

# Statistical Analysis of Bayes Optimal Subset Ranking

David Cossock and Tong Zhang

**Abstract**—The ranking problem has become increasingly important in modern applications of statistical methods in automated decision making systems. In particular, we consider a formulation of the statistical ranking problem which we call subset ranking, and focus on the discounted cumulated gain (DCG) criterion that measures the quality of items near the top of the rank-list. Similar to error minimization for binary classification, direct optimization of natural ranking criteria such as DCG leads to a nonconvex optimization problems that can be NP-hard. Therefore, a computationally more tractable approach is needed. We present bounds that relate the approximate optimization of DCG to the approximate minimization of certain regression errors. These bounds justify the use of convex learning formulations for solving the subset ranking problem. The resulting estimation methods are not conventional, in that we focus on the estimation quality in the top-portion of the rank-list. We further investigate the asymptotic statistical behavior of these formulations. Under appropriate conditions, the consistency of the estimation schemes with respect to the DCG metric can be derived.

**Index Terms**—Bayes optimal, consistency, convex surrogate, ranking.

## I. INTRODUCTION

WE consider the general ranking problem, where a computer system is required to rank a set of items based on a given input. In such applications, the system often needs to present only a few top ranked items to the user. Therefore, the quality of the system output is determined by the performance near the top of its rank-list.

Ranking is especially important in electronic commerce and many internet applications, where personalization and information based decision making are critical to the success of such businesses. The decision making process can often be posed as a problem of selecting top candidates from a set of potential alternatives, leading to a conditional ranking problem. For example, in a recommender system, the computer is asked to choose a few items a user is most likely to buy based on the user's profile and buying history. The selected items will then be presented to the user as recommendations. Another important example that affects millions of people every day is the Internet search problem, where the user presents a query to the search engine, and the search engine then selects a few web pages that are most relevant to the query from the whole web. The quality of a search engine is largely determined by the top-ranked, or

highest ranked results the search engine can display on the first page. Internet search is the main motivation of this theoretical study, although the model presented here can be useful for many other applications. For example, another ranking problem is ad placement in a web page (either search result, or some content page) according to revenue-generating potential.

Although there has been much theoretical investigation of the ranking problem in recent years, many authors have only considered the global ranking problem, where a single ranking function is used to order a fixed set of items. However, in web search, a different ranking of web pages is needed for each different query. That is, we want to find a ranking function that is conditioned on (or local to) the query. We may look at the problem from an equivalent point of view: in web search, instead of considering many different ranking functions of web pages (one for each query), we may consider a single ranking function that depends on the (query, web page) pair. Given a query  $q$ , we only need to rank the subset of all possible (query, web page) pairs restricted to query  $q$ . Moreover, we may have a preprocessing step to filter out documents unlikely to be relevant to the query  $q$ , so that the subset of (query, web page) pairs to be ranked is further reduced. We call this framework *subset ranking*. A formal mathematical definition will be given in Section III.

For web search (and many other ranking problems), we are only interested in the quality of the top choices; the evaluation of the system output is different from many traditional error metrics such as classification error. In this setting, a useful figure of merit should focus on the top portion of the rank-list. This characteristic of ranking problems has not been carefully explored in earlier studies (except for a recent paper [24], which also touched on this issue). The purpose of this paper is to develop some theoretical results for converting a ranking problem into a convex optimization problem that can be efficiently solved. The resulting formulation focuses on the quality of the top ranked results. The theory can be regarded as an extension of related theory for convex risk minimization formulations for classification, which has drawn much attention recently in the statistical learning literature [3], [18], [26], [27], [30], [31].

Due to our motivation from web search, this paper focuses on ranking problems that measure quality only in the top portion of the rank-list. However, it is important to note that in some other applications, global ranking criteria such as Spearman rank correlation and Kendall's  $\tau$  metric are used. Such global metrics have been investigated in a number of papers in recent years.

We organize the paper as follows. Section II discusses earlier work in statistics and machine learning on global and pair-wise ranking. Section III introduces the subset ranking problem. We define two ranking metrics: one is the discounted cumulated gain (DCG) measure which we focus on in this paper, and the other is a measure that counts the number of correctly ranked

Manuscript received September 5, 2007; revised April 18, 2008. Current version published October 22, 2008. This work was supported in part by the National Science Foundation under Grant DMS-0706805. Part of the work was performed when the second author was at Yahoo Inc.

D. Cossock is with Yahoo Inc., Sunnyvale, CA 94089-0703 USA (e-mail: dcossock@yahoo-inc.com).

T. Zhang is with Rutgers–State University of New Jersey, Piscataway NJ 08854 USA (e-mail: tzhang@stat.rutgers.edu).

Communicated by P. L. Bartlett, Associate Editor for Pattern Recognition, Statistical Learning and Inference.

Digital Object Identifier 10.1109/TIT.2008.929939

pairs. The latter has been studied recently by several authors in the context of pair-wise preference learning. The order induced by the regression function is shown to be optimal with respect to both metrics. Section IV introduces some basic estimation methods for ranking. Although this paper focuses on the least squares regression based formulation, we also briefly discuss other possible loss functions approximating the optimal order here. The later sections develop bounds and consistency arguments for a faster convergence rate to optimal ranking than is possible with naive (uniform) regression. Section V contains the main theoretical results in this paper, where we show that the approximate minimization of certain regression errors leads to the approximate optimization of the ranking metrics defined earlier. This implies that asymptotically the nonconvex ranking problem can be solved using regression methods that are convex. Section VI presents the regression learning formulation derived from the theoretical results in Section V. Similar methods are currently used to optimize Yahoo’s production search engine. Section VII studies the generalization and asymptotic behavior of regression learning, where we focus on the  $L_2$ -regularization approach. Together with earlier theoretical results, we can establish the consistency of regression-based ranking under appropriate conditions.

## II. RANKING AND PAIR-WISE PREFERENCE LEARNING

In the standard (global) ranking problem, a set of items is ranked relative to each other according to a single criterion. The goal is to learn a ranking (or linear ordering) for all items from a small set of training items (with partially defined preference relations among them), so that the remaining items can be ranked according to the same criterion. However, as we have explained in the introduction, this paper considers subset ranking, where only a subset of items are ranked according to an input  $q$  (representing the query in the web search example).

Since our motivation is document retrieval, we will follow the convention of using  $q$  to represent query and  $p$  to represent pages to be retrieved. Further discussion on ranking in the document retrieval domain and some mathematical formulations can be found in [17], [21], although we use different notations here.

In the context of document retrieval, we may consider ranking as a prediction problem. The traditional prediction problem in statistical machine learning assumes that we observe an input vector  $q \in \mathcal{Q}$ , so as to predict an unobserved output  $p \in \mathcal{P}$ . However, in a ranking problem, if we assume  $\mathcal{P} = \{1, \dots, m\}$  contains  $m$  possible values, then instead of predicting a value in  $\mathcal{P}$ , we predict a permutation of  $\mathcal{P}$  that gives an optimal ordering of  $\mathcal{P}$ . That is, if we denote by  $\mathcal{P}!$  the set of permutations of  $\mathcal{P}$ , then the goal is to predict an output in  $\mathcal{P}!$ . There are two fundamental issues: first, how to measure the quality of ranking; second, how to learn a good ranking procedure from historical data.

At first glance, it may seem that we can simply cast the ranking problem as an ordinary prediction problem where the output space becomes  $\mathcal{P}!$ . However, the number of permutations in  $\mathcal{P}!$  is  $m!$ , which can be extremely large even for small

$m$ . Therefore, it is not practical to solve the ranking problem directly without imposing certain structures on the search space. Moreover, in practice, given a training point  $q \in \mathcal{Q}$ , we are generally not given an optimal permutation in  $\mathcal{P}!$  as the observed output. Instead, we may observe another form of output from which an optimal ranking can be inferred, but it may also contain extra information. For example, in web search, we may observe relevance score or click-through-rate of a web page with respect to a query, either using human editorial judgment, or from search-logs. The training procedure should take advantage of such information.

A common method of obtaining an optimal permutation in  $\mathcal{P}!$  is via a scoring function which maps a pair  $(q, p)$  in  $\mathcal{Q} \times \mathcal{P}$  to a real valued number  $r(q, p)$ . For each  $q$ , the predicted permutation in  $\mathcal{P}!$  induced by this scoring function is defined as the ordering of  $p \in \mathcal{P}$  sorted with nonincreasing value  $r(q, p)$ . This is the method we will focus on in this paper. In subset ranking, we consider an input space  $\mathcal{X} = \mathcal{Q} \times \mathcal{P}$ , and the goal is to learn a scoring function  $r(q, p)$  on  $\mathcal{X}$ , so that the items  $(q, p)$  can be ranked using this function.

Although the ranking problem has received considerable interest in machine learning due to many important applications in information processing systems, the problem has not been extensively studied in the traditional statistical literature. A relevant statistical model is ordinal regression [20]. In this model, we are still interested in predicting a single output. We consider the input space  $\mathcal{X} = \mathcal{Q} \times \mathcal{P}$ , and for each  $x = (q, p)$ , we observe an output value  $y \in \mathcal{Y}$ . Moreover, we assume that the values in  $\mathcal{Y} = \{1, \dots, L\}$  are ordered, and the cumulative probability  $P(y \leq j|x)$  ( $j = 1, \dots, L$ ) has the form  $\gamma(P(y \leq j|x)) = \theta_j + g_\beta(x)$ . In this model, both  $\gamma(\cdot)$  and  $g_\beta(\cdot)$  have known functional forms, and  $\theta$  and  $\beta$  are model parameters.

Note that the ordinal regression model induces a stochastic preference relationship on the input space  $\mathcal{X}$ . Consider two samples  $(x_1, y_1)$  and  $(x_2, y_2)$  on  $\mathcal{X} \times \mathcal{Y}$ . We say  $x_1 \prec x_2$  if and only if  $y_1 < y_2$ . This is a classification problem that takes a pair of inputs  $x_1$  and  $x_2$  and tries to predict whether  $x_1 \prec x_2$  or not (that is, whether the corresponding outputs satisfy  $y_1 < y_2$  or not). In this formulation, the optimal prediction rule to minimize classification error is induced by the ordering of  $g_\beta(x)$  on  $\mathcal{X}$  because if  $g_\beta(x_1) < g_\beta(x_2)$ , then  $P(y_1 < y_2|x_1, x_2) > 0.5$  (based on the ordinal regression model), which is consistent with the Bayes rule. Motivated by this observation, a support vector machine (SVM) ranking method is proposed in [15]. The idea is to reformulate ordinal regression as a model to learn a preference relation on the input space  $\mathcal{X}$ , which can be learned using pair-wise classification. Given the parameter  $\hat{\beta}$  learned from training data, the scoring function is simply  $r(q, p) = g_{\hat{\beta}}(x)$ , where  $x = (q, p)$ .

The pair-wise preference learning model has become the major ranking topic in the machine learning literature. For example, in addition to SVM, a similar method based on AdaBoost is proposed in [12]. The idea was also used in optimizing the Microsoft web search system [6].

A number of researchers have proposed theoretical analysis of ranking based on the pair-wise ranking model. The criterion is to minimize the error of pair-wise preference prediction when

we draw two pairs  $x_1$  and  $x_2$  randomly from the input space  $\mathcal{X}$ . That is, given a scoring function  $g : \mathcal{X} \rightarrow \mathcal{R}$ , the ranking loss is

$$\begin{aligned} & \mathbf{E}_{(X_1, Y_1)} \mathbf{E}_{(X_2, Y_2)} [I(Y_1 < Y_2)I(g(X_1) \geq g(X_2)) \\ & \quad + I(Y_1 > Y_2)I(g(X_1) \leq g(X_2))] \\ & = \mathbf{E}_{X_1, X_2} [P(Y_1 < Y_2 | X_1, X_2)I(g(X_1) \geq g(X_2)) \\ & \quad + P(Y_1 > Y_2 | X_1, X_2)I(g(X_1) \leq g(X_2))] \quad (1) \end{aligned}$$

where  $I(\cdot)$  denotes the indicator function. For binary output  $y = 0, 1$ , it is known that this metric is equivalent to the AUC measure (area under the receiver operating characteristic curve (ROC)) for binary classifiers up to a scaling, and it is closely related to the Mann-Whitney-Wilcoxon statistic [14]. In the literature, theoretical analysis has focused mainly on this ranking criterion (for example, see [1], [2], [8], [23]).

