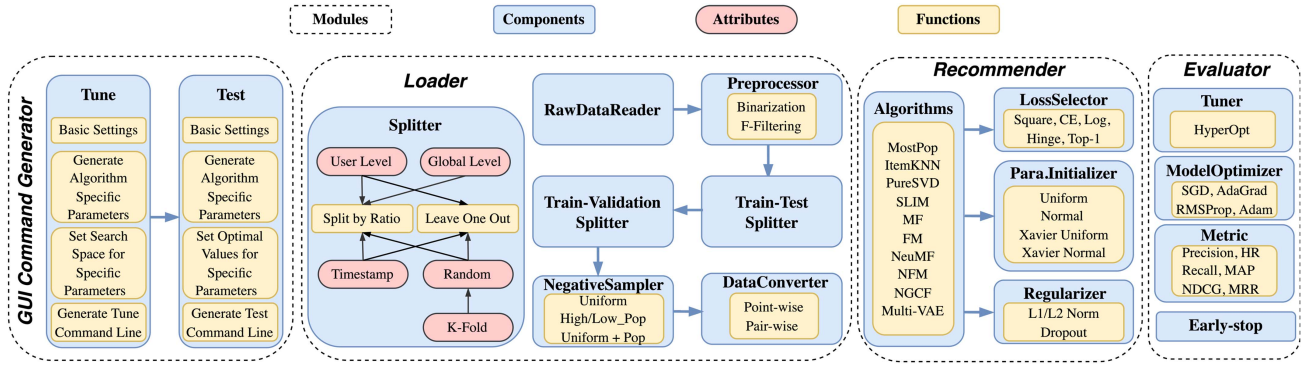


Fig. 4. The overall structure of DaisyRec 2.0, composed of four components, i.e., GUI Command Generator, Loader, Recommender, and Evaluator.

```
python hpo_tuner.py --tune_epochs=30 --problem_type=point --algo_name=multi-vae
--dataset=ml-1m --prepro=origin --test_method=lloo --val_method=lloo --loss_type=CL
--tune_pack="[\"batch_size\":[32,512],[32,64,128,512], \"choice\"], \"kl_reg\":
[0, 1, 0.1, \"float\"]]"
```

Fig. 5. An example of the generated tune command for Multi-VAE.

on the test set. Fig. 6 depicts the final results, where SLIM is omitted due to its extremely high computational complexity on large-scale datasets, which is unable to complete in a reasonable amount of time; and NGCF on Yelp and AMZe under origin view is also omitted because of the same reason (see Section IV-B3). Due to the space limitation, we only report the results on NDCG@10.

Overall, three different trends can be observed from the results: (1) the performance of different baselines keeps relatively stable on ML-1 M with varied settings; (2) on Book-X, Yelp and AMZe, the performance of all baselines generally climbs up; and (3) an obvious performance drop is observed on LastFM and Epinions. The most probable explanation is that although the density of the datasets increases (origin \rightarrow 5-filter \rightarrow 10-filter) as shown in Table II, the average length of the training sets for each user keeps stable on ML-1 M (86); increases on Book-X, Yelp and AMZe; and decreases on LastFM (39 \rightarrow 30 \rightarrow 27) and Epinions (35 \rightarrow 23 \rightarrow 20), as depicted by Table III. The more training data per user, the better a model can be trained, meaning that the more accurate performance can be achieved, and *vice versa*.

Regarding the performance of different baselines, several major findings can be noted as below. (1) Regarding MMs, MostPop performs the worst in most cases, showing the importance of personalization in recommendation; and ItemKNN is defeated by LFM and DLM, indicating the superiority of LFM and DLM on effective recommendation. However, on ML-1 M, the performance of MostPop exceeds that of ItemKNN, PureSVD and even BPRMF, demonstrating the potential of popularity in effective recommendation; and on LastFM, ItemKNN achieves the best performance compared with LFM and DLM. This implies that, the neighborhood-based idea, though simple, could be absorbed by LFM and DLM to further improve the recommendation accuracy [189]. (2) W.r.t. the three LFM, BPRMF generally performs better than PureSVD but worse

than BPRFM. Although PureSVD is simple – directly applying conventional singular value decomposition on the user-item interaction matrix, it sometimes achieves comparable and even better performance in comparison with BPRMF and BPRFM (see LastFM, Book-X and Yelp with 5/10-filter views). (3) For the four DLMs (i.e., NeuMF, NFM, NGCF and Multi-VAE), their performance varies across different datasets. For instance, NeuMF obtains the highest accuracy on Epinions; NFM is the winner on AMZe; NGCF defeats the others on both LastFM and Yelp; and Multi-VAE achieves extraordinary results on ML-1 M. However, they generally perform comparably to BPRFM across all datasets, and sometimes even worse than BPRFM (e.g., ML-1 M and Book-X). Besides, on LastFM, they even underperform ItemKNN. This is consistent with the previous findings [6] that DLMs are not always better than traditional methods with well-tuned parameters, but mostly cost much more in training as verified by Table V.

2) *Impacts of Dataset Splitting Methods*: We now test the impacts of different data splitting methods on the recommendation performance. To this end, we compare *random-* and *time-aware split-by-ratio* (i.e., **RSBR** vs. **TSBR**) at global level with $\rho = 80\%$ as well as *random-* and *time-aware leave-one-out* (i.e., **RLOO** vs. **TLOO**) on the 10-filter view. Note that we adopt Bayesian HyperOpt to perform hyper-parameter optimization w.r.t. NDCG@10 for each baseline under each data splitting method on each dataset for 30 trails. Meanwhile, LastFM and Book-X do not have the timestamp information, so their results on TSBR and TLOO for each baseline are omitted.

Fig. 7 displays the results of ten baselines on the six datasets. First of all, we can clearly observe that the performance of TSBR/SBR (split-by-ratio) is generally better than that of TLOO/LOO (leave-one-out). This could be largely affected by the different settings on the test procedure. To be specific, as introduced in Section III-B3, to improve the test efficiency, we randomly sample negative items for each user to ensure her test candidates to be 1,000, and then rank all test items among the 1,000 items w.r.t. both SBR and LOO. However, the number of positive items in the test set of SBR (> 1) is normally larger than that of LOO ($= 1$), thus leading to a higher accuracy of SBR. Second, baselines with RSBR/RLOO outperform those with TSBR/TLOO, especially on Epinions. The reason behind is that compared with random-aware split, time-aware split poses

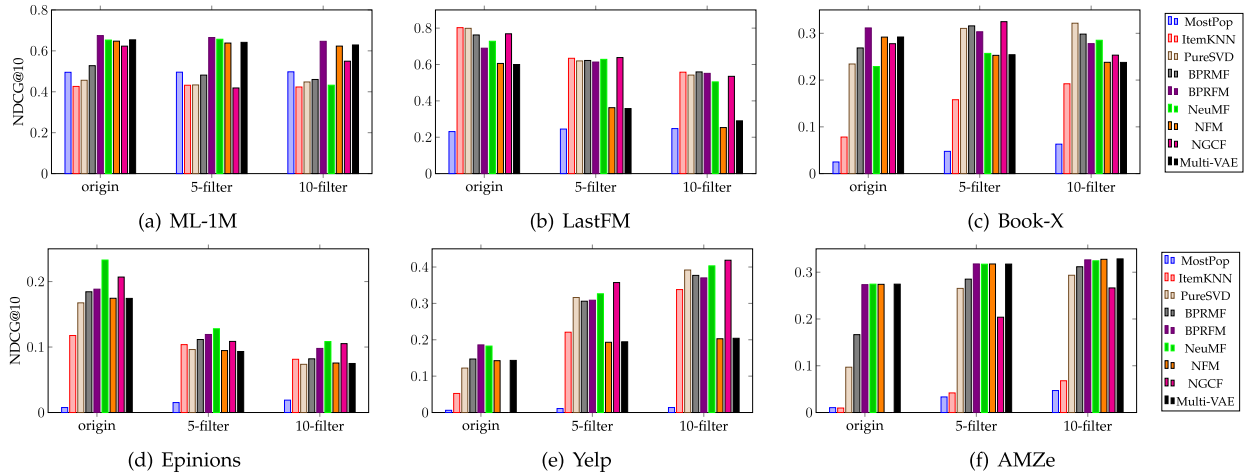


Fig. 6. Performance of baselines w.r.t. time-aware split-by-ratio on the six datasets across origin, 5- and 10-filter settings.

TABLE V
BASELINE COMPARISONS ON TRAINING TIME W.R.T. TIME-AWARE SPLIT-BY-RATIO ON THE 10-FILTER VIEW (SECONDS).

	MMs		LFMs				DLMs			
	MostPop	ItemKNN	PureSVD	BPRMF	BPRFM	SLIM	NeuMF	NFM	NGCF	Multi-VAE
ML-1M	0.0048	22.882	0.6808	2286.0	1847.4	62.924	9627.6	1682.6	2216.2	58.197
Lastfm	0.0010	2.4913	0.4728	246.35	103.35	5.1933	95.156	1178.6	538.06	31.807
Book-X	0.0055	29.013	3.1499	506.89	2129.0	860.71	5250.0	895.42	6251.6	101.00
Epinions	0.0070	25.847	2.8042	1497.0	3525.4	3283.7	4265.6	3387.0	9453.9	177.75
Yelp	0.0314	244.80	16.284	1566.4	6882.0	63990	5931.2	2459.4	38351	1388.6
AMZe	0.0167	121.96	1.7262	491.56	1258.4	18273	8625.5	2290.9	10099	1497.3

