

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## Quantifying consensus of rankings based on $q$ -support patterns

Zhengui Xue, Zhiwei Lin\*, Hui Wang, Sally McClean

*School of Computing, University of Ulster, United Kingdom*

---

### Abstract

Rankings, representing preferences over a set of candidates, are widely used in many applications, e.g., group decision making and information retrieval. Rankings may be obtained by different agents (humans or systems). It is often necessary to evaluate consensus of obtained rankings from multiple agents, as a measure of consensus provides insights into the rankings. Moreover, a consensus measure could provide a quantitative basis for comparing groups and for improving a ranking system. Existing studies on consensus measurement are insufficient, since they did not evaluate consensus among most rankings or consensus with respect to specific preference patterns. In this paper, a novel consensus quantifying approach, without the use of correlation or distance functions as in existing studies of consensus, is proposed based on the concept of  $q$ -support patterns, which represent the commonality embedded in a set of rankings. A pattern is regarded as a  $q$ -support pattern if it is included by at least  $q$  rankings in the ranking set. A method for detecting outliers in a set of rankings is naturally derived from the proposed consensus quantifying approach. Experimental studies are conducted to demonstrate the effectiveness of the proposed approach.

*Keywords:* Rankings, consensus, support patterns, outlier detection

---

\*Corresponding author

*Email addresses:* zhenguixue@gmail.com (Zhengui Xue), z.lin@ulster.ac.uk (Zhiwei Lin), h.wang@ulster.ac.uk (Hui Wang), si.mcclean@ulster.ac.uk (Sally McClean)

## 1. Introduction

Extensive studies have been carried out in social science to measure group cohesion, in order to gain insight into the factors affecting group cohesion and further promote higher group consistency (see, e.g., [21, 7, 35, 10]). In artificial intelligence, rankings have been widely used to represent the preferences of agents (humans or systems) over a set of candidates in many information systems, such as group decision making [27, 34, 43] and information retrieval [22, 29, 33]. It is important to evaluate the degree to which the rankings obtained by different agents agree, as it would help to understand the obtained rankings. Quantifying the *consensus* of the obtained rankings can provide an accurate evaluation about the overall agreement. It is also a quantitative indicator for comparing consensus between groups (e.g., two sets of rankings) [4] or for further improving the ranking systems. For example, in group decision making, if the consensus score is extremely low, it is necessary for the experts to adjust their rankings in order to reach an agreement [27]. However, to the best of our knowledge, there are only a few existing studies [1, 2, 4, 13, 16, 17] on consensus evaluation for a set of rankings.

In the literature, rank *correlation* and *distance* functions, such as Kendall's  $\tau$  [26] and Spearman's  $\rho$  [37], are used to measure the correlation and disagreement of two rankings. Kendall's  $\tau$  measures the correlation of two rankings by considering their concordant and discordant pairs, and Spearman's  $\rho$  evaluates the rank correlation by taking into account the positions of the items in two rankings. The Kemeny distance [25] is extended to measure pairwise disagreements in two rankings. For a set with more than two rankings, the related concepts are *consensus* and *diversity* of rankings. Consensus is also used interchangeably for *cohesiveness* [2]. Existing approaches measure consensus of rankings by considering the similarity of preferences in a group based on rank correlation functions. One typical approach as discussed in [2] is to calculate the similarity for each pair of rankings based on correlation functions, such as Kendall's  $\tau$  and Spearman's  $\rho$ , and then aggregate the obtained results. *Diversity* and *consensus* are considered as two opposite concepts about rankings in social choice theory [24]. Research was carried out to measure the diversity of a ranking set based on distance

1  
2  
3  
4  
5  
6  
7  
8  
9 functions (see [16]). These existing studies are not sufficient in evaluating the overall  
10 consensus of a ranking set. It is difficult to use the rank correlation or distance functions  
11 based approaches to completely quantify the level of consensus for a set of rankings.  
12 As pointed out in [12], the pairwise comparison reflects the degree of commonality in  
13 two rankings, and consequently the aggregated result of the pairwise comparisons is  
14 not informative enough to tell the degree to which the ranking set agrees. In reality, it  
15 is often the case that certain preference patterns are embedded in most of the rankings  
16 obtained for a task. The existing work cannot tell the degree to which preferences over  
17 candidates are shared by the majority of the rankings. In addition, they did not provide  
18 a solution to identifying the majority of rankings in order to filter irrelevant results in  
19 the ranking set, which could play an important role in modern information systems.  
20 For instance, in query expansion [6], it is reasonable to expand a query for ‘film’ to  
21 its relevant query ‘movie’ in the context of entertainment, but not in the context of ‘a  
22 thin coat or layer’. It is impossible to manually check if the expansions from the source  
23 query are consistent as there is no ground truth available and moreover the meaning of  
24 the queries may evolve from time to time (e.g, ‘apple’ in fruit context to the context  
25 of cooperation). Therefore, using the rankings obtained from the expansions to under-  
26 stand the extent to which the query expansions provide high level of consistency is key  
27 to provide good search results.

