

Learning sparse mixtures of rankings from noisy information

Anindya De*
Northwestern University
anindya@eecs.northwestern.edu

Ryan O’Donnell†
Carnegie Mellon University
odonnell@cs.cmu.edu

Rocco A. Servedio‡
Columbia University
rocco@cs.columbia.edu

November 6, 2018

Abstract

We study the problem of learning an unknown mixture of k rankings over n elements, given access to noisy samples drawn from the unknown mixture. We consider a range of different noise models, including natural variants of the “heat kernel” noise framework and the Mallows model. For each of these noise models we give an algorithm which, under mild assumptions, learns the unknown mixture to high accuracy and runs in $n^{O(\log k)}$ time. The best previous algorithms for closely related problems have running times which are exponential in k .

*Supported by NSF grant CCF-1814706

†Supported by NSF grants CCF-1618679 and CCF-1717606

‡Supported by NSF grants CCF-1563155 and CCF-1814873 and by the Simons Collaboration on Algorithms and Geometry. This material is based upon work supported by the National Science Foundation under grant numbers listed above. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

1 Introduction

This paper considers the following natural scenario: there is a large heterogeneous population which consists of k disjoint subgroups, and for each subgroup there is a “central preference order” specifying a ranking over a fixed set of n items (equivalently, specifying a permutation in the symmetric group \mathbb{S}_n). For each $i \in \{1, \dots, k\}$, the preference order of each individual in subgroup i is assumed to be a noisy version of the central preference order (the permutation corresponding to subgroup i). A natural learning task which arises in this scenario is the following: given access to the preference order of randomly selected members of the population, is it possible to learn the central preference orders of the k sub-populations, as well as the relative sizes of these k sub-populations within the overall population?

Worst-case formulations of the above problem typically tend to be (difficult) variants of the feedback arc set problem, which is known to be NP-complete [GJ79]. In view of the practical importance of problems of this sort, though, there has been considerable recent research interest in studying various generative models corresponding to the above scenario (we discuss some of the recent work which is most closely related to our results in [Section 1.3](#)). In this paper we will model the above general problem schema as follows: The k “central preference orders” of the subgroups are given by k unknown permutations $\sigma_1, \dots, \sigma_k \in \mathbb{S}_n$. The fraction of the population belonging to the i -th subgroup, for $1 \leq i \leq k$, is given by an unknown $w_i \geq 0$ (so $w_1 + \dots + w_k = 1$). Finally, the noise is modeled by some family of distributions $\{\mathcal{K}_\theta\}$, where each distribution \mathcal{K}_θ is supported on \mathbb{S}_n , and the preference order of a random individual in the i -th subgroup is given by $\pi\sigma_i$, where $\pi \sim \mathcal{K}_\theta$. Here θ is a model parameter capturing the “noise rate” (we will have much more to say about this for each of the specific noise models we consider below). The learning task is to recover the central rankings $\sigma_1, \dots, \sigma_k$ and their proportions w_1, \dots, w_k , given access to preference orders of randomly chosen individuals from the population. In other words, each sample provided to the learner is independently generated by first choosing a random permutation σ , where σ is chosen to be σ_i with probability w_i ; then independently drawing a random $\pi \sim \mathcal{K}_\theta$; and finally, providing the learner with the permutation $\pi\sigma \in \mathbb{S}_n$. Let $f : \mathbb{S}_n \rightarrow \mathbb{R}^{\geq 0}$ denote the function which is w_i at σ_i and 0 otherwise. With this notation, we write “ $\mathcal{K}_\theta * f$ ” to denote the distribution over noisy samples described above, and our goal is to approximately recover f given such noisy samples. The reader may verify that the distribution defined by $\pi\sigma$ is precisely given by the group convolution $\mathcal{K}_\theta * f$ (and hence the notation).

1.1 The noise models that we consider

We consider a range of different noise models, corresponding to different choices for the parametric family $\{\mathcal{K}_\theta\}$, and for each one we give an efficient algorithm for recovering the population in the presence of that kind of noise. In this subsection we detail the three specific noise models that we will work with (though as we discuss later, our general mode of analysis could be applied to other noise models as well).

(A.) **Symmetric noise.** In the *symmetric noise* model, the parametric family of distributions over \mathbb{S}_n is denoted $\{\mathcal{S}_{\bar{p}}\}_{\bar{p} \in \Delta^n}$. Given a vector $\bar{p} = (p_0, \dots, p_n) \in \Delta^n$ (so each $p_i \geq 0$ and $\sum_{i=0}^n p_i = 1$), a draw of $\pi \sim \mathcal{S}_{\bar{p}}$ is obtained as follows:

1. Choose $0 \leq j \leq n$, where value j is chosen with probability p_j .

2. Choose a uniformly random subset $\mathbf{A} \subseteq [n]$ of size exactly j . Draw $\boldsymbol{\pi}$ uniformly from $\mathbb{S}_{\mathbf{A}}$; in other words, $\boldsymbol{\pi}$ is a uniformly random permutation over the set \mathcal{A} and is the identity permutation on elements in $[n] \setminus \mathcal{A}$. (We denote this uniform distribution over $\mathbb{S}_{\mathbf{A}}$ by $\mathbb{U}_{\mathbf{A}}$.)

Note that in this model, if the noise vector \bar{p} has $p_n = 1$, then every draw from $\mathcal{S}_{\bar{p}} * f$ is a uniform random permutation and there is no useful information available to the learner.

In order to define the next two noise models that we consider, let us recall the notion of a *right-invariant* metric on \mathbb{S}_n . Such a metric $d(\cdot, \cdot)$ is one that satisfies $d(\sigma, \pi) = d(\sigma\tau, \pi\tau)$ for all $\sigma, \pi, \tau \in \mathbb{S}_n$. We note that a metric is right-invariant if and only if it is invariant under relabeling of the items $1, \dots, n$, and that most metrics considered in the literature satisfy this condition (see [KV10, Dia88b] for discussions of this point). In this paper, for technical convenience we restrict our attention to the metric $d(\cdot, \cdot)$ being the *Cayley distance* over \mathbb{S}_n (though see Section 1.5 for a discussion of how our methods and results could potentially be generalized to other right-invariant metrics):

Definition 1.1. Let G be the undirected graph with vertex set \mathbb{S}_n and an edge between permutations σ and π if there is a transposition τ such that $\sigma = \tau \cdot \pi$. The *Cayley distance* over \mathbb{S}_n is the metric induced by this graph; in other words, $d(\pi, \sigma) = t$ where t is the smallest value such that there are transpositions τ_1, \dots, τ_t satisfying $\sigma = \tau_1 \cdots \tau_t \pi$.

Now we are ready to define the next two parameterized families of noise distributions that we consider. We note that each of the noise distributions \mathcal{K} considered below has the natural property that $\Pr_{\boldsymbol{\pi} \sim \mathcal{K}}[\boldsymbol{\pi} = \pi]$ decreases with $d(\pi, e)$ where e is the identity distribution.

(B.) **Heat kernel random walk under Cayley distance.** Let \mathcal{L} be the Laplacian of the graph G from Definition 1.1. Given a “temperature” parameter $t \in \mathbb{R}^+$, the *heat kernel* is the $n! \times n!$ matrix $H_t = e^{-t\mathcal{L}}$. It is well known that H_t is the transition matrix of the random walk induced by choosing a Poisson-distributed time parameter $\mathbf{T} \sim \text{Poi}(t)$ and then taking \mathbf{T} steps of a uniform random walk in the graph G . With this motivation, we define the *heat kernel noise model* as follows: the parametric family of distributions is $\{\mathcal{H}_t\}_{t \in \mathbb{R}^+}$, where the probability weight that \mathcal{H}_t assigns to permutation π is the probability that the above-described random walk, starting at the identity permutation $e \in \mathbb{S}_n$, reaches π . (Observe that higher temperature parameters t correspond to higher rates of noise. More precisely, it is well known that the mixing time of a uniform random walk on G is $\Theta(n \log n)$ steps, so if t grows larger than $n \log n$ then the distribution \mathcal{H}_t converges rapidly to the uniform distribution on \mathbb{S}_n ; see [DS81] for detailed results along these lines.) We note that these probability distributions (or more precisely, the associated heat kernel H_t) have been previously studied in the context of learning rankings, see e.g. [KL02, KB10, JV18]. In some of this work, a different underlying distance measure was used over \mathbb{S}_n rather than the Cayley distance; see our discussion of related work in Section 1.3.

(C.) **Mallows-type model under Cayley distance (Cayley-Mallows / Ewens model).** While the heat kernel noise model arises naturally from an analyst’s perspective, a somewhat different model, called the *Mallows model*, has been more popular in the statistics and machine learning literature. The Mallows model is defined using the “Kendall τ -distance” $K(\cdot, \cdot)$ between permutations (defined in Section 1.3) rather than the Cayley distance $d(\cdot, \cdot)$; the Mallows model with parameter $\theta > 0$ assigns probability weight $e^{-\theta K(\pi, e)} / Z_K(\theta)$ to the permutation π , where $Z_K(\theta) = \sum_{\pi \in \mathbb{S}_n} e^{-\theta K(\pi, e)}$ is a normalizing constant. As proposed by Fligner and Verducci [FV86],

it is natural to consider generalizations of the Mallows model in which other distance measures take the place of the Kendall τ -distance. The model which we consider is one in which the Cayley distance is used as the distance measure; so given $\theta > 0$, the noise distribution \mathcal{M}_θ which we consider assigns weight $e^{-\theta d(\pi, e)}/Z(\theta)$ to each permutation $\pi \in \mathbb{S}_n$, where $Z(\theta) = \sum_{\pi \in \mathbb{S}_n} e^{-\theta d(\pi, e)}$ is a normalizing constant. In fact, this noise model was already proposed in 1972 by W. Ewens in the context of population genetics [Ewe72] and has been intensively studied in that field (we note that [Ewe72] has been cited more than 2000 times according to Google Scholar). To align our terminology with the strand of research in machine learning and theoretical computer science which deals with the Mallows model, in the rest of this paper we refer to \mathcal{M}_θ as the *Cayley-Mallows* model. For the same reason, we will also refer to the usual Mallows model (with the Kendall τ -distance) as the *Kendall-Mallows* model. We observe that for the Cayley-Mallows model \mathcal{M}_θ , in contrast with the heat kernel noise model now *smaller* values of θ correspond to higher levels of noise, and that when $\theta = 0$ the distribution \mathcal{M}_θ is simply the uniform distribution over \mathbb{S}_n and there is no useful information available to the learner.

1.2 Our results

For each of the noise models defined above, we give algorithms which, under a mild technical assumption (that no mixing weight w_i is too small), provably recover the unknown central rankings $\sigma_1, \dots, \sigma_k$ and associated mixing weights w_1, \dots, w_k up to high accuracy. A notable feature of our results is that the sample and running time dependence is only *quasipolynomial* in the number of elements n and the number of sub-populations k ; as we detail in [Section 1.3](#) below, this is in contrast with recent results for similar problems in which the dependence on k is exponential.

Below we give detailed statements of our results. The following notation and terminology will be used in these statements: for f a distribution over \mathbb{S}_n (or any function from \mathbb{S}_n to \mathbb{R}) we write $\text{supp}(f)$ to denote the set of permutations $\sigma \in \mathbb{S}_n$ that have $f(\sigma) \neq 0$. For a given noise model \mathcal{K} , we write “ $\mathcal{K} * f$ ” to denote the distribution over noisy samples that is provided to the learning algorithm as described earlier. Given two functions $f, g : \mathbb{S}_n \rightarrow \mathbb{R}$, we write “ $\|f - g\|_1$ ” to denote $\sum_{\pi \in \mathbb{S}_n} |f(\pi) - g(\pi)|$, the ℓ_1 distance between f and g . If f and g are both distributions then we write $d_{\text{TV}}(f, g)$ to denote the *total variation distance* between f and g , which is $\frac{1}{2}\|f - g\|_1$. Finally, if f is a distribution over \mathbb{S}_n in which $f(\sigma) > \varepsilon$ for every σ such that $f(\sigma) > 0$, we say that f is ε -heavy.

Learning from noisy rankings: Positive and negative results. Our first algorithmic result is for the symmetric noise model (A) defined earlier. [Theorem 1.2](#), stated below, gives an efficient algorithm as long as the vector \bar{p} is “not too extreme” (i.e. not too biased towards putting almost all of its weight on large values very close to n):

Theorem 1.2 (Algorithm for symmetric noise). *There is an algorithm with the following guarantee: Let f be an unknown ε -heavy distribution over \mathbb{S}_n with $|\text{supp}(f)| \leq k$. Let $\bar{p} = (p_0, \dots, p_n) \in \Delta^n$ be such that*

$$\sum_{j=0}^{n-\log k} p_j \geq \frac{1}{n^{O(\log k)}}.$$

*Given \bar{p} , the value of $\varepsilon > 0$, a confidence parameter $\delta > 0$, and access to random samples from $\mathcal{S}_{\bar{p}} * f$, the algorithm runs in time $\text{poly}(n^{\log k}, 1/\varepsilon, \log(1/\delta))$ and with probability $1 - \delta$ outputs a distribution $g : \mathbb{S}_n \rightarrow \mathbb{R}$ such that $d_{\text{TV}}(f, g) \leq \varepsilon$.*

Our second algorithmic result, which is similar in spirit to [Theorem 1.2](#), is for the heat kernel noise model:

Theorem 1.3 (Algorithm for heat kernel noise). *There is an algorithm with the following guarantee: Let f be an unknown ε -heavy distribution over \mathbb{S}_n with $|\text{supp}(f)| \leq k$. Let $t \in \mathbb{R}^+$ be any value that is $O(n \log n)$. Given t , the value of $\varepsilon > 0$, a confidence parameter $\delta > 0$, and access to random samples from $\mathcal{H}_t * f$, the algorithm runs in time $\text{poly}(n^{\log k}, 1/\varepsilon, \log(1/\delta))$ and with probability $1 - \delta$ outputs a distribution $g : \mathbb{S}_n \rightarrow \mathbb{R}$ such that $d_{\text{TV}}(f, g) \leq \varepsilon$.*

Recalling that the uniform random walk on the Cayley graph of \mathbb{S}_n mixes in $\Theta(n \log n)$ steps, we see that the algorithm of [Theorem 1.3](#) is able to handle quite high levels of noise and still run quite efficiently (in quasi-polynomial time).

