
Ranking Distributions based on Noisy Sorting

Adil El Mesaoudi-Paul¹ Eyke Hüllermeier¹ Róbert Busa-Fekete²

Abstract

We propose a new statistical model for ranking data, i.e., a new family of probability distributions on permutations. Our model is inspired by the idea of a data-generating process in the form of a noisy sorting procedure, in which deterministic comparisons between pairs of items are replaced by Bernoulli trials. The probability of producing a certain ranking as a result then essentially depends on the Bernoulli parameters, which can be interpreted as pairwise preferences. We show that our model can be written in closed form if insertion sort is used as sorting algorithm and can be characterized recursively if quick sort is used, and propose a maximum likelihood approach for parameter estimation. We also introduce a generalization of the model, in which the constraints on pairwise preferences are relaxed, and for which maximum likelihood estimation can be carried out based on a variation of the generalized iterative scaling algorithm. Experimentally, we show that the models perform very well in terms of goodness of fit, compared to existing models for ranking data.

1. Introduction

The analysis of ranking data has a long tradition in statistics, and corresponding methods have been used in various fields of application, such as psychology and the social sciences (Marden, 1996). More recently, applications in information retrieval (Liu, 2009) and machine learning (Fürnkranz & Hüllermeier, 2010) have caused a renewed interest in the analysis of rankings and related statistical tools, such as probability distributions on rankings.

In contrast to probability distributions on the reals, the number of parametric distributions on rankings (permutations

of a fixed size) is rather limited. The most popular models are Mallows (Mallows, 1957) and Plackett-Luce (Plackett, 1975; Luce, 1959), and to a lesser extent Babington Smith (Babington-Smith, 1950). In this paper, we add another class of probability distributions to this repertoire.

Our model is inspired by the idea of a data-generating process in the form of a noisy sorting procedure (Biernacki & Jacques, 2013), that is, the idea that a ranking is produced as the result of a sorting process, in which comparisons are not deterministic but dependant on chance. More specifically, comparisons between pairs of items are modelled as Bernoulli trials, with the Bernoulli parameters representing pairwise preferences. While these preferences obey certain consistency constraints for our basic model, we also introduce a generalization for which these constraints are relaxed. For two sorting algorithms, insertion sort and quick sort, we show that the former model can be written in closed form, and that the latter has a recursive characterization.

In addition to proposing the models themselves, we address the problem of parameter estimation based on sample data. More specifically, we devise procedures for efficient maximum likelihood estimation. In an experimental study, we assess the performance of our models in terms of goodness of fit on a large number of real-world data sets.

The rest of the paper is organized as follows. In the next section, we introduce notation and recall the basic families of probability distributions on rankings. Our new model classes are introduced in Section 3, their instantiation for specific sorting algorithms is discussed in Section 4, and the problem of parameter estimation is addressed in Section 5. Experimental results are presented in Section 6, prior to concluding the paper in Section 7.

2. Probability Distributions on Rankings

Consider a fixed set $O = \{o_1, \dots, o_K\}$ of K choice alternatives (objects/options/items). We identify a ranking over O with a permutation $\pi \in \mathbb{S}_K$, where \mathbb{S}_K denotes the collection of permutations on $[K] = \{1, \dots, K\}$. Thus, each π is a mapping $[K] \rightarrow [K]$, such that $\pi(k)$ denotes the position of the k^{th} item o_k in the associated ranking. With each ranking π , we associate an ordering π^{-1} , where $\pi^{-1}(j)$ is the index of the item on position j . To simplify notation, we

¹Heinz Nixdorf Institute and Department of Computer Science, Paderborn University, Germany ²Yahoo Research, New York, USA. Correspondence to: Adil El Mesaoudi-Paul <adil.paul@upb.de>.

shall denote by π both a ranking and the associated ordering, writing the former in brackets and the latter in parentheses. For example, $\pi = [2, 3, 1]$, $\pi = (3, 1, 2)$, as well as the function π defined by $\pi(1) = 2, \pi(2) = 1, \pi(3) = 1$, all denote the ranking in which o_3 is at the top, o_1 in the middle, and o_2 on the last position.

2.1. Mallows Distribution and Extensions

The Mallows model (MM) (Mallows, 1957) belongs to the exponential family of distributions and is parametrized by a reference ranking τ and a dispersion parameter ϕ :

$$\mathbb{P}_{\tau, \phi}(\pi) = \frac{1}{C(\phi)} \exp(-\phi D(\pi, \tau)),$$

where $D(\pi, \tau)$ is the Kendall distance (the number of pairwise inversions between π and τ) and $C(\phi)$ a normalization constant. Thus, Mallows is a distance-based model: the probability of a ranking π decreases with increasing distance from τ , which is the mode of the distribution.

The generalized Mallows model (GMM) (Fligner & Verducci, 1986) is an extension of the MM model, which has $K - 1$ dispersion parameters $\phi_1, \dots, \phi_{K-1}$. Each of the latter affects one specific position in the ranking, thereby allowing permutations at the same distance from the reference ranking to have different probabilities. The probability of a ranking π according to the GMM model is given by

$$\mathbb{P}_{\tau, \phi}(\pi) = \frac{1}{C(\phi)} \exp\left(-\sum_{j=1}^{K-1} \phi_j V_j(\pi)\right),$$

where $V_j(\pi) = \sum_{i>j} \mathbb{1}[\pi^{-1}(i) < \pi^{-1}(j)]$ is the number of inversions for item o_j in π with respect to the identity permutation¹. As such, the GMM model uses an insertion procedure in its generative process, in which a ranking is generated by iteratively inserting elements according to the reference ranking into a list. The probability of inserting an element into a specific position is controlled by the inversion distance and the ϕ -parameters.

