



Offline Evaluation of Ranking Policies with Click Models

Shuai Li

The Chinese University of Hong Kong

Joint work with

Yasin Abbasi-Yadkori (Adobe Research)

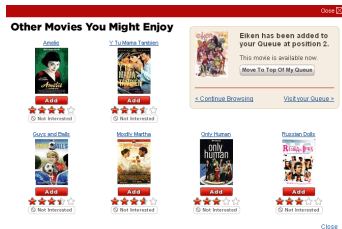
Branislav Kveton (Google Research, was in Adobe Research)

S. Muthukrishnan (Rutgers University)

Vishwa Vinay (Adobe Research)

Zheng Wen (Adobe Research)

Motivation



Amazon, Facebook, Netflix

Production Policy



Number of clicks: $f(\) = 1$

- How can we know $f(\zeta) = 2$?

Hypothetical Policy ζ



Number of clicks: $f(\zeta) = 2$



Directly Implementing New Policy ζ

Risks for directly implementing new policy ζ

- Expensive
- Uses a portion of live users and poor policy might harm user experience
- Not replicable



Can we know $f(\zeta) = 2$ without directly implementing it?

Offline Evaluation!

Setting

- Lists: $\mathcal{Y} = (\check{Y}_1; \dots; \check{Y}_T)$

- The value of list \mathcal{Y} with the click realization \check{U} :

$$f(\mathcal{Y}; \check{U}) = \prod_{\check{U}=1}^{\check{U}} (\check{Y}_{\check{U}}; \check{U})$$

- The value of a policy ζ :

$$f(\mathcal{Q}) = E_{\zeta(j)} [f(\mathcal{Y}; \check{U})]$$

- Suppose the clicks depend only on (item, position) pairs
- The CTR of putting item \check{Y} at \check{U} th position under context \check{U}_j is

$$f(\check{Y}; \check{U}_j)$$

- The expected value of

$$f(\check{Y}) = \sum_{\check{U}=1}^{\check{X}} f(\check{Y}_{\check{U}}; \check{U}_j)$$

Setting (Problem Definition)

Logged dataset $k = f(\cdot; \cdot; \cdot)_{g_{=1}^{\bar{a}}}$

- At each time
 - The environment draws context and click realizations
 - The learner observes and selects according to policy
 - The environment reveals $f(\cdot; \cdot; \cdot)_{g_{U=1}^I}$

Objective

- To design statistically efficient estimators based on logged dataset for any ranking policy

Challenge

- The number of different lists is exponential in I

- Direct Method

$$\hat{f}(\zeta) = \frac{1}{\tilde{a}} \sum_{j=1}^{\tilde{a}} \sum_{\tilde{y}} \sum_{\tilde{u}=1}^{\tilde{a}} \zeta(\tilde{y}; \tilde{u}_j) \wedge (\tilde{y}; \tilde{u}_j)$$

- Can be used to evaluate any policy
- Unstable when the number of observations for some item is small
- No theoretical guarantee for known computationally efficient method for some click models

- Importance sampling (for list level)

$$\begin{aligned}
 f(Q) &= E_{\zeta} [f(\cdot)] \\
 &= E_{\zeta} \left[f(\cdot) \frac{(\cdot)}{(\cdot)} \right] \\
 &= E_{\zeta} \left[f(\cdot) \frac{\zeta(\cdot)}{(\cdot)} \right]
 \end{aligned}$$

- List estimator

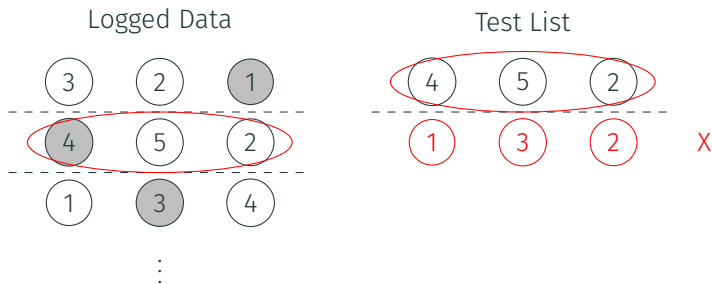
$$\hat{f}_L(Q) = \frac{1}{|j|k|} \sum_{(\cdot; \cdot) \in \mathcal{Z}_k} f(\cdot; \cdot) \min_{\zeta} \frac{\zeta(\cdot; \cdot)}{\hat{(\cdot; \cdot)}}; Q$$

