

Adversarial Attacks for Recommendations

Wenqi Fan

The Hong Kong Polytechnic University

<https://wenqifan03.github.io>, wenqifan@polyu.edu.hk

Tutorial website: <https://advanced-recommender-systems.github.io/ijcai2021-tutorial/>



Adversarial Attacks on Deep Learning



Classified as panda

Small adversarial noise

Classified as gibbon

x

ϵ

x'

Attacks can happen in Recommender Systems



NEWS Menu

Business | Market Data | New Economy | New Tech Economy | Companies | Entrepreneurship | Technology of Business | Business of Sport | Global Education | Economy | Global Car Industry

Amazon 'flooded by fake five-star reviews' - Which? report

16 April 2019

GOV.UK Menu

Coronavirus (COVID-19) Guidance and support

Home > Competition

Press release

Facebook and eBay pledge to combat trading in fake reviews

Following action from the CMA, Facebook and eBay have committed to combatting the trade of fake and misleading reviews on their sites.

From: [Competition and Markets Authority](#)

Published 8 January 2020

“More than three-quarters of people are influenced by reviews when they shop online.”

Understand how attacks can be performed

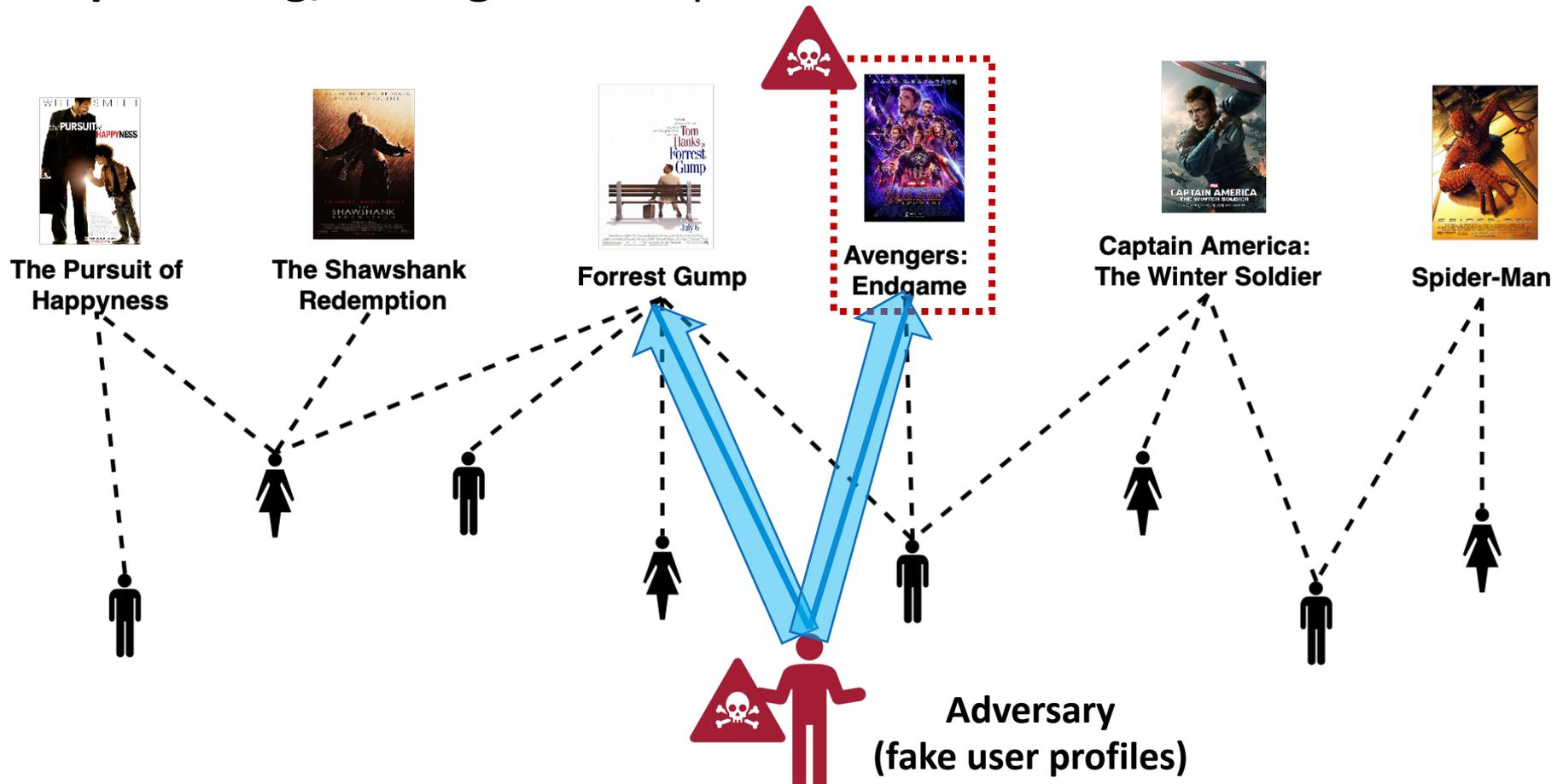


Defend against potential adversarial attacks

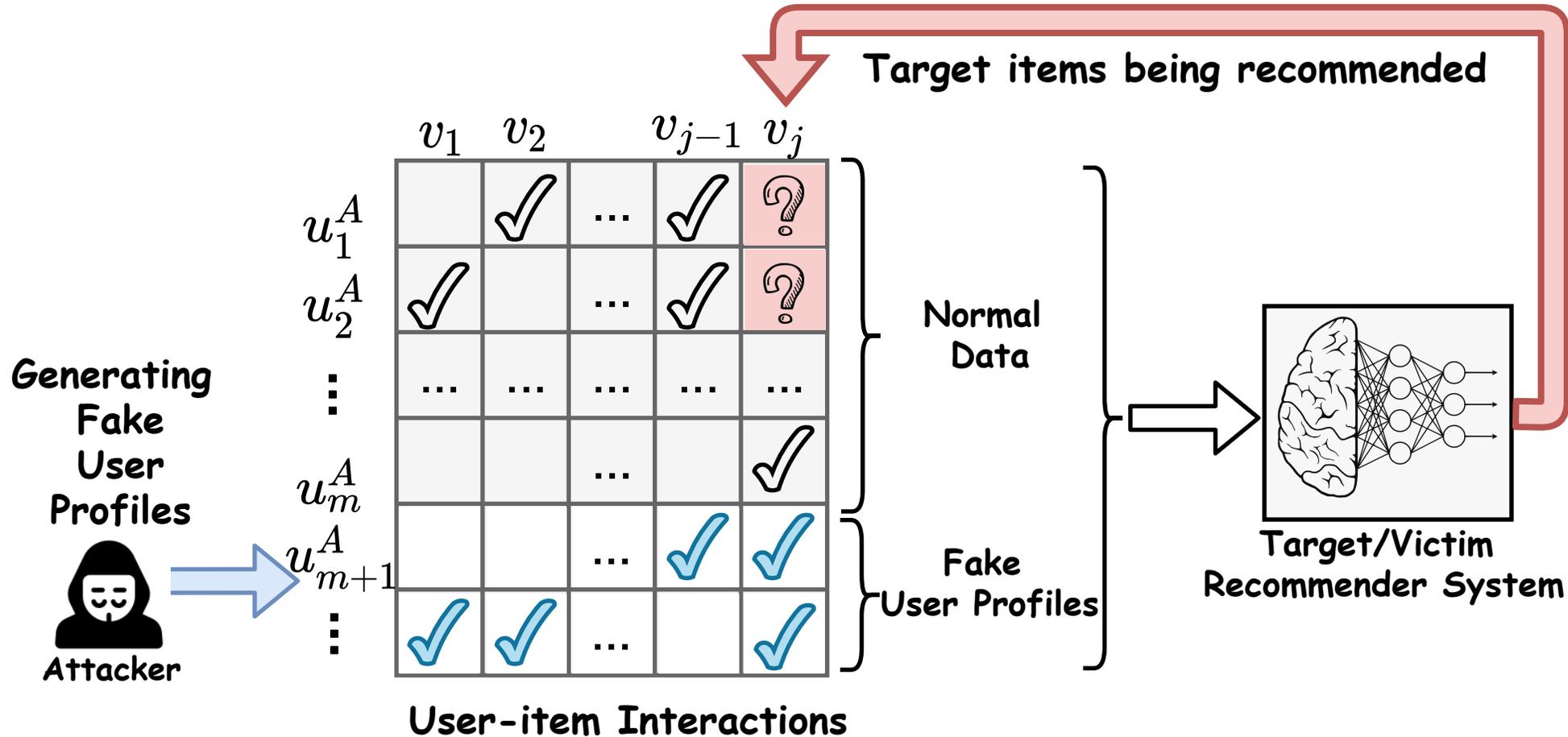


Attacks can happen in Recommender Systems

- Security (Attacking) in Recommender Systems
 - **Data poisoning/shilling attacks:** promote/demote a set of items



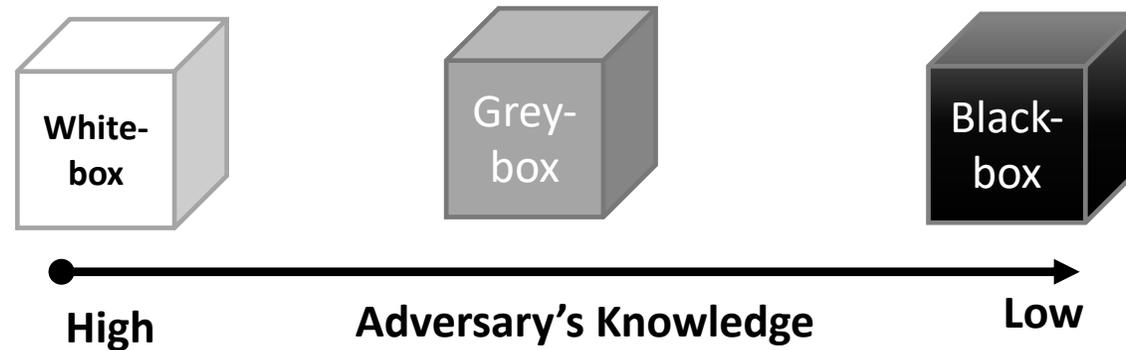
A General Attacking Framework



Attack settings

White/grey-box attacks vs. Black-box attacks.

- have **full/partial knowledge** of the victim model/have **no knowledge**.



Targeted Attacks vs. Non-Targeted Attacks.

- attack **specific target** items / hurt the overall recommendation performance.

