

# **Causal Understanding of Fake News Dissemination on Social Media**

Speaker: Zi-Xin, Chen

Advisor: Jia-Ling, Koh

Date: 2021/12/14

# Introduction

# Motivation

- False or misleading information disguised in news articles to mislead consumers has raised serious concerns, demanding novel approaches to understanding fake news dissemination
- Great effort can be seen in computational fake news detection, but less is known about what user attributes **cause** some users to share fake news

# Goal

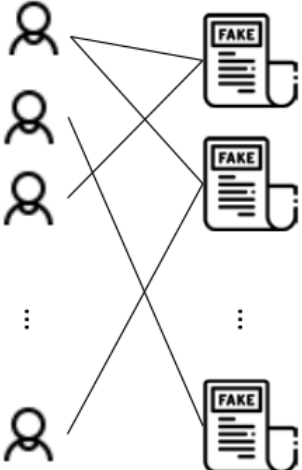
1. Propose a principled approach to alleviating **selection bias** in fake news dissemination.
2. Then consider the learned **unbiased** fake news sharing behavior as the surrogate confounder that can fully capture the causal links between user attributes and user susceptibility.
3. Understand what user attributes potentially **cause** users to share fake news

# Input and Output

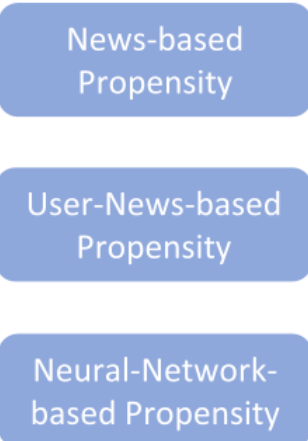
- Input: Fake news sharing behavior
- Output: User attributes that cause the sharing behavior

# Overview Framework

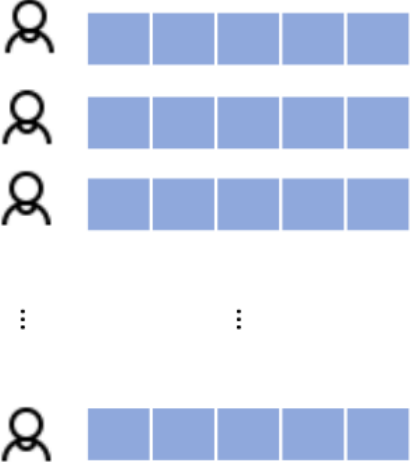
① Fake news dissemination



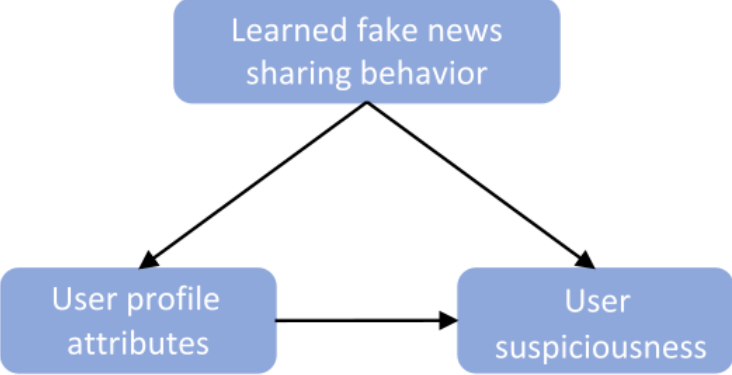
② IPS reweighting



③ Unbiased fake news sharing behavior

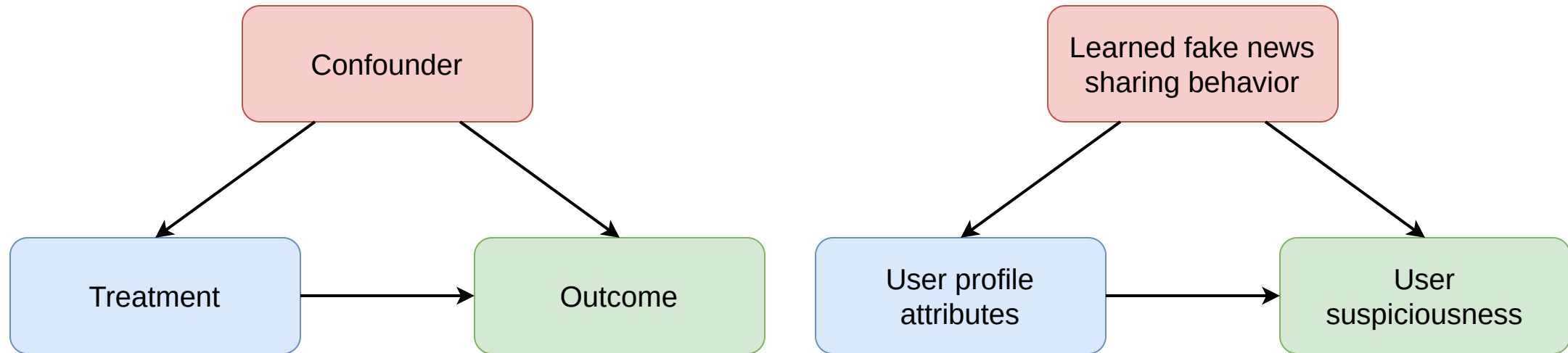


④ (Potentially) causal attributes identification



# Causal Inference

Confounder: Variables that cause spurious associations between treatments and outcome



# Method



# Problem Statement

Users

$$\mathcal{U} = \{1, 2, \dots, u, \dots, U\}$$

Fake News

$$\mathcal{C} = \{1, 2, \dots, i, \dots, N\}$$

# Problem Statement

Interactions between user  $u$  and fake news  $i$

$$Y_{ui} \in \mathcal{Y}$$

$$Y_{ui} = \begin{cases} 1, & \text{if } u \text{ spreads } i \\ 0, & \text{else} \end{cases}$$

$Y_{ui} = 0$  can be interpreted as either  $u$  is not interested in  $i$  or  $u$  did not observe  $i$ .