Our third positive result, for the Cayley-Mallows model, displays an intriguing qualitative difference from [Theorems 1.2](#) and [1.3](#). To state our result, let us define the function $\text{dist} : \mathbb{R}^+ \times \mathbb{N} \rightarrow \mathbb{R}^+$ as follows:

$$\text{dist}(\theta, \ell) := \min_{j \in \{1, \dots, \ell\}} |e^\theta - j|,$$

so $\text{dist}(\theta, \ell)$ measures the minimum distance between e^θ and any integer in $\{1, \dots, \ell\}$. [Theorem 1.4](#) gives an algorithm which can be quite efficient for the Cayley-Mallows noise model if the noise parameter θ is such that $\text{dist}(\theta, \log k)$ is not too small:

Theorem 1.4 (Algorithm for the Cayley-Mallows model). *There is an algorithm with the following guarantee: Let f be an unknown ε -heavy distribution over \mathbb{S}_n with $|\text{supp}(f)| \leq k$. Given $\theta > 0$, the value of $\varepsilon > 0$, a confidence parameter $\delta > 0$, and access to random samples from $\mathcal{M}_\theta * f$, the algorithm runs in time $\text{poly}(n^{\log k}, 1/\varepsilon, \log(1/\delta), \text{dist}(\theta, \log k)^{-\sqrt{\log k}})$ and with probability $1 - \delta$ outputs a distribution $g : \mathbb{S}_n \rightarrow \mathbb{R}$ such that $d_{\text{TV}}(f, g) \leq \varepsilon$.*

As alluded to earlier, as θ approaches 0 the difficulty of learning in the \mathcal{M}_θ noise model increases (and indeed learning becomes impossible at $\theta = 0$); since for small θ we have $\text{dist}(\theta, \ell) \approx \theta$, this is accounted for by the $\text{dist}(\theta, \log k)^{-\sqrt{\log k}}$ factor in our running time bound above. However, for larger values of θ the $\text{dist}(\theta, \log k)^{-\sqrt{\log k}}$ dependence may strike the reader as an unnatural artifact of our analysis: is it really hard to learn when θ is very close to $\ln 2 \approx 0.63147$, easy when θ is very close to $\ln 2.5 \approx 0.91629$, and hard again when θ is very close to $\ln 3 \approx 1.09861$? Perhaps surprisingly, the answer is yes: it turns out that the $\text{dist}(\cdot, \cdot)$ parameter captures a *fundamental barrier* to learning in the Cayley-Mallows model. We establish this by proving the following lower bound for the Cayley-Mallows model, which shows that a dependence on dist as in [Theorem 1.4](#) is in fact inherent in the problem:

Theorem 1.5. *Given $j \in \mathbb{N}$, there are infinitely many values of k and $m = m(k)$ such that the following holds: Let $\theta > 0$ be such that $|e^\theta - j| \leq \eta \leq 1/2$, and let A be any algorithm which, when given access to random samples from $\mathcal{M}_\theta * f$ where f is a distribution over \mathbb{S}_m with $|\text{supp}(f)| \leq k$, with probability at least 0.51 outputs a distribution h over \mathbb{S}_m that has $d_{\text{TV}}(f, h) \leq 0.99$. Then A must use $\eta^{-\Omega\left(\sqrt{\frac{\log k}{\log \log k}}\right)}$ samples.*

1.3 Relation to prior work

Starting with the work of Mallows [[Mal57](#)], there is a rich line of work in machine learning and statistics on probabilistic models of ranking data, see e.g. [[Mar14](#), [LL02](#), [BOB07](#), [MM09](#), [MC10](#),

LB11]. In order to describe the prior works which are most relevant to our paper, it will be useful for us to define the *Kendall-Mallows* model (referred to in the literature just as the Mallows model) in slightly more detail than we gave earlier. Introduced by Mallows [Mal57], the Kendall-Mallows model is quite similar to the Cayley-Mallows model that we consider — it is specified by a parametric family of distributions $\{\mathcal{M}_{\tau,\theta}\}_{\theta \in \mathbb{R}^+}$ and a central permutation $\sigma \in \mathbb{S}_n$, and a draw from the model is generated as follows: sample $\pi \sim \mathcal{M}_{\tau,\theta}$ and output $\pi \cdot \sigma$. The distribution $\mathcal{M}_{\tau,\theta}$ assigns probability weight $e^{-\theta K(\pi,e)}/Z_K(\theta)$ to the permutation π where $Z_K(\theta) = \sum_{\pi \in \mathbb{S}_n} e^{-\theta K(\pi,e)}$ is the normalizing constant and $K(\cdot, \cdot)$ is the Kendall τ -distance (defined next):

Definition 1.6. The *Kendall τ -distance* $K : \mathbb{S}_n \times \mathbb{S}_n \rightarrow \mathbb{R}^{\geq 0}$ is a distance metric on \mathbb{S}_n defined as

$$K(\pi, \pi') = \{(i, j) : i < j \text{ and } ((\pi(i) < \pi(j)) \oplus (\pi'(i) < \pi'(j)) = 1)\}$$

In other words, $K(\pi, \pi')$ is the number of inversions between π and π' . Like the Cayley distance, the Kendall τ -distance is also a right-invariant metric. Another equivalent way to define $K(\cdot, \cdot)$ is to consider the undirected graph on \mathbb{S}_n where vertices π_1, π_2 share an edge if and only if $\pi_1 = \tau \cdot \pi_2$ where τ is an *adjacent transposition* — in other words, $\tau = (i, i+1)$ for some $1 \leq i < n$. Then $K(\cdot, \cdot)$ is defined as the shortest path metric on this graph. From this perspective, the difference between the Kendall τ -distance and the Cayley distance is that the former only allows adjacent transpositions while the latter allows all transpositions.

Learning mixture models: As mentioned earlier, probabilistic models of ranking data have been studied extensively in probability, statistics and machine learning. Models that have been considered in this context include the Kendall-Mallows model [Mal57, LB11, MPPB07, GP18], the Cayley-Mallows model (and generalizations of it) [FV86, MM03, Muk16, DH92, Dia88a, Ewe72] and the heat kernel random walk model [KL02, KB10, JV18], among others. In contrast, within theoretical computer science interest in probabilistic models of ranking data is somewhat more recent, and the best-studied model in this community is the Kendall-Mallows model. Braverman and Mossel [BM08] initiated this study and (among other results) gave an efficient algorithm to recover a single Kendall-Mallows model from random samples. The question of learning mixtures of k Kendall-Mallows models was raised soon thereafter, and Awasthi *et al.* [ABSV14] gave an efficient algorithm for the case $k = 2$. We note two key distinctions between our work and that of [ABSV14]: (i) our results apply to the Cayley-Mallows model rather than the Kendall-Mallows model, and (ii) the work of [ABSV14] allows for the two components in the mixture to have two different noise parameters θ_1 and θ_2 whereas our mixture models allow for only one noise parameter θ across all the components.

Very recently, Liu and Moitra [LM18] extended the result of [ABSV14] to any constant k . In particular, the running time of the [LM18] algorithm scales as $n^{\text{poly}(k)}$. It is interesting to contrast our results with those of [LM18]. Besides the obvious difference in the models treated (namely Kendall-Mallows in [LM18] versus Cayley-Mallows in this paper), another significant difference is that our running time scales only quasipolynomially in k versus exponentially in k for [LM18]. (In fact, [LM18] shows that an exponential dependence on k is necessary for the problem they consider.) Another difference is that their algorithm allows each mixture component to have a different noise parameter θ_i whereas our result requires the same noise parameter θ across the mixture components. We observe that one curious feature of the algorithm of [LM18] is the following: When all the noise parameters $\{\theta_i\}_{1 \leq i \leq k}$ are *well-separated* (meaning that for all $i \neq j$, $|\theta_i - \theta_j| \geq \gamma$), then the running

time of [LM18] can be improved to $\text{poly}(n) \cdot 2^{\text{poly}(k)}$. This suggests that the case when all θ_i are the same might be the hardest for the Liu-Moitra [LM18] algorithm.

Finally, we note that while the analysis in this paper does not immediately extend to the Kendall-Mallows model (see Section 1.5 for more details), we point out that there is a sense in which the Kendall-Mallows and Cayley-Mallows models are fundamentally incomparable. This is because, while the results of [LM18] show that mixtures of Kendall-Mallows models are identifiable whenever each $\theta_i \neq 1$, Theorem 1.5 shows that mixtures of Cayley-Mallows models are not identifiable at various larger values of θ such as $\ln 2, \ln 3, \dots$, even when all of the noise parameters are the same value θ which is provided to the algorithm.

1.4 Our techniques

A key notion for our algorithmic approach is that of the *marginal* of a distribution f over \mathbb{S}_n :

Definition 1.7. Fix $f : \mathbb{S}_n \rightarrow [0, 1]$ to be some distribution over \mathbb{S}_n . Let $t \in \{1, \dots, n\}$, let $\bar{i} = (i_1, \dots, i_t)$ be a vector of t distinct elements of $\{1, \dots, n\}$ and likewise $\bar{j} = (j_1, \dots, j_t)$. We say the (\bar{i}, \bar{j}) -*marginal* of f is the probability

$$\Pr_{\sigma \sim f}[\sigma(i_1) = j_1 \text{ and } \dots \text{ and } \sigma(i_t) = j_t]$$

that for all $\ell = 1, \dots, t$, the i_ℓ -th element of a random σ drawn from f is j_ℓ . When \bar{i} and \bar{j} are of length t we refer to such a probability as a t -way *marginal* of f .

The first key ingredient of our approach for learning from noisy rankings is a reduction from the problem of learning f (the unknown distribution supported on k rankings $\sigma_1, \dots, \sigma_k$) given access to samples from $\mathcal{K} * f$, to the problem of estimating t -way marginals (for a not-too-large value of t). More precisely, in Section 2 we give an algorithm which, given the ability to efficiently estimate t -way marginals of f , efficiently computes a high-accuracy approximation for an unknown ε -heavy distribution f with support size at most k (see Theorem 2.1). This algorithm builds on ideas in the population recovery literature, suitably extended to the domain \mathbb{S}_n rather than $\{0, 1\}^n$.

With the above-described reduction in hand, in order to obtain a positive result for a specific noise model \mathcal{K} the remaining task is to develop an algorithm A_{marginal} which, given access to noisy samples from $\mathcal{K} * f$, can reliably estimate the required marginals. In Section 3 we show that if the noise distribution \mathcal{K} (a distribution over \mathbb{S}_n) is efficiently samplable, then given samples from $\mathcal{K} * f$, the time required to estimate the required marginals essentially depends on the minimum, over a certain set of matrices arising from the Fourier transform (over the symmetric group \mathbb{S}_n) of the noise distribution, of the minimum singular value of the matrix. (See Theorem 3.1 for a detailed statement.) At this point, we have reduced the algorithmic problem of obtaining a learning algorithm for a particular noise model to the analytic task of lower bounding the relevant singular values. We carry out the required analyses on a noise-model-by-noise-model basis in Sections 4, 5, and 6. These analyses employ ideas and results from the representation theory of the symmetric group and its connections to enumerative combinatorics; we give a brief overview of the necessary background in Appendix A.

To establish our lower bound for the Cayley-Mallows model, Theorem 1.5, we exhibit two distributions f_1 and f_2 over the symmetric group such that the distributions of noisy rankings $\mathcal{M}_\theta * f_1$ and $\mathcal{M}_\theta * f_2$ have very small statistical distance from each other. Not surprisingly, the inspiration

for this construction also comes from the representation theory of the symmetric group; more precisely, the two above-mentioned distributions are obtained from the character (over the symmetric group) corresponding to a particular carefully chosen partition of $[n]$. A crucial ingredient in the proof is the fact that characters of the symmetric group are rational-valued functions, and hence any character can be split into a positive part and a negative part; details are given in [Section 8](#).

1.5 Discussion and future work

In this paper we have considered three particular noise models — symmetric noise, heat kernel noise, and Cayley-Mallows noise — and given efficient algorithms for these noise models. Looking beyond these specific noise models, though, our approach provides a general framework for obtaining algorithms for learning mixtures of noisy rankings. Indeed, for essentially any efficiently samplable noise distribution \mathcal{K} , given access to samples from $\mathcal{K} * f$ our approach reduces the algorithmic problem of learning f to the analytic problem of lower bounding the minimum singular values of matrices arising from the Fourier transform of \mathcal{K} (see [Theorem 3.1](#)). We believe that this technique may be useful in a broader range of contexts, e.g. to obtain results analogous to ours for the original Kendall-Mallows model or for other noise models.

As is made clear in [Sections 4, 5, and 6](#), the representation-theoretic analysis that we require for our noise models is facilitated by the fact that each of the noise distributions considered in those sections is a *class function* (in other words, the value of the distribution on a given input permutation depends only on the cycle structure of the permutation). Extending the kinds of analyses that we perform to other noise models which are not class functions is a technical challenge that we leave for future work.

2 Algorithmic recovery of sparse functions

The main result of this section is the reduction alluded to in [Section 1.4](#). In more detail, we give an algorithm which, given the ability to efficiently estimate t -way marginals, efficiently computes a high-accuracy approximation for an unknown ε -heavy distribution f with support size at most k :

Theorem 2.1. *Let f be an unknown ε -heavy distribution over \mathbb{S}_n with $|\text{supp}(f)| \leq k$. Suppose there is an algorithm A_{marginal} with the following property: given as input a value $\delta > 0$ and two vectors $\bar{i} = (i_1, \dots, i_t)$ and $\bar{j} = (j_1, \dots, j_t)$ each composed of t distinct elements of $\{1, \dots, n\}$, algorithm A_{marginal} runs in time $T(\delta, t, k, n)$ and outputs an additively $\pm\delta$ -accurate estimate of the (\bar{i}, \bar{j}) -marginal of f (recall [Definition 1.7](#)). Then there is an algorithm A_{learn} with the following property: given the value of ε , algorithm A_{learn} runs in time $\text{poly}(n/\varepsilon, n^{\log k}) \cdot T(\frac{\varepsilon}{2k^{O(\log k)}}, 2 \log k, k^2, n)$ and returns a function $g : \mathbb{S}_n \rightarrow \mathbb{R}^+$ such that $\|f - g\|_1 \leq \varepsilon$.*

Looking ahead, given [Theorem 2.1](#), in order to obtain a positive result for a specific noise model \mathcal{K} the remaining task is to develop an algorithm A_{marginal} which, given access to noisy samples from $\mathcal{K} * f$, can reliably estimate the required marginals. The algorithm is given in [Section 3](#) and the detailed analyses establishing its efficiency for each of the noise models (by bounding minimum singular values of certain matrices arising from each specific noise distribution) is given in [Sections 4, 5, and 6](#).

2.1 A useful structural result

The following structural result on functions from \mathbb{S}_n to \mathbb{R} with small support will be useful for us:

Claim 2.2 (Small-support functions are correlated with juntas). *Fix $1 \leq \ell \leq n$ and let $g : [n]^\ell \rightarrow \mathbb{R}$ be such that $\|g\|_1 = 1$ and $|\text{supp}(g)| \leq k$. There is a subset $U \subseteq [n]$ and a list of values $\alpha_1, \dots, \alpha_{|U|} \in [n]$ such that $|U| \leq \log k$ and*

$$\left| \sum_{x \in [n]^\ell} g(x) \cdot \mathbf{1}[x_i = \alpha_i \text{ for all } i \in U] \right| \geq k^{-O(\log k)}. \quad (1)$$

Claim 2.2 is reminiscent of analogous structural results for functions over $\{0, 1\}^\ell$ which are implicit in the work of [WY12] (specifically, Theorem 1.5 of that work), and indeed **Claim 2.2** can be proved by following the techniques of [WY12]. Michael Saks [Sak18] has communicated to us an alternative, and arguably simpler, argument for the relevant structural result over $\{0, 1\}^\ell$; here we follow that alternative argument (extending it in the essentially obvious way to the domain $[n]^\ell$ rather than $\{0, 1\}^\ell$).

Proof. Let the support of g be $S \subseteq [n]^\ell$. Note that since $|S| \leq k$, there must exist some set of $k' := \min\{k, \ell\}$ coordinates such that any two elements of S differ in at least one of those coordinates. Without loss of generality, we assume that this set is the first k' coordinates $\{1, \dots, k'\}$.

We prove **Claim 2.2** by analyzing an iterative process that iterates over the coordinates $1, \dots, k'$. At the beginning of the process, we initialize a set $\text{Coord}_{\text{live}}$ of “live coordinates” to be $[k']$, initialize a set Constr of constraints to be initially empty, and initialize a set S_{live} of “live support elements” to be the entire support S of g . We will see that the iterative process maintains the following invariants:

- (I1) The coordinates in $\text{Coord}_{\text{live}}$ are sufficient to distinguish between the elements in S_{live} , i.e. any two distinct strings in S_{live} have distinct projections onto the coordinates in $\text{Coord}_{\text{live}}$;
- (I2) The only elements of S that satisfy all the constraints in Constr are the elements of S_{live} .

