

Decomposition and Interleaving for Variance Reduction of Post-click Metrics

Kojiro Iizuka
Gunosy Inc. / University of Tsukuba
iizuka.kojiro@gmail.com

Yoshifumi Seki
Gunosy Inc.
yoshifumi.seki@gunosy.com

Makoto P. Kato
University of Tsukuba / JST, PRESTO
mpkato@acm.org

ABSTRACT

In this study, we propose an efficient method for comparing the post-click metric (e.g., dwell time and conversion rate) of multiple rankings in online experiments. The proposed method involves (1) the decomposition of the post-click metric measurement of a ranking into a click model estimation and a post-click metric measurement of each item in the ranking, and (2) interleaving of multiple rankings to produce a single ranking that preferentially exposes items possessing a high population variance. The decomposition of the post-click metric measurement enables the free layout of items in a ranking and focuses on the measurement of the post-click metric of each item in the multiple rankings. The interleaving of multiple rankings reduces the sample variance of the items possessing a high population variance by optimizing a ranking to be presented to the users so that those items received more samples of the post-click metric. In addition, we provide a proof that the proposed method leads to the minimization of the evaluation error in the ranking comparison and propose two practical techniques to stabilize the online experiment. We performed a comprehensive simulation experiment and a real service setting experiment. The experimental results revealed that (1) the proposed method outperformed existing methods in terms of efficiency and accuracy, and the performance was especially remarkable when the input rankings shared many items, and (2) the two stabilization techniques successfully improved the evaluation accuracy and efficiency.

CCS CONCEPTS

• Information systems → Retrieval efficiency; Retrieval effectiveness.

KEYWORDS

interleaving, online evaluation, post-click metrics

ACM Reference Format:

Kojiro Iizuka, Yoshifumi Seki, and Makoto P. Kato. 2021. Decomposition and Interleaving for Variance Reduction of Post-click Metrics. In *Proceedings of the 2021 ACM SIGIR International Conference on the Theory of Information Retrieval (ICTIR '21)*, July 11, 2021, Virtual Event, Canada. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3471158.3472235>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICTIR '21, July 11, 2021, Virtual Event, Canada

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8611-1/21/07...\$15.00

<https://doi.org/10.1145/3471158.3472235>

1 INTRODUCTION

Online controlled experiments are conducted daily to evaluate search and recommendation algorithms. A/B testing is a common approach that compares two different outcomes by showing them to two different user groups. A/B testing is applicable for a variety of purposes, including measurement of click-based metrics (e.g., click-through rate (CTR)), and measurement of *post-click metrics* (e.g., music listening time [12], news reading time [18, 33], and the number of reservations [14]). As post-click metrics are closely related to user satisfaction and the sales of services, the measurement of post-click metrics is particularly important for the continuous improvement of algorithms in search and recommender systems.

However, there are some challenges in measuring post-click metrics in online experiments. Suppose that a news article is ranked at the bottom of a ranking, for which users spend a significantly different length of time to read. Generally, low-ranked items are infrequently clicked. Thus, the post-click metric results in high variance of the sample mean. Moreover, some types of post-click metrics highly depending on users are of high *population* variance by nature, and also likely to result in a high variance of the sample mean. High variance in each item naturally leads to high variance for the whole ranking. This high variance can prevent online experiments from efficiently discriminating between competitive rankings.

This study proposes an efficient method for comparing the multiple rankings of post-click metrics in online experiments. The key ideas of the proposed method are (1) the decomposition of the post-click metric measurement of a ranking into a click model estimation and post-click metric measurement of each item in the ranking, and (2) interleaving multiple rankings to produce a single ranking that preferentially exposes items with high population variances. The decomposition of the post-click metric measurement permits free layout of items in a ranking and focus on the measurement of the post-click metric of each item in multiple rankings. The interleaving of multiple rankings reduces the variance of the sample mean for items with a high population variance by optimizing the ranking to be presented so that those items received more samples of the post-click metric. In this paper, the method is referred to as the *Decomposition and Interleaving for Reducing the Variance of post-click metrics (DIRV)*. DIRV uses interleaving to evaluate post-click metrics for the variance reduction, which is a major distinction from interleaving designed to evaluate click-based metrics.

In addition to the proposal of DIRV, we have the following theoretical and technical contributions in this work. First, we proved that the ranking optimization by the DIRV leads to minimization of the evaluation error in the ranking comparison. Second, we propose two techniques to stabilize the evaluation by the DIRV. The first technique predicts the population variance of each item based

on the observed samples and item features for stabilizing the ranking optimization at the beginning of the online experiment. The second technique corrects the systematic error between the estimated post-click metrics and the ground-truth post-click metrics.

We performed comprehensive experiments in simulation as well as real service settings. The experimental results revealed that (1) the proposed method outperformed existing methods in terms of efficiency and accuracy, and the performance was especially remarkable when the input rankings shared many items, (2) the two stabilization techniques successfully improved the evaluation accuracy and efficiency.

The major contributions of this study are as follows:

- To efficiently compare post-click metrics of multiple rankings, we proposed an interleaving method (DIRV) that decomposes the post-click metric measurement and preferentially exposes items with high population variance to minimize the evaluation error.
- We provided a theoretical guarantee that the DIRV ranking optimization minimizes the evaluation error in the ranking comparison.
- We proposed two techniques to stabilize the evaluation by DIRV and demonstrated that these techniques were empirically effective.
- We extensively evaluated DIRV using both simulation and real service settings. The results demonstrated its high accuracy and efficiency.
- We published the real service data used in our experiments. This data can be used to validate our study and be used in future research on user modeling and bandit algorithms.

The remainder of this paper is organized as follows: Section 2 provides an overview of the related work. Section 3 describes the problem setting in this study. Section 4 explains the proposed method. Sections 5 and 6 report the experimental settings and discuss the experimental results, respectively. Finally, Section 7 presents the conclusions and future work on this topic.

2 RELATED WORK

Online evaluations are the basis of data-driven decision-making [10]. A typical online evaluation method is A/B testing [17], while there are also alternatives such as interleaving and bandit algorithms.

Interleaving (or *multileaving*, which refers to interleaving of more than two rankings) is a method used to increase the efficiency of online evaluation. Interleaving was reported to be 10 to 100 times more efficient than A/B testing [5, 24]. Schuth et al. proposed two multileaving methods called team draft multileaving (TDM) and optimized multileaving (OM) [28]. Recently, sample-scored-only multileaving (SOSM) [3] and pairwise preference multileaving (PPM) [19] were developed as more scalable multileaving methods. While PPM designed to evaluate click-based metrics is the state-of-the-art method in the interleaving, it is not trivial how to modify PPM to evaluate post-click metrics. Specifically, Schuth et al. found that a simple extension of the TDM resulted in a low accuracy in the non-click-based metrics (e.g., time to click metrics) [27]. Our approach differs from previous interleaving methods because we formulated a new *credit function*, which is used to aggregate

user feedback in the interleaving, as the expectation of post-click metrics for accurate evaluation.

