

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

---

# Learning Mixtures of Ranking Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

This work concerns learning probabilistic models for ranking data in a heterogeneous population. The specific problem we study is learning the parameters of a *Mallows Mixture Model*. Despite being widely studied, current heuristics for this problem do not have theoretical guarantees and can get stuck in bad local optima. We present the first polynomial time algorithm which provably learns the parameters of a mixture of two Mallows models. A key component of our algorithm is a novel use of tensor decomposition techniques to learn the top- $k$  prefix in both the rankings. Before this work, even the question of *identifiability* in the case of a mixture of two Mallows models was unresolved.

## 1 Introduction

Probabilistic modeling of ranking data is an extensively studied problem with a rich body of past work [1, 2, 3, 4, 5, 6, 7, 8, 9]. Ranking using such models has applications in a variety of areas ranging from understanding user preferences in electoral systems and social choice theory, to more modern learning tasks in online web search, crowd-sourcing and recommendation systems. Traditionally, models for generating ranking data consider a homogeneous group of users with a *central ranking* (permutation)  $\pi^*$  over a set of  $n$  elements or alternatives. (For instance,  $\pi^*$  might correspond to a “ground-truth ranking” over a set of movies.) Each individual user generates her own ranking as a noisy version of this one central ranking and independently from other users. The most popular ranking model of choice is the *Mallows model* [1], where in addition to  $\pi^*$  there is also a scaling parameter  $\phi \in (0, 1)$ . Each user picks her ranking  $\pi$  w.p. proportional to  $\phi^{d_{kt}(\pi, \pi^*)}$  where  $d_{kt}(\cdot)$  denotes the Kendall-Tau distance between permutations (see Section 2).<sup>1</sup> We denote such a model as  $\mathcal{M}_n(\phi, \pi^*)$ .

The Mallows model and its generalizations have received much attention from the statistics, political science and machine learning communities, relating this probabilistic model to the long-studied work about voting and social choice [10, 11]. From a machine learning perspective, the problem is to find the parameters of the model — the central permutation  $\pi^*$  and the scaling parameter  $\phi$ , using independent samples from the distribution. There is a large body of work [4, 6, 5, 7, 12] providing efficient algorithms for learning the parameters of a Mallows model.

In many scenarios, however, the population is heterogeneous with multiple groups of people, each with their own central ranking [2]. For instance, when ranking movies, the population may be divided into two groups corresponding to men and women; with men ranking movies with one underlying central permutation, and women ranking movies with another underlying central permutation. This naturally motivates the problem of learning a *mixture* of multiple Mallows models for rankings, a problem that has received significant attention [8, 13, 3, 4]. Heuristics like the EM algorithm have

---

<sup>1</sup>In fact, it was shown [1] that this model is the result of the following simple (inefficient) algorithm: rank every pair of elements randomly and independently s.t. with probability  $\frac{1}{1+\phi}$  they agree with  $\pi^*$  and with probability  $\frac{\phi}{1+\phi}$  they don't; if all  $\binom{n}{2}$  pairs agree on a single ranking – output this ranking, otherwise resample.

054 been applied to learn the model parameters of a mixture of Mallows models [8]. The problem has  
 055 also been studied under distributional assumptions over the parameters, e.g. weights derived from  
 056 a Dirichlet distribution [13]. However, unlike the case of a single Mallows model, algorithms with  
 057 provable guarantees have remained elusive for this problem.

058 In this work we give the *first polynomial time algorithm* that *provably* learns a mixture of two  
 059 Mallows models. The input to our algorithm consists of i.i.d random rankings (samples), with  
 060 each ranking drawn with probability  $w_1$  from a Mallows model  $\mathcal{M}_n(\phi_1, \pi_1)$ , and with probability  
 061  $w_2 (= 1 - w_1)$  from a different model  $\mathcal{M}_n(\phi_2, \pi_2)$ .

062 **Informal Theorem.** *Given sufficiently many i.i.d samples drawn from a mixture of two Mallows*  
 063 *models, we can learn the central permutations  $\pi_1, \pi_2$  exactly and parameters  $\phi_1, \phi_2, w_1, w_2$  upto*  
 064  *$\epsilon$ -accuracy in time  $\text{poly}(n, (\min\{w_1, w_2\})^{-1}, \frac{1}{\phi_1(1-\phi_1)}, \frac{1}{\phi_2(1-\phi_2)}, \epsilon^{-1})$ .*

066 It is worth mentioning that, to the best of our knowledge, prior to this work even the question of iden-  
 067 tifiability was unresolved for a mixture of two Mallows models; given infinitely many i.i.d. samples  
 068 generated from a mixture of two distinct Mallows models with parameters  $\{w_1, \phi_1, \pi_1, w_2, \phi_2, \pi_2\}$   
 069 (with  $\pi_1 \neq \pi_2$  or  $\phi_1 \neq \phi_2$ ), could there be a different set of parameters  $\{w'_1, \phi'_1, \pi'_1, w'_2, \phi'_2, \pi'_2\}$   
 070 which explains the data just as well. Our result shows that this is not the case and the mixture is  
 071 uniquely identifiable given polynomially many samples.

072 **Intuition and a Naïve First Attempt.** It should be evident that having access to sufficiently many  
 073 random samples allows one to learn a single Mallows model. Let the elements in the permutations  
 074 be denoted as  $\{e_1, e_2, \dots, e_n\}$ . In a single Mallows model, the probability of element  $e_i$  going to  
 075 position  $j$  (for all  $j = 1, 2, \dots, n$ ) drops off exponentially as one goes farther from the true position  
 076 of  $e_i$  [12]. So by assigning each  $e_i$  the most frequent position in our sample, we can find the central  
 077 ranking  $\pi^*$ .

