

Debias the Black-box: A Fair Ranking Framework via Knowledge Distillation

Zhitao Zhu^{1,2}, Shijing Si^{1,3}, Jianzong Wang^{1(✉)}, Yaodong Yang^{4(✉)}, and Jing Xiao¹

¹ Ping An Technology (Shenzhen) Co., Ltd., Shenzhen, China
andyzzt@mail.ustc.edu.cn, shijing.si@outlook.com, jzwang@188.com,
xiaojing661@pingan.com.cn

² IAT, University of Science and Technology of China, Hefei, China

³ School of Economics and Finance, Shanghai International Studies University,
Shanghai, China

⁴ Institute for AI, Peking University, Beijing, China
yaodong.yang@pku.edu.cn

Abstract. Deep neural networks can capture the intricate interaction history information between queries and documents, because of their many complicated nonlinear units, allowing them to provide correct search recommendations. However, service providers frequently face more complex obstacles in real-world circumstances, such as deployment cost constraints and fairness requirements. Knowledge distillation, which transfers the knowledge of a well-trained complex model (teacher) to a simple model (student), has been proposed to alleviate the former concern, but the best current distillation methods focus only on how to make the student model imitate the predictions of the teacher model. To better facilitate the application of deep models, we propose a fair information retrieval framework based on knowledge distillation. This framework can improve the exposure based fairness of models while considerably decreasing model size. Our extensive experiments on three huge datasets show that our proposed framework can reduce the model size to a minimum of 1% of its original size while maintaining its black-box state. It also improves fairness performance by 15%~46% while keeping a high level of recommendation effectiveness.

Keywords: Information Retrieval · Knowledge distillation · Fairness · Learning to rank · Exposure

1 Introduction

Information Retrieval (IR) systems are nowadays one of the most pervasive techniques in a variety of industries. The sophisticated architectures and growing

Co-corresponding authors: Jianzong Wang (jzwang@188.com) and Yaodong Yang (yaodong.yang@pku.edu.cn).

data scale of application scenarios cause the size of models to increase rapidly. Large models tend to capture complicated interactions between queries and documents, yielding increased performance at the expense of increased computational time and memory phase.

To tackle the difficulty of applying such large models to web-scale and real-time platforms, a few recent works [8, 10, 11, 23] have applied Knowledge Distillation (KD) [6] to IR. Most ranking distillation approaches, however, focus on the balance of prediction performance and computing speed. Little attention has been paid to the fairness of ranking models during the distillation process [20]. Ranking systems employ deterministic models to assign an individual score to each item, and then sort items in descending order of their assigned scores to obtain rankings. That calculation pattern is succinct and intuitive, yet, leads to unfairness in the distribution of exposure. Exposure represents the expected number of people who will check an item. User behavior is affected by position bias: they are less likely to check items at lower positions. As illustrated in Figure 1, this results in the allocation of user attention being disproportionate to the rankings on the recommendation list. A small difference in relevance can

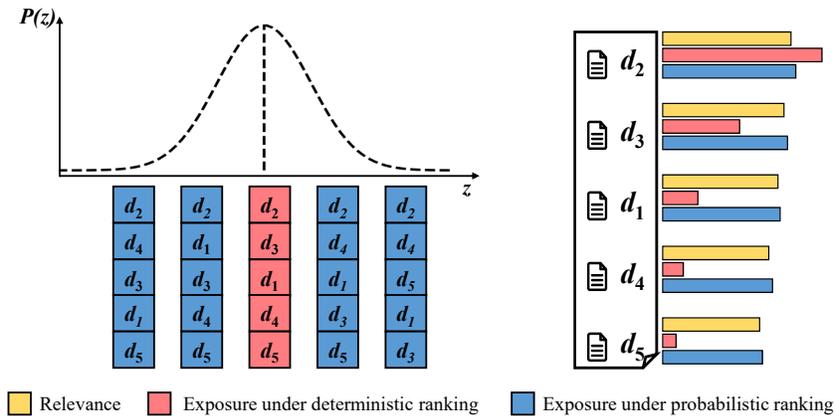


Fig. 1. A simple ranking method that relies only on ordering relevance scores makes the ratio between exposure and relevance disproportionate, magnifying subtle gaps in relevance. In contrast, a probabilistic ranking model maintains a positive relationship between exposure and relevance. z denotes possible permutations.

make a world of difference in exposure as a result of the winner-take-all [20] allocation of exposure. Under this circumstance, the unfairness between candidates is significantly amplified. Furthermore, while the combination of deep learning with ranking models provides a fresh leap forward, its lack of interpretability and increased reliance on data may introduce endogenous biases. The trade-off between model performance and fairness has become a huge issue for ranking models.

