

## THE ASYMPTOTICS OF RANKING ALGORITHMS

BY JOHN C. DUCHI<sup>\*,†</sup>, LESTER MACKEY<sup>\*,†</sup> AND MICHAEL I. JORDAN<sup>†</sup>

*University of California, Berkeley and Stanford University*

We consider the predictive problem of supervised ranking, where the task is to rank sets of candidate items returned in response to queries. Although there exist statistical procedures that come with guarantees of consistency in this setting, these procedures require that individuals provide a complete ranking of all items, which is rarely feasible in practice. Instead, individuals routinely provide partial preference information, such as pairwise comparisons of items, and more practical approaches to ranking have aimed at modeling this partial preference data directly. As we show, however, such an approach raises serious theoretical challenges. Indeed, we demonstrate that many commonly used surrogate losses for pairwise comparison data do not yield consistency; surprisingly, we show inconsistency even in low-noise settings. With these negative results as motivation, we present a new approach to supervised ranking based on aggregation of partial preferences and we develop  $U$ -statistic-based empirical risk minimization procedures. We present an asymptotic analysis of these new procedures, showing that they yield consistency results that parallel those available for classification. We complement our theoretical results with an experiment studying the new procedures in a large-scale web-ranking task.

**1. Introduction.** Recent years have seen significant developments in the theory of classification, most notably binary classification, where strong theoretical results are available that quantify rates of convergence and shed light on qualitative aspects of the problem [44, 3]. Extensions to multi-class classification have also been explored, and connections to the theory of regression are increasingly well understood, so that overall a satisfactory theory of supervised machine learning has begun to emerge [43, 39].

In many real-world problems in which labels or responses are available, however, the problem is not merely to classify or predict a real-valued response, but rather to list a set of items in order. The theory of supervised

---

<sup>\*</sup>Supported by DARPA through the National Defense Science and Engineering Graduate Fellowship Program (NDSEG).

<sup>†</sup>Partially supported by the U.S. Army Research Laboratory and the U.S. Army Research Office under contract/grant number W911NF-11-1-0391.

*AMS 2000 subject classifications:* 62F07, 62F12, 68Q32, 62C99

*Keywords and phrases:* Ranking, consistency, Fisher consistency, asymptotics, rank aggregation,  $U$ -statistics

learning cannot be considered complete until it also provides a treatment of such *ranking* problems. For example, in information retrieval, the goal is to rank a set of documents in order of relevance to a user’s search query; in medicine, the object is often that of ranking drugs in order of probable curative outcomes for a given disease; and in recommendation or advertising systems, the aim is to present a set of products in order of a customer’s willingness to purchase or consume. In each example, the goal is to order a set of items in accordance with the preferences of an individual or population. While such problems are often converted to classification problems for simplicity (for example, a document is classified as “relevant” or not), decision makers frequently require the ranks (for example, a search engine must place documents in a particular ordering on the page). Despite its ubiquity, our statistical understanding of ranking falls short of our understanding of classification and regression. Our aim here is to characterize the statistical behavior of computationally tractable inference procedures for ranking under natural data-generating mechanisms.

We consider a general decision-theoretic formulation of the *supervised ranking problem* in which preference data are drawn i.i.d. from an unknown distribution, and where each datum consists of a *query*,  $Q \in \mathcal{Q}$ , and a *preference judgment*,  $Y \in \mathcal{Y}$ , over a set  $M_Q$  of candidate items that are available based on the query  $Q$ . The exact nature of the query and preference judgment will depend on the ranking context. In the setting of information retrieval, for example, each datum corresponds to a user issuing a natural language query and expressing a preference by selecting or clicking on zero or more of the returned results.

The statistical task is to discover a function that provides a query-specific ordering of items that best respects the observed preferences. This query-indexed setting is especially natural for tasks like information retrieval in which a different ranking of webpages is needed for each natural language query.

Following existing literature, we estimate a *scoring function*  $f : \mathcal{Q} \rightarrow \mathbb{R}^m$ , where  $f(q)$  assigns a score to each of  $m$  candidate items for the query  $q$ , and the results are ranked according to their scores [23, 21]. Throughout the paper, we adopt a decision-theoretic perspective and assume that given a query-judgment pair  $(Q, Y)$ , we evaluate the scoring function  $f$  via a loss  $L(f(Q), Y)$ . The goal is to choose the  $f$  minimizing the risk

$$(1) \quad R(f) := \mathbb{E}[L(f(Q), Y)].$$

While minimizing the risk (1) directly is in general intractable, researchers in machine learning and information retrieval have developed surrogate loss

functions that yield procedures for selecting  $f$ . Unfortunately, as we show, extant procedures fail to solve the ranking problem under reasonable data generating mechanisms. The goal in the remainder of the paper is to explain this failure and to propose a novel solution strategy based on preference aggregation.

Let us begin to elucidate the shortcomings of current approaches to ranking. One main problem lies in their unrealistic assumptions about available data. The losses proposed and most commonly used for evaluation in the information retrieval literature [29, 26] have a common form, generally referred to as (Normalized) Discounted Cumulative Gain ((N)DCG). The NDCG family requires that the preference judgements  $Y$  associated with the datum  $(Q, Y)$  be a vector  $Y \in \mathbb{R}^m$  of *relevance scores* for the entire set of items; that is,  $Y_j$  denotes the real-valued relevance of item  $j$  to the query  $Q$ . While having complete preference information makes it possible to design procedures that asymptotically minimize NDCG losses [e.g., 12], in practice such complete preferences are unrealistic: they are expensive to collect and difficult to trust. In biological applications, evaluating the effects of all drugs involved in a study—or all doses—on a single subject is infeasible. In web search, users click on only one or two results: no feedback is available for most items. Even when practical and ethical considerations do not rule out collecting complete preference information from participants in a study, a long line of psychological work has highlighted the inconsistency with which humans assign numerical values to multiple objects [e.g., 38, 40, 30].

The inherent practical difficulties that arise in using losses based on relevance scores has led other researchers to propose loss functions that are suitable for *partial preference data* [27, 21, 16]. Such data arise naturally in a number of real-world situations; for example, a patient’s prognosis may improve or deteriorate after administration of treatment, competitions and sporting matches provide paired results, and shoppers at a store purchase one item but not others. Moreover, the psychological literature shows that human beings are quite good at performing pairwise distinctions and forming relative judgments [see, e.g., 36, and references therein].

More formally, let  $\alpha := f(Q) \in \mathbb{R}^m$  denote the vector of predicted scores for each item associated with query  $Q$ . If a preference  $Y$  indicates that item  $i$  is preferred to  $j$  then the natural associated loss is the zero-one loss  $L(\alpha, Y) = 1(\alpha_i \leq \alpha_j)$ . Minimizing such a loss is well known to be computationally intractable; nonetheless, the classification literature [43, 44, 3, 39] has shown that it is possible to design convex Fisher-consistent surrogate losses for the 0-1 loss in classification settings and has linked Fisher consistency to consistency. By reduction to classification, similar consistency

results are possible in certain bipartite or binary ranking scenarios [10]. One might therefore hope to make use of these surrogate losses in the ranking setting to obtain similar guarantees. Unfortunately, however, this hope is not borne out; as we illustrate in Section 3, it is generally computationally intractable to minimize any Fisher-consistent loss for ranking, and even in favorable low-noise cases, convex surrogates that yield Fisher consistency for binary classification fail to be Fisher-consistent for ranking.

We find ourselves at an impasse: existing methods based on practical data-collection strategies do not yield a satisfactory theory, and those methods that do have theoretical justification are not practical. Our approach to this difficulty is to take a new approach to supervised ranking problems in which partial preference data are aggregated before being used for estimation. The point of departure for this approach is the notion of *rank aggregation* [e.g., 20], which has a long history in voting [15], social choice theory [11, 2], and statistics [41, 28]. In Section 2 we discuss some of the ways in which partial preference data can be aggregated, and we propose a new family of  $U$ -statistic-based loss functions that are computationally tractable. Sections 3 and 4 present a theoretical analysis of procedures based on these loss functions, establishing their consistency. We provide a further discussion of practical rank aggregation strategies in Section 5 and present experimental results in Section 6. Section 7 contains our conclusions, with proofs deferred to appendices.

**2. Ranking with rank aggregation.** We begin by considering several ways in which partial preference data arise in practice. We then turn to a formal treatment of our aggregation-based strategy for supervised ranking.

1. *Paired comparison data.* Data in which an individual judges one item to be preferred over another in the context of a query are common. Competitions and sporting matches, where each pairwise comparison may be accompanied by a magnitude such as a difference of scores, naturally generate such data. In practice, a single individual will not provide feedback for all possible pairwise comparisons, and we do not assume transitivity among the observed preferences for an individual. Thus it is natural to model the pairwise preference judgment space  $\mathcal{Y}$  as the set of weighted directed graphs on  $m$  nodes.
2. *Selection data.* A ubiquitous source of partial preference information is the selection behavior of a user presented with a small set of potentially ordered items. For example, in response to a search query, a web search engine presents an ordered list of webpages and records the URL a user clicks on, and a store records inventory and tracks the items customers

purchase. Such selections provide partial information: that a user or customer prefers one item to others presented.

3. *Partial orders.* An individual may also provide preference feedback in terms of a partial ordering over a set of candidates or items. In the context of elections, for example, each preference judgment  $Y \in \mathcal{Y}$  specifies a partial order  $\prec_Y$  over candidates such that candidate  $i$  is preferred to candidate  $j$  whenever  $i \prec_Y j$ . A partial order need not specify a preference between every pair of items.

Using these examples as motivation, we wish to develop a formal treatment of ranking based on aggregation. To provide intuition for the framework presented in the remainder of this section, let us consider a simple aggregation strategy appropriate for the case of paired comparison data. Let each relevance judgment  $Y \in \mathcal{Y}$  be a weighted adjacency matrix where the  $(i, j)$ th entry expresses a preference for item  $i$  over  $j$  whenever this entry is non-zero. In this case, a natural aggregation strategy is to average all observed adjacency matrices for a fixed query. Specifically, for a set of adjacency matrices  $\{Y_l\}_{l=1}^k$  representing user preferences for a given query, we form the average  $(1/k) \sum_{l=1}^k Y_l$ . As  $k \rightarrow \infty$ , the average adjacency matrix captures the mean population preferences, and we thereby obtain complete preference information over the  $m$  items.

This averaging of partial preferences is one example of a general class of aggregation strategies that form the basis of our theoretical framework. To formalize this notion, we modify the loss formulation slightly and hereafter assume that the loss function  $L$  is a mapping  $\mathbb{R}^m \times \mathcal{S} \rightarrow \mathbb{R}$ , where  $\mathcal{S}$  is a problem-specific *structure space*. We further assume the existence of a series of *structure functions*,  $s_k : \mathcal{Y}^k \rightarrow \mathcal{S}$ , that map sets of preference judgments  $\{Y_j\}$  into  $\mathcal{S}$ . The loss  $L$  depends on the preference feedback  $(Y_1, \dots, Y_k)$  for a given query only via the structure  $s_k(Y_1, \dots, Y_k)$ . In the example of the previous paragraph,  $\mathcal{S}$  is the set of  $m \times m$  adjacency matrices, and  $s_k(Y_1, \dots, Y_k) = (1/k) \sum_{l=1}^k Y_l$ . A typical loss for this setting is the pairwise loss [21, 27]

$$L(\alpha, s(Y_1, \dots, Y_k)) = L(\alpha, A) = \sum_{i < j} A_{ij} 1(\alpha_i \leq \alpha_j) + \sum_{i > j} A_{ij} 1(\alpha_i < \alpha_j),$$

where  $\alpha$  is a set of scores, and  $A = s_k(Y_1, \dots, Y_k)$  is the average adjacency matrix with entries  $A_{ij}$ . In Section 5, we provide other examples of structure functions for different data collection mechanisms and losses. Hereafter, we abbreviate  $s_k(Y_1, \dots, Y_k)$  as  $s(Y_1, \dots, Y_k)$  whenever the input length  $k$  is clear from context.

To meaningfully characterize the asymptotics of inference procedures, we make a mild assumption on the limiting behavior of the structure functions.

