

User Rankings from Comparisons: Learning Permutations in High Dimensions

Ioannis Mitliagkas, Aditya Gopalan, Constantine Caramanis, Sriram Vishwanath

Department of Electrical and Computer Engineering

The University of Texas at Austin

Austin, TX 78712

ioannis@utexas.edu, {gadit,caramanis}@mail.utexas.edu, sriram@austin.utexas.edu

Abstract—We consider the problem of learning users’ preferential orderings for a set of items when only a limited number of pairwise comparisons of items from users is available. This problem is relevant in large collaborative recommender systems where overall rankings of users for objects need to be predicted using partial information from simple pairwise item preferences from chosen users. We consider two natural schemes of obtaining pairwise item orderings – random and active (or intelligent) sampling. Under both these schemes, assuming that the users’ orderings are constrained in number, we develop efficient, low-complexity algorithms that reconstruct all the orderings with provably order-optimal sample complexities. Finally, our algorithms are shown to outperform a matrix completion based approach in terms of sample and computational requirements in numerical experiments.

I. INTRODUCTION

Modeling and understanding user ratings based on structure is a recent but well-studied discipline. In this setting, we have n products and m users, and our goal is to determine the overall rating-matrix – which is comprised of ratings each user for each product. The main issue though, is that users only provide us with a subset (possibly random) of ratings, and we must now attempt to learn the remainder of the matrix entries. To this end, structure plays a key role, and low-rank structure is particularly useful in helping complete the overall matrix [1], [2], [3].

In many scenarios, however, the ultimate goal is to understand user *ranking*, with ratings merely being a stepping stone along the way. In other words, we are interested in determining the order in which each of the users would like these products. For example, if the n products were movies, ranking reflects each user’s preference of movies using an integer ordering, with ties broken randomly. Similarly, whenever we have multiple products/brands of the same type (whether they be toasters, washers or restaurants), a rank-ordering of them proves to be an effective representation of their relative merits. Intuitively, a raw rating of 7 out of 10 in the absence of any other information is potentially useless

(Is 7/10 good? or bad? or average?). In standardized tests such as SAT and LSAT, relative performance is captured using a percentile ranking which has now become the gold standard for admissions. Thus, the ranking of a product relative to its peers is valuable information for ultimate user consumption.

As mentioned before, given all ratings, rankings can be obtained by sorting the ratings for each user. However, finding user ratings first and then transforming them to rankings is indirect, and may require much more information and structure than the problem setting allows. This is due to the fact that the range of user ratings can be quite subjective. For example, given a rating scale of 0 to 10, a user can pick her top rating to be 6 and least favorable rating as 4, limiting the actual range of values significantly. As a result, the actual ranking of a product can be considerably different from what the rating indicates when taken out-of-context. In other words, two users with identical rankings of products can have a very different set of ratings. In some settings such as funding-proposal rating, coursework grading etc., we frequently observe a rating/grade *inflation*, where the range of ratings associated with the work being assessed is skewed in favor of a less-punitive scale.

The subjectivity of ratings provided by users also negatively impacts low-rank structure – the basis for the effectiveness of powerful matrix completion techniques in predicting missing ratings.

Motivating Example: Even if all m users in the system have exactly the same ranking for all products, their choices of real-valued ratings can result in a rating matrix that is full rank. Without loss of generality, we can assume the common ranking to be $[1, \dots, n]$. If each user were to generate n real numbers uniformly over $[0, 10]$ and then sort them in descending order, the resulting $m \times n$ matrix will be full rank with high probability. Intuitive justification for low-rank matrix completion techniques originates from the fact that user preferences have only a few degrees of freedom. However, with significant user subjectivity, we expect rankings to cap-

ture similarities in user preferences more effectively than ratings.

Consequently, learning the rankings of a collection of users directly is of primary interest. Indeed, as Weimer et al. [4] argue, “Rating algorithms solve the wrong problem, and one that is actually harder: The absolute value of the rating for an item is highly biased for different users, while the ranking is far less prone to this problem.”

Much of existing work on learning rankings of objects deals with learning a single, “globally appropriate” ordering using preferences from training examples, to minimize a suitable notion of loss ([5], [6]). These include the popular *learning-to-rank* approaches [7], [8] and graph-based learning techniques [9], [10], and on-line permutation learning algorithms and frameworks [11], [12], [13]. Related work on sorting with noise or sorting partially ordered sets can be found in [14], [15], [16].

When a collection of orderings from users is to be learnt, such methods could ideally be applied in a sequential, decoupled fashion to deduce the orderings. However, structure among user orderings, if present, can potentially be exploited to learn the orderings with savings in sample complexity. Researchers have noted that rankings in a population of users often exhibit forms of “low-dimensional” structure – to paraphrase Jagabathula and Shah [17], “Irrespective of the number of candidates in an election, the number of distinct vote rankings that prevail in the population are likely to be few, considering a small set of ‘issues’ influences ranking patterns over candidates.” This inspires the following question when jointly estimating users’ rankings of objects: How can structure among user orderings be effectively leveraged to learn orderings with significantly less effort?

In this work, we study the problem of learning a collection of permutations chosen by m users for n items using only pairwise ordering information. Pairwise sampling asks a user to compare two specified items each time, and is not only a natural choice for attempting to deduce ordering information, but also easy to implement in practical systems. We consider the learning problem under both *random* (i.e., algorithm-independent) and *active* (i.e., algorithm-dependent) pairwise sampling schemes. As a reasonable structural constraint on the space of user permutations, we assume a stochastic model in which the users pick permutations uniformly from a pool of r possible orderings.

For both the random and active sampling schemes, we design efficient, low-complexity algorithms that can reconstruct all the users’ orderings with a guaranteed number of pairwise samples, with high probability. Moreover, we establish, using information-theoretic techniques and concentration results, that the sample-complexity of our

algorithms matches lower bounds on the number of pairwise samples needed by any procedure to learn permutations with high probability, when m , n , and r are large. This shows that these reconstruction algorithms are *order-optimal* – in the sense of sample complexity – for learning users’ rankings from pairwise comparisons. The superior performance of our algorithms for the task of learning user orderings is also borne out in practice in the results of numerical experiments that we report.

Organization: The remainder of the paper is organized as follows. We describe the setup for the problem of learning users’ orderings from pairwise comparisons in Section II. In Section III, we present our algorithms to infer users’ orderings, state performance guarantees and converse results for the learning problem, and discuss the implications of our results. Section IV contains the complete arguments required to prove our main results. Section V presents numerical results for the performance of our approach compared to that of a matrix completion based technique to solve the same problem.

Notation: We let $[n]$ denote the set of all integers from 1 to n . We denote the symmetric group on $[n]$ by \mathcal{S}_n . A permutation $\pi \in \mathcal{S}_n$ is a bijection on $[n]$, and $\pi(i)$ represents the rank of object i . Throughout this paper, we use $N \triangleq \binom{n}{2} = n(n-1)/2$ to denote the number of distinct pairs $(i, j) \in [n] \times [n]$, $i < j$. We can also represent a permutation $\pi \in \mathcal{S}_n$ by a $n \times n$ matrix \mathbf{P}_π such that $P_\pi(i, j) = -1$ if $\pi(i) > \pi(j)$ and $P_\pi(i, j) = +1$ otherwise. Since \mathbf{P}_π is skew-symmetric, a more practical representation is the stacking of its upper triangular entries into a vector $\mathbf{p}_\pi \in \{-1, 1\}^N$. There is a trivial bijection between the two representations, so we use them interchangeably. Throughout, the phrase “with high probability” is used to mean with probability at least $1 - cn^{-1}$ for constant $c > 0$.

II. LEARNING USERS’ ORDERINGS: SETUP

Consider the setup where each one of m users totally orders a set of n objects; we denote the resulting permutation of user $k \in [m]$ by $\pi_k \in \mathcal{S}_n$. The goal is to recover all of these permutations with a small number of *pairwise ordering samples*, i.e. how a user relatively orders a specified pair of objects, from each user. Specifically, let $\mathbf{M} = [\mathbf{p}_{\pi_1} \ \mathbf{p}_{\pi_2} \ \dots \ \mathbf{p}_{\pi_m}]$ be the $N \times m$ matrix of pairwise orderings for all users. The *sampling set* $\Omega \subseteq [N] \times [m]$ denotes the indices of entries of \mathbf{M} we sample, $\mathbf{M}(\Omega)$ denotes the set of all samples acquired, and $s = |\Omega|$ is the number of acquired samples. Sampling can be performed either uniformly at random (*random sampling*) or arbitrarily and adaptively by the algorithm (*active sampling*). In this setup, we are interested in

- Quantifying the minimum *sample complexity* of the learning problem, i.e., the number of samples

required to infer all the users' permutations with high probability, and

- Developing *efficient algorithms* that are *optimal for sample-complexity*, i.e., that successfully recover all permutations drawing the minimum number of samples required.

Model for User Permutations: Without further assumptions on the permutations π_k that all the users choose, the problem of learning all the π_k is in general decoupled. This renders unnecessary anything other than a sequential, independent approach to learn each permutation with pairwise samples. The problem of learning a collection of orderings becomes interesting when we impose structure on these orderings, since we can then hope to exploit the resulting “coupling” between user ordering behavior.

In practice, as noted in the introduction, item orderings across a population of users are likely to be much fewer than all the $n!$ permutations in \mathcal{S}_n . This can be attributed primarily to a small set of underlying “features” that essentially drive the users' preferences. We consider a natural structural model where each user picks her permutation uniformly at random and independently from a common pool of randomly selected permutations. Specifically, we impose a “low-dimensionality” constraint as follows:

Assumption: There exists a set of r permutations $\{\rho_1, \rho_2, \dots, \rho_r\}$, where each ρ_j is drawn independently and uniformly at random from \mathcal{S}_n . Each π_k is drawn from the ρ_j independently and uniformly at random, i.e., $\mathbb{P}(\pi_k = \rho_j) = 1/r \quad \forall k \in [m], j \in [r]$.

We remark that in correspondence with the matrix-completion literature, the assumption above makes the ± 1 matrix of pairwise orderings across all users ($\mathbf{M} = [\mathbf{p}_{\pi_1} \ \mathbf{p}_{\pi_2} \ \dots \ \mathbf{p}_{\pi_m}]$) at most rank r , and thus may be viewed as a surrogate to “low-rank” structure in our permutation-learning setup. The setup is characterized completely by the triple (n, m, r) , and our algorithms and results are expressed chiefly in terms of these parameters.

III. ALGORITHMS, MAIN RESULTS AND IMPLICATIONS

In this section, we present algorithms for recovering (and sampling when permitted) all the permutations under both the random and active sampling models. For each case, we provide rigorous analytical guarantees on the number of samples sufficient for our algorithms to exactly recover all permutations with high probability. This is followed by matching converse results, using information-theoretic source-coding techniques, that establish fundamental lower bounds on the sample-complexity required by *any* algorithm to learn the permutations with a significant probability. We discuss

the implications of our results and comment on their consequences.