28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39 This paper studies the consensus degree of a ranking set from a different perspective  
40 to provide a full picture on the degree to which a set of rankings mutually agree. A novel  
41 framework is proposed to analyze consensus of rankings by considering the common  
42 patterns embedded in a ranking set. A new concept of  $q$ -support patterns is introduced  
43 to represent how common patterns are embedded in rankings, by which the preferences  
44 of a group over candidates can be expressed at a subtle and fine-grained level. A pattern  
45 is regarded as a  $q$ -support pattern if it is included by at least  $q$  rankings in the ranking  
46 set. Thus, a  $q$ -support pattern represents the partial coverage of the pattern by rankings,  
47 where the integer  $q$  can be specified as needed when a ranking system is evaluated. The  
48 consensus degree of rankings is quantified based on  $q$ -support patterns. Compared with  
49 the existing work based on correlation or distance functions, this new approach gives a  
50 finer characterization and quantification of the commonality embedded in the rankings.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

The contributions of this paper are: (1) a new representation of the commonality within a set of rankings,  $q$ -support pattern, is proposed; (2) a new framework (non-distance and non-correlation) for quantifying consensus with  $q$ -support patterns is introduced; (3) an efficient algorithm is developed to calculate consensus scores and characterize the set of  $q$ -support patterns; (4) consensus scores are defined for each ranking to reflect its relationship with the other rankings, which can be used to detect outliers in a ranking set; (5) extensive experiments have been conducted to show the effectiveness and usefulness of the proposed approach.

The rest of the paper is organized as follows. In Section 2 related work on the comparison of two rankings and the measure of consensus and diversity of a ranking set is reviewed. In Section 3 the  $q$ -support pattern of rankings is formulated and consensus scores are defined based on it. An algorithm is then introduced to calculate ranking consensus. In Section 4 weighted consensus scores are defined. In Section 5, an outlier detection method is developed. Section 6 gives experimental studies to evaluate the proposed approach. Section 7 concludes this paper.

## 2. Related work

**Rank correlation and distance functions.** Historically developed by Maurice Kendall in 1938 [26], Kendall’s  $\tau$  measures the correlation between two rankings by considering the numbers of pairwise items ranked in same orders and in opposite orders. Suppose that we consider rankings over candidates  $\{\sigma_1, \sigma_2, \dots, \sigma_n\}$ . A ranking is an ordered list in which items in higher positions are more preferred than items in lower positions. Let  $\pi(\cdot, \cdot)$  be the position function. The function  $\pi(\sigma_i, \mathbf{r}_l)$  returns the position of item  $\sigma_i$  in ranking  $\mathbf{r}_l$ . Kendall’s  $\tau$  for two rankings  $\mathbf{r}_l$  and  $\mathbf{r}_z$  is

$$\tau(\mathbf{r}_l, \mathbf{r}_z) = \frac{\sum_{\substack{i,j \in \{1, \dots, n\} \\ i < j}} \text{sgn}(\pi(\sigma_i, \mathbf{r}_l) - \pi(\sigma_j, \mathbf{r}_l)) \text{sgn}(\pi(\sigma_i, \mathbf{r}_z) - \pi(\sigma_j, \mathbf{r}_z))}{n(n-1)/2}.$$

This coefficient is in the range  $-1 \leq \tau(\mathbf{r}_l, \mathbf{r}_z) \leq 1$ , where value 1 corresponds to the case that the two rankings are in the same order and value  $-1$  indicates that one ranking is in the reverse order of the other.

Spearman's  $\rho$  proposed by Charles Spearman in 1904 [37] is defined based on the positions of each item in two rankings as follows

$$\rho(\mathbf{r}_l, \mathbf{r}_z) = \frac{\sum_{i=1}^n (\pi(\sigma_i, \mathbf{r}_l) - \bar{\pi}_l)(\pi(\sigma_i, \mathbf{r}_z) - \bar{\pi}_z)}{\sqrt{\sum_{i=1}^n (\pi(\sigma_i, \mathbf{r}_l) - \bar{\pi}_l)^2 \sum_{i=1}^n (\pi(\sigma_i, \mathbf{r}_z) - \bar{\pi}_z)^2}},$$

where  $\bar{\pi}_l = \frac{1}{n} \sum_{i=1}^n \pi(\sigma_i, \mathbf{r}_l)$  and  $\bar{\pi}_z = \frac{1}{n} \sum_{i=1}^n \pi(\sigma_i, \mathbf{r}_z)$ . Similarly, this coefficient satisfies  $-1 \leq \rho(\mathbf{r}_l, \mathbf{r}_z) \leq 1$ .

These rank correlation functions do not take into account the varying relevance of ranked items in different positions. They are not suitable for evaluating the rankings where items at the top of a ranking are much more important than those at the bottom [15]. Further studies on weighted rank correlation were carried out extensively based on these two functions [8, 23, 28, 36, 39, 41, 42]. More reasonable variants of rank correlation functions were also proposed in the literature [14, 19, 20, 38].

Distance metrics have been used to analyze ranking data. One of the most widely used distance functions to measure rankings is the Kemeny distance [25]. It is defined as the sum of pairs where the ranking preferences disagree. One can refer to [3, 31, 11] for more information about the commonly used distance metrics.

**Measuring consensus and diversity of rankings.** For a ranking set with the number of rankings greater than two, work [4] is known as the first study to define a consensus measure as a function mapping linear orders (i.e, rankings without ties) to a number between 0 and 1. Kendall's coefficient of concordance was introduced in [4] as a measure of consensus of a ranking set. Given a set of rankings  $\mathcal{R} = \{\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N\}$  over candidates  $\{\sigma_1, \sigma_2, \dots, \sigma_n\}$ , the total positions of the candidates in all rankings need to be calculated first, which are  $\sum_{l=1}^N \pi(\sigma_i, \mathbf{r}_l), i = 1, \dots, n$ . Kendall's coefficient of concordance is defined based on the deviations of the total positions from their mean as

$$W = \frac{12}{N^2(n^3 - n)} \sum_{i=1}^n \left( \sum_{l=1}^N \pi(\sigma_i, \mathbf{r}_l) - \frac{1}{n} \sum_{i=1}^n \sum_{l=1}^N \pi(\sigma_i, \mathbf{r}_l) \right)^2,$$

where the term  $\frac{12}{N^2(n^3 - n)}$  is for normalization.