078 The above mentioned intuition suggests the following clustering based approach to learn a mixture  
 079 of two Mallows models — look at the distribution of the positions where element  $e_i$  appears. If the  
 080 distribution has 2 clearly separated “peaks” then they will correspond to the positions of  $e_i$  in the  
 081 central permutations. Now, dividing the samples according to  $e_i$  being ranked in a high or a low  
 082 position is likely to give us two pure (or almost pure) subsamples, each one coming from a single  
 083 Mallows model. We can then learn the individual models separately. More generally, this strategy  
 084 works when the two underlying permutations  $\pi_1$  and  $\pi_2$  are far apart which can be formulated as  
 085 a *separation condition*.<sup>2</sup> Indeed, the above-mentioned intuition works only under strong separator  
 086 conditions: otherwise, the observation regarding the distribution of positions of element  $e_i$  is no  
 087 longer true<sup>3</sup>. For example, if  $\pi_1$  ranks  $e_i$  in position  $k$  and  $\pi_2$  ranks  $e_i$  in position  $k + 2$ , it is likely  
 088 that the most frequent position of  $e_i$  is  $k + 1$ , which differs from  $e_i$ ’s position in either permutations!

089 **Handling arbitrary permutations.** Learning mixture models under no separation requirements is  
 090 a challenging task. To the best of our knowledge, the only polynomial time algorithm known is  
 091 for the case of a mixture of a constant number of Gaussians [17, 18]. Other works, like the recent  
 092 developments that use tensor based methods for learning mixture models without distance-based  
 093 separation condition [19, 20, 21] still require non-degeneracy conditions and/or work for specific  
 094 sub cases (e.g. spherical Gaussians).

095 These sophisticated tensor methods form a key component in our algorithm for learning a mixture  
 096 of two Mallows models. This is non-trivial as learning over rankings poses challenges which are  
 097 not present in other widely studied problems such as mixture of Gaussians. For the case of Gaus-  
 098 sians, spectral techniques have been extremely successful [22, 16, 19, 21]. Such techniques rely on  
 099 estimating the covariances and higher order moments in terms of the model parameters to detect  
 100 structure and dependencies. On the other hand, in the mixture of Mallows models problem there is  
 101 no “*natural*” notion of a *second/third moment*. A key contribution of our work is defining analogous  
 102 notions of moments which can be represented succinctly in terms of the model parameters. As we  
 103 later show, this allows us to use tensor based techniques to get a good starting solution.

105 <sup>2</sup>Identifying a permutation  $\pi$  over  $n$  elements with a  $n$ -dimensional vector  $(\pi(i))_i$ , this separation condition  
 106 can be roughly stated as  $\|\pi_1 - \pi_2\|_\infty = \tilde{\Omega}\left(\frac{1}{\min\{w_1, w_2\}} \cdot \frac{1}{\min\{\log(1/\phi_1), \log(1/\phi_2)\}}\right)$ .

107 <sup>3</sup>Much like how other mixture models are solvable under separation conditions, see [14, 15, 16].

**Overview of Techniques.** One key difficulty in arguing about the Mallows model is the lack of closed form expressions for basic propositions like “the probability that the  $i$ -th element of  $\pi^*$  is ranked in position  $j$ .” Our first observation is that the distribution of a given element appearing at the top, i.e. the first position, behaves nicely. Given an element  $e$  whose rank in the central ranking  $\pi^*$  is  $i$ , the probability that a ranking sampled from a Mallows model ranks  $e$  as the first element is  $\propto \phi^{i-1}$ . A length  $n$  vector consisting of these probabilities is what we define as the *first moment vector* of the Mallows model. Clearly by sorting the coordinate of the first moment vector, one can recover the underlying central permutation and estimate  $\phi$ . Going a step further, consider any two elements which are in positions  $i, j$  respectively in  $\pi^*$ . We show that the probability that a ranking sampled from a Mallows model ranks  $\{i, j\}$  in (any of the  $2!$  possible ordering of) the first two positions is  $\propto f(\phi)\phi^{i+j-2}$ . We call the  $n \times n$  matrix of these probabilities as the *second moment matrix* of the model (analogous to the covariance matrix). Similarly, we define the 3rd moment tensor as the probability that any 3 elements appear in positions  $\{1, 2, 3\}$ . We show in the next section that in the case of a mixture of two Mallows models, the 3rd moment tensor defined this way has a rank-2 decomposition, with each rank-1 term corresponds to the first moment vector of each of two Mallows models. This motivates us to use tensor-based techniques to estimate the first moment vectors of the two Mallows models, thus learning the models’ parameters.

The above mentioned strategy would work if one had access to infinitely many samples from the mixture model. But notice that the probabilities in the first-moment vectors decay exponentially, so by using polynomially many samples we can only recover a prefix of length  $\sim \log_{1/\phi} n$  from both rankings. This forms the first part of our algorithm which outputs good estimates of the mixture weights, scaling parameters  $\phi_1, \phi_2$  and prefixes of a certain size from both the rankings. Armed with  $w_1, w_2$  and these two prefixes we next proceed to recover the full permutations  $\pi_1$  and  $\pi_2$ . In order to do this, we take two new fresh batches of samples. On the first batch, we estimate the probability that element  $e$  appears in position  $j$  for all  $e$  and  $j$ . On the second batch, which is noticeably larger than the first, we estimate the probability that  $e$  appears in position  $j$  *conditioned on a carefully chosen element  $e^*$  appearing as the first element*. We show that this conditioning is *almost* equivalent to sampling from the same mixture model but with rescaled weights  $w_1'$  and  $w_2'$ . The two estimations allow us to set a system of two linear equations in two variables:  $f^{(1)}(e \rightarrow j)$  – the probability of element  $e$  appearing in position  $j$  in  $\pi_1$ , and  $f^{(2)}(e \rightarrow j)$  – the same probability for  $\pi_2$ . Solving this linear system we find the position of  $e$  in each permutation.

