
Incorporating Bias-aware Margins into Contrastive Loss for Collaborative Filtering

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Collaborative filtering (CF) models easily suffer from popularity bias, which makes
2 recommendation deviate from users' actual preferences. However, most current
3 debiasing strategies are prone to playing a trade-off game between head and tail
4 performance, thus inevitably degrading the overall recommendation accuracy. To
5 reduce the negative impact of popularity bias on CF models, we incorporate Bias-
6 aware margins into Contrastive loss and propose a simple yet effective **BC Loss**,
7 where the margin tailors quantitatively to the bias degree of each user-item interac-
8 tion. We investigate the geometric interpretation of BC loss, then further visualize
9 and theoretically prove that it simultaneously learns better head and tail represen-
10 tations by encouraging the compactness of similar users/items and enlarging the
11 dispersion of dissimilar users/items. Over six benchmark datasets, we use BC
12 loss to optimize two high-performing CF models. On various evaluation settings
13 (*i.e.*, imbalanced/balanced, temporal split, fully-observed unbiased, tail/head test
14 evaluations), BC loss outperforms the state-of-the-art debiasing and non-debiasing
15 methods with remarkable improvements. Considering the theoretical guarantee
16 and empirical success of BC loss, we advocate using it not just as a debiasing
17 strategy, but also as a standard loss in recommender models. Codes are available at
18 <https://anonymous.4open.science/r/BC-Loss-8764/model.py>.

19 1 Introduction

20 At the core of leading collaborative filtering (CF) models is the learning of high-quality representations
21 of users and items from historical interactions. However, most CF models easily suffer from the
22 popularity bias issue in the interaction data [1, 2, 3, 4]. Specifically, the training data distribution is
23 typically long-tailed, *e.g.*, a few head items occupy most of the interactions, whereas the majority
24 of tail items are unpopular and receive little attention. The CF models built upon the imbalanced
25 data are prone to learn the popularity bias and even amplify it by over-recommending head items and
26 under-recommending tail items. As a result, the popularity bias causes the biased representations
27 with poor generalization ability, making recommendations deviate from users' actual preferences.

28 Motivated by concerns of popularity bias, studies on debiasing have been conducted to lift the tail
29 performance. Unfortunately, most prevalent debiasing strategies focus on the trade-off between
30 head and tail evaluations (see Table 3), including post-processing re-ranking [5, 6, 7, 8, 9], balanced
31 training loss [10, 11, 12, 9], sample re-weighting [13, 14, 15, 16, 17, 18], and head bias removal
32 by causal inference [19, 20, 21, 22]. Worse still, many of them hold some assumptions that are
33 infeasible in practice, such as the balanced test distribution is known in advance to guide the
34 hyperparameters' adjustment [23, 22], or a small unbiased data is present to train the unbiased model
35 [24, 19]. Consequently, they pursue improvements on tail items but exacerbate the performance

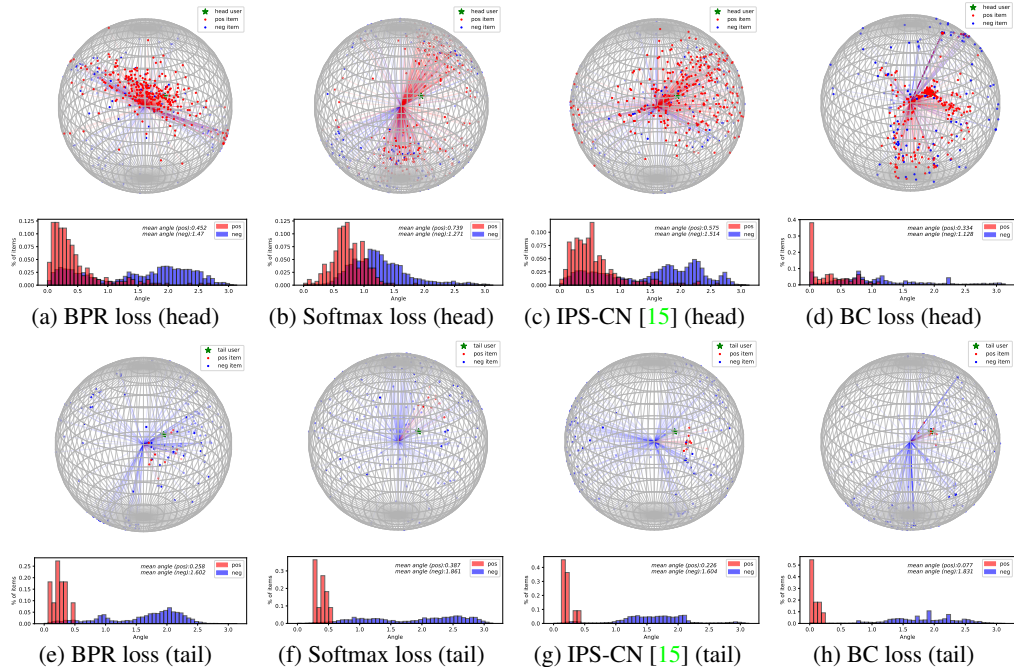


Figure 1: Visualizations of item representations learned by LightGCN [25] on Yelp2018 [25], where subfigures (a-d)/(e-h) depict the identical head/tail user as a green star, while the red and blue points denote positive and negative items, respectively. In each subfigure, the first row presents the 3D item representations projected on the unit sphere, while the second row shows the angle distribution of items *w.r.t.* the specific user and the statistics of mean angles. Compared to other losses, BC loss learns better head representations (*cf.* with the smallest mean positive angle, the vast majority of positive items fall into the group closest to the user) and tail representations (*cf.* a clear margin exists between positive and negative items for the tail user). **BC loss learns a more reasonable representation distribution that is locally clustered and globally separated.** See more details in Appendix A.1.

36 sacrifice of head items, leading to a severe overall performance drop. The trade-off between the head
 37 and tail evaluations results in suboptimal representations, which derails the generalization ability.

38 In this paper, we conjecture that an ideal debiasing strategy should learn high-quality head and tail
 39 representations with powerful discrimination and generalization abilities, rather than playing a trade-
 40 off game between the head and tail performance. Here we follow the prior studies [15, 23, 13, 14]
 41 to focus on one key ingredient in representation learning: the loss function. Figure 1 depicts the
 42 item representations, which is optimized via two non-debiasing losses (BPR [26] and Softmax [27])
 43 and one debiasing loss (IPS-CN [15]). Wherein, representation discrimination is reflected in how
 44 well the positive items of a user are apart from the negatives. Our insights are: (1) For a user, the
 45 non-debiasing losses are inadequate to discriminate his/her positive and negative items well, since
 46 their representations are largely overlapped as Figures 4a and 4b show; (2) Although IPS-CN achieves
 47 better discrimination power in the tail group than BPR (*cf.* positive items get smaller angles to the
 48 ego user in Figure 4g, as compared to Figure 4e), it gets worse discrimination ability in the head (*cf.*
 49 positive items hold larger angles to the ego user in Figure 4c, as compared to Figure 4a).

50 Towards this end, we incorporate Bias-aware margins into Contrastive Loss and devise a simple yet
 51 effective **BC Loss** to guide the head and tail representation learning of CF models. Specifically, we
 52 first employ a bias degree extractor to quantify the influence of interaction-wise popularity — that
 53 is, how well an interaction is predicted, when only popularity information of the target user and
 54 item is used. Interactions involving inactive users and unpopular items often align with lower bias
 55 degrees, indicating that popularity fails to reflect user preference faithfully. In contrast, interactions
 56 with active users and popular items are spurred by the popularity information, thus easily inclining to
 57 high bias degrees. We then move on to train the CF model by converting the bias degrees into the
 58 angular margins between user and item representations. If the bias degree is low, we impose a larger
 59 margin to strongly squeeze the tightness of representations. In contrast, if the bias degree is large, we
 60 exert a small or vanishing margin to reduce the influences of biased representations. Through this

61 way, for each ego user’s representation, BC quantitatively controls its bias-aware margins with item
 62 representations — adaptively intensifying the representation similarity among positive items, while
 63 diluting that among negative items. Benefiting from stringent and discriminative representations, BC
 64 loss significantly improves both head and tail performance.

