



Double Correction Framework for Denoising Recommendation

Zhuangzhuang He^{*†}
School of Computer Science and
Information Engineering, Hefei
University of Technology
hyicheng223@gmail.com

Peijie Sun
Department of Computer Science and
Technology, Tsinghua University
sun.hfut@gmail.com

Jinqi Gong
Department of Mathematics,
University of Macau
eggmangong@gmail.com

Yifan Wang[†]
Department of Computer Science and
Technology, Tsinghua University
yf-wang21@mails.tsinghua.edu.cn

Le Wu[‡]
School of Computer Science and
Information Engineering, Hefei
University of Technology
Institute of Dataspace,
Hefei Comprehensive National
Science Center
lewu.ustc@gmail.com

Richang Hong
School of Computer Science and
Information Engineering, Hefei
University of Technology
hongrc.hfut@gmail.com

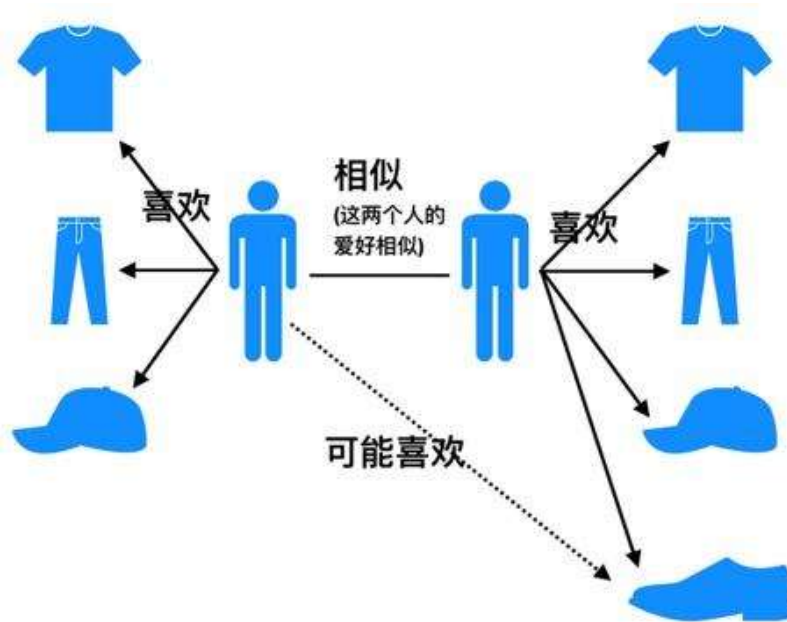
Yonghui Yang
School of Computer Science and
Information Engineering, Hefei
University of Technology
yyh.hfut@gmail.com

Haoyue Bai
School of Computer Science and
Information Engineering, Hefei
University of Technology
baihaoyue621@gmail.com

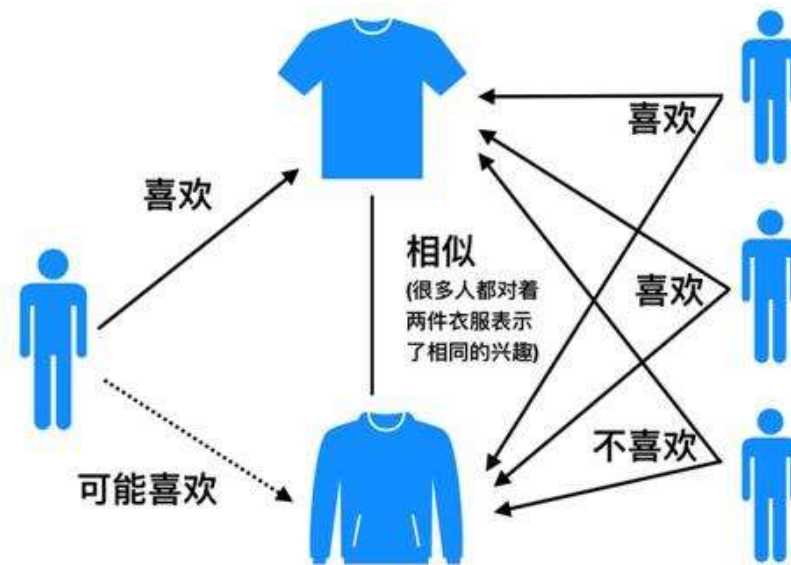
Min Zhang[‡]
Department of Computer Science and
Technology, Tsinghua University
z-m@tsinghua.edu.cn