# Problem Statement

Users have  $m$  profile attributes denoted by

$$\text{matrix } A = (A_1, A_2, \dots, A_m).$$

Susceptibility to spread fake news of user  $u$

$$B \in (0, 1]$$

Each user  $u$  is associated with an outcome  $B$

# Problem Statement

## Tasks

- Fake News Sharing Behavior Learning
  - Model the fake news dissemination process
  - Learn fake news sharing behavior  $U$  under selection biases
- Causal User Attributes Identification
  - Identify user attributes that potentially cause users to spread fake news and estimate the effects.

# Modeling Fake News Dissemination

Interestingness

$$R_{ui} \in \{0, 1\}$$

Exposure

$$O_{ui} \in \{0, 1\}$$

# Modeling Fake News Dissemination

Assume that a user spreads fake news iff s/he is both exposed to and interested in it

$$Y_{ui} = O_{ui} \cdot R_{ui}$$

$$P(Y_{ui} = 1) = P(O_{ui} = 1) \cdot P(R_{ur} = 1) = \theta_{ui} \cdot r_{ui}, \theta_{ui} > 0; r_{ui} > 0; \forall Y_{ui} \in \mathcal{Y}$$

$$\mathcal{D}_{pair} = \mathcal{U} \times \mathcal{C} \times \mathcal{C}$$

is the set of all observed(positive) interactions  $(u, i)$  and all unobserved(negative) interactions  $(u, j)$

# Modeling Fake News Dissemination

Optimizing the pairwise BPR loss

$$\mathcal{L}_{ideal}(\hat{S}) = \frac{1}{|\mathcal{D}_{pair}|} \sum_{(u,i,j) \in \mathcal{D}_{pair}} r_{ui}(1 - r_{uj})\ell(\hat{S}_{uij})$$

where  $\hat{S}_{ui}$  is the difference between the predicted scores of fake news  $i$  and  $j$ , and  $\ell = -\ln(\sigma(\cdot))$  represents the local loss for the triplet  $(u, i, j)$

# Learning Unbiased Sharing Behavior

## Inverse Probability Weighting (IPS)

Reweighting mechanism by assigning larger weights to news that is **less** likely to be observed

## Propensity Score

The propensity score of user  $u$  being exposed to news  $i$  is

$$\theta_{ui} = P(O_{ui} = 1) = P(Y_{ui} = 1 | R_{ui} = 1)$$



# Learning Unbiased Sharing Behavior

## Unbiased estimator

$$\mathcal{L}_{unbiased}(\hat{S}) = \frac{1}{|\mathcal{D}_{pair}|} \sum_{(u,i,j) \in \mathcal{D}_{pair}} \frac{Y_{ui}}{\theta_{ui}} \left(1 - \frac{Y_{uj}}{\theta_{uj}}\right) \ell(\hat{S}_{uij})$$

## Proposition

$$\mathbb{E}[\hat{\mathcal{L}}_{unbiased}(\hat{S})] = \mathcal{L}_{ideal}(\hat{S})$$

## News-based Propensity

$$P_{\text{news}} = \hat{\theta}_{,i}^{\text{news}} = \left( \frac{\sum_{u \in \mathcal{U}} Y_{ui}}{\max_{i \in \mathcal{C}} \sum_{u \in \mathcal{U}} Y_{ui}} \right)^\eta$$

## User-News-based Propensity

$$P_{\text{user}} = \hat{\theta}_{u,i}^{\text{user}} = \left( \frac{\sum_{u \in \mathcal{U}} Y_{ui} \cdot F_u}{\max_{i \in \mathcal{C}} \sum_{u \in \mathcal{U}} Y_{ui} \cdot F_u} \right)^\eta$$

## Neural-Network-based Propensity

$$P_{\text{neural}} = \hat{\theta}_{,i}^{\text{neural}} = \sigma(e_i).$$

# Identifying Causal User Attributes

User  $u$ 's attribute

$$a = (a_1, a_2, \dots, a_m), a \in A$$

User susceptibility

$$B_u \in (0, 1]$$

# Identifying Causal User Attributes

Assume a large portion of news a user has shared is fake, more *susceptible* s/he is to share fake news.

