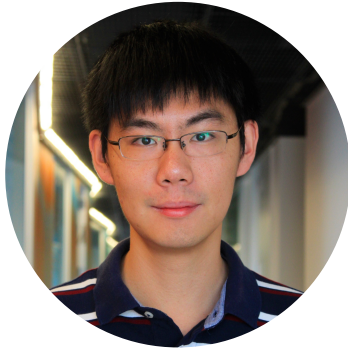


# Unbiased Offline Recommender Evaluation for Missing-Not-At-Random Implicit Feedback



Longqi Yang



Yin Cui



Yuan Xuan



Chenyang Wang



Serge Belongie



Deborah Estrin



**CORNELL  
TECH**



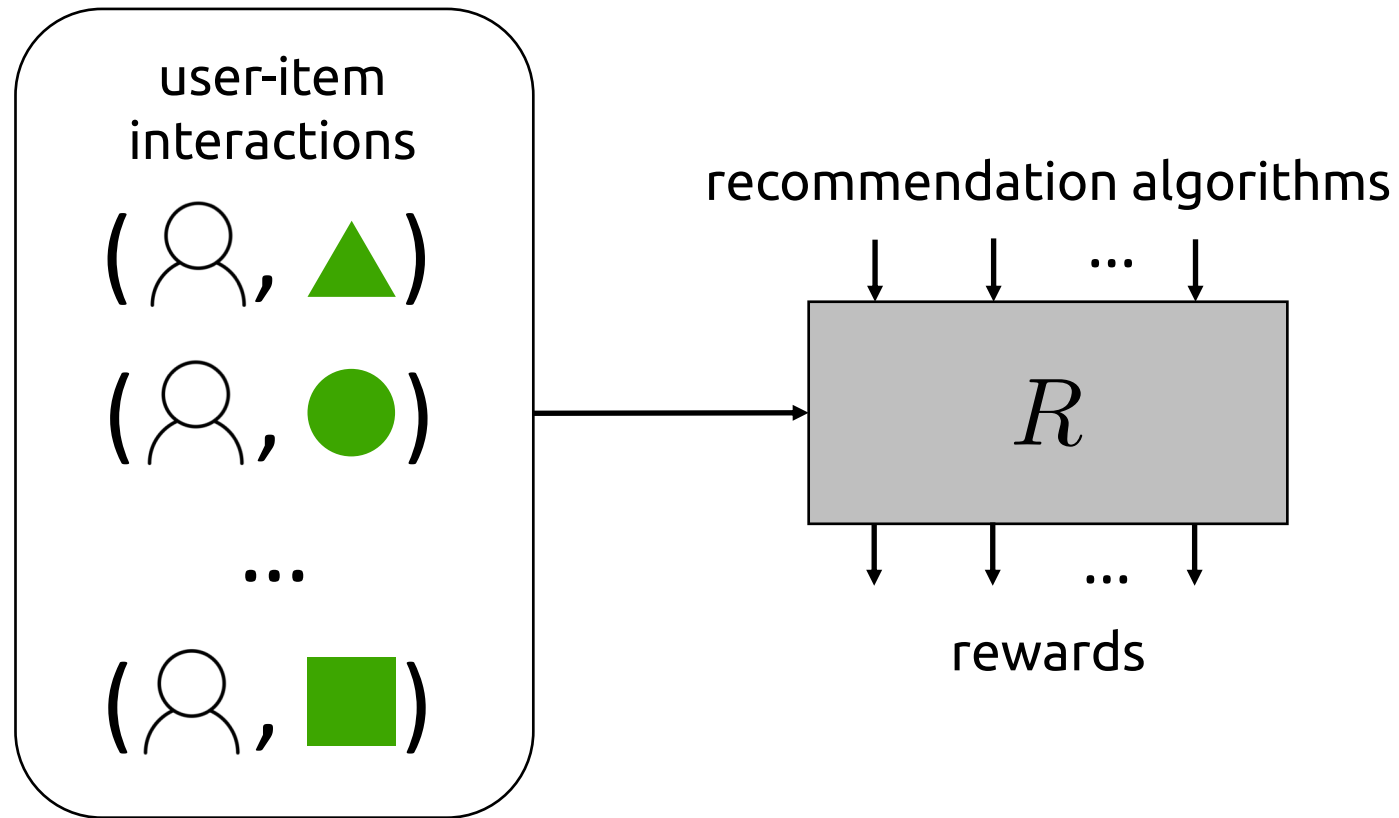
Cornell CIS  
**Computer Science**

Funders:

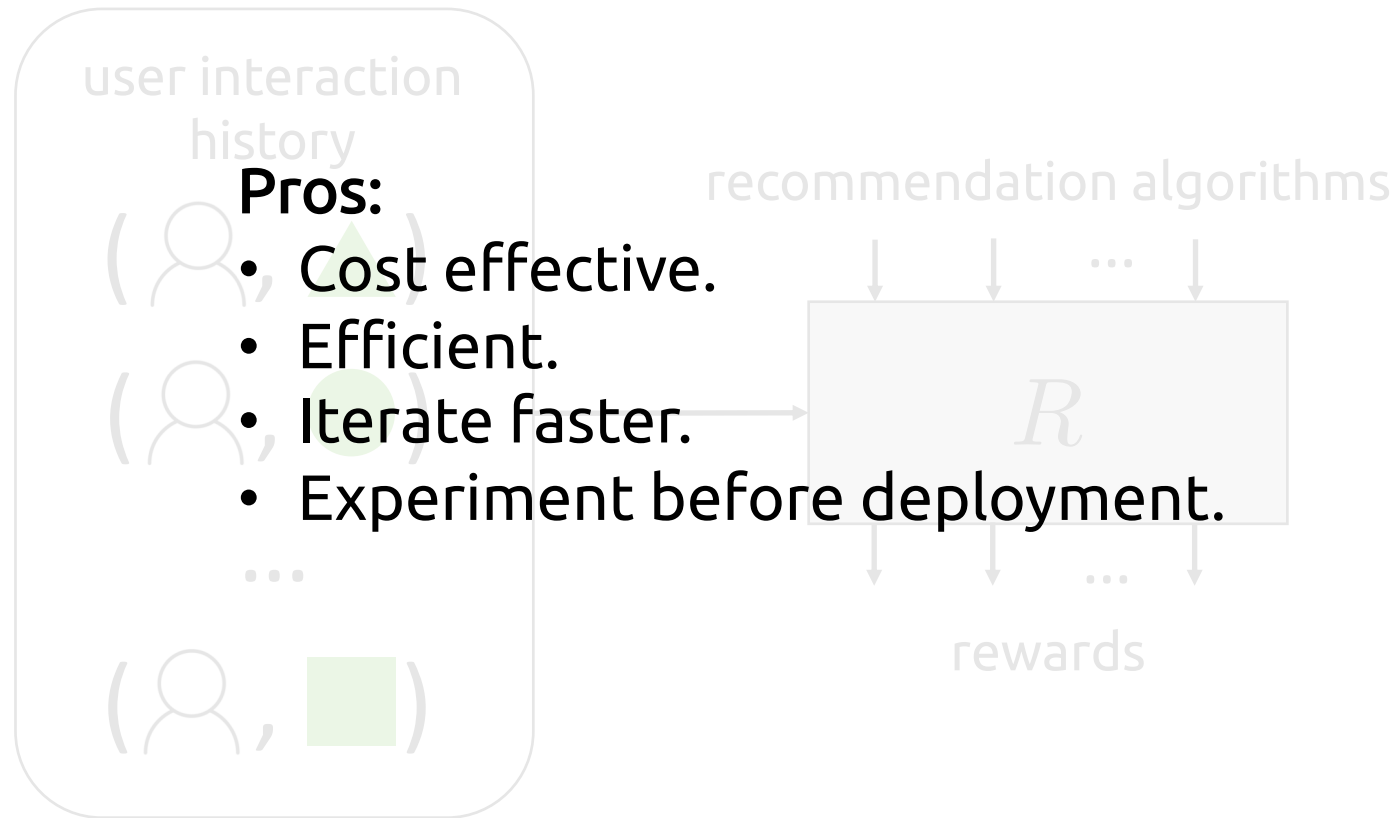


**Oath:**  
A Verizon company

