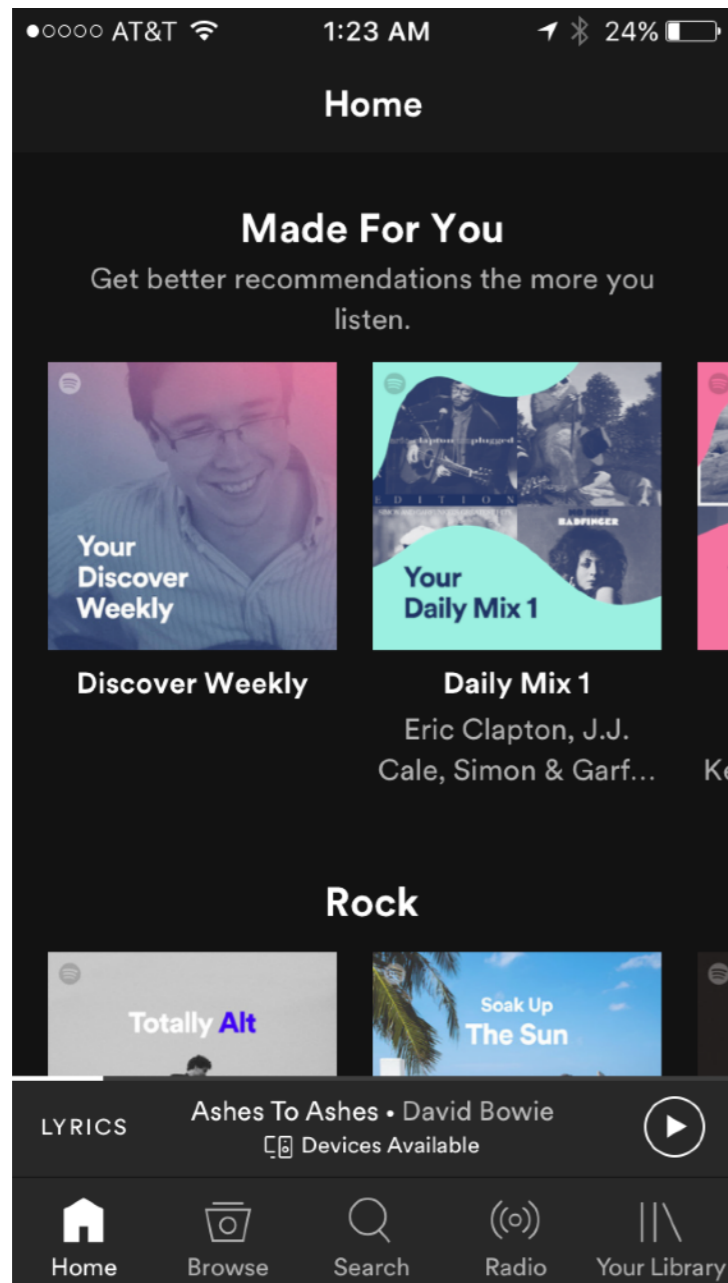# Explore, Exploit, and Explain: Personalizing Explainable Recommendations with Bandits

James McInerney, Ben Lacker, Samantha Hansen, Karl Higley, Hugues Bouchard, Alois Gruson, Rishabh Mehrotra

email: jamesm@spotify.com

# Research question: how to explore-exploit over explainable recommendations?
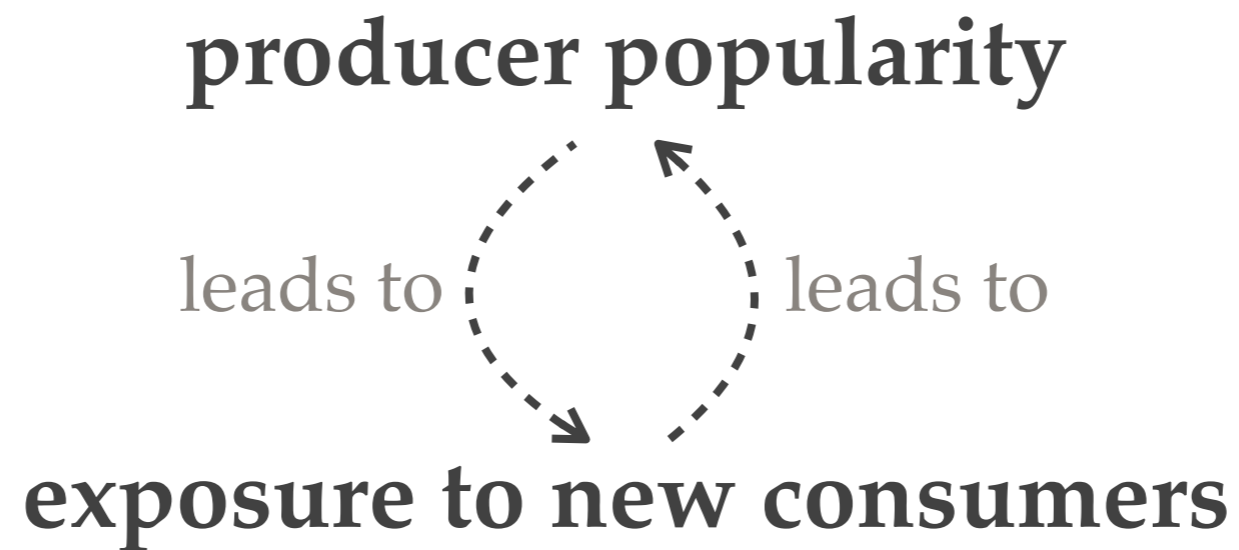


- e.g. home page of Spotify

- items arranged into shelves, each shelf has an explanation for the associated recommendation

# Outline

1. Pareto principle for content producers

2. a causal diagnosis of filter bubbles in recommendation

3. contextual bandits for recommendation

4. explained recommendations

5. introducing Bart (<u>ban</u>dits for <u>r</u>ecommendations as <u>t</u>reatments)

6. offline and online experiments on homepage data

7. conclusions & future work

# A small number of producers dominate consumption in culture

**producer popularity**

leads to · · ·  · · · leads to

**exposure to new consumers**

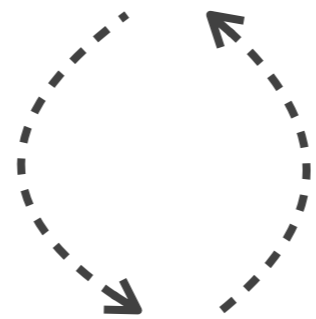# A small number of producers dominate consumption in culture

e.g. musicians, authors, actors

**producer popularity**

leads to    leads to

**exposure to new consumers**

# A small number of producers dominate consumption in culture

e.g. musicians, authors, actors
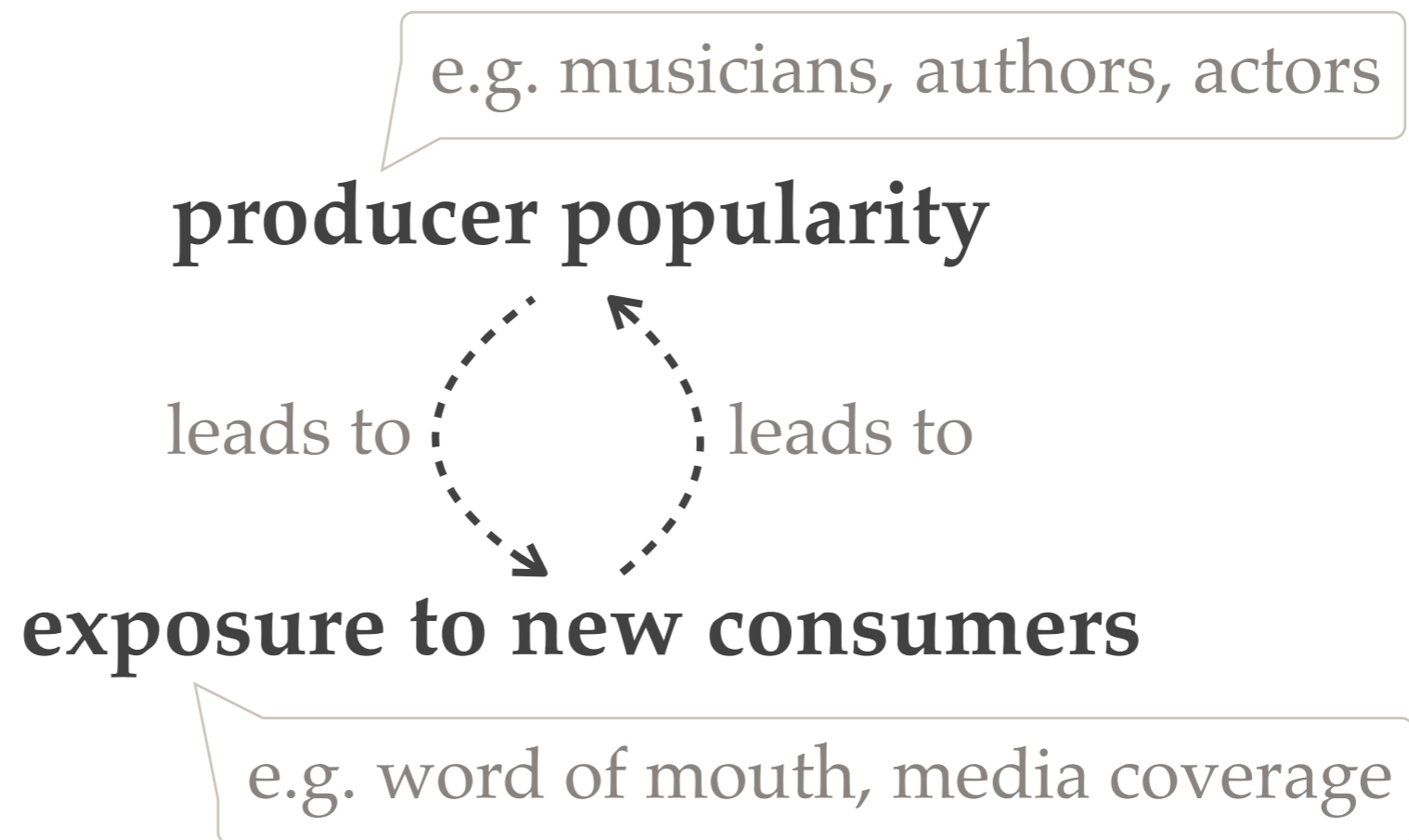
**producer popularity**

