

Aggregating partial rankings with applications to peer grading in massive online open courses*

Ioannis Caragiannis[†] George A. Krimpas[‡] Alexandros A. Voudouris[§]

Abstract

We investigate the potential of using ordinal peer grading for the evaluation of students in massive online open courses (MOOCs). According to such grading schemes, each student receives a few assignments (by other students) which she has to rank. Then, a global ranking (possibly translated into numerical scores) is produced by combining the individual ones. This is a novel application area for social choice concepts and methods where the important problem to be solved is as follows: how should the assignments be distributed so that the collected individual rankings can be easily merged into a global one that is as close as possible to the ranking that represents the relative performance of the students in the assignment? Our main theoretical result suggests that using very simple ways to distribute the assignments so that each student has to rank only k of them, a Borda-like aggregation method can recover a $1 - \mathcal{O}(1/k)$ fraction of the true ranking when each student correctly ranks the assignments she receives. Experimental results strengthen our analysis further and also demonstrate that the same method is extremely robust even when students have imperfect capabilities as graders. We believe that our results provide strong evidence that ordinal peer grading can be a highly effective and scalable solution for evaluation in MOOCs.

1 Introduction

Massive online open courses (MOOCs) such as Coursera and EdX have currently become a trend and have attracted significant funding from VCs and support from leading academics. Their vision is to use the Internet and provide (to huge numbers of students) an educational experience that is typical in courses targeted to small audiences in top-class Universities. Whether MOOCs will become the next big business over the Internet strongly depends on whether they will satisfy the fundamental need for easy and cheap access to high quality education without restrictions. An apparent bottleneck for their full deployment and success is the fact that assessment and grading with the classical means is extremely costly. A typical approach is to use closed type questions in exams or assignments so that grading can be done automatically. This is highly unsatisfactory when, as part of a course, one would like to evaluate the students' ability of proving a mathematical statement, or expressing their critical thinking over an issue, or even demonstrating their creative writing skills. Evaluating this ability is inherently a *human computation* task [14].

An approach that has been proposed is to outsource the grading task to the students participating in the exam or assignment themselves; for example, they can be required to grade (a small number of) their

*This work was partially supported by the European Social Fund and Greek national funds through the research funding program Thales on "Algorithmic Game Theory", by the Caratheodory research grant E.114 from the University of Patras, and by the COST Action IC1205 on "Computational Social Choice".

[†]Computer Technology Institute "Diophantus" & Department of Computer Engineering and Informatics, University of Patras, 26504 Rion, Greece. Email: caragian@ceid.upatras.gr

[‡]Department of Computer Engineering and Informatics, University of Patras, 26504 Rion, Greece. Email: krimpas@ceid.upatras.gr

[§]Department of Computer Engineering and Informatics, University of Patras, 26504 Rion, Greece. Email: voudouris@ceid.upatras.gr

peers’ assignments as part of their own assignment [22]. Of course, allowing the students to grade using cardinal scores is risky; they are not experienced in assessing their peers’ performance in absolute terms¹ and they may have strong incentives to assign low scores to everybody in order to increase their own relative success in the assignment. An alternative that sounds feasible is to ask each student to provide a ranking of a small number of her peers’ assignments and then compute a global ranking by merging the partial ones; this is known as *ordinal peer grading* (e.g., see [23, 24]). Can this global ranking be in accordance to the objective comparison of students in terms of their performance in the assignment? Which are the necessary methods for this computation? And how accurate can this global ranking be? In this paper, we address these questions and provide both conceptual and technical answers.

Merging individual rankings into a global one is the main goal of voting rules from *social choice theory*, where a set of voters provide full rankings over the available alternatives and a voting rule has to transform this input into a winning alternative or an aggregate ranking of the alternatives. At first glance, ordinal peer grading seems to be a natural application area for classical voting theory. Interestingly, its particular characteristics deviate from those usually studied in the voting literature. First, each voter is also an alternative. This is a rare assumption in social choice in works that focus mostly on incentives issues (e.g., see [2, 13]). Second, the input consists of partial rankings over small subsets of alternatives. The closest such approach in social choice is known as preference elicitation [7] where simple queries are asked to each voter about their preferences; for example, in top- k elicitation [10], each voter provides the partial ranking of the k alternatives she likes the most. The (complexity) effects of using only partial rankings in voting have been studied under the possible and necessary winner problems (e.g., see [27]). An important characteristic of ordinal peer grading is that the partial rankings have the same size and that each assignment is given to the same number of graders. And finally, there is an objective way to assess the ordinal peer grading outcome by comparing it to the objective comparison of the students in terms of performance in the assignment. This is close in spirit to recent approaches that use voting in order to *learn* a ground truth [5, 6], such as a winning alternative or an underlying true ranking. In our work, we deviate from these studies as well since we aim to learn the ground truth only approximately. So, ordinal peer grading is a setting where ideas and analysis techniques from human computation, voting, and learning are blended together in novel ways.

In particular, our model uses a grading scheme that asks each student to rank the assignments of k other students. For fairness reasons, we restrict ourselves to grading schemes that distribute each assignment to exactly k students. Unlike recent studies [23, 24], we investigate the potential of applying ordinal peer grading exclusively, i.e., without involving any professionals in grading. We assume that there is an underlying true (strict) ranking of the assignments (the ground truth) and we would like to recover correctly an as high as possible fraction of it using input from the students. We have two scenarios that determine the input. In the first one, we assume that, after the students have submitted their assignments, the instructor announces indicative solutions and grading instructions. Here, we make the simplifying assumption that each student grades the assignments in her bundle consistently to the ground truth (perfect grading). In a second scenario that is also assumed in [23], we assume that grading is performed without any guidance by the instructor. Here, the natural assumption is that the quality of a student determines both her performance in the assignment and her grading ability. We have mostly focused on simple rank aggregation rules such as the adaptation of the classical Borda count [4], where the partial ranking provided by each grader is interpreted as follows: k points are given to the assignment ranked first, $k - 1$ points to the one ranked second, and so on. The global ranking is then computed by ordering the assignments in decreasing order of these Borda scores. We have also considered more aggregation rules which are described in detail in Sections 2 and 4.

Our technical contributions can be summarized as follows. In Section 3, we present a theoretical analysis of Borda when the partial rankings on input are consistent to the ground truth. We prove that using any way to distribute k assignments per student, Borda recovers correctly an expected fraction of

¹This is in contrast to the main assumption behind the reviewing systems that are used in academic conferences.

$1 - \mathcal{O}(1/\sqrt{k})$ of the pairwise relations in the ground truth. If the distribution of the assignments has some particularly desired simple structure, an even better guarantee of $1 - \mathcal{O}(1/k)$ is obtained. The independence of these results from the number of students is rather surprising. Our proofs exploit the beautiful theory of *martingales* in order to cope with dependencies between random variables that are involved in the analysis. To the best of our knowledge, this is the first application of martingales in social choice. We also present extensive experiments with Borda and other aggregation rules (in Section 4). Our findings further justify the robustness of Borda, even in the scenario of imperfect grading. For example, Borda is shown to recover more than 88% of the ground truth by distributing 8 assignments per student (with students having highly varying grading capabilities). Here, we borrow ideas from recent studies on voting and learning (e.g., [5]) and use noise models for the generation of random partial rankings whose distance from the ground truth depends probabilistically on the quality of the graders. En route, we provide some intuition about the problem (in Section 2). We conclude with a discussion of (the many) possible extensions of our work in Section 5.