ASSUMPTION A. *Fix a query  $Q = q$ . Let the sequence  $Y_1, Y_2 \dots$  be drawn i.i.d. conditional on  $q$  and define the random variables  $S_k := s(Y_1, \dots, Y_k)$ . If  $\mu_q^k$  denotes the distribution of  $S_k$ , there exists a limiting law  $\mu_q$  such that*

$$\mu_q^k \xrightarrow{d} \mu_q \quad \text{as } k \rightarrow \infty.$$

For example, the averaging structure function satisfies Assumption A so long as  $\mathbb{E}[|Y_{ij}| \mid Q] < \infty$  with probability 1. Aside from the requirements of Assumption A, we allow arbitrary aggregation within the structure function.

In addition, our main assumption on the loss function  $L$  is as follows:

ASSUMPTION B. *The loss function  $L : \mathbb{R}^m \times \mathcal{S} \rightarrow \mathbb{R}$  is bounded in  $[0, 1]$ , and, for any fixed vector  $\alpha \in \mathbb{R}^m$ ,  $L(\alpha, \cdot)$  is continuous in the topology of  $\mathcal{S}$ .*

Having stated assumptions on the asymptotics of the structure function  $s$  and the loss  $L$ , we can turn to a discussion of the risk functions that guide our design of inference procedures. We begin with the pointwise conditional risk, which maps predicted scores and a measure  $\mu$  on  $\mathcal{S}$  to  $[0, 1]$ :

$$(2) \quad \ell : \mathbb{R}^m \times \mathcal{M}(\mathcal{S}) \rightarrow [0, 1] \quad \text{where} \quad \ell(\alpha, \mu) = \int L(\alpha, s) d\mu(s).$$

Here  $\mathcal{M}(\mathcal{S})$  denotes the closure of the subset of probability measures on the set  $\mathcal{S}$  for which  $\ell$  is defined. For any query  $q$  and  $\alpha \in \mathbb{R}^m$ , we have  $\lim_i \ell(\alpha, \mu_q^i) = \ell(\alpha, \mu_q)$  by the definition of convergence in distribution. This convergence motivates our decision-theoretic approach.

Our goal in ranking is thus to minimize the risk

$$(3) \quad R(f) := \sum_q p_q \ell(f(q), \mu_q),$$

where  $p_q$  denotes the probability that the query  $Q = q$  is issued. The risk of the scoring function  $f$  can also be obtained in the limit as the number of preference judgments for each query goes to infinity:

$$(4) \quad R(f) = \lim_k \mathbb{E}[L(f(Q), s(Y_1, \dots, Y_k))] = \lim_k \sum_q p_q \ell(f(q), \mu_q^k).$$

That the limiting expectation (4) is equal to the risk (3) follows from the definition of weak convergence.

We face two main difficulties in the study of the minimization of the risk (3). The first difficulty is that of *Fisher consistency* mentioned previously: since  $L$  may be non-smooth in the function  $f$  and is typically intractable to minimize, when will the minimization of a tractable surrogate lead to the minimization of the loss (3)? We provide a precise formulation of and answer to this question in Section 3. In addition, we demonstrate the inconsistency of many commonly used pairwise ranking surrogates and show that aggregation leads to tractable Fisher consistent inference procedures for both complete and partial data losses.

The second difficulty is that of *consistency*: for a given Fisher consistent surrogate for the risk (3), are there tractable statistical procedures that converge to a minimizer of the risk? Yes: in Section 4, we develop a new family of aggregation losses based on  $U$ -statistics of increasing order, showing that uniform laws of large numbers hold for the resulting M-estimators.

**3. Fisher consistency of surrogate risk minimization.** In this section, we formally define the Fisher consistency of a surrogate loss and give general necessary and sufficient conditions for consistency to hold for losses satisfying Assumption B. To begin, we assume that the space  $\mathcal{Q}$  of queries is countable (or finite) and thus bijective with  $\mathbb{N}$ . Recalling the definition (3) of the risk and the pointwise conditional risk (2), we define the Bayes risk for  $R$  as the minimal risk over all measurable functions  $f : \mathcal{Q} \rightarrow \mathbb{R}^m$ :

$$R^* := \inf_f R(f) = \sum_q p_q \inf_{\alpha \in \mathbb{R}^m} \ell(\alpha, \mu_q).$$

The second equality follows because  $\mathcal{Q}$  is countable and the infimum is taken over all measurable functions.

Since it is infeasible to minimize the risk (3) directly, we consider a bounded-below surrogate  $\varphi$  to minimize in place of  $L$ . For each structure  $s \in \mathcal{S}$ , we write  $\varphi(\cdot, s) : \mathbb{R}^m \rightarrow \mathbb{R}_+$ , and we assume that for  $\alpha \in \mathbb{R}^m$ , the function  $s \mapsto \varphi(\alpha, s)$  is continuous with respect to the topology on  $\mathcal{S}$ . We then define the conditional  $\varphi$ -risk as

$$(5) \quad \ell_\varphi(\alpha, \mu) := \int_{\mathcal{S}} \varphi(\alpha, s) d\mu(s)$$

and the asymptotic  $\varphi$ -risk of the function  $f$  as

$$(6) \quad R_\varphi(f) := \sum_q p_q \ell_\varphi(f(q), \mu_q),$$

whenever each  $\ell_\varphi(f(q), \mu_q)$  exists (otherwise  $R_\varphi(f) = +\infty$ ). The optimal  $\varphi$ -risk is defined to be  $R_\varphi^* := \inf_f R_\varphi(f)$ , and throughout we make the

assumption that there exist measurable  $f$  such that  $R_\varphi(f) < +\infty$  so that  $R_\varphi^*$  is finite. The following is our general notion of consistency.

DEFINITION 1. *The surrogate loss  $\varphi$  is Fisher-consistent for the loss  $L$  if for any  $\{p_q\}$  and probability measures  $\mu_q \in \mathcal{M}(\mathcal{S})$ , the convergence*

$$R_\varphi(f_n) \rightarrow R_\varphi^* \quad \text{implies} \quad R(f_n) \rightarrow R^*.$$

To achieve more actionable risk bounds and to more accurately compare surrogate risks, we also draw upon a uniform statement of consistency:

DEFINITION 2. *The surrogate loss  $\varphi$  is uniformly consistent for the loss  $L$  if for any  $\epsilon > 0$ , there exists a  $\delta(\epsilon) > 0$  such that for any  $\{p_q\}$  and probability measures  $\mu_q \in \mathcal{M}(\mathcal{S})$ ,*

$$(7) \quad R_\varphi(f) < R_\varphi^* + \delta(\epsilon) \quad \text{implies} \quad R(f) < R^* + \epsilon.$$

The bound (7) is equivalent to the assertion that there exists a non-decreasing function  $\zeta$  such that  $\zeta(0) = 0$  and  $R(f) - R^* \leq \zeta(R_\varphi(f) - R_\varphi^*)$ . Bounds of this form have been completely characterized in the case of binary classification [3], and Steinwart [39] has given necessary and sufficient conditions for uniform consistency to hold. We now turn to analyzing conditions under which a surrogate loss  $\varphi$  is consistent for ranking.

3.1. *General theory.* The main approach in establishing conditions for the surrogate risk consistency in Definition 1 is to move from global conditions for consistency to local, pointwise consistency. Following the treatment of Steinwart [39], we begin by defining a function measuring the discriminating ability of the surrogate  $\varphi$ :

$$(8) \quad H(\epsilon) := \inf_{\mu \in \mathcal{M}(\mathcal{S}), \alpha} \left\{ \ell_\varphi(\alpha, \mu) - \inf_{\alpha'} \ell_\varphi(\alpha', \mu) \mid \ell(\alpha, \mu) - \inf_{\alpha'} \ell(\alpha', \mu) \geq \epsilon \right\}.$$

This function is familiar from work on surrogate risk consistency in classification [3] and measures surrogate risk suboptimality as a function of risk suboptimality. A reasonable conditional  $\varphi$ -risk will declare a set of scores  $\alpha \in \mathbb{R}^m$  suboptimal whenever the conditional risk  $\ell$  declares them suboptimal. This corresponds to  $H(\epsilon) > 0$  whenever  $\epsilon > 0$ , and we call any loss satisfying this condition *pointwise consistent*.

From these definitions, we can conclude the following consistency result, which is analogous to the results of [39]. For completeness, we provide a proof in supplementary Appendix 8.1.

PROPOSITION 1. *Let  $\varphi : \mathbb{R}^m \times \mathcal{S} \rightarrow \mathbb{R}_+$  be a bounded-below loss function such that for some  $f$ ,  $R_\varphi(f) < +\infty$ . Then  $\varphi$  is pointwise consistent if and only if the uniform consistency definition (7) holds.*

Proposition 1 makes it clear that pointwise consistency for general measures  $\mu$  on the set of structures  $\mathcal{S}$  is a stronger condition than that of consistency in Definition 1. In some situations, however, it is possible to connect the weaker surrogate risk consistency of Definition 1 with uniform consistency and pointwise consistency. Ranking problems with appropriate choices of the space  $\mathcal{S}$  give rise to such connections. Indeed, consider the following:

ASSUMPTION C. *The space of possible structures  $\mathcal{S}$  is finite, and the loss  $L$  is discrete, meaning that it takes on only finitely many values.*

Binary and multiclass classification provide examples of settings in which Assumption C is appropriate, since the set of structures  $\mathcal{S}$  is the set of class labels, and  $L$  is usually a version of the 0-1 loss. We also sometimes make a weaker version of Assumption C:

ASSUMPTION C'. *The (topological) space of possible structures  $\mathcal{S}$  is compact, and there exists a partition  $\mathcal{A}_1, \dots, \mathcal{A}_d$  of  $\mathbb{R}^m$  such that for any  $s \in \mathcal{S}$ ,*

$$L(\alpha, s) = L(\alpha', s) \quad \text{whenever } \alpha, \alpha' \in \mathcal{A}_i.$$

Assumption C' may be more natural in ranking settings than Assumption C. The compactness assumption holds, for example, if  $\mathcal{S} \subset \mathbb{R}^m$  and is closed and bounded, such as in our pairwise aggregation example in Section 2. Losses  $L$  that depend only on the relative order of the coordinate values of  $\alpha \in \mathbb{R}^m$ —common in ranking problems—provide a collection of examples for which the partitioning condition holds.

Under Assumption C or C', we can provide a definition of local consistency that is often more user-friendly than pointwise consistency (8):

DEFINITION 3. *Let  $\varphi$  be a bounded-below surrogate loss such that  $\varphi(\cdot, s)$  is continuous for all  $s \in \mathcal{S}$ . The function  $\varphi$  is structure-consistent with respect to the loss  $L$  if for all  $\mu \in \mathcal{M}(\mathcal{S})$ ,*

$$\ell_\varphi^*(\mu) := \inf_{\alpha} \ell_\varphi(\alpha, \mu) < \inf_{\alpha'} \left\{ \ell_\varphi(\alpha, \mu) \mid \alpha \notin \underset{\alpha'}{\operatorname{argmin}} \ell(\alpha', \mu) \right\}.$$

Definition 3 describes the set of loss functions  $\varphi$  satisfying the intuitively desirable property that the surrogate  $\varphi$  cannot be minimized if the scores

$\alpha \in \mathbb{R}^m$  are restricted to not minimize the loss  $L$ . As we see presently, Definition 3 captures exactly what it means for a surrogate loss  $\varphi$  to be consistent when one of Assumptions C or C' holds. Moreover, the set of consistent surrogates coincides with the set of uniformly consistent surrogates in this case. The following theorem formally states this result; we give a proof in supplementary Appendix 8.2.