García-Lapresta and Pérez-Román [16] extended the work [4] by considering weak

1  
2  
3  
4  
5  
6  
7  
8  
9 orders (i.e., ranking with ties). A measure based on a weighted Kemeny distance was  
10 introduced. In [2], it was discussed that one prominent approach of constructing a con-  
11 sensus or diversity measure is to make pairwise comparisons of the rankings with a rank  
12 correlation or distance function, such as the functions introduced in the above section,  
13  
105 and then aggregate the comparison results. Thus, two key issues with this approach are  
14 the choice of a proper pairwise comparison metric and the utilization of an aggregation  
15 method. Kendall's  $\tau$  was used to compare the similarity of each pair of rankings in  
16 [2], and the consensus measure of a ranking set was constructed by taking the average  
17  
18 of the comparison results. Studies with more reasonable similarity or distance metrics  
19  
20 were carried out in [1, 17, 18, 13]. Karpov [24] considered to aggregate the comparison  
21  
22 of the comparison results. Studies with more reasonable similarity or distance metrics  
23  
24 were carried out in [1, 17, 18, 13]. Karpov [24] considered to aggregate the comparison  
25  
26 results with a geometric mean aggregator.

27 Although these studies discussed different aspects of consensus measures, they are  
28 still inefficient in the assessment of overall consensus of a ranking set. In information  
29  
30 systems, it is often the case that certain preference patterns are embedded in most of  
31  
32 the rankings. The existing studies based on rank correlation and distance functions did  
33  
34 not provide a full picture about this kind of common patterns. They cannot quantify the  
35  
36 degree to which preferences over candidates are shared by the majority of the rankings.  
37  
38 To solve this problem, this paper proposes a concept of  $q$ -support patterns to represent  
120 the commonality in a ranking set and the consensus is quantified based on the  $q$ -support  
39  
40 patterns.

### 41 42 **3. Quantifying consensus with $q$ -support patterns**

43  
44 This section first defines the  $q$ -support patterns and the consensus scores of a rank-  
45  
46 ing set. Then, an algorithm is presented to calculate the consensus scores by utilizing  
47  
48 matrices to represent the  $q$ -support patterns.  
125

#### 49 50 *3.1. $q$ -support patterns*

51  
52 Let  $\mathcal{C} = \{\sigma_1, \sigma_2, \dots, \sigma_n\}$  be a set of  $n$  candidates to be ranked. A ranking  $\mathbf{r}_l =$   
53  
54  $(r_{l_1}, r_{l_2}, \dots, r_{l_m})$  is an ordered list in which item  $r_{l_i} \in \mathcal{C}$  is more preferred than item  
55  
56  $r_{l_j} \in \mathcal{C}$  for  $i < j$ . Given two items  $\sigma_x$  and  $\sigma_y \in \mathcal{C}$ , if there exists  $i \leq j$  such that  
57  
58

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

130  $r_{l_i} = \sigma_x$  and  $r_{l_j} = \sigma_y$ , we write  $\sigma_x\sigma_y \sqsubset \mathbf{r}_l$ ; otherwise  $\sigma_x\sigma_y \not\sqsubset \mathbf{r}_l$ . Specially, if  
 $\sigma_x = \sigma_y$ ,  $\sigma_x\sigma_x \sqsubset \mathbf{r}_l$  simply means that item  $\sigma_x$  is included in ranking  $\mathbf{r}_l$ , also written  
as  $\sigma_x \sqsubset \mathbf{r}_l$ .

It is usually the case that most of the rankings obtained for a task share certain commonality. Suppose that there is a set of rankings  $\mathcal{R} = \{\mathbf{r}_1 = (a, b, c, d, e, f), \mathbf{r}_2 = (b, a, c, d, e, f), \mathbf{r}_3 = (a, b, c, e, d, f), \mathbf{r}_4 = (c, b, d, e, f, g)\}$ . It can be seen that item  $a$  and the pairwise items  $bc$  are common patterns for most of the rankings, but not for all the rankings in  $\mathcal{R}$  (e.g.,  $bc \sqsubset \mathbf{r}_1, bc \sqsubset \mathbf{r}_2, bc \sqsubset \mathbf{r}_3$ , but  $bc \not\sqsubset \mathbf{r}_4$ ). These patterns, partially included in a set of rankings, show the **extent** to which the rankings agree. Therefore, it is necessary to consider these patterns to understand the consensus level in a set of rankings. As such, we define the following  $q$ -support patterns for a ranking set.