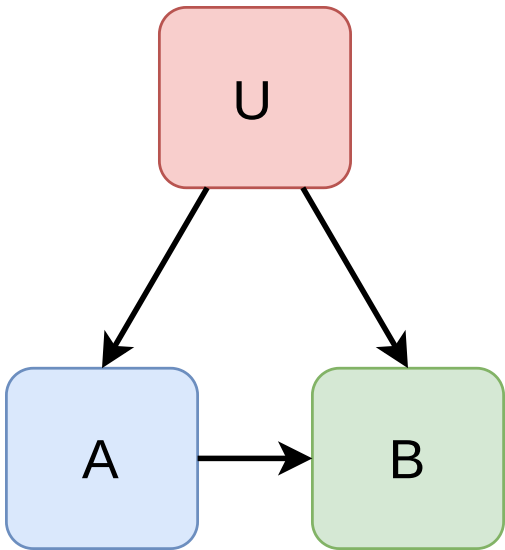
$$B_u = n_{\text{fake}}^u / (n_{\text{fake}}^u + n_{\text{true}}^u)$$

where  $n_{\text{fake}}^u$  is the number of fake news  $u$  has shared.

# Identifying Causal User Attributes

## Causal Model

$$B_u = \beta^\top a_u + \gamma^\top U_u$$



# Experiment

# Dataset

Dataset	#Real	#Fake	#Total	#User
PolitiFact	624	432	1,056	110,127
GossipCop	16,817	5,323	22,140	194,788

- PolitiFact: Fake or real, are provided by journalists and domain experts.
- GossipCop: The fact-checking evaluation results came from the rating scores on the GossipCop website.

