



# Lightweight Unbiased Multi-teacher Ensemble for Review-based Recommendation

Guipeng Xv<sup>\*†</sup>

School of Informatics, Xiamen University  
Xiamen, China  
xvguipeng.xgp@alibaba-inc.com

Xinyi Liu<sup>\*</sup>

School of Informatics, Xiamen University  
Xiamen, China  
xinyiliu@stu.xmu.edu.cn

Chen Lin<sup>‡</sup>

School of Informatics, Xiamen University  
Xiamen, China  
chenlin@xmu.edu.cn

Hui Li

School of Informatics, Xiamen University  
Xiamen, China  
hui@xmu.edu.cn

Chenliang Li

School of Cyber Science and Engineering, Wuhan University  
Wuhan, China  
cllee@whu.edu.cn

Zhenhua Huang

School of Computer Science, South China Normal University  
Guangzhou, China  
huangzhenhua@m.scnu.edu.cn

## ABSTRACT

Review-based recommender systems (RRS) have received an increasing interest since reviews greatly enhance recommendation quality and interpretability. However, existing RRS suffer from high computational complexity, biased recommendation and poor generalization. The three problems make them inadequate to handle real recommendation scenarios. Previous studies address each issue separately, while none of them consider solving three problems together under a unified framework. This paper presents LUME (a Lightweight Unbiased Multi-teacher Ensemble) for RRS. LUME is a novel framework that addresses the three problems simultaneously. LUME uses multi-teacher ensemble and debiased knowledge distillation to aggregate knowledge from multiple pretrained RRS, and generates a small, unbiased student recommender which generalizes better. Extensive experiments on various real-world benchmarks demonstrate that LUME successfully tackles the three problems and has superior performance than state-of-the-art RRS and knowledge distillation based RS.

## CCS CONCEPTS

• Information systems → Recommender systems.

## KEYWORDS

review-based recommender systems, bias in recommender systems, knowledge distillation

<sup>\*</sup>Both authors contributed equally to this research.

<sup>†</sup>Work done during internship at Alibaba Group.

<sup>‡</sup>Corresponding author. Supported by the Natural Science Foundation of China (No. 61972328), Alibaba Group through Alibaba Innovative Research program.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

CIKM '22, October 17–21, 2022, Atlanta, GA, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9236-5/22/10...\$15.00

<https://doi.org/10.1145/3511808.3557629>

## ACM Reference Format:

Guipeng Xv, Xinyi Liu, Chen Lin, Hui Li, Chenliang Li, and Zhenhua Huang. 2022. Lightweight Unbiased Multi-teacher Ensemble for Review-based Recommendation. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management (CIKM '22)*, October 17–21, 2022, Atlanta, GA, USA. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3511808.3557629>

## 1 INTRODUCTION

Online reviews are valuable feedback in recommender systems, as they provide explanations on various aspects of a product and guide users towards purchase. Due to the prevalence of online reviewing sites, Review-based Recommender Systems (RRS) have attracted a great amount of attention [1, 2, 11, 17, 25, 31, 35].

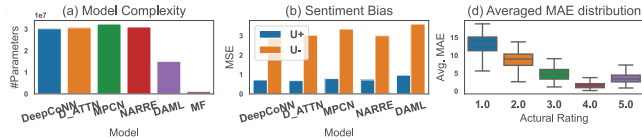
Although existing RRS provide high-quality and interpretable recommendations, they still suffer from issues such as **high computational complexity, biased recommendations, and poor generalization**. (1) As shown in Fig. 1 (a), since state-of-the-art RRS typically adopt deep neural networks to analyze review contents, they have much more parameters than the non-review-based, shallow recommendation models such as Matrix Factorization (MF) [12]. (2) As shown in Fig. 1(b), RRS are recently found to show sentiment bias [16], i.e., they generate more significant errors on critical users who write fewer positive reviews than on positive users who post more positive reviews. (3) As shown in Fig. 1(c), RRS make inaccurate (i.e., large median MAE) and unreliable (i.e., large variance) predictions on low-rating reviews, which are with insufficient training samples but more valuable than high-rating reviews [22]. These issues severely affect the efficiency in terms of inference time and storage cost, and the quality and fairness of recommendations, when RRS are deployed in practice.