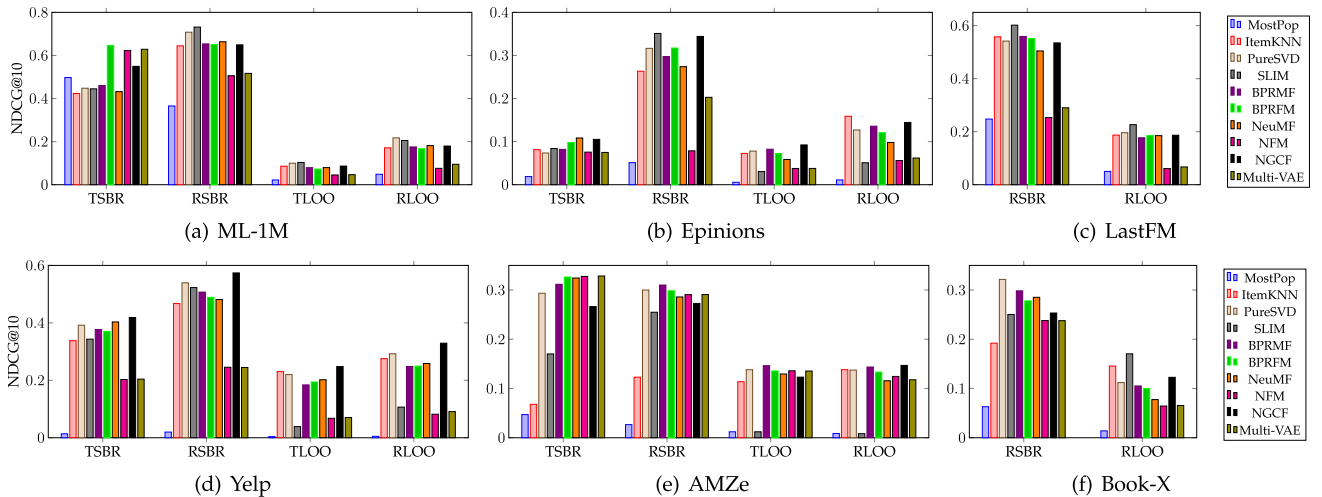


Fig. 7. Performance of baselines w.r.t. 10-filter on the six datasets with different data splitting methods.

a stronger constraint on the pattern of training and test data, thus increasing the training difficulty. However, this is more close to the real recommendation scenario, which strives to infer future by history. Our study also implies that the empirical results disclosed in previous studies using RSBR might be overestimated compared to those for real-world scenarios.

3) *Complexity of Comparison Baselines*: Table V shows the training time for the ten baselines on the six datasets with optimal hyper-parameters found by Bayesian HyperOpt on 10-filter view

via split-by-ratio. Note that, the optimal batch size may differ for different baselines, which may also affect the training time. All the experiments are executed on an Nvidia V100 GPU with 32 GiB memory, each running is paired with 11 Intel(R) Xeon(R) Platinum 8260 CPU (2.4 GHz) sharing 40 GiB memory.

According to Table V, we can note that *MostPop* is the fastest one in training, as it merely ranks all the items by the calculated popularity. *PureSVD* is the runner-up with time complexity $\mathcal{O}(\min\{m^2 f, n^2 f\})$, where m, n, f are the number

of users, items and singular values, respectively. Compared with other LFM and DLMs, it achieves a better balance between time complexity and ranking accuracy. Particularly, it performs comparably and sometimes even better than, e.g., BPRMF and NeuMF on LastFM, as depicted by Fig. 6, while its training time is hundreds or thousands times less than that of BPRMF and NeuMF as shown in Table V. Although the training efficiency of ItemKNN ranks third among all baselines with 10-filter setting, the time cost quadratically increases with origin setting due to its time complexity $\mathcal{O}(mn^2)$. Besides, the similarity matrix also takes up huge memory, for example, on the original AMZe ($n \approx 10^6$), it will cost $(64 \text{ bit} * 10^6 * 10^6)/10^{12} = 64 \text{ T}$ to save the similarity matrix. To ease this issue, we only keep the top-100 similar items for each target item in the memory.

The training time of BPRMF and BPRFM is comparable, where the time complexity for both methods is $\mathcal{O}(|\mathcal{R}|d)$, where \mathcal{R} is the total number of observed feedback and d is the dimension of latent factors. Similar to ItemKNN, the time cost of SLIM with 10-filter setting is acceptable, while it tremendously increases with origin setting due to its time complexity $\mathcal{O}(|R|n)$. Even with the 10-filter view, it takes the longest training time among all baselines on the two large datasets (i.e., Yelp and AMZe). Meanwhile, it also suffers from the huge memory cost issue because of the learned item similarity matrix. Hence, both ItemKNN and SLIM are not scalable for large-scale datasets. Regarding the four DLMs (i.e., NeuMF, NFM, NGCF and Multi-VAE), NFM and Multi-VAE are usually more efficient than NeuMF and NGCF regarding time complexity. Although DLMs yield comparable performance with LFM, they generally cost much more training time, especially on larger datasets. For example, on AMZe, the training time of NeuMF and NGCF is around 20 times larger than that of BPRMF.

4) *Correlations of Evaluation Metrics:* As discussed in Section IV-B1, we adopt Bayesian HyperOpt to perform hyper-parameter optimization for 30 trials via optimizing NDCG@10. However, six metrics are used in our study, namely Precision, Recall, HR, MAP, MRR and NDCG. The best hyper-parameter settings for optimal NDCG does not necessarily guarantee optimum w.r.t. the other five metrics. Hence, we study the correlation of different metrics when their respective optimums are achieved. In particular, for each baseline on each dataset with 10-filter view, the Bayesian HyperOpt executes 30 trials; we thus have 30 entries for the validation performance of the baseline correspondingly, where each entry includes the results on the six metrics, e.g., [Precision: 0.24; Recall: 0.07; HR: 0.57; MAP: 0.17; MRR: 0.76; NDCG: 0.42]. Due to the optimal results for the six metrics may not achieve simultaneously, we select the optimal one among the 30 entries for each metric, and ultimately obtain six entries, where each entry records the best result on the corresponding metric.

Based on the six selected entries of each method per dataset, we pair-wisely calculate and record the times that any two of them (e.g., NDCG and HR) can achieve their best results simultaneously entry by entry. For example, given the optimal entry for NDCG, we will check whether the rest five metrics (e.g., HR) in this entry are optimal or not. If yes, we will add one at the corresponding position (NDCG, HR) of the correlation

matrix; otherwise 0. The same rule is applied to the optimal entries for the other five metrics. Except MostPop, as it does not have any hyper-parameters, we accumulate the results of nine baselines across the six datasets ($9*6=54$), and ultimately obtain their correlation matrices regarding time-/random-aware split-by-ratio and leave-one-out as illustrated by Figures 8(a-d), where all values are normalized into the range of $[0,1]$ (divided by 54), and a darker color indicates a stronger correlation, that is, a higher probability of two metrics achieving their best results in the meanwhile.

The results help verify our argument that best hyper-parameter settings for optimal NDCG cannot ensure optimal results for all the other five metrics. Several detailed findings can be noted. (1) The correlation matrix is asymmetrical, for instance, the correlation for (NDCG, HR, 0.69) is higher than (HR, NDCG, 0.61) as shown in Fig. 8(a). That is to say, the probability of a model with best NDCG to achieve the best HR is higher than that of a model with best HR to reap the optimal NDCG. (2) Similar trends can be noticed within a same base data splitting method (i.e., TSBR and RSBR; TLOO and RLOO) no matter whether the timestamp information is considered or not; whilst the patterns are different across different base splitting fashions (i.e., TSBR/RSBR and TLOO/RLOO). (3) Regarding TSBR/RSBR, the best MRR and MAP are more easily to be guaranteed concurrently, e.g., (MRR, MAP, 0.81) and (MAP, MRR, 0.80) in Fig. 8(a); and (MRR, MAP, 0.89) and (MAP, MRR, 0.78) in Fig. 8(b). (4) For TLOO/RLOO, Precision, Recall and HR are more likely to be optimized simultaneously; and MAP, MRR and NDCG have a higher probability to reach their peaks together. This is mainly due to only one positive item inside the test set for each user; consequently, Recall and HR are equivalent, which is positively correlated with Precision; meanwhile, MAP, MRR and NDCG are also positively correlated with each other.

Additionally, we examine the Kendall's correlation [83] among metrics in terms of indicating recommendation performance on the ten baselines across the six datasets under 10-filter view with different data splitting methods. The results are depicted in Figures 8(e-h), where a darker color (a stronger correlation) implies that the metrics produce more identical ranking. We find that (1) different from Figures 8(a-d), the Kendall's correlation matrix is symmetrical; (2) similarly, the trends are consistent within a same base data splitting method, e.g., Figures 8(g-h), while vary slightly across different base splitting ways, e.g., Figures 8(e) and 8(g); and (3) a common observation across Figures 8(e-h) is that MAP, MRR and NDCG are more likely to generate consistent ranking. Besides, for TSBR/RSBR, (Precision, NDCG) and (Recall, HR) show a fairly strong correlation, while w.r.t. TLOO/RLOO, (Precision, Recall, HR) exhibits obvious correlation, which is also caused by the single positive item inside the test set for each user as explained previously. In summary, a convincing and solid evaluation should be performed w.r.t. more diverse metrics.

C. Analysis on Model-Dependent Hyper-Factors

1) *Impacts of Loss Functions:* To examine the impacts of different objective functions, we adopt the optimal parameters

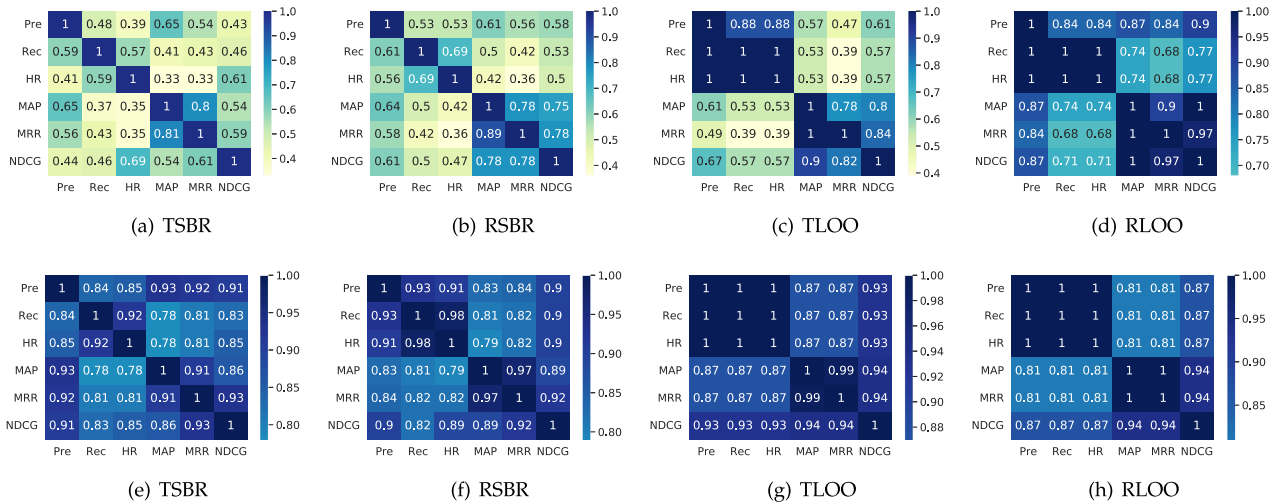


Fig. 8. The correlations of evaluation metrics w.r.t. different data splitting methods on 10-filter. ‘Pre’ and ‘Rec’ are Precision and Recall, respectively.