Trade-off between bias and variance

Empirical distribution over logged data

List Estimator - Disadvantages

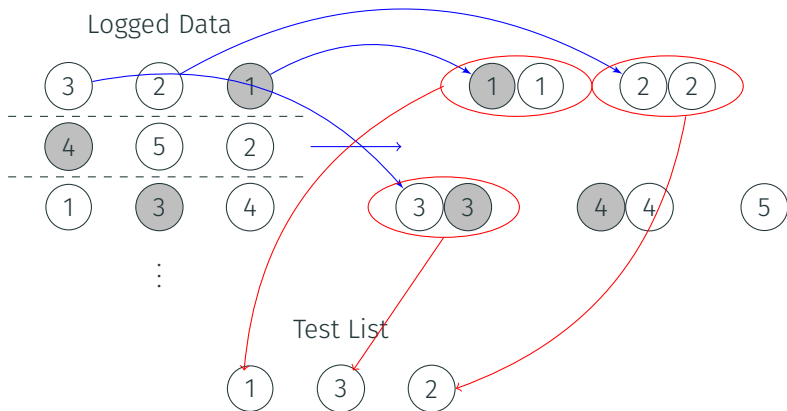
- Disadvantages
 - Have to match the exact lists
 - The number of lists is exponential in I , thus $\hat{p}(j)$ is very small



Click No Click

Document-Based Click Model (DCTR)

- $(\hat{Y}; \hat{U}_j)$ only depends on item \hat{Y} , for any context



- DCTR: $(\tilde{Y}; \tilde{U}j)$ only depends on item \tilde{Y} for any context

$$\hat{f}_i(Q) = \frac{1}{jkj} \sum_{(; ;)_{2k} \tilde{U}=1}^{\times \times} (\tilde{Y}; \tilde{U}) \min \frac{\zeta(\tilde{Y}\tilde{U}j)}{\wedge(\tilde{Y}\tilde{U}j)}; Q$$

- List estimator

$$\hat{f}_L(Q) = \frac{1}{jkj} \sum_{(; ;)_{2k} \tilde{U}=1}^{\times \times} (\tilde{Y}; \tilde{U}) \min \frac{\zeta(j)}{\wedge(j)}; Q$$

Estimators for Click Models

Click Model	Assumption	Estimator
Random	$(\check{Y}; \hat{U}_j)$ constant	\hat{f}_R
Rank-Based	$(\check{Y}; \hat{U}_j)$ only depends on position \hat{U}	\hat{f}_R
Document-Based	$(\check{Y}; \hat{U}_j)$ only depends on item \check{Y}	\hat{f}_I
Position-Based	$(\check{Y}; \hat{U}_j) = (\check{Y}_j) \cdot \check{U}_j$	\hat{f}_{PBM}
Item-Position	$(\check{Y}; \hat{U}_j)$	\hat{f}_{IP}

Proposition (Unbiased in a larger class of policies)

$q_{\mathcal{C}} \cdot \dot{y} \dot{u} - \dot{u}^3 \cdot \dot{y} \hat{\mathcal{L}} \ddot{Y} \hat{e} \ddot{y} \ddot{Y} \dot{u} \quad \tilde{a} - \hat{\mathcal{E}} \dot{y} \cdot \hat{\mathcal{L}} \ddot{Y} \ddot{Y} \hat{u} \dot{u} - \ddot{Y} \dot{y} \hat{e} \hat{A} \ddot{e} \hat{\mathcal{L}} \hat{\mathcal{E}} \dot{y}$
 $\mathcal{C} \ddot{Y} \hat{a} \hat{\mathcal{L}} \dot{y} \hat{\mathcal{L}} \ddot{Y} \hat{e} \dot{u}$

Proposition (Lower bias in estimating policy)

$q_{\mathcal{C}} \cdot \dot{y} \dot{u} - \dot{u}^3 \cdot \dot{y} \hat{\mathcal{L}} \ddot{Y} \hat{e} \ddot{y} \mathcal{C} \ddot{Y} \cdot \hat{\mathcal{L}} \dot{u} - \hat{\mathcal{E}} \dot{y} \mathcal{C} \ddot{Y} \hat{a} \hat{\mathcal{L}} \dot{y} \hat{\mathcal{L}} \ddot{Y} \hat{e} \dot{u}$

Proposition (Better guarantee for policy optimization)

$q_{\mathcal{C}} \cdot \dot{y} \ddot{e} \hat{\mathcal{L}} \hat{\mathcal{E}} \hat{a}^3 \dot{y} \dot{u} - \dot{u}^3 \cdot \dot{y} \hat{\mathcal{L}} \ddot{Y} \hat{e} \ddot{y} \mathcal{C} \ddot{Y} \cdot \dot{y} \cdot \dot{u}$
 $\mathcal{C} \hat{e} \dot{u} \hat{\mathcal{E}} \ddot{Y} \hat{u} \hat{Y} \ddot{a} \cdot \dot{y} \mathcal{C} \ddot{Y} \hat{a} \mathcal{C} \ddot{Y} \dot{y} \hat{\mathcal{L}} \ddot{Y} \hat{e} \dot{u}$