**Definition 1** ( $q$ -support patterns). *Consider a set of  $N$  rankings  $\mathcal{R} = \{\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N\}$  over candidate set  $\mathcal{C} = \{\sigma_1, \sigma_2, \dots, \sigma_n\}$ . For  $\sigma_x$  and  $\sigma_y \in \mathcal{C}$ , we have the following subset  $\mathcal{R}'(\sigma_x, \sigma_y) \subseteq \mathcal{R}$*

$$\mathcal{R}'(\sigma_x, \sigma_y) = \{\mathbf{r}_z | \sigma_x\sigma_y \sqsubset \mathbf{r}_z, \mathbf{r}_z \in \mathcal{R}\}. \quad (1)$$

Let  $q \in (0, N]$  be an integer. The pattern  $\sigma_x\sigma_y$  is a  $q$ -support pattern of  $\mathcal{R}$ , denoted by  $\sigma_x\sigma_y \stackrel{q}{\sqsubset} \mathcal{R}$ , if the size of  $\mathcal{R}'(\sigma_x, \sigma_y)$  satisfies  $|\mathcal{R}'(\sigma_x, \sigma_y)| \geq q$ ; otherwise  $\sigma_x\sigma_y \not\stackrel{q}{\sqsubset} \mathcal{R}$ . If  $\sigma_x = \sigma_y$ ,  $\sigma_x\sigma_x \stackrel{q}{\sqsubset} \mathcal{R}$  indicates that item  $\sigma_x$  is a single  $q$ -support item of  $\mathcal{R}$ , also written as  $\sigma_x \stackrel{q}{\sqsubset} \mathcal{R}$ .

The notation  $\sigma_x\sigma_y \stackrel{q}{\sqsubset} \mathcal{R}$  means that  $\sigma_x\sigma_y$  occurs in at least  $q$  rankings in  $\mathcal{R}$ . We use  $\mathbf{S}_1(q)$  and  $\mathbf{S}_2(q)$  to respectively denote the set of the single  $q$ -support items and the set of the pairwise  $q$ -support patterns, i.e.,

$$\mathbf{S}_1(q) = \{\sigma_x | \sigma_x \stackrel{q}{\sqsubset} \mathcal{R}, \sigma_x \in \mathcal{C}\} \quad (2)$$

$$\mathbf{S}_2(q) = \{\sigma_x\sigma_y | \sigma_x\sigma_y \stackrel{q}{\sqsubset} \mathcal{R}, \sigma_x \neq \sigma_y, \sigma_x \in \mathcal{C}, \sigma_y \in \mathcal{C}\}. \quad (3)$$

The set  $\mathbf{S}_1(q)$  is important in the evaluation of incomplete rankings, where not all the candidates under consideration are ranked in the rankings. It gives the items with more preferences among the candidates, which are ranked in at least  $q$  rankings. The set  $\mathbf{S}_2(q)$  collects the preference orders embedded in at least  $q$  rankings.

### 3.2. Consensus scores

The  $q$ -support patterns describe how common patterns are embedded in rankings. This section first defines individual consensus scores for a ranking  $\mathbf{r}_l \in \mathcal{R}$  based on the  $q$ -support patterns. Then, the overall consensus scores are introduced for the ranking set  $\mathcal{R}$ . The relative consensus degree that a ranking  $\mathbf{r}_l$  shares with the others can be revealed by the individual and the overall consensus scores. In Section 5, it shows that this information can be used in the detection of an outlier from a ranking set.

The following individual consensus scores are defined for a ranking  $\mathbf{r}_l$ .

**Definition 2** (Individual consensus scores). *For a ranking  $\mathbf{r}_l = (r_{l_1}, r_{l_2}, \dots, r_{l_m}) \in \mathcal{R}$ , the sets of the single  $q$ -support items and the pairwise  $q$ -support patterns are defined as*

$$\mathbf{S}_1^{\mathbf{r}_l}(q) = \{r_{l_i} | r_{l_i} \stackrel{q}{\sqsubseteq} \mathcal{R}, i \in \{1, 2, \dots, m\}\} \quad (4)$$

$$\mathbf{S}_2^{\mathbf{r}_l}(q) = \{r_{l_i} r_{l_j} | r_{l_i} r_{l_j} \stackrel{q}{\sqsubseteq} \mathcal{R}, i, j \in \{1, 2, \dots, m\}, i < j\}. \quad (5)$$

The individual consensus scores of  $\mathbf{r}_l$  are

$$\kappa_1^{\mathbf{r}_l}(q) = \frac{1}{N_1^{\mathbf{r}_l}} |\mathbf{S}_1^{\mathbf{r}_l}(q)| \quad (6)$$

$$\kappa_2^{\mathbf{r}_l}(q) = \frac{1}{N_2^{\mathbf{r}_l}} |\mathbf{S}_2^{\mathbf{r}_l}(q)|, \quad (7)$$

where  $N_1^{\mathbf{r}_l} = m$  and  $N_2^{\mathbf{r}_l} = \frac{m(m-1)}{2}$  respectively represent the number of the ranked items and the number of the pairwise patterns of  $\mathbf{r}_l$ .