65 Furthermore, BC loss has three desirable advantages. First, it has a clear geometric interpretation,
 66 as illustrated in Figure 2. Second, it brings forth a simple but effective mechanism of hard example
 67 mining (See Appendix A.2). Third, we theoretically reveal that BC loss tends to learn a low-entropy
 68 cluster for positive pairs (*e.g.*, compactness of matched users and items) and a high-entropy space
 69 for negative pairs (*e.g.*, dispersion of unmatched users and items) (See Theorem 1). Considering the
 70 theoretical guarantee and empirical effectiveness, we argue that BC loss is not only promising to
 71 alleviate popularity bias, but also suitable as a standard learning strategy in CF.

72 2 Preliminary of Collaborative Filtering (CF)

73 **Task Formulation.** Personalized recommendation is retrieving a subset of items from a large catalog
 74 to match user preference. Here we consider a typical scenario, collaborative filtering (CF) with
 75 implicit feedback [28], which can be framed as a top- N recommendation problem. Let $\mathcal{O}^+ =$
 76 $\{(u, i) | y_{ui} = 1\}$ be the historical interactions between users \mathcal{U} and items \mathcal{I} , where $y_{ui} = 1$ indicates
 77 that user $u \in \mathcal{U}$ has adopted item $i \in \mathcal{I}$ before. Our goal is to optimize a CF model $\hat{y} : \mathcal{U} \times \mathcal{I} \rightarrow \mathbb{R}$
 78 that latches on user preference towards items.

79 **Modeling Scheme.** Scrutinizing leading CF models [26, 25, 29, 30], we systematize the common
 80 paradigm as a combination of three modules: user encoder $\psi(\cdot)$, item encoder $\phi(\cdot)$, and similarity
 81 function $s(\cdot)$. Formally, we depict one CF model as $\hat{y}(u, i) = s(\psi(u), \phi(i))$, where $\psi : \mathcal{U} \rightarrow \mathbb{R}^d$
 82 and $\phi : \mathcal{I} \rightarrow \mathbb{R}^d$ encode the identity (ID) information of user u and item i into d -dimensional
 83 representations, respectively; $s : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ measures the similarity between user and item
 84 representations. In literature, there are various choices of encoders and similarity functions:

- 85 • Common encoders roughly fall into three groups: ID-based (*e.g.*, MF [26, 29], NMF [31], CMN
 86 [32]), history-based (*e.g.*, SVD++ [29], FISM [33], MultVAE [30]), and graph-based (*e.g.*, GCMC
 87 [34], PinSage [35], LightGCN [25]) fashions. Here we select two high-performing encoders, MF
 88 and LightGCN, as the backbone models being optimized.
- 89 • The widely-used similarity functions include dot product [26], cosine similarity [36], and neural
 90 networks [31]. As suggested in the recent study [36], cosine similarity is a simple yet effective
 91 and efficient similarity function in CF models, having achieved strong performance. For better
 92 interpretation, we take a geometric view and denote it by:

$$s(\psi(u), \phi(i)) = \frac{\psi(u)^\top \phi(i)}{\|\psi(u)\| \cdot \|\phi(i)\|} \doteq \cos(\hat{\theta}_{ui}), \quad (1)$$

93 in which $\hat{\theta}_{ui}$ is the angle between the user representation $\psi(u)$ and item representation $\phi(i)$.

94 **Learning Strategy.** To optimize the model parameters, CF models mostly frame the top- N rec-
 95 ommendation problem into a supervised learning task, and resort to one of three classical learning
 96 strategies: pointwise loss (*e.g.*, binary cross-entropy [37], mean square error [29]), pairwise loss
 97 (*e.g.*, BPR [26], WARP [38]), and softmax loss [28]. Among them, pointwise and pairwise losses are
 98 long-standing and widely-adopted objective functions in CF. However, extensive studies [9, 1, 39]
 99 have analytically and empirically confirmed that using pointwise or pairwise loss is prone to propagate
 100 more information towards the head user-item pairs, which amplifies popularity bias.

101 Softmax loss is much less explored in CF than its application in other domains like CV [40, 41].
 102 Recent studies [36, 42, 43, 44, 45] find that it inherently conducts hard example mining over multiple
 103 negatives and aligns well with the ranking metric, thus attracting a surge of interest in recommendation.
 104 Hence, we cast the minimization of softmax loss [27] as the representative learning strategy:

$$\mathcal{L}_0 = - \sum_{(u,i) \in \mathcal{O}^+} \log \frac{\exp(\cos(\hat{\theta}_{ui})/\tau)}{\exp(\cos(\hat{\theta}_{ui})/\tau) + \sum_{j \in \mathcal{N}_u} \exp(\cos(\hat{\theta}_{uj})/\tau)}, \quad (2)$$

105 where $(u, i) \in \mathcal{O}^+$ is one observed interaction of user u , while $\mathcal{N}_u = \{j | y_{uj} = 0\}$ is the set of
 106 sampled unobserved items that u did not interact with before; τ is the hyper-parameter known as

107 the temperature in softmax [46]. Nonetheless, modifying softmax loss to enhance the discriminative
 108 power of representations and alleviate the popularity bias remains largely unexplored. Therefore,
 109 our work aims to devise a more generic and broadly-applicable variant of softmax loss for CF tasks,
 110 which can improve the long-tail performance fundamentally.

111 3 Methodology of BC Loss

112 On the basis of softmax loss, we devise our BC loss and present its desirable characteristics.

113 3.1 Popularity Bias Extractor

114 Before mitigating popularity bias, we need to quantify the influence of popularity bias on a single
 115 user-item pair. One straightforward solution is to compare the performance difference between
 116 the biased and unbiased evaluations. However, this is not feasible as the unbiased data is usually
 117 unavailable in practice. Statistical metrics of popularity could be a reasonable proxy of the biased
 118 information, such as user popularity statistics $p_u \in \mathbb{P}$ (*i.e.*, the number of historical items that user
 119 u has interacted with before) and item popularity statistics $p_i \in \mathbb{P}$ (*i.e.*, the number of observed
 120 interactions that item i is involved in). If the impact of the interaction between u and i can be
 121 captured well based solely on such statistics, the model is susceptible to exploiting popularity bias for
 122 prediction. Hence, we argue that the popularity-only prediction will delineate the influence of bias.

123 Towards this end, we first train an additional module, termed popularity bias extractor, which only
 124 takes the popularity statistics as input to make prediction. Similar to the modeling of CF (*cf.* Section
 125 2), the bias extractor is formulated as a function $\hat{y}_b : \mathbb{P} \times \mathbb{P} \rightarrow \mathbb{R}$:

$$\hat{y}_b(p_u, p_i) = s(\psi_b(p_u), \phi_b(p_i)) \doteq \cos(\hat{\xi}_{ui}), \quad (3)$$

126 where the user popularity encoder $\psi_b : \mathbb{P} \rightarrow \mathbb{R}^d$ and the item popularity encoder $\phi_b : \mathbb{P} \rightarrow \mathbb{R}^d$
 127 map the popularity statistics of user u and item i into d -dimensional popularity embeddings $\psi_b(p_u)$
 128 and $\phi_b(p_i)$, respectively; $s : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is the cosine similarity function between popularity
 129 embeddings (*cf.* Equation (1)). $\hat{\xi}_{ui}$ is the angle between $\psi_b(p_u)$ and $\phi_b(p_i)$.

