



浙江大学  
ZHEJIANG UNIVERSITY

# 如何找新Idea?

---- 以推荐系统纠偏为例

陈佳伟

浙江大学计算机学院

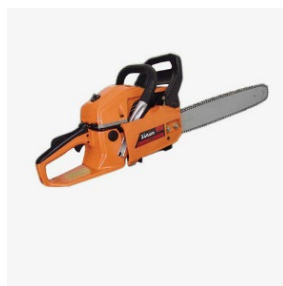
Sep 2024

# 什么是idea?

大部分Idea: 用适当方法  
解决有用的问题



方法:



问题:

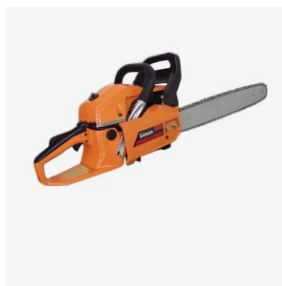


# 怎么找idea? --- 更好的模型(借鉴)

现有方法:



观察到:



新的idea:



为何方法更  
适合这个问题?

# 怎么找idea? --- 更好的模型(改进)

现有方法:



发现: 柄太短了

针对问题,  
挖掘现有方  
法的缺陷

新的idea:

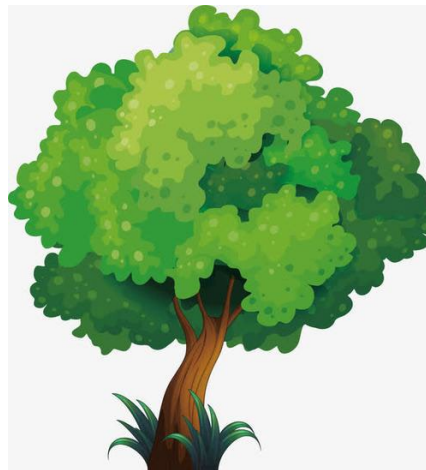


# 怎么找idea? --- 新的问题

现有方法:



新的idea:



阐明新问题的  
意义，分  
析其性质

# 怎么找idea? --- 分析型 (方法或问题的性质)



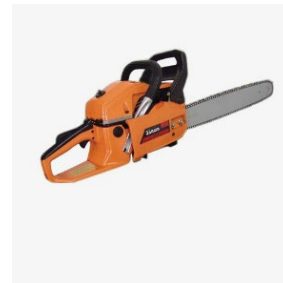
深入研究斧头(方法)的性质



深入研究树(问题)的性质

# 怎么找idea? --- 怎么做?

掌握更多的方法:

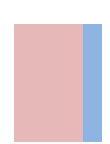


关注其他领域的  
方法

深入理解问题:

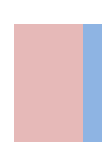


对该领域全  
面调研



# 写好一篇文章





# 写好一篇文章

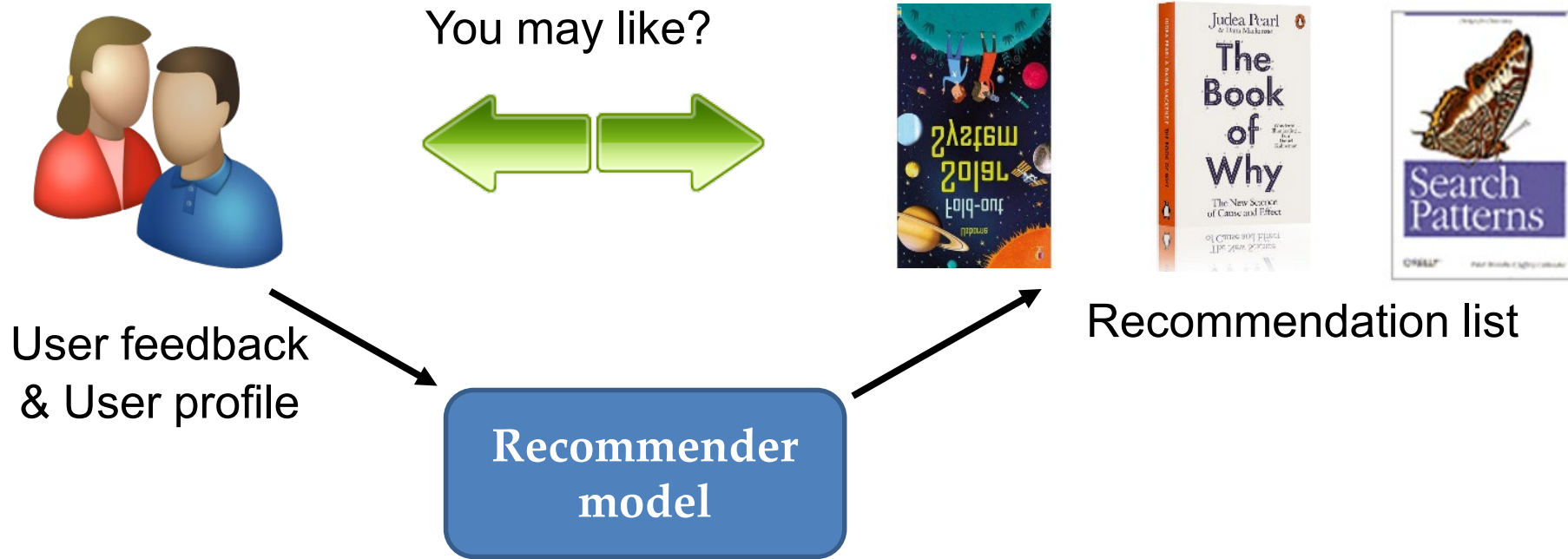
## **1. Bias in Recommendation**

- Recommender Systems
- Bias in Recommendation
- The Influences of Bias

## **2. My Recent Debiasing Work**

# Recommender System

□ Recommender system (RS) helps address information overload

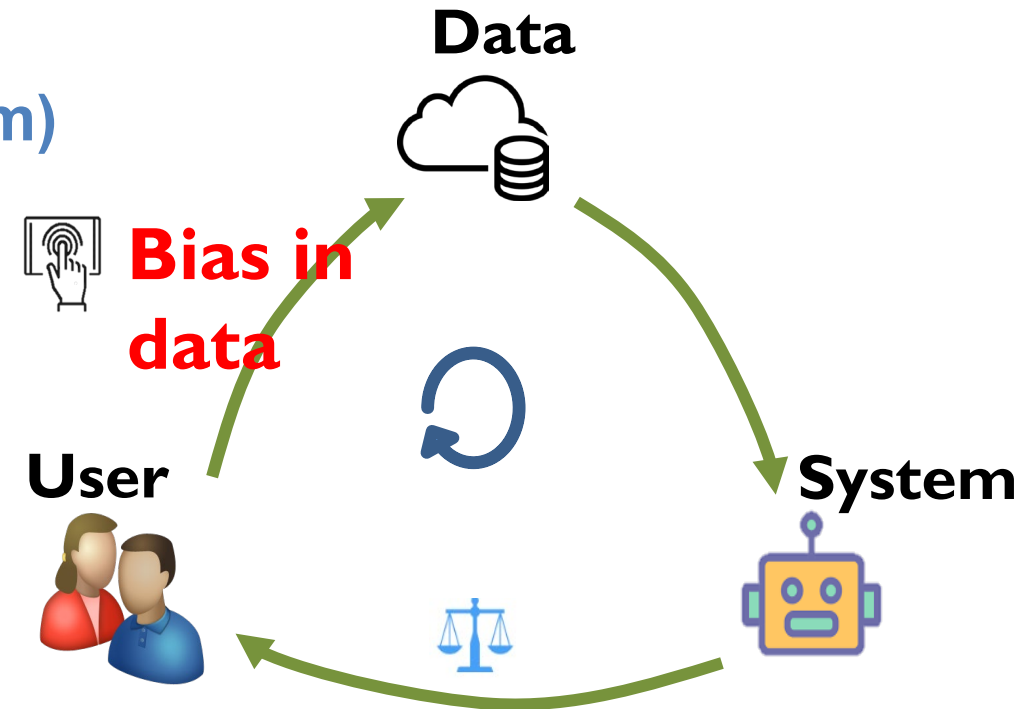


- RS captures user preference and provides **personalized information filtering**
- **However, bias widely exists in RS**

# Bias in Recommendation

## □ Bias in Data (collecting)

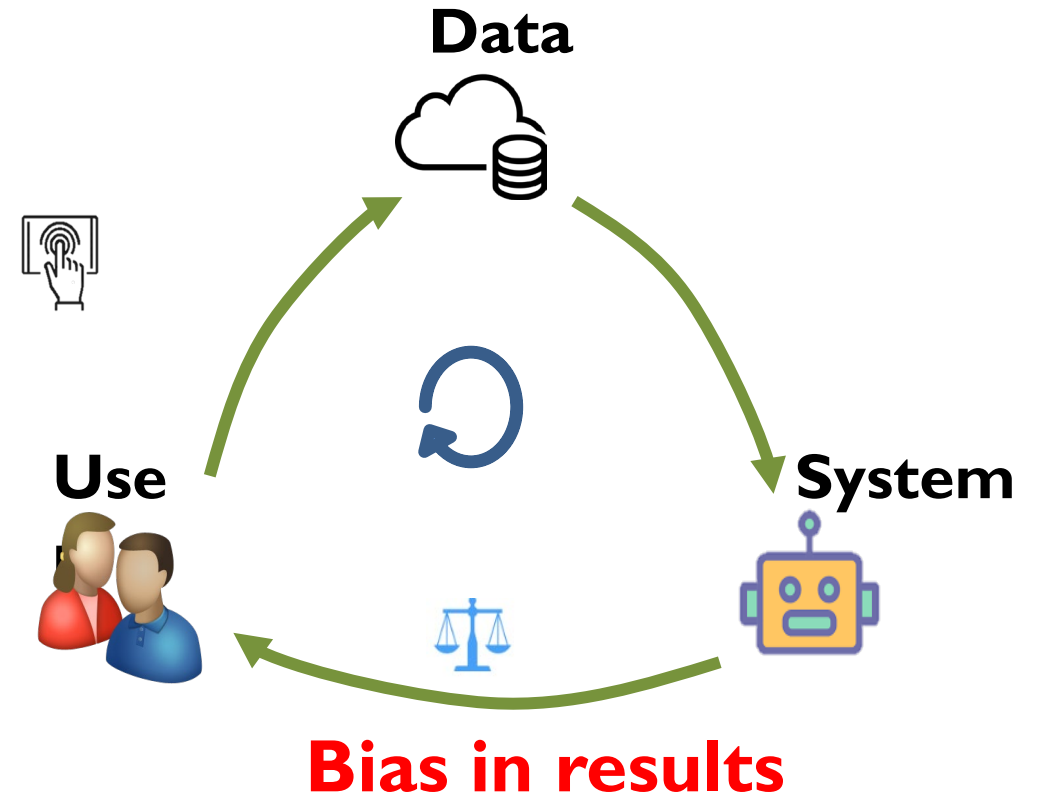
- Data is observational rather than experimental (i.e., missing-not-at-random)
- Labels do not faithfully reflect user preference
- Affected by many factors:
  - Exposure mechanism
  - Public opinions
  - Display position
  - .....
- The collected data deviates from user true preference



# Bias in Recommendation

## □ Bias in Results (training & serving)

- Models do not only **inherit** the bias in data but also **amplify/generate** bias
- Models naturally bias towards dominated groups
- E.g. Popularity bias, mainstream bias



# Bias is Evil

## □ Economic perspective

- Decreasing recommendation accuracy
- Hurting user experience and satisfaction
- Causing the losses of users

- **Example: The clickbait problem.**

Items with interesting title (but boring content) may get more exposure opportunity. But users do not like them!



# Bias is Evil

## □ Social perspective



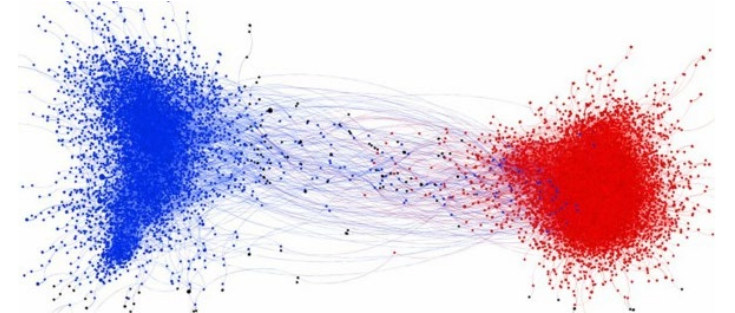
### Unfairness

- E.g., job recommendation [1]:
- software developer → man
  - registered nurse → woman



