

# A General Framework for Counterfactual Learning-to-Rank

Aman Agarwal  
Cornell University  
Ithaca, NY  
aman@cs.cornell.edu

Kenta Takatsu  
Cornell University  
Ithaca, NY  
kt426@cornell.edu

Ivan Zaitsev  
Cornell University  
Ithaca, NY  
iz44@cornell.edu

Thorsten Joachims  
Cornell University  
Ithaca, NY  
tj@cs.cornell.edu

## ABSTRACT

Implicit feedback (e.g., click, dwell time) is an attractive source of training data for Learning-to-Rank, but its naive use leads to learning results that are distorted by presentation bias. For the special case of optimizing average rank for linear ranking functions, however, the recently developed SVM-PropRank method has shown that counterfactual inference techniques can be used to provably overcome the distorting effect of presentation bias. Going beyond this special case, this paper provides a general and theoretically rigorous framework for counterfactual learning-to-rank that enables unbiased training for a broad class of additive ranking metrics (e.g., Discounted Cumulative Gain (DCG)) as well as a broad class of models (e.g., deep networks). Specifically, we derive a relaxation for propensity-weighted rank-based metrics which is subdifferentiable and thus suitable for gradient-based optimization. We demonstrate the effectiveness of this general approach by instantiating two new learning methods. One is a new type of unbiased SVM that optimizes DCG – called SVM PropDCG –, and we show how the resulting optimization problem can be solved via the Convex Concave Procedure (CCP). The other is Deep PropDCG, where the ranking function can be an arbitrary deep network. In addition to the theoretical support, we empirically find that SVM PropDCG significantly outperforms existing linear rankers in terms of DCG. Moreover, the ability to train non-linear ranking functions via Deep PropDCG further improves performance.

## KEYWORDS

Learning to rank, presentation bias, counterfactual inference

### ACM Reference Format:

Aman Agarwal, Kenta Takatsu, Ivan Zaitsev, and Thorsten Joachims. 2019. A General Framework for Counterfactual Learning-to-Rank. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '19)*, July 21–25, 2019, Paris, France. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3331184.3331202>

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](mailto:permissions@acm.org).

*SIGIR '19*, July 21–25, 2019, Paris, France

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6172-9/19/07...\$15.00

<https://doi.org/10.1145/3331184.3331202>

## 1 INTRODUCTION

Implicit feedback from user behavior is an attractive source of data in many information retrieval (IR) systems, especially ranking applications where collecting relevance annotations from experts can be economically infeasible or even impossible (e.g., personal collection search, intranet search, scholarly search). While implicit feedback is often abundant, cheap, timely, user-centric, and routinely logged, it suffers from inherent biases. For example, the position of a result in a search ranking strongly affects how likely it is to be seen by a user and thus clicked. So, naively using click data as a relevance signal leads to sub-optimal performance.

A counterfactual inference approach for learning-to-rank (LTR) from logged implicit feedback was recently developed to deal with such biases [15]. This method provides a rigorous approach to unbiased learning despite biased data and overcomes the limitations of alternative bias-mitigation strategies. In particular, it does not require the same query to be seen multiple times as necessary for most generative click models, and it does not introduce alternate biases like treating clicks as preferences between clicked and skipped documents.

The key technique in counterfactual learning is to incorporate the propensity of obtaining a particular training example into an Empirical Risk Minimization (ERM) objective that is provably unbiased [28]. While it was shown that this is possible for learning to rank, existing theoretical support is limited to linear ranking functions and optimizing average rank of the relevant documents as objective [15]. In this paper, we generalize the counterfactual LTR framework to a broad class of additive IR metrics as well as non-linear deep models. Specifically, we show that any IR metric that is the sum of individual document relevances weighted by some function of document rank can be directly optimized via Propensity-Weighted ERM. Moreover, if an IR metric meets the mild requirement that the rank weighting function is monotone, we show that a hinge-loss upper-bounding technique enables learning of a broad class of differentiable models (e.g. deep networks).

To demonstrate the effectiveness of the general framework, we fully develop two learning-to-rank methods that optimize the Discounted Cumulative Gain (DCG) metric. The first is SVM PropDCG, which generalizes a Ranking SVM to directly optimize a bound on DCG from biased click data. The resulting optimization problem is no longer convex, and we show how to find a local optimum using the Convex Concave Procedure (CCP). In CCP, several iterations of convex sub-problems are solved. In the case of SVM PropDCG,

these convex sub-problems have the convenient property of being a Quadratic Program analogous to a generalized Ranking SVM. This allows the CCP to work by invoking an existing and fast SVM solver in each iteration until convergence. The second method we develop is Deep PropDCG, which further generalizes the approach to deep networks as non-linear ranking functions. Deep PropDCG also optimizes a bound on the DCG, and we show how the resulting optimization problem can be solved via stochastic gradient descent for any network architecture that shares neural network weights across candidate documents for the same query.

In addition to the theoretical derivation and the justification, we also empirically evaluate the effectiveness of both SVM PropDCG and Deep PropDCG, especially in comparison to the existing SVM PropRank method [15]. We find that SVM PropDCG performs significantly better than SVM PropRank in terms of DCG, and that it is robust to varying degrees of bias, noise and propensity-model misspecification. In our experiments, CCP convergence was typically achieved quickly within three to five iterations. For Deep PropDCG, the results show that DCG performance is further improved compared to SVM PropDCG when using a neural network, thus demonstrating that the counterfactual learning approach can effectively train non-linear ranking functions. SVM PropDCG and Deep PropDCG are also seen to outperform LambdaRank in terms of DCG.

## 2 RELATED WORK

Generative click models are a popular approach for explaining the bias in user behavior and for extracting relevance labels for learning. For example, in the cascade model [9] users are assumed to sequentially go down a ranking and click on a document, thus revealing preferences between clicked and skipped documents. Learning from these relative preferences lowers the impact of some biases [13]. Other click models ([3, 7, 9], also see [8]) train to maximize log-likelihood of observed clicks, where relevance is modeled as a latent variable that is inferred over multiple instances of the same query. In contrast, the counterfactual framework [15] does not require latent-variable inference and repeat queries, but allows directly incorporating click feedback into the learning-to-rank algorithm in a principled and unbiased way, thus allowing the direct optimization of ranking performance over the natural query distribution.

The counterfactual approach uses inverse propensity score (IPS) weighting, originally employed in survey sampling [11] and causal inference from observational studies [24], but more recently also in whole page optimization [34], IR evaluation with manual judgments [25], and recommender evaluation [17, 26]. This approach is similar in spirit to [32], where propensity-weighting is used to correct for selection bias when discarding queries without clicks during learning-to-rank.

Recently, inspired by the IPS correction approach for unbiased LTR, some algorithms (Ai et al. [2], Hu et al. [12]) have been proposed that jointly estimate the propensities and learn the ranking function. However, this requires an accurate relevance model to succeed, which is at least as hard as the LTR from biased feedback problem in question. Moreover, the two-step approach of propensity estimation followed by training an unbiased ranker allows direct optimization of any chosen target ranking metric independent of

Metric	$\lambda(\text{rank})$
<i>AvgRank</i>	rank
<i>DCG</i>	$-1/\log(1 + \text{rank})$
<i>Prec@k</i>	$-1_{\text{rank} \leq k}/k$
<i>RBP-p</i> [20]	$-(1-p)/p^{\text{rank}}$

**Table 1: Some popular linearly decomposable IR metrics that can be directly optimized by Propensity-Weighted ERM.  $\lambda(r)$  is the rank weighting function.**

the propensity estimation step.