leads to          leads to
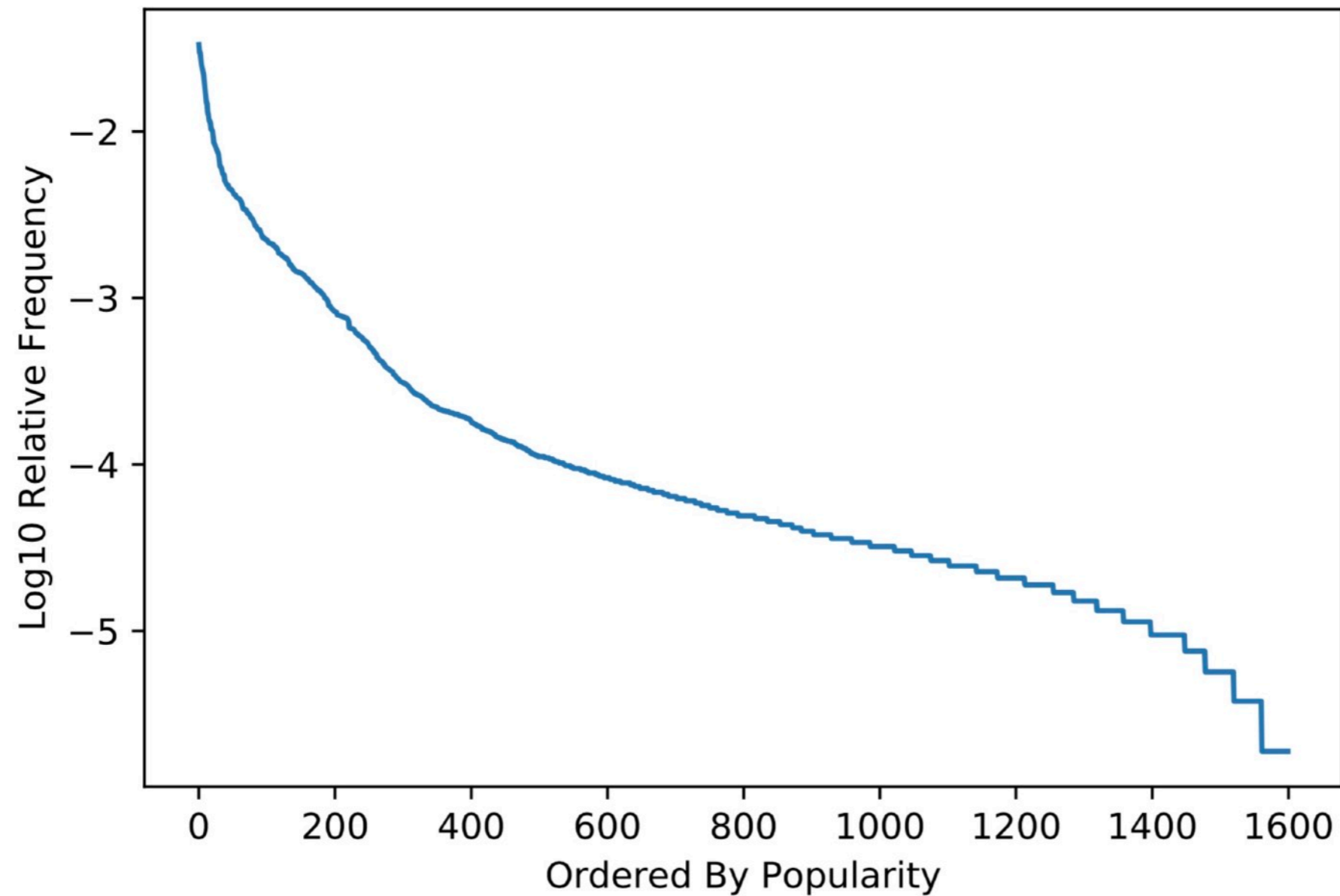
**exposure to new consumers**

e.g. word of mouth, media coverage

# A small number of producers dominate consumption in culture

e.g. musicians, authors, actors

**producer popularity**

leads to        leads to

**exposure to new consumers**
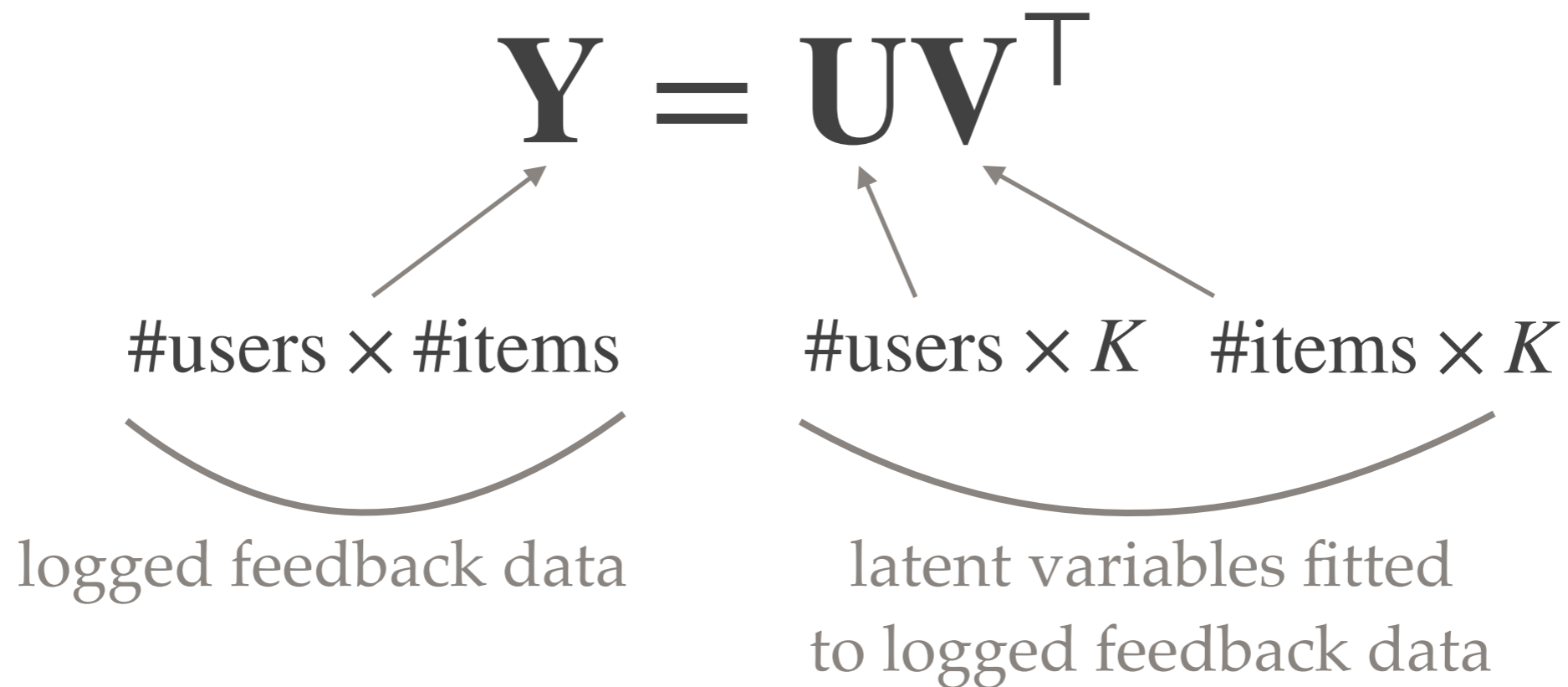
e.g. word of mouth, media coverage

- known as the Matthew effect or Pareto principle [Juran, 1937]

# A small number of producers dominate consumption in culture

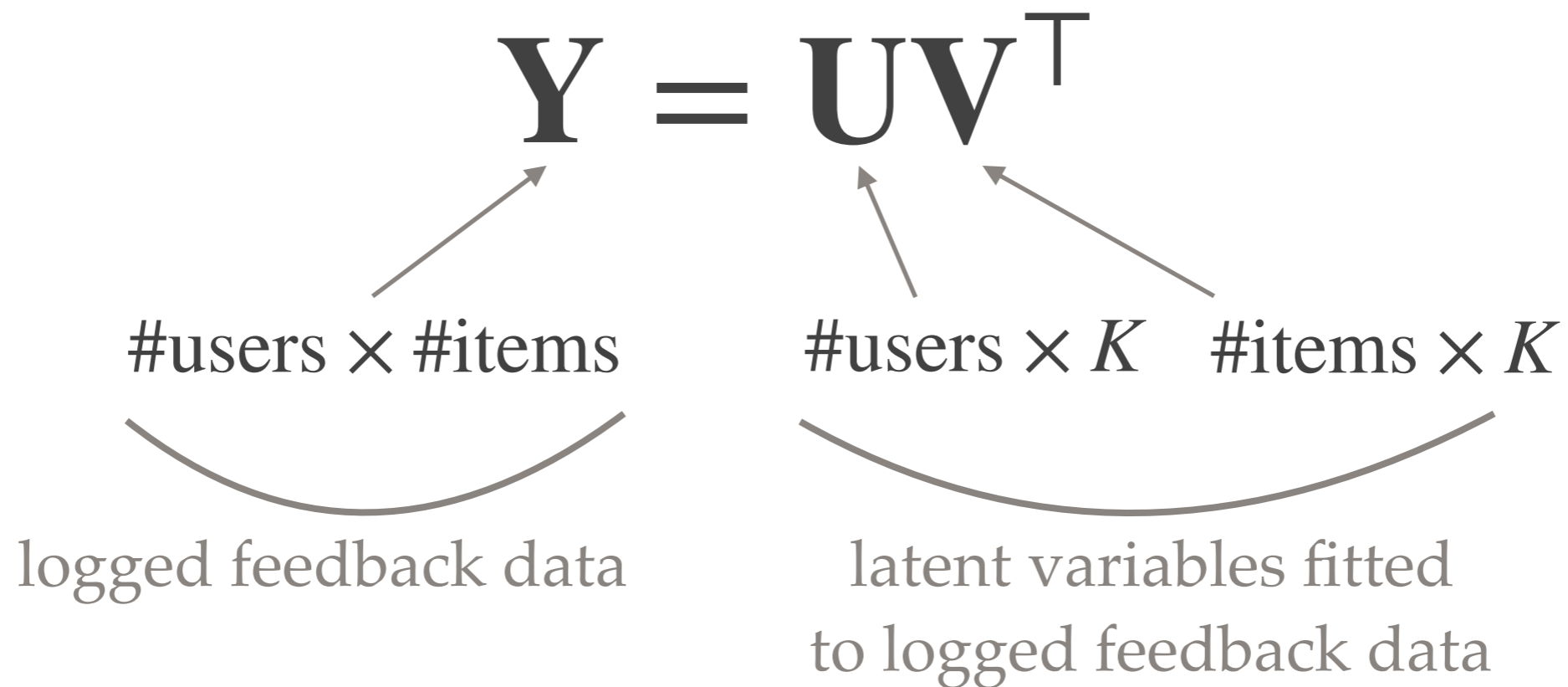# Collaborative filtering perpetuates the Pareto principle

e.g. matrix factorization

$$\mathbf{Y} = \mathbf{U}\mathbf{V}^{\top}$$

#users $\times$ #items          #users $\times K$    #items $\times K$

logged feedback data          latent variables fitted
                              to logged feedback data

# Collaborative filtering perpetuates the Pareto principle

e.g. matrix factorization

$$\mathbf{Y} = \mathbf{U}\mathbf{V}^{\top}$$

#users × #items        #users × $K$    #items × $K$

logged feedback data        latent variables fitted
to logged feedback data

- in general: collaborative filtering engines use implicit feedback data from users to learn a model of user preferences