Before presenting the iterative process we need to define some pertinent quantities. For each coordinate $j \in \text{Coord}_{\text{live}}$ and each index $\alpha \in [n]$, we define

$$\text{Wt}(j, \alpha) := \sum_{x \in S_{\text{live}}: x_j = \alpha} |g(x)|,$$

the *weight* under g of the live support elements x that have $x_j = \alpha$, and we define

$$\text{Num}(j, \alpha) := |\{x \in S_{\text{live}} : x_j = \alpha\}|,$$

the number of live support elements x that have $x_j = \alpha$ (note that $\text{Num}(j, \alpha)$ has nothing to do with g). It will also be useful to have notation for fractional versions of each of these quantities, so we define

$$\text{FracWt}(j, \alpha) := \frac{\text{Wt}(j, \alpha)}{\sum_{x \in S_{\text{live}}} |g(x)|} \quad \text{and} \quad \text{Frac}(j, \alpha) := \frac{\text{Num}(j, \alpha)}{|S_{\text{live}}|}$$

Note that for any $j \in \text{Coord}_{\text{live}}$ we have that $\sum_{\alpha} \text{Num}(j, \alpha) = |S_{\text{live}}|$, or equivalently $\sum_{\alpha} \text{Frac}(j, \alpha) = 1$.

For each coordinate $j \in \text{Coord}_{\text{live}}$, we write $\text{MAJ}(j)$ to denote the element $\beta \in [n]$ which is such that $\text{Num}(j, \beta) \geq \text{Num}(j, \alpha)$ for all $\alpha \in [n]$ (we break ties arbitrarily). Finally, we let $\text{FracWtMaj}(j) = \text{FracWt}(j, \text{MAJ}(j))$.

Now we are ready to present the iterative process:

1. If every $j \in \text{Coord}_{\text{live}}$ has $\text{FracWtMaj}(j) > 1 - \frac{1}{10k'}$ ¹, then halt the process. Otherwise, let j be any element of $\text{Coord}_{\text{live}}$ for which $\text{FracWtMaj}(j) \leq 1 - \frac{1}{10k'}$.
2. For this coordinate j , choose $\alpha \in [n]$ which maximizes the ratio $\frac{\text{FracWt}(j, \alpha)}{\text{Frac}(j, \alpha)}$ (or equivalently, maximizes $\frac{\text{FracWt}(j, \alpha)}{\text{Num}(j, \alpha)}$) subject to $\text{Frac}(j, \alpha) \neq 0$ and $\alpha \neq \text{MAJ}(j)$.
3. Add the constraint $x_j = \alpha$ to Constr , remove j from $\text{Coord}_{\text{live}}$, and remove all x such that $x_j \neq \alpha$ from S_{live} . Go to Step 1.

When the iterative process ends, suppose that the set Constr is $\{x_{j_1} = \alpha_1, \dots, x_{j_\ell} = \alpha_\ell\}$. Then we claim that [Equation \(1\)](#) holds for $U = \{j_1, \dots, j_\ell\}$.

To argue this, we first observe that both invariants (I1) and (I2) are clearly maintained by each round of the iterative process. We next observe that each time a pair (j, α) is processed in Step 3, it holds that $\text{Frac}(j, \alpha) \leq \frac{1}{2}$, and hence each round shrinks S_{live} by a factor of at least 2. Thus, after $\log k$ steps, the set S_{live} must be of size at most 1 and hence the process must halt. (Note that the claimed bound $|U| \leq \log k$ follows from the fact that the process runs for at most $\log k$ stages.)

Next, note that when the process halts, by a union bound over the at most k' coordinates in $\text{Coord}_{\text{live}}$ it holds that

$$\sum_{x \in S_{\text{live}}: x_j = \text{MAJ}(j) \text{ for all } j \in \text{Coord}_{\text{live}}} |g(x)| \geq \frac{9}{10} \cdot \sum_{x \in S_{\text{live}}} |g(x)|.$$

On the other hand, by the first invariant (I1), the cardinality of the set $\{x \in S_{\text{live}} : x_j = \text{MAJ}(j) \text{ for all } j \in \text{Coord}_{\text{live}}\}$ is precisely 1. This immediately implies that almost all of the weight of g , across elements of S_{live} , is on a single element; more precisely, that

$$\left| \sum_{x \in S_{\text{live}}} g(x) \right| \geq \frac{4}{5} \cdot \sum_{x \in S_{\text{live}}} |g(x)|,$$

from which it follows that

$$\left| \sum_{x \in [n]^\ell} g(x) \cdot \mathbb{1}[x_i = \alpha_i \text{ for all } i \in U] \right| \geq \frac{4}{5} \cdot \sum_{x \in S_{\text{live}}} |g(x)|. \quad (2)$$

So to establish [Equation \(1\)](#), it remains only to establish a lower bound on $\sum_{x \in S_{\text{live}}} |g(x)|$ when the process terminates. To do this, let us suppose that the process runs for T steps where in the

¹Note that this means almost all of the weight under g of the live support elements is on elements that all agree with the majority value on coordinate j . Note further that if $\text{Coord}_{\text{live}}$ is empty then this condition trivially holds.

t^{th} step the coordinate chosen is j_t . Now, at any stage t , we have

$$\frac{\sum_{\beta \in \text{Coord}_{\text{live}}: \beta \neq \text{MAJ}(j_t)} \text{FracWt}(j_t, \beta)}{\sum_{\beta \in \text{Coord}_{\text{live}}: \beta \neq \text{MAJ}(j_t)} \text{Frac}(j_t, \beta)} \geq \frac{1}{10k'}.$$

(because the denominator is at most 1 and since the process does not terminate, the numerator is at least $\frac{1}{10k}$). As a result, we get that if the constraint chosen at time t is $x_{j_t} = \alpha_t$, then

$$\frac{\text{FracWt}(j_t, \alpha_t)}{\text{Frac}(j_t, \alpha_t)} \geq \frac{1}{10k'}. \quad (3)$$

By [Equation \(3\)](#), when the process halts we have

$$\sum_{x \in S_{\text{live}}} |g(x)| = \prod_{t=1}^T \text{FracWt}(j_t, \alpha_t) \geq \frac{1}{(10k')^T} \prod_{t=1}^T \text{Frac}(j_t, \alpha_t).$$

But since at least one element remains, we have that $\prod_{t=1}^T \text{Frac}(j_t, \alpha_t) \geq \frac{1}{k}$, and since $T \leq \log k$, we conclude (recalling that $k' \leq k$) that

$$\sum_{x \in S_{\text{live}}} |g(x)| \geq k^{-O(\log k)}.$$

Combining with [\(2\)](#), this yields the claim. \square

2.2 Proof of [Theorem 2.1](#)

The idea of the proof is quite similar to the algorithmic component of several recent works on population recovery [[MS13](#), [WY12](#), [LZ15](#), [DST16](#)]. Given any function $f : \mathbb{S}_n \rightarrow \mathbb{R}$ and any integer $i \in \{1, \dots, n\}$, we define the function $f_i : [n]^i \rightarrow \mathbb{R}$ as follows:

$$f_i(x_1, \dots, x_i) := \sum_{\sigma \in \mathbb{S}_n} f(\sigma) \cdot \mathbb{1}[\sigma(1) = x_1 \wedge \dots \wedge \sigma(i) = x_i]. \quad (4)$$

At a high level, the algorithm A_{learn} of [Theorem 2.1](#) works in stages, by successively reconstructing f_0, \dots, f_n . In each stage it uses the procedure described in the following claim, which says that high-accuracy approximations of the $(\log k)$ -marginals *together with the support of f_ℓ* (or a not-too-large superset of it) suffices to reconstruct f_ℓ :

Claim 2.3. *Let f_ℓ be an unknown distribution over $[n]^\ell$ supported on a given set S of size k . There is an algorithm $A_{\text{one-stage}}$ which has the following guarantee: The algorithm is given as input $\delta > 0$, and parameters $\beta_{J,y}$ (for every set $J \subseteq [\ell]$ of size at most $\log k$ and every $y \in [n]^J$) which satisfy*

$$\left| \beta_{J,y} - \sum_{x \in S} f(x) \cdot \mathbb{1}[x_i = y_i \text{ for all } i \in J] \right| \leq \delta.$$

$A_{\text{one-stage}}$ runs in time $\text{poly}(n, \ell^{\log k})$ and outputs a function $\tilde{f} : [n]^\ell \rightarrow [0, 1]$ such that $\|f - \tilde{f}\|_1 \leq \delta \cdot k^{O(\log k)}$.

Proof. We consider a linear program which has a variable s_x for each $x \in S$ (representing the probability that f puts on x) and is defined by the following constraints:

1. $s_x \geq 0$ and $\sum_{x \in S} s_x = 1$.
2. For each $J \subseteq [\ell]$ of size at most $\log k$ and each $y \in [n]^J$, include the constraint

$$\left| \beta_{J,y} - \sum_{x \in S} s_x \cdot \mathbf{1}[x_i = y_i \text{ for all } i \in J] \right| \leq \delta. \quad (5)$$

Algorithm $A_{\text{one-stage}}$ sets up and solves the above linear program (this can clearly be done in time $\text{poly}(n, \ell^{\log k})$). We observe that the linear program is feasible since by definition $s_x = f_\ell(x)$ is a feasible solution. To prove the claim it suffices to show that every feasible solution is ℓ_1 -close to f_ℓ ; so let $f^*(x)$ denote any other feasible solution to the linear program, and let η denote $\|f^* - f_\ell\|_1$. Define $h(x) = f^*(x) - f_\ell(x)$, so $\|h\|_1 = \eta$. By [Claim 2.2](#), we have that there is a subset $J \subseteq [\ell]$ of size at most $\log k$ and a $y \in [n]^J$ such that

$$\left| \sum_x h(x) \cdot \mathbf{1}[x_i = y_i \text{ for all } i \in J] \right| \geq \eta \cdot k^{-O(\log k)}. \quad (6)$$

On the other hand, since both $f_\ell(x)$ and $f^*(x)$ are feasible solutions to the linear program, by the triangle inequality it must be the case that

$$\left| \sum_x h(x) \cdot \mathbf{1}[x_i = y_i \text{ for all } i \in J] \right| \leq 2\delta. \quad (7)$$

Equations [6](#) and [2.2](#) together give the desired upper bound on η , and the claim is proved. \square

Essentially the only remaining ingredient required to prove [Theorem 2.1](#) is a procedure to find (a not-too-large superset of) the support of f . This is given by the following claim, which inductively uses the algorithm $A_{\text{one-stage}}$ to successively construct suitable (approximations of) the support sets for f_1, \dots, f_n .

Claim 2.4. *Under the assumptions of [Theorem 2.1](#), there is an algorithm A_{support} with the following property: given as input a value $\delta > 0$, algorithm A_{support} runs in time $\text{poly}(n/\varepsilon, n^{\log k}) \cdot T(\frac{\varepsilon}{2k^{O(\log k)}}, 2 \log k, k^2, n)$ and for each $\ell = 1, \dots, n$ outputs a set $S'_{(\ell)}$ of size at most k which contains the support of f_ℓ .*

Proof. The algorithm A_{support} works inductively, where at the start of stage ℓ (in which it will construct the set $S'_{(\ell)}$) it is assumed to have a set $S'_{(\ell-1)}$ with $|S'_{(\ell-1)}| \leq k$ which contains the support of $f_{\ell-1}$. (Note that at the start of the first stage $\ell = 1$ this holds trivially since f_0 trivially has empty support).

Let us describe the execution of the ℓ -th stage of A_{support} . For $1 \leq \ell \leq n$, we define the set $S_{\text{marg},\ell}$ as follows:

$$S_{\text{marg},\ell} = \left\{ t : \sum_{\sigma \in \mathbb{S}_n} f(\sigma) \cdot \mathbf{1}[\sigma(\ell) = t] > 0 \right\}.$$

Observe that in time $\text{poly}(n/\varepsilon) \cdot T(\frac{\varepsilon}{4}, 1, k, n)$, we can compute $f(\sigma) \cdot \mathbb{1}[\sigma(\ell) = t]$ up to error $\pm\varepsilon/4$ (denote this estimate by $\beta_{\ell,t}$) for all $1 \leq t \leq n$. Since f is ε -heavy, we have that

$$t \in S_{\text{marg},\ell} \text{ implies } \beta_{\ell,t} \geq \frac{3\varepsilon}{4} \quad \text{and} \quad t \notin S_{\text{marg},\ell} \text{ implies } \beta_{\ell,t} \leq \frac{\varepsilon}{4}.$$

Consequently, we can compute the set $S_{\text{marg},\ell}$ in time $\text{poly}(n/\varepsilon) \cdot T(\frac{\varepsilon}{4}, 1, k, n)$. The final observation is that the set $S_{(\ell)}^*$ (of cardinality at most k^2) obtained by appending each final ℓ -th character from $S_{\text{marg},\ell}$ to each element of $S'_{(\ell-1)}$ must contain the support $S_{(\ell)}$ of f_ℓ . Set $\delta = \frac{\varepsilon}{2k^{O(\log k)}}$; by the assumption of [Theorem 2.1](#), in time $T(\frac{\varepsilon}{2k^{O(\log k)}}, 2 \log k, k^2, n)$ it is possible to obtain additively $\pm\delta$ -accurate estimates of each of the $(2 \log k)$ -way marginals of f_ℓ . In the ℓ -th stage, algorithm A_{support} runs $A_{\text{one-stage}}$ using $S_{(\ell)}^*$ and these estimates of the marginals; by [Claim 2.3](#), this takes time $\text{poly}(n/\varepsilon, n^{\log k})$ and yields a function $\tilde{f}_\ell : [n]^\ell \rightarrow [0, 1]$ such that $\|f_\ell - \tilde{f}_\ell\|_1 \leq \frac{\delta}{2k^{O(\log k)}} \cdot k^{O(\log k)} = \varepsilon/4$. Since by assumption f is ε -heavy, it follows that any element x in the support of \tilde{f}_ℓ such that $\tilde{f}_\ell(x) \leq \varepsilon/4$ must not be in the support of f_ℓ ; so the algorithm removes all such elements x from $S_{(\ell)}^*$ to obtain the set $S'_{(\ell)}$. This resulting $S'_{(\ell)}$ is precisely the support of f_ℓ , and is clearly of size at most k . \square

Finally, the overall algorithm A_{learn} works by running A_{support} to get the set $S' = S'_{(n)}$ of size at most k which is the support of $f_n = f$, and then uses S' and the algorithm A_{marginal} from the assumptions of [Theorem 2.1](#)) to run algorithm $A_{\text{one-stage}}$ and obtain the required ε -accurate approximator g of f . This concludes the proof of [Theorem 2.1](#).

3 Computing limited way marginals from noisy samples

Recall that the noisy ranking learning problems we consider are of the following sort: There is a known noise distribution \mathcal{K} supported on \mathbb{S}_n , and an unknown k -sparse ε -heavy distribution $f : \mathbb{S}_n \rightarrow [0, 1]$. Each sample provided to the learning algorithm is generated by the following probabilistic process: independent draws of $\pi \sim \mathcal{K}$ and $\sigma \sim f$ are obtained, and the sample given to the learner is $(\pi\sigma) \in \mathbb{S}_n$. By the reduction established in [Theorem 2.1](#), in order to give an algorithm that learns the distribution f in the presence of a particular kind of noise \mathcal{K} , it suffices to give an algorithm that can efficiently estimate t -way marginals given samples $\pi\sigma \sim \mathcal{K} * f$.

The main result of this section, [Theorem 3.1](#), gives such an algorithm. Before stating the theorem we need some terminology and notation and we need to recall some necessary background from representation theory of the symmetric group (see [Appendix A](#) for a detailed overview of all of the required background).

First, let \mathcal{K} be a distribution over \mathbb{S}_n (which should be thought of as a noise distribution as described earlier). We say that \mathcal{K} is *efficiently samplable* if there is a $\text{poly}(n)$ -time randomized algorithm which takes no input and, each time it is invoked, returns an independent draw of $\pi \sim \mathcal{K}$.

