McInerney, James, Benjamin Lacker, Samantha Hansen, Karl Higley, Hugues Bouchard, Alois Gruson, and Rishabh Mehrotra. "Explore, exploit, and explain: personalizing explainable recommendations with bandits." In Proceedings of the 12th ACM Conference on Recommender Systems, pp. 31-39. ACM, 2018.
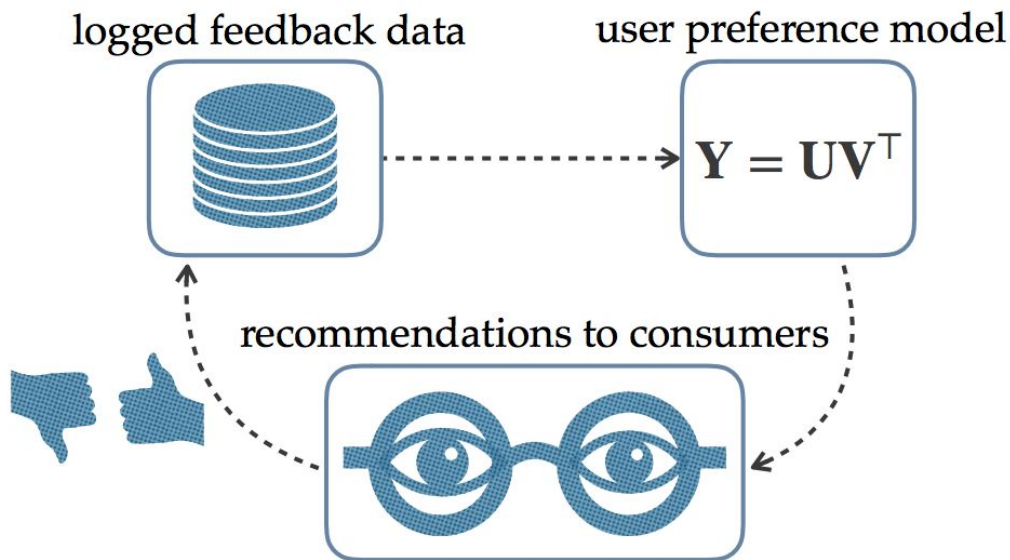
Slides

- Spotify

- Factorization Machine

- Epsilon-greedy

- IPS(inverse propensity score)

# Factorization Machine

$$r^{(2)}(j, e, x) = \sigma\left(\theta_{\text{global}} + \theta^\top x' + \sum_{a=1}^{D} \sum_{b>a}^{D} v_a^\top v_b x_a' x_b'\right) \quad (3)$$

# 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

# Exploration-Exploitation Policy & Propensity Scoring

$$\pi_c^{\text{item}}(j \mid x, e) = \begin{cases} (1 - \epsilon) + \frac{\epsilon}{|f(e,u)|}, & \text{if } j = j^*, j \in f(e,x) \\ \frac{\epsilon}{|f(e,u)|}, & \text{if } j \neq j^*, j \in f(e,x) \\ 0, & \text{otherwise.} \end{cases}$$