Adversarial Attacks



■ White-box Attacks

- Data Poisoning Attacks on Factorization-Based Collaborative Filtering (NIPS'16)

■ Grey-box Attacks

- Revisiting Adversarially Learned Injection Attacks Against Recommender Systems (RecSys'20)
- Adversarial Attacks on an Oblivious Recommender (RecSys'19)

■ Black-box Attacks

- CopyAttack: Attacking Black-box Recommendations via Copying Cross-domain User Profiles (ICDE'21)
- PoisonRec: An Adaptive Data Poisoning Framework for Attacking Black-box Recommender Systems (ICDE'20)

Adversarial Attacks



■ White-box Attacks

- **Data Poisoning Attacks on Factorization-Based Collaborative Filtering (NIPS'16)**

■ Grey-box Attacks

- Revisiting Adversarially Learned Injection Attacks Against Recommender Systems (RecSys'20)
- Adversarial Attacks on an Oblivious Recommender (RecSys'19)

■ Black-box Attacks

- CopyAttack: Attacking Black-box Recommendations via Copying Cross-domain User Profiles (ICDE'21)
- PoisonRec: An Adaptive Data Poisoning Framework for Attacking Black-box Recommender Systems (ICDE'20)

Preliminaries



- Collaborative Filtering:

- Given data. $\mathbf{M} \in \mathbb{R}^{m \times n}$, $\Omega = \{(i, j) : \mathbf{M}_{ij} \text{ is observed}\}$
- Goal: matrix completion

$$\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \|\mathcal{R}_\Omega(\mathbf{M} - \mathbf{X})\|_F^2, \quad s.t. \text{rank}(\mathbf{X}) \leq k.$$

- Alternating minimization:

$$\min_{\mathbf{U} \in \mathbb{R}^{m \times k}, \mathbf{V} \in \mathbb{R}^{n \times k}} \{ \|\mathcal{R}_\Omega(\mathbf{M} - \mathbf{UV}^\top)\|_F^2 + 2\lambda_U \|\mathbf{U}\|_F^2 + 2\lambda_V \|\mathbf{V}\|_F^2 \}$$

Attacking Formulation



- Inject malicious users $\widetilde{\mathbf{M}} \in \mathbb{R}^{m' \times n}$

- The CF formulations will be:

$$\Theta_{\lambda}(\widetilde{\mathbf{M}}; \mathbf{M}) = \arg \min_{\mathbf{U}, \widetilde{\mathbf{U}}, \mathbf{V}} \|\mathcal{R}_{\Omega}(\mathbf{M} - \mathbf{U}\mathbf{V}^{\top})\|_F^2 + \|\mathcal{R}_{\widetilde{\Omega}}(\widetilde{\mathbf{M}} - \widetilde{\mathbf{U}}\mathbf{V}^{\top})\|_F^2 + 2\lambda_U(\|\mathbf{U}\|_F^2 + \|\widetilde{\mathbf{U}}\|_F^2) + 2\lambda_V\|\mathbf{V}\|_F^2$$

- Goal : $\widehat{\mathbf{M}}^* \in \operatorname{argmax}_{\widetilde{\mathbf{M}} \in \mathbb{M}} R(\widehat{\mathbf{M}}(\Theta_{\lambda}(\widetilde{\mathbf{M}}; \mathbf{M})), \mathbf{M})$

- Solution: Projected gradient ascent (PGA)

$$\widehat{\mathbf{M}}^{(t+1)} = \operatorname{Proj}_{\mathbb{M}} \left(\widehat{\mathbf{M}}^{(t)} + s_t \cdot \nabla_{\widetilde{\mathbf{M}}} R(\widehat{\mathbf{M}}, \mathbf{M}) \right)$$

Adversarial Attacks



■ White-box Attacks

- Data Poisoning Attacks on Factorization-Based Collaborative Filtering (NIPS'16)

■ Grey-box Attacks

- **Revisiting Adversarially Learned Injection Attacks Against Recommender Systems (RecSys'20)**
- Adversarial Attacks on an Oblivious Recommender (RecSys'19)

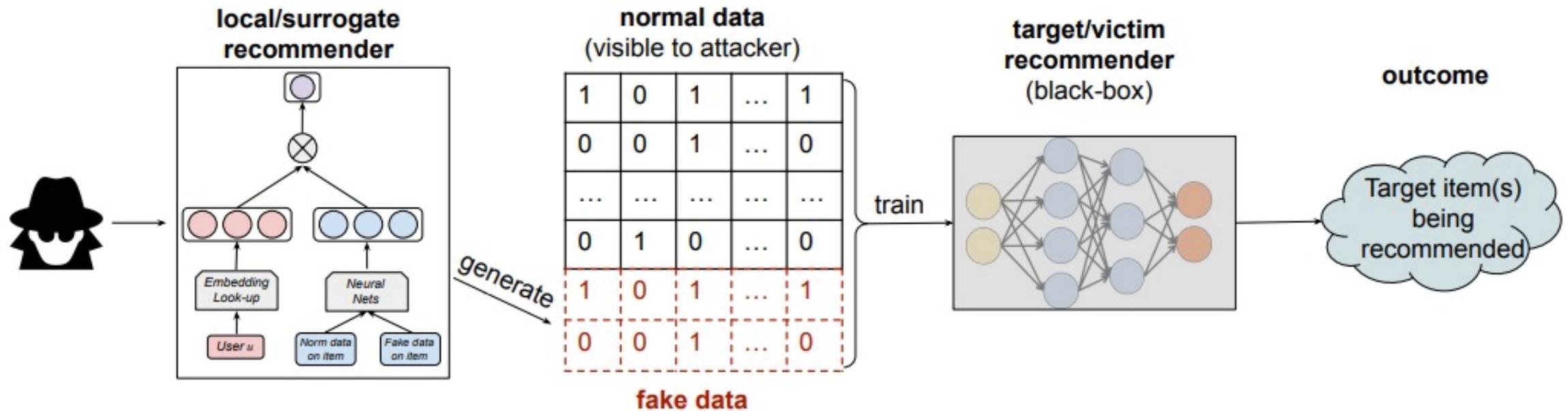
■ Black-box Attacks

- CopyAttack: Attacking Black-box Recommendations via Copying Cross-domain User Profiles (ICDE'21)
- PoisonRec: An Adaptive Data Poisoning Framework for Attacking Black-box Recommender Systems (ICDE'20)