A. Learning with Random Pairwise Samples

Suppose that the set of samples Ω is obtained by uniform sampling with replacement from $[N] \times [m]$, i.e., the set of all (object pair, user) combinations. This models the case where, for instance, every user is asked to independently provide pairwise comparisons for a uniformly randomly chosen set of object pairs. The problem is then to use these results to deduce the users' orderings of all the objects. We introduce our first algorithm (Algorithm 1) to learn the permutations given s randomly drawn samples, and show that it recovers all the orderings with high probability given a sufficient number of random samples. In this description, we denote the sampling set by $\Omega_s \subset [N] \times [m]$ to indicate its size s and use $\Omega_{s,u} \subset [N]$ to denote the positions (object pairs) sampled from user $u \in [m]$.

In essence, Algorithm 1 uses the s pairwise samples to first separate each pair of users if there is any *discrepancy* in their sampled comparisons. A discrepancy occurs between two users u and k if their sampled orderings for a pair of objects (i, j) disagree, i.e., $(i, j) \in \Omega_{s,u} \cap \Omega_{s,k}$ and $\mathbf{M}(\Omega_{s,u}, u)_{(i,j)} \neq \mathbf{M}(\Omega_{s,k}, k)_{(i,j)}$. Having “clustered” the users' permutations thus, the algorithm proceeds to completely learn the (presumably correctly clustered) permutations by collecting all pairwise samples from users belonging to each cluster and topologically sorting the resulting Directed Acyclic Graph (DAG).

Our first result concerns the sample-complexity of Algorithm 1:

Theorem 1 (Random Sampling, Algorithm 1). *Suppose $r = \theta(m^\gamma)$ for a fixed $\gamma > 0$. Algorithm 1 recovers all permutations correctly, with high probability, when the number of random samples s is at least $\max\{(12/\gamma)m \log r, 2rN \log n\}$.*

Proof sketch: The two terms given in the bound of Theorem 1 above quantify separately the sample-complexities needed to successfully complete both steps of the algorithm, i.e. *clustering* and *learning*. We first establish a concentration result for the pairwise Hamming distance between two distinct permutations drawn uniformly at random. This allows us to show that $O(\log r)$ random pairwise comparisons per permutation are sufficient to distinguish them. Alongside, for any fixed permutation, we identify a necessary and sufficient condition to exactly learn the permutation, from pairwise samples, in terms of the unique Hamiltonian path of the digraph induced by the permutation. This is used together with a coupon-collecting argument to show that the permutation-learning stage of Algorithm 1 requires $O(N \log n)$ random samples per cluster (i.e. for each

Algorithm 1

Input: Set of sampled positions $\Omega_s \subset [N] \times [m]$ and samples $\mathbf{M}(\Omega_s) \in \{-1, +1\}^s$.

Output: Permutations of all users $(\hat{\pi}_k)_{k=1}^m \in \mathcal{S}_n$.

Stage 1: Clustering:

- 1) Set \mathcal{C} to be an empty collection of clusters.
- 2) Set $\mathcal{U} \leftarrow [m]$, to be the set of all unclustered users.
- 3) If $\mathcal{U} = \emptyset$, go to Stage 2.
- 4) Let $u \leftarrow \min_{k \in \mathcal{U}} k$, set $\mathcal{U} \leftarrow \mathcal{U} \setminus u$ and $\mathcal{L} \leftarrow \{u\}$.
- 5) For every $k \in \mathcal{U}$
 - If $\mathbf{M}(\Omega_{s,u} \cap \Omega_{s,k}, u) = \mathbf{M}(\Omega_{s,u} \cap \Omega_{s,k}, k)$ then set $\mathcal{U} \leftarrow \mathcal{U} \setminus k$ and $\mathcal{L} \leftarrow \mathcal{L} \cup \{k\}$.
- 6) Set $\mathcal{C} \leftarrow \mathcal{C} \cup \{\mathcal{L}\}$ and go to Step 3.

Stage 2: Permutation Learning:

For every cluster $\mathcal{L} \in \mathcal{C}$,

- 1) Let $\Omega_{\mathcal{L}} \leftarrow \bigcup_{k \in \mathcal{L}} \Omega_{s,k}$
 - 2) Let $G = (V, E)$ denote a directed graph, with vertex set $V = [n]$ and edge set $E = \emptyset$.
 - 3) For every sample position p in $\Omega_{\mathcal{L}}$, drawn from user k and corresponding to object pair (i, j)
 - if $\mathbf{M}(p, k) = -1$ then $E \leftarrow E \cup \{(i, j)\}$; else $E \leftarrow E \cup \{(j, i)\}$
 - 4) Set $\hat{\rho}_{\mathcal{L}} \leftarrow \text{TopologicalSorting}(G)$
 - 5) Set $\hat{\pi}_k \leftarrow \hat{\rho}_{\mathcal{L}}$ for all $k \in \mathcal{L}$.
-

ρ_j) to completely infer the cluster. Putting together these estimates gives the theorem. The full details of the proof are provided in Section IV-A.

On the other hand, we establish a converse result on the minimum number of samples needed for successful permutation recovery. For this purpose, consider a general algorithm \mathcal{A} that takes as input s pairwise random samples $\mathbf{M}(\Omega_s)$ and maps it to its output: a possibly random estimate $\hat{\mathbf{M}} \triangleq [p_{\hat{\pi}_1} \ p_{\hat{\pi}_2} \ \dots \ p_{\hat{\pi}_m}]$ of all the permutations, one for each user. We denote the probability of successful reconstruction with \mathcal{A} on s samples by $P_{\text{succ}} = P_{\text{succ}}(\mathcal{A}, s) = \mathbb{P}[\hat{\mathbf{M}} = \mathbf{M}]$.

Theorem 2 (Random Sampling, Lower Bound on Sample Complexity). *For any algorithm \mathcal{A} , if $s < \max\{(m-r) \log r, rN \log n\}$, then $P_{\text{succ}} \rightarrow 0$ as $n \rightarrow \infty$.*

Proof sketch: For the sake of contradiction, assume that an algorithm can learn the orderings with fewer than the claimed number of samples. It follows that this results in an algorithm that performs the easier clustering and permutation learning (given clustering information) tasks with those samples. Hence, it suffices to prove separate converse results for these two stages. A key contribution of this work is to prove a strong converse theorem for clustering. This is carried out by first extending the information-theoretic notion of typicality to

“clusterings”, and using a source-coding proof technique that results in a converse theorem. For the permutation learning stage, an important step is to show that the number of random pairwise samples needed to learn a single permutation with high probability is $\Omega(N \log n)$. For this purpose, the necessary and sufficient Hamiltonian path condition – from the proof of Theorem 1 – that characterizes when a permutation is learnt can be used along with concentration estimates for the lower tail of the coupon-collector problem to prove the result. We defer details of the full proof to Section IV-B.

Note that Theorem 2 is a *strong converse*, i.e., it states that *any* algorithm fails to recover all the users’ item orderings correctly with *overwhelming probability* when the number of random samples drawn is below a threshold.

Implications of Theorems 1 and 2:

- When the number of users m is relatively small, viz. $m = O(rN \log n) = O(rn^2 \log n)$, Algorithm 1 succeeds overwhelmingly with $O(rN \log n)$ samples according to Theorem 1. At the same time, with our standing assumption that $r = \theta(m^\gamma)$ for a fixed γ , the converse Theorem 2 forces at least $\Omega(rN \log n)$ samples to be drawn for correct reconstruction. Thus, in this regime, Algorithm 1 is *order-optimal* for the sample-complexity of the problem. Further, it demands on an average $\frac{r}{m} N \log n$ random samples from each user. If $\gamma < 1$ additionally, this means that each user needs to contribute a *vanishing number* of pairwise comparisons for successful recovery. This represents a significant gain compared to decoupling the learning problem across users and reconstructing each permutation independently (whose net sample complexity is $mN \log n$).
- In general, the best reconstruction algorithm for any number of samples is information-theoretically specified to be the *Maximum-Likelihood (ML)* reconstruction algorithm, i.e., the algorithm that outputs a set of user permutations that maximizes the *a posteriori* probability of permutations given sampled observations. Solving the maximum likelihood problem requires performing a potentially hard combinatorial optimization over the space of all possible user ordering patterns – a computationally infeasible task. However, in the above regime with a relatively small number of users, it is remarkable that the efficient and simple Algorithm 1 achieves the same sample complexity as the ML algorithm for the permutation-learning problem.

B. Learning with Active Pairwise Samples

In many application scenarios, it is often desirable (and possible) to *actively* query users for comparisons

of objects. Thus, the choice of samples could be more intelligent, and we can hope to accomplish the learning task with a *smaller number* of carefully chosen pairwise samples than if we took uncontrolled random pairwise samples. Here, we indeed show that this is the case, and provide a joint sampling and permutation-learning algorithm (Algorithm 2) that is both (a) order-optimal across all learning algorithms, and (b) requires fewer samples than its random-sampling counterpart.

Algorithm 2

Input: Pairwise representation matrix \mathbf{M} ; Number of samples s the algorithm is allowed to use.

Output: Permutations of all users $(\hat{\pi}_k)_{k=1}^m \in \mathcal{S}_n$.

Stage 1: Clustering

- 1) Set $\Omega_{\mathcal{C}}$ to be a random subset of $[N]$, of size $\min\{c \log r, s\}$
- 2) Set \mathcal{C} to be an empty collection of clusters.
- 3) Set $\mathcal{U} \leftarrow [m]$, to be the set of all unclustered users.
- 4) If $\mathcal{U} = \emptyset$, go to Stage 2.
- 5) Let $u \leftarrow \min_{k \in \mathcal{U}} k$, set $\mathcal{U} \leftarrow \mathcal{U} \setminus u$ and $\mathcal{L} \leftarrow \{u\}$.
- 6) For every $k \in \mathcal{U}$
 - If $\mathbf{M}(\Omega_{\mathcal{C}}, u) = \mathbf{M}(\Omega_{\mathcal{C}}, k)$ then set $\mathcal{U} \leftarrow \mathcal{U} \setminus k$ and $\mathcal{L} \leftarrow \mathcal{L} \cup \{k\}$.
- 7) Set $\mathcal{C} \leftarrow \mathcal{C} \cup \{\mathcal{L}\}$ and go to Step 3.

Stage 2: Permutation Learning through Sorting

For every cluster $\mathcal{L} \in \mathcal{C}$,

- 1) Use a sorting algorithm to sort the $[n]$ objects in cluster \mathcal{L} balancing the sample load across all users in the cluster.
 - 2) If sample budget s is reached before completion, stop and declare failure.
-

Algorithm 2 follows the same basic outline of operation as Algorithm 1, i.e., working by clustering and learning each clustered permutation. The key departure here is that it is free to specify pairs of items that it wants compared by certain users, and so it uses this flexibility to cluster and learn permutations faster. Specifically, it first picks a random *common* set of $c \log r$ pairs of objects that it asks *all* r users to order, and uses these samples to cluster users' putative permutations. Once clustering is accomplished, the algorithm pretends that each cluster is a single ordering and attempts to learn the ordering using a standard sorting algorithm on n items (we use Quicksort in our implementation, Section V) that issues pairwise queries to essentially "complete" the permutation. Following the same outline as with Algorithm 1, we first bound the sample complexity of Algorithm 2 as follows.

Theorem 3 (Active Sampling, Algorithm 2). *Algorithm 2 correctly recovers all permutations, with high probabil-*

ity, taking $s = O(m \log r + rn \log n)$ pairwise samples.