### Filter Bubble

- E.g., Tittytainment
- Entertainment videos/games
  - Addictive



### Polarization

- E.g., political polarization [1]:
- Democrats vs Republicans

**Debiasing is vitally important!**

## 1. Bias in Recommendation

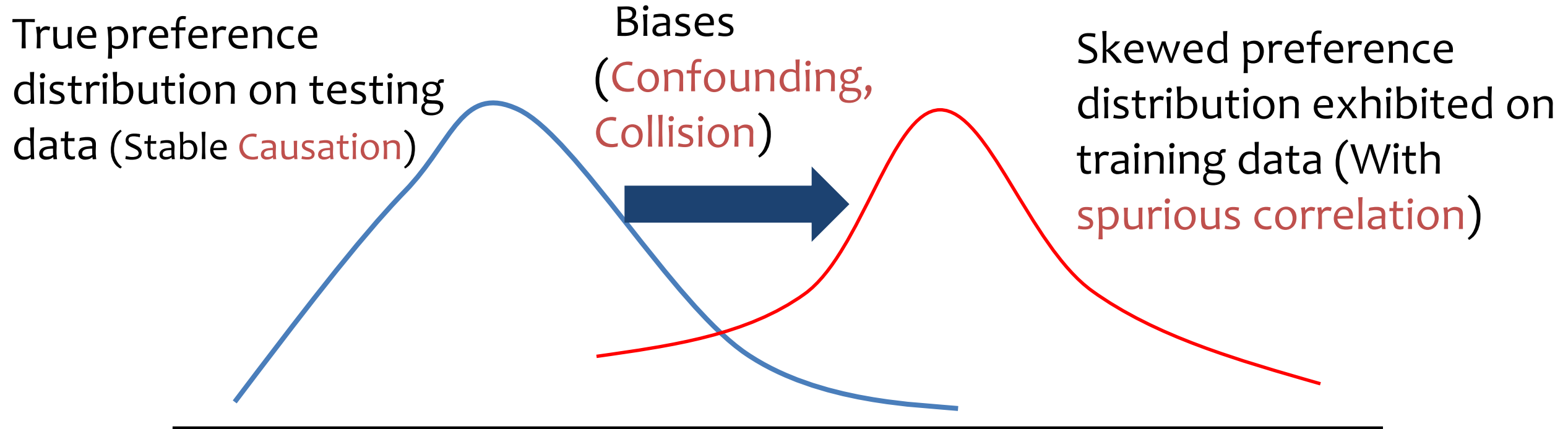
## 2. Our Recent Debiasing Work

- Survey on Recommendation Bias
- MACR (借鉴)
- TIDE (改进)
- UnKD (新问题)
- Adap- $\tau$  (分析)



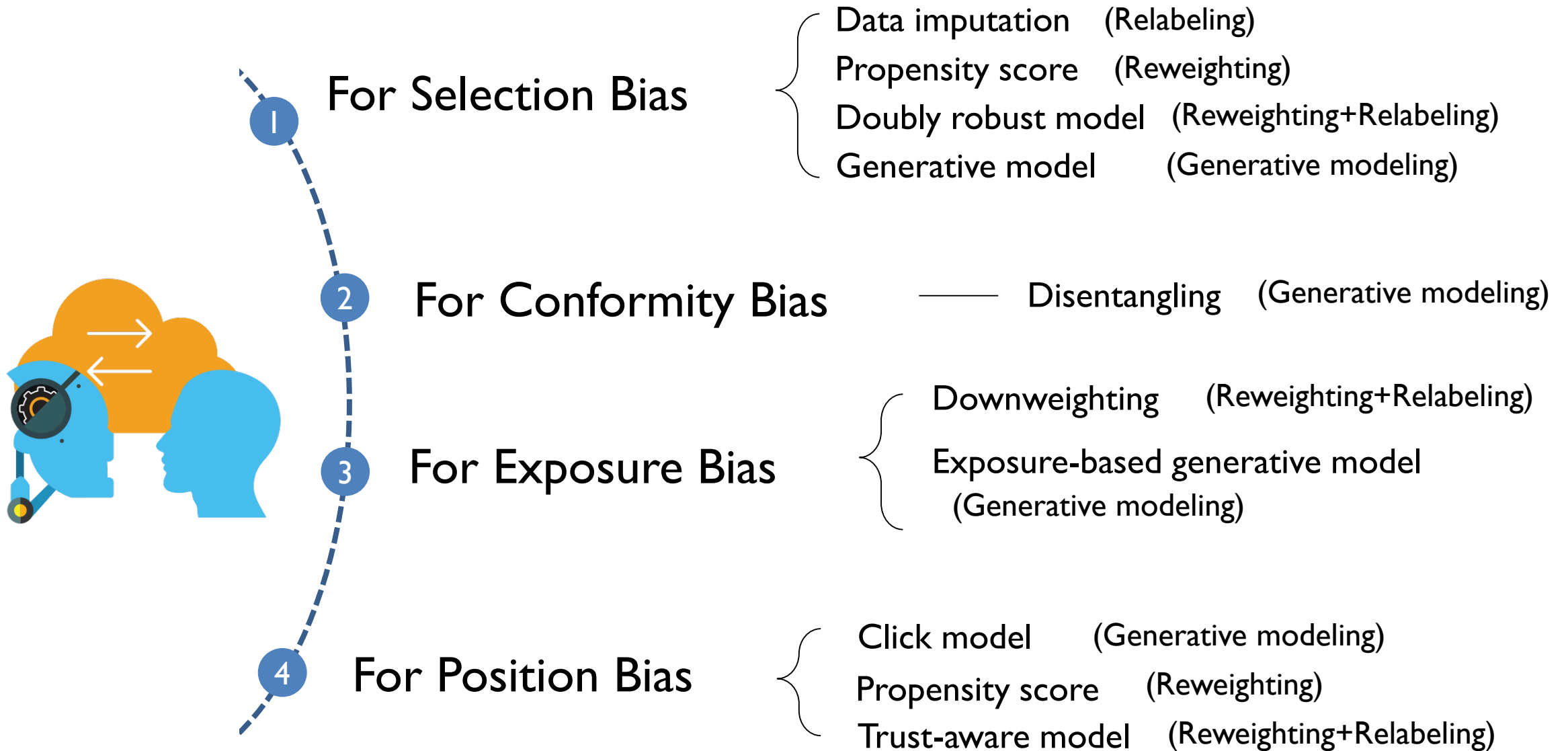
# Survey on Recommendation Bias(对问题总结和归纳)

- Data-driven methods would learn skewed user preference:



- Data-driven methods may infer **spurious correlations**, which are deviated from reflecting user true preference, and lack interpretation.