Next, we recall that a *partition* λ of the natural number n (written “ $\lambda \vdash n$ ”) is a vector of natural numbers $(\lambda_1, \dots, \lambda_k)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k > 0$ and $\lambda_1 + \dots + \lambda_k = n$ (see [Appendix A.2](#) for more detail). For two partitions λ and μ of n , we say that μ *dominates* λ , written $\mu \triangleright \lambda$, if $\sum_{j \leq i} \mu_j \geq \sum_{j \leq i} \lambda_j$ for all $i > 0$ (see [Definition A.11](#)). Given any $\lambda \vdash n$, let $\text{Up}(\lambda)$ denote the set of all partitions $\mu \vdash n$ such that $\mu \triangleright \lambda$.

We recall that a *representation* of the symmetric group \mathbb{S}_n is a group homomorphism from \mathbb{S}_n to $\mathbb{C}^{m \times m}$ (see [Appendix A](#)). We further recall that for each partition $\lambda \vdash n$ there is a corresponding *irreducible* representation, denoted ρ_λ (see [Appendix A.2](#)). For a matrix M we write $\sigma_{\min}(M)$ to denote the smallest singular value of M . Given a partition $\lambda \vdash n$ we define the value $\sigma_{\min, \text{Up}(\lambda), \mathcal{K}}$ to be

$$\sigma_{\min, \text{Up}(\lambda), \mathcal{K}} := \min_{\mu \in \text{Up}(\lambda)} \sigma_{\min}(\widehat{\mathcal{K}}(\rho_\mu)), \quad (8)$$

the smallest singular value across all Fourier coefficients of the noise distribution of irreducible representations corresponding to partitions that dominate λ . (We recall that the Fourier coefficients of functions over the symmetric group, and indeed over any finite group, are matrices; see [Appendix A.2](#).)

Finally, for $0 \leq \ell \leq n - 1$ we define the partition $\lambda_{\text{hook}, \ell} \vdash n$ to be

$$\lambda_{\text{hook}, \ell} := (n - \ell, 1, \dots, 1).$$

Now we can state the main result of this section:

Theorem 3.1. *Let \mathcal{K} be an efficiently samplable distribution over \mathbb{S}_n . Let f be an unknown distribution over \mathbb{S}_n . There is an algorithm A_{marginal} with the following properties: A_{marginal} receives as input a parameter $\delta > 0$, a confidence parameter $\tau > 0$, a pair of ℓ -tuples $\bar{i} = (i_1, \dots, i_\ell) \in [n]^\ell$, $\bar{j} = (j_1, \dots, j_\ell) \in [n]^\ell$ each composed of ℓ distinct elements, and has access to random samples from $\mathcal{K} * f$. Algorithm A_{marginal} runs in time $\text{poly}\left(\binom{n}{\ell}, \delta^{-1}, \sigma_{\min, \text{Up}(\lambda_{\text{hook}, \ell}), \mathcal{K}}^{-1}, \log(1/\tau)\right)$ and outputs a value $\kappa_{\bar{i}, \bar{j}}$ which with probability at least $1 - \tau$ is a $\pm\delta$ -accurate estimate of the (\bar{i}, \bar{j}) -marginal of f .*

We will use the following claim to prove [Theorem 3.1](#):

Claim 3.2. *Let $\rho : \mathbb{S}_n \rightarrow \mathbb{C}^{m \times m}$ be any unitary representation of \mathbb{S}_n , let \mathcal{K} be any efficiently samplable distribution over \mathbb{S}_n , and let σ_{\min} denote the smallest singular value of $\widehat{\mathcal{K}}(\rho)$. Let f be an unknown distribution over \mathbb{S}_n . There is an algorithm which, given random samples from $\mathcal{K} * f$ and an error parameter $0 < \delta < 1$, runs in time $\text{poly}(m, n, \sigma_{\min}^{-1}, \delta^{-1})$ and with high probability outputs a matrix $M_{f, \rho}$ such that $\|M_{f, \rho} - \widehat{f}(\rho)\| \leq \delta$.*

Proof. Let $\eta_1, \eta_2 > 0$ denote two error parameters that will be fixed later. Since f is a distribution, the Fourier coefficient $\widehat{f}(\rho)$ is equal to $\mathbf{E}_{\sigma \sim f}[\rho(\sigma)]$. Consequently, since \mathcal{K} is assumed to be efficiently samplable and the algorithm is given samples from $\mathcal{K} * f$, by sampling from \mathcal{K} and from $\mathcal{K} * f$ it is straightforward to obtain matrices M_1, M_2 in time $\text{poly}(m, n, \log(1/\tau))$ which with probability $1 - \tau$ satisfy

$$\|M_1 - \widehat{\mathcal{K}}(\rho)\|_2 \leq \eta_1 \text{ and } \|M_2 - \widehat{\mathcal{K} * f}(\rho)\|_2 \leq \eta_2.$$

Now we recall the following matrix perturbation inequality (see [Theorem 2.2](#) of [\[Ste77\]](#)):

Lemma 3.3. *Let $A \in \mathbb{R}^{n \times n}$ be a non-singular matrix and further let $\Delta A \in \mathbb{R}^{n \times n}$ be such that $\|\Delta A\|_2 \cdot \|A^{-1}\|_2 < 1$. Then $A + \Delta A$ is non-singular. Further, if $\gamma = 1 - \|A^{-1}\|_2 \|\Delta A\|_2$, then*

$$\|A^{-1} - (A + \Delta A)^{-1}\|_2 \leq \frac{\|A^{-1}\|_2^2 \|\Delta A\|_2}{\gamma}.$$

Let us now set the error parameters η_1 and η_2 as follows (recall that $\delta < 1$):

$$\eta_1 = \min \left\{ \frac{\delta \cdot \sigma_{\min}^2}{4}, \frac{\delta \cdot \sigma_{\min}}{4} \right\} \quad \text{and} \quad \eta_2 = \min \left\{ \frac{\delta \cdot \sigma_{\min}}{4}, 1 \right\}. \quad (9)$$

Applying Lemma 3.3 with $\widehat{\mathcal{K}}(\rho)$ in place of A and $M_1 - \widehat{\mathcal{K}}(\rho)$ in place of ΔA , using (9) (more precisely, the upper bound $\eta_1 \leq \delta \cdot \sigma_{\min}^2/4$ in the numerator and the upper bound $\eta_1 \leq \delta \cdot \sigma_{\min}/4$ in the denominator) we get that

$$\|M_1^{-1} - \widehat{\mathcal{K}}(\rho)^{-1}\|_2 \leq \frac{\|\widehat{\mathcal{K}}(\rho)^{-1}\|_2^2 \cdot \|M_1 - \widehat{\mathcal{K}}(\rho)\|_2}{1 - \|\widehat{\mathcal{K}}(\rho)^{-1}\|_2 \cdot \|M_1 - \widehat{\mathcal{K}}(\rho)\|_2} \leq \frac{\delta}{3}. \quad (10)$$

Now using $\widehat{\mathcal{K}} * \widehat{f}(\rho) = \widehat{\mathcal{K}}(\rho) \cdot \widehat{f}(\rho)$, we get

$$\begin{aligned} \|M_1^{-1} \cdot M_2 - \widehat{f}(\rho)\|_2 &= \|M_1^{-1} \cdot M_2 - \widehat{\mathcal{K}}(\rho)^{-1} \cdot \widehat{\mathcal{K}} * \widehat{f}(\rho)\|_2 \\ &\leq \|M_1^{-1} \cdot M_2 - M_1^{-1} \cdot \widehat{\mathcal{K}} * \widehat{f}(\rho)\|_2 + \|M_1^{-1} \cdot \widehat{\mathcal{K}} * \widehat{f}(\rho) - \widehat{\mathcal{K}}(\rho)^{-1} \cdot \widehat{\mathcal{K}} * \widehat{f}(\rho)\|_2 \\ &\leq \|M_1^{-1}\|_2 \cdot \|M_2 - \widehat{\mathcal{K}} * \widehat{f}(\rho)\|_2 + \|M_1^{-1} - \widehat{\mathcal{K}}(\rho)^{-1}\|_2 \cdot \|\widehat{\mathcal{K}} * \widehat{f}(\rho)\|_2 \\ &\leq \|M_1^{-1}\|_2 \cdot \eta_2 + \|\widehat{\mathcal{K}} * \widehat{f}(\rho)\|_2 \cdot \frac{\delta}{3}. \quad (\text{using (10)}) \\ &\leq \eta_2 \left(\|\widehat{\mathcal{K}}(\rho)^{-1}\|_2 + \|M_1^{-1} - \widehat{\mathcal{K}}(\rho)^{-1}\|_2 \right) + \|\widehat{\mathcal{K}} * \widehat{f}(\rho)\|_2 \cdot \frac{\delta}{3}. \quad (\text{using (10)}) \\ &\leq \sigma_{\min}^{-1} \cdot \eta_2 + \frac{\delta}{3} \cdot \eta_2 + \|\widehat{\mathcal{K}} * \widehat{f}(\rho)\|_2 \cdot \frac{\delta}{3}. \end{aligned} \quad (11)$$

Next we use the following fact, which is an easy consequence of the triangle inequality and the assumption that ρ is unitary:

Fact 3.4. *Let $\rho : \mathbb{S}_n \rightarrow \mathbb{C}^{m \times m}$ be a unitary representation and let $g : \mathbb{S}_n \rightarrow \mathbb{R}^+$. Then we have that $\|\widehat{g}(\rho)\|_2 \leq \|g\|_1$.*

Combining this fact with (11) and (9), since $\|\mathcal{K} * f\|_1 = 1$, we get that

$$\|M_1^{-1} \cdot M_2 - \widehat{f}(\rho)\|_2 \leq \sigma_{\min}^{-1} \cdot \eta_2 + \frac{\delta}{3} \cdot \eta_2 + \frac{\delta}{3} \leq \frac{\delta}{4} + \frac{\delta}{3} + \frac{\delta}{3} < \delta.$$

This concludes the proof of Claim 3.2. \square

With Claim 3.2 in hand we are ready to prove Theorem 3.1:

Proof of Theorem 3.1. Let $\tau_{\lambda_{\text{hook},\ell}}$ be the permutation representation corresponding to the partition $\lambda_{\text{hook},\ell}$; for conciseness we subsequently write ρ for $\tau_{\lambda_{\text{hook},\ell}}$. Definition A.10 immediately gives that the dimension of ρ is $\binom{n}{\ell}$. Observe that ρ is a unitary representation. Let σ_{\min} denote the smallest singular value of $\widehat{\mathcal{K}}(\rho)$; applying Claim 3.2, we get an algorithm running in time $\text{poly}(\binom{n}{\ell}, \sigma_{\min}^{-1}, \delta)$ which outputs a matrix $M_{f,\rho}$ such that $\|M_{f,\rho} - \widehat{f}(\rho)\| \leq \delta$. Next, we observe that the Young tableaux corresponding to the partition $\lambda_{\text{hook},\ell}$ (which, recalling Definition A.10, index the rows and columns of $\rho(\cdot)$) correspond precisely to ordered t -tuples of distinct entries of $[n]$. If $Y_{\lambda_{\text{hook},\ell},i} = \bar{i}$ and $Y_{\lambda_{\text{hook},\ell},j} = \bar{j}$, then it follows that

$$\widehat{f}(\rho)(i, j) = \sum_{\sigma \in \mathbb{S}_n} f(\sigma) \cdot \mathbf{1}[f(i_1) = j_1 \text{ and } \dots \text{ and } f(i_\ell) = j_\ell],$$

which is the (\bar{i}, \bar{j}) -marginal of f as desired; so the output of the algorithm is $M_{f,\rho}(\bar{i}, \bar{j})$.

To finish the correctness argument it remains only to argue that σ_{\min}^{-1} is at most $\text{poly}(\sigma_{\min, \text{Up}(\lambda_{\text{hook}, \ell)}}^{-1})$. To see that this is indeed the case, we observe that by [Theorem A.12](#), the permutation representation $\tau_{\lambda_{\text{hook}, \ell}}$ block diagonalizes into a direct sum of irreducible representations ρ_{μ} where each μ belongs to $\text{Up}(\lambda_{\text{hook}, \ell})$. This finishes the proof of [Theorem 3.1](#). \square

3.1 Efficient samplability of our noise distributions

In order to apply [Theorem 3.1](#) to a particular noise distribution \mathcal{K} we need to confirm that \mathcal{K} is efficiently samplable; we now do this for each of the three noise models that we consider. It is immediate from the definition that it is straightforward (given \bar{p}) to efficiently generate a random σ drawn from the symmetric noise distribution $\mathcal{S}_{\bar{p}}$, and the same is true for the heat kernel noise distribution \mathcal{H}_t .

For the generalized Mallows model \mathcal{M}_{θ} , the characterization $\Pr_{\sigma \sim \mathcal{M}_{\theta}}[\sigma = \pi] = e^{-\theta d(\pi, e)}/Z(\theta)$ given earlier does not directly yield an efficient sampling algorithm, since it may be hard to compute or approximate the normalizing factor $Z(\theta) = \sum_{\pi \in \mathbb{S}_n} e^{-\theta d(\pi, e)}$. Instead, we recall (see e.g. Section 2.1 of [\[DS98\]](#)) that the Metropolis algorithm can be used to efficiently perform a random walk on \mathbb{S}_n whose unique stationary distribution is the generalized Mallows distribution \mathcal{M}_{θ} . (Each step of the random walk can be carried out efficiently because it is computationally easy to compute the Cayley distance between two permutations: if π is the permutation that brings σ to τ , then the Cayley distance $d(\sigma, \tau)$ is $n - \text{cycles}(\pi)$ where $\text{cycles}(\pi)$ is the number of cycles in π .) It is known (see e.g. Theorem 2 of [\[DH92\]](#)) that this random walk has rapid convergence, and consequently it is indeed possible to sample efficiently from \mathcal{M}_{θ} (up to an exponentially small statistical distance which can be ignored in our applications since our algorithms use a sub-exponential number of samples).

4 Representations of symmetric noise

In this section we establish lower bounds on the smallest singular value for the relevant matrices corresponding to “symmetric noise” $\mathcal{S}_{\bar{p}}$ on \mathbb{S}_n . In more detail, the main result of this section is the following lower bound:

Lemma 4.1. *Let $\ell \in \{1, \dots, n\}$ and let $\bar{p} = (p_0, \dots, p_n) \in \Delta^n$ (i.e. \bar{p} is a non-negative vector whose entries sum to 1) which is such that*

$$\sum_{j=0}^{n-\ell} p_j \geq \kappa.$$

Then (recalling [Equation \(8\)](#)) we have that

$$\sigma_{\min, \text{Up}(\lambda_{\text{hook}, \ell}), \mathcal{S}_{\bar{p}}} \geq \frac{\kappa}{n^{\ell}}. \tag{12}$$

4.1 Setup

To analyze the smallest singular value of $\widehat{\mathcal{S}}_{\bar{p}}(\rho_\mu)$ (as required by the definition of $\sigma_{\min, \mathbf{Up}(\lambda_{\text{hook}}, \ell), \mathcal{S}_{\bar{p}}}$), we start by observing that symmetric noise is a *class function* (meaning that it is invariant under conjugation, see [Definition A.6](#)):

Claim 4.2. *For any vector $\bar{p} = (p_0, \dots, p_n) \in \Delta^n$, the distribution $\mathcal{S}_{\bar{p}}$ (viewed as a function from \mathbb{S}_n to $[0, 1]$) is a class function (i.e. $\mathcal{S}_{\bar{p}}(\pi) = \mathcal{S}_{\bar{p}}(\tau\pi\tau^{-1})$ for every $\pi, \tau \in \mathbb{S}_n$).*

Proof. For $0 \leq j \leq n$, let \bar{e}_j denote the vector in \mathbb{R}^{n+1} which has a 1 in the j -th position and a 0 in every other position. By linearity, to prove [Claim 4.2](#) it suffices to prove that $\mathcal{S}_{\bar{e}_j}$ is invariant under conjugation for every j ; to establish this, it suffices to show that $\mathcal{S}_{\bar{e}_j}$ is invariant under conjugation by any transposition τ . By symmetry, it suffices to consider the transposition $\tau = (1, 2)$.