Proof sketch: The arguments used to prove Theorem 3 follow the same outline as those for Theorem 1, viz. estimating the number of samples sufficient to perform the two steps of clustering users and learning the clusters. Concentration properties for the Hamming distance between randomly chosen permutations are employed for the estimate of $O(\log r)$ *random common object pairs* which guarantees successful clustering for all users. This is followed by observing that actively sampling and learning a single cluster/permutation is equivalent to sorting items using pairwise comparisons, and standard tail bounds for the performance of a standard sorting algorithm such as Quicksort show that $O(n \log n)$ pairwise samples suffice to learn each cluster with high probability. Putting together these estimates gives the promised bound. The full proof details are provided in Section IV-C.

For the case of active sampling, we provide a matching converse theorem – in the same spirit as Theorem 2 – for the sample complexity of *any* algorithm that is *free* to draw any pairwise samples from the users. Recall from the random sampling scenario, that we denote the probability of successful reconstruction, using algorithm \mathcal{A} on s samples, by $P_{\text{succ}} = P_{\text{succ}}(\mathcal{A}, s)$.

Theorem 4 (Active Sampling, Lower Bound on Sample Complexity). *For any active-sampling algorithm \mathcal{A} , if $s < \max\{(m - r) \log r, rn \log n\}$, then $P_{\text{succ}} \rightarrow 0$ as $n \rightarrow \infty$.*

Proof sketch: The proof for this active sampling converse theorem uses the same high-level outline as that for Theorem 2. The first part – for the clustering stage – is the same information-theoretic source coding argument as in the proof of Theorem 2. For the second part, we use a converse argument for jointly learning a set of r distinct permutations, which essentially generalizes the $\Omega(n \log n)$ pairwise comparison result to sort a set of n item. We reproduce the full details of the proof in Section IV-D.

Implications of Theorems 3 and 4:

- Algorithm 2 achieves *perfect reconstruction* with an *order optimal* number of samples (i.e. $O(m \log r + rn \log n)$). In other words, distinguishing users on the basis of a few ($O(m \log r)$) common pairwise comparisons decouples the overall learning problem *tightly* into r independent “cluster-learning” or sorting problems.
- Compared to the sample complexity of learning with random samples (Theorems 1, 2), Algorithm 2 exhibits a saving in sample complexity of the order of n . This can be directly attributed to the gain in “collaboratively sorting” users clustered together

as the “same”, in the second phase of the algorithm. Also, the sample complexity of Algorithm 2 translates into an average of $\frac{rn}{m} \log n$ samples per user – a gain of the order of r/mn over trying to reconstruct all permutations independent of each other.

Remark. *It is worth noticing that, even if the order samples are q -ary (i.e. full orderings of subsets of size q) instead of pairwise samples, for constant q , the order-wise behaviour of the sample complexities does not change.*

IV. PROOFS OF MAIN RESULTS

A. Proof of Theorem 1

In this section we provide the proof for the correctness of Algorithm 1 in detail. First, we identify a necessary and sufficient condition (Lemma 1) for the exact recovery of a single permutation and use it to provide an upper bound on the number of samples sufficient for recovery (Lemma 2). We then define a useful metric and characterize a concentration on distance between permutation pairs drawn from a uniform prior (Lemma 4). These tools are in turn used to prove that $O(\log r)$ samples per user are enough to distinguish two distinct permutations drawn uniformly at random (Lemma 7). Finally, all these intermediate results are used to prove Theorem 1 in its generality.

Lemma 1. *Let $\pi \in \mathcal{S}_n$ denote the true permutation we draw samples from. In order to recover the fully ordered set, we need all $n - 1$ samples in the set*

$$\mathcal{H}_\pi \triangleq \{(i, j) : i = \pi^{-1}(r), j = \pi^{-1}(r + 1), \\ \text{for } r = 1, \dots, n - 1\}.$$

That is, we need to sample all edges on the unique Hamiltonian path on the directed graph induced by permutation π . Furthermore, this set of samples is sufficient for exact recovery.

Proof: Let $(i, j) \in \mathcal{H}_\pi$ and $(i, j) \notin \Omega_s$, i.e. one pair of consecutive objects in the true ordering is not sampled. Considering i, j are consecutive, any indirect comparison through samples (i, k) and (k, j) is inconclusive, since either $\{\pi(i) < \pi(k) \text{ and } \pi(j) < \pi(k)\}$ or $\{\pi(i) > \pi(k) \text{ and } \pi(j) > \pi(k)\}$.

Conversely, if all pairs in \mathcal{H}_π are sampled, the unique ordering implied by the samples can be found by a simple topological sorting algorithm. ■

Lemma 2 (Simple achievability). *Given $s > 2N \log n$ samples, there exists an algorithm that recovers any permutation $\pi \in \mathcal{S}_n$ with high probability.*

Proof: By Lemma 1, sampling all edges on the unique Hamiltonian path, \mathcal{H}_π , is sufficient. Hence, the

probability of error is given by

$$P_e = \mathbb{P}(\mathcal{H}_\pi \not\subseteq \Omega_s) = \mathbb{P}\left(\bigcup_{(i,j) \in \mathcal{H}_\pi} \{(i, j) \notin \Omega_s\}\right) \\ \leq \sum_{(i,j) \in \mathcal{H}_\pi} \mathbb{P}((i, j) \notin \Omega_s) = \sum_{(i,j) \in \mathcal{H}_\pi} \left(1 - \frac{1}{N}\right)^s \\ \leq ne^{-\frac{s}{N}} = n^{-1}.$$

■

The complexity of discerning a couple of permutations, $\pi_1, \pi_2 \in \mathcal{S}_n$, depends on the magnitude of their difference. Since samples are pairwise orderings, a meaningful notion of distance is the following metric.

Definition 3 (Permutation distance metric). *We define the distance of a couple of permutations $\pi_1, \pi_2 \in \mathcal{S}_n$, to be the Hamming distance of the vector pairwise representations $\mathbf{p}_{\pi_1}, \mathbf{p}_{\pi_2} \in \{-1, +1\}^N$.*

$$d(\mathbf{p}_{\pi_1}, \mathbf{p}_{\pi_2}) \triangleq \sum_{i=1}^N \mathbb{1}[p_{\pi_1}(i) \neq p_{\pi_2}(i)]$$

We use this distance metric throughout this work.

Lemma 4 (Typical permutation distance). *Let $\pi_1, \pi_2 \in \mathcal{S}_n$ be drawn independently uniformly at random from the symmetric group. The distance between the two permutations is concentrated around $N/2$,*

$$\mathbb{P}(|d(\mathbf{p}_{\pi_1}, \mathbf{p}_{\pi_2}) - N/2| > \sqrt{n^{2+\gamma}}) \leq 2e^{-n^\gamma}, \quad \text{for } \gamma > 0.$$

Proof: First, by the linearity of expectation we get the expected distance to be

$$\mathbb{E}[d(\mathbf{p}_{\pi_1}, \mathbf{p}_{\pi_2})] = \mathbb{E}\left[\sum_{i=1}^N \mathbb{1}[p_{\pi_1}(i) \neq p_{\pi_2}(i)]\right] \\ = \sum_{i=1}^N \mathbb{P}(p_{\pi_1}(i) \neq p_{\pi_2}(i)) \\ = N/2.$$

To show concentration around this expected value, we construct a Doob martingale and use the Azuma-Hoeffding inequality. Define $X_i \triangleq \mathbb{1}\{p_{\pi_1}(i) \neq p_{\pi_2}(i)\}$ and let

$$Z_k = \mathbb{E}_{X_{k+1}, X_{k+2}, \dots}[d(\mathbf{p}_{\pi_1}, \mathbf{p}_{\pi_2}) | X_1, \dots, X_k],$$

and, specifically, $Z_0 = \mathbb{E}[d(\mathbf{p}_{\pi_1}, \mathbf{p}_{\pi_2})]$. Notice that, surely, $|Z_k - Z_{k-1}| \leq 1$. Then, by the Azuma-Hoeffding

inequality,

$$\begin{aligned}
\mathbb{P}(|d(\mathbf{p}_{\pi_1}, \mathbf{p}_{\pi_2}) - N/2| > \sqrt{n^{2+\gamma}}) \\
&= \mathbb{P}(|Z_N - Z_0| > \sqrt{n^{2+\gamma}}) \\
&\leq 2 \exp\left(\frac{-(\sqrt{n^{2+\gamma}})^2}{2 \sum_{k=1}^N 1^2}\right) \\
&= 2 \exp\left(\frac{-n^{2+\gamma}}{n(n-1)}\right) \\
&\approx 2e^{-n^\gamma}.
\end{aligned}$$

Lemma 5 (Number of occupied bins). *Consider the balls-and-bins problem, where m balls are thrown independently and uniformly at random into n bins. Let Z denote the number of occupied bins after the end of the process. If $m < n/2$, then for any $0 < \delta < 1$,*

$$\mathbb{P}(Z \leq (1-\delta)m) \leq (1-\delta)m \left(\frac{e^{\frac{\delta-1}{\delta}} n}{(1-\delta)m}\right)^{-\delta m}.$$

Proof: Let L_i denote the load of bin i after the end of the process. Since our main focus is the asymptotic behaviour, as both n and m grow, we ignore the integer constraints on the cardinality of bin set S .

$$\begin{aligned}
\mathbb{P}(Z \leq (1-\delta)m) \\
&= \mathbb{P}\left(\bigcup_{S \subset [n], |S| > n-(1-\delta)m} \bigcap_{i \in S} \{L_i = 0\}\right) \\
&\stackrel{(a)}{\leq} \sum_{S \subset [n], |S| > n-(1-\delta)m} \mathbb{P}\left(\bigcap_{i \in S} \{L_i = 0\}\right) \\
&= \sum_{S \subset [n], |S| > n-(1-\delta)m} \left(\frac{n-|S|}{n}\right)^m \\
&\stackrel{(b)}{\leq} \sum_{S \subset [n], |S| > n-(1-\delta)m} \left(\frac{n-n+(1-\delta)m}{n}\right)^m \\
&= \sum_{k=n-(1-\delta)m+1}^n \binom{n}{k} \left(\frac{(1-\delta)m}{n}\right)^m \\
&\stackrel{(c)}{\leq} (1-\delta)m \binom{n}{n-(1-\delta)m} \left(\frac{(1-\delta)m}{n}\right)^m \\
&= (1-\delta)m \binom{n}{(1-\delta)m} \left(\frac{(1-\delta)m}{n}\right)^m \\
&\stackrel{(d)}{\leq} (1-\delta)m \left(\frac{ne}{(1-\delta)m}\right)^{(1-\delta)m} \left(\frac{(1-\delta)m}{n}\right)^m \\
&= (1-\delta)m \left(\frac{e^{\frac{\delta-1}{\delta}} n}{(1-\delta)m}\right)^{-\delta m}.
\end{aligned}$$

Here, step (a) follows from a union bound and step (b) from the monotonicity of the summands with respect to

$|S|$. For step (c) notice that, since by assumption $m < n/2$, k is strictly greater than $n/2$ and in this range the binomial coefficient is decreasing in k . Finally, step (d) stems from a simple use of the bound $\binom{n}{k} \leq \left(\frac{ne}{k}\right)^k$. ■