The pair-wise preference learning model has some limitations. First, although the criterion in (1) measures the global pair-wise ranking quality, it is not the best metric to evaluate some practical ranking systems such as those in web search. In such applications, a system does not need to rank all data-pairs, but only a subset of them each time. Furthermore, the topmost positions are of primary importance. Another issue with the pair-wise preference learning model is that the scoring function is usually learned by minimizing a convex relaxation of the pair-wise classification error, similar to large margin classification. However, if the preference relationship is noisy (that is, not all preference relations can be satisfied by a single ranking order), then an important question that should be addressed is whether such a learning algorithm leads to a Bayes optimal ranking function in the large sample limit. Unfortunately for general risk minimization formulations, the problem is quite complex and difficult if the decision rule is induced by a single-variable scoring function of the form  $r(x)$ . The regression formulations considered in this paper have simpler (although less general) forms, which can be more easily analyzed.

The problem of Bayes optimality in the pair-wise learning model was partially investigated in [8], but with a decision rule of a general form  $r(x_1, x_2)$ : we predict  $x_1 \prec x_2$  if  $r(x_1, x_2) < 0$ . To our knowledge, this method is not widely used in practice because a naive application can lead to contradictions: we may predict  $r(x_1, x_2) < 0$ ,  $r(x_2, x_3) < 0$ , and  $r(x_3, x_1) < 0$ . Therefore, in order to use such a method effectively for ranking, there needs to be a mechanism to resolve such contradictions (see [9]). For example, one possibility is to define a scoring function  $f(x) = \sum_{x'} r(x, x')$ , and rank the data accordingly. Another possibility is to use a sorting method (such as quick-sort) directly with the comparison function given by  $r(x_1, x_2)$ . In order to show that such contradiction resolution methods are well behaved asymptotically, it is necessary to analyze the corresponding error. Although useful, such analysis is beyond the scope of this paper.

It is worth mentioning that there are more elaborate ideas in the preference learning framework that can address some limitations of global pair-wise ranking. One proposal that involves only partially defined preference relations has been considered in [25]. Their methods allow optimization of a convex surrogate that can be more tightly coupled with the goals in

ranking, while still enabling fast optimization. The methods are directly motivated by considering the noiseless situation where all preference relationships can be fully satisfied. If the preference relationships cannot be fully satisfied, the behavior of their methods is not well-understood. In contrast, the regression formulations considered in this paper are easier to analyze. Our results imply that methods we propose are well-behaved even when the ranking problem contains noise (that is, we cannot find a ranking function that preserves all preference relationships).

### III. SUBSET RANKING MODEL

The global pair-wise preference learning model in Section II has some limitations. In this paper, we shall describe a model more relevant to some practical ranking systems such as web search.

#### A. Problem Definition

The problem of web search is to rank web pages (denoted by  $p$ ) based on a query  $q$ . As explained in the Introduction, one may consider this problem as learning a global ranking function that orders (query, web page) pairs. However, given a query  $q$ , we only need to rank the subset of all (query, web page) pairs that are consistent with the query  $q$ . Our goal is to rank the (query, web page) pairs according to how relevant the web page is to the query. As mentioned in Section II, this is achieved by taking  $(q, p)$  as input, and estimate a scoring function  $r(q, p)$  as output, so that a more relevant pair gets a higher score.

In practical applications, the standard input to a machine learning algorithm is usually not the original input, but a feature vector created from the input. Therefore, from now on, we denote the space of observable feature vectors as  $\mathcal{X}$ . In web search, feature vectors are constructed from  $(q, p)$ . For each  $x \in \mathcal{X}$ , we also observe a nonnegative real-valued variable  $y$  that measures the quality of  $x$  (this corresponds to the relevance score of the (query, web page) pair represented by  $x$ ). Denote by  $\mathcal{S}$  the set of all finite subsets of  $\mathcal{X}$  that may possibly contain elements that are redundant (which happens when two (query, web page) pairs map to the same feature vector). Note that sets with repeated memberships are some time referred to as *multiset* in the literature. We do not differentiate set and multiset here, as this detail is not essential in the paper. In subset ranking, each instance is an ordered finite subset  $S = \{x_1, \dots, x_m\} \in \mathcal{S}$  (i.e., the set of (query, web page) pairs consistent with the given query  $q$  in the web search example). Note that the actual order of the items in the set is of no importance; the numerical subscripts are for notational purposes only, so that permutations can be more conveniently defined.

In our model, we randomly draw an ordered subset  $S = \{x_1, \dots, x_m\} \in \mathcal{S}$  consisting of feature vectors  $x_j$  in  $\mathcal{X}$ ; at the same time, we are given a set of real-valued grades  $\{y_j\} = \{y_1, \dots, y_m\}$  such that for each  $j$ ,  $y_j$  corresponds to  $x_j$  (representing its relevance score). Whether the size of the set  $m$  should be a random variable has no importance in our analysis. In this paper, we assume that it is fixed for simplicity.

Based on the observed subset  $S = \{x_1, \dots, x_m\}$ , the system is required to output an ordering (ranking) of the items in the set. Using our notation, this ordering can be represented as a permutation  $J = [j_1, \dots, j_m]$  of  $[1, \dots, m]$ . Our goal is to

produce a permutation such that  $y_{j_i}$  is in decreasing order for  $i = 1, \dots, m$ . A quality criterion is used to evaluate the system produced ranking of subset  $S$  based on the observed grades  $\{y_j\}$ . Our goal is to maximize the expected value of this quality criterion over random draws of  $S$  and  $\{y_j\}$  from an underlying distribution  $D$ . A formal definition is given below, where for notational simplicity, we state it with fixed  $m$ .

*Definition 1:* In subset ranking, we are given a set of items  $\mathcal{X}$ . Let  $\mathcal{S} = \mathcal{X}^m$  be the set of (ordered) finite subsets of  $\mathcal{X}$  of cardinality  $m$ . Let  $Q(J, \{y_j\})$  be a real-valued quality function, where we denote by  $J$  a permutation of  $[1, \dots, m]$ , and  $\{y_j\} = [y_1, \dots, y_m] \in \mathbb{R}^m$  a vector of  $m$  grades. A ranking function  $r(S)$  maps an ordered subset  $S \in \mathcal{S}$  to a permutation  $J$ . Assume that  $S = \{x_1, \dots, x_m\}$  and  $\{y_j\}$  are drawn randomly from an underlying distribution  $D$ , which is invariant under any simultaneous permutation of  $(x_j, y_j)$  for  $j = 1, \dots, m$  (that is, the ordering of  $S$  is of no importance). Our goal is to find a ranking function  $r(S)$  to maximize the expected subset ranking quality

$$\mathbf{E}_{(S, \{y_j\}) \sim D} Q(r(S), \{y_j\}).$$

The machine learning problem is to estimate  $r(S)$  from a set of  $n$  training pairs  $\{(S, \{y_j\})\}$  that are independently drawn from  $D$ .

In practical applications, each available position  $i$  can be associated with a weight  $c_i$  that measures the importance of that position. Now, given the grades  $y_j (j = 1, \dots, m)$ , a very natural measure of the rank-list  $J = [j_1, \dots, j_m]$ 's quality is the following weighted sum:

$$DCG(J, [y_j]) = \sum_{i=1}^m c_i y_{j_i}.$$

We assume that  $\{c_i\}$  is a predefined sequence of nonincreasing nonnegative discount factors that may or may not depend on  $S$ . This metric, described in [16] as a DCG, is one of the preferred metrics used in the evaluation of internet search systems, including the production system of Yahoo and that of Microsoft [6]. In this context, a typical choice of  $c_i$  is to set  $c_i = 1/\log(1+i)$  when  $i \leq k$  and  $c_i = 0$  when  $i > k$  for some  $k$ . One may also use other choices, such as letting  $c_i$  be the probability of user viewing (or clicking) the result at position  $i$ .

Although introduced in the information retrieval (IR) context and applied to web search, the DCG criterion can be useful for many other ranking applications such as recommender systems. By choosing a decaying sequence of  $c_i$ , this measure naturally focuses on the quality of the top portion of the rank-list. This is in contrast with the pair-wise error criterion in (1), which does not distinguish the top portion of the rank-list from the bottom portion.

For the DCG criterion, the goal is to train a ranking function  $r$  that can take a subset  $S \in \mathcal{S}$  as input, and produce an output permutation  $J = r(S)$  such that the expected DCG is as large as possible

$$DCG(r) = \mathbf{E}_S DCG(r, S) \quad (2)$$

where

$$DCG(r, S) = \sum_{i=1}^m c_i \mathbf{E}_{y_{j_i} | (x_{j_i}, S)} y_{j_i}. \quad (3)$$

Note that in the above equation, we let  $r(S) = [j_1, \dots, j_m]$ . We also assume that  $x_{j_i}$  is the  $j_i$ th item in  $S$ , as in Definition 1. The notation  $\mathbf{E}_{y_{j_i} | (x_{j_i}, S)} y_{j_i}$  is the expectation of the score  $y_{j_i}$  corresponding to the item  $x_{j_i} \in S$ . We use  $x_{j_i}$  instead of  $j_i$  to indicate the position of  $y_{j_i}$  in order to emphasize that the ordering information in  $S$  is not important (see Definition 1). That is,  $y_{j_i}$  depends on  $x_{j_i}$  instead of the index  $j_i$  itself. Conditioned on the value  $x_{j_i}$ , the distribution of  $y_{j_i}$  is not affected when we reorder  $S$ .

The global pair-wise preference learning metric (1) can be adapted to the subset ranking setting. We may consider the following weighted total of correctly ranked pairs minus incorrectly ranked pairs:

$$\mathbf{T}(J, [y_j]) = \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (y_{j_i} - y_{j_{i'}}).$$

If the output label  $y_i$  takes binary values, and the subset  $S = \mathcal{X}$  is global (we may assume that it is finite), then this metric is equivalent to (1). Although we pay special attention to the DCG metric, we shall also include some analysis of the  $\mathbf{T}$  criterion for completeness.

Similar to (2) and (3), we can define the following quantities:

$$\mathbf{T}(r) = \mathbf{E}_S \mathbf{T}(r, S) \quad (4)$$

where if we let  $r(S) = [j_1, \dots, j_m]$ , then

$$\mathbf{T}(r, S) = \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (\mathbf{E}_{y_{j_i} | (x_{j_i}, S)} y_{j_i} - \mathbf{E}_{y_{j_{i'}} | (x_{j_{i'}}, S)} y_{j_{i'}}). \quad (5)$$

Similar to the concept of Bayes classifier in classification, we can define the Bayes ranking function that optimizes the  $DCG$  and  $\mathbf{T}$  measures. Based on the conditional formulations in (3) and (5), we have the following result.

*Theorem 1:* Given a set  $S \in \mathcal{S}$ , for each  $x_j \in S$ , we define the Bayes-scoring function as

$$f_B(x_j, S) = \mathbf{E}_{y_j | (x_j, S)} y_j.$$

An optimal Bayes ranking function  $r_B(S)$  that maximizes (5) returns a rank list  $J = [j_1, \dots, j_m]$  such that  $f_B(x_{j_1}, S) \geq f_B(x_{j_2}, S) \geq \dots \geq f_B(x_{j_m}, S)$ . An optimal Bayes ranking function  $r'_B(S)$  that maximizes (3) returns a rank list  $J = [j_1, \dots, j_m]$  such that  $c_k > c_{k'}$  implies that  $f_B(x_{j_k}, S) > f_B(x_{j_{k'}}, S)$ .