130 We then minimize the following softmax loss to optimize the popularity bias extractor:

$$\mathcal{L}_b = - \sum_{(u,i) \in \mathcal{O}^+} \log \frac{\exp(\cos(\hat{\xi}_{ui})/\tau)}{\exp(\cos(\hat{\xi}_{ui})/\tau) + \sum_{j \in \mathcal{N}_u} \exp(\cos(\hat{\xi}_{uj})/\tau)}. \quad (4)$$

131 This optimization enforces the extractor to reconstruct the historical interactions and makes the
 132 reconstruction reflect the interaction-wise bias degree. Furthermore, we can distinguish hard inter-
 133 actions based on the bias degree, *i.e.*, the interactions that can be hardly predicted by popularity
 134 statistics ought to be more informative for representation learning in the target CF model. In a
 135 nutshell, the popularity bias extractor underscores the bias degree of each user-item interaction, which
 136 substantively reflects how hard it is to be predicted.

137 3.2 BC Loss

138 We move on to devise a new BC loss for the target CF model. Our BC loss stems from softmax loss but
 139 converts the interaction-bias degrees into the bias-aware angular margins among the representations
 140 to enhance the discriminative power of representations. Our BC loss is:

$$\mathcal{L}_{\text{BC}} = - \sum_{(u,i) \in \mathcal{O}^+} \log \frac{\exp(\cos(\hat{\theta}_{ui} + M_{ui})/\tau)}{\exp(\cos(\hat{\theta}_{ui} + M_{ui})/\tau) + \sum_{j \in \mathcal{N}_u} \exp(\cos(\hat{\theta}_{uj})/\tau)}, \quad (5)$$

141 where M_{ui} is the bias-aware angular margin for the interaction (u, i) defined as:

$$M_{ui} = \min\{\hat{\xi}_{ui}, \pi - \hat{\theta}_{ui}\} \quad (6)$$

142 where $\hat{\xi}_{ui}$ is derived from the popularity bias extractor (*cf.* Equation (3)), and $\pi - \hat{\theta}_{ui}$ is the upper
 143 bound to restrict $\cos(\cdot + M_{ui})$ to be a monotonically decreasing function. Intuitively, if a user-item
 144 pair (u, i) is the hard interaction that can hardly be reconstructed by its popularity statistics, it holds a

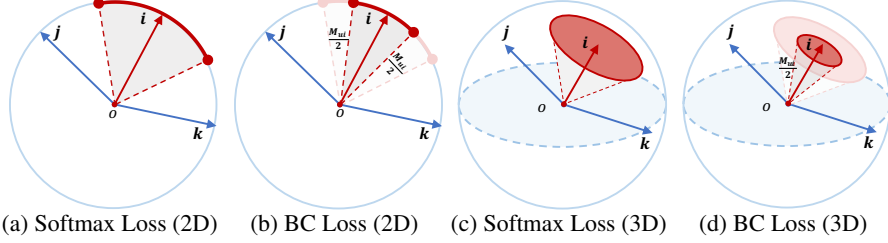


Figure 2: Geometric Interpretation of softmax loss and BC loss in 2D and 3D hypersphere. The dark red region indicates the discriminative user constraint, while the light red region is for comparison.

145 high value of $\hat{\xi}_{ui}$ and leads to a high value of M_{ui} ; henceforward, BC loss imposes the large angular
 146 margin M_{ui} between the negative item j and positive item i and optimizes the representations of user
 147 u and item i to lower $\hat{\xi}_{ui}$. See more details and analyses in Section 4.

148 It is noted that BC loss is extremely easy to implement in recommendation tasks, which only needs to
 149 revise several lines of code. Moreover, compared with softmax loss, BC loss only adds negligible
 150 computational complexity during training but achieves more discriminative representations. Hence,
 151 we recommend to use BC loss not only as a debiasing strategy to alleviate the popularity bias, but
 152 also as a standard loss in recommender models to enhance the discriminative power. Note that the
 153 modeling of M_{ui} is worth exploring, such as the more complex version $M_{ui} = \min\{\lambda \cdot \hat{\xi}_{ui}, \pi - \hat{\theta}_{ui}\}$
 154 where λ controls the strength of the bias-margin. Meanwhile, carefully designing a monotonically
 155 decreasing function helps to get rid of the upper bound restriction. We will leave the exploration of
 156 bias-margin in future work.

157 4 Analyses of BC Loss

158 We analyze desirable characteristics of BC loss. Specifically, we start by presenting its geometric inter-
 159 pretation, and then show its theoretical properties *w.r.t.* compactness and dispersion of representations.
 160 The hard mining mechanism of BC loss is discussed in Appendix A.2.

161 4.1 Geometric Interpretation

162 Here we probe into the ranking criteria of softmax loss and BC loss, from the geometric perspective.
 163 To simplify the geometric interpretation, we analyze one user u with one observed item i and
 164 only two unobserved items j and k . Then the posterior probabilities obtained by softmax loss are:
 165
$$\frac{\exp(\cos(\hat{\theta}_{ui})/\tau)}{\exp(\cos(\hat{\theta}_{ui})/\tau) + \exp(\cos(\hat{\theta}_{uj})/\tau) + \exp(\cos(\hat{\theta}_{uk})/\tau)}$$
. During training, softmax loss encourages the ranking
 166 criteria $\hat{\theta}_{ui} < \hat{\theta}_{uj}$ and $\hat{\theta}_{ui} < \hat{\theta}_{uk}$ to model the basic assumption that the observed interaction (u, i)
 167 indicates more positive cues of user preference than the unobserved interactions (u, j) and (u, k) .

168 Intuitively, to make the ranking criteria more stringent, we can impose an angular margin M_{ui} on it
 169 and establish a new criteria $\hat{\theta}_{ui} + M_{ui} < \hat{\theta}_{uj}$ and $\hat{\theta}_{ui} + M_{ui} < \hat{\theta}_{uk}$. Directly formulating this idea
 170 arrives at the posterior probabilities of BC loss:
$$\frac{\exp(\cos(\hat{\theta}_{ui} + M_{ui})/\tau)}{\exp(\cos(\hat{\theta}_{ui} + M_{ui})/\tau) + \exp(\cos(\hat{\theta}_{uj})/\tau) + \exp(\cos(\hat{\theta}_{uk})/\tau)}$$
.
 171 Obviously, BC loss is more rigorous about the ranking assumption compared with softmax loss. See
 172 Appendix A.2 for more detailed explanations.

173 We then depict the geometric interpretation and comparison of softmax loss and BC loss in Figure 2.
 174 Assume the learned representations of i, j , and k are given, and softmax and BC losses are optimized
 175 to the same value. In softmax loss, the constraint boundaries for correctly ranking user u 's preference
 176 are $\hat{\theta}_{ui} = \hat{\theta}_{uj}$ and $\hat{\theta}_{ui} = \hat{\theta}_{uk}$; whereas, in BC loss, the constraint boundaries are $\hat{\theta}_{ui} + M_{ui} = \hat{\theta}_{uj}$
 177 and $\hat{\theta}_{ui} + M_{ui} = \hat{\theta}_{uk}$. Geometrically, from softmax loss (*cf.* Figure 2c) to BC loss (*cf.* Figure 2d), it
 178 is a more stringent circle-like region on the unit sphere in the 3D case. Further enlarging the margin
 179 M_{ui} will lead to a smaller hypercycle-like region, which is an explicit discriminative constraint on a
 180 manifold. As a result, limited constraint regions squeeze the tightness of similar items and encourages
 181 the separation of dissimilar items. Moreover, with the increase of representation dimension, BC loss
 182 has more restricted learning requirements, exponentially decreasing the area of constraint regions for
 183 correct ranking, and becomes progressively powerful to learn discriminative representations.