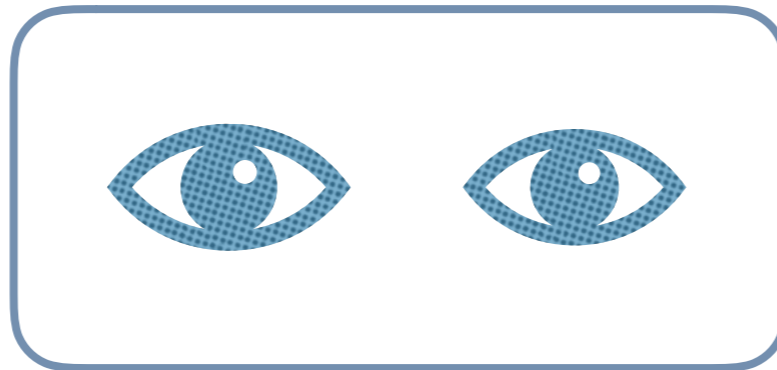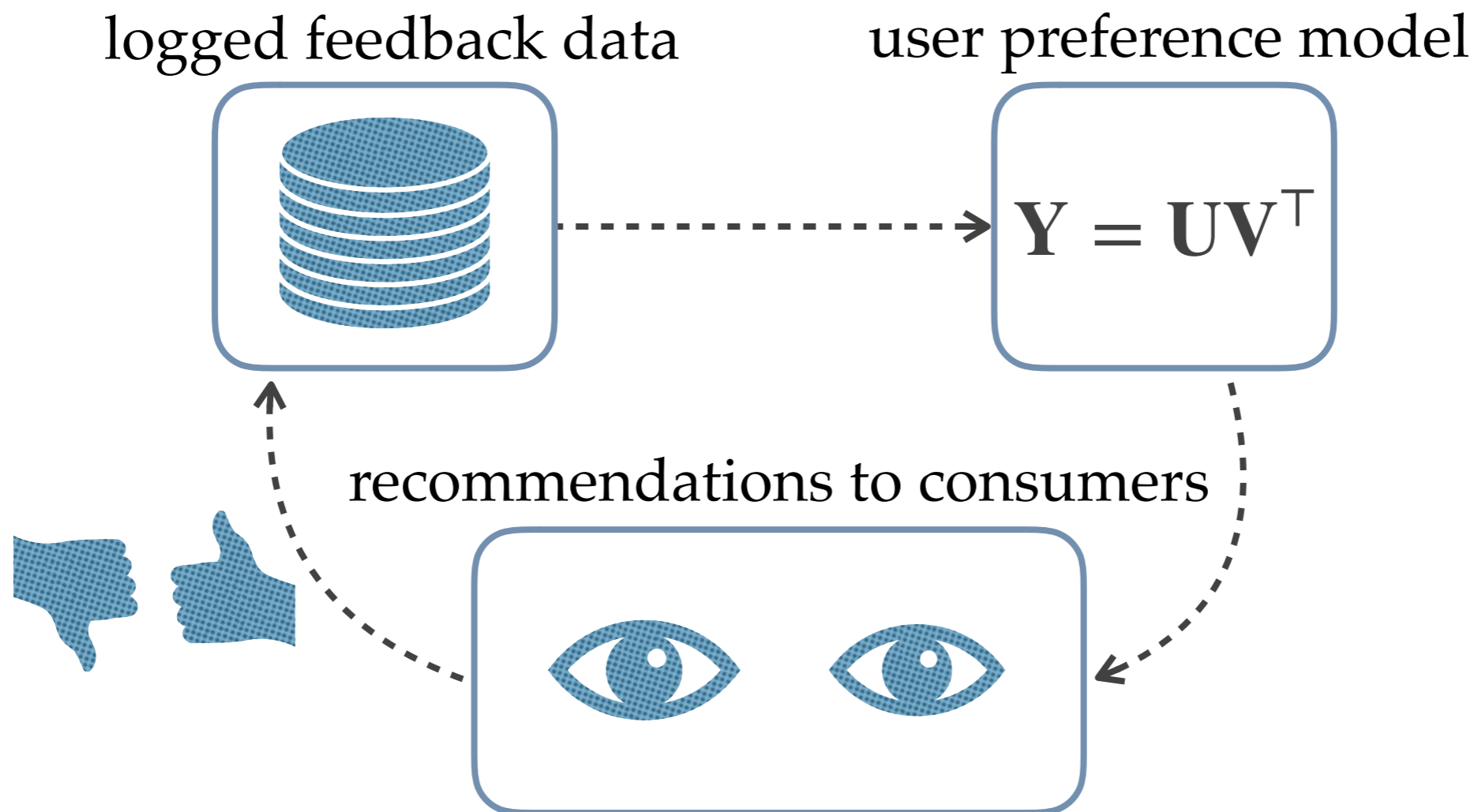# Collaborative filtering perpetuates the Pareto principle

logged feedback data

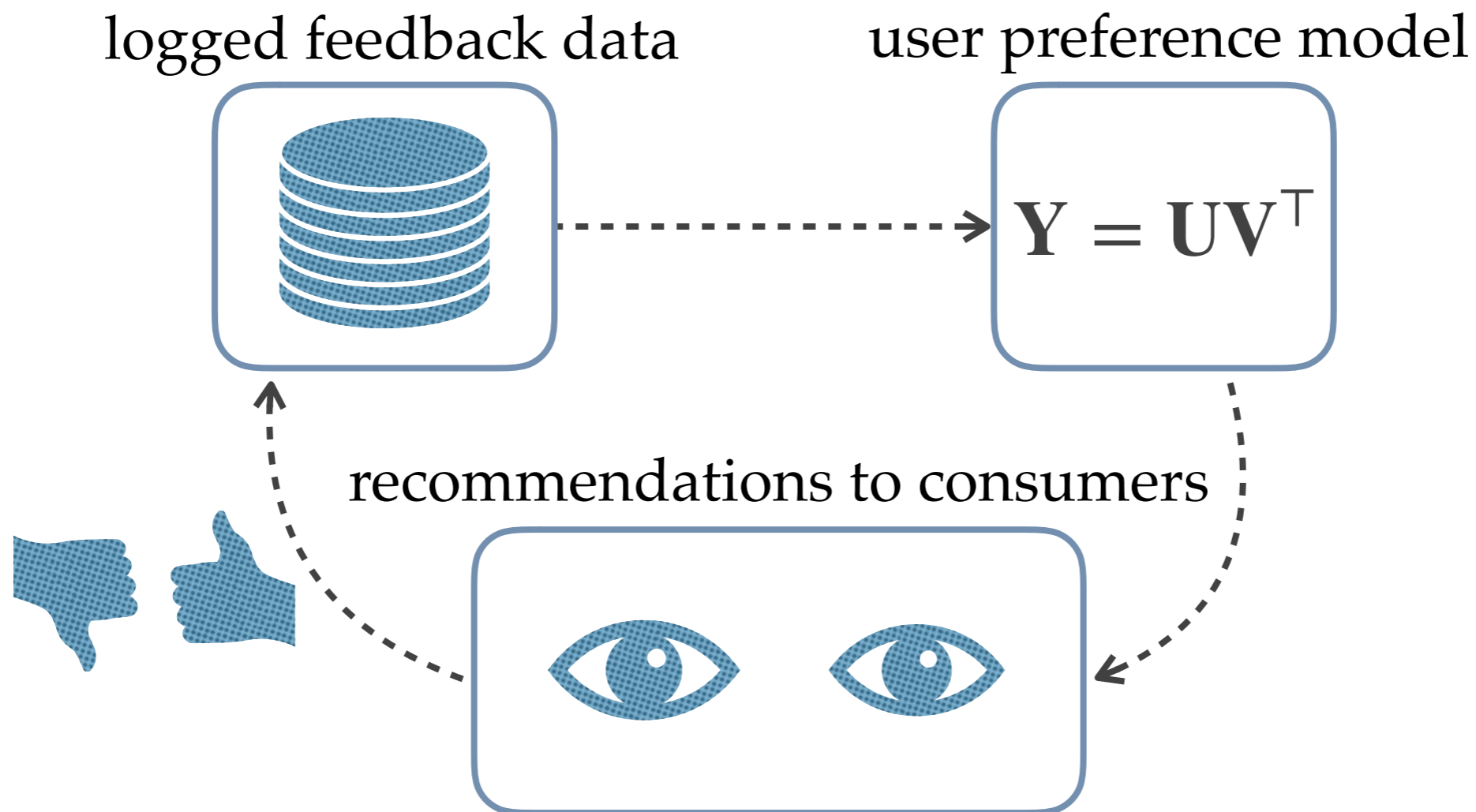user preference model

$$\mathbf{Y} = \mathbf{U}\mathbf{V}^{\top}$$

recommendations to consumers
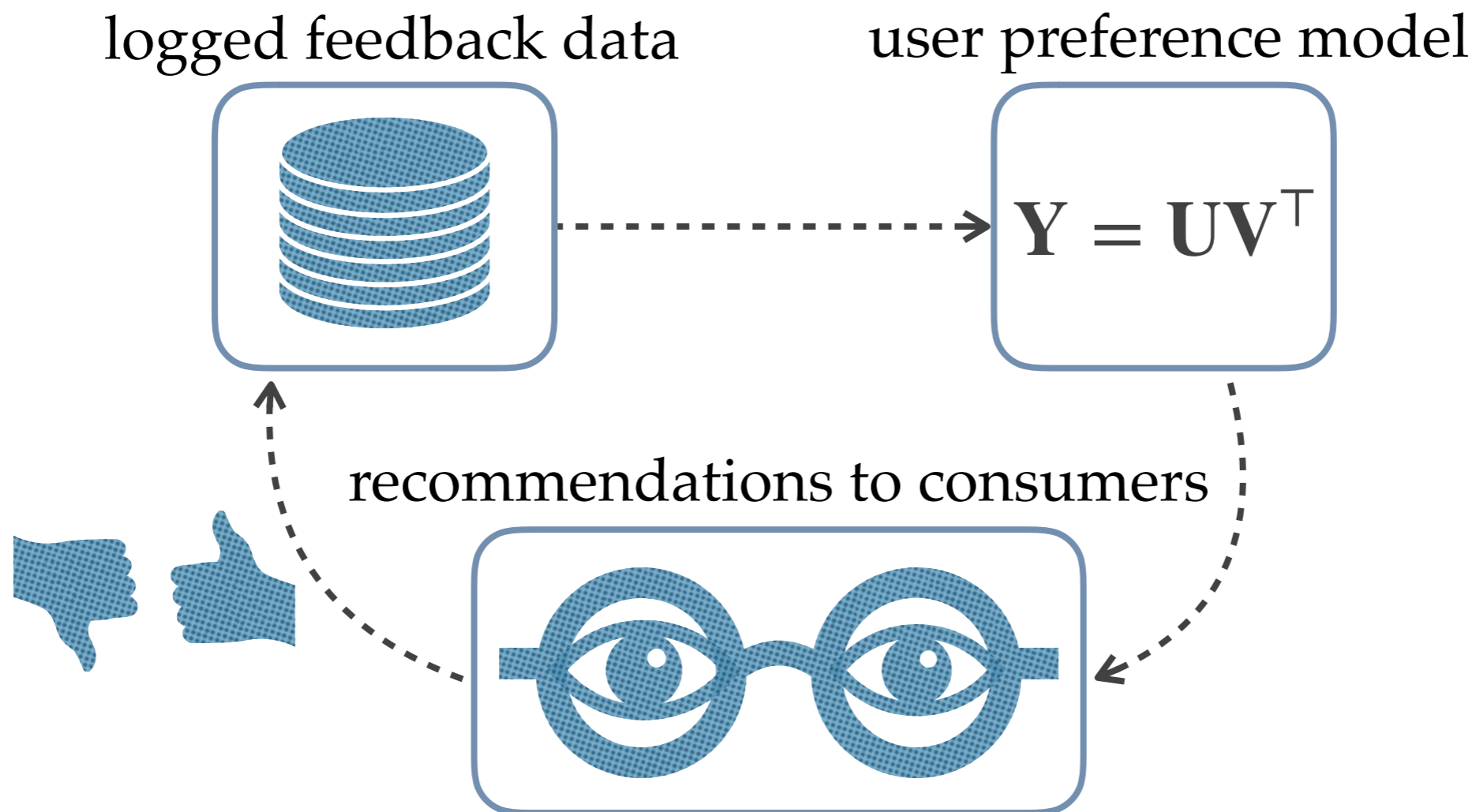
# Collaborative filtering perpetuates the Pareto principle

logged feedback data

user preference model

$$\mathbf{Y} = \mathbf{U}\mathbf{V}^{\top}$$

recommendations to consumers

# Collaborative filtering perpetuates the Pareto principle

logged feedback data          user preference model



$$\mathbf{Y} = \mathbf{U}\mathbf{V}^\top$$

recommendations to consumers

"How Algorithmic Confounding in Recommendation Systems Increases Homogeneity and Decreases Utility" (Chaney et al. 2017)

"Modeling User Exposure in Recommendation" (Liang et al. 2016)

# Collaborative filtering perpetuates the Pareto principle



logged feedback data

user preference model

$$\mathbf{Y} = \mathbf{U}\mathbf{V}^{\top}$$

recommendations to consumers

"How Algorithmic Confounding in Recommendation Systems Increases Homogeneity and Decreases Utility" (Chaney et al. 2017)

"Modeling User Exposure in Recommendation" (Liang et al. 2016)

# Standard collaborative filtering methods are limited because they can only exploit or ignore

**recommender system relevance certainty**

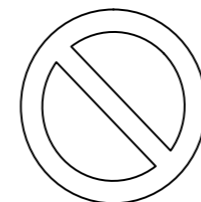|                                          | Low certainty | High certainty |
|------------------------------------------|---------------|----------------|
| **ground truth item relevance** — Low relevance | Sometimes Exploit / Sometimes Ignore | Ignore |
| **ground truth item relevance** — High relevance | Sometimes Exploit / Sometimes Ignore | Exploit |

# Standard collaborative filtering methods are limited because they can only exploit or ignore

**recommender system relevance certainty**

|  | Low certainty | High certainty |
|---|---|---|
| **Low relevance** | ⭐ Sometimes Exploit / 🚫 Sometimes Ignore | 🚫 Ignore |
| **High relevance** | ⭐ Sometimes Exploit / 🚫 Sometimes Ignore | ⭐ Exploit |