Threat Model

Attacker's Goal: promote certain items availability of being recommended

Attacker's knowledge: fully (partial) observable dataset



How to attack a RecSys: A bi-level optimization problem



- Step 1: Train surrogate model

normal data
(visible to attacker)

1	0	1	...	1
0	0	1	...	0
...
0	1	0	...	0

fake data

1	0	1	...	1
0	0	1	...	0

Training Recommender System

$$\theta^* = \arg \min_{\theta} (\mathcal{L}_{\text{train}}(\mathbf{X}, \mathbf{R}_{\theta}) + \mathcal{L}_{\text{train}}(\hat{\mathbf{X}}, \hat{\mathbf{R}}_{\theta}))$$

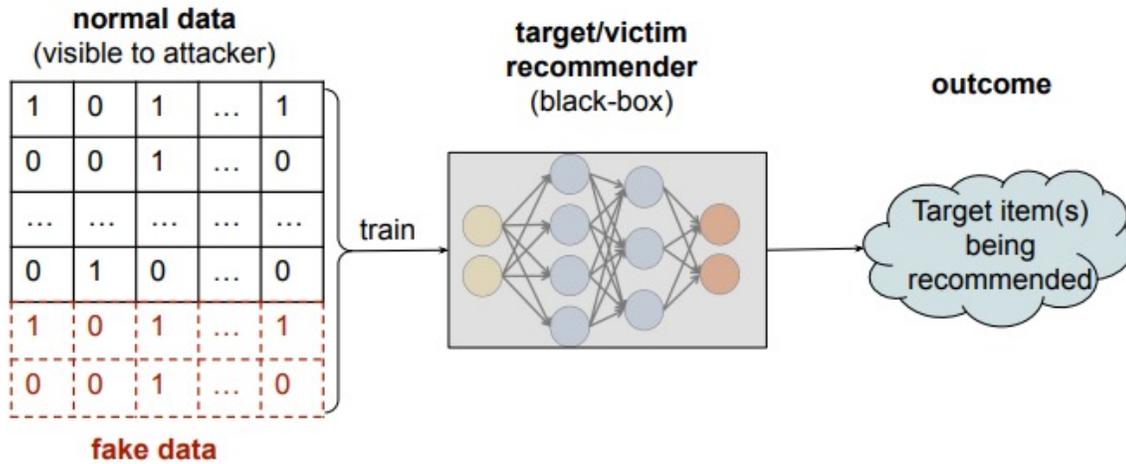
Normal Data Injected fake data

where $\hat{\mathbf{X}}$ is the fake rating matrix, θ^* is the parameters of the surrogate model

How to attack a RecSys: A bi-level optimization problem



- Step 2: Evaluate the malicious goal after fake data are consumed



Adversarial objective
(defined on prediction on normal data)

$$\min_{\hat{X}} \mathcal{L}_{\text{adv}}(\mathbf{R}_{\theta^*}) = - \sum_{u \in \mathcal{U}} \log \left(\frac{\exp(r_{uk})}{\sum_{i \in \mathcal{I}} \exp(r_{ui})} \right).$$

subject to $\theta^* = \arg \min_{\theta} (\mathcal{L}_{\text{train}}(\mathbf{X}, \mathbf{R}_{\theta}) + \mathcal{L}_{\text{train}}(\hat{\mathbf{X}}, \hat{\mathbf{R}}_{\theta}))$

Well-trained surrogate model parameters

Solving the bi-level optimization: gradient-based



Algorithm 1 Learning fake user data with Gradient Descent

- 1: **Input:** Normal user data \mathbf{X} ; learning rate for inner and outer objective: α and η ; max iteration for inner and outer objective: L and T .
 - 2: **Output:** Learned fake user data for malicious goal.
 - 3: Initialize fake data $\widehat{\mathbf{X}}^{(0)}$ and surrogate model parameters $\theta^{(0)}$
 - 4: **for** $t = 1$ to T **do**
 - 5: **for** $l = 1$ to L **do**
 - 6: Optimize inner objective with SGD: $\theta^{(l)} \leftarrow \theta^{(l-1)} - \alpha \cdot \nabla_{\theta} (\mathcal{L}_{\text{train}}(\mathbf{X}, \mathbf{R}_{\theta^{(l-1)}}) + \mathcal{L}_{\text{train}}(\widehat{\mathbf{X}}^{(t)}, \widehat{\mathbf{R}}_{\theta^{(l-1)}}))$
 - 7: **end for**
 - 8: Evaluate $\mathcal{L}_{\text{adv}}(\mathbf{R}_{\theta^{(L)}})$ and compute gradients $\nabla_{\widehat{\mathbf{X}}} \mathcal{L}_{\text{adv}}$
 - 9: Update fake data: $\widehat{\mathbf{X}}^{(t)} = \text{Proj}_{\Lambda}(\widehat{\mathbf{X}}^{(t-1)} - \eta \cdot \nabla_{\widehat{\mathbf{X}}} \mathcal{L}_{\text{adv}})$
 - 10: **end for**
 - 11: **Return:** $\widehat{\mathbf{X}}^{(T)}$
-

Train surrogate model
based on new fake data
Obtain gradient and
update fake data