Table 1: Overall debiasing performance comparison in balanced and imbalanced test sets.

	Tencent				Amazon-book				Alibaba-iFashion			
	Balanced		Imbalanced		Balanced		Imbalanced		Balanced		Imbalanced	
	Recall	NDCG	Recall	NDCG	Recall	NDCG	Recall	NDCG	Recall	NDCG	Recall	NDCG
MostPop	0.0002	0.0002	0.0384	0.0208	0.0001	0.0001	0.0102	0.0063	0.0003	0.0001	0.0212	0.0084
MF	0.0052	0.0040	0.0982	0.0643	0.0109	0.0103	0.0856	0.0638	0.0056	0.0028	0.0843	0.0411
+ IPS-CN	0.0075	0.0058	0.0686	0.0421	0.0132	0.0123	0.0765	0.0554	0.0050	0.0027	0.0551	0.0255
+ CausE	0.0056	0.0043	0.0687	0.0468	0.0115	0.0105	0.0720	0.0551	0.0005	0.0003	0.0185	0.0086
+ sam+reg	0.0070	0.0054	0.0406	0.0266	0.0141	0.0132	0.0599	0.0443	0.0067	0.0032	0.0305	0.0146
+ MACR	0.0067	0.0046	0.0326	0.0241	0.0181	0.0146	0.0292	0.0229	0.0086	0.0041	0.0650	0.0331
+ BC Loss	0.0087*	0.0068*	0.1298*	0.0904*	0.0221*	0.0202*	0.1198*	0.0948*	0.0095*	0.0048*	0.0967*	0.0487*
Imp. %	16.0%	17.2%	32.2%	40.1%	22.1%	38.4%	40.0%	49.6%	10.5%	17.1%	14.7%	18.5%
LightGCN	0.0055	0.0042	0.1065	0.0712	0.0123	0.0116	0.0941	0.0724	0.0036	0.0017	0.0660	0.0322
+ IPS-CN	0.0072	0.0054	0.0900	0.0599	0.0148	0.0136	0.0836	0.0639	0.0038	0.0017	0.0658	0.0317
+ CausE	0.0055	0.0040	0.0966	0.0665	0.0134	0.0121	0.0926	0.0717	0.0029	0.0013	0.0449	0.0221
+ sam+reg	0.0076	0.0056	0.0653	0.0436	0.0157	0.0149	0.0773	0.0600	0.0056	0.0027	0.0502	0.0252
+ MACR	0.0075	0.0050	0.0731	0.0532	0.0183	0.0153	0.0767	0.0600	0.0033	0.0015	0.0475	0.0238
+ BC Loss	0.0095*	0.0073*	0.1194*	0.0832*	0.0257*	0.0227*	0.1123*	0.0903*	0.0077*	0.0037*	0.0992*	0.0510*
Imp. %	25.0%	30.1%	12.1%	16.9%	40.4%	48.4%	19.3%	24.7%	37.5%	37.0%	50.3%	58.4%

184 4.2 Theoretical Properties

185 BC loss improves head and tail representation learning by enforcing the compactness of matched
 186 users and items, while imposing the dispersion of unmatched users and items. See detailed proof in
 187 Appendix A.3.

188 **Theorem 1.** Let $\mathbf{v}_u \doteq \psi(u)$, $\mathbf{v}_i \doteq \phi(i)$, and $\mathbf{c}_u = \frac{1}{|\mathcal{P}_u|} \sum_{i \in \mathcal{P}_u} \mathbf{v}_i$, $\mathbf{c}_i = \frac{1}{|\mathcal{P}_i|} \sum_{u \in \mathcal{P}_i} \mathbf{v}_u$, where
 189 $\mathcal{P}_u = \{i | y_{ui} = 1\}$ and $\mathcal{N}_u = \{i | y_{ui} = 0\}$ are the sets of user u 's positive and negative items,
 190 respectively; $\mathcal{P}_i = \{u | y_{ui} = 1\}$ is the set of item i 's positive users. Assuming the representations of
 191 users and items are normalized, the minimization of BC loss is equivalent to minimizing a compactness
 192 part and a dispersion part simultaneously:

$$\mathcal{L}_{BC} \geq \underbrace{\sum_{u \in \mathcal{U}} \|\mathbf{v}_u - \mathbf{c}_u\|^2}_{\text{Compactness part}} + \underbrace{\sum_{i \in \mathcal{I}} \|\mathbf{v}_i - \mathbf{c}_i\|^2 - \sum_{u \in \mathcal{U}} \sum_{j \in \mathcal{N}_u} \|\mathbf{v}_u - \mathbf{v}_j\|^2}_{\text{Dispersion part}} \propto \underbrace{H(\mathbf{V}|Y)}_{\text{Compactness}} - \underbrace{H(\mathbf{V})}_{\text{Dispersion}}. \quad (7)$$

193 **Discussion.** \mathbf{c}_u is the averaged representations of all items that u has interacted with, which describes
 194 u 's interest; similarly, \mathbf{c}_i profiles its user group. For the compactness part, BC loss forces the user's
 195 positive items to be user-centric and vice versa. From the entropy perspective, compactness part tends
 196 to learn a low-entropy cluster for positive interactions, *i.e.*, high compactness for similar users and
 197 items. For the dispersion part, for users and items from unobserved interactions, BC loss maximizes
 198 the pairwise euclidean distance between their representations and encourages them to be distant from
 199 each other; Hence, from the entropy viewpoint, dispersion part levers the spread of representations to
 200 learn a high-entropy representation space, *i.e.*, large separation degree for dissimilar users and items.

201 5 Experiments

202 We aim to answer the following research questions:

- 203 • **RQ1:** How does BC Loss perform compared with debiasing strategies in various evaluations?
- 204 • **RQ2:** Does BC loss cause the trade-off between head and tail performance?
- 205 • **RQ3:** What are the impacts of the components (*e.g.*, temperature, margin) on BC Loss?

206 **Baselines & Datasets.** SOTA debiasing strategies in various research lines are compared: sample
 207 re-weighting (IPS-CN [15]), bias removal by causal inference (MACR [22], CausE [19]), and
 208 regularization-based framework (sam+reg [12]). Extensive experiments are conducted on six real-
 209 world benchmark datasets: Tencent [47], Amazon-Book [48], Alibaba-iFashion [49], Yelp2018 [25],
 210 Douban Movie [50], and KuaiRec [51]. For comprehensive comparisons, almost all standard test
 211 distributions in CF are covered in the experiments: balanced test set [22, 23, 24], randomly selected
 212 imbalanced test set [10, 52], temporal split test set [20, 21, 12], and unbiased test set [13, 51, 53].
 213 See more experiments on KuaiRec [51] for unbiased test evaluation in Appendix B.3.

214 5.1 Performance Comparison (RQ1)

215 5.1.1 Evaluations on Imbalanced and Balanced Test Sets.

216 **Motivation.** Many prevalent debiasing methods assume that test distribution is known in advance
217 [22, 23, 10], *i.e.*, the validation set has similar distribution with the test set. Moreover, only an
218 imbalanced or balanced test set is evaluated. However, in real-world applications, the test distributions
219 are usually unavailable and can even reverse the prior in the training distribution. We conjecture that
220 a good debiasing recommender is required to perform well on both imbalanced and balanced test
221 distributions. In our settings, no information about the balanced test is provided in advance.