2 Problem statement, terminology and notation

Let \mathcal{A} denote a universe of n elements. A collection \mathcal{B} of subsets of \mathcal{A} is called a *grading scheme* with parameters n and $k \leq n$ (or (n, k) -grading scheme) if \mathcal{B} consists of n subsets of \mathcal{A} called *bundles*, each bundle has size k , and each element of \mathcal{A} belongs to exactly k subsets of \mathcal{B} . To see the relation to peer grading, we can view the elements of the universe \mathcal{A} as the n papers of students participating in an assignment. Each bundle contains k papers that will be graded by a distinct student. Of course, we require that no student will grade her own paper. This can be easily achieved by a matching computation.²

Alternatively, we can represent the (n, k) -grading scheme with a bipartite graph $G = (U, V, E)$ which we will call (n, k) -*bundle graph*. The set of nodes U has size n and contains a distinct node for each element of \mathcal{A} . The set of nodes V has size n too and contains a node for each bundle of \mathcal{B} . The set of edges E contains an edge (u, v) connecting node $u \in U$ with node $v \in V$ if and only if the element corresponding to node u belongs to the bundle corresponding to node v . Clearly, an (n, k) -bundle graph is k -regular. Actually, every k -regular bipartite graph has the same number n of nodes in both bipartition sides and be used as an (n, k) -bundle graph.

A *partial ranking* \succ_b associated with a bundle $b \in \mathcal{B}$ is simply a ranking of the elements b contains. We remark that \succ_b is undefined for elements not belonging to \mathcal{B} . A *profile* is simply the collection that contains the partial ranking \succ_b for each bundle b of \mathcal{B} . An *aggregation rule* takes as input a profile of partial rankings and computes a complete ranking of all elements. A typical example is the following rule that extends Borda count from classical voting theory. Each element gets a score from each appearance in a partial ranking. The *Borda score* of an element is then the sum of the scores from all partial rankings. Within each partial ranking, a score of k is given to the element that is ranked first, a score of $k - 1$ to the element that is ranked second, and so on. The final complete ranking is computed by sorting the elements in decreasing order in terms of their Borda scores. We will use the term *Borda* to refer to this aggregation rule. Even though one can think of several different ways to resolve ties, we simply ignore ties in our theoretical analysis (Section 3) and use uniformly random tie-breaking in our experiments (Section 4).

We have also considered another aggregation rule which we call *Random Serial Dictatorship* (RSD). The term is inspired by the well-known mechanism for house allocation markets [1]. A complete ranking is computed gradually starting from an initially empty one. In a first *serial phase*, the partial rankings are considered in a random order. When considering a partial ranking, we copy to the global one all the pairwise relations that do not contradict (i.e., do not form cycles with) relations copied earlier. When all

²Indeed, for every student i , there are $n - k$ bundles that do not contain her paper. Then, the bipartite graph that represents the information about the bundles that a student is allowed to grade is regular and, by Hall's matching theorem, has a perfect matching. This matching can be used to assign bundles of papers to students.

partial rankings have been considered, the global partial ranking is augmented by the pairwise relations implied due to transitivity (e.g., the pairwise relations $x \succ y$ and $y \succ z$ copied from two partial rankings imply that $x \succ z$ as well). Then, we use a second *random completion phase* to complete the global ranking as follows. In each step, we pick a random pair of elements whose relation has not been decided so far. We make this decision randomly and update all pairwise relations that this decision and the existing ones imply due to transitivity. We continue this way until all pairwise relations have been decided.

We are now ready to give the statement of the problem that we consider more formally. In general, we would like to use the grading schemes and aggregation rules in order to *learn* an unknown *ground truth*, i.e., a ranking of the elements representing their relative quality. A first question is whether the ground truth can be learnt *with certainty* when the partial rankings are consistent to it. In other words, we ask for an *order-revealing* grading scheme (and a corresponding order-revealing bundle graph) which defines the bundles in such a way that the partial rankings contain enough information so that all pairwise relations in the ground truth can be recovered with certainty. Unfortunately, order-revealing grading schemes have severe limitations. In particular, they should have the following too demanding property: for every pair of elements, there should be some bundle that contains both of them.³ Indeed, let \mathcal{B} be an order-revealing grading scheme over a universe \mathcal{A} of n elements and assume that there are two elements x and y so that no bundle contains both x and y . Now, consider a ranking \succ that has x and y in the first two positions and let \succ' be the ranking that differs from \succ only in the order of x and y . Clearly, the partial rankings within the bundles are identical in both cases and, as a result, there is no way to identify whether the ground truth is the ranking \succ or the ranking \succ' . Notice that the above property implies that RSD combined with order-revealing grading schemes recovers the ground truth with certainty (and does not have to run the random completion phase). This is not the case for Borda unless any two elements co-exist in the same number of bundles (like in the bundle graphs constructed below).

Clearly, the maximum number of elements that belong to a bundle with x is $k(k-1)$ and this number should be at least $n-1$ if we want x to belong to some bundle with every other element. This immediately implies that order-revealing grading schemes should have bundles of size $\Omega(\sqrt{n})$. In sharp contrast to this disappointing observation, we will see that the goal of *approximate* order-revealing grading schemes is a very feasible one and leads to effective and scalable grading solutions in theory and practice. Interestingly, many of our findings make use of bundle graphs that are order-revealing; this is why we have included the following explicit construction of order-revealing grading schemes for particular values of the parameters n and k here.

Let $p \geq 1$ be a prime and let \mathcal{A} be a universe with $n = p^2 + p + 1$ elements. We will construct the grading scheme \mathcal{B} in which each bundle has size exactly $k = p + 1$. Observe that these values for n and k satisfy the lower-bound condition mentioned above with equality. Rename the elements of \mathcal{A} as $\mathcal{A} = \{u\} \cup \{v_i | i = 0, \dots, p-1\} \cup \{w_{i,j} | i = 0, \dots, p-1, j = 0, \dots, p-1\}$ and define the bundles of \mathcal{B} as follows:

- $F = \{u, v_0, v_1, \dots, v_{p-1}\}$;
- For $i = 0, \dots, p-1$, $R_i = \{u\} \cup \{w_{i,j} | j = 0, \dots, p-1\}$;
- For $i = 0, \dots, p-1$ and $s = 0, \dots, p-1$, $C_{i,s} = \{v_s\} \cup \{w_{j,(i+j \cdot s) \bmod p} | j = 0, \dots, p-1\}$.

An order-revealing $(7, 3)$ -bundle graph is depicted in Figure 1; it represents the following grading scheme \mathcal{B} . The underlying universe is $\mathcal{A} = \{1, 2, 3, 4, 5, 6, 7\}$ and \mathcal{B} has the following seven 3-sized bundles: $\{1, 2, 3\}$, $\{1, 4, 5\}$, $\{1, 6, 7\}$, $\{2, 4, 6\}$, $\{2, 5, 7\}$, $\{3, 4, 7\}$, and $\{3, 5, 6\}$. The numbering of

³This property essentially asks for a k -regular bipartite graph of diameter at most 3. Our order-revealing bundle graphs are known as Moore bipartite graphs, i.e., they are the smallest bipartite graphs of degree at least k and of diameter at most 3; see [18] for a detailed survey on the degree-diameter problem.

nodes in set V indicates an assignment of bundles to students for grading and, hence, nodes with the same number are not adjacent.

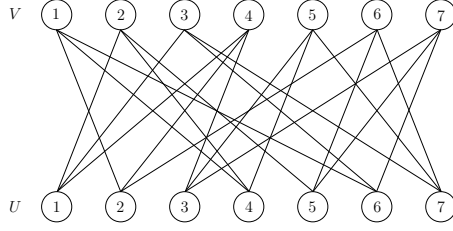


Figure 1: An order-revealing $(7, 3)$ -bundle graph.

We prove the correctness of our construction using basic facts from number theory.