We observe that $\mathcal{S}_{\bar{e}_j}$ is a uniform average of \mathbb{U}_A over all $\binom{[n]}{j}$ subsets A of $[n]$ of size exactly j . Now we consider two cases: the first is that $|A \cap \{1, 2\}|$ is 0 or 2. In this case it is easy to see that \mathbb{U}_A does not change under conjugation by the transposition $(1, 2)$. The remaining case is that $|A \cap \{1, 2\}| = 1$; in this case it is easy to see that conjugation by $(1, 2)$ converts \mathbb{U}_A into $\mathbb{U}_{A \Delta \{1, 2\}}$. Since the collection of size- j sets A with $A \cap \{1, 2\} = \{1\}$ are in 1-1 correspondence with the collection of size- j sets A with $A \cap \{1, 2\} = \{2\}$, it follows that $\mathcal{S}_{\bar{e}_j}$ is invariant under conjugation by $\tau = (1, 2)$, and the proof is complete. \square

Before stating the next lemma we remind the reader that for partitions $\mu \vdash m, \lambda \vdash n$ where $m \leq n$, we write $\text{Paths}(\mu, \lambda)$ to denote the number of paths from μ to λ in Young's lattice (see [Appendix A.2](#) and [Theorem A.15](#)). We write Triv_j to denote the trivial partition (j) of j .

Lemma 4.3. *Let $\lambda \vdash n$ and let ρ_λ be the corresponding irreducible representation of \mathbb{S}_n . Given $\bar{p} = (p_0, \dots, p_n) \in \Delta^n$, we have that*

$$\widehat{\mathcal{S}}_{\bar{p}}(\rho_\lambda) = c(\bar{p}, \lambda) \cdot \text{Id} \quad \text{where} \quad c(\bar{p}, \lambda) := \frac{\sum_{j=0}^n p_j \cdot \text{Paths}(\text{Triv}_j, \lambda)}{\dim(\rho_\lambda)}. \quad (13)$$

Proof. By [Claim 4.2](#), we have that $\mathcal{S}_{\bar{p}}$ is a class function, so we may apply [Lemma A.9](#) to conclude that

$$\widehat{\mathcal{S}}_{\bar{p}}(\rho_\lambda) = c(\bar{p}, \lambda) \cdot \text{Id},$$

where

$$c(\bar{p}, \lambda) = \frac{1}{\dim(\rho_\lambda)} \cdot \left(\sum_{\sigma \in \mathbb{S}_n} \mathcal{S}_{\bar{p}}(\sigma) \cdot \chi_\lambda(\sigma) \right)$$

and χ_λ denotes the character of the irreducible representation ρ_λ . Thus it remains to show that $\sum_{\sigma \in \mathbb{S}_n} \mathcal{S}_{\bar{p}}(\sigma) \cdot \chi_\lambda(\sigma)$ is equal to the numerator of [Equation \(13\)](#). By definition of $\mathcal{S}_{\bar{p}}$, we have that

$$\sum_{\sigma \in \mathbb{S}_n} \mathcal{S}_{\bar{p}}(\sigma) \cdot \chi_\lambda(\sigma) = \sum_{0 \leq j \leq n} \bar{p}_j \mathbf{E}_{\mathcal{A}: |\mathcal{A}|=j} \mathbf{E}_{\sigma \in \mathbb{U}_{\mathcal{A}}} \chi_\lambda(\sigma). \quad (14)$$

We proceed to analyze $\mathbf{E}_{\sigma \in \mathbb{U}_{\mathcal{A}}} \chi_\lambda(\sigma)$. Let $\rho_\lambda^{\mathcal{A}}$ denote the representation ρ_λ restricted to the subgroup $\mathbb{S}_{\mathcal{A}}$. By [Theorem A.15](#), the representation $\rho_\lambda^{\mathcal{A}}$ splits as follows:

$$\rho_\lambda^{\mathcal{A}} = \bigoplus_{\mu \vdash |\mathcal{A}|} \text{Paths}(\mu, \lambda) \rho_\mu.$$

Thus, we have that

$$\mathbf{E}_{\sigma \in \mathbb{U}_{\mathcal{A}}} \chi_{\lambda}(\sigma) = \sum_{\mu \vdash |\mathcal{A}|} \text{Paths}(\mu, \lambda) \mathbf{E}_{\sigma \in \mathbb{U}_{\mathcal{A}}} \chi_{\mu}(\sigma) = \text{Paths}(\text{Triv}_{|\mathcal{A}|}, \lambda).$$

The second equality follows from that fact that if μ is a non-trivial partition of $|\mathcal{A}|$ then $\mathbf{E}_{\sigma \in \mathbb{U}_{\mathcal{A}}} \chi_{\mu}(\sigma) = 0$, while if $\mu = \text{Triv}_{|\mathcal{A}|}$ then $\mathbf{E}_{\sigma \in \mathbb{U}_{\mathcal{A}}} \chi_{\mu}(\sigma) = 1$. Plugging this into (14) we get that $\sum_{\sigma \in \mathbb{S}_n} \mathcal{S}_{\bar{p}}(\sigma) \cdot \chi_{\lambda}(\sigma) = \sum_{j=0}^n p_j \cdot \text{Paths}(\text{Triv}_j, \lambda)$, and the lemma is proved. \square

4.2 Proof of Lemma 4.1

We recall from Equation (8) that

$$\sigma_{\min, \text{Up}(\lambda_{\text{hook}, \ell}), \mathcal{S}_{\bar{p}}} := \min_{\mu \in \text{Up}(\lambda_{\text{hook}, \ell})} \sigma_{\min}(\widehat{\mathcal{S}}_{\bar{p}}(\rho_{\mu})).$$

Fix any $\mu \in \text{Up}(\lambda_{\text{hook}, \ell})$, so μ is a partition of n of the form $(n - \ell', \ell_2, \dots, \ell_r)$ where $\ell' \leq \ell$. By Lemma 4.3 we have that the smallest singular value of $\widehat{\mathcal{S}}_{\bar{p}}(\rho_{\mu})$ is

$$c(\bar{p}, \mu) := \frac{\sum_{j=0}^n p_j \cdot \text{Paths}(\text{Triv}_j, \mu)}{\dim(\rho_{\mu})}. \quad (15)$$

To upper bound $\dim(\rho_{\mu})$, we observe that

$$\dim(\rho_{\mu}) \leq \dim(\tau_{\mu}) = \binom{n}{n - \ell', \ell_2, \dots, \ell_r} \leq \frac{n!}{(n - \ell')!} \leq n^{\ell'} \leq n^{\ell},$$

where the first inequality is by Theorem A.12. For the numerator, we observe that if $j \leq n - \ell$ then there is at least one path in the Young lattice from Triv_j to μ , so under the assumptions of Lemma 4.1 the numerator of Equation (15) is at least κ . This proves the lemma. \square

5 Representations of heat kernel noise

In this section, analogous to Section 4, we lower bound Equation (8) when the noise distribution \mathcal{K} is \mathcal{H}_t , corresponding to “heat kernel noise” at temperature parameter t :

Lemma 5.1. *Let $t \geq 1$ and let $\ell \in \{1, \dots, cn\}$ for some suitably small universal constant $c > 0$. Then we have that*

$$\sigma_{\min, \text{Up}(\lambda_{\text{hook}, \ell}), \mathcal{H}_t} \geq \frac{1}{2} \cdot e^{-O(\ell t)/n}. \quad (16)$$

5.1 Setup

Let $\text{trans} : \mathbb{S}_n \rightarrow [0, 1]$ be the following probability distribution over \mathbb{S}_n :

$$\text{trans}(\pi) = \begin{cases} 1/n & \text{if } \pi \text{ is the identity,} \\ 2/n^2 & \text{if } \pi \text{ is a transposition,} \\ 0 & \text{otherwise.} \end{cases}$$

Since $\text{trans}(\pi)$ depends only on the cycle structure of π , the function $\text{trans}(\cdot)$ is a class function. Fix any $\mu \in \text{Up}(\lambda_{\text{hook}, \ell})$, so μ is a partition of n of the form (μ_1, \dots, μ_r) where $\mu_1 \geq n - \ell$. As in the proof of [Lemma 4.3](#) we may apply [Lemma A.9](#) to conclude that

$$\widehat{\text{trans}}(\rho_\mu) = c_{\text{trans}, \mu} \cdot \text{Id}$$

for some constant $c_{\text{trans}, \mu}$. By Corollary 1 of Diaconis and Shahshahani [[DS81](#)], we have that

$$c_{\text{trans}, \mu} = \frac{1}{n} + \frac{n-1}{n} \cdot \frac{\chi_\mu(\tau)}{\dim(\rho_\mu)}, \quad (17)$$

where as before χ_μ denotes the character of the irreducible representation ρ_μ and τ is any transposition. [[DS81](#)] further shows that for ρ_μ an irreducible representation of \mathbb{S}_n with μ as above and τ any transposition, it holds that

$$\frac{\chi_\mu(\tau)}{\dim(\rho_\mu)} = \frac{1}{n(n-1)} \cdot \sum_{j=1}^r (\mu_j - j)(\mu_j - j + 1) - j(j-1). \quad (18)$$

In our setting we have

$$(18) \geq \frac{(n-\ell)(n-\ell-1)}{n(n-1)} + \frac{1}{n(n-1)} \sum_{j=2}^r (\mu_j - j)(\mu_j - j + 1) - j(j-1). \quad (19)$$

where the inequality holds because $\mu_1 \geq n - \ell$. Now, we observe that for each summand in [Equation \(19\)](#), we have

$$\begin{aligned} (\mu_j - j)(\mu_j - j + 1) - j(j-1) &= \mu_j^2 - \mu_j(2j-1) \\ &\geq -\mu_j(2j-1) \\ &\geq \frac{-\ell}{j-1} \cdot (2j-1) \geq -3\ell. \end{aligned}$$

The second inequality above holds because $\mu_2 + \dots + \mu_j \leq \ell$ and the μ_j 's are non-increasing, so $\mu_j \leq \frac{\ell}{j-1}$. Since $r-1 \leq \ell$, this means that

$$(18) \geq \frac{(n-\ell)(n-\ell-1)}{n(n-1)} - \frac{3\ell^2}{n(n-1)} \geq 1 - \frac{O(\ell)}{n},$$

and recalling [Equation \(17\)](#) we get that

$$1 \geq c_{\text{trans}, \mu} \geq 1 - \frac{O(\ell)}{n}. \quad (20)$$

5.2 Proof of [Lemma 5.1](#)

As in [Section 4](#) we recall from [Equation \(8\)](#) that

$$\sigma_{\min, \text{Up}(\lambda_{\text{hook}, \ell}), \mathcal{H}_t} := \min_{\mu \in \text{Up}(\lambda_{\text{hook}, \ell})} \sigma_{\min}(\widehat{\mathcal{H}}_t(\rho_\mu)),$$

Fix any $\mu \in \text{Up}(\lambda_{\text{hook},\ell})$ (so μ is a partition of n of the form (μ_1, \dots, μ_r) where $\mu_1 \geq n - \ell$). We recall that the function $\mathcal{H}_t : \mathbb{S}_n \rightarrow [0, 1]$ is defined by

$$\mathcal{H}_t = \sum_{j=0}^{\infty} \Pr_{\mathbf{T} \sim \text{Poi}(t)}[\mathbf{T} = j] (\text{trans})^j,$$

where “ $(\text{trans})^T$ ” denotes T -fold convolution of trans . Since convolution corresponds to multiplication of Fourier coefficients, this gives that

$$\widehat{\mathcal{H}}_t(\rho_\mu) = \mathbf{c}(t, \mu) \cdot \text{Id}, \text{ where } \mathbf{c}(t, \mu) := \sum_{j=0}^{\infty} \Pr_{\mathbf{T} \sim \text{Poi}(t)}[\mathbf{T} = j] (\mathbf{c}_{\text{trans}, \mu})^j. \quad (21)$$

Recalling [Cho94] that the median of the Poisson distribution $\text{Poi}(t)$ is at most $t + 1/3$, we get that

$$\mathbf{c}(t, \mu) \geq \frac{1}{2} \cdot (\mathbf{c}_{\text{trans}, \mu})^{t+1/3} \geq \frac{1}{2} \cdot e^{-O(\ell t)/n},$$

(where the second inequality uses $\ell \leq cn$ and $t \geq 1$), and the lemma is proved. \square

6 Representations of Cayley-Mallows model noise

In this section we lower bound Equation (8) when the noise distribution \mathcal{K} is \mathcal{M}_θ , corresponding to the Cayley-Mallows noise model with parameter θ :

Lemma 6.1. *Let $\theta > 0$, let $\ell \in \{1, \dots, n\}$, and let $\eta := \text{dist}(\theta, \ell) = \min_{j \in \{1, \dots, \ell\}} |e^\theta - j|$. Then (recalling Equation (8)) we have that*

$$\sigma_{\min, \text{Up}(\mu_{\text{hook}, \ell}), \mathcal{M}_\theta} \geq (2n)^{-\ell} \eta^{2\sqrt{\ell}}. \quad (22)$$

Similar to the previous two sections, Lemma 6.1 follows immediately from the following lower bound on singular values of certain irreducible representations:

Lemma 6.2. *Let μ be a partition of n of the form (μ_1, \dots, μ_r) where $\mu_1 \geq n - \ell$. Let $\theta > 0$ and let $\eta := \text{dist}(\theta, \ell) = \min_{j \in \{1, \dots, \ell\}} |e^\theta - j|$. Then we have that*

$$\widehat{\mathcal{M}}_\theta(\rho_\mu) = c_{\mu, \theta} \cdot \text{Id} \text{ where } |c_{\mu, \theta}| \geq (2n)^{-\ell} \eta^{2\sqrt{\ell}}.$$

To prove Lemma 6.2, we will need the notions of *content* and *hook length* for boxes in a Young diagram:

Definition 6.3. Let μ be a partition $\mu \vdash n$. The *hook length* of a box u in the Young diagram for μ , denoted by $h(u)$, is the sum

$$(\# \text{ of boxes to the right of } u \text{ in its row}) + (\# \text{ of boxes below } u \text{ in its column}) + 1 \text{ (for } u \text{ itself)}.$$

The *content* $c(u)$ of a box u is $c(u) := j - i$, where j is its column number (from the left, starting with column 1) and i is its row number (from the top, starting with row 1).



Figure 1: On the left is a Young diagram in which each box has been labeled with its hook length; on the right is a Young diagram in which each box has been labeled with its content.

The left portion of **Figure 1** depicts a Young diagram annotated with the hook lengths of each of its boxes. The right portion of **Figure 1** depicts the same Young diagram annotated with the contents of each of its boxes.

We will need the following technical result to prove Lemma 6.2:

Lemma 6.4. *Let $\mu \vdash n$ and let χ_μ be the corresponding character in \mathbb{S}_n . For any $q \in \mathbb{R}$,*

$$\frac{1}{n!} \sum_{\sigma \in \mathbb{S}_n} \chi_\mu(\sigma) \cdot q^{\text{cycles}(\sigma)} = \prod_{u \in \mu} \frac{q + c(u)}{h(u)},$$

where the subscript “ $u \in \mu$ ” means that u ranges over all the boxes in the Young diagram corresponding to μ .

Proof. The above identity is given as Exercise 7.50 in Stanley’s book [Sta99]. For the sake of completeness, we provide the proof here.

