

Comparing and Aggregating Rankings with Ties

Ronald Fagin^{*}
IBM Almaden

Ravi Kumar^{*}
IBM Almaden

Mohammad Mahdian[†]
MIT

D. Sivakumar^{*}
IBM Almaden

Erik Vee[‡]
University of Washington

ABSTRACT

Rank aggregation has recently been proposed as a useful abstraction that has several applications, including meta-search, synthesizing rank functions from multiple indices, similarity search, and classification. In database applications (catalog searches, fielded searches, parametric searches, etc.), the rankings are produced by sorting an underlying database according to various fields. Typically, there are a number of fields that each have very few distinct values, and hence the corresponding rankings have many ties in them. Known methods for rank aggregation are poorly suited to this context, and the difficulties can be traced back to the fact that we do not have sound mathematical principles to compare two *partial rankings*, that is, rankings that allow ties.

In this work, we provide a comprehensive picture of how to compare partial rankings. We propose several metrics to compare partial rankings, present algorithms that efficiently compute them, and prove that they are within constant multiples of each other. Based on these concepts, we formulate aggregation problems for partial rankings, and develop a highly efficient algorithm to compute the top few elements of a near-optimal aggregation of multiple partial rankings. In a model of access that is suitable for databases, our algorithm reads essentially as few elements of each partial ranking as are necessary to determine the winner(s).

^{*}IBM Almaden Research Center, Department K53/B2, 650 Harry Road, San Jose, CA 95120, USA. email: {fagin,ravi,siva}@almaden.ibm.com

[†]CSAIL, MIT, 200 Technology Square, Cambridge, MA 02139, USA. Part of this work was done while visiting IBM Almaden. email: mahdian@theory.lcs.mit.edu

[‡]Dept. of Computer Science and Engineering, University of Washington, Seattle, WA 98195, USA. Supported by NSF grant CCR-0098066. Part of this work was done while visiting IBM Almaden. email: env@cs.washington.edu

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage, and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PODS 2004 June 14-16, 2004, Paris, France.

Copyright 2004 ACM 1-58113-858-X/04/06 . . . \$5.00.

1. INTRODUCTION

Rank aggregation is the problem of combining several ranked lists of objects in a robust way to produce a single ranking of the objects. This problem has a long and interesting history that goes back at least two centuries. While the philosophical aspects of rank aggregation have been debated extensively during this period, the mathematics of rank aggregation has gained more attention in the last eighty years, and the computational aspects are still within the purview of active research.

In computer science, rank aggregation has proved to be a useful and powerful paradigm in several applications including meta-search [7, 20, 18, 17, 1, 16], combining experts [3], synthesizing rank functions from multiple indices [8], biological databases [19], similarity search [10], and classification [16, 10]. An important contribution of the work of [7] is to adopt and highlight the merits of a proposal of Kemeny's for performing rank aggregation: namely, given multiple rankings, find a ranking whose total *Kendall tau distance* to the given rankings is minimized. The Kendall tau distance between two rankings (permutations) is defined as the number of pairwise disagreements between the two rankings; it is easy to see that it is a metric on the space of permutations. While this formulation is mathematically crisp, applications usually demand additional flexibility from the aggregation algorithm. For example, in aggregating search results, we are faced with the problem that we have access only to the top few elements of the rankings. In [7], this issue was addressed by suitably modifying the heuristics for full rank aggregation, but without providing any mathematical justification. This situation was later remedied in [9], where notions of "near metrics" were introduced, and a robust, unified class of near metrics was identified to compare "top k lists". This allowed the formulation of appropriate rank aggregation problems for top k lists, and the design of efficient (approximation) algorithms for these problems.

Challenges for rank aggregation in databases. While the extensive work in economics and computer science provide a mathematical basis for aggregation of full/top k rankings, the context of database-centric applications poses two formidable challenges for rank aggregation. We outline these next.

(1) In many applications of rank aggregation, there is an underlying database of records that are first ranked in several ways; typically, each ranked list is produced when the user specifies some criterion to rank (and/or filter) the records according to one of the attributes in the schema. Once the records are sorted in different ways, an aggregation algorithm combines the ranked lists to produce the final output. Common examples are catalog searches, fielded/parametric searches, and "advanced search" options.

For example, in online commerce, users often state their preferences for products according to various criteria. In a database of restaurants (e.g., www.dine.com), it is common to rank the restaurants based on the user's preferences for cuisine, driving distance, price, star ratings, etc.; in airline reservations (e.g., www.travelocity.com), it is common to rank flight plans by price, airline preference, number of connections, flight times, etc. Other examples include searching for an article in scientific bibliography databases (e.g., www.ams.org/mathscinet) using preference criteria on attributes such as title, year of publication, number of citations, etc.; searching for a protein in biological databases (e.g., www.rcsb.org/pdb) based on attributes like chain type, compound information, experimental technique, resolution, etc.; and searching for NSF awards (www.nsf.gov/verity/srchawdf.htm) based on attributes such as award amount, start date, etc.

While many database attributes are usually numeric, there are attributes that are inherently non-numeric. For instance, in the restaurant selection example above, "type of cuisine" is a non-numeric attribute. The number of distinct values in such non-numeric attributes is often very small. Therefore, when one sorts according to values this attribute can take, the resulting rank ordering of the objects is not a permutation any more; it is an ordering with ties, also known as a *partial ranking*. Notice that partial rankings could result even for numeric attributes. For example, in travel reservations, the field "number of connections" is a numeric attribute, but usually has no more than four values. Furthermore, the user may not be interested in using the complete range of values for a numeric attribute even if the database might permit it. For instance, in the restaurant selection example, even though distance is numeric, the user might wish to treat any distance up to ten miles to be the same in his/her preference.

Thus, the first main feature of rank aggregation in database applications is that, due to preference criteria on few-valued attributes, we need to deal with partial rankings rather than full rankings. While it is possible to treat this issue heuristically by arbitrarily ordering the tied elements to produce a full ranking, we seek ways that are mathematically more well-founded.

(2) In database-centric applications, we are often interested in only the top few answers of the aggregation. Certainly, this is the case with all the above examples. This feature leads to the quest for algorithms that quickly obtain the top result(s) of aggregation, perhaps in sub-linear time, without even having to read each ranking in its entirety. This issue was addressed in [10], where an aggregation heuristic based on median rank values was studied. This median rank aggregation has the nice property that it admits an *instance-optimal algorithm* in the sense of [11] under a model of access that is relevant for databases. The instance optimality feature is shared neither by the more sophisticated heuristics in [7] based on matchings and Markov chains, nor by a natural heuristic based on average ranks. Also, median is robust, as it mitigates the effect of outliers.

Since the applications we focus on in this paper are database-centric, it is tempting to try to adapt the median-based algorithm for aggregating partial rankings. However, there are two obstacles to such an attempt. First of all, though the median rank aggregation algorithm was argued to be heuristically good as an aggregation algorithm, nothing provable was known about its efficacy. In particular, it was not known if median rank aggregation produced an approximately optimal aggregation with respect to the Kendall distance. Secondly, the median rank aggregation algorithm was proposed in [7, 10] assuming that the inputs are permutations. Consequently, it is not clear if the algorithm would perform well, even in a heuristic sense, when the inputs are partial rankings.

To summarize, the aggregation of partial rankings is an important problem in the context of many database applications and it is useful to develop algorithms that quickly obtain the top few results of the aggregation. The single main obstacle is that we do not have sound mathematical principles to compare two partial rankings; this is precisely what we study in this paper. Our main contribution is a comprehensive solution to comparing and aggregating partial rankings.

Summary of our contributions. We define four metrics between partial rankings. These are obtained by suitably generalizing the Kendall tau distance and the Spearman footrule distance on permutations (cf. [5]) in two different ways. In both approaches, to compare two partial rankings, we compare the two sets of full rankings obtained from the partial rankings by breaking ties in all possible ways. A classical way (cf. [4]) to compare two sets in a metric space is the well-known method of using the Hausdorff distance between the sets (see Section 3.2). The drawback of using the Hausdorff extensions of Kendall tau and Spearman footrule is that they are less intuitive. Our second method to compare the two sets avoids this pitfall, and is based on succinctly summarizing the two sets by compact vectors—their "profiles"—and applying the L_1 distance between the profile vectors. By definition, these metrics admit efficient computation, and furthermore, they are extremely intuitive and quite natural. These metrics are defined and discussed in Section 3.

While the metrics obtained through profiles can be efficiently computed, the Hausdorff metrics are max-min over exponentially large sets and it is not at all obvious a priori if they can be computed efficiently as well. We solve this problem by first obtaining a complete characterization of how the Hausdorff distance is achieved between two partial rankings (for both Kendall tau and Spearman footrule versions). Specifically, we show how to efficiently construct full rankings from the given partial rankings so that the Kendall/Spearman distance between the full rankings is equal to the corresponding Hausdorff distances between the partial rankings. These characterizations enable us to compute the Hausdorff distances efficiently; furthermore, while the proofs of the characterizations are technically quite intricate, the resulting algorithms are extremely simple. The computational aspects of the metrics are discussed in Section 4.

Having four metrics on partial rankings is good news, but exactly which one should a practitioner use to compare partial rankings? Furthermore, which one is best suited to formulating an aggregation problem for partial rankings? Our summary answer to these questions is that the exact choice doesn't matter much. Namely, following the lead of [9], we define two metrics to be *equivalent* if they are within constant multiples of each other. This notion was inspired by the Diaconis-Graham inequality [6], which says that the Kendall tau distance and the Spearman footrule distance are within a factor of two of each other. We show that all four of our metrics are equivalent in this sense. It is easy to show that the Hausdorff versions of the Kendall tau distance and the Spearman footrule distance are equivalent; the equivalence between the Hausdorff and the profile versions of the Kendall tau metric is also simple to establish. Proving equivalence for the profile metrics turns out to be rather tricky, and requires us to uncover considerable structure inside partial rankings. We present these equivalence results in Section 5.

Finally, we turn to algorithms that obtain the top few answers when we aggregate partial rankings. Here we fully reap the benefit of having defined four distinct metrics on partial rankings and having established their equivalence with much technical maneu-

vering. Namely, if we care primarily about aggregations that are approximately optimal with respect to a metric, we now have four viewpoints from which to attack the problem! Thus, a constant factor approximation algorithm for aggregation with respect to one metric is automatically a constant factor approximation algorithm for aggregation with respect to all the other metrics. It turns out that an algorithm that is based on the median rank algorithm [7, 10] lends itself naturally to efficient aggregation with respect to the profile version of the Spearman footrule metric. We show that the algorithm derived from median ranks is a constant factor approximation algorithm with respect to this metric. Also, as mentioned before, being a median-based algorithm, our algorithm reads only as few elements of each partial ranking as possible in order to determine the winner(s) of the aggregation—in this aspect, the algorithm is extremely database-friendly and practical.

By the equivalence outlined above, it follows that the median rank algorithm gives an approximation algorithm for rank aggregation with respect to all our metrics. It also vindicates the use of median in [7, 10]. These results are presented in Section 6.

Related work. Kendall [15] defined two variations of the Kendall tau distance for partial rankings of which one is a normalized version of the Kendall tau distance through profiles. Baggerly [2] defined two versions of the Spearman footrule distance for partial rankings of which one is similar to our Spearman footrule metric through profiles. However, neither work proceeds significantly beyond simply providing the definition. For top k lists, which are special case of partial rankings, Critchlow [4] defined Hausdorff versions of Kendall tau and Spearman footrule distances and Fagin et al. [9] studied further properties of these metrics. Goodman and Kruskal [12] proposed an approach for comparing partial rankings, which was recently utilized [13] for evaluating strategies for similarity search on the Web. A serious disadvantage of Goodman and Kruskal’s approach is that it is not always defined (this problem did not arise in the application of [13]).