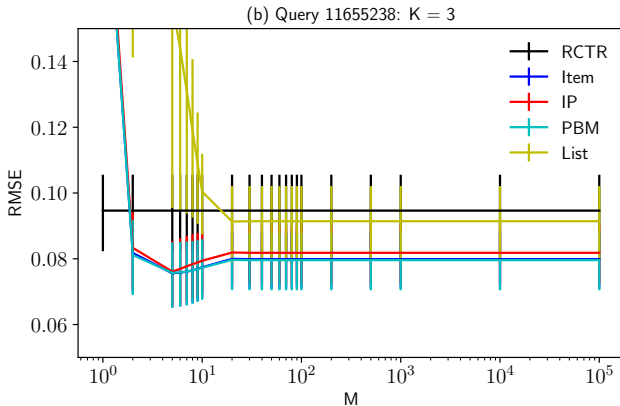
Personalized Web Search Challenge

Re-rank web documents using personal preferences

\$9,000 · 194 teams · 5 years ago

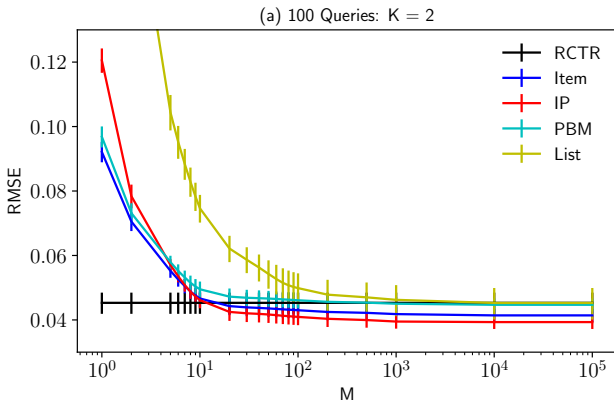
- Recorded over 27 days
- Each record contains
 - A query ID
 - The day when the query occurs
 - 10 displayed item as a response to the query
 - The corresponding click indicators of each displayed items
- **Logged dataset k**
 - Any record except day 3
 - $\hat{\cdot}$ is the empirical distribution over k
- **Evaluation policy ζ**
 - Records of day 3
 - ζ is the empirical distribution over these records
 - $f(\zeta)$ is the average CTR for these records

Experiments - Example Query with $I = 3$



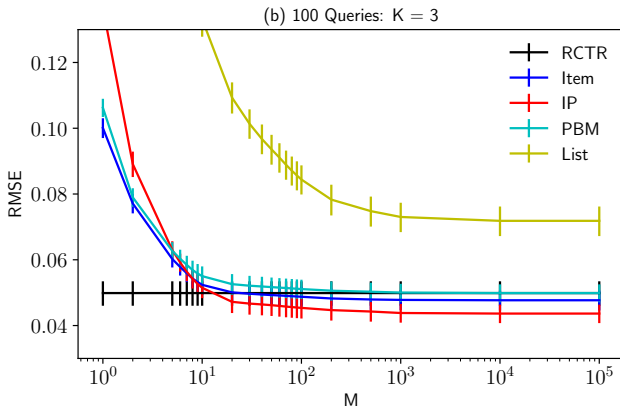
- Structured estimators better
- Tuning of Q matters

Experiments - 100 Most Frequent Queries with $I = 2$



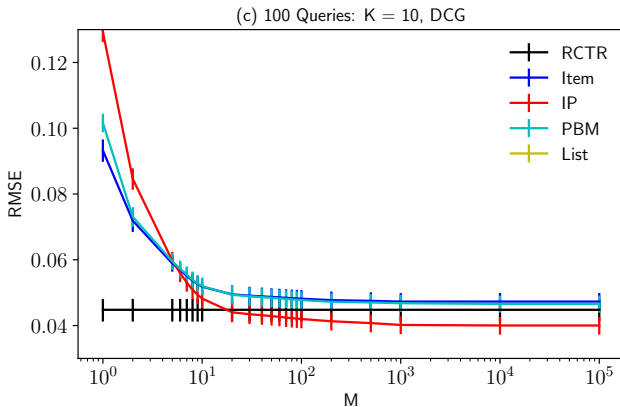
- IP estimator improves 18% over list estimator
- IP estimator improves 13% over RCTR estimator

Experiments - 100 Most Frequent Queries with $I = 3$



- IP estimator improves 46% over list estimator
- IP estimator improves 13% over RCTR estimator




Experiments - 100 Most Frequent Queries with $I = 10$, DCG



- IP estimator improves **82%** over list estimator
- IP estimator improves **11%** over RCTR estimator

Conclusions

- We propose various estimators for the expected number of clicks on lists generated by ranking policies that leverage the structure of click models
- We prove that our estimators are better than the unstructured list estimators
 - Less biased
 - Better guarantees for policy optimization
- Our estimators consistently outperform the list estimator in experiments

-  T. Joachims, A. Swaminathan, and T. Schnabel.
Unbiased learning-to-rank with biased feedback.
In *Proceedings of the 31st International Conference on Machine Learning (ICML 2017)*, pages 781–789. ACM, 2017.
-  A. Strehl, J. Langford, L. Li, and S. M. Kakade.
Learning from logged implicit exploration data.
In *Proceedings of the 31st International Conference on Machine Learning (ICML 2017)*, pages 2217–2225, 2010.
-  A. Swaminathan, A. Krishnamurthy, A. Agarwal, M. Dudik, J. Langford, D. Jose, and I. Zitouni.
Off-policy evaluation for slate recommendation.
In *Proceedings of the 31st International Conference on Machine Learning (ICML 2017)*, pages 3632–3642, 2017.