The above description contains most of the core ideas involved in the algorithm. We need two additional components. First, notice that the 3rd moment tensor is not well defined for triplets  $(i, j, k)$ , when  $i, j, k$  are not all distinct and hence cannot be estimated from sampled data. To get around this barrier we consider a random partition of our element-set into 3 disjoint subsets. The actual tensor we work with consists only of triplets  $(i, j, k)$  where the indices belong to different partitions. Secondly, we have to handle the case where tensor based-technique fails, i.e. when the 3rd moment tensor isn’t full-rank. This is a *degenerate case*. Typically, tensor based approaches for other problems cannot handle such degenerate cases. However, in the case of the Mallows mixture model, we show that such a degenerate case provides a lot of useful information about the problem. In particular, it must hold that  $\phi_1 = \phi_2$ ,<sup>4</sup> and  $\pi_1$  and  $\pi_2$  are fairly close — one is almost a cyclic shift of the other. To show this we use a characterization of the when the tensor decomposition is unique (for tensors of rank 2), and we handle such degenerate cases separately. Altogether, we find the mixture model’s parameters with no non-degeneracy conditions.

**Lower bound under the pairwise access model.** Given that a single Mallows model can be learned using only pairwise comparisons, a very restricted access to each sample, it is natural to ask, “Is it possible to learn a mixture of Mallows models from pairwise queries?”. This next example shows that we cannot hope to do this even for a mixture of two Mallows models. Fix some  $\phi$  and  $\pi$  and assume our sample is taken using mixing weights of  $w_1 = w_2 = \frac{1}{2}$  from the two Mallows models  $\mathcal{M}_n(\phi, \pi)$  and  $\mathcal{M}_n(\phi, \text{rev}(\pi))$ , where  $\text{rev}(\pi)$  indicates the reverse permutation (the first element of  $\pi$  is the last of  $\text{rev}(\pi)$ , the second is the next-to-last, etc.) . Consider two elements,  $e$  and  $e'$ . Using only pairwise comparisons, we have that it is just as likely to rank  $e > e'$  as it is to rank  $e' > e$  and so this case cannot be learned regardless of the sample size.

<sup>4</sup>Or rather:  $|\phi_1 - \phi_2| \ll \epsilon$  where  $\epsilon$  is our accuracy parameter.

162 **3-wise queries** We would also like to stress that our algorithm does not need full access to the  
 163 sampled rankings and instead will work with access to certain 3-wise queries. Observe that the first  
 164 part of our algorithm, where we recover the top elements in each of the two central permutations,  
 165 only uses access to the top 3 elements in each sample. In that sense, we replace the pairwise query  
 166 “do you prefer  $e$  to  $e'$ ?” with a 3-wise query: “what are your top 3 choices?” Furthermore, the  
 167 second part of the algorithm (where we solve a set of 2 linear equations) can be altered to support  
 168 3-wise queries of the (admittedly, somewhat unnatural) form “if  $e^*$  is your top choice, do you prefer  
 169  $e$  to  $e'$ ?” For ease of exposition, we will assume full-access to the sampled rankings for the rest of  
 170 the paper.

171 **Future Directions.** Several interesting directions come out of this work. A natural next step is to  
 172 generalize our results to learn a mixture of  $k$  Mallows models for  $k > 2$ . We believe that most  
 173 of these techniques can be extended to design algorithms that take  $\text{poly}(n, 1/\epsilon)^k$  time. It would  
 174 also be interesting to get algorithms for learning a mixture of  $k$  Mallows models which run in time  
 175  $\text{poly}(k, n)$ , perhaps in an appropriate smoothed analysis setting [23] or under other non-degeneracy  
 176 assumptions. Perhaps, more importantly, our result indicates that tensor based methods which have  
 177 been very popular for learning problems, might also be a powerful tool for tackling ranking-related  
 178 problems in the fields of machine learning, voting and social choice.

179 **Organization.** In Section 2 we give the formal definition of the Mallow model and of the problem  
 180 statement, as well as some useful facts about the Mallow model. Our algorithm and its numerous  
 181 subroutines are detailed in Section 3. In Section 4 we experimentally compare our algorithm with a  
 182 popular EM based approach for the problem. We conclude with future directions in Section ?? . The  
 183 complete details of our algorithms and proofs are included in the supplementary material.

## 185 2 Notations and Properties of the Mallows Model

186  
 187 Let  $U_n = \{e_1, e_2, \dots, e_n\}$  be a set of  $n$  distinct elements. We represent permutations over the  
 188 elements in  $U_n$  through their indices  $[n]$ . (E.g.,  $\pi = (n, n-1, \dots, 1)$  represents the permutation  
 189  $(e_n, e_{n-1}, \dots, e_1)$ .) Let  $\text{pos}_\pi(e_i) = \pi^{-1}(i)$  refer to the position of  $e_i$  in the permutation  $\pi$ . We  
 190 omit the subscript  $\pi$  when the permutation  $\pi$  is clear from context. For any two permutations  $\pi, \pi'$   
 191 we denote  $d_{\text{kt}}(\pi, \pi')$  as the Kendall-Tau distance [24] between them (number of pairwise inversions  
 192 between  $\pi, \pi'$ ). Given some  $\phi \in (0, 1)$  we denote  $Z_i(\phi) = \frac{1-\phi^i}{1-\phi}$ , and partition function  $Z_{[n]}(\phi) =$   
 193  $\sum_\pi \phi^{d_{\text{kt}}(\pi, \pi_0)} = \prod_{i=1}^n Z_i(\phi)$  (see Section ?? in supplementary material).

194 **Definition 2.1.** [Mallows model  $(\mathcal{M}_n(\phi, \pi_0))$ .] Given a permutation  $\pi_0$  on  $[n]$  and a parameter  
 195  $\phi \in (0, 1)$ ,<sup>5</sup> a Mallows model is a permutation generation process that returns permutation  $\pi$  w.p.

$$196 \Pr(\pi) = \phi^{d_{\text{kt}}(\pi, \pi_0)} / Z_{[n]}(\phi)$$

199 In Section ?? of the supplementary material we show many useful properties of the Mallows model  
 200 which we use repeatedly throughout this work. Even though these properties are deferred to the  
 201 supplementary material, we believe they provide an insight to Mallows model, and we advise the  
 202 reader to go through them. We proceed with the main definition.