**Definition 3** (Overall consensus scores). *For a ranking set  $\mathcal{R}$  with the individual consensus scores defined as (6) and (7), the overall consensus scores of  $\mathcal{R}$  are*

$$\bar{\kappa}_1(q) = \frac{1}{N} \sum_{l=1}^N \kappa_1^{\mathbf{r}_l}(q) \quad (8)$$

$$\bar{\kappa}_2(q) = \frac{1}{N} \sum_{l=1}^N \kappa_2^{\mathbf{r}_l}(q). \quad (9)$$

The individual consensus scores measure the proportions of the preference patterns of  $\mathbf{r}_l$  embedded in at least  $q$  rankings, where  $\kappa_1^{\mathbf{r}_l}(q)$  measures consensus in terms of single  $q$ -support items and  $\kappa_2^{\mathbf{r}_l}(q)$  measures consensus in terms of pairwise  $q$ -support



patterns. The overall consensus scores give the average proportions and they are used to evaluate the consensus degree of a whole ranking set. Note that a  $q$ -support pattern depicts the commonality embedded in at least  $q$  rankings in a ranking set. The choice of  $q$  in the consensus evaluation depends on the specific need in the evaluation of ranking data. For example, many information systems may expect that ranked patterns are supported by at least half of the experts, and the value of  $q$  can be set to  $\lceil \frac{N}{2} \rceil$  for this case. In addition, by studying the consensus degree based on different values of  $q$ , a more comprehensive understanding about the ranking set can be obtained, as different values of  $q$  reflect the extents of different partial coverage of the patterns embedded in rankings.

The consensus scores have the following property.

**Property 1.** The overall consensus scores satisfy

$$0 \leq \bar{\kappa}_1(q) \leq 1 \quad (10)$$

$$0 \leq \bar{\kappa}_2(q) \leq 1. \quad (11)$$

The score  $\bar{\kappa}_1(q) = 0$  if and only if arbitrary  $q$  rankings in  $\mathcal{R}$  share no common item, and  $\bar{\kappa}_1(q) = 1$  if and only if every ranked item of all the rankings is shared by at least  $q$  rankings. Similarly,  $\bar{\kappa}_2(q) = 0$  if and only if arbitrary  $q$  rankings in  $\mathcal{R}$  share no common pairwise pattern, and  $\bar{\kappa}_2(q) = 1$  if and only if every pairwise preference pattern of all the rankings is embedded in at least  $q$  rankings.

### 3.3. An efficient algorithm for quantifying consensus

In this section, a matrix representation is introduced to represent the  $q$ -support patterns, as shown in Theorem 1, which implies an algorithm for calculating the consensus scores.

**Theorem 1.** Consider a set of  $N$  rankings  $\mathcal{R} = \{\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N\}$  over candidates  $\mathcal{C} = \{\sigma_1, \sigma_2, \dots, \sigma_n\}$ . For a ranking  $\mathbf{r}_l = (r_{l_1}, r_{l_2}, \dots, r_{l_m}) \in \mathcal{R}$  and  $\forall \mathbf{r}_z = (r_{z_1}, r_{z_2}, \dots, r_{z_u}) \in \mathcal{R}$ , with the position function

$$\pi(r_{l_i}, \mathbf{r}_z) = \begin{cases} 0, & \text{if } r_{l_i} \notin \mathbf{r}_z \\ p, & \text{if } r_{l_i} = r_{z_p} \end{cases} \quad (12)$$

and the Heaviside function

$$H(x) = \begin{cases} 1, & \text{if } x > 0 \\ 0, & \text{otherwise,} \end{cases} \quad (13)$$

we define

$$f(r_{l_i}, r_{l_j}) = \begin{cases} \sum_{z=1}^N H(\pi(r_{l_i}, \mathbf{r}_z)), & \text{if } i = j \\ \sum_{z=1}^N H(\pi(r_{l_j}, \mathbf{r}_z) - \pi(r_{l_i}, \mathbf{r}_z)) H(\pi(r_{l_i}, \mathbf{r}_z)), & \text{otherwise} \end{cases} \quad (14)$$

and matrix  $\mathbf{A}^{\mathbf{r}_l} = (A^{\mathbf{r}_l}[j, i]) \in \mathbb{R}^{m \times m}$  as

$$A^{\mathbf{r}_l}[j, i] = \begin{cases} 1, & \text{if } i \leq j \text{ and } f(r_{l_i}, r_{l_j}) \geq q \\ 0, & \text{otherwise.} \end{cases} \quad (15)$$

Then, we have

$$\kappa_1^{\mathbf{r}_l}(q) = \frac{1}{N_1^{\mathbf{r}_l}} \text{tr}(\mathbf{A}^{\mathbf{r}_l}) \quad (16)$$

$$\kappa_2^{\mathbf{r}_l}(q) = \frac{1}{N_2^{\mathbf{r}_l}} (\mathbf{e}^T \mathbf{A}^{\mathbf{r}_l} \mathbf{e} - \text{tr}(\mathbf{A}^{\mathbf{r}_l})), \quad (17)$$

where  $\mathbf{e} = [1, 1, \dots, 1]^T$  is an  $m$ -row vector of all ones.

*Proof.* By (12), it can be known that  $\pi(r_{l_i}, \mathbf{r}_z)$  gives the position of item  $r_{l_i}$  in  $\mathbf{r}_z$ . From the definition of  $f(r_{l_i}, r_{l_j})$ , it can be seen that  $f(r_{l_i}, r_{l_j})$  counts the number of rankings  $\forall \mathbf{r}_z \in \mathcal{R}$  satisfying  $r_{l_i} r_{l_j} \sqsubset \mathbf{r}_z$ . Thus, the entry  $A^{\mathbf{r}_l}[j, i] = 1$  represents  $r_{l_i} r_{l_j} \stackrel{q}{\sqsubset} \mathcal{R}$ . Moreover, note that  $\mathbf{e}^T \mathbf{A}^{\mathbf{r}_l} \mathbf{e}$  gives the sum of the all entries in matrix  $\mathbf{A}^{\mathbf{r}_l}$ .