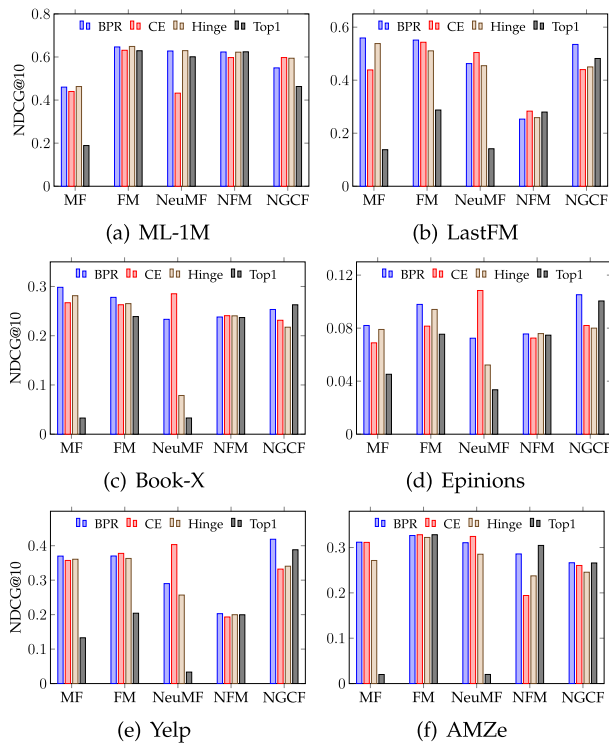


Fig. 9. Performance of baselines w.r.t. time-aware split-by-ratio on 10-filter view across the six datasets with different loss functions.

for the baselines found on 10-filter view with time-aware split-by-ratio in Section IV-B1, and only vary objective functions for MF, FM, NeuMF, NFM and NGCF. The results are depicted in Fig. 9, where BPR (pair-wise log loss), CE (point-wise cross entropy loss), Hinge (pair-wise hinge loss) and Top1 (pair-wise top1 loss) correspond to $\mathcal{L}_{pai} + f_{ll}$, $\mathcal{L}_{poi} + f_{cl}$, $\mathcal{L}_{pai} + f_{hl}$ and $\mathcal{L}_{pai} + f_{tl}$ in Table IV, respectively. Several conclusions can be drawn. As a whole, for different baselines on the six datasets, (1) BPR loss generally achieves the best performance; (2) CE loss and Hinge loss perform comparably; and (3) Top1 loss

possesses the largest performance variation. From the perspective of different baselines, (1) MF and FM usually achieve the best performance with BPR loss; (2) NeuMF performs better with CE loss in most cases; (3) NFM is relatively less sensitive to different losses; and (4) NGCF generally obtains better accuracy with either BPR or Top1 loss.

2) *Impacts of Negative Sampling Strategies:* We now explore the impact of different negative samplers, i.e., uniform (U), high-popularity (HP), low-popularity (LP), uniform+high-popularity (U+HP) and uniform+low-popularity (U+LP) on BPRMF, BPRFM, NeuMF, NFM and NGCF across the six datasets under 10-filter view with time-aware split-by-ratio. To this end, we only vary negative samplers for the baselines while keeping other parameters fixed. First, the uniform sampler, though simple, achieves comparable performance in comparison with popularity samplers (HP and LP) as illustrated in Fig. 10. Intuitively, users may not tend to buy the less popular items, that is, the items with low popularity are more likely to be the negative items for users. However, it is overturned by the empirical results. Second, U+HP and U+LP samplers are generally defeated by U/HP/LP samplers. However, there are some exceptions, e.g., BPRMF on ML-1M and NeuMF on Yelp. Lastly, U+LP exceeds U+HP in most cases, indicating that generally the popular items have a lower probability to be negative items than the less popular ones.

3) *Impacts of Parameter Initializers:* To study the impact of different parameter initializers, we compare the results of six baselines (BPRMF, BPRFM, NeuMF, NFM, NGCF and MultiVAE) across the six datasets under 10-filter view with time-aware split-by-ratio. Specifically, for BPRMF and BPRFM, we adopt uniform ($\alpha = 1$) and normal distribution ($\sigma = 0.01$) for initialization; while for the rest four deep learning baselines, we utilize Xavier uniform and normal distribution for initialization. As depicted in Fig. 11, we can note that (1) for the two LFM, i.e., BPRMF and BPRFM, initializer with normal distribution dramatically beats that with uniform distribution; and (2) for the four DLMs, some baselines (e.g., NGCF) gain better accuracy with uniform distribution than normal distribution on the six

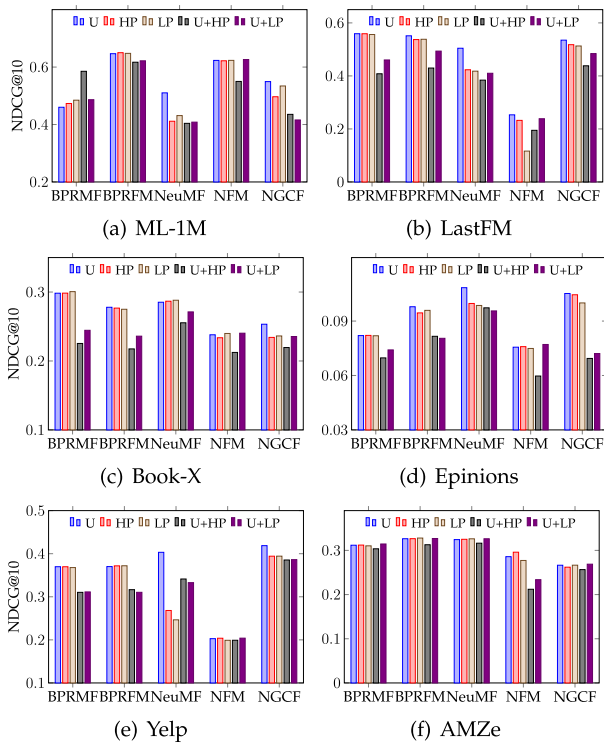


Fig. 10. Performance of baselines w.r.t. time-aware split-by-ratio on 10-filter view across the six datasets with different sampling strategies.

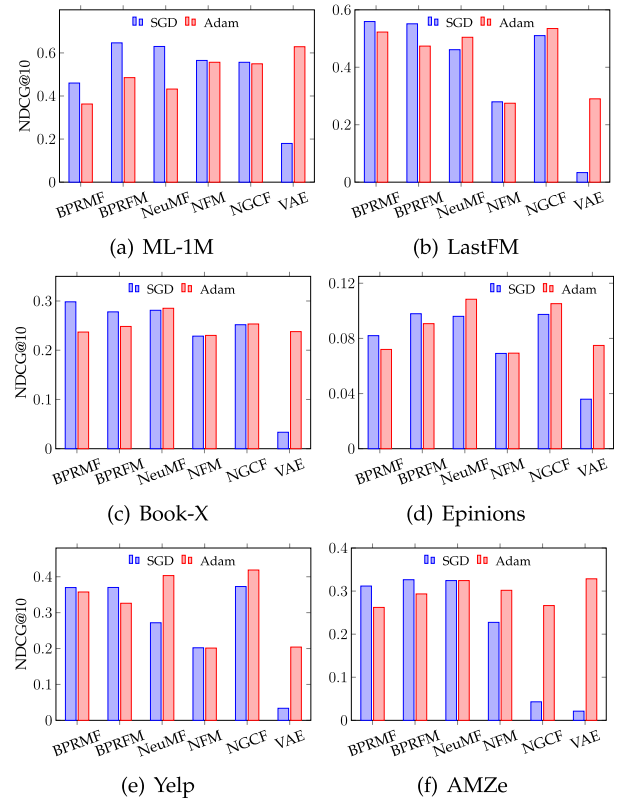


Fig. 12. Performance of baselines w.r.t. time-aware split-by-ratio on 10-filter view across the six datasets with different optimizers.

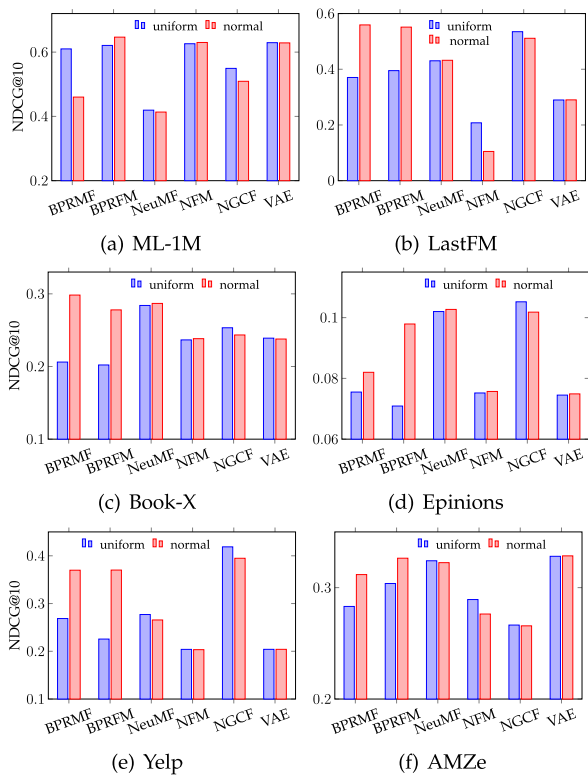


Fig. 11. Performance of baselines w.r.t. time-aware split-by-ratio on 10-filter view across the six datasets with different initializers.

datasets; while some perform comparably with the two types of initializers, e.g., Multi-VAE, across the six datasets. In a nutshell, different parameter initializers produce different recommendation performance. With a proper initializer, the LFM may easily beat DLMs, for example, BPRMF defeats DLMs on LastFM and Book-X.