**THEOREM 1.** *Let  $\varphi : \mathbb{R}^m \times \mathcal{S} \rightarrow \mathbb{R}_+$  satisfy  $R_\varphi(f) < +\infty$  for some measurable  $f$ . If Assumption C holds, then*

- (a) *If  $\varphi$  is structure consistent (Definition 3), then  $\varphi$  is uniformly consistent for the loss  $L$  (Definition 2).*
- (b) *If  $\varphi$  is Fisher-consistent for the loss  $L$  (Definition 1), then  $\varphi$  is structure consistent.*

*If the function  $\varphi(\cdot, s)$  is convex for  $s \in \mathcal{S}$ , and for  $\mu \in \mathcal{M}(\mathcal{S})$  the conditional risk  $\ell_\varphi(\alpha, \mu) \rightarrow \infty$  as  $\|\alpha\| \rightarrow \infty$ , then Assumption C' implies (a) and (b).*

Theorem 1 shows that as long as Assumption C holds, pointwise consistency, structure consistency, and both uniform and non-uniform surrogate loss consistency coincide. These four also coincide under the weaker Assumption C' so long as the surrogate is 0-coercive, which is not restrictive in practice. As a final note, we recall a result due to Steinwart [39], which gives general necessary and sufficient conditions for the consistency in Definition 1 to hold. We begin by giving a weaker version of the suboptimality function (8) that depends on  $\mu$ :

$$(9) \quad H(\epsilon, \mu) := \inf_{\alpha} \left\{ \ell_\varphi(\alpha, \mu) - \inf_{\alpha'} \ell_\varphi(\alpha', \mu) \mid \ell(\alpha, \mu) - \inf_{\alpha'} \ell(\alpha', \mu) \geq \epsilon \right\}.$$

**PROPOSITION 2** (Steinwart [39], Theorems 2.8 and 3.3). *The suboptimality function (9) satisfies  $H(\epsilon, \mu_q) > 0$  for any  $\epsilon > 0$  and  $\mu_q$  with  $q \in \mathcal{Q}$  and  $p_q > 0$  if and only if  $\varphi$  is Fisher-consistent for the loss  $L$  (Definition 1).*

We remark that as a corollary of this result, any structure-consistent surrogate loss  $\varphi$  (in the sense of Definition 3) is consistent for the loss  $L$  whenever the conditional risk  $\ell(\alpha, \mu)$  has finite range, so that  $\alpha \notin \operatorname{argmin}_{\alpha'} \ell(\alpha', \mu) \neq \emptyset$  implies the existence of an  $\epsilon > 0$  such that  $\ell(\alpha, \mu) - \inf_{\alpha'} \ell(\alpha, \mu) \geq \epsilon$ .

**3.2. The difficulty of consistency for ranking.** We now turn to the question of whether there exist structure-consistent ranking losses. In a preliminary version of this work [18], we focused on the practical setting of

learning from pairwise preference data and demonstrated that many popular ranking surrogates are inconsistent for standard pairwise ranking losses. We review and generalize our main inconsistency results here, noting that while the losses considered use pairwise preferences, they perform no aggregation. Their theoretically poor performance provides motivation for the aggregation strategies proposed in this work; we explore the connections in Section 5 (focusing on pairwise losses in Section 5.3). We provide proofs of our inconsistency results in supplementary Appendix 10.

To place ourselves in the general structural setting of the paper, we consider the structure function  $s(Y_1, \dots, Y_k) = Y_1$  which performs no aggregation for all  $k$ , and we let  $Y$  denote the weighted adjacency matrix of a directed acyclic graph (DAG)  $G$ , so that  $Y_{ij}$  is the weight of the directed edge ( $i \rightarrow j$ ) in the graph  $G$ . We consider a pairwise loss that imposes a separate penalty for each misordered pair of results:

$$(10) \quad L(\alpha, Y) = \sum_{i < j} Y_{ij} 1(\alpha_i \leq \alpha_j) + \sum_{i > j} Y_{ij} 1(\alpha_i < \alpha_j),$$

where the cases  $i < j$  and  $i > j$  are distinguished to avoid doubly penalizing  $1(\alpha_i = \alpha_j)$ . When pairwise preference judgments are available, use of such losses is common. Indeed, this loss generalizes the disagreement error described by Dekel et al. [16] and is similar to losses used by Joachims [27]. If we define  $Y_{ij}^\mu := \int Y_{ij} d\mu(Y)$ , then

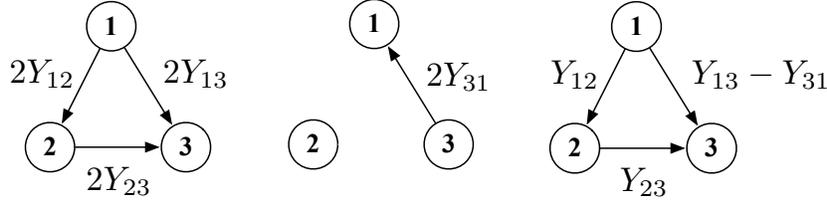
$$(11) \quad \ell(\alpha, \mu) = \sum_{i < j} Y_{ij}^\mu 1(\alpha_i \leq \alpha_j) + \sum_{i > j} Y_{ij}^\mu 1(\alpha_i < \alpha_j).$$

We assume that the number of nodes in any graph  $G$  (or, equivalently, the number of results returned by any query) is bounded by a finite constant  $M$ . Hence, the conditional risk (11) has a finite range; if there are a finite number of preference labels  $Y$  or the set of weights is compact, Assumptions C or C' are satisfied, whence Theorem 1 applies.

**3.2.1. General inconsistency.** Let the set  $P$  denote the complexity class of problems solvable in polynomial time and  $NP$  denote the class of non-deterministic polynomial time problems [see, e.g., 25]. Our first inconsistency result (see also [18, Lemma 7]) is that unless  $P = NP$  (a widely doubted proposition), any loss that is tractable to minimize cannot be a consistent surrogate for the loss (10) and its associated risk.

**PROPOSITION 3.** *Finding an  $\alpha$  minimizing  $\ell$  is NP-hard.*

In particular, most convex functions are minimizable to an accuracy of  $\epsilon$  in time polynomial in the dimension of the problem times a multiple of  $\log \frac{1}{\epsilon}$ ,



**Fig 1.** The two leftmost DAGs occur with probability  $\frac{1}{2}$ , yielding the difference graph  $G_\mu$  at right, assuming  $Y_{23} > Y_{32}$ .

known as poly-logarithmic time [4]. Since any  $\alpha$  minimizing  $\ell_\varphi(\alpha, \mu)$  must minimize  $\ell(\alpha, \mu)$  for a consistent surrogate  $\varphi$ , and  $\ell(\cdot, \mu)$  has a finite range (so that optimizing  $\ell_\varphi$  to a fixed  $\epsilon$  accuracy is sufficient), convex surrogate losses are inconsistent for the pairwise loss (10) unless  $P = NP$ .

**3.2.2. Low-noise inconsistency.** We now turn to showing that, surprisingly, many common convex surrogates are inconsistent even in low-noise settings in which it is easy to find an  $\alpha$  minimizing  $\ell(\alpha, \mu)$ . (Weaker versions of the results in this section appeared in our preliminary paper [18].) Inspecting the loss definition (10), a natural choice for a surrogate loss is one of the form [23, 21, 16]

$$(12) \quad \varphi(\alpha, Y) = \sum_{i,j} h(Y_{ij}) \phi(\alpha_i - \alpha_j),$$

where  $\phi \geq 0$  is a convex function, and  $h$  is a some function of the penalties  $Y_{ij}$ . This surrogate implicitly uses the structure function  $s(Y_1, \dots, Y_k) = Y_1$  and performs no preference aggregation. The conditional surrogate risk is thus  $\ell_\varphi(\alpha, \mu) = \sum_{i \neq j} h_{ij} \phi(\alpha_i - \alpha_j)$ , where  $h_{ij} := \int h(Y_{ij}) d\mu(Y)$ . Surrogates of the form (12) are convenient in margin-based binary classification, where the complete description by Bartlett, Jordan, and McAuliffe [3] shows  $\phi$  is Fisher-consistent if and only if it is differentiable at 0 with  $\phi'(0) < 0$ .

We now precisely define our low-noise setting. For any measure  $\mu$  on a space  $\mathcal{Y}$  of adjacency matrices, let the directed graph  $G_\mu$  be the *difference graph*, that is, the graph with edge weights  $\max\{Y_{ij}^\mu - Y_{ji}^\mu, 0\}$  on edges  $(i \rightarrow j)$ , where  $Y_{ij}^\mu = \int Y_{ij} d\mu(Y)$ . Then we say that the edge  $(i \rightarrow j) \notin G_\mu$  if  $Y_{ij}^\mu \leq Y_{ji}^\mu$  (see Figure 1). We define the following low-noise condition based on self-reinforcement of edges in the difference graph.

**DEFINITION 4.** *The measure  $\mu$  on a set  $\mathcal{Y}$  of adjacency matrices is low-noise when the corresponding difference graph  $G_\mu$  satisfies the following reverse triangle inequality: whenever there is an edge  $(i \rightarrow j)$  and an edge*

( $j \rightarrow k$ ) in  $G_\mu$ , then the weight  $Y_{ik}^\mu - Y_{ki}^\mu$  on the edge ( $i \rightarrow k$ ) is greater than or equal to the path weight  $Y_{ij}^\mu - Y_{ji}^\mu + Y_{jk}^\mu - Y_{kj}^\mu$  on the path ( $i \rightarrow j \rightarrow k$ ).

If  $\mu$  satisfies Definition 4, its difference graph  $G_\mu$  is a DAG. Indeed, the definition ensures that all global preference information in  $G_\mu$  (the sum of weights along any path) conforms with and reinforces local preference information (the weight on a single edge). Hence, we would expect any reasonable ranking method to be consistent in this setting. Nevertheless, typical pair-wise surrogate losses are inconsistent in this low-noise setting (see also the weaker Theorem 11 in our preliminary work [18]):

**THEOREM 2.** *Let  $\varphi$  be a loss of the form (12) and assume  $h(0) = 0$ . If  $\phi$  is convex, then even in the low-noise setting of Definition 4 the loss  $\varphi$  is not structure-consistent.*

Given the difficulties we encounter using losses of the form (12), it is reasonable to consider a reformulation of the surrogate. A natural alternative is a margin-based loss, which encodes a desire to separate ranking scores by large margins dependent on the preferences in a graph. Similar losses have been proposed, e.g., by Shashua and Levin [37]. The next result, shows that convex margin-based losses are also inconsistent, even in low-noise settings. (See also the weaker Theorem 12 of our preliminary work [18].)

**THEOREM 3.** *Let  $\varphi$  be a loss of the form*

$$(13) \quad \varphi(\alpha, Y) = \sum_{i,j:Y_{ij}>0} \phi(\alpha_i - \alpha_j - h(Y_{ij})),$$

where  $h : \mathbb{R} \rightarrow \mathbb{R}$ . If  $\phi$  is convex, then even in the low-noise setting of Definition 4 the loss  $\varphi$  is not structure-consistent.

**3.3. Achieving consistency.** Although Section 3.2 suggests an inherent difficulty in the development of tractable losses for ranking, tractable consistency is in fact achievable if one has access to *complete* preference data. We review a few of the known results here, showing how they follow from the consistency guarantees in Section 3.1, deriving some new consistency guarantees for the complete data setting (we defer all proofs to supplementary Appendix 8). As we have argued, these results are of limited practical value per se, since complete preference judgements are typically unavailable or untrustworthy, but, as we show in Sections 4 and 5, they can be combined with aggregation strategies to yield procedures that are both practical and come with consistency guarantees.

We first define the Normalized Discounted Cumulative Gain (NDCG) family of complete data losses. Such losses are common in applications like web search, since they penalize ranking errors at the top of a ranked list more heavily than errors farther down the list. Let  $s \in \mathcal{S} \subseteq \mathbb{R}^m$  be a vector of relevance scores and  $\alpha \in \mathbb{R}^m$  be a vector of predicted scores. Define  $\pi_\alpha$  to be the permutation associated with  $\alpha$ , so that  $\pi_\alpha(j)$  is the rank of item  $j$  in the ordering induced by  $\alpha$ . Following Ravikumar et al. [34], a general class of NDCG loss functions can be defined as follows:

$$(14) \quad L(\alpha, s) = 1 - \frac{1}{Z(s)} \sum_{j=1}^m \frac{G(s_j)}{F(\pi_\alpha(j))}, \quad Z(s) = \max_{\alpha'} \sum_{j=1}^m \frac{G(s_j)}{F(\pi_{\alpha'}(j))},$$

where  $G$  and  $F$  are functions monotonically increasing in their arguments. By inspection,  $L \in [0, 1]$ , and we remark that the standard NDCG criterion [26] uses  $G(s_j) = 2^{s_j} - 1$  and  $F(j) = \log(1 + j)$ . The ‘‘precision at  $k$ ’’ loss [29] can also be written in the form (14), where  $G(s_j) = s_j$  (assuming that  $s_j \geq 0$ ) and  $F(j) = 1$  for  $j \leq k$  and  $F(j) = +\infty$  otherwise, which measures the relevance of the top  $k$  items given by the vector  $\alpha$ . This form generalizes standard forms of precision, which assume  $s_j \in \{0, 1\}$ .