KDD 2024

Introduction

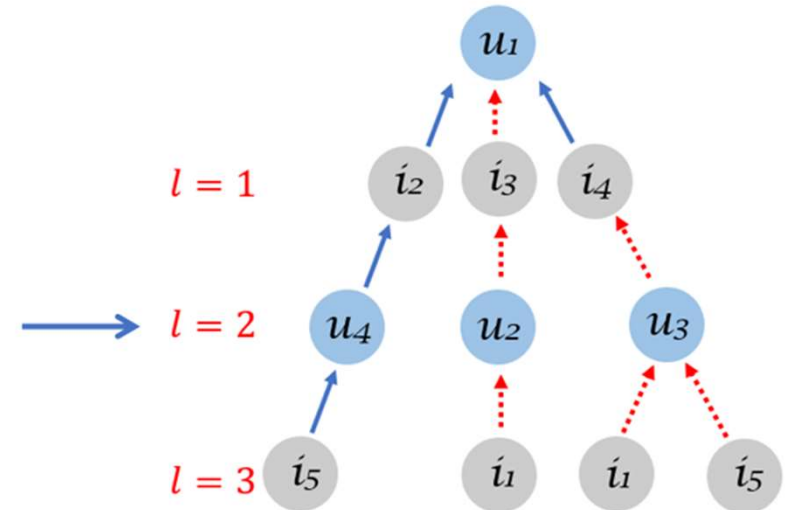
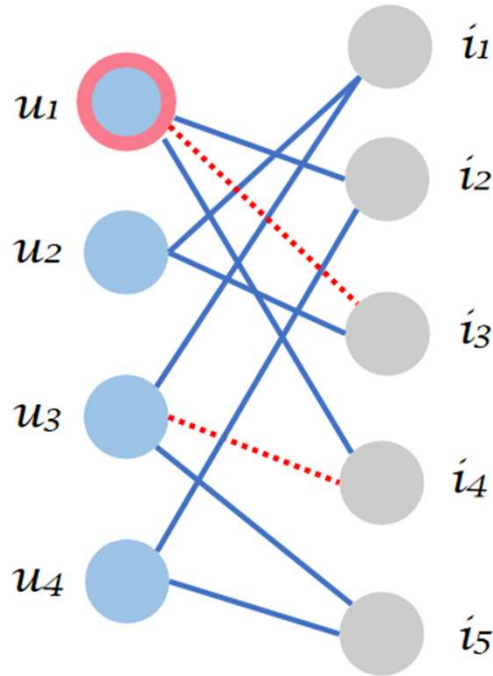


“人以群分”的基于用户的协同过滤



“物以类聚”的基于物品的协同过滤

Introduction



噪声来源:

1. 用户的误触, 错误点击
2. 用户代为购买
3. 恶意生成无效数据的攻击

GNN中的噪声传递:

噪声信息顺着邻域不断传递, 对推荐的结果造成很大影响

Introduction



Truncated Loss: 基于一个动态调整的阈值，将loss比较高的样本loss直接置0，使得这些样本在当前训练轮不参与模型更新。

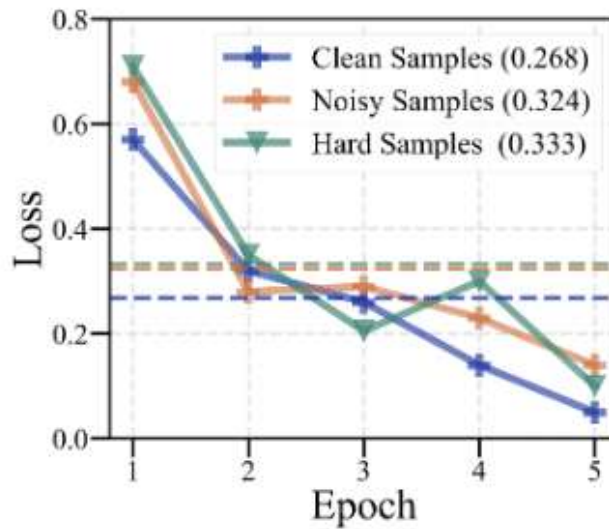
$$\mathcal{L}_{T-CE}(u, i) = \begin{cases} 0, & \mathcal{L}_{CE}(u, i) > \tau \wedge \bar{y}_{ui} = 1 \\ \mathcal{L}_{CE}(u, i), & \text{otherwise,} \end{cases}$$