For any $\bar{t} = (t_1, \dots, t_n)$, we define the polynomial

$$a_{\bar{t}}(x_1, \dots, x_n) := \det \begin{bmatrix} x_1^{t_1} & x_2^{t_1} & x_3^{t_1} & \dots & x_n^{t_1} \\ x_1^{t_2} & x_2^{t_2} & x_3^{t_2} & \dots & x_n^{t_2} \\ \dots & \dots & \dots & \dots & \dots \\ x_1^{t_n} & x_2^{t_n} & x_3^{t_n} & \dots & x_n^{t_n} \end{bmatrix}.$$

Given any partition $\mu \vdash n$, we now define the Schur polynomial $s_\mu(x_1, \dots, x_n)$ as follows: Define $\bar{t}_\mu = (\mu_1 + n - 1, \dots, \mu_n + 0)$ and $\bar{t}_0 = (n - 1, \dots, 0)$. Then,

$$s_\mu(x_1, \dots, x_n) := \frac{a_{\bar{t}_\mu}(x_1, \dots, x_n)}{a_{\bar{t}_0}(x_1, \dots, x_n)}.$$

The denominator is just the Vandermonde determinant of the variables (x_1, \dots, x_n) . As the polynomial $a_{\bar{t}_\mu}(x_1, \dots, x_n)$ is alternating, it follows that $s_\mu(x_1, \dots, x_n)$ is a polynomial (as opposed to a rational function) and further, it is symmetric.

The following is a fundamental fact connecting Schur polynomials and cycles: For any $0 \leq k \leq n$,

$$s_\mu(\underbrace{1, \dots, 1}_k, \underbrace{0, \dots, 0}_{n-k}) = \sum_{\sigma \in \mathbb{S}_n} \frac{1}{n!} \cdot \chi_\mu(\sigma) \cdot k^{\text{cycles}(\sigma)} \quad (23)$$

(see equation 7.78 in [Sta99]). On the other hand, there are known explicit formulas for evaluations of the Schur polynomial at specific inputs. In particular, Corollary 7.21.4 of [Sta99] states that

$$s_\mu(\underbrace{1, \dots, 1}_k, \underbrace{0, \dots, 0}_{n-k}) = \prod_{u \in \mu} \frac{k + c(u)}{h(u)}. \quad (24)$$

Combining (23) and (24), we get that for any $0 \leq k \leq n$, we have

$$\frac{1}{n!} \sum_{\sigma \in \mathbb{S}_n} \chi_\mu(\sigma) \cdot k^{\text{cycles}(\sigma)} = \prod_{u \in \mu} \frac{k + c(u)}{h(u)}.$$

However, note that both the left and the right hand sides can be seen as polynomials of degree at most n in the variable k . Since they agree at $n + 1$ values $k = 0, \dots, n$, they must be identical as formal functions. This concludes the proof. \square

Proof of Lemma 6.2. Recall that the distribution \mathcal{M}_θ over \mathbb{S}_n is defined by $\mathcal{M}_\theta(\pi) = e^{-\theta d(\pi, e)} / Z(\theta)$, where $Z(\theta) = \sum_{\pi \in \mathbb{S}_n} e^{-\theta d(\pi, e)}$ is a normalizing constant. Since the Cayley distance $d(\sigma, \tau)$ is equal to $n - \text{cycles}(\sigma^{-1}\tau)$, where $\text{cycles}(\pi)$ is the number of cycles in π , we have that

$$\mathcal{M}_\theta(\pi) = \frac{e^{\theta \cdot \text{cycles}(\pi)}}{C}, \quad \text{where } C = \sum_{\pi \in \mathbb{S}_n} e^{\theta \cdot \text{cycles}(\pi)}.$$

Since the $\text{cycles}(\cdot)$ function is a class function so is \mathcal{M}_θ , so we can apply Lemma A.9 and we get that $\widehat{\mathcal{M}}_\theta(\rho_\mu) = c_{\mu, \theta} \cdot \text{Id}$, where

$$c_{\mu, \theta} = \frac{\sum_{\sigma \in \mathbb{S}_n} \mathcal{M}_\theta(\sigma) \cdot \chi_\mu(\sigma)}{\dim(\rho_\mu)} = \frac{\sum_{\sigma \in \mathbb{S}_n} e^{\theta \cdot \text{cycles}(\sigma)} \cdot \chi_\mu(\sigma)}{\dim(\rho_\mu) \cdot (\sum_{\sigma \in \mathbb{S}_n} e^{\theta \cdot \text{cycles}(\sigma)})} = \frac{\sum_{\sigma \in \mathbb{S}_n} q^{\text{cycles}(\sigma)} \cdot \chi_\mu(\sigma)}{\dim(\rho_\mu) \cdot (\sum_{\sigma \in \mathbb{S}_n} q^{\text{cycles}(\sigma)})}, \quad \text{where } q := e^\theta.$$

We re-express the numerator by applying Lemma 6.4 to get

$$\sum_{\sigma \in \mathbb{S}_n} q^{\text{cycles}(\sigma)} \cdot \chi_\mu(\sigma) = n! \cdot \prod_{u \in \mu} \frac{q + c(u)}{h(u)}. \quad (25)$$

To analyze the denominator of $c_{\mu, \theta}$, applying Lemma 6.4 to the trivial partition $\text{Triv}_n = (n)$ of n (the character of which is identically 1), we get that

$$\sum_{\sigma \in \mathbb{S}_n} q^{\text{cycles}(\sigma)} = n! \cdot \prod_{u \in \text{Triv}_n} \frac{q + c(u)}{h(u)} = q(q + 1) \cdots (q + n - 1). \quad (26)$$

For the rest of the denominator, we recall the following well-known fact about the dimension of irreducible representations of the symmetric group:

Fact 6.5 (Hook length formula, see e.g. Theorem 3.41 of [Mél17]). *For $\mu \vdash n$, $\dim(\rho_\mu) = \frac{n!}{\prod_{u \in \mu} h(u)}$.*

Combining (25), (26) and Fact 6.5, we get

$$c_{\mu, \theta} = \frac{\prod_{u \in \mu} (q + c(u))}{q(q + 1) \cdots (q + n - 1)}. \quad (27)$$

Let \mathcal{A} denote the set consisting of the cells of the Young diagram of μ which are not in the first row. Since $n - \mu_1 = \ell'$ for some $\ell' \leq \ell$, the above expression simplifies to

$$c_{\mu,\theta} = \frac{\prod_{u \in \mathcal{A}} (q + c(u))}{(q + n - \ell') \cdots (q + n - 1)}. \quad (28)$$

To bound this ratio, first observe that both the numerator and denominator are ℓ' -way products. There are two possibilities now:

1. **Case 1:** $q \geq \ell + 1$. In this case we observe that each cell $u \in \mathcal{A}$ satisfies $c(u) \geq -\ell' \geq -\ell$. Thus $c_{\mu,\theta}$ can be expressed as a product of ℓ' many fractions, each of which is at least $\frac{q-\ell}{q+n-1} \geq \frac{1}{\ell+n}$. This implies that

$$c_{\mu,\theta} \geq \left(\frac{1}{n + \ell} \right)^{\ell'} \geq (2n)^{-\ell}.$$

2. **Case 2:** $q \leq \ell$. In this case, the denominator of [Equation \(28\)](#) is at most $(2n)^\ell$. To lower bound the numerator, observe that for every cell u of \mathcal{A} , the value of $c(u)$ is an integer in $\{-\ell, \dots, \ell\}$. Let j_0 and j_1 denote the two values in $\{-\ell, \dots, \ell\}$ for which $|q - j|$ achieves its smallest value η and its next smallest value (note that these values are equal if $\eta = 1/2$). Next, we observe that at most $\sqrt{\ell}$ many cells of \mathcal{A} have content equal to any given fixed integer value. Since j_0 and j_1 are the only possible values of $j \in \{-\ell, \dots, \ell\}$ for which $|q + j| < 1$, it follows that

$$\prod_{u \in \mathcal{A}} |(q + c(u))| \geq \left(\prod_{u \in \mathcal{A}: c(u)=j_0} |(q + c(u))| \right) \cdot \left(\prod_{u \in \mathcal{A}: c(u)=j_1} |(q + c(u))| \right) \geq \eta^{2\sqrt{\ell}}.$$

This finishes the proof. □

7 Our positive results for noisy rankings: Putting the pieces together

In this brief section we put all the pieces together to obtain our main positive results, [Theorems 1.2, 1.3](#) and [1.4](#), for the symmetric, heat kernel, and generalized Mallows noise models respectively.

Symmetric noise. Under the assumptions of [Theorem 1.2](#) (that $\sum_{j=0}^{n-\log k} p_j \geq \frac{1}{n^{O(\log k)}}$), taking $\ell = \log k$ in [Lemma 4.1](#), we have that $\sigma_{\min, \text{Up}(\lambda_{\text{hook}, \log k}), \mathcal{S}_{\bar{p}}} \geq \frac{1}{n^{O(\log k)}}$. Since (as discussed in [Section 3.1](#)) $\mathcal{S}_{\bar{p}}$ is efficiently samplable given \bar{p} , by [Theorem 3.1](#) in time $\text{poly}(n^{\log k}, 1/\delta, \log(1/\tau))$ with probability $1 - \tau$ it is possible to obtain $\pm\delta$ -accurate estimates of all of the $(\log k)$ -way marginals of f . Setting $\delta = \frac{\epsilon}{2k^{O(\log k)}}$ and applying [Theorem 2.1](#), we get [Theorem 1.2](#).

Heat kernel noise. First observe that we may assume that the temperature parameter t is at least 1 (since otherwise it is easy to artificially add noise to achieve $t = 1$). Under the assumptions of [Theorem 1.3](#) (that $t = O(n \log n)$), taking $\ell = \log k$ in [Lemma 5.1](#), we have that

$\sigma_{\min, \text{Up}(\lambda_{\text{hook}, \log k}), \mathcal{H}_t} \geq \frac{1}{n^{O(\log k)}}$. **Theorem 1.3** follows as in the previous paragraph (this time using the efficient samplability of \mathcal{H}_t given t).

Cayley-Mallows noise. Under the assumptions of **Theorem 1.4**, taking $\ell = \log k$ in **Lemma 6.1** we get that $\sigma_{\min, \text{Up}(\lambda_{\text{hook}, \log k}), \mathcal{M}_\theta} \geq \frac{1}{n^{O(\log k)}} \cdot \text{dist}(\theta, \log k)^{2\sqrt{\log k}}$. **Theorem 1.4** follows as in the previous paragraph (this time using the efficient samplability of \mathcal{M}_θ given θ).

8 Lower bound for Cayley-Mallows models

Recall that because of the $\text{poly}(\text{dist}(\theta, \log k)^{-\sqrt{\log k}})$ dependence in **Theorem 1.4**, the algorithm of that theorem is inefficient if e^θ is very close to an integer. In this section we prove **Theorem 1.5**, which establishes that *any algorithm* for learning in the presence of Cayley-Mallows noise *must* be inefficient if e^θ is very close to an integer.

8.1 A key technical result

The following lemma is at the heart of our lower bound. It shows that if e^θ is close to an integer, then any partition μ of $n \geq m$ which extends a particular partition λ_{sq} of m must be such that the Fourier coefficient $\widehat{\mathcal{M}_\theta}(\rho_\mu)$ of Cayley-Mallows noise has small singular values.

Lemma 8.1. *Let λ_{sq} denote the partition (t, \dots, t) of $m = t(t + j)$ whose Young tableau is a rectangle with $t + j$ rows and t columns. Let $\theta > 0$ be such that $|e^\theta - j| \leq \eta$ where $\eta \leq 1/2$. Let $n \geq m$, $\mu \vdash n$ and $\lambda_{\text{sq}} \uparrow \mu$ (recall **Definition A.13**). Then*

$$\widehat{\mathcal{M}_\theta}(\rho_\mu) = c_{\mu, \theta} \cdot \text{Id}, \quad \text{where } c_{\mu, \theta} \leq \eta^t.$$

Here ρ_μ denotes the irreducible representation of \mathbb{S}_n corresponding to the partition μ .

Proof. Let $\mu = (\mu_1, \dots, \mu_r)$. By **Lemma 6.2**, we have that

$$\widehat{\mathcal{M}_\theta}(\rho_\mu) = c_{\mu, \theta} \cdot \text{Id},$$

where **Equation (28)** gives the precise value of $c_{\mu, \theta}$ as

$$c_{\mu, \theta} = \frac{\prod_{u \in \mathcal{A}} (q + c(u))}{\prod_{u \in \mathcal{B}} (q + c(u))}, \quad \text{where } q = e^\theta. \quad (29)$$

Here \mathcal{A} denotes the set of cells of the Young diagram of μ which are not in the first row and \mathcal{B} denotes the rightmost $n - \mu_1$ many cells in the Young diagram of the trivial partition $\text{Triv}_n = (n)$. Note that in this lemma, we are trying to upper bound **Equation (29)** whereas **Lemma 6.2** was about lower bounding this quantity.

To upper bound **Equation (29)**, we first observe that there is an obvious bijection $\Phi : \mathcal{A} \rightarrow \mathcal{B}$ such that if $\Phi(u) = v$, then $c(v) > |c(u)| > 0$.

Next, let $\mathcal{A}_{-j} \subset \mathcal{A}$ be $\mathcal{A} := \{(r, s) : s - r = j \text{ and } (r, s) \in \mathcal{A}\}$. Since $\lambda_{\text{sq}} \uparrow \mu$, it follows that $|\mathcal{A}_{-j}| \geq t$. As a result, we can upper bound $c_{\mu, \theta}$ as follows:

$$\begin{aligned} c_{\mu, \theta} &= \frac{\prod_{u \in \mathcal{A}} (q + c(u))}{\prod_{u \in \mathcal{B}} (q + c(u))} = \prod_{u \in \mathcal{A}} \frac{q + c(u)}{q + c(\Phi(u))} = \left(\prod_{u \in \mathcal{A}_{-j}} \frac{q + c(u)}{q + c(\Phi(u))} \right) \left(\prod_{u \in \mathcal{A} \setminus \mathcal{A}_{-j}} \frac{q + c(u)}{q + c(\Phi(u))} \right) \\ &\leq \prod_{u \in \mathcal{A}_{-j}} q + c(u) \quad (\text{using } c(\Phi(u)) > |c(u)| > 0 \text{ and } q > 0) \\ &\leq \eta^t. \end{aligned}$$

□

8.2 Proof of Theorem 1.5

Theorem 1.5 is an immediate consequence of the following result. It shows that if e^θ is close to an integer j , then it may be statistically impossible to learn a distribution f supported on k rankings without using many samples from $\mathcal{M}_\theta * f$:

Theorem 8.2. *Given $j \in \mathbb{N}$, there are infinitely many values of k and $m = m(k) \approx \frac{\log k}{\log \log k}$ such that the following holds: there are two distributions f_1, f_2 over \mathbb{S}_m with the following properties:*

1. $d_{\text{TV}}(f_1, f_2) = 1$ (i.e. the distributions f_1 and f_2 have disjoint support);
2. $|\text{supp}(f_1)|, |\text{supp}(f_2)| \leq k$;
3. For any $\theta > 0$ such that $|e^\theta - j| \leq \eta \leq 1/2$, we have that $d_{\text{TV}}(\mathcal{M}_\theta * f_1, \mathcal{M}_\theta * f_2) \leq 2 \cdot \eta^\Theta \left(\sqrt{\frac{\log k}{\log \log k}} \right)$.

Proof. Let $t \geq j$ be any integer, let $m = t(t + j)$, and let $k = m!$. We first construct the two distributions f_1, f_2 over \mathbb{S}_m and argue that properties (1) and (2) hold.