Therefore, the result of (16) and (17) can be further obtained based on Definition 2.  $\square$

The matrix  $\mathbf{A}^{\mathbf{r}_l}$  provides a proper representation of the  $q$ -support patterns in  $\mathbf{r}_l$ . This representation can further facilitate the analysis of the commonality that individual rankings share with the others. Based on Theorem 1, we develop Algorithm 1 to calculate the consensus scores and characterize the  $q$ -support patterns more efficiently.

In Algorithm 1, when a ranking  $\mathbf{r}_l$  is considered, for a pattern  $r_{l_i} r_{l_j}$  embedded in the ranking, there is no need to judge if the pattern is a  $q$ -support pattern by checking all the rankings in some cases. Suppose  $q = \lceil \frac{2N}{3} \rceil$ , which means that we consider

1  
2  
3  
4  
5  
6  
7  
8  
9  $r_{l_i}r_{l_j}$  as a common pattern if it is contained by at least two third of the rankings.

10 For  $\mathbf{r}_l$ , if  $r_{l_i}r_{l_j}$  is a  $\lceil \frac{2N}{3} \rceil$ -support pattern, it must be included by one of the rankings

11  
12 210  $\mathbf{r}_x \in \{\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_{\lfloor \frac{N}{3} \rfloor + 1}\}$ . If  $r_{l_i}r_{l_j}$  is not included by one of the first  $\lfloor \frac{N}{3} \rfloor + 1$  rankings  
13 of the ranking set,  $r_{l_i}r_{l_j}$  cannot be a  $q$ -support pattern and  $A^{r_l}[j, i]$  should be zero.

14 Thus, we do not need to calculate  $A^{r_l}[j, i]$  by always checking all the rankings. Line

15 7 in Algorithm 1 checks if  $r_{l_i}r_{l_j}$  of  $\mathbf{r}_l$  is included by a ranking  $\mathbf{r}_x$  for which matrix

16  $\mathbf{A}^{r_x}$  has already been constructed. If the number of the rankings whose corresponding

17 215 matrix is not constructed is greater than  $q$ , we look for  $\mathbf{r}_x$  in the previously considered

18 rankings  $\{\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_{l-1}\}$ ; Otherwise if the number of rankings not yet considered is

19 less than  $q$ , we only need to check if there is an  $\mathbf{r}_x$  in the first  $N - q + 1$  rankings. As

20 shown on Lines 8 and 9, if  $r_{l_i}r_{l_j}$  has been considered in a constructed matrix for  $\mathbf{r}_x$ ,

21 it is not necessary to recalculate the corresponding entry of the current matrix  $\mathbf{A}^{r_l}$  and

22 220 the entry is equal to that of  $\mathbf{A}^{r_x}$  corresponding to the pattern. Otherwise, as on Line

23 10, only when the number of the rankings  $\{\mathbf{r}_l, \mathbf{r}_{l+1}, \dots, \mathbf{r}_N\}$  is no less than  $q$ , pattern

24  $r_{l_i}r_{l_j}$  has the possibility to be a  $q$ -support pattern and we need to check the remaining

25 rankings to see if the pattern is a  $q$ -support ranking. In this way, the computation cost

26 can be significantly reduced. From Lines 11 to 16,  $f(r_{l_i}, r_{l_j})$  accumulates the number

27 225 of rankings containing  $r_{l_i}r_{l_j}$ . To further improve the computation efficiency, the sum

28 of  $f(r_{l_i}, r_{l_j})$  and the number of the remaining rankings not yet considered is checked

29 during the accumulation process. If it is less than  $q$ , then  $r_{l_i}r_{l_j}$  has no chance to be a

30  $q$ -support pattern and there is no need to check if the remaining rankings contain  $r_{l_i}r_{l_j}$ .

31 In Algorithm 1, for the case that  $r_{l_i}r_{l_j}$  has no chance to be a  $q$ -support pattern,  $A^{r_l}[j, i]$

32 230 keeps the initialized value, i.e., zero.

33  
34 The following example shows how the matrix representation can be used to evaluate  
35 the ranking consensus.

36  
37 **Example 1.** Consider a set of rankings  $\mathcal{R} = \{\mathbf{r}_1 = (a, b, c, d, e, f), \mathbf{r}_2 = (b, c, d, e, f, a), \mathbf{r}_3 =$   
38  $(b, d, a, g, h, f), \mathbf{r}_4 = (b, a, c, d, f, e)\}$  over candidates  $\{a, b, c, d, e, f, g, h\}$ , and let  $q = 3$ .

We have

$$\mathbf{A}^{\mathbf{r}_1} = \begin{matrix} & a & b & c & d & e & f \\ \begin{matrix} a \\ b \\ c \\ d \\ e \\ f \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 & 0 & 1 \end{pmatrix} \end{matrix}, \mathbf{A}^{\mathbf{r}_2} = \begin{matrix} & b & c & d & e & f & a \\ \begin{matrix} b \\ c \\ d \\ e \\ f \\ a \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \end{matrix}$$

$$\mathbf{A}^{\mathbf{r}_3} = \begin{matrix} & b & d & a & g & h & f \\ \begin{matrix} b \\ d \\ a \\ g \\ h \\ f \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 \end{pmatrix} \end{matrix}, \mathbf{A}^{\mathbf{r}_4} = \begin{matrix} & b & a & c & d & f & e \\ \begin{matrix} b \\ a \\ c \\ d \\ f \\ e \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 \end{pmatrix} \end{matrix}.$$