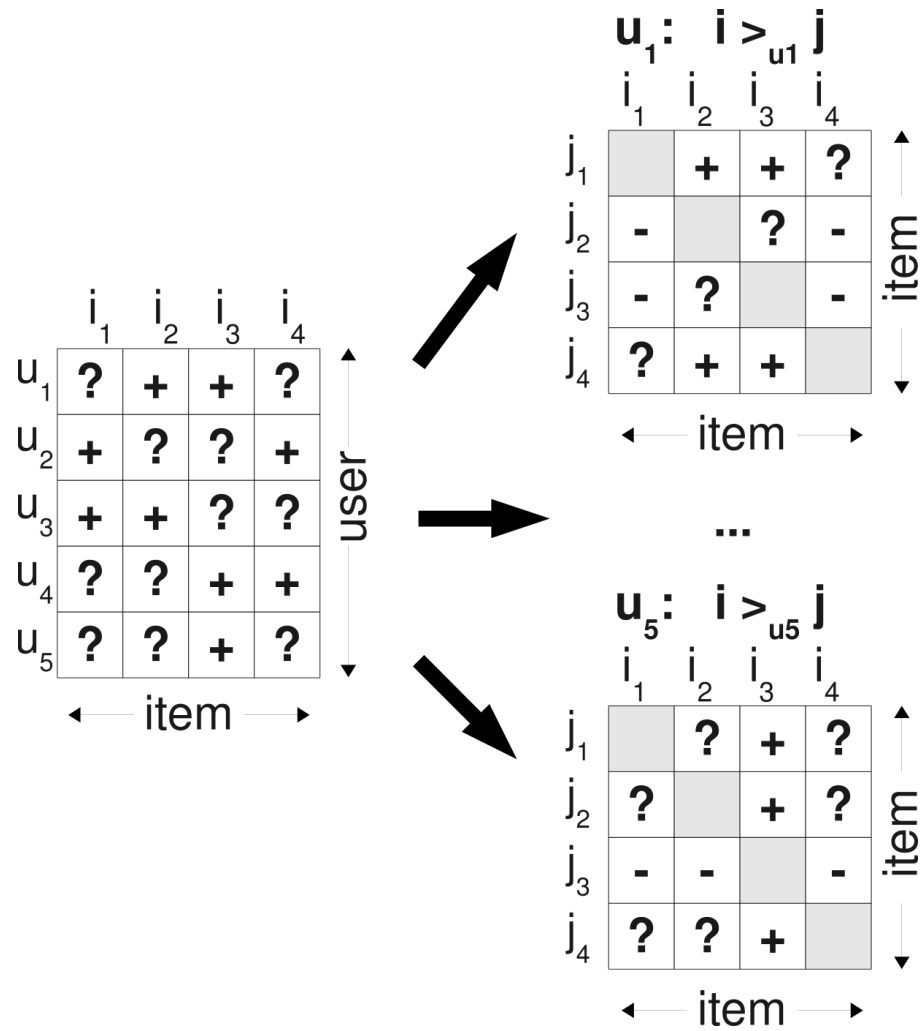
# Baseline

Closely related to recommender system

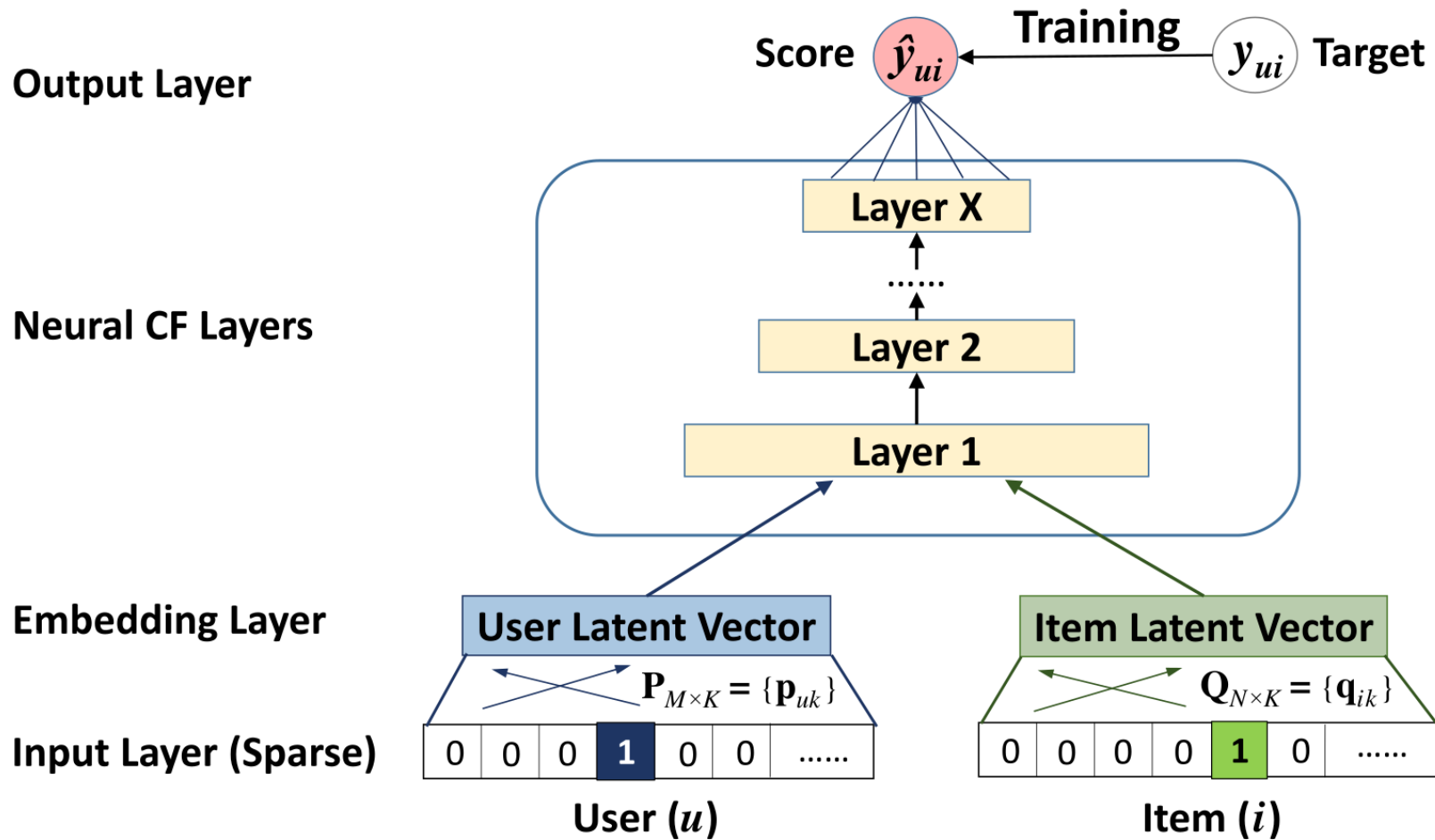
- **BPRMF** (Bayesian Personalized Ranking for Matrix Factorization)
- **NCF** (Neural Collaborative Filtering Model)



# BPRMF



# NCF



# Evaluation Metrics

- Recall@K

$$\text{Recall@K} = \frac{\# \text{ interesting fake news @K.}}{\text{Total \# interesting fake news}}$$

- NDCG@K

$$\text{NDCG@K} = \frac{\text{DCG@K}}{\text{IDCG@K}}, \text{DCG@K} = \sum_{i=1}^K \frac{2^{\text{interest}_i} - 1}{\log_2(i + 1)}$$

# Evaluation on Fake News Dissemination

**Table 2: Performance comparisons w.r.t. predicting fake news to be shared using data *PolitiFact* and base model BPRMF (%).  $p < 0.05$ .**

**(a) Recall@K with K=20,40,60,80.**

K	20	40	60	80
BPRMF	12.36	22.18	31.10	39.51
BPRMF-N	14.45 <sup>↑16.9%</sup>	25.11 <sup>↑13.2%</sup>	34.34 <sup>↑10.4%</sup>	42.72 <sup>↑8.1%</sup>
BPRMF-U	14.78 <sup>↑19.6%</sup>	25.65 <sup>↑15.6%</sup>	34.91 <sup>↑12.2%</sup>	<b>43.63</b> <sup>↑10.4%</sup>
BPRMF-Neu	<b>14.90</b> <sup>↑20.6%</sup>	<b>25.83</b> <sup>↑16.5%</sup>	<b>35.13</b> <sup>↑13.0%</sup>	43.55 <sup>↑10.2%</sup>

**(b) NDCG@K with K=20,40,60,80.**

K	20	40	60	80
BPRMF	5.33	7.51	9.22	10.71
BPRMF-N	6.39 <sup>↑19.9%</sup>	8.73 <sup>↑16.2%</sup>	10.49 <sup>↑13.8%</sup>	11.97 <sup>↑11.8%</sup>
BPRMF-U	<b>6.54</b> <sup>↑22.7%</sup>	8.92 <sup>↑18.8%</sup>	10.69 <sup>↑15.9%</sup>	<b>12.21</b> <sup>↑14.0%</sup>
BPRMF-Neu	6.53 <sup>↑22.5%</sup>	<b>8.93</b> <sup>↑18.9%</sup>	<b>10.71</b> <sup>↑16.2%</sup>	12.19 <sup>↑13.8%</sup>

**Table 3: Performance comparisons w.r.t. predicting fake news to be shared using data *PolitiFact* and base model NCF (%).  $p < 0.05$ .**