$$\epsilon(T) = \min(\alpha T, \epsilon_{max})$$

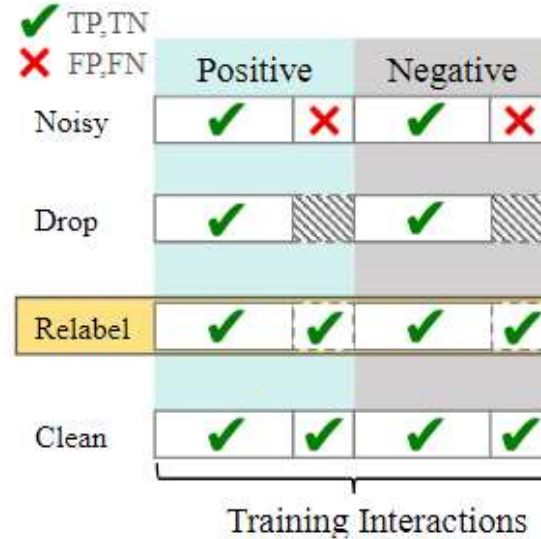
Problem



- (1) 忽视了损失值与噪声之间可能并非高度相关
- (2) 简单地丢弃样本可能导致数据更稀疏



(a) Illustration of unstable losses (dotted line represents the mean of losses).



(b) Different strategies in their ideal cases.

Problem

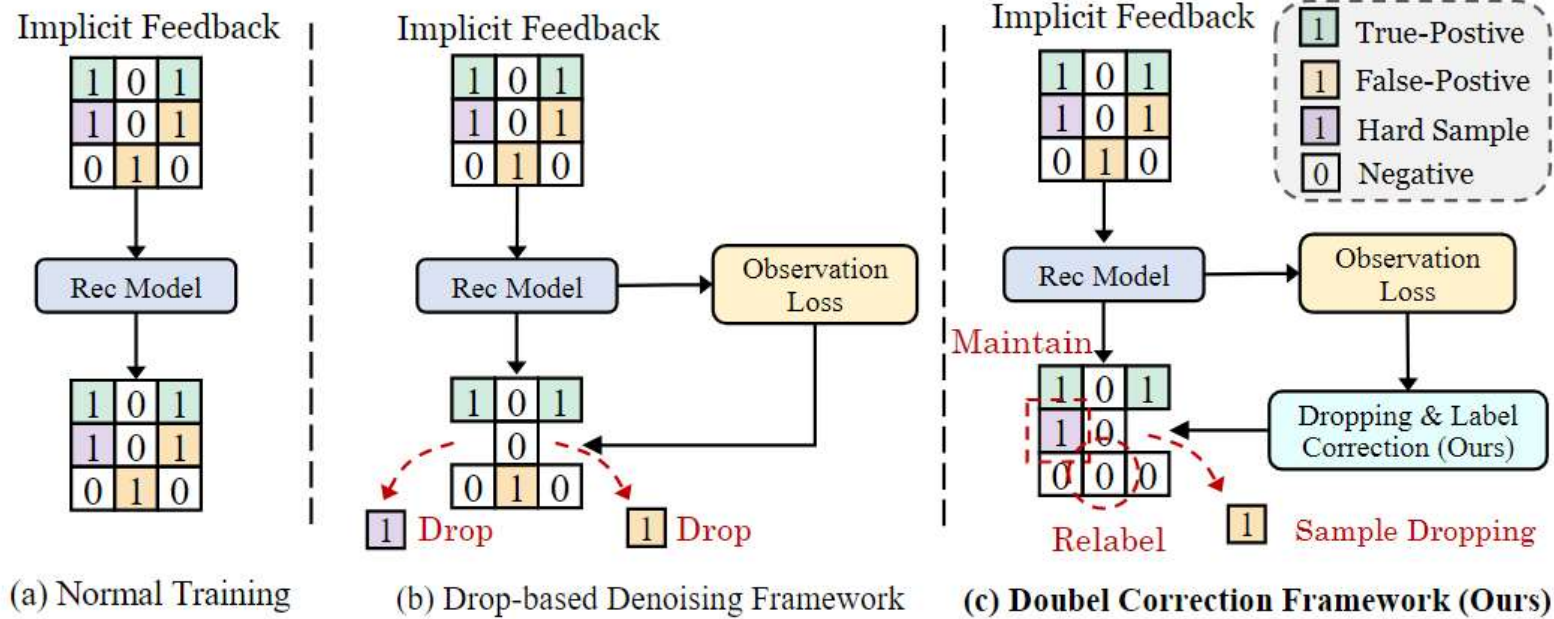
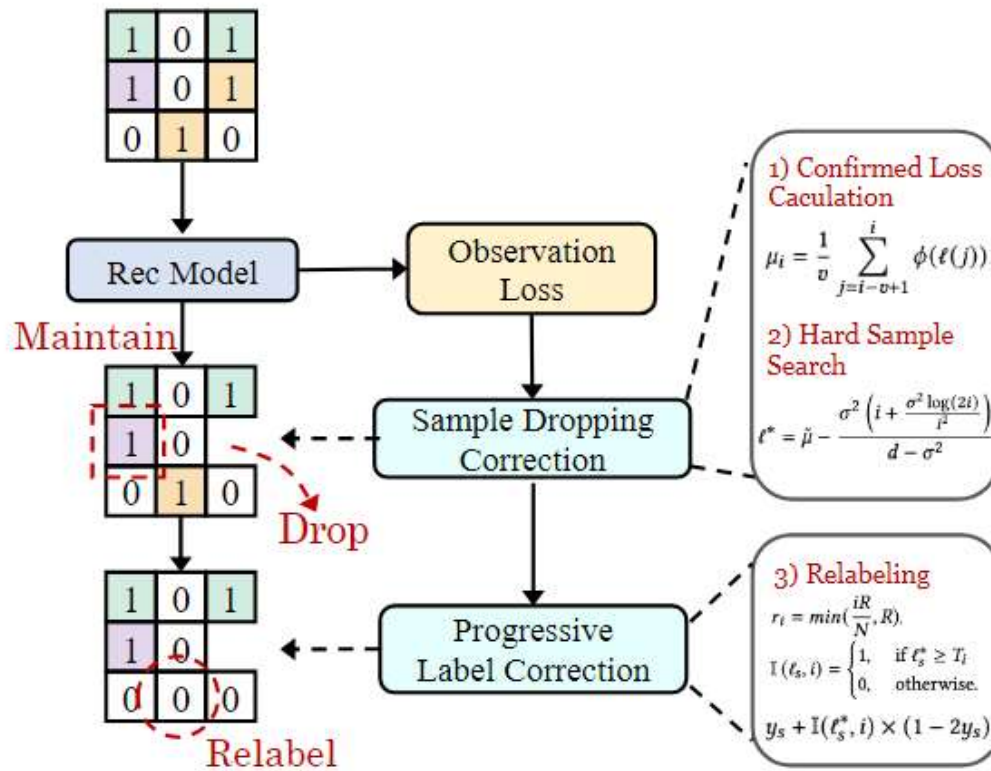