While our focus is on directly optimizing ranking performance in the implicit feedback partial-information setting, several approaches have been proposed for the same task in the full-information supervised setting, i.e. when the relevances of all the documents in the training set are known. A common strategy is to use some smoothed version of the ranking metric for optimization, as seen in SoftRank [30] and others [6, 14, 35, 36]. In particular, SoftRank optimizes the expected performance metric over the distribution of rankings induced by smoothed scores, which come from a normal distribution centered at the query-document mean scores predicted by a neural net. This procedure is computationally expensive with an  $O(n^3)$  dependence on the number of documents for a query. In contrast, our approach employs an upper bound on the performance metric, whose structure makes it amenable to the Convex Concave Procedure for efficient optimization, as well as adaptable to non-linear ranking functions via deep networks.

Finally, several works exist [4, 5, 23, 30] that have proposed neural network architectures for learning-to-rank. We do not focus on a specific network architecture in this paper, but instead propose a new training criterion for learning-to-rank from implicit feedback that in principle allows unbiased network training for a large class of architectures.

## 3 UNBIASED ESTIMATION OF RANK-BASED IR METRICS

We begin by developing a counterfactual learning framework that covers the full class of linearly decomposable metrics as defined below (e.g. DCG). This extends [15] which was limited to the Average Rank metric. Suppose we are given a sample  $X$  of i.i.d. query instances  $\mathbf{x}_i \sim P(\mathbf{x})$ ,  $i \in [N]$ . A query instance can include personalized and contextual information about the user in addition to the query string. For each query instance  $\mathbf{x}_i$ , let  $r_i(y)$  denote the user-specific relevance of result  $y$  for instance  $\mathbf{x}_i$ . For simplicity, assume that relevances are binary,  $r_i(y) \in \{0, 1\}$ . In the following, we consider the class of additive ranking performance metrics, which includes any metric that can be expressed as

$$\Delta(\mathbf{y}|\mathbf{x}_i, r_i) = \sum_{y \in \mathbf{y}} \lambda(\text{rank}(y|\mathbf{y})) \cdot r_i(y). \quad (1)$$

$\mathbf{y}$  denotes a ranking of results, and  $\lambda(\cdot)$  can be any weighting function that depends on the rank  $\text{rank}(y|\mathbf{y})$  of document  $y$  in ranking  $\mathbf{y}$ . A broad range of commonly used ranking metrics falls into this class, and Table 1 lists some of them. For instance, setting  $\lambda(\text{rank}) = \text{rank}$  gives the sum of relevant ranks metric (also called average rank when normalized) considered in [15], and  $\lambda(\text{rank}) = \frac{-1}{\log(1+\text{rank})}$

gives the DCG metric. Note that we consider negative values whenever necessary to make the notation consistent with risk minimization.

A ranking system  $S$  maps a query instance  $\mathbf{x}_i$  to a ranking  $\mathbf{y}$ . Aggregating the losses of individual rankings over the query distribution, we can define the overall *risk* (e.g., the expected DCG) of a system as

$$R(S) = \int \Delta(S(\mathbf{x})|\mathbf{x}, r) dP(\mathbf{x}, r). \quad (2)$$

A key problem when working with implicit feedback data is that we cannot assume that all relevances  $r_i$  are observed. In particular, while a click (or a sufficiently long dwell time) provides a noisy indicator of positive relevance in the presented ranking  $\tilde{\mathbf{y}}_i$ , a missing click does not necessarily indicate lack of relevance as the user may not have observed that result. From a machine learning perspective, this implies that we are in a partial-information setting, which we will deal with by explicitly modeling missingness in addition to relevance. Let  $o_i \sim P(o|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i)$  denote the 0/1 vector indicating which relevance values are revealed. While  $o_i$  is not necessarily fully observed either, we can now model its distribution, which we will find below is sufficient for unbiased learning despite the missing data. In particular, the *propensity* of observing  $r_i(y)$  for query instance  $\mathbf{x}_i$  given presented ranking  $\tilde{\mathbf{y}}$  is then defined as  $Q(o_i(y) = 1|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i)$ .

Using this counterfactual setup, an unbiased estimate of  $\Delta(\mathbf{y}|\mathbf{x}_i, r_i)$  for any ranking  $\mathbf{y}$  can be obtained via IPS weighting

$$\hat{\Delta}_{IPS}(\mathbf{y}|\mathbf{x}_i, \tilde{\mathbf{y}}_i, o_i) = \sum_{\substack{y: o_i(y)=1 \\ \wedge r_i(y)=1}} \frac{\lambda(\text{rank}(\mathbf{y}|\mathbf{y}))}{Q(o_i(y) = 1|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i)}. \quad (3)$$

This is an unbiased estimator since,

$$\begin{aligned} & \mathbb{E}_{o_i}[\hat{\Delta}_{IPS}(\mathbf{y}|\mathbf{x}_i, \tilde{\mathbf{y}}_i, o_i)] \\ &= \mathbb{E}_{o_i} \left[ \sum_{y: o_i(y)=1} \frac{\lambda(\text{rank}(\mathbf{y}|\mathbf{y})) \cdot r_i(y)}{Q(o_i(y) = 1|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i)} \right] \\ &= \sum_{y \in \mathbf{y}} \mathbb{E}_{o_i} \left[ \frac{o_i(y) \cdot \lambda(\text{rank}(\mathbf{y}|\mathbf{y})) \cdot r_i(y)}{Q(o_i(y) = 1|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i)} \right] \\ &= \sum_{y \in \mathbf{y}} \frac{Q(o_i(y) = 1|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i) \cdot \lambda(\text{rank}(\mathbf{y}|\mathbf{y})) \cdot r_i(y)}{Q(o_i(y) = 1|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i)} \\ &= \sum_{y \in \mathbf{y}} \lambda(\text{rank}(\mathbf{y}|\mathbf{y})) r_i(y) \\ &= \Delta(\mathbf{y}|\mathbf{x}_i, r_i), \end{aligned} \quad (4)$$

assuming  $Q(o_i(y) = 1|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i) > 0$  for all  $\mathbf{y}$  that are relevant  $r_i(y) = 1$ . The above proof is a generalized version of the one in [15] for the Average Rank metric. Note that the estimator in Equation (3) sums only over the results where the feedback is observed (i.e.,  $o_i(y) = 1$ ) and positive (i.e.,  $r_i(y) = 1$ ), which means that we do not have to disambiguate whether lack of positive feedback (e.g., the lack of a click) is due to a lack of relevance or due to missing the observation (e.g., result not relevant vs. not viewed).

Using this unbiased estimate of the loss function, we get an unbiased estimate of the risk of a ranking system  $S$

$$\hat{R}_{IPS}(S) = \frac{1}{N} \sum_{i=1}^N \sum_{\substack{y: o_i(y)=1 \\ \wedge r_i(y)=1}} \frac{\lambda(\text{rank}(y|S(\mathbf{x}_i)))}{Q(o_i(y) = 1|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i)}. \quad (5)$$

Note that the propensities  $Q(o_i(y) = 1|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i)$  are generally unknown, and must be estimated based on some model of user behavior. Practical approaches to estimating the propensities are given in [1, 10, 15, 33].

## 4 UNBIASED EMPIRICAL RISK MINIMIZATION FOR LTR

The propensity-weighted empirical risk from Equation (5) can be used to perform Empirical Risk Minimization (ERM)

$$\hat{S} = \underset{S \in \mathcal{S}}{\text{argmin}} \{ \hat{R}_{IPS}(S) \}.$$

Under the standard uniform convergence conditions [31], the unbiasedness of the risk estimate implies consistency in the sense that given enough training data, the learning algorithm is guaranteed to find the best system in  $\mathcal{S}$ . We have thus obtained a theoretically justified training objective for learning-to-rank with additive metrics like DCG. However, it remains to be shown that this training objective can be implemented in efficient and practical learning methods. This section shows that this is indeed possible for a generalization of Ranking SVMs and for deep networks as ranking functions.