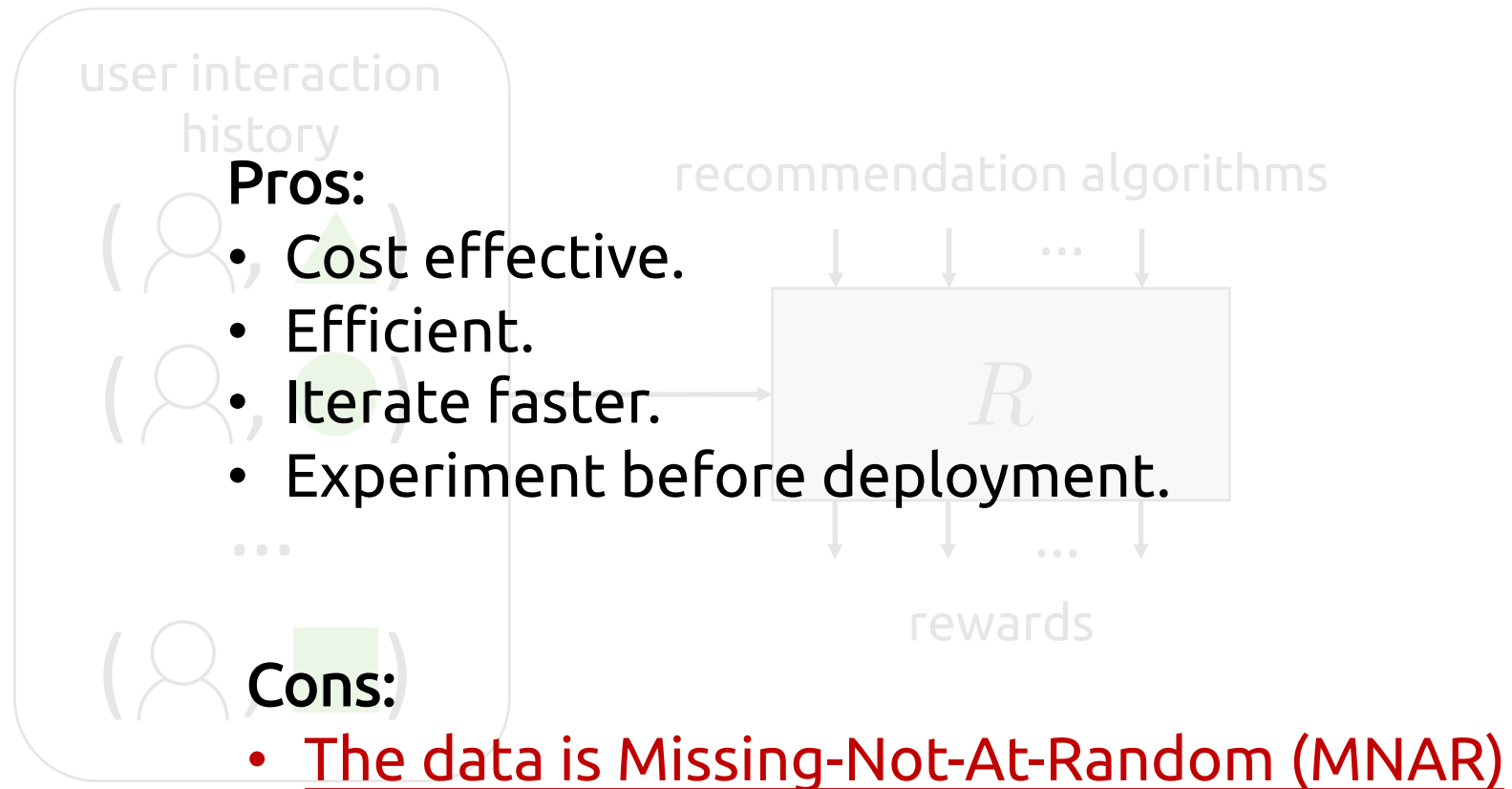
# Offline Evaluation of Recommendation Algorithm