222 **Data Splits.** The models are identical across both imbalanced and balanced evaluations. The test
223 distribution in the balanced evaluation is uniform, *i.e.*, randomly sample 15% of interactions with
224 equal probability *w.r.t.* items. Besides, the test splits for the imbalanced test are similarly long-
225 tailed like the train and validation sets, *i.e.*, randomly split the remaining interactions into training,
226 validation, and imbalanced test sets (60% : 10% : 15%).

227 **Results.** Table 1 reports the comparison of performance in imbalanced and balanced test evaluations.
228 The best performing methods are bold and starred, while the strongest baselines are underlined;
229 Imp.% measures the relative improvements of BC loss over the strongest baselines. We observe that:

230 • **BC loss significantly outperforms the state-of-the-art baselines in both balanced and imbal-**
231 **anced evaluations across all datasets.** In particular, it achieves consistent improvements over the
232 best debiasing baselines and original CF models by 12.1% ~ 58.4%. This clearly demonstrates
233 that BC loss not only effectively alleviates the amplification of popularity bias but also improves
234 the discriminative power of representations. Moreover, Table 5 shows the computational costs of
235 all methods. Compared to the backbone models, BC loss only adds negligible time complexity.

236 • **Debiasing baselines sacrifice the imbalanced performance and perform inconsistently across**
237 **datasets.** Debiasing strategies generally achieve higher balanced performance at the expense of a
238 large imbalanced performance drop. Specifically, the strongest baselines over all imbalanced test
239 sets are the original CF models. Worse still, as the degree of data sparsity increases, some debiasing
240 methods fail to quantify the popularity bias and limit their bias removal ability. For example, in the
241 sparsest Alibaba-iFashion dataset, the results of MF+IPS-CN, MF+CausE, LightGCN+MACR, and
242 LightGCN+CausE on the balanced evaluation are lower than original CF models (MF or Light-
243 GCN). In contrast, benefiting from popularity bias-aware margin, BC loss can learn discriminative
244 representations that accomplish more profound user and item understanding, leading to higher head
245 and tail recommendation quality.

246 5.1.2 Evaluations on Temporal Split Test Set

247 **Motivation.**

248 In real applications, popularity bias
249 dynamically changes over time. Here
250 we consider temporal split test eval-
251 uation on Douban Movie where the
252 historical interactions are sliced into
253 the training, validation, and test sets
254 (70%:10%:20%) according to the
255 timestamps.

Table 2: The performance comparison on Douban dataset.

	MF			LightGCN		
	HR	Recall	NDCG	HR	Recall	NDCG
Backbone	<u>0.2924</u>	<u>0.0294</u>	<u>0.0472</u>	0.3543	0.0313	0.0602
+ IPS-CN	0.2514	0.0174	0.0324	0.3212	0.0261	0.0502
+ CausE	0.2725	0.0203	0.0376	0.3403	0.0275	0.0514
+ sam+reg	0.2826	0.0191	0.0390	0.2944	0.0252	0.0488
+ MACR	0.1084	0.0087	0.0163	0.3127	0.0271	0.0519
+ BC loss	0.3742*	0.0324*	0.0601*	0.3562*	0.0346*	0.0652*
Imp. %	28.0%	10.2%	27.3%	0.5%	10.4%	8.3%

256 **Results.** As Table 2 shows, BC loss is steadily superior to all baselines *w.r.t.* all metrics on Douban
257 Movie. For instance, it achieves significant improvements over the MF and LightGCN backbones
258 *w.r.t.* Recall@20 by 10.2% and 10.4%, respectively. This validates that BC loss endows the backbone
259 models with better robustness against the popularity distribution shift and alleviates the negative in-
260 fluence of popularity bias. Surprisingly, none of the debiasing baselines could maintain a comparable
261 performance to the backbones. We ascribe the failure to their preconceived idea of tail items, which
262 possibly change over time.

Table 3: The performance evaluations of head, mid, and tail on Tencent dataset.

	Balanced NDCG@20				Imbalanced NDCG@20			
	Tail	Mid	Head	Overall	Tail	Mid	Head	Overall
MF	0.00004	0.00097	0.01250	0.00402	0.00021	0.00197	0.06837	0.06431
+ IPS-CN	0.00009 ^{+125%}	0.00212 ^{+119%}	0.01684 ^{+35%}	0.00575 ^{+43%}	0.00056 ^{+167%}	0.00401 ^{+104%}	0.04439 ^{-35%}	0.04205 ^{-35%}
+ CausE	0.00008 ^{+100%}	0.00149 ^{+54%}	0.01168 ^{-7%}	0.00430 ^{+7%}	0.00038 ^{+81%}	0.00253 ^{+28%}	0.04876 ^{-29%}	0.04680 ^{-27%}
+ sam-reg	0.00006 ^{+50%}	0.00135 ^{+39%}	0.01573 ^{+26%}	0.00535 ^{+33%}	0.00011 ^{-48%}	0.00281 ^{+43%}	0.02850 ^{-58%}	0.02661 ^{-59%}
+ MACR	0.00188^{+4600%}	0.00521^{+437%}	0.00555 ^{-56%}	0.00456 ^{+13%}	0.00370^{+1662%}	0.00615 ^{+212%}	0.02748 ^{-60%}	0.02413 ^{-62%}
+ BC loss	0.00024 ^{+500%}	0.00355 ^{+266%}	0.01831^{+46%}	0.00680^{+69%}	0.00142 ^{+576%}	0.00712^{+261%}	0.09552^{+40%}	0.09040^{+41%}
LightGCN	0.00025	0.00193	0.01136	0.00417	0.00094	0.00391	0.07561	0.07121
+ IPS-CN	0.00140 ^{+460%}	0.00241 ^{+25%}	0.01560 ^{+37%}	0.00544 ^{+30%}	0.00109 ^{+16%}	0.00522 ^{+34%}	0.06333 ^{-16%}	0.05993 ^{-16%}
+ CausE	0.00006 ^{-76%}	0.00138 ^{-29%}	0.01177 ^{+4%}	0.00403 ^{-3%}	0.00040 ^{-57%}	0.00279 ^{-29%}	0.06996 ^{-7%}	0.06650 ^{-7%}
+ sam-reg	0.00006 ^{-76%}	0.00120 ^{-38%}	0.01727 ^{+52%}	0.00560 ^{+34%}	0.00024 ^{-74%}	0.00253 ^{-35%}	0.04647 ^{-39%}	0.04355 ^{-39%}
+ MACR	0.00287^{+1048%}	0.00461^{+139%}	0.00454 ^{-60%}	0.00501 ^{+20%}	0.00389^{+313%}	0.00635^{+62%}	0.04058 ^{-46%}	0.05323 ^{-25%}
+ BC loss	0.00057 ^{+128%}	0.00321 ^{+66%}	0.01943^{+71%}	0.00730^{+75%}	0.00125 ^{+33%}	0.00516 ^{+32%}	0.08823^{+17%}	0.08320^{+17%}

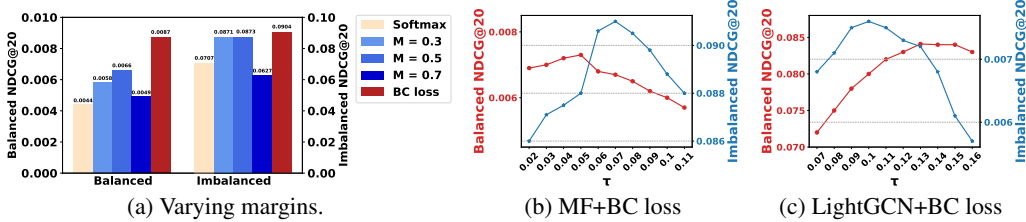


Figure 3: (a) Comparisons with a varying margin; (b-c) Temperature τ sensitivity analysis on Tencent.