To analyze the consistency of surrogate losses for the NDCG family (14), we begin by computing the loss  $\ell(\alpha, \mu)$ , then state a corollary to Proposition 2. First, we observe that for any  $\mu \in \mathcal{M}(\mathcal{S})$ ,

$$\ell(\alpha, \mu) = 1 - \sum_{j=1}^m \frac{1}{F(\pi_\alpha(j))} \int \frac{G(s_j)}{Z(s)} d\mu(s).$$

Since the function  $F$  is increasing in its argument, minimizing  $\ell(\alpha, \mu)$  corresponds to choosing any vector  $\alpha$  whose values  $\alpha_j$  obey the same order as the  $m$  points  $\int G(s_j)/Z(s) d\mu(s)$ . In particular, the range of  $\ell$  is finite for any  $\mu$  since it depends only on the permutation induced by  $\alpha$ , so we have

COROLLARY 1. *Define the set*

$$(15) \quad A(\mu) = \left\{ \alpha \in \mathbb{R}^m \mid \alpha_j > \alpha_l \text{ when } \int \frac{G(s_j)}{Z(s)} d\mu(s) > \int \frac{G(s_l)}{Z(s)} d\mu(s) \right\}.$$

*A surrogate loss  $\varphi$  is Fisher-consistent for the NDCG family (14) if and only if for all  $\mu \in \mathcal{M}(\mathcal{S})$ ,*

$$\inf_{\alpha} \left\{ \ell_{\varphi}(\alpha, \mu) - \inf_{\alpha'} \ell_{\varphi}(\alpha', \mu) \mid \alpha \notin A(\mu) \right\} > 0.$$

Corollary 1 recovers the main flavor of the consistency results in the papers of Ravikumar et al. [34] and Buffoni et al. [6]. The surrogate  $\varphi$  is consistent if

and only if it preserves the order of the integrated terms  $\int G(s_j)/Z(s)d\mu(s)$ , that is, any sequence  $\alpha_n$  tending toward the infimum of  $\ell_\varphi(\alpha, \mu)$  must satisfy  $\alpha_n \in A(\mu)$  for large enough  $n$ . Zhang [43] presents several examples of such losses; as a corollary to his Theorem 5 [also noted by 6], the loss

$$\varphi(\alpha, s) := \sum_{j=1}^m \frac{G(s_j)}{Z(s)} \sum_{l=1}^m \phi(\alpha_l - \alpha_j)$$

is convex and structure-consistent (in the sense of Definition 3) whenever  $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$  is non-increasing, differentiable, and satisfies  $\phi'(0) < 0$ . The papers [34, 6] contain more examples and a deeper study of NDCG losses. To extend Corollary 1 to a uniform result, we note that if  $G(s_j) > 0$  for all  $j$  and  $\mathcal{S}$  is compact, then  $\varphi$  is 0-coercive, whence Theorem 1 implies that structure consistency coincides with uniform consistency.

Another family of loss functions is based on a cascade model of user behavior [8]. These losses model dependency among items or results by assuming that a user scans an ordered list of results from top to bottom and selects the first satisfactory result, where satisfaction is determined independently at each position. The form of such expected reciprocal rank (ERR) losses is

$$(16) \quad L(\alpha, s) = 1 - \sum_{i=1}^m \frac{1}{F(i)} G(s_{\pi_\alpha(i)}) \prod_{j=1}^{i-1} (1 - G(s_{\pi_\alpha(j)})),$$

where  $G : \mathbb{R} \rightarrow [0, 1]$  is a non-decreasing function that indicates the prior probability that a result with score  $s_j$  is selected, and  $F : \mathbb{N} \rightarrow [1, \infty)$  is an increasing function that more heavily weights the first items. The ERR family also satisfies  $L \in [0, 1]$ , and empirically correlates well with user satisfaction in ranking tasks [8].

Computing the expected conditional risk  $\ell(\alpha, \mu)$  for general  $\mu \in \mathcal{M}(\mathcal{S})$  is difficult, but we can compute it when  $\mu$  is a product measure over  $s_1, \dots, s_m$ . Indeed, in this case, we have

$$\begin{aligned} \ell(\alpha, \mu) &= 1 - \sum_{i=1}^m \frac{1}{F(i)} \int G(s_{\pi_\alpha(i)}) \prod_{j=1}^{i-1} (1 - G(s_{\pi_\alpha(j)})) d\mu(s) \\ &= 1 - \sum_{i=1}^m \frac{1}{F(i)} \mathbb{E}_\mu[G(s_{\pi_\alpha(i)})] \prod_{j=1}^{i-1} (1 - \mathbb{E}_\mu[G(s_{\pi_\alpha(j)})]). \end{aligned}$$

When one believes that the values  $G(s_i)$  represent the a priori relevance of the result  $i$ , this independence assumption is not unreasonable, and indeed, in Section 5 we provide examples in which it holds. Regardless, we see that  $\ell(\alpha, \mu)$  depends only on the permutation  $\pi_\alpha$ , and we can compute the min-

imizers of the conditional risk for the ERR family (16) using the following lemma, whose proof we provide in supplementary Appendix 8.3.

LEMMA 1. *Let  $p_i = \mathbb{E}_\mu[G(s_i)]$ . The permutation  $\pi$  minimizing  $\ell(\alpha, \mu)$  is in decreasing order of the  $p_i$ .*

Lemma 1 shows that an order-preserving property is necessary and sufficient for the consistency of a surrogate  $\varphi$  for the ERR family (16), as it was for the NDCG family (14). To see this, we apply a variant of Corollary 1 where  $A(\mu)$  as defined in Eq. (15) is replaced with the set

$$A(\mu) = \left\{ \alpha \in \mathbb{R}^m \mid \alpha_j > \alpha_l \text{ whenever } \int G(s_j) d\mu(s) > \int G(s_l) d\mu(s) \right\}.$$

Theorem 5 of [43] implies that  $\varphi(\alpha, s) = \sum_{j=1}^m G(s_j) \sum_{l=1}^m \phi(\alpha_l - \alpha_j)$  is a consistent surrogate when  $\phi$  is convex, differentiable, and non-increasing with  $\phi'(0) < 0$ . Theorem 1 also yields an equivalence between structure and uniform consistency under suitable conditions on  $\mathcal{S}$ .

Before concluding this section, we make a final remark, which has bearing on the aggregation strategies we discuss in Section 5. We have assumed that the structure spaces  $\mathcal{S}$  for the NDCG (14) and ERR (16) loss families consist of real-valued relevance scores. This is certainly not necessary. In some situations, it may be more beneficial to think of  $s \in \mathcal{S}$  as simply an ordered list of the results or as a directed acyclic graph over  $\{1, \dots, m\}$ . We can then apply a transformation  $r : \mathcal{S} \rightarrow \mathbb{R}^m$  to get relevance scores, using those in the losses (14) and (16). This has the advantage of causing  $\mathcal{S}$  to be finite, so Theorem 1 applies and there exists a non-decreasing function  $\zeta$  with  $\zeta(0) = 0$  such that for any distribution and any measurable  $f$ ,

$$R(f) - R^* \leq \zeta(R_\varphi(f) - R_\varphi^*).$$

**4. Uniform laws and asymptotic consistency.** In Section 3, we gave examples of losses based on readily available pairwise data but for which Fisher-consistent tractable surrogates do not exist. The existence of Fisher-consistent tractable surrogates for other forms of data, as in Section 3.3, suggests that aggregation of pairwise and partial data into more complete data structures, such as lists or scores, makes the problem easier. However, it is not obvious how to design statistical procedures based on aggregation. In this section, we formally define a class of suitable estimators that permit us to take advantage of the weak convergence of Assumption A and show that uniform laws of large numbers hold for our surrogate losses. This means that we can indeed asymptotically minimize the risk (3) as desired.

Our aim is to develop an empirical analogue of the population surrogate risk (6) that converges uniformly to the population risk under minimal assumptions on the loss  $\varphi$  and structure function  $s$ . Given a dataset  $\{(Q_i, Y_i)\}_{i=1}^n$  with  $(Q_i, Y_i) \in \mathcal{Q} \times \mathcal{Y}$ , we begin by defining, for each query  $q$ , the batch of data belonging to the query,  $\mathcal{B}(q) = \{i \in \{1, \dots, n\} \mid Q_i = q\}$ , and the empirical count of the query,  $\hat{n}_q = |\mathcal{B}(q)|$ . As a first attempt at developing an empirical objective, we might consider an empirical surrogate risk based on complete aggregation over the batch of data belonging to each query:

$$(17) \quad \frac{1}{n} \sum_q \hat{n}_q \varphi \left( f(q), s(\{Y_{i_1}, \dots, Y_{i_{\hat{n}_q}} \mid i_j \in \mathcal{B}(q)\}) \right).$$

While we would expect this risk to converge uniformly when  $\varphi$  is a sufficiently smooth function of its structure argument, the analysis of the complete aggregation risk requires overly detailed knowledge of the surrogate  $\varphi$  and the structure function  $s$ .

To develop a more broadly applicable statistical procedure, we instead consider an empirical surrogate based on  $U$ -statistics. By trading off the nearness of an order- $k$   $U$ -statistic to an i.i.d. sample and the nearness of the limiting structure distribution  $\mu_q$  to a structure  $s(Y_1, \dots, Y_k)$  aggregated over  $k$  draws, we can obtain consistency under mild assumptions on  $\varphi$  and  $s$ . More specifically, for each query  $q$ , we consider the surrogate loss

$$(18) \quad \binom{\hat{n}_q}{k}^{-1} \sum_{\substack{i_1 < \dots < i_k, \\ i_j \in \mathcal{B}(q)}} \varphi(f(q), s(Y_{i_1}, \dots, Y_{i_k})).$$

When  $\hat{n}_q < k$ , we adopt the convention  $\binom{\hat{n}_q}{k} = 1$  and the above sum becomes the single term  $\varphi(f(q), s(Y_{i_1}, \dots, Y_{i_{\hat{n}_q}}))$  with  $\{i_1, \dots, i_{\hat{n}_q}\} = \mathcal{B}(q)$  as in the expression (17). Hence, our  $U$ -statistic loss recovers the complete aggregation loss (17) when  $k = \infty$ .

An alternative formulation to loss (18) might consist of  $\lceil |\mathcal{B}(q)|/k \rceil$  aggregation terms per query, with each query-preference pair appearing in a single term. However, the instability of such a strategy is high: a change in the ordering of the data or a substitution of queries could have a large effect on the final estimator. The  $U$ -statistic (18) grants robustness to such perturbations in the data. Moreover, by choosing the right rate of increase of the aggregation order  $k$  as a function of  $n$ , we obtain consistent procedures for a broad class of surrogates  $\varphi$  and structures  $s$ .

We associate with the surrogate loss (18) a surrogate empirical risk which

weights each query by its empirical probability of appearance:

$$(19) \quad \widehat{R}_{\varphi,n}(f) := \frac{1}{n} \sum_q \widehat{n}_q \binom{\widehat{n}_q}{k}^{-1} \sum_{\substack{i_1 < \dots < i_k, \\ i_j \in \mathcal{B}(q)}} \varphi(f(q), s(Y_{i_1}, \dots, Y_{i_k})).$$

Let  $\mathbb{P}_n$  denote the probability distribution of the queries given that the dataset consists of a total of  $n$  samples. Then by iteration of expectation and Fubini's theorem, the surrogate risk (19) is an unbiased estimate of the population quantity

$$(20) \quad R_{\varphi,n}(f) := \sum_q \left[ \sum_{l=1}^n l \mathbb{P}_n(\widehat{n}_q = l) \mathbb{E}[\varphi(f(Q), s(Y_1, \dots, Y_{l \wedge k})) \mid Q = q] \right].$$

It remains to establish a uniform law of large numbers guaranteeing the convergence of the empirical risk (19) to the target population risk (6). Under suitable conditions such as those of Section 3, this ensures the asymptotic consistency of computationally tractable statistical procedures. Hereafter, we assume that we have a non-decreasing sequence of function classes  $\mathcal{F}_n$ , where any  $f \in \mathcal{F}_n$  is a scoring function for queries, mapping  $f : \mathcal{Q} \rightarrow \mathbb{R}^m$  and giving scores to the (at most  $m$ ) results for each query  $q \in \mathcal{Q}$ . Our goal is to give sufficient conditions for the convergence in probability

$$(21) \quad \sup_{f \in \mathcal{F}_n} \left| \widehat{R}_{\varphi,n}(f) - R_{\varphi}(f) \right| \xrightarrow{p} 0 \quad \text{as } n \rightarrow \infty.$$

While we do not provide fully general conditions under which the convergence (21) occurs, we provide representative, checkable conditions sufficient for convergence. At a high level, to establish (21), we control the uniform difference between the expectations  $R_{\varphi,n}(f)$  and  $R_{\varphi}(f)$  and bound the distance between the empirical risk  $\widehat{R}_{\varphi,n}$  and its expectation  $R_{\varphi,n}$ , via covering number arguments. We now specify assumptions under which our results hold, deferring all proofs to the supplementary Appendix 9.