203 **Definition 2.2.** [Mallows Mixture model  $w_1 \mathcal{M}_n(\phi_1, \pi_1) \oplus w_2 \mathcal{M}_n(\phi_2, \pi_2)$ .] Given parameters  
 204  $w_1, w_2 \in (0, 1)$  s.t.  $w_1 + w_2 = 1$ , parameters  $\phi_1, \phi_2 \in (0, 1)$  and two permutations  $\pi_1, \pi_2$ , we call  
 205 a mixture of two Mallows models to be the process that with probability  $w_1$  generates a permutation  
 206 from  $\mathcal{M}(\phi_1, \pi_1)$  and with probability  $w_2$  generates a permutation from  $\mathcal{M}(\phi_2, \pi_2)$ .

207 Our next definition is crucial for our application of tensor decomposition techniques.

208 **Definition 2.3.** [Representative vectors.] The representative vector of a Mallows model is a vector  
 209 where for every  $i \in [n]$ , the  $i$ th-coordinate is  $\phi^{\text{pos}_\pi(e_i)-1} / Z_n$ .

211 The expression  $\phi^{\text{pos}_\pi(e_i)-1} / Z_n$  is precisely the probability that a permutation generated by a model  
 212  $\mathcal{M}_n(\phi, \pi)$  ranks element  $e_i$  at the first position (proof deferred to the supplementary material).  
 213 Given that our focus is on learning a mixture of two Mallows models  $\mathcal{M}_n(\phi_1, \pi_1)$  and  $\mathcal{M}_n(\phi_2, \pi_2)$ ,  
 214 we denote  $x$  as the representative vector of the first model, and  $y$  as the representative vector of the

215 <sup>5</sup>It is also common to parameterize using  $\beta \in \mathbb{R}^+$  where  $\phi = e^{-\beta}$ . For small  $\beta$  we have  $(1 - \phi) \approx \beta$ .

latter. Note that retrieving the vectors  $x$  and  $y$  exactly implies that we can learn the permutations  $\pi_1$  and  $\pi_2$  and the values of  $\phi_1, \phi_2$ .

**Tensors:** Given two vectors  $u \in \mathbb{R}^{n_1}, v \in \mathbb{R}^{n_2}$ , we define  $u \otimes v \in \mathbb{R}^{n_1 \times n_2}$  as the matrix  $uv^T$ . Given also  $z \in \mathbb{R}^{n_3}$  then  $u \otimes v \otimes z$  denotes the 3-tensor (of rank-1) whose  $(i, j, k)$ -th coordinate is  $u_i v_j z_k$ . A tensor  $T \in \mathbb{R}^{n_1 \times n_2 \times n_3}$  has a rank- $r$  decomposition if  $T$  can be expressed as  $\sum_{i \in [r]} u_i \otimes v_i \otimes z_i$  where  $u_i \in \mathbb{R}^{n_1}, v_i \in \mathbb{R}^{n_2}, z_i \in \mathbb{R}^{n_3}$ . Given two vectors  $u, v \in \mathbb{R}^n$ , we use  $(u; v)$  to denote the  $n \times 2$  matrix that is obtained with  $u$  and  $v$  as columns.

We now define first, second and third order statistics (frequencies) that serve as our proxies for the first, second and third order moments.

**Definition 2.4.** [Moments] Given a Mallows mixture model, we denote for every  $i, j, k \in [n]$

- $P_i = \Pr(\text{pos}(e_i) = 1)$  is the probability that element  $e_i$  is ranked at the first position
- $P_{ij} = \Pr(\text{pos}(\{e_i, e_j\}) = \{1, 2\})$ , is the probability that  $e_i, e_j$  are ranked at the first two positions (in any order)
- $P_{ijk} = \Pr(\text{pos}(\{e_i, e_j, e_k\}) = \{1, 2, 3\})$  is the probability that  $e_i, e_j, e_k$  are ranked at the first three positions (in any order).

For convenience, let  $P$  represent the set of quantities  $(P_i, P_{ij}, P_{ijk})_{1 \leq i < j < k \leq n}$ . These can be estimated up to any inverse polynomial accuracy using only polynomial samples. The following simple, yet crucial lemma relates  $P$  to the vectors  $x$  and  $y$ , and demonstrates why these statistics and representative vectors are ideal for tensor decomposition.

**Lemma 2.5.** Given a mixture  $w_1 \mathcal{M}(\phi_1, \pi_1) \oplus w_2 \mathcal{M}(\phi_2, \pi_2)$  let  $x, y$  and  $P$  be as defined above.

1. For any  $i$  it holds that  $P_i = w_1 x_i + w_2 y_i$ .
2. Denote  $c_2(\phi) = \frac{Z_n(\phi)}{Z_{n-1}(\phi)} \frac{1+\phi}{\phi}$ . Then for any  $i \neq j$  it holds that  $P_{ij} = w_1 c_2(\phi_1) x_i x_j + w_2 c_2(\phi_2) y_i y_j$ .
3. Denote  $c_3(\phi) = \frac{Z_n^2(\phi)}{Z_{n-1}(\phi) Z_{n-2}(\phi)} \frac{1+2\phi+2\phi^2+\phi^3}{\phi^3}$ . Then for any distinct  $i, j, k$  it holds that  $P_{ijk} = w_1 c_3(\phi_1) x_i x_j x_k + w_2 c_3(\phi_2) y_i y_j y_k$ .

Clearly, if  $i = j$  then  $P_{ij} = 0$ , and if  $i, j, k$  are not all distinct then  $P_{ijk} = 0$ .

In addition, in Lemma ?? we prove the bounds  $c_2(\phi) = O(1/\phi)$  and  $c_3(\phi) = O(\phi^{-3})$ .