263 5.2 Head, Mid, & Tail Performance (RQ2)

264 **Motivation.** To further evaluate whether BC loss lifts the tail performance by inevitably sacrificing
 265 the head performance, we divide the test set of Tencent into three subgroups, according to the
 266 interaction number of each item: head (popular items that are in the top third), mid (normal items
 267 in the middle), and tail (unpopular items in the bottom third). Most previous studies focus on
 268 average NDCG@20 for evaluation, especially balanced test evaluations [22, 23]. However, average
 269 metrics could be insufficient to reflect the performance of each subgroup. A trivial solution to
 270 achieve high performance is promoting the rankings of low-popularity items in the recommendations.
 271 In this case, only the average metrics are not reliable on the balanced test. Therefore, we report
 272 the performance of individual subgroups on both balanced and imbalanced test sets for a more
 273 comprehensive comparison.

274 **Results.** Table 3 shows the evaluations of the head, mid, and tail subgroups. The red and blue
 275 numbers in percentage separately refer to the improvement and decline of each method relative to the
 276 original CF model (MF or LightGCN). We find that:

- 277 • **BC loss is the only method that consistently yields remarkable improvements in every sub-**
 278 **group.** With a closer look at the head evaluation, BC loss shows its ability to learn more discrimi-
 279 native representations for popular items across imbalanced and balanced settings. In particular, it
 280 achieves significant improvements over MF and LightGCN *w.r.t.* head NDCG by 40% and 17% in
 281 the imbalanced test evaluation, respectively. We attribute improvements to the usage of bias-aware
 282 margin, which boosts the recommendation quality for the tail and head items.
- 283 • **As the performance comparison among subgraphs in the imbalanced scenario shows, the**
 284 **baselines enhance the tail performance but sacrifice the head performance.** Specifically,
 285 these baselines hardly maintain the head performance and show a clear trade-off trend between
 286 the head and tail performance. Taking MACR as an example, although the great improvement
 287 (+1662%) over MF is achieved in the tail subgraph, it brings in the dramatic drop (-62%) in the
 288 head subgraph, which lowers the overall performance by a big drop (-60%). Here we ascribe the
 289 trade-off to blindly promoting the rankings of tail items for matched and unmatched, rather than
 290 improving the discriminative power of representations.

291 5.3 Study on BC Loss (RQ3)

292 **Effect of Bias-aware Margin.** Figure 3a displays the performance on balanced and imbalanced test
 293 sets on Tencent among softmax loss, BC loss with constant margin M [40], and BC loss with adaptive
 294 bias-aware margin. BC loss achieves the best performance, illustrating that bias-aware margin indeed
 295 is effective at reducing popularity bias and learning high-quality representations.

296 **Effect of Temperature** τ . BC loss has one hyperparameter to tune — temperature τ in Equation
297 (5). In Figure 3b and 3c, both balanced and imbalanced evaluations exhibit the concave unimodal
298 functions of τ , where the curves reach the peak almost synchronously in a small range of τ . For
299 example, MF+BC loss gets the best performance when $\tau = 0.05$ and $\tau = 0.07$ in balanced and
300 imbalanced settings, respectively; We observe similar trends on other datasets and skip them due to
301 the space limit. This justifies that BC loss does not suffer from the trade-off between the balanced and
302 imbalanced evaluations and improves the generalization without sacrificing the head performance.

303 6 Related Work

304 Prevalent popularity debiasing strategies in CF roughly fall into four research lines.

305 **Post-processing re-ranking methods** [5, 6, 7, 8, 9] are applied to the output of the recommender
306 system without changing the representations of users and items. The purposes of modifying the
307 ranking of models can be various: Calibration [5] ensures that the past interests proportions of
308 users are expected to maintain at the same level; RankALS [6] aims to increase the diversification
309 of recommendation; FPC [8] investigates the popularity bias in the dynamic recommendation by
310 rescaling the predicted scores.

311 **Regularization-based frameworks** [10, 11, 12, 9] explore the use of regularization that provides a
312 tunable mechanism for controlling the trade-off between recommendation accuracy and coverage.
313 The difference among these methods is the design of penalty terms: ALS+Reg [11] defines intra-list
314 distance as the penalty to achieve the fair recommendation; ESAM [10] introduces the attribute
315 correlation alignment, center-clustering, and self-training regularization to learn good feature rep-
316 resentations; sam-reg [12] regularizes the biased correlation between user-item relevance and item
317 popularity; Reg [9] decouples the item popularity with the model preference predictions.

318 **Sample re-weighting methods** [13, 14, 15, 16, 17, 18], also known as Inverse Propensity Score (IPS)
319 view the item popularity in the training set as the propensity score and exploit its inverse to re-weight
320 loss of each instance. To address the high variance of re-weighted loss, many of them [15, 14] further
321 employ normalization or smoothing penalty to attain a more stable output. However, the unreliability
322 of methods is due to their measurement of the propensity score, leveraging the item frequency but
323 failing to consider interaction-wise popularity bias.

324 **Bias removal by causal inference methods** [19, 24, 23, 20, 21, 22], getting inspiration from the
325 recent success of counterfactual inference, specify the role of popularity bias in assumed causal
326 graphs and mitigate the bias effect on the prediction. However, the causal structure is heuristically
327 assumed based on the author’s understanding, without any theoretical guarantee.

328 BC loss opens up a possibility of conventional debiasing methods in CF that mitigate the popularity
329 bias by enhancing the discriminative power. Recent studies, boosting the discriminative feature
330 spaces by modified softmax loss are mainly discussed in face recognition, where a constant margin is
331 added [40] to better classify. We transfer it in CF and compare it with BC loss in Figure 3a.

332 7 Conclusion

333 Despite the great success in collaborative filtering, today’s popularity debiasing methods are still
334 far from being able to improve the recommendation quality. In this work, we proposed a simple yet
335 effective BC loss, utilizing popularity bias-aware margin to eliminate the popularity bias. Grounded
336 by theoretical proof, clear geometric interpretation and real-world visualization study, BC loss boosts
337 the head and tail performance by learning a more discriminative representation space. Extensive
338 experiments verify that the remarkable improvement in head and tail evaluations on various test sets
339 indeed comes from the better representation rather than simply catering to the tail.

340 The limitations of BC loss are in three respects, which will be addressed in future work: 1) the model-
341 ing of bias-aware margin is worth exploring, which could significantly influence the performance of
342 BC loss, 2) multiple important biases, such as exposure and selection bias, are not considered, and 3)
343 more experiments comparing BC loss to standard CF losses (*e.g.*, cross-entropy, WARP) are needed
344 to further demonstrate the power of BC loss in regular recommendation tasks. We believe that this
345 work provides a potential research direction to diagnose the debiasing of long-tail ranking and will
346 inspire more works.