Variance reduction is a technique commonly used to improve efficiency in online evaluations. [21] used a boosted decision tree regression to reduce the variance by matching similar users. The efficiency was improved in [10] by utilizing pre-experiment data through variance reduction. [32] implemented stratification for variance reduction. These studies reduced the variance for each tested group in a bucket-level evaluation like A/B testing. In contrast to bucket-level variance reduction, we reduce the item-level variances in post-click values based on the variance of the post-click values differing greatly for each item. By reducing item-level variances, we were able to improve the evaluation efficiency. Oosterhuis and de Rijke [20] applied variance reduction to an efficient evaluation of CTR. However, it was unclear how to apply this method to a post-click metric, which we focus on in this work. The proposed method in [20] relies heavily on the examination probability and has only been validated by simulation-based experiments. We found that the evaluation accuracy of the post-click metrics was highly compromised when there was an estimation error in the examination probability of the real service data. We developed a stabilization technique to address this problem. The stabilization technique reduces the systematic error caused by the estimation error of the examination probability.

Recently, several off-policy evaluation methods have been developed [15, 16, 31]. The goal of off-policy evaluation is to estimate the performance of a policy from the data generated by another policy(ies) in reinforcement learning [11]. Specifically, Saito [25] proposed an unbiased estimator for post-click metrics. Our method is different in that we interleave rankings dynamically so that the estimation error is minimized by the variance reduction in an on-policy manner.

The multi-armed bandit problem is a problem in which a fixed limited set of resources must be allocated between competing choices in a way that maximizes the expected gain [2, 4]. Algorithms used for bandit problems are similar to the proposed method of evaluating multiple rankings in this study. However, the objective of a bandit algorithm is different. The objective of a bandit algorithm is to maximize the total gain during a specific period or to identify the best arm. The objective of our study is to evaluate the superiority or inferiority of each pair of rankings. The evaluation result for each pair cannot be obtained from a bandit algorithm. However, the evaluation result for each pair may help decision-makers choose the best ranking in terms of effectiveness and the cost of producing the rankings.

3 PROBLEM SETTING

The primary goal of this study is to compare multiple *input rankings* $R = \{r_1, r_2, \dots\}$ in terms of a specific post-click metric in an online experiment. A user is shown a ranking in response to her/his query, clicks the items in the ranking, and consumes the clicked items. The post-click metric of an item can only be measured after the item has been clicked by the user. To compare rankings, we define the post-click metric of a ranking as the expectation of the

post-click metric of items in ranking r_i :

$$E[x|r_i] = \int_x xP(x|r_i)dx \quad (1)$$

where x is a random variable representing a value of the post-click metric, and $P(x|r_i)$ is the probability of observing x when a user interacts with the ranking r_i . $E[x|r_i]$ can be interpreted as how effective the ranking is to trigger users' behaviors (e.g., conversion) quantified by the post-click metric. Thus, $E[x|r_i]$ can be used to evaluate rankings based on users' post-click behaviors.

Following [3, 19, 26], we set the goal of our study to efficiently estimate the pairwise preference between rankings. The pairwise preference between rankings is defined as a matrix $\mathbf{P} \in \mathbb{R}^{|R| \times |R|}$, where $P_{i,j}$ indicates the difference between the expected post-click metric for ranking r_i and $r_j \in R$, that is,

$$P_{i,j} = E[x|r_i] - E[x|r_j].$$

A positive value of $P_{i,j}$ indicates that the ranking r_i is superior to r_j . The binary error E_{bin} is commonly used error metric to evaluate \mathbf{P} [3, 19, 26]. Letting $\bar{P}_{i,j}$ be the preference estimated by a certain method, the binary error of this estimate is defined as follows:

$$E_{\text{bin}} = \frac{1}{|R|(|R|-1)} \sum_{r_i, r_j \in R} \text{sgn}(P_{i,j}) \neq \text{sgn}(\bar{P}_{i,j}), \quad (2)$$

where $\text{sgn}(\cdot)$ returns -1 for negative values, 1 for positive values, and 0 otherwise. The operator \neq returns 1 whenever the signs are unequal.

A straightforward approach for estimating $P_{i,j}$ (or obtaining $\bar{P}_{i,j}$) is to conduct A/B testing with the rankings r_i and r_j , take the mean of x for each of the rankings, and then approximate $E[x|r_i]$ and $E[x|r_j]$ using these means. We discuss possible improvements over A/B testing and introduce our proposed method in the next section.

4 METHODOLOGY

We propose a method named *Decomposition and Interleaving for Reducing the Variance of post-click metrics (DIRV)* for efficient evaluation based on post-click metrics. We first discuss possible improvements over A/B testing by deeply look at the expectation of a post-click metric and introduce DIRV as our solution. We then show that the optimization by DIRV leads to the reduction of the binary error defined in Equation (2). Finally, we propose two techniques to stabilize evaluations by DIRV in actual applications.

4.1 Decomposition

We deeply look at Equation (1) by decomposing it with some assumptions and discuss possible improvements to reduce the evaluation error.

In Equation (1), the value of a post-click metric x can be observed for each clicked item in ranking r_i . Thus, $P(x|r_i)$ can be obtained by marginalizing all of the items in r_i , with the assumption that x depends solely on d , but not on the ranking r_i :

$$P(x|r_i) = \sum_{d \in r_i} P(x|d)P(c_d = 1|r_i) \quad (3)$$

where c_d is a random binary variable indicating an event of click on d , $P(c_d = 1|r_i)$ is the click probability of d in the ranking r_i , and $P(x|d)$ is the probability of observing x at d .

The independence assumption between x and r_i allows the following decomposition of the expectation of the post-click metric:

$$\begin{aligned} E[x|r_i] &= \int_x x \left\{ \sum_{d \in r_i} P(x|d)P(c_d = 1|r_i) \right\} dx \\ &= \sum_{d \in r_i} \int_x xP(x|d)P(c_d = 1|r_i)dx \\ &= \sum_{d \in r_i} P(c_d = 1|r_i) \int_x xP(x|d)dx = \sum_{d \in r_i} P(c_d = 1|r_i)E[x|d] \end{aligned} \quad (4)$$

where $E[x|d]$ is the expectation of the post-click metric for item d .

This equation immediately suggests that an unbiased estimate of $E[x|r_i]$ can be obtained by:

$$\bar{E}[x|r_i] = \sum_{d \in r_i} \bar{P}(c_d = 1|r_i)\bar{E}[x|d] \quad (5)$$

where $\bar{P}(c_d = 1|r_i)$ and $\bar{E}[x|d]$ are unbiased estimates for $P(c_d = 1|r_i)$ and $E[x|d]$, respectively. A standard approach to obtain these two estimates is to use the sample means of the random variables c_d and x for $\bar{P}(c_d = 1|r_i)$ and $\bar{E}[x|d]$, respectively. This decomposition can potentially achieve higher efficiency than A/B testing, because the samples for item d in *any* rankings in R can be used to estimate $E[x|d]$. Thus, we can increase the sample size for items that are shared by multiple rankings. Rankings are likely to include identical items, especially when differences between similar rankings (e.g., the same algorithm with different parameters) are evaluated.

While only Equation (5) enables better estimation of $E[x|d]$ by increasing the sample size, the use of the sample mean is not efficient enough especially when (1) the population variance of x is high and (2) the sample size of x is small due to small $P(c_d = 1|r_i)$. Both cases are likely to lead to high variance of $\bar{E}[x|d]$, resulting in high variance of $\bar{E}[x|r_i]$. Consequently, a large binary error in the ranking comparison is obtained. Both cases can be alleviated by prioritizing the exposure of items that satisfy these conditions for increasing the sample size. However, simple manipulation of a ranking prevents us from correctly estimating $P(c_d = 1|r_i)$ since this probability depends on the position of item d in ranking r_i .