Definition 6 (Clustering error). *Let $\pi_1, \pi_2 \in \mathcal{S}_n$, where $\pi_1 \neq \pi_2$, and $\Omega_1, \Omega_2 \subset [N]$ be the multisets of positions that have been sampled from \mathbf{p}_{π_1} and \mathbf{p}_{π_2} respectively. We define as clustering error the event when the two vectors agree on all common positions sampled, i.e.,*

$$\mathcal{E}_c(\pi_1, \pi_2, \Omega_1, \Omega_2) \triangleq \{\mathbf{p}_{\pi_1}(\Omega_1 \cap \Omega_2) = \mathbf{p}_{\pi_2}(\Omega_1 \cap \Omega_2)\}.$$

Lemma 7. *Let π_1, π_2 be drawn from \mathcal{S}_n independently uniformly at random. We independently and uniformly at random draw two multisets of samples with replacement from the two vectors \mathbf{p}_{π_1} and \mathbf{p}_{π_2} . The multisets of sampled indices, of sizes s_1 and s_2 respectively, are denoted by $\Omega_{1,s_1}, \Omega_{2,s_2} \subset [N]$. If $s_1 = s_2 = c \log r$, for constants $c > 4\beta > 0$, then, for large enough n and r , the probability of a clustering error is upper bounded as follows.*

$$P_e \triangleq \mathbb{P}(\mathcal{E}_c(\pi_1, \pi_2, \Omega_{1,s_1}, \Omega_{2,s_2})) \leq 2e^{-n} + 3r^{-\beta}.$$

Proof: Denote $d(\mathbf{p}_{\pi_1}, \mathbf{p}_{\pi_2})$ by d and the conflicting sample set of \mathbf{p}_{π_1} and \mathbf{p}_{π_2} by

$$\Delta = \{i : \mathbf{p}_{\pi_1}(i) \neq \mathbf{p}_{\pi_2}(i)\},$$

where $|\Delta| = d$. Then, the probability of clustering error is $\mathbb{P}(\Delta \cap \Omega_{1,s_1} \cap \Omega_{2,s_2} = \emptyset)$. We define $X \triangleq |\Delta \cap \Omega_{1,s_1}| \sim \text{Bin}(s_1, p)$, where $p = d/N$ and $\mathbb{E}[X] = s_1 d/N$. Notice that, since sampling is done with replacement and Ω_{1,s_1} is a multiset, we also consider $\Delta \cap \Omega_{1,s_1}$ to be a multiset, which justifies the Binomial model for its cardinality.

Now let Z denote the number of distinct samples in the multiset $\Delta \cap \Omega_{1,s_1}$. To facilitate exposition, let us define the following events:

$$\mathcal{E} = \{\Delta \cap \Omega_{1,s_1} \cap \Omega_{2,s_2} = \emptyset\},$$

$$\mathcal{A} = \left\{|\Delta \cap \Omega_{1,s_1}| \leq (1-\delta_1) \frac{s_1 d}{N}\right\}$$

and

$$\mathcal{B} = \{Z \leq (1-\delta_2) |\Delta \cap \Omega_{1,s_1}|\},$$

for $0 < \delta_1, \delta_2 < 1$. We upper bound the probability that none of the samples in Ω_{2,s_2} intersect with the samples in $\Delta \cap \Omega_{1,s_1}$.

$$\begin{aligned}
&\mathbb{P}(\Delta \cap \Omega_{1,s_1} \cap \Omega_{2,s_2} = \emptyset) \\
&= \mathbb{P}(\mathcal{E}|\mathcal{A})\mathbb{P}(\mathcal{A}) + \mathbb{P}(\mathcal{E}|\mathcal{A}^c)\mathbb{P}(\mathcal{A}^c) \\
&\leq \mathbb{P}(\mathcal{A}) + \mathbb{P}(\mathcal{E}|\mathcal{A}^c) \\
&= \mathbb{P}(\mathcal{A}) + \mathbb{P}(\mathcal{E}|\mathcal{B}, \mathcal{A}^c)\mathbb{P}(\mathcal{B}|\mathcal{A}^c) \\
&\quad + \mathbb{P}(\mathcal{E}|\mathcal{B}^c, \mathcal{A}^c)\mathbb{P}(\mathcal{B}^c|\mathcal{A}^c) \\
&\leq \mathbb{P}(\mathcal{A}) + \mathbb{P}(\mathcal{B}|\mathcal{A}^c) + \mathbb{P}(\mathcal{E}|\mathcal{B}^c, \mathcal{A}^c) \quad (1)
\end{aligned}$$

Using the Chernoff bound,

$$\mathbb{P}(X \leq (1 - \delta_1)\mathbb{E}[X]) \leq \exp\{-\delta_1^2\mathbb{E}[X]/2\}$$

for a binomially distributed X , we can upper bound the first term like this:

$$\begin{aligned} \mathbb{P}(\mathcal{A}) &= \mathbb{P}\left(|\Delta \cap \Omega_{1,s_1}| \leq (1 - \delta_1)\frac{s_1 d}{N}\right) \\ &\leq \exp\left\{-\delta_1^2 \frac{dc \log r}{2N}\right\} = r^{-c\delta_1^2 \frac{d}{2N}}. \end{aligned}$$

For the second term in (1), since surely $|\Delta \cap \Omega_{1,s_1}| < s_1$, by Lemma 5,

$$\begin{aligned} \mathbb{P}(\mathcal{B}|\mathcal{A}^c) &= \mathbb{P}\left(Z \leq (1 - \delta_2)|\Delta \cap \Omega_{1,s_1}| \mid \mathcal{A}^c\right) \\ &\leq \mathbb{P}(Z \leq (1 - \delta_2)s_1) \\ &\leq (1 - \delta_2)c \log r \left(\frac{e^{\frac{\delta_2-1}{\delta_2}} d}{(1 - \delta_2)c \log r}\right)^{-\delta_2 c \log r}. \end{aligned}$$

For the final term in (1),

$$\begin{aligned} \mathbb{P}(\mathcal{E}|\mathcal{B}^c, \mathcal{A}^c) &= \mathbb{P}\left(\mathcal{E} \mid Z > (1 - \delta_2)(1 - \delta_1)\frac{s_1 d}{N}\right) \\ &\leq \left(1 - (1 - \delta_2)(1 - \delta_1)\frac{s_1 d}{N}\right)^{s_2} \\ &\leq \exp\left\{-(1 - \delta_2)(1 - \delta_1)\frac{s_1 s_2 d}{N}\right\} \\ &= \exp\left\{-(1 - \delta_2)(1 - \delta_1)\frac{dc^2 \log^2 r}{N}\right\}. \end{aligned}$$

Let $\underline{d} \triangleq N/2 - n\sqrt{n}$. By Lemma 4, $\mathbb{P}(d < \underline{d}) \leq 2e^{-n}$. Now, since the bound on the right hand side of (1) is decreasing in d ,

$$\begin{aligned} \mathbb{P}(\mathcal{E}|d \geq \underline{d}) &\leq r^{-c\delta_1^2 \frac{d}{2N}} \\ &+ (1 - \delta_2)c \log r \left(\frac{e^{\frac{\delta_2-1}{\delta_2}} \underline{d}}{(1 - \delta_2)c \log r}\right)^{-\delta_2 c \log r} \\ &+ \exp\left\{-(1 - \delta_2)(1 - \delta_1)\frac{dc^2 \log^2 r}{N}\right\}. \end{aligned}$$

Straightforward calculations yield that, for any $c > 4\beta$, there exists n' large enough such that the first term decays at least as fast as $r^{-\beta}$ and there exists n'' large enough such that

$$\mathbb{P}(\mathcal{E}|d > \underline{d}) \leq 3r^{-\beta}, \quad \forall n > n''.$$

Finally,

$$\begin{aligned} \mathbb{P}(\mathcal{E}) &= \mathbb{P}(\mathcal{E}|d < \underline{d})\mathbb{P}(d < \underline{d}) + \mathbb{P}(\mathcal{E}|d \geq \underline{d})\mathbb{P}(d \geq \underline{d}) \\ &\leq \mathbb{P}(d < \underline{d}) + \mathbb{P}(\mathcal{E}|d \geq \underline{d}) \\ &\leq 2e^{-n} + 3r^{-\beta}. \end{aligned}$$

Having collected all the necessary lemmata and intermediate results we, finally, turn to proving Theorem 1.

Theorem. Suppose $r = \theta(m^\gamma)$ for a fixed $\gamma > 0$. Given $s > \max\left\{cm \log r, \frac{2rN \log n}{(1-\delta_1)(1-\delta_2)}\right\}$ random samples, for $c > 12/\gamma$ and $0 < \delta_1, \delta_2 < 1$, Algorithm 1 fully recovers all m permutations correctly, in the sense that $\{\hat{\pi}_k = \pi_k, \forall k \in [m]\}$ with high probability.

Proof: By Lemma 7, all comparisons between distinct permutations at Step 5 of the clustering stage, are bound to reveal a difference between \mathbf{p}_{π_u} and \mathbf{p}_{π_k} with high probability. A simple union bound over all user pairs gives us a guarantee that the clustering stage will succeed with high probability. Let $\mathcal{E}_{u,k}$ denote the event that the sampled values from users u and k do not reveal a difference, i.e.,

$$\mathcal{E}_{u,k} \triangleq \{\mathbf{M}(\Omega_{s,u} \cap \Omega_{s,k}, u) = \mathbf{M}(\Omega_{s,u} \cap \Omega_{s,k}, k)\}.$$

The probability that at least one error occurs in Step 5 is upper bounded by,

$$\begin{aligned} \mathbb{P}\left(\bigcup_{\substack{(u,k) \\ \pi_u \neq \pi_k}} \mathcal{E}_{u,k}\right) &\leq \sum_{\substack{(u,k) \in [m]^2 \\ \pi_u \neq \pi_k}} \mathbb{P}(\mathcal{E}_{u,k}) \\ &\leq 2m^2 e^{-n} + 3m^2 r^{-c/4} \\ &\leq 2m^2 e^{-n} + 3c' m^{2-c\gamma/4}, \end{aligned}$$

where in the second inequality we use the result from Lemma 7 and in the last one, the assumption $r = \theta(m^\gamma)$. When $m^2 = O(e^n/m)$, for any $c > 12/\gamma$ this bound decays as m^{-1} for large enough m .