2. PRELIMINARIES

Bucket orders. A *bucket order* is, intuitively, a linear order with ties. More formally, a bucket order is a transitive binary relation \triangleleft for which there are sets $\mathcal{B}_1, \dots, \mathcal{B}_t$ (the *buckets*) that form a partition of the domain such that $x \triangleleft y$ if and only if there are i, j with $i < j$ such $x \in \mathcal{B}_i$ and $y \in \mathcal{B}_j$. If $x \in \mathcal{B}_i$, we may refer to \mathcal{B}_i as the *bucket of x* . We may say that bucket \mathcal{B}_i *precedes* bucket \mathcal{B}_j if $i < j$. Thus, $x \triangleleft y$ if and only if the bucket of x precedes the bucket of y . We think of the members of a given bucket as “tied”. A linear order is a bucket order where every bucket is of size 1. We now define the *position* of bucket \mathcal{B} , denoted $\text{pos}(\mathcal{B})$. Let $\mathcal{B}_1, \dots, \mathcal{B}_t$ be the buckets in order (so that bucket \mathcal{B}_i precedes bucket \mathcal{B}_j when $i < j$). Then $\text{pos}(\mathcal{B}_i) = (\sum_{j < i} |\mathcal{B}_j|) + (|\mathcal{B}_i| + 1)/2$. Intuitively, $\text{pos}(\mathcal{B}_i)$ is the average location within bucket \mathcal{B}_i .

Partial ranking. Just as we can associate a ranking with a linear order (i.e., permutation), we associate a *partial ranking* σ with each bucket order, by letting $\sigma(x) = \text{pos}(\mathcal{B})$ when \mathcal{B} is the bucket of x . We refer to a partial ranking associated with a linear order as a *full ranking*. When it is not otherwise specified, we assume that all partial rankings have the same domain, denoted D . We say that x is *ahead of y in σ* if $\sigma(x) < \sigma(y)$. We say that x and y are *tied in σ* if $\sigma(x) = \sigma(y)$. When we speak of the buckets of a partial ranking, we are referring to the buckets of the corresponding bucket order.

We define a *top k list* to be a partial ranking where the top k buckets are singletons, representing the top k elements, and the

bottom bucket contains all other members of the domain. Note that in [9] there is no bottom bucket in a top k list. This is because in [9] each top k list has its own domain of size k , unlike our scenario where there is a fixed domain.

Given a partial ranking σ with domain D , we define its *reverse*, denoted σ^R , in the expected way. That is, for all $d \in D$, let $\sigma^R(d) = |D| + 1 - \sigma(d)$.

We also define the notion of *swapping* in the normal way. If $a, b \in D$, then *swapping a and b in σ* produces a new order σ' where $\sigma'(a) = \sigma(b)$, $\sigma'(b) = \sigma(a)$, and $\sigma'(d) = \sigma(d)$ for all $d \in D - \{a, b\}$.

Refinements of partial rankings. Given two partial rankings σ and τ , both with domain D , we say that σ is a *refinement* of τ and write $\sigma \succeq \tau$ if the following holds: for all $i, j \in D$, we have $\sigma(i) < \sigma(j)$ whenever $\tau(i) < \tau(j)$. Notice that when $\tau(i) = \tau(j)$, there is no order forced on σ . When σ is a full ranking, we say that σ is a *full refinement* of τ . Given two partial rankings σ and τ , both with domain D , we frequently make use of a particular refinement of σ in which ties are broken according to τ . Define the τ -*refinement* of σ , denoted $\tau * \sigma$, to be the refinement of σ with the following properties. For all $i, j \in D$, if $\sigma(i) = \sigma(j)$ and $\tau(i) < \tau(j)$, then $\tau * \sigma(i) < \tau * \sigma(j)$. If $\sigma(i) = \sigma(j)$ and $\tau(i) = \tau(j)$, then $\tau * \sigma(i) = \tau * \sigma(j)$. Notice that when τ is in fact a full ranking, then $\tau * \sigma$ is also a full ranking. Also note that $*$ is an associative operation, so that if ρ is a partial ranking with domain D , it makes sense to talk about $\rho * \tau * \sigma$.

Notation. When f and g are functions with the same domain D , we denote the L_1 distance between f and g by $L_1(f, g)$. Thus, $L_1(f, g) = \sum_{i \in D} |f(i) - g(i)|$.

2.1 Metrics, near metrics, equivalence classes

A binary function d is called *symmetric* if $d(x, y) = d(y, x)$ for all x, y in the domain, and is called *regular* if $d(x, y) = 0$ if and only if $x = y$. A *distance measure* is a nonnegative, symmetric, regular binary function. A *metric* is a distance measure d that satisfies the *triangle inequality* $d(x, z) \leq d(x, y) + d(y, z)$ for all x, y, z in the domain.

The definitions and results in this section were derived in [9], in the context of comparing top k lists. Two seemingly different notions of a “near metric” were defined in [9]: their first notion of near metric is based on “relaxing” the polygonal inequality that a metric is supposed to satisfy.

DEFINITION 1 (NEAR METRIC). *A distance measure on partial rankings with domain D is a near metric if there is a constant c , independent of the size of D , such that the distance measure satisfies the relaxed polygonal inequality: $d(x, z) \leq c(d(x, x_1) + d(x_1, x_2) + \dots + d(x_{n-1}, z))$ for all $n > 1$ and $x, z, x_1, \dots, x_{n-1} \in D$.*

It makes sense to say that the constant c is independent of the size of D when, as in [9], each of the distance measures considered is actually a family, parameterized by D . We need to make an assumption that c is independent of the size of D , since otherwise we are simply considering distance measures over finite domains, where there is always such a constant c .

The other notion of near metric given in [9] is based on bounding the distance measure above and below by positive constant multiples of a metric. It was shown that both the notions of near metrics coincide.¹ This theorem inspired a definition of what it means for

¹This result would not hold if instead of relaxing the polygonal inequality, we simply relaxed the triangle inequality.

a distance measure to be “almost” a metric, and a robust notion of “similar” or “equivalent” distance measures. We modify the definitions in [9] slightly to fit our scenario, where there is a fixed domain D .

DEFINITION 2 (EQUIVALENT DISTANCE MEASURES). *Two distance measures d and d' between partial rankings with domain D are equivalent if there are positive constants c_1 and c_2 , independent of the size of D , such that $c_1 d'(\sigma_1, \sigma_2) \leq d(\sigma_1, \sigma_2) \leq c_2 d'(\sigma_1, \sigma_2)$, for every pair σ_1, σ_2 of partial rankings.*

It is clear that the above definition leads to an equivalence relation (i.e., reflexive, symmetric, and transitive). It follows from [9] that a distance measure is equivalent to a metric if and only if it is a near metric.

2.2 Metrics on full rankings

The study of metrics on full rankings is classical (cf. [14, 5]). We now review two well-known notions of metrics on full rankings, namely the Kendall tau distance and the Spearman footrule distance.

Let σ_1, σ_2 be two full rankings with domain D . The *Spearman footrule distance* is simply the L_1 distance $L_1(\sigma_1, \sigma_2)$. The definition of the Kendall tau distance requires a little more work.

Let $\mathcal{P} = \{\{i, j\} \mid i \neq j \text{ and } i, j \in D\}$ be the set of unordered pairs of distinct elements. The *Kendall tau distance* between full rankings is defined as follows. For each pair $\{i, j\} \in \mathcal{P}$ of distinct members of D , if i and j are in the same order in σ_1 and σ_2 , then let the penalty $\bar{K}_{i,j}(\sigma_1, \sigma_2) = 0$; and if i and j are in the opposite order (such as i being ahead of j in σ_1 and j being ahead of i in σ_2), then let $\bar{K}_{i,j}(\sigma_1, \sigma_2) = 1$. The Kendall tau distance is given by $K(\sigma_1, \sigma_2) = \sum_{\{i,j\} \in \mathcal{P}} \bar{K}_{i,j}(\sigma_1, \sigma_2)$. The Kendall tau distance turns out to be equal to the number of exchanges needed in a bubble sort to convert one full ranking to the other.

Diaconis and Graham [6] proved a classical result, which states that for every two full rankings σ_1, σ_2 ,

$$K(\sigma_1, \sigma_2) \leq F(\sigma_1, \sigma_2) \leq 2K(\sigma_1, \sigma_2). \quad (1)$$

In other words, Kendall tau and Spearman footrule are equivalent metrics for full rankings.

3. COMPARING PARTIAL RANKINGS

In this section we define the distance between partial rankings. The first set of metrics is based on profile vectors (Section 3.1) and the second set is based on the Hausdorff distance (Section 3.2). Section 3.3 compares these metrics (when the partial rankings are top k lists) with the distance measures for top k lists that are developed in [9].

3.1 Metrics based on profiles

Let σ_1, σ_2 be two partial rankings with domain D . We now define a family of generalizations of the Kendall tau distance to partial rankings. These are based on a generalization [9] of the Kendall tau distance to top k lists.

Let p be a fixed parameter, $0 \leq p \leq 1$. Similar to our definition of $\bar{K}_{i,j}(\sigma_1, \sigma_2)$ for full rankings σ_1, σ_2 , we define a penalty $\bar{K}_{i,j}^{(p)}(\sigma_1, \sigma_2)$ for partial rankings σ_1, σ_2 for $\{i, j\} \in \mathcal{P}$. There are three cases.

Case 1: i and j are in different buckets in both σ_1 and σ_2 . If i and j are in the same order in σ_1 and σ_2 (such as $\sigma_1(i) > \sigma_1(j)$ and $\sigma_2(i) > \sigma_2(j)$) then let $\bar{K}_{i,j}^{(p)}(\sigma_1, \sigma_2) = 0$; this corresponds to “no penalty” for $\{i, j\}$. If i and j are in the opposite order in σ_1

and σ_2 (such as $\sigma_1(i) > \sigma_1(j)$ and $\sigma_2(i) < \sigma_2(j)$) then let the penalty $\bar{K}_{i,j}^{(p)}(\sigma_1, \sigma_2) = 1$.

Case 2: i and j are in the same bucket in both σ_1 and σ_2 . We then let the penalty $\bar{K}_{i,j}^{(p)}(\sigma_1, \sigma_2) = 0$. Intuitively, both partial rankings agree that i and j are tied.

Case 3: i and j are in the same bucket in one of the partial rankings σ_1 and σ_2 , but in different buckets in the other partial ranking. In this case, we let the penalty $\bar{K}_{i,j}^{(p)}(\sigma_1, \sigma_2) = p$.

Based on these cases, define $K^{(p)}$, the *Kendall distance with penalty parameter p* , as follows:

$$K^{(p)}(\sigma_1, \sigma_2) = \sum_{\{i,j\} \in \mathcal{P}} \bar{K}_{i,j}^{(p)}(\sigma_1, \sigma_2).$$

We now discuss our choice of penalty in Cases 2 and 3. If the penalty value in Case 2 were strictly positive, then $K^{(p)}$ would not even be a distance measure, since we could have $K^{(p)}(\sigma, \sigma) > 0$. If $p = 0$ in Case 3, then we could have $K^{(p)}(\sigma_1, \sigma_2) = 0$ even with $\sigma_1 \neq \sigma_2$, and so again $K^{(p)}$ would not even be a distance measure. As for other choices of p , it turns out that $K^{(p)}$ is a metric for $p \in [1/2, 1]$, and is a near metric for $p \in (0, 1/2)$. The proof is deferred to the full version of the paper.

For the rest of the paper, we focus on the natural case $p = 1/2$, since it corresponds to an “average” penalty for two elements i and j that are tied in one partial ranking but not in the other partial ranking. We denote $K^{(1/2)}$ by K_{prof} , since, as we now show, there is an alternative but equivalent definition in terms of a “profile”.