Meek & Meila (2014) further extend the GMM to the recursive inversion model (RIM), which is able to capture a hierarchical structure on the items. Instead of inserting single items, complete subsequences are merged in a recursive manner, preserving the order within each subsequence. The model is specified by a binary recursive decomposition of the items represented by a structure τ , and the number of inversions is controlled by a parameter θ_i associated with each merge operation. By representing a RIM as a binary tree, where the leaves correspond to the items and the internal vertices \mathcal{I} to the parameters θ_i , the probability of a

ranking $\tau(\theta)$ becomes proportional to

$$\prod_{i \in \mathcal{I}} \exp(-\theta_i v_i(\pi, \pi_\tau)),$$

where $v_i(\pi, \pi_\tau)$ is the number of inversions at vertex i of $\tau(\theta)$ for the ranking π .

2.2. Plackett-Luce Distribution

The Plackett-Luce (PL) model (Plackett, 1975; Luce, 1959) is parametrized by a vector $\theta = (\theta_1, \theta_2, \dots, \theta_K) \in \mathbb{R}_+^K$. Each θ_i can be interpreted as the weight or “strength” of the option o_i . The probability assigned by the PL model to a ranking represented by a permutation $\pi \in \mathbb{S}_K$ is given by

$$\mathbb{P}_\theta(\pi) = \prod_{i=1}^K \frac{\theta_{\pi^{-1}(i)}}{\theta_{\pi^{-1}(i)} + \theta_{\pi^{-1}(i+1)} + \dots + \theta_{\pi^{-1}(K)}} \quad (1)$$

The product on the right-hand side of (1) is the probability of producing the ranking π in a *stagewise* process: First, the item on the first position is selected, then the item on the second position, and so forth. In each step, the probability of an item to be chosen next is proportional to its weight. Consequently, items with a higher weight tend to occupy higher positions. In particular, the most probable ranking (i.e., the mode of the PL distribution) is simply obtained by sorting the items in decreasing order of their weight:

$$\tau = \operatorname{argmax}_{\pi \in \mathbb{S}_K} \mathbb{P}_\theta(\pi) = \operatorname{argsort}\{\theta_1, \dots, \theta_K\} \quad (2)$$

2.3. Babington Smith Distribution

The Babington Smith (BS) model is defined as follows (Babington-Smith, 1950):

$$\mathbb{P}_\theta(\pi) = \frac{1}{C(\theta)} \prod_{1 \leq i < j \leq K} p_{\pi^{-1}(i), \pi^{-1}(j)}, \quad (3)$$

where $p_{i,j}$ is the probability to observe a preference $o_i \succ o_j$ in a direct comparison between o_i and o_j , and $C(\theta)$ is a normalization constant. Thus, the parametrization θ of the BS model consists of all pairwise probabilities $p_{i,j} = 1 - p_{j,i}$, $1 \leq i < j \leq K$.

The BS distribution results from the following “trial and error” data-generating process: First, the order of each pair of objects o_i and o_j is determined independently at random (as a result of a Bernoulli trial, i.e., by flipping a coin with bias $p_{i,j}$). Then, in case all pairwise comparisons form a consistent ranking, this ranking is adopted, otherwise the first step is repeated.

2.4. Comparison

The previous models can naturally be distinguished in terms of their parametrization. The Mallows model is quite restricted and not very flexible. It has one degree of freedom

¹ $\mathbb{1}[\cdot]$ maps true predicates to 1 and false predicates to 0.

to determine the location of the distribution (the reference ranking), and another parameter to determine the spread (comparable, for example, to the normal distribution on the reals). The PL model is more flexible (for example, see (Cheng et al., 2012) for a comparison of the expressivity of Mallows and PL), with a number of parameters that is linear in the number of items. BS has an even richer parametrization, the size of which grows quadratically with the number of items.

From a preference modeling point of view, the parametrizations of PL and BS are both quite natural: PL specifies the strength of each option individually, whereas BS takes pairwise comparisons as a point of departure. Thus, while PL implies relatively strong consistency properties, such as strong stochastic transitivity, BS principally allows for preferential cycles.

3. Ranking Distributions based on Sorting

PL and BS can both be interpreted in terms of an underlying data-generating process, in which a ranking is produced as the result of a specific stochastic process. However, especially in the case of BS, the “cognitive plausibility” of the process is questionable: It is difficult to imagine that a ranking of items is indeed produced by repeating the full set of stochastic pairwise comparisons, independently of each other, till reaching consistency (especially since most of such repetitions will be futile).

As an arguably more plausible assumption, one could imagine that a ranking is the result of a (noisy) sorting procedure. Indeed, when people produce a ranking, they often apply some kind of sorting process, in which items are compared only if necessary. This idea has recently been put forward by Biernacki and Jacques (Biernacki & Jacques, 2013), and provides the main point of departure for our contribution.

A sorting algorithm puts objects stored in a list in a certain order, based on pairwise comparisons between these objects. Most often, the objects to be sorted are numbers, and the pairwise comparison is determined based on some binary relation, for example, the \leq relation for increasing and \geq relation for decreasing order. Note that the list submitted as input to a (deterministic) sorting algorithm does not affect its output, but it does have an influence on its time complexity, and on the pairs of items that are compared. Therefore, the time complexity of sorting algorithms is often analyzed under the assumption of a uniform distribution over the possible inputs (average time complexity analysis).

3.1. Insertion Sort Rank Data Model

The model by Biernacki and Jacques (Biernacki & Jacques, 2013), called Insertion Sort Rank data (ISR) model, is specified by a reference ranking τ and real parameter p , very

much like the Mallows model. The former corresponds to the “correct” ranking, i.e., the mode of the distribution, and $p \in [0.5, 1]$ is the noise parameter that controls the peakedness of the distribution. More specifically, the following assumption is made: A sorting algorithm (insertion sort) is run on an initial ordering π , and whenever two items o_i and o_j are compared, the “right” outcome (consistent with τ) is produced with probability p (hence the “wrong” outcome with probability $1 - p$).