**(a) Recall@K with K=20,40,60,80.**

K	20	40	60	80
NCF	9.59	18.45	27.30	36.33
NCF-N	<b>10.42</b> ↑8.7%	<b>19.34</b> ↑4.8%	28.58↑4.7%	37.07↑2.0%
NCF-U	10.29↑7.3%	19.29↑4.6%	27.34↑0.1%	34.87↓4.0%
NCF-Neu	10.20↑6.4%	19.11↑3.6%	<b>28.74</b> ↑5.3%	<b>38.39</b> ↑5.7%

**(b) NDCG@K with K=20,40,60,80.**

K	20	40	60	80
NCF	3.72	5.66	7.35	8.94
NCF-N	4.13↑11.2%	6.09↑7.6%	<b>7.85</b> ↑6.8%	9.36↑4.7%
NCF-U	<b>4.19</b> ↑12.6%	<b>6.18</b> ↑9.2%	7.75↑5.4%	9.10↑1.8%
NCF-Neu	4.04↑8.6%	5.99↑5.8%	7.82↑6.4%	<b>9.52</b> ↑6.5%

**Table 4: Performance comparisons w.r.t. predicting fake news to be shared using data *GossipCop* and base model BPRMF (%).  $p < 0.05$ .**

**(a) Recall@K with K=20,40,60,80.**

K	20	40	60	80
BPRMF	13.31	16.38	18.77	20.8
BPRMF-N	14.92 <sup>↑12.2%</sup>	17.61 <sup>↑7.5%</sup>	19.70 <sup>↑5.0%</sup>	21.52 <sup>↑3.5%</sup>
BPRMF-U	14.97 <sup>↑12.6%</sup>	17.70 <sup>↑8.1%</sup>	19.73 <sup>↑5.1%</sup>	21.58 <sup>↑3.8%</sup>
BPRMF-Neu	<b>15.72</b> <sup>↑18.2%</sup>	<b>18.76</b> <sup>↑14.5%</sup>	<b>21.03</b> <sup>↑12.0%</sup>	<b>22.96</b> <sup>↑10.4%</sup>

**(b) NDCG@K with K=20,40,60,80.**

K	20	40	60	80
BPRMF	10.52	11.32	11.86	12.30
BPRMF-N	12.38 <sup>↑17.7%</sup>	13.11 <sup>↑15.8%</sup>	13.60 <sup>↑14.7%</sup>	13.97 <sup>↑13.6%</sup>
BPRMF-U	12.22 <sup>↑16.2%</sup>	12.95 <sup>↑14.4%</sup>	13.42 <sup>↑13.2%</sup>	13.81 <sup>↑12.3%</sup>
BPRMF-Neu	<b>12.74</b> <sup>↑21.1%</sup>	<b>13.56</b> <sup>↑19.8%</sup>	<b>14.08</b> <sup>↑18.7%</sup>	<b>14.49</b> <sup>↑17.8%</sup>

**Table 5: Performance comparisons w.r.t. predicting fake news to be shared using data *GossipCop* and base model NCF (%).  $p < 0.05$ .**

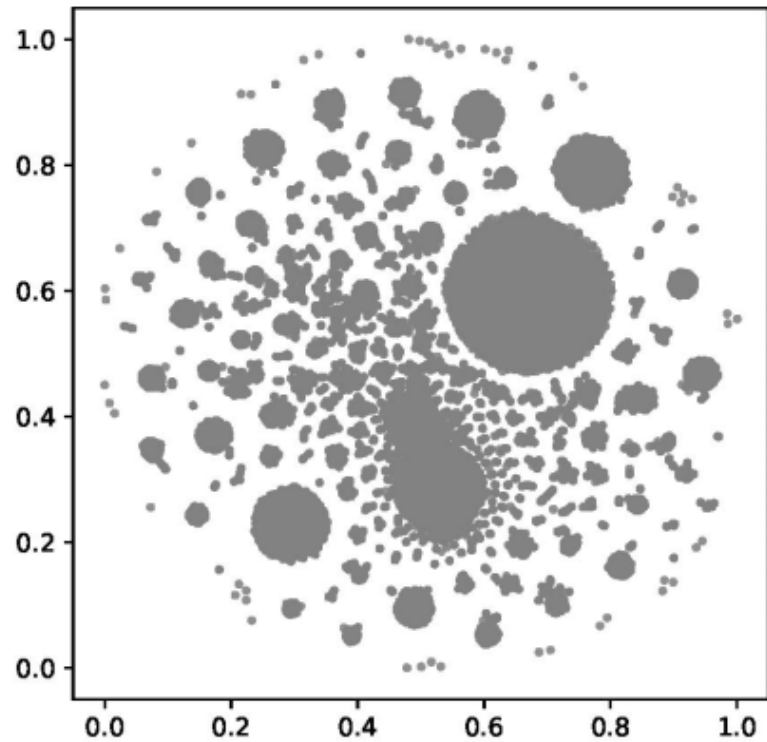
**(a) Recall@K with K=20,40,60,80.**

K	20	40	60	80
NCF	5.87	8.01	9.72	11.63
NCF-N	7.59 <sup>↑29.3%</sup>	9.50 <sup>↑18.6%</sup>	11.22 <sup>↑15.4%</sup>	12.74 <sup>↑9.5%</sup>
NCF-U	<b>8.99</b> <sup>↑53.2%</sup>	<b>10.93</b> <sup>↑36.5%</sup>	<b>12.73</b> <sup>↑31.0%</sup>	<b>14.42</b> <sup>↑24.0%</sup>
NCF-Neu	8.36 <sup>↑42.4%</sup>	10.53 <sup>↑31.5%</sup>	12.39 <sup>↑27.5%</sup>	13.97 <sup>↑20.1%</sup>