Repeat until
converge

Limitations



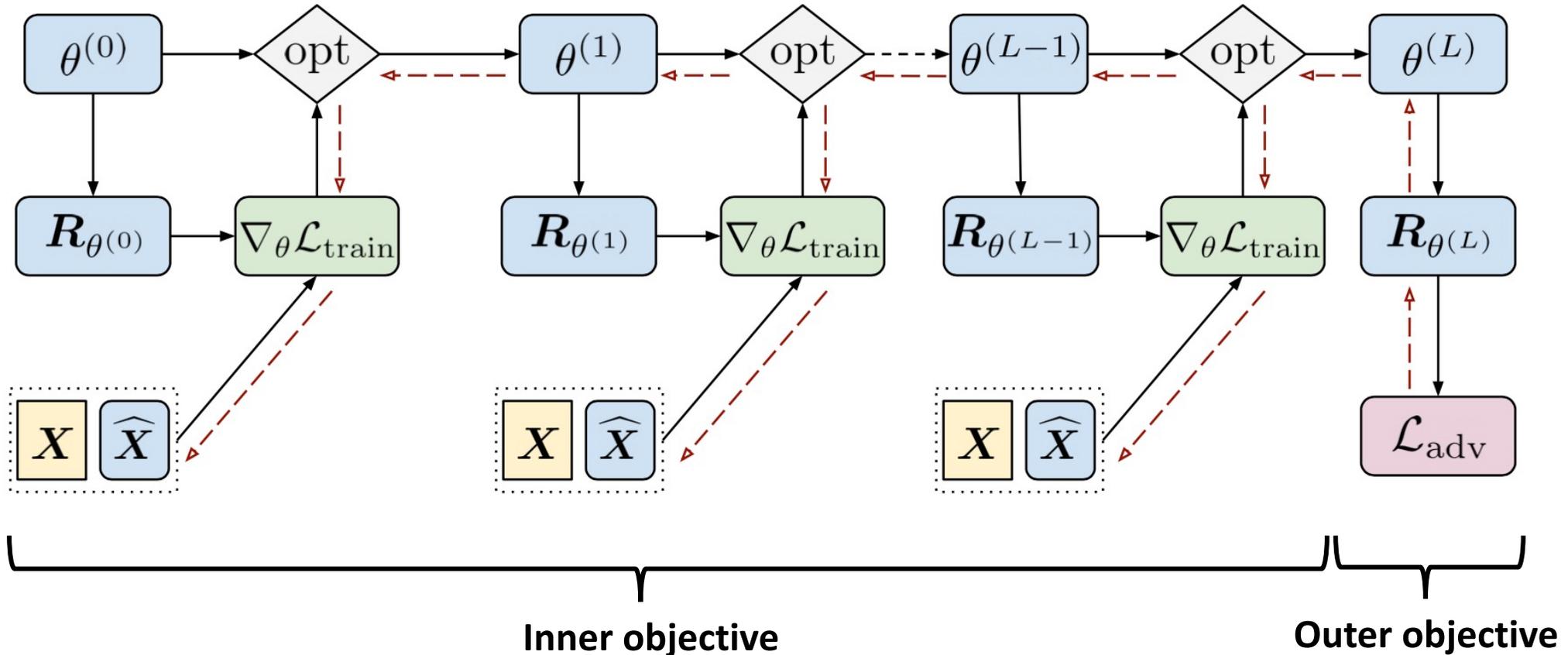
How to obtain the desired gradients $\nabla_{\hat{X}} \mathcal{L}_{\text{adv}}$

Lacking exactness in gradient computation

$$\nabla_{\hat{X}} \mathcal{L}_{\text{adv}} = \frac{\partial \mathcal{L}_{\text{adv}}}{\partial \hat{X}} + \underbrace{\frac{\partial \mathcal{L}_{\text{adv}}}{\partial \theta^*} \cdot \frac{\partial \theta^*}{\partial \hat{X}}}_{\text{ignored}}$$

Computational graph

Exact Solution



Adversarial Attacks



■ White-box Attacks

- Data Poisoning Attacks on Factorization-Based Collaborative Filtering (NIPS'16)

■ Grey-box Attacks

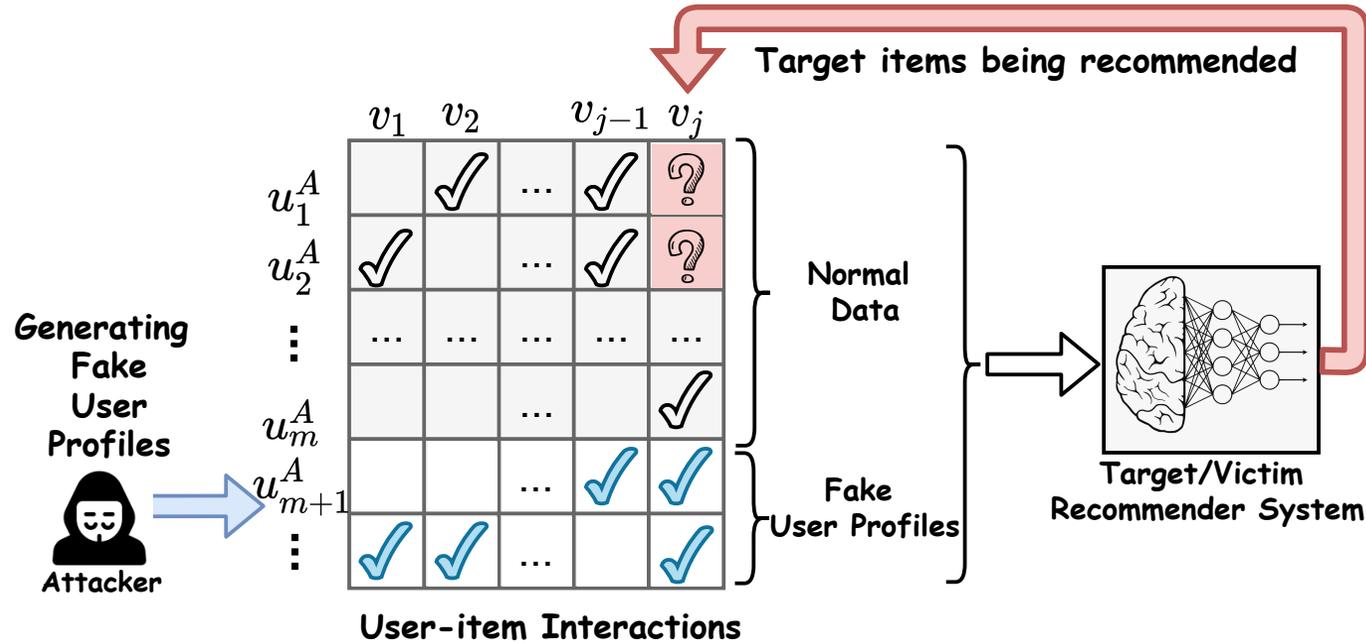
- Revisiting Adversarially Learned Injection Attacks Against Recommender Systems (RecSys'20)
- Adversarial Attacks on an Oblivious Recommender (RecSys'19)