Let $\mathcal{O} = \{(i, j) : i \neq j \text{ and } i, j \in D\}$ be the set of ordered pairs of distinct elements in the domain D . Let σ be a partial ranking (as usual, with domain D). For $(i, j) \in \mathcal{O}$, define p_{ij} to be $1/4$ if $\sigma(i) < \sigma(j)$, to be 0 if $\sigma(i) = \sigma(j)$, and to be $-1/4$ if $\sigma(i) > \sigma(j)$. Define the *K -profile* of σ to be the vector $\langle p_{ij} : (i, j) \in \mathcal{O} \rangle$. It is straightforward to verify that $K_{\text{prof}}(\sigma_1, \sigma_2)$ is simply the L_1 distance between the K -profiles of σ_1 and σ_2 .²

It is clear how to generalize the Spearman footrule distance to partial rankings—we simply take it to be $L_1(\sigma_1, \sigma_2)$, just as before. We refer to this value as $F_{\text{prof}}(\sigma_1, \sigma_2)$, for reasons we now explain. Let us define the *F -profile* of a partial ranking σ to be the vector of values $\sigma(i)$. So the F -profile is indexed by D , whereas the K -profile is indexed by \mathcal{O} . Just as the K_{prof} value of two partial rankings (or of the corresponding bucket orders) is the L_1 distance between their K -profiles, the F_{prof} value of two partial rankings (or of the corresponding bucket orders) is the L_1 distance between their F -profiles. Since F_{prof} and K_{prof} are L_1 distances, they are automatically metrics.

3.2 The Hausdorff metrics

Let A and B be finite sets of objects and let d be a metric of distances between objects. The *Hausdorff distance* between A and B is given by

$$d_{\text{Haus}}(A, B) = \max \left\{ \max_{\gamma_1 \in A} \min_{\gamma_2 \in B} d(\gamma_1, \gamma_2), \max_{\gamma_2 \in B} \min_{\gamma_1 \in A} d(\gamma_1, \gamma_2) \right\}. \quad (2)$$

Although this looks fairly nonintuitive, it is actually quite natural, as we now explain. The quantity $\min_{\gamma_2 \in B} d(\gamma_1, \gamma_2)$ is the distance between γ_1 and the set B . Therefore, the quantity $\max_{\gamma_1 \in A} \min_{\gamma_2 \in B} d(\gamma_1, \gamma_2)$ is the maximal distance

²Each pair $\{i, j\}$ with $i \neq j$ is counted twice, once as (i, j) and once as (j, i) . This is why the values of p_{ij} are $1/4, 0$, and $-1/4$ rather than $1/2, 0$, and $-1/2$.

of a member of A from the set B . Similarly, the quantity $\max_{\gamma_2 \in B} \min_{\gamma_1 \in A} d(\gamma_1, \gamma_2)$ is the maximal distance of a member of B from the set A . Therefore, the Hausdorff distance between A and B is the maximal distance of a member of A or B from the other set. Thus, A and B are within Hausdorff distance s of each other precisely if every member of A and B is within distance s of some member of the other set. The Hausdorff distance is well known to be a metric.

Critchlow [4] used the Hausdorff distance to define a metric between top k lists. We generalize his construction to give a metric between partial rankings. Given a metric d that gives the distance $d(\gamma_1, \gamma_2)$ between full rankings γ_1 and γ_2 , define the distance between partial rankings σ_1 and σ_2 to be

$$\max \left\{ \max_{\gamma_1 \succeq \sigma_1} \min_{\gamma_2 \succeq \sigma_2} d(\gamma_1, \gamma_2), \max_{\gamma_2 \succeq \sigma_2} \min_{\gamma_1 \succeq \sigma_1} d(\gamma_1, \gamma_2) \right\}, \quad (3)$$

where γ_1 and γ_2 are full rankings. In particular, when d is the footrule distance, this gives us the metric F_{Haus} between partial rankings, and when d is the Kendall distance, this gives us the metric K_{Haus} between partial rankings. Both F_{Haus} and K_{Haus} are indeed metrics, since they are special cases of the Hausdorff distance.

3.3 Discussion

Metrics on partial rankings naturally induce metrics on top k lists. We now compare the metrics on top k lists that are induced by our metrics on partial rankings with the distance measures on top k lists that were introduced in [9]. Recall that for us, a top k list is a partial ranking consisting of k singleton buckets, followed by a bottom bucket of size $|D| - k$. However, in [9], a top k list is a bijection of a domain onto $\{1, \dots, k\}$. Let σ and τ be top k lists (of our form). Define the *active domain* for σ, τ to be the union of the elements in the top k buckets of σ and the elements in the top k buckets of τ . In order to make our scenario compatible with the scenario of [9], we assume during our comparison that the domain D equals the active domain for σ, τ . Our definitions of $K^{(p)}$, F_{Haus} , and K_{Haus} are then exactly the same in the two scenarios. (Unlike earlier, even the case $p = 0$ gives a distance measure, since the unpleasant situation where $K^{(0)}(\sigma_1, \sigma_2) = 0$ even though $\sigma_1 \neq \sigma_2$ does not arise for top k lists σ_1 and σ_2 .) In spite of this, $K^{(p)}$, F_{Haus} , and K_{Haus} are only near metrics in [9] in spite of being metrics for us. This is because, in [9], the active domain varies depending on which pair of top k lists is being compared.

Our definition of $K_{\text{prof}}(\sigma, \tau)$ is equivalent to the definition of $K_{\text{avg}}(\sigma, \tau)$ in [9], namely the average value of $K(\sigma, \tau)$ over all full rankings σ, τ where $\sigma \succeq \sigma$ and $\tau \succeq \tau$. It is interesting to note that if σ and τ were not top k lists but arbitrary partial rankings, then K_{avg} would not be a distance measure, since $K_{\text{avg}}(\sigma, \sigma)$ can be strictly positive if σ is an arbitrary partial ranking.

Let ℓ be a real number greater than k . The *footrule distance with location parameter* ℓ , denoted $F^{(\ell)}$, is defined in [9] to be obtained, intuitively, by treating each element that is not among the top k elements as if it were in position ℓ , and then taking the L_1 distance. More formally, let σ and τ be top k lists (of our form). Define the function f_σ with domain D by letting $f_\sigma(i) = \sigma(i)$ if $1 \leq \sigma(i) \leq k$, and $f_\sigma(i) = \ell$ otherwise. Similarly, define the function f_τ with domain D by letting $f_\tau(i) = \tau(i)$ if $1 \leq \tau(i) \leq k$, and $f_\tau(i) = \ell$ otherwise. Then $F^{(\ell)}(\sigma, \tau)$ is defined to be $L_1(f_\tau, f_\sigma)$. It is straightforward to verify that $F_{\text{prof}}(\sigma, \tau) = F^{(\ell)}(\sigma, \tau)$ for $\ell = (|D| + k + 1)/2$.

4. COMPUTING THE METRICS

It is clear from the definition that both K_{prof} and F_{prof} can be computed in polynomial time. In this section we show how to compute the Hausdorff metrics K_{Haus} and F_{Haus} in polynomial time. We make use of these results later to prove that all of our metrics are in the same equivalence class. (Note that once we show in Section 5 that all the metrics are equivalent, then it follows that both the Hausdorff metrics can be approximated in polynomial time by computing the profile metrics.)

First, we state a fact that we use often (proof omitted).

LEMMA 3. *Suppose $a \leq b$ and $c \leq d$. Then $|a - c| + |b - d| \leq |a - d| + |b - c|$.*

We next show a simple lemma.

LEMMA 4. *Let π be a full ranking, and let σ be a partial ranking. Suppose that $\pi \neq \sigma$. Then there exist i, j such that $\pi(j) = \pi(i) + 1$ while $\sigma(j) \leq \sigma(i)$. If σ is in fact a full ranking, then $\sigma(j) < \sigma(i)$.*

PROOF. Order the elements of the domain $D = \{d_1, \dots, d_{|D|}\}$ so that $\pi(d_1) < \pi(d_2) < \dots < \pi(d_{|D|})$. If $\sigma(d_\ell) < \sigma(d_{\ell+1})$ for all ℓ , then we would have $K_{\text{prof}}(\sigma, \pi) = 0$, contradicting the fact that $\pi \neq \sigma$. Hence, there must be some ℓ for which $\sigma(d_{\ell+1}) \leq \sigma(d_\ell)$. Setting $i = d_\ell$ and $j = d_{\ell+1}$ gives us the lemma.

If σ is a full ranking, then $\sigma(j) \neq \sigma(i)$, showing $\sigma(j) < \sigma(i)$. \square

The next two lemmas are towards obtaining a characterization of the Hausdorff distance.

LEMMA 5. *Let σ be a full ranking, and let τ be a partial ranking. Then the quantity $F(\sigma, \tau)$, taken over all full refinements $\tau \succeq \tau$, is minimized for $\tau = \sigma * \tau$. Similarly, the quantity $K(\sigma, \tau)$, taken over all full refinements $\tau \succeq \tau$, is minimized for $\tau = \sigma * \tau$.*

PROOF. First, note that if $\tau \succeq \tau$ then there is a full ranking π such that $\tau = \pi * \tau$. We show that $F(\sigma, \sigma * \tau) \leq F(\sigma, \pi * \tau)$ and $K(\sigma, \sigma * \tau) \leq K(\sigma, \pi * \tau)$ for every full ranking π . The lemma will then follow. Let

$$U = \{\pi \mid \pi \text{ is a full ranking and } F(\sigma, \sigma * \tau) > F(\sigma, \pi * \tau)\},$$

$$V = \{\pi \mid \pi \text{ is a full ranking and } K(\sigma, \sigma * \tau) > K(\sigma, \pi * \tau)\},$$

and let $S = U \cup V$. If S is empty, then we are done. So suppose not. Over all full rankings $\pi \in S$, choose π to be the full ranking that minimizes $K(\sigma, \pi)$.

Since $\pi \neq \sigma$, Lemma 4 guarantees that we can find a pair i, j such that $\pi(j) = \pi(i) + 1$, but $\sigma(j) < \sigma(i)$. Produce π' by swapping i and j in π . Clearly, π' has one fewer inversion with respect to σ than π does. Hence, $K(\sigma, \pi') < K(\sigma, \pi)$. We show that $\pi' \in S$, thus giving a contradiction.

If i and j are in different buckets for τ , then $\pi' * \tau = \pi * \tau$. Hence, $F(\sigma, \pi' * \tau) = F(\sigma, \pi * \tau)$ and $K(\sigma, \pi' * \tau) = K(\sigma, \pi * \tau)$. So if $\pi \in U$, then $\pi' \in U$ as well. Similarly, if $\pi \in V$, then $\pi' \in V$. In either case, $\pi' \in S$.

On the other hand, assume that i and j are in the same bucket for τ . Then $\pi' * \tau(i) = \pi * \tau(j)$ and $\pi' * \tau(j) = \pi * \tau(i)$. Furthermore, since $\pi(i) < \pi(j)$ and i and j are in the same bucket, we have $\pi * \tau(i) < \pi * \tau(j)$, while $\sigma(j) < \sigma(i)$.

Either $\pi \in U$ or $\pi \in V$. First, consider the case where $\pi \in U$. Substituting $a = \pi * \tau(i)$, $b = \pi * \tau(j)$, $c = \sigma(j)$, $d = \sigma(i)$ in Lemma 3, we have

$$\begin{aligned} & |\pi' * \tau(j) - \sigma(j)| + |\pi' * \tau(i) - \sigma(i)| \\ &= |\pi * \tau(i) - \sigma(j)| + |\pi * \tau(j) - \sigma(i)| \\ &\leq |\pi * \tau(i) - \sigma(i)| + |\pi * \tau(j) - \sigma(j)| \end{aligned}$$

We also have $|\pi' * \tau(d) - \sigma(d)| = |\pi * \tau(d) - \sigma(d)|$ for all $d \in D - \{i, j\}$ since $\pi' * \tau$ and $\pi * \tau$ agree everywhere but at i and j . Summing, we have $F(\sigma, \pi' * \tau) \leq F(\sigma, \pi * \tau)$. Since $\pi \in U$, then $F(\sigma, \pi * \tau) < F(\sigma, \sigma * \tau)$. So $\pi' \in U$ by transitivity.