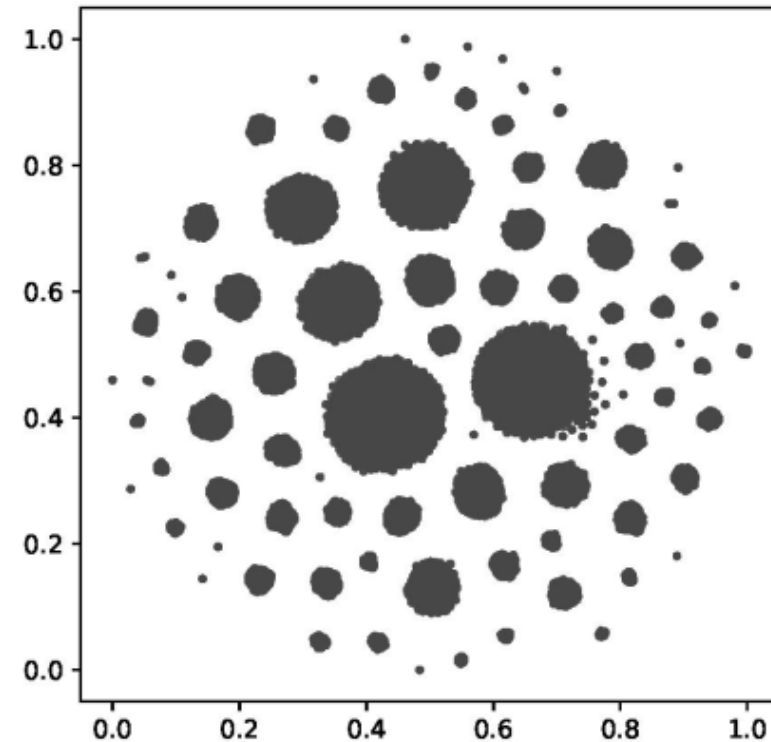
**(b) NDCG@K with K=20,40,60,80.**

K	20	40	60	80
NCF	4.41	4.97	5.37	5.77
NCF-N	5.96 <sup>↑35.1%</sup>	6.50 <sup>↑30.8%</sup>	6.91 <sup>↑28.7%</sup>	7.23 <sup>↑25.3%</sup>
NCF-U	<b>7.36</b> <sup>↑66.9%</sup>	<b>7.91</b> <sup>↑59.2%</sup>	<b>8.33</b> <sup>↑55.1%</sup>	<b>8.68</b> <sup>↑50.4%</sup>
NCF-Neu	6.53 <sup>↑48.1%</sup>	7.14 <sup>↑43.7%</sup>	7.57 <sup>↑41.0%</sup>	7.91 <sup>↑37.1%</sup>





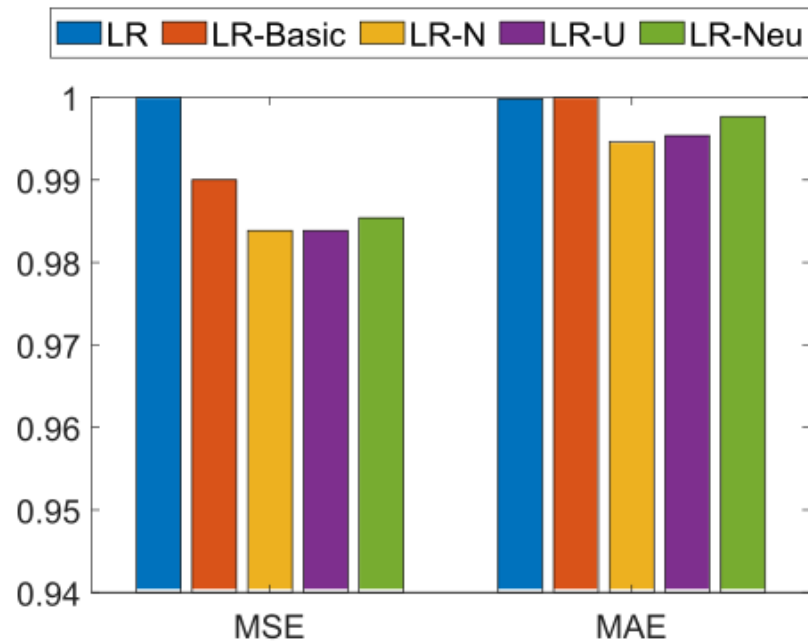
**(a) Fake news sharing behavior.**  
**Silhouette Coefficient=-0.124**