The number of samples drawn from user k is $S_k \sim \text{Bin}(s, 1/m)$ and the number of users in cluster j is $U_j \sim \text{Bin}(m, 1/r)$. Let L_j denote the number of samples drawn from all users in cluster j . A sure lower bound for $L_{\min} \triangleq \min_{j=1}^r L_j$ is

$$L_{\min} \geq U_{\min} S_{\min},$$

where $U_{\min} \triangleq \min_{j=1}^r U_j$ and $S_{\min} \triangleq \min_{k=1}^m S_k$. Simple Chernoff and union bounds give us a lower bound on U_{\min} , for $0 < \delta_1 < 1$,

$$\mathbb{P}\left(U_{\min} \leq (1 - \delta_1)\frac{m}{r}\right) \leq re^{-\frac{m\delta_1^2}{2r}},$$

and on S_{\min} , for $0 < \delta_2 < 1$,

$$\mathbb{P}\left(S_{\min} \leq (1 - \delta_2)\frac{s}{m}\right) \leq me^{-\frac{s\delta_2^2}{2m}}.$$

Then by a union bound,

$$\begin{aligned} \mathbb{P}\left(L_{\min} \leq (1 - \delta_1)\frac{m}{r}(1 - \delta_2)\frac{s}{m}\right) \\ \leq re^{-\frac{m\delta_1^2}{2r}} + me^{-\frac{s\delta_2^2}{2m}}, \end{aligned}$$

and setting $s = \frac{2rN \log n}{(1-\delta_1)(1-\delta_2)}$,

$$\mathbb{P}(L_{\min} \leq 2N \log n) \leq r e^{-\frac{m\delta_1^2}{2r}} + m e^{-\frac{s\delta_2^2}{2m}},$$

which decays at least as fast as r^{-1} when given the prescribed number of samples. The result follows as a consequence of Lemma 2. This proves that both stages of Algorithm 1 succeed with high probability. ■

B. Proof of Theorem 2

In this section we provide the proof of the strong converse theorem for the random sampling complexity in detail. We first prove a converse theorem for the sample complexity of learning a single permutation (Lemma 8). To that end, we show that, given an insufficient number of samples, the random sampling process will miss an increasingly large number of critical samples leading to a probability of learning error that approaches 1. Next, we prove a strong converse theorem for the task of clustering (Lemma 14) by extending the information theoretic definitions of typicality to colourings and clusterings and using a popular source coding converse theorem technique. Eventually, we use these intermediate results in the proof for Theorem 2.

Lemma 8 (Strong converse for single permutation learning). *For any recovery algorithm $\hat{\pi} = g(\Omega_s, P_\pi(\Omega_s))$ that uses $s < N(\log n - \log \log n - 1 - \sqrt{d\pi^2 \log n/3})$ random samples, the probability of error, $P_e = \mathbb{P}(\hat{\pi} \neq \pi)$, tends to 1 as $n \rightarrow \infty$.*

Proof: By Lemma 1, $\{\hat{\pi} \neq \pi\} = \{\mathcal{H}_\pi \not\subseteq \Omega_s\}$. Let T denote the number of samples drawn until $\mathcal{H}_\pi \subseteq \Omega_s$. This is an instance of the classic coupon collector problem, where $T = T_1 + T_2 + \dots + T_{n-1}$ for $T_k \sim \text{Geom}(p_k)$, and $p_k = \frac{n-k}{N}$, and all T_k 's are independent. Further, define the partial sum $T^{(l)} = \sum_{k=1}^l T_k$ for $l \in [n-1]$. We use a Chernoff bound to show a concentration of $T^{(l)}$ around the mean. Let $S_i \triangleq T_i - \mathbb{E}T_i$ and $S^{(l)} \triangleq \sum_{k=1}^l S_k$. For any $t > 0$,

$$\begin{aligned} \mathbb{P}(T^{(l)} \leq \mathbb{E}T^{(l)} - cN) &= \mathbb{P}(S^{(l)} \leq -cN) \\ &= \mathbb{P}(e^{-tS^{(l)}} \geq e^{tcN}) \\ &\leq e^{-tcN} \mathbb{E}[e^{-tS^{(l)}}]. \end{aligned}$$

Now, to bound the moment generating function (MGF) of $S^{(l)}$,

$$\begin{aligned} \mathbb{E}[e^{-tS^{(l)}}] &= \prod_{k=1}^l \mathbb{E}[e^{-t(T_k - 1/p_k)}] = \prod_{k=1}^l e^{t/p_k} \mathbb{E}[e^{-tT_k}] \\ &= \prod_{k=1}^l \frac{e^{t/p_k}}{(e^t - 1)/p_k + 1} \leq \prod_{k=1}^l \frac{e^{t/p_k}}{t/p_k + 1} \\ &\leq e^{\frac{1}{2}t^2 \sum_{k=1}^l p_k^{-2}} \leq e^{\frac{1}{2}t^2 \frac{N^2 \pi^2}{6}}. \end{aligned}$$

$$\Rightarrow \mathbb{P}(T^{(l)} \leq \mathbb{E}[T^{(l)}] - cN) \leq e^{\frac{1}{2}(N\pi t)^2 - tcN}$$

Optimizing over t yields

$$\mathbb{P}(T^{(l)} \leq \mathbb{E}[T^{(l)}] - cN) \leq e^{-\frac{3c^2}{\pi^2}}. \quad (2)$$

Let us set $l = l(n) = n - \lfloor \log n \rfloor$, so that for n large enough, we have

$$\begin{aligned} \mathbb{E}[T^{(l(n))}] &= \sum_{k=1}^{l(n)} 1/p_k = N \sum_{k=1}^{l(n)} \frac{1}{n-k} \\ &= N(H_{n-1} - H_{n-l(n)-1}) = N(H_{n-1} - H_{\lfloor \log n \rfloor - 1}) \\ &\geq N(\log(n-1) - \log(\lfloor \log n \rfloor - 1) - 1) \\ &\geq N(\log(n-1) - \log(\log n - 1) - 1) \\ &\geq N(\log(n/\log n) - 1). \end{aligned}$$

Using the bound (2) with our choice of $l = l(n)$, and with $c = c(n) = \sqrt{d\pi^2 \log n/3}$ for a fixed $d > 0$, we can write

$$\begin{aligned} \mathbb{P}(T^{(l)} \leq s) &\leq \mathbb{P}(T^{(l)}) \\ &\leq N(\log n - \log \log n - 1) - cN \\ &\leq \mathbb{P}(T^{(l)} \leq \mathbb{E}[T^{(l)}] - cN) \leq e^{-\frac{3c^2}{\pi^2}} = n^{-d}. \end{aligned}$$

Since $\{T^{(l)} \leq s\} = \{|\mathcal{H}_\pi \setminus \Omega_s| \leq n - 1 - l\}$, this means that

$$\mathbb{P}(|\mathcal{H}_\pi \setminus \Omega_s| \geq \lfloor \log n \rfloor) \geq 1 - n^{-d}. \quad (3)$$

A permutation-learning algorithm g needing s samples, in its most general form, maps the input of sampled pairs $\Omega_s \subseteq [N]$ and their values $P_\pi(\Omega_s) \in \{1, +1\}$ to a probability distribution on \mathcal{S}_n . The probability of successful learning for g is then $P_g \triangleq \mathbb{P}(g(\Omega_s, P_\pi(\Omega_s)) = \pi)$. A standard information-theoretic argument yields the following fact: the highest probability of success over all permutation-learning algorithms is attained by a Maximum Likelihood (ML) learning algorithm g^* given by

$$g^*(\omega_s, p_s) = \arg \max_{\bar{\pi} \in \mathcal{S}_n} \mathbb{P}(\pi = \bar{\pi} \mid \Omega_s = \omega_s, P_\pi(\Omega_s) = p_s),$$

with the corresponding success probability being

$$\begin{aligned} P_{g^*} &= \mathbb{P}(g^*(\Omega_s, P_{\bar{\pi}}(\Omega_s)) = \pi) \\ &= \mathbb{E} \left[\max_{\bar{\pi} \in \mathcal{S}_n} \mathbb{P}(\pi = \bar{\pi} \mid \Omega_s, P_{\bar{\pi}}(\Omega_s)) \right]. \quad (4) \end{aligned}$$

We show that as $n \rightarrow \infty$, $P_{g^*} \rightarrow 0$, completing the proof of the theorem. This is accomplished, in turn, via the following claims:

Lemma 9. *Let $\bar{\pi} \in \mathcal{S}_n$, $\omega_s \subset [N]$ and p_s be such that $P_{\bar{\pi}}(\omega_s) = p_s$. If $|\mathcal{H}_{\bar{\pi}} \setminus \omega_s| \geq k$, then $P_{\bar{\pi}'}(\omega_s) = p_s$ for at least 2^k distinct permutations $\bar{\pi}'$ in \mathcal{S}_n .*

Proof: This follows from the observation that by assigning any combination of ± 1 values to the entries

of $\bar{\pi}$ in the positions $\mathcal{H}_{\bar{\pi}} \setminus \omega_s$ of $P_{\bar{\pi}}$, and keeping all other entries unchanged, we get a valid permutation (pairwise representation). ■

Lemma 10. *If $\bar{\pi}_1, \bar{\pi}_2 \in \mathcal{S}_n$, ω_s and p_s are such that $P_{\bar{\pi}_1}(\omega_s) = P_{\bar{\pi}_2}(\omega_s) = p_s$, then*

$$\begin{aligned} & \mathbb{P}\left(\pi = \bar{\pi}_1 \mid \Omega_s = \omega_s, P_\pi(\Omega_s) = p_s\right) \\ &= \mathbb{P}\left(\pi = \bar{\pi}_2 \mid \Omega_s = \omega_s, P_\pi(\Omega_s) = p_s\right). \end{aligned}$$

Proof: We have,

$$\begin{aligned} & \mathbb{P}\left(\pi = \bar{\pi}_1 \mid \Omega_s = \omega_s, P_\pi(\Omega_s) = p_s\right) \\ &= \frac{\mathbb{P}\left(\Omega_s = \omega_s, P_\pi(\Omega_s) = p_s \mid \pi = \bar{\pi}_1\right) \mathbb{P}\left(\pi = \bar{\pi}_1\right)}{\mathbb{P}\left(\Omega_s = \omega_s, P_\pi(\Omega_s) = p_s\right)} \\ &= \frac{\mathbb{P}\left(\Omega_s = \omega_s\right) \mathbb{P}\left(\pi = \bar{\pi}_1\right)}{\mathbb{P}\left(\Omega_s = \omega_s, P_\pi(\Omega_s) = p_s\right)} \\ &= \frac{\mathbb{P}\left(\Omega_s = \omega_s\right) \mathbb{P}\left(\pi = \bar{\pi}_2\right)}{\mathbb{P}\left(\Omega_s = \omega_s, P_\pi(\Omega_s) = p_s\right)} \\ &= \mathbb{P}\left(\pi = \bar{\pi}_2 \mid \Omega_s = \omega_s, P_\pi(\Omega_s) = p_s\right). \end{aligned}$$

Combining the results of the two lemmata above, we have that if ω_s and p_s are such that there exists $\bar{\pi} \in \mathcal{S}_n$ satisfying $P_{\bar{\pi}}(\omega_s) = p_s$ and $|\mathcal{H}_{\bar{\pi}} \setminus \omega_s| \geq k$, then

$$\begin{aligned} & \max_{\hat{\pi} \in \mathcal{S}_n} \mathbb{P}\left(\pi = \hat{\pi} \mid \Omega_s = \omega_s, P_\pi(\Omega_s) = p_s\right) \leq 2^{-k} \\ \Rightarrow & \mathbb{1}_{\{|\mathcal{H}_{\bar{\pi}} \setminus \omega_s| \geq k\}} \cdot \max_{\hat{\pi} \in \mathcal{S}_n} \mathbb{P}\left(\pi = \hat{\pi} \mid \Omega_s = \omega_s, P_\pi(\Omega_s) = p_s\right) \leq 2^{-k}. \end{aligned}$$

Setting $k = \lfloor \log n \rfloor$ in the above, and using (3) and (4) gives

$$\begin{aligned} P_{g^*} &= \mathbb{E} \left[\max_{\hat{\pi} \in \mathcal{S}_n} \mathbb{P}\left(\pi = \hat{\pi} \mid \Omega_s = \omega_s, P_\pi(\Omega_s) = p_s\right) \right] \\ &\leq 2^{-\log n} + n^{-d} \rightarrow_{n \rightarrow \infty} 0 \end{aligned}$$

as promised. ■

Learning the full set of permutations intuitively involves a step of clustering, where users that adhere to the same permutation are grouped in the first stage of the algorithm. We argue in the proof for Theorem 2 that an algorithm, that cannot – even as a by-product – cluster the users, cannot perform the learning task at hand.