**Partitioning Indices:** Given a partition of  $[n]$  into  $S_a, S_b, S_c$ , let  $x^{(a)}, y^{(a)}$  be the representative vectors  $x, y$  restricted to the indices (rows) in  $S_a$  (similarly for  $S_b, S_c$ ). Then the 3-tensor

$$T^{(abc)} \equiv (P_{ijk})_{i \in S_a, j \in S_b, k \in S_c} = w_1 c_3(\phi_1) x^{(a)} \otimes x^{(b)} \otimes x^{(c)} + w_2 c_3(\phi_2) y^{(a)} \otimes y^{(b)} \otimes y^{(c)}.$$

This tensor has a rank-2 decomposition, with one rank-1 term for each Mallows model. Finally for convenience we define the matrix  $M = (x; y)$ , and similarly define the matrices  $M_a = (x^{(a)}; y^{(a)})$ ,  $M_b = (x^{(b)}; y^{(b)})$ ,  $M_c = (x^{(c)}; y^{(c)})$ .

**Error Dependency and Error Polynomials.** Our algorithm gives an estimate of the parameters  $w, \phi$  that we learn in the first stage, and we use these estimates to figure out the entire central rankings in the second stage. The following lemma (proof deferred to the supplementary material) essentially allows us to assume instead of estimations, we have access to the true values of  $w$  and  $\phi$ .

**Lemma 2.6.** For every  $\delta > 0$  there exists a function  $f(n, \phi, \delta)$  s.t. for every  $n, \phi$  and  $\hat{\phi}$  satisfying  $|\phi - \hat{\phi}| < \frac{\delta}{f(n, \phi, \delta)}$  we have that the total-variation distance satisfies  $\|\mathcal{M}(\phi, \pi) - \mathcal{M}(\hat{\phi}, \pi)\|_{\text{TV}} \leq \delta$ .

For the ease of presentation, we do not optimize constants or polynomial factors in all parameters. In our analysis, we show how our algorithm is robust (in a polynomial sense) to errors in various statistics, to prove that we can learn with polynomial samples. However, the simplification when there are no errors (infinite samples) still carries many of the main ideas in the algorithm — this in fact shows the identifiability of the model, which was not known previously.

### 3 Algorithm Overview

---

**Algorithm 1** LEARN MIXTURES OF TWO MALLOWS MODELS, **Input:** a set  $\mathcal{S}$  of  $N$  samples from  $w_1\mathcal{M}(\phi_1, \pi_1) \oplus w_2\mathcal{M}(\phi_2, \pi_2)$ , Accuracy parameters  $\epsilon, \epsilon_2$ .

---

1. Let  $\hat{P}$  be the empirical estimate of  $P$  on samples in  $\mathcal{S}$ .
  2. Repeat  $O(\log n)$  times:
    - (a) Partition  $[n]$  randomly into  $S_a, S_b$  and  $S_c$ . Let  $T^{(abc)} = (\hat{P}_{ijk})_{i \in S_a, j \in S_b, k \in S_c}$ .
    - (b) Run TENSOR-DECOMP from [25, 23] to get a decomposition of  $T^{(abc)} = u^{(a)} \otimes u^{(b)} \otimes u^{(c)} + v^{(a)} \otimes v^{(b)} \otimes v^{(c)}$ .
    - (c) If  $\min\{\sigma_2(u^{(a)}; v^{(a)}), \sigma_2(u^{(b)}; v^{(b)}), \sigma_2(u^{(c)}; v^{(c)})\} > \epsilon_2$   
(In the *non-degenerate* case these matrices are far from being rank-1 matrices in the sense that their least singular value is bounded away from 0.)
      - i. Obtain parameter estimates  $(\hat{w}_1, \hat{w}_2, \hat{\phi}_1, \hat{\phi}_2)$  and prefixes of the central rankings  $\pi_1', \pi_2'$  from INFER-TOP-K( $\hat{P}, M'_a, M'_b, M'_c$ ).
      - ii. Use RECOVER-REST to find the full central rankings  $\hat{\pi}_1, \hat{\pi}_2$ .  
Return SUCCESS and output  $(\hat{w}_1, \hat{w}_2, \hat{\phi}_1, \hat{\phi}_2, \hat{\pi}_1, \hat{\pi}_2)$ .
  3. Run HANDLE DEGENERATE CASES ( $\hat{P}$ ).
- 

Our algorithm (Algorithm 1) has two main components. First we invoke a decomposition algorithm [25, 23] over the tensor  $T^{(abc)}$ , and retrieve approximations of the two Mallows models' representative vectors which in turn allow us to approximate the weight parameters  $w_1, w_2$ , scale parameters  $\phi_1, \phi_2$ , and the top few elements in each central ranking. We then use the inferred parameters to recover the entire rankings  $\pi_1$  and  $\pi_2$ . Should the tensor-decomposition fail, we invoke a special procedure to handle such degenerate cases. Altogether, our algorithm has the following guarantee.

**Theorem 3.1.** *Let  $w_1\mathcal{M}(\phi_1, \pi_1) \oplus w_2\mathcal{M}(\phi_2, \pi_2)$  be a mixture of two Mallows models and let  $w_{\min} = \min\{w_1, w_2\}$  and  $\phi_{\max} = \max\{\phi_1, \phi_2\}$  and similarly  $\phi_{\min} = \min\{\phi_1, \phi_2\}$ . Denote  $\epsilon_0 = \frac{w_{\min}^2(1-\phi_{\max})^{10}}{16n^{22}\phi_{\max}^2}$ . Then, given any  $0 < \epsilon < \epsilon_0$ , suitably small  $\epsilon_2 = \text{poly}(\frac{1}{n}, \epsilon, \phi_{\min}, w_{\min})$  and  $N = \text{poly}\left(n, \frac{1}{\min\{\epsilon, \epsilon_0\}}, \frac{1}{\phi_1(1-\phi_1)}, \frac{1}{\phi_2(1-\phi_2)}, \frac{1}{w_1}, \frac{1}{w_2}\right)$  i.i.d samples from the mixture model, Algorithm 1 recovers, in poly-time and with probability  $\geq 1 - n^{-3}$ , the model's parameters up to  $\epsilon$ -accuracy.*

Next we detail the various subroutines of the algorithm, and give an overview of the analysis for each subroutine. The full analysis is given in the supplementary material.