Now consider the case where $\pi \in V$. By our choice, $\pi(j) = \pi(i) + 1$. Hence, $\pi * \tau(j) = \pi * \tau(i) + 1$ since i and j are in the same bucket of τ . Similarly, $\pi' * \tau(i) = \pi' * \tau(j) + 1$. And as we noted earlier, $\pi * \tau$ and $\pi' * \tau$ agree everywhere except at i and j . In other words, $\pi' * \tau$ is just $\pi * \tau$, with the adjacent elements i and j swapped. Since $\sigma(i) > \sigma(j)$ we see that $\pi' * \tau$ has exactly one fewer inversion with respect to σ than $\pi * \tau$ does. That is, $K(\sigma, \pi' * \tau) < K(\sigma, \pi * \tau)$. Since $\pi \in V$, we have $K(\sigma, \pi * \tau) < K(\sigma, \sigma * \tau)$. So π' must be in V as well, by transitivity.

In either case, we have produced a $\pi' \in S$ such that $K(\sigma, \pi') < K(\sigma, \pi)$, contradicting the minimality of π . Hence, S must have been empty, as we wanted. \square

LEMMA 6. *Let σ and τ be partial rankings, and let ρ be any full ranking. Then the quantity $F(\sigma, \sigma * \tau)$, taken over all full refinements $\sigma \succeq \sigma$, is maximized when $\sigma = \rho * \tau^R * \sigma$. Similarly, the quantity $K(\sigma, \sigma * \tau)$, taken over all full refinements $\sigma \succeq \sigma$, is maximized when $\sigma = \rho * \tau^R * \sigma$.*

PROOF. First, note that for any full refinement $\sigma \succeq \sigma$, there is some full ranking π , such that $\sigma = \pi * \sigma$. We show that for all full rankings π that

$$\begin{aligned} F(\rho * \tau^R * \sigma, \rho * \tau^R * \sigma * \tau) &\geq F(\pi * \sigma, \pi * \sigma * \tau) \\ \text{and } K(\rho * \tau^R * \sigma, \rho * \tau^R * \sigma * \tau) &\geq K(\pi * \sigma, \pi * \sigma * \tau) \end{aligned}$$

The lemma will then follow.

Let $U = \{\text{full } \pi \mid F(\rho * \tau^R * \sigma, \rho * \tau^R * \sigma * \tau) < F(\pi * \sigma, \pi * \sigma * \tau)\}$, let $V = \{\text{full } \pi \mid K(\rho * \tau^R * \sigma, \rho * \tau^R * \sigma * \tau) < K(\pi * \sigma, \pi * \sigma * \tau)\}$, and let $S = U \cup V$. If S is empty, then we are done. So suppose not. Over all full rankings $\pi \in S$, choose π to be the full ranking that minimizes $K(\rho * \tau^R, \pi)$.

Since $\pi \neq \rho * \tau^R$, Lemma 4 guarantees that we can find a pair i, j such that $\pi(j) = \pi(i) + 1$, but $\rho * \tau^R(j) < \rho * \tau^R(i)$. Produce π' by swapping i and j . Clearly, π' has one fewer inversion with respect to $\rho * \tau^R$ than π does. That is, $K(\rho * \tau^R, \pi') < K(\rho * \tau^R, \pi)$. We now show that $\pi' \in S$, producing a contradiction.

If i and j are in different buckets for σ , then $\pi' * \sigma = \pi * \sigma$. Hence, $F(\pi' * \sigma, \pi' * \sigma * \tau) = F(\pi * \sigma, \pi * \sigma * \tau)$ and $K(\pi' * \sigma, \pi' * \sigma * \tau) = K(\pi * \sigma, \pi * \sigma * \tau)$. So if $\pi \in U$, then $\pi' \in U$. Similarly, if $\pi \in V$, then $\pi' \in V$. Hence, π must be in S .

Likewise, if i and j are in the same bucket for both σ and τ , then swapping i and j in π swaps their positions in both $\pi * \sigma * \tau$ and $\pi * \sigma$. So again, we see $F(\pi' * \sigma, \pi' * \sigma * \tau) = F(\pi * \sigma, \pi * \sigma * \tau)$ and $K(\pi' * \sigma, \pi' * \sigma * \tau) = K(\pi * \sigma, \pi * \sigma * \tau)$. As before, $\pi' \in S$.

Now, consider the case when i and j are in the same bucket for σ , but in different buckets for τ . First of all, $\pi' * \sigma$ is just $\pi * \sigma$ with i and j swapped since i and j are in the same bucket for σ ; further, notice that i and j are adjacent in $\pi * \sigma$. Second, $\pi' * \sigma * \tau = \pi * \sigma * \tau$ since i and j are in different buckets for τ .

Since $\pi(i) < \pi(j)$, we have $\pi * \sigma(i) < \pi * \sigma(j)$. Further, $\tau(i) < \tau(j)$ since $\rho * \tau^R(j) < \rho * \tau^R(i)$ and $\rho * \tau^R$ is a refinement of the reverse of τ . Hence, $\pi * \sigma * \tau(i) < \pi * \sigma * \tau(j)$. We have two cases to consider. Either $\pi \in U$ or $\pi \in V$.

Let us first examine the case that $\pi \in U$. Substituting $a = \pi * \sigma(i)$, $b = \pi * \sigma(j)$, $c = \pi * \sigma * \tau(i)$, $d = \pi * \sigma * \tau(j)$, in Lemma 3 gives us

$$\begin{aligned} &|\pi * \sigma(i) - \pi * \sigma * \tau(i)| + |\pi * \sigma(j) - \pi * \sigma * \tau(j)| \\ &\leq |\pi * \sigma(i) - \pi * \sigma * \tau(j)| + |\pi * \sigma(j) - \pi * \sigma * \tau(i)| \\ &= |\pi' * \sigma(j) - \pi' * \sigma * \tau(j)| + |\pi' * \sigma(i) - \pi' * \sigma * \tau(i)| \end{aligned}$$

We also have that $|\pi' * \sigma(d) - \pi' * \sigma * \tau(d)| = |\pi * \sigma(d) - \pi * \sigma * \tau(d)|$ for all $d \in D - \{i, j\}$. Summing over all d , we obtain $F(\pi * \sigma, \pi * \sigma * \tau) \leq F(\pi' * \sigma, \pi' * \sigma * \tau)$. Since $\pi \in U$, we have that $F(\rho * \tau^R * \sigma, \rho * \tau^R * \sigma * \tau) < F(\pi * \sigma, \pi * \sigma * \tau)$. Hence, $\pi' \in U$ by transitivity.

We now examine the case that $\pi \in V$. From above, we see that $\pi' * \sigma * \tau = \pi * \sigma * \tau$, while $\pi' * \sigma$ and $\pi * \sigma$ differ only by swapping the adjacent elements i and j . Since $\pi' * \sigma(i) > \pi' * \sigma(j)$ while $\pi' * \sigma * \tau(i) < \pi' * \sigma * \tau(j)$, we see that there is exactly one more inversion between $\pi' * \sigma$ and $\pi' * \sigma * \tau$ than between $\pi * \sigma$ and $\pi * \sigma * \tau$. That is, $K(\pi' * \sigma, \pi' * \sigma * \tau) < K(\pi * \sigma, \pi * \sigma * \tau)$. By our assumption, $\pi \in V$, hence $K(\rho * \tau^R * \sigma, \rho * \tau^R * \sigma * \tau) < K(\pi * \sigma, \pi * \sigma * \tau)$. It follows that $\pi' \in V$.

So in each case, we have produced a $\pi' \in S$ such that $K(\rho * \tau^R, \pi') < K(\rho * \tau^R, \pi)$, contradicting the minimality of π . Hence, S must have been empty, as we wanted. \square

Combining the previous two lemmas, we obtain a complete characterization of the Hausdorff distance.

THEOREM 7. *Let σ and τ be partial rankings, let σ^R be the reverse of σ , and let τ^R be the reverse of τ . Let ρ be any full ranking. Then*

$$\begin{aligned} F_{\text{Haus}}(\sigma, \tau) &= \max\{ F(\rho * \tau^R * \sigma, \rho * \sigma * \tau), \\ &\quad F(\rho * \tau * \sigma, \rho * \sigma^R * \tau) \} \\ K_{\text{Haus}}(\sigma, \tau) &= \max\{ K(\rho * \tau^R * \sigma, \rho * \sigma * \tau), \\ &\quad K(\rho * \tau * \sigma, \rho * \sigma^R * \tau) \} \end{aligned}$$

PROOF. We prove it for F_{Haus} . The proof for K_{Haus} is analogous. Recall that

$$F_{\text{Haus}}(\sigma, \tau) = \max\left\{ \max_{\sigma} \min_{\tau} F(\sigma, \tau), \max_{\tau} \min_{\sigma} F(\sigma, \tau) \right\}$$

where throughout this proof, σ and τ range through all full refinements of σ and τ , respectively. We show $\max_{\sigma} \min_{\tau} F(\sigma, \tau) = F(\rho * \tau^R * \sigma, \rho * \sigma * \tau)$. The fact that $\max_{\tau} \min_{\sigma} F(\sigma, \tau) = F(\rho * \tau * \sigma, \rho * \sigma^R * \tau)$ follows similarly.

Think of $\sigma \succeq \sigma$ as fixed. Then by Lemma 5, the quantity $F(\sigma, \tau)$, where τ ranges over all full refinements of τ , is minimized when $\tau = \sigma * \tau$. That is, $\min_{\tau} F(\sigma, \tau) = F(\sigma, \sigma * \tau)$.

By Lemma 6, the quantity $F(\sigma, \sigma * \tau)$, where σ ranges over all full refinements of σ , is maximized when $\sigma = \rho * \tau^R * \sigma$. Hence, $\max_{\sigma} \min_{\tau} F(\sigma, \tau) = F(\rho * \tau^R * \sigma, \rho * \tau^R * \sigma * \tau)$. Since $\rho * \tau^R * \sigma * \tau = \rho * \sigma * \tau$, we have $\max_{\sigma} \min_{\tau} F(\sigma, \tau) = F(\rho * \tau^R * \sigma, \rho * \sigma * \tau)$, as we wanted. \square

Let σ and τ be partial rankings. Theorem 7 gives us a simple algorithm for computing $F_{\text{Haus}}(\sigma, \tau)$ and $K_{\text{Haus}}(\sigma, \tau)$: we simply pick an arbitrary full ranking ρ and do the computations given in Theorem 7. Let $\sigma_1 = \rho * \tau^R * \sigma$, let $\tau_1 = \rho * \sigma * \tau$, let $\sigma_2 = \rho * \tau * \sigma$, and let $\tau_2 = \rho * \sigma^R * \tau$. Theorem 7 tells us that $F_{\text{Haus}}(\sigma, \tau) = \max\{F(\sigma_1, \tau_1), F(\sigma_2, \tau_2)\}$ and $K_{\text{Haus}}(\sigma, \tau) = \max\{K(\sigma_1, \tau_1), K(\sigma_2, \tau_2)\}$. It is interesting that the same pairs, namely (σ_1, τ_1) and (σ_2, τ_2) are the candidates for exhibiting the Hausdorff distance for both F and K . Note that the only role that the arbitrary full ranking ρ plays is to arbitrarily break ties (in the same way for σ and τ) for pairs (i, j) of distinct elements that are in the same bucket in both σ and τ . A way to describe the pair (σ_1, τ_1) intuitively is: break the ties in σ based on the reverse of the ordering in τ , break the ties in τ based on the ordering in σ , and break any remaining ties arbitrarily (but in the same way in both). A similar description applies to the pair (σ_2, τ_2) .