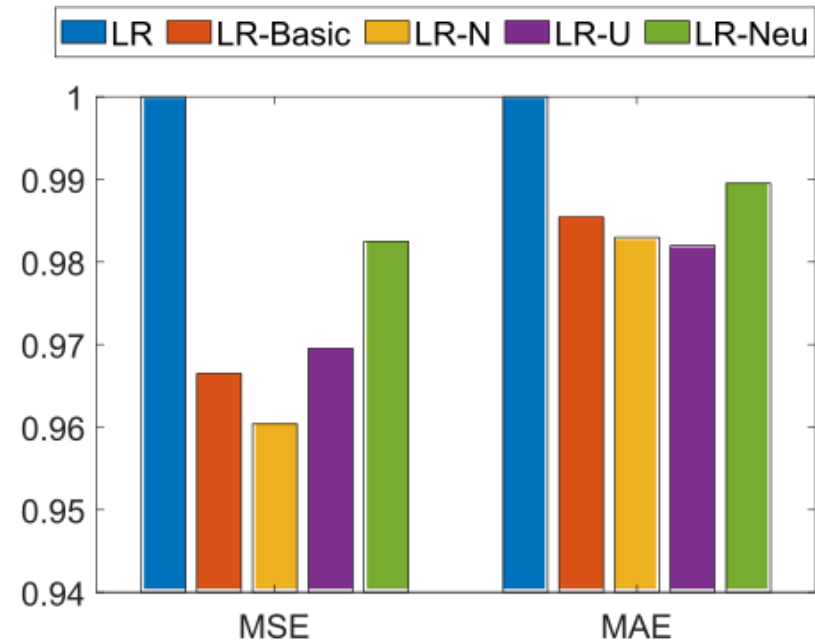
**(b) True news sharing behavior.**  
**Silhouette Coefficient=0.903**

**Figure 2: Behavior comparisons using 2-D t-SNE visualizations.**

# Evaluation on Identifying Causal User

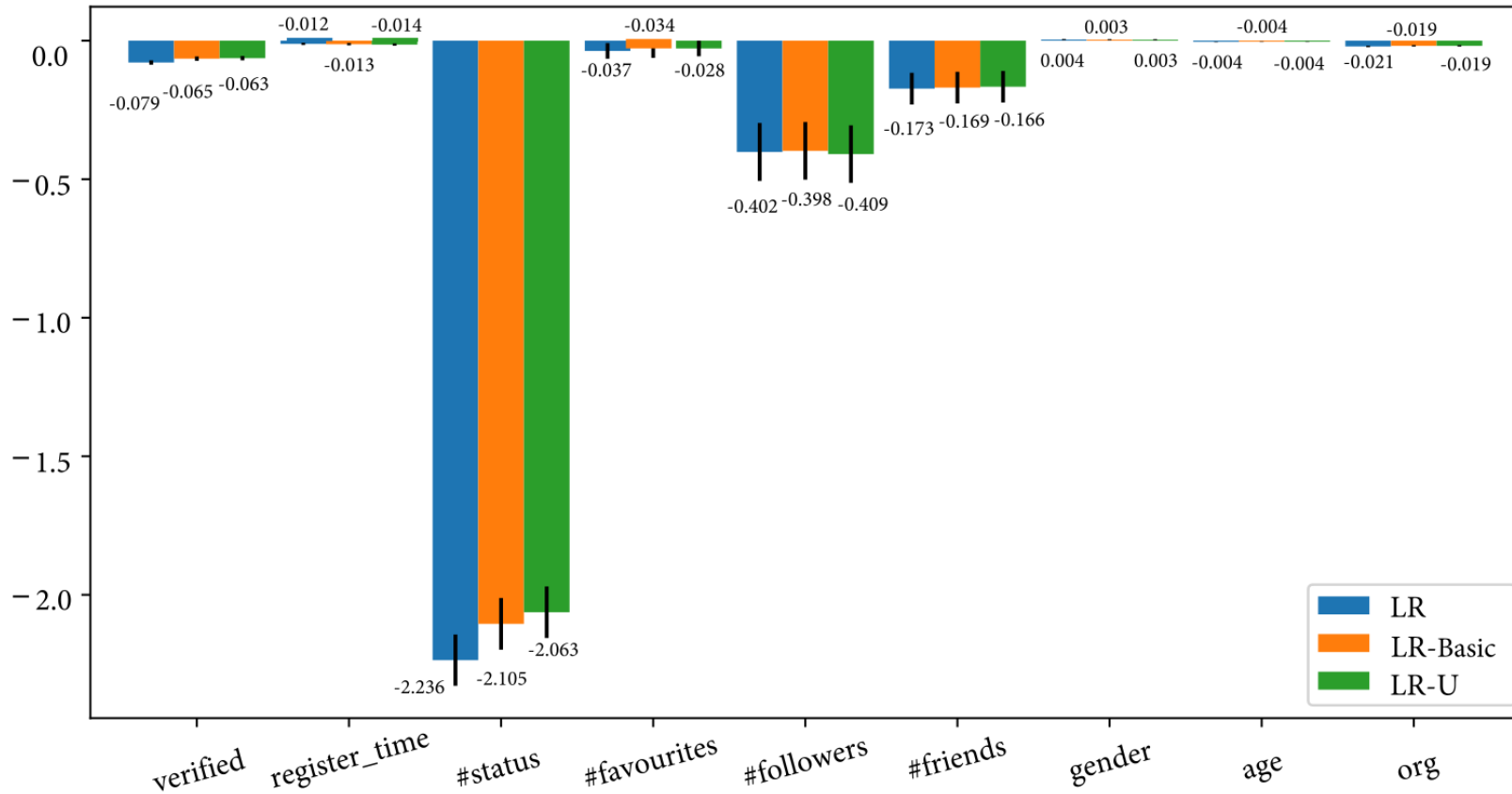


(a) *PolitiFact*.