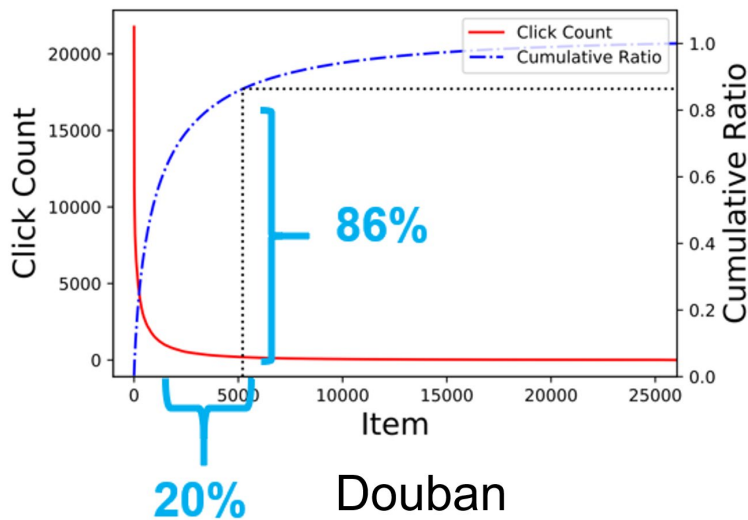
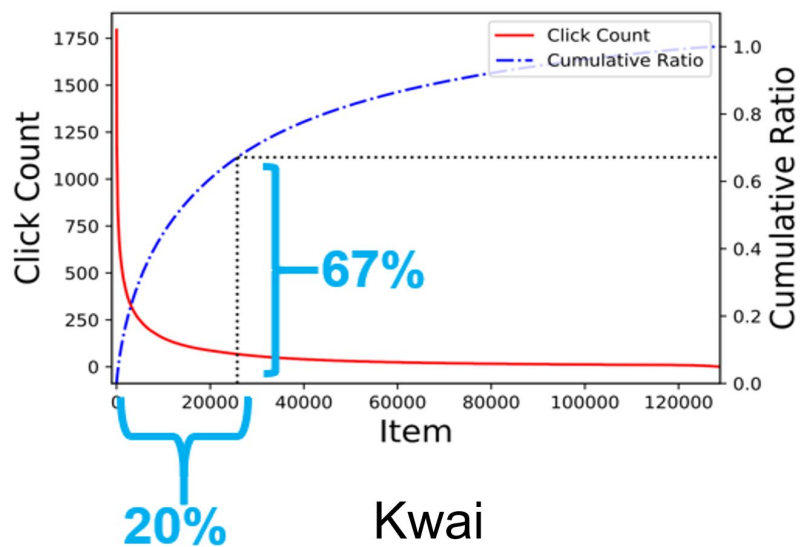
# • Debiasing Strategies Overview



# Popularity Bias

## □ Popularity bias in recommender system

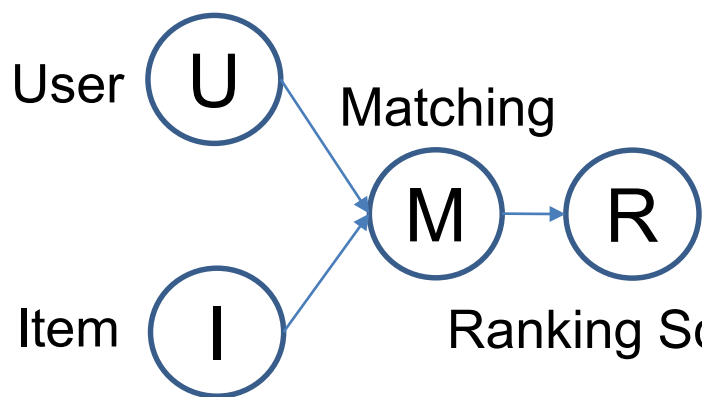
- **Popularity bias:** Popular items are recommended even **more frequently** than they would warrant
  - Long-tail distribution, Matthew effect
  - Harmful for personality and accuracy



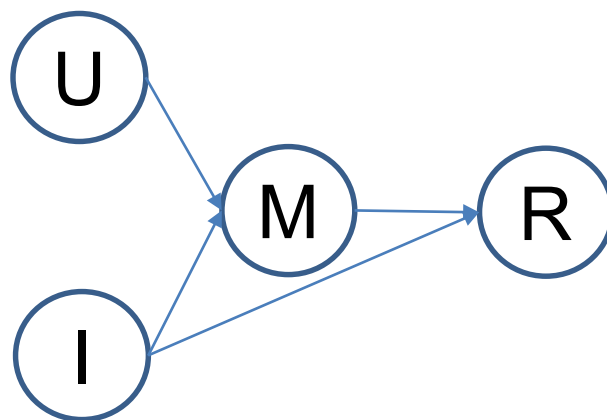
Long-tail distribution

# MACR: Model-Agnostic Counterfactual Reasoning (借鉴)

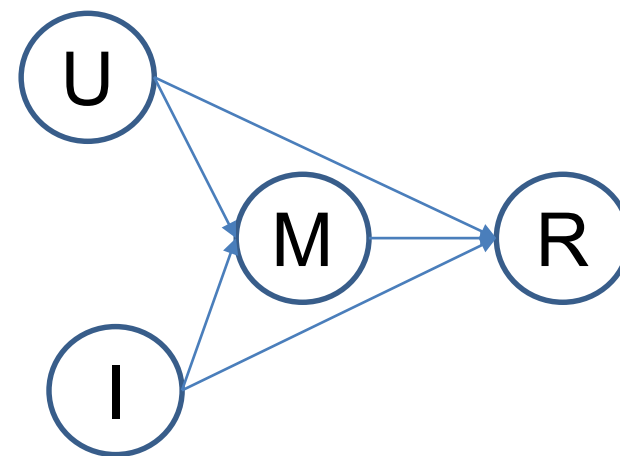
## □ Causal View of Popularity Bias



Common Recommender  
User-Item Matching



Popularity bias modeling:  
Incorporating item popularity



User-specific modeling:  
Incorporating item popularity  
& user activity

- Edge  $I \rightarrow R$  captures **popularity bias**.
- Edge  $U \rightarrow R$  captures the **user' sensitivity to popularity**.

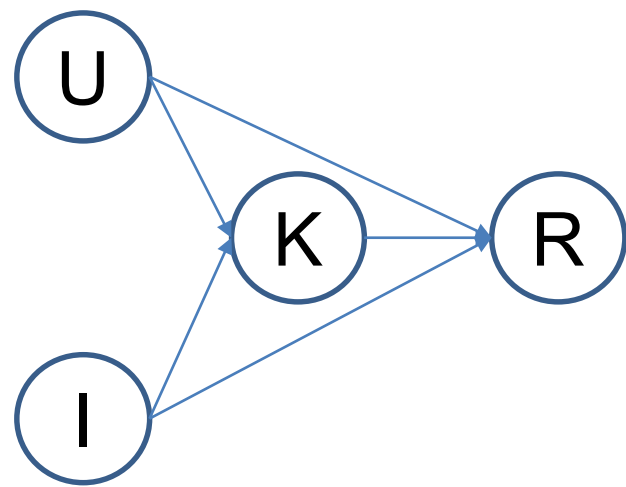
## □ Solution Idea:

- Train a recommender based on the causal graph via a **multi-task learning**
- Perform **counterfactual inference** to eliminate **popularity bias**