Lemma 1. *The above construction yields an order-revealing grading scheme.*

Proof. Clearly, the above grading scheme consists of $n = p^2 + p + 1$ bundles of size $k = p + 1$. Also, observe that each element belongs to exactly $p + 1$ bundles. Indeed, element u belongs to sets F and R_i for $i = 0, 1, \dots, p - 1$. Element v_s belongs to sets F and $C_{i,s}$ for $i = 0, 1, \dots, p - 1$. Element $w_{i,j}$ belongs to sets R_i and $C_{t,s}$ such that $j = (t + i \cdot s) \bmod p$.

We complete the proof by showing that for every pair $x, y \in \mathcal{A}$, there exists a bundle that contains both x and y . This is clearly true if one of x and y is u or if both x and y belong to F or to some R_i , for $i = 0, \dots, p - 1$. So, there are two more cases to be considered. First, assume that $x = v_s$ for $s \in \{0, \dots, p - 1\}$ and $y = w_{j,\ell}$. Then, there exists an $i \in \{0, \dots, p - 1\}$ such that $i + j \cdot s = \ell \pmod{p}$ and, hence, both x and y belong to set $C_{i,s}$. It remains to consider the case where $x = w_{i_1,j_1}$ and $y = w_{i_2,j_2}$ with $0 \leq i_1 < i_2 \leq p - 1$. Then, there exists a unique $s \in \{0, \dots, p - 1\}$ such that $(i_2 - i_1) \cdot s = (j_2 - j_1) \pmod{p}$. This follows from the facts that p is prime and that any linear equation of the form $a \cdot z = b \pmod{n}$ has $\gcd(a, n)$ solutions if and only if $\gcd(a, n)$ divides b . Now, set $i = (j_1 - i_1 \cdot s) \pmod{p}$ and observe that $i_1 = (i + j_1 \cdot s) \pmod{p}$ and $i_2 = (i + j_2 \cdot s) \pmod{p}$. Hence, both x and y belong to $C_{i,s}$ and the proof is complete. \square

We now relax our requirements and seek for an approximate order-revealing grading scheme. Our aim is to use a bundle graph of simple structure and of very low (i.e., independent of n) degree and still be able to correctly recover a high fraction of the $\binom{n}{2}$ pairwise relations in the ground truth. Our grading schemes will be randomized in the sense that we will always randomly permute the elements before associating them to nodes of set U of the bundle graph; let $\pi : U \rightarrow \mathcal{A}$ denote this bijection (or permutation). Sometimes, in our experiments, the bundle graphs we use are themselves random. Much of our work (i.e., our theoretical analysis in Section 3 as well as the first among the two sets of experiments reported in Section 4) has focused on the scenario where the partial rankings are consistent to the ground truth. Our second set of experiments in Section 4 uses partial rankings that deviate from the ground truth according to a noise model.

3 Analysis of Borda

In this section, we present our theoretical results. We assume that the (n, k) -bundle graph $G = (U, V, E)$ has $k \geq 3$ and $n \geq 3k(k - 1) + 2$. These are technical assumptions that do not affect the applicability of our results; recall that, in practice, we would like n and k to be huge and very small, respectively. Surprisingly, Borda correctly recovers a very large fraction of the ground truth as the next statement suggests.

Theorem 2. *When Borda is applied on partial rankings that are consistent to the ground truth, the expected fraction of correctly recovered pairwise relations is at least $1 - \mathcal{O}(1/k)$ when the (n, k) -bundle graph has girth at least 6, and at least $1 - \mathcal{O}(1/\sqrt{k})$ in general.*

We prove this theorem by relating the performance of Borda only to the degree k and on a quantity $\eta(G)$ that characterizes the structure of the bundle graph. For the definition of $\eta(G)$, we need some notation; this will be heavily used throughout this section. Given two nodes u, v of U , we use $\lambda_{u,v}$ to denote their common neighbourhood in V , i.e., $\lambda_{u,v} = |N(u) \cap N(v)|$. Observe that $\sum_{v \in U \setminus \{u\}} \lambda_{u,v} = k(k-1)$ since G is k -regular. Also, we define the quantity $\theta_{u,v}$ as $\theta_{u,v} = 4 \sum_{z \in N(N(u,v)) \setminus \{u,v\}} (\lambda_{u,z} + \lambda_{v,z})^2$. Then,

$$\eta(G) = \frac{1}{n(n-1)} \sum_{u,v \in U} \sqrt{\theta_{u,v}},$$

where the sum runs over all ordered pairs of u, v in U .

Intuitively, the quantity $\eta(G)$ is small when, on average, the common neighbourhood between pairs of nodes is small. The extreme case is when the common neighbourhood consists of a single node; in this case, the graph has girth⁴ at least 6. The next lemma provides upper bounds on $\eta(G)$ that will be useful later.

Lemma 3. *For every k -regular bipartite graph G , $\eta(G) \leq \sqrt{8k(k-1)(4k-3)}$. Every k -regular bipartite graph G of girth at least 6 has $\eta(G) \leq 4\sqrt{k(k-1)}$.*

Proof. Consider two nodes $u, v \in U$ of an arbitrary k -regular bipartite graph. We will show that $\theta_{u,v}$ is at most $8k(k-1)(4k-3)$. Consider the sets of nodes $N(u) \cap N(v)$, $N(u) \setminus N(v)$, and $N(v) \setminus N(u)$, and the edges connecting these nodes to $N(N(u,v)) \setminus \{u,v\}$. Each edge from a node of $N(u) \cap N(v)$ to a node $z \in N(N(u,v)) \setminus \{u,v\}$ contributes 2 to the quantity $\lambda_{u,z} + \lambda_{v,z}$, which can be up to $2k$. Hence, each edge from a node of $N(u) \cap N(v)$ to a node $z \in N(N(u,v)) \setminus \{u,v\}$ contributes at most $(2k)^2 - (2k-2)^2 = 8k-4$ to the quantity $(\lambda_{u,z} + \lambda_{v,z})^2$ and there are $|N(u) \cap N(v)|(k-2)$ such edges. Similarly, each edge from a node of $N(u) \setminus N(v)$ and $N(v) \setminus N(u)$ to a node $z \in N(N(u,v)) \setminus \{u,v\}$ contributes 1 to the quantity $\lambda_{u,z} + \lambda_{v,z}$, which can be up to $2k-1$. Hence, each edge from a node of $N(u) \setminus N(v)$ or $N(v) \setminus N(u)$ to a node $z \in N(N(u,v)) \setminus \{u,v\}$ contributes at most $(2k-1)^2 - (2k-2)^2 = 4k-3$ to the quantity $(\lambda_{u,z} + \lambda_{v,z})^2$ and there are $2(k - |N(u) \cap N(v)|)(k-1)$ such edges. So, $\theta_{u,v}$ is bounded by 4 times the total contributions to quantities $(\lambda_{u,z} + \lambda_{v,z})^2$ by the edges between $N(u, v)$ and $N(N(u, v)) \setminus \{u, v\}$, i.e., by $4(|N(u) \cap N(v)|(k-2)(8k-4) + 2(k - |N(u) \cap N(v)|)(k-1)(4k-3)) \leq 8k(k-1)(4k-3)$.