In the literature, model compression [5, 14, 24, 28], debiasing [3, 30] and generalization [15] have been studied for RS. However, existing works address these problems separately. Furthermore, the three issues are correlated, e.g., model complexity and generalization ability, generalization ability to low-rating reviews and bias against critical users. Lacking consideration of any of the above problems will result in sub-optimal, ineffective and/or inefficient RRS that are inadequate to handle real recommendation scenarios.

Inspired by recent advances on Knowledge Distillation (KD) [8], we propose Lightweight Unbiased Multi-teacher Ensemble (LUME)

for the review-based recommendation task. LUME first captures high quality, generalizable common knowledge shared within multiple teacher RRS, and then trains a lightweight student model and mitigates the biases via a KD process. Experiments on various real-world RS datasets verify the superiority of LUME to make consistent, high-quality and unbiased review-based recommendations.

In summary, our contributions are three-fold. (1) We design a novel framework, LUME, which simultaneously alleviates high computational complexity, bias and poor generalization of RRS. (2) Unlike most existing KD-based RS that only learn from one teacher model, LUME compresses and accelerates multiple teachers by fusing common knowledge and adapting it to the student model. (3) The KD process of LUME is specially designed to handle biases



**Figure 1: (a) Number of parameters in five state-of-the-art RRS (i.e., DeepCoNN [35], MPCN [25], NARRE [2], DAML [17] and D\_ATT [11]) and conventional MF (with embedding/factor size = 50). (b) RRS produce much higher Mean Square Error (MSE) on critical users  $\mathcal{U}^-$  than on positive users  $\mathcal{U}^+$ .  $\mathcal{U}^-$ ,  $\mathcal{U}^+$  are decided based on the sentiment scores of reviews as in [16]. (c) Averaged MAE (i.e., MAE averaged over different RRS on every sample) on one-star rating samples has a large variance and a large median value.**

## 2 RELATED WORK

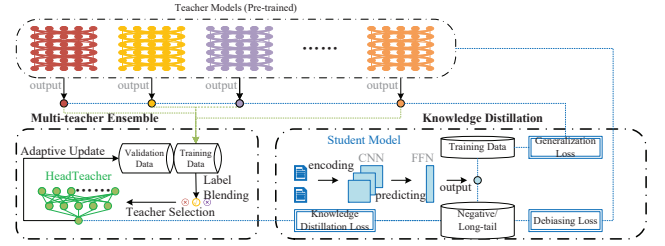
**Review-based Recommender Systems (RRS).** Traditional RRS have utilized latent semantic analysis [34], LDA [27] or latent factor model [20] to model reviews and provide better recommendations. Recently, deep learning based techniques such as CNN [35]MLP [11], LSTM [23, 26], Autoencoder [36], and the attention mechanism [2, 17, 25] have significantly facilitated the development of RRS [31].

**Biases in RS.** Several biases have been observed in RS [4], including selection bias [19], conformity bias [18], position bias [9], popularity bias [32], and exposure bias [29]. The ubiquitous sentiment bias problem [16] in RRS is hard to handle. To mitigate biases in RS, numerous debiasing methods [4] adopt causal inference methods [29], or regularization methods [3, 16].

**Knowledge Distillation (KD) for RS.** A number of recent studies have investigated KD [8, 10] in RS, where a small student model learns to rank items [14, 24] from one large teacher models to reduce model complexity. The ranking distillation framework can be enhanced by a three-player game where a discriminator is introduced to learn the true data distribution from the teacher [28]. The student and the teacher can learn from each other simultaneously [13, 33] to enhance the interpretability or the performance of RS. For RRS, an adversarial distillation framework [5] is proposed to make review predictions.