Figure 2: Illustration comparison: (a) Normal training model without denoising, (b) Denoising model with drop strategy, (c) Double correction framework for denoising recommendation (Ours).

Method



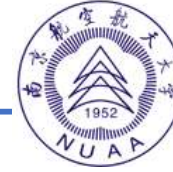
DCF

1. Sample Dropping Correction

- Confirmed loss calculation
- Cautious hard sample search

2. Progressive label correction

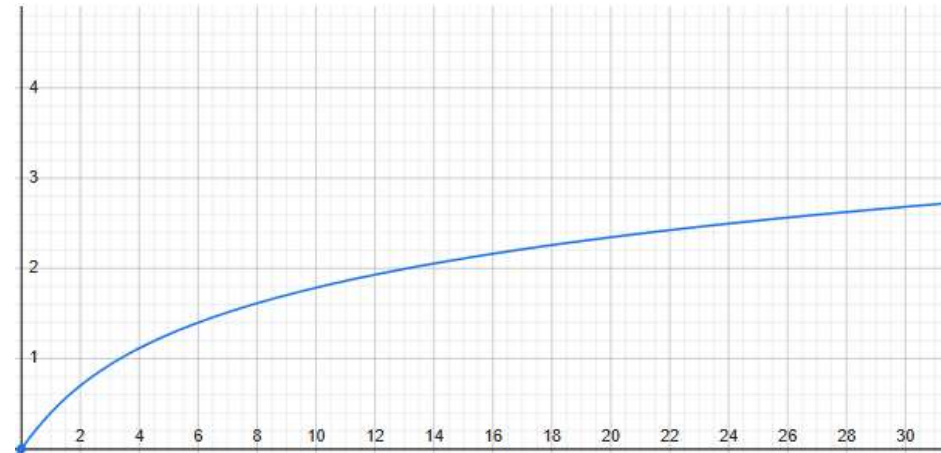
Method



1.1 Confirmed loss calculation

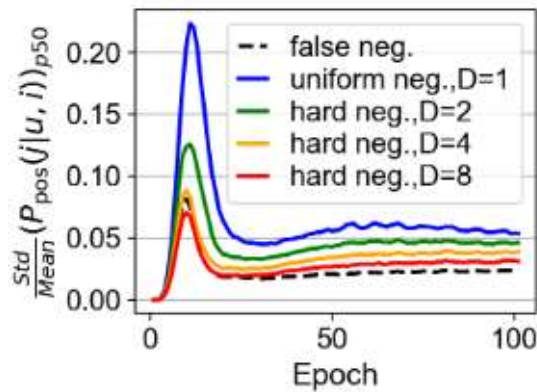
$$\mu_i = \frac{1}{v} \sum_{j=i-v+1}^i \phi(\ell(j)),$$