The algorithm’s probability to terminate with a ranking σ obviously depends on the initial ordering π , which is a latent variable of the model. To get rid of this influence, the initialization is “averaged out”, i.e., an expectation is taken over all initial rankings. Assuming a uniform distribution for π , we thus obtain

$$\mathbb{P}(\sigma | \tau, p) = \frac{1}{K!} \sum_{\pi \in \mathbb{S}_K} \mathbb{P}(\sigma | \pi, \tau, p) . \quad (4)$$

This model can also be written as follows:

$$\begin{aligned} \mathbb{P}(\sigma | \mathbf{P}) &= \frac{1}{C'(\mathbf{P})} \sum_{\pi \in \mathbb{S}_K} \mathbb{P}(\sigma | \pi, \mathbf{P}) , \\ &= \frac{1}{C'(\mathbf{P})} \sum_{\pi \in \mathbb{S}_K} \prod_{i=1}^K \prod_{j \neq i} p_{i,j}^{d_{i,j}^{\sigma, \pi}} , \end{aligned} \quad (5)$$

where \mathbf{P} is a $K \times K$ matrix $\mathbf{P} = [p_{i,j}]_{1 \leq i, j \leq K}$ with entries

$$p_{i,j} = p_{i,j}(\tau, p) = \begin{cases} p & \text{if } \tau(o_i) < \tau(o_j) \\ 1 - p & \text{if } \tau(o_i) > \tau(o_j) \end{cases} . \quad (6)$$

That is, the matrix \mathbf{P} is uniquely determined by τ and p (and vice versa). Moreover, for rankings $\sigma, \pi \in \mathbb{S}_K$,

$$\mathbf{D}^{\sigma, \pi} = [d_{i,j}^{\sigma, \pi}]_{1 \leq i, j \leq K}$$

is a binary matrix with entries $d_{i,j}^{\sigma, \pi} = 1$ if the sorting algorithm, given π as initial ordering and producing σ as output, has compared o_i to o_j with a win for o_i , and to 0 otherwise. Finally,

$$C'(\mathbf{P}) = \sum_{\sigma \in \mathbb{S}_K} \sum_{\pi \in \mathbb{S}_K} \prod_{i=1}^K \prod_{j \neq i} p_{i,j}^{d_{i,j}^{\sigma, \pi}}$$

is the normalization constant.

Biernacki and Jacques tackle the problem of estimating the parameters of the model, τ and p , using the maximum likelihood principle. To this end, they adopt a latent variable interpretation of the model and propose an EM algorithm.

3.2. The Conjunctive Noisy Sorting Model

Our model is a modification of (5), which looks very similar at first sight: Instead of averaging out the influence of the

initial ranking π in an additive way, by aggregating the probabilities $\mathbb{P}(\sigma | \pi, \tau, p)$ with an arithmetic mean, we apply the product as an aggregation function:

$$\mathbb{P}_{\mathcal{A}}(\sigma | \tau, p) \propto \prod_{\pi \in \mathbb{S}_K} \mathbb{P}_{\mathcal{A}}(\sigma | \pi, \tau, p),$$

where \mathcal{A} is the underlying sorting algorithm. As for the latter, one may of course consider algorithms other than insertion sort. Indeed, any pairwise-comparison-based sorting algorithm can in principle be extended to a noisy sorting model by using stochastic pairwise comparisons (Braverman & Mossel, 2008; 2009). In Section 4, we will instantiate our model for two algorithms, insertion sort and quick sort.

There are different motivations for the above modification. First, as will be seen, the multiplicative variant has appealing mathematical properties and can be handled a bit more easily. Second, the model can also be motivated intuitively. The product is a *conjunctive* aggregation function (Grabisch et al., 2009), and combining probabilities in a conjunctive way is in agreement with standard (deterministic) sorting, where the “correct” output ordering σ is obtained regardless of the initial ordering π , that is, as a result for *all* initial orderings π . Therefore, we call our model the Conjunctive Noisy Sorting (CNS) model.

Recall the definition of the matrix \mathbf{P} with entries (6), which is in one-to-one relationship with the model parameters τ and p . With this notation, the CNS model can also be written as follows:

$$\begin{aligned} \mathbb{P}_{\mathcal{A}}(\sigma | \mathbf{P}) &= \frac{1}{C(\mathbf{P})} \prod_{\pi \in \mathbb{S}_K} \mathbb{P}(\sigma | \pi, \mathbf{P}), \\ &= \frac{1}{C(\mathbf{P})} \prod_{\pi \in \mathbb{S}_K} \prod_{i=1}^K \prod_{j \neq i} p_{i,j}^{d_{i,j}^{\sigma, \pi}}, \end{aligned} \quad (7)$$

where $C(\mathbf{P}) = \sum_{\sigma \in \mathbb{S}_K} \prod_{\pi \in \mathbb{S}_K} \prod_{i=1}^K \prod_{j \neq i} p_{i,j}^{d_{i,j}^{\sigma, \pi}}$. To simplify this representation, we introduce

$$\mathbf{D}^{\sigma} = [d_{i,j}^{\sigma}]_{1 \leq i, j \leq K}, \quad d_{i,j}^{\sigma} = \sum_{\pi \in \mathbb{S}_K} d_{i,j}^{\sigma, \pi}.$$

With this notation, the model becomes

$$\mathbb{P}_{\mathcal{A}}(\sigma | \mathbf{P}) = \frac{1}{C(\mathbf{P})} \prod_{i=1}^K \prod_{j \neq i} p_{i,j}^{d_{i,j}^{\sigma}}, \quad (8)$$

where

$$C(\mathbf{P}) = \sum_{\sigma \in \mathbb{S}_K} \prod_{i=1}^K \prod_{j \neq i} p_{i,j}^{d_{i,j}^{\sigma}}. \quad (9)$$

In Section 4, explicit expressions for the exponents $d_{i,j}^{\sigma}$ will be provided for two instantiations of the model (insertion sort and quick sort). Note that, since $p_{i,j} = p$ or $p_{i,j} = 1-p$

in (9), the normalization constant can be written as a power series in p :

$$C(p) = \sum_{j=0}^{e(K)} \alpha_j^{(K)} \cdot p^j$$

The degree $e(K)$ and the coefficients $\alpha_j^{(K)}$ are specific to K but can be precomputed.