By (16) and (17), the following result can be obtained

| $l$                          | 1    | 2    | 3    | 4    |
|------------------------------|------|------|------|------|
| $\kappa_1^{\mathbf{r}_1}(3)$ | 1.00 | 1.00 | 0.67 | 1.00 |
| $\kappa_2^{\mathbf{r}_1}(3)$ | 0.67 | 0.67 | 0.33 | 0.73 |

The overall consensus scores are

$$\bar{\kappa}_1(3) = 0.92, \bar{\kappa}_2(3) = 0.60.$$

235 Since  $A^{\mathbf{r}_l}[j, i]$  represents if  $r_{l_i} r_{l_j}$  is a  $q$ -support pattern, it can be known  $\mathbf{S}_1^{\mathbf{r}_1}(3) =$   
 240  $\{a, b, c, d, e, f\}$ ,  $\mathbf{S}_2^{\mathbf{r}_1}(3) = \{af, bc, bd, be, bf, cd, ce, cf, de, df\}$ ,  $\mathbf{S}_1^{\mathbf{r}_2}(3) = \{b, c, d, e, f, a\}$ ,  
 $\mathbf{S}_2^{\mathbf{r}_2}(3) = \{bc, bd, be, bf, ba, cd, ce, cf, de, df\}$ ,  $\mathbf{S}_1^{\mathbf{r}_3}(3) = \{b, d, a, f\}$ ,  $\mathbf{S}_2^{\mathbf{r}_3}(3) = \{bd, ba,$   
 $bf, df, af\}$ ,  $\mathbf{S}_1^{\mathbf{r}_4}(3) = \{b, a, c, d, f, e\}$ ,  $\mathbf{S}_2^{\mathbf{r}_4}(3) = \{ba, bc, bd, bf, be, af, cd, cf, ce, df, de\}$ .  
 Furthermore, the sets of the  $q$ -support patterns of the whole ranking set are  $\mathbf{S}_1(3) =$   
 $\mathbf{S}_1^{\mathbf{r}_1}(3) \cup \mathbf{S}_1^{\mathbf{r}_2}(3) \cup \mathbf{S}_1^{\mathbf{r}_3}(3) \cup \mathbf{S}_1^{\mathbf{r}_4}(3) = \{a, b, c, d, e, f\}$ ,  $\mathbf{S}_2(3) = \mathbf{S}_2^{\mathbf{r}_1}(3) \cup \mathbf{S}_2^{\mathbf{r}_2}(3) \cup$   
 $\mathbf{S}_2^{\mathbf{r}_3}(3) \cup \mathbf{S}_2^{\mathbf{r}_4}(3) = \{af, ba, bc, bd, be, bf, cd, ce, cf, de, df\}$ .

---

**Algorithm 1: Quantifying consensus with matrix representation**

---

**Data:** A set of rankings  $\mathcal{R}$ , the value of  $q$

**Result:**  $\kappa_1^{r_l}(q), \kappa_2^{r_l}(q), l = 1, 2, \dots, N; \mathbf{S}_1^{r_l}(q), \mathbf{S}_2^{r_l}(q); \bar{\kappa}_1(q), \bar{\kappa}_2(q); \mathbf{S}_1(q), \mathbf{S}_2(q)$

```
1 Initialize  $\mathbf{A}^{r_l}, l = 1, 2, \dots, N$  with zero matrices
2 for  $l = 1$  to  $N$  do
3    $m \leftarrow$  Length of  $\mathbf{r}_l$ 
4   for  $i = 1$  to  $m$  do
5     for  $j = i$  to  $m$  do
6        $f(r_{l_i}, r_{l_j}) = 0$ 
7       if  $(l > 1, N - l + 1 \geq q, \exists x \in [1, l - 1])$  or
8          $(l > 1, N - l + 1 < q, \exists x \in [1, N - q + 1])$  such that  $r_{l_i} r_{l_j} \sqsubset \mathbf{r}_x$  then
9          $A^{r_l}[j, i] = A^{r_x}[\pi(r_{l_j}, \mathbf{r}_x), \pi(r_{l_i}, \mathbf{r}_x)]$ 
10        continue
11      else if  $N - l + 1 \geq q$  then
12        for  $z = l$  to  $N$  do
13          Calculate  $\pi(r_{l_i}, \mathbf{r}_z), \pi(r_{l_j}, \mathbf{r}_z)$  by (12)
14          Calculate  $f(r_{l_i}, r_{l_j}) + =$ 
15            
$$\begin{cases} H(\pi(r_{l_i}, \mathbf{r}_z)), & \text{if } i = j \\ H(\pi(r_{l_j}, \mathbf{r}_z) - \pi(r_{l_i}, \mathbf{r}_z))H(\pi(r_{l_i}, \mathbf{r}_z)), & \text{otherwise} \end{cases}$$

16          if  $N - z + f(r_{l_i}, r_{l_j}) < q$  then
17            break
18        end
19      end
20      if  $f(r_{l_i}, r_{l_j}) \geq q$  then
21         $A^{r_l}[j, i] = 1$ 
22      end
23    end
24  end
25 Calculate  $\kappa_1^{r_l}(q), \kappa_2^{r_l}(q)$  by (16) and (17)
26 Get  $\mathbf{S}_1^{r_l}(q), \mathbf{S}_2^{r_l}(q)$  based on  $\mathbf{A}^{r_l}$ 
27 end
28 Calculate  $\bar{\kappa}_1(q), \bar{\kappa}_2(q)$  by (8) and (9)
29 Get  $\mathbf{S}_1(q), \mathbf{S}_2(q)$  by  $\mathbf{S}_1(q) = \bigcup_{r_l \in \mathcal{R}} \mathbf{S}_1^{r_l}(q), \mathbf{S}_2(q) = \bigcup_{r_l \in \mathcal{R}} \mathbf{S}_2^{r_l}(q)$ 
30 return  $\{\kappa_1^{r_l}(q), \kappa_2^{r_l}(q), l = 1, 2, \dots, N; \mathbf{S}_1^{r_l}(q), \mathbf{S}_2^{r_l}(q); \bar{\kappa}_1(q), \bar{\kappa}_2(q); \mathbf{S}_1(q), \mathbf{S}_2(q)\}$ 
```