$$\phi(\ell) = \log(1 + \ell + \ell^2/2)$$



1. l 很大时, $\phi(l)$ 增长的很慢, 减小极端值的影响
2. l 较小时, $\phi(l)$ 和 l 近似为线性关系, 保留了原始的损失值

1.2 Cautious Hard Samples Search



(d) Comparing $\frac{\text{Std}}{\text{Mean}} P_{\text{pos}}$.

Theorem 1. Let $Z_n = \{z_1, \dots, z_n\}$ be an observation set with mean μ_z and variance σ^2 . By exploiting the non-decreasing damping function $\phi(z) = \log(1 + z + z^2/2)$. For any $\epsilon > 0$, we have

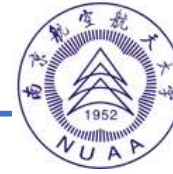
$$\left| \frac{1}{n} \sum_{i=1}^n \phi(z_i) - \mu_z \right| \leq \frac{\sigma^2 \left(n + \frac{\sigma^2 \log(\epsilon^{-1})}{n^2} \right)}{n - \sigma^2}, \quad (3)$$

with probability at least $1 - 2\epsilon$.

噪声样本的方差相对困难样本更低

使用置信区间的下界作为筛选困难样本的标准

Method



目标：寻找损失均值的置信上下界

$$\tilde{\mu}_z = \frac{1}{n} \sum_{i=1}^n \phi(z_i)$$
$$\tilde{\mu}_z^- \leq \tilde{\mu}_z \leq \tilde{\mu}_z^+$$

(1) 构建估计器

$$r(\tilde{\mu}_z) = \sum_{i=1}^n \phi[\alpha(z_i - \tilde{\mu}_z)] = 0.$$

其中 α 是正的实数参数

(2) 调整归一化

$$r(\theta) = \frac{1}{\alpha n} \sum_{i=1}^n \phi[\alpha(z_i - \theta)], \theta \in \mathbb{R}.$$

(3) 不等式推导

$$\mathbb{E}\{\exp[\alpha n r(\theta)]\} \leq \left\{ 1 + \alpha(\mu_z - \theta) + \frac{\alpha^2}{2} [\alpha^2 + (\mu_z - \theta)^2] \right\}^n$$
$$\leq \exp \left\{ n\alpha(\mu_z - \theta) + \frac{n\alpha^2}{2} [\alpha^2 + (\mu_z - \theta)^2] \right\}.$$

(4) 上下界计算

定义参数 θ_- 和 θ_+ , 使得 $r(\theta_-) > 0, r(\theta_+) < 0$

$$B_-(\theta) = \mu_z - \theta - \alpha [\sigma^2 + (\mu_z - \theta)^2] - \frac{\log(\epsilon^{-1})}{\alpha n}$$

$$B_+(\theta) = \mu_z - \theta + \alpha [\sigma^2 + (\mu_z - \theta)^2] + \frac{\log(\epsilon^{-1})}{\alpha n}$$

$$P(r(\theta) > B_-(\theta)) \geq 1 - \epsilon \quad \text{and} \quad P(r(\theta) < B_+(\theta)) \geq 1 - \epsilon.$$

根据切比雪夫不等式, 假设

$$4\alpha^2\sigma^2 + \frac{4\log(\epsilon^{-1})}{n} \leq 1.$$

置信区间, 以 $1 - 2\epsilon$ 的概率 $\tilde{\mu}_z^- \leq \tilde{\mu}_z \leq \tilde{\mu}_z^+$

$$\tilde{\mu}_z^- \geq \mu_z - \frac{\alpha\sigma^2 + \frac{\log(\epsilon^{-1})}{\alpha n}}{\alpha - 1} \quad \tilde{\mu}_z^+ \leq \mu_z + \frac{\alpha\sigma^2 + \frac{\log(\epsilon^{-1})}{\alpha n}}{\alpha - 1}.$$