*Proof:* Without confusion, we use the same notation  $J = [j_1, \dots, j_m]$  to denote the rank list from an optimal Bayes ranking function  $r_B(S)$  for (5), or from  $r'_B(S)$  for (3). Given any  $k, k' \in \{1, \dots, m\}$ , we consider a modified ranking function  $r''_B(S)$  that returns a rank list  $J'$  derived from  $J$  by switching  $j_k$  and  $j_{k'}$ . That is, we define  $J' = [j'_1, \dots, j'_m]$ , where  $j'_i = j_i$  when  $i \neq k, k'$ , and  $j'_k = j_{k'}$ , and  $j'_{k'} = j_k$ .

Due to the optimality of Bayes ranking function, we have  $\mathbf{T}(J, S) \geq \mathbf{T}(J', S)$  for (5) and  $\mathbf{DCG}(J, S) \geq \mathbf{DCG}(J', S)$  for (3).

We consider the  $\mathbf{T}$ -criterion first. In order to show that  $f_B(x_{j_k}, S)$  is in descending order, we only need to show that  $f_B(x_{j_{k+1}}, S) \leq f_B(x_{j_k}, S)$  for  $k = 1, \dots, m-1$ . Consider  $k' = k+1$ , it is easy to check that

$$\mathbf{T}(J', S) - \mathbf{T}(J, S) = 4(f_B(x_{j_{k+1}}, S) - f_B(x_{j_k}, S)) / (m(m-1)).$$

Therefore,  $\mathbf{T}(J', S) \leq \mathbf{T}(J, S)$  implies that  $f_B(x_{j_{k+1}}, S) \leq f_B(x_{j_k}, S)$ . This proves the first claim of the theorem. Note that the reverse is also true because all ranking orders with  $f_B(x_{j_k}, S)$  in descending order have the same  $\mathbf{T}$ -value.

Now we consider the  $\mathbf{DCG}$ -criterion. Given any  $k, k'$ , we have

$$\mathbf{DCG}(J', S) - \mathbf{DCG}(J, S) = (c_k - c_{k'}) (f_B(x_{j_{k'}}, S) - f_B(x_{j_k}, S)).$$

Now  $c_k > c_{k'}$  and  $\mathbf{DCG}(J', S) \leq \mathbf{DCG}(J, S)$  implies that  $f_B(x_{j_k}, S) \geq f_B(x_{j_{k'}}, S)$ . This proves the second claim. Note that the reverse is also true because all such ranking orders the same  $\mathbf{DCG}$ -value.  $\square$

The result indicates that the optimal ranking can be induced by a single variable scoring function of the form  $f(x, S) : \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{R}$ . In this paper,  $f(x, S)$  is meaningful only in the subset of  $\mathcal{X} \times \mathcal{S}$  with  $x \in S$ . However, for notational simplicity, we consider the domain of  $f(x, S)$  as  $\mathcal{X} \times \mathcal{S}$ , where we simply let  $f(x, S) = 0$  when  $x \notin S$ . We use a specific function form  $f(x, S)$  to emphasize that the order of  $S$  is of no importance. In fact, in practical applications, we usually consider a function  $f(x, S)$  that is independent of the ordering of  $S$ . That is, if  $S'$  contains the same elements as  $S$  in a different order, then  $f(x, S') = f(x, S)$ .

Due to the dependency of conditional probability of  $y$  on  $S$ , the optimal scoring function also depends on  $S$ . Therefore, a complete solution of the subset ranking problem can be difficult when  $m$  is large. In practice, in order to remove the set dependency, we may consider an encoding of  $S$  into a set-dependent feature vector  $g(S)$ , and then incorporate this information into the feature vector  $x$ . If we are able to find such salient features, then the augmented feature vector contains all necessary information of  $S$ . Given the augmented  $x_j \in S$ , we may then assume that  $y_j$  is conditionally independent of  $S$ . That is,  $f_B(x_j, S) = \mathbf{E}_{y_j | (x_j, S)} y_j = \mathbf{E}_{y_j | x_j} y_j = f_B(x_j)$ , which removes the dependence on  $S$ . Note that this explicit formulation allows predictive features which are functions of the result set  $S$ , in addition to those which depend exclusively on  $x_j$ .

One may also consider the combinatorial problem of finding a global set independent scoring function  $f(x_j)$  that approximately maximizes the DCG (i.e., approximately preserves the ranking order of  $f_B(x_j, S)$ ). In the general case, this problem is computationally difficult. To see this, we may consider for simplicity that  $\mathcal{X}$  is finite; moreover, each subset only contains two elements, and one is preferred over the other (deterministically). Now in the subset learning model, such a preference relationship  $x \prec x'$  of two elements  $x, x' \in \mathcal{X}$  can be denoted by a directed edge from  $x$  to  $x'$ . In this setting, to find a global scoring function that approximates the optimal set dependent

Bayes scoring rule is equivalent to finding a maximum subgraph that is acyclic. In general, this problem is known to be NP-hard as well as APX-hard [11]: the class APX consists of problems having an approximation to within  $1+c$  of the optimum for some  $c$ . If any APX-hard problem admits a polynomial time approximation scheme, then  $P = NP$ .

## B. Web-Search Example

Web search is a concrete example of subset ranking. In this application, a user submits a query  $q$ , and expects the search engine to return a rank-list of web pages  $\{p_j\}$  such that a more relevant page is placed before a less relevant page. In a typical internet search engine, the system takes a query and uses a simple ranking formula for the initial filtering, which limits the set of web pages to an initial pool  $\{p_j\}$  of size  $m$  (e.g.,  $m = 100000$ ).

After this initial ranking, the system utilizes a more complicated second stage ranking process, which reorders the pool. This critical stage is the focus of this paper. This step generates a feature vector  $x_j$  for each page  $p_j$  in the initial pool, using information about the query  $q$ , page  $p_j$ , query/page matching relationships, or the result pool itself. The feature vector can encode various types of information, such as the query length, characteristics of the pool, number of query terms that match in the title or body of  $p_j$ , web linkage of  $p_j$ , etc. The set of all possible feature vectors  $x_j$  is  $\mathcal{X}$ . The ranking algorithm only observes a list of feature vectors  $\{x_1, \dots, x_m\}$  with each  $x_j \in \mathcal{X}$ . A human editor is presented with a pair  $(q, p_j)$  and assigns a score  $s_j$  on a scale, e.g., 1–5 (least relevant to highly relevant). The corresponding target value  $y_j$  is defined as a transformation of  $s_j$ ,<sup>1</sup> which maps the grade into the interval  $[0, 1]$ . Another possible choice of  $y_j$  is to normalize it by multiplying each  $y_j$  by a factor such that the optimal DCG is no more than one.

## IV. RISK MINIMIZATION BASED ESTIMATION METHODS

From the previous section, we know that the optimal scoring function is the conditional expectation of the grades  $y$ . We investigate some basic estimation methods for conditional expectation learning.

### A. Relation to Multicategory Classification

The subset ranking problem is a generalization of multicategory classification. In the latter case, we observe an input  $x_0$ , and are interested in classifying it into one of the  $m$  classes. Let the output value be  $k \in \{1, \dots, m\}$ . We encode the input  $x_0$  into  $m$  feature vectors  $\{x_1, \dots, x_m\}$ , where  $x_i = [0, \dots, 0, x_0, 0, \dots, 0]$  with the  $i$ th component being  $x_0$ , and the other components are zeros. We then encode the output  $k$  into  $m$  values  $\{y_j\}$  such that  $y_k = 1$  and  $y_j = 0$  for  $j \neq k$ . In this setting, we try to find a scoring function  $f$  such that  $f(x_k) > f(x_j)$  for  $j \neq k$ . Consider the DCG criterion with  $c_1 = 1$  and  $c_j = 0$  when  $j > 1$ . Then the classification accuracy is given by the corresponding DCG.

Given any multicategory classification algorithm, one may use it to solve the subset ranking problem (with nonnegative  $y_j$  value) as follows. Consider a sample  $S = [x_1, \dots, x_m]$  as input, and a set of outputs  $\{y_j\}$ . We randomly draw  $k$  from 1 to

<sup>1</sup>For example, the formula  $(2^{s_j} - 1) / (2^5 - 1)$  is used in [6]. Yahoo uses a different transformation based on empirical user surveys.

$m$  according to the distribution  $y_k / \sum_j y_j$ . We then form another sample with weight  $\sum_j y_j$ , which has the vector  $\bar{S} = [x_1, \dots, x_m]$  (where order is important) as input, and label  $y' = k \in \{1, \dots, m\}$  as output. This changes the problem formulation into multicategory classification. Since the conditional expectation can be expressed as

$$\mathbf{E}_{y_k|(x_k, S)} y_k = P(y' = k|S) \mathbf{E}_{\{y_j\}|S} \sum_j y_j;$$

the order induced by the scoring function  $\mathbf{E}_{y_k|(x_k, S)} y_k$  is the same as that induced by  $P(y' = k|S)$ . Therefore, a multicategory classification solver that estimates conditional probability can be used to solve the subset ranking problem. In particular, if we consider a risk minimization based multicategory classification solver for  $m$ -class problem [27], [30] of the following form:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \Phi(f(X_i), Y_i)$$

where  $(X_i, Y_i)$  are training points with  $Y_i \in \{1, \dots, m\}$ ,  $\mathcal{F}$  is a vector function class that takes values in  $R^m$ , and  $\Phi$  is some risk functional. Then for ranking with training points  $(\bar{S}_i, \{y_{i,1}, \dots, y_{i,m}\})$  and  $\bar{S}_i = [x_{i,1}, \dots, x_{i,m}]$ , the corresponding learning method becomes

$$\hat{f} = \arg \min_{f \in \bar{\mathcal{F}}} \sum_{i=1}^n \sum_{j=1}^m y_{i,j} \Phi(f(\bar{S}_i), j)$$

where the function space  $\bar{\mathcal{F}}$  contains a subset of functions  $\{f(\bar{S}) : \mathcal{X}^m \rightarrow R^m\}$  of the form

$$f(\bar{S}) = [f(x_1, S), \dots, f(x_m, S)]$$

where  $S = \{x_1, \dots, x_m\}$  is an unordered set. A concrete example is maximum entropy (multicategory logistic regression) which employs the following loss function:

$$\Phi(f(\bar{S}), j) = -f(x_j, S) + \ln \sum_{k=1}^m e^{f(x_k, S)}.$$

### B. Regression-Based Learning

Since in ranking problems  $y_{i,j}$  can take values other than 0 or 1, we can have more general formulations than multicategory classification. In particular, we may consider variations of the following regression-based learning method to train a scoring function in  $\mathcal{F} \subset \{\mathcal{X} \times \mathcal{S} \rightarrow R\}$ :

$$\begin{aligned} \hat{f} &= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \sum_{j=1}^m \phi(f(x_{i,j}, S_i), y_{i,j}) \\ S_i &= \{x_{i,1}, \dots, x_{i,m}\} \in \mathcal{S} \end{aligned} \quad (6)$$

where we assume that

$$\phi(a, b) = \phi_0(a) + \phi_1(a)b + \phi_2(b).$$

The estimation formulation is decoupled for each element  $x_{i,j}$  in a subset  $S_i$ , which makes the problem easier to solve.

In this method, each training point  $((x_{i,j}, S_i), y_{i,j})$  is treated as a single sample (for  $i = 1, \dots, n$  and  $j = 1, \dots, m$ ). The population version of the risk function is

$$\mathbf{E}_S \sum_{x \in S} [\phi_0(f(x, S)) + \phi_1(f(x, S)) \mathbf{E}_{y|(x, S)} y + \mathbf{E}_{y|(x, S)} \phi_2(y)].$$

This implies that the optimal population solution is a function that minimizes

$$\phi_0(f(x, S)) + \phi_1(f(x, S)) \mathbf{E}_{y|(x, S)} y$$

which is a function of  $\mathbf{E}_{y|(x, S)} y$ . Therefore, the estimation method in (6) leads to an estimator of conditional expectation with a reasonable choice of  $\phi_0(\cdot)$  and  $\phi_1(\cdot)$ .