4) *Impacts of Model Optimizers*: We further investigate the impacts of different optimizers on the final recommendation performance. In particular, we vary optimizers (i.e., SGD, and Adam) for the six baselines (i.e., BPRMF, BPRFM, NeuMF, NFM, NGCF and Multi-VAE) on the six datasets under 10-filter view with time-aware split-by-ratio. The results are presented in Fig. 12, where we observe that a better performance is achieved via SGD in comparison with Adam for LFM (i.e., BPRMF and BPRFM); whereas Adam generally outperforms SGD regarding DLM (i.e., NeuMF, NFM, NGCF and Multi-VAE).

5) *Impacts of Strategies to Avoid Over-Fitting*: As illustrated in Section III-C5, regularization term, dropout and early-stop mechanism are widely adopted to avoid over-fitting. To verify their impacts, we compare the results of six baselines (BPRMF, BPRFM, NeuMF, NFM, NGCF and Multi-VAE) across the six datasets under 10-filter view with time-aware split-by-ratio by removing these strategies. In particular, +all, -L2, -dropout and -ES respectively indicate the baseline with all over-fitting prevention strategies, variant without L2 regularization term, variant without dropout (only for deep learning baseline), and variant without early-stop. The results are displayed in Fig. 13,

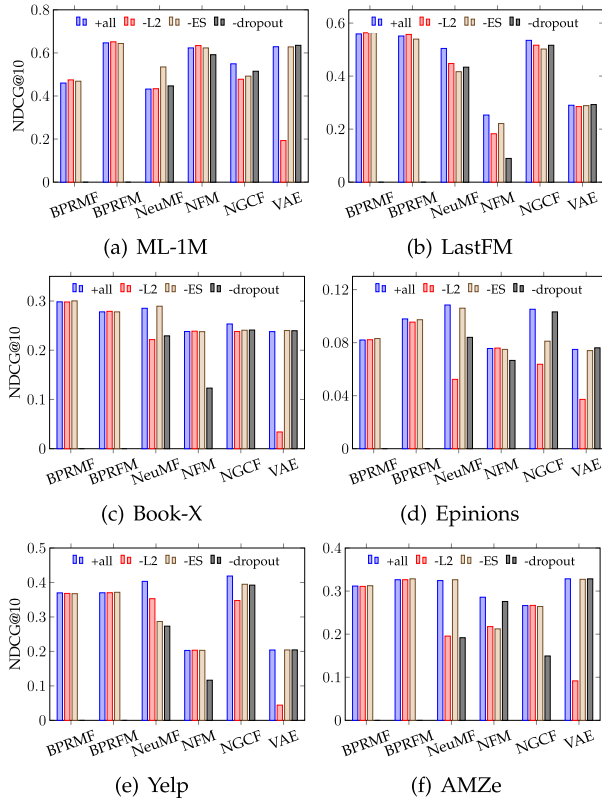


Fig. 13. Performance of baselines w.r.t. time-aware split-by-ratio on 10-filter across the six datasets with different strategies to avoid over-fitting.

where several observations can be noted. First, the over-fitting prevention strategies generally facilitate to enhance the recommendation accuracy to some extent for all baselines across the six datasets. However, there are a few exceptions, e.g., NeuMF on ML-1 M, indicating some of these strategies may also lead to the under-fitting issue occasionally. Second, the impact of these strategies is more significant on DLMs (e.g., NeuMF) than LFM (e.g., BPRMF). Lastly, the performance of DLMs may be remarkably affected by a certain strategy, e.g., NFM is heavily affected by dropout, whilst a major impact of L2 regularization term on Multi-VAE can be observed.

6) *Impacts of Hyper-Parameter Tuning Strategies:* As illustrated in Section III-C6, a validation set should be held out for hyper-parameter tuning to avoid data leakage. To investigate its impact, we compare with the results by directly tuning hyper-parameters on the test set under 10-filter view with time-aware split-by-ratio. For simplicity, we only select four representative baselines on four datasets as displayed in Fig. 14. Accordingly, we can easily find that in most cases directly tuning hyper-parameters on the test set indeed guarantees a better performance compared with tuning hyper-parameters on the validation set. This implies that the empirical results reported in existing studies without a validation set might be overestimated.

D. Summary and Further Discussion

For ease of reading, Table VI summarizes the most important findings w.r.t. the impact of different hyper-factors on

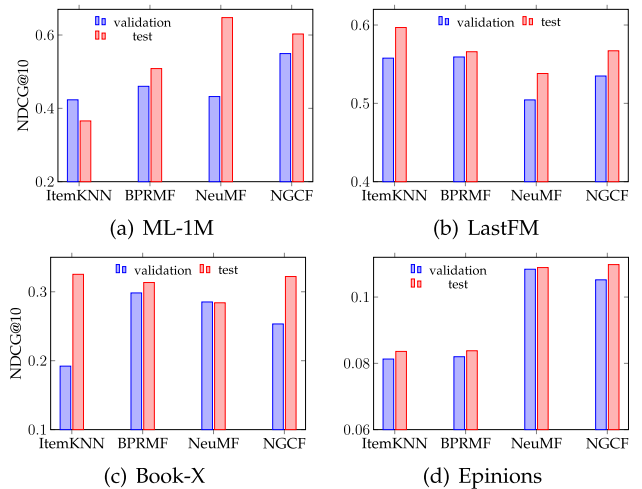


Fig. 14. Performance of baselines w.r.t. time-aware split-by-ratio on 10-filter across the six datasets by tuning on validation and test sets.

recommendation evaluation in Sections IV-B and C. These findings seek to provide invaluable instructions and guidance for both researchers and practitioners in the area of recommender systems. Moreover, we further perform the horizontal comparison to examine which factors affect more for the evaluation. To this end, regarding each hyper-factor, we calculate and average the relative performance gap (i.e., $(\text{best_result} - \text{worst_result}) / \text{worst_result}$) for each baseline across all the datasets shown in Table VII, where ‘DP, DS, LF, NS, MI, MO, OP, HT’ are short for Dataset Pre-processing, Dataset Splitting, Loss Function, Negative Sampling, Model Initializer, Model Optimizer, Over-fitting Prevention and Hyper-parameter Tuning, respectively. Note that a larger relative performance gap indicates a higher impact of the corresponding hyper-factor. As such, the result highlighted in bold in each row suggests the corresponding hyper-factor has the highest impact; ‘-’ denotes the result is not available; and the last row ‘Average’ implies the overall impact generated by each hyper-factor regardless of the baselines. According to the table, we find that (1) the impacts of different hyper-factors vary a lot across different baselines, for instance, the most impactful hyper-factor for BPRMF and NeuMF is ‘Loss Function,’ whereas ‘Model Optimizer’ affects Multi-VAE the most; and (2) as a whole, ‘Dataset Splitting,’ ‘Loss Function’ and ‘Model Optimizer’ are the top-3 impactful hyper-factors in recommendation evaluation based on our empirical study.

V. BENCHMARKING RECOMMENDATION

A. Standardized Procedures

Section III shows the hyper-factors in recommendation evaluation, and their impacts are empirically analyzed in Section IV. To achieve a rigorous evaluation, the *mixed mode* discussed in Section III-D is encouraged to be adopted. Accordingly, we propose a series of standardized procedures and correspondingly call for endeavors of all researchers, aiming to effectively

TABLE VI
SUMMARY OF IMPORTANT FINDINGS IN SECTION IV-B AND SECTION IV-C

Section	Hyper-factors	Important Findings
4.2.1	Dataset Pre-processing	The more training data per user, the better a model can be trained; The performance of different methods varies a lot across different datasets, e.g., traditional ones can defeat DLMs on some cases.
4.2.2	Dataset Splitting	SBR and LOO can achieve different performance, which relies on the number of negative samples in the test procedure; Baselines with random-aware split outperform those with time-aware split; Empirical results disclosed in previous studies using RSBR might be overestimated compared to those for real-world scenarios.
4.2.3	Comparison Baseline	PureSVD achieves a better balance between time complexity and ranking accuracy among all selected baselines; Both ItemKNN and SLIM are not scalable for large-scale datasets due to the high time and memory cost; DLMs yield comparable performance with LFMs, but possess higher complexity.
4.2.4	Evaluation Metric	Best parameter settings for optimizing one metric, e.g., NDCG, cannot guarantee optimal results for other metrics; For SBR, the best MRR and MAP are more easily to be guaranteed concurrently; For LOO, Precision, Recall and HR are more likely to be optimized simultaneously; MAP, MRR and NDCG have a higher probability to reach their peaks together.
4.3.1	Loss Function	BPR loss generally performs the best; CE and Hinge loss perform comparably; Top1 loss has the largest performance variation; Different baselines usually achieve their best performance with different loss functions.
4.3.2	Negative Sampling	The uniform sampler, though simple, achieves comparable performance with popularity samplers; Popular items have a lower probability to be negative items than the less popular ones.
4.3.3	Parameter Initializer	For LFMs, the initializer with normal distribution beats that with uniform distribution; For DLMs, some gain better accuracy with uniform distribution; while some perform comparably with the two types of initializers; Different initializers produce different performance; with a proper initializer, LFMs may easily defeat DLMs.
4.3.4	Model Optimizer	For LFMs, better performance is achieved via SGD compared with Adam; For DLMs, Adam generally outperforms SGD.
4.3.5	Over-fitting Prevention	Over-fitting prevention strategies (regularization term, dropout and early-stop) generally help enhance the accuracy; The impact of these over-fitting prevention strategies is more significant on DLMs than LFMs; The performance of DLMs may be remarkably affected by a certain strategy, e.g., dropout for NFM and L2 term for Multi-VAE.
4.3.6	Hyper-Parameter Tuning	Directly tuning hyper-parameters on the test set gains a better performance compared with tuning them on the validation set; Empirical results reported in existing studies without a validation set might be overestimated.