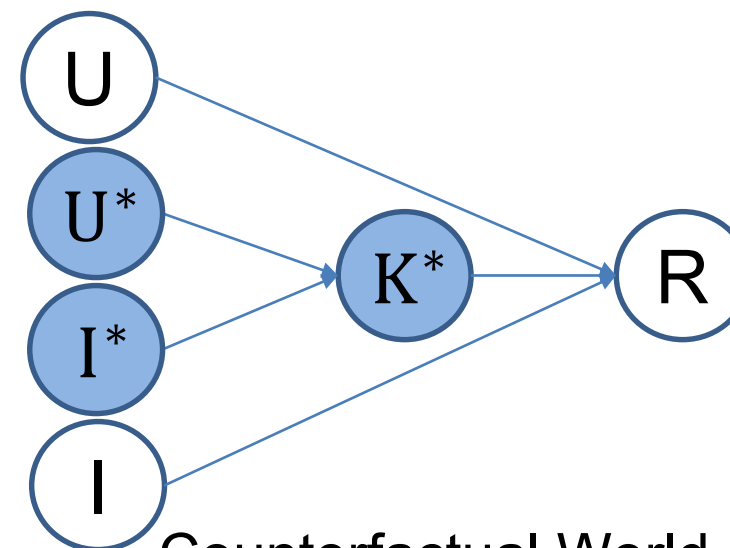
*Question to answer: what would the prediction be if there were only popularity bias?*<sup>20</sup>

# MACR: Model-Agnostic Counterfactual Reasoning

- Counterfactual Inference to Remove Bias



Factual World  
(original prediction)



Counterfactual World  
(block matching to capture bias)

$$TIE = TE - NDE = Y(U = u, I = i, K = K_{u,i}) - Y(U = u, I = i, K = K_{u^*,i^*})$$

Total Indirect  
Effect

Total  
Effect

Natural  
Direct

Factual world

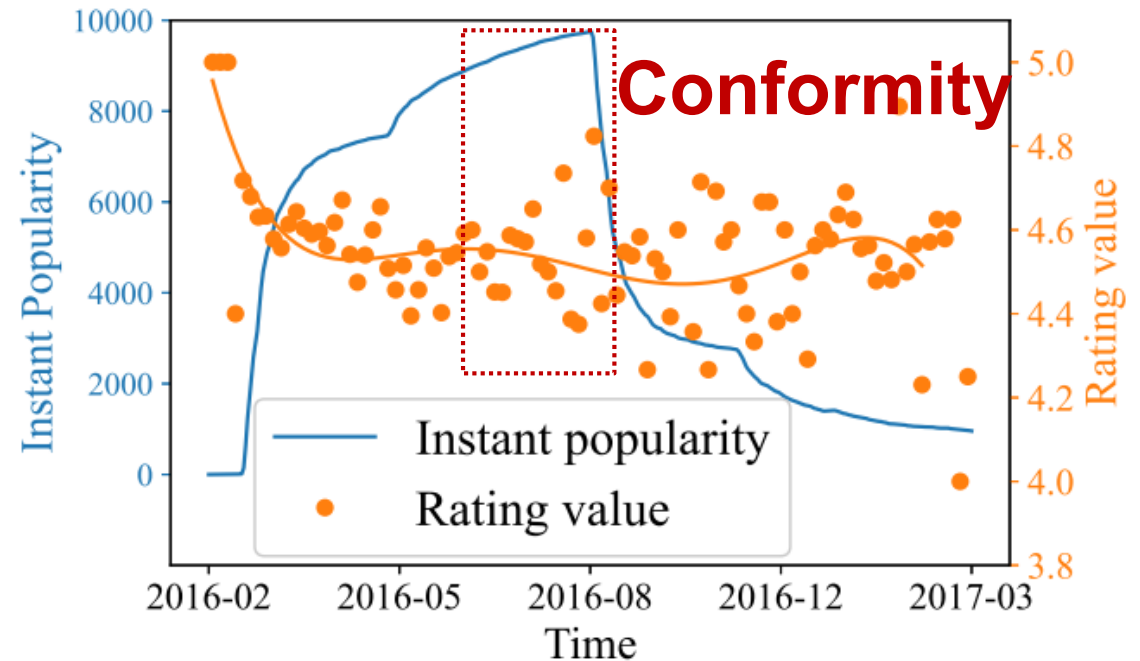
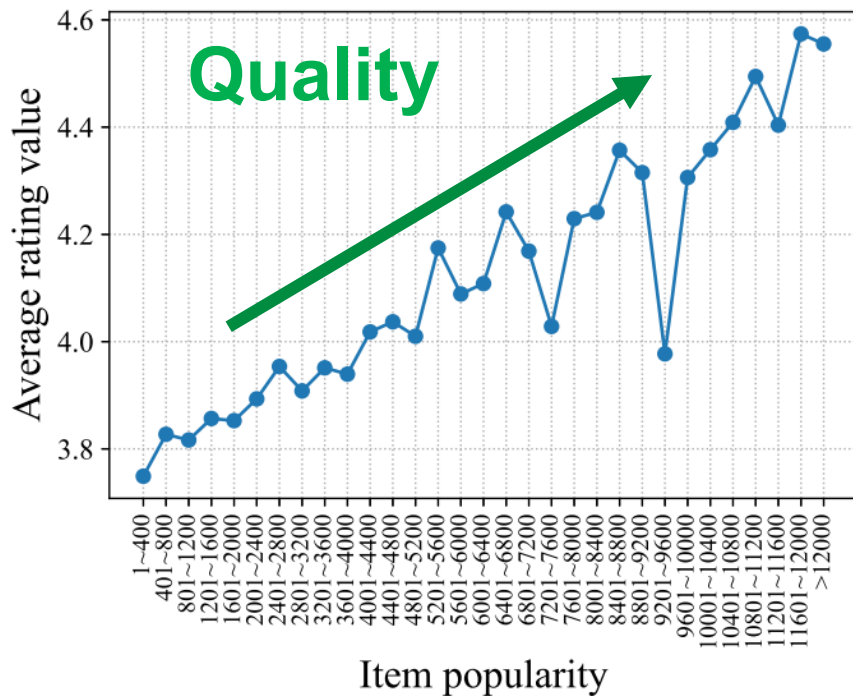
Counterfactual world

**Inference** with  $TIE = \hat{y}_k \times \sigma(\hat{y}_i) \times \sigma(\hat{y}_u) - c \times \sigma(\hat{y}_i) \times \sigma(\hat{y}_u)$

# TIDE: Disentangling Benign and Harmful Bias (改进)

- **Conflicting Observation:**

- The more **popular** an item is, the larger **average rating value** the item tends to have (**positive correlation**).
- From the temporal view, for a large proportion of items, the **rating value** exhibits **negative correlation** with the item **popularity** at that time



# TIDE: Disentangling Benign and Harmful Bias

## Time-aware DisEntangled framework(TIDE)

- Main challenge: Lack of explicit signal for disentanglement

## Quality is static: $I \rightarrow Q \rightarrow Y$

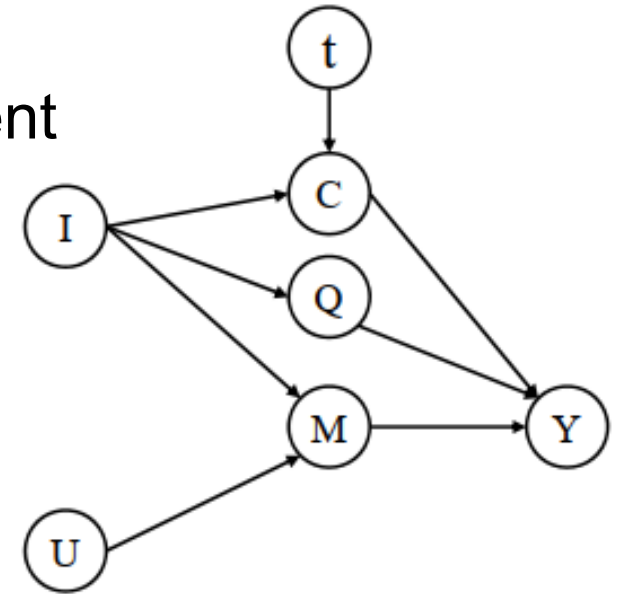
- Quality has **stable** influence on users' behavior

## Conformity is dynamic: $(I, t) \rightarrow C \rightarrow Y$

- Conformity is **time-sensitive**, since recent interactions should have stronger influence.