A simple example is the least squares method, where we pick  $\phi_0(a) = a^2$ ,  $\phi_1(a) = -2a$ , and  $\phi_2(b) = b^2$ . That is, the learning method (6) becomes least squares estimation

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \sum_{j=1}^m (f(x_{i,j}, S_i) - y_{i,j})^2. \quad (7)$$

This method, and some essential variations which we will introduce later, will be the focus of our analysis.

It was shown in [7] that the only loss function with conditional expectation as the minimizer (for an arbitrary conditional distribution of  $y$ ) is least squares. However, for practical purposes, we only need to estimate a monotonic transformation of the conditional expectation. For this purpose, we can have additional loss functions of the form (6). In particular, let  $\phi_0(a)$  be an arbitrary convex function such that  $\phi'_0(a)$  is a monotone increasing function of  $a$ , then we may simply take the function  $\phi(a, b) = \phi_0(a) - ab$  in (6). The optimal population solution is uniquely determined by  $\phi'_0(f(x, S)) = \mathbf{E}_{y|(x, S)} y$ . A simple example is  $\phi_0(a) = a^4/4$  such that the population optimal solution is  $f(x, S) = (\mathbf{E}_{y|(x, S)} y)^{1/3}$ . Clearly, such a transformation does not affect ranking. Moreover, in many ranking problems, the range of  $y$  is bounded. It is known that additional loss functions can be used for computing the conditional expectation. As a simple example, if we assume that  $y \in [0, 1]$ , then the following modified least squares can be used:

$$\begin{aligned} \hat{f} &= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \sum_{j=1}^m [(1 - y_{i,j}) \max(0, f(x_{i,j}, S_i))^2 \\ &\quad + y_{i,j} \max(0, 1 - f(x_{i,j}, S_i))^2]. \end{aligned} \quad (8)$$

One may replace this with other loss functions used for binary classification that estimate conditional probability, such as those discussed in [31]. Although such general formulations might be interesting for certain applications, advantages over the simpler least squares loss of (7) are not completely certain, and they are more complicated to deal with. Therefore, we will not consider such general formulations in this paper, but rather focus on adapting the least squares method in (7) to the ranking problems. As we shall see, nontrivial modifications of (7) are necessary to optimize system performance near the top of rank-list.

### C. Pair-Wise Preference Learning

A popular idea in the recent machine learning literature is to pose the ranking problem as a pair-wise preference relationship

learning problem (see Section II). Using this idea, the scoring function for subset ranking can be trained by the following method:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \sum_{(j,j') \in E_i} \phi(f(x_{i,j}, S_i), f(x_{i,j'}, S_i); y_{i,j}, y_{i,j'}) \quad (9)$$

where each  $E_i = \{(j, j') : y_{i,j} < y_{i,j'}\}$  is a subset of  $\{1, \dots, m\} \times \{1, \dots, m\}$ . For example, we may use a nonincreasing monotone function  $\phi_0$  and let

$$\phi(a_1, a_2; b_1, b_2) = \phi_0((a_2 - a_1) - (b_2 - b_1))$$

or

$$\phi(a_1, a_2; b_1, b_2) = (b_2 - b_1)\phi_0(a_2 - a_1).$$

Example loss functions include SVM loss  $\phi_0(x) = \max(0, 1 - x)$  and AdaBoost loss  $\phi_0(x) = \exp(-x)$  (see [12], [15], [24]).

The approach works well if the ranking problem is noise-free (that is,  $y_{i,j}$  is deterministic). However, one difficulty with this approach is that if  $y_{i,j}$  is stochastic, then the corresponding population estimator from (9) may not be Bayes optimal, unless a more complicated scheme such as [8] is used. It will be interesting to investigate the error of such an approach, but the analysis is beyond the scope of this paper.

One argument used by the advocates of the pair-wise learning formulation is that we do not have to learn an absolute grade judgment (or its expectation), but rather only the relative judgment that one item is better than another. In essence, this means that for each subset  $S$ , if we shift each judgment by a constant, the ranking is not affected. If invariance with respect to a set-dependent judgment shift is a desirable property, then it can be incorporated into the regression-based model [28]. For example, we may introduce an explicit set dependent shift feature (which is rank-preserving) into (6)

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \min_{b_i \in \mathbb{R}} \sum_{j=1}^m \phi(f(x_{i,j}, S_i) + b_i, y_{i,j}).$$

In particular, for least squares, we have the following method:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \min_{b_i \in \mathbb{R}} \sum_{j=1}^m (f(x_{i,j}, S_i) + b_i - y_{i,j})^2. \quad (10)$$

More generally, we may also introduce more sophisticated set-dependent features (such as set-dependent scaling factors) and even hierarchical set-dependent models into the regression formulation. This general approach can remove many limitations of the standard regression method when compared to the pair-wise approach.

## V. CONVEX SURROGATE BOUNDS

The subset ranking problem defined in Section III is combinatorial in nature, which is very difficult to solve. Since the optimal Bayes ranking rule is given by conditional expectation, in Section IV, we discussed various formulations to estimate the conditional expectation. In particular, we are interested in

least squares regression based methods. In this context, a natural question to ask is if a scoring function approximately minimizes regression error, how well it can optimize ranking metrics such as DCG or  $\mathbf{T}$ . This section provides some theoretical results that relate the optimization of the ranking metrics defined in Section III to the minimization of regression errors. This allows us to design appropriate convex learning formulations that improve the simple least squares methods in (7) and (10).

*Definition 2:* A scoring function  $f(x, S)$  maps each  $x \in S$  to a real valued score. It induces a ranking function  $r_f$ , which ranks elements  $\{x_j\}$  of  $S$  in descending order of  $f(x_j)$ .

We are interested in bounding the **DCG** performance of  $r_f$  compared with that of  $f_B$ . This can be regarded as extensions of Theorem 1 that motivate regression based learning.

*Theorem 2:* Let  $f(x, S)$  be a real-valued scoring function, which induces a ranking function  $r_f$ . Consider pair  $p, q \in [1, \infty]$  such that  $1/p + 1/q = 1$ . We have the following relationship for each  $S = \{x_1, \dots, x_m\}$ :

$$\begin{aligned} & \mathbf{DCG}(r_B, S) - \mathbf{DCG}(r_f, S) \\ & \leq \left( 2 \sum_{i=1}^m c_i^p \right)^{1/p} \left( \sum_{j=1}^m |f(x_j, S) - f_B(x_j, S)|^q \right)^{1/q}. \end{aligned}$$

*Proof:* Let  $S = \{x_1, \dots, x_m\}$ ,  $r_f(S) = J = [j_1, \dots, j_m]$ , and  $r_B(S) = J_B = [j_1^*, \dots, j_m^*]$ . We have

$$\begin{aligned} & \mathbf{DCG}(r_f, S) \\ & = \sum_{i=1}^m c_i f_B(x_{j_i}, S) \\ & = \sum_{i=1}^m c_i f(x_{j_i}, S) + \sum_{i=1}^m c_i (f_B(x_{j_i}, S) - f(x_{j_i}, S)) \\ & \geq \sum_{i=1}^m c_i f(x_{j_i^*}, S) + \sum_{i=1}^m c_i (f_B(x_{j_i}, S) - f(x_{j_i}, S)) \\ & = \sum_{i=1}^m c_i f_B(x_{j_i^*}, S) + \sum_{i=1}^m c_i (f(x_{j_i^*}, S) - f_B(x_{j_i^*}, S)) \\ & \quad + \sum_{i=1}^m c_i (f_B(x_{j_i}, S) - f(x_{j_i}, S)) \\ & \geq \mathbf{DCG}(r_B, S) - \sum_{i=1}^m c_i (f(x_{j_i}, S) - f_B(x_{j_i}, S))_+ \\ & \quad - \sum_{i=1}^m c_i (f_B(x_{j_i^*}, S) - f(x_{j_i^*}, S))_+ \\ & \geq \mathbf{DCG}(r_B, S) - \left( 2 \sum_{i=1}^m c_i^p \right)^{1/p} \\ & \quad \times \left( \sum_{j=1}^m |f(x_j, S) - f_B(x_j, S)|^q \right)^{1/q} \end{aligned}$$

where we used the notation  $(z)_+ = \max(0, z)$ . The first inequality in the above derivation is a direct consequence of the

definition of  $J = r_f(S)$ , which implies that  $J = [j_1, \dots, j_m]$  achieves the maximum value of  $\sum_{i=1}^m c_i f(x_{j_i}, S)$  among all possible permutations  $[j_1, \dots, j_m]$  of  $[1, \dots, m]$  (see proof of Theorem 1). The last inequality is due to Hölder's inequality.  $\square$

The above theorem shows that the DCG criterion can be bounded through regression error. Although the theorem applies to any arbitrary pair of  $p$  and  $q$  such that  $1/p + 1/q = 1$ , the most useful case is with  $p = q = 2$ . This is because in this case, the problem of minimizing  $\sum_{j=1}^m (f(x_j, S) - f_B(x_j, S))^2$  can be directly achieved using least squares regression in (7). If regression error goes to zero, then the resulting ranking converges to the optimal DCG. Similarly, we can show the following result for the  $\mathbf{T}$  criterion.

*Theorem 3:* Let  $f(x, S)$  be a real-valued scoring function, which induces a ranking function  $r_f$ . We have the following relationship for each  $S = \{x_1, \dots, x_m\}$ :

$$\begin{aligned} \mathbf{T}(r'_B, S) - \mathbf{T}(r_f, S) &\leq \frac{4}{\sqrt{m}} \left( \sum_{j=1}^m (f(x_j, S) - f_B(x_j, S))^2 \right)^{1/2} \end{aligned}$$

where  $r'_B$  is an optimal Bayes ranking function for the  $\mathbf{T}$  criterion which is characterized in Theorem 1.

*Proof:* Let  $S = \{x_1, \dots, x_m\}$ ,  $r_f(S) = J = [j_1, \dots, j_m]$ , and  $r'_B(S) = J_B = [j_1^*, \dots, j_m^*]$ . We have

$$\begin{aligned} \mathbf{T}(r_f, S) &= \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (f_B(x_{j_i}, S) - f_B(x_{j_{i'}}, S)) \\ &= \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (f(x_{j_i}, S) - f(x_{j_{i'}}, S)) \\ &\quad - \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m [(f(x_{j_i}, S) - f_B(x_{j_i}, S)) \\ &\quad - (f(x_{j_{i'}}, S) - f_B(x_{j_{i'}}, S))] \\ &\geq \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (f(x_{j_i}, S) - f(x_{j_{i'}}, S)) \\ &\quad - \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m [|f(x_{j_i}, S) \\ &\quad - f_B(x_{j_i}, S)| + |f(x_{j_{i'}}, S) - f_B(x_{j_{i'}}, S)|]. \end{aligned}$$