## 3 OUR METHOD: LUME

As shown in Fig. 2, LUME mainly consists of two parts. Given a set of pre-trained RRS (i.e., teachers), LUME first learns a *HeadTeacher* model to fuse the knowledge from multiple teachers and further



**Figure 2: Overview of LUME.**

improves the quality of common knowledge, using three steps: label blending, teacher selection, and adaptive model update (Sec. 3.1). Then, LUME trains a student model using the guidance from the HeadTeacher, mitigates biases, and strengthens generalization in a KD process (Sec. 3.2).

### 3.1 Multi-teacher Ensemble

Different from existing KD-based RS [5, 14, 24, 28] that leverage a single-teacher architecture, LUME uses a multi-teacher architecture, so that the student will not be easily misled by a single teacher if the teacher performs poorly in some cases. A natural approach to fuse multiple teachers is to use an ensemble model  $M^e$ , parameterized by  $\Theta^e$  to integrate the predictions of multiple teachers, i.e., ensemble learning [6]. However, as illustrated in Sec. 1: different RRS produce inconsistent predictions on hard cases, which introduces noise to the KD process. It is questionable whether abnormal predictions from some teachers should be incorporated into the ensemble model. To overcome the above problem, we propose the *multi-teacher ensemble* to generate the HeadTeacher.

Suppose that we have a number of teacher models, where each teacher model  $t \in \mathcal{T}$  gives the prediction  $\hat{X}_{u,i}^t$  for the rating  $X_{u,i}$  of a user  $u \in \mathcal{U}$  on an item  $i \in \mathcal{I}$ . The teacher models are independently pre-trained on the training set  $\mathcal{DS}$ , and they are fixed during the training phase of the later KD process.

We first use a **label blending** step which traverses the training set  $\mathcal{DS}$  and removes low-quality teacher predictions to fuse outputs from multiple teachers. A label  $l(t, u, i)$  is assigned for each teacher  $t$  on every prediction  $\hat{X}_{u,i}^t$  to indicate whether the prediction should be utilized in training the HeadTeacher. If the deviation between the prediction and the actual rating, i.e.,  $|\hat{X}_{u,i}^t - X_{u,i}|$ , is larger than a predefined threshold  $\zeta$ ,  $\hat{X}_{u,i}^t$  will be considered as abnormal and it will not benefit the ensemble learning.

**Teacher selection.** Then, the HeadTeacher takes the output of each teacher model  $\hat{X}_{u,i}^t$ , if  $l(t, u, i)$  equals 1, and make a fused prediction. The HeadTeacher uses a two-layer feed-forward network (FFN). In the first layer, predictions from teachers are aggregated to generate the probabilities of different rating values:

$$\mathbf{z}_{u,i} = \mathbf{w}_1 \left( \text{concate}(l(t, u, i)\hat{X}_{u,i}^t, t \in \mathcal{T}) \right) + \mathbf{b}_1, \quad (1)$$

where  $\text{concate}(\cdot) \in \mathcal{R}^{|\mathcal{T}| \times 1}$  is a concatenated vector,  $\mathbf{z}_{u,i} \in \mathcal{R}^{5 \times 1}$  indicates the probability distribution of ratings.  $\mathbf{w}_1 \in \mathcal{R}^{5 \times |\mathcal{T}|}$  and  $\mathbf{b}_1 \in \mathcal{R}^{5 \times 1}$  are learnable weight vector and bias vector, respectively. In the second layer, different rating values are aggregated to form the predicted rating:

$$\hat{X}_{u,i}^e = \mathbf{w}_2^T \mathbf{z}_{u,i} + b_2, \quad (2)$$

where  $\mathbf{w}_2 \in \mathbb{R}^{5 \times 1}$  and  $b_2$  are learnable parameters.

**Adaptive model update.** Besides, we use a subset of the testing data as a validation set  $\mathcal{DV}$  to improve the generalization of LUME. We derive the gradient of the HeadTeacher parameters in the validation set and carry a small number of trials to update the ensemble model. The motivation is similar to model-agnostic meta-learning [7]: Since RRS will be updated using a gradient-based method on new data (including low-rating reviews) that they can not learn well (i.e., poor generalization), LUME is designed to find model parameters that are sensitive to new data so that small changes in model parameters will produce large improvements on the loss function.