## User interest: $(U, I) \rightarrow M \rightarrow Y$

- User and item's matching score, can be Implemented by various recommendation models, such as MF, LightGCN, etc.



(a) Causal graph of our TIDE.

U: User    I: Item  
t: time    C: conformity  
Q: Quality    Y: Prediction  
M: Matching score

# TIDE: Disentangling Benign and Harmful Bias

## □ Training Stage:

- ⊙ **Popularity** comes from **Quality** and **Conformity**
- ⊙ Prediction with **Popularity** and **matching score**

$$\hat{y}_{ui}^t = \text{Tanh}(q_i + c_i^t) \times \text{Softplus}(m_{ui})$$

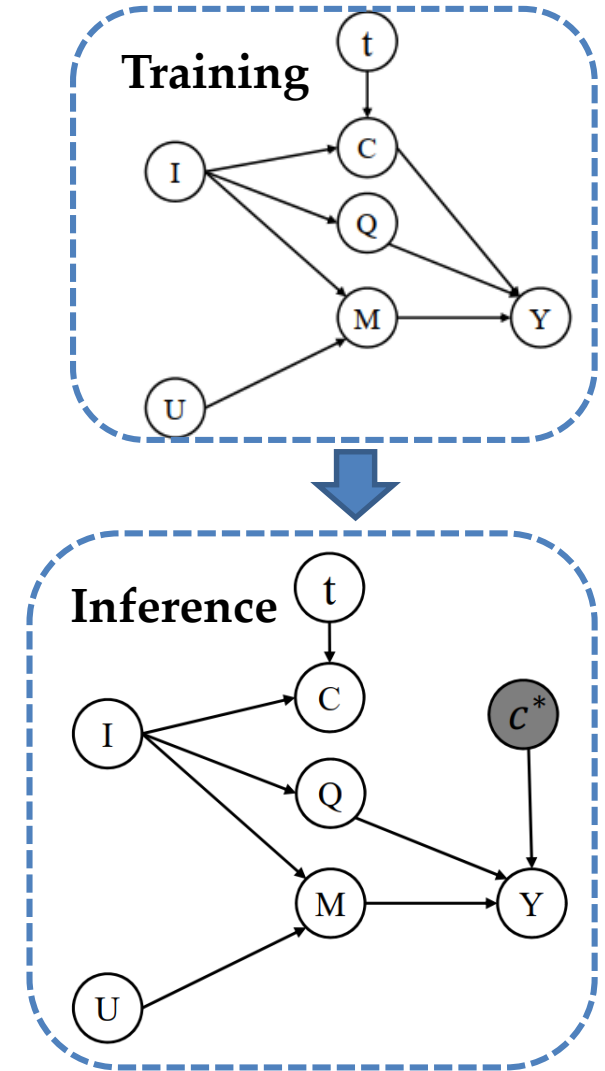
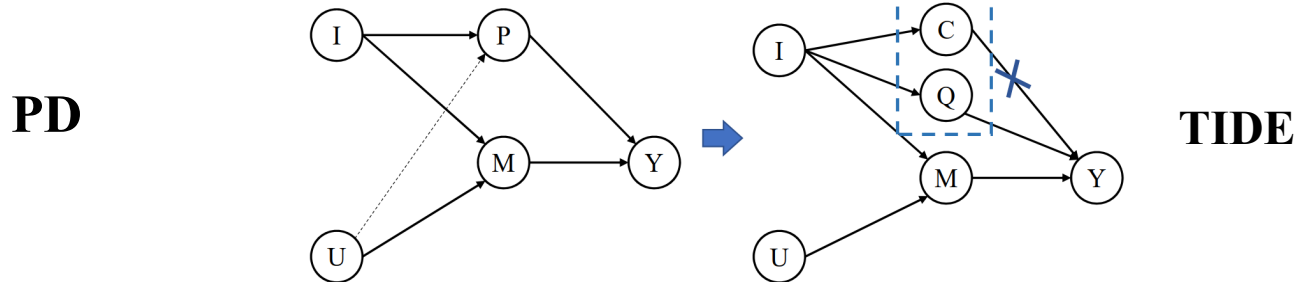
## □ Inference Stage:

- ⊙ Intervention: set  $c$  as reference vector  $c^*$  (e.g., zero) during inference to **remove** the **improper effect from C to Y**.

$$\hat{y}_{ui}^* = \tanh(q_i + c^*) \times \text{Softplus}(m_{ui})$$

## □ Comparison with PD

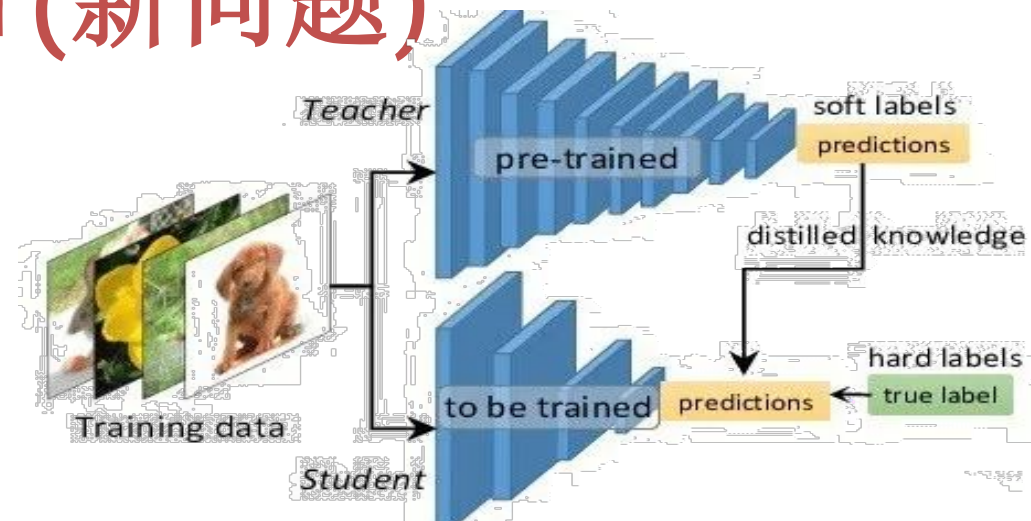
- ⊙ TIDE further conduct disentanglement of popularity bias