3.3. The Generalized Conjunctive Noisy Sorting Model

CNS is a relatively simple model, comparable to Mallows and ISR in terms of its parametrization. The distribution has a single mode at τ , and all pairwise preferences are consistent with this reference. Here, we consider a more general model, which subsumes the CNS model as a special case, and in which these assumptions are relaxed. More specifically, like in the BS model, pairwise preferences $p_{i,j}$ are allowed to be defined independently for each pair of objects o_i and o_j , and are not assumed to obey any consistency conditions. Thus, we assume a noisy sorting procedure in which, whenever the comparison of objects o_i and o_j is required, a coin with success probability $p_{i,j}$ is flipped, and the outcome of this Bernoulli experiment determines the order of the two elements: o_i is preferred to o_j if the outcome is 1, and o_j is preferred to o_i otherwise. We furthermore assume that all pairwise comparisons are independent of each other, and that $p_{i,j} = 1 - p_{j,i}$ for all $i, j \in [K]$. We summarize the probabilities $p_{i,j}$ in the matrix $\mathbf{P} \in [0, 1]^{K \times K}$, which constitutes the parametrization of the model, referred to as Generalized Conjunctive Noisy Sorting (GCNS) model.

Together with a sorting algorithm \mathcal{A} and an initial ordering π , GCNS defines a distribution $\mathbb{P}_{\mathcal{A}}(\cdot | \pi, \mathbf{P})$ over \mathbb{S}_K . Thus, for each ranking $\sigma \in \mathbb{S}_K$, $\mathbb{P}_{\mathcal{A}}(\sigma | \pi, \mathbf{P})$ is the probability to end up with σ when applying \mathcal{A} to the input π , and comparing items o_i and o_j according to $p_{i,j}$. Again, we eliminate the latent variable π via conjunctive aggregation:

$$\mathbb{P}_{\mathcal{A}}(\sigma | \mathbf{P}) \propto \prod_{\pi \in \mathbb{S}_K} \mathbb{P}_{\mathcal{A}}(\sigma | \pi, \mathbf{P})$$

To obtain a more compact representation, we introduce binary matrices $\mathbf{D}^{\sigma, \pi} = [d_{i,j}^{\sigma, \pi}]_{1 \leq i, j \leq K}$ for rankings $\sigma, \pi \in \mathbb{S}_K$, where the entry $d_{i,j}^{\sigma, \pi}$ in $\mathbf{D}^{\sigma, \pi}$ is set to 1 if the sorting algorithm \mathcal{A} , given π as initial ordering and producing σ as output, has compared o_i to o_j with a win for o_i , and to 0 otherwise, and the matrices

$$\mathbf{D}^{\sigma} = [d_{i,j}^{\sigma}]_{1 \leq i, j \leq K}, \quad d_{i,j}^{\sigma} = \sum_{\pi \in \mathbb{S}_K} d_{i,j}^{\sigma, \pi}. \quad (10)$$

We shall consider only such sorting algorithms for which all these matrices are well-defined (which means that, given σ and π , it is clear whether and how o_i and o_j have been compared); this includes insertion sort and quick sort, amongst

others. Using this notation, the GCNS model can be written as follows:

$$\mathbb{P}_{\mathcal{A}}(\sigma | \mathbf{P}) = \frac{1}{C(\mathbf{P})} \prod_{i=1}^K \prod_{j \neq i} p_{i,j}^{d_{i,j}^{\sigma}}, \quad (11)$$

where

$$C(\mathbf{P}) = \sum_{\sigma \in \mathbb{S}_K} \prod_{i=1}^K \prod_{j \neq i} p_{i,j}^{d_{i,j}^{\sigma}}. \quad (12)$$

Based on (11), one can see that GCNS is a special case of the log-linear model over the symmetric group, because the log of the probabilities can be written as a linear function of the logarithm of the parameters. Note that the key quantity in the model is \mathbf{D}^{σ} , which we shall compute in a closed form when insertion sort is used as sorting algorithm, and characterize recursively when quick sort is used.

Extreme probabilities $p_{i,j} \in \{0, 1\}$ may cause problems in the case of inconsistencies, such as preferential cycles $p_{1,2} = p_{2,3} = p_{3,1} = 1$, which are not excluded in our general model. Applying a sorting algorithm \mathcal{A} to some $\mathbf{P} \in \{0, 1\}^{K \times K}$, an initial ordering π will be turned into an ordering σ with probability 1, i.e., $\mathbb{P}_{\mathcal{A}}(\sigma | \pi, \mathbf{P}) = 1$ and $\mathbb{P}_{\mathcal{A}}(\sigma' | \pi, \mathbf{P}) = 0$ for all $\sigma' \neq \sigma$. Then, unless the same σ is produced for all initial orderings π , which is unlikely in the case of inconsistencies, the product $\prod_{\pi} \mathbb{P}_{\mathcal{A}}(\sigma | \pi, \mathbf{P})$ will vanish for all σ , which means that (11) is no longer well-defined. Therefore, we subsequently exclude extreme probabilities and assume $0 < p_{i,j} < 1$ for all $i, j \in [K]$.

Observation 1. *Assuming that $p_{i,j} > 0$ for all $i, j \in [K]$, the model (11) is well-defined in the sense that $C(\mathbf{P}) > 0$; moreover, $\mathbb{P}_{\mathcal{A}}(\sigma | \mathbf{P}) > 0$ for all $\sigma \in \mathbb{S}_K$.*

3.4. Connection to BS

Our model has an interesting connection to BS. The latter is parametrized by the same probability matrix \mathbf{P} , specifying probabilities $p_{i,j} = 1 - p_{j,i}$ for each pair of objects o_i, o_j . Moreover, with a normalizing constant $C''(\mathbf{P})$, it can be written as follows:

$$\mathbb{P}(\sigma | \mathbf{P}) = \frac{1}{C''(\mathbf{P})} \prod_{i=1}^K \prod_{j \neq i} p_{i,j}^{d_{i,j}^{\sigma}},$$

where $d_{i,j}^{\sigma} = 1$ if $\sigma(i) < \sigma(j)$ and $= 0$ otherwise. Comparing this expression with (11), it can be seen that BS has exactly the same structure as our model. The key difference concerns the values $d_{i,j}^{\sigma}$, which can be seen as weights specifying the importance of the comparison between o_i and o_j . In BS, $d_{i,j}^{\sigma} \in \{0, 1\}$ and $d_{i,j}^{\sigma} + d_{j,i}^{\sigma} \equiv 1$, which means that each pair has the same importance. In our model, where a pair (o_i, o_j) can be more or less relevant when producing a ranking σ with a sorting algorithms \mathcal{A} , more general (integer) values are possible.