Now, assume that the graph has girth at least 6; this means that $\lambda_{u,z} + \lambda_{v,z} \leq 2$ for any node $z \in N(N(u,v)) \setminus \{u,v\}$, otherwise z would be in a 4-cycle with either u or v . We will show that $\theta_{u,v} \leq 16k(k-1)$. Each node $z \in N(N(u,v)) \setminus \{u,v\}$ can be adjacent to either one node of $N(u) \cap N(v)$ or (exclusive) to at most one node of $N(u) \setminus N(v)$ and at most one node of $N(v) \setminus N(u)$. Among the nodes in $N(N(u,v)) \setminus \{u,v\}$, denote by D_2 the ones that are adjacent to one node from $N(u) \setminus N(v)$ and to one node from $N(v) \setminus N(u)$. So, any node z that is among the $|N(u) \cap N(v)|(k-2)$ neighbours of $N(u) \cap N(v)$ in $N(N(u,v)) \setminus \{u,v\}$ or belongs to D_2 has $\lambda_{u,z} + \lambda_{v,z} = 2$. Any node z among the remaining $2(k - |N(u) \cap N(v)|)(k-1) - 2|D_2|$ nodes of $N(N(u,v)) \setminus \{u,v\}$ has $\lambda_{u,z} + \lambda_{v,z} = 1$. Hence, $\theta_{u,v} = 4(4(|N(u) \cap N(v)|(k-2) + |D_2|) + 2(k - |N(u) \cap N(v)|)(k-1) - 2|D_2|) = 8k(k-1) + 8(|N(u) \cap N(v)|(k-2) + |D_2|) - 8|N(u) \cap N(v)|$. The second part of the lemma follows by observing that the quantity $|N(u) \cap N(v)|(k-2) + |D_2|$ is the number of nodes z of $N(N(u,v)) \setminus \{u,v\}$ with $\lambda_{u,z} + \lambda_{v,z} = 2$ which cannot be higher than $k(k-1)$. \square

The important step in the proof of Theorem 2 is to focus on two elements a_r and a_q with ranks (positions) $r < q$ in the ground truth and to bound from above the probability that the difference in

⁴The girth of a graph is the length of its smallest cycle.

their Borda scores is inconsistent to their rank difference. This will require to take care of several subtle dependencies among the random variables involved. We will do so by exploiting the beautiful theory of martingales and a well-known tail inequality about them. The necessary background from martingale theory is presented below; the interested reader can refer to the textbooks [19] and [20] for an introduction to martingales and their applications.

Definition 4. A sequence of random variables Z_0, Z_1, \dots, Z_m is a *martingale* with respect to a second sequence of random variables X_1, X_2, \dots, X_m if for every $i = 1, \dots, m$, it holds that $\mathbb{E}[Z_i | X_1, \dots, X_i] = Z_{i-1}$.

The next definition provides a general way to define martingales associated with *any* random variable and was first used by Doob [8].

Definition 5. Consider a random variable W and a sequence of random variables X_1, \dots, X_m . Then, the sequence of random variables Z_0, \dots, Z_m such that $Z_0 = \mathbb{E}[W]$ and $Z_i = \mathbb{E}[W | X_1, \dots, X_i]$ for every $i = 1, \dots, m$ is a martingale, called a *Doob martingale*.

We can now present a powerful tail inequality for martingales that is known as Azuma-Hoeffding inequality (see Azuma [3] and Hoeffding [12]).

Lemma 6 (Azuma-Hoeffding inequality). *Let Z_0, Z_1, \dots, Z_m be a martingale with $|Z_i - Z_{i-1}| \leq c_i$ for $i = 1, \dots, m$. Then, for all $t \geq 0$, it holds that*

$$\Pr[Z_m - Z_0 \leq -t] \leq \exp\left(-\frac{t^2}{2 \sum_{i=1}^m c_i^2}\right).$$

We are now ready to show that the probability that the Borda score of a high-rank element is larger than the Borda score of a low-rank element is small. Importantly, it turns out that this probability decreases exponentially in terms of the rank difference. We will first study such phenomena under particular conditions on our bijection π .

Lemma 7. *Let $u, v \in U$, and consider the two elements $a_r, a_q \in \mathcal{A}$ with ranks $r < q$ in the ground truth. Let $W_{r,q}$ be the random variable denoting the difference of the Borda score of a_r minus the Borda score of a_q and let $\Gamma_{u,v}^{r,q}$ be the event that $\pi(u) = a_r$ and $\pi(v) = a_q$. Then,*

$$\mathbb{E}[W_{r,q} | \Gamma_{u,v}^{r,q}] = (k(k-1) - \lambda_{u,v}) \frac{q-r-1}{n-2} + \lambda_{u,v}$$

and

$$\Pr[W_{r,q} \leq 0 | \Gamma_{u,v}^{r,q}] \leq \exp\left(-\frac{\mathbb{E}[W_{r,q} | \Gamma_{u,v}^{r,q}]^2}{2\theta_{u,v}}\right).$$

Proof. We begin the proof by computing the expectation of the Borda scores. Element a_r gets one point for each bundle it belongs to plus one additional point for each appearance of an element with rank higher than r in the bundles a_r belongs to. Assuming that $\pi(u) = a_r$ and $\pi(v) = a_q$, there are $\lambda_{u,v}$ appearances of a_q in the bundles of a_r and $k(k-1) - \lambda_{u,v}$ appearances of elements different than a_r and a_q ; each of them has probability $\frac{n-r-1}{n-2}$ to have higher rank than r . Hence, the expected Borda score of element a_r is $k + (k(k-1) - \lambda_{u,v}) \frac{n-r-1}{n-2} + \lambda_{u,v}$. Similarly, element a_q gets one point for each bundle it belongs to plus one additional point for each appearance of an element with rank higher than q . There are $k(k-1) - \lambda_{u,v}$ appearances of elements different than a_r and a_q in bundles of a_q and each of them has rank higher than q with probability $\frac{n-q}{n-2}$. Hence, the expected Borda score of element a_q is $k + (k(k-1) - \lambda_{u,v}) \frac{n-q}{n-2}$, and the expectation of the difference $W_{r,q}$ is indeed

$$\mathbb{E}[W_{r,q} | \Gamma_{u,v}^{r,q}] = (k(k-1) - \lambda_{u,v}) \frac{q-r-1}{n-2} + \lambda_{u,v}.$$

Given $\Gamma_{u,v}^{r,q}$, define $S = N(N(u, v)) \setminus \{u, v\}$ to be the set of nodes in G that are at distance exactly 2 from u or v (not including u and v); notice that $|S| \leq 2k(k-1)$. Now, consider an arbitrary ordering $o : [|S|] \rightarrow S$ of the nodes of S and let X_i be the random variable denoting the rank of the element $\pi(o(i))$. Using the random variables X_i and the random variable $W_{r,q}$, we define the Doob martingale $Z_0, Z_1, \dots, Z_{|S|}$ such that $Z_0 = \mathbb{E}[W_{r,q} | \Gamma_{u,v}^{r,q}]$ and $Z_i = \mathbb{E}[W_{r,q} | \Gamma_{u,v}^{r,q}, X_1, \dots, X_i]$ (hence, given $\Gamma_{u,v}^{r,q}$, $W_{r,q} = Z_{|S|}$). The next technical lemma bounds the difference $|Z_i - Z_{i-1}|$ for $i = 1, \dots, |S|$.

Lemma 8. *For every $i = 1, \dots, |S|$, it holds that $|Z_i - Z_{i-1}| \leq 2(\lambda_{u,o(i)} + \lambda_{v,o(i)})$.*

Proof. Throughout this proof, all random variables and probabilities are conditioned on the event $\Gamma_{u,v}^{r,q}$, even if, in order to simplify notation, we do not explicitly write so.