Hence, we introduce a click model for further decomposing $\bar{P}(c_d = 1|r_i)\bar{E}[x|d]$ in the summation. We assume the examination hypothesis [7] that the item click can be decomposed into two variables: examination e_d and attraction a_d . This enables us to decompose the unbiased estimate $\bar{P}(c_d = 1|r_i)$ as follows:

$$\bar{P}(c_d = 1|r_i) = \bar{P}(e_d = 1|r_i)\bar{P}(a_d = 1|d), \quad (6)$$

where $\bar{P}(e_d = 1|r_i)$ can be estimated by a position-based click model or a cascade click model [7, 9]. The position-based click model assumes that the examination probability depends solely on the rank of the item in ranking r : i.e., $\bar{P}(e_d = 1|r_i) = g(\text{rank}(d, r))$, where g is a function taking a rank to return a probability, and $\text{rank}(d, r)$ is the rank of item d in ranking r . This assumption allows us to avoid estimating the probabilities specific to a particular ranking: i.e., $\bar{P}(e_d = 1|r_i)$. The estimations of $E[x|d]$, $P(e_d = 1|r_i)$ and $P(a_d = 1|r_i)$ are detailed in Section 5.3. Note that the simple assumption for the click model might sacrifice the evaluation error for efficiency. Thus,

in Section 4.4, we introduce a technique correcting the systematic error between the assumed and actual click models.

In summary, the decomposition of Equation (1) increases the sample size of x for estimating $E[x|d]$ and provides greater flexibility about the ranking presented to the users. Now one can freely produce and present a ranking containing items from input rankings R , and estimate each of $P(a_d = 1|d)$ and $E[x|d]$ based on the users' interactions with the presented ranking, in order to obtain $\bar{E}[x|r_i]$. In the next subsection, we explain what ranking should be presented to the users for minimizing the evaluation error.

4.2 Interleaving

Our proposed DIRV method is designed to minimize the variance of a post-click metric by interleaving input rankings to produce a single ranking to be presented to the users. The interleaved ranking preferentially exposes items with high population variance. The method attempts to receive more samples (or clicks) for these items to reduce the variance of the input rankings.

As shown in Appendix A, the variance $V[\bar{E}[x|r_i]]$ can be defined as a monotonically decreasing function ϕ on the number of impressions of the item d (denoted by n_d^i) and the number of clicks on the item d (denoted by n_d^c):

$$V[\bar{E}[x|r_i]] = \sum_{d \in r_i} \phi_{d,r_i}(n_d^i, n_d^c). \quad (7)$$

Our proposed DIRV method determines a ranking o produced by interleaving input rankings. The ranking o minimizes the summation of $V[\bar{E}[x|r_i]]$ over all the rankings in R when o is exposed to a user:

$$\min_o f(o)$$

where

$$f(o) = \sum_{r_i \in R} V[\bar{E}[x|r_i]] = \sum_{r_i \in R} \sum_{d \in r_i} \phi_{d,r_i}(n_d^i, n_d^c), \quad (8)$$

$$\dot{n}_d^i = n_d^i + 1, \quad \dot{n}_d^c = n_d^c + E[c_d|o]. \quad (9)$$

\dot{n}_d^i indicates the expected sample size after the interleaved ranking o is presented to a user. The number of impressions, n_d^i , is incremented by one each time the ranking o is presented. The number of clicks for item d , n_d^c , is incremented by $E[c_d|o] = P(c_d = 1|o)$, which indicates the expected number of clicks on item d in ranking o for a single impression.

As a brute-force search is infeasible to determine the optimal ranking o , we employ a greedy algorithm based on Equation (8). The greedy algorithm starts with an empty ranking r , and repeats appending an item to r that maximizes the difference between the current ranking and the ranking with the new item in terms of the variance of the sample mean:

$$\max_{d \in D \setminus r} \sum_{r_i \in R} (\phi_{d,r_i}(n_d^i, n_d^c) - \phi_{d,r_i}(\dot{n}_d^i, \dot{n}_d^c)), \quad (10)$$

where

$$\dot{n}_d^i = n_d^i + 1, \quad \dot{n}_d^c = n_d^c + E[c_d|r \oplus d]. \quad (11)$$

and D is a set of all items in R and $r \oplus d$ is the ranking r with d appended to the bottom. This greedy algorithm repeatedly finds the item to minimize the variance when it is appended. The algorithm stops when the ranking reaches a predefined depth.

4.3 Theoretical Justification

We provide theoretical justification that the minimization of the summation of $V[\bar{E}[x|r_i]]$ over R leads to the minimization of the expected binary error $E[E_{\text{bin}}]$ defined in Equation (2).

THEOREM 4.1.

$$\sum_{r_i \in R} V[\bar{E}[x|r_i]] \quad (12)$$

constitutes the upper bound of $E[E_{\text{bin}}]$.

PROOF. Let $\mu_i = E[x|r_i]$ and $\bar{\mu}_i = \bar{E}[x|r_i]$. Without a loss of generality, we assume that $\mu_i > \mu_j$. The binary error probability $P(\text{sgn}(P_{i,j}) \neq \text{sgn}(\bar{P}_{i,j}))$ can be interpreted as the probability of estimating a higher value for μ_j than that for μ_i , that is, $P(\text{sgn}(P_{i,j}) \neq \text{sgn}(\bar{P}_{i,j})) = P(\bar{\mu}_j > \bar{\mu}_i) = P(\bar{\mu}_j - \bar{\mu}_i > 0)$. Letting $\Delta_{ij} = \mu_j - \mu_i$ and $\bar{\Delta}_{ij} = \bar{\mu}_j - \bar{\mu}_i$, and using Chebyshev's inequality ($P(|x - \mu| > k) \leq \sigma^2/k^2$), this probability can be bounded as follows:

$$\begin{aligned} P(\text{sgn}(P_{i,j}) \neq \text{sgn}(\bar{P}_{i,j})) &= P(\bar{\Delta}_{ij} - \Delta_{ij} > -\Delta_{ij}) \\ &\leq P(|\bar{\Delta}_{ij} - \Delta_{ij}| > -\Delta_{ij}) \\ &\leq \frac{V[\bar{\Delta}_{ij}]}{\Delta_{ij}^2} = \frac{V[\bar{\mu}_i] + V[\bar{\mu}_j]}{\Delta_{ij}^2} \end{aligned} \quad (13)$$

Note that $P(|x| > y) = P(x > y) + P(x < -y) \geq P(x > y)$, and $V[x - y] = V[x] + V[y]$ if x and y are assumed independent.

Therefore, the expected binary error is bounded as follows:

$$\begin{aligned} E[E_{\text{bin}}] &= \frac{1}{|R|(|R|-1)} \sum_{r_i, r_j \in R} P(\text{sgn}(P_{i,j}) \neq \text{sgn}(\bar{P}_{i,j})) \\ &\leq \frac{1}{|R|(|R|-1)} \sum_{r_i, r_j \in R} \frac{V[\bar{\mu}_i] + V[\bar{\mu}_j]}{\Delta_{ij}^2} \\ &\leq \frac{1}{C|R|} \sum_{r_i \in R} V[\bar{\mu}_i] = \frac{1}{C|R|} \sum_{r_i \in R} V[\bar{E}[x|r_i]] \end{aligned} \quad (14)$$

where C is the smallest value for Δ_{ij}^2 ($r_i, r_j \in R$). Note that $V[\bar{\mu}_i]$ appears $|R|-1$ times in the summation, which cancels out $|R|-1$ in the denominator. \square

This theorem suggests that the upper bound of $E[E_{\text{bin}}]$ becomes lower if the variance of $\bar{E}[x|r_i]$ is reduced. This is the theoretical basis of our proposal and is formally expressed as follows:

$$E[E_{\text{bin}}] \rightarrow 0 \left(\sum_{r_i \in R} V[\bar{E}[x|r_i]] \rightarrow 0 \right). \quad (15)$$

4.4 Stabilization Techniques

There still remain some practical problems associated with the application of DIRV to an online experiment. These problems are discussed in this section.