Consider a dataset of  $n$  examples of the following form. For each query-result pair  $(\mathbf{x}_i, y_i)$  that is clicked, let  $q_i = Q(o_i(y) = 1|\mathbf{x}_i, \tilde{\mathbf{y}}_i, r_i)$  be the propensity of the click according to a click propensity model such as the Position-Based Model [15, 33]. We also record the candidate set  $Y_i$  of all results for query  $\mathbf{x}_i$ . Note that each click generates a separate training example, even if multiple clicks occur for the same query.

Given this propensity-scored click data, we would like to learn a scoring function  $f(\mathbf{x}, y)$ . Such a scoring function  $f$  naturally specifies a ranking system  $S$  by sorting candidate results  $Y$  for a given query  $\mathbf{x}$  by their scores.

$$S_f(\mathbf{x}) \equiv \text{argsort}_Y \{ f(\mathbf{x}, y) \} \quad (6)$$

Since  $\text{rank}(y|S_f(\mathbf{x}))$  of a result is a discontinuous step function of the score, tractable learning algorithms typically optimize a substitute loss that is (sub-)differentiable [13, 30, 36]. Following this route, we now derive a tractable substitute for the empirical risk of (5) in terms of the scoring function. This is achieved by the following hinge-loss upper bound [15] on the rank

$$\begin{aligned} \text{rank}(y_i|\mathbf{y}) - 1 &= \sum_{\substack{y \in Y_i \\ y \neq y_i}} \mathbb{1}_{f(\mathbf{x}_i, y) - f(\mathbf{x}_i, y_i) > 0} \\ &\leq \sum_{\substack{y \in Y_i \\ y \neq y_i}} \max(1 - (f(\mathbf{x}_i, y_i) - f(\mathbf{x}_i, y)), 0). \end{aligned}$$

Using this upper bound, we can also get a bound for any IR metric that can be expressed through a monotonically increasing weighting function  $\lambda(r)$  of the rank. Note that this monotonicity condition is satisfied by all the metrics in Table 1. By rearranging terms and

applying the weighting function  $\lambda(r)$ , we have

$$\lambda(\text{rank}(y_i|\mathbf{y})) \leq \lambda \left( 1 + \sum_{\substack{y \in Y_i \\ y \neq y_i}} \max(1 - (f(\mathbf{x}_i, y_i) - f(\mathbf{x}_i, y)), 0) \right).$$

This provides the following continuous and subdifferentiable upper bound  $\hat{R}_{IPS}^{\text{hinge}}(f)$  on the propensity-weighted risk estimate of (5).

$$\begin{aligned} \hat{R}_{IPS}(S_f) &\leq \hat{R}_{IPS}^{\text{hinge}}(f) \\ &= \frac{1}{n} \sum_{i=1}^n \frac{1}{q_i} \lambda \left( 1 + \sum_{\substack{y \in Y_i \\ y \neq y_i}} \max(1 - (f(\mathbf{x}_i, y_i) - f(\mathbf{x}_i, y)), 0) \right) \end{aligned} \quad (7)$$

Focusing on the DCG metric, we show in the following how this upper bound can be optimized for linear as well as non-linear neural network scoring functions. For the general class of additive IR metrics, the optimization depends on the properties of the weighting function  $\lambda(r)$ , and we highlight them wherever appropriate.

#### 4.1 SVM PropDCG

The following derives an SVM-style method, called SVM PropDCG, for learning a linear scoring function  $f(\mathbf{x}, y) = \mathbf{w} \cdot \phi(\mathbf{x}, y)$ , where  $\mathbf{w}$  is a weight vector and  $\phi(\mathbf{x}, y)$  is a feature vector describing the match between query  $\mathbf{x}$  and result  $y$ . For such linear ranking functions – which are widely used in Ranking SVMs [13] and many other learning-to-rank methods [19] –, the propensity-weighted ERM bound from Equation (7) can be expressed as the following SVM-type optimization problem.

$$\begin{aligned} \hat{\mathbf{w}} &= \underset{\mathbf{w}, \xi}{\text{argmin}} \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + \frac{C}{n} \sum_{i=1}^n \frac{1}{q_i} \lambda \left( \sum_{y \in Y_i} \xi_{iy} + 1 \right) \\ \text{s.t.} \quad &\forall y \in Y_1 \setminus \{y_1\} : \mathbf{w} \cdot [\phi(\mathbf{x}_1, y_1) - \phi(\mathbf{x}_1, y)] \geq 1 - \xi_{1y} \\ &\vdots \\ &\forall y \in Y_n \setminus \{y_n\} : \mathbf{w} \cdot [\phi(\mathbf{x}_n, y_n) - \phi(\mathbf{x}_n, y)] \geq 1 - \xi_{ny} \\ &\forall i \forall y : \xi_{iy} \geq 0 \end{aligned}$$

$C$  is a regularization parameter. The training objective optimizes the  $\mathcal{L}_2$ -regularized hinge-loss upper bound on the empirical risk estimate (7). This upper bound holds since for any feasible  $(\mathbf{w}, \xi)$  and any monotonically increasing weighting function  $\lambda(r)$

$$\begin{aligned} &\lambda \left( 1 + \sum_{\substack{y \in Y_i \\ y \neq y_i}} \max(1 - (f(\mathbf{x}_i, y_i) - f(\mathbf{x}_i, y)), 0) \right) \\ &= \lambda \left( 1 + \sum_{\substack{y \in Y_i \\ y \neq y_i}} \max(1 - \mathbf{w} \cdot [\phi(\mathbf{x}_i, y_i) - \phi(\mathbf{x}_i, y)], 0) \right) \leq \lambda \left( 1 + \sum_{y \in Y_i} \xi_{iy} \right). \end{aligned}$$

As shown in [15], for the special case of using the sum of relevant ranks as the metric to optimize, i.e.  $\lambda(r) = r$ , this SVM optimization problem is a convex Quadratic Program which can be solved efficiently using standard SVM solvers, like SVM-rank [14], via a one-slack formulation.

Moving to the case of DCG as the training metric via the weighting function  $\lambda(r) = \frac{1}{\log(1+r)}$ , we get the following optimization problem for SVM PropDCG

$$\begin{aligned} \hat{\mathbf{w}} &= \underset{\mathbf{w}, \xi}{\text{argmin}} \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \frac{C}{n} \sum_{i=1}^n \frac{1}{q_i} \frac{1}{\log(\sum_{y \in Y_i} \xi_{iy} + 2)} \\ \text{s.t.} \quad &\forall j \forall y \in Y_i \setminus \{y_i\} : \mathbf{w} \cdot [\phi(\mathbf{x}_i, y_i) - \phi(\mathbf{x}_i, y)] \geq 1 - \xi_{iy} \\ &\forall j \forall y : \xi_{iy} \geq 0. \end{aligned}$$

This optimization problem is no longer a convex Quadratic Program. However, all constraints are still linear inequalities in the variables  $\mathbf{w}$  and  $\xi$ , and the objective can be expressed as the difference of two convex functions  $h$  and  $g$ . Let  $h(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$  and  $g(\xi) = \frac{C}{n} \sum_{j=1}^n \frac{1}{q_j} \frac{1}{\log(\sum_{y \in Y_i} \xi_{iy} + 2)}$ . Then the function  $h$  is the  $\mathcal{L}_2$  norm of the vector  $\mathbf{w}$  and is thus a convex function. As for the function  $g$ , the function  $k : x \mapsto \frac{1}{\log x}$  is convex as it is the composition of a convex decreasing function ( $x \mapsto \frac{1}{x}$ ) with a concave function ( $x \mapsto \log x$ ). So, since the sum of affine transformations of a convex function is convex,  $g$  is convex.