By simplifying the last term, we have

$$\begin{aligned} \mathbf{T}(r_f, S) &\geq \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (f(x_{j_i}, S) - f(x_{j_{i'}}, S)) \\ &\quad - \frac{2}{m} \sum_{i=1}^m |f(x_i, S) - f_B(x_i, S)| \\ &\geq \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (f(x_{j_i^*}, S) - f(x_{j_{i'}^*}, S)) \\ &\quad - \frac{2}{m} \sum_{i=1}^m |f(x_i, S) - f_B(x_i, S)| \end{aligned}$$

$$\begin{aligned} &= \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (f_B(x_{j_i^*}, S) - f_B(x_{j_{i'}^*}, S)) \\ &\quad - \frac{2}{m} \sum_{i=1}^m |f(x_i, S) - f_B(x_i, S)| \\ &\quad - \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m [(f_B(x_{j_i^*}, S) - f(x_{j_i^*}, S)) \\ &\quad - (f_B(x_{j_{i'}^*}, S) - f(x_{j_{i'}^*}, S))] \\ &\geq \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (f_B(x_{j_i^*}, S) - f_B(x_{j_{i'}^*}, S)) \\ &\quad - \frac{2}{m} \sum_{i=1}^m |f(x_i, S) - f_B(x_i, S)| \\ &\quad - \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m [|f_B(x_{j_i^*}, S) - f(x_{j_i^*}, S)| \\ &\quad + |f_B(x_{j_{i'}^*}, S) - f(x_{j_{i'}^*}, S)|] \\ &= \frac{2}{m(m-1)} \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (f_B(x_{j_i^*}, S) - f_B(x_{j_{i'}^*}, S)) \\ &\quad - \frac{4}{m} \sum_{i=1}^m |f(x_i, S) - f_B(x_i, S)| \\ &= \mathbf{T}(r'_B, S) - \frac{4}{m} \sum_{i=1}^m |f(x_i, S) - f_B(x_i, S)|. \end{aligned}$$

The first inequality is a simplification of the last inequality from the previous chain of derivations. The second inequality in the above derivation is a direct consequence of the definition of  $J = r_f(S)$ , which implies that  $J = [j_1, \dots, j_m]$  achieves the maximum value of

$$\mathbf{T}(J, [f(x_i, S)]) = \sum_{i=1}^{m-1} \sum_{i'=i+1}^m (f(x_{j_i}, S) - f(x_{j_{i'}}, S))$$

among all possible permutations of  $[1, \dots, m]$  (see proof of Theorem 1). Now, Jensen's inequality implies that

$$\begin{aligned} \frac{4}{m} \sum_{i=1}^m |f(x_i, S) - f_B(x_i, S)| &\leq \frac{4}{\sqrt{m}} \left( \sum_{i=1}^m (f(x_i, S) - f_B(x_i, S))^2 \right)^{1/2}. \end{aligned}$$

Combining this estimate with the previous inequality, we obtain the desired bound.  $\square$

The above approximation bounds imply that least square regression can be used to learn the optimal ranking functions. The approximation error converges to zero when  $f$  converges to  $f_B$  in  $L_2$ . However, in general, requiring  $f$  to converge to  $f_B$  in  $L_2$  is not necessary. More importantly, in real applications, we are often only interested in the top portion of the rank-list. Our bounds should reflect this practical consideration. Assume that the coefficients  $c_i$  in the DCG criterion decay fast, so that  $\sum_i c_i$  is bounded (independent of  $m$ ). In this case, we may pick  $p = 1$  and  $q = \infty$  in Theorem 2. If  $\sup_j |f(x_j, S) - f_B(x_j, S)|$  is

small, then we obtain a better bound than the least squares error bound with  $p = q = 1/2$  which depends on  $m$ .

However, we cannot ensure that  $\sup_j |f(x_j, S) - f_B(x_j, S)|$  is small using the simple least squares estimation in (7). Therefore, in the following, we develop a more refined bound for the DCG metric, which will then be used to motivate practical learning methods that improve on the simple least squares method.

*Theorem 4:* Let  $f(x, S)$  be a real-valued scoring function, which induces a ranking function  $r_f$ . Given  $S = \{x_1, \dots, x_m\}$ , let the optimal ranking order be  $J_B = [j_1^*, \dots, j_m^*]$ , where  $f_B(x_{j_i^*})$  is arranged in nonincreasing order. Assume that  $c_i = 0$  for all  $i > k$ . Then we have the following relationship for all  $\gamma \in (0, 1)$ ,  $p, q \geq 1$ , such that  $1/p + 1/q = 1$ ,  $u > 0$ , and subset  $K \subset \{1, \dots, m\}$  that contains  $j_1^*, \dots, j_k^*$ :

$$\begin{aligned} & \mathbf{DCG}(r_B, S) - \mathbf{DCG}(r_f, S) \\ & \leq C_p(\gamma, u) \left( \sum_{j \in K} |f(x_j, S) - f_B(x_j, S)|^q \right. \\ & \quad \left. + u \sup_{j \notin K} (f(x_j, S) - f'_B(x_j, S))_+^q \right)^{1/q} \end{aligned}$$

where  $(z)_+ = \max(z, 0)$ ,  $M = f_B(x_{j_k^*}, S)$ , and

$$\begin{aligned} C_p(\gamma, u) &= \frac{1}{1-\gamma} \left( 2 \sum_{i=1}^k c_i^p + u^{-p/q} \left( \sum_{i=1}^k c_i \right)^p \right)^{1/p} \\ f'_B(x_j, S) &= f_B(x_j, S) + \gamma(M - f_B(x_j, S))_+. \end{aligned}$$

*Proof:* Let  $S = \{x_1, \dots, x_m\}$ ,  $r_f(S) = J = [j_1, \dots, j_m]$ , and  $r_B(S) = J_B = [j_1^*, \dots, j_m^*]$ . Since  $(f_B(x_{j_i^*}, S) - M)_+$  is in descending order of  $i$ ,  $[j_i^*]$  achieves the maximum of  $\sum_{i=1}^n c_i (f_B(x_{j_i}, S) - M)_+$  among all possible permutations  $[j_i]$  of  $[1, \dots, m]$ . Therefore, using  $c_i = 0$  ( $i > k$ ), we have

$$\begin{aligned} & \sum_{i=1}^m c_i ((f_B(x_{j_i^*}, S) - M) - (f_B(x_{j_i}, S) - M)_+) \\ &= \sum_{i=1}^m c_i ((f_B(x_{j_i^*}, S) - M)_+ - (f_B(x_{j_i}, S) - M)_+) \\ &\geq 0. \end{aligned}$$

This implies that

$$\begin{aligned} & \sum_{i=1}^m c_i ((f_B(x_{j_i^*}, S) - M) - (f_B(x_{j_i}, S) - M)_+) \\ & \leq \frac{1}{1-\gamma} \sum_{i=1}^m c_i ((f_B(x_{j_i^*}, S) - M) - (f_B(x_{j_i}, S) - M)_+). \end{aligned}$$

Therefore

$$\begin{aligned} & \mathbf{DCG}(r_B, S) - \mathbf{DCG}(r_f, S) \\ &= \sum_{i=1}^m c_i ((f_B(x_{j_i^*}, S) - M) - (f_B(x_{j_i}, S) - M)) \\ &= \sum_{i=1}^m c_i ((f_B(x_{j_i^*}, S) - M) - (f_B(x_{j_i}, S) - M)_+) \end{aligned}$$

$$\begin{aligned} & + \sum_{i=1}^m c_i (M - f_B(x_{j_i}, S))_+ \\ & \leq \frac{1}{1-\gamma} \left[ \sum_{i=1}^m c_i ((f_B(x_{j_i^*}, S) - M) \right. \\ & \quad \left. - (f_B(x_{j_i}, S) - M)_+) \right. \\ & \quad \left. + (1-\gamma) \sum_{i=1}^m c_i (M - f_B(x_{j_i}, S))_+ \right] \\ &= \frac{1}{1-\gamma} \left[ \sum_{i=1}^m c_i ((f_B(x_{j_i^*}, S) - M) \right. \\ & \quad \left. - (f_B(x_{j_i}, S) - M)) \right. \\ & \quad \left. - \gamma \sum_{i=1}^m c_i (M - f_B(x_{j_i}, S))_+ \right]. \end{aligned}$$

By using the definition of  $f'_B$  to simplify the last inequality, we obtain

$$\begin{aligned} & \mathbf{DCG}(r_B, S) - \mathbf{DCG}(r_f, S) \\ & \leq \frac{1}{1-\gamma} \left[ \sum_{i=1}^m c_i ((f_B(x_{j_i^*}, S) - M) - (f'_B(x_{j_i}, S) - M)) \right] \\ & \leq \frac{1}{1-\gamma} \left( \sum_{i=1}^m c_i (f_B(x_{j_i^*}, S) - f(x_{j_i^*}, S)) \right. \\ & \quad \left. - \sum_{i=1}^m c_i (f'_B(x_{j_i}, S) - f(x_{j_i}, S)) \right) \\ & \leq \frac{1}{1-\gamma} \left( \sum_{i=1}^m c_i (f_B(x_{j_i^*}, S) - f(x_{j_i^*}, S))_+ \right. \\ & \quad \left. + \sum_{i=1}^m c_i (f(x_{j_i}, S) - f'_B(x_{j_i}, S))_+ \right) \\ & \leq \frac{1}{1-\gamma} \left( \sum_{i=1}^k c_i (f_B(x_{j_i^*}, S) - f(x_{j_i^*}, S))_+ \right. \\ & \quad \left. + \sum_{i=1}^k c_i (f(x_{j_i}, S) - f'_B(x_{j_i}, S))_+ I(j_i \in K) \right. \\ & \quad \left. + \left( \sum_{i=1}^k c_i \right) \sup_{j \notin K} (f(x_j, S) - f'_B(x_j, S))_+ \right) \\ & \leq \frac{1}{1-\gamma} \left( \left( 2 \sum_{i=1}^k c_i^p \right)^{1/p} \left( \sum_{j \in K} (f_B(x_j, S) - f(x_j, S))_+^q \right. \right. \\ & \quad \left. \left. + \sum_{j \notin K} (f(x_j, S) - f'_B(x_j, S))_+^q \right)^{1/q} \right. \\ & \quad \left. + \left( \sum_{i=1}^k c_i \right) \sup_{j \notin K} (f(x_j, S) - f'_B(x_j, S))_+ \right) \\ & \leq \frac{1}{1-\gamma} \left( \left( 2 \sum_{i=1}^k c_i^p \right)^{1/p} \left( \sum_{j \in K} |f_B(x_j, S) - f(x_j, S)|^q \right)^{1/q} \right. \\ & \quad \left. + \sum_{i=1}^k c_i \sup_{j \notin K} (f(x_j, S) - f'_B(x_j, S))_+ \right). \end{aligned}$$

In the preceding derivation, the second inequality uses the fact that  $J = [j_i]$  achieves the maximum value of  $\sum_{i=1}^m c_i f(x_{j_i}, S)$  among all possible permutations  $[j_i]$  of  $[1, \dots, m]$ . Hölder's inequality has been applied to obtain the second to the last inequality. The last inequality uses the fact  $f'_B(x_j, S) \geq f_B(x_j, S)$ , implying that

$$(f(x_j, S) - f'_B(x_j, S))_+^q \leq (f(x_j, S) - f_B(x_j, S))_+^q.$$

From the last inequality, we can apply the Hölder's inequality again to obtain the desired bound.  $\square$

If  $f_B(x_j, S) \geq 0$  for all  $x_j \in S$ , then  $f'_B(x_j, S) \geq \gamma f_B(x_{j_k^*}, S)$ . Therefore, in this case we have (with  $p = q = 2$ )

$$\begin{aligned} & DCG(r_B, S) - DCG(r_f, S) \\ & \leq C_2(\gamma, u) \left( \sum_{j \in K} |f(x_j, S) - f_B(x_j, S)|^2 \right. \\ & \quad \left. + u \sup_{j \notin K} (f(x_j, S) - \gamma M)_+^2 \right)^{1/2}. \quad (11) \end{aligned}$$

Intuitively, the bound says the following: if we can find a set  $K$  such that we are certain that an item  $j \notin K$  is irrelevant, then we only need to estimate the quality scores of the items  $j \in K$  reliably (using least squares regression), while making sure that the score  $f(x_j, S)$  for any irrelevant item  $j \notin K$  is no more than  $\gamma M$  (which is a much easier task than estimating the expected quality score). If  $|K| \ll m$ , then the bound is a significant improvement over the standard least squares bound in Theorem 2 because the right-hand side depends only on the sum of least squared error over  $|K|$  error terms instead of  $m$  error terms. Note that in standard least squares, we try to estimate the quality scores uniformly well. Although idealized, this assumption is quite reasonable for web search because given a query  $q$ , most pages can be reliably determined to be irrelevant. The remaining pages, which we are uncertain about, form the set  $K$  that is much smaller than the total number of pages  $m$ . This bound motivates an importance weighted regression formulation which we consider in Section VI. In practical implementations, we do not have to set  $\gamma M$  as in the theorem, but rather regard it as a tuning parameter.