Definition 11 (*r*-colouring). *We call a mapping $L : [m] \rightarrow [r]$ an *r*-colouring.*

Definition 12 (*r*-clustering). *Consider two *r*-colourings L_1 and L_2 and define the equivalence relation*

$$\begin{aligned} L_1 &\sim_C L_2 \\ \Leftrightarrow & \{L_1(i) = L_1(j) \text{ iff } L_2(i) = L_2(j), \forall i, j \in [m]\}. \end{aligned}$$

The two colourings are equivalent if their colour groups are exactly the same, i.e., they cluster elements in $[m]$ in exactly the same way. Let \mathcal{T} denote the quotient group

*of all colourings by \sim_C . An element $C \in \mathcal{T}$, is an equivalence class of colourings, also denoted by $\langle L \rangle_C$ for any colouring $L \in C$, and we call it an *r*-clustering.*

Definition 13 (True clustering). *We denote by $C(\mathbf{M}) \in \mathcal{T}$ the true *r*-clustering induced by the pairwise representation matrix \mathbf{M} .*

Lemma 14 (Clustering converse). *For any clustering algorithm, given by $\hat{C} = g(\mathbf{M})$, $\hat{C} \in \mathcal{T}$, that uses $s < m(\log r - \zeta) - r \log r$ samples, for any $\zeta > 0$, the probability of error, $P_e \triangleq \mathbb{P}(\hat{C} \neq C(\mathbf{M}))$, goes to 1 as $n \rightarrow \infty$.*

Proof: In the system model described in Section II, each of m users is assigned to a cluster i.i.d., using a uniform distribution in $[r]$. More generally, we can assume that the cluster indices q_i are drawn i.i.d. from some arbitrary distribution $p(q_1^m) = \prod_{i=1}^m p(q_i)$. Our analysis uses the notion of strong typicality ([18], [19]) and goes along the lines of a converse theorem for source coding in [19].

■ **Definition 15** (Strongly typical *r*-colouring). *Let L be a random *r*-colouring, where colours are drawn from $[r]$ i.i.d. according to $p(\cdot)$, and define $N(q; L)$ to be the number of occurrences of the colour $q \in [r]$ in L . We define the set of δ -strongly typical *r*-colourings with respect to $p(q)$ to be*

$$T_{p(q), \delta}^m = \left\{ L : \sum_{q=1}^r \left| \frac{1}{m} N(q; L) - p(q) \right| \leq \delta \right\}. \quad (5)$$

Standard typicality results give us the following bounds on the cardinality of typical sets and the probability of typical *r*-colourings.

Lemma 16 (Theorem 6.2, [19]). *There exists $\epsilon > 0$ such that $\epsilon \rightarrow 0$ as $\delta \rightarrow 0$ and if $L \in T_{p(q), \delta}^m$,*

$$2^{-m(H(p)+\epsilon)} \leq p(L) \leq 2^{-m(H(p)-\epsilon)}. \quad (6)$$

Also, for m sufficiently large,

$$\mathbb{P}(L \in T_{p(q), \delta}^m) > 1 - \epsilon, \quad (7)$$

$$(1 - \delta) 2^{m(H(p)-\epsilon)} \leq |T_{p(x), \delta}^m| \leq 2^{m(H(p)+\epsilon)}. \quad (8)$$

*An interesting fact about members of this typical set is that for any permutation $\sigma \in \mathcal{S}_r$, i.e. a recolouring that preserves colour groups, and any *r*-colouring L ,*

$$L \in T_{p(q), \delta}^m \Leftrightarrow \sigma(L) \in T_{p(\sigma^{-1}(q)), \delta}^m,$$

where $\sigma(L) \triangleq [\sigma(L(1)), \sigma(L(2)), \dots, \sigma(L(m))]$. More specifically, in our current setup, $p(q)$ is uniform in $[r]$ so, $p(\sigma^{-1}(q)) = p(q)$ for all $q \in [r]$ and

$$L \in T_{p(q), \delta}^m \Leftrightarrow \sigma(L) \in T_{p(q), \delta}^m.$$

As a direct consequence, $L \in T_{p(q),\delta}^m \Rightarrow L' \in T_{p(q),\delta}^m$, for all $L' \in \langle L \rangle_C$, and we can give a simple definition for strongly typical r -clusterings.

Definition 17 (Strongly typical r -clustering). *Under uniform $p(q)$, let L be a δ -strongly typical r -colouring. We call the equivalence class $\langle L \rangle_C$ a δ -strongly typical r -clustering, and denote the collection of all δ -strongly typical r -clusterings by \mathcal{T}_δ^m .*

The collection \mathcal{T}_δ^m is essentially a partition of all δ -strongly typical r -colourings and otherwise contains the exact same elements. Hence from (7),

$$\mathbb{P}(C \in \mathcal{T}_\delta^m) = \mathbb{P}(L \in T_{p(q),\delta}^m) > 1 - \epsilon.$$

Now, to calculate the cardinality of this collection, notice that, by (5), for any $\delta < 1/r$, any δ -strongly typical r -colouring is onto, that is, it maps to all values in $[r]$. In this case, since $\langle L \rangle_C$ consists of all possible colour permutations $\sigma(L)$ for some $\sigma \in \mathcal{S}_r$, $|\langle L \rangle_C| = r!$. From this fact and (6), for $C \in \mathcal{T}_\delta^m$, we get

$$\begin{aligned} p(C) &= \sum_{L \in C} p(L) \leq r! 2^{-m(H(p)-\epsilon)} \\ &\leq 2^{-m(\log r - \epsilon) + r \log r}. \end{aligned} \quad (9)$$

Now consider an arbitrary deterministic clustering algorithm $\hat{C} = g(\mathbf{M})$ that uses $s < m(\log r - \zeta) - r \log r$ binary samples from the matrix \mathbf{M} , for some constant $\zeta > 0$. It can map each of the 2^s possible unique inputs to either strongly typical r -clusterings, or non-strongly typical r -clusterings. Using (9), the probability mass within \mathcal{T}_δ^m covered by the outputs of the algorithm is at most

$$2^s 2^{-m(\log r - \epsilon) + r \log r} = 2^{-m(\zeta - \epsilon)}$$

and the total probability covered is, by (7), at most

$$2^{-m(\zeta - \epsilon)} + \mathbb{P}(C \notin \mathcal{T}_\delta^m) \leq 2^{-m(\zeta - \epsilon)} + \epsilon.$$

It follows that any mass in \mathcal{T} not covered by the algorithm gives a lower bound on the probability of error, hence

$$\mathbb{P}_e > 1 - (2^{-m(\zeta - \epsilon)} + \epsilon)$$

which goes to 1 for any positive value of ζ . \blacksquare

Lemma 18 (Learning converse). *For any recovery algorithm given by $(\hat{\pi}_k)_{k=1}^m = g(\Omega_s, \mathbf{M}(\Omega_s))$, that uses $s < rN(\log n - \log \log n - 1 - \sqrt{d\pi^2 \log n/3})$ random samples, the probability of error, $P_e = \mathbb{P}(\bigcup_{k=1}^m \{\hat{\pi}_k \neq \pi_k\})$, goes to 1 as $n \rightarrow \infty$.*

Proof: Let L_j be the number of samples drawn from users adhering to permutation ρ_j , for $j \in [r]$. Since $\sum_{j=1}^r L_j = s$, there exists surely an j' such that, $L_{j'} \leq s/r$. Therefore,

$$L_{j'} < N(\log n - \log \log n - 1 - \sqrt{d\pi^2 \log n/3})$$

and from the result in Lemma 8, the probability of learning $\rho_{j'}$ goes to zero. That is, there exists user $k \in [m]$ such that $\pi_k = \rho_{j'}$, so $\mathbb{P}(\hat{\pi}_k \neq \pi_k) = \mathbb{P}(\hat{\rho}_{j'} \neq \rho_{j'}) \rightarrow 1$ and

$$P_e = \mathbb{P}\left(\bigcup_{u=1}^m \{\hat{\pi}_u \neq \pi_u\}\right) \geq \mathbb{P}(\hat{\pi}_k \neq \pi_k) \rightarrow 1. \quad \blacksquare$$

Theorem. *For any recovery algorithm given by $(\hat{\pi}_k)_{k=1}^m = g(\Omega_s, \mathbf{M}(\Omega_s))$, that uses*

$$s < \max \left\{ rN \left(\log n - \log \log n - 1 - \sqrt{d\pi^2 \log n/3} \right), \right. \\ \left. (m - r) \log r \right\}$$

random samples, the probability of error, $P_e = \mathbb{P}(\bigcup_{k=1}^m \{\hat{\pi}_k \neq \pi_k\})$, goes to 1 as $n \rightarrow \infty$.

Proof: The first term is directly explained by Lemma 18. To justify the second term we give the following argument. Consider an algorithm $g(\Omega_s, \mathbf{M}(\Omega_s))$ that uses $s < (m - r) \log r$ samples and suppose it sustains a probability of error uniformly bounded away from 1. That is, $P_e \leq 1 - \delta$, for $\delta > 0$ and for all $n > 0$. Then the clustering induced by the output of the algorithm is correct with probability at least δ , or equivalently,

$$\mathbb{P}\left(\hat{C}(g(\Omega_s, \mathbf{M}(\Omega_s))) \neq C(\mathbf{M})\right) \leq 1 - \delta, \quad \forall n.$$

This result directly contradicts Lemma 14 and the assumption that δ is strictly positive must be false. This concludes the second argument. \blacksquare

C. Proof of Theorem 3

In this section we prove correctness of Algorithm 2. We use results already established in Section IV-A as stepping stones to prove correctness of the clustering step and known properties of sorting algorithms to prove correctness of the learning step.