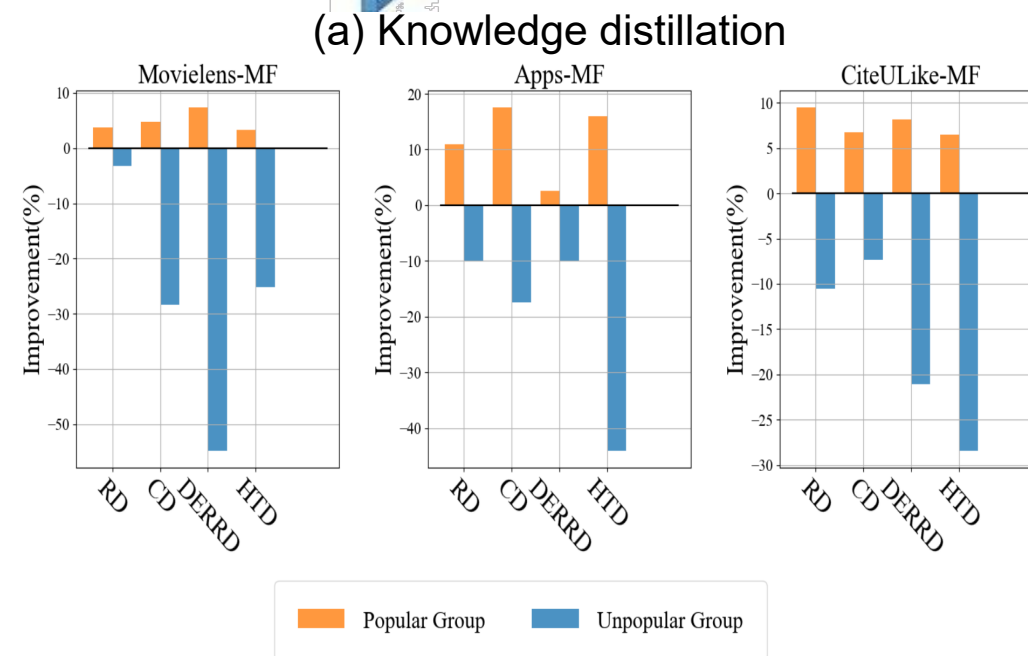
# UnKD: Unbiased Distillation (新问题)

□ **Knowledge distillation:** Transfer knowledge from large models to small models.



□ **Existing knowledge distillation:**

- Performance is mainly improved in the popular group
- The performance of the unpopular group was severely degraded.



(b) The improvement ratio of different distillation methods.

# UnKD: Unbiased Distillation

- **Conditional total effect** : For a particular user  $u$ ,

$$TE_i = Y_{i|u} - Y_{i^*|u} \quad ; \quad i^* \text{ is the benchmark situation}$$

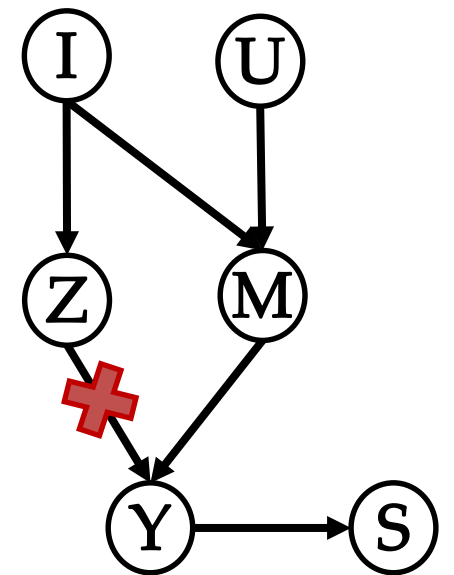
- **Adverse causal effect** :  $PEZ_i = Y_{i^*,Z_i|u} - Y_{i^*|u}$

- **Eliminating bias** :  $PEM_i = TE_i - PEZ_i = Y_{i|u} - Y_{i^*,Z_i|u}$

- **For any two items  $i$  and  $j$  with highly similar popularity**, we

have  $Z_i \approx Z_j$ , so the equation  $Y_{i^*,Z_i|u} = Y_{i^*,Z_j|u}$  almost holds. Then, we have:

$$Y_{i|u} > Y_{j|u} \Leftrightarrow Y_{i|u} - Y_{i^*,Z_i|u} > Y_{j|u} - Y_{i^*,Z_j|u} \Leftrightarrow PEM_i > PEM_j$$



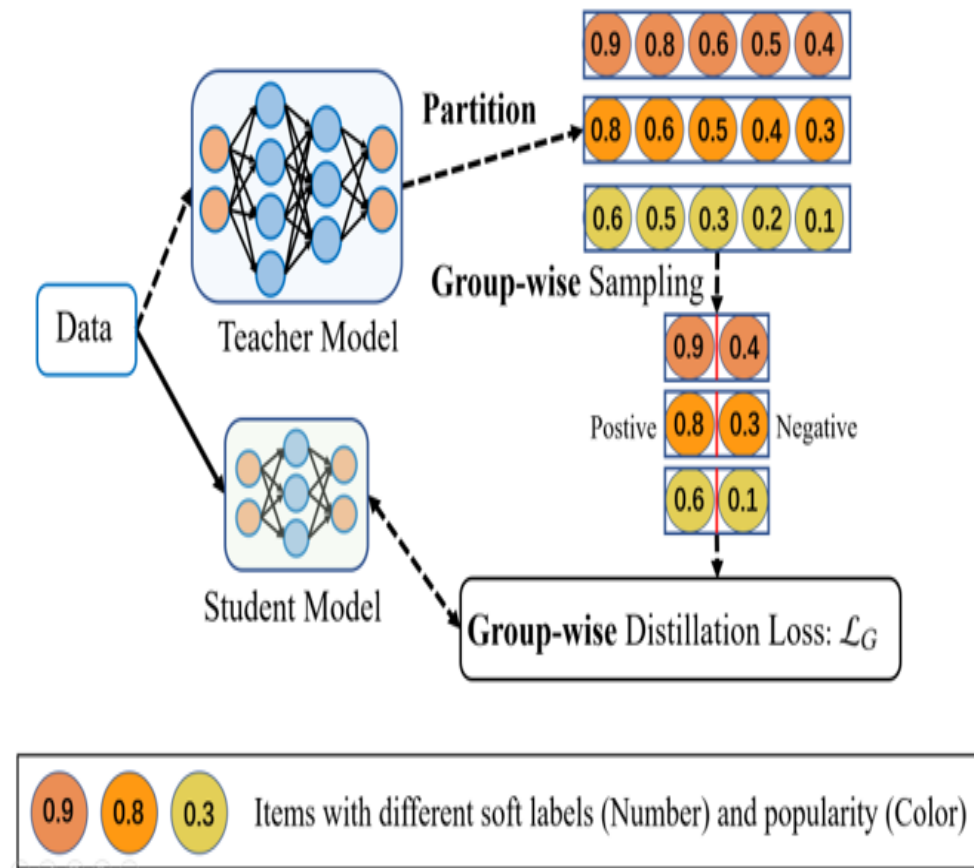
$U$ : user;  
 $I$ : item;  
 $M$ : affinity score;  
 $Z$ : popularity;  
 $Y$ : soft label ;  
 $S$ : student