**ground truth item relevance**

# Standard collaborative filtering methods are limited in that they can only exploit or ignore

- e.g. two items, A and B, with the same click rate = 0.1

**observed implicit feedback for item A**

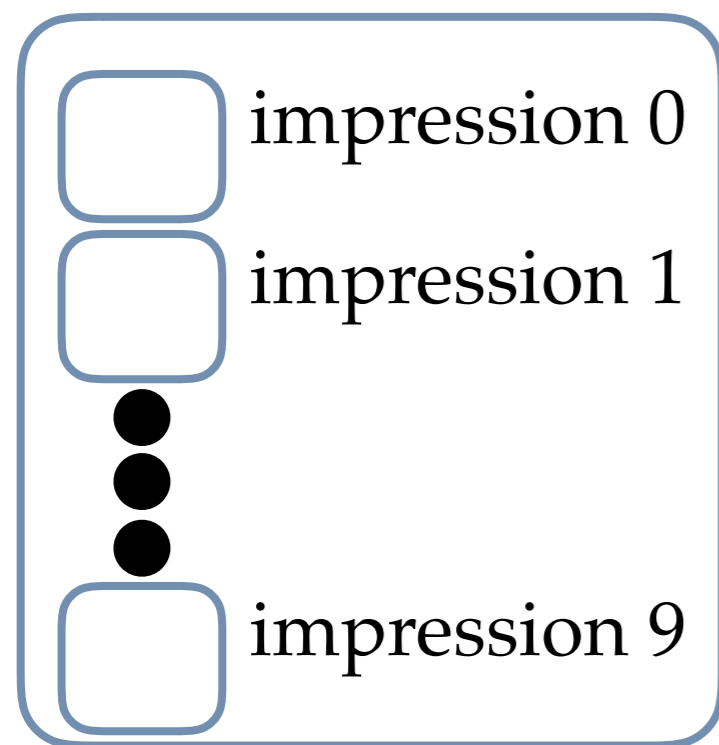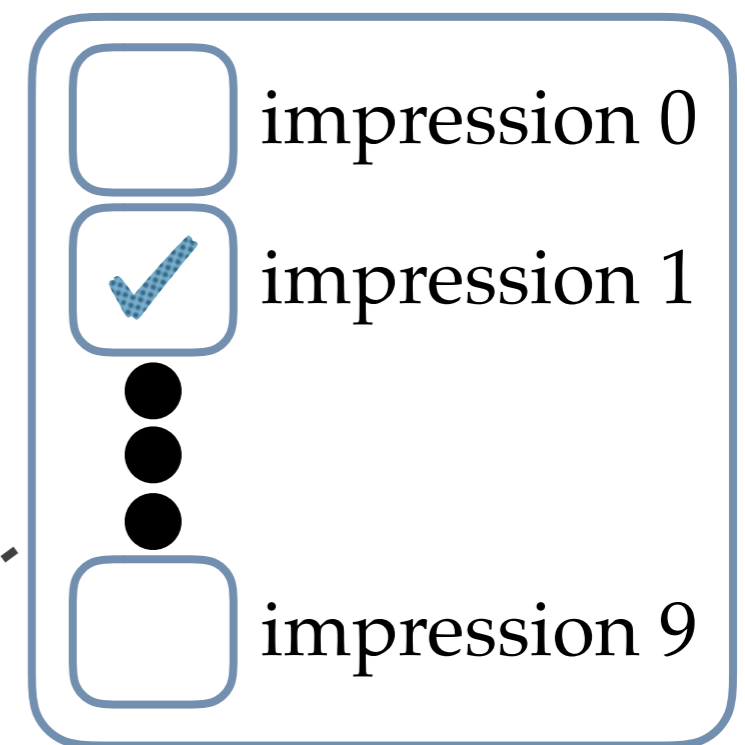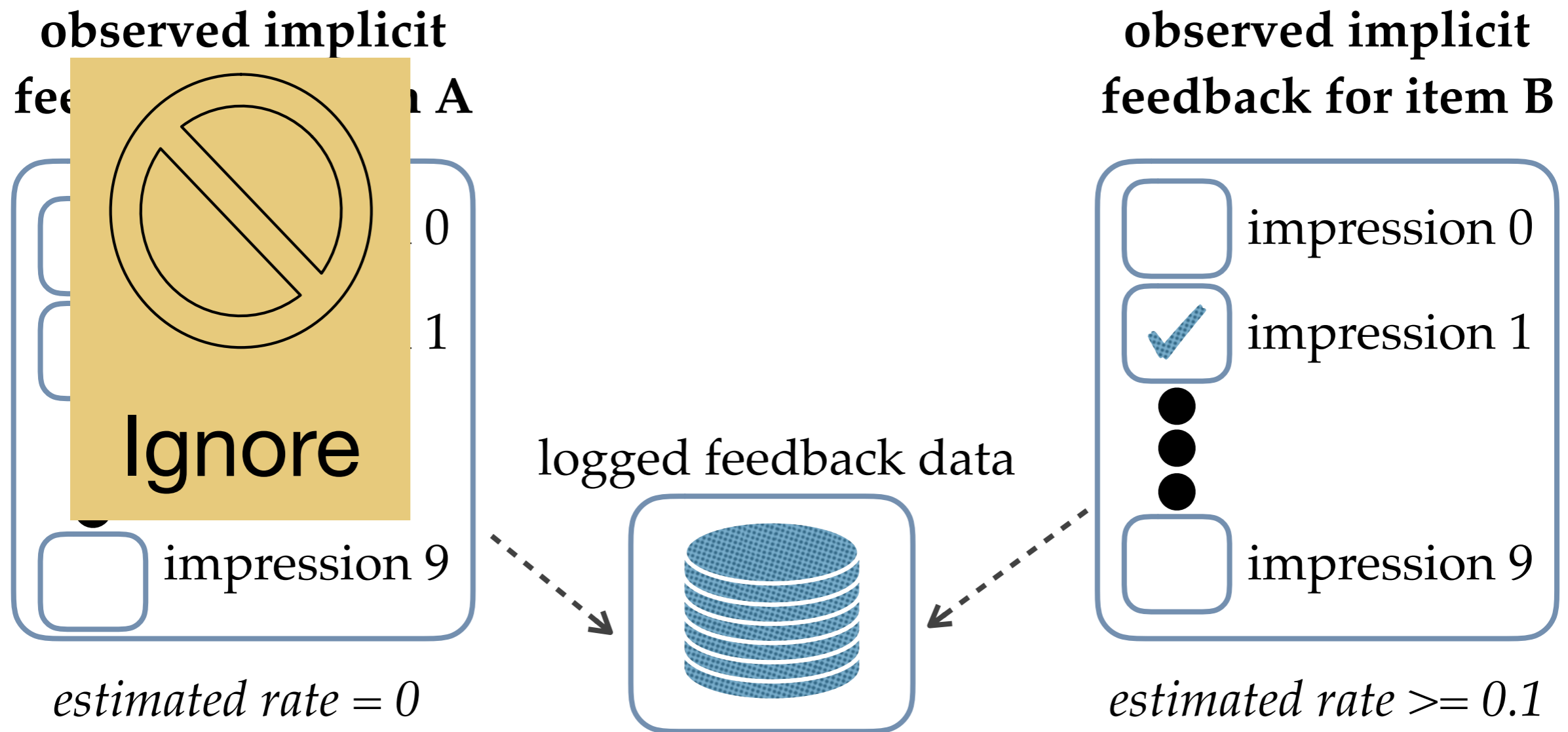**observed implicit feedback for item B**
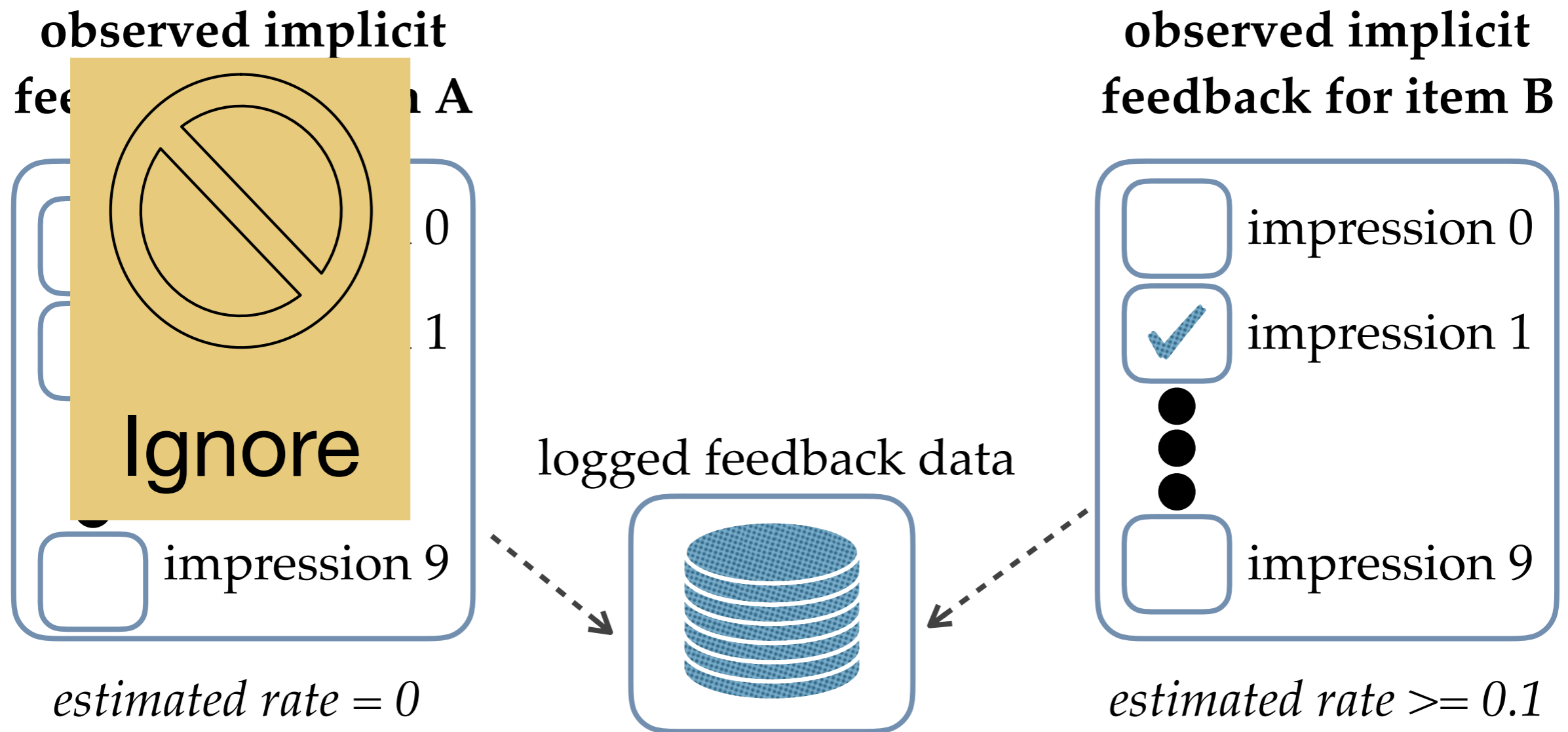
logged feedback data

# Standard collaborative filtering methods are limited in that they can only <u>exploit</u> or <u>ignore</u>

- e.g. two items, A and B, with the same click rate = 0.1



**observed implicit feedback for item A**

impression 0

impression 1

impression 9

logged feedback data

**observed implicit feedback for item B**

impression 0

✓ impression 1

impression 9

# Standard collaborative filtering methods are limited in that they can only <u>exploit</u> or <u>ignore</u>

- e.g. two items, A and B, with the same click rate = 0.1

**observed implicit feedback for item A**

**observed implicit feedback for item B**

impression 0

impression 1

impression 9

logged feedback data

impression 0

✓ impression 1

impression 9

*estimated rate = 0*

*estimated rate >= 0.1*

# Standard collaborative filtering methods are limited in that they can only <u>exploit</u> or <u>ignore</u>

- e.g. two items, A and B, with the same click rate = 0.1

**observed implicit feedback for item A**

**observed implicit feedback for item B**

impression 0

impression 1

Ignore

impression 9

logged feedback data

impression 0

✓ impression 1

impression 9

*estimated rate = 0*

*estimated rate >= 0.1*

# Standard collaborative filtering methods are limited in that they can only exploit or ignore

- e.g. two items, A and B, with the same click rate = 0.1

**observed implicit feedback for item A**

impression 0

impression 1

Ignore

impression 9

*estimated rate = 0*

logged feedback data

**observed implicit feedback for item B**

impression 0

✓ impression 1

impression 9

*estimated rate >= 0.1*

- the estimated performance will be identical only 31.3% of the time

# Randomized controlled trials

**Charles Sanders Peirce**

*"At the beginning […] the pack was well shuffled, and, the operator and subject having taken their places, the operator was governed by the color of the successive cards in choosing whether he should first diminish the weight and then increase it, or vice versa."*