Without loss of generality, we assume that  $p_q$ , the true probability of seeing the query  $q$ , is non-increasing in the query index  $q$ . First, we describe the tails of the query distribution:

**ASSUMPTION D.** *There exist constants  $\beta > 0$  and  $K_1 > 0$  such that  $p_q \leq K_1 q^{-\beta-1}$  for all  $q$ , that is,  $p_q = \mathcal{O}(q^{-\beta-1})$ .*

Infinite sets of queries  $\mathcal{Q}$  are reasonable, since search engines, for example, receive a large volume of entirely new queries each day. Our arguments also apply when  $\mathcal{Q}$  is finite, in which case we can take  $\beta \uparrow \infty$ .

Our second main assumption concerns the behavior of the surrogate loss  $\varphi$  over the function class  $\mathcal{F}_n$ , which we assume is contained in a normed space with norm  $\|\cdot\|$ .

ASSUMPTION E (Bounded Lipschitz Losses). *The surrogate loss function  $\varphi$  is bounded and Lipschitz continuous over  $\mathcal{F}_n$ : for any  $s \in \mathcal{S}$ , any  $f, f_1, f_2 \in \mathcal{F}_n$ , and any  $q \in \mathcal{Q}$ , there exist constants  $B_n$  and  $L_n < \infty$  such that*

$$0 \leq \varphi(f(q), s) \leq B_n$$

and

$$|\varphi(f_1(q), s) - \varphi(f_2(q), s)| \leq L_n \|f_1 - f_2\|.$$

This assumption is satisfied whenever  $\varphi(\cdot, s)$  is convex and  $\mathcal{F}_n$  is compact (and contained in the interior of the domain of  $\varphi(\cdot, s)$ ) [24]. Our final assumption gives control over the sizes of the function classes  $\mathcal{F}_n$  as measured by their covering numbers. (The  $\epsilon$ -covering number of  $\mathcal{F}$  is the smallest  $N$  such that there are  $f^i, i \leq N$ , such that  $\min_i \|f^i - f\| \leq \epsilon$  for any  $f \in \mathcal{F}$ .)

ASSUMPTION F. *For all  $\epsilon > 0$ ,  $\mathcal{F}_n$  has  $\epsilon$ -covering number  $N(\epsilon, n) < \infty$ .*

With these assumptions in place, we give a few representative conditions that enable us to guarantee uniform convergence (21). Roughly, these conditions control the interaction between the size of the function classes  $\mathcal{F}_n$  and the order  $k$  of aggregation used with  $n$  samples. To that end, we let the aggregation order  $k_n$  grow with  $n$ . In stating the conditions, we make use of the shorthand  $\mathbb{E}_q[\varphi(f(q), s(Y_{1:k}))]$  for  $\mathbb{E}[\varphi(f(Q), s(Y_1, \dots, Y_k)) \mid Q = q]$ .

CONDITION I. There exists a  $\rho > 0$  and constant  $C$  such that for all  $q \in \mathcal{Q}$ ,  $n \in \mathbb{N}$ ,  $k \in \mathbb{N}$ , and  $f \in \mathcal{F}_n$ ,

$$\left| \mathbb{E}_q[\varphi(f(q), s(Y_1, \dots, Y_k))] - \lim_{k'} \mathbb{E}_q[\varphi(f(q), s(Y_1, \dots, Y_{k'}))] \right| \leq C B_n k^{-\rho}.$$

Additionally, the sequences  $B_n$  and  $k_n$  satisfy  $B_n = o(k_n^\rho)$ .

This condition is not unreasonable; when  $\varphi$  and  $s$  are suitably continuous, we expect  $\rho \geq \frac{1}{2}$ . We also consider an alternative covering number condition.

CONDITION I'. The sequences  $\epsilon_n$  and  $k_n$  and an  $\epsilon_n$ -cover  $\mathcal{F}_n^1, \dots, \mathcal{F}_n^{N(\epsilon_n, n)}$  of  $\mathcal{F}_n$  can be chosen such that

$$\max_{i \in [N(\epsilon_n, n)]} \inf_{f \in \mathcal{F}_n^i} \left| R_\varphi(f) - \sum_q p_q \mathbb{E}_q[\varphi(f(q), s(Y_1, \dots, Y_{k_n}))] \right| + 2L_n \epsilon_n \rightarrow 0.$$

Condition **I'** is weaker than Condition **I**, since it does not require uniform convergence over  $q \in \mathcal{Q}$ . If the function class  $\mathcal{F}$  is fixed for all  $n$ , then the weak convergence of  $s(Y_1, \dots, Y_k)$  as in Assumption **A** guarantees Condition **I'**, since  $N(\epsilon, n) = N(\epsilon, n') < \infty$  and we may take  $\epsilon$  arbitrarily small. We require one additional condition, which relates the growth of  $k_n$ ,  $B_n$ , and the function classes  $\mathcal{F}_n$  more directly.

CONDITION II. The sequences  $k_n$  and  $B_n$  satisfies  $k_n B_n^{\frac{1+\beta}{\beta}} = o(n)$ . Additionally, for any fixed  $\epsilon > 0$ , the sequences satisfy

$$k_n B_n \left[ \log N \left( \frac{\epsilon}{4L_n}, n \right) \right]^{\frac{1}{2}} = o(\sqrt{n}).$$

By inspection, Condition **II** is satisfied for any  $k_n = o(\sqrt{n})$  if the function classes  $\mathcal{F}_n$  are fixed for all  $n$ . Similarly, if for all  $k \geq k_0$ ,  $s(Y_1, \dots, Y_k) = s(Y_1, \dots, Y_{k_0})$ , so  $s$  depends only on its first  $k_0$  arguments, Condition **II** holds whenever  $\max\{B_n^{(1+\beta)/\beta}, B_n^2 \log N(\epsilon/4L_n, n)\} = o(n)$ . If the function classes  $\mathcal{F}_n$  consist of linear functionals represented by vectors  $\theta \in \mathbb{R}^{d_n}$  in a ball of some finite radius, then  $\log N(\epsilon, n) \approx d_n \log \epsilon^{-1}$ , which means that Condition **II** roughly requires  $k_n \sqrt{d_n/n} \rightarrow 0$  as  $n \rightarrow \infty$ . Modulo the factor  $k_n$ , this condition is familiar for its necessity for convergence of parametric statistical problems.

The conditions in place, we come to our main result on the convergence of our  $U$ -statistic-based empirical loss minimization procedures.

THEOREM 4. *Assume Condition **I** or **I'** and additionally assume the growth Condition **II**. Under Assumptions **D**, **E**, and **F**, we obtain*

$$\sup_{f \in \mathcal{F}_n} \left| \widehat{R}_{\varphi, n}(f) - R_{\varphi}(f) \right| \xrightarrow{P} 0.$$

We remark in passing that if Condition **II** holds, except that the  $o(\sqrt{n})$  bound is replaced by  $O(n^{-\rho})$  for some  $\rho < \frac{1}{2}$ , the conclusion of Theorem 4 can be strengthened to both convergence almost surely and in expectation.

By inspection, Theorem 4 provides our desired convergence guarantee (21). By combining the Fisher-consistent loss families outlined in Section 3.3 with the consistency guarantees provided by Theorem 4, it is thus possible to design statistical procedures that are both computationally tractable—minimizing only convex risks—and asymptotically consistent.

**5. Rank aggregation strategies.** In this section, we give several examples of practical strategies for aggregating disparate user preferences under our framework. Motivated by the statistical advantages of complete preference data highlighted in Section 3.3, we first present strategies for constructing complete vectors of relevance scores from pairwise preference data. We then discuss a model for the selection or “click” data that arises in web search and information retrieval and show that maximum likelihood estimation under this model allows for consistent ranking. We conclude this section with a brief overview of structured aggregation strategies.

5.1. *Recovering scores from pairwise preferences.* Here we treat partial preference observations as noisy evidence of an underlying complete ranking and attempt to achieve consistency with respect to a complete preference data loss. We consider three methods that take as input pairwise preferences and output a relevance score vector  $s \in \mathbb{R}^m$ . Such procedures fit naturally into our ranking-with-aggregation framework: the results in Section 3.3 and Section 4 show that a Fisher-consistent loss is consistent for the limiting distribution of the scores  $s$  produced by the aggregation procedure. Thus, it is the responsibility of the statistician—the designer of an aggregation procedure—to determine whether the scores accurately reflect the judgments of the population. We present our first example in some detail to show how aggregation of pairwise judgments can lead to consistency in our framework, following with brief descriptions of alternate aggregation strategies. For an introduction to the design of aggregation strategies for pairwise data, see Tsukida and Gupta [42] as well as the book by David [14].

*Thurstone-Mosteller least squares and skew-symmetric scoring.* The first aggregation strategy constructs a relevance score vector  $s$  in two phases. First, it aggregates a sequence of observed preference judgments  $Y_i \in \mathcal{Y}$ , provided in any form, into a skew-symmetric matrix  $A \in \mathbb{R}^{m \times m}$  satisfying  $A = -A^\top$ . Each entry  $A_{ij}$  encodes the extent to which item  $i$  is preferred to item  $j$ . Given such a skew-symmetric matrix, Thurstone and Mosteller [31] recommend deriving a score vector  $s$  such that  $s_i - s_j \approx A_{ij}$ . In practice, one may not observe preference information for every pair of results, so we define a masking matrix  $\Omega \in \{0, 1\}^{m \times m}$  with  $\Omega = \Omega^\top$ ,  $\Omega_{ii} = 1$ , and  $\Omega_{ij} = 1$  if and only if preference information has been observed for the pair  $i \neq j$ . Letting  $\circ$  denote the Hadamard product, a natural objective for selecting scores [e.g., 22] is the least squares objective

$$(22) \quad \underset{x: x^\top \mathbf{1} = 0}{\text{minimize}} \quad \frac{1}{4} \sum_{i,j} \Omega_{ij} (A_{ij} - (x_i - x_j))^2 = \frac{1}{4} \left\| \Omega \circ (A - (\mathbf{1}x^\top - x\mathbf{1}^\top)) \right\|_{\text{Fr}}^2.$$

The gradient of the objective (22) is

$$D_\Omega x - (\Omega \circ A)\mathbb{1} - \Omega x \quad \text{where} \quad D_\Omega := \text{diag}(\Omega\mathbb{1}).$$

Setting  $s = (D_\Omega - \Omega)^\dagger(\Omega \circ A)\mathbb{1}$  yields the solution to the minimization problem (22), since  $D_\Omega - \Omega$  is an unnormalized graph Laplacian matrix [9], and therefore  $\mathbb{1}^\top s = \mathbb{1}^\top (D_\Omega - \Omega)^\dagger(\Omega \circ A)\mathbb{1} = 0$ .

If  $\Omega = \mathbb{1}\mathbb{1}^\top$ , so that all pairwise preferences are observed, then the eigenvalue decomposition of  $D_\Omega - \Omega = mI - \mathbb{1}\mathbb{1}^\top$  can be computed explicitly as  $V\Sigma V^\top$ , where  $V$  is any orthonormal matrix whose first column is  $1/\sqrt{m}$ , and  $\Sigma$  is a diagonal matrix with entries 0 (once) and  $m$  repeated  $m-1$  times. Thus, letting  $x_A$  and  $x_B$  denote solutions to the minimization problem (22) with different skew-symmetric matrices  $A$  and  $B$  and noting that  $A\mathbb{1} \perp \mathbb{1}$  since  $\mathbb{1}^\top A\mathbb{1} = 0$ , we have the Lipschitz continuity of the solutions  $s$  in  $A$ :

$$\|x_A - x_B\|_2^2 = \left\| (mI - \mathbb{1}\mathbb{1}^\top)^\dagger (A - B)\mathbb{1} \right\|_2^2 = \frac{1}{m^2} \|(A - B)\mathbb{1}\|_2^2 \leq \frac{1}{m} \|A - B\|_2^2.$$

Similarly, when  $\Omega$  is fixed, the score structure  $s$  is likewise Lipschitz in  $A$  for any norm  $\|\cdot\|$  on skew-symmetric matrices.