# Offline Evaluation of Recommendation Algorithm

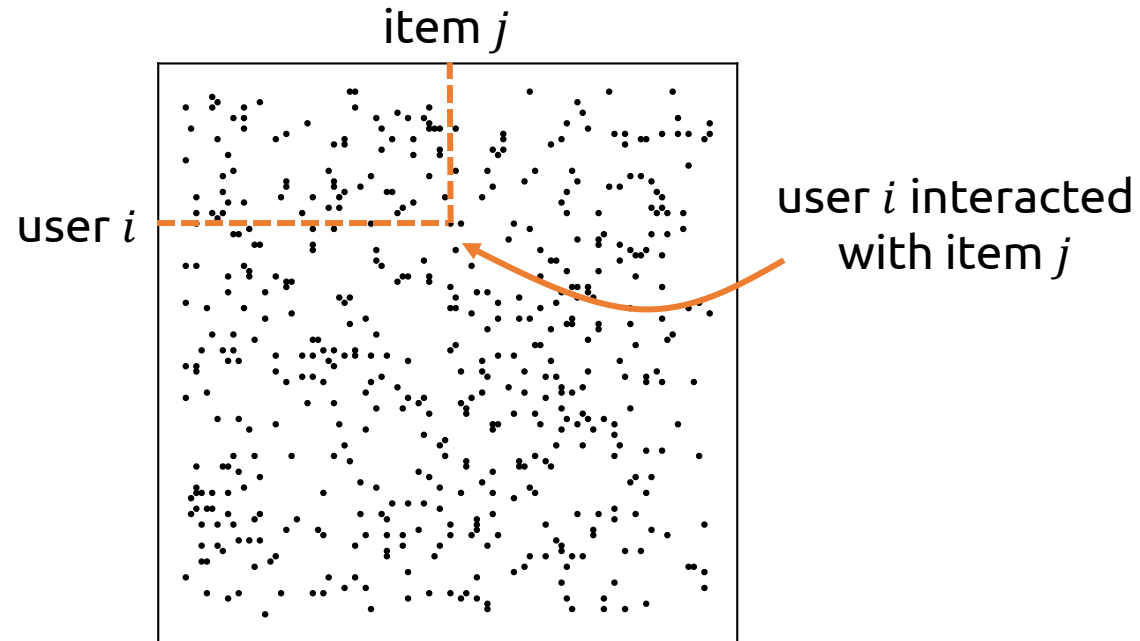


# Offline Evaluation of Recommendation Algorithm

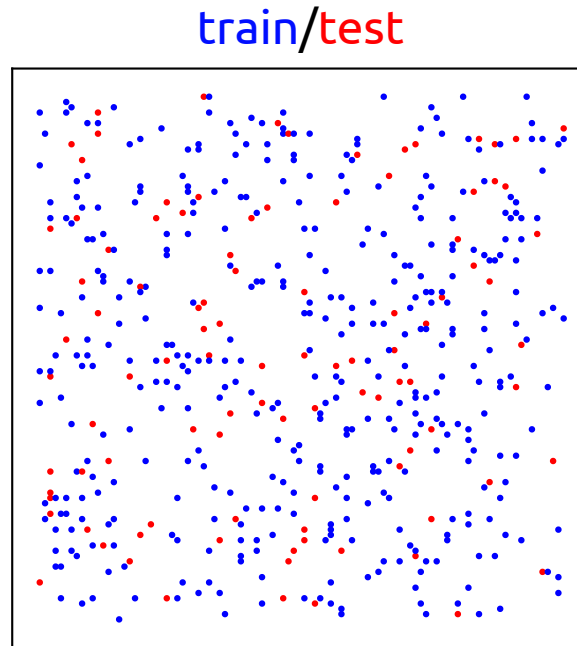




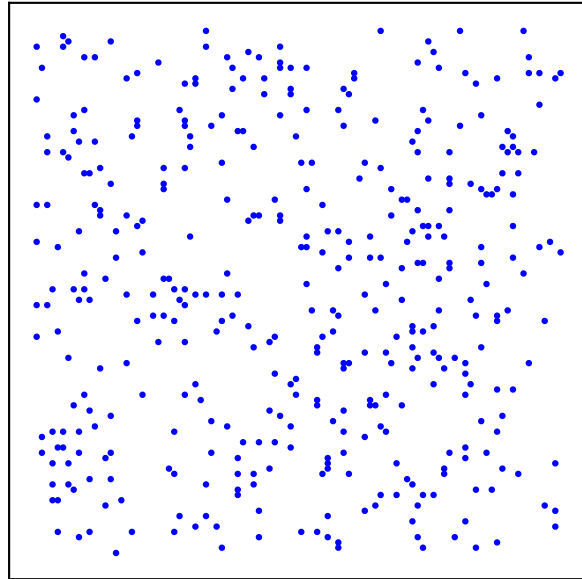
# Offline Evaluation procedure



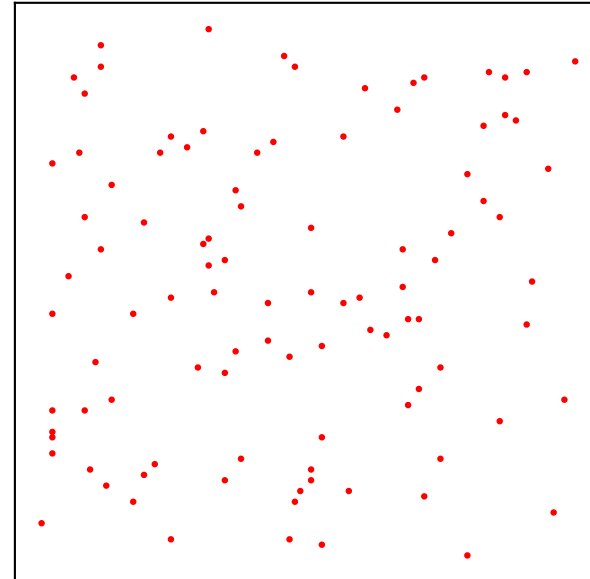
# Offline Evaluation procedure



# Offline Evaluation procedure

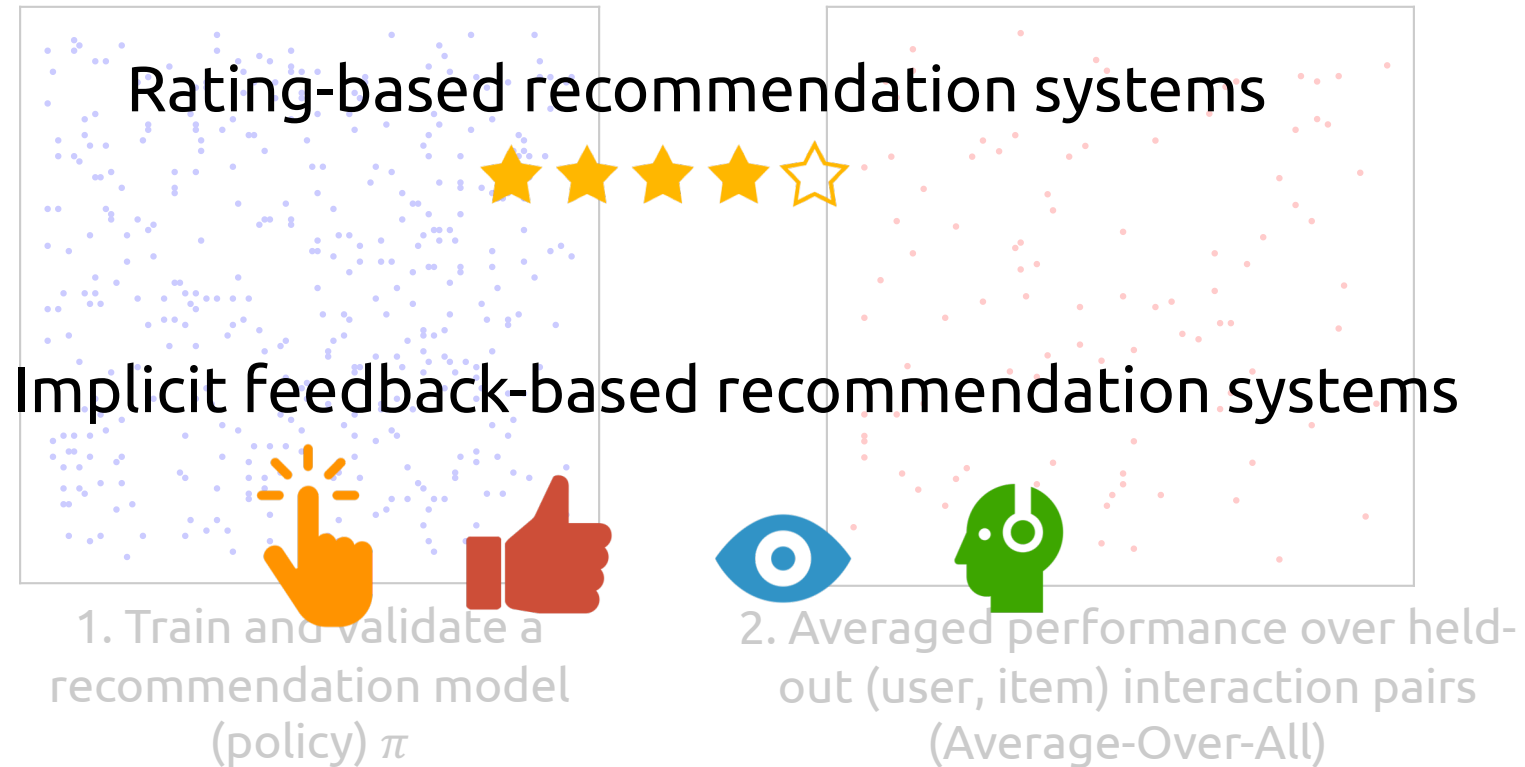


1. Train and validate a recommendation model



2. Averaged performance over held-out (user, item) interaction pairs (Average-Over-All)

# Offline Evaluation procedure



Previous work: **Average-Over-All** is **biased** for rating-based recommendation systems, because ratings are **MNAR**  
[Marlin et al. 09], [Schnabel et al. 16], [Steck 10], [Steck 11], and [Steck 13]