Let $\lambda_{\text{sq}} \vdash m$ be the partition whose Young tableau is a rectangle with $t + j$ rows and t columns. Let us consider the character $\chi_{\text{sq}} : \mathbb{S}_m \rightarrow \mathbb{Q}$ corresponding to the partition λ_{sq} . By **Fact A.16** we have that χ_{sq} is rational valued, and by **Theorem A.8** we have that $\sum_{\sigma \in \mathbb{S}_n} \chi_{\text{sq}}(\sigma) = 0$. Thus, we have that

$$\sum_{\sigma \in \mathbb{S}_n} |\chi_{\text{sq}}(\sigma)| \cdot \mathbf{1}_{\chi_{\text{sq}}(\sigma) > 0} = \sum_{\sigma \in \mathbb{S}_n} |\chi_{\text{sq}}(\sigma)| \cdot \mathbf{1}_{\chi_{\text{sq}}(\sigma) < 0} =: C_{\text{sq}} \quad (30)$$

for some C_{sq} (which is nonzero again by **Theorem A.8**). We now define distributions f_1 and f_2 over \mathbb{S}_m as

$$f_1(\sigma) = \begin{cases} \frac{1}{C_{\text{sq}}} \cdot \chi_{\text{sq}}(\sigma) & \text{if } \chi_{\text{sq}}(\sigma) > 0 \\ 0 & \text{otherwise,} \end{cases} \quad f_2(\sigma) = \begin{cases} \frac{-1}{C_{\text{sq}}} \cdot \chi_{\text{sq}}(\sigma) & \text{if } \chi_{\text{sq}}(\sigma) < 0 \\ 0 & \text{otherwise.} \end{cases}$$

From their definitions and **Equation (30)** it is immediate that f_1 and f_2 are distributions over \mathbb{S}_m which have disjoint support. Since $|\mathbb{S}_m| = k$, this gives items 1 and 2 of the theorem.

To prove the third item, observe (recalling the comment immediately after **Definition A.7**) that the function $g : \mathbb{S}_m \rightarrow \mathbb{C}$, defined as $g(\sigma) := f_1(\sigma) - f_2(\sigma) = \frac{1}{C_{\text{sq}}} \cdot \chi_{\text{sq}}(\sigma)$, is a class function. Choose

any partition $\lambda \vdash m$ and the corresponding irreducible representation ρ_λ of \mathbb{S}_m . By applying Lemma A.9, we have that

$$\widehat{g}(\rho_\lambda) = c_\lambda \cdot \text{Id} \quad \text{where} \quad c_\lambda = \frac{\sum_{\sigma \in \mathbb{S}_m} g(\sigma) \cdot \chi_\lambda(\sigma)}{\dim(\rho_\lambda)}. \quad (31)$$

We analyze the multiplier c_λ by noting that

$$\begin{aligned} c_\lambda &= \frac{\sum_{\sigma \in \mathbb{S}_m} g(\sigma) \cdot \chi_\lambda(\sigma)}{\dim(\rho_\lambda)} = \frac{\sum_{\sigma \in \mathbb{S}_m} \chi_{\text{sq}}(\sigma) \cdot \chi_\lambda(\sigma)}{\dim(\rho_\lambda) \cdot C_{\text{sq}}} \\ &= \frac{m! \cdot \mathbb{1}[\lambda = \lambda_{\text{sq}}]}{\dim(\rho_\lambda) \cdot C_{\text{sq}}} \quad \text{using Theorem A.8.} \end{aligned} \quad (32)$$

Thus, we have

$$\begin{aligned} \|\mathcal{M}_\theta * f_1 - \mathcal{M}_\theta * f_2\|_1 &= \sum_{\sigma \in \mathbb{S}_m} |\mathcal{M}_\theta * f_1(\sigma) - \mathcal{M}_\theta * f_2(\sigma)| \\ &= \sum_{\sigma \in \mathbb{S}_m} |\mathcal{M}_\theta * g(\sigma)| \quad (\text{linearity and } g = f_1 - f_2) \\ &= \frac{1}{m!} \sum_{\sigma \in \mathbb{S}_m} \left| \sum_{\mu \vdash m} \dim(\rho_\mu) \text{Tr}[\widehat{\mathcal{M}}_\theta * g(\rho_\mu) \rho_\mu(\sigma^{-1})] \right| \\ &\quad (\text{Definition A.5, inverse Fourier transform of } \mathcal{M}_\theta * g) \\ &= \frac{1}{m!} \sum_{\sigma \in \mathbb{S}_m} \left| \sum_{\mu \vdash m} \dim(\rho_\mu) \text{Tr}[\widehat{\mathcal{M}}_\theta(\rho_\mu) \widehat{g}(\rho_\mu) \rho_\mu(\sigma^{-1})] \right| \quad (\text{convolution identity}) \\ &= \frac{1}{\dim(\rho_{\lambda_{\text{sq}}}) \cdot C_{\text{sq}}} \sum_{\sigma \in \mathbb{S}_m} \left| \dim(\rho_{\lambda_{\text{sq}}}) \text{Tr}[\widehat{\mathcal{M}}_\theta(\rho_{\lambda_{\text{sq}}}) \rho_{\lambda_{\text{sq}}}(\sigma^{-1})] \right| \\ &\quad (\text{Equations 31 and 32}) \\ &= \frac{1}{C_{\text{sq}}} \sum_{\sigma \in \mathbb{S}_m} \left| \text{Tr}[\widehat{\mathcal{M}}_\theta(\rho_{\lambda_{\text{sq}}}) \rho_{\lambda_{\text{sq}}}(\sigma^{-1})] \right| \end{aligned} \quad (33)$$

To deal with $\widehat{\mathcal{M}}_\theta(\rho_{\lambda_{\text{sq}}})$, we apply Lemma 8.1. In particular, by setting $n = m$ and $\mu = \lambda_{\text{sq}}$ in Lemma 8.1, we get that

$$\widehat{\mathcal{M}}_\theta(\rho_{\lambda_{\text{sq}}}) = c_{\lambda_{\text{sq}}, \theta} \cdot \text{Id},$$

where $|c_{\lambda_{\text{sq}}, \theta}| \leq \eta^t$, and we thus get that

$$\|\mathcal{M}_\theta * f_1 - \mathcal{M}_\theta * f_2\|_1 \leq \frac{\eta^t}{C_{\text{sq}}} \cdot \sum_{\sigma \in \mathbb{S}_m} |\text{Tr}[\rho_{\lambda_{\text{sq}}}(\sigma^{-1})]| = \frac{\eta^t}{C_{\text{sq}}} \cdot \sum_{\sigma \in \mathbb{S}_m} |\chi_{\text{sq}}(\sigma^{-1})|. \quad (34)$$

Finally, recalling that

$$C_{\text{sq}} = \frac{\sum_{\sigma \in \mathbb{S}_n} |\chi_{\text{sq}}(\sigma)|}{2},$$

we get that the RHS of Equation (34) is $2\eta^t$. Recalling that $t \geq \sqrt{m/2}$, the theorem is proved. \square

A Basics of representation theory over the symmetric group

Representation theory of the symmetric group \mathbb{S}_n is at the technical core of this paper. In this appendix we briefly review the definitions and results that we require, starting first with general groups and then specializing to \mathbb{S}_n as necessary. See Curtis and Reiner [CR66] (or many other sources) for an extensive reference on representation theory of finite groups and James [Jam06] or Méliot [Mél17] for an extensive reference on representation theory of S_n .

A.1 General groups

We start by recalling the definition of a representation:

Definition A.1. For any group G , a *representation* $\rho : G \rightarrow \mathbb{C}^{m \times m}$ is a group homomorphism, i.e. a function from G to $\mathbb{C}^{m \times m}$ that satisfies $\rho(g) \cdot \rho(h) = \rho(g \cdot h)$ for all $g, h \in G$. The *dimension* of such a representation ρ is m .

In this paper, unless otherwise mentioned, all representations ρ are *unitary* – in other words, for every $g \in G$, $\rho(g)$ is a unitary matrix. Over finite groups, any representation can be made unitary by applying a similarity transformation; by this we mean that if ρ is a representation, then there is an invertible matrix Z such that the new map $\tilde{\rho}$ defined as $\tilde{\rho}(g) = Z^{-1} \cdot \rho(g) \cdot Z$ is a unitary representation. (The reader should verify that as long as Z is invertible, the map $\tilde{\rho}$ is always a representation if ρ is a representation.) Two such representations ρ and $\tilde{\rho}$ are said to be *equivalent*.

Next we recall the notion of an irreducible representation:

Definition A.2. A representation $\rho : G \rightarrow \mathbb{C}^{m \times m}$ is said to be *reducible* if there exists a proper subspace V of \mathbb{C}^m such that $\rho(g) \cdot V \subseteq V$ for all $g \in G$. If there is no such proper subspace V , then ρ is said to be *irreducible*.

It is well known that any finite group has only finitely many irreducible representations, up to the above notion of equivalence, and that every representation of a finite group G can be written as a direct sum of irreducible representations:

Theorem A.3 (Maschke’s theorem, see e.g. Theorem 1.3 of [Mél17]). *For G a finite group, there is a finite set of distinct irreducible representations $\{\rho_1, \dots, \rho_r\}$ such that for any representation $\rho : G \rightarrow \mathbb{C}^{m \times m}$, there is a invertible transformation $Z \in \mathbb{C}^{m \times m}$ such that $Z^{-1}\rho Z$ is block diagonal where each block is one of $\{\rho_1, \dots, \rho_r\}$. In other words, $Z^{-1}\rho Z$ is equal to the direct sum $\bigoplus_{\ell=1}^M \mu_\ell$ where each μ_ℓ is an element of $\{\rho_1, \dots, \rho_r\}$.*

We remind the reader that elements g, h in a group G are said to be *conjugates* if there is an element $t \in G$ such that $tgt^{-1} = h$. Define $\text{Cl}(g)$, the *conjugacy class* of g , to be $\{h : h \text{ is conjugate to } g\}$; it is easy to see that the different conjugacy classes form a partition of G .

We recall some very standard facts about irreducible representations:

Theorem A.4 (see e.g. Theorem 2.3.1 of [GW10]). *Let G be a finite group and let $\{\rho_1, \dots, \rho_r\}$ be the set of its irreducible representations, where $\rho_i : G \rightarrow \mathbb{C}^{d_i \times d_i}$. Then*

1. $\sum_{i=1}^r d_i^2 = |G|$.
2. *The number of conjugacy classes is equal to r , the number of distinct irreducible representations.*

3. For $1 \leq s, t \leq d_i$, let $\rho_{i,s,t} : G \rightarrow \mathbb{C}$ be the (s, t) entry of $\rho_i(g)$. Then, for $1 \leq i_1, i_2 \leq r$, $1 \leq s_1, t_1 \leq d_{i_1}$ and $1 \leq s_2, t_2 \leq d_{i_2}$

$$\mathbf{E}_{g \in G} [\rho_{i_1, s_1, t_1}(g) \cdot \overline{\rho_{i_2, s_2, t_2}(g)}] = \begin{cases} \frac{1}{d_{i_1}} & \text{if } i_1 = i_2, s_1 = s_2 \text{ and } t_1 = t_2 \\ 0 & \text{otherwise} \end{cases}$$

4. The representations ρ_1, \dots, ρ_r are unitary.

A restatement of (3) above is that the functions $\{\rho_{i,s,t}(\cdot)\}$ are orthogonal. Combining this with $\sum_{i=1}^r d_i^2 = |G|$ (given by (1)), we get that the functions $\{\rho_{i,s,t}\}_{1 \leq i \leq r, 1 \leq s, t \leq d_i}$ form an orthogonal basis for \mathbb{C}^G .

With an orthonormal basis for the set of complex-valued functions on G in hand (in other words, a basis for the group algebra $\mathbb{C}[G]$), we are ready to define the *Fourier transform* of a function $f : G \rightarrow \mathbb{C}$:

Definition A.5. Let G be a finite group with irreducible representations given by $\{\rho_1, \dots, \rho_r\}$ and let $f : G \rightarrow \mathbb{C}$. The *Fourier transform* of f is given by matrices $\widehat{f}(\rho_1), \dots, \widehat{f}(\rho_r)$, where

$$\widehat{f}(\rho_i) = \sum_{g \in G} f(g) \cdot \rho_i(g).$$

The inverse transform is given by

$$f(g) = \frac{1}{|G|} \sum_{i=1}^r \dim(\rho_i) \text{Tr}[\widehat{f}(\rho_i) \rho_i(g^{-1})].$$

Parseval's identity states that for any f as above, we have

$$\sum_{i=1}^r \|\widehat{f}(\rho_i)\|_F^2 = |G| \cdot \sum_{g \in G} |f(g)|^2. \quad (35)$$

We next recall the definition of characters and class functions for a group G .

Definition A.6. Given a finite group G , a function $f : G \rightarrow \mathbb{C}$ is said to be a *class function* of G if $f(g)$ only depends on the conjugacy class of g , i.e. $f(g) = f(hgh^{-1})$ for every $h \in G$.

Definition A.7. The *character* $\chi_\rho : G \rightarrow \mathbb{C}$ corresponding to a representation $\rho : G \rightarrow \mathbb{C}^{m \times m}$ is given by $\chi_\rho(g) := \text{Tr}(\rho(g))$.

We observe that $\chi_\rho(\cdot)$ is a class function of G , and that if ρ and $\tilde{\rho}$ are unitarily equivalent, then $\chi_\rho(\cdot) = \chi_{\tilde{\rho}}(\cdot)$. We recall some standard facts about characters and class functions:

Theorem A.8. Let G be a finite group and let $\{\rho_1, \dots, \rho_r\}$ be its set of irreducible representations. Let $\chi_{\rho_1}, \dots, \chi_{\rho_r}$ be the corresponding characters. Then we have:

1. [Schur's lemma] $\mathbf{E}_{g \in G} [\chi_{\rho_i}(g) \cdot \overline{\chi_{\rho_j}(g)}] = \delta_{i,j}$.
2. The functions $\{\chi_{\rho_i}(\cdot)\}_{1 \leq i \leq r}$ forms an orthonormal basis for all class functions of G .

The following (standard) claim shows that the Fourier transform of any class function f is a diagonal matrix (in fact, a scalar multiple of the identity matrix):

Lemma A.9. *Let $f : G \rightarrow \mathbb{C}$ be a class function and let $\rho : G \rightarrow \mathbb{C}^{m \times m}$ be an irreducible representation of G . Then $\widehat{f}(\rho) = c \cdot \text{Id}$ where $c = \frac{\sum_{g \in G} f(g)\chi_\rho(g)}{m}$ and Id is the identity matrix.*

Proof. Choose any $h \in G$, and observe that

$$\begin{aligned} \rho(h) \cdot \widehat{f}(\rho) &= \rho(h) \cdot \left(\sum_{g \in G} f(g)\rho(g) \right) \\ &= \rho(h) \cdot \left(\sum_{g \in G} f(h^{-1}gh)\rho(h^{-1}gh) \right) = \rho(h) \cdot \left(\sum_{g \in G} f(g)\rho(h^{-1}gh) \right) \\ &= \rho(h) \cdot \rho(h^{-1}) \cdot \left(\sum_{g \in G} f(g)\rho(g) \right) \cdot \rho(h) = \widehat{f}(\rho) \cdot \rho(h). \end{aligned}$$

As a consequence of Schur's lemma, we have that if a matrix A is such that $A \cdot \rho(h) = \rho(h) \cdot A$ for all $h \in G$, then $A = c \cdot \text{Id}$. Thus, we get that $\widehat{f}(\rho) = c \cdot \text{Id}$. The lemma follows by taking trace on both sides. \square

A.2 Representation theory of the symmetric group

Representation theory of the symmetric group has many applications to algebra, combinatorics and statistical physics and has been intensively studied (as mentioned earlier, see e.g. [Jam06, M el17] for detailed treatments). Below we only recall a few basics which we will need.

The first notion we require is that of a *Young diagram*. Consider a partition $\lambda = (\lambda_1, \dots, \lambda_k)$ of n where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k > 0$ and $\lambda_1 + \dots + \lambda_k = n$. We indicate that λ is such a partition by writing " $\lambda \vdash n$." The Young diagram corresponding to such a partition λ is a two-dimensional left-justified array of empty cells in which the i^{th} row has λ_i cells. See the left portion of Figure 2 for an example of a Young diagram. A *Young tableau* corresponding to a partition λ is obtained by filling in the n cells of the Young diagram with the elements of $[n]$, using each element exactly once, where the ordering within rows of the Young diagram is irrelevant.