In this research, we propose a general Fairness-aware Ranking Distillation framework (FRD) for improving the fairness of top- K ranking models. FRD eliminates the average disparity of exposures documents received, and it is a direct approach irrelevant to specific protected attributes. Our framework takes advantage of ranking distillation, so it can build on top of large well-trained ranking models without extra retraining cost. The main contributions of this work can be summarized as follows:

- We adopt a ranking distillation framework with sub-regional treatment that retains the top rankings while penalizing items ranked lower by teachers. Since only the soft label results of the teacher model are used, our framework can be effectively applied to the black-box models.
- We identify the great potential of ranking distillation and propose that implementing a personalized fairness correction in the process of ranking distillation can avoid the vast additional computational costs necessary to achieve fairness purposes directly on complex models.
- By applying the latest optimization method of the PL ranking model, we have further reduced the computational cost of fairness correction significantly.

2 Related Work

2.1 Knowledge Distillation

The initial KD method transfers knowledge through the softmax output of a teacher model. Some subsequent work has extended the scope of statistics used for matching, such as intermediate feature responses [3], gradient [22], and distribution [7], etc. Other work, on the other hand, opted for a more subtle structural design. For example, DE-RRD [8] uses distillation experts to learn the middle layer representation space mapping function of Teacher while additionally using the rank order given by Teacher for rank matching; DCD [12] assigns more training resource to instances that were not correctly predicted by Student; MiniLM [24] proposes an assistant mechanism to distill only the last layer of the self-attentive matrix and the value-value matrix of the pre-trained model. But the potential of KD for other objectives remains largely unexplored.

2.2 Fairness in rankings

Despite the growing impact of online information systems on our society and economy, the fairness of rankings has been a relatively under-explored area. In the existing work, some consider the equity of groups with respect to a set of categorical sensitive attributes (or features) to be ranked according to the principle of population parity, which can be further divided into group fairness [9] and individual fairness [15]. But there are other works such as [4, 20] that argue that the fairness of ranking systems corresponds to how they assign exposure based on the merits of individual items or item groups. These works specify and enforce fairness constraints that explicitly link relevance to exposure in expectation or amortized over a set of queries [21].

3 Methodology

In this section, we illustrate FRD, our universal fairness-aware ranking distillation framework, which is made up of ranking distillation and a fairness penalty. We elaborate on these two components in this section and present the algorithm for FRD.

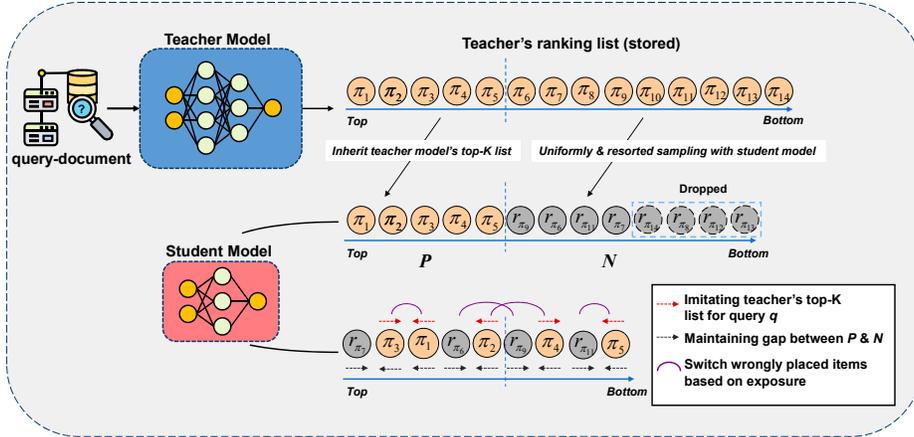


Fig. 2. Illustration of FRD framework. Student studies to rearrange the items into two sets to (a). maintain the same top- K list as the teacher; (b). penalize negative set items that are ranked too high; (c). ensure that items expected to get more exposure are placed higher up the list by pairwise comparison

3.1 Ranking Distillation

Ranking distillation attempts to mimic the teacher ranking model by transferring its knowledge to the student. For we only care about the top performers, considering the scores of all query-document pairs is undoubtedly inefficient and wasteful. A straightforward approach is to use only the information of top- K items and discard items ranked low by the teacher like [23]. But this approach leads to performance degradation because some useful information is discarded. It is undesirable to discard directly or incorporate wholesale into the training for candidate items that are not on the teacher model recommendation list.