Theorem. *There exist constants c_1, c_2 such that, using $s > c_1 m \log r + c_2 r n \log n$ samples, Algorithm 2 fully recovers all m permutations correctly, in the sense that $\mathbb{P}(\hat{\pi}_k = \pi_k, \forall k)$ tends to 1.*

Proof: We first analyze the probability of a clustering error; i.e., the event that a mistake happens in Step 6 of Stage 1 of our algorithm. Consider the case when $\pi_u \neq \pi_k$. By Lemma 4, the event

$$\mathcal{D} \triangleq \{d(\mathbf{p}_{\pi_u}, \mathbf{p}_{\pi_k}) < N/2 - n\sqrt{n}\}$$

happens with probability at most $2e^{-n}$. Then, the probability that for $\pi_u \neq \pi_k$, the sampled positions match

perfectly can be upper bounded as follows.

$$\begin{aligned}
& \mathbb{P}(\mathbf{M}(\Omega_C, u) = \mathbf{M}(\Omega_C, k) | \pi_u \neq \pi_k) \\
&= \mathbb{P}(\mathbf{M}(\Omega_C, u) = \mathbf{M}(\Omega_C, k) | \pi_u \neq \pi_k, \mathcal{D}) \mathbb{P}(\mathcal{D}) \\
&\quad + \mathbb{P}(\mathbf{M}(\Omega_C, u) = \mathbf{M}(\Omega_C, k) | \pi_u \neq \pi_k, \mathcal{D}^c) \mathbb{P}(\mathcal{D}^c) \\
&\leq 1 \cdot 2e^{-n} + \left(1 - \frac{N/2 - n\sqrt{n}}{N}\right)^{c_1 \log r} \cdot 1 \\
&\leq 2e^{-n} + \exp\left\{-c_1 \log r \frac{N/2 - n\sqrt{n}}{N}\right\} \\
&= 2e^{-n} + r^{-\frac{c_1}{2} + \frac{2\sqrt{n}}{n-1}}.
\end{aligned} \tag{10}$$

Let us now define the event, that the sampled values from users u and k do not reveal a difference, as

$$\begin{aligned}
\mathcal{E}_{u,k} &= \{\mathbf{M}(\Omega_C, u) = \mathbf{M}(\Omega_C, k)\} \\
&= \{\mathbf{p}_{\pi_u}(\Omega_C) \neq \mathbf{p}_{\pi_k}(\Omega_C)\}.
\end{aligned}$$

The probability that at least one error occurs in Step 6 is upper bounded by,

$$\begin{aligned}
& \mathbb{P}\left(\bigcup_{\substack{(u,k) \\ \pi_u \neq \pi_k}} \mathcal{E}_{u,k}\right) = \mathbb{P}\left(\bigcup_{\substack{(i,l) \in [r]^2 \\ i \neq l}} \bigcup_{\substack{(u,k) \\ \pi_u = \rho_i, \pi_k = \rho_l}} \mathcal{E}_{u,k}\right) \\
&= \mathbb{P}\left(\bigcup_{\substack{(i,l) \in [r]^2 \\ i \neq l}} \bigcup_{\substack{(u,k) \\ \pi_u = \rho_i, \pi_k = \rho_l}} \{\mathbf{p}_{\pi_u}(\Omega_C) \neq \mathbf{p}_{\pi_k}(\Omega_C)\}\right) \\
&= \mathbb{P}\left(\bigcup_{\substack{(i,l) \in [r]^2 \\ i \neq l}} \{\mathbf{p}_{\rho_i}(\Omega_C) \neq \mathbf{p}_{\rho_l}(\Omega_C)\}\right) \\
&\leq \sum_{\substack{(i,l) \in [r]^2 \\ i \neq l}} \mathbb{P}(\mathbf{p}_{\rho_i}(\Omega_C) \neq \mathbf{p}_{\rho_l}(\Omega_C)) \\
&\leq 2r^2 e^{-n} + r^{2 - \frac{c_1}{2} + \frac{2\sqrt{n}}{n-1}},
\end{aligned}$$

where in the last step we use the result from (10). For any $c_1 > 6$ this bound decays as r^{-1} for large enough n .

The second stage of the algorithm consists of r sorting tasks. The sample complexity depends on the sorting algorithm used. The usual candidate, Quicksort, is famous for its low complexity. The authors in [20] argue the existence of constants $c_2(k)$, for every k , such that the sample complexity of Randomized Quicksort is at most $c_2(k)n \log n$ with probability at least $1 - \frac{1}{n^k}$. In the regime where $r = \theta(n^\gamma)$, this result along with a union bound over the r sorting tasks establish that $c_2(k)rn \log n$ samples are sufficient for learning with high probability. In other regimes, a different sorting

algorithm with better worst case guarantees can be used (e.g. Mergesort) to give complexity of the same order. ■

D. Proof of Theorem 4

In this section we provide the proof for the strong converse theorem on the sample complexity in the active sampling scenario. We reuse the clustering converse argument from the proof of Theorem 2 (Lemma 14) and provide a new information theoretic argument for the intrinsic complexity of the joint learning task.

Theorem. For any recovery algorithm given by $(\hat{\pi}_k)_{k=1}^m = g(\mathbf{M})$, and using

$$s < \max\{(m-r) \log r, crn \log n\}$$

samples for any $0 < c < 1$, the probability of error, $P_e = \mathbb{P}(\bigcup_{k=1}^m \{\hat{\pi}_k \neq \pi_k\})$, goes to 1 as $n \rightarrow \infty$.

Proof: The first term is related to the clustering ability of an algorithm and follows from an argument identical to the one given in the proof for Theorem 2 hinging on Lemma 14. We do not repeat it here.

By a similar argument, any algorithm that cannot identify the common pool of permutations, $\{\rho_j\}_{j=1}^r$ (see Section II), cannot learn all user permutations correctly. To see this, whenever $\hat{\pi}_k = \pi_k$ for all k , the set $\{\hat{\rho}_j\}_{j=1}^r$ consisting of the r unique elements in $\{\hat{\pi}_k\}_{k=1}^m$, is exactly $\{\rho_j\}_{j=1}^r$ (modulo a relabeling). Hence, it suffices to show that for any algorithm, that is given by $\{\hat{\rho}_j\}_{j=1}^r = \tilde{g}(\mathbf{M})$ and uses $s < crn \log n$ samples, $\mathbb{P}(\{\hat{\rho}_j\}_{j=1}^r \neq \{\rho_j\}_{j=1}^r)$ goes to 1 as $n \rightarrow \infty$.

Uniformity in the choice of the ρ_j 's and a union bound tell us that, when $r = o(\sqrt{n!})$, all ρ_j 's are distinct with high probability – a birthday paradox type result. This probability is

$$\mathbb{P}(\rho_j \neq \rho_k, \forall j \neq k) \geq 1 - \frac{r(r-1)}{2n!} = 1 - \epsilon$$

for $\epsilon \rightarrow 0$ as $n \rightarrow \infty$. In this case, there are $\binom{n!}{r}$ distinct permutation pool sets possible; all with equal probability, at least $(1 - \epsilon)/\binom{n!}{r}$ and at most $1/\binom{n!}{r}$. Now consider any deterministic algorithm that maps the 2^s different inputs to different common permutation sets. Following an argument similar to Lemma 14, we argue that the probability mass covered within the set of all solutions is at most

$$\begin{aligned}
2^s \frac{1}{\binom{n!}{r}} + \epsilon &\leq 2^{crn \log_2 n} 2^{-rn \log_2 n + \frac{rn}{\ln 2}} + \epsilon \\
&= 2^{(c-1)rn \log_2 n + \frac{rn}{\ln 2}} + \epsilon
\end{aligned}$$

and any mass not covered results in an error. Hence,

$$\mathbb{P}(\{\hat{\rho}_j\}_{j=1}^r \neq \{\rho_j\}_{j=1}^r) \geq 1 - 2^{(c-1)rn \log_2 n + \frac{rn}{\ln 2}} - \epsilon$$

and, for any $c < 1$, the probability of error goes to 1 as $n \rightarrow \infty$. This concludes the proof. ■

V. EXPERIMENTAL RESULTS

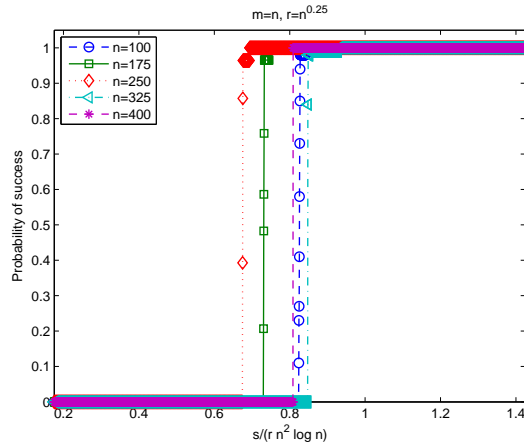
This section describes results of numerical experiments on synthetic data following the stochastic model for user permutations introduced in this work. For the random sampling case, Matrix Completion (MC) is an attractive choice of algorithm to hope to recover users' permutations, since (a) we essentially get to see partial entries from the ± 1 matrix \mathbf{M} of pairwise representations of the permutations (note, though, that the $O(n^2)$ pairwise representation is over-complete), and (b) if the users have at most r distinct permutations among themselves, the rank of \mathbf{M} is at most r . Hence, for the random sampling problem, we compare the performance of our algorithms – both in terms of sample complexity and running time – with an Augmented Lagrangian Method version of Matrix Completion (ALM-MC, relevant code from [21]). Finally, we also present an overall comparison of sample complexities across both random and active sampling cases and algorithms. This helps to put both sampling methods in perspective, and also illustrates the order-wise gains when the learning algorithm is allowed to sample pairwise orderings from users at will. All the routines run in MATLAB on a 2.4 GHz desktop computer system with 4 GB of memory.

A. Random Sampling: Algorithm 1 and ALM-MC

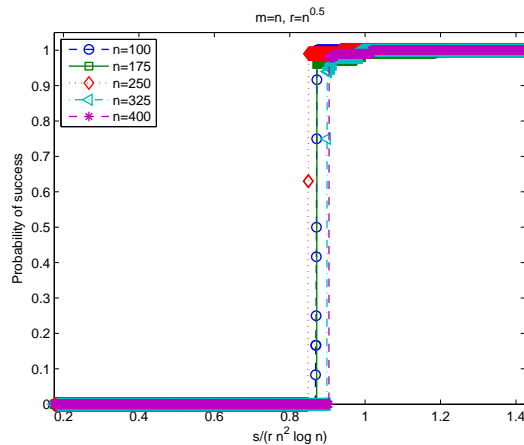
For our random sampling experiments we set $m = n$, varying from 100 to 400, and consider two scaling regimes of r : $r = m^{0.25}$ and $r = m^{0.5}$. We generate random permutations for all m users by first picking r random permutations and then assigning users randomly to these permutations.

Figure 1 shows how the reconstruction success probability of Algorithm 1 scales with the normalized number $s/(rn^2 \log n)$ of random pairwise samples drawn from users. Each of the colored curves represents a fixed value of n , the number of items, and depicts the success probability of Algorithm 1 as the number of random samples s is varied. We observe a sharp “phase transition” effect for the probability of success at this normalized scale of samples – the “correct” normalization suggested by Theorems 1 and 2.

We plot the corresponding probabilities of success for ALM-MC (matrix completion on the pairwise ± 1 matrix) in Figure 2. An important point here is that due to our stochastic model for user permutations, the matrix of pairwise ordering representations of users is at most rank r , and the number of samples required to complete this matrix (and recover all orderings) is $O(rn^{1.2} \log n) = O(rn^{2.4} \log n)$ when incoherence is constant, according to Candès and Recht [2]. Thus, we use this normalizing factor for the number of samples in our plots.



(a)



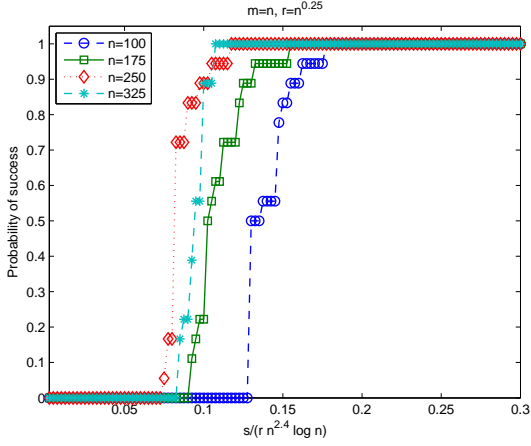
(b)

Fig. 1. Probability of success vs. number of random samples s normalized by $rn^2 \log n$ for Algorithm 1, (a) $r = m^{0.25}$, $m = n$, (b) $r = m^{0.5}$, $m = n$.