The bound in Theorem 4 can still be refined. However, the resulting inequalities will become more complicated. Therefore, we will not include such bounds in this paper. Similar to Theorem 4, such refined bounds show that we do not have to estimate conditional expectation uniformly well. We present a simple example as illustration.

*Proposition 1:* Consider  $m = 3$  and  $S = \{x_1, x_2, x_3\}$ . Let  $c_1 = 2, c_2 = 1, c_3 = 0$ , and  $f_B(x_1, S) = 1, f_B(x_2, S) = f_B(x_3, S) = 0$ . Let  $f(x, S)$  be a real-valued scoring function, which induces a ranking function  $r_f$ . Then

$$\begin{aligned} & DCG(r_B, S) - DCG(r_f, S) \\ & \leq 2|f(x_1, S) - f_B(x_1, S)| + |f(x_2, S) - f_B(x_2, S)| \\ & \quad + |f(x_3, S) - f_B(x_3, S)|. \end{aligned}$$

The coefficients on the right-hand side cannot be improved.

*Proof:* Note that  $f$  is suboptimal only when either  $f(x_3, S) \geq f(x_1, S)$  or when  $f(x_3, S) \geq f(x_2, S)$ . This gives the following bound:

$$\begin{aligned} & DCG(r_B, S) - DCG(r_f, S) \\ & \leq I(f(x_2, S) \geq f(x_1, S)) + I(f(x_3, S) \geq f(x_1, S)) \\ & \leq I(|f(x_2, S) - f_B(x_2, S)| + |f(x_1, S) - f_B(x_1, S)| \geq 1) \\ & \quad + I(|f(x_3, S) - f_B(x_3, S)| + |f(x_1, S) - f_B(x_1, S)| \geq 1) \\ & \leq [|f(x_2, S) - f_B(x_2, S)| + |f(x_1, S) - f_B(x_1, S)|] \\ & \quad + [|f(x_3, S) - f_B(x_3, S)| + |f(x_1, S) - f_B(x_1, S)|] \\ & = 2|f(x_1, S) - f_B(x_1, S)| + |f(x_2, S) - f_B(x_2, S)| \\ & \quad + |f(x_3, S) - f_B(x_3, S)|. \end{aligned}$$

In the above,  $I(\cdot)$  is the set indicator function. To see that the coefficients cannot be improved, we simply note that the bound is tight when either  $f(x_1, S) = f(x_2, S) = f(x_3, S) = 0$ , or when  $f(x_1, S) = f(x_2, S) = 1$  and  $f(x_3, S) = 0$ , or when  $f(x_2, S) = 0$  and  $f(x_1, S) = f(x_3, S) = 1$ .  $\square$

The Proposition does not make the same assumption as in Theorem 4, where we assume that there are many irrelevant items that can be reliably identified. The proposition implies that even in more general situations, we should not weight all errors equally. In this example, getting  $x_1$  right is more important than getting  $x_2$  or  $x_3$  right. Conceptually, Theorem 4 and Proposition 1 show the following.

- Since we are interested in the top portion of the rank-list, we only need to estimate the top rated items accurately, while preventing the bottom items from being overestimated (the conditional expectations do not have to be estimated accurately).
- For ranking purposes, some points are more important than other. Therefore, we should bias our learning method to produce more accurate conditional expectation estimation at the more important points.

## VI. IMPORTANCE WEIGHTED REGRESSION

The key message from the analysis in Section V is that we do not have to estimate the conditional expectations equally well for all items. In particular, since we are interested in the top portion of the rank-list, Theorem 4 implies that we need to estimate the top portion more accurately than the bottom portion.

Motivated by this analysis, we consider a regression-based training method to solve the DCG optimization problem but weight different points differently according to their importance. We shall not discuss the implementation details for modeling the function  $f(x, S)$ , which is beyond the scope of this paper.

Let  $\mathcal{F}$  be a function space that contains functions  $\mathcal{X} \times \mathcal{S} \rightarrow R$ . We draw  $n$  sets  $S_1, \dots, S_n$  randomly, where  $S_i = \{x_{i,1}, \dots, x_{i,m}\}$ , with the corresponding grades  $\{y_{i,j}\} = \{y_{i,1}, \dots, y_{i,m}\}$  (in this notation,  $i$  is fixed and the set  $\{y_{i,j}\}$  ranges over  $j = 1, \dots, m$ ). Based on Theorem 2, the simple least squares regression (7) can be applied. However, this direct regression method is not adequate for many practical problems such as web search, for which there are many items to rank (that is,  $m$  is large) but only the top ranked pages are important.

This is because the method pays equal attention to relevant and irrelevant pages. In reality, one should pay more attention to the top-ranked (relevant) pages. The grades of lower rank pages do not need to be estimated accurately, as long as we do not overestimate them so that these pages appear in the top ranked positions.

The above mentioned intuition can be captured by Theorem 4 and Proposition 1, which motivate the following alternative training method:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n L(f, S_i, \{y_{i,j}\}) \quad (12)$$

where for  $S = \{x_1, \dots, x_m\}$ , with the corresponding  $\{y_j\}$ , we have the following importance weighted regression loss as in (11):

$$L(f, S, \{y_j\}) = \sum_{j=1}^m w(x_j, S)(f(x_j, S) - y_j)^2 + u \sup_j w'(x_j, S)(f(x_j, S) - \delta(S))_+^2 \quad (13)$$

where  $u$  is a nonnegative parameter. Throughout this section, an abstract notion of weight function  $w(x_j, S)$  is employed. This may equivalently represent a reweighting of either the training sample probability distribution or the loss function. The number  $\delta(S)$  is a set-dependent threshold that corresponds to  $\gamma M$  in (11). However, since  $\gamma M$  is not known, in practice,  $\delta(S)$  may be considered as a tuning parameter, and a good choice may be determined with heuristics or cross-validation.

A variation of this method is used to optimize the production system of Yahoo's internet search engine. The detailed implementation and parameter choices cannot be disclosed here.<sup>2</sup> It is also irrelevant for the purpose of this paper. However, in the following, we shall briefly explain the intuition behind (13) using Theorem 4, and some practical considerations.

It is useful to mention that the distribution of web search relevance is heavily skewed in the sense that relevance is a "rare-event." In fact, empirically, we observe that the distribution of a regression-based scoring function  $f$  approximately fits a probability model with exponential decay of the form  $P(f > t) < \exp(-at + b)$ . The significance is that the improved rate of convergence to optimal DCG over naive uniform regression is due to both the relative importance of top-ranked documents in the DCG cost function and the exponentially small probability of the relevant documents. In order to boost the presence of relevant documents that are of maximal importance to DCG, the readers should bear in mind that importance sampling can be used as a component in the weighting scheme  $w(x_j, S)$ .

The weight function  $w(x_j, S)$  in (13) is chosen so that it focuses only on the most important examples (the weight is set to zero for pages that we know are irrelevant). This part of the formulation corresponds to the first part of the bound in Theorem 4 (in that case, we choose  $w(x_j, S)$  to be one for the top part of the example with index set  $K$ , and zero otherwise). The usefulness of nonuniform weighting is also demonstrated in Proposition 1. The specific choice of the weight function requires various engineering considerations that are not important for the purpose of this paper. In general, if there are many items with

similar grades, then it is beneficial to give each of the similar items a smaller weight. In the second part of (13), we choose  $w'(x_j, S)$  so that it focuses on the examples not covered by  $w(x_j, S)$ . In particular, it only covers those data points  $x_j$  that are low-ranked with high confidence. We choose  $\delta(S)$  to be a small threshold that can be regarded as a lower bound of  $\gamma M$  in (11). An important observation is that although  $m$  is often very large, the number of points so that  $w(x_j, S)$  is nonzero is often small. Moreover,  $(f(x_j, S) - \delta(S))_+$  is not zero only when  $f(x_j, S) \geq \delta(S)$ . In practice, the number of these points is usually small (that is, most irrelevant pages will be predicted as irrelevant). Therefore, the formulation completely ignores those low-ranked data points such that  $f(x_j, S) \leq \delta(S)$ . This makes the learning procedure computationally efficient even when  $m$  is large. The analogy here is support vector machines, where only the support vectors are useful in the learning formulation. One can completely ignore samples corresponding to non support vectors.

In the practical implementation of (13), we can use an iterative refinement scheme, where we start with a small number of samples to be included in the first part of (13), and then put the low-ranked points into the second part of (13) only when their ranking scores exceed  $\delta(S)$ . In fact, one may also put these points into the first part of (13), so that the second part always has zero values (which makes the implementation simpler). In this sense, the formulation in (13) suggests a selective sampling scheme, in which we pay special attention to important and highly ranked data points, while completely ignoring most of the low ranked data points. In this regard, with appropriately chosen  $w(x, S)$ , the second part of (13) can be completely ignored.

The empirical risk minimization method in (12) approximately minimizes the following criterion:

$$Q(f) = \mathbf{E}_S L(f, S) \quad (14)$$

where

$$\begin{aligned} L(f, S) &= \mathbf{E}_{\{y_j\}|S} L(f, S, \{y_j\}) \\ &= \sum_{j=1}^m w(x_j, S) \mathbf{E}_{y_j|(x_j, S)} (f(x_j, S) - y_j)^2 \\ &\quad + u \sup_j w'(x_j, S) (f(x_j, S) - \delta(S))_+^2. \end{aligned}$$

The following theorem shows that under appropriate assumptions, approximate minimization of (14) leads to the approximate optimization of DCG. For clarity, the assumptions are idealized. For example, in practice the condition  $\delta(S) \leq \gamma f_B(x_{j_k^*}, S)$  may be violated for some  $S$ . However, as long as it holds for most  $S$ , the consequence of the theorem is still valid approximately. In this regard, the theorem itself should only be considered as a formal justification of (12) under idealized assumptions that specify good parameter choices in (12). The method itself may still yield good performance when some of the assumptions fail.

*Theorem 5:* Assume that  $c_i = 0$  for all  $i > k$ . Assume the following conditions hold for each  $S = \{x_1, \dots, x_m\}$ .

- Let the optimal ranking order be  $J_B = [j_1^*, \dots, j_m^*]$ , where  $f_B(x_{j_k^*})$  is arranged in nonincreasing order.
- For all  $x_j \in S$ ,  $f_B(x_j, S) \geq 0$ .

<sup>2</sup>Some aspects of the implementation were covered in [10].

- There exists  $\gamma \in [0, 1)$  such that  $\delta(S) \leq \gamma f_B(x_{j_k}^*, S)$ , where  $\delta(S)$  is an appropriately chosen set-dependent threshold in (13).
- For all  $f_B(x_j, S) > \delta(S)$ , we have  $w(x_j, S) \geq 1$ .
- Let  $w'(x_j, S) = I(w(x_j, S) < 1)$ .

Then the following results hold.

- A function  $f_*$  minimizes (14) if  $f_*(x_j, S) = f_B(x_j, S)$  when  $w(x_j, S) > 0$  and  $f_*(x_j, S) \leq \delta(S)$  otherwise.
- For all  $f$ , let  $r_f$  be the induced ranking function. Let  $r_B$  be the optimal Bayes ranking function, we have:

$$DCG(r_B) - DCG(r_f) \leq C(\gamma, u)(Q(f) - Q(f_*))^{1/2}.$$

*Proof:* Note that if  $f_B(x_j, S) > \delta(S)$ , then  $w(x_j, S) \geq 1$  and  $w'(x_j, S) = 0$ . Therefore, the minimizer  $f_*(x_j, S)$  should minimize  $\mathbf{E}_{y_j|(x_j, S)}(f(x_j, S) - y_j)^2$ , achieved at  $f_*(x_j, S) = f_B(x_j, S)$ . If  $f_B(x_j, S) \leq \delta(S)$ , then there are two cases.