### 3.2 Knowledge Distillation

Formally, the student model is denoted as  $M^s$ , parameterized with  $\Theta^s$ , which makes predictions  $\hat{X}_{u,i}^s = M_{\Theta^s}(\mathbf{P}_u, \mathbf{Q}_i)$  for each user profile  $\mathbf{P}_u$  and item profile  $\mathbf{Q}_i$ . We construct a user profile  $\mathbf{P}_u$  by concatenating all reviews written by user  $u$ . An embedding vector is used to represent each review token, and thus a user profile is defined as  $\mathbf{P}_u \in \mathcal{R}^{N_u \times N_w}$ , where  $N_u$  is the maximal number of reviews that LUME includes in a user profile, and  $N_w$  is the number of the tokens that LUME considers for each review from its beginning. Similarly, we construct an item profile  $\mathbf{Q}_i$  by concatenating all reviews written on item  $i$ .

The architecture of the student model, as most RRS models, consists of an encoding module that learns feature representations of textual reviews and a prediction module. The goal of the student model in LUME is to make it as lightweight as possible. We experimentally find that Convolutional Neural Network (CNN), as an encoding module, generates stable performance. To reduce the computational complexity, we use the same CNN module for both user profiles and item profiles. The prediction module in the student model is a one-layer FFN that predicts the ratings in one to five stars.

The student optimizes a combined loss that helps the student mimic the behavior of the HeadTeacher via a *teacher distillation loss*  $\mathcal{L}_t$ , *student loss*  $\mathcal{L}_s$ , and mitigates biases via a *debiasing loss*  $\mathcal{L}_x$ , and strengthens generalization via a *generalization losses*  $\mathcal{L}_g$ . The overall loss for training the student model is defined as:

$$\mathcal{L} = \lambda_t \mathcal{L}_t + \lambda_s \mathcal{L}_s + \lambda_x \mathcal{L}_x + \lambda_g \mathcal{L}_g, \quad (3)$$

where  $\lambda_t, \lambda_s, \lambda_x, \lambda_g$  are loss weights.

**Teacher distillation loss.** Recall that the HeadTeacher contains two layers, where the output of the first layer (i.e., logits  $\mathbf{z}_{u,i}$  in Eq. 1) carries ensemble knowledge from different teachers, by predicting the probabilities of one to five rating stars, i.e.,  $\mathbf{z}_{u,i,c} = \Pr(\mathbf{X}_{u,i} = c), c \in \{1, 2, 3, 4, 5\}$ . However, the student outputs numerical rating values instead of discrete rating categories. Thus, the cross-entropy loss used in many KD systems [8] is infeasible for RRS.

To transfer the ensemble knowledge in  $M^e$  to  $M^s$ , LUME uses the logits as supervision signals and optimizes the MSE loss between logits and the student model's output as the teacher loss  $\mathcal{L}_t$ :

$$\mathcal{L}_t = \sum_{u \in \mathcal{U}, i \in \mathcal{I}, \mathbf{X}_{u,i} \neq 0} \left( \sum_c c \cdot \mathbf{z}_{u,i,c} - \hat{X}_{u,i}^s \right)^2, \quad (4)$$

where  $c = \{1, 2, 3, 4, 5\}$  refers to discrete ratings in RS,  $\mathbf{z}_{u,i,c}$  is the logit output from the first layer of HeadTeacher on the neuron for  $c$  (Eq. 1).

**Student loss.**  $\mathcal{L}_s$  in LUME is defined between the ground truth rating value  $\mathbf{X}_{u,i}$  and the output of the student to encourage the student to make accurate predictions:

$$\mathcal{L}_s = \sum_{u \in \mathcal{U}, i \in \mathcal{I}, \mathbf{X}_{u,i} \neq 0} (\hat{X}_{u,i}^s - \mathbf{X}_{u,i})^2. \quad (5)$$