Interestingly, if the BS model is restricted such that $p_{i,j} = p$ if $\tau(i) < \tau(j)$ and $p_{i,j} = 1 - p$ if $\tau(i) > \tau(j)$, for a fixed $\tau \in \mathbb{S}_K$ and probability p , it reduces to the Mallows model (Mallows, 1957). For exactly the same restriction, GCNS reduces to CNS. Roughly speaking, Mallows is to BS what CNS is to GCNS. Moreover, since GCNS can be seen as a ‘‘sorting variant’’ of BS, CNS can also be seen as a ‘‘sorting variant’’ of Mallows. This is another strong motivation of our model.

4. Instantiations of the Ranking Model

To make the definition of our models complete, we make use of two sorting algorithms \mathcal{A} : insertion sort, denoted by \mathcal{I} , and quick sort, denoted by \mathcal{Q} . For insertion sort algorithm, we show that (8) and (11) can be written in closed form, and for quick sort algorithm, we show that they can be characterized in a recursive way.

4.1. Insertion Sort

In (stochastic) insertion sort, we start with an empty ordering, in which all K objects are inserted one by one, in the order determined by the initial ranking π . In the l^{th} iteration, we are given a partial ordering $(o_{i(1)}, \dots, o_{i(l)})$ of $l < K$ objects and insert another object o . To this end, o is first compared with $o_{i(1)}$, then with $o_{i(2)}$, and so forth. It is inserted as position j if $o_{i(j)}$ is the first item that loses its comparison with o ; in case o is beaten by all l items, it is put on position $l + 1$.

Thus, in stochastic insertion sort, we produce an output ranking σ from an initial ranking π by comparing only a subset of all possible pairs of items. Note that the output of the noisy sorting procedure is a random ordering that depends on the success probabilities $\mathbf{P} = [p_{i,j}]$, and also on the initial ordering π , i.e., the order in which items are inserted. The following example elaborates on this dependence.

Example 1. *Consider insertion sort with two different initial orderings $\pi = (o_1, o_2, o_3)$ and $\pi' = (o_3, o_2, o_1)$. Let the pairwise probabilities be $p_{1,2} = 1/4$, and $p_{1,3} = p_{2,3} = 1/2$. Now let us compute the probability of observing $\sigma = (o_1, o_2, o_3)$. Starting from π , we first insert o_1 , then o_2 , and finally o_3 . Ending with σ is thus only possible if o_1 has beaten o_2 in the first comparison, o_1 has also beaten o_3 , and o_2 has beaten o_3 . Therefore, the probability of observing σ is proportional to $p_{1,2}p_{1,3}p_{2,3} = 0.0625$. Starting from π' , σ is produced with fewer comparisons, and the same probability is proportional to $p_{1,2}p_{2,3} = 0.125$.*

Now, we are going to focus on $\mathbf{D}^{\sigma} = \sum_{\pi \in \mathbb{S}_K} \mathbf{D}^{\sigma, \pi}$, where the sum is elementwise. The following observation allows us to compute \mathbf{D}^{σ} in a concise way.

Lemma 1. *Assume that $\sigma_{id} = (o_1, \dots, o_K)$. Then, for*

insertion sort, the matrix $D^{\sigma_{id}}$ is given by

$$d_{i,j}^{\sigma_{id}} = \begin{cases} \frac{K!}{2} & \text{if } i < j \\ \binom{K}{b_{i,j}+2} (K - b_{i,j} - 2)! b_{i,j}! & \text{if } j < i \\ 0 & \text{otherwise} \end{cases},$$

where $b_{i,j} = i - j - 1$. Furthermore, for any $\sigma \in \mathbb{S}_K$, we have $\mathbf{D}^\sigma = \mathbf{B}^\sigma \mathbf{D}^{\sigma_{id}} \mathbf{B}^{\sigma^T}$, where \mathbf{B}^σ is the permutation matrix that corresponds to σ .

Proof. The first case is easy to verify, since insertion sort with an initial ordering π compares two objects only in case they are concordant (in the same order) in σ_{id} and π . The number of such orderings is $\frac{K!}{2}$.

The second case is more involved, since one needs to calculate all initial orders $\pi \in \mathbb{S}_K$ in which a pair of objects, say o_i and o_j , are discordant and compared to each other in the course of the sorting procedure. Assume that insertion sort is run with π as initial order, and the output is σ_{id} . It is easy to see that object o_i and o_j are not compared if there is a third object o_k , which is between o_j and o_i with respect to σ_{id} , and also between o_j and o_i in π ; Figure 1 illustrates such a configuration of objects. Therefore, the number of orderings for which o_i and o_j are compared is equal to $\binom{K}{b_{i,j}+2} (K - b_{i,j} - 2)!$, because o_i and o_j and the items between them have to be ordered such that o_i is the first, o_j is the second, and all items between o_j and o_i with respect to σ_{id} precede them in π . In addition, the items between o_i and o_j can be permuted arbitrarily in the initial order, which results in the term $b_{i,j}!$.

The last claim can be verified based on the fact that the argument above holds for an arbitrary permutation of objects. This concludes the proof. \square

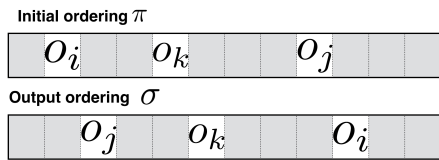


Figure 1. An initial ordering π and output ordering for which insertion sort does not compare o_i and o_j .

4.2. Quick Sort

The quick sort algorithm is inherently random due to the random choice of the pivot item. We make use of a derandomized version by taking as pivot the item in the middle of the initial ordering (i.e., the item on position $\lceil K/2 \rceil$ for an ordering with K items).