Previous work: **Average-Over-All** is **biased** for rating-based recommendation systems, because ratings are **MNAR**  
[Marlin et al. 09], [Schnabel et al. 16], [Steck 10], [Steck 11], and [Steck 13]

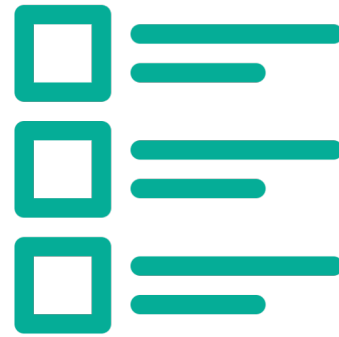
Previous work: **Average-Over-All** is **unbiased** for implicit feedback-based recommendation systems, because implicit feedback is **missing uniformly at random**.  
[Lim 15]

This work: Average-Over-All is **biased** for implicit feedback-based recommendation systems, because implicit feedback is **NOT missing uniformly at random**.

This work: Average-Over-All is **biased** for implicit feedback-based recommendation systems, because implicit feedback is **NOT missing uniformly at random**.



trending



recommendation

Popularity bias (Users are more likely to be exposed to popular items)



# A Hypothetical Example

	Popular Items		Long-tail Items
# of liked items (over all items)	1	:	10
# of liked items (over observations)	10	:	1
Algorithm 1 Performance	0.8		0
Algorithm 2 Performance	0.75		0.75

# A Hypothetical Example

Popular Items

Long-tail Items

# of liked items  
(over all items)

1

:

10

# of liked items  
(over observations)

10

:

1

Algorithm 1  
Performance

0.8

0

Algorithm 2  
Performance

0.75

0.75

# A Hypothetical Example

Popular Items

Long-tail Items

# of liked items  
(over all items)

1

:

10

# of liked items  
(over observations)

10

:

1

Algorithm 1  
Performance

0.8

0

Algorithm 2  
Performance

0.75

0.75

# A Hypothetical Example

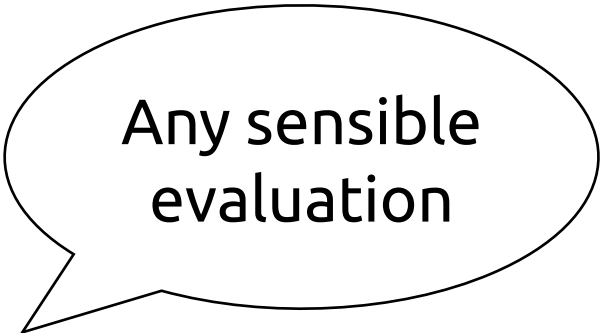
	Popular Items		Long-tail Items
# of liked items (over all items)	1	:	10
# of liked items (over observations)	10	:	1
Algorithm 1 Performance	0.8		0
Algorithm 2 Performance	0.75		0.75

# A Hypothetical Example

	Popular Items		Long-tail Items
# of liked items (over all items)	1	:	10
# of liked items (over observations)	10	:	1
Algorithm 1 Performance	0.8		0
Algorithm 2 Performance	0.75		0.75

# A Hypothetical Example

	Popular Items		Long-tail Items
# of liked items (over all items)	1	:	10
# of liked items (over observations)	10	:	1
Algorithm 1 Performance			0
Algorithm 2 Performance	0.75		0.75



# A Hypothetical Example


	Popular Items		Long-tail Items
# of liked items (over all items)	1	:	10
# of liked items (over observations)	0	:	1
Algorithm 1 Performance	0.8		0
Algorithm 2 Performance	0.75		0.75

Average-Over-All

# Formalize Reward $R$

Item rankings predicted by an algorithm

Ideal evaluation:  $R(\hat{Z}) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{1}{|\mathcal{S}_u|} \sum_{i \in \mathcal{S}_u} c(\hat{Z}_{u,i})$