**Debiasing loss.** In the following, we use the sentiment bias, which exists in most RRS [16], as the example to illustrate how LUME mitigate biases. The idea can be generalized to other biases (e.g., popularity bias). Intuitively, to reduce sentiment bias, the student model must be enhanced to provide better predictions on negative users/items. We propose  $E_o(\mathcal{S}, t)$  to evaluate teacher model  $t$ , based on how much the embedding vectors of negative items spread out in the batch containing samples  $\mathcal{S}$ :  $E_o(\mathcal{S}, t) = \sum_{\mathbf{X}_{u,i} \in \mathcal{S} \& i \in \mathcal{I}^-} \|\mathbf{e}_i^t - \mathbf{e}^t(\mathcal{S})\|_2^2$ , where  $\mathbf{e}^t(\mathcal{S})$  is the mean embedding vector in the set  $\mathcal{S}$ . When the best teacher model  $x$ , in terms of the smallest  $E_o$  is selected, we can use the output of  $x$  (i.e.,  $\hat{X}_{u,i}^x$ ) to guide the student model and reduce sentiment bias on negative items via the following debiasing loss:

$$\mathcal{L}_x = \sum_{u \in \mathcal{U}, i \in \mathcal{I}^-, \mathbf{X}_{u,i} \neq 0} (\hat{X}_{u,i}^x - \hat{X}_{u,i}^s)^2. \quad (6)$$

**Generalization losses.** The Generalization loss  $\mathcal{L}_g$  is defined as  $\mathcal{L}_g = \lambda_y \mathcal{L}_y + \lambda_z \mathcal{L}_z$ . If teacher models do not agree with each other, we increase the uncertainty of student model's output. We first select samples  $\mathcal{O}$  in the batch (i.e.,  $\mathbf{X}_{u,i} \in \mathcal{S}$ ) using the following evaluation function:  $E_o(u, i) = \sum_{t \in \mathcal{T}} \sum_{\mathbf{X}_{u,i} \in \mathcal{S}} (\hat{X}_{u,i}^t - \mathbf{X}_{u,i})^2$ , where  $\mathbf{X}_{u,i}$  is the average output of all teacher models for the sample  $\mathbf{X}_{u,i}$ . If the variance of teacher model outputs (i.e.,  $E_o(u, i)$ ) is large, LUME uses the entropy-based regularizer  $\mathcal{L}_y$  to increase the uncertainty of the final output:

$$\mathcal{L}_y = \sum_{u \in \mathcal{U}, i \in \mathcal{I}^-, E_o(u, i) > \phi} \sum_{c=1}^5 p(u, i, c) \log p(u, i, c), \quad (7)$$

where  $\phi$  denotes a predefined threshold to judge whether teachers agree or not. Simply connecting a FFN layer with softmax to the prediction layer of the student, we can obtain  $p(u, i, c) = \Pr(\mathbf{X}_{u,i} = c)$ , which denotes the probability that user  $u$  gives item  $i$  a rating of  $c$ , where  $0 \leq p(u, i, c) \leq 1$ ,  $\sum_c p(u, i, c) = 1$ , and  $c \in \{1, 2, 3, 4, 5\}$ .

To further enhance the generalization on low-value ratings, we present the error function  $E_g(\mathcal{S}, t)$ , to evaluate whether a teacher model  $t$  provides unbiased predictions on low ratings in a set of ratings  $\mathcal{S}$ :  $E_g(\mathcal{S}, t) = \sum_{\mathbf{X}_{u,i} \in \mathcal{S} \& \mathbf{X}_{u,i} < 3} (\hat{X}_{u,i}^t - \mathbf{X}_{u,i})^2$ . When the best teacher model  $z$ , in terms of the smallest  $E_g$  is selected, we can use the output of  $z$  (i.e.,  $\hat{X}_{u,i}^z$ ) to strengthen the student model's performance on low-value ratings:

$$\mathcal{L}_z = \sum_{u \in \mathcal{U}, i \in \mathcal{I}, \mathbf{X}_{u,i} \neq 0} (\hat{X}_{u,i}^z - \hat{X}_{u,i}^s)^2. \quad (8)$$