The algorithm we have described for computing $F_{\text{Haus}}(\sigma, \tau)$ and $K_{\text{Haus}}(\sigma, \tau)$ is based on creating pairs (σ_1, τ_1) and (σ_2, τ_2) ,

one of which must exhibit the Hausdorff distance. The next proposition gives a direct algorithm for computing $K_{\text{Haus}}(\sigma, \tau)$, that we make use of later.

PROPOSITION 8. *Let σ and τ be partial rankings. Let S be the set of pairs $\{i, j\}$ of distinct elements such that i and j appear in the same bucket of σ but in different buckets of τ , let T be the set of pairs $\{i, j\}$ of distinct elements such that i and j appear in the same bucket of τ but in different buckets of σ , and let U be the set of pairs $\{i, j\}$ of distinct elements that are in different buckets of both σ and τ and are in a different order in σ and τ . Then $K_{\text{Haus}}(\sigma, \tau) = |U| + \max\{|S|, |T|\}$.*

PROOF. As before, let $\sigma_1 = \rho * \tau^{\text{R}} * \sigma$, let $\tau_1 = \rho * \sigma * \tau$, let $\sigma_2 = \rho * \tau * \sigma$, and let $\tau_2 = \rho * \sigma^{\text{R}} * \tau$. It is straightforward to see that the set of pairs $\{i, j\}$ of distinct elements that are in a different order in σ_1 and τ_1 is exactly the union of the disjoint sets U and S . Therefore, $K(\sigma_1, \tau_1) = |U| + |S|$. Identically, the set of pairs $\{i, j\}$ of distinct elements that are in a different order in σ_2 and τ_2 is exactly the union of the disjoint sets U and T , and hence $K(\sigma_2, \tau_2) = |U| + |T|$. But by Theorem 7, we know that $K_{\text{Haus}}(\sigma, \tau) = \max\{K(\sigma_1, \tau_1), K(\sigma_2, \tau_2)\} = \max\{|U| + |S|, |U| + |T|\}$. The result follows immediately. \square

5. EQUIVALENCE BETWEEN METRICS

In this section we show that all of our metrics are in the same equivalence class. The following theorem is proved in three parts, as Theorem 11 (Section 5.1), Theorem 15 (Section 5.2), and Theorem 16 (Section 5.3).

THEOREM 9. *The metrics F_{prof} , K_{prof} , F_{Haus} , and K_{Haus} are all in the same equivalence class.*

As we discussed earlier, the above theorem shows that our metrics are quite robust. The equivalence will come in handy when we design aggregation algorithms for partial rankings in Section 6.

5.1 Equivalence of F_{Haus} and K_{Haus}

In this section, we prove the simple result that the Diaconis–Graham inequalities (1) extend to F_{Haus} and K_{Haus} . We begin with a lemma. In this lemma, for metric d , we define d_{Haus} as in (2), and similarly for metric d' .

LEMMA 10. *Assume that d and d' are metrics where there is a constant c such that $d \leq c \cdot d'$. Then $d_{\text{Haus}} \leq c \cdot d'_{\text{Haus}}$.*

PROOF. Let A and B be as in (2). Assume without loss of generality that $d_{\text{Haus}}(A, B) = \max_{\gamma_1 \in A} \min_{\gamma_2 \in B} d(\gamma_1, \gamma_2)$. Find γ_1 in A that maximizes $\min_{\gamma_2 \in B} d(\gamma_1, \gamma_2)$, and γ_2 in B that minimizes $d(\gamma_1, \gamma_2)$. Therefore, $d_{\text{Haus}}(A, B) = d(\gamma_1, \gamma_2)$. Find γ'_2 in B that minimizes $d'(\gamma_1, \gamma'_2)$. (There is such a γ'_2 since by assumption on the definition of Hausdorff distance, A and B are finite sets.) Then $d_{\text{Haus}}(A, B) = d(\gamma_1, \gamma_2) \leq d(\gamma_1, \gamma'_2)$, since γ_2 minimizes $d(\gamma_1, \gamma_2)$. Also $d(\gamma_1, \gamma'_2) \leq c \cdot d'(\gamma_1, \gamma'_2)$, by assumption on d and d' . Finally $c \cdot d'(\gamma_1, \gamma'_2) \leq c \cdot d'_{\text{Haus}}(A, B)$, by definition of d'_{Haus} and the fact that γ'_2 minimizes $d'(\gamma_1, \gamma'_2)$. Putting these inequalities together, we obtain $d_{\text{Haus}}(A, B) \leq c \cdot d'_{\text{Haus}}(A, B)$, which completes the proof. \square

We can now show that the Diaconis–Graham inequalities (1) extend to F_{Haus} and K_{Haus} .

THEOREM 11. *Let σ_1 and σ_2 be partial rankings. Then $K_{\text{Haus}}(\sigma_1, \sigma_2) \leq F_{\text{Haus}}(\sigma_1, \sigma_2) \leq 2K_{\text{Haus}}(\sigma_1, \sigma_2)$.*

PROOF. The first inequality $K_{\text{Haus}}(\sigma_1, \sigma_2) \leq F_{\text{Haus}}(\sigma_1, \sigma_2)$ follows from the first Diaconis–Graham inequality $K(\sigma_1, \sigma_2) \leq F(\sigma_1, \sigma_2)$ and Lemma 10, where we let the roles of d, d' and c be played by K, F , and 1 respectively. The second inequality $F_{\text{Haus}}(\sigma_1, \sigma_2) \leq 2K_{\text{Haus}}(\sigma_1, \sigma_2)$ follows from the second Diaconis–Graham inequality $F(\sigma_1, \sigma_2) \leq 2K(\sigma_1, \sigma_2)$ and Lemma 10, where we let the roles of d, d' and c be played by F, K , and 2 respectively. \square

5.2 Equivalence of F_{prof} and K_{prof}

In order to generalize the Diaconis–Graham inequalities to F_{prof} and K_{prof} , we convert a pair of partial rankings into full rankings in such a way that both the F_{prof} and K_{prof} distances between the partial rankings is precisely 4 times the F and K distances between the full rankings, respectively. Given a partial ranking, σ , with domain D , produce a duplicate set, $D^\sharp = \{i^\sharp : i \in D\}$. Further, produce a new partial ranking, σ^\sharp , with domain $D \cup D^\sharp$ defined by $\sigma^\sharp(i) = \sigma^\sharp(i^\sharp) = 2\sigma(i) - 1/2$ for all $i \in D$.

It is easy to see that σ^\sharp is a well-defined partial ranking. Further, it is not hard to check that for any partial ranking τ ,

$$\begin{aligned} F_{\text{prof}}(\sigma^\sharp, \tau^\sharp) &= 4F_{\text{prof}}(\sigma, \tau) \\ K_{\text{prof}}(\sigma^\sharp, \tau^\sharp) &= 4K_{\text{prof}}(\sigma, \tau) \end{aligned}$$

In order for us to prove our theorem, we still need to convert σ^\sharp from a partial ranking to a full ranking. For any full ranking π with domain D , define a full ranking π^\sharp with domain $D \cup D^\sharp$ as follows:

$$\begin{aligned} \pi^\sharp(d) &= \pi(d) \text{ for all } d \in D \\ \pi^\sharp(d^\sharp) &= 2|D| + 1 - \pi(d) \text{ for all } d \in D \end{aligned}$$

so that π^\sharp ranks elements of D in the same order as π , elements of D^\sharp in the reverse order of π , and all elements of D before all elements of D^\sharp .

We define $\sigma_\pi = \pi^\sharp * (\sigma^\sharp)$. For instance, suppose \mathcal{B} is a bucket of σ^\sharp containing the items $a, b, c, a^\sharp, b^\sharp, c^\sharp$, and suppose that π orders the items $\pi(a) < \pi(b) < \pi(c)$. Then σ_π will contain the sequence $a, b, c, c^\sharp, b^\sharp, a^\sharp$. Also notice that $\frac{1}{2}(\sigma_\pi(a) + \sigma_\pi(a^\sharp)) = \frac{1}{2}(\sigma_\pi(b) + \sigma_\pi(b^\sharp)) = \frac{1}{2}(\sigma_\pi(c) + \sigma_\pi(c^\sharp)) = \text{pos}(\mathcal{B})$. In fact, because of this “reflected-duplicate” property, we see that in general, for any $d \in D$,

$$\frac{1}{2}(\sigma_\pi(d) + \sigma_\pi(d^\sharp)) = \sigma^\sharp(d) = \sigma^\sharp(d^\sharp) = 2\sigma(d) - 1/2 \quad (4)$$

The following lemma shows that no matter what order π we choose, the Kendall distance between σ_π and τ_π is exactly 4 times the K_{prof} distance between σ and τ .

LEMMA 12. *Let σ, τ be partial rankings, and let π be any full ranking on the same domain. Then $K(\sigma_\pi, \tau_\pi) = 4K_{\text{prof}}(\sigma, \tau)$.*

PROOF. By cases. \square

Notice that Lemma 12 holds for any choice of π . The analogous statement is not true for F_{prof} . In that case, we need to choose π specifically for the pair of partial rankings we are given. In particular, we need to avoid a property we call “nesting.”

Given fixed σ, τ , we say that an element $d \in D$ is *nested* with respect to π if either

$$\begin{aligned} [\sigma_\pi(d), \sigma_\pi(d^\sharp)] &\sqsubset [\tau_\pi(d), \tau_\pi(d^\sharp)] \\ \text{or } [\tau_\pi(d), \tau_\pi(d^\sharp)] &\sqsubset [\sigma_\pi(d), \sigma_\pi(d^\sharp)] \end{aligned}$$

where the notation $[s, t] \sqsubset [u, v]$ for integers s, t, u, v means that $[s, t] \subseteq [u, v]$ and $s \neq u, t \neq v$. It is sometimes convenient to write $[u, v] \supset [s, t]$ for $[s, t] \sqsubset [u, v]$.

The following lemma shows us why we want to avoid nesting.

LEMMA 13. *Given partial rankings σ, τ and full ranking π , suppose that there are no elements that are nested with respect to π . Then $F(\sigma_\pi, \tau_\pi) = 4F_{\text{prof}}(\sigma, \tau)$.*

PROOF. Let $d \in D$. Since d is not nested with respect to π , either

$$\begin{aligned} & \sigma_\pi(d) \leq \tau_\pi(d) \text{ and } \sigma_\pi(d^\sharp) \leq \tau_\pi(d^\sharp) \\ \text{or} & \quad \sigma_\pi(d) \geq \tau_\pi(d) \text{ and } \sigma_\pi(d^\sharp) \geq \tau_\pi(d^\sharp) \end{aligned}$$

In either case, we see

$$\begin{aligned} & |\sigma_\pi(d) - \tau_\pi(d)| + |\sigma_\pi(d^\sharp) - \tau_\pi(d^\sharp)| \\ &= |\sigma_\pi(d) - \tau_\pi(d) + \sigma_\pi(d^\sharp) - \tau_\pi(d^\sharp)| \end{aligned}$$

But recall that $\frac{1}{2}(\sigma_\pi(d) + \sigma_\pi(d^\sharp)) = 2\sigma(d) - 1/2$ and similarly for τ_π . Substituting gives us

$$|\sigma_\pi(d) - \tau_\pi(d)| + |\sigma_\pi(d^\sharp) - \tau_\pi(d^\sharp)| = 4|\sigma(d) - \tau(d)|$$

Hence,

$$\begin{aligned} F(\sigma_\pi, \tau_\pi) &= \sum_{d \in D} (|\sigma_\pi(d) - \tau_\pi(d)| + |\sigma_\pi(d^\sharp) - \tau_\pi(d^\sharp)|) \\ &= \sum_{d \in D} 4|\sigma(d) - \tau(d)| \\ &= 4F_{\text{prof}}(\sigma, \tau). \quad \square \end{aligned}$$

In the proof of the following lemma, we show that in fact, there is always a full ranking π with no nested elements.