TABLE VII
THE RELATIVE PERFORMANCE GAP BETWEEN THE BEST AND WORST RESULT OF DIFFERENT BASELINES REGARDING VARIOUS HYPER-FACTORS

	DP	DS	LF	NS	MI	MO	OP	HT
MostPop	131%	762%	-	-	-	-	-	-
ItemKNN	232%	225%	-	-	-	-	-	24%
PureSVD	107%	261%	-	-	-	-	-	-
BPRMF	73%	280%	496%	23%	31%	16%	1%	5%
BPRFM	42%	302%	38%	18%	32%	16%	1%	-
NeuMF	63%	311%	653%	25%	1%	20%	50%	14%
NFM	60%	399%	14%	36%	18%	6%	68%	-
NGCF	55%	245%	23%	22%	4%	91%	32%	12%
Multi-VAE	55%	452%	-	-	0.24%	615%	261%	-
Average	91%	360%	245%	25%	14%	127%	69%	14%

enhance the standardization of recommendation evaluation. Regarding model-independent hyper-factors, five procedures are recommended.

- It is impossible to evaluate recommenders on all public datasets covering each domain. However, at least one widely-used dataset discussed in Section III-B1 should be considered, especially for the papers evaluated on the private datasets (e.g., confidential data from commercial companies). Otherwise, the results could not be easily reproduced by the subsequent studies.
- Section IV-B1 verifies that different data pre-processing strategies impact the performance. Besides origin view, 5- and 10-filter views are recommended to ease the data sparsity issue, and a clear description on data pre-processing details is indispensable.
- For data splitting methods, both time-aware split-by-ratio and time-aware leave-one-out are recommended. With timestamp, the real recommendation scenario will be better simulated. W.r.t. split-by-ratio, both global- and user-level work well and $\rho = 80\%$ is recommended for a more feasible and convenient comparison.

- The representative baselines with different types (MMs, LFMs and DLMs) in Section III-B4 are recommended to be selected and compared. As shown in Section IV-B1, the performance of different types of baselines vary a lot in different scenarios, that is, the MMs (e.g., MostPop) and simple LFMs (e.g., PureSVD) sometimes even perform better than DLMs (e.g., NeuMF). The more diverse baselines are compared, the more comprehensive and reliable the evaluation is.
- At least two of the six discussed metrics in Section III-B5 should be adopted, where one (e.g., Precision) measures whether a test item is present on the top-N recommendation list, and the other (e.g., NDCG) measures the ranking positions of the recommended items.

With respect to model-dependent hyper-factors, there are also five procedures recommended as below.

- For a fair comparison, it is better to evaluate all methods with a same type of objective functions and thus better positioning a proposed method's contributions.
- All the compared methods should adopt the same negative sampler, except the papers with the goal of proposing or studying different negative sampling strategies.
- The parameter initializer and model optimizer should be consistent across all compared methods as demonstrated in Section IV-C3 and Section IV-C4.
- The same basic overfitting prevention strategies should be applied to all compared methods, except the methods with specially-designed strategies, e.g., the message dropout in NGCF [102].
- With regards to the hyper-parameter tuning, a nested validation is mandatory, that is, retaining partial (e.g., 10%) training data as validation set. Bayesian HyperOpt, as a more intelligent parameter searching strategy, is recommended, and the search space should be kept the same for

the shared parameters of different baselines. The number of trails (we set 30 in this study by following [6]) may be increased for further performance improvements. Most importantly, the optimal parameter settings should be well reported for reproduction.

Meanwhile, the source codes and datasets for each publication should be available for reproduction [190]. The conference venues could make them as necessities, measure the quality, and even require a short code demonstration along with each accepted paper during the conference.

B. Performance of Baselines

With the goal of providing a better reference for fair comparison, Tables 12-17 (Appendix, available in the online supplemental material) show the performance of ten baselines across six metrics on the six datasets under three different views (i.e., origin, 5-filter and 10-filter) with time-aware split-by-ratio ($N = 10$). Due to space limitation, other results (e.g., $N = 1, 5, 20, 30, 50$ and other data splitting methods) are on our GitHub. All optimal hyper-parameters are found by Bayesian HyperOpt to optimize NDCG@10 for 30 trials (see Section IV-B1), and the corresponding detailed parameter settings are shown in Tables 18-21.

Based on the results, several major observations can be noted. **(1)** BPRFM achieves the best performance on ML-1 M across all views. **(2)** Regarding LastFM, ItemKNN/NGCF performs the best on origin and 5-filter views, while SLIM achieves the best performance on 10-filter view. **(3)** For Book-X, BPRFM/NGCF and PureSVD/NGCF respectively beat other baselines on origin and 5-filter views, and PureSVD is the winner on 10-filter view. **(4)** W.r.t. Epinions, NeuMF obtains the highest accuracy on origin view; and NeuMF/NGCF helps reach the best performance on 5- and 10-filter views. **(5)** On Yelp, the best performance on the origin, 5- and 10-filter views are respectively gained by BPRFM, NGCF and NGCF. **(6)** With regards to AMZe, NeuMF/Multi-VAE defeat other baselines on origin view; BPRFM/NFM obtains the optimal results on 5-filter view; and Multi-VAE is the top method on 10-filter view.

VI. CONCLUSION

This paper aims to benchmark recommendation for reproducible evaluation and fair comparison from the angles of both practical theory analysis and empirical study. Regarding theory analysis, 141 recommendation papers published in the four recent years (2017-2020) from eight top tier conferences have been systematically reviewed, whereby we define and extract the hyper-factors affecting recommendation evaluation, classified into model-independent (e.g., dataset splitting methods) and -dependent (e.g., loss function design) factors. Accordingly, different modes for rigorous evaluation are defined and discussed in-depth. To support the empirical study, a user-friendly Python toolkit – DaisyRec 2.0 has been released and updated by seamlessly accommodating the extracted hyper-factors. Thereby, the impacts of different hyper-factors on evaluation are then empirically examined and comprehensively analyzed. Lastly, we create benchmarks for rigorous evaluation by proposing standardized

procedures and providing the performance of ten well-tuned state-of-the-art algorithms on six widely-used datasets across six metrics as a reference for later study. For the future work, we plan to deepen our investigation by, for example, diving into more diverse (e.g., session/sequential-aware) recommendation tasks, and more evaluation metrics (e.g., diversity, novelty and serendipity).

REFERENCES

- [1] Z. Sun et al., "Research commentary on recommendations with side information: A survey and research directions," *Electron. Commerce Res. Appl.*, vol. 37, 2019, Art. no. 100879.
- [2] A. Said and A. Bellogín, "Comparative recommender system evaluation: Benchmarking recommendation frameworks," in *Proc. ACM Conf. Recommender Syst.*, 2014, pp. 129–136.
- [3] A. Said and A. Bellogín, "Rival: A toolkit to foster reproducibility in recommender system evaluation," in *Proc. ACM Conf. Recommender Syst.*, 2014, pp. 371–372.
- [4] Z. Sun et al., "Are we evaluating rigorously? Benchmarking recommendation for reproducible evaluation and fair comparison," in *Proc. ACM Conf. Recommender Syst.*, 2020, pp. 23–32.
- [5] S. Rendle et al., "On the difficulty of evaluating baselines: A study on recommender systems," 2019, *arXiv:1905.01395*.
- [6] M. F. Dacrema et al., "Are we really making much progress? A worrying analysis of recent neural recommendation approaches," in *Proc. ACM Conf. Recommender Syst.*, 2019, pp. 101–109.
- [7] B. Sarwar et al., "Item-based collaborative filtering recommendation algorithms," in *Proc. World Wide Web Conf.*, 2001, pp. 285–295.
- [8] Z. Sun et al., "Exploiting both vertical and horizontal dimensions of feature hierarchy for effective recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2017, pp. 189–195.
- [9] D. Li et al., "ERMMA: Expected risk minimization for matrix approximation-based recommender systems," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2017, pp. 1403–1409.
- [10] L. Yu et al., "Walkranker: A unified pairwise ranking model with multiple relations for item recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2018, pp. 2596–2603.
- [11] M. Wang et al., "Collaborative filtering with social exposure: A modular approach to social recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2018, pp. 2516–2523.
- [12] T. D. T. Do and L. Cao, "Coupled poisson factorization integrated with user/item metadata for modeling popular and sparse ratings in scalable recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2018, pp. 2918–2925.
- [13] J. Zhang et al., "Hierarchical reinforcement learning for course recommendation in MOOCs," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2019, pp. 435–442.
- [14] C. Wang et al., "CAMO: A collaborative ranking method for content based recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2019, pp. 5224–5231.
- [15] C. Lin et al., "Non-compensatory psychological models for recommender systems," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2019, pp. 4304–4311.
- [16] J. Li et al., "From zero-shot learning to cold-start recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2019, pp. 4189–4196.
- [17] C. Liu et al., "Discrete social recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2019, pp. 208–215.
- [18] Z.-H. Deng et al., "DeepCF: A unified framework of representation learning and matching function learning in recommender system," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2019, pp. 61–68.
- [19] L. Hu et al., "HERS: Modeling influential contexts with heterogeneous relations for sparse and cold-start recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2019, pp. 3830–3837.
- [20] X. Wang et al., "Explainable reasoning over knowledge graphs for recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2019, pp. 5329–5336.
- [21] T. Shen et al., "PEIA: Personality and emotion integrated attentive model for music recommendation on social media platforms," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2020, pp. 206–213.
- [22] Q. Zhu et al., "A knowledge-aware attentional reasoning network for recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2020, pp. 6999–7006.