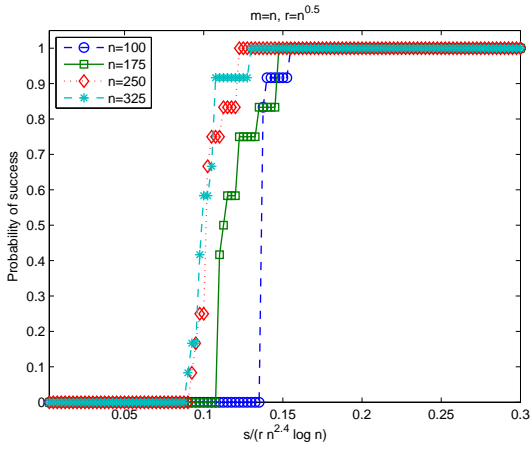
Note the similar phase transition for the success probability for ALM-MC as the number of drawn samples varies. Also, since the phase transition occurs at the $n^{2.4}$ timescale instead of the n^2 timescale for Algorithm 1 (Figure 1), *ALM-MC requires order-wise $n^{0.4}$ more samples than Algorithm 1 to succeed*. Not only is this concordant with the lower bound on sample complexity given by Theorem 2, but it also demonstrates the order-wise superior performance of Algorithm 1 to solve the permutation-learning problem from random samples.

B. Active Sampling: Algorithm 2

We plot the success probability of Algorithm 2 which draws active samples (Figure 3), for the same regime (m, n, r) as in the random sampling case. Here, as indicated by Theorems 3 and 4, the right scale of normalization for the number of samples taken is $rn \log n$ – an n -fold improvement over reconstruction with random



(a)



(b)

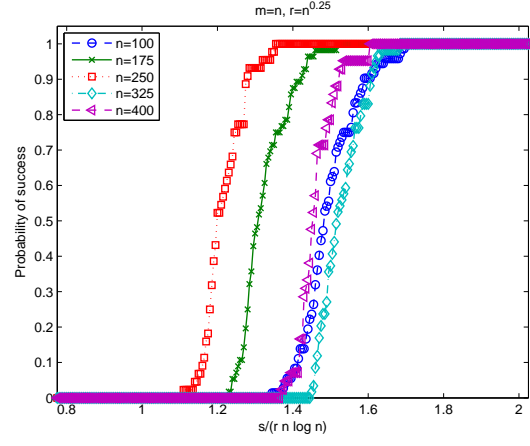
Fig. 2. Probability of success vs. number of random samples s normalized by $rn^{2.4} \log n$ for Augmented Lagrange Multiplier-based Matrix Completion (ALM-MC), (a) $r = m^{0.25}, m = n$, (b) $r = m^{0.5}, m = n$.

sampling. The phase transition for Algorithm 2's success probability is clearly visible in Figure 3.

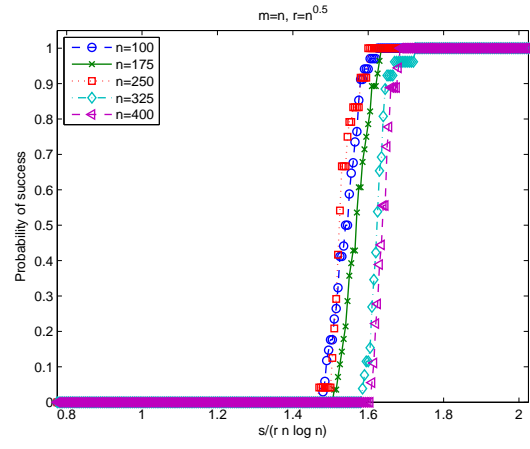
C. Overall Comparison: Sample Complexities and Running Time

To put all the algorithms considered so far in perspective, we compare and contrast their sample complexities and running times together.

Figure 4(a) compares the number of samples at the success probability phase transition (s) against the problem size (n) for Algorithms 1 and 2 and ALM-MC. Algorithm 2 dramatically outperforms its random sampling counterparts, lending support to the active sampling model of attempting to learn user rankings. On the other hand, in the random sampling case, though ALM-MC (matrix completion) fares better than the order-optimal Algorithm 1, we suspect that this is due to overheads



(a)



(b)

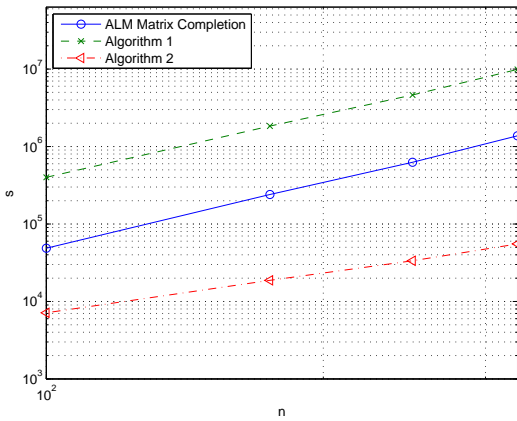
Fig. 3. Probability of success vs. number of samples s normalized by $rn \log n$ for Algorithm 2 (active sampling), (a) $r = m^{0.25}, m = n$, (b) $r = m^{0.5}, m = n$.

occurring at low problem sizes and observe that the curves for ALM-MC and Algorithm 1 are projected to cross over at larger sizes of n .

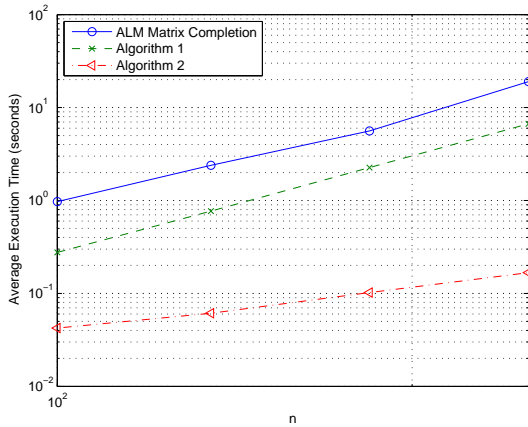
In Figure 4(b) shows the running times on a 2.4 GHz CPU with 4 GB of memory, for our implementations of Algorithms 1 and 2, and the off-the-shelf implementation of ALM-MC. Here, *the standard ALM matrix completion algorithm is outperformed by both Algorithms 1 and 2, illustrating the gain in computational efficiency that these algorithms offer.*

VI. CONCLUSION

We considered the problem of learning a collection of users' permutations of items using just partial pairwise comparisons. Both random and active/intelligent sampling schemes were separately considered. In both cases, we developed efficient algorithms that reconstruct the permutations with a guaranteed sample complexity, and



(a) Experimental sample-complexities vs. n for Algorithms 1, 2 and ALM-MC, $m = n$ and $r = n^{0.5}$.



(b) Experimental execution time in seconds vs. n for Algorithms 1, 2 and ALM-MC, $m = n$ and $r = n^{0.5}$.

Fig. 4. Experimental results on sample and time complexity of all algorithms.

using corresponding lower bounds on sample complexity showed that these algorithms are order-optimal, additionally with an order-wise performance improvement when sampling actively. Moreover, there is a significant gain when solving the problem jointly compared to learning each permutation individually. Experiments were carried out that validated the performance benefits of the algorithms we presented, and in many cases showed their superiority over traditional matrix-completion approaches.

REFERENCES

- [1] R. H. Keshavan, A. Montanari, and S. Oh, "Matrix completion from a few entries," *IEEE Transactions on Information Theory*, vol. 56, pp. 2980–2998, 2010.
- [2] E. J. Candès and B. Recht, "Exact matrix completion via convex optimization," *Found. Comput. Math.*, vol. 9, pp. 717–772, December 2009.
- [3] B. Recht, "A simpler approach to matrix completion," *CoRR*, vol. abs/0910.0651, 2009.

- [4] M. Weimer, A. Karatzoglou, Q. V. Le, and A. Smola, "CoFiRank, Maximum Margin Matrix Factorization for Collaborative Ranking," in *Advances in Neural Information Processing Systems*, vol. 20, 2007.
- [5] B. Zelinka, "Distances between partially ordered sets," *Mathematica Bohemica*, vol. 118, no. 2, pp. 167–170, 1993.
- [6] A. Haviar and B. Bystrica, "A metric on a system of ordered sets," *Mathematica Bohemica*, vol. 121, no. 2, pp. 123–131, 1996.
- [7] Y. Freund, R. Iyer, R. E. Schapire, and Y. Singer, "An efficient boosting algorithm for combining preferences," *J. Mach. Learn. Res.*, vol. 4, pp. 933–969, December 2003.
- [8] W. W. Cohen, R. E. Schapire, and Y. Singer, "Learning to order things," *J. Artif. Intell. Res. (JAIR)*, vol. 10, pp. 243–270, 1999.
- [9] S. Agarwal, "Learning to rank on graphs," *Machine Learning*, vol. 81, no. 3, pp. 333–357, 2010.
- [10] S. Rajaram and S. Agarwal, "Generalization bounds for k -partite ranking," in *Proceedings of the NIPS-2005 Workshop on Learning to Rank*, 2005.
- [11] D. Helmbold and M. Warmuth, "Learning permutations with exponential weights," in *Proceedings of the 20th annual conference on Learning theory*. Springer-Verlag, 2007, pp. 469–483.
- [12] G. Fung, R. Rosales, and B. Krishnapuram, "Learning rankings via convex hull separation," *Advances in Neural Information Processing Systems*, vol. 18, p. 395, 2006.
- [13] G. Lebanon and J. Lafferty, "Cranking: Combining rankings using conditional probability models on permutations," in *Proceedings of the Nineteenth International Conference on Machine Learning*. Morgan Kaufmann Publishers Inc., 2002, pp. 363–370.
- [14] M. Ajtai, V. Feldman, A. Hassidim, and J. Nelson, "Sorting and selection with imprecise comparisons," *Automata, Languages and Programming*, pp. 37–48, 2009.
- [15] U. Feige, P. Raghavan, D. Peleg, and E. Upfal, "Computing with noisy information," *SIAM J. Comput.*, vol. 23, no. 5, pp. 1001–1018, 1994.
- [16] C. Daskalakis, R. Karp, E. Mossel, S. Riesenfeld, and E. Verbin, "Sorting and selection in posets," in *Proceedings of the twentieth Annual ACM-SIAM Symposium on Discrete Algorithms*. Society for Industrial and Applied Mathematics, 2009, pp. 392–401.
- [17] S. Jagabathula and D. Shah, "Inferring rankings under constrained sensing," in *Advances in Neural Information Processing Systems*, 2008, pp. 753–760.
- [18] T. Cover and J. Thomas, "Elements of information theory 2nd edition," 2006.
- [19] R. Yeung, *Information theory and network coding*. Springer Verlag, 2008.
- [20] D. Dubhashi, A. Panconesi, and C. U. Press, *Concentration of measure for the analysis of randomized algorithms*. Cambridge University Press, 2009.
- [21] Z. Lin, M. Chen, L. Wu, and Y. Ma, "The augmented Lagrange multiplier method for exact recovery of corrupted low-rank matrices," *Arxiv preprint arXiv:1009.5055*, 2010.