On Small Differences in Sensation,
C. S. Peirce & J. Jastrow (1885)

# Randomized controlled trials

**Charles Sanders Peirce**

*"At the beginning […] the pack was well shuffled, and, the operator and subject having taken their places, the operator was governed by the color of the successive cards in choosing whether he should first diminish the weight and then increase it, or vice versa."*
On Small Differences in Sensation,
C. S. Peirce & J. Jastrow (1885)

In recommendation: uniform random items

# Let's restart from the basic ideal of randomized controlled trials

logged feedback data

recommendations to consumers

# Let's restart from the basic ideal of randomized controlled trials

logged feedback data

recommendations to consumers

uniform random items

# Let's restart from the basic ideal of randomized controlled trials

logged feedback data

user preference model

$$\mathbf{Y} = \mathbf{U}\mathbf{V}^\top$$

recommendations to consumers

uniform random items

# Let's restart from the basic ideal of randomized controlled trials

logged feedback data

user preference model

$$\mathbf{Y} = \mathbf{U}\mathbf{V}^\top$$

recommendations to consumers

uniform random items

# Let's restart from the basic ideal of randomized controlled trials

✓ =

# Let's restart from the basic ideal of randomized controlled trials

$$\checkmark = \mathbb{E}_{X, A \sim \text{Uniform}(\mathscr{A}), Y}[\log p_\theta(Y|A, X)]$$

# Let's restart from the basic ideal of randomized controlled trials

✓ $= \mathbb{E}_{X, A \sim \text{Uniform}(\mathscr{A}), Y}[\log p_\theta(Y|A, X)]$

"choose a model and train it on data how you like"

# Let's restart from the basic ideal of randomized controlled trials

"train on the right data"

$$\checkmark = \mathbb{E}_{X,A\sim\text{Uniform}(\mathscr{A}),Y}[\log p_\theta(Y|A,X)]$$

# Let's restart from the basic ideal of randomized controlled trials

"train on the right data"

$$\checkmark = \mathbb{E}_{X, A \sim \text{Uniform}(\mathscr{A}), Y}[\log p_\theta(Y | A, X)]$$

random item recommended     set of all items     model parameters     context

# But we don't want to just recommend random stuff all the time ⚡⚡⚡

# But we don't want to just recommend random stuff all the time ⚡⚡⚡

- Enter exploration-exploitation  [Sutton & Barto, 1998]

# But we don't want to just recommend random stuff all the time ⚡⚡⚡

- Enter <u>exploration-exploitation</u> [Sutton & Barto, 1998]

**recommender system relevance certainty**



|  | Low certainty | High certainty |
|---|---|---|
| **Low relevance** | Explore | Ignore |
| **High relevance** | Explore | Exploit |

*(y-axis label: ground truth item relevance)*

# But we don't want to just recommend random stuff all the time ⚡⚡⚡

- Enter <u>exploration-exploitation</u>  [Sutton & Barto, 1998]

**recommender system relevance certainty**

|  | Low certainty | High certainty |
|---|---|---|
| **Low relevance** | Explore | Ignore |
| **High relevance** | Explore | Exploit |

*ground truth item relevance*

- When the recommender is certain it has a bad item, it ignores it.

- When the recommender is certain it has a good item, it recommends it.

# But we don't want to just recommend random stuff all the time ⚡⚡⚡

- ## Enter <u>exploration-exploitation</u> [Sutton & Barto, 1998]

**recommender system relevance certainty**

|  | Low certainty | High certainty |
|---|---|---|
| **Low relevance** | Explore | Ignore |
| **High relevance** | Explore | Exploit |

*ground truth item relevance*

- When the recommender is certain it has a bad item, it ignores it.

- When the recommender is certain it has a good item, it recommends it.

# How to balance exploration and exploitation?

# How to balance exploration and exploitation?

- the central question of contextual multi-armed bandits

- standard methods include epsilon-greedy, Thompson sampling, and upper confidence bounds

# How to balance exploration and exploitation?

- the central question of contextual multi-armed bandits

- standard methods include epsilon-greedy, Thompson sampling, and upper confidence bounds

**ε-greedy**

best action A*
under preference
model prediction

$$\pi(A \mid X) = \begin{cases} (1 - \epsilon) + \dfrac{\epsilon}{|\mathscr{A}|} & \text{when } A = A* \\[2em] \dfrac{\epsilon}{|\mathscr{A}|} & \text{otherwise} \end{cases}$$

# How to balance exploration and exploitation?

- the central question of contextual multi-armed bandits

- standard methods include epsilon-greedy, Thompson sampling, and upper confidence bounds

**ε-greedy**

best action A*
under preference
model prediction

exploration parameter (when fixed -> crude exploitation; can also decay over time)

$$\pi(A \,|\, X) = \begin{cases} (1 - \epsilon) + \dfrac{\epsilon}{|\mathscr{A}|} & \text{when } A = A* \\ \dfrac{\epsilon}{|\mathscr{A}|} & \text{otherwise} \end{cases}$$

# Research question: how to explore-exploit over explainable recommendations?

# Research question: how to explore-exploit over explainable recommendations?

# Research question: how to explore-exploit over explainable recommendations?



$card_1$  $shelf_1$

$card_1$  $shelf_2$

$card_1$  $shelf_3$

# Research question: how to explore-exploit over explainable recommendations?



$card_1$   $shelf_1$

# Research question: how to explore-exploit over explainable recommendations?



$shelf_1$: $card_1$ $card_2$ $card_3$ $card_4$ $card_5$

$shelf_2$: $card_1$ $card_2$ $card_3$ $card_4$ $card_5$

$shelf_3$: $card_1$ $card_2$ $card_3$ $card_4$ $card_5$

# Research question: how to explore-exploit over explainable recommendations?



| | | | | |
|---|---|---|---|---|
| $card_1$ | $card_2$ | $card_3$ | $card_4$ | $card_5$ |

$shelf_1$

| | | | | |
|---|---|---|---|---|
| $card_1$ | $card_2$ | $card_3$ | $card_4$ | $card_5$ |

$shelf_2$

| | | | | |
|---|---|---|---|---|
| $card_1$ | $card_2$ | $card_3$ | $card_4$ | $card_5$ |

$shelf_3$

naively, the bandit has to try every possible combination of item and explanation many times before being able to exploit the best combinations

# Bart

- Bart (bandits for recommendations as treatments) consists of:
  - a user preference model conditioned on the context
  - a ranking procedure + propensities
  - a training procedure

For details, see our new publication "Explore, Exploit, Explain" at RecSys
www.jamesmc.com/s/BartRecSys.pdf

# Bart

- Bart (<u>ban</u>dits for <u>r</u>ecommendations as <u>t</u>reatments) consists of:
    - a user preference model conditioned on the context
    - a ranking procedure + propensities
    - a training procedure

factorization machine capturing interactions between features in a parameter efficient manner [Rendle, 2010]

For details, see our new publication "Explore, Exploit, Explain" at RecSys
www.jamesmc.com/s/BartRecSys.pdf

# Bart

- Bart (bandits for recommendations as treatments) consists of:

  - a user preference model conditioned on the context

  - a ranking procedure + propensities

  - a training procedure

anything we know about the user and item, including region, age group, recent listening patterns, time of day

factorization machine capturing interactions between features in a parameter efficient manner [Rendle, 2010]

For details, see our new publication "Explore, Exploit, Explain" at RecSys
www.jamesmc.com/s/BartRecSys.pdf

# Bart

- Bart (<u>ban</u>dits for <u>r</u>ecommendations as <u>t</u>reatments) consists of:

  - a user preference model conditioned on the context

  - a ranking procedure + propensities

  - a training procedure

counterfactual maximum likelihood
[Joachims & Swaminathan, 2016]

anything we know about the user and item, including region, age group, recent listening patterns, time of day

factorization machine capturing interactions between features in a parameter efficient manner [Rendle, 2010]

For details, see our new publication "Explore, Exploit, Explain" at RecSys
www.jamesmc.com/s/BartRecSys.pdf

# Bart

- Bart (bandits for recommendations as treatments) consists of:

    - a user preference model conditioned on the context
    - a <mark>ranking procedure + propensities</mark>
    - a training procedure

counterfactual maximum
likelihood
[Joachims & Swaminathan, 2016]

factorization machine capturing
interactions between features in a
parameter efficient manner [Rendle, 2010]

anything we know about the
user and item, including region,
age group, recent listening
patterns, time of day

For details, see our new publication "Explore, Exploit, Explain" at RecSys
www.jamesmc.com/s/BartRecSys.pdf

# Ranking Procedure

Let's make our lives easy: aim to train user preference model on logged impressions assumed independent given context.

| impression_id | card_id | shelf_id | context | streamed? |
|---|---|---|---|---|
| 0 | 101 | 0 | Stockholm | No |
| 1 | 3 | 0 | Stockholm | Yes |
| 2 | 45 | 1 | Stockholm | No |
| 3 | 99 | 1 | New York | No |
| 4 | 11 | 0 | New York | Yes |

# Ranking Procedure

Let's make our lives easy: aim to train user preference model on logged impressions assumed independent given context.

| impression_id | card_id | shelf_id | context | streamed? |
|---|---|---|---|---|
| 0 | 101 | 0 | Stockholm | No |
| 1 | 3 | 0 | Stockholm | Yes |
| 2 | 45 | 1 | Stockholm | No |
| 3 | 99 | 1 | New York | No |
| 4 | 11 | 0 | New York | Yes |

What set of bandit assumptions lead to this procedure?

# Ranking procedure

# Ranking procedure

## Assumptions of shelf browsing model

Horizontal scrolling

# Ranking procedure

## Assumptions of shelf browsing model

### Horizontal scrolling

#### User awareness

# Ranking procedure

## Assumptions of shelf browsing model

### Horizontal scrolling

User awareness

card-1

# Ranking procedure

## Assumptions of shelf browsing model

### Horizontal scrolling

User awareness

card-1

→ *next card*

# Ranking procedure

## Assumptions of shelf browsing model
### Horizontal scrolling

User awareness

# Ranking procedure

## Assumptions of shelf browsing model

### Horizontal scrolling

User awareness

```
┌ ─ ─ ─ ─ ─ ─ ─ ─ ┐          ┌ ─ ─ ─ ─ ─ ─ ─ ─ ┐
│  ┌──────────┐   │          │  ┌──────────┐   │
│  │          │   │          │  │          │   │
│  │  card-1  │ ──┼──▷──────┼─▷│  card-2  │   │
│  │          │   │          │  │          │   │
│  └──────────┘   │ next card│  └──────────┘   │
└ ─ ─ ─ ─ ─ ─ ─ ─ ┘          └ ─ ─ ─ ─ ─ ─ ─ ─ ┘
```