- [23] C. Chen et al., "Efficient heterogeneous collaborative filtering without negative sampling for recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2020, pp. 19–26.
- [24] G. Guo et al., "Leveraging title-abstract attentive semantics for paper recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2020, pp. 67–74.
- [25] M. Li et al., "Symmetric metric learning with adaptive margin for recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2020, pp. 4634–4641.
- [26] Y. Xu et al., "Multi-feature discrete collaborative filtering for fast cold-start recommendation," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2020, pp. 270–278.
- [27] J. Chen et al., "Fast adaptively weighted matrix factorization for recommendation with implicit feedback," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2020, pp. 3470–3477.
- [28] D. D. Le and H. W. Lauw, "Stochastically robust personalized ranking for LSH recommendation retrieval," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2020, pp. 4594–4601.
- [29] C. Wang et al., "SetRank: A setwise Bayesian approach for collaborative ranking from implicit feedback," in *Proc. Conf. Assoc. Advance. Artif. Intell.*, 2020, pp. 6127–6136.
- [30] Y. Zhang et al., "Joint representation learning for top-n recommendation with heterogeneous information sources," in *Proc. Conf. Inf. Knowl. Manage.*, 2017, pp. 1449–1458.
- [31] D. D. Le and H. W. Lauw, "Indexable Bayesian personalized ranking for efficient top-k recommendation," in *Proc. Conf. Inf. Knowl. Manage.*, 2017, pp. 1389–1398.
- [32] W. Pei et al., "Interacting attention-gated recurrent networks for recommendation," in *Proc. Conf. Inf. Knowl. Manage.*, 2017, pp. 1459–1468.
- [33] L. Mei et al., "An attentive interaction network for context-aware recommendations," in *Proc. Conf. Inf. Knowl. Manage.*, 2018, pp. 157–166.
- [34] H. Wang et al., "RippleNet: Propagating user preferences on the knowledge graph for recommender systems," in *Proc. Conf. Inf. Knowl. Manage.*, 2018, pp. 417–426.
- [35] T. Tran et al., "Regularizing matrix factorization with user and item embeddings for recommendation," in *Proc. Conf. Inf. Knowl. Manage.*, 2018, pp. 687–696.
- [36] J. Ma et al., "DBRec: Dual-bridging recommendation via discovering latent groups," in *Proc. Conf. Inf. Knowl. Manage.*, 2019, pp. 1513–1522.
- [37] W.-C. Kang and J. McAuley, "Candidate generation with binary codes for large-scale Top-N recommendation," in *Proc. Conf. Inf. Knowl. Manage.*, 2019, pp. 1523–1532.
- [38] F. Xu et al., "Relation-aware graph convolutional networks for agent-initiated social e-commerce recommendation," in *Proc. Conf. Inf. Knowl. Manage.*, 2019, pp. 529–538.
- [39] B. Chang et al., "Learning graph-based geographical latent representation for point-of-interest recommendation," in *Proc. Conf. Inf. Knowl. Manage.*, 2020, pp. 135–144.
- [40] B. Chen et al., "TGCN: Tag graph convolutional network for tag-aware recommendation," in *Proc. Conf. Inf. Knowl. Manage.*, 2020, pp. 155–164.
- [41] D. Lee et al., "News recommendation with topic-enriched knowledge graphs," in *Proc. Conf. Inf. Knowl. Manage.*, 2020, pp. 695–704.
- [42] R. Sun et al., "Multi-modal knowledge graphs for recommender systems," in *Proc. Conf. Inf. Knowl. Manage.*, 2020, pp. 1405–1414.
- [43] S. Kang et al., "DE-RRD: A knowledge distillation framework for recommender system," in *Proc. Conf. Inf. Knowl. Manage.*, 2020, pp. 605–614.
- [44] Y.-N. Chuang et al., "TPR: Text-aware preference ranking for recommender systems," in *Proc. Conf. Inf. Knowl. Manage.*, 2020, pp. 215–224.
- [45] Z. Xu et al., "E-commerce recommendation with weighted expected utility," in *Proc. Conf. Inf. Knowl. Manage.*, 2020, pp. 1695–1704.
- [46] Y. Wang et al., "DisenHAN: Disentangled heterogeneous graph attention network for recommendation," in *Proc. Conf. Inf. Knowl. Manage.*, 2020, pp. 1605–1614.
- [47] Y. Xian et al., "CAFE: Coarse-to-fine knowledge graph reasoning for e-commerce recommendation," in *Proc. Conf. Inf. Knowl. Manage.*, 2020, pp. 1645–1654.
- [48] F. Yuan et al., "Exploring missing interactions: A convolutional generative adversarial network for collaborative filtering," in *Proc. Conf. Inf. Knowl. Manage.*, 2020, pp. 1773–1782.
- [49] F. Zhao and Y. Guo, "Learning discriminative recommendation systems with side information," in *Proc. Int. Joint Conf. Artif. Intell.*, 2017, pp. 3469–3475.
- [50] Z. Sun et al., "MRLR: Multi-level representation learning for personalized ranking in recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2017, pp. 2807–2813.
- [51] H.-J. Xue et al., "Deep matrix factorization models for recommender systems," in *Proc. Int. Joint Conf. Artif. Intell.*, 2017, pp. 3203–3209.
- [52] Y. Liu et al., "Dynamic Bayesian logistic matrix factorization for recommendation with implicit feedback," in *Proc. Int. Joint Conf. Artif. Intell.*, 2018, pp. 3463–3469.
- [53] Z. Wang et al., "Matrix completion with preference ranking for Top-N recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2018, pp. 3585–3591.
- [54] W. Zhao et al., "PLASTIC: Prioritize long and short-term information in top-n recommendation using adversarial training," in *Proc. Int. Joint Conf. Artif. Intell.*, 2018, pp. 3676–3682.
- [55] J. Ding et al., "Improving implicit recommender systems with view data," in *Proc. Int. Joint Conf. Artif. Intell.*, 2018, pp. 3343–3349.
- [56] W. Cheng et al., "DELTA: A dual-embedding based deep latent factor model for recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2018, pp. 3329–3335.
- [57] H. Liu et al., "Discrete factorization machines for fast feature-based recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2018, pp. 3449–3455.
- [58] X. Xin et al., "CFM: Convolutional factorization machines for context-aware recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2019, pp. 3926–3932.
- [59] J. Jiang et al., "Convolutional Gaussian embeddings for personalized recommendation with uncertainty," in *Proc. Int. Joint Conf. Artif. Intell.*, 2019, pp. 2642–2648.
- [60] G. Guo et al., "Discrete trust-aware matrix factorization for fast recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2019, pp. 1380–1386.
- [61] Y. Xu et al., "Learning shared vertex representation in heterogeneous graphs with convolutional networks for recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2019, pp. 4620–4626.
- [62] S. Zhang et al., "Quaternion collaborative filtering for recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2019, pp. 4313–4319.
- [63] W. Fan et al., "Deep adversarial social recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2019, pp. 1351–1357.
- [64] Z. Wang et al., "Unified embedding model over heterogeneous information network for personalized recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2019, pp. 3813–3819.
- [65] P. Han et al., "Contextualized point-of-interest recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2020, pp. 2484–2490.
- [66] R. Xie et al., "Internal and contextual attention network for cold-start multi-channel matching in recommendation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2020, pp. 2732–2738.
- [67] H. Chen and J. Li, "Neural tensor model for learning multi-aspect factors in recommender systems," in *Proc. Int. Joint Conf. Artif. Intell.*, 2020, pp. 2449–2455.
- [68] R. Liu et al., "Hypernews: Simultaneous news recommendation and active-time prediction via a double-task deep neural network," in *Proc. Int. Joint Conf. Artif. Intell.*, 2020, pp. 3487–3493.
- [69] X. Li and J. She, "Collaborative variational autoencoder for recommender systems," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2017, pp. 305–314.
- [70] H. Zhu et al., "Learning tree-based deep model for recommender systems," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2018, pp. 1079–1088.
- [71] E. Christakopoulou and G. Karypis, "Local latent space models for Top-N recommendation," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2018, pp. 1235–1243.
- [72] B. Hu et al., "Leveraging meta-path based context for Top-N recommendation with a neural co-attention model," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2018, pp. 1531–1540.
- [73] X. Wang et al., "KGAT: Knowledge graph attention network for recommendation," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2019, pp. 950–958.
- [74] X. Tang et al., "AKUPM: Attention-enhanced knowledge-aware user preference model for recommendation," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2019, pp. 1891–1899.
- [75] J. Zhao et al., "IntentGC: A scalable graph convolution framework fusing heterogeneous information for recommendation," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2019, pp. 2347–2357.
- [76] H. Wang et al., "Knowledge-aware graph neural networks with label smoothness regularization for recommender systems," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2019, pp. 968–977.
- [77] Y. Chen et al., "LambdaOpt: Learn to regularize recommender models in finer levels," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2019, pp. 978–986.