4.4.1 Feature-based Variance Prediction. The first problem is the variance estimation of the post-click metric $V[x|d]$. This variance is included in the variance of $\bar{E}[x|d]$ (see Appendix A) and is required to be estimated to produce interleaved rankings. However, the estimation can be inaccurate, especially when the sample size is small. As demonstrated in our experiments, this problem can cause degradation of the estimation efficiency.

To overcome this problem, we propose using a regression model that predicts the variance of the post-click metric based on item features (e.g., titles and categories). As a result, we can use two estimates for the variance of the post-click metric $V[x|d]$, namely, the unbiased variance obtained from the observed values (denoted by $\bar{V}[x|d]$), and the predicted variance estimated by a regression model (denoted by $\hat{V}[x|d]$).

The larger value between $\bar{V}[x|d]$ and $\hat{V}[x|d]$ (i.e., $\max(\bar{V}[x|d], \hat{V}[x|d])$) is used as the estimate of $V[x|d]$. This idea is similar to the *clipping technique* used in [6, 13, 29]. As DIRV prioritizes items with a high variance in the interleaved ranking, there are opportunities to correct the estimation errors for items whose variance of the post-click are overestimated. In contrast, underestimation of the variance is more problematic because there are only limited opportunities to correct the error. Therefore, we use a higher variance to avoid underestimating the variance. Our experiments demonstrated that this technique is effective for efficient evaluation.

4.4.2 Systematic Error Correction. The second technique reduces the systematic error between the actual and assumed click models. Although the click probability $P(c_d = 1|r_i)$ is efficiently estimated with the click model, its estimation accuracy may not be sufficient as the click model does not precisely reflect real user behaviors.

To alleviate this problem, we estimate the click probability $P(c_d = 1|r_i)$ by the click model as well as a model-agnostic estimate (i.e., the sample mean of clicks). We present each ranking r_i with a small probability and observe the clicks on d in r_i . The clickthrough rate (i.e., $\hat{P}(c_d = 1|r_i) = n_{d,r_i}^c/n_{r_i}^i$) is used together with $\bar{P}(c_d = 1|r_i)$ in Equation (6), where n_{d,r_i}^c is the number of clicks on item d in the ranking r_i and $n_{r_i}^i$ is the number of impressions of the ranking r_i . The two estimates are linearly combined for approximating $P(c_d = 1|r_i)$:

$$P(c_d = 1|r_i) \approx \theta_i \bar{P}(c_d = 1|r_i) + (1 - \theta_i) \hat{P}(c_d = 1|r_i) \quad (16)$$

where we set the parameter θ_i to a value decreasing as the sample size increases (i.e., $1/(n_{r_i}^i + 1)^{0.5}$) to weight the click model when $n_{r_i}^i$ is small.

To present each ranking r_i with a small probability, we incorporate another objective function $g(o)$ into the policy for generating rankings. $g(o)$ is defined as follows:

$$g(o) = \sum_{r_i \in R} \theta_i \sum_{d \in r_i} \phi_{d,r_i}(\hat{n}_{r_i}^i, \hat{n}_{d,r_i}^c) \quad (17)$$

where

$$\hat{n}_{r_i}^i = n_{r_i}^i + \text{sgn}(r_i = o), \quad \hat{n}_{d,r_i}^c = n_{d,r_i}^c + \text{sgn}(r_i = o)P(c_d = 1|o) \quad (18)$$

and $\text{sgn}(r_i = o)$ returns 1 if all of the items are identical in r_i and o , otherwise 0. In the last procedure for generating rankings, we select the ranking o that minimizes $f(o) + \gamma g(o)$ from the set $\{o^*\} \cup R$, where o^* is a greedy solution in Equation (8), and γ is a hyperparameter. Although this technique is not justified theoretically, it was shown to be effective by our experiments.

5 EXPERIMENT SETUP

We describe the experimental settings to answer the following research questions (RQs):

- **RQ1:** Can DIRV identify preferences between rankings more accurately and efficiently than other methods?

Table 1: Experimental Setting Summary

	Simulation-based	Real Service
Input ranking	generated	raw data
Systematic error	-	exists
Click	partially generated	raw data
Post-click	partially generated	raw data

- **RQ2:** How does the variance prediction technique affect the evaluation efficiency?
- **RQ3:** How does the error correction technique affect the evaluation accuracy?

We have two experimental settings for comprehensive validation, namely, the simulation-based setting and the real service setting. A summary of the experimental settings is included in Table 1. The simulation-based setting is based on the traditional setting but is extended for the post-click evaluation. The real service setting is designed to validate methods close to the actual online evaluation by use of raw user feedback about the clicks and post-click values obtained from the raw ranking impressions. The real service dataset¹ are described in Section 5.2.

5.1 Simulation-based Settings

5.1.1 Evaluation Procedure. With the exception of post-click behavior, we simulated user behavior using the same four steps as in the previous interleaved evaluations [23, 28]:

- (1) **Impression:** A list of ranked items is displayed for a user.
- (2) **Click:** The user decides whether to click an item.
- (3) **Post-click Behavior:** If the user clicks the item, then the user consumes it, and post-click values are produced.
- (4) **Session End:** The user ends the session.

The details of the user’s behaviors are as follows:

First, we fix a query by uniformly sampling from a static dataset. Then, each method generates rankings from the query and displays them to the users. Whether a user clicks on an item depends on the item’s click probability.

After clicking on the item, the user is assumed to consume the item. Then, we can measure the value of a post-click metric (e.g., dwell time). The post-click value is assumed to follow a particular distribution. For example, the occurrence of a conversion in EC dataset follows a Bernoulli distribution. Throughout the experiment, we assumed a cascade click model [7], in which users view results from top to bottom and leave as soon as they see a worthwhile item. Each experiment involved a sequence of 10,000 simulated user impressions. All runs were repeated 30 times for each dataset and parameter setting.

5.1.2 Datasets. Three datasets were used in the simulation-based settings. The first dataset (called News) was generated from our news media service named Gunosy. The second dataset (called LETOR) is the learning to rank dataset. These two datasets were used to measure the dwell time that varies for each user. The third dataset (called EC) is an artificially generated dataset used to simulate e-commerce product purchases that has a predefined post-click value for the purchase price of the product.

¹It is public available at <https://github.com/koiizukag/DIRV>.

News Dataset. This dataset comes from a real-world, mobile news application. The News dataset contains the dwell time and the click probability. We generated this dataset using the implicit feedback from a single day. Each selected article had more than 20,000 raw post-click values (i.e., the dwell time for each user). This dataset contains 50 articles, which is sufficient to generate rankings. The ground truth for click probability and dwell time for each article was calculated from all of the user logs.