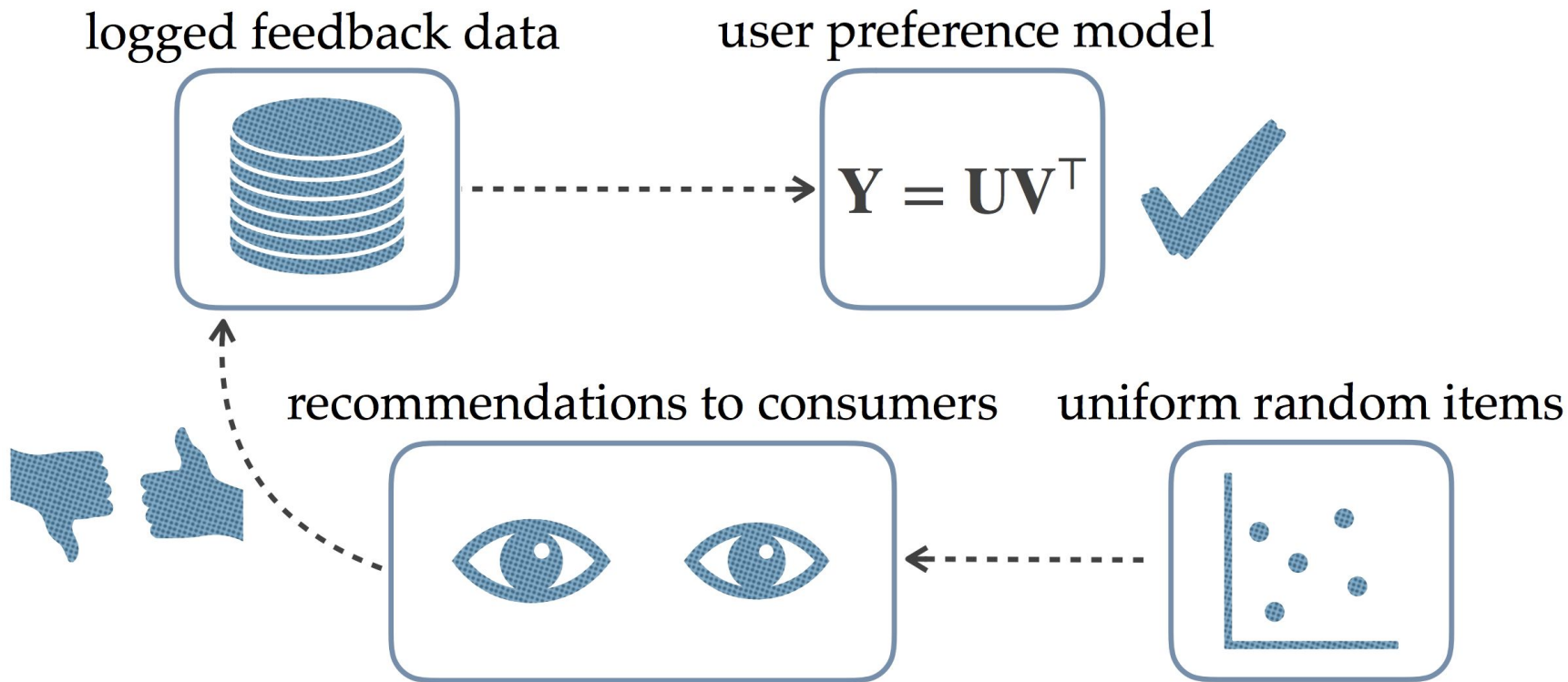$$\text{where } j^* = \arg_{j_1} \max r(j_1, e, x) \tag{7}$$

$$\pi_c^{\text{expl.}}(e \mid x, j) = \begin{cases} (1 - \epsilon) + \frac{\epsilon}{|\mathcal{E}|}, & \text{if } e = e^*, j \in f(e,x) \\ \frac{\epsilon}{|\mathcal{E}|}, & \text{if } e \neq e^*, j \in f(e,x) \\ 0, & \text{otherwise.} \end{cases}$$

$$\text{where } e^* = \arg_{e_1} \max r(j, e_1, x). \tag{8}$$

# 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

# Off-Policy Training

$$\hat{\theta}, \hat{v} = \arg_{\theta,v} \max \mathbb{E}_{A\sim\text{Uniform}(\cdot)}[\mathbb{E}_{X,R}[\log p_{\theta,v}(R|A,X)]] \qquad (5)$$
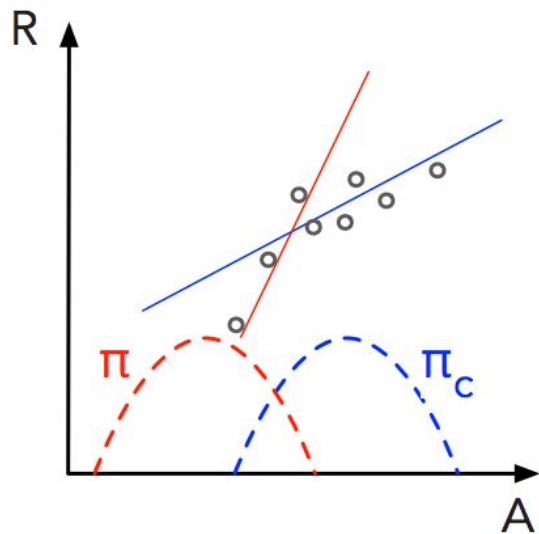
$$\approx \arg_{\theta,v} \max \frac{1}{N} \sum_{n=1}^{N} \frac{\text{Uniform}(a_n)}{\pi_c(a_n)} \log p_{\theta,v}(r_n|a_n,x_n). \qquad (6)$$