# Ranking procedure

## Assumptions of shelf browsing model

### Horizontal scrolling

User awareness

# Ranking procedure

## Assumptions of shelf browsing model
### Horizontal scrolling

User awareness

# Ranking procedure

## Assumptions of shelf browsing model
### Horizontal scrolling

User awareness

# Ranking procedure with bandit

## Horizontal scrolling

User awareness

# Ranking procedure with bandit

## Horizontal scrolling

User awareness

**Candidate set:**

card-1 card-2

card-3 $M_1$

# Ranking procedure with bandit

## Horizontal scrolling

User awareness

card-1 card-2

card-3 $M_1$

**Candidate set:**

**Action select:** $\mathrm{card}_1 \sim \pi_{s,r}(M_1)$

# Ranking procedure with bandit

## Horizontal scrolling

User awareness

card-1

**Candidate set:**

card-1  card-2

card-3  $M_1$

**Action select:**  $\mathrm{card}_1 \sim \pi_{s,r}(M_1)$

# Ranking procedure with bandit

## Horizontal scrolling

User awareness

card-1

→ *next card*

**Candidate set:**

card-1  card-2

card-3  $M_1$

**Action select:**   $\mathrm{card}_1 \sim \pi_{s,r}(M_1)$

# Ranking procedure with bandit

## Horizontal scrolling

User awareness

card-1 → *next card*

**Candidate set:** card-1 card-2 card-3 $M_1$

**Action select:** $\mathrm{card}_1 \sim \pi_{s,r}(M_1)$

# Ranking procedure with bandit

## Horizontal scrolling

User awareness

card-1

→ *next card*

**Candidate set:**

card-1 card-2 card-3 $M_1$

card-2 card-3 $M_2$

**Action select:** $\mathrm{card}_1 \sim \pi_{s,r}(M_1)$

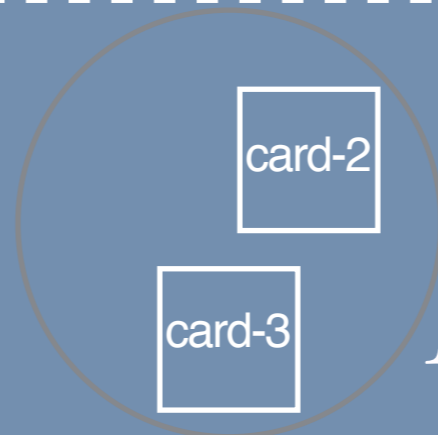# Ranking procedure with bandit

## Horizontal scrolling

User awareness



**Candidate set:**

$M_1$

$M_2$

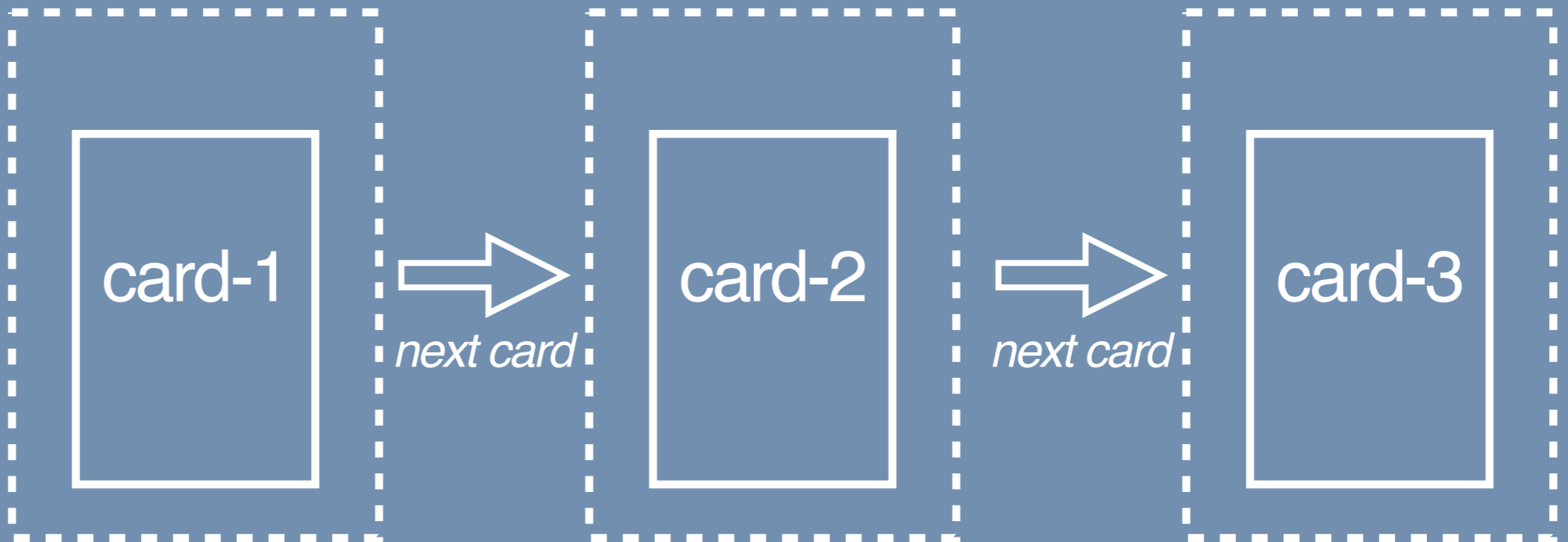**Action select:** $\text{card}_1 \sim \pi_{s,r}(M_1)$ $\quad \text{card}_2 \sim \pi_{s,r}(M_2)$

# Ranking procedure with bandit

## Horizontal scrolling

User awareness



**Candidate set:**

$M_1$

$M_2$

**Action select:** $\mathrm{card}_1 \sim \pi_{s,r}(M_1)$  $\mathrm{card}_2 \sim \pi_{s,r}(M_2)$

# Ranking procedure with bandit

## Horizontal scrolling

User awareness



**Candidate set:**

$M_1$  $M_2$

**Action select:** $\text{card}_1 \sim \pi_{s,r}(M_1)$  $\text{card}_2 \sim \pi_{s,r}(M_2)$
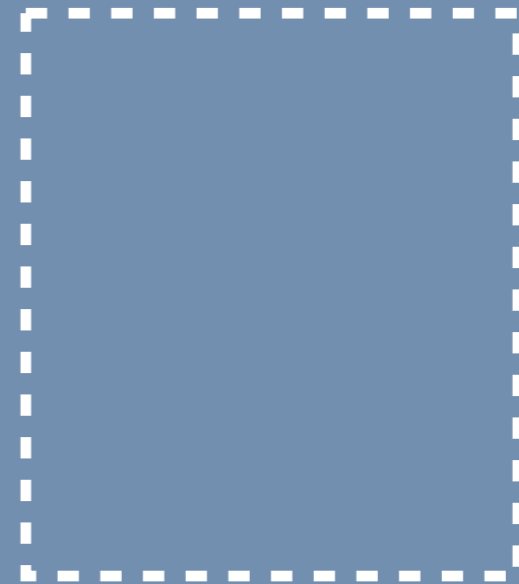
# Ranking procedure with bandit

## Horizontal scrolling

User awareness



**Candidate set:**

**Action select:** $\text{card}_1 \sim \pi_{s,r}(M_1)$ $\quad$ $\text{card}_2 \sim \pi_{s,r}(M_2)$

# Ranking procedure with bandit

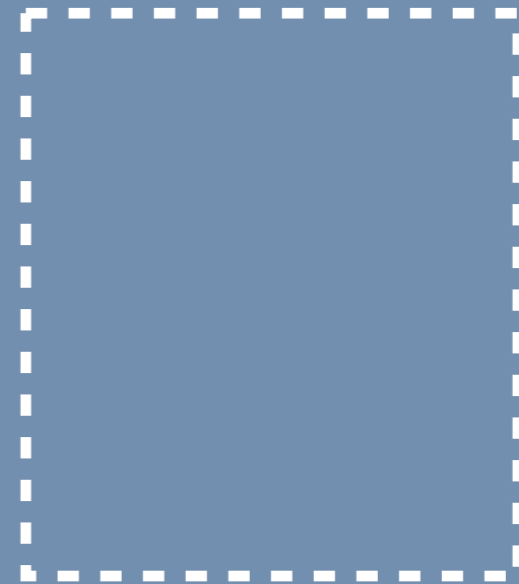## Horizontal scrolling

User awareness



**Candidate set:**

$M_1$ $M_2$ $M_3$

**Action select:** $\mathrm{card}_1 \sim \pi_{s,r}(M_1)$ $\mathrm{card}_2 \sim \pi_{s,r}(M_2)$

# Ranking procedure with bandit

## Horizontal scrolling

User awareness



**Candidate set:**

**Action select:** $\text{card}_1 \sim \pi_{s,r}(M_1)$ $\qquad$ $\text{card}_2 \sim \pi_{s,r}(M_2)$ $\qquad$ $\text{card}_3 \sim \pi_{s,r}(M_3)$

# Ranking procedure with bandit

## Horizontal scrolling

User awareness



**Candidate set:**

$M_1$ $M_2$ $M_3$

**Action select:** $\mathrm{card}_1 \sim \pi_{s,r}(M_1)$ $\mathrm{card}_2 \sim \pi_{s,r}(M_2)$ $\mathrm{card}_3 \sim \pi_{s,r}(M_3)$

# Ranking procedure with bandit

## Horizontal scrolling

User awareness



**Candidate set:**

$M_1$
$M_2$
$M_3$

repeat for each shelf

**Action select:** $\mathrm{card}_1 \sim \pi_{s,r}(M_1)$ $\mathrm{card}_2 \sim \pi_{s,r}(M_2)$ $\mathrm{card}_3 \sim \pi_{s,r}(M_3)$