- $w(x_j, S) > 0$ ,  $f_*(x_j, S)$  should minimize  $\mathbf{E}_{y_j|(x_j, S)}(f(x_j, S) - y_j)^2$ , achieved at

$$f_*(x_j, S) = f_B(x_j, S).$$

- $w(x_j, S) = 0$ ,  $f_*(x_j, S)$  should minimize  $\mathbf{E}_{y_j|(x_j, S)}(f(x_j, S) - \delta(S))_+^2$ , achieved at

$$f_*(x_j, S) \leq \delta(S).$$

This proves the first claim.

For each  $S$ , denote by  $K$  the set of  $x_j$  such that  $w'(x_j, S) = 0$ . The second claim follows from the following derivation:

$$\begin{aligned} & Q(f) - Q(f_*) \\ &= \mathbf{E}_S(L(f, S) - L(f_*, S)) \\ &= \mathbf{E}_S \left[ \sum_{j=1}^k w(x_j, S)(f(x_j, S) - f_B(x_j, S))^2 \right. \\ &\quad \left. + u \sup_j w'(x_j, S)(f(x_j, S) - \delta(S))_+^2 \right] \\ &\geq \mathbf{E}_S \left[ \sum_{j \in K} (f_B(x_j, S) - f(x_j, S))_+^2 \right. \\ &\quad \left. + u \sup_{j \notin K} (f(x_j, S) - \delta(S))_+^2 \right] \\ &\geq \mathbf{E}_S \left[ \sum_{j \in K} (f_B(x_j, S) - f(x_j, S))_+^2 \right. \\ &\quad \left. + u \sup_{j \notin K} (f(x_j, S) - f'_B(x_j, S))_+^2 \right] \\ &\geq \mathbf{E}_S (DCG(r_B, S) - DCG(r_f, S))^2 C(\gamma, u)^{-2} \\ &\geq (DCG(r_B) - DCG(r_f))^2 C(\gamma, u)^{-2}. \end{aligned}$$

Note that the second inequality follows from  $f'_B(x_j, S) \geq \gamma f_B(x_{j_k}^*, S) \geq \delta(S)$ , and the third inequality follows from Theorem 4.  $\square$

## VII. ASYMPTOTIC ANALYSIS

In this section, we analyze the asymptotic statistical performance of (12). The analysis depends on the underlying func-

tion class  $\mathcal{F}$ . In the literature, one often employs a linear function class with appropriate regularization condition, such as  $L_1$  or  $L_2$  regularization for the linear weight coefficients. Yahoo's machine learning ranking system employs the gradient boosting method described in [13], which is closely related to  $L_1$  regularization, analyzed in [4], [18], [19]. Although the consistency of boosting for the standard least squares regression is known (for example, see [5], [32]), such analysis does not deal with the situation that  $m$  is large and thus is not suitable for analyzing the ranking problem considered in this paper.

In this section, we will consider a linear function class with  $L_2$  regularization, which is closely related to kernel methods. We employ a relatively simple stability analysis which is suitable for  $L_2$  regularization. Our result does not depend on  $m$  explicitly, which is important for large-scale ranking problems such as web search. Although similar results can be obtained for  $L_1$  regularization or gradient boosting, the analysis will become much more complicated.

For  $L_2$  regularization, we consider a feature map  $\psi : \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{H}$ , where  $\mathcal{H}$  is a vector space. We denote by  $w^T v$  the  $L_2$  inner product of  $w$  and  $v$  in  $\mathcal{H}$ . The function class  $\mathcal{F}$  considered here is of the following form:

$$\{\beta^T \psi(x, S); \beta \in \mathcal{H}, \beta^T \beta \leq A^2\} \subset \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{R} \quad (15)$$

where the complexity is controlled by  $L_2$  regularization of the weight vector  $\beta^T \beta \leq A^2$ . We use  $(S_i = \{x_{i,1}, \dots, x_{i,m}\}, \{y_{i,j}\})$  to indicate a sample point indexed by  $i$ . As before, we let  $\{y_{i,j}\} = \{y_{i,1}, \dots, y_{i,m}\}$ . Note that for each sample  $i$ , we do not need to assume that  $y_{i,j}$  are independently generated for different  $j$ . Using (15), the importance weighted regression in (12) becomes the following regularized empirical risk minimization method:

$$\begin{aligned} f_{\hat{\beta}}(x, S) &= \hat{\beta}^T \psi(x, S), \\ \hat{\beta} &= \arg \min_{\beta \in \mathcal{H}} \left[ \frac{1}{n} \sum_{i=1}^n L(\beta, S_i, \{y_{i,j}\}) + \lambda \beta^T \beta \right], \\ L(\beta, S, \{y_j\}) &= \sum_{j=1}^m w(x_j, S)(\beta^T \psi(x_j, S) - y_j)^2 \\ &\quad + u \sup_j w'(x_j, S)(\beta^T \psi(x_j, S) - \delta(S))_+^2. \end{aligned} \quad (16)$$

In this method, we replace the hard regularization in (15) with tuning parameter  $A$  by soft regularization with tuning parameter  $\lambda$ , which is computationally more convenient.

The following result is an expected generalization bound for the  $L_2$ -regularized empirical risk minimization method (16), which uses the stability analysis in [29]. The bound in the theorem compares the performance of a finite sample statistical estimator to that of the optimal Bayes estimator. Such a bound is generally referred to as *oracle inequality* in the literature. The proof is in Appendix.

*Theorem 6:* Let  $M = \sup_{x,S} \|\psi(x, S)\|_2$  and

$$W = \sup_S \left[ \sum_{x_j \in S} w(x_j, S) + u \sup_{x_j \in S} w'(x_j, S) \right].$$

Let  $f_{\hat{\beta}}$  be the estimator defined in (16). Then we have

$$\mathbf{E}_{\{S_i, \{y_{i,j}\}\}_{i=1}^n} Q(f_{\hat{\beta}}) \leq \left(1 + \frac{WM^2}{\sqrt{2}\lambda n}\right)^2 \inf_{\beta \in \mathcal{H}} [Q(f_{\beta}) + \lambda \beta^T \beta].$$

We have paid special attention to the properties of (16). In particular, the quantity  $W$  is usually much smaller than  $m$ , which is large for web search applications. The point we would like to emphasize here is that even though the number  $m$  is large, the estimation complexity is only affected by the top portion of the rank-list. If the estimation of the lowest ranked items is relatively easy (as is generally the case), then the learning complexity does not depend on the majority of items near the bottom of the rank-list.

We can combine Theorems 5 and 6, giving the following bound.

*Theorem 7:* Suppose the conditions in Theorems 5 and 6 hold with  $f_*$  minimizing (14). Letting  $\hat{f} = f_{\hat{\beta}}$ , we get the expression at the bottom of the page.

*Proof:* From Theorem 5, we obtain

$$\begin{aligned} DCG(r_B) - \mathbf{E}_{\{S_i, \{y_{i,j}\}\}_{i=1}^n} DCG(r_{\hat{f}}) \\ \leq C(\gamma, u) \mathbf{E}_{\{S_i, \{y_{i,j}\}\}_{i=1}^n} (Q(\hat{f}) - Q(f_*))^{1/2} \\ \leq C(\gamma, u) [\mathbf{E}_{\{S_i, \{y_{i,j}\}\}_{i=1}^n} Q(f_{\hat{\beta}}) - Q(f_*)]^{1/2}. \end{aligned}$$

The second inequality is a consequence of Jensen's inequality. Now by applying Theorem 6, we obtain the desired bound.  $\square$

The theorem implies that if  $Q(f_*) = \inf_{\beta \in \mathcal{H}} Q(f_{\beta})$ , then as  $n \rightarrow \infty$ , we can let  $\lambda \rightarrow 0$  and  $\lambda n \rightarrow \infty$  so that the second term on the right-hand side vanishes in the large sample limit. Therefore, asymptotically, we can achieve the optimal DCG score. This implies the consistency of regression-based learning methods for the DCG criterion. Moreover, the rate of convergence does not depend on  $m$ , but rather the relatively small quantities  $W$  and  $M$ .

## VIII. CONCLUSION

Ranking problems have many important real world applications. Although various formulations have been investigated in the literature, most theoretical results are concerned with global ranking using the pair-wise AUC criterion. Motivated by applications such as web search, we introduced the subset ranking problem, and focus on the DCG criterion that measures the quality of the top-ranked items.

We derived bounds that relate the optimization of DCG scores to the minimization of convex regression errors. In our analysis, it is essential to weight samples differently according to their importance. These bounds are used to motivate modifications of least squares regression methods that focus on the top portion of the rank-list. In addition to conceptual advantages, these

methods have significant computational advantages over standard regression methods because only a small number of items contribute to the solution. This means that they are computationally efficient to solve. The implementation of these methods can be achieved through appropriate selective sampling procedures. Moreover, we showed that the expected generalization performance of the system does not depend on  $m$ . Instead, it only depends on the estimation quality of the top-ranked items. Again this is important for many practical applications.

Results obtained here are closely related to the theoretical analysis for solving classification methods using convex optimization formulations. Our theoretical results show that the regression approach provides a solid basis for solving the subset ranking problem. The practical value of such methods is also significant. In Yahoo's case, substantial improvement of DCG has been achieved after the deployment of a machine learning based ranking system.

Although the DCG criterion is difficult to optimize directly, it is a natural metric for ranking. The investigation of convex surrogate formulations provides a systematic approach to developing efficient machine learning methods for solving this difficult problem. This paper shows that with appropriate features, importance weighted regression methods can produce the optimal scoring function in the large sample limit. However, regression methods proposed in this paper may not necessarily be optimal algorithms for learning ranking functions. Other methods, such as pair-wise preference learning in Section II, can also be effective. It will be interesting to investigate such alternative formulations using a similar analysis.

## APPENDIX

We shall introduce the following notation: let  $Z_n = \{(S_i, \{y_{i,j}\}) : i = 1, \dots, n\}$ . Let  $\hat{\beta}(Z_n)$  be the solution of (16) and  $\hat{\beta}(Z_{n+1})$  be the solution using training data  $Z_{n+1}$

$$\hat{\beta}(Z_{n+1}) = \arg \min_{\beta \in \mathcal{H}} \left[ \frac{1}{n} \sum_{i=1}^{n+1} L(\beta, S_i, \{y_{i,j}\}) + \lambda \beta^T \beta \right].$$

We have the following stability lemma in [29], which can be stated with our notation as follows.

*Lemma 1:* The following inequality holds:

$$\begin{aligned} \|\hat{\beta}(Z_n) - \hat{\beta}(Z_{n+1})\|_2 \\ \leq \frac{1}{2\lambda n} \left\| \frac{\partial}{\partial \beta} L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\}) \right\|_2 \end{aligned}$$

where  $\frac{\partial}{\partial \beta} L(\beta, S, \{y_j\})$  denotes a subgradient of  $L$  with respect to  $\beta$ .