- [78] J. Jin et al., "An efficient neighborhood-based interaction model for recommendation on heterogeneous graph," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2020, pp. 75–84.
- [79] C. Ma et al., "Probabilistic metric learning with adaptive margin for top-K recommendation," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2020, pp. 1036–1044.
- [80] S. Ji et al., "Dual channel hypergraph collaborative filtering," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2020, pp. 2020–2029.
- [81] R. Otunba et al., "MPR: Multi-objective pairwise ranking," in *Proc. ACM Conf. Recommender Syst.*, 2017, pp. 170–178.
- [82] D. Rafailidis and F. Crestani, "Learning to rank with trust and distrust in recommender systems," in *Proc. ACM Conf. Recommender Syst.*, 2017, pp. 5–13.
- [83] D. Valcarce et al., "On the robustness and discriminative power of information retrieval metrics for top-N recommendation," in *Proc. ACM Conf. Recommender Syst.*, 2018, pp. 260–268.
- [84] Z. Sun et al., "Recurrent knowledge graph embedding for effective recommendation," in *Proc. ACM Conf. Recommender Syst.*, 2018, pp. 297–305.
- [85] L. Zheng et al., "Spectral collaborative filtering," in *Proc. ACM Conf. Recommender Syst.*, 2018, pp. 311–319.
- [86] S. Ouyang et al., "Asymmetric Bayesian personalized ranking for one-class collaborative filtering," in *Proc. ACM Conf. Recommender Syst.*, 2019, pp. 373–377.
- [87] H. Liu et al., "Deep generative ranking for personalized recommendation," in *Proc. ACM Conf. Recommender Syst.*, 2019, pp. 34–42.
- [88] A. N. Nikolakopoulos et al., "Personalized diffusions for top-n recommendation," in *Proc. ACM Conf. Recommender Syst.*, 2019, pp. 260–268.
- [89] F. S. d. Costa and P. Dolog, "Collective embedding for neural context-aware recommender systems," in *Proc. ACM Conf. Recommender Syst.*, 2019, pp. 201–209.
- [90] E. Frolov and I. Oseledets, "HybridSVD: When collaborative information is not enough," in *Proc. ACM Conf. Recommender Syst.*, 2019, pp. 331–339.
- [91] E. Elahi et al., "Variational low rank multinomials for collaborative filtering with side-information," in *Proc. ACM Conf. Recommender Syst.*, 2019, pp. 340–347.
- [92] Y. Zhang et al., "Content-collaborative disentanglement representation learning for enhanced recommendation," in *Proc. ACM Conf. Recommender Syst.*, 2020, pp. 43–52.
- [93] D. Liu et al., "KRED: Knowledge-aware document representation for news recommendations," in *Proc. ACM Conf. Recommender Syst.*, 2020, pp. 200–209.
- [94] J. P. Zhou et al., "TAFA: Two-headed attention fused autoencoder for context-aware recommendations," in *Proc. ACM Conf. Recommender Syst.*, 2020, pp. 338–347.
- [95] J. Chen et al., "Attentive collaborative filtering: Multimedia recommendation with item-and component-level attention," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2017, pp. 335–344.
- [96] X. He et al., "Adversarial personalized ranking for recommendation," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2018, pp. 355–364.
- [97] T. Ebesu et al., "Collaborative memory network for recommendation systems," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2018, pp. 515–524.
- [98] Q. Xu et al., "GraphCAR: Content-aware multimedia recommendation with graph autoencoder," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2018, pp. 981–984.
- [99] R. Canameres and P. Castells, "Should I follow the crowd?: A probabilistic analysis of the effectiveness of popularity in recommender systems," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2018, pp. 415–424.
- [100] W. Wang et al., "Streaming ranking based recommender systems," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2018, pp. 525–534.
- [101] Y. Chen et al., "Bayesian personalized feature interaction selection for factorization machines," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2019, pp. 665–674.
- [102] X. Wang et al., "Neural graph collaborative filtering," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2019, pp. 165–174.
- [103] G. Wu et al., "Noise contrastive estimation for one-class collaborative filtering," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2019, pp. 135–144.
- [104] X. Xin et al., "Relational collaborative filtering: Modeling multiple item relations for recommendation," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2019, pp. 125–134.
- [105] Z. Wang et al., "CKAN: Collaborative knowledge-aware attentive network for recommender systems," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 219–228.
- [106] C. Chen et al., "Jointly non-sampling learning for knowledge graph enhanced recommendation," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 189–198.
- [107] J. Gong et al., "Attentional graph convolutional networks for knowledge concept recommendation in MOOCs in a heterogeneous view," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 79–88.
- [108] C. Hansen et al., "Content-aware neural hashing for cold-start recommendation," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 971–980.
- [109] X. He et al., "LightGCN: Simplifying and powering graph convolution network for recommendation," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 639–648.
- [110] S. Shi et al., "Beyond user embedding matrix: Learning to hash for modeling large-scale users in recommendation," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 319–328.
- [111] C.-Y. Tai et al., "MVIN: Learning multiview items for recommendation," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 99–108.
- [112] L. Wu et al., "Joint item recommendation and attribute inference: An adaptive graph convolutional network approach," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 679–688.
- [113] D.-K. Chae et al., "AR-CF: Augmenting virtual users and items in collaborative filtering for addressing cold-start problems," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 1251–1260.
- [114] X. Wang et al., "Disentangled graph collaborative filtering," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 1001–1010.
- [115] L. Zou et al., "Neural interactive collaborative filtering," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 749–758.
- [116] J. Sun et al., "Neighbor interaction aware graph convolution networks for recommendation," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2020, pp. 1289–1298.
- [117] Q. Zhao et al., "Multi-product utility maximization for economic recommendation," in *Proc. Int. Conf. Web Search Data Mining*, 2017, pp. 435–443.
- [118] Y. Zhang et al., "Discrete deep learning for fast content-aware recommendation," in *Proc. Int. Conf. Web Search Data Mining*, 2018, pp. 717–726.
- [119] Z. Jiang et al., "Recommendation in heterogeneous information networks based on generalized random walk model and Bayesian personalized ranking," in *Proc. Int. Conf. Web Search Data Mining*, 2018, pp. 288–296.
- [120] W. Niu et al., "Neural personalized ranking for image recommendation," in *Proc. Int. Conf. Web Search Data Mining*, 2018, pp. 423–431.
- [121] C. Ma et al., "Gated attentive-autoencoder for content-aware recommendation," in *Proc. Int. Conf. Web Search Data Mining*, 2019, pp. 519–527.
- [122] A. N. Nikolakopoulos and G. Karypis, "RecWalk: Nearly uncoupled random walks for top-n recommendation," in *Proc. Int. Conf. Web Search Data Mining*, 2019, pp. 150–158.
- [123] C. Chen et al., "Social attentional memory network: Modeling aspect- and friend-level differences in recommendation," in *Proc. Int. Conf. Web Search Data Mining*, 2019, pp. 177–185.
- [124] D. Liu et al., "Spiral of silence in recommender systems," in *Proc. Int. Conf. Web Search Data Mining*, 2019, pp. 222–230.
- [125] H. Steck et al., "ADMM SLIM: Sparse recommendations for many users," in *Proc. Int. Conf. Web Search Data Mining*, 2020, pp. 555–563.
- [126] R. Li et al., "Adversarial learning to compare: Self-attentive prospective customer recommendation in location based social networks," in *Proc. Int. Conf. Web Search Data Mining*, 2020, pp. 349–357.
- [127] F. Liu et al., "End-to-end deep reinforcement learning based recommendation with supervised embedding," in *Proc. Int. Conf. Web Search Data Mining*, 2020, pp. 384–392.
- [128] Y. Gu et al., "Hierarchical user profiling for e-commerce recommender systems," in *Proc. Int. Conf. Web Search Data Mining*, 2020, pp. 223–231.
- [129] J. Wang et al., "Key opinion leaders in recommendation systems: Opinion elicitation and diffusion," in *Proc. Int. Conf. Web Search Data Mining*, 2020, pp. 636–644.
- [130] C. Sun et al., "LARA: Attribute-to-feature adversarial learning for new-item recommendation," in *Proc. Int. Conf. Web Search Data Mining*, 2020, pp. 582–590.
- [131] H. Zamani and W. B. Croft, "Learning a joint search and recommendation model from user-item interactions," in *Proc. Int. Conf. Web Search Data Mining*, 2020, pp. 717–725.