**The TENSOR-DECOMP Procedure.** This procedure is a straight-forward invocation of the algorithm detailed in [25, 23]. This algorithm uses spectral methods to retrieve the two vectors generating the rank-2 tensor  $T^{(abc)}$ . This technique works when all factor matrices  $M_a = (x^{(a)}; y^{(a)})$ ,  $M_b = (x^{(b)}; y^{(b)})$ ,  $M_c = (x^{(c)}; y^{(c)})$  are well-conditioned. We note that any algorithm that decomposes non-symmetric tensors which have well-conditioned factor matrices, can be used here as a black box.

**Lemma 3.2 (Full rank case).** *In the conditions of Theorem 3.1, suppose our algorithm picks some partition  $S_a, S_b, S_c$  such that the matrices  $M_a, M_b, M_c$  are all well-conditioned — i.e. have  $\sigma_2(M_a), \sigma_2(M_b), \sigma_2(M_c) \geq \epsilon'_2 \geq \text{poly}(\frac{1}{n}, \epsilon, \epsilon_2, w_1, w_2)$  then with high probability, Algorithm TENSORDECOMP of [25] finds  $M'_a = (u^{(a)}; v^{(a)})$ ,  $M'_b = (u^{(b)}; v^{(b)})$ ,  $M'_c = (u^{(c)}; v^{(c)})$  such that for any  $\tau \in \{a, b, c\}$ , we have  $u^{(\tau)} = \alpha_\tau x^{(\tau)} + z_1^{(\tau)}$  and  $v^{(\tau)} = \beta_\tau y^{(\tau)} + z_2^{(\tau)}$ ; with  $\|z_1^{(\tau)}\|, \|z_2^{(\tau)}\| \leq \text{poly}(\frac{1}{n}, \epsilon, \epsilon_2, w_{\min})$  and,  $\sigma_2(M'_\tau) > \epsilon_2$  for  $\tau \in \{a, b, c\}$ .*

**The INFER-TOP-K procedure.** This procedure uses the output of the tensor-decomposition to retrieve the weights,  $\phi$ 's and the representative vectors. In order to convert  $u^{(a)}, u^{(b)}, u^{(c)}$  into an approximation of  $x^{(a)}, x^{(b)}, x^{(c)}$  (and similarly with  $v^{(a)}, v^{(b)}, v^{(c)}$  and  $y^{(a)}, y^{(b)}, y^{(c)}$ ), we need to find a good approximation of the scalars  $\alpha_a, \alpha_b, \alpha_c$ . This is done by solving a certain linear system. This also allows us to estimate  $\hat{w}_1, \hat{w}_2$ . Given our approximation of  $x$ , it is easy to find  $\phi_1$  and the

top first elements of  $\pi_1$  — we sort the coordinates of  $x$ , setting  $\pi'_1$  to be the first elements in the sorted vector, and  $\phi_1$  as the ratio between any two adjacent entries in the sorted vector. We refer the reader to Section ?? in the supplementary material for full details.

**The RECOVER-REST procedure.** The algorithm for recovering the remaining entries of the central permutations (Algorithm 2) is more involved.

---

**Algorithm 2** RECOVER-REST, **Input:** a set  $\mathcal{S}$  of  $N$  samples from  $w_1\mathcal{M}(\phi_1, \pi_1) \oplus w_2\mathcal{M}(\phi_2, \pi_2)$ , parameters  $\hat{w}_1, \hat{w}_2, \hat{\phi}_1, \hat{\phi}_2$  and initial permutations  $\hat{\pi}_1, \hat{\pi}_2$ , and accuracy parameter  $\epsilon$ .

---

1. For elements in  $\hat{\pi}_1$  and  $\hat{\pi}_2$ , compute representative vectors  $\hat{x}$  and  $\hat{y}$  using estimates  $\hat{\phi}_1$  and  $\hat{\phi}_2$ .
  2. Let  $|\hat{\pi}_1| = r_1, |\hat{\pi}_2| = r_2$  and wlog  $r_1 \geq r_2$ .  
If there exists an element  $e_i$  such that  $\text{pos}_{\hat{\pi}_1}(e_i) > r_1$  and  $\text{pos}_{\hat{\pi}_2}(e_i) < r_2/2$  (or in the symmetric case), then:  
Let  $\mathcal{S}_1$  be the subsample with  $e_i$  ranked in the first position.
    - (a) Learn a single Mallows model on  $\mathcal{S}_1$  to find  $\hat{\pi}_1$ . Given  $\hat{\pi}_1$  use dynamic programming to find  $\hat{\pi}_2$
  3. Let  $e_{i^*}$  be the first element in  $\hat{\pi}_1$  such that  $|\hat{x}(e_{i^*}) - \hat{y}(e_{i^*})| > \epsilon$ . Define  $\hat{w}'_1 = \frac{1}{1 + \frac{\hat{w}_2 \hat{y}(e_{i^*})}{\hat{w}_1 \hat{x}(e_{i^*})}}$  and  $\hat{w}'_2 = 1 - \hat{w}'_1$ . Let  $\mathcal{S}_1$  be the subsample with  $e_{i^*}$  ranked at the first position.
  4. For each  $e_i$  that doesn't appear in either  $\hat{\pi}_1$  or  $\hat{\pi}_2$  and any possible position  $j$  it might belong to
    - (a) Use  $\mathcal{S}$  to estimate  $\hat{f}_{i,j} = \Pr(e_i \text{ goes to position } j)$ , and  $\mathcal{S}_1$  to estimate  $\hat{f}_{e_i, j|e_{i^*} \rightarrow 1} = \Pr(e_i \text{ goes to position } j | e_{i^*} \rightarrow 1)$ .
    - (b) Solve the system
 
$$\hat{f}_{i,j} = \hat{w}_1 f^{(1)}(i \rightarrow j) + \hat{w}_2 f^{(2)}(i \rightarrow j) \quad (1)$$

$$\hat{f}_{i, j|e_{i^*} \rightarrow 1} = \hat{w}'_1 f^{(1)}(i \rightarrow j) + \hat{w}'_2 f^{(2)}(i \rightarrow j) \quad (2)$$
  5. To complete  $\hat{\pi}_1$  assign each  $e_i$  to position  $\arg \max_j \{f^{(1)}(i \rightarrow j)\}$ . Similarly complete  $\hat{\pi}_2$  using  $f^{(2)}(i \rightarrow j)$ . Return the two permutations.
- 