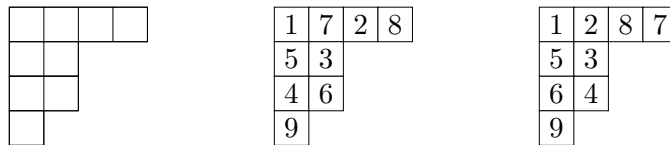


Figure 2: On the left is the Young diagram for the partition $\lambda = (4, 2, 2, 1)$. The middle and right present two equivalent Young tableaux for $(\{1, 7, 2, 8\}, \{5, 3\}, \{4, 6\}, \{9\})$.

For each partition $\lambda = (\lambda_1, \dots, \lambda_k)$ of n , there is an associated representation, denoted τ_λ , which we now define. Let $N_\lambda = \binom{n}{\lambda_1, \dots, \lambda_k}$ be the number of Young tableaux corresponding to partition λ , and let $Y_{\lambda,1}, \dots, Y_{\lambda,N_\lambda}$ be an enumeration of these tableaux in some order.

Definition A.10. The *permutation representation* τ_λ corresponding to λ is defined as follows: For each $g \in \mathbb{S}_n$, $\tau_\lambda(g)$ is the $N_\lambda \times N_\lambda$ matrix (where we view rows and columns as indexed by Young tableaux corresponding to λ) which has $\tau_\lambda(g)(i, j) = 1$ iff $Y_{\lambda,i}$ maps to $Y_{\lambda,j}$ under the action of g .

It is easy to check that $\tau_\lambda : \mathbb{S}_n \rightarrow \mathbb{C}^{N_\lambda \times N_\lambda}$ as defined above is indeed a representation. In fact, since the range of τ_λ is always a permutation matrix, τ_λ is also a unitary representation.

It turns that for $\lambda \neq (n)$, the permutation representation τ_λ is not an irreducible representation. However, it also turns out that all of the irreducible representations of \mathbb{S}_n can be obtained from the permutation representations. To explain this, we need to define a partial order over partitions of n :

Definition A.11. For two partitions λ and μ of n , we say that λ *dominates* μ , written $\lambda \triangleright \mu$, if $\sum_{j \leq i} \lambda_j \geq \sum_{j \leq i} \mu_j$ for all $i > 0$. The partial order defined by \triangleright is said to be the *dominance order* over the partitions (equivalently, Young diagrams) of n .

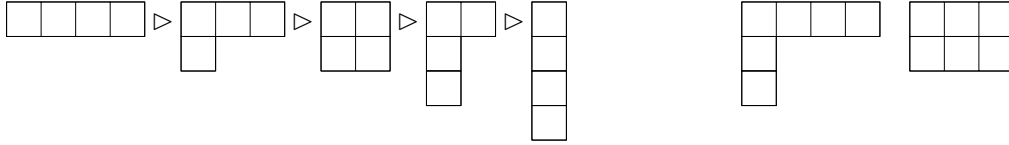


Figure 3: The left part of the picture depicts the dominance order across the partitions of 4; it happens to be the case that the dominance order is a total order across the partitions of 4. This is not true in general; as depicted on the right, the two partitions $(4, 1, 1)$ and $(3, 3)$ of 6 are incomparable under the dominance order.

The next result explains how the irreducible representations of \mathbb{S}_n can be obtained from the representations $\{\tau_\lambda\}_{\lambda \vdash n}$:

Theorem A.12 (James submodule theorem, see e.g. Theorem 3.34 of [Mél17]). *The irreducible representations of \mathbb{S}_n are in one-to-one correspondence with the partitions $\lambda \vdash n$; we denote the irreducible representation corresponding to λ by ρ_λ . In particular, when $\lambda = (n)$, then ρ_λ is the trivial irreducible representation (which maps each $g \in G$ to 1). Moreover, each permutation representation τ_λ is a direct sum of irreducible representations corresponding to partitions which dominate λ , i.e.*

$$\tau_\lambda = \bigoplus_{\mu \triangleright \lambda} \bigoplus_{\ell=1}^{K_{\lambda,\mu}} \rho_\mu.$$

Here the $K_{\lambda,\mu}$'s are non-negative integers, known as the Kostka numbers, which are such that $K_{\lambda,\lambda} = 1$.

A.2.1 Restrictions of irreducible representations

Fix $\lambda \vdash n$ and consider the irreducible representation ρ_λ of \mathbb{S}_n . For any $m \leq n$, \mathbb{S}_m can be viewed as the subgroup of \mathbb{S}_n where elements $\{m+1, \dots, n\}$ are fixed. Hence ρ_λ can also be viewed as a representation of \mathbb{S}_n ; this representation of \mathbb{S}_m is written ρ_λ^m and is called the *restriction* of ρ_λ to

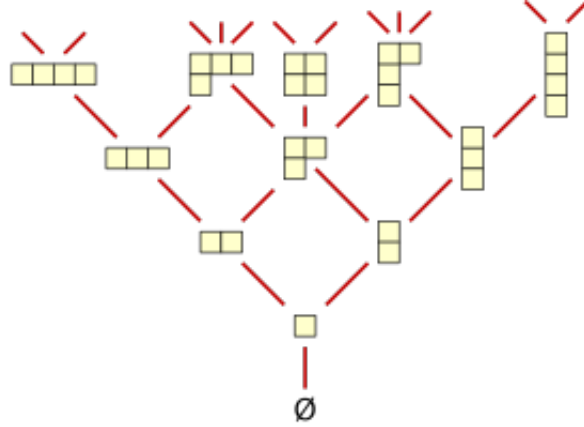


Figure 4: The first five levels of Young's lattice.

\mathbb{S}_m . Note that ρ_λ^m may not be an *irreducible* representation of \mathbb{S}_m . By [Theorem A.3](#), we have that ρ_λ^m is equivalent to some direct sum

$$\bigoplus_{\mu \vdash m} M_{\lambda, \mu} \rho_\mu,$$

in which there are $M_{\lambda, \mu}$ many copies of ρ_μ , for some non-negative integers $M_{\lambda, \mu}$. These integers are given by the so-called “branching rule” on Young's lattice, which we now describe.

Definition A.13. *Young's lattice* is the partially ordered set of Young diagrams in which the partial order is given by inclusion in the following sense: given partitions μ and λ , we write “ $\mu \uparrow \lambda$ ” if λ can be obtained by adding one box to μ (in such a way that λ is a valid partition, of course). If there are partitions μ^1, \dots, μ^r such that $\mu^1 \uparrow \mu^2 \uparrow \dots \uparrow \mu^r$, we write “ $\mu^1 \uparrow \mu^r$.”

It is convenient to draw Young's lattice in such a way that the n -th level contains all and only the Young diagrams with n boxes. The diagram in [Figure 4](#) depicts the first five levels of Young's lattice.

The next result, known as the “branching rule,” states that for $\lambda \vdash n$, ρ_λ splits into a direct sum of ρ_μ over all $\mu \uparrow \lambda$ when ρ_λ is restricted to \mathbb{S}_{n-1} :

Lemma A.14 (Branching rule). *Let λ be a partition of n and let ρ_λ be the corresponding irreducible representation of \mathbb{S}_n . Then ρ_λ^{n-1} , the restriction of ρ_λ to \mathbb{S}_{n-1} , is equivalent to*

$$\bigoplus_{\mu \vdash n-1: \mu \uparrow \lambda} \rho_\mu.$$

By applying [Lemma A.14](#) inductively we get a complete description of how ρ_λ splits when it is restricted to any \mathbb{S}_m , $m < n$:

Theorem A.15. *Let $\lambda \vdash n$ and let ρ_λ be the corresponding irreducible representation of \mathbb{S}_n . For $m < n$ we have that ρ_λ^m , the restriction of ρ_λ to \mathbb{S}_m , is equivalent to*

$$\bigoplus_{\mu \vdash m} \text{Paths}(\mu, \lambda) \rho_\mu,$$

where $\text{Paths}(\mu, \lambda)$ denotes the number of paths in Young's lattice from μ to λ .

Irreducible characters of the symmetric group. Finally, we recall the following fundamental fact (which is a consequence, e.g., of the Murnaghan-Nakayama rule) which we will use:

Fact A.16. [see e.g. Theorem 3.10 in [Mél17]] Let $\chi : \mathbb{S}_m \rightarrow \mathbb{C}$ be a character of \mathbb{S}_m . Then in fact χ is \mathbb{Q} -valued.

Acknowledgments

We thank Mike Saks for allowing us to include his proof of Claim 2.2 here. We also thank Vic Reiner and Yuval Roichman for answering several questions about representation theory. Anindya is grateful to Aravindan Vijayaraghavan for many useful discussions about ranking models.

References

- [ABSV14] P. Awasthi, A. Blum, O. Sheffet, and A. Vijayaraghavan. Learning mixtures of ranking models. In *Advances in Neural Information Processing Systems*, pages 2609–2617, 2014. [1.3](#)
- [BM08] M. Braverman and E. Mossel. Noisy sorting without resampling. In *Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 268–276, 2008. [1.3](#)
- [BOB07] L. Busse, P. Orbanz, and J. Buhmann. Cluster analysis of heterogeneous rank data. In *Proceedings of the 24th ICML*, pages 113–120, 2007. [1.3](#)
- [Cho94] K. P. Choi. On the medians of gamma distributions and an equation of Ramanujan. *Proc. Amer. Math. Soc.*, 121:245–251, 1994. [5.2](#)
- [CR66] C. Curtis and I. Reiner. *Representation theory of finite groups and associative algebras*, volume 356. American Mathematical Society, 1966. [A](#)
- [DH92] Persi Diaconis and Phil Hanlon. Eigen Analysis for Some Examples of the Metropolis Algorithm. *Contemporary Mathematics*, 138:99–117, 1992. [1.3](#), [3.1](#)
- [Dia88a] P. Diaconis. Group representations in probability and statistics. *Lecture Notes-Monograph Series*, 11:i–192, 1988. [1.3](#)
- [Dia88b] Persi Diaconis. *Chapter 6: Metrics on Groups, and Their Statistical Uses*, volume Volume 11 of *Lecture Notes-Monograph Series*, pages 102–130. Institute of Mathematical Statistics, 1988. [1.1](#)
- [DS81] Persi Diaconis and Mehrdad Shahshahani. Generating a Random Permutation with Random Transpositions. *Z. Wahrscheinlichkeitstheorie verw. Gebiete*, 57:159–179, 1981. [1.1](#), [5.1](#), [5.1](#)
- [DS98] Persi Diaconis and Laurent Saloff-Coste. What do we know about the metropolis algorithm? *J. Comput. Syst. Sci.*, 57(1):20–36, 1998. [3.1](#)

- [DST16] A. De, M. Saks, and S. Tang. Noisy population recovery in polynomial time. In *2016 Foundations of Computer Science*, pages 675–684. IEEE, 2016. [2.2](#)
- [Ewe72] W. Ewens. The sampling theory of selectively neutral alleles. *Theoretical Population Biology*, 3:87–112, 1972. [1.1](#), [1.3](#)
- [FV86] M. Fligner and J. Verducci. Distance based ranking models. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 359–369, 1986. [1.1](#), [1.3](#)
- [GJ79] Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman, 1979. [1](#)
- [GP18] A. Gladkikh and R. Peled. On the cycle structure of mallows permutations. 46(2):1114–1169, 03 2018. [1.3](#)
- [GW10] B. Green and A. Wigderson. Lecture notes for the 22nd McGill Invitational Workshop on Computational Complexity. 2010. [A.4](#)
- [Jam06] Gordon Douglas James. *The representation theory of the symmetric groups*, volume 682. Springer, 2006. [A](#), [A.2](#)
- [JV18] Y. Jiao and J. Vert. The Kendall and Mallows kernels for permutations. *IEEE transactions on pattern analysis and machine intelligence*, 40(7):1755–1769, 2018. [1.1](#), [1.3](#)
- [KB10] R. Kondor and M. Barbosa. Ranking with Kernels in Fourier space. In *COLT 2010*, pages 451–463, 2010. [1.1](#), [1.3](#)
- [KL02] R. Kondor and J. Lafferty. Diffusion kernels on graphs and other discrete structures. In *Machine Learning, Proceedings of the 19th International Conference (ICML 2002)*, 2002. [1.1](#), [1.3](#)
- [KV10] R. Kumar and S. Vassilvitskii. Generalized distances between rankings. In *WWW*, pages 571–580, 2010. [1.1](#)
- [LB11] T. Lu and C. Boutilier. Learning Mallows models with pairwise preferences. In *Proceedings of the 28th ICML*, pages 145–152, 2011. [1.3](#), [1.3](#)
- [LL02] G. Lebanon and J. Lafferty. Cranking: Combining rankings using conditional probability models on permutations. In *Proceedings of the Nineteenth International Conference on Machine Learning*, pages 363–370, 2002. [1.3](#)
- [LM18] A. Liu and A. Moitra. Efficiently Learning Mixtures of Mallows Models. In *Proceedings of FOCS, 2018*, 2018. [1.3](#)
- [LZ15] Shachar Lovett and Jiapeng Zhang. Improved noisy population recovery, and reverse Bonami-Beckner inequality for sparse functions. In *Proceedings of the 47th Annual ACM Symposium on Theory of Computing*, pages 137–142, 2015. [2.2](#)
- [Mal57] C. Mallows. Non-null ranking models. I. *Biometrika*, 44(1/2):114–130, 1957. [1.3](#), [1.3](#)
- [Mar14] J. Marden. *Analyzing and modeling rank data*. Chapman and Hall/CRC, 2014. [1.3](#)

- [MC10] M. Meilă and H. Chen. Dirichlet process mixtures of generalized mallows models. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pages 358–367, 2010. [1.3](#)
- [Mél17] P. Méliot. *Representation theory of symmetric groups*. Chapman and Hall/CRC, 2017. [6.5](#), [A](#), [A.3](#), [A.2](#), [A.12](#), [A.16](#)
- [MM03] T. Murphy and D. Martin. Mixtures of distance-based models for ranking data. *Computational statistics & data analysis*, 41(3-4):645–655, 2003. [1.3](#)
- [MM09] B. Mandhani and M. Meila. Tractable search for learning exponential models of rankings. In *Artificial Intelligence and Statistics*, pages 392–399, 2009. [1.3](#)
- [MPPB07] M. Meilă, K. Phadnis, A. Patterson, and J. Bilmes. Consensus ranking under the exponential model. In *Proceedings of the Twenty-Third Conference on Uncertainty in Artificial Intelligence*, pages 285–294, 2007. [1.3](#)
- [MS13] Ankur Moitra and Michael Saks. A polynomial time algorithm for lossy population recovery. In *2013 Foundations of Computer Science*, pages 110–116. IEEE, 2013. [2.2](#)
- [Muk16] S. Mukherjee. Estimation in exponential families on permutations. *The Annals of Statistics*, 44(2):853–875, 2016. [1.3](#)
- [Sak18] M. Saks. Personal communication, 2018. [2.1](#)
- [Sta99] Richard P. Stanley. *Enumerative Combinatorics: Volume 2*. Cambridge University Press, 1999. [6](#), [6](#)
- [Ste77] G. W. Stewart. On the Perturbation of Pseudo-Inverses, Projections and Linear Least Squares Problems. *SIAM Review*, 19(4):634–662, 1977. [3](#)
- [WY12] Avi Wigderson and Amir Yehudayoff. Population Recovery and Partial Identification. In *53rd Annual IEEE Symposium on Foundations of Computer Science*, pages 390–399, 2012. [2.1](#), [2.2](#)