LEMMA 14. *Let σ, τ be partial rankings. Then there exists a full ranking π on the same domain such that $F(\sigma_\pi, \tau_\pi) = 4F_{\text{prof}}(\sigma, \tau)$.*

PROOF. We produce a full ranking π that has no nested elements. For any full ranking ρ , we say its *first nest* is $\min_d \pi(d)$, where d is allowed to range over all nested elements of ρ ; we say its first nest is ∞ if ρ has no nests. Choose π so that its first nest is as large as possible.

If π has no nested elements, then we are done. Otherwise, let a be the element such that $\pi(a)$ is the first nest of π . By definition, a is nested. Without loss of generality, assume that $[\sigma_\pi(a), \sigma_\pi(a^\sharp)] \supset [\tau_\pi(a), \tau_\pi(a^\sharp)]$. We find $b \in D$ so that $\pi(a) < \pi(b)$, and swapping a and b in π will leave b unnested. To this end, let

$$S_1 = \left\{ d \in D - \{a\} \mid [\sigma_\pi(a), \sigma_\pi(a^\sharp)] \supset [\sigma_\pi(d), \sigma_\pi(d^\sharp)] \right\}$$

$$S_2 = \left\{ d \in D - \{a\} \mid [\sigma_\pi(a), \sigma_\pi(a^\sharp)] \supset [\tau_\pi(d), \tau_\pi(d^\sharp)] \right\}$$

Choose $b \in S_1 - S_2$. To see such a b exists, note that $|S_1| = \frac{1}{2} |[\sigma_\pi(a), \sigma_\pi(a^\sharp)]| - 1$, while $|S_2| \leq \frac{1}{2} |[\sigma_\pi(a), \sigma_\pi(a^\sharp)]| - 2$, since $[\sigma_\pi(a), \sigma_\pi(a^\sharp)] \supset [\tau_\pi(a), \tau_\pi(a^\sharp)]$ but a is not counted in S_2 . Note that $b \in S_1$ implies a and b are in the same bucket for σ . It further implies that $\pi(a) < \pi(b)$.

Furthermore, a and b are in different buckets for τ . To see this, suppose that a and b were in the same bucket for τ . Then since $\pi(a) < \pi(b)$, we would have $\tau_\pi(a) < \tau_\pi(b)$ and $\tau_\pi(a^\sharp) > \tau_\pi(b^\sharp)$. That is, $[\tau_\pi(a), \tau_\pi(a^\sharp)] \supset [\tau_\pi(b), \tau_\pi(b^\sharp)]$. But a is nested, so by our assumption, $[\sigma_\pi(a), \sigma_\pi(a^\sharp)] \supset [\tau_\pi(a), \tau_\pi(a^\sharp)] \supset [\tau_\pi(b), \tau_\pi(b^\sharp)]$. This contradicts the fact that $b \notin S_2$. Hence, a and b must be in different buckets for τ .

Now, produce π' by swapping a and b in π . Since $\pi(a) < \pi(b)$, we see $\pi'(b) = \pi(a) < \pi(b) = \pi'(a)$. We wish to prove that the

first nest for π' is larger than the first nest for π , giving a contradiction. We do so by showing that b is unnested for π' and further, that d is unnested for π' for all $d \in D$ such that $\pi'(d) < \pi'(b) = \pi(a)$. In order to prove this, we need to examine the affect of swapping a and b in π .

To this end, consider a bucket \mathcal{B} of σ . Let $\pi|_{\mathcal{B}}$ denote the order that π induces on \mathcal{B} . Since $\pi'(d) = \pi(d)$ for all d such that $\pi(d) < \pi(a)$, we see that $\pi'|_{\mathcal{B}}(d) = \pi|_{\mathcal{B}}(d)$ for all such d . Hence, $\sigma_{\pi'}(d) = \sigma_\pi(d)$ and $\sigma_{\pi'}(d^\sharp) = \sigma_\pi(d^\sharp)$ for all such d . Therefore, for all d such that $\pi(d) < \pi(a)$

$$[\sigma_{\pi'}(d), \sigma_{\pi'}(d^\sharp)] = [\sigma_\pi(d), \sigma_\pi(d^\sharp)] \quad (5)$$

Let \mathcal{B} be the bucket of σ that contains a and b . Then $\pi'|_{\mathcal{B}}$ is just $\pi|_{\mathcal{B}}$ with a and b swapped. So $\pi'|_{\mathcal{B}}(b) = \pi|_{\mathcal{B}}(a)$. Hence, $\sigma_{\pi'}(b) = \sigma_\pi(a)$ and $\sigma_{\pi'}(b^\sharp) = \sigma_\pi(a^\sharp)$. That is,

$$[\sigma_{\pi'}(b), \sigma_{\pi'}(b^\sharp)] = [\sigma_\pi(a), \sigma_\pi(a^\sharp)] \quad (6)$$

We now consider a bucket \mathcal{B} of τ . Arguing as we did for buckets of σ , we have that for all d such that $\pi(d) < \pi(a)$,

$$[\tau_{\pi'}(d), \tau_{\pi'}(d^\sharp)] = [\tau_\pi(d), \tau_\pi(d^\sharp)] \quad (7)$$

Now, let \mathcal{B} be the bucket of τ that contains a . Since π and π' differ only by swapping a and b , and $\pi'(a) > \pi(a)$, we see that $\pi'|_{\mathcal{B}}(a) \geq \pi|_{\mathcal{B}}(a)$. Hence, $\tau_{\pi'}(a) \geq \tau_\pi(a)$ and $\tau_{\pi'}(a^\sharp) \leq \tau_\pi(a^\sharp)$. That is,

$$[\tau_{\pi'}(a), \tau_{\pi'}(a^\sharp)] \subseteq [\tau_\pi(a), \tau_\pi(a^\sharp)] \quad (8)$$

Finally, let \mathcal{B} be the bucket of τ that contains b . Since π and π' differ only by swapping a and b , and $\pi'(b) < \pi(b)$, we see that $\pi'|_{\mathcal{B}}(b) \leq \pi|_{\mathcal{B}}(b)$. Hence, $\tau_{\pi'}(b) \leq \tau_\pi(b)$ and $\tau_{\pi'}(b^\sharp) \geq \tau_\pi(b^\sharp)$. That is,

$$[\tau_{\pi'}(b), \tau_{\pi'}(b^\sharp)] \supseteq [\tau_\pi(b), \tau_\pi(b^\sharp)] \quad (9)$$

We are now ready to prove the lemma. From (5) and (7), we see that d remains unnested for all d such that $\pi'(d) < \pi(a) = \pi'(b)$. So we only need to show that b is unnested for π' to finish the proof.

If b were nested for π' , then either $[\sigma_{\pi'}(b), \sigma_{\pi'}(b^\sharp)] \supset [\tau_{\pi'}(b), \tau_{\pi'}(b^\sharp)]$ or $[\tau_{\pi'}(b), \tau_{\pi'}(b^\sharp)] \supset [\sigma_{\pi'}(b), \sigma_{\pi'}(b^\sharp)]$. First, suppose that $[\sigma_{\pi'}(b), \sigma_{\pi'}(b^\sharp)] \supset [\tau_{\pi'}(b), \tau_{\pi'}(b^\sharp)]$. Then

$$\begin{aligned} [\sigma_\pi(a), \sigma_\pi(a^\sharp)] &= [\sigma_{\pi'}(b), \sigma_{\pi'}(b^\sharp)] \text{ from (6)} \\ &\supset [\tau_{\pi'}(b), \tau_{\pi'}(b^\sharp)] \text{ by supposition} \\ &\supseteq [\tau_\pi(b), \tau_\pi(b^\sharp)] \text{ from (9)} \end{aligned}$$

But this contradicts the fact that $b \notin S_2$. Now, suppose that $[\tau_{\pi'}(b), \tau_{\pi'}(b^\sharp)] \supset [\sigma_{\pi'}(b), \sigma_{\pi'}(b^\sharp)]$. Then

$$\begin{aligned} [\tau_{\pi'}(b), \tau_{\pi'}(b^\sharp)] &\supset [\sigma_{\pi'}(b), \sigma_{\pi'}(b^\sharp)] \text{ by supposition} \\ &= [\sigma_\pi(a), \sigma_\pi(a^\sharp)] \text{ from (6)} \\ &\supset [\tau_\pi(a), \tau_\pi(a^\sharp)] \text{ } a \text{ is nested, by assumption} \\ &\supseteq [\tau_{\pi'}(a), \tau_{\pi'}(a^\sharp)] \text{ from (8)} \end{aligned}$$

But this implies that a and b are in the same bucket for τ , a contradiction. Hence, b must not be nested for π' .

Hence, if any element d is nested for π' , it must be the case that $\pi'(d) > \pi'(b) = \pi(a)$. That is, the first nest for π' is larger than the first nest for π , contradicting our choice of π . Therefore, π must have had no nested elements. By Lemma 13, $F(\sigma_\pi, \tau_\pi) = 4F_{\text{prof}}(\sigma, \tau)$, as we wanted. \square

Putting these two lemmas together, we conclude the following.

THEOREM 15. *Let σ and τ be partial rankings. Then $K_{\text{prof}}(\sigma, \tau) \leq F_{\text{prof}}(\sigma, \tau) \leq 2K_{\text{prof}}(\sigma, \tau)$.*

PROOF. Given σ and τ , let π be the full ranking guaranteed in Lemma 14. Then we have

$$\begin{aligned} K_{\text{prof}}(\sigma, \tau) &= 4K(\sigma_\pi, \tau_\pi) \text{ by Lemma 12} \\ &\leq 4F(\sigma_\pi, \tau_\pi) \text{ from Diaconis–Graham} \\ &= F_{\text{prof}}(\sigma, \tau) \text{ by Lemma 14} \end{aligned}$$

And similarly,

$$\begin{aligned} F_{\text{prof}}(\sigma, \tau) &= 4F(\sigma_\pi, \tau_\pi) \text{ by Lemma 14} \\ &\leq 8K(\sigma_\pi, \tau_\pi) \text{ from Diaconis–Graham} \\ &= 2K_{\text{prof}}(\sigma, \tau) \text{ by Lemma 12. } \quad \square \end{aligned}$$

5.3 Equivalence of K_{Haus} and K_{prof}

We now show that K_{Haus} and K_{prof} are in the same equivalence class.

THEOREM 16. *Let σ_1 and σ_2 be partial rankings. Then $K_{\text{prof}}(\sigma_1, \sigma_2) \leq K_{\text{Haus}}(\sigma_1, \sigma_2) \leq 2K_{\text{prof}}(\sigma_1, \sigma_2)$.*

PROOF. As in Proposition 8 (but where we let σ_1 play the role of σ and σ_2 play the role of τ), let S be the set of pairs $\{i, j\}$ of distinct elements such that i and j appear in the same bucket of σ_1 but in different buckets of σ_2 , let T be the set of pairs $\{i, j\}$ of distinct elements such that i and j appear in the same bucket of σ_2 but in different buckets of σ_1 , and let U be the set of pairs $\{i, j\}$ of distinct elements that are in different buckets of both σ_1 and σ_2 and are in a different order in σ_1 and σ_2 . By Proposition 8, we know that $K_{\text{Haus}}(\sigma_1, \sigma_2) = |U| + \max\{|S|, |T|\}$. It follows from the definition of K_{prof} that $K_{\text{prof}}(\sigma_1, \sigma_2) = |U| + \frac{1}{2}|S| + \frac{1}{2}|T|$. The theorem now follows from the straightforward inequalities $|U| + \frac{1}{2}|S| + \frac{1}{2}|T| \leq |U| + \max\{|S|, |T|\} \leq 2(|U| + \frac{1}{2}|S| + \frac{1}{2}|T|)$. \square

This concludes the proof that all our metrics are in the same equivalence class.

6. AGGREGATING PARTIAL RANKINGS

In this section, we show how to use median-based algorithms to achieve constant factor approximations for rank aggregation. Since all of our metrics are in the same equivalence class, a constant factor approximation for rank aggregation according to one of our metrics is a constant factor approximation for rank aggregation according to all of our metrics. Our theorems will be stated in terms of L_1 , which gives the F_{prof} metric.