For every node $w \in S$, denote by $\mu_{u,v,w} = |N(u) \cap N(v) \cap N(w)|$ the number of common neighbours between u, v , and w . We can now express $W_{r,q}$ using the following observations: the Borda score difference

- increases for each appearance of element a_q in the same bundle with a_r ,
- for each appearance of element $\pi(o(j))$ in a bundle containing a_r but not a_q provided that the rank of $\pi(o(j))$ is higher than r , and
- for each appearance of an element $\pi(o(j))$ in a bundle containing both a_r and a_q provided that the rank of $\pi(o(j))$ is between r and q , and
- decreases for each appearance of element $\pi(o(j))$ in a bundle containing a_q but not a_r provided that the rank of $\pi(o(j))$ is higher than q .

Using our notation $\lambda_{u,v}$ and $\mu_{u,v,o(j)}$, we have

$$\begin{aligned} W_{r,q} &= \lambda_{u,v} + \sum_{j=1}^{|S|} (\lambda_{u,o(j)} - \mu_{u,v,o(j)}) \mathbb{1}\{X_j > r\} + \sum_{j=1}^{|S|} \mu_{u,v,o(j)} \mathbb{1}\{r < X_j < q\} \\ &\quad - \sum_{j=1}^{|S|} (\lambda_{v,o(j)} - \mu_{u,v,o(j)}) \mathbb{1}\{X_j > q\} \\ &= \lambda_{u,v} + \sum_{j=1}^{|S|} (\lambda_{u,o(j)} \mathbb{1}\{X_j > r\} - \lambda_{v,o(j)} \mathbb{1}\{X_j > q\}). \end{aligned}$$

Denoting by \mathbf{X}_i the sequence X_1, \dots, X_i , we have that the difference $|Z_i - Z_{i-1}|$ is

$$\begin{aligned} Z_i - Z_{i-1} &= \sum_{j=i}^{|S|} \lambda_{u,o(j)} (\Pr[X_j > r | \mathbf{X}_i] - \Pr[X_j > r | \mathbf{X}_{i-1}]) \\ &\quad - \sum_{j=i}^{|S|} \lambda_{v,o(j)} (\Pr[X_j > q | \mathbf{X}_i] - \Pr[X_j > q | \mathbf{X}_{i-1}]). \end{aligned} \tag{1}$$

Once the values of X_1, \dots, X_{i-1} are determined, let x and y be the number of available ranks from $[n] \setminus \{r, q, X_1, \dots, X_{i-1}\}$ that are between r and q and higher than q , respectively. Hence, for $j = i, \dots, |S|$, we have

$$\begin{aligned} \Pr[X_j > r | \mathbf{X}_{i-1}] &= \frac{x+y}{n-i-1}, \\ \Pr[X_j > q | \mathbf{X}_{i-1}] &= \frac{y}{n-i-1}, \end{aligned}$$

and for $j = i + 1, \dots, |S|$, we have

$$\Pr[X_j > r | \mathbf{X}_i] = \frac{x + y - \mathbb{1}\{X_i > r\}}{n - i - 2},$$

$$\Pr[X_j > q | \mathbf{X}_i] = \frac{y - \mathbb{1}\{X_i > q\}}{n - i - 2}.$$

Now, (1) yields

$$\begin{aligned} Z_i - Z_{i-1} &= \lambda_{u,o(i)} \left(\mathbb{1}\{X_i > r\} - \frac{x + y}{n - i - 1} \right) - \lambda_{v,o(j)} \left(\mathbb{1}\{X_i > q\} - \frac{y}{n - i - 1} \right) \\ &+ \sum_{j=i+1}^{|S|} \lambda_{u,o(j)} \left(\frac{x + y - \mathbb{1}\{X_i > r\}}{n - i - 2} - \frac{x + y}{n - i - 1} \right) \\ &+ \sum_{j=i+1}^{|S|} \lambda_{v,o(j)} \left(\frac{y - \mathbb{1}\{X_i > q\}}{n - i - 2} - \frac{y}{n - i - 1} \right) \\ &= \left(\lambda_{u,o(i)} - \frac{\sum_{j=i+1}^{|S|} \lambda_{u,o(j)}}{n - i - 2} \right) \left(\mathbb{1}\{X_i > r\} - \frac{x + y}{n - i - 1} \right) \\ &+ \left(\lambda_{v,o(i)} - \frac{\sum_{j=i+1}^{|S|} \lambda_{v,o(j)}}{n - i - 2} \right) \left(\frac{y}{n - i - 1} - \mathbb{1}\{X_i > q\} \right) \end{aligned}$$

The second and fourth parenthesis in the above expression are obviously between -1 and 1 . Recall that $\sum_{j=i+1}^{|S|} \lambda_{u,o(j)} \leq k(k-1)$ and $\sum_{j=i+1}^{|S|} \lambda_{v,o(j)} \leq k(k-1)$. Also, by the definition of S , $\lambda_{u,o(i)} + \lambda_{v,o(i)} \geq 1$, for every $i = 1, \dots, |S|$. Combined with our assumption that $n \geq 3k(k-1) + 2$, these properties imply that the first parenthesis is between $-\max\{\lambda_{u,o(i)}, 1\}$ and $\max\{\lambda_{u,o(i)}, 1\}$, and the third one is between $-\max\{\lambda_{v,o(i)}, 1\}$ and $\max\{\lambda_{v,o(i)}, 1\}$. The lemma follows since $|\max\{\lambda_{u,o(i)}, 1\} + \max\{\lambda_{v,o(i)}, 1\}| \leq 2(\lambda_{u,o(i)} + \lambda_{v,o(i)})$. \square

Lemma 7 then follows by applying the Azuma-Hoeffding inequality (Lemma 6) with $t = \mathbb{E}[W_{r,q} | \Gamma_{u,v}^{r,q}]$ and using Lemma 8 to bound the difference $|Z_i - Z_{i-1}|$. \square

The proof of Theorem 2 can now be completed using Lemmas 3 and 7.

Proof of Theorem 2. Consider the pair of elements with true ranks r and q so that $r < q$. The correct pairwise relation between the two elements will be recovered when the Borda score of the low-rank element is higher than the Borda score of the high-rank one (there is the additional case where the two elements are tied and the tie is resolve in favour of the low-rank element but we will ignore this case; this will only make our result stronger). Again, $W_{r,q}$ will be the random variable denoting the difference between the Borda scores of the low- and high-rank elements. Then, by Lemma 7 the probability that the relation between the elements with ranks r and q is correctly recovered is

$$\begin{aligned} \Pr[W_{r,q} > 0] &= 1 - \sum_{u,v \in U} (\Pr[W_{r,q} \leq 0 | \Gamma_{u,v}^{r,q}] \Pr[\Gamma_{u,v}^{r,q}]) \\ &\geq 1 - \frac{1}{n(n-1)} \sum_{u,v \in U} \exp\left(-\frac{\mathbb{E}[W_{r,q} | \Gamma_{u,v}^{r,q}]^2}{2\theta_{u,v}}\right) \\ &= 1 - \frac{1}{n(n-1)} \sum_{u,v \in U} e^{-(\beta(u,v)y(q-r) + \delta(u,v))^2}, \end{aligned}$$

where $\beta(u, v) = \frac{k(k-1)-\lambda_{u,v}}{\sqrt{2\theta_{u,v}}}$, $\delta(u, v) = \frac{\lambda_{u,v}}{\sqrt{2\theta_{u,v}}}$, and $y(t) = \frac{t-1}{n-2}$. Now, denoting the expected number of correctly recovered pairwise relations by C , we have

$$\begin{aligned}
C &= \sum_{r=1}^{n-1} \sum_{q=r+1}^n \Pr[W_{r,q} > 0] \\
&\geq \sum_{r=1}^{n-1} \sum_{q=r+1}^n \left(1 - \frac{1}{n(n-1)} \sum_{u,v \in U} e^{-(\beta(u,v)y(q-r)+\delta(u,v))^2} \right) \\
&= \frac{n(n-1)}{2} - \frac{1}{n(n-1)} \sum_{u,v \in U} \sum_{d=1}^{n-1} (n-d) e^{-(\beta(u,v)y(d)+\delta(u,v))^2} \\
&\geq \frac{n(n-1)}{2} - \sum_{u,v \in U} \int_0^1 (1-y) e^{-(\beta(u,v)y+\delta(u,v))^2} dy.
\end{aligned}$$

We will estimate the (Gaussian) integral using the following claim.