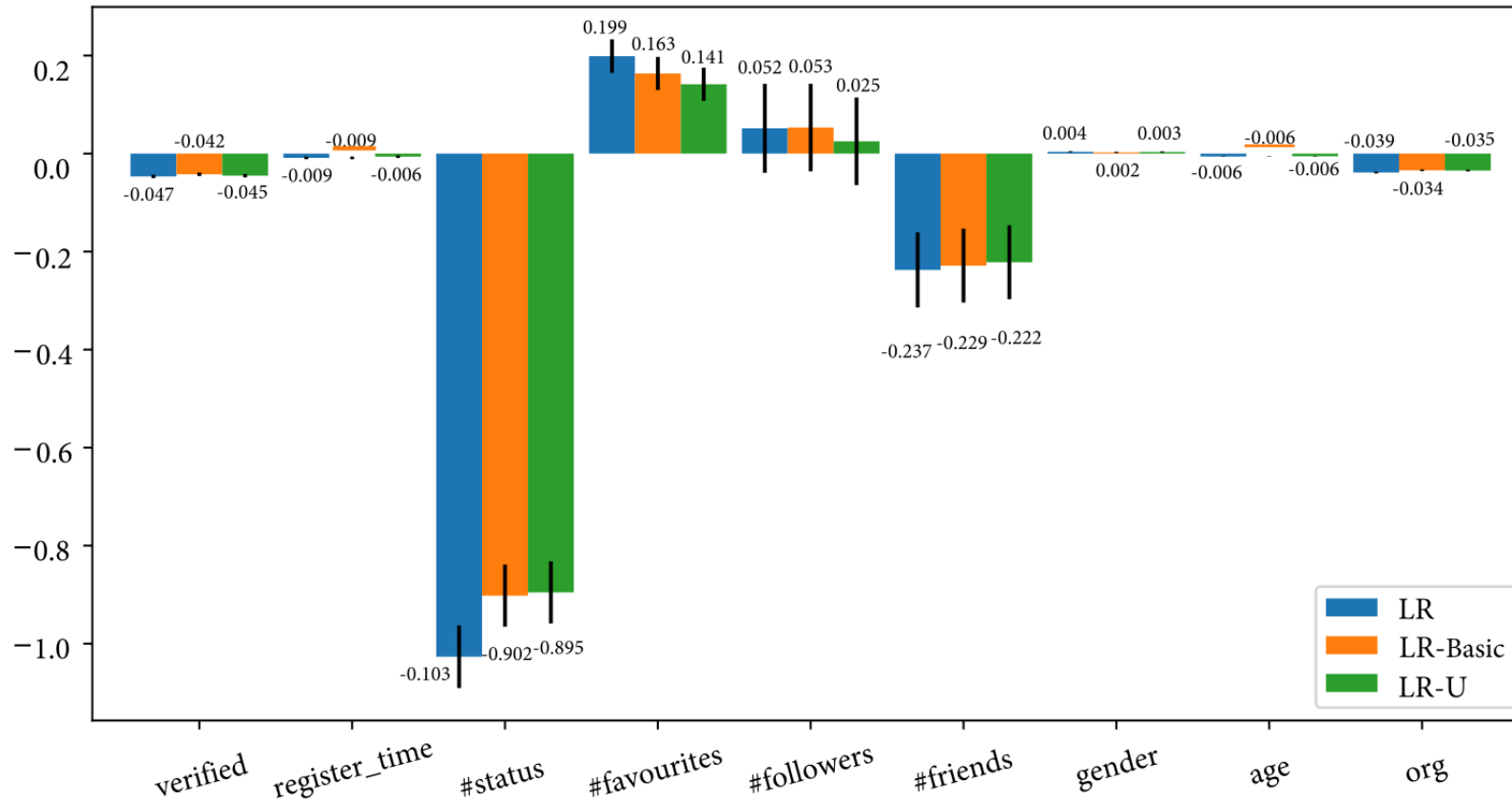


(b) *GossipCop*.

**Figure 3: Performance comparisons w.r.t. predicting user susceptibility using both datasets.  $y$ -axis denotes relative results.**



**Figure 4: *PolitiFact*: Effects comparisons w.r.t. each potential causal user attribute. All the results are statistically significant.**



**Figure 5: *GossipCop*: Effects comparisons w.r.t. each potential causal user attribute. All the results except for that of *#followers* are statistically significant.**

# Conclusion

1. IPS-weighted models can learn unbiased embeddings of fake news sharing behavior that lead to more accurate predictions of fake news that users will share and user susceptibility
2. The identified causal attributes show that ***verified, statuses count, friends count, and org*** relate significantly with user susceptibility to share fake news.

# References

- <https://arxiv.org/pdf/2010.10580.pdf>
- <https://github.com/GitHubLuCheng/Causal-Understanding-of-Fake-News-Dissemination>
- <https://towardsdatascience.com/causality-an-introduction-f8a3f6ac4c4a>
- <https://towardsdatascience.com/causal-inference-962ae97cefda>
- <https://towardsdatascience.com/causal-discovery-6858f9af6dcb>
- <https://towardsdatascience.com/propensity-scores-and-inverse-probability-weighting-in-causal-inference-97aa53f3b6ce>