■ Black-box Attacks

- **CopyAttack: Attacking Black-box Recommendations via Copying Cross-domain User Profiles (ICDE'21)**
- PoisonRec: An Adaptive Data Poisoning Framework for Attacking Black-box Recommender Systems (ICDE'20)

Challenges

- Challenges in existing attacking methods:
 - Less "realistic" user profiles (easily detected)



Solution

■ Cross-domain Information

- Share a lot of items
- Users from these platforms with **similar functionalities** also share similar behavior patterns/preferences.



The image displays four e-commerce platforms: Taobao.com, JD.com, Amazon, and Best Buy. Each platform shows listings for the iPhone 12 series. The Taobao and JD.com listings are on the left, Amazon is in the middle, and Best Buy is on the right. A large 'VS' graphic with lightning bolts is placed between the Chinese and American sites, indicating a comparison or cross-domain analysis.

Taobao listings include:

- iPhone 12 官方正品 苹果新品12 iph one12全新5G手机 (¥5699)
- Apple/苹果 iPhone (¥6799)
- 【限时享6期免息】Apple, 苹果 iPhone 12 全网通5G新品 (¥6799)

JD.com listings include:

- Apple iPhone 12 (128GB, 6.1英寸) (¥6799.00)
- Apple iPhone 12 Pro Max (256GB, 6.7英寸) (¥10099.00)

Amazon listings include:

- iPhone 12. Hello 5G. (New Apple iPhone 12 (256GB, Green) Locked + Carrier Subscription, \$31)
- New Apple iPhone 12 (128GB, Blue) Locked + Carrier Subscription (prime, \$31)

Best Buy listings include:

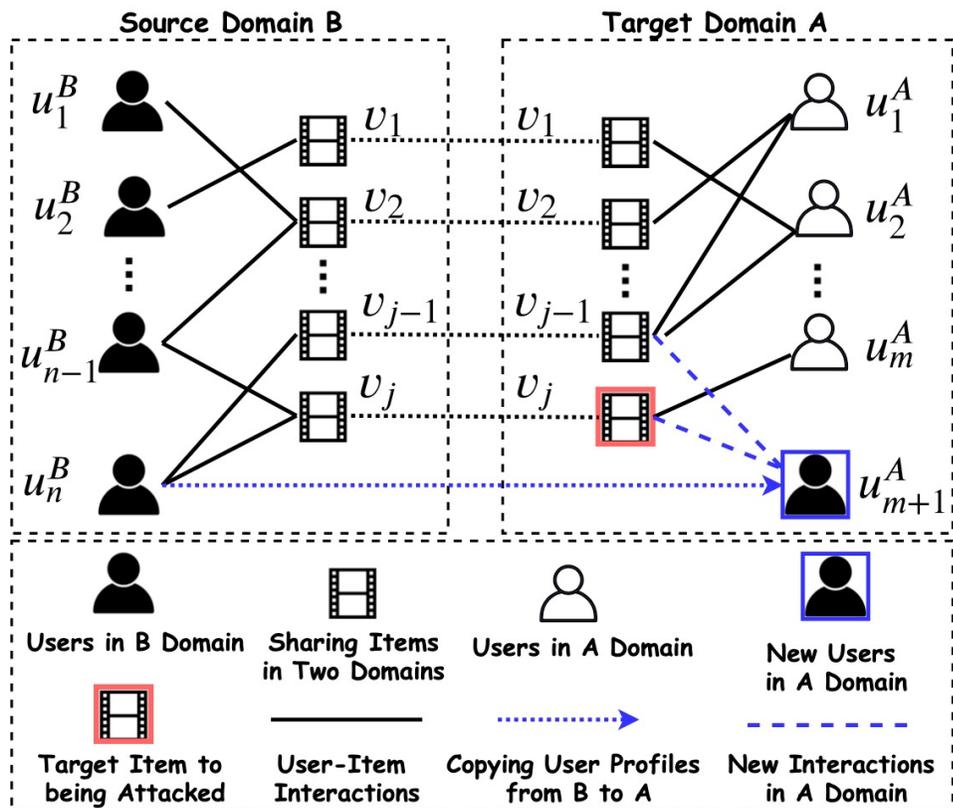
- iPhone 12 mini, iPhone 12, iPhone 12 Pro Max (Now available)

Solution

Challenges in existing attacking methods:

- Less "realistic" user profiles (easily detected)

 **Copy cross-domain users with real profiles from other domains**



Challenges



- **Challenges in existing attacking methods:**

- **Less "realistic" user profiles (easily detected)**

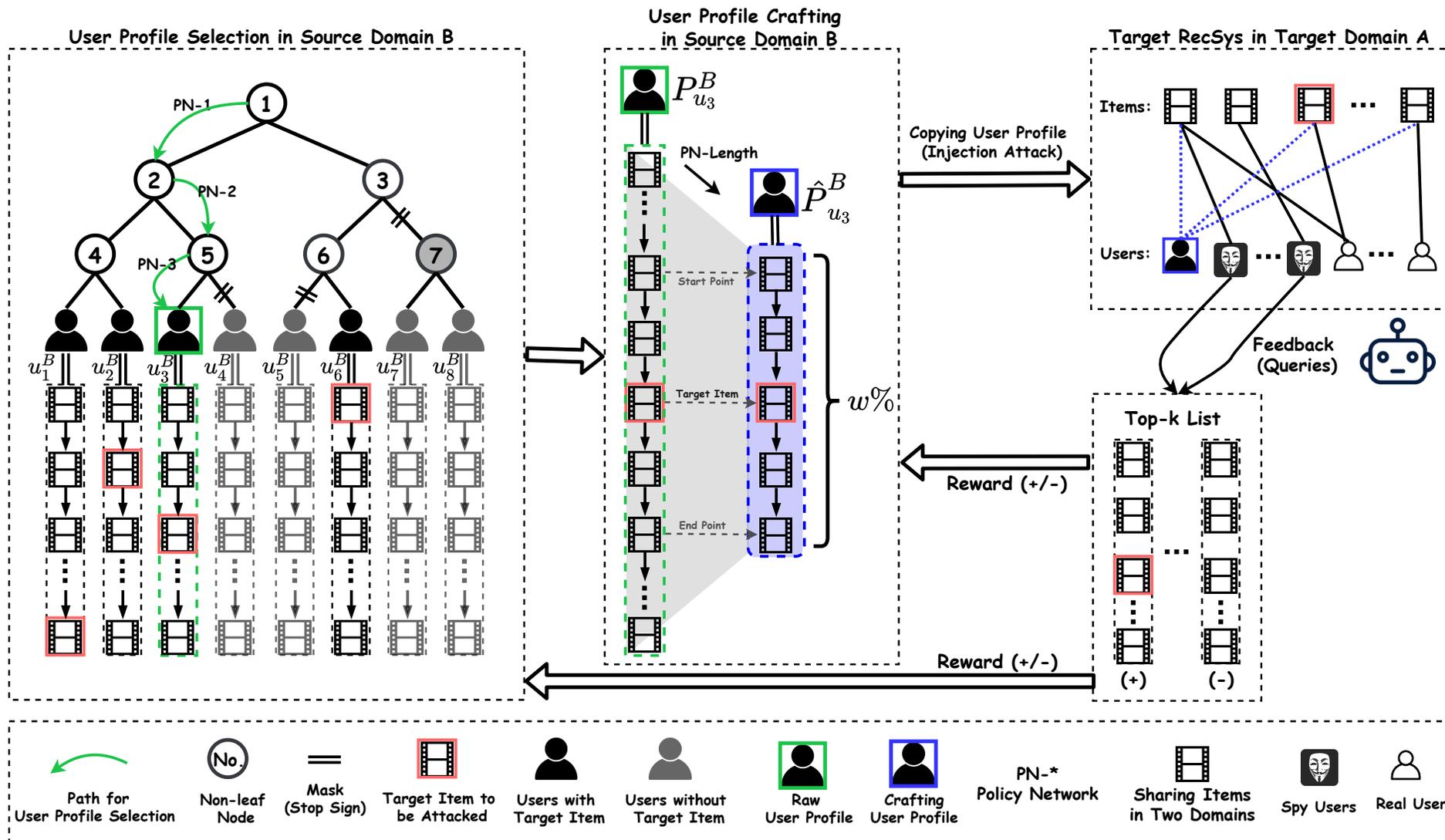
-  Cross-domain Information

- **White/Grey-box setting** (i.e., model architecture and parameters, and datasets)
 - impossible and unrealistic (**privacy and security**)

- **Black-box setting**

-  ➤ Reinforcement Learning (RL) -- Query Feedback (Reward)

CopyAttack



User Profile Selection

- User Profile Selection

- Construct hierarchical clustering tree
- **Masking** Mechanism - specific target items
- Hierarchical-structure Policy Gradient