Such an optimization problem is called a convex-concave problem<sup>1</sup> and a local optimum can be obtained efficiently via the Convex-Concave Procedure (CCP) [18]. At a high level, the procedure works by repeatedly approximating the second convex function with its first order Taylor expansion which makes the optimization problem convex in each iteration. The Taylor expansion is first done at some chosen initial point in the feasible region, and then the solution of the convex problem in a particular iteration is used as the Taylor approximation point for the next iteration. It can be shown that this procedure converges to a local optimum [18].

Concretely, let  $\mathbf{w}^k, \xi^k$  be the solution in the  $k^{\text{th}}$  iteration. Then, we have the Taylor approximation

$$\begin{aligned} \hat{g}(\xi; \xi^k) &= g(\xi^k) + \nabla g(\xi^k)^T (\xi - \xi^k) \\ &= g(\xi^k) - \frac{C}{n} \sum_{j=1}^n \frac{1}{q_j} \frac{\sum_{y \in Y_i} \xi_{iy} - \xi_{iy}^k}{\left( \sum_{y \in Y_i} \xi_{iy}^k + 2 \right) \log^2 \left( \sum_{y \in Y_i} \xi_{iy}^k + 2 \right)} \end{aligned}$$

Letting  $q'_i = q_i \left( \sum_{y \in Y_i} \xi_{iy}^k + 2 \right) \log^2 \left( \sum_{y \in Y_i} \xi_{iy}^k + 2 \right)$ , and dropping the additive constant terms from  $\hat{g}$ , we get the following convex program that needs to be solved in each CCP iteration.

$$\begin{aligned} &\underset{\mathbf{w}, \xi}{\text{argmin}} \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + \frac{C}{n} \sum_{i=1}^n \frac{1}{q'_i} \sum_{y \in Y_i} \xi_{iy} \\ \text{s.t.} \quad &\forall i \forall y \in Y_i \setminus \{y_i\} : \mathbf{w} \cdot [\phi(\mathbf{x}_i, y_i) - \phi(\mathbf{x}_i, y)] \geq 1 - \xi_{iy} \\ &\forall i \forall y : \xi_{iy} \geq 0 \end{aligned}$$

Observe that this problem is of the same form as SVM PropRank, the Propensity Ranking SVM for the average rank metric, i.e.  $\lambda(r) = r$  (with the caveat that  $q'_i$  are not propensities). This nifty feature allows us to solve the convex problem in each iteration of the CCP using the fast solver for SVM PropRank provided in [15]. In our

<sup>1</sup>More generally, the inequality constraints can also be convex-concave and not just convex

experiments, CCP convergence was achieved within a few iterations – as detailed in the empirical section. For other IR metrics, the complexity and feasibility of the above Ranking SVM optimization procedure will depend on the form of the target IR metric. In particular, if the rank weighting function  $\lambda(r)$  is convex, it may be solved directly as a convex program. If  $\lambda(r)$  is concave, then the CCP may be employed as shown for the DCG metric above.

An attractive theoretical property of SVM-style methods is the ability to switch from linear to non-linear functions via the Kernel trick. In principle, kernelization can be applied to SVM PropDCG as is evident from the representer theorem [27]. Specifically, by taking the Lagrange dual, the problem can be kernelized analogous to [13]. While it can be shown that the dual is convex and strong duality holds, it is not clear that the optimization problem has a convenient and compact form that can be efficiently solved in practice. Even for the special case of the average rank metric,  $\lambda(r) = r$ , the associated kernel matrix  $K_{iy, jy'}$  has a size equal to the total number of candidates  $\sum_{i=1}^n |Y_i|^2$  squared, making the kernelization approach computationally infeasible or challenging at best. We therefore explore a different route for extending our approach to non-linear scoring functions in the following.

## 4.2 Deep PropDCG

Since moving to non-linear ranking functions through SVM kernelization is challenging, we instead explore deep networks as a class of non-linear scoring functions. Specifically, we replace the linear scoring function  $f(\mathbf{x}, y) = \mathbf{w} \cdot \phi(\mathbf{x}, y)$  with a neural network

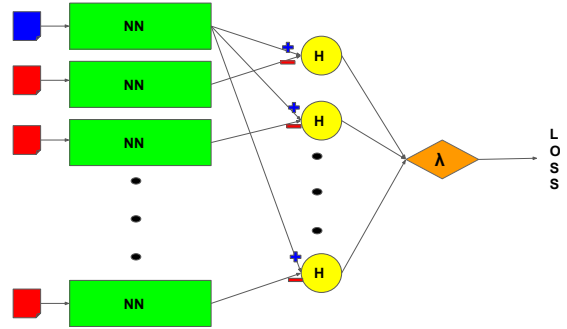
$$f(\mathbf{x}, y) = NN_{\mathbf{w}}[\phi(\mathbf{x}, y)] \quad (8)$$

This network is generally non-linear in both the weights  $\mathbf{w}$  and the features  $\phi(\mathbf{x}, y)$ . However, this does not affect the validity of the hinge-loss upper bound from Equation (7), which now takes the form

$$\frac{1}{n} \sum_{j=1}^n \frac{1}{q_i} \lambda \left( 1 + \sum_{\substack{y \in Y_i \\ y \neq y_i}} \max(1 - (NN_{\mathbf{w}}[\phi(\mathbf{x}_i, y_i)] - NN_{\mathbf{w}}[\phi(\mathbf{x}_i, y)]), 0) \right)$$

During training, we need to minimize this function with respect to the network parameters  $\mathbf{w}$ . Unlike in the case of SVM PropDCG, this function can no longer be expressed as the difference of a convex and a concave function, since  $NN_{\mathbf{w}}[\phi(\mathbf{x}_i, y_i)]$  is neither convex nor concave in general. Nevertheless, the empirical success of optimizing non-convex  $NN_{\mathbf{w}}[\phi(\mathbf{x}_i, y_i)]$  via gradient descent to a local optimum is well documented, and we will use this approach in the following. This is possible since the training objective is subdifferentiable as long as the weighting function  $\lambda(r)$  is differentiable. However, the non-linearity of  $\lambda(r)$  adds a challenge in applying *stochastic* gradient descent methods to our training objective, since the objective no longer decomposes into a sum over all  $(\mathbf{x}_i, y)$  as in standard network training. We discuss in the following how to handle this situation to arrive at an efficient stochastic-gradient procedure.

For concreteness, we again focus on the case of optimizing DCG via  $\lambda(r) = \frac{-1}{\log(1+r)}$ . In particular, plugging in the weighting function



**Figure 1: Deep PropDCG schema for computing the loss from one query instance. The blue document is the positive (clicked) result, and the red documents are the other candidates. The neural net NN is used to compute document scores for each set of candidate features. Pairs of scores are passed through the hinge node, and then finally the weighting function is applied as shown.**

for DCG, we get the Deep PropDCG minimization objective

$$\frac{1}{n} \sum_{j=1}^n \frac{-1}{q_i} \log^{-1} \left( 2 + \sum_{\substack{y \in Y_i \\ y \neq y_i}} \max(1 - (NN_{\mathbf{w}}[\phi(\mathbf{x}_i, y_i)] - NN_{\mathbf{w}}[\phi(\mathbf{x}_i, y)]), 0) \right)$$

to which a regularization term can be added (our implementation uses weight decay).

Since the weighting function ties together the hinge losses from pairs of documents in a non-linear way, stochastic gradient descent (SGD) is not directly feasible at the level of individual documents. In the case of DCG, since the rank weighting function is concave, one possible workaround is a Majorization-Minimization scheme [28] (akin to CCP): upper bound the loss function with a linear Taylor approximation at the current neural net weights, perform SGD at the level of document pairs  $(y_i, y)$  to update the weights, and repeat until convergence.