LETOR Dataset. We used eight publicly available LETOR datasets with varying sizes and representing different search tasks [19, 28]. Each dataset consisted of a set of queries and a set of items corresponding to each query. Feature representations and relevance labels were provided for each item and query pair. The datasets had three labels: unrelated (0), relevant (1), and highly relevant (2). The LETOR dataset was broken down by task. Most of the tasks were from the TREC Web Tracks from 2003 to 2008 [8, 22, 30]. These TREC Web Tracks were HP2003, HP2004, NP2003, NP2004, TD2003, and TD2004. Each TREC Web Track contained 50–150 queries. The OHSUMED dataset was based on a search engine’s query log of the MEDLINE abstract database and contained 106 queries. MQ2007 was based on the million query track [1] and consisted of 1,700 queries.

We generated click probability and post-click values because these datasets do not contain these values. We assumed that a user would read an article for a long time if the article was relevant to the user. To reflect this, we generated dwell time values using an exponential distribution with the rate parameters set to $\frac{1}{\lambda_d} \sim (\text{relevance} + 1) \cdot \text{uniform}(1, 20)$ for each item. We set the click probability to $P(a_d = 1|d) \sim \min((\text{relevance} + 1) \cdot \text{uniform}(0.0, 0.5), 1)$.

EC Dataset. We generated a dataset to simulate a product purchase in an e-commerce setting. This dataset can be easily reproduced based on the parameters described below. In this dataset, we randomly set the click probability to $P(a_d = 1|d) \sim \text{uniform}(0.0, 0.5)$, $\text{price}(d) \sim \text{uniform}(1, 1000)$, and $\text{conversion_rate}(d) \sim \text{uniform}(0, 0.5)$. The relationship between the conversion rate and the price is $E[x|d] = \text{conversion_rate}(d) \cdot \text{price}(d)$. We generated 50 items using these parameters.

5.1.3 Input Ranking. In the LETOR datasets, we generated *input rankings* by sorting items with features used in the existing interleaving experiments [19, 28]. We used the BM25, TF.IDF, TF, IDF, and LMIR.JM features for MQ2007. For the other datasets, we used the BM25, TF.IDF, LMIR.JM, and Hyperlink features. First, we randomly selected 20 items to obtain the same number of item candidates for each dataset. We then sorted ten items by each feature. As a result, we generated $|R|= 5$ rankings of length $|r_i|= 10$, which were similar to the settings reported in [28].

For the News and EC datasets, we generated input rankings by prioritizing items with positive feedback. We assumed that News and EC companies are likely to pay more attention to the user’s feedback in their ranking algorithm. The algorithm for generating input rankings first retrieved the top- k items ranked in descending order of $P(a_d = 1|d)E[x|d]$. This algorithm was used to approximate the level of user satisfaction. Next, we randomly retrieved items that had not been selected and appended them to the input

ranking. We continued this process until the ranking reached a certain length. Finally, we randomly shuffled the order of each input ranking randomly. We used different values for k in our experiments, which we called the *item duplication rate*. As a result, we generated $|R|= 5$ input rankings of length $|r_i|= 10$, which are similar values to the setting reported by [28].

5.2 Real Service Settings

Simulation-based experiments are easily conducted. However, their settings are limited in the validation of their assumptions. In the proposed method, estimating the expected post-click metric relies on the assumption that post-click behavior is independent of other behaviors by the user. However, this assumption may not always be true. For example, in the news service, the dwell time on a clicked article may depend on the dwell time of another clicked article when the two articles contain similar or redundant content. The simulation-based experiments were based on user simulation, which assumes the post-click behavior is independent. Moreover, the assumed click model may have been wrong or contained estimation errors. Therefore, the validity of these assumptions was not tested by the simulation-based experiments. We generated a real service dataset using the rankings presented to real users to validate the proposed method under conditions closer to actual user behaviors in an online evaluation.

5.2.1 Real Service Dataset. The evaluation of interleaving requires 1) *input rankings* to be used as the input, 2) generation of *interleaved rankings* from the *input rankings* for presentation to the users, and 3) *user behaviors* associated with the interleaved rankings. The interleaving method generates interleaved rankings from the input rankings using various policies for generating the rankings. To evaluate methods in this setting, we generated interleaved rankings by identifying *all possible combinations* of the items that the input rankings could generate and collected user behavior associated with the interleaved rankings from the service log. As they cover all possible interleaved rankings, one can obtain real statistics observed for any interleaved rankings without actually presenting them.

The data were collected in a news application. The news application was the same application used in the simulation-based News dataset. This dataset consisted of a set of queries, a set of input rankings and the interleaved rankings with user behaviors corresponding to each query. User behaviors data contained click or not for each article in the ranking and dwell-time for each article that the user clicked. The query was composed of the topic and day requested by the user. We used a social topic that was one of the most popular topics in this service. The ranking was personalized and differed for each user and query pair. The top three rankings with the highest number of impressions for each day were selected as the input ranking. The rankings were fixed at three in length to reduce the possible number of item combinations, since the number of interleaved rankings could be explosively large depending on the number of items in the input rankings. We generated 30 pairs of queries and their respective responses for a million-scale ranking impression in a one-month period.

Table 2: Evaluation accuracy for the News, EC, and LETOR datasets in the simulation-based setting. The binary error E_{bin} of all the comparison methods after 10,000 impressions on comparisons of $|R|=5$ rankings was averaged over 30 times per dataset and parameter. The best performances are noted in bold.

Duplication ratio (%)	News					EC				
	0	20	40	60	80	0	20	40	60	80
A/B Testing	0.100	0.190	0.095	0.140	0.160	0.015	0.040	0.070	0.075	0.185
TDM	0.425	0.465	0.470	0.500	0.615	0.055	0.170	0.200	0.215	0.300
DIRV w/o Var Pred	0.150	0.280	0.115	0.100	0.110	0.350	0.280	0.275	0.320	0.335
DIRV w/o Err Corr	0.075	0.095	0.070	0.035	0.070	0.035	0.045	0.055	0.030	0.050
DIRV	0.095	0.165	0.080	0.045	0.075	0.015	0.035	0.045	0.055	0.050

Duplication ratio (%)	LETOR							
	HP2003	HP2004	TD2003	TD2004	NP2003	NP2004	MQ2007	OHSUMED
A/B Testing	0.085	0.105	0.070	0.110	0.095	0.080	0.100	0.085
TDM	0.360	0.255	0.295	0.345	0.285	0.330	0.355	0.295
DIRV w/o Var Pred	0.070	0.055	0.060	0.060	0.075	0.055	0.065	0.035
DIRV w/o Err Corr	0.020	0.045	0.035	0.025	0.030	0.030	0.015	0.040
DIRV	0.030	0.015	0.035	0.040	0.045	0.035	0.055	0.020

The real service setting has advantages and disadvantages over simulation-based settings. The real service dataset was limited by the number of items in the input rankings. In addition, interleaved rankings did not cover all combinations of the items in the input rankings. On the other hand, this real service setting did allow us to test the assumptions of the model close to the actual online evaluation.

5.2.2 Experimental Runs. The ground truth of the preference in E_{bin} was calculated from half of the total data included for each query using A/B testing logic. The method was validated using the other half of the data. The experimental runs consisted of 5,000 ranking impressions for each of the 30 queries. These procedures were iterated 30 times.