A variety of procedures are available for aggregating pairwise comparison data  $Y_i \in \mathcal{Y}$  into a skew-symmetric matrix  $A$ . One example, the Bradley-Terry-Luce (BTL) model [5], is based upon empirical log-odds ratios. Specifically, assume that  $Y_i \in \mathcal{Y}$  are pairwise comparisons of the form  $j \succ l$ , meaning item  $j$  is preferred to item  $l$ . Then we can set

$$A_{jl} = \log \frac{\widehat{\mathbb{P}}(j \succ l) + c}{\widehat{\mathbb{P}}(j \prec l) + c} \quad \text{for observed pairs } j, l,$$

where  $\widehat{\mathbb{P}}$  denotes the empirical distribution over  $\{Y_1, \dots, Y_k\}$  and  $c > 0$  is a smoothing parameter.

Since the proposed structure  $s$  is a continuous function of the skew-symmetric matrix  $A$ , the limiting distribution  $\mu$  is a point mass whenever  $A$  converges almost surely, as it does in the BTL model. If aggregation is carried out using only a finite number of preferences rather than letting  $k$  approach  $\infty$  with  $n$ , then  $\mu$  converges to a non-degenerate distribution. Theorem 1 grants uniform consistency since the score space  $\mathcal{S}$  is finite.

*Borda count and budgeted aggregation.* The Borda count [15] provides a computationally efficient method for computing scores from election results. In a general election setting, the procedure counts the number of times that a particular item was rated as the best, second best, and so on. Given a skew-symmetric matrix  $A$  representing the outcomes of elections, the Borda count assigns the scores  $s = A\mathbb{1}$ . As above, a skew-symmetric matrix  $A$

can be constructed from input preferences  $\{Y_1, \dots, Y_k\}$ , and the choice of this first-level aggregation can greatly affect the resulting rankings. Ammar and Shah [1] suggest that if one has limited computational budget and only pairwise preference information then one should assign to item  $j$  the score

$$s_j = \frac{1}{m-1} \sum_{l \neq j} \widehat{\mathbb{P}}(j \succ l),$$

which estimates of the probability of winning an election against an opponent chosen uniformly. This is equivalent to the Borda count when we choose  $A_{jl} = \widehat{\mathbb{P}}(j \succ l) - \widehat{\mathbb{P}}(j \prec l)$  as the entries in the skew-symmetric aggregate  $A$ .

*Principal eigenvector method.* Saaty [35] describes the principal eigenvector method, which begins by forming a reciprocal matrix  $A \in \mathbb{R}^{m \times m}$ , with positive entries  $A_{ij} = (A_{ji})^{-1}$ , from pairwise comparison judgments. Here  $A_{ij}$  encodes a multiplicative preference for item  $i$  over item  $j$ ; the idea is that ratios preserve preference strength [35]. To generate  $A$ , one may use, for example, smoothed empirical ratios  $A_{jl} = \frac{\widehat{\mathbb{P}}(j \succ l) + c}{\widehat{\mathbb{P}}(j \prec l) + c}$ . Saaty recommends finding a vector  $s$  so that  $s_i/s_j \approx A_{ij}$ , suggesting using the Perron vector of the matrix, that is, the first eigenvector of  $A$ .

5.2. *Cascade models for selection data.* Cascade models [13, 8] explain the behavior of a user presented with an ordered list of items, for example from a web search. In a cascade model, a user considers results in the presented order and selects the first to satisfy him or her, and the model assumes the result  $l$  satisfies a user with probability  $p_l$ , independently of previous items in the list. It is natural to express a variety of ranking losses, including the expected reciprocal rank (ERR) family (16), as expected disutility under a cascade model, but computation and optimization of these losses require knowledge of the satisfaction probabilities  $p_l$ . When the satisfaction probabilities are unknown, Chapelle et al. [8] recommend plugging in those values  $p_l$  that maximize the likelihood of observed click data. Here we show that risk consistency for the ERR family is simple to characterize when scores are estimated via maximum likelihood.

To this end, fix a query  $q$  and let each affiliated preference judgment  $Y_i$  consist of a triple  $(m_i, \pi_i, c_i)$ , where  $m_i$  is the number of results presented to the user,  $\pi_i$  is the order of the presented results, which maps positions  $\{1, \dots, m_i\}$  to the full result set  $\{1, \dots, m\}$ , and  $c_i \in \{1, \dots, m_i + 1\}$  is the position clicked on by the user ( $m_i + 1$  if the user chooses nothing). The

likelihood  $g$  of an i.i.d. sequence  $\{Y_1, \dots, Y_k\}$  under a cascade model  $p$  is

$$g(p, \{Y_1, \dots, Y_k\}) = \prod_{i=1}^k p^{1(c_i \leq m_i)} \prod_{j=1}^{c_i-1} (1 - p_{\pi_i(j)}),$$

and the maximum likelihood estimator of the satisfaction probabilities has the closed form

$$\hat{p}_l(Y_1, \dots, Y_k) = \frac{\sum_{i=1}^k \mathbf{1}(\pi_i(c_i) = l)}{\sum_{i=1}^k \sum_{j=1}^{c_i} \mathbf{1}(\pi_i(j) = l)}.$$

To incorporate this maximum likelihood aggregation procedure into our framework, we define the structure function  $s$  to be the vector

$$s(Y_1, \dots, Y_k) := \hat{p}(Y_1, \dots, Y_k) \in \mathbb{R}^m$$

of maximum likelihood probabilities, and we take as our loss  $L$  any member of the ERR family (16). The strong law of large numbers implies the a.s. convergence of  $\hat{p}$  to a vector  $p \in [0, 1]^m$ , so that the limiting law  $\mu_q(\{p\}) = 1$ . Since  $\mu_q$  is a product measure over  $[0, 1]^m$ , Lemma 1 implies that any  $\alpha$  inducing the same ordering over results as  $p$  minimizes the conditional ERR risk  $\ell(\alpha, \mu)$ . By application of Theorems 1 (or Proposition 2) and 4, it is possible to asymptotically minimize the Expected Reciprocal Rank by aggregation.

**5.3. Structured aggregation.** Our framework can leverage aggregation procedures [see, e.g., 20] that map input preferences into representations of combinatorial objects. Consider the setting of Sec. 3.2, in which each observed preference judgment  $Y$  is the weighted adjacency matrix of a directed acyclic graph, our loss of interest  $L$  is the edgewise indicator loss (10), and our candidate surrogate losses have the form (18). Theorems 2 and 3 establish that risk consistency is not generally attainable when  $s(Y_1, \dots, Y_k) = Y_1$ . In certain cases, aggregation can recover consistency. Indeed, define

$$s(Y_1, \dots, Y_k) := \frac{1}{k} \sum_{i=1}^k Y_i,$$

the average of the input adjacency matrices. For an i.i.d. sequence  $Y_1, Y_2, \dots$  associated with a given query  $q$ , we have  $s(Y_1, \dots, Y_n) \xrightarrow{a.s.} \mathbb{E}(Y \mid Q = q)$  by the strong law of large numbers, and hence the asymptotic surrogate risk

$$R_\varphi(f) = \sum_q p_q \int \varphi(f(q), s) d\mu_q(s) = \sum_q p_q \varphi(f(q), \mathbb{E}(Y \mid Q = q)).$$

Recalling the conditional pairwise risk (11), we can rewrite the risk as

$$\begin{aligned} R(f) &= \sum_q p_q \left[ \sum_{i < j} Y_{ij}^{\mu_q} 1(f_i(q) \leq f_j(q)) + \sum_{i > j} Y_{ij}^{\mu_q} 1(f_i(q) < f_j(q)) \right] \\ &= \sum_q p_q \sum_{i > j} \mathbb{E}[Y_{ij} | Q = q] + \sum_q p_q \sum_{i < j} \mathbb{E}[Y_{ij} - Y_{ji} | Q = q] 1(f_i(q) \leq f_j(q)). \end{aligned}$$

The discussion immediately following Proposition 2 shows that any consistent surrogate  $\varphi$  must be bounded away from its minimum for  $\alpha \notin \operatorname{argmin}_{\alpha'} \ell_\varphi(\alpha', \mu)$ . Since the limiting distribution  $\mu$  is a point mass at some adjacency matrix  $s$  for each  $q$ , a surrogate loss  $\varphi$  is consistent if and only if

$$\inf_{\alpha} \left\{ \varphi(\alpha, s) - \inf_{\alpha'} \varphi(\alpha', s) \mid \alpha \notin \operatorname{argmin}_{\alpha'} L(\alpha', s) \right\} > 0.$$

In the important special case when the difference graph  $G_\mu$  associated with  $\mathbb{E}[Y | Q = q]$  is a DAG for each query  $q$  (recall Section 3.2.2), structure consistency is obtained if for each  $\alpha^* \in \operatorname{argmin}_{\alpha} \varphi(\alpha, s)$ ,  $\operatorname{sign}(\alpha_i^* - \alpha_j^*) = \operatorname{sign}(s_{ij} - s_{ji})$  for each pair of results  $i, j$ . As an example, in this setting

$$(23) \quad \varphi(\alpha, s) := \sum_{i,j} [s_{ij} - s_{ji}]_+ \phi(\alpha_i - \alpha_j)$$

is consistent when  $\phi$  is non-increasing, convex, and has derivative  $\phi'(0) < 0$ .

The Fisher-consistent loss (23) is similar to the inconsistent losses (12) considered in Section 3.2, but the coefficients adjoining each  $\phi(\alpha_i - \alpha_j)$  summand exhibit a key difference. While the inconsistent losses employ coefficients based solely on the average  $i \rightarrow j$  weight  $s_{ij}$ , the consistent loss coefficients are nonlinear functions of the edge weight differences  $s_{ij} - s_{ji}$ : they are precisely the edge weights of the difference graph  $G_\mu$  introduced Section 3.2.2. Since at least one of the two coefficients  $[s_{ij} - s_{ji}]_+$  and  $[s_{ji} - s_{ij}]_+$  is always zero, the loss (23) penalizes misordering either edge  $i \rightarrow j$  or  $j \rightarrow i$ . This contrasts with the inconsistent surrogates of Section 3.2, which simultaneously associate non-zero convex losses with opposing edges  $i \rightarrow j$  and  $j \rightarrow i$ . Note also that our argument for the consistency of the loss (23) does not require Definition 4's low-noise assumption: consistency holds under the weaker condition that, on average, a population's preferences are acyclic.

**6. Experimental Study and Implementation.** In this section, we describe strategies for solving the convex programs that emerge from our aggregation approach to ranking and demonstrate the empirical utility of our proposed procedures. We begin with a broad description of implementation strategies, which we follow with our specific experiments.

6.1. *Minimizing the empirical risk.* At first glance, the empirical risk (19) appears difficult to minimize, since the number of terms grows exponentially in the level of aggregation  $k$ . Fortunately, we may leverage techniques from the stochastic optimization literature [32, 17] to minimize the risk (19) in time linear in  $k$  and independent of  $n$ . Let us consider minimizing a function of the form

$$(24) \quad \widehat{R}_N(f) := \frac{1}{N} \sum_{i=1}^N \varphi(f, s^i) + \Phi(f),$$

where  $\{s^i\}_{i=1}^N$  is some collection of data,  $\varphi(\cdot, s)$  is convex in its first argument, and  $\Phi$  is a convex regularizing function (possibly zero).