While this Majorization-Minimization scheme in analogy to the SVM approach is possible also for deep networks, we chose a different approach for the reasons given below. In particular, given the success of stochastic-gradient training of deep networks in other settings, we directly perform stochastic-gradient updates at the level of query instances, not individual  $(\mathbf{x}_i, y)$ . At the level of query instances, the objective does decompose linearly such that any subsample of query instances can provide an unbiased gradient estimate. Note that this approach works for any differentiable weighting function  $\lambda(r)$ , does not require any alternating approximations as in Majorization-Minimization, and processes each candidate document  $y$  including the clicked document  $y_i$  only once in one SGD step.

For SGD at the level of query instances, a forward pass of the neural network – with the current weights fixed – must be performed on each document  $y$  in candidate set  $Y_i$  in order to compute the loss from training instance  $(\mathbf{x}_i, y_i)$ . Since the number of documents in each candidate set varies, this is best achieved by processing

each input instance (including the corresponding candidate set) as a (variable-length) sequence so that the neural net weights are effectively shared across candidate documents for the same query instance.

This process is most easily understood via the network architecture illustrated in Figure 1. The scoring function  $NN_w[\phi(x_i, y_i)]$  is replicated for each result in the candidate set using shared weights  $w$ . In addition there is a hinge-loss node  $H(u, v) = \max(1 - (u - v), 0)$  that combines the score of the clicked result with each other result in the candidate set  $Y_i$ . For each such pair  $(y_i, y)$ , the corresponding hinge-loss node computes its contribution  $h_j$  to the upper bound on the rank. The result of the hinge-loss nodes then feeds into a single weighting node  $\Lambda(\vec{h}) = \lambda \left(1 + \sum_j h_j\right)$  that computes the overall bound on the rank and applies the weighting function. The result is the loss of that particular query instance.

Note that we have outlined a very general method which is agnostic about the size and architecture of the neural network. As a proof-of-concept, we achieved superior empirical results over a linear scoring function even with a simple two layer neural network, as seen in Section 5.8. We conjecture that DCG performance may be enhanced further with deeper, more specialized networks. Moreover, in principle, the hinge-loss nodes can be replaced with nodes that compute any other differentiable loss function that provides an upper bound on the rank without fundamental changes to the SGD algorithm.

## 5 EMPIRICAL EVALUATION

While the derivation of SVM PropDCG and Deep PropDCG has provided a theoretical justification for both methods, it still remains to show whether this theoretical argument translates to improved empirical performance. To this effect, the following empirical evaluation addresses three key questions.

First, we investigate whether directly optimizing DCG improves performance as compared to baseline methods, in particular, SVM PropRank as the most relevant method for unbiased LTR from implicit feedback, as well as LambdaRank, a common strong non-linear LTR method. Comparing SVM PropDCG to SVM PropRank is particularly revealing about the importance of direct DCG optimization, since both methods are linear SVMs and employ the same software machinery for the Quadratic Programs involved, thus eliminating any confounding factors. We also experimentally analyze the CCP optimization procedure to see whether SVM PropDCG is practical and efficient. Second, we explore the robustness of the generalized counterfactual LTR approach to noisy feedback, the severity of the presentation bias, and misspecification of the propensity model. And, finally, we compare the DCG performance of Deep PropDCG with a simple two layer neural network against the linear SVM PropDCG to understand to what extent non-linear models can be trained effectively using the generalized counterfactual LTR approach.

### 5.1 Setup

We conducted experiments on synthetic click data derived from two major LTR datasets, the Yahoo Learning to Rank Challenge corpus and LETOR4.0 [22]. LETOR4.0 contains two separate corpora: MQ2007 and MQ2008. Since MQ2008 is significantly smaller than

Dataset	# Avg. train clicks	# Train queries	# Features
Yahoo	173,986	20,274	699
LETOR4.0	25,870	1,484	46

**Table 2: Properties of the two benchmark datasets.**

Model	Avg. DCG (Yahoo)	Avg. DCG (LETOR4.0)
SVM Rank	0.6223 $\pm$ 8e-4	0.6841 $\pm$ 2e-3
LambdaRank	0.6435 $\pm$ 4e-4	0.6915 $\pm$ 4e-3
SVM PropRank	0.6410 $\pm$ 1e-3	0.7004 $\pm$ 1e-2
SVM PropDCG	0.6468 $\pm$ 2e-3	0.7043 $\pm$ 1e-2
Deep PropDCG	<b>0.6517 <math>\pm</math> 4e-4</b>	<b>0.7244 <math>\pm</math> 4e-3</b>

**Table 3: Performance comparison of different methods on two benchmark datasets ( $\eta = 1, \epsilon_- = 0.1, \epsilon_+ = 1$ ).**

Yahoo Learning to Rank Challenge, with only 784 queries, we follow the data augmentation approach proposed in [21], combining the MQ2007 and MQ2008 train sets for training and using the MQ2008 validation and test sets for validation and testing respectively.

Our experiment setup matches [15] for the sake of consistency and reproducibility. Briefly, the training and validation click data were generated from the respective full-information datasets (with relevances binarized) by simulating the position-based click model. Following [15], we use propensities that decay with presented rank of the result as  $p_r = \left(\frac{1}{r}\right)^\eta$ . The rankings that generate the clicks are given by a “production ranker” which was a conventional Ranking SVM trained on 1 percent of the full-information training data. The parameter  $\eta$  controls the severity of bias, with higher values causing greater position bias.

We also introduced noise into the clicks by allowing some irrelevant documents to be clicked. Specifically, an irrelevant document ranked at position  $r$  by the production ranker is clicked with probability  $p_r$  times  $\epsilon_-$ . When not mentioned otherwise, we used the parameters  $\eta = 1, \epsilon_- = 0.1$  and  $\epsilon_+ = 1$ , which is consistent with the setup used in [15]. Other bias profiles are also explored in the following.

Both the SVM PropRank and SVM PropDCG models were trained and cross-validated to pick the regularization constant  $C$ . For cross-validation, we use the partial feedback data in the validation set and select based on the IPS estimate of the DCG [29]. The performance of the models is reported on the binarized fully labeled test set which is never used for training or validation.

### 5.2 How do SVM PropDCG and Deep PropDCG compare against baselines?

We begin the empirical evaluation by comparing our counterfactual LTR methods against standard methods that follow a conventional ERM approach, namely LambdaRank and SVM-Rank. We generate synthetic click data using the procedure describe above, iterating over the training set 10 times for the Yahoo dataset and 100 times for MQ2008. This process was repeated over 6 independent runs, and we report the average performance along with the standard deviation over these runs. The regularization constant  $C$  for all SVM methods was picked based on the average DCG performance across the validation click data sampled over the 6 runs. Table 2

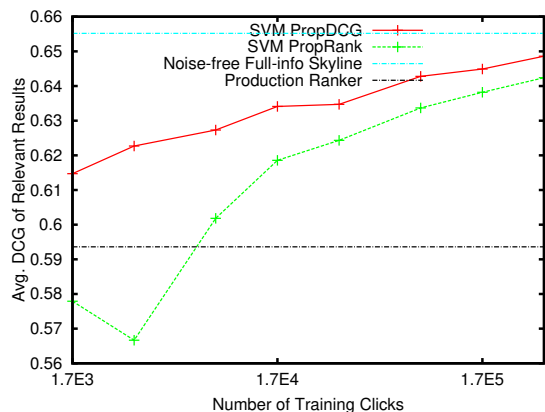


Figure 2: Test set Avg DCG performance for SVM PropDCG and SVM PropRank ( $\eta = 1, \epsilon_- = 0.1$ )

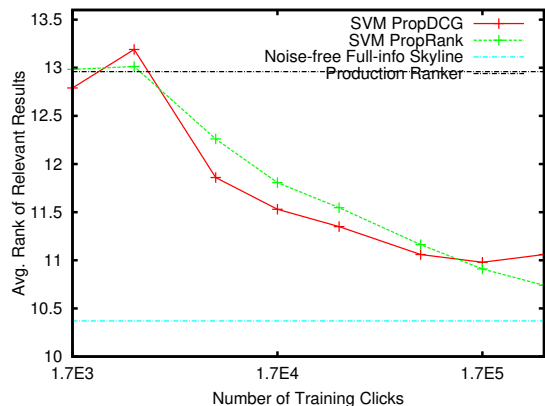


Figure 3: Test set Avg Rank performance for SVM PropDCG and SVM PropRank ( $\eta = 1, \epsilon_- = 0.1$ )

shows the average number of clicks along with other information about the training sets.