5.3 Parameter Estimation

Each parameter was estimated as follows: $E[x|d]$ was estimated by the sample mean of the observed post-click value x for item d over all rankings. $P(e_{d_j} = 1|r_i)$ that was the examination probability of j -th item in the ranking r_i was estimated assuming the cascade click model as $\bar{P}(e_{d_j} = 1|r_i) = \prod_{k=1}^{j-1} (1 - \bar{P}(c_{d_k} = 1|d_k))$ where $\bar{P}(c_{d_k} = 1|d_k) = \bar{P}(e_{d_k} = 1|d_k)\bar{P}(a_{d_k} = 1|d_k)$. $P(a_d = 1|d)$ was estimated by $\bar{P}(a_d = 1|d) = n_d^c/n_d^e$ where n_d^c is the number of clicks and n_d^e is the number of presentations of item d over all of the rankings. We note that n_d^e is incremented for the items between the 1st position and the last position of the click in the ranking for each ranking impression. If there were no clicks in the ranking, n_d^e for each item was increased.

5.4 Variance Prediction

We used the News and real service datasets for studying the variance prediction because the EC and LETOR datasets had no ground-truth for the variance. We generated a dataset that included each article’s title length and category, as well as the media source and the sample variance for ground-truth that was calculated from the

Table 3: Features and importance

Feature	Importance
Category ID to which the article belongs	879
Supplier ID of the article	2,342
Content length of the article	1,854
Title length of the article	1,045

dwelling time of all users. A tree-based training model with a gradient-boosting framework² was used for prediction. We used the features detailed in Table 3 for prediction. The training epoch was set to 1,000. The root mean square error was used for the loss function. The early stopping parameter was 10. The other parameters were set to default values.

Figure 4 shows the plots of the predicted variance and the actual variance. A Pearson correlation value of 0.76 was achieved by only using meta-features. Feature importance is shown in Table 3.

In the EC and LETOR datasets, we artificially generated the predicted variance. Specifically, uniformly randomized noise was added to the population variance $V[x|d]$ to generate a predicted $\hat{V}[x|d]$ that satisfies $\hat{V}[x|d] \leq 2V[x|d]$.

5.5 Comparison Methods

We used five methods for the evaluation: A/B Testing, modified TDM, DIRV w/o Var Pred, DIRV w/o Err Corr, and DIRV. PPM is designed to evaluate click-based metrics and considered as the state-of-the-art method for interleaving [19]. However, it is not trivial to modify PPM to evaluate post-click metrics. As previously discussed, the modified TDM [27] is the only method designed to evaluate non-click-based metrics using interleaving. Based on [27], we modified the credit function to be the product of the post-click value and the original credit function of the TDM. DIRV w/o Var Pred is a DIRV method without the variance prediction, while DIRV w/o Err Corr is a DIRV method without the error correction. DIRV is a method that has both stabilization techniques. The hyper-parameter γ was set to 1.0.

²<https://github.com/microsoft/LightGBM>

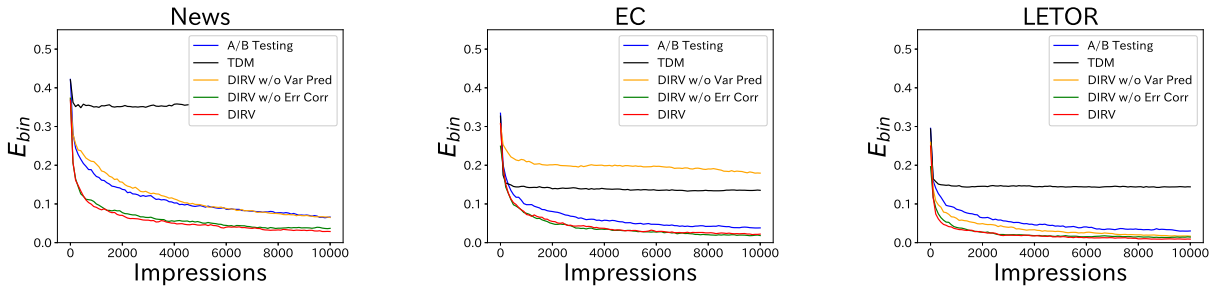


Figure 1: E_{bin} averaged over the number of impressions of the News, EC, and LETOR datasets. For all the datasets, DIRV or DIRV w/o Err Corr had the lowest E_{bin} for each impression.

TDM could not be used in the real service setting because there are cases where the rankings generated by the TDM might not exist in the real service data, as discussed in 5.2. Therefore, TDM was only used for the simulation-based experiments. In the DIRV methods for the real service setting, we generated rankings that minimize variance using $f(o)$ and $g(o)$ from the candidates of the interleaved rankings.

6 RESULTS AND DISCUSSION

In this section, we discuss the experimental results for answering the research questions posed in section 5.

6.1 RQ1: Can DIRV identify preferences between rankings more efficiently and accurately than comparison methods?

6.1.1 Efficiency. Figure 1 shows the efficiency results for the simulation-based setting for the News, LETOR, and EC datasets. Figure 2 shows the efficiency for the real service setting. In Figures 1 and 2, the x-axis represents the number of impressions, and the y-axis represents E_{bin} . For all of the datasets, DIRV had the lowest E_{bin} for each impression. In contrast to the simulation-based setting, DIRV w/o Err Corr had higher E_{bin} after 1,000 impressions compared to the results from the A/B testing in the real service setting. TDM stopped decreasing E_{bin} during the early stage in all datasets.

Among the evaluated methods, DIRV achieved the highest efficiency by reducing the variances. DIRV was designed to reduce variance by (1) aggregating post-clicks from different rankings that leads to an increased sample size for each item and (2) exposing the item that has high variance. DIRV successfully reduced the variance, which led to high efficiency (as shown in Figure 3).

6.1.2 Accuracy. Table 2 details the evaluation accuracy results. The number closest to 0.0 for each parameter and dataset are highlighted in bold. In the LETOR dataset, there were eight E_{bin} results for each dataset. In the News and EC datasets, there were five E_{bin} results based on the item duplication ratios (which ranged from 0% to 80%, in 20% increments). Figure 2 shows the accuracy of the real service setting in the last impression (i.e., at 5,000 impressions).

Overall, DIRV outperformed the existing methods for all of the datasets. A/B testing resulted in the lowest E_{bin} in the EC dataset with duplication ratios of 0%. We designed the estimator using

the expectation for post-click metrics. DIRV achieved high accuracy compared with TDM due to the use of our estimator. It is remarkable that our method performed well when many items were shared among the input rankings (i.e., high duplication ratio). Similar to the results from [27], our results demonstrated that TDM is difficult to extend for aggregating continuous values like post-click values.

Regarding RQ1, DIRV outperformed the existing methods in efficiency and accuracy. The performance was especially remarkable when many items were shared among the input rankings.

6.2 RQ2: How does the variance prediction technique affect the evaluation efficiency?

Figure 3 details the variance reduction on the News dataset in the simulation setting. The results show that DIRV achieved the lowest variance among the evaluated methods. DIRV w/o Var Pred decreased the variance slowly compared to DIRV, especially for a small number of impressions. TDM had the highest variance for each impression. These trends were also observed in the other datasets.