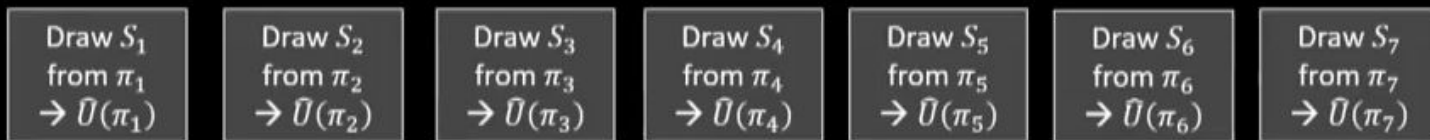IPS(inverse propensity score)

# Off-Policy Training



**Figure 3: Off-policy training fits a reward function that best fits the input points with respect to the target policy $\pi$ using input points generated by the collection policy $\pi_c$.**
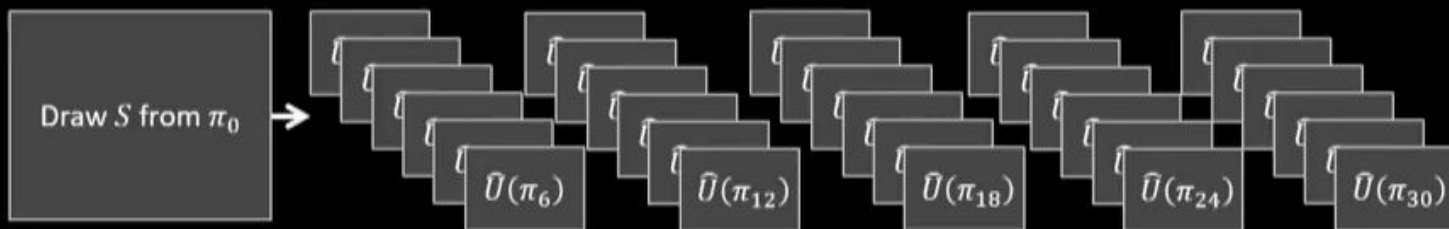
# Causal Inference Recommendation Papers

- Schnabel, Tobias, Adith Swaminathan, Ashudeep Singh, Navin Chandak, and Thorsten Joachims. "Recommendations as treatments: Debiasing learning and evaluation." arXiv preprint arXiv:1602.05352 (2016).
  - Propensity-Scored Matrix Factorization (= IPS + MF)
- Liang, Dawen, Laurent Charlin, and David M. Blei. "Causal Inference for Recommendation." (2016).
  - IPS + MF
- Bonner, Stephen, and Flavian Vasile. "Causal embeddings for recommendation." In Proceedings of the 12th ACM Conference on Recommender Systems, pp. 104-112. ACM, 2018.
  - Criteo, MF + Counterfactual Risk Minimization(CRM)
- Gilotte, Alexandre, Clément Calauzènes, Thomas Nedelec, Alexandre Abraham, and Simon Dollé. "Offline A/B testing for Recommender Systems." In Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, pp. 198-206. ACM, 2018.
  - Criteo, IPS
- Joachims, Thorsten, Adith Swaminathan, and Maarten de Rijke. "Deep learning with logged bandit feedback." (2018).
  - Unbiased estimate of risk => IPS
  - Partian Information => Variacne Control
  - Propensity Overfitting => SNIPS(self-normalized IPS estimator)

# Evaluating Online Metrics Offline



- Online: On-policy A/B Test

  Draw $S_1$ from $\pi_1$ → $\hat{U}(\pi_1)$ | Draw $S_2$ from $\pi_2$ → $\hat{U}(\pi_2)$ | Draw $S_3$ from $\pi_3$ → $\hat{U}(\pi_3)$ | Draw $S_4$ from $\pi_4$ → $\hat{U}(\pi_4)$ | Draw $S_5$ from $\pi_5$ → $\hat{U}(\pi_5)$ | Draw $S_6$ from $\pi_6$ → $\hat{U}(\pi_6)$ | Draw $S_7$ from $\pi_7$ → $\hat{U}(\pi_7)$