# Ranking procedure with bandit

## Vertical scrolling

User awareness

# Ranking procedure with bandit

## Vertical scrolling

### Candidate set

shelf-1 shelf-2

$L_1$

### User awareness

# Ranking procedure with bandit

## Vertical scrolling

Candidate set
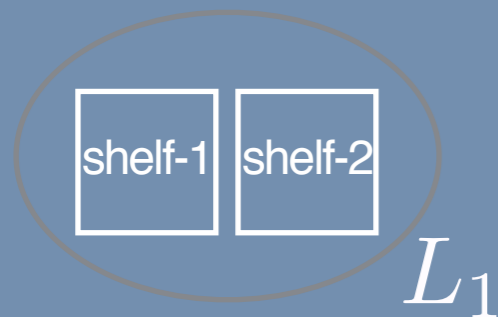
Action select

User awareness

shelf-1 shelf-2

$L_1$

$$\mathrm{shelf}_1 \sim \pi_{s,r'}(L_1)$$

# Ranking procedure with bandit

## Vertical scrolling

Candidate set

Action select

User awareness

$$\text{shelf}_1 \sim \pi_{s,r'}(L_1)$$

shelf-1 shelf-2

$L_1$

**shelf-1**

card-1,1

# Ranking procedure with bandit

## Vertical scrolling

### Candidate set

shelf-1 | shelf-2

$L_1$

### Action select

$$\text{shelf}_1 \sim \pi_{s,r'}(L_1)$$

### User awareness
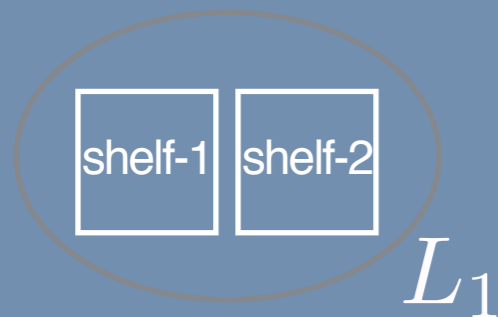
**shelf-1**

card-1,1

*next card*

*next shelf*

# Ranking procedure with bandit

## Vertical scrolling

Candidate set

Action select

User awareness

shelf-1 shelf-2

$L_1$

$$\text{shelf}_1 \sim \pi_{s,r'}(L_1)$$

**shelf-1**

card-1,1

*next card*

*next shelf*

# Ranking procedure with bandit

## Vertical scrolling

**Candidate set**

**Action select**

**User awareness**

shelf-1 shelf-2

$L_1$

$$\text{shelf}_1 \sim \pi_{s,r'}(L_1)$$

**shelf-1**

card-1,1

⟹ *next card*

⟱ *next shelf*

shelf-2

$L_2$

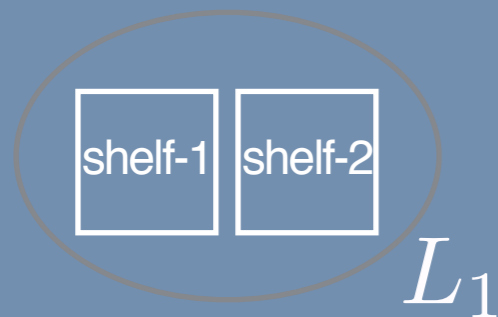# Ranking procedure with bandit

## Vertical scrolling

**Candidate set**



$L_1$

**Action select**

$$\mathrm{shelf}_1 \sim \pi_{s,r'}(L_1)$$

$$\mathrm{shelf}_2 \sim \pi_{s,r'}(L_2)$$

**User awareness**

**shelf-1**

card-1,1

*next card*

*next shelf*

$L_2$

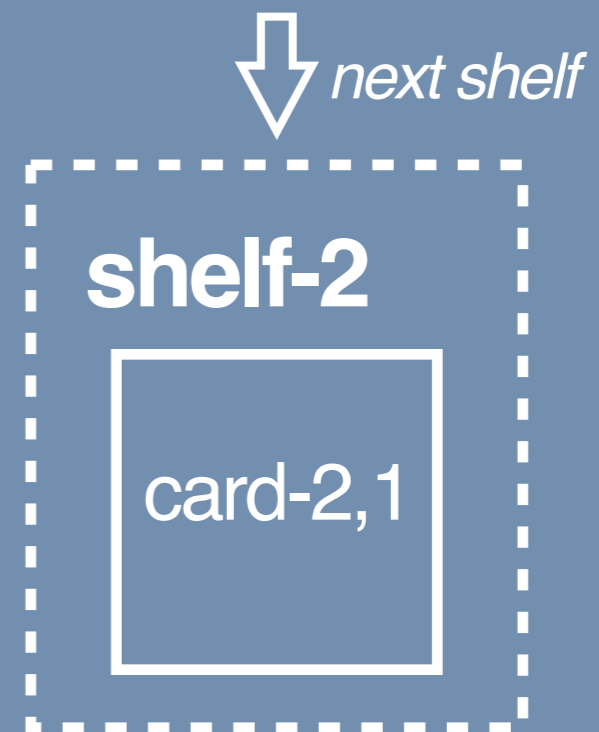# Ranking procedure with bandit

## Vertical scrolling

**Candidate set**



$$L_1$$

$$L_2$$

**Action select**

$$\text{shelf}_1 \sim \pi_{s,r'}(L_1)$$

$$\text{shelf}_2 \sim \pi_{s,r'}(L_2)$$

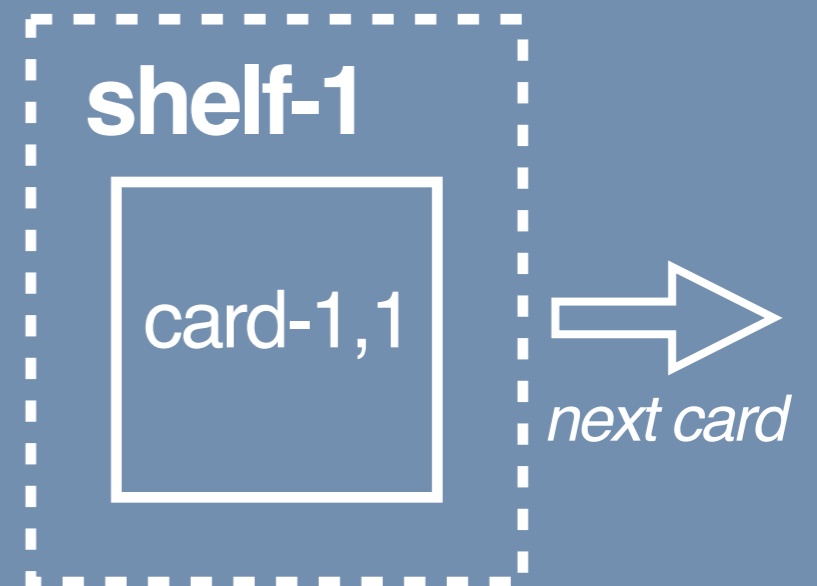**User awareness**

**shelf-1**

card-1,1

*next card*

*next shelf*

**shelf-2**

card-2,1

# Ranking procedure with bandit

## Vertical scrolling

| Candidate set | Action select | User awareness |
|---|---|---|

**shelf-1**, **shelf-2** $L_1$

$$\text{shelf}_1 \sim \pi_{s,r'}(L_1)$$

**shelf-1**

card-1,1

$\Rightarrow$ *next card*

$\Downarrow$ *next shelf*

**shelf-2** $L_2$

$$\text{shelf}_2 \sim \pi_{s,r'}(L_2)$$

**shelf-2**

card-2,1
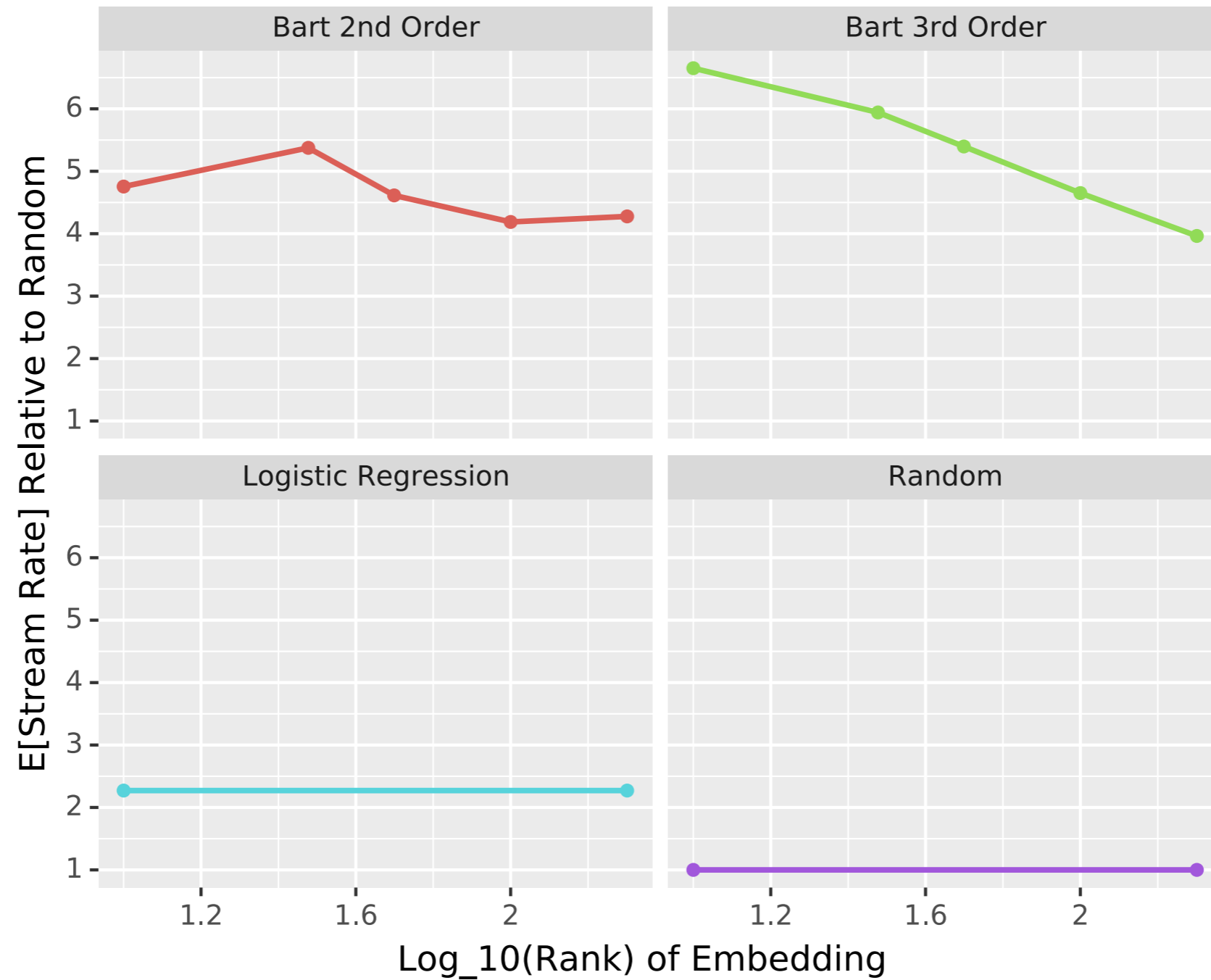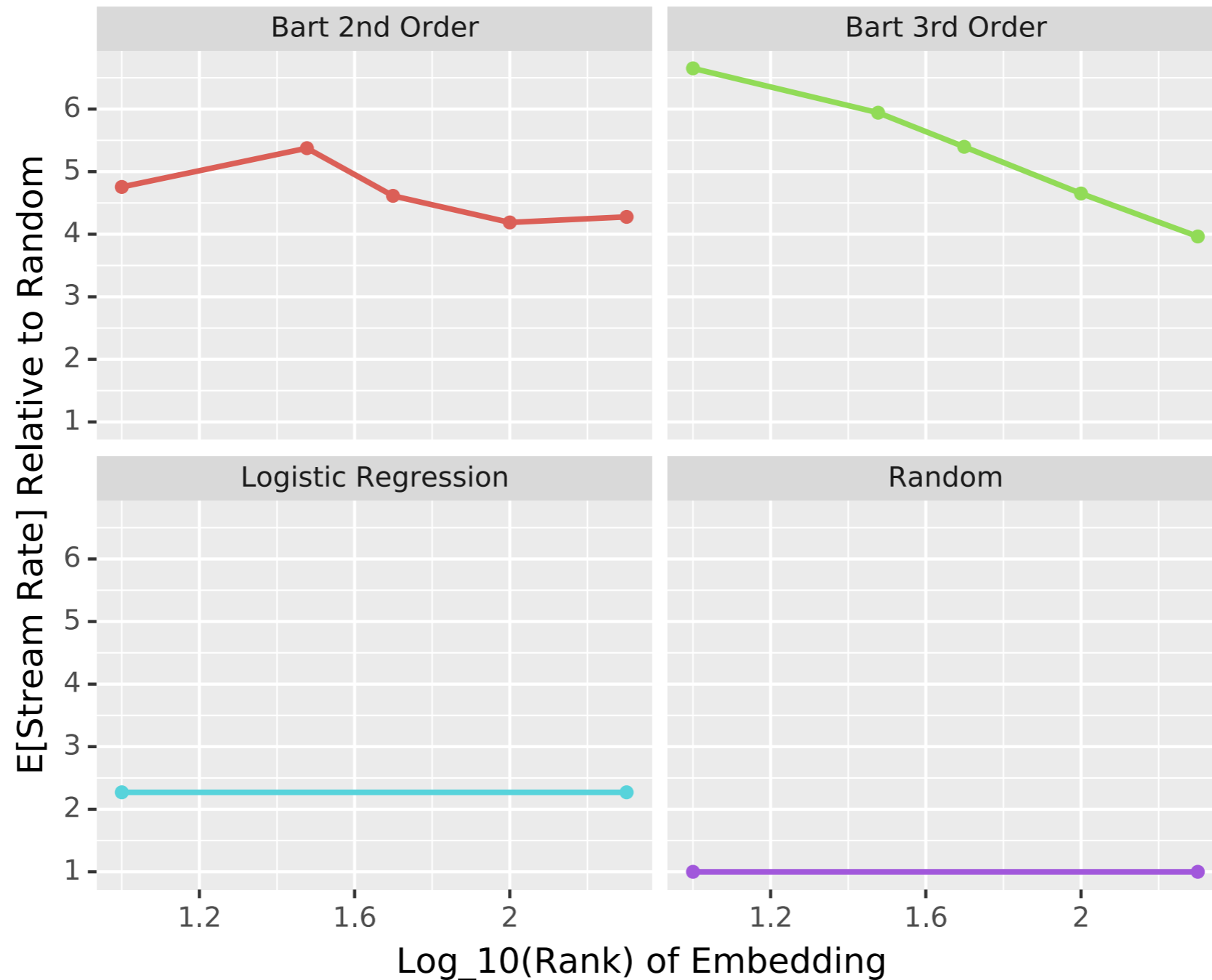
etc.