Figure 3 shows that the predicted variance contributed to reducing the variance in small impressions. In contrast, the variance of some items was underestimated in DIRV w/o Var Pred and could not obtain a sufficient number of exposures during the early stage. This resulted in reducing the variance in the post-click metrics slowly. The variance prediction technique reduced the variance and improved efficiency. In TDM, some items at the bottom of the input ranking was not selected for the interleaved ranking, which led to high variance.

Figure 2 shows that the variance prediction did not contribute to the efficiency in the real service setting. This is because the total number of items was small in this setting. Thus, there were fewer benefits from manipulating the order of the items using the predicted variance.

Regarding RQ2, the variance prediction techniques contributed to reducing the variance, which led to improved efficiency.

6.3 RQ3: How does the error correction technique affect the evaluation accuracy?

Figure 2 shows the accuracy of the real service setting at the end of the evaluation (i.e., at 5,000 impressions). The results show that the accuracy was close to 0.0 for DIRV and A/B testing methods.

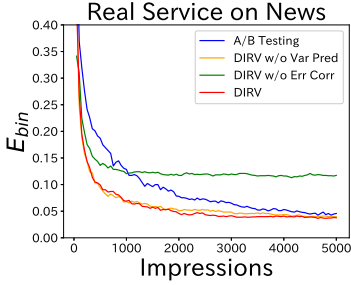


Figure 2: E_{bin} averaged over the number of impressions of the real service News dataset. DIRV or DIRV w/o Var Pred had the lowest E_{bin} for each impression.

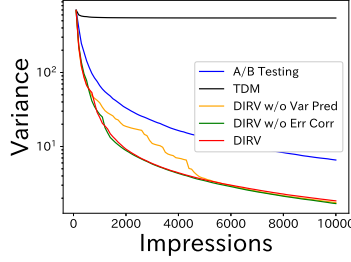


Figure 3: Variance versus the number of impressions in the simulation-based News dataset. The result shows that DIRV had lower variance for each impression, especially for a small number of impressions.

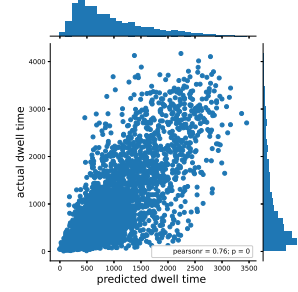


Figure 4: Visualization of the predicted variance. Each dot corresponds to the predicted and actual variance values for each News article. The Pearson’s correlation value was 0.76.

In contrast, the DIRV w/o Err Corr method stopped decreasing the E_{bin} after 1,000 impressions.

The results show that the error was reduced as the number of impressions increased in the real service setting by introducing an error correction technique. In particular, we successfully confirmed that the evaluation error converged around at 0.0 using the error correction technique, which reduced the systematic error as the experiment progressed. This result supported the assumption that the post-click behavior was independent of other user behavior in this real service setting.

The difference between the assumed cascade click model and the actual user behavior in the real service setting resulted in a systematic error. Figure 2 illustrated the error in the results of the DIRV w/o Err Corr method. In the simulation-based setting, we assumed that the actual user behavior also obeyed the cascade click model. This assumption led to the same accuracy between the DIRV and DIRV w/o Err Corr methods in Table 2.

Regarding RQ3, we reduced the systematic error using the error correction technique. The error correction technique led to improved accuracy when an estimation error existed in the click model.

7 CONCLUSION AND FUTURE WORK

In this study, we proposed a method for accurate and efficient post-click evaluation using interleaving. First, we introduced the click model to aggregate the post-click values. Then, we showed that minimizing the variance of post-click metrics leads to a reduction in the evaluation error. Next, we provided a policy to generate rankings to reduce the variance. Finally, we proposed two stabilization techniques to support the proposed method.

To evaluate this method, we conducted comprehensive experiments with both simulation-based and real service settings. The experimental results from this study indicated that 1) the proposed method outperformed existing methods in efficiency and accuracy and the performance was especially remarkable when many items were shared among the input rankings, 2) the variance prediction

techniques contributed to reducing the variance, which led to improved efficiency, and 3) we could successfully reduce the systematic error using the error correction technique, which led to improved accuracy when there existed an estimation error in the click model.

In this study, we focused on building an online evaluating method. In the future, we plan to extend our framework to include a bandit algorithm in ranking optimization.

APPENDIX

A VARIANCE OF $\bar{E}[x|r_i]$

Noting that $V[x + y] = V[x] + V[y]$ if x and y are independent, we obtain:

$$V[\bar{E}[x|r_i]] = V\left[\sum_{d \in r_i} \bar{P}(c_d = 1|r_i) \bar{E}[x|d]\right] = \sum_{d \in r_i} V\left[\bar{P}(c_d = 1|r_i) \bar{E}[x|d]\right]$$

Using $V[xy] = V[x]V[y] + E[x]^2V[y] + V[x]E[y]^2$ (x and y are independent), the variance for each item is obtained as follows:

$$\begin{aligned} V\left[\bar{P}(c_d = 1|r_i) \bar{E}[x|d]\right] &= V[\bar{P}(c_d = 1|r_i)]V[\bar{E}[x|d]] \\ &+ E[\bar{P}(c_d = 1|r_i)]^2V[\bar{E}[x|d]] \\ &+ E[\bar{E}[x|d]]^2V[\bar{P}(c_d = 1|r_i)] \end{aligned}$$

Since $V[\bar{x}] = \sigma^2/n$ (\bar{x} is the sample mean, σ^2 is the population variance, and n is the sample size), $V[\bar{P}(c_d = 1|r_i)]$ can be computed as $V[\bar{P}(c_d = 1|r_i)] = V[P(c_d = 1|d)]/n_d^i$, where n_d^i is the number of impressions of item d . We can also obtain $V[\bar{E}[x|d]] = V[x|d]/n_d^c$ where n_d^c is the number of clicks on item d .

Finally, we express the variance of $\bar{E}[x|r_i]$ as a function of the samples sizes n_d^i and n_d^c :

$$\begin{aligned} V[\bar{E}[x|r_i]] &= \sum_{d \in r_i} \left\{ \frac{V[P(c_d = 1|d)]}{n_d^i} \frac{V[x|d]}{n_d^c} + E[\bar{P}(c_d = 1|r_i)]^2 \frac{V[x|d]}{n_d^c} \right. \\ &\left. + E[\bar{E}[x|d]]^2 \frac{V[P(c_d = 1|d)]}{n_d^i} \right\} := \sum_{d \in r_i} \phi_{d,r_i}(n_d^i, n_d^c) \end{aligned}$$

The function ϕ_{d,r_i} monotonically decreases for n_d^i and n_d^c when x is always positive. As $\bar{P}(c_d = 1|r_i)$ and $\bar{E}[x|d]$ are unbiased estimators, $E[\bar{P}(c_d = 1|r_i)]$ and $E[\bar{E}[x|d]]$ can be approximated by

$\bar{P}(c_d = 1|r_i)$ and $\bar{E}[x|d]$ with a sufficient number of samples. Assuming that $P(a_d|d)$ follows a Bernoulli distribution, we can also approximate $V[P(c_d = 1|d)]$ by $\bar{P}(c_d = 1|r_i)(1 - \bar{P}(c_d = 1|r_i))$.