- Offline: Off-policy Counterfactual Estimates

  Draw $S$ from $\pi_0$ → $\hat{U}(\pi_6)$ | $\hat{U}(\pi_{12})$ | $\hat{U}(\pi_{18})$ | $\hat{U}(\pi_{24})$ | $\hat{U}(\pi_{30})$
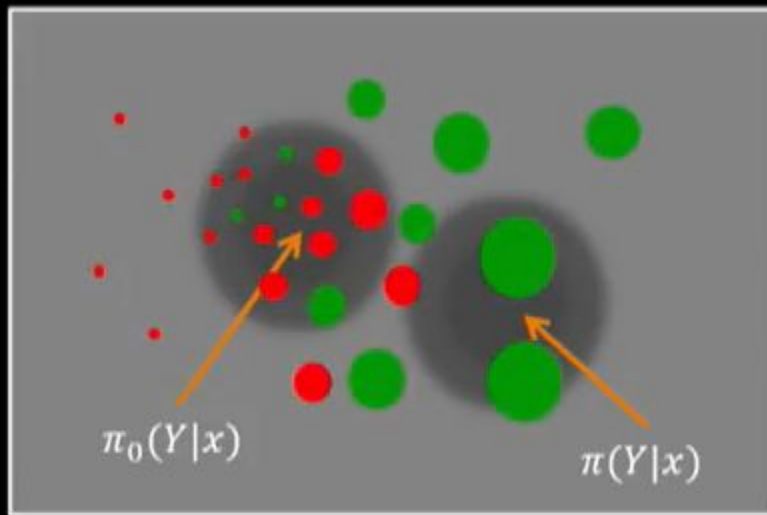
# Off-Policy Risk Evaluation



Given $S = \left( (x_1, y_1, \delta_1), \ldots, (x_n, y_n, \delta_n) \right)$ collected under $\pi_0$,

$$\hat{R}(\pi) = \frac{1}{n} \sum_{i=1}^{n} \delta_i \frac{\pi(y_i|x_i)}{\pi_0(y_i|x_i)}$$

Propensity $p_i$

$\pi_0(Y|x)$

$\pi(Y|x)$

→ Unbiased estimate of risk, if propensity nonzero everywhere (where it matters).

# Partial Information Empirical Risk Minimization



- Training

$$\hat{\pi} := \mathrm{argmin}_{\pi \in H} \sum_i^n \frac{\pi(y_i|x_i)}{p_i} \delta_i$$

# Variance Control

$$R(\pi) \leq \hat{R}(\pi) + O\left(\sqrt{\widehat{Var}(\pi)/n}\right) + O(C)$$

Unbiased Estimator

Variance Control

Capacity Control

$$\hat{R}(\pi) = \widehat{Mean}\left(\frac{\pi(y_i|x_i)}{p_i}\delta_i\right)$$

$$\widehat{Var}(\pi) = \widehat{Var}\left(\frac{\pi(y_i|x_i)}{p_i}\delta_i\right)$$

# Problem: Propensity Overfitting



- Example
  - Training sample with losses:
  - Which $\pi(y|x)$ minimize IPS?

$$R(\pi) = \min_{\pi \in H} \frac{1}{n} \sum_{i}^{n} \frac{\pi(y_i|x_i)}{p_i} \delta_i$$

# EQUIVARIANT COUNTERFACTUAL RISK MINIMIZATION

$$\hat{R}_{SNIPS}(\pi_w) \;=\; \frac{\frac{1}{n}\sum_{i=1}^{n}\delta_i\frac{\pi_w(y_i|x_i)}{\pi_0(y_i|x_i)}}{\frac{1}{n}\sum_{i=1}^{n}\frac{\pi_w(y_i|x_i)}{\pi_0(y_i|x_i)}}.$$

# TRAINING ALGORITHM

$$w = \underset{w \in \mathfrak{R}^N}{\text{argmin}} \left[ \hat{R}^{SNIPS}(w) + \lambda_1 \sqrt{\widehat{Var}\left(\hat{R}^{SNIPS}(w)\right)} + \lambda_2 ||w||^2 \right]$$

# More about causal inference

- Ryan Tibshirani's lecture notes on causal inference
  - http://www.stat.cmu.edu/~larry/=sml/Causation.pdf
- Hernán, Miguel A. and James M. Robins. 2012. Causal Inference. Forthcoming, Cambridge University Press.
  - https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/