In noisy quick sort, we start by picking a pivot element o_p from the elements $\{o_1, \dots, o_K\}$ to be ordered. We then

proceed with the partition operation, in which we construct two sub-orderings, one collecting the items that lost and the other one the items that won the pairwise comparison with the pivot element. The same process is repeated with each of the two sub-orderings (unless a sub-ordering reduces to a single item), eventually producing a complete ordering of the items. Again, we note that the output of the noisy sorting model based on quick sort depends on the pairwise probabilities \mathbf{P} and the initial ordering π . This dependence is illustrated in the following example.

Example 2. Consider the quick sort algorithm with two different initial orderings $\pi = (o_1, o_2, o_3)$ and $\pi' = (o_2, o_1, o_3)$. Let the pairwise probabilities be given as in Example 1. It is easy to see that the probability of observing $\sigma = (o_1, o_2, o_3)$ starting from π is proportional to $p_{1,2}p_{2,3} = 0.125$. When starting with π' , this probability is proportional to $p_{1,2}p_{1,3}p_{2,3} = 0.0625$.

The following lemma gives a recursive expression of \mathbf{D}^σ for the case of quick sort as a sorting algorithm.

Lemma 2. Assume that $\sigma_{id} = (o_1, \dots, o_K)$. Then, for quick sort, the matrix $D^{\sigma_{id}}$ is given by

$$d_{i,j}^{\sigma_{id}} = (K-1)! \left[\sum_{k=1}^i Q(k, K, i, j) + \sum_{k=j}^K Q(1, k, i, j) \right],$$

where

$$Q(\ell, u, i, j) = \begin{cases} 0 & \text{if } i < p = \lceil \frac{i+j}{2} \rceil < j \\ 2 & \text{if } p = i \text{ or } p = j \\ Q(1, p-1, i, j) & \text{if } j < p \\ Q(p+1, u, i, j) & \text{if } i > p \end{cases}$$

Furthermore, for any $\sigma \in \mathbb{S}_K$, we have $\mathbf{D}^\sigma = \mathbf{B}^\sigma \mathbf{D}^{\sigma_{id}} \mathbf{B}^{\sigma^T}$, where \mathbf{B}^σ is the permutation matrix that corresponds to σ .

Proof. The first case of the recursion $Q(\ell, u, i, j)$ follows from the fact that neither o_i nor o_j is chosen as pivot, in which case they will not be compared any more. In the second case, either o_i or o_j is chosen as pivot, in which case they will be compared. Otherwise, the recursion corresponds to the quick sort recursion. \square

5. Parameter Estimation

In this section, we address the problem of parameter estimation, i.e., the question of how to fit our models to a given sample $\mathcal{D} = \{\sigma_1, \dots, \sigma_n\} \subset \mathbb{S}_K$ using the principle of maximum likelihood (ML) estimation.

5.1. The CNS Model

Given a set of observations $\{\sigma_1, \dots, \sigma_n\}$, the ML estimation consists of solving the following constrained optimiza-

tion problem:

$$\begin{aligned} \max_{\mathbf{P}^\tau} \quad & \sum_{\ell=1}^n \sum_{i=1}^K \sum_{j \neq i} d_{i,j}^{\sigma_\ell, \tau} \log p_{i,j}^\tau - n \log C(\mathbf{P}^\tau) \quad (13) \\ \text{s. t.} \quad & p_{i,j}^\tau = \begin{cases} p & \text{if } \tau(i) < \tau(j) \\ 1-p & \text{if } \tau(i) > \tau(j) \end{cases} \quad \forall i, j \in [K], i \neq j \end{aligned}$$

Recall that \mathbf{P}^τ is equivalently represented by the reference order τ and the probability p , i.e., the maximization in the above problem is over these two parameters.

We tackle the problem with simple hill-climbing search for τ in the discrete space \mathbb{S}_K , initialized with the Borda ranking (i.e., sorting items according to their average rank in the data). The neighborhood of an ordering is defined as the set of all orderings that can be obtained by a swap of two adjacent items. For a fixed τ , the optimization problem (13) reduces to a simple one-dimensional problem:

$$\begin{aligned} \max_p \quad & \sum_{\ell=1}^n \left[\sum_{\tau(i) < \tau(j)} d_{i,j}^{\sigma_\ell, \tau} \log p + \sum_{\tau(i) > \tau(j)} d_{i,j}^{\sigma_\ell, \tau} \log(1-p) \right] \\ & - n \log C(\mathbf{P}^\tau) \\ \text{s. t.} \quad & p \in [0.5, 1] \quad (14) \end{aligned}$$

This problem is convex (the distribution belongs to the exponential family) and can be solved numerically, for example by means of the golden section method.

In each iteration of the algorithm, the best candidate solution (τ, p) in the neighborhood of the current best solution is adopted, and the search stops if no improvement is possible anymore.

5.2. The GCNS Model

The GCNS model (11) is parametrized by \mathbf{P} . Here, the maximum likelihood (ML) principle cannot be applied directly, because the normalizing factor $C(\mathbf{P})$ in (12) cannot be written in a closed form in terms of the model parameters. Therefore, we opt for using the generalized iterative scaling (GIS) procedure (Darroch & Ratcliff, 1972), an iterative method for estimating the probabilities in a log-linear model. Given a set of observations $\{\sigma_1, \dots, \sigma_n\}$, ML estimation amounts to solving the following constrained optimization problem:

$$\begin{aligned} \max_{\mathbf{P}} \quad & \sum_{\ell=1}^n \sum_{i=1}^K \sum_{j \neq i} d_{i,j}^{\sigma_\ell} \log p_{i,j} - n \log C(\mathbf{P}) \quad (15) \\ \text{s. t.} \quad & p_{i,j} + p_{j,i} = 1, \quad \forall i, j \in [K], i \neq j. \end{aligned}$$

Let $\sigma_{\downarrow j}$ denote the j^{th} ranking according to some fixed ordering over \mathbb{S}_K (e.g. Lehmer code). With $f_j = \#\{i \in [n] : \sigma_i = \sigma_{\downarrow j}\}$, the empirical frequencies corresponding

to the probabilities of all possible permutations, the GIS procedure seeks to find a parameter estimate \mathbf{P}' for which

$$\sum_{\ell=1}^{K!} p'_\ell d_{i,j}^{\sigma_{\downarrow \ell}} = \sum_{\ell=1}^{K!} \widehat{p}_\ell d_{i,j}^{\sigma_{\downarrow \ell}} \quad (16)$$

for all $i \neq j$, where $p'_\ell = \mathbb{P}_{\mathcal{A}}(\sigma_{\downarrow \ell} | \mathbf{P}')$.