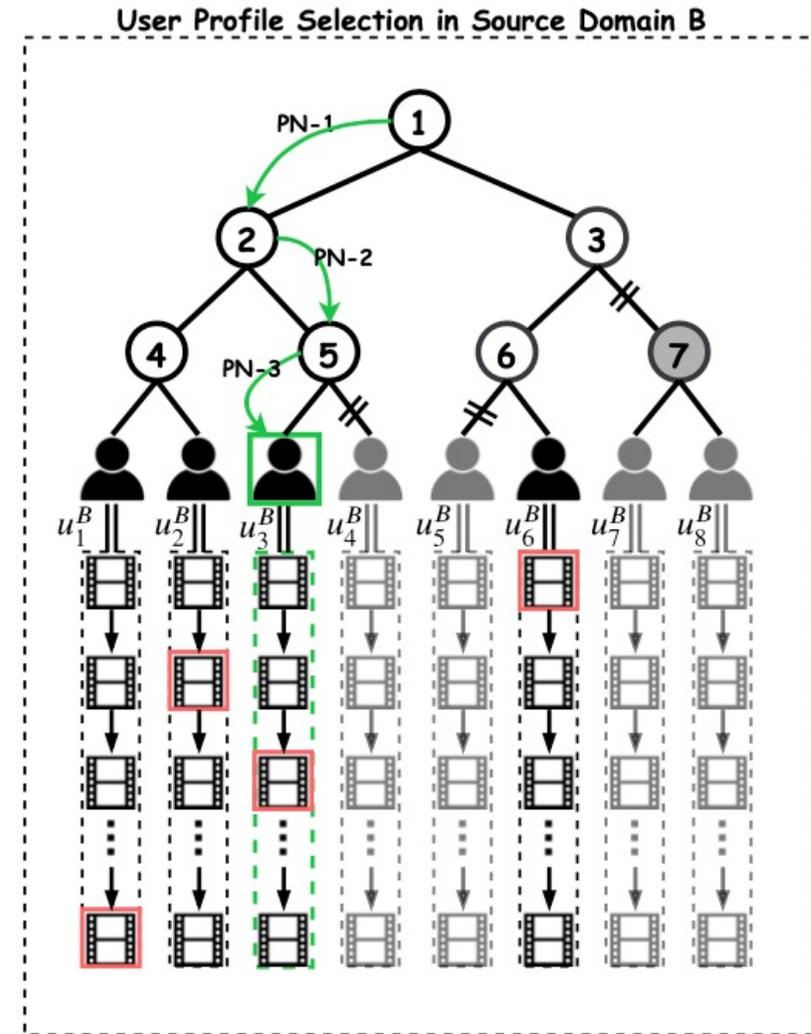
$$a_t^u = \{a_{[t,1]}^u, a_{[t,2]}^u, \dots, a_{[t,d]}^u\}$$

$$\begin{aligned} p^u(a_t^u | s_t^u) &= \prod_{d=1}^d p_d^u(a_{[t,d]}^u | \cdot, s_t^u) \\ &= p_d^u(a_{[t,d]}^u | s_t^u) \cdot p_{d-1}^u(a_{[t,d-1]}^u | s_t^u) \cdots p_1^u(a_{[t,1]}^u | s_t^u) \end{aligned}$$

$$\mathbf{x}_{v_*} = RNN(\mathcal{U}_t^{B \rightarrow A}),$$

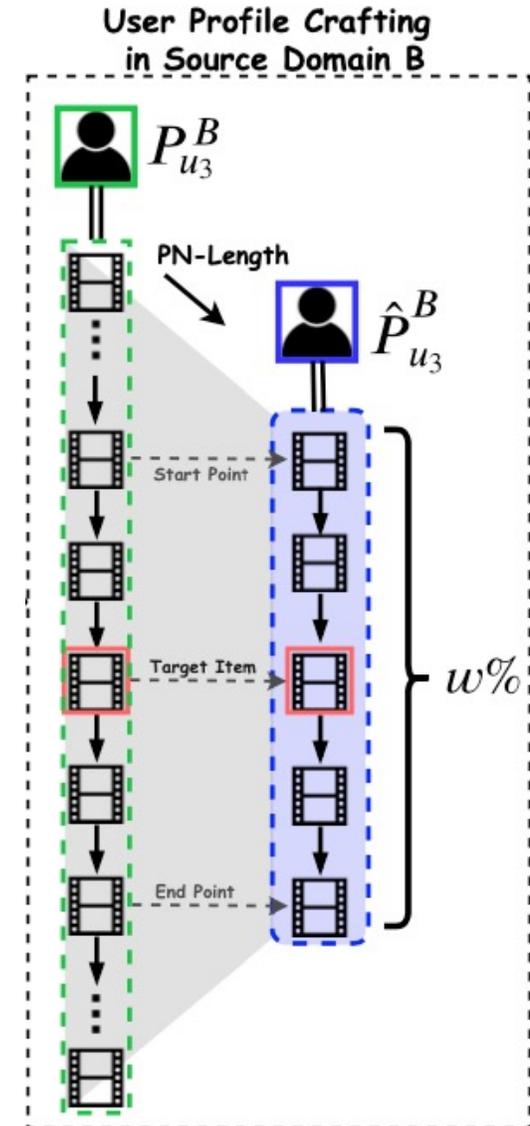
$$p_i^u(\cdot | s_t^u) = \text{softmax}(MLP([\mathbf{q}_{v_*}^B \oplus \mathbf{x}_{v_*}] | \theta_i^u))$$

Time Complexity : $\mathcal{O}(|\mathcal{U}^B|) \longrightarrow \mathcal{O}(d \times |\mathcal{U}^B|^{1/d})$



User Profile Crafting

- User Profile Crafting
 - Clipping operation to craft the raw user profiles
 - $W = \{10\%, 20\%, 30\%, 40\%, 50\%, 60\%, 70\%, 80\%, 90\%, 100\%\}$
 - Sequential patterns (forward/backward)



Example:

$$P_{u_i}^B = \{v_1 \rightarrow v_2 \rightarrow v_3 \rightarrow v_4 \rightarrow v_{5*} \rightarrow v_6 \rightarrow v_7 \rightarrow v_8 \rightarrow v_9 \rightarrow v_{10}\}$$

w = 50%

$$\hat{P}_{u_i}^B = \{v_3 \rightarrow v_4 \rightarrow v_{5*} \rightarrow v_6 \rightarrow v_7\}$$

$$p^l(\cdot | s_t^l) = \text{softmax}(\text{MLP}([\mathbf{p}_i^B \oplus \mathbf{q}_{v_*}^B] | \theta^l))$$