As a representative for non-linear LTR methods that use a conventional ERM approach, we also conducted experiments with LambdaRank as one of the most popular tree-based rankers. We use the LightGBM implementation [16]. During training, LambdaRank optimizes Normalized Discounted Cumulative Gain (NDCG). Since LambdaRank is a full-information method, we used clicks as relevance labels, i.e. all clicked documents as relevant and all non-clicked documents as irrelevant. The hyperparameters for LambdaRank, namely learning rate and the number of leaves were tuned based on the average DCG of clicked documents in the validation sets. More specifically, we performed a grid search to finetune learning rate from 0.001 to 0.1 and the number of leaves from 2 to 256. After tuning, we selected the learning rate to be 0.1, and the number of leaves to be 64 for the Yahoo dataset and 4 for MQ2008. We also made sure each split does not use more than 50% of the input features.

As shown in in Table 3, the counterfactual ERM approach via IPS weighting and directly optimizing for the target metric DCG yield superior results for SVM PropDCG and Deep PropDCG. The best results on both benchmarks are achieved by Deep PropDCG,

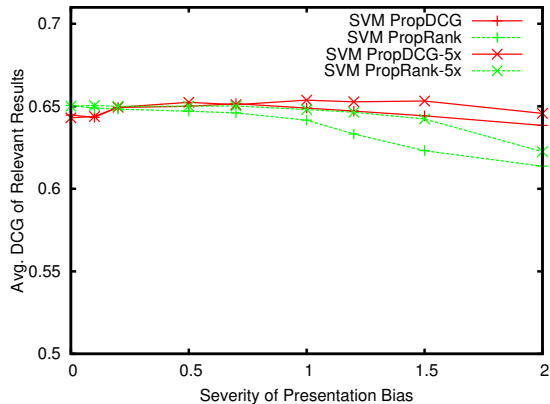


Figure 4: Test set Avg DCG performance for SVM PropDCG and SVM PropRank as presentation bias becomes more severe in terms of  $\eta$  ( $n = 45K$  and  $n = 225K, \epsilon_- = 0$ ).

which learns a two-layer neural network ranker. We conjecture that more sophisticated network architectures can further improve performance.

### 5.3 How does ranking performance scale with training set size?

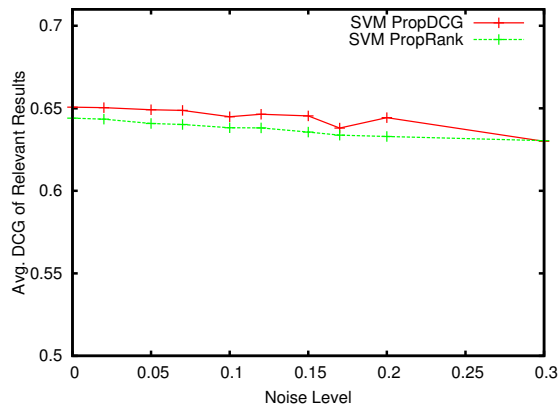
Next, we explore how the test-set ranking performance changes as the learning algorithm is given more and more click data. The resulting learning curves are given in Figures 2 and 3. The click data has presentation bias with  $\eta = 1$  and noise with  $\epsilon_- = 0.1$ . For small datasets, results are averaged over 3 draws of the click data. Both curves show the performance of the Production Ranker used to generate the click data, and the SVM skyline performance trained on the full-information training set. Ideally, rankers trained on click data should outperform the production ranker and approach the skyline performance.

Figure 2 shows that the DCG performance of both SVM PropDCG and SVM PropRank. As expected, both improve with increasing amounts of click data. Moreover, SVM PropDCG performs substantially better than the baseline SVM PropRank in maximizing test set DCG.

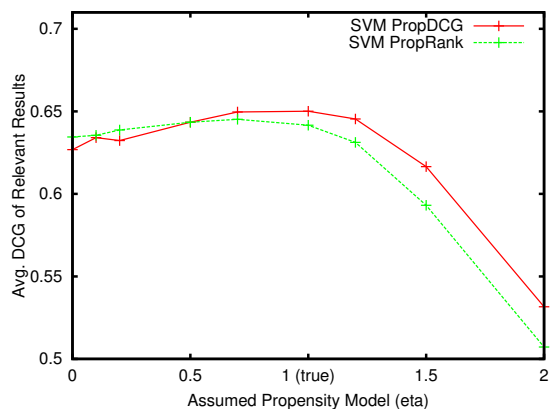
More surprisingly, Figure 3 shows both methods perform comparably in minimizing the average rank metric, with SVM PropDCG slightly better at smaller amounts of data and SVM PropRank better at larger amounts. We conjecture that this is due the variance-limiting effect of the DCG weights in SVM PropDCG when substituting the propensity weights  $q_i$  with the new constants  $q'_i$  in the SVM PropDCG CCP iterations. This serves as implicit variance control in the IPS estimator similar to clipping [15] by preventing propensity weights from getting too big. Since variance dominates estimation error at small amounts of data and bias dominates at large amounts, our conjecture is consistent with the observed trend.

### 5.4 How much presentation bias can be tolerated?

We now vary the severity of the presentation bias via  $\eta$  – higher values leading to click propensities more skewed to the top positions



**Figure 5: Test set Avg DCG performance for SVM PropDCG and SVM PropRank as the noise level increases in terms of  $\epsilon$  ( $n = 170K$ ,  $\eta = 1$ ).**



**Figure 6: Test set Avg DCG performance for SVM PropDCG and SVM PropRank as propensities are misspecified (true  $\eta = 1$ ,  $n = 170K$ ,  $\epsilon_- = 0.1$ ).**

– to understand its impact on the learning algorithm. Figure 4 shows the impact on DCG performance for both methods. We report performance for two training set sizes that differ by a factor of 5 (noise  $\epsilon_- = 0$ ). We see that SVM PropDCG is at least as robust to the severity of bias as SVM PropRank. In fact, SVM PropRank’s performance degrades more at high bias than that of SVM PropDCG, further supporting the conjecture that the DCG weighting in SVM PropDCG provides improved variance control which is especially beneficial when propensity weights are large. Furthermore, as also noted for SVM PropRank in [15], increasing the amount of training data by a factor of 5 improves performance of both methods due to variance reduction, which is an advantage that unbiased learning methods have over those that optimize a biased objective.

### 5.5 How robust is SVM PropDCG to noise?

Figure 5 shows the impact of noise on DCG performance, as noise levels in terms of  $\epsilon_-$  increase from 0 to 0.3. The latter results in click data where 59.8% of all clicks are on irrelevant documents. As expected, performance degrades for both methods as noise in-

creases. However, there is no evidence that SVM PropDCG is less robust to noise than the baseline SVM PropRank.

### 5.6 How robust is SVM PropDCG to misspecified propensities?

So far all experiments have had access to the true propensities that generated the synthetic click data. However, in real-world settings propensities need to be estimated and are necessarily subject to modeling assumptions. So, we evaluate the robustness of the learning algorithm to propensity misspecification.