Algorithm 2 first attempts to find a pivot — an element  $e_i$  which appears at a fairly high rank in one permutation, yet does not appear in the other prefix  $\hat{\pi}_2$ . Let  $E_{e_i}$  be the event that a permutation ranks  $e_i$  at the first position. As  $e_i$  is a pivot, then  $\Pr_{\mathcal{M}_1}(E_{e_i})$  is noticeable whereas  $\Pr_{\mathcal{M}_2}(E_{e_i})$  is negligible. Hence, conditioning on  $e_i$  appearing at the first position leaves us with a subsample in which all sampled rankings are generated from the first model. This subsample allows us to easily retrieve the rest of  $\pi_1$ . Given  $\pi_1$ , the rest of  $\pi_2$  can be recovered using a dynamic programming procedure. Refer to the supplementary material for details.

The more interesting case is when no such pivot exists, i.e., when the two prefixes of  $\pi_1$  and  $\pi_2$  contain almost the same elements. Yet, since we invoke RECOVER-REST after successfully calling TENSOR-DECOMP, it must hold that the distance between the obtained representative vectors  $\hat{x}$  and  $\hat{y}$  is noticeably large. Hence some element  $e_{i^*}$  satisfies  $|\hat{x}(e_{i^*}) - \hat{y}(e_{i^*})| > \epsilon$ , and we proceed by setting up a linear system.

Let  $f^{(1)}(i \rightarrow j)$  be the probability that element  $e_i$  goes to position  $j$  according to Mallows Model  $\mathcal{M}_1$  (and similarly  $f^{(2)}(i \rightarrow j)$  for model  $\mathcal{M}_2$ ). To find the complete rankings, we measure appropriate statistics to set up a system of linear equations to calculate  $f^{(1)}(i \rightarrow j)$  and  $f^{(2)}(i \rightarrow j)$  up to inverse polynomial accuracy. The largest of these values  $\{f^{(1)}(i \rightarrow j)\}$  corresponds to the position of  $e_i$  in the central ranking of  $\mathcal{M}_1$ .

To compute the values  $\{f^{(r)}(i \rightarrow j)\}_{r=1,2}$  we consider  $f^{(1)}(i \rightarrow j|e_{i^*} \rightarrow 1)$  — the probability that  $e_i$  is ranked at the  $j$ th position conditioned on the element  $e_{i^*}$  ranking first according to  $\mathcal{M}_1$  (and resp. for  $\mathcal{M}_2$ ). Using  $w'_1$  and  $w'_2$  as in Algorithm 2, it holds that

$$\Pr(e_i \rightarrow j | e_{i^*} \rightarrow 1) = w'_1 f^{(1)}(i \rightarrow j | e_{i^*} \rightarrow 1) + w'_2 f^{(2)}(i \rightarrow j | e_{i^*} \rightarrow 1).$$

We need to relate  $f^{(r)}(i \rightarrow j | e_{i^*} \rightarrow 1)$  to  $f^{(r)}(i \rightarrow j)$ . Indeed Lemma ?? shows that  $\Pr(e_i \rightarrow j | e_{i^*} \rightarrow 1)$  is an *almost* linear equations in the two unknowns. We show that if  $e_{i^*}$  is

ranked above  $e_i$  in the central permutation, then for some small  $\delta$  it holds that

$$\Pr(e_i \rightarrow j | e_{i^*} \rightarrow 1) = w'_1 f^{(1)}(i \rightarrow j) + w'_2 f^{(2)}(i \rightarrow j) \pm \delta$$

We refer the reader to Section ?? in the supplementary material for full details.

**The HANDLE-DEGENERATE-CASES procedure.** We call a mixture model  $w_1 \mathcal{M}(\phi_1, \pi_1) \oplus w_2 \mathcal{M}(\phi_2, \pi_2)$  *degenerate* if the parameters of the two Mallows models are equal, and the edit distance between the prefixes of the two central rankings is at most two.<sup>6</sup> We show that unless  $w_1 \mathcal{M}(\phi_1, \pi_1) \oplus w_2 \mathcal{M}(\phi_2, \pi_2)$  is degenerate, then a random partition  $(S_a, S_b, S_c)$  is likely to satisfy the requirements of Lemma 3.2 and result in a successful execution of the TENSOR-DECOMP procedure. Therefore, when TENSOR-DECOMP repeatedly fail, we deduce our model is indeed degenerate. To show this, we use a complete characterization of the uniqueness of decompositions of rank 2, along with some very useful properties of random partitions. In such degenerate cases, we find the two prefixes and then remove the elements in the prefixes from  $U$ , and recurse on the remaining elements. We refer the reader to Section ?? in the supplementary material for full details.

## 4 Experiments

**Goal.** The main contribution of our paper is devising an algorithm that provably learns any mixture of two Mallows models, in comparison to previously existing unproven heuristics. But it could be that such heuristics do well in practice. Our goal is to show that indeed — in certain settings heuristics do not succeed in learning a mixture of Mallows models, as opposed to our algorithm which does succeed.