Claim 9. *Let $\beta > 0$ and $\delta \geq 0$. Then, $\int_0^1 (1-y) e^{-(\beta y+\delta)^2} dy \leq \frac{\beta+\delta}{2\beta^2} \sqrt{\pi}$.*

Proof. Denote by $\operatorname{erf}(y) = \frac{2}{\sqrt{\pi}} \int_0^y e^{-t^2} dt$ the error function. Then, we can verify by tedious calculations that

$$\begin{aligned}
\int_0^1 (1-y) e^{-(\beta y+\delta)^2} dy &= \frac{\beta+\delta}{2\beta^2} \sqrt{\pi} (\operatorname{erf}(\beta+\delta) - \operatorname{erf}(\beta)) + \frac{1}{2\beta^2} (e^{-(\beta+\delta)^2} - e^{-\beta^2}) \\
&\leq \frac{\beta+\delta}{2\beta^2} \sqrt{\pi},
\end{aligned}$$

where the inequality follows since the error function $\operatorname{erf}(y)$ takes values in $[0, 1]$ when $y \geq 0$. \square

Now, we use Claim 9 and the facts $\beta(u, v) = \frac{k(k-1)-\lambda_{u,v}}{\sqrt{2\theta_{u,v}}} \leq \frac{k(k-2)}{\sqrt{2\theta_{u,v}}}$ and $\delta(u, v) = \frac{\lambda_{u,v}}{\sqrt{2\theta_{u,v}}}$ to obtain

$$\begin{aligned}
C &\geq \frac{n(n-1)}{2} - \sum_{u,v \in U} \frac{\beta(u, v) + \delta(u, v)}{2\beta(u, v)^2} \sqrt{\pi} \\
&\geq \frac{n(n-1)}{2} - \frac{k-1}{k(k-2)^2} \sqrt{\frac{\pi}{2}} \sum_{u,v \in U} \sqrt{\theta_{u,v}} \\
&= \frac{n(n-1)}{2} \left(1 - \frac{k-1}{k(k-2)^2} \sqrt{2\pi} \eta(G) \right).
\end{aligned}$$

Now, the theorem follows by Lemma 3. Recall that $\eta(G)$ is at most $\sqrt{8k(k-1)(4k-3)}$ for every k -regular bipartite graph G and at most $4\sqrt{k(k-1)}$ when G has girth at least 6. Using the assumption that $k \geq 3$, we obtain that the rightmost parenthesis in the above expression becomes at least $1 - \frac{48\sqrt{2\pi}}{\sqrt{k}}$ and $1 - \frac{16\sqrt{3\pi}}{k}$, respectively. \square

4 Experimental evaluation

We now describe two sets of experiments that we have conducted.⁵ In the first one, we have studied perfect grading with Borda and RSD. We have considered three different types of bundle graphs. The

⁵All experiments presented in this section have been conducted in an Intel 12-core i7 machine with 32Gb of RAM running Windows 7. Our methods have been implemented in Matlab R2013a.

first type is that of random k -regular bipartite graphs. We build these graphs by picking k perfect matchings in the complete bipartite graph $K_{n,n}$ as follows. For each node of $K_{n,n}$ in –say– the upper⁶ node side, we select one edge among its incident ones uniformly at random. We remove this edge from $K_{n,n}$ and continue for the remaining nodes; this defines a random perfect matching. We repeat the above procedure k times. If a node at the upper side becomes isolated before the completion of the above procedure, we repeat from scratch. Otherwise, the set of edges that have been removed constitutes the bundle graph. The second type of graphs consists of many components of small girth-6 graphs. For $k = p + 1$, where p is a prime, we use the k -regular bipartite graph with $k^2 - k + 1$ nodes per side whose construction is described in Section 2 and which was proved to be order-revealing in Lemma 1. The bundle graph consists of multiple disconnected copies of this graph. Similarly, the third type of bundle graphs contains copies of the complete bipartite graph $K_{k,k}$ (possibly, containing one small non-complete k -regular bipartite graph if k does not divide n). The selection of highly disconnected bundle graphs is intentional; these graphs are in a sense extreme (within their category) and can challenge our methods.

Table 1 depicts the data (percentage of correctly recovered pairwise relations) from the execution of Borda and RSD on 18 distinct triplets of graph type and values⁷ for the parameters n and k . The data in the column labelled “random k -regular” show the average performance of Borda and RSD using 50 random bundle graphs. A different random permutation is used each time in order to assign elements to nodes. For graphs of the second and third type, one graph is used for each pair of values for n and k . For example, the data entries in the columns labeled “girth-6” and “copies of $K_{k,k}$ ” in the line with $k = 3$ and $n = 1001$ correspond to the performance of Borda and RSD on a girth-6 bundle graph which consists of 143 copies of the (7, 3)-bundle graph of Figure 1, and on a third-type graph that consists of 332 copies of $K_{3,3}$ and one more 3-regular graph with 5 nodes per side. Again, the data are average performance values from 50 executions; in each execution, a different random assignment of the elements to the nodes of the bundle graph is used.

| graph | | random k -regular | | girth-6 | | copies of $K_{k,k}$ | |
|-------|------|---------------------|------|---------|------|---------------------|------|
| k | n | Borda | RSD | Borda | RSD | Borda | RSD |
| 2 | 1002 | 73.3 | 62.7 | 73.5 | 60.3 | 66.8 | 56.8 |
| 3 | 1001 | 83.0 | 77.2 | 83.2 | 66.0 | 73.1 | 60.2 |
| 4 | 1001 | 87.5 | 86.8 | 87.7 | 68.7 | 77.1 | 62.2 |
| 6 | 1023 | 92.0 | 94.6 | 92.1 | 72.7 | 81.6 | 65.2 |
| 8 | 1026 | 94.2 | 97.2 | 94.1 | 72.8 | 84.3 | 66.5 |
| 12 | 1064 | 96.3 | 98.9 | 96.6 | 76.0 | 87.3 | 68.5 |

Table 1: Performance of Borda and RSD with perfect grading on different bundle graphs of similar size.

The results for Borda complement our theoretical analysis from Section 3. Indeed, the Borda-columns with bundle graphs of the second and third type indicate that the fraction of correctly recovered pairwise relations follows patterns of $1 - \mathcal{O}(1/k)$ and $1 - \mathcal{O}(1/\sqrt{k})$, respectively. Interestingly, the constants hidden in the \mathcal{O} notation are significantly smaller than the theoretical constants $16\sqrt{3\pi}$ and $48\sqrt{2\pi}$, respectively. The results from the execution of Borda on random bundle graphs shows a pattern of $1 - \mathcal{O}(1/k)$ as well, albeit with a slightly higher constant hidden in the \mathcal{O} notation. We believe that this can be proved by extending our analysis in Section 3. Even though we have not managed to prove that the quantity $\eta(G)$ is $\mathcal{O}(k^2)$ for these graphs, we strongly believe that this is the case.

⁶Consider the graph with a bipartition into an upper and lower set of nodes like in Figure 1.

⁷In all experiments reported here, n equals or is very close to 1000. This is because the results are essentially identical when significantly higher values of n are used (up to 10,000) and since the value of 1000 has allowed us to complete our experiments in a reasonable time frame.