347 **References**

- 348 [1] Himan Abdollahpouri, Masoud Mansoury, Robin Burke, and Bamshad Mobasher. The un-
349 fairness of popularity bias in recommendation. In *RMSE@RecSys*, volume 2440 of *CEUR*
350 *Workshop Proceedings*, 2019.
- 351 [2] Rocío Cañamares and Pablo Castells. Should I follow the crowd?: A probabilistic analysis of
352 the effectiveness of popularity in recommender systems. In *SIGIR*, pages 415–424, 2018.
- 353 [3] Yitong Ji, Aixin Sun, Jie Zhang, and Chenliang Li. A re-visit of the popularity baseline in
354 recommender systems. In *SIGIR*, pages 1749–1752, 2020.
- 355 [4] Sirui Yao and Bert Huang. Beyond parity: Fairness objectives for collaborative filtering. In
356 *NIPS*, pages 2921–2930, 2017.
- 357 [5] Harald Steck. Calibrated recommendations. In *RecSys*, pages 154–162. ACM, 2018.
- 358 [6] Himan Abdollahpouri, Robin Burke, and Bamshad Mobasher. Managing popularity bias in
359 recommender systems with personalized re-ranking. In *FLAIRS Conference*, pages 413–418,
360 2019.
- 361 [7] Harald Steck. Collaborative filtering via high-dimensional regression. *CoRR*, abs/1904.13033,
362 2019.
- 363 [8] Ziwei Zhu, Yun He, Xing Zhao, and James Caverlee. Popularity bias in dynamic recommenda-
364 tion. In *KDD*, pages 2439–2449, 2021.
- 365 [9] Ziwei Zhu, Yun He, Xing Zhao, Yin Zhang, Jianling Wang, and James Caverlee. Popularity-
366 opportunity bias in collaborative filtering. In *WSDM*, pages 85–93. ACM, 2021.
- 367 [10] Zhihong Chen, Rong Xiao, Chenliang Li, Gangfeng Ye, Haochuan Sun, and Hongbo Deng.
368 ESAM: discriminative domain adaptation with non-displayed items to improve long-tail perfor-
369 mance. In *SIGIR*, pages 579–588, 2020.
- 370 [11] Himan Abdollahpouri, Robin Burke, and Bamshad Mobasher. Controlling popularity bias in
371 learning-to-rank recommendation. In *RecSys*, pages 42–46, 2017.
- 372 [12] Ludovico Boratto, Gianni Fenu, and Mirko Marras. Connecting user and item perspectives in
373 popularity debiasing for collaborative recommendation. *Inf. Process. Manag.*, 58(1):102387,
374 2021.
- 375 [13] Tobias Schnabel, Adith Swaminathan, Ashudeep Singh, Navin Chandak, and Thorsten Joachims.
376 Recommendations as treatments: Debiasing learning and evaluation. In *ICML*, pages 1670–
377 1679, 2016.
- 378 [14] Léon Bottou, Jonas Peters, Joaquin Quiñonero Candela, Denis Xavier Charles, Max Chickering,
379 Elon Portugaly, Dipankar Ray, Patrice Y. Simard, and Ed Snelson. Counterfactual reasoning and
380 learning systems: the example of computational advertising. *J. Mach. Learn. Res.*, 14(1):3207–
381 3260, 2013.
- 382 [15] Alois Gruson, Praveen Chandar, Christophe Charbuillet, James McInerney, Samantha Hansen,
383 Damien Tardieu, and Ben Carterette. Offline evaluation to make decisions about playlist
384 recommendation algorithms. In *WSDM*, pages 420–428, 2019.
- 385 [16] Thorsten Joachims, Adith Swaminathan, and Tobias Schnabel. Unbiased learning-to-rank with
386 biased feedback. In *IJCAI*, pages 5284–5288. ijcai.org, 2018.
- 387 [17] Longqi Yang, Yin Cui, Yuan Xuan, Chenyang Wang, Serge J. Belongie, and Deborah Estrin.
388 Unbiased offline recommender evaluation for missing-not-at-random implicit feedback. In
389 *RecSys*, pages 279–287. ACM, 2018.
- 390 [18] Yuta Saito, Suguru Yaginuma, Yuta Nishino, Hayato Sakata, and Kazuhide Nakata. Unbiased
391 recommender learning from missing-not-at-random implicit feedback. In *WSDM*, pages 501–
392 509, 2020.

- 393 [19] Stephen Bonner and Flavian Vasile. Causal embeddings for recommendation. In *RecSys*, pages
394 104–112, 2018.
- 395 [20] Yang Zhang, Fuli Feng, Xiangnan He, Tianxin Wei, Chonggang Song, Guohui Ling, and
396 Yongdong Zhang. Causal intervention for leveraging popularity bias in recommendation. In
397 *SIGIR*, pages 11–20, 2021.
- 398 [21] Wenjie Wang, Fuli Feng, Xiangnan He, Xiang Wang, and Tat-Seng Chua. Deconfounded
399 recommendation for alleviating bias amplification. In *KDD*, pages 1717–1725. ACM, 2021.
- 400 [22] Tianxin Wei, Fuli Feng, Jiawei Chen, Ziwei Wu, Jinfeng Yi, and Xiangnan He. Model-agnostic
401 counterfactual reasoning for eliminating popularity bias in recommender system. In *KDD*,
402 pages 1791–1800, 2021.
- 403 [23] Yu Zheng, Chen Gao, Xiang Li, Xiangnan He, Yong Li, and Depeng Jin. Disentangling
404 user interest and conformity for recommendation with causal embedding. In *WWW*, pages
405 2980–2991, 2021.
- 406 [24] Dugang Liu, Pengxiang Cheng, Zhenhua Dong, Xiuqiang He, Weike Pan, and Zhong Ming. A
407 general knowledge distillation framework for counterfactual recommendation via uniform data.
408 In *SIGIR*, pages 831–840, 2020.
- 409 [25] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yong-Dong Zhang, and Meng Wang. Lightgcn:
410 Simplifying and powering graph convolution network for recommendation. In *SIGIR*, pages
411 639–648. ACM, 2020.
- 412 [26] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. BPR:
413 bayesian personalized ranking from implicit feedback. *CoRR*, abs/1205.2618, 2012.
- 414 [27] Yoshua Bengio and Jean-Sébastien Senecal. Quick training of probabilistic neural nets by
415 importance sampling. In *AISTATS*, 2003.
- 416 [28] Steffen Rendle. Item recommendation from implicit feedback. *CoRR*, abs/2101.08769, 2021.
- 417 [29] Yehuda Koren. Factorization meets the neighborhood: a multifaceted collaborative filtering
418 model. In *KDD*, pages 426–434. ACM, 2008.
- 419 [30] Dawen Liang, Rahul G. Krishnan, Matthew D. Hoffman, and Tony Jebara. Variational autoen-
420 coders for collaborative filtering. In *WWW*, pages 689–698. ACM, 2018.
- 421 [31] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural
422 collaborative filtering. In *WWW*, pages 173–182. ACM, 2017.
- 423 [32] Travis Ebesu, Bin Shen, and Yi Fang. Collaborative memory network for recommendation
424 systems. In *SIGIR*, pages 515–524. ACM, 2018.
- 425 [33] Santosh Kabbur, Xia Ning, and George Karypis. FISM: factored item similarity models for
426 top-n recommender systems. In *KDD*, pages 659–667. ACM, 2013.
- 427 [34] Rianne van den Berg, Thomas N. Kipf, and Max Welling. Graph convolutional matrix comple-
428 tion. In *KDD*, 2018.
- 429 [35] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, and Jure
430 Leskovec. Graph convolutional neural networks for web-scale recommender systems. In *KDD*,
431 pages 974–983. ACM, 2018.
- 432 [36] Steffen Rendle, Walid Krichene, Li Zhang, and John R. Anderson. Neural collaborative filtering
433 vs. matrix factorization revisited. In *RecSys*, pages 240–248. ACM, 2020.
- 434 [37] Ying Shan, T. Ryan Hoens, Jian Jiao, Haijing Wang, Dong Yu, and J. C. Mao. Deep crossing:
435 Web-scale modeling without manually crafted combinatorial features. In *KDD*, pages 255–262.
436 ACM, 2016.
- 437 [38] Jason Weston, Samy Bengio, and Nicolas Usunier. WSABIE: scaling up to large vocabulary
438 image annotation. In *IJCAI*, pages 2764–2770. IJCAI/AAAI, 2011.