1). For items recommended by Teacher, student network learns to match all their orders to get fine-grained scoring results.

2). Candidates not recommended by the teacher model are also internally differentiated: items with very low scores do not help maintain a distinct score boundary, and the model should be allowed to learn more information from candidates close to the boundary.

FRD treats Teacher's top- K candidates as positive set P to help the student imitate top performers' order and samples out negative set N to keep less relevant documents at lower end. Inspired by [19], two categorical ranking losses

$\mathcal{L}_P(t, s, P)$ and $\mathcal{L}_N(s, P, N)$ are employed to allocate computing resources rationally according to importance, where $t, s \in \mathbb{R}^A$ represent the score list generated by Teacher and Student respectively. For the top-K list containing more refined sorting information of f^T , we construct a softmax based cross-entropy loss to minimize the difference between student’s and teacher’s top-p lists (positive set P). It is worth mentioning that here we relax K to p in Algorithm 1 to give more items recommendation opportunities.

$$\mathcal{L}_P(t, s, P) = \sum_{i \in P} \left\{ -\frac{e^{t_i}}{\delta_P^t} \cdot \log \frac{e^{s_i}}{\delta_P^s} \right\} \text{ where } \delta_P^t = \sum_{j \in P} e^{t_j} \quad (1)$$

Later, for the large negative sample set, the binary or pairwise comparison loss functions like Equation (2) with less complexity can be employed since we only need to discriminate the rest items without considering the precise ranking. A negative sampling strategy is adopted to further reduce the number of negative candidates. A compressed set C is first sampled out with a categorical distribution Q , and then score with the student model to pick the top n documents to form the negative set N . This process brings a huge reduction in the calculation amount but no loss of effect.

$$\begin{aligned} \mathcal{L}_{N_b}(s, N) &= \sum_{i \in N} \log(1 + e^{s_i}) && \text{(Binary)} \\ \mathcal{L}_{N_p}(s, P, N) &= \sum_{i \in N} \sum_{j \in P} \log(1 + e^{s_i - s_j}) && \text{(Pairwise)} \end{aligned} \quad (2)$$

This categorical distillation strategy helps us concentrate finite computational power on the most useful retrieved list displayed to users while making good use of the remaining items.

3.2 Exposure-based Fair Ranking

It is widely accepted that an item’s position in the rankings has a critical impact on its exposure and financial success. Yet, surprisingly, the algorithms used to learn these rankings are often blind to their impact on items. The key goal of our amelioration is to promote exposure-based fairness during distillation. Based on the exposure estimating approach proposed below, we can explicitly specify how exposure is allocated such as making exposure proportional to relevance.

In recent years, the popularity of the Plackett-Luce (PL) ranking model [14] has increased. It is more efficient to use the PL ranking model to model the probabilistic distribution over rankings directly. Especially for fair distributions of attention exposure, multiple lines of previous works have found that PL ranking models may transfer the deterministic recommendation pattern into reflecting the relevance by the probability of occurrence [4, 20, 21], which gives an approximately equal possibility of being the top-item to candidates with slightly lower relevance indicators.

In the meantime, we follow the definition of exposure-based fairness: the items with higher ranking should achieve more exposure in practical environment. Thus we can constraint the difference between existing arrangement and

counterfactual arrangement by minimize their exchanged exposure difference. For brevity we denote $\varepsilon_d = \varepsilon(q, d)$ as the exposure an item d receives under a probabilistic model π :

$$\varepsilon_d = \sum_{z \in \pi} \pi(z) \sum_{k=1}^K \theta_k \mathbb{I}[z_k = d] \quad (3)$$

where z denotes a possible permutation and rank weight θ_k indicates the probability that a user examines the item at rank k (e.g. $\frac{1}{k}$). For each query, this statistic calculates the averaged disparity between every two items and takes the average disparity across all item pairs:

$$\mathcal{L}_{Fair}(s, P) = \frac{1}{|P|(|P|-1)} \sum_{d \in P} \sum_{d' \in P} (\varepsilon_{d'} \mathcal{R}_d - \varepsilon_d \mathcal{R}_{d'})^2 \quad (4)$$

Here the relevance of a document \mathcal{R}_d is set as a transformation of its label and the exposures are estimated by d 's total possibly gained rank weight [16] calculated through a continued product of following softmax principle:

$$\pi(d|z_{1:k-1,A}) = \frac{e^{s(d)} \cdot \mathbb{I}[d \notin z_{1:k-1}]}{\sum_{d' \in A \setminus z_{1:k-1}} e^{s(d')}} \quad (5)$$

where $s(d)$ denotes the score that the student model assigns to document d .