Figure 6 shows the performance of SVM PropDCG and SVM PropRank when the training data is generated with  $\eta = 1$ , but the propensities used in learning are misspecified according to the  $\eta$  on the x-axis. The results show that SVM PropDCG is at least as robust to misspecified propensities as SVM PropRank. Both methods degrade considerably in the high bias regime when small propensities are underestimated – this is often tackled by clipping [15]. It is worth noting that SVM PropDCG performs better than SVM PropRank when misspecification leads to propensities that are underestimated, further strengthening the implicit variance control conjecture for SVM PropDCG discussed above.

### 5.7 How well does the CCP converge?

Next, we consider the computational efficiency of employing the CCP optimization procedure for training SVM PropDCG. Recall that the SVM PropDCG objective is an upper bound on the regularized (negative) DCG IPS estimate. It is optimized via CCP which repeatedly solves convex subproblems using the SVM PropRank solver until the objective value converges.

In Figure 7, optimization progress vs number of iterations as indicated by the change in objective value as well as the training DCG SNIPS estimate [29] is shown for 17K training clicks and the full range of regularization parameter  $C$  used in validation. The figure shows that the objective value usually converges in 3-5 iterations, a phenomenon observed in our experiments for other amounts of training data as well. In fact, the convergence tends to take slightly fewer iterations for larger amounts of data. The figure also shows that progress in objective is well-tracked with progress in the training DCG estimate, which suggests that the objective is a suitable upper bound for DCG optimization.

It is worth noting that restarting the optimizer across multiple CCP iterations can be substantially less time consuming than the initial solution that SVM PropRank computes. Since only the coefficients of the Quadratic Program change, the data does not need to be reloaded and the optimizer can be warm-started for quicker convergence in subsequent CCP iterations.

### 5.8 When does the non-linear model improve over the linear model?

We have seen that SVM PropDCG optimizes DCG better than SVM PropRank, and that it is a robust method across a wide range of biases and noise levels. Now we explore if performance can be improved further by introducing non-linearity via neural networks. Since the point of this paper is not a specific deep architecture but



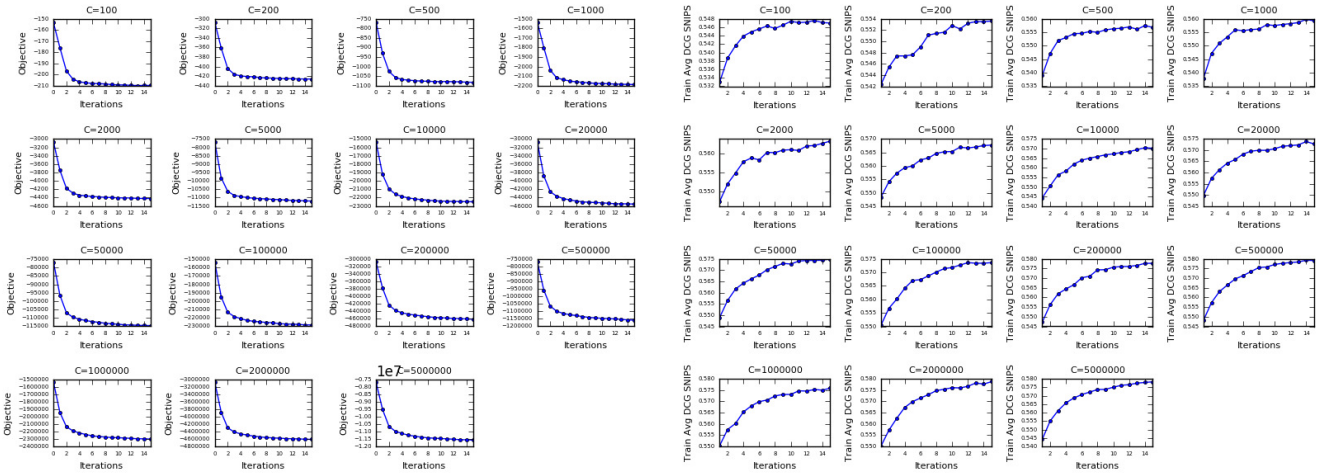


Figure 7: Optimization progress with respect to the number of CCP iterations. The objective value is shown in the left plots, and the training set DCG estimate on the right plots. Each plot corresponds to a particular value of regularization constant  $C$  ( $n = 17K$ ,  $\eta = 1$ ,  $\epsilon_- = 0.1$ ).

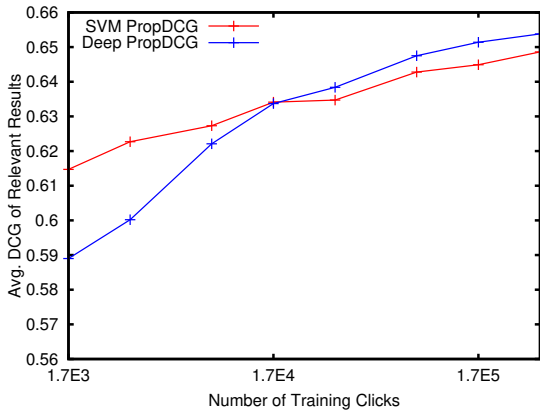


Figure 8: Test set Avg DCG performance for SVM PropDCG and Deep PropDCG ( $\eta = 1$ ,  $\epsilon_- = 0.1$ )

a novel training objective, we used a simple two-layer neural network with 200 hidden units and sigmoid activation. We expect that specialized deep architectures will further improve performance.

Figure 8 shows that Deep PropDCG achieves improved DCG compared to the linear SVM PropDCG given enough training data. For small amounts of training data, the linear model performs better, which is to be expected given the greater robustness to overfitting of linear models.

We also expect improved performance from tuning the hyperparameters of Deep PropDCG. In fact, we only used default parameters for Deep PropDCG, while we optimized the hyperparameters of SVM PropDCG on the validation set. In particular, Adam was used for stochastic gradient descent with weight decay regularizer at  $10^{-6}$ , minibatch size of 1000 documents and 750 epochs. The learning rate began at  $10^{-6}$  for the first 300 epochs, dropping by one order of magnitude in the next 200 epochs and another order of

magnitude in the remaining epochs. We did not try any other hyperparameter settings and these settings were held fixed across varying amounts of training data.

## 6 CONCLUSION

In this paper, we proposed a counterfactual learning-to-rank framework that is broad enough to cover a broad class of additive IR metrics as well as non-linear deep network models. Based on the generalized framework, we developed the SVM PropDCG and Deep PropDCG methods that optimize DCG via the Convex-Concave Procedure (CCP) and stochastic gradient descent respectively. We found empirically that SVM PropDCG performs better than SVM PropRank in terms of DCG, that it is robust to a substantial amount of presentation bias, noise and propensity misspecification, and that it can be optimized efficiently. DCG was improved further by using a neural network in Deep PropDCG.

There are many directions for future work. First, it is open for which other ranking metrics it is possible to develop efficient and effective methods using the generalized counterfactual framework. Second, the general counterfactual learning approach may also provide unbiased learning objectives for other settings beyond ranking, like full-page optimization and browsing-based retrieval tasks. Finally, it is an open question whether non-differentiable (e.g. tree-based) ranking models can be trained in the counterfactual framework as well.

## 7 ACKNOWLEDGMENTS

This research was supported in part by NSF Awards IIS-1615706 and IIS-1513692, an Amazon Research Award, and the Criteo Faculty Research Award program. All content represents the opinion of the authors, which is not necessarily shared or endorsed by their respective employers and/or sponsors.