RSD has poor performance on bundle graphs of the second and third type. This can be easily explained by recalling that these bundle graphs consist of small connected components. Even though all pairwise relations between elements assigned to nodes of the same component are correctly recovered, the vast majority of the pairwise relations are between elements that are assigned to different components. The probability that such a relation will be recovered correctly is only $1/2$. This explains the small percentages in the second and third RSD-columns.

In contrast, the first RSD-column (for random bundle graphs) shows a very interesting pattern. RSD is clearly worse than Borda for values of k up to 4 and becomes better as k increases further. Actually, this is more apparent in Figure 2 where Borda and RSD are compared in (n, k) -bundle graphs for all values of k from 2 up to 25 (and $n = 1000$). Each data point in Figure 2 corresponds to the average performance among 50 executions. Here, we can again recognize the $1 - \mathcal{O}(1/k)$ pattern for Borda that was observed in Table 1 and we further conjecture an even better pattern of $1 - \mathcal{O}(1/k^2)$ for RSD. Proving such a statement formally seems to be a challenging task.

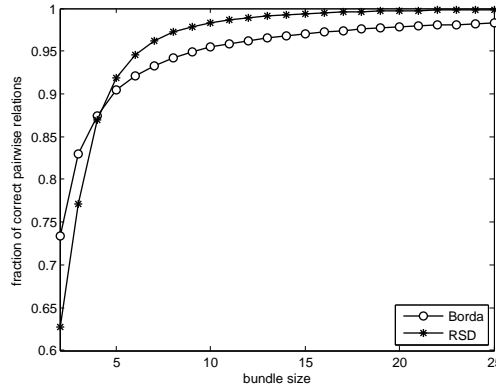


Figure 2: Borda vs. RSD with perfect grading and bundle size ranging from 2 to 25.

In a second set of experiments, we have studied imperfect grading. Now, we do not assume that the partial rankings are consistent to the ground truth any more. Instead, we have implemented generators of noisy rankings that may differ from the ground truth. In particular, we assume that each student has a quality that affects her position in the ground truth but also her ability to grade. First, the ground truth is the ranking of the elements in decreasing order of quality. Then, the ability of a student to rank the elements in a bundle depends on her quality q and is modelled by the following process. For every pair of elements a and b in the bundle that is ranked as $a \succ b$ in the ground truth, decide the correct pairwise relation with probability q and the opposite relation with probability $1 - q$. If this process creates a circular pairwise relation, we repeat the whole process from scratch. Otherwise, the output induces a ranking in the obvious way; this ranking is the one computed by the student. Clearly, a student of quality 1 will always produce a ranking that is consistent to the ground truth while a student of quality $1/2$ will produce a totally random ranking. This model was proposed by Condorcet in the 18th century; today, it is known as the Mallows model [17].

In our experiments, we use different noise levels that indicate the range of student qualities. For example, a noise level of 30% means that the qualities of the students are drawn uniformly at random from the interval $[0.7, 1]$. We use random bundle graphs for different values of k and besides Borda and RSD, we have also consider Markov chain-based aggregation methods. Dwork et al. [9] have studied a series of such methods; we describe the most powerful among them (even though we have experimented with a lot of variations of all the methods presented in [9]), which is known as MC4. MC4 defines a Markov chain (or random walk) over the elements and ranks them in decreasing order of their probabilities in the stationary distribution of this chain. The transition matrix of the Markov chain is defined as follows: when at an element a , pick an element b uniformly at random; if the number of partial rankings where

b is ranked above a is higher than the number of partial rankings where a is ranked above b , we have a transition to element b , otherwise we stay with element a .

Table 2 presents experimental data from the execution of Borda, RSD, and MC4 with random bundles for different values of the bundle size parameter and noise levels ranging from 50% to perfect grading. RSD has poor performance for high noise levels and small values of k . For non-zero noise levels, Borda has the best performance. MC4 and RSD are good choices only in the case of perfect grading, with RSD outperforming MC4 for the high values of $k = 8$ and 12. Overall, our experiments suggest that Borda is extremely robust. Note that there are some values missing from Table 2; this is due to the (exponential-time) implementation of Mallows generator which “takes forever” to come up with a set of non-circular pairwise relations that induces a ranking, when both k and the noise level are high.

| | $k = 5$ | | | $k = 8$ | | | $k = 12$ | | |
|-------------|---------|------|------|---------|------|------|----------|-------|-------|
| noise level | Borda | RSD | MC4 | Borda | RSD | MC4 | Borda | RSD | MC4 |
| 50 | 81.6 | 70.2 | 78.4 | 88.3 | 74.0 | 84.3 | ###.# | ###.# | ###.# |
| 40 | 84.9 | 75.1 | 81.2 | 91.1 | 80.1 | 86.5 | ###.# | ###.# | ###.# |
| 30 | 87.1 | 80.0 | 83.7 | 92.6 | 85.4 | 88.3 | ###.# | ###.# | ###.# |
| 20 | 88.6 | 84.2 | 86.0 | 93.5 | 89.6 | 89.8 | 95.5 | 92.2 | 92.6 |
| 10 | 89.6 | 88.4 | 88.8 | 93.9 | 93.2 | 91.2 | 96.1 | 95.7 | 93.6 |
| 0 | 90.4 | 92.0 | 92.7 | 94.2 | 97.2 | 96.4 | 96.2 | 98.9 | 97.8 |

Table 2: Performance of Borda, RSD, and MC4 with random bundle graphs of size 1000 and noise levels ranging from 50% to perfect grading.

We conclude by examining how sharply concentrated around the expectations the outcomes of the above experiments are. In Figure 3, we have plotted the fractions of correctly recovered pairwise relations obtained by Borda and RSD in the two extreme cases of perfect grading and noise level of 50%. Each figure contains data from 500 executions (a random bundle graph and a random element-to-node assignment defines each execution) with $n = 1000$ and $k = 8$. The spread of fractions of correctly recovered pairwise relations achieved by Borda is almost the same in both cases. In contrast, RSD has a very high spread when the noise level is high (observe the long and narrow form of the left plot in Figure 3) while it is only marginally better than Borda in the perfect grading case. In conclusion, Borda appears to be robust with respect to this metric as well.

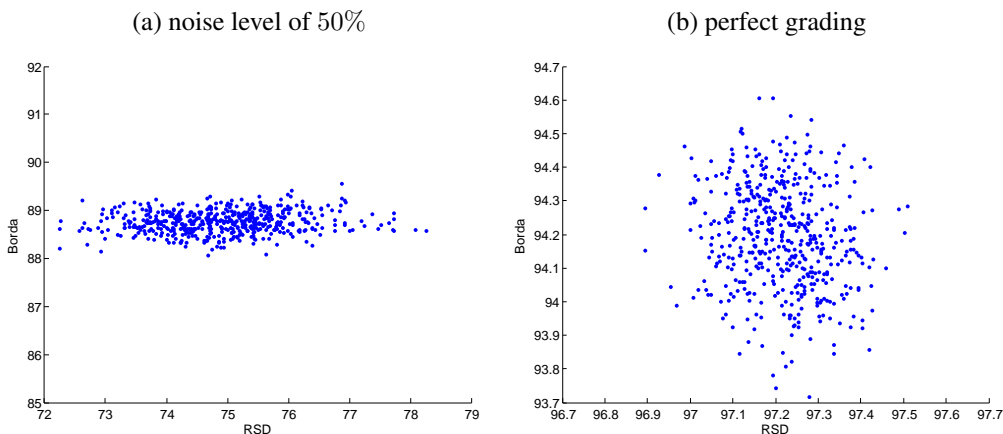


Figure 3: A comparison of Borda and RSD in 500 executions for two different noise levels ($n = 1000$, $k = 8$).