Calculating the gradient of a PL ranking model, on the other hand, necessitates iterating through every conceivable ranking that the model may generate, i.e., every possible permutation. In reality, this computational impossibility is overcome by estimating the gradient using model rankings as samples [17]. Employing the estimator to the gradient of PL model proposed in [16], we can optimize the exposure-based unfairness metrics simultaneously during knowledge distillation. This Metric (4) is designed to measure the difference in reward if the exposures of two items were swapped. Substituting the derivative of this metric with respect to ε_d , we can employ the below-mentioned estimator [16] to estimate its gradient and then update the PL model by back propagation, m denotes the ranking merit and an entire ranking is sampled with Gumble Softmax [5] trick rather than calculating any of the actual probabilities [1].

$$\begin{aligned} \frac{\delta}{\delta m} \text{Metric}(q) &= \sum_{d \in P} \left[\frac{\delta}{\delta m} m(d) \right] \mathbb{E}_z \left[\left(\sum_{k=\text{rank}(d,z)+1}^K \theta_k \left[\frac{\delta \text{Metric}(q)}{\delta \varepsilon_{z_k}} \right] \right. \right. \\ &\quad \left. \left. + \sum_{i=1}^{\text{rank}(d,z)} \pi(d|z_{1:k-1}) (\theta_k \left[\frac{\delta \text{Metric}(q)}{\delta \varepsilon_{z_k}} \right]) \right] \right] \end{aligned} \quad (6)$$

In total, the loss of FRD is

$$\mathcal{L} = \mathcal{L}_P(t, s, P) + \mathcal{L}_N(s, P, N) + \lambda \mathcal{L}_{Fair}(s, P) \quad (7)$$

where λ is the tuning parameter specific to the dataset. This framework is model-agnostic as it only requires the teacher network's scoring data, making it easy to distill black box models. It is also able to debias ranking models against many attributes as it does not require any side information.

Algorithm 1: FRD

- 1: Student model f^S and well-trained Teacher model f^T .
Hyper-parameters: Positive set size $p \in [K, a]$, Compressed set size $c \leq a - p$, Negative set size $n \leq c$, Learning rate η , Epochs R .
 - 2: **for** each $r \in [0, R - 1]$ **do**
 - 3: Uniformly randomly select an query x , denote its candidate set of a documents as A with stored teacher scores t
 - 4: /* Indices of vector's largest p elements */
 Compute positive set $P = Top_p(t)$
 - 5: Randomly sample the compressed set C of c items with $Q(A - P)$
 - 6: Compute negative set $N = Top_n(f^S(C))$
 - 7: Compute student scores $s = f^S(P \cup N)$
 - 8: Update f^S by optimizing the total loss in Equation (7)
 - 9: **end for**
-

4 Experiments

In this section, we show empirical performance for FRD. We compare our method with the other baselines on two aspects: 1). ranking performance through imitating teacher; 2). unfairness correction level during distillation.

4.1 Datasets and Training Configurations

Datasets. We conduct our experiments on three public datasets: MSLR-WEB30K (fold 1) [18], Yahoo LETOR [2] and Istella LETOR full dataset [13], which are all benchmark datasets for large-scale experiments on the efficiency and scalability of LTR solutions. The basic statistics of these datasets can be found in [19].

Teacher/Student Models. To better control the ratio of distillation, we employ the Multi-Layered Perceptron (MLP) as the architecture for both the teacher and student networks. Specifically, the teacher model is a 3 layer FC-BN-ReLU model with hidden units of sizes 1024, 512, 256 for all three datasets and the student models are chosen by varying the proportion of parameter counts of the teacher and student models.

Evaluation Protocols. Following [19, 23], we measure the ranking performance of student models in terms of the normalized discounted cumulative gain (NDCG) at top 5 & 10, which is commonly used for ranking tasks. For the fairness of a ranking model, we utilize the averaged squared disparity in Equation (4) over all queries as the disparity metric.