# Formalize Reward $R$

Item rankings predicted by an algorithm

Predicted ranking of item  $i$  for user  $u$

Items liked by user  $u$  among the entire item set

scoring metric

Reward for  $(u, i)$  pair

Ideal evaluation: 
$$R(\hat{Z}) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{1}{|\mathcal{S}_u|} \sum_{i \in \mathcal{S}_u} c(\hat{Z}_{u,i})$$

# Formalize Reward $R$

Item rankings predicted by an algorithm

Predicted ranking of item  $i$  for user  $u$

Items liked by user  $u$  among the entire item set

scoring metric

Reward for user  $u$

Ideal evaluation: 
$$R(\hat{Z}) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \left[ \frac{1}{|\mathcal{S}_u|} \sum_{i \in \mathcal{S}_u} c(\hat{Z}_{u,i}) \right]$$

# Formalize Reward $R$

Item rankings predicted by an algorithm

Predicted ranking of item  $i$  for user  $u$

Ideal evaluation:  $R(\hat{Z}) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{1}{|\mathcal{S}_u|} \sum_{i \in \mathcal{S}_u} c(\hat{Z}_{u,i})$

Items liked by user  $u$  among the entire item set

scoring metric


Reward for the algorithm

The diagram shows the formula  $R(\hat{Z}) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{1}{|\mathcal{S}_u|} \sum_{i \in \mathcal{S}_u} c(\hat{Z}_{u,i})$  enclosed in a dashed green box. An orange arrow points from the text 'Item rankings predicted by an algorithm' to the  $\hat{Z}$  in the formula. Another orange arrow points from 'Predicted ranking of item  $i$  for user  $u$ ' to  $\hat{Z}_{u,i}$ . A third orange arrow points from 'Items liked by user  $u$  among the entire item set' to  $\mathcal{S}_u$ . A fourth orange arrow points from 'scoring metric' to  $c(\hat{Z}_{u,i})$ . A green arrow points from 'Reward for the algorithm' to the entire dashed green box.

# Formalize Reward $R$

Average-Over-All:  $\hat{R}_{\text{AOA}}(\hat{Z}) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{1}{|\mathcal{S}_u^*|} \sum_{i \in \mathcal{S}_u^*} c(\hat{Z}_{u,i})$

Items liked by user  $u$  (observed)



# Formalize Bias

$$\mathbb{E}_O \left[ \hat{R}_{\text{AOA}}(\hat{Z}) \right] \neq R(\hat{Z})$$

  $O_{u,i} = 1$  if  $(u, i)$  is observed, and  $O_{u,i} = 0$  otherwise

$$O_{u,i} \sim \mathcal{B}(1, P_{u,i})$$

# Inverse-Propensity-Scoring (IPS)

$$\hat{R}_{\text{AOA}}(\hat{Z}) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \left[ \frac{1}{|\mathcal{S}_u^*|} \sum_{i \in \mathcal{S}_u^*} c(\hat{Z}_{u,i}) \right]$$
$$\hat{R}_{\text{IPS}}(\hat{Z}|P) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \left[ \frac{1}{|\mathcal{S}_u|} \sum_{i \in \mathcal{S}_u^*} \frac{c(\hat{Z}_{u,i})}{P_{u,i}} \right]$$

# Inverse-Propensity-Scoring (IPS)


$$\hat{R}_{\text{AOA}}(\hat{Z}) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \left[ \frac{1}{|\mathcal{S}_u^*|} \sum_{i \in \mathcal{S}_u^*} c(\hat{Z}_{u,i}) \right]$$
$$\hat{R}_{\text{IPS}}(\hat{Z}|P) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \left[ \frac{1}{|\mathcal{S}_u|} \sum_{i \in \mathcal{S}_u^*} \frac{c(\hat{Z}_{u,i})}{P_{u,i}} \right]$$

$$\mathbb{E}_O \left[ \hat{R}_{\text{IPS}}(\hat{Z}|P) \right] = R(\hat{Z})$$

# Self-Normalized Inverse-Propensity-Scoring (SNIPS)

[Swaminathan et al.15]

$$\hat{R}_{\text{IPS}}(\hat{Z}|P) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \left[ \frac{1}{|\mathcal{S}_u|} \right] \sum_{i \in \mathcal{S}_u^*} \frac{c(\hat{Z}_{u,i})}{P_{u,i}}$$


$$\hat{R}_{\text{SNIPS}}(\hat{Z}|P) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \left[ \frac{1}{\sum_{i \in \mathcal{S}_u^*} \frac{1}{P_{u,i}}} \right] \sum_{i \in \mathcal{S}_u^*} \frac{c(\hat{Z}_{u,i})}{P_{u,i}}$$