# UnKD: Unbiased Distillation

□ UnKD consists of the following three steps:

- **Group partition.** Items are divided into  $K$  groups according to their popularity.
- **Group-wise Sampling.** Sample a set of  $\mathcal{S}_{ug}$  positive and negative item pairs for each group  $g(i^+, i^-)$ .
- **Group-wise Learning.** The student model is trained with a group distillation loss:

$$\mathcal{L}_G = - \sum_u \frac{1}{|\mathcal{U}|} \sum_{g \in \mathcal{G}} \sum_{(i^+, i^- \in \mathcal{S}_{ug})} \log \sigma(\mathbf{e}_u^T \mathbf{e}_{i^+} - \mathbf{e}_u^T \mathbf{e}_{i^-})$$



(a) UnKD Structure

# Adap- $\tau$ : Embedding Normalization(发现+分析问题)

## □ Harm of Un-normalized embedding:

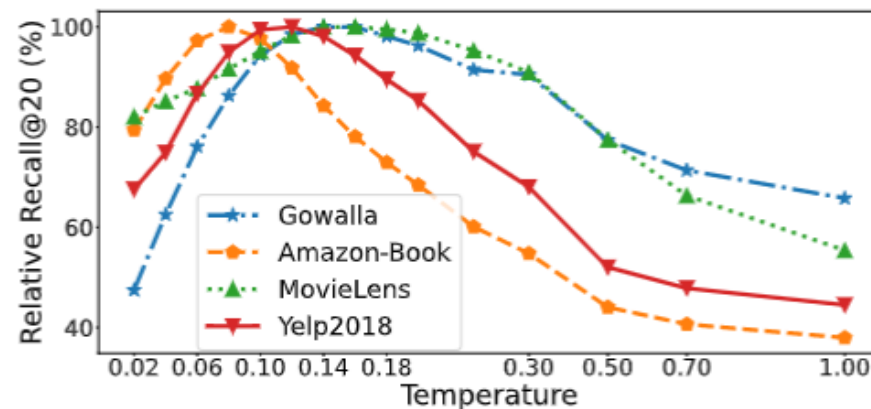
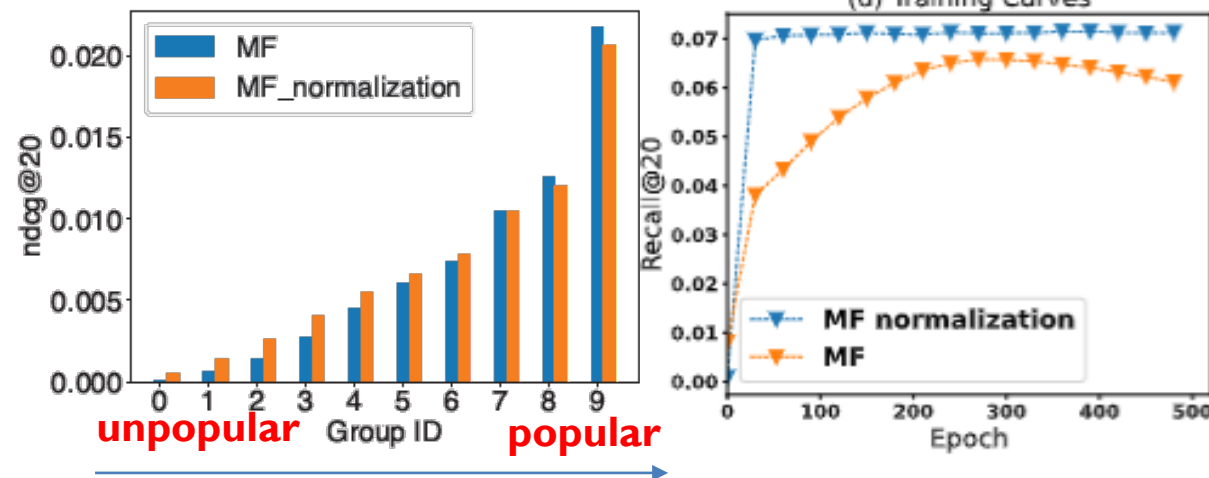
- Aggravates **popularity bias**
- Hurts training **convergence**

## □ Challenge of embedding normalization:

- highly sensitive to the hyper-parameter **temperature  $\tau$** .
- Finding proper  $\tau$  is difficult.

inner product:  $f_{u,i} = e_u^T \cdot e_i$

normalization:  $f_{u,i} = \frac{e_u^T \cdot e_i}{\|e_u\| \cdot \|e_i\|} \cdot \frac{1}{\tau}$



# Adap- $\tau$ : Embedding Normalization

## □ Understanding roles of temperature:

- R1: Avoiding from **gradient vanishment**  
=> temperature should adapt to the data and model
- R2: **Balancing contributions** from hard instances and easy instances  
=> temperature could be fined-grained (personalized) for flexible adjustment
- **We propose adaptive and fine-grained temperature:**

- **Adaptive:** maximizing the cumulated gradient magnitude:

$$\tau_0 \approx \frac{\mu_+ - \mu}{\log(\frac{nm}{2|D|})}$$

$\mu_+$ : average score of positive instances

$\mu$ : average score of all instances

$n$ : the number of users

$m$ : the number of items

$D$ : the number of positive instances

- **Fine-grained:** monitoring the loss for each user:

$$\tau_u^* = \tau_0 \cdot \exp(\mathbb{W}(\max(-\frac{1}{e}, \frac{L(u) - m_u}{2\beta})))$$

$L(u)$ : the loss for the user  $u$

$m_u$ : mean of  $L(u)$

$\mathbb{W}$ : lambert-W function

# Conclusion and Future Work

- **Idea? 合适的方法+有用的问题**
  - 合适的方法：借鉴，改进
  - 有用的问题：寻找新的问题（可解决的）
  - 分析型文章：研究方法<sub>和问题的性质</sub>
- **Future direction: 以动态图表征为例**
  - 合适的方法：自监督学习、Transformer、改进效率
  - 有用的问题：鲁棒性、可解释性、隐私保护、动态图应用
  - 分析型：表征的几何结构？



THANK YOU!