## 4 EXPERIMENTS

Experiments are conducted on four Amazon datasets [21] and Yelp dataset. We apply 5-core pre-processing on Yelp to make sure each user/item has at least five ratings. We use 8:1:1 training/validation/test split. Five state-of-the-art RRS models are used as teacher models and competitors: DeepCoNN [35], MPCN [25], NARRE [2], DAML [17], D\_ATTEN [11]. Other baselines include simple RRS and state-of-the-art KD-based RS: (1) CNN: we train a student network with a CNN encoding module and a FFN prediction layer via the student loss in Eq. 5. This baseline does not



## REFERENCES

- [1] Charu C. Aggarwal. 2016. *Recommender Systems - The Textbook*. Springer.
- [2] Chong Chen, Min Zhang, Yiqun Liu, and Shaoping Ma. 2018. Neural Attentional Rating Regression with Review-level Explanations. In *WWW*. 1583–1592.
- [3] Jiawei Chen, Hande Dong, Yang Qiu, and Xiangnan He. 2021. AutoDebias: Learning to Debias for Recommendation. In *SIGIR*. 21–30.
- [4] Jiawei Chen, Hande Dong, Xiang Wang, Fuli Feng, Meng Wang, and Xiangnan He. 2020. Bias and Debias in Recommender System: A Survey and Future Directions. *arXiv Preprint (2020)*. <https://arxiv.org/abs/2010.03240>
- [5] Xu Chen, Yongfeng Zhang, Hongteng Xu, Zheng Qin, and Hongyuan Zha. 2019. Adversarial Distillation for Efficient Recommendation with External Knowledge. *ACM Trans. Inf. Syst.* 37, 1 (2019), 12:1–12:28.
- [6] Thomas G. Dietterich. 2000. Ensemble Methods in Machine Learning. In *Multiple Classifier Systems*, Vol. 1857. 1–15.
- [7] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. In *ICML*, Vol. 70. 1126–1135.
- [8] Jianping Gou, Baosheng Yu, Stephen J. Maybank, and Dacheng Tao. 2021. Knowledge Distillation: A Survey. *Int. J. Comput. Vis.* 129, 6 (2021), 1789–1819.
- [9] Huifeng Guo, Jinkai Yu, Qing Liu, Ruiming Tang, and Yuzhou Zhang. 2019. PAL: a position-bias aware learning framework for CTR prediction in live recommender systems. In *RecSys*. 452–456.
- [10] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the Knowledge in a Neural Network. *arXiv Preprint (2015)*. <https://arxiv.org/abs/1503.02531>
- [11] Dongmin Hyun, Chanyoung Park, Min-Chul Yang, Ilhyeon Song, Jung-Tae Lee, and Hwanjo Yu. 2018. Review Sentiment-Guided Scalable Deep Recommender System. In *SIGIR*. 965–968.
- [12] Yehuda Koren, Robert M. Bell, and Chris Volinsky. 2009. Matrix Factorization Techniques for Recommender Systems. *Computer* 42, 8 (2009), 30–37.
- [13] Wonbin Kweon, SeongKu Kang, and Hwanjo Yu. 2021. Bidirectional Distillation for Top-K Recommender System. In *WWW*. 3861–3871.
- [14] Jae-woong Lee, Minjin Choi, Jongwuk Lee, and Hyunjung Shim. 2019. Collaborative Distillation for Top-N Recommendation. In *ICDM*. 369–378.
- [15] Ruihui Li, Xian Wu, Xiushi Chen, and Wei Wang. 2020. Few-Shot Learning for New User Recommendation in Location-based Social Networks. In *WWW*. 2472–2478.
- [16] Chen Lin, Xinyi Liu, Guipeng Xv, and Hui Li. 2021. Mitigating Sentiment Bias for Recommender Systems. In *SIGIR*. 31–40.
- [17] Donghua Liu, Jing Li, Bo Du, Jun Chang, and Rong Gao. 2019. DAML: Dual Attention Mutual Learning between Ratings and Reviews for Item Recommendation. In *KDD*. 344–352.