Method



目标：寻找损失均值的置信上下界

(5)最终结果

$$|\tilde{\mu}_z - \mu_z| \leq \frac{\sigma^2 \left(n + \frac{\sigma^2 \log(\epsilon^{-1})}{n^2} \right)}{n - \sigma^2}.$$

$$\text{let } \epsilon = \frac{1}{2i},$$

$$l^* = \tilde{\mu} - \frac{\sigma^2 \left(i + \frac{\sigma^2 \log(2i)}{i^2} \right)}{d - \sigma^2},$$

d 表示和样本未被丢弃的次数;

置信区间下界 l^* 依赖于样本方差 σ^2 、当前迭代次数 i 和样本空间大小 d

Method



2 Progressive Label Correction

重标记的比例

$$r_i = \min\left(\frac{iR}{O}, R\right),$$

R: 最终重新标注的百分比

O: 第O轮训练后, 重新标注比例保持不变

选取重标记样本

$$\mathbb{I}(\ell_s^*, i) = \begin{cases} 1, & \text{if } \ell_s^* \geq T_i \\ 0, & \text{otherwise.} \end{cases}$$

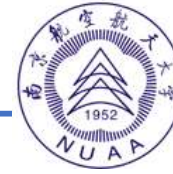
确定重标记样本的阈值

$$T_i = \lceil \lfloor B(1 - r_i) \rfloor \rceil.$$

标签翻转

$$y'_s = y_s + \mathbb{I}(\ell_s^*, i) \times (1 - 2y_s),$$

Method



Input : Training set D , Maximum Epochs N , Time interval v , Relabel Ratio R ;

Output: Trained Model \mathcal{M} .

```
1 Initialize model  $\mathcal{M}$  with random weights;
2 Initialize empty loss history  $L$  for each sample in  $D$ ;
3 for  $i = 1$  to  $N$  do
4   for each mini-batch  $B$  from  $D$  do
5     // Sample Dropping Correction
6     for each sample  $s$  in  $B$  do
7       Compute loss  $\ell$  using model  $\mathcal{M}$ ;
8       Append  $\ell$  to loss history  $L[s]$ ;
9       if  $\text{length}(L[s]) > v$  then
10        Remove samples in  $L[s]$  with time interval
11        greater than  $v$ ;
12      else
13        Compute mean loss  $\mu$  using Eq. (2);
14        Compute lower bound  $\ell^*$  using Eq. (4);
15      end
16    end
17    Update model  $\mathcal{M}$  using lower bound  $\ell^*$ ;
18  end
19 end
```

```
17 // Progressive Label Correction
18 Compute relabel ratio  $r$  using Eq. (5);
19 for each sample  $s$  in  $D$  do
20   if  $\ell^* \geq T_i$  ( $T_i$  is computed based on  $r$  and  $\ell^*$ ) then
21     Flip label  $y_s$  using Eq. (7);
22   else
23     Keep label  $y_s$ ;
24   end
25 end
26 Training the model with the corrected samples at the
27 next epoch;
28 return Trained Model  $\mathcal{M}$ ;
```