- 439 [39] Xianfeng Tang, Huaxiu Yao, Yiwei Sun, Yiqi Wang, Jiliang Tang, Charu Aggarwal, Prasenjit Mi-
440 tra, and Suhang Wang. Investigating and mitigating degree-related biases in graph convoltuional
441 networks. In *CIKM*, pages 1435–1444, 2020.
- 442 [40] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular
443 margin loss for deep face recognition. In *CVPR*, pages 4690–4699. Computer Vision Foundation
444 / IEEE, 2019.
- 445 [41] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface:
446 Deep hypersphere embedding for face recognition. In *CVPR*, pages 6738–6746, 2017.
- 447 [42] Chang Zhou, Jianxin Ma, Jianwei Zhang, Jingren Zhou, and Hongxia Yang. Contrastive
448 learning for debiased candidate generation in large-scale recommender systems. In *KDD*, pages
449 3985–3995. ACM, 2021.
- 450 [43] Kun Zhou, Hui Wang, Wayne Xin Zhao, Yutao Zhu, Sirui Wang, Fuzheng Zhang, Zhongyuan
451 Wang, and Ji-Rong Wen. S3-rec: Self-supervised learning for sequential recommendation with
452 mutual information maximization. In *CIKM*, pages 1893–1902. ACM, 2020.
- 453 [44] Defu Lian, Qi Liu, and Enhong Chen. Personalized ranking with importance sampling. In
454 *WWW*, pages 1093–1103, 2020.
- 455 [45] Xinyang Yi, Ji Yang, Lichan Hong, Derek Zhiyuan Cheng, Lukasz Heldt, Aditee Kumthekar,
456 Zhe Zhao, Li Wei, and Ed H. Chi. Sampling-bias-corrected neural modeling for large corpus
457 item recommendations. In *RecSys*, pages 269–277, 2019.
- 458 [46] Feng Wang and Huaping Liu. Understanding the behaviour of contrastive loss. In *CVPR*, pages
459 2495–2504. Computer Vision Foundation / IEEE, 2021.
- 460 [47] Fajie Yuan, Xiangnan He, Haochuan Jiang, Guibing Guo, Jian Xiong, Zhezhaoh Xu, and Yilin
461 Xiong. Future data helps training: Modeling future contexts for session-based recommendation.
462 In *Proceedings of The Web Conference 2020*, pages 303–313, 2020.
- 463 [48] Ruining He and Julian J. McAuley. Ups and downs: Modeling the visual evolution of fashion
464 trends with one-class collaborative filtering. In *WWW*, pages 507–517. ACM, 2016.
- 465 [49] Wen Chen, Pipei Huang, Jiaming Xu, Xin Guo, Cheng Guo, Fei Sun, Chao Li, Andreas
466 Pfadler, Huan Zhao, and Binqiang Zhao. POG: personalized outfit generation for fashion
467 recommendation at alibaba ifashion. In *KDD*, pages 2662–2670, 2019.
- 468 [50] Weiping Song, Zhiping Xiao, Yifan Wang, Laurent Charlin, Ming Zhang, and Jian Tang.
469 Session-based social recommendation via dynamic graph attention networks. In *WSDM*, pages
470 555–563. ACM, 2019.
- 471 [51] Chongming Gao, Shijun Li, Wenqiang Lei, Biao Li, Peng Jiang, Jiawei Chen, Xiangnan He,
472 Jiaxin Mao, and Tat-Seng Chua. KuaiRec: A fully-observed dataset for recommender systems.
473 *CoRR*, abs/2202.10842, 2022.
- 474 [52] Ziwei Zhu, Jianling Wang, and James Caverlee. Measuring and mitigating item under-
475 recommendation bias in personalized ranking systems. In *SIGIR*, pages 449–458. ACM,
476 2020.
- 477 [53] Benjamin M. Marlin and Richard S. Zemel. Collaborative prediction and ranking with non-
478 random missing data. In *RecSys*, pages 5–12, 2009.
- 479 [54] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie.
480 Self-supervised graph learning for recommendation. In *SIGIR*, pages 726–735. ACM, 2021.
- 481 [55] Malik Boudiaf, Jérôme Rony, Imtiaz Masud Ziko, Eric Granger, Marco Pedersoli, Pablo
482 Piantanida, and Ismail Ben Ayed. A unifying mutual information view of metric learning:
483 Cross-entropy vs. pairwise losses. In *ECCV (6)*, volume 12351 of *Lecture Notes in Computer
484 Science*, pages 548–564. Springer, 2020.
- 485 [56] Meihong Wang and Fei Sha. Information theoretical clustering via semidefinite programming.
486 In *AISTATS*, volume 15 of *JMLR Proceedings*, pages 761–769. JMLR.org, 2011.

- 487 [57] Xiang Wang, Xiangnan He, Yixin Cao, Meng Liu, and Tat-Seng Chua. KGAT: knowledge
488 graph attention network for recommendation. In *KDD*, pages 950–958. ACM, 2019.
- 489 [58] Dawen Liang, Laurent Charlin, and David M Blei. Causal inference for recommendation. In
490 *Causation: Foundation to Application, Workshop at UAI. AUAI*, 2016.
- 491 [59] Zaiqiao Meng, Richard McCreadie, Craig Macdonald, and Iadh Ounis. Exploring data splitting
492 strategies for the evaluation of recommendation models. In *RecSys*, pages 681–686. ACM,
493 2020.
- 494 [60] Walid Krichene and Steffen Rendle. On sampled metrics for item recommendation. In *KDD*,
495 pages 1748–1757, 2020.
- 496 [61] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*
497 (*Poster*), 2015.
- 498 [62] Kelong Mao, Jieming Zhu, Jinpeng Wang, Quanyu Dai, Zhenhua Dong, Xi Xiao, and Xiuqiang
499 He. Simplex: A simple and strong baseline for collaborative filtering. In *CIKM*, 2021.
- 500 [63] Jiancan Wu, Xiang Wang, Xingyu Gao, Jiawei Chen, Hongcheng Fu, Tianyu Qiu, and Xiangnan
501 He. On the effectiveness of sampled softmax loss for item recommendation. *CoRR*, 2022.

502 **Checklist**

- 503 1. For all authors...
- 504 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
505 contributions and scope? [Yes]
- 506 (b) Did you describe the limitations of your work? [Yes] See Section 7.
- 507 (c) Did you discuss any potential negative societal impacts of your work? [No] This paper
508 proposes a novel debiasing algorithm for recommendation system, which does not have
509 any negative societal impacts.
- 510 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
511 them? [Yes]
- 512 2. If you are including theoretical results...
- 513 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 514 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix A.3.
- 515 3. If you ran experiments...
- 516 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
517 imental results (either in the supplemental material or as a URL)? [Yes] We include
518 code by URL in abstract.
- 519 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
520 were chosen)? [Yes] See Appendix B.1.
- 521 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
522 ments multiple times)? [Yes]
- 523 (d) Did you include the total amount of compute and the type of resources used (e.g., type
524 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix B.1.
- 525 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 526 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 527 (b) Did you mention the license of the assets? [Yes]
- 528 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 529 (d) Did you discuss whether and how consent was obtained from people whose data you’re
530 using/curating? [N/A]
- 531 (e) Did you discuss whether the data you are using/curating contains personally identifiable
532 information or offensive content? [N/A]
- 533 5. If you used crowdsourcing or conducted research with human subjects...
- 534 (a) Did you include the full text of instructions given to participants and screenshots, if
535 applicable? [N/A]
- 536 (b) Did you describe any potential participant risks, with links to Institutional Review
537 Board (IRB) approvals, if applicable? [N/A]
- 538 (c) Did you include the estimated hourly wage paid to participants and the total amount
539 spent on participant compensation? [N/A]