Observe that the GIS procedure can be adapted to produce the parameters of the log-linear model instead of the probabilities $p_{i,j}$ (Malouf, 2002). In addition, we note that GIS requires the computation of a vector of length $K!$, a very costly operation that will be tackled based on a Monte Carlo-based approximation technique. Further, based on Lemma 1 and 2, it is easy to see that the sum of exponents is constant for every ordering in case of both insertion and quick sort, that is $\sum_{i=1}^K \sum_{i \neq j} d_{i,j}^\sigma = B_K$ for all $\sigma \in \mathbb{S}_K$.

According to (Darroch & Ratcliff, 1972), \mathbf{P}' in (16) is the (unconstrained) ML estimate for \mathbf{P} . In our case, however, the constraints $p_{i,j} + p_{j,i} = 1$ in (15) need to be taken into account. Therefore, we accompany each update step in the GIS procedure with a projection step, which ensures that the estimated parameters satisfy the constraints. One update step of the iterative procedure for the parameter estimation thus can be written as

$$p_{i,j}^{(n+1)} = \Pi \left(p_{i,j}^{(n)} + \delta^{(n)} \right),$$

where

$$\delta^{(n)} = \log \left(\frac{\sum_{\ell=1}^{K!} \widehat{p}_\ell d_{i,j}^{\sigma_{\downarrow \ell}}}{\sum_{\ell=1}^{K!} p_{i,j}^{(n)} d_{i,j}^{\sigma_{\downarrow \ell}}} \right)^{\frac{1}{B_K}},$$

and $\Pi(x)$ denotes the least-squares projection of $x = (x_{i,j}, x_{j,i})$, given by

$$\underset{y \in \mathbb{R}_+^2}{\text{argmin}} \quad \|x - y\|^2 \quad \text{s. t.} \quad y_{i,j} + y_{j,i} = 1, \quad (17)$$

which can be determined analytically.

5.3. Sampling

The model (11) can be sampled by using MCMC based on the fact that one can compute the acceptance ratio as

$$\log \frac{\mathbb{P}_{\mathcal{A}}(\sigma | \mathbf{P})}{\mathbb{P}_{\mathcal{A}}(\sigma' | \mathbf{P})} = \sum_{i=1}^K \sum_{j \neq i} (d_{i,j}^\sigma - d_{i,j}^{\sigma'}) \log p_{i,j}.$$

This allows us to make use of the Metropolis-Hastings (MH) algorithm. We use Mallows (Mallows, 1957) as proposal distribution. The pseudo-code of the sampling is given in Algorithm 1. The reference ranking of the Mallows model $\mathbb{P}(\cdot | \phi, \sigma)$, denoted by σ , is always set to the current ranking σ_{i-1} (see line 5). In this case, it is easy to verify that the stationary distribution of the Markov chain is indeed

$\mathbb{P}_{\mathcal{A}}(\sigma | \mathbf{P})$, because the Mallows model is symmetric in the sense that $\mathbb{P}(\sigma | \phi, \sigma') = \mathbb{P}(\sigma' | \phi, \sigma)$, and assigns positive probability to every ranking when $\phi > 0$. Therefore, the detailed balance condition is satisfied, and the ergodicity of the chain is also ensured.

Algorithm 1 Metropolis-Hastings with Mallows proposal

```

1: procedure MH( $T, \phi$ )
2:   Select initial ordering  $\sigma_0$ 
3:    $\mathcal{D} = \emptyset$ 
4:   for  $i = 1 \rightarrow T$  do
5:      $\sigma_i \sim \mathbb{P}(\cdot | \phi, \sigma_{i-1}) \triangleright$  Proposal from Mallows
6:      $q_i \leftarrow \sum_{i=1}^K \sum_{j \neq i}^K (d_{i,j}^{\sigma_i} - d_{i,j}^{\sigma_{i-1}}) \log p_{i,j}$ 
7:     Accept  $\sigma_i$  with probability  $\min(1, \exp(q_i))$ 
8:      $\mathcal{D} = \mathcal{D} \cup \{\sigma_i\}$ 
9:   return  $\mathcal{D}$ 
    
```

6. Experiments

To investigate the performance of our new model and the effectiveness of parameter estimation, we conducted experiments on 213 real-world data sets from the PrefLib repository (<http://www.preflib.org>). These data sets originate from different domains, ranging from actual elections over movie rankings to competitor rankings from various sporting competitions. The number of items varies between 3 and 10 (details are summarized in the supplementary material).

All models are fit using maximum likelihood estimation, and Kullback-Leibler (KL) divergence between an empirical distribution and its estimation is used as a measure of the goodness of fit. In a first setting, we fit the models to the entire data, while in a second setting, we only fit to half of the data and determine divergence on the other half (averaging over 20 random splits).

In a first experiment, we compare ISR with our new variant CNS, with both insertion and quick sort as underlying sorting algorithms, using MM as an additional baseline. The ISR, CNS, and MM models are comparable in terms of their parametrization. A summary of the results in terms of win/tie/loss statistics is given in Table 1 (while the complete results can be found in the supplementary material). As can be seen, CNS shows a very strong performance, especially with insertion sort as a sorting algorithm.

In a second experiment, we compare CNS with its generalization GCNS, again with insertion and quick sort as underlying sorting algorithms in both models. The results in Table 1 clearly show that GCNS leads to better approximations. This is hardly surprising, given that GCNS has more parameters and therefore allows for fitting distributions in a more flexible way. Again, an instantiation with insertion sort seems to be preferable to the use of quick sort.