**Baseline.** In this section we compare our algorithm with a popular EM based algorithm of [13] for learning a mixture of two Mallows model. EM based heuristics are the most popular way to learn a mixture of Mallows models. The EM algorithm starts with a random guess for the two central permutations. At each iteration, the algorithm assigns scores to each data point  $i$ . These scores reflect the probability that point  $i$  was generated from permutation one or two. This is the expectation step. In the maximization step, the algorithm performs a local search around the current central permutations to find new permutations which minimize the average distance from the sampled points weighted according to the assigned scores. We also implemented a simpler version of our algorithm which is designed to handle non-degenerate cases (which is the more interesting case in practice).

**Setting.** We ran both the algorithms on synthetic data comprising of rankings of size up to 10. For every size  $n$ , ranging from 1 to 10, we generated two random rankings at various distances  $d$ , ranging from 0 to  $\binom{n}{2}$ . For each choice of  $n$  and  $d$ , we fix one of the rankings to be  $[1 : n]$ . The other ranking is generated as follows: Any other ranking is uniquely defined by the number of inversions that each element participates in. Hence we choose  $n$  random numbers that sum up to  $d$ . These numbers specify the number of inversions for each element and hence a unique ranking at a distance  $d$  from  $[1 : n]$ . For each such combination, we generated data from a mixture model using the two rankings as central permutations and using random values for mixture weights and  $\phi$  parameters.

**Evaluation Metric.** For each value of  $n$  and  $d$ , we ran both the algorithm 20 times and counted the fraction of instances on which they return the exact rankings in the mixture.

**Results.** The results of the experiment for rankings of size  $n = 10$  are tabulated below. The performance of the EM algorithm is bad if the two rankings are close to each other and it improves as the rankings get far away. On the other hand, our algorithm is able to recover the true rankings even at very close distances. As the rankings get slightly farther, our algorithm recovers the true rankings all the time. Similar performance was observed for other values of  $n$  as well.

<sup>6</sup>I.e., by changing the positions of at most two elements in  $\pi_1$  we retrieve  $\pi_2$ .

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

distance between rankings	success rate of EM	success rate of our algorithm
0	0%	10%
2	0%	10%
4	0%	40%
8	10%	70%
16	30%	60 %
24	30%	100%
30	60%	100%
35	60%	100%
40	80%	100%
45	60%	100%

## References

- [1] C. L. Mallows. Non-null ranking models i. *Biometrika*, 44(1-2), 1957.
- [2] John I. Marden. *Analyzing and Modeling Rank Data*. Chapman & Hall, 1995.
- [3] Guy Lebanon and John Lafferty. Cranking: Combining rankings using conditional probability models on permutations. In *ICML*, 2002.
- [4] Thomas Brendan Murphy and Donal Martin. Mixtures of distance-based models for ranking data. *Computational Statistics and Data Analysis*, 41, 2003.
- [5] Marina Meila, Kapil Phadnis, Arthur Patterson, and Jeff Bilmes. Consensus ranking under the exponential model. Technical report, UAI, 2007.
- [6] Ludwig M. Busse, Peter Orbanz, and Joachim M. Buhmann. Cluster analysis of heterogeneous rank data. In *ICML, ICML '07*, 2007.
- [7] Bhushan Mandhani and Marina Meila. Tractable search for learning exponential models of rankings. *Journal of Machine Learning Research - Proceedings Track*, 5, 2009.
- [8] Tyler Lu and Craig Boutilier. Learning mallows models with pairwise preferences. In *ICML*, 2011.
- [9] Joel Oren, Yuval Filmus, and Craig Boutilier. Efficient vote elicitation under candidate uncertainty. *JCAI*, 2013.
- [10] H Peyton Young. Condorcet’s theory of voting. *The American Political Science Review*, 1988.
- [11] Persi Diaconis. *Group representations in probability and statistics*. Institute of Mathematical Statistics, 1988.
- [12] Mark Braverman and Elchanan Mossel. Sorting from noisy information. *CoRR*, abs/0910.1191, 2009.
- [13] Marina Meila and Harr Chen. Dirichlet process mixtures of generalized mallows models. In *UAI*, 2010.
- [14] Sanjoy Dasgupta. Learning mixtures of gaussians. In *FOCS*, 1999.
- [15] Sanjeev Arora and Ravi Kannan. Learning mixtures of arbitrary gaussians. In *STOC*, 2001.
- [16] Dimitris Achlioptas and Frank McSherry. On spectral learning of mixtures of distributions. In *COLT*, 2005.
- [17] Adam Tauman Kalai, Ankur Moitra, and Gregory Valiant. Efficiently learning mixtures of two gaussians. In *STOC*, STOC ’10, 2010.
- [18] A. Moitra and G. Valiant. Settling the polynomial learnability of mixtures of gaussians. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, 2010.
- [19] Anima Anandkumar, Rong Ge, Daniel Hsu, Sham M. Kakade, and Matus Telgarsky. Tensor decompositions for learning latent variable models. *CoRR*, abs/1210.7559, 2012.
- [20] Animashree Anandkumar, Daniel Hsu, and Sham M. Kakade. A method of moments for mixture models and hidden markov models. In *COLT*, 2012.
- [21] Daniel Hsu and Sham M. Kakade. Learning mixtures of spherical gaussians: moment methods and spectral decompositions. In *ITCS, ITCS ’13*, 2013.
- [22] Santosh Vempala and Grant Wang. A spectral algorithm for learning mixture models. *J. Comput. Syst. Sci.*, 68(4), 2004.
- [23] Aditya Bhaskara, Moses Charikar, Ankur Moitra, and Aravindan Vijayaraghavan. Smoothed analysis of tensor decompositions. In *Symposium on the Theory of Computing (STOC)*, 2014.
- [24] M. G. Kendall. *Biometrika*, 30(1/2), 1938.
- [25] Aditya Bhaskara, Moses Charikar, and Aravindan Vijayaraghavan. Uniqueness of tensor decompositions with applications to polynomial identifiability. *CoRR*, abs/1304.8087, 2013.