# Experimental evaluation

- we collected <u>randomized</u> recommendation data
- offline experiments:
  - counterfactual estimation of A/B test performance using importance sampling reweighting
- online A/B test experiments

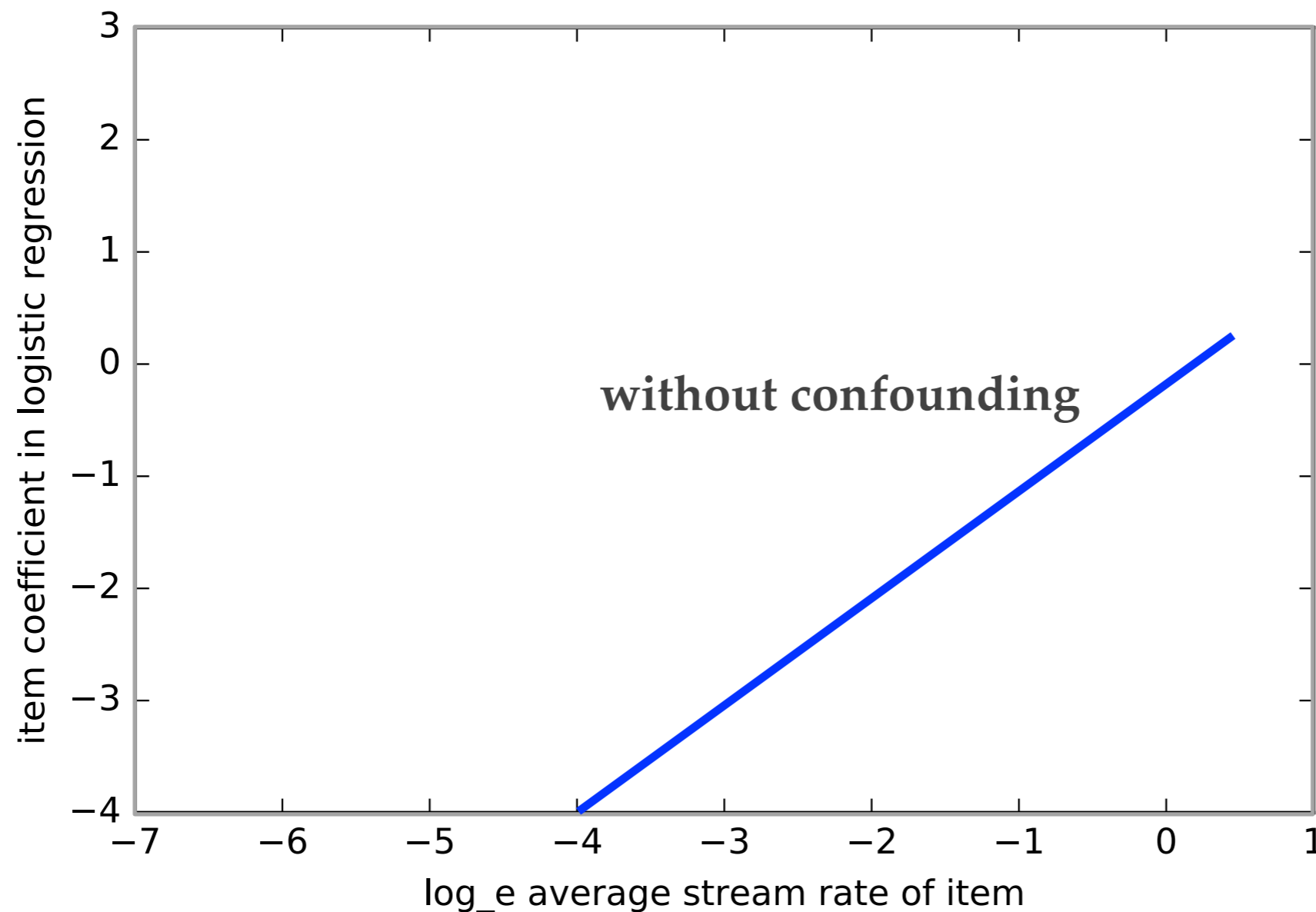# Offline experiments

# Offline experiments

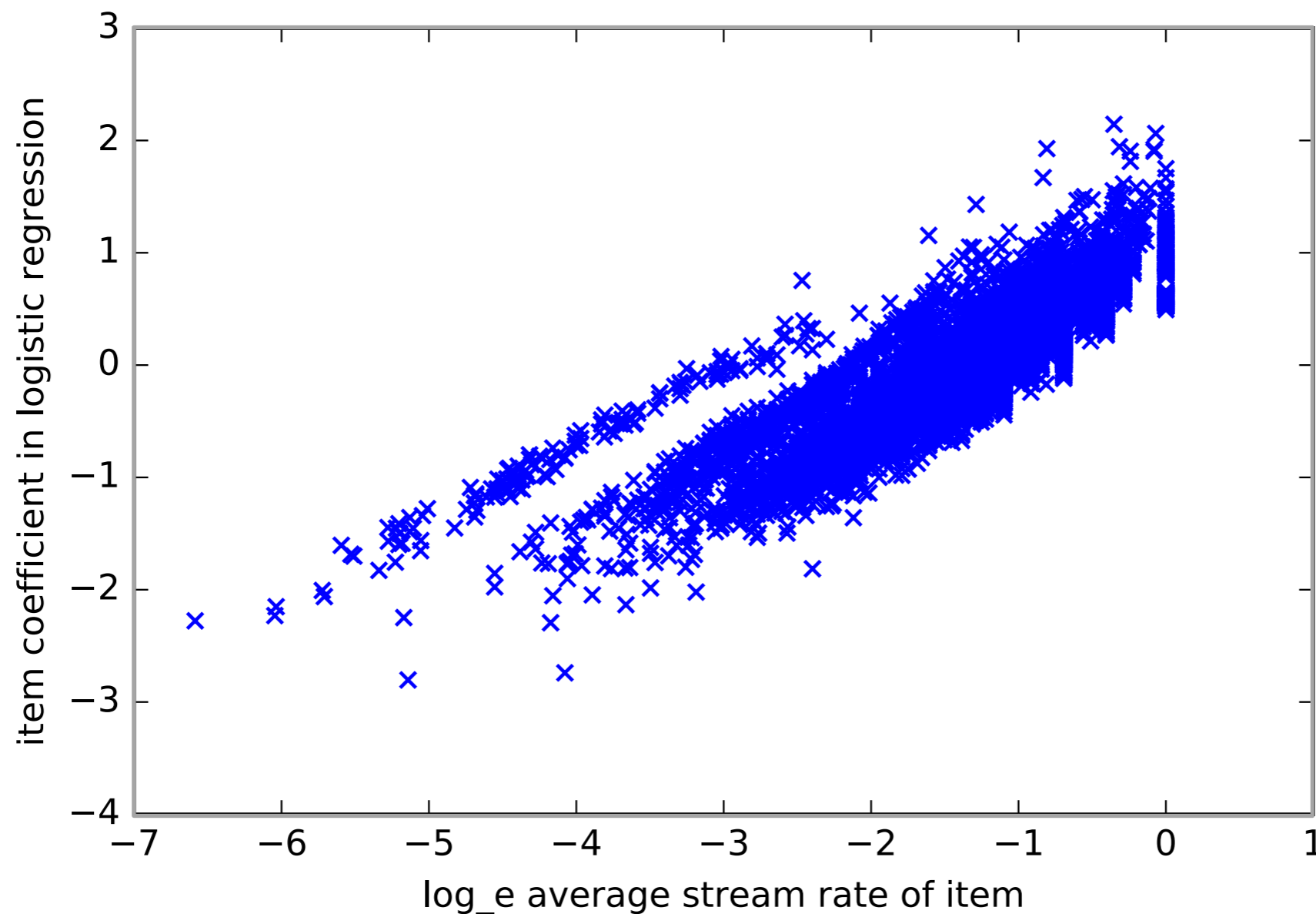

(similar conclusions as NDCG@10 for the metric)

# Offline experiments

- how does the empirical stream rate of an item relate to its stream rate controlling for other factors?

# Offline experiments

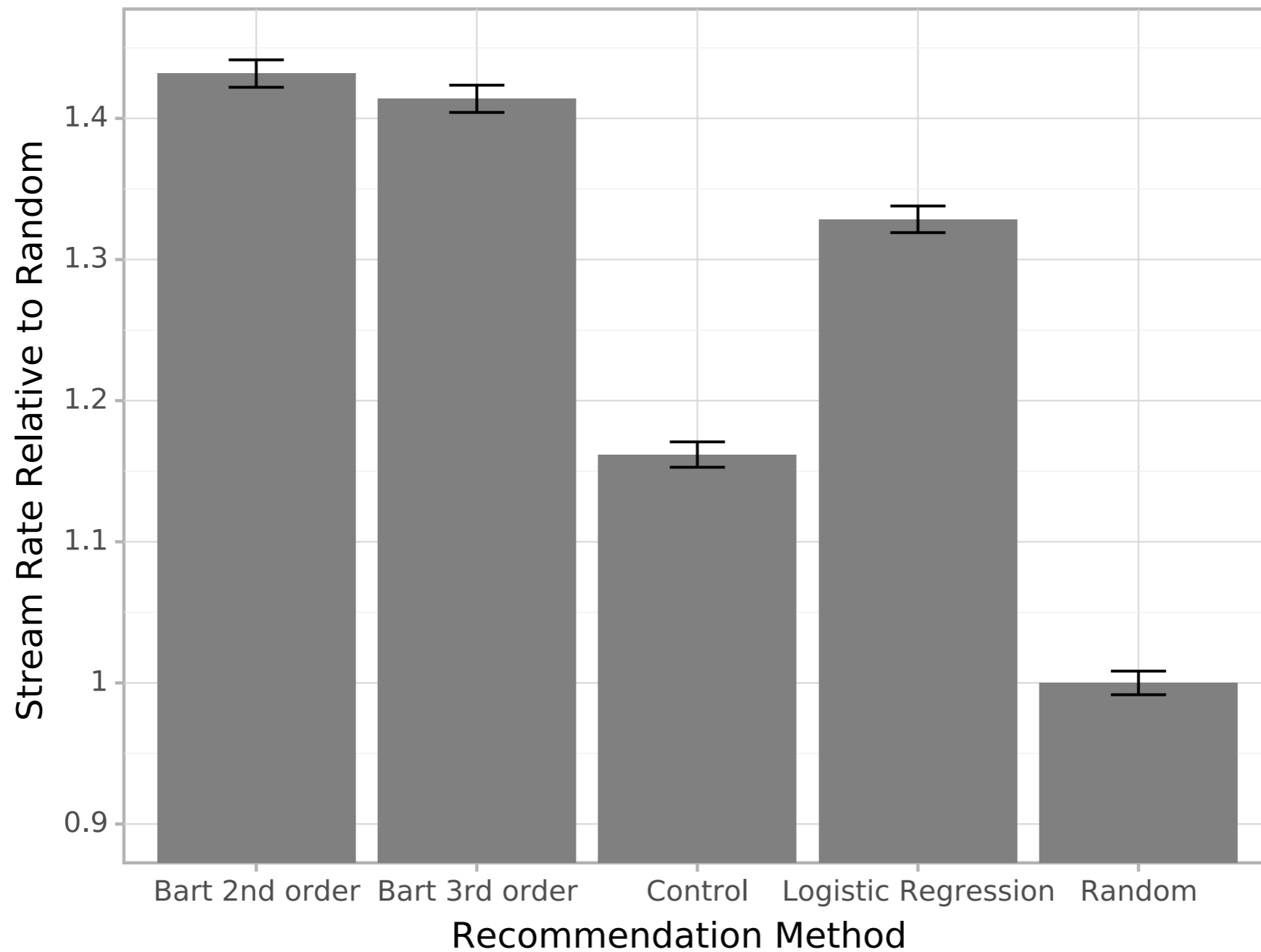- how does the empirical stream rate of an item relate to its stream rate controlling for other factors?

# Offline experiments

- how does the empirical stream rate of an item relate to its stream rate controlling for other factors?

# Bart limitations and future work

- user preference model:
  - assumes independence of impression outcomes
  - attempts to estimate absolute reward, competitive pairwise model closer to how humans judge items
  - maximizes our defined reward, does it approximate user satisfaction?
- ranking model not defined to promote diversity
- exploration-exploitation over a candidate set not the full item set

# Is bandits a good idea for your problem?

Things to consider:

- <u>confounding</u>: are you training a model using data collected with another model?

  - consider counterfactual evaluation on its own; less need to explore/exploit

- <u>auto-confounding</u>: are you repeatedly training a model using data generated by the same model?

  - consider counterfactual evaluation and explore/exploit

# Thank you, any questions?

email: jamesm@spotify.com