Experiment



南京航空航天大学

Nanjing University of Aeronautics and Astronautics

RQ1: performance

Dataset		Adressa				MovieLens				Yelp			
Base Model	Method	R@5	R@20	N@5	N@20	R@5	R@20	N@5	N@20	R@5	R@20	N@5	N@20
GMF	Normal	0.1257	0.2126	0.0908	0.1210	0.0355	0.1196	0.0482	0.0715	0.0149	0.0431	0.0152	0.0243
	WBPR	0.1258	0.2131	0.0912	0.1212	0.0357	0.1199	0.0486	0.0718	0.0148	0.0434	0.0155	0.0246
	WRMF	0.1263	0.2132	0.0911	0.1214	0.0363	0.1201	0.0491	0.0722	0.0144	0.0446	0.0140	0.0242
	T-CE	0.1251	0.2155	0.0913	0.1228	0.0374	0.1202	0.0509	0.0738	0.0143	0.0447	0.0143	0.0244
	DeCA	0.1232	0.2147	0.0866	0.1181	0.0364	0.1075	0.0423	0.0631	0.0141	0.0442	0.0139	0.0236
	BOD	0.1237	0.2153	0.0892	0.1193	0.0366	0.1083	0.0474	0.0685	0.0151	0.0436	0.0157	0.0247
	DCF (Ours)		0.1296*	0.2183*	0.0938*	0.1254*	0.0427*	0.1175	0.0543*	0.0743*	0.0155**	0.0458*	0.0158
NeuMF	Normal	0.1909	0.3078	0.1427	0.1851	0.0439	0.1084	0.0516	0.0724	0.0123	0.0386	0.0123	0.0210
	WBPR	0.1903	0.3082	0.1428	0.1848	0.0426	0.1132	0.0504	0.0735	0.0104	0.0384	0.0108	0.0193
	WRMF	0.1922	0.3084	0.1424	0.1852	0.0418	0.1180	0.0512	0.0729	0.0119	0.0378	0.0121	0.0198
	T-CE	0.1880	0.3080	0.1410	0.1847	0.0366	0.1065	0.0482	0.0680	0.0108	0.0383	0.0105	0.0190
	DeCA	0.1870	0.3076	0.1402	0.1804	0.0327	0.0990	0.0388	0.0590	0.0103	0.0381	0.0101	0.0182
	BOD	0.1890	0.3082	0.1414	0.1828	0.0375	0.0134	0.0489	0.0703	0.0126	0.0389	0.0119	0.0215
	DCF (Ours)		0.1979*	0.3134*	0.1439*	0.1853	0.0513*	0.1210*	0.0642*	0.0816*	0.0132*	0.0411*	0.0132*
NGCF	Normal	0.1235	0.2257	0.0934	0.1291	0.0335	0.1015	0.0452	0.0634	0.0172	0.0495	0.0174	0.0273
	WBPR	0.1237	0.2252	0.0936	0.1289	0.0331	0.1013	0.0446	0.0637	0.0166	0.0481	0.0164	0.0267
	WRMF	0.1244	0.2255	0.0942	0.1304	0.0334	0.1011	0.0449	0.0639	0.0169	0.0486	0.0167	0.0270
	T-CE	0.1260	0.2270	0.0959	0.1313	0.0335	0.1014	0.0450	0.0635	0.0173	0.0497	0.0173	0.0270
	DeCA	0.1172	0.2235	0.0846	0.1037	0.0318	0.0973	0.0436	0.0627	0.0169	0.0464	0.0166	0.0268
	BOD	0.1212	0.2246	0.0901	0.1265	0.0321	0.1008	0.0437	0.0633	0.0174	0.0492	0.0169	0.0274
	DCF (Ours)		0.1267*	0.2275**	0.0970*	0.1321*	0.0353*	0.1037*	0.0468*	0.0647*	0.0180*	0.0503**	0.0179**
LightGCN	Normal	0.1236	0.2257	0.0933	0.1289	0.0347	0.1029	0.0457	0.0642	0.0185	0.0514	0.0183	0.0291
	WBPR	0.1239	0.2253	0.0914	0.1295	0.0353	0.1043	0.0460	0.0638	0.0182	0.0514	0.0178	0.0282
	WRMF	0.1242	0.2260	0.0928	0.1298	0.0357	0.1046	0.0464	0.0650	0.0181	0.0515	0.0180	0.0283
	T-CE	0.1261	0.2261	0.0956	0.1306	0.0290	0.1020	0.0409	0.0612	0.0185	0.0516	0.0183	0.0290
	DeCA	0.1185	0.2251	0.0859	0.1038	0.0347	0.0987	0.0440	0.0640	0.0176	0.0503	0.0172	0.0273
	BOD	0.1225	0.2254	0.0902	0.1287	0.0325	0.1014	0.0443	0.0638	0.0182	0.0513	0.0177	0.0282
	DCF (Ours)		0.1258	0.2274*	0.0961**	0.1315*	0.0365*	0.1050	0.0472*	0.0659*	0.0192*	0.0523*	0.0187**

Experiment



南京航空航天大学

Nanjing University of Aeronautics and Astronautics

RQ2: each component's contribution

Table 4: The effect of each components on DCF.

Method	R@5	R@10	N@5	N@10
T-CE	0.0374	0.0734	0.0509	0.0591
DCF _{CL}	0.0432(15.5%)	0.0751(2.3%)	0.0527(3.5%)	0.0600(1.5%)
DCF _{HS}	0.0435(16.3%)	0.0772(5.2%)	0.0535(5.1%)	0.0610(3.2%)
DCF _{LC}	0.0430(15.0%)	0.0756(3.0%)	0.0529(3.9%)	0.0604(2.2%)
DCF _{CL+HS}	0.0436(16.6%)	0.0743(1.2%)	0.0531(4.3%)	0.0596(0.8%)
DCF _{CL+LC}	0.0459(22.7%)	0.0763(4.0%)	0.0549(7.9%)	0.0612(3.6%)
DCF _{HS+LC}	0.0454(21.4%)	0.0779(6.1%)	0.0545(7.1%)	0.0613(3.7%)
DCF _{ALL}	0.0471(25.9%)	0.0789(7.5%)	0.0553(8.6%)	0.0621(5.1%)