Note that from simple subgradient algebra in [22], we know that a subgradient of  $\sup_j L_j(\beta)$  for a convex function  $L_j(\beta)$  can be written as  $\sum_j \alpha_j \partial L_j(\beta) / \partial \beta$ , where  $\sum_j \alpha_j \leq 1$  and

$$\mathbf{E}_{\{S_i, \{y_{i,j}\}\}_{i=1}^n} DCG(r_{\hat{f}}) \geq DCG(r_B) - C(\gamma, u) \left[ \left(1 + \frac{WM^2}{\sqrt{2}\lambda n}\right)^2 \inf_{\beta \in \mathcal{H}} (Q(f_{\beta}) + \lambda \beta^T \beta) - Q(f_*) \right]^{1/2}.$$

$$\begin{aligned} & \left\| \frac{\partial}{\partial \beta} L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\}) \right\|_2^2 \\ &= 2 \left\| \sum_{j=1}^m w(x_{n+1,j}, S_{n+1}) (\beta^T \psi(x_{n+1,j}, S_{n+1}) - y_{n+1,j}) \right. \\ & \quad \cdot \psi(x_{n+1,j}, S_{n+1}) + u \sum_{j=1}^m \alpha_j w'(x_{n+1,j}, S_{n+1}) (\beta^T \psi(x_{n+1,j}, S_{n+1}) - \delta(S_{n+1})) + \psi(x_{n+1,j}, S_{n+1}) \left. \right\|_2^2 \\ &\leq 2L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\}) \left( \sum_{j=1}^m (w(x_{n+1,j}, S_{n+1}) + u \alpha_j w'(x_{n+1,j}, S_{n+1})) \|\psi(x_{n+1,j}, S_{n+1})\|_2^2 \right) \\ &\leq 2L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\}) \left( \sum_{j=1}^m w(x_{n+1,j}, S_{n+1}) + u \sup_j w'(x_{n+1,j}, S_{n+1}) \right) M^2 \end{aligned}$$

$\alpha_j \geq 0$ . Therefore, we can find  $\alpha_j \geq 0$  and  $\sum_j \alpha_j \leq 1$  as shown in the expressions at the top of the page, where the first inequality in the derivation is a direct application of Cauchy–Schwartz inequality. Now by applying Lemma 1 with  $\delta\beta = \hat{\beta}(Z_n) - \hat{\beta}(Z_{n+1})$ , we can derive the following chain of inequalities. The first inequality uses the following form of Jensen’s inequality

$$\sum_i \rho_i (a_i + b_i)^2 \leq \left[ \left( \sum_i \rho_i a_i^2 \right)^{1/2} + \left( \sum_i \rho_i b_i^2 \right)^{1/2} \right]^2;$$

the third inequality is due to  $(a+b)^2 \leq (1+s)a^2 + (1+s^{-1})b^2$  (where  $s > 0$ ); the fourth inequality uses Lemma 1.

$$\begin{aligned} & L(\hat{\beta}(Z_n), S_{n+1}, \{y_{n+1,j}\}) \\ &= L(\hat{\beta}(Z_{n+1}) + \delta\beta, S_{n+1}, \{y_{n+1,j}\}) \\ &\leq \left[ L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\})^{1/2} \right. \\ & \quad \left. + \left( \sum_{j=1}^m w(x_{n+1,j}, S_{n+1}) |\delta\beta^T \psi(x_{n+1,j}, S_{n+1})|^2 \right. \right. \\ & \quad \left. \left. + u \sup_j w'(x_j, S) |\delta\beta^T \psi(x_{n+1,j}, S_{n+1})|^2 \right)^{1/2} \right]^2 \\ &\leq \left[ L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\})^{1/2} \right. \\ & \quad \left. + W(S_{n+1})^{1/2} M \|\delta\beta\|_2 \right]^2 \\ &\leq (1+s)L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\}) \\ & \quad + (1+s^{-1})W(S_{n+1}) \|\delta\beta\|_2^2 M^2 \\ &\leq (1+s)L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\}) + (1+s^{-1}) \\ & \quad \cdot W(S_{n+1}) \frac{M^2}{4\lambda^2 n^2} \left\| \frac{\partial}{\partial \beta} L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\}) \right\|_2^2 \\ &\leq L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\}) \\ & \quad \cdot \left[ (1+s) + (1+s^{-1})W(S_{n+1})^2 \frac{M^4}{2\lambda^2 n^2} \right] \end{aligned}$$

where we define

$$W(S) = \sum_{x_j \in S} w(x_j, S) + u \sup_{x_j \in S} w'(x_j, S).$$

By optimizing over  $s$ , we obtain

$$\begin{aligned} & L(\hat{\beta}(Z_n), S_{n+1}, \{y_{n+1,j}\}) \\ &\leq \left( 1 + \frac{W(S_{n+1})M^2}{\sqrt{2}\lambda n} \right)^2 L(\hat{\beta}(Z_{n+1}), S_{n+1}, \{y_{n+1,j}\}). \end{aligned}$$

Now denote by  $Z_{n+1}^{(i)}$  the training data obtained from  $Z_{n+1}$  by removing the  $i$ th datum  $(S_i, \{y_{i,j}\})$ , and let  $\hat{\beta}(Z_{n+1}^{(i)})$  be the solution of (16) with  $Z_n$  replaced by  $Z_{n+1}^{(i)}$ , then we have

$$\begin{aligned} & \sum_{i=1}^{n+1} L(\hat{\beta}(Z_{n+1}^{(i)}), S_i, \{y_{i,j}\}) \\ &\leq \left( 1 + \frac{WM^2}{\sqrt{2}\lambda n} \right)^2 \sum_{i=1}^{n+1} L(\hat{\beta}(Z_{n+1}), S_i, \{y_{i,j}\}) \\ &\leq \left( 1 + \frac{WM^2}{\sqrt{2}\lambda n} \right)^2 \\ & \quad \cdot \inf_{\beta \in \mathcal{H}} \left[ \sum_{i=1}^{n+1} L(\beta, S_i, \{y_{i,j}\}) + \lambda(n+1)\beta^T \beta \right]. \end{aligned}$$

The second inequality is due to the definition of  $\hat{\beta}(Z_{n+1})$ , which minimizes the regularized empirical loss with respect to the augmented training data  $Z_{n+1}$  over  $\beta \in \mathcal{H}$ . Now, in order to obtain the desired bound, we simply take expectation with respect to  $Z_{n+1}$  on both sides.

REFERENCES

- [1] S. Agarwal, T. Graepel, R. Herbrich, S. Har-Peled, and D. Roth, “Generalization bounds for the area under the ROC curve,” *J. Machine Learning Res.*, vol. 6, pp. 393–425, 2005.
- [2] S. Agarwal and D. Roth, “Learnability of bipartite ranking functions,” in *Proc. 18th Annu. Conf. Learning Theory*, Bertinoro, Italy, Jun. 2005, pp. 16–31.
- [3] P. Bartlett, M. Jordan, and J. McAuliffe, “Convexity, classification, and risk bounds,” *J. Amer. Statist. Assoc.*, vol. 101, no. 473, pp. 138–156, 2006.
- [4] G. Blanchard, G. Lugosi, and N. Vayatis, “On the rate of convergence of regularized boosting classifiers,” *J. Machine Learning Res.*, vol. 4, pp. 861–894, 2003.

- [5] P. Bühlmann, "Boosting for high-dimensional linear models," *Ann. Statist.*, vol. 34, pp. 559–583, 2006.
- [6] C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. Hullender, "Learning to rank using gradient descent," in *Proc. Int. Conf. machine Learning (ICML'05)*, Bonn, Germany, Aug. 2005, pp. 89–96.
- [7] A. Caponnetto, "A Note on the Role of Squared Loss in Regression," CBCL, MIT, Cambridge, MA, 2005, Tech. Rep..
- [8] S. Clemencon, G. Lugosi, and N. Vayatis, "Ranking and scoring using empirical risk minimization," in *Proc. 18th Annu. Conf. Learning Theory*, Bartinoro, Italy, Jun. 2005, pp. 1–15.
- [9] W. W. Cohen, R. E. Schapire, and Y. Singer, "Learning to order things," *JAIR*, vol. 10, pp. 243–270, 1999.
- [10] D. Cossock, "Method and Apparatus for Machine Learning a Document Relevance Function," US Patent 7197497, 2007.
- [11] P. Crescenzi and V. Kann, "A Compendium of NP Optimization Problems," Dipartimento di Scienze dell'Informazione, Università di Roma La Sapienza, Rome, Italy, 1995 [Online]. Available: <http://www.nada.kth.se/~viggo/wwwcompendium>
- [12] Y. Freund, R. Iyer, R. E. Schapire, and Y. Singer, "An efficient boosting algorithm for combining preferences," *J. Machine Learning Res.*, vol. 4, pp. 933–969, 2003.
- [13] J. Friedman, "Greedy function approximation: A gradient boosting machine," *Ann. Statist.*, vol. 29, pp. 1189–1232, 2001.
- [14] J. Hanley and B. McNeil, "The meaning and use of the area under a receiver operating characteristic (ROC) curve," *Radiology*, pp. 29–36, 1982.
- [15] R. Herbrich, T. Graepel, and K. Obermayer, "Large margin rank boundaries for ordinal regression," in *Advances in Large Margin Classifiers*, B. S. A. Smola, P. Bartlett, and D. Schuurmans, Eds. Cambridge, MA: MIT Press, 2000, pp. 115–132.
- [16] K. Jarvelin and J. Kekalainen, "IR evaluation methods for retrieving highly relevant documents," in *Proc. 23rd. Annu. Int. ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR'00)*, Athens, Greece, Jul. 2000, pp. 41–48.
- [17] T. Joachims, "Optimizing search engines using clickthrough data," in *Proc. ACM Conf. Knowledge Discovery and Data Mining (KDD)*, Edmonton, AB, Canada, Jul 2002, pp. 133–142.
- [18] G. Lugosi and N. Vayatis, "On the Bayes-risk consistency of regularized boosting methods," *Ann. Statist.*, vol. 32, pp. 30–55, 2004.
- [19] S. Mannor, R. Meir, and T. Zhang, "Greedy algorithms for classification—Consistency, convergence rates, and adaptivity," *J. Machine Learning Res.*, vol. 4, pp. 713–741, 2003.
- [20] P. McCullagh and J. A. Nelder, *Generalized Linear Models*. London, U.K.: Chapman & Hall, 1989.
- [21] F. Radlinski and T. Joachims, "Query chains: Learning to rank from implicit feedback," in *Proc. 11th ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, Chicago, IL, Aug. 2005, pp. 239–248.
- [22] R. T. Rockafellar, *Convex Analysis*. Princeton, NJ: Princeton Univ. Press, 1970.
- [23] S. Rosset, "Model selection via the AUC," in *Proc. 21st Int. Conf. Machine Learning*, Banff, AB, Canada, Jul. 2004.
- [24] C. Rudin, "Ranking with a p-norm push," in *Proc. 19th Annu. Conf. Learning Theory*, Pittsburgh, PA, Jun. 2006, pp. 589–604.
- [25] S. Shalev-Shwartz and Y. Singer, "Efficient learning of label ranking by soft projections onto polyhedra," *J. Machine Learning Res.*, vol. 7, pp. 1567–1599, 2006.
- [26] I. Steinwart, "Support vector machines are universally consistent," *J. Complexity*, vol. 18, pp. 768–791, 2002.
- [27] A. Tewari and P. Bartlett, "On the consistency of multiclass classification methods," in *Proc. 18th Annu. Conf. Learning Theory*, Bertinoro, Italy, Jun. 2005, pp. 143–157.
- [28] H. Zha, Z. Zheng, H. Fu, and G. Sun, "Incorporating query difference for learning retrieval functions in world wide web search," in *CIKM '06: Proc. 15th ACM Int. Conf. Information and Knowledge Management*, Arlington, VA, 2006, pp. 307–316, ISBN 1-59593-433-2.
- [29] T. Zhang, "Leave-one-out bounds for kernel methods," *Neural Computation*, vol. 15, pp. 1397–1437, 2003.
- [30] T. Zhang, "Statistical analysis of some multi-category large margin classification methods," *J. Machine Learning Res.*, vol. 5, pp. 1225–1251, 2004.
- [31] T. Zhang, "Statistical behavior and consistency of classification methods based on convex risk minimization," *Ann. Statist.*, vol. 32, pp. 56–85, 2004.
- [32] T. Zhang and B. Yu, "Boosting with early stopping: Convergence and consistency," *Ann. Statist.*, vol. 33, pp. 1538–1579, 2005.