Duchi and Singer [17], using ideas similar to those of Nemirovski et al. [32], develop a specialized stochastic gradient descent method for minimizing composite objectives of the form (24) (providing a unified and sharper analysis in the subsequent paper [19]). Such methods maintain a parameter  $f^t$ , which is assumed to live in convex subset  $\mathcal{F}$  of a Hilbert space with inner product  $\langle \cdot, \cdot \rangle$ , and iteratively update  $f^t$  as follows. At iteration  $t$ , an index  $i_t \in [N]$  is chosen uniformly at random and the gradient  $\nabla_f \varphi(f^t, s^{i_t})$  is computed at  $f^t$ . The parameter  $f$  is then updated via

$$(25) \quad f^{t+1} = \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \langle f, \nabla \varphi(f^t, s^{i_t}) \rangle + \Phi(f) + \frac{1}{2\eta_t} \|f - f^t\|^2 \right\},$$

where  $\eta_t > 0$  is an iteration-dependent stepsize and  $\|\cdot\|$  denotes the Hilbert norm. The convergence guarantees of the update (25) are well-understood [32, 17, 19]. Define  $\bar{f}^T = (1/T) \sum_{t=1}^T f^t$  to be the average parameter after  $T$  iterations. If the function  $\widehat{R}_n$  is strongly convex—meaning it has at least quadratic curvature—the step-size choice  $\eta_t \propto 1/t$  gives

$$\mathbb{E} \left[ \widehat{R}_N(\bar{f}^T) \right] - \inf_{f \in \mathcal{F}} \widehat{R}_N(f) = \mathcal{O} \left( \frac{1}{T} \right),$$

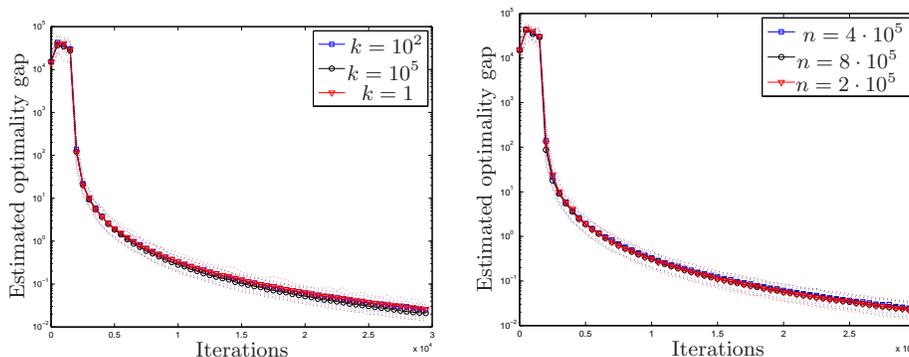
where the expectation is taken with respect to the indices  $i_t$  chosen during each iteration of the algorithm. In the convex case (without assuming any stronger properties than convexity), the step-size choice  $\eta_t \propto 1/\sqrt{t}$  yields

$$\mathbb{E} \left[ \widehat{R}_N(\bar{f}^T) \right] - \inf_{f \in \mathcal{F}} \widehat{R}_N(f) = \mathcal{O} \left( \frac{1}{\sqrt{T}} \right).$$

These guarantees also hold with high probability [19].

Neither of the convergence guarantees  $1/T$  or  $1/\sqrt{T}$  depends on the number  $N$  of terms in the stochastic objective (24). As a consequence, we can apply the composite stochastic gradient method (25) directly to the empirical risk (19): we sample a query  $q$  with probability  $\widehat{n}_q/n$ , after which we

uniformly sample one of the  $\binom{\hat{n}_q}{k}$  collections  $\{i_1, \dots, i_k\}$  of  $k$  indices associated with query  $q$ , and we then perform the gradient update (25) using the gradient sample  $\nabla\varphi(f^t, s(Y_{i_1}, \dots, Y_{i_k}))$ . This stochastic gradient scheme means that we can minimize the empirical risk in a number of iterations independent of both  $n$  and  $k$ ; the run-time behavior of the method scales independently of  $n$  and depends on  $k$  only so much as computing an instantaneous gradient  $\nabla\varphi(f, s(Y_1, \dots, Y_k))$  increases with  $k$ .



**Fig 2.** Timing experiments for different values of  $k$  and  $n$  when applying the method (25). The horizontal axes are the number of stochastic gradient iterations, the vertical axes are the estimated optimality gap for the empirical surrogate risk. Left: varying amount of aggregation  $k$ , fixed  $n = 4 \cdot 10^5$ . Right: varying total number of samples  $n$ , fixed  $k = 10^2$ .

In Figure 2, we show empirical evidence that the stochastic method (25) works as described. In particular, we minimize the empirical  $U$ -statistic-based risk (19) with the loss (28) we employ in our experiments in the next section. In each plot in Figure 2, we give an estimated optimality gap,  $\widehat{R}_{\varphi,n}(f^t) - \inf_{f \in \mathcal{F}} \widehat{R}_{\varphi,n}(f)$ , as a function of  $t$ , the number of iterations. As in the section to follow,  $\mathcal{F}$  consists of linear functionals parameterized by a vector  $\theta \in \mathbb{R}^d$  with  $d = 136$ . To estimate  $\inf_{f \in \mathcal{F}}$ , we perform 100,000 updates of the procedure (25), then estimate  $\inf_{f \in \mathcal{F}} \widehat{R}_{\varphi,n}(f)$  using the output predictor  $\widehat{f}$  evaluated on an additional (independent) 50,000 samples (the number of terms in the true objective is too large to evaluate). To estimate the risk  $\widehat{R}_{\varphi,n}(f^t)$ , we use a moving average of the previous 100 sampled losses  $\varphi(f^\tau, s^{i_\tau})$  for  $\tau \in \{t - 99, \dots, t\}$ , which is an unbiased estimate of an upper bound on the empirical risk  $\widehat{R}_{\varphi,n}(f^t)$  (see, e.g. [7]). We perform the experiment 20 times and plot averages as well as 90% confidence intervals. As predicted by our theoretical results, the number of iterations to attain a particular accuracy is essentially independent of  $n$  and  $k$ ; all the plots lie on

one another.

6.2. *Experimental evaluation.* To perform our experimental evaluation, we use a subset of the Microsoft Learning to Rank Web10K dataset [33], which consists of 10,000 web searches (queries) issued to the Microsoft Bing search engine, a set of approximately 100 potential results for each query, and a relevance score  $r \in \mathbb{R}$  associated with each query/result pair. A query/result pair is represented by a  $d = 136$ -dimensional feature vector of standard document-retrieval features.

To understand the benefits of aggregation and consistency in the presence of partial preference data, we generate pairwise data from the observed query/result pairs, so that we know the true asymptotic generating distribution. We adopt a loss  $L$  from the NDCG-family (14) and compare three surrogate losses: a consistent regression surrogate based on aggregation, an inconsistent but commonly used pairwise logistic loss [16], and a consistent loss that requires access to complete preference data [34]. Recalling the NDCG score (14) of a prediction vector  $\alpha \in \mathbb{R}^m$  for scores  $s \in \mathbb{R}^m$  (where  $\pi_\alpha$  is the permutation induced by  $\alpha$ ), we have the loss

$$L(\alpha, s) = 1 - \frac{1}{Z(s)} \sum_{j=1}^m \frac{G(s(j))}{F(\pi_\alpha(j))},$$

where  $Z(s)$  is the normalizing value for the NDCG score, and  $F(\cdot)$  and  $G(\cdot)$  are increasing functions.

Given a set of queries  $q$  and relevance scores  $r_i \in \mathbb{R}$ , we generate  $n$  pairwise preference observations according to a Bradley-Terry-Luce (BTL) model [5]. That is, for each observation, we choose a query  $q$  uniformly at random and then select a uniformly random pair  $(i, j)$  of results to compare. The pair is ordered as  $i \succ j$  (item  $i$  is preferred to  $j$ ) with probability  $p_{ij}$ , and  $j \succ i$  with probability  $1 - p_{ij} = p_{ji}$ , where

$$(26) \quad p_{ij} = \frac{\exp(r_i - r_j)}{1 + \exp(r_i - r_j)},$$

for  $r_i$  and  $r_j$  the respective relevances of results  $i$  and  $j$  under query  $q$ .

We define our structure functions  $s_k$  as score vectors in  $\mathbb{R}^m$ , where given a set of  $k$  preference pairs, the score for item  $i$  is

$$s_k(i) = \frac{1}{m-1} \sum_{j \neq i} \log \frac{\widehat{\mathbb{P}}(j \prec i)}{\widehat{\mathbb{P}}(j \succ i)},$$

the average empirical log-odds of result  $i$  being preferred to any other result.

Under the BTL model (26), as  $k \rightarrow \infty$  the structural score converges to

$$(27) \quad s(i) = \frac{1}{m-1} \sum_{j \neq i} [\log(1 + \exp(r_i - r_j)) - \log(1 + \exp(r_j - r_i))].$$

In our setting we may thus evaluate the asymptotic NDCG risk of a scoring function  $f$  by computing the asymptotic scores (27). In addition, Corollary 1 shows that if all minimizers of a loss obey the ordering of the values

$$\int_{\mathcal{S}} \frac{G(s(j))}{Z(s)} d\mu(s), \quad j \in \{1, \dots, m\}$$

then the loss is Fisher-consistent. A well-known example [12, 34] of such a loss is the least-squares loss, where the regression labels are  $G(s_j)/Z(s)$ :

$$(28) \quad \varphi(\alpha, s) = \frac{1}{2m} \sum_{j=1}^m \left( \alpha_j - \frac{G(s(j))}{Z(s)} \right)^2.$$

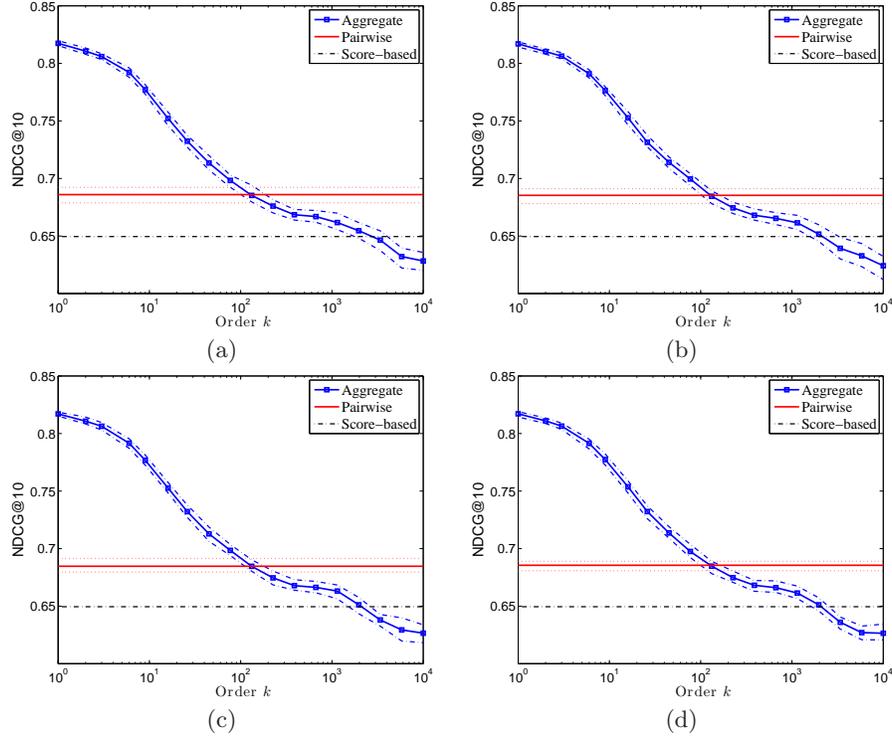
We compare the least-squares aggregation loss with a pairwise logistic loss natural for the pairwise data generated according to the BTL model (26). Specifically, given a sample pair with  $i \succ j$ , the logistic surrogate loss is

$$(29) \quad \varphi(\alpha, i \succ j) = \log(1 + \exp(\alpha_j - \alpha_i)),$$

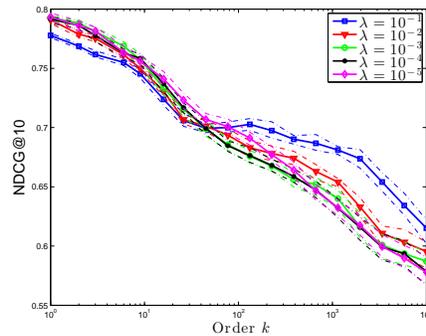
which is equivalent or similar to previous losses used for pairwise data in the ranking literature [27, 16]. For completeness, we also compare with a consistent surrogate that requires access to complete preference information in the form of the asymptotic structure scores (27). Following Ravikumar et al. [34], we obtain such a surrogate by granting the regression loss (28) direct access to the asymptotic structure scores. Note that such a construction would be infeasible in any true pairwise data setting.