Table 1. Win/tie/loss statistics for the first (above) and second (below) experiment (first line/first setting, second line/second setting).

	CNS _I	CNS _Q	ISR	MM
CNS _I	—	197/0/16	197/0/16	170/0/43
	—	204/0/9	191/0/22	167/0/46
CNS _Q	16/0/197	—	165/0/48	143/0/70
	9/0/204	—	153/0/60	139/0/74
ISR	16/0/197	48/0/165	—	60/1/152
	22/0/191	60/0/153	—	57/0/156
MM	43/0/170	70/0/143	152/1/60	—
	46/0/167	74/0/139	156/0/57	—

	CNS _I	CNS _Q	GCNS _I	GCNS _Q
CNS _I	—	197/0/16	0/1/212	65/0/148
	—	204/0/9	25/0/188	81/0/132
CNS _Q	16/0/197	—	4/0/209	8/0/205
	9/0/204	—	10/0/203	18/0/195
GCNS _I	212/1/0	209/0/4	—	170/0/43
	188/0/25	203/0/10	—	169/0/44
GCNS _Q	148/0/65	205/0/8	43/0/170	—
	132/0/81	195/0/18	44/0/169	—

7. Conclusion and Future Work

Adopting the idea of a data-generating process in the form of a noisy sorting procedure, we proposed a variant of a parametrized probability distribution on rankings as recently proposed by Biernacki and Jacques (Biernacki & Jacques, 2013), as well as a generalization that is more flexible and makes less stringent coherence assumptions. Our models have an intuitive interpretation, exhibit convenient mathematical properties, and seem to fit empirical data very well. For two sorting algorithms, insertion sort and quick sort, we developed parameter estimation techniques based on a closed-form expression of the likelihood function for the former, and a recursive characterization of it for the latter. Experimentally, insertion sort leads to better performance.

In future work, we plan to consider other sorting algorithms, such as merge sort and heap sort. Another direction worth to investigate is the analysis of algebraic properties of our models using tools from computational algebraic geometry (Geiger et al., 2006); such properties may simplify the handling of the model and help to further improve efficiency of parameter estimation. Last but not least, we are also interested in using the model for other machine learning problems, in which distributions on rankings are needed, such as learning to rank (Ailon et al., 2005; Ailon, 2008; Cao et al., 2007) and multi-armed bandits (Busa-Fekete & Hüllermeier, 2014; Szörényi et al., 2015).

Acknowledgements

The authors gratefully acknowledge financial support by the Germany Research Foundation (DFG).

References

- Ailon, N. Reconciling Real Scores with Binary Comparisons: A New Logistic Based Model for Ranking. In *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, pp. 25–32, 2008.
- Ailon, N., Charikar, M., and Newman, A. Aggregating Inconsistent Information: Ranking and Clustering. In *Proceedings of the Annual ACM Symposium on Theory of Computing (STOC)*, pp. 684–693, 2005.
- Babington-Smith, B. Discussion of Professor Ross’s Paper. *Journal of the Royal Statistical Society B*, 12:153–162, 1950.
- Biernacki, C. and Jacques, J. A Generative Model for Rank Data Based on Insertion Sort Algorithm. *Computational Statistics & Data Analysis*, 58(15):162–176, 2013.
- Braverman, M. and Mossel, E. Noisy Sorting Without Resampling. In *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 268–276, 2008.
- Braverman, M. and Mossel, E. Sorting from Noisy Information. *CoRR*, abs/0910.1191, 2009.
- Busa-Fekete, R. and Hüllermeier, E. A Survey of Preference-based Online Learning with Bandit Algorithms. In *Proceedings of Conference on Learning Theory (COLT)*, pp. 18–39, 2014.
- Cao, Z., Qin, T., Liu, T.-Y., Tsai, M.-F., and Li, H. Learning to Rank: From Pairwise Approach to Listwise Approach. In *Proceedings of International Conference on Machine Learning (ICML)*, pp. 129–136, 2007.
- Cheng, W., Hüllermeier, E., Waegeman, W., and Welker, V. Label Ranking with Partial Abstention Based on Thresholded Probabilistic Models. In *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, pp. 2501–2509, 2012.
- Darroch, J. N. and Ratcliff, D. Generalized Iterative Scaling for Log-Linear Models. *The Annals of Mathematical Statistics*, 43(5):1470–1480, 1972.
- Fligner, M. A. and Verducci, J. S. Distance Based Ranking Models. *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 359–369, 1986.
- Fürnkranz, J. and Hüllermeier, E. Preference Learning: An Introduction. In *Preference learning*, pp. 1–17. Springer, 2010.
- Geiger, D., Meek, C., and Sturmfels, B. On the Toric Algebra of Graphical Models. *The Annals of Statistics*, 34(3):1463–1492, 2006.
- Grabisch, M., Marichal, J.-L., Mesiar, R., and Pap, E. *Aggregation Functions*. Cambridge University Press, 2009.
- Liu, T.-Y. Learning to Rank for Information Retrieval. *Foundations and Trends in Information Retrieval*, 3(3):225–331, 2009.
- Luce, R. D. *Individual Choice Behavior: A Theoretical Analysis*. Wiley, 1959.
- Mallows, C. Non-null Ranking Models. *Biometrika*, 44(1):114–130, 1957.
- Malouf, R. A Comparison of Algorithms for Maximum Entropy Parameter Estimation. In *Proceedings of Conference on Natural Language Learning (COLING)*, pp. 1–7, 2002.
- Marden, J. I. *Analyzing and Modeling Rank Data*. CRC Press, 1996.
- Meek, C. and Meila, M. Recursive Inversion Models for Permutations. In *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, pp. 631–639, 2014.
- Plackett, R. The Analysis of Permutations. *Applied Statistics*, 24:193–202, 1975.
- Szörényi, B., Busa-Fekete, R., Paul, A., and Hüllermeier, E. Online Rank Elicitation for Plackett-Luce: A Dueling Bandits Approach. In *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, pp. 604–612, 2015.