Given a list a_1, \dots, a_m of numbers, we let $\text{median}(a_1, \dots, a_m)$ be the set of values that would typically be taken as the median of the list (note that if m is odd, it is a set containing just one number). More precisely, suppose that the a_i 's are relabeled so that $a_1 \leq \dots \leq a_m$. Then $\text{median}(a_1, \dots, a_m)$ is the set $\left\{a_{\frac{m+1}{2}}\right\}$ when m is odd, and the set $\left\{a_{\frac{m}{2}}, a_{\frac{m}{2}+1}, (a_{\frac{m}{2}} + a_{\frac{m}{2}+1})/2\right\}$ when m is even. Given a list f_1, \dots, f_m of functions, each mapping $D \rightarrow \mathcal{R}$, we abuse notation slightly and define $\text{median}(f_1, \dots, f_m)$ to be the set of valid median functions. More precisely, we define $\text{median}(f_1, \dots, f_m)$ to be

$$\{f : D \rightarrow \mathcal{R} \mid f(d) \in \text{median}(f_1(d), \dots, f_m(d)), \forall d \in D\}$$

We shall show that aggregation methods based on ranking objects according to their median ranks, and breaking ties arbitrarily when needed, provide constant factor approximations for rank aggregation.

We now recall the merits of median as an aggregation operator in the context of databases, as discussed in [10]. In [10], the median rank aggregation algorithm was implemented by using two cursors for each attribute to implicitly rank the database objects with respect to the query without having to sort for every query. This ensures that the data is accessed in a localized and pre-defined order, without any random access or extra storage, thereby permitting extremely efficient implementations. In fact, this algorithm was shown to be instance-optimal [11]—among the class of algorithms that access the input rankings in sequential order, this algorithm is the best possible algorithm (to within a constant factor) on every input instance. Our results show that we automatically inherit the benefits of the median rank aggregation even for partial rankings.

To see the simplicity of the whole algorithm, here is an instantiation to obtain the top element: access each of the input (partial) rankings, one element at a time, until some database object is seen in more than $m/2$ (i.e., more than half the number of the inputs) times; output this object as the top result of the aggregation.

We now show that the median based algorithm achieves constant factor approximations to the Spearman footrule distance in the following three interesting scenarios; the scenarios are based on what kind of inputs/outputs (partial rankings, full rankings, or top k lists) the aggregation algorithm handles. The proof of these results are presented in a generalized form in Section 6.2.

From partial ranking to top k lists. When the inputs are partial rankings, we show that the median aggregation algorithm produces a top k list that is within a factor of three of the optimum top k list (in fact, we need to run the median aggregation algorithm only long enough to output the first k objects). Noting that a full ranking is actually a top $|D|$ list, we see that this also implies that the median aggregation algorithm produces a near-optimal full ranking.

Given a function $f : D \rightarrow \mathcal{R}$, it naturally defines a partial ranking, denoted \hat{f} , as follows: for all $i, j \in D$, if $f(i) < f(j)$, then set $\hat{f}(i) < \hat{f}(j)$; if $f(i) = f(j)$, then set $\hat{f}(i) = \hat{f}(j)$.

THEOREM 17. *Let $\sigma_1, \sigma_2, \dots, \sigma_m$ be partial rankings. Assume $f \in \text{median}(\sigma_1, \dots, \sigma_m)$, and let σ be a top k list of \hat{f} where ties are broken arbitrarily. Then for every top k list τ ,*

$$\sum_{i=1}^m L_1(\sigma, \sigma_i) \leq 3 \sum_{i=1}^m L_1(\tau, \sigma_i).$$

We remark that the output top k list of Theorem 17 satisfies an even stronger notion of optimality—it is the top k list of a near-optimal partial ranking.

From full rankings to full rankings. When the inputs themselves are full rankings and the output is required to be a full ranking, we show that the median aggregation algorithm gives a near-optimal full ranking, with an approximation factor of two. Note that the question of whether the median based algorithm is a constant factor approximation with respect to Kendall distance was left open in [7, 10]; our work answers this question in the affirmative.

THEOREM 18. *Let $\sigma_1, \dots, \sigma_m$ be full rankings. Assume $f \in \text{median}(\sigma_1, \dots, \sigma_m)$, and let σ be a refinement of \hat{f} where ties are broken arbitrarily. Then $\sum_{i=1}^m L_1(\sigma, \sigma_i) \leq 2 \sum_{i=1}^m L_1(\tau, \sigma_i)$ for every partial ranking τ .*

From partial rankings to partial rankings. In the above discussion, we assumed that the final goal of aggregation is to produce a full ranking (or top k list) that is good when compared against other full rankings (or top k lists). In some applications, it may be desirable (and sufficient) for the aggregation to obtain a partial ranking,

but then the partial ranking should compare well against all partial rankings (and not just full rankings or top k lists). Based on the fact that \hat{f} gives a near-optimal aggregation (after suitable tie-breaking) in both the top k case (Theorem 17) and in the case of full rankings (Theorem 18), we would certainly guess that \hat{f} also gives a near-optimal aggregation among partial rankings. Surprisingly, this is not the case, as the next proposition says.

PROPOSITION 19. *For each constant c , there are partial rankings $\sigma_1, \dots, \sigma_m$ and a partial ranking σ such that for each $f \in \text{median}(\sigma_1, \dots, \sigma_m)$,*

$$\sum_{i=1}^m L_1(\hat{f}, \sigma_i) > c \sum_{i=1}^m L_1(\sigma, \sigma_i).$$

In fact, we can have $c = \Omega(|D|)$ in Proposition 19, where D is the domain of the partial rankings.

Even though \hat{f} does not give us a near-optimal partial ranking, we can obtain another near-optimal partial ranking that is based on median rank. Unfortunately, the algorithm we use to obtain this near-optimal partial ranking cannot be branded database-friendly, as it is based on dynamic programming (described in Section 6.3).

THEOREM 20. *Let $\sigma_1, \dots, \sigma_m$ be partial rankings, and assume $f \in \text{median}(\sigma_1, \dots, \sigma_m)$. Suppose that f^\dagger is a partial ranking such that for all partial rankings τ , we have $L_1(f^\dagger, f) \leq L_1(\tau, f)$. Then for every partial ranking σ , we have*

$$\sum_{i=1}^m L_1(f^\dagger, \sigma_i) \leq 2 \sum_{i=1}^m L_1(\sigma, \sigma_i).$$

Furthermore, an f^\dagger that satisfies $L_1(f^\dagger, f) \leq L_1(\tau, f)$ for all τ can be computed in $O(|D|^2)$ time by dynamic programming.

The rest of this section is devoted to proving these theorems. Section 6.1 develops some basic machinery needed for the proofs. Section 6.2 presents a unified proof of Theorems 17, 18, and 20. Section 6.3 describes the dynamic programming algorithm.

6.1 Basic machinery

We first develop some basic machinery for the proof. The following lemma, which is folklore and was previously noted in [10], shows the importance of the median function for rank aggregation. Basically, it says that median is the best function for minimizing L_1 -norm quantities.

LEMMA 21. *Let f_1, f_2, \dots, f_m be functions mapping $D \rightarrow \mathcal{R}$. Assume $f \in \text{median}(f_1, \dots, f_m)$. Then for every function $g : D \rightarrow \mathcal{R}$,*

$$\sum_{i=1}^m L_1(f, f_i) \leq \sum_{i=1}^m L_1(g, f_i).$$

The following lemma appears to be folklore. Note that Lemma 3 is, in fact, a special case of this lemma.

LEMMA 22. *If A and B are two multisets of numbers of the same size, and the cost of matching $a \in A$ to $b \in B$ is $|a - b|$, then the order-preserving perfect matching (i.e., the matching that matches the i -th largest element of A to the i -th largest element of B) is a minimum cost perfect matching between A and B .*

Given functions $f : D \rightarrow \mathcal{R}$ and $g : D \rightarrow \mathcal{R}$, we say that f and g are *consistent* with each other if there is no pair $i, j \in D$ such that $f(i) < f(j)$ and $g(i) > g(j)$. We now show that this notion is symmetric in the role of f and g . Assume that f and g are

consistent with each other, and there is a pair $i, j \in D$ such that $g(i) < g(j)$ and $f(i) > f(j)$. By reversing the roles of i and j , we obtain a contradiction to the fact that f and g are consistent with each other. Although, as we just showed, this relationship is symmetric, it is not transitive, since the constant function is consistent with all other functions. We define $\langle f \rangle$ to be the set of all *partial rankings* that are consistent with f .

Types. Let $\mathcal{B}_1, \dots, \mathcal{B}_t$ be the buckets of the partial ranking σ in order (thus, $\text{pos}(\mathcal{B}_i) < \text{pos}(\mathcal{B}_j)$ when $i < j$). We define the *type* of σ to be the sequence $|\mathcal{B}_1|, |\mathcal{B}_2|, \dots, |\mathcal{B}_t|$, and denote it by $\text{type}(\sigma)$. For example, if σ is a full ranking, then $\text{type}(\sigma)$ is the sequence $1, 1, \dots, 1$ with the number 1 appearing $|D|$ times; a *top k list* is a partial ranking σ where $\text{type}(\sigma)$ is the sequence $1, 1, \dots, 1, |D| - k$, with the number 1 appearing before $|D| - k$ a total of k times. Given a type α , define $\langle f \rangle_\alpha$ to be the subset of $\langle f \rangle$ consisting of partial rankings with type α .

LEMMA 23. *Let $f : D \rightarrow \mathcal{R}$, let α be a type, and suppose $\sigma \in \langle f \rangle_\alpha$. Then $L_1(\sigma, f) \leq L_1(\tau, f)$ for all partial rankings τ such that $\text{type}(\tau) = \alpha$.*

PROOF. Consider the multisets $A = \{\{\sigma(x) : x \in D\}\}$ and $B = \{\{f(x) : x \in D\}\}$. It is clear from the definition of partial rankings and types that every partial ranking τ of type α corresponds to a perfect matching between D and A . Since there is a one-to-one correspondence between D and B defined by f , every such τ also corresponds to a perfect matching between A and B . Furthermore, the cost of this perfect matching (assuming that the cost of matching $a \in A$ with $b \in B$ is $|a - b|$) is precisely $L_1(\tau, f)$. Thus, by Lemma 22, the minimum value of $L_1(\tau, f)$ is achieved when τ is consistent with f , that is, when it belongs to $\langle f \rangle_\alpha$. Hence, $L_1(\sigma, f) \leq L_1(\tau, f)$. \square

LEMMA 24. *Let $f : D \rightarrow \mathcal{R}$, and let \hat{f} be the partial ranking associated with it. Let σ be a refinement of \hat{f} . Then for every full ranking τ , we have $L_1(\sigma, f) \leq L_1(\tau, f)$.*

PROOF. Assume $\sigma \succeq \hat{f}$, and let σ be a full ranking that is a refinement of σ . We shall show that $L_1(\sigma, f) \leq L_1(\tau, f)$, which implies that $L_1(\sigma, f) \leq L_1(\tau, f)$ by Lemma 23 (both have the same type, namely $1, \dots, 1$ with 1 repeated $|D|$ times).