# Estimating Propensity Scores

**Factor:** Popularity bias (Users are more likely to be exposed to popular items)

## Assumptions:

- User-independence assumption  $P_{u,i} = P(O_{u,i} = 1) = P(O_{*,i} = 1) = P_{*,i}$

- Two-steps assumption  $P_{*,i} = P_{*,i}^{\text{select}} \cdot P_{*,i}^{\text{interact|select}}$

- User preference is not affected by item presentation

$$P_{*,i}^{\text{interact|select}} = P_{*,i}^{\text{interact}}$$

# Estimating Propensity Scores

Popularity bias model [Steck 11]:

$$\hat{P}_{*,i}^{\text{select}} \propto (n_i^*)^\gamma$$



Observed item  
popularity

# Estimating Propensity Scores

Popularity bias model [Steck 11]:

$$\hat{P}_{*,i}^{\text{select}} \propto (n_i^*)^\gamma$$

Estimated from  
known online content  
serving policy

$$\hat{P}_{*,i} \propto (n_i^*)^{\left(\frac{\gamma+1}{2}\right)}$$

# Measuring bias in recommender evaluation (Yahoo! music rating dataset)

Mean Absolute Error (MAE), Recall

Model	Average-Over-All	$R_{SNIPS}$ ( $\gamma = 1.5$ )	$R_{SNIPS}$ ( $\gamma = 2.0$ )	$R_{SNIPS}$ ( $\gamma = 2.5$ )	$R_{SNIPS}$ ( $\gamma = 3.0$ )
U-CML	0.401	0.270	0.260	0.253	0.248
A-CML	0.399	0.274	0.264	0.258	0.253
BPR	0.380	0.275	0.268	0.258	0.258
PMF	0.386	0.267	0.268	0.258	0.248

$R_{SNIPS}$  produces  
significantly lower  
MAE

# Measuring bias in recommender evaluation (Yahoo! music rating dataset)

Mean Absolute Error (MAE), Recall

Model	Average- Recall@10	RSNIPS ( $\lambda = 0.01$ )	RSNIPS ( $\lambda = 0.001$ )	RSNIPS ( $\lambda = 0.0001$ )	RSNIPS ( $\lambda = 3.0$ )
U-CML	0.381	0.272	0.269	0.259	0.248
A-CML	0.399	0.274	0.261	0.258	0.253
BPR	0.380	0.275	0.269	0.258	0.258
PMF	0.386	0.267			0.258

The accuracy of recommending popular items is a significant **overestimation** of the true recommendation performance

RSNIPS produces significantly lower MAE

Please come to our poster or refer to our paper for:

- Proofs
- Experimental details.
- More experiments.
- Deeper analysis of the unbiased evaluator.

# Conclusions and Future Work

$$\mathbb{E}_O \left[ \hat{R}_{\text{IPS}}(\hat{Z}|P) \right] = R(\hat{Z})$$

$$\hat{R}_{\text{SNIPS}}(\hat{Z}|P) = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{1}{\sum_{i \in \mathcal{S}_u^*} \frac{1}{P_{u,i}}} \sum_{i \in \mathcal{S}_u^*} \frac{c(\hat{Z}_{u,i})}{P_{u,i}}$$

- Understanding variance of evaluators.
- Propensity estimation (e.g., incorporate auxiliary user and item information).
- Debias training of recommendation systems (e.g., [Liang et al. 16]).



<http://www.openrec.ai>

Github link, documents, and tutorials

## Longqi Yang

Ph.D. candidate

Computer Science, Cornell Tech, Cornell University

Email: [ylongqi@cs.cornell.edu](mailto:ylongqi@cs.cornell.edu)

Web: [bit.ly/longqi](http://bit.ly/longqi)

Twitter: [@ylongqi](https://twitter.com/ylongqi)

Connected Experiences Lab

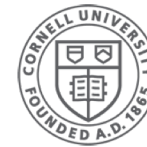
<http://cx.jacobs.cornell.edu/>

Small Data Lab

<http://smalldata.io/>



**CORNELL  
TECH**



Cornell CIS  
**Computer Science**

Funders:



**Oath:**  
A Verizon company