5 Discussion

Let us conclude by discussing some aspects of our work and possible future directions. Even though our analysis of Borda is targeted to the perfect grading case, we believe that our martingale-based arguments could be extended to handle imperfect grading under the Mallows noise model that we use in our experiments. This requires taking care of even more dependencies but we are confident that martingale theory will be useful here as well. We plan to consider extending our analysis in this direction in follow-up versions of this work.

Besides Borda, we have attempted a theoretical analysis of RSD as well. Here, our starting point has been to exploit the developments in the degree/diameter problem [18] and use a diameter-5 low-degree bipartite graph as a bundle graph. The important property this graph has is that for every pair of nodes u and v of the node set U , these nodes either have a common neighbour in V or there is another (intermediate) node z in U that has a common neighbour with u and another common neighbour with v . Hence, in the perfect grading case, the pairwise relation between the elements a and b that are assigned to nodes u and v can be indirectly learnt during the serial phase through the pairwise relations of a and b with the element c that is assigned to node z , provided that c is ranked between a and b in the true ranking. Furthermore, if the bundle graph had more than one disjoint paths between any pair of nodes in U (and more than one intermediates for any pair of nodes), the probability that the relation between two elements can be learnt correctly would be very high, provided that these elements have a relatively large rank difference in the true ranking. Unfortunately, even though some theoretical guarantees can indeed be formally proved in this way, the bundle graphs required have degree that strongly depends on the number of elements. So, this approach fails to explain the performance of RSD that we observed experimentally. Instead, one should reason about pairwise relations that can be learnt indirectly through long chains of intermediate elements. Unfortunately, exploiting such arguments seems elusive at this point.

In our experimental work, we have implemented and tested many more aggregation rules than the ones presented in Section 4. These include rules that put more weight on the partial rankings of low-rank (i.e., good) students. Such rules are usually defined using Markov chains that are variations of PageRank [21] (such as the PeerRank method in [26]), where the idea is that the confidence about the quality of a student depends on the performance of her graders (and this is reflected in the definition of the transition matrix of the Markov chain). Unfortunately, we have not observed any significant improvement compared to the rules considered in Section 4. We believe that this can be explained by the fact that k is a small constant.

In future work, we would also like to consider more realistic noise models that generalize Mallows (see, e.g., [11, 16, 23]) and ranking models that are inherently associated with cardinal utilities such as the generalized random utility model of Soufiani et al. [25] (see also the book of [15] and the references therein). Of course, it is important to perform real-world experiments (with students in the classroom or with participants in real MOOCs, if possible) in order to justify our methods and determine the noise model that is closest to practice.

Acknowledgements. We would like to thank Stratis Gallopoulos and Steve Vavasis for discussions on early stages of this work, and Panagiotis Kanellopoulos and Nisarg Shah for technical comments and remarks.

References

- [1] A. Abdulkadiroğlu and T. Sönmez. Random serial dictatorship and the core from random endowments in house allocation problems. *Econometrica*, 66:689–702, 1998.

- [2] N. Alon, F. A. Fischer, A. D. Procaccia, and M. Tennenholtz. Sum of us: strategyproof selection from the selectors. In *Proceedings of the 13th Conference on Theoretical Aspects of Rationality and Knowledge (TARK)*, pages 101–110, 2011.
- [3] K. Azuma. Weighted sums of certain dependent random variables. *Tôhoku Mathematical Journal*, 19:357–367, 1967.
- [4] J. C. Borda. Mémoire sur les élections au scrutin. *Histoire de l'Académie Royale des Sciences*, 1781.
- [5] I. Caragiannis, A. D. Procaccia, and N. Shah. When do noisy votes reveal the truth? In *Proceedings of the 14th ACM conference on Electronic commerce (EC)*, pages 143–160, 2013.
- [6] F. Chierichetti and J. M. Kleinberg. Voting with limited information and many alternatives. In *Proceedings of the 23rd Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1036–1055, 2012.
- [7] V. Conitzer and T. Sandholm. Vote elicitation: complexity and strategy-proofness. In *Proceedings of the 18th AAAI Conference on Artificial Intelligence (AAAI)*, pages 392–397, 2002.
- [8] J. L. Doob. *Stochastic Processes*. Wiley and Sons, New York, 1953.
- [9] C. Dwork, R. Kumar, M. Naor, and D. Sivakumar. Rank aggregation methods for the web. In *Proceedings of the 10th International World Wide Web Conference (WWW)*, pages 613–622, 2001.
- [10] Y. Filmus and J. Oren. Efficient voting via the top- k elicitation scheme: a probabilistic approach. In *Proceedings of the 15th Conference on Economics and Computation (EC)*, pages 295–312, 2014.
- [11] M. A. Fligner and J. S. Verducci. Distance based ranking models. *Journal of the Royal Statistical Society B*, 48(3):359–369, 1986.
- [12] W. Hoeffding. Probability inequalities for the sum of bounded random variables. *Journal of the American Statistical Association*, 58:13–30, 1963.
- [13] R. Holzman and H. Moulin. Impartial nominations for a prize. *Econometrica*, 81:173–196, 2013.
- [14] E. Law and L. von Ahn. *Human Computation*. Synthesis Lecture on Artificial Intelligence and Machine Learning. Morgan & Claypool, 2011.
- [15] T. Y. Liu. *Learning to Rank for Information Retrieval*. Springer-Verlag, 2011.
- [16] T. Lu and C. Boutilier. Learning Mallows models with pairwise preferences. In *Proceedings of the 28th International Conference on Machine Learning (ICML)*, pages 145–152, 2011.
- [17] C. L. Mallows. Non-null ranking models. *Biometrika*, 44:114–130, 1957.
- [18] M. Miller and J. Sirán. Moore graphs and beyond: a survey of the degree/diameter problem. *The Electronic Journal of Combinatorics*, Dynamic Survey, 2013.
- [19] M. Mitzenmacher and E. Upfal. *Probability and Computing – Randomized Algorithms and Probabilistic Analysis*. Cambridge University Press, 2005.
- [20] R. Motwani and P. Raghavan. *Randomized Algorithms*. Cambridge University Press, New York, 1995.

- [21] L. Page, S. Brin, R. Motwani, and T. Winograd. The PageRank citation ranking: bringing order to the web. Stanford InfoLab technical report, 1999.
- [22] C. Piech, J. Huang, Z. Chen, C. Do, A. Ng, and D. Koller. Tuned models of peer assessment in MOOCs. In *Proceedings of the 6th International Conference on Educational Data Mining (EDM)*, pages 153–160, 2013.
- [23] K. Raman and T. Joachims. Methods for ordinal peer grading. In *Proceedings of the 20th ACM Conference on Knowledge Discovery and Data Mining (KDD)*, page forthcoming, 2014.
- [24] N. B. Shah, J. K. Bradley, A. Parekh, M. Wainwright, and K. Ramchandran. A case for ordinal peer-evaluation in MOOCs. In *Neural Information Processing Systems (NIPS): Workshop on Data Driven Education*, 2013.
- [25] H. A. Soufiani, D. C. Parkes, and L. Xia. Preference elicitation for general random utility models. In *Proceedings of the 29th Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 596–605, 2013.
- [26] T. Walsh. The PeerRank method for peer assessment. In *Proceedings of the 21st European Conference on Artificial Intelligence (ECAI)*, pages 909–914, 2014.
- [27] L. Xia and V. Conitzer. Determining possible and necessary winners given partial orders. *Journal of Artificial Intelligence Research*, 41:25–67, 2011.