To this end, let \mathcal{B} be a bucket of σ . Since σ is a refinement of \hat{f} , we see that f is constant over all $i \in \mathcal{B}$; call this value $f_{\mathcal{B}}$. Since σ is a refinement of σ , it follows that $\sum_{i \in \mathcal{B}} \sigma(i) = |\mathcal{B}| \cdot \text{pos}(\mathcal{B})$. So we have

$$\begin{aligned} \sum_{i \in \mathcal{B}} |\sigma(i) - f_{\mathcal{B}}| &\geq \left| \sum_{i \in \mathcal{B}} (\sigma(i) - f_{\mathcal{B}}) \right| \\ &= |\mathcal{B}| \cdot |\text{pos}(\mathcal{B}) - f_{\mathcal{B}}| \\ &= \sum_{i \in \mathcal{B}} |\sigma(i) - f_{\mathcal{B}}|. \end{aligned}$$

Summing the above over all buckets of σ , we see that $L_1(\sigma, f) \geq L_1(\tau, f)$, as we wanted. \square

6.2 Unified proof of Theorems 17, 18, and 20

Theorems 17, 18, and 20 are special cases of the following theorem, as we will show.

THEOREM 25. *Let f_1, \dots, f_m be functions mapping $D \rightarrow \mathcal{R}$, and assume $f \in \text{median}(f_1, \dots, f_m)$. Also, let S be a set of functions (for instance, the set of top k lists, or the set of partial*

rankings). Suppose that f' is a function such that for all functions $g \in S$, we have $L_1(f', f) \leq L_1(g, f)$. Then for all functions $g \in S$, we have

$$\sum_{i=1}^m L_1(f', f_i) \leq 3 \sum_{i=1}^m L_1(g, f_i).$$

If the functions $f_1, \dots, f_m \in S$, then we have for all functions h (not necessarily in S),

$$\sum_{i=1}^m L_1(f', f_i) \leq 2 \sum_{i=1}^m L_1(h, f_i).$$

PROOF.

$$\begin{aligned} \sum_{i=1}^m L_1(f', f_i) &\leq \sum_{i=1}^m (L_1(f', f) + L_1(f, f_i)) \text{ by } \Delta \text{ ineq.} \\ &\leq \sum_{i=1}^m (L_1(g, f) + L_1(f, f_i)) \text{ by assumption} \\ &\leq \sum_{i=1}^m (L_1(g, f_i) + L_1(f_i, f) + L_1(f, f_i)) \text{ by } \Delta \text{ ineq.} \\ &\leq 3 \sum_{i=1}^m L_1(g, f_i) \text{ by Lemma 21} \end{aligned}$$

As for the second part,

$$\begin{aligned} \sum_{i=1}^m L_1(f', f_i) &\leq \sum_{i=1}^m (L_1(f', f) + L_1(f, f_i)) \text{ by } \Delta \text{ ineq.} \\ &\leq \sum_{i=1}^m 2L_1(f, f_i) \text{ by assumption, since each } f_i \in S \\ &\leq 2 \sum_{i=1}^m L_1(h, f_i) \text{ by Lemma 21. } \quad \square \end{aligned}$$

If S is the set of top k lists, then combining Lemma 23 and Theorem 25 gives us Theorem 17. Setting S to be the set of all full rankings in Theorem 25, and using Lemma 24, Theorem 18 follows. Finally, setting S to be the set of all partial rankings in Theorem 25, Theorem 20 is immediate once we are given f^\dagger .

6.3 The dynamic programming algorithm

Let $|D| = n$. We now describe an algorithm that given a function $f \in \text{median}(f_1, \dots, f_m)$ for partial orders f_1, \dots, f_m , finds a partial ranking f^\dagger so that $L_1(f^\dagger, f)$ is minimized. (Note that the algorithm does not actually need f to be a median function.) Using the recurrence we define, it is easy to produce an algorithm running in time $O(n^2)$ if we are allowed to use $O(n^2)$ space. If we make the additional assumption that $2f(i)$ is integral for all i , then we have an algorithm that runs in linear space and time $O(n^2)$. Note that this assumption is not very restrictive, since the median function for a set of partial orders will always satisfy this when the median does not average two values (For instance, if we have a set of m values $a_1 \leq a_2 \leq \dots \leq a_m$, we take the median value to be $a_{\lfloor \frac{m+1}{2} \rfloor}$.)

Suppose without loss of generality that $f(1) \leq f(2) \leq \dots \leq f(n)$. Let π be the usual total order on $1, 2, \dots, n$. Let τ be the partial ranking that minimizes $L_1(\tau, f)$, and let α be its type. By Lemma 23, we know that if $f^\dagger \in \langle \pi \rangle_\alpha$ then $L_1(f^\dagger, f) \leq L_1(\tau, f)$. Hence, to find an optimal partial ranking, we simply need to find an optimal type. We can determine such an optimal type using dynamic programming.

To do so, we first need several definitions. For any i, j with $0 \leq i < j \leq n$, we define

$$c(i, j) = \sum_{\ell=i+1}^j \left| f(\ell) - \frac{i+j+1}{2} \right|$$

To motivate our definition of $c(i, j)$, imagine that we alter the type of π so that there is a bucket starting at $i+1$ and going until j . Then the position of that bucket is $\frac{i+j+1}{2}$, and the distance between that bucket and f (on the values $\{i+1, i+2, \dots, j\}$) is precisely $c(i, j)$.

In general, let S be a sequence of integers s_0, s_1, \dots, s_t where $s_0 < s_1 < \dots < s_t$. (Throughout the rest of our discussion, all sequences will be integral.) Then we define

$$c(S) = \sum_{\ell=0}^{t-1} c(s_\ell, s_{\ell+1})$$

Intuitively, we think of each s_ℓ as marking a point where one bucket ends and the next begins. Put another way, suppose $s_0 = 0$ and $s_t = n$. Then let β be the type defined by the sequence $s_1 - s_0, s_2 - s_1, \dots, s_t - s_{t-1}$. We see immediately that if $\tau \in \langle \pi \rangle_\beta$ then

$$L_1(f, \tau) = \sum_{\ell=0}^{t-1} c(s_\ell, s_{\ell+1}) = c(S)$$

Conversely, it is not hard to see similarly that every type corresponds to an increasing integral sequence whose first element is 0 and whose last is n . Hence, to find an optimal type, our dynamic programming algorithm will instead find the corresponding sequence that minimizes $c(\cdot)$.

To this end, we find $n+1$ different integral sequences $S_0, S_1, S_2, \dots, S_n$. For all $j > 0$, the sequence S_j will have the property that its first element is 0, and its last element is j . Our goal is to have $c(S_j)$ minimal over all such sequences.

We define $S_0 = 0$, and recursively define $S_j = S_{i_0, j}$, where $i_0 = \text{argmin}_i [c(S_i) + c(i, j)]$. Then we have the following.

LEMMA 26. *Let S_0, S_1, \dots, S_n be defined as above. Then for all j and for all strictly increasing sequences S'_j that start with 0 and end with j , we have $c(S'_j) \geq c(S_j)$.*

PROOF. We proceed by induction. The case $j = 0$ is trivially true. So assume that $j > 0$ and that our claim is true for all indices smaller than j .

Now, let S'_j be a strictly increasing sequence starting with 0 and ending with j . Suppose the penultimate element of S'_j is i . Then there is a strictly increasing sequence S'_i ending with i such that $S'_j = S'_i, j$. By definition, $c(S'_j) = c(S'_i) + c(i, j)$. But by induction, $c(S'_i) \geq c(S_{i_0, i})$. Hence, $c(S'_j) \geq c(S_{i_0, i}) + c(i, j) \geq c(S_j)$. \square

Given the recurrence relation, it is a simple matter to calculate S_n . Since $c(i, j)$ can be calculated in $O(n)$ time for all i, j , we see there is a simple algorithm to calculate S_n in time $O(n^3)$. However, we can in fact calculate $c(i, j)$ in amortized $O(1)$ time. In the case where we do not have memory restrictions, we simply utilize the following recurrence:

$$c(i-1, j+1) = c(i, j) + \left| f(i-1) - \frac{i+j}{2} \right| + \left| f(j+1) - \frac{i+j}{2} \right|$$

Using this, we can calculate $c(i, j)$ for all i, j in $O(n^2)$ time, but $O(n^2)$ space.

If $2f(i)$ is integral for all i , then we can calculate S_n in linear space and $O(n^2)$ time using a slightly more complicated algorithm. The details are omitted.

7. CONCLUSIONS

In this paper we consider metrics between partial rankings, motivated by need for such metrics in various database applications. We define four intuitive and natural metrics between partial rankings. We obtain efficient polynomial time algorithms to compute these metrics. We also show that these metrics are all within constant multiples of each other. Armed with this, we obtain a constant factor approximation algorithm for aggregation with respect to each of the metrics by obtaining a constant factor approximation algorithm with respect to just one of them. Our algorithm is based on median rank and admits very efficient database-friendly implementations.

8. REFERENCES

- [1] J. A. Aslam and M. Montague. Models for metasearch. In *Proceedings of the 24th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 276–284, 2001.
- [2] K. A. Baggerly. *Visual Estimation of Structure in Ranked Data*. PhD thesis, Rice University, 1995.
- [3] W. W. Cohen, R. E. Schapire, and Y. Singer. Learning to order things. *Journal of Artificial Intelligence Research*, 10:243–270, 1999.
- [4] D. E. Critchlow. *Metric Methods for Analyzing Partially Ranked Data*. Number 34 in Lecture Notes in Statistics. Springer-Verlag, 1980.
- [5] P. Diaconis. *Group Representation in Probability and Statistics*. Number 11 in IMS Lecture Series. Institute of Mathematical Statistics, 1988.
- [6] P. Diaconis and R. Graham. Spearman’s footrule as a measure of disarray. *Journal of the Royal Statistical Society, Series B*, 39(2):262–268, 1977.
- [7] C. Dwork, R. Kumar, M. Naor, and D. Sivakumar. Rank aggregation methods for the web. In *Proceedings of the 10th International World Wide Web Conference*, pages 613–622, 2001.
- [8] R. Fagin, R. Kumar, K. McCurley, J. Novak, D. Sivakumar, J. Tomlin, and D. Williamson. Searching the workplace web. In *Proceedings of the 12th International World Wide Web Conference*, pages 366–375, 2003.
- [9] R. Fagin, R. Kumar, and D. Sivakumar. Comparing top k lists. In *Proceedings of the 14th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 28–36, 2003. Full version in *SIAM Journal on Discrete Mathematics*, 17(1):134–160, 2003.
- [10] R. Fagin, R. Kumar, and D. Sivakumar. Efficient similarity search and classification via rank aggregation. In *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data*, pages 301–312, 2003.
- [11] R. Fagin, A. Lotem, and M. Naor. Optimal aggregation algorithms for middleware. In *Proceedings of the 20th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*, pages 102–113, 2001. Full version in *Journal of Computer and System Sciences*, 66(4):614–656, 2003.
- [12] L. A. Goodman and W. H. Kruskal. Measures of association for cross classification. *Journal of the American Statistical Association*, 49:732–764, 1954.
- [13] T. H. Haveliwala, A. Gionis, D. Klein, and P. Indyk. Evaluating strategies for similarity search on the web. In *Proceedings of the 11th International World Wide Web Conference*, pages 432–442, 2002.
- [14] M. Kendall and J. D. Gibbons. *Rank Correlation Methods*. Edward Arnold, 1990.
- [15] M. G. Kendall. The treatment of ties in ranking problems. *Biometrika*, 33(3):239–251, 1945.
- [16] G. Lebanon and J. D. Lafferty. Cranking: Combining rankings using conditional probability models on permutations. In *Proceedings of the 19th International Conference on Machine Learning*, pages 363–370, 2002.
- [17] M. Montague and J. A. Aslam. Condorcet fusion for improved retrieval. In *Proceedings of the 11th International Conference on Information and Knowledge Management*, pages 538–548, 2002.
- [18] M. E. Renda and U. Straccia. Web metasearch: Rank vs. score based rank aggregation methods. In *Proceedings of the 18th Annual Symposium on Applied Computing*, pages 841–846, 2003.
- [19] J. Sese and S. Morishita. Rank aggregation method for biological databases. *Genome Informatics*, 12:506–507, 2001.
- [20] R. R. Yager and V. Kreinovich. On how to merge sorted lists coming from different web search tools. *Soft Computing Research Journal*, 3:83–88, 1999.