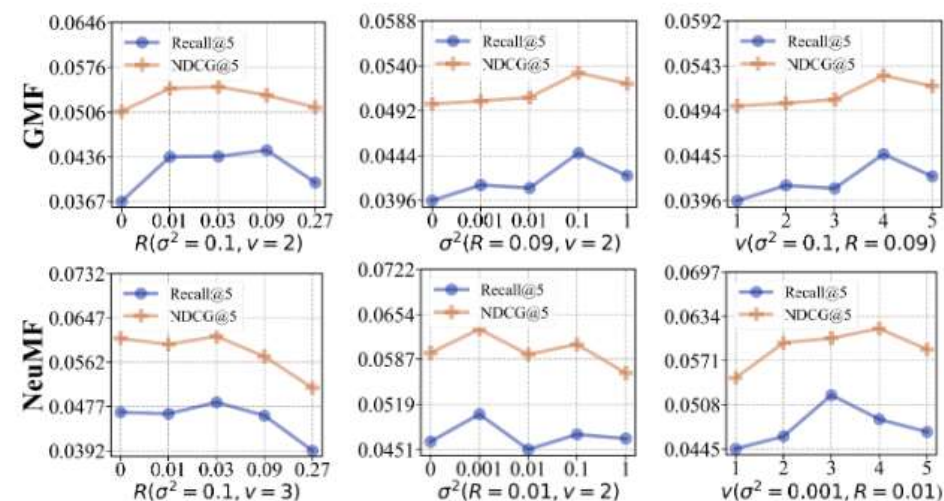


Figure 4: Impact of the relabel ratio R , search discretion level σ^2 and time interval v .

Experiment

RQ3: hard samples

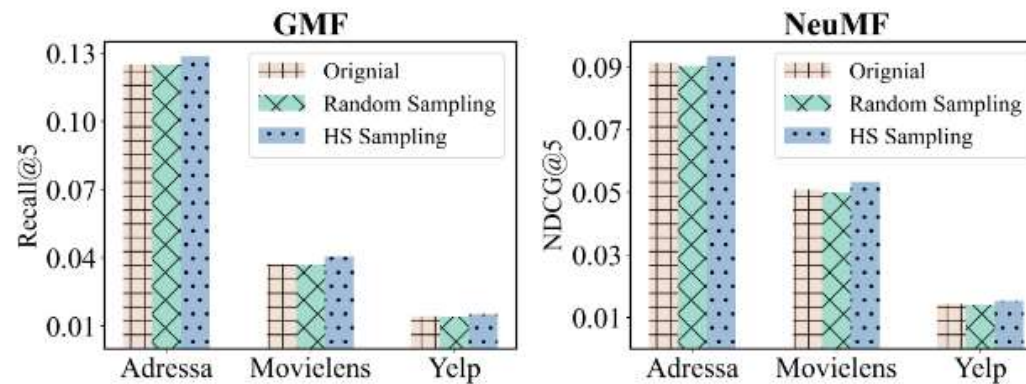


Figure 5: Comparative experiments on three datasets with two backbones validate the effectiveness of our hard sample search strategy to improve model performance.

Experiment

RQ4: progressive strategy

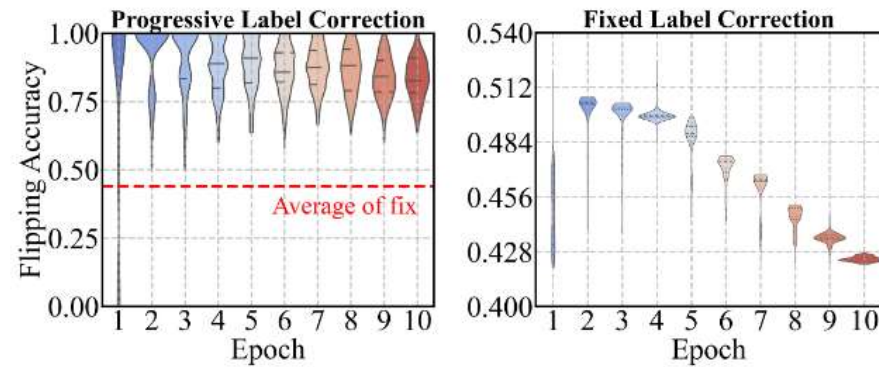


Figure 6: Comparison of flip accuracy between progressive label correction and fixed. For a clear presentation, we use a violin plot here. Additionally, we mark the average flip accuracy of fixed with a red line to clearly highlight the superiority of our progressive strategy.