- [18] Yiming Liu, Xuezhai Cao, and Yong Yu. 2016. Are You Influenced by Others When Rating?: Improve Rating Prediction by Conformity Modeling. In *RecSys*. 269–272.
- [19] Benjamin M. Marlin, Richard S. Zemel, Sam T. Roweis, and Malcolm Slaney. 2007. Collaborative Filtering and the Missing at Random Assumption. In *UAI*. 267–275.
- [20] Julian J. McAuley and Jure Leskovec. 2013. Hidden factors and hidden topics: understanding rating dimensions with review text. In *RecSys*. 165–172.
- [21] Julian J. McAuley, Rahul Pandey, and Jure Leskovec. 2015. Inferring Networks of Substitutable and Complementary Products. In *KDD*. 785–794.
- [22] Sangwon Park and Juan L. Nicolau. 2015. Asymmetric effects of online consumer reviews. *Annals of Tourism Research* 50 (2015), 67–83. Issue 2015.
- [23] Peijie Sun, Le Wu, Kun Zhang, Yanjie Fu, Richang Hong, and Meng Wang. 2020. Dual Learning for Explainable Recommendation: Towards Unifying User Preference Prediction and Review Generation. In *WWW*. 837–847.
- [24] Jiayi Tang and Ke Wang. 2018. Ranking Distillation: Learning Compact Ranking Models With High Performance for Recommender System. In *KDD*. 2289–2298.
- [25] Yi Tay, Anh Tuan Luu, and Siu Cheung Hui. 2018. Multi-Pointer Co-Attention Networks for Recommendation. In *KDD*. 2309–2318.
- [26] Quoc-Tuan Truong and Hady W. Lauw. 2019. Multimodal Review Generation for Recommender Systems. In *WWW*. 1864–1874.
- [27] Hongning Wang, Yue Lu, and ChengXiang Zhai. 2011. Latent aspect rating analysis without aspect keyword supervision. In *KDD*. 618–626.
- [28] Xiaojie Wang, Rui Zhang, Yu Sun, and Jianzhong Qi. 2018. KDGAN: Knowledge Distillation with Generative Adversarial Networks. In *NeurIPS*. 783–794.
- [29] Yixin Wang, Dawen Liang, Laurent Charlin, and David M. Blei. 2020. Causal Inference for Recommender Systems. In *RecSys*. 426–431.
- [30] Jae woong Lee, Seongmin Park, and Jongwuk Lee. 2021. Dual Unbiased Recommender Learning for Implicit Feedback. In *SIGIR*. 1647–1651.
- [31] Le Wu, Xiangnan He, Xiang Wang, Kun Zhang, and Meng Wang. 2021. A Survey on Neural Recommendation: From Collaborative Filtering to Content and Context Enriched Recommendation. *arXiv Preprint (2021)*. <https://arxiv.org/abs/2104.13030>
- [32] Yang Zhang, Fuli Feng, Xiangnan He, Tianxin Wei, Chonggang Song, Guohui Ling, and Yongdong Zhang. 2021. Causal Intervention for Leveraging Popularity Bias in Recommendation. In *SIGIR*. 11–20.
- [33] Yuan Zhang, Xiaoran Xu, Hanning Zhou, and Yan Zhang. 2020. Distilling Structured Knowledge into Embeddings for Explainable and Accurate Recommendation. In *WSDM*. 735–743.
- [34] Yin Zhang, Yueting Zhuang, Jiangqin Wu, and Liang Zhang. 2009. Applying probabilistic latent semantic analysis to multi-criteria recommender system. *AI Commun.* 22, 2 (2009), 97–107.
- [35] Lei Zheng, Vahid Noroozi, and Philip S. Yu. 2017. Joint Deep Modeling of Users and Items Using Reviews for Recommendation. In *WSDM*. 425–434.
- [36] Jin Peng Zhou, Zhao Yue Cheng, Felipe Pérez, and Maksims Volkovs. 2020. TAFA: Two-headed Attention Fused Autoencoder for Context-Aware Recommendations. In *RecSys*. 338–347.
- [37] Jieming Zhu, Jinyang Liu, Weiqi Li, Jincan Lai, Xiuqiang He, Liang Chen, and Zibin Zheng. 2020. *Ensembled CTR Prediction via Knowledge Distillation*. 2941–2958.