REFERENCES

- [1] James Allan, Ben Carterette, Javed A. Aslam, Virgil Pavlu, Blagovest Dachev, and Evangelos Kanoulas. 2007. *Million query track 2007 overview*. Technical Report. University of Massachusetts Amherst.
- [2] Djallel Bouneffouf and Irina Rish. 2019. A survey on practical applications of multi-armed and contextual bandits. *arXiv preprint arXiv:1904.10040* (2019).
- [3] Brian Brost, Ingemar J Cox, Yevgeny Seldin, and Christina Lioma. 2016. An improved multileaving algorithm for online ranker evaluation. In *Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 745–748.
- [4] Giuseppe Burtini, Jason Loeppky, and Ramon Lawrence. 2015. A survey of online experiment design with the stochastic multi-armed bandit. *arXiv preprint arXiv:1510.00757* (2015).
- [5] Olivier Chapelle, Thorsten Joachims, Filip Radlinski, and Yisong Yue. 2012. Large-scale validation and analysis of interleaved search evaluation. *ACM Trans. Inf. Syst.* 30, 1 (2012), 1–41.
- [6] Minmin Chen, Alex Beutel, Paul Covington, Sagar Jain, Francois Belletti, and Ed H Chi. 2019. Top-k off-policy correction for a REINFORCE recommender system. In *Proceedings of the 12th ACM International Conference on Web Search and Data Mining*. 456–464.
- [7] Aleksandr Chuklin, Ilya Markov, and Maarten De Rijke. 2015. *Click models for web search*. Morgan & Claypool Publishers.
- [8] Charles L. Clarke, Nick Craswell, and Ian Soboroff. 2009. *Overview of the TREC 2009 web track*. Technical Report. University of Waterloo.
- [9] Nick Craswell, Onno Zoeter, Michael Taylor, and Bill Ramsey. 2008. An experimental comparison of click position-bias models. In *Proceedings of the 1st ACM International Conference on Web Search and Data Mining*. 87–94.
- [10] Alex Deng, Ya Xu, Ron Kohavi, and Toby Walker. 2013. Improving the sensitivity of online controlled experiments by utilizing pre-experiment data. In *Proceedings of the 6th ACM International Conference on Web Search and Data Mining*. 123–132.
- [11] Mehrdad Farajtabar, Yinlam Chow, and Mohammad Ghavamzadeh. 2018. More robust doubly robust off-policy evaluation. In *International Conference on Machine Learning*. PMLR, 1447–1456.
- [12] Jean Garcia-Gathright, Brian St. Thomas, Christine Hosey, Zahra Nazari, and Fernando Diaz. 2018. Understanding and evaluating user satisfaction with music discovery. In *The 41st International ACM SIGIR Conference on Research and Development in Information Retrieval*. 55–64.
- [13] Alexandre Gilotte, Clément Calauzènes, Thomas Nedelec, Alexandre Abraham, and Simon Dollé. 2018. Offline A/B testing for recommender systems. In *Proceedings of the 11th ACM International Conference on Web Search and Data Mining*. 198–206.
- [14] Mihajlo Grbovic and Haibin Cheng. 2018. Real-time personalization using embeddings for search ranking at Airbnb. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 311–320.
- [15] Alois Gruson, Praveen Chandar, Christophe Charbuillet, James McInerney, Samantha Hansen, Damien Tardieu, and Ben Carterette. 2019. Offline evaluation to make decisions about playlist recommendation algorithms. In *Proceedings of the 12th ACM International Conference on Web Search and Data Mining*. 420–428.
- [16] Thorsten Joachims, Adith Swaminathan, and Tobias Schnabel. 2017. Unbiased learning-to-rank with biased feedback. In *Proceedings of the 10th ACM International Conference on Web Search and Data Mining*. 781–789.
- [17] Ron Kohavi, Alex Deng, Brian Frasca, Toby Walker, Ya Xu, and Nils Pohlmann. 2013. Online controlled experiments at large scale. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1168–1176.
- [18] Shumpei Okura, Yukihiko Tagami, Shingo Ono, and Akira Tajima. 2017. Embedding-based news recommendation for millions of users. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1933–1942.
- [19] Harrie Oosterhuis and Maarten de Rijke. 2017. Sensitive and scalable online evaluation with theoretical guarantees. In *Proceedings of the 26th ACM International Conference on Information and Knowledge Management*. 77–86.
- [20] Harrie Oosterhuis and Maarten de Rijke. 2020. Taking the Counterfactual Online: Efficient and Unbiased Online Evaluation for Ranking. In *Proceedings of the 2020 ACM SIGIR on International Conference on Theory of Information Retrieval*. 137–144.
- [21] Alexey Poyarkov, Alexey Drutsa, Andrey Khalyavin, Gleb Gusev, and Pavel Serdyukov. 2016. Boosted decision tree regression adjustment for variance reduction in online controlled experiments. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 235–244.
- [22] Tao Qin, Tie-Yan Liu, Jun Xu, and Hang Li. 2016. LETOR: A benchmark collection for research on learning to rank for information retrieval.
- [23] Filip Radlinski and Nick Craswell. 2013. Optimized interleaving for online retrieval evaluation. In *Proceedings of the 6th ACM International Conference on Web Search and Data Mining*. 245–254.
- [24] Filip Radlinski, Madhu Kurup, and Thorsten Joachims. 2008. How does click-through data reflect retrieval quality?. In *Proceedings of the 17th ACM International Conference on Information and Knowledge Management*. 43–52.
- [25] Yuta Saito. 2020. Doubly robust estimator for ranking metrics with post-click conversions. In *Proceedings of the 14th International ACM Conference on Recommender Systems*. 92–100.
- [26] Anne Schuth, Robert-Jan Bruinjtjes, Fritjof Buüttner, Joost van Doorn, Carla Groenland, Harrie Oosterhuis, Cong-Nguyen Tran, Bas Veeling, Jos van der Velde, Roger Wechsler, et al. 2015. Probabilistic multileave for online retrieval evaluation. In *Proceedings of the 38th international ACM SIGIR Conference on Research and Development in Information Retrieval*. 955–958.
- [27] Anne Schuth, Katja Hofmann, and Filip Radlinski. 2015. Predicting search satisfaction metrics with interleaved comparisons. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 463–472.
- [28] Anne Schuth, Floor Sietsma, Shimon Whiteson, Damien Lefortier, and Maarten de Rijke. 2014. Multileaved comparisons for fast online evaluation. In *Proceedings of the 23rd ACM International Conference on Information and Knowledge Management*. 71–80.
- [29] Yi Su, Lequn Wang, Michele Santacatterina, and Thorsten Joachims. 2019. Cab: Continuous adaptive blending for policy evaluation and learning. In *International Conference on Machine Learning*. 6005–6014.
- [30] Ellen M Voorhees and Donna Harman. 2003. Overview of TREC 2003. In *Trec*. 1–13.
- [31] Xuanhui Wang, Nadav Golbandi, Michael Bendersky, Donald Metzler, and Marc Najork. 2018. Position bias estimation for unbiased learning to rank in personal search. In *Proceedings of the 11th ACM International Conference on Web Search and Data Mining*. 610–618.
- [32] Huizhi Xie and Juliette Aurisset. 2016. Improving the sensitivity of online controlled experiments: Case studies at Netflix. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 645–654.
- [33] Xing Yi, Liangjie Hong, Erheng Zhong, Nanthan Nan Liu, and Suju Rajan. 2014. Beyond clicks: Dwell time for personalization. In *Proceedings of the 8th International ACM Conference on Recommender Systems*. 113–120.

This figure "sample-franklin.png" is available in "png" format from:

<http://arxiv.org/ps/2306.10024v1>