4.2 Ranking and Fairness Performance

The distillation and fairness performances of teacher and student models on three datasets are presented in Table 1. The student model is a two-layered FC-BN-ReLU model of hidden units of sizes 28, 28. We utilize two kinds of distillation

loss, the pairwise and binary loss, shown in methods ending with ‘P’ and ‘B’, respectively.

FRD consistently outperforms RankDistil and the teacher model in fairness. The averaged disparity of FRD methods decreases by 15%~46%, compared with the teacher model. The student models of RankDistil also exhibit a slight improvement in fairness, which is caused by the reduced size of the model. Note that we attribute to the differences unfairness values due to the different composition of the datasets, for example, each query in the Istella test set corresponds to an average of 319 documents, while for the Yahoo dataset, the number is 24. Deeper reasons remain to be explored, but for now, this fairness metric is not applicable to cross-dataset comparisons. With regards to the ranking perfor-

Table 1. Performance of various distillation methods on three datasets. Methods ending with ‘B’ indicate that binary loss is used for distillation, and ‘P’ denotes pairwise distillation loss. ‘Disp.’ denotes the averaged disparity over all queries.

Method/Data	MSLR Web30K			Yahoo LETOR			Istella		
	N ₅ ↑	N ₁₀ ↑	Disp. ↓	N ₅ ↑	N ₁₀ ↑	Disp. ↓	N ₅ ↑	N ₁₀ ↑	Disp. ↓
Teacher	0.467	0.488	0.080	0.722	0.764	0.647	0.628	0.684	0.013
Ranking Distillation [23]	0.428	0.451	0.080	0.684	0.731	0.641	0.490	0.525	0.013
RankDistil-B [19]	0.423	0.457	0.075	0.683	0.725	0.621	0.490	0.525	0.011
RankDistil-P [19]	0.457	0.479	0.076	0.710	<u>0.750</u>	0.625	0.481	0.528	0.012
FRD-B(<i>ours</i>)	0.420	0.452	0.063	0.678	0.721	0.551	0.483	0.520	0.008
FRD-P(<i>ours</i>)	<u>0.455</u>	<u>0.478</u>	0.059	<u>0.706</u>	0.751	0.550	<u>0.487</u>	<u>0.526</u>	0.007

mance, student models trained with pairwise loss significantly outperform binary loss. Both RankDistil and FRD with pairwise loss yield similar performances to the teacher model on MSLR Web30K and Yahoo LETOR datasets. On Istella dataset, student models perform significantly worse than the teacher model due to their lack of capacity to capture enough patterns from the huge dataset.

4.3 Performance of Student Models versus Size

Table 2 displays the performance of pairwise distillation methods on two datasets while varying the size of student models from 50% to 1% of Teacher. We vary the size of student models by adjusting the number of hidden neuron units. Fewer parameters means faster inference.

Both RankDistil and FRD methods with pairwise loss are robust to the size of student models, as the NDCG₅ scores of the 1% student models are comparable to the teachers. FRD can reduce the average squared disparity significantly faster than RankDistil, which indicates that the fairness penalty plays an essential role in alleviating the disparity in ranking. The distillation performance of FRD has some decrease compared to RankDistil, but not more than 0.06%, which is negligible considering the improvement in fair performance.

Table 2. Performance of distillation methods while varying the size of student models. The ‘NDCG’ is evaluated on top-5 documents, and ‘Disp.’ is the averaged disparity over all queries.

Models	Size	MSLR Web30K		Yahoo LETOR	
		N ₅ ↑	Disp. ↓	N ₅ ↑	Disp. ↓
Teacher	100%	0.467	0.080	0.722	0.647
RankDistil-P	50%	0.435	0.079	0.721	0.635
FRD-P	50%	0.433	0.069	0.720	0.578
RankDistil-P	10%	0.431	0.078	0.721	0.612
FRD-P	10%	0.428	0.064	0.719	0.563
RankDistil-P	1%	0.422	0.074	0.715	0.589
FRD-P	1%	0.421	0.059	0.714	0.550

5 Conclusion

In this paper, we propose FRD, a fairness-aware ranking distillation framework that leverages the potential of knowledge distillation. It can inherit the excellent ranking retrieval capabilities of the teacher model in the black-box state, while reducing the model size to a minimum of one percent. We pioneer the use of KD in achieving the exposure fairness. Our experiments show that the FRD with a bias correction strategy can achieve a significant reduction in model size and a large improvement in the reasonable distribution of exposure.