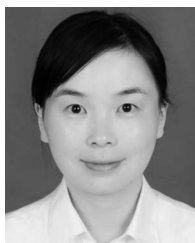
- [132] I. Shenbin et al., "RecVAE: A new variational autoencoder for top-n recommendations with implicit feedback," in *Proc. Int. Conf. Web Search Data Mining*, 2020, pp. 528–536.
- [133] C.-K. Hsieh et al., "Collaborative metric learning," in *Proc. World Wide Web Conf.*, 2017, pp. 193–201.
- [134] X. He and T.-S. Chua, "Neural factorization machines for sparse predictive analytics," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2017, pp. 355–364.
- [135] W. Yu et al., "Aesthetic-based clothing recommendation," in *Proc. World Wide Web Conf.*, 2018, pp. 649–658.
- [136] D. Liang et al., "Variational autoencoders for collaborative filtering," in *Proc. World Wide Web Conf.*, 2018, pp. 689–698.
- [137] H. Wang et al., "Multi-task feature learning for knowledge graph enhanced recommendation," in *Proc. World Wide Web Conf.*, 2019, pp. 2000–2010.
- [138] W. Ma et al., "Jointly learning explainable rules for recommendation with knowledge graph," in *Proc. World Wide Web Conf.*, 2019, pp. 1210–1221.
- [139] Y. Cao et al., "Unifying knowledge graph learning and recommendation: Towards a better understanding of user preferences," in *Proc. World Wide Web Conf.*, 2019, pp. 151–161.
- [140] T. Tran et al., "Signed distance-based deep memory recommender," in *Proc. World Wide Web Conf.*, 2019, pp. 1841–1852.
- [141] H. Wang et al., "Knowledge graph convolutional networks for recommender systems," in *Proc. World Wide Web Conf.*, 2019, pp. 3307–3313.
- [142] C.-M. Chen et al., "Collaborative similarity embedding for recommender systems," in *Proc. World Wide Web Conf.*, 2019, pp. 2637–2643.
- [143] F. Khawar et al., "Learning the structure of auto-encoding recommenders," in *Proc. World Wide Web Conf.*, 2020, pp. 519–529.
- [144] H. Liu et al., "Deep global and local generative model for recommendation," in *Proc. World Wide Web Conf.*, 2020, pp. 551–561.
- [145] A. Javari et al., "Weakly supervised attention for hashtag recommendation using graph data," in *Proc. World Wide Web Conf.*, 2020, pp. 1038–1048.
- [146] C. Wang et al., "Personalized employee training course recommendation with career development awareness," in *Proc. World Wide Web Conf.*, 2020, pp. 1648–1659.
- [147] Q. Tan et al., "Learning to hash with graph neural networks for recommender systems," in *Proc. World Wide Web Conf.*, 2020, pp. 1988–1998.
- [148] C. Chen et al., "Efficient non-sampling factorization machines for optimal context-aware recommendation," in *Proc. World Wide Web Conf.*, 2020, pp. 2400–2410.
- [149] Y. Koren et al., "Matrix factorization techniques for recommender systems," *Computer*, vol. 42, no. 8, pp. 30–37, 2009.
- [150] S. Rendle et al., "BPR: Bayesian personalized ranking from implicit feedback," in *Proc. Conf. Uncertainty Artif. Intell.*, 2009, pp. 452–461.
- [151] B. Hidasi et al., "Parallel recurrent neural network architectures for feature-rich session-based recommendations," in *Proc. ACM Conf. Recommender Syst.*, 2016, pp. 241–248.
- [152] G. Guo et al., "LibRec: A Java library for recommender system," in *Proc. 23rd Conf. User Modelling Adapt. Personalization Posters Demos Late-Breaking Results Workshop*, 2015, pp. 38–45.
- [153] S. Zhang et al., "DeepRec: An open-source toolkit for deep learning based recommendation," in 2019, *arXiv:1905.10536*.
- [154] W. X. Zhao et al., "RecBole: Towards a unified, comprehensive and efficient framework for recommendation algorithms," in *Proc. Conf. Inf. Knowl. Manage.*, 2021, pp. 4653–4664.
- [155] P. Cremonesi et al., "Performance of recommender algorithms on top-n recommendation tasks," in *Proc. ACM Conf. Recommender Syst.*, 2010, pp. 39–46.
- [156] O. S. Collaboration et al., "Estimating the reproducibility of psychological science," *Science*, vol. 349, no. 6251, 2015, Art. no. aac4716.
- [157] M. Baker, "Reproducibility crisis," *Nature*, vol. 533, no. 26, pp. 353–66, 2016.
- [158] M. R. Munafò et al., "A manifesto for reproducible science," *Nat. Hum. Behav.*, vol. 1, no. 1, pp. 1–9, 2017.
- [159] N. Ferro and D. Kelly, "SIGIR initiative to implement ACM artifact review and badging," vol. 52, no. 1, pp. 4–10, 2018.
- [160] J. Freire et al., "Report from Dagstuhl seminar 16041: Reproducibility of data-oriented experiments in e-science," *Dagstuhl Rep.*, vol. 6, no. 1, pp. 108–159, 2016.
- [161] R. Clancy et al., "Overview of the 2019 open-source IR replicability challenge (OSIRRC 2019)," in *Proc. 42nd Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2019, pp. 1–7.
- [162] J. Pineau et al., "Improving reproducibility in machine learning research: A report from the NeurIPS 2019 reproducibility program," *J. Mach. Learn. Res.*, vol. 22, pp. 7459–7478, 2021.
- [163] A. Hanbury et al., *Advances in Information Retrieval: 37th European Conference on IR Research, ECIR 2015, Vienna, Austria, March 29–April 2, 2015. Proceedings*, vol. 9022. Berlin, Germany: Springer, 2015.
- [164] *Proc. 29th ACM Int. Conf. Multimedia*, New York, NY, USA, 2021.
- [165] A. Hotho et al., *The Semantic Web—ISWC 2021: 20th International Semantic Web Conference, ISWC 2021, Virtual Event, October 24–28, 2021, Proceedings*, vol. 12922. Berlin, Germany: Springer, 2021.
- [166] *Proc. 14th ACM Conf. Recommender Syst.*, New York, NY, USA, 2020.
- [167] V. W. Anelli et al., "Elliot: A comprehensive and rigorous framework for reproducible recommender systems evaluation," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2021, pp. 2405–2414.
- [168] J. Beel et al., "Towards reproducibility in recommender-systems research," *User Model. User-Adapted Interact.*, vol. 26, no. 1, pp. 69–101, 2016.
- [169] N. Sachdeva, C.-J. Wu, and J. McAuley, "On sampling collaborative filtering datasets," in *Proc. Int. Conf. Web Search Data Mining*, 2022, pp. 842–850.
- [170] S. M. McNee et al., "Being accurate is not enough: How accuracy metrics have hurt recommender systems," in *Proc. Conf. Extended Abstr. Hum. Factors Comput. Syst.*, 2006, pp. 1097–1101.
- [171] Z. Gantner, S. Rendle, C. Freudenthaler, and L. Schmidt-Thieme, "MyMediaLite: A free recommender system library," in *Proc. ACM Conf. Recommender Syst.*, 2011, pp. 305–308.
- [172] N. Hug, "Surprise: A Python library for recommender systems," *J. Open Source Softw.*, vol. 5, no. 52, 2020, Art. no. 2174.
- [173] I. Cantador et al., "The 2nd workshop on information heterogeneity and fusion in recommender systems (HetRec)," in *Proc. ACM Conf. Recommender Syst.*, 2011, pp. 387–388.
- [174] J. Tang et al., "eTrust: Understanding trust evolution in an online world," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2012, pp. 253–261.
- [175] C.-N. Ziegler et al., "Improving recommendation lists through topic diversification," in *Proc. World Wide Web Conf.*, 2005, pp. 22–32.
- [176] S. Rendle, "Factorization machines," in *Proc. Int. Conf. Des. Minings*, 2010, pp. 995–1000.
- [177] Y. Hu et al., "Collaborative filtering for implicit feedback datasets," in *Proc. Int. Conf. Des. Minings*, 2008, pp. 263–272.
- [178] X. Ning and G. Karypis, "SLIM: Sparse linear methods for top-N recommender systems," in *Proc. Int. Conf. Des. Minings*, 2011, pp. 497–506.
- [179] X. He et al., "Fast matrix factorization for online recommendation with implicit feedback," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2016, pp. 549–558.
- [180] H. Zhang et al., "Discrete collaborative filtering," in *Proc. Int. ACM SIGIR Conf. Res. Develop. Informat. Retrieval*, 2016, pp. 325–334.
- [181] F. Zhang et al., "Collaborative knowledge base embedding for recommender systems," in *Proc. Int. Conf. Knowl. Discov. Data Mining*, 2016, pp. 353–362.
- [182] Y. Wu et al., "Collaborative denoising auto-encoders for top-N recommender systems," in *Proc. Int. Conf. Web Search Data Mining*, 2016, pp. 153–162.
- [183] Q. Zhao et al., "Interpreting user inaction in recommender systems," in *Proc. ACM Conf. Recommender Syst.*, 2018, pp. 40–48.
- [184] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2010, pp. 249–256.
- [185] N. Srivastava et al., "Dropout: A simple way to prevent neural networks from overfitting," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [186] J. Bergstra and Y. Bengio, "Random search for hyper-parameter optimization," *J. Mach. Learn. Res.*, vol. 13, no. 1, pp. 281–305, 2012.
- [187] J. Snoek et al., "Practical Bayesian optimization of machine learning algorithms," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2012, pp. 2960–2968.
- [188] L. Yang et al., "OpenRec: A modular framework for extensible and adaptable recommendation algorithms," in *Proc. Int. Conf. Web Search Data Mining*, 2018, pp. 664–672.
- [189] D. Jannach and M. Ludewig, "When recurrent neural networks meet the neighborhood for session-based recommendation," in *Proc. ACM Conf. Recommender Syst.*, 2017, pp. 306–310.
- [190] E. Raff, "A step toward quantifying independently reproducible machine learning research," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2019, pp. 5485–5495.
- [191] B. McFee et al., "The million song dataset challenge," in *Proc. World Wide Web Conf.*, 2012, pp. 909–916.



Zhu Sun received the PhD degree from Nanyang Technological University, Singapore, in 2018. Her main research topic is recommender systems. Her research has been published in leading conferences and journals (e.g., IJCAI, AAAI, SIGIR, CIKM, RecSys, *IEEE Transactions on Knowledge and Data Engineering* and *IEEE Transactions on Neural Networks and Learning Systems*). She is the AE with ECRA journal, and PC Member for KDD, SIGIR, IJCAI, AAAI, CIKM, RecSys, IUI and UMAP, etc.



Hongyang Liu is currently working toward the master's degree with the School of Information Science and Engineering, Yanshan University, China. His research topic includes machine learning and recommender systems. He mainly focused on applying deep learning techniques (e.g., reinforcement learning and generative adversarial networks) to improve the performance of recommender systems in the related areas.



Hui Fang received the PhD degree from Nanyang Technological University, Singapore. She is an associate professor with the Shanghai University of Finance and Economics, China. Her main research topic is personalized machine learning, including trust/link prediction in online communities, and recommender systems. She has published papers in leading conferences (e.g., IJCAI, AAAI and AAMAS), and journals (e.g., *Journal of Artificial Intelligence*, *IEEE Transactions on Knowledge and Data Engineering*, *ACM Transactions on Information Systems* and *Decision Support Systems*). She is the SE of the ECRA journal, and serves as a PC Board of IJCAI, and (Senior) PC Member for UMAP, IJCAI, AAAI and AAMAS, etc.



Di Yu received the bachelor's degree in computer science and master's degree in management science from Shanghai University of Finance and Economics, China in 2017 and in 2019, respectively. He is currently working toward the master's degree with Singapore Management University. His research interests include personalized recommendation and FinTech.



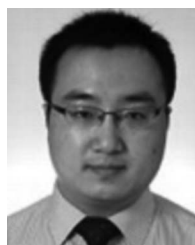
Jie Yang is assistant professor with the Web Information Systems group of Delft University of Technology (TU Delft). Before joining TU Delft, he worked as a machine learning scientist with Amazon and a senior researcher with the eXascale Infolab, University of Fribourg. His research focuses on human-centered AI for Web-scale information systems, aiming at leveraging the joint power of human and machine intelligence for understanding and making use of data in large-scale information systems.



Yew-Soon Ong (Fellow, IEEE) received the PhD degree from the University of Southampton, U.K., in 2003. He is president's chair professor in Computer Science with Nanyang Technological University (NTU), and holds the position of chief artificial intelligence scientist of A*STAR, Singapore. At NTU, he serves as co-director of the Singtel-NTU Cognitive & Artificial Intelligence Joint Lab. His research interest is in artificial and computational intelligence. He is founding EIC of *IEEE Transactions on Emerging Topics in Computational Intelligence* and AE of *IEEE Transactions on Neural Networks and Learning Systems*, *IEEE Transactions on Cybernetics*, *IEEE Transactions on Artificial Intelligence* and others. He has received several IEEE outstanding paper awards and was listed as a Thomson Reuters highly cited Researcher and among the World's Most Influential Scientific Minds.



Xinghua Qu received the PhD degree from Nanyang Technological University, Singapore in 2022. Currently he is a research scientist with Bytedance AI Lab, Singapore. His primary research interest includes machine learning and optimisation. Recently his research interests mainly focus on the adversarial robustness of diverse machine learning areas, including deep reinforcement learning, computer vision, and speech recognition.



Jie Zhang received the PhD degree from the University of Waterloo, Canada, in 2009. He is currently an associate professor with the School of Computer Science and Engineering, Nanyang Technological University and Singapore Institute of Manufacturing Technology, Singapore. He was a recipient of the Alumni Gold Medal at the 2009 Convocation Ceremony, which is awarded once a year to honor the top PhD graduate from the University of Waterloo. During his PhD study, he held the prestigious NSERC Alexander Graham Bell Canada Graduate Scholarship rewarded for top PhD students across Canada. His papers have been published by top journals and conferences and received several best paper awards. He is also active in serving research communities.