Having described our sampling procedure, aggregation strategy, and loss functions, we now describe our model. We let  $x_i^q$  denote the feature vector for the  $i$ th result from query  $q$ , and we model the scoring function  $f(q)_i = \langle \theta, x_i^q \rangle$  for a vector  $\theta \in \mathbb{R}^d$ . For the regression loss (28), we minimize the  $U$ -statistic-based empirical risk (19) over a variety of orders  $k$ , while for the pairwise logistic loss (29), we minimize the empirical risk over all pairs sampled according to the BTL model (26). We regularize our estimates by adding  $\Phi(\theta) = (\lambda/2) \|\theta\|_2^2$  to the objective minimized, and we use the specialized stochastic method (25) to minimize the empirical risk.

Our goals in the experiments are to understand the behavior of the empirical risk minimizer as the order  $k$  of the aggregating statistic is varied and to evaluate the extent to which aggregation improves the estimated scoring function. A secondary concern is to verify that the method is insensitive to



**Fig 3.** NDCG risk and 95% confidence intervals for  $\theta$  estimated using the logistic pairwise loss (29) and the  $U$ -statistic empirical risk with  $\varphi$  chosen to be regression loss (28). The horizontal axis of each plot is the order  $k$  of the aggregation in the  $U$ -statistic (19), the vertical axis is the NDCG risk, and each plot corresponds to a different number  $n$  of samples. (a)  $n = 2 \cdot 10^5$  (b)  $n = 4 \cdot 10^5$  (c)  $n = 8 \cdot 10^5$  (d)  $n = 1.6 \cdot 10^6$ .



**Fig 4.** NDCG risk and 95% confidence intervals for  $\theta$  estimated using the  $U$ -statistic empirical risk (19) with  $\varphi$  chosen as the regression loss (28) under various choices of the regularization parameter,  $\lambda$ .

the amount  $\lambda$  of regularization performed on  $\theta$ . We run each experiment 50 times and report confidence intervals based on those 50 experiments.

Let  $\theta_{n,k}^{\text{reg}}$  denote the estimate of  $\theta$  obtained from minimizing the empirical risk (19) with the regression loss (28) on  $n$  samples with aggregation order  $k$ , let  $\theta_n^{\text{log}}$  denote the estimate of  $\theta$  obtained from minimizing the empirical pairwise logistic loss (29), and let  $\theta^{\text{full}}$  denote the estimate of  $\theta$  obtain from minimizing the empirical risk with surrogate loss (28) using the asymptotic structure scores (27) directly. Then each plot of Figure 3 displays the risk  $R(\theta_{n,k}^{\text{reg}})$  as a function of the aggregation order  $k$ , using  $R(\theta_n^{\text{log}})$  and  $R(\theta^{\text{full}})$  as references. The four plots in the figure correspond to different numbers  $n$  of sample pairs.

Broadly, the four plots in Figure 3 match our theoretical results. Consistently across the plots, we see that for small  $k$ , it appears there is not sufficient aggregation in the regression-loss-based empirical risk, and for such small  $k$  the pairwise logistic loss is better. However, as the order of aggregation  $k$  grows, the risk performance of  $\theta_{n,k}^{\text{reg}}$  improves. In addition, with larger sample sizes  $n$ , the difference between the risk of  $\theta_n^{\text{log}}$  and  $\theta_{n,k}^{\text{reg}}$  becomes more pronounced. The second salient feature of the plots is a moderate flattening of the risk  $R(\theta_{n,k}^{\text{reg}})$  and widening of the confidence interval for large values of  $k$ . This seems consistent with the estimation error guarantees in Propositions 5 and 6, where the order  $k$  being large has an evidently detrimental effect. Interestingly, however, large values of  $k$  still yield significant improvements over  $R(\theta_n^{\text{log}})$ . For very large  $k$ , the improved performance of  $\theta_{n,k}^{\text{reg}}$  over  $\theta^{\text{full}}$  is a consequence of sampling artifacts and the fact that we use a finite dimensional representation. (By using sufficiently many dimensions  $d$ , the estimator  $\theta^{\text{full}}$  attains zero risk by matching the asymptotic scores (27) directly.)

Figure 4 displays the risk  $R(\theta_{n,k}^{\text{reg}})$  for  $n = 800000$  pairs,  $k = 100$ , and multiple values of the regularization multiplier  $\lambda$  on  $\|\theta\|_2^2$ . The results, which are consistent across many choices of  $n$ , suggest that minimization of the aggregated empirical risk (19) is robust to the choice of regularization multiplier.

**7. Conclusions.** In this paper, we demonstrated both the difficulty and the feasibility of designing consistent, practicable procedures for ranking. By giving necessary and sufficient conditions for the Fisher-consistency of ranking algorithms, we proved that many natural ranking procedures based on surrogate losses are inconsistent, even in low-noise settings. To address this inconsistency while accommodating the incomplete nature of typical ranking data, we proposed a new family of surrogate losses, based on  $U$ -statistics, that aggregate disparate partial preferences. We showed how our

losses can fruitfully leverage any well behaved rank aggregation procedure and demonstrated their empirical benefits over more standard surrogates in a series of ranking experiments.

Our work thus takes a step toward bringing the consistency literature for ranking in line with that for classification, and we anticipate several directions of further development. First, it would be interesting to formulate low-noise conditions under which faster rates of convergence are possible for ranking risk minimization (see, e.g., the work of Cl  men  on et al. [10], which focuses on the minimization of a single pairwise loss). Additionally, it may be interesting to study structure functions  $s$  that yield non-point distributions  $\mu$  as the number of arguments  $k$  grows to infinity. For example, would scaling the Thurstone-Mosteller least-squares solutions (22) by  $\sqrt{k}$ —to achieve asymptotic normality—induce greater robustness in the empirical minimizer of the  $U$ -statistic risk (19)? Finally, exploring tractable formulations of other supervised learning problems in which label data is naturally incomplete could be fruitful.

## References.

- [1] A. Ammar and D. Shah. Ranking: compare, don't score. In *The 49th Allerton Conference on Communication, Control, and Computing*, 2011.
- [2] K. J. Arrow. *Social Choice and Individual Values*. Wiley, New York, 1951.
- [3] P. L. Bartlett, M. I. Jordan, and J. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101:138–156, 2006.
- [4] A. Ben-Tal and A. Nemirovski. *Lectures on Modern Convex Optimization*. SIAM, 2001.
- [5] R. A. Bradley and M. E. Terry. Rank analysis of incomplete block designs: I. The method of paired comparisons. *Biometrika*, 39(3/4), 1952.
- [6] D. Buffoni, C. Calauzenes, P. Gallinari, and N. Usunier. Learning scoring functions with order-preserving losses and standardized supervision. In *Proceedings of the 28th International Conference on Machine Learning*, 2011.
- [7] N. Cesa-Bianchi, A. Conconi, and C. Gentile. On the generalization ability of on-line learning algorithms. In *Advances in Neural Information Processing Systems 14*, pages 359–366, 2002.
- [8] O. Chapelle, D. Metzler, Y. Zhang, and P. Grinspan. Expected reciprocal rank for graded relevance. In *Conference on Information and Knowledge Management*, 2009.
- [9] F. R. K. Chung. *Spectral Graph Theory*. AMS, 1998.
- [10] S. Cl  men  on, G. Lugosi, and N. Vayatis. Ranking and empirical minimization of  $u$ -statistics. *Annals of Statistics*, 36(2):844–874, 2008.
- [11] N. Condorcet. *Essai sur l'Application de l'Analyse    la Probabilit   des D  cisions Rendues    la Pluralit   des Voix*. Paris, 1785.
- [12] D. Cossack and T. Zhang. Statistical analysis of Bayes optimal subset ranking. *IEEE Transactions on Information Theory*, 16:1274–1286, 2008.
- [13] N. Craswell, O. Zoeter, M. J. Taylor, and B. Ramsey. An experimental comparison of click position-bias models. In *Web Search and Data Mining (WSDM)*, pages 87–94, 2008.
- [14] H. A. David. *The Method of Paired Comparisons*. Charles Griffin & Company, 1969.

- [15] J. C. de Borda. *Memoire sur les Elections au Scrutin*. Histoire de l'Academie Royale des Sciences, Paris, 1781.
- [16] O. Dekel, C. Manning, and Y. Singer. Log-linear models for label ranking. In *Advances in Neural Information Processing Systems 16*, 2004.
- [17] J. C. Duchi and Y. Singer. Efficient online and batch learning using forward-backward splitting. *Journal of Machine Learning Research*, 10:2873–2898, 2009.
- [18] J. C. Duchi, L. Mackey, and M. I. Jordan. On the consistency of ranking algorithms. In *Proceedings of the 27th International Conference on Machine Learning*, 2010.
- [19] J. C. Duchi, S. Shalev-Shwartz, Y. Singer, and A. Tewari. Composite objective mirror descent. In *Proceedings of the Twenty Third Annual Conference on Computational Learning Theory*, 2010.
- [20] C. Dwork, R. Kumar, M. Naor, and D. Sivakumar. Rank aggregation methods for the web. In *Proceedings of the Tenth International World Wide Web Conference (WWW10)*, 2001.
- [21] Y. Freund, R. Iyer, R. E. Schapire, and Y. Singer. Efficient boosting algorithms for combining preferences. *Journal of Machine Learning Research*, 4:933–969, 2003.
- [22] H. Gulliksen. A least squares method for paired comparisons with incomplete data. *Psychometrika*, 21:125–134, 1956.
- [23] R. Herbrich, T. Graepel, and K. Obermayer. Large margin rank boundaries for ordinal regression. In *Advances in Large Margin Classifiers*. MIT Press, 2000.
- [24] J. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms I & II*. Springer, New York, 1996.
- [25] J. Hopcroft and J. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley, Reading, MA, 1979.
- [26] K. Järvelin and J. Kekäläinen. Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems*, 20(4):422–446, 2002.
- [27] T. Joachims. Optimizing search engines using clickthrough data. In *Proceedings of the ACM Conference on Knowledge Discovery and Data Mining*, 2002.
- [28] C. L. Mallows. Non-null ranking models. *Biometrika*, 44(1/2):pp. 114–130, 1957.
- [29] C. Manning, P. Raghavan, and H. Schütze. *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- [30] G. Miller. The magic number seven, plus or minus two: Some limits on our capacity for processing information. *Psychology Review*, 63:81–97, 1956.
- [31] F. Mosteller. Remarks on the method of paired comparisons: I. The least squares solution assuming equal standard deviations and equal correlations. *Psychometrika*, 16(1):3–9, 1951.
- [32] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- [33] T. Qin, T.-Y. Liu, W. Ding, J. Xu, and H. Li. Microsoft learning to rank datasets. URL <http://research.microsoft.com/en-us/projects/mslr/>, accessed March 1, 2012, 2011.
- [34] P. Ravikumar, A. Tewari, and E. Yang. On NDCG consistency of listwise ranking methods. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 2011.
- [35] T. L. Saaty. Decision making with the AHP: why is the principal eigenvector necessary. *European Journal of Operational Research*, 145:85–91, 2003.
- [36] T. L. Saaty. Relative measurement and its generalization in decision making. *Review of the Royal Spanish Academy of Sciences, Series A, Mathematics*, 102(2):251–318, 2008.

- [37] A. Shashua and A. Levin. Ranking with large margin principle: Two approaches. In *Advances in Neural Information Processing Systems 15*, 2002.
- [38] R. Shiffrin and R. Nosofsky. Seven plus or minus two: a commentary on capacity limitations. *Psychological Review*, 101(2):357–361, 1994.
- [39] I. Steinwart. How to compare different loss functions. *Constructive Approximation*, 26:225–287, 2007.
- [40] N. Stewart, G. Brown, and N. Chater. Absolute identification by relative judgment. *Psychological Review*, 112(4):881–911, 2005.
- [41] L. L. Thurstone. A law of comparative judgment. *Psychological Review*, 34(4):273–286, 1927.
- [42] K. Tsukida and M. R. Gupta. How to analyze paired comparison data. Technical Report UWEEETR-2011-0004, University of Washington Department of Electrical Engineering, 2011.
- [43] T. Zhang. Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5:1225–1251, 2004.
- [44] T. Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *Annals of Statistics*, 32:56–85, 2004.

DEPARTMENTS OF EECS AND STATISTICS  
UNIVERSITY OF CALIFORNIA, BERKELEY  
E-MAIL: [jduchi@cs.berkeley.edu](mailto:jduchi@cs.berkeley.edu)  
[jordan@stat.berkeley.edu](mailto:jordan@stat.berkeley.edu)

DEPARTMENT OF STATISTICS  
STANFORD UNIVERSITY  
E-MAIL: [lmackey@stanford.edu](mailto:lmackey@stanford.edu)