Acknowledgements Zhitao Zhu and Shijing Si contributed equally to this work. Co-corresponding authors: Jianzong Wang and Yaodong Yang . This work is supported by the Key Research and Development Program of Guangdong Province under grant No.2021B0101400003.

References

1. Bruch, S., Han, S., Bendersky, M., Najork, M.: A stochastic treatment of learning to rank scoring functions. In: WSDM. pp. 61–69. ACM (2020)
2. Chapelle, O., Chang, Y.: Yahoo! learning to rank challenge overview. In: Proceedings of the Yahoo! Learning to Rank Challenge, held at ICML. JMLR Proceedings, vol. 14, pp. 1–24. JMLR.org (2011)
3. Denton, E.L., Zaremba, W., Bruna, J., LeCun, Y., Fergus, R.: Exploiting linear structure within convolutional networks for efficient evaluation. In: NeuIPS. pp. 1269–1277 (2014)
4. Diaz, F., Mitra, B., Ekstrand, M.D., Biega, A.J., Carterette, B.: Evaluating stochastic rankings with expected exposure. In: CIKM. pp. 275–284. ACM (2020)
5. Gumbel, E.J.: Statistical theory of extreme values and some practical applications: a series of lectures, vol. 33. US Government Printing Office (1954)
6. Hinton, G.E., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. In: NeuIPS. pp. 1–9. MIT Press (2015)
7. Huang, Z., Wang, N.: Like what you like: Knowledge distill via neuron selectivity transfer. CoRR [abs/1707.01219](#) (2017)

8. Kang, S., Hwang, J., Kweon, W., Yu, H.: DE-RRD: A knowledge distillation framework for recommender system. In: CIKM. pp. 605–614. ACM (2020)
9. Kusner, M.J., Loftus, J.R., Russell, C., Silva, R.: Counterfactual fairness. In: NeurIPS. pp. 4066–4076 (2017)
10. Kweon, W., Kang, S., Yu, H.: Bidirectional distillation for top-k recommender system. In: WWW. pp. 3861–3871. ACM / IW3C2 (2021)
11. Lee, J., Choi, M., Lee, J., Shim, H.: Collaborative distillation for top-n recommendation. In: ICDM. pp. 369–378. IEEE (2019)
12. Lee, Y., Kim, K.E.: Dual correction strategy for ranking distillation in top-n recommender system. In: CIKM. pp. 3186–3190. ACM (2021)
13. Lucchese, C., Nardini, F.M., Orlando, S., Perego, R., Tonellotto, N., Venturini, R.: Exploiting CPU SIMD extensions to speed-up document scoring with tree ensembles. In: SIGIR. pp. 833–836. ACM (2016)
14. Luce, R.D.: Individual choice behavior: A theoretical analysis. Courier Corporation (2012)
15. Maity, S., Xue, S., Yurochkin, M., Sun, Y.: Statistical inference for individual fairness. In: ICLR. pp. 1–19. OpenReview.net (2021)
16. Oosterhuis, H.: Computationally efficient optimization of plackett-luce ranking models for relevance and fairness. In: SIGIR. pp. 1023–1032. ACM (2021)
17. Oosterhuis, H., de Rijke, M.: Unifying online and counterfactual learning to rank: A novel counterfactual estimator that effectively utilizes online interventions. In: WSDM. pp. 463–471. ACM (2021)
18. Qin, T., Liu, T.: Introducing LETOR 4.0 datasets. CoRR **abs/1306.2597**, 1–6 (2013)
19. Reddi, S., Pasumarthi, R.K., Menon, A., Rawat, A.S., Yu, F., Kim, S., Veit, A., Kumar, S.: Rankdistil: Knowledge distillation for ranking. In: AISTATS. pp. 2368–2376. PMLR, JMLR (2021)
20. Singh, A., Joachims, T.: Fairness of exposure in rankings. In: SIGKDD. pp. 2219–2228. ACM (2018)
21. Singh, A., Joachims, T.: Policy learning for fairness in ranking. In: NeurIPS. pp. 5427–5437. MIT Press (2019)
22. Srinivas, S., Fleuret, F.: Knowledge transfer with jacobian matching. In: ICML. Proceedings of Machine Learning Research, vol. 80, pp. 4730–4738. PMLR (2018)
23. Tang, J., Wang, K.: Ranking distillation: Learning compact ranking models with high performance for recommender system. In: SIGKDD. pp. 2289–2298. ACM (2018)
24. Wang, W., Wei, F., Dong, L., Bao, H., Yang, N., Zhou, M.: Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. In: NeurIPS (2020)