## REFERENCES

- [1] Aman Agarwal, Ivan Zaitsev, Xuanhui Wang, Cheng Li, Marc Najork, and Thorsten Joachims. 2019. Estimating Position Bias without Intrusive Interventions. In *International Conference on Web Search and Data Mining (WSDM)*. 474–482.
- [2] Qingyao Ai, Keping Bi, Cheng Luo, Jiafeng Guo, and W. Bruce Croft. 2018. Unbiased Learning to Rank with Unbiased Propensity Estimation. In *The 41st International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*. ACM, New York, NY, USA, 385–394. <https://doi.org/10.1145/3209978.3209986>
- [3] Alexey Borisov, Ilya Markov, Maarten de Rijke, and Pavel Serdyukov. 2016. A Neural Click Model for Web Search. In *Proceedings of the 25th International Conference on World Wide Web (WWW)*. 531–541.
- [4] Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Hullender. 2005. Learning to Rank Using Gradient Descent. In *Proceedings of the 22nd International Conference on Machine Learning (ICML)*. ACM, New York, NY, USA, 89–96.
- [5] Christopher J Burges, Robert Ragno, and Quoc V Le. 2007. Learning to rank with nonsmooth cost functions. In *Advances in Neural Information Processing Systems (NeurIPS)*. 193–200.
- [6] Olivier Chapelle and Mingrui Wu. 2010. Gradient Descent Optimization of Smoothed Information Retrieval Metrics. *Information Retrieval* 13, 3 (June 2010), 216–235. <https://doi.org/10.1007/s10791-009-9110-3>
- [7] Olivier Chapelle and Ya Zhang. 2009. A dynamic bayesian network click model for web search ranking. In *International Conference on World Wide Web (WWW)*. ACM, 1–10.
- [8] Aleksandr Chuklin, Ilya Markov, and Maarten de Rijke. 2015. *Click Models for Web Search*. Morgan & Claypool Publishers.
- [9] Nick Craswell, Onno Zoeter, Michael Taylor, and Bill Ramsey. 2008. An Experimental Comparison of Click Position-bias Models. In *International Conference on Web Search and Data Mining (WSDM)*. ACM, 87–94.
- [10] Zhichong Fang, A. Agarwal, and T. Joachims. 2019. Intervention Harvesting for Context-Dependent Examination-Bias Estimation. In *ACM Conference on Research and Development in Information Retrieval (SIGIR)*.
- [11] D. G. Horvitz and D. J. Thompson. 1952. A Generalization of Sampling Without Replacement from a Finite Universe. *J. Amer. Statist. Assoc.* 47, 260 (1952), 663–685.
- [12] Ziniu Hu, Yang Wang, Qu Peng, and Hang Li. 2018. A Novel Algorithm for Unbiased Learning to Rank. (2018). [arXiv:cs.LR/1809.05818](https://arxiv.org/abs/1809.05818)
- [13] T. Joachims. 2002. Optimizing Search Engines Using Clickthrough Data. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*. 133–142.
- [14] T. Joachims, T. Finley, and Chun-Nam Yu. 2009. Cutting-Plane Training of Structural SVMs. *Machine Learning* 77, 1 (2009), 27–59.
- [15] Thorsten Joachims, Adith Swaminathan, and Tobias Schnabel. 2017. Unbiased Learning-to-Rank with Biased Feedback. In *ACM International Conference on Web Search and Data Mining (WSDM)*. ACM, New York, NY, USA, 781–789.
- [16] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In *Advances in Neural Information Processing Systems (NeurIPS)*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.). 3146–3154.
- [17] Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. 2011. Unbiased Offline Evaluation of Contextual-bandit-based News Article Recommendation Algorithms. In *International Conference on Web Search and Data Mining (WSDM)*. 297–306.
- [18] Thomas Lipp and Stephen Boyd. 2016. Variations and extension of the convex-concave procedure. *Optimization and Engineering* 17, 2 (2016), 263–287.
- [19] Tie-Yan Liu. 2009. Learning to Rank for Information Retrieval. *Foundations and Trends in Information Retrieval* 3, 3 (2009), 225–331.
- [20] Alistair Moffat and Justin Zobel. 2008. Rank-biased Precision for Measurement of Retrieval Effectiveness. *ACM Transactions on Information Systems (TOIS)* 27, 1 (2008), 2:1–2:27.
- [21] Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, Jingfang Xu, and Xueqi Cheng. 2017. DeepRank: A New Deep Architecture for Relevance Ranking in Information Retrieval. In *ACM Conference on Information and Knowledge Management (CIKM)*. ACM, 257–266.
- [22] Tao Qin and Tie-Yan Liu. 2013. Introducing LETOR 4.0 Datasets. *CoRR* abs/1306.2597 (2013).
- [23] L. Rigutini, T. Papini, M. Maggini, and F. Scarselli. 2011. SortNet: Learning to Rank by a Neural Preference Function. *IEEE Transactions on Neural Networks* 22, 9 (Sept 2011), 1368–1380.
- [24] Paul R. Rosenbaum and Donald B. Rubin. 1983. The central role of the propensity score in observational studies for causal effects. *Biometrika* 70, 1 (1983), 41–55.
- [25] T. Schnabel, A. Swaminathan, P. Frazier, and T. Joachims. 2016. Unbiased Comparative Evaluation of Ranking Functions. In *ACM International Conference on the Theory of Information Retrieval (ICTIR)*.
- [26] T. Schnabel, A. Swaminathan, A. Singh, N. Chandak, and T. Joachims. 2016. Recommendations as Treatments: Debiasing Learning and Evaluation. In *International Conference on Machine Learning (ICML)*.
- [27] B. Schoelkopf and A. J. Smola. 2002. *Learning with Kernels*. The MIT Press, Cambridge, MA.
- [28] A. Swaminathan and T. Joachims. 2015. Batch Learning from Logged Bandit Feedback through Counterfactual Risk Minimization. *Journal of Machine Learning Research (JMLR)* 16 (Sep 2015), 1731–1755.
- [29] A. Swaminathan and T. Joachims. 2015. The Self-Normalized Estimator for Counterfactual Learning. In *Neural Information Processing Systems (NeurIPS)*.
- [30] Michael Taylor, John Guiver, Stephen Robertson, and Tom Minka. 2008. SoftRank: Optimizing Non-smooth Rank Metrics. In *ACM International Conference on Web Search and Data Mining (WSDM)*. ACM, New York, NY, USA.
- [31] V. Vapnik. 1998. *Statistical Learning Theory*. Wiley, Chichester, GB.
- [32] Xuanhui Wang, Michael Bendersky, Donald Metzler, and Marc Najork. 2016. Learning to Rank with Selection Bias in Personal Search. In *ACM Conference on Research and Development in Information Retrieval (SIGIR)*. ACM.
- [33] Xuanhui Wang, Nadav Golbandi, Michael Bendersky, Donald Metzler, and Marc Najork. 2018. Position Bias Estimation for Unbiased Learning to Rank in Personal Search. In *ACM International Conference on Web Search and Data Mining (WSDM)*.
- [34] Yue Wang, Dawei Yin, Luo Jie, Pengyuan Wang, Makoto Yamada, Yi Chang, and Qiaozhu Mei. 2016. Beyond Ranking: Optimizing Whole-Page Presentation. In *Proceedings of the Ninth ACM International Conference on Web Search and Data Mining (WSDM)*. 103–112.
- [35] Mingrui Wu, Yi Chang, Zhaohui Zheng, and Hongyuan Zha. 2009. Smoothing DCG for Learning to Rank: A Novel Approach Using Smoothed Hinge Functions. In *ACM Conference on Information and Knowledge Management (CIKM)*. ACM, New York, NY, USA, 1923–1926. <https://doi.org/10.1145/1645953.1646266>
- [36] Yisong Yue, T. Finley, F. Radlinski, and T. Joachims. 2007. A Support Vector Method for Optimizing Average Precision. In *ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*. 271–278.