---

#### 4. Quantifying consensus with consideration of positions and position gaps

The rank positions of an item and the position gaps of pairwise items may be significantly different in a ranking set. Consider the items  $a$  and  $f$  in Example [1](#). The rank positions of item  $a$  are  $\pi(a, \mathbf{r}_1) = 1, \pi(a, \mathbf{r}_2) = 6, \pi(a, \mathbf{r}_3) = 3, \pi(a, \mathbf{r}_4) = 2$  and the position gaps of the two items are  $\pi(f, \mathbf{r}_1) - \pi(a, \mathbf{r}_1) = 5, \pi(f, \mathbf{r}_3) - \pi(a, \mathbf{r}_3) = 3, \pi(f, \mathbf{r}_4) - \pi(a, \mathbf{r}_4) = 3$ . These differences influence the ranking consensus. However, the consensus scores defined in the previous section only involve the existence of  $q$ -support patterns. To reflect the importance of these position and gap information, the following definition presents an extension to [\(6\)](#) and [\(7\)](#) for quantifying consensus of a ranking set more effectively.

**Definition 4** (Weighted individual consensus scores). *The weighted consensus scores of ranking  $\mathbf{r}_l \in \mathcal{R}$  are*

$$\kappa_1^{\mathbf{r}_l}(q) = \frac{1}{N_1^{\mathbf{r}_l}} \sum_{r_{l_i} \in \mathbf{S}_1^{\mathbf{r}_l}(q)} \gamma^{h(r_{l_i}, \mathbf{r}_l)} \quad (18)$$

$$\kappa_2^{\mathbf{r}_l}(q) = \frac{1}{N_2^{\mathbf{r}_l}} \sum_{r_{l_i}, r_{l_j} \in \mathbf{S}_2^{\mathbf{r}_l}(q)} \lambda^{d(r_{l_i}, r_{l_j}, \mathbf{r}_l)}, \quad (19)$$

where the constants  $0 < \gamma \leq 1$  and  $0 < \lambda \leq 1$  are the weights,  $h(r_{l_i}, \mathbf{r}_l)$  is the deviation of the position of  $r_{l_i}$  in  $\mathbf{r}_l$  from its average position in the ranking set, and  $d(r_{l_i}, r_{l_j}, \mathbf{r}_l)$  is the deviation of the position gaps between  $r_{l_i}$  and  $r_{l_j}$  in  $\mathbf{r}_l$  from the average.

The deviations  $h(r_{l_i}, \mathbf{r}_l)$  and  $d(r_{l_i}, r_{l_j}, \mathbf{r}_l)$  are calculated as follows. For ranking  $\mathbf{r}_l \in \mathcal{R}$ , we have the sets  $\mathbf{S}_1^{\mathbf{r}_l}(q)$  and  $\mathbf{S}_2^{\mathbf{r}_l}(q)$  of the  $q$ -support patterns defined as [\(4\)](#) and [\(5\)](#), the function  $f(r_{l_i}, r_{l_j})$  in the form of [\(14\)](#), and the subset  $\mathcal{R}'(r_{l_i}, r_{l_j})$  of  $\mathcal{R}$  containing pattern  $r_{l_i}, r_{l_j}$  as [\(1\)](#). The average position of item  $r_{l_i}$  in the ranking set is defined as

$$\bar{\pi}(r_{l_i}) = \frac{1}{f(r_{l_i}, r_{l_i})} \sum_{\mathbf{r}_z \in \mathcal{R}'(r_{l_i}, r_{l_i})} \pi(r_{l_i}, \mathbf{r}_z). \quad (20)$$

The deviation  $h(r_{l_i}, \mathbf{r}_l)$  is

$$h(r_{l_i}, \mathbf{r}_l) = |\pi(r_{l_i}, \mathbf{r}_l) - \bar{\pi}(r_{l_i})|. \quad (21)$$

The position gap between  $r_{l_i}$  and  $r_{l_j}$  in ranking  $\mathbf{r}_z$  is

$$\omega(r_{l_i}, r_{l_j}, \mathbf{r}_z) = \pi(r_{l_j}, \mathbf{r}_z) - \pi(r_{l_i}, \mathbf{r}_z). \quad (22)$$

The average position gap of  $r_{l_i}$  and  $r_{l_j}$  in the ranking set is defined as

$$\bar{\omega}(r_{l_i}, r_{l_j}) = \frac{1}{f(r_{l_i}, r_{l_j})} \sum_{\mathbf{r}_z \in \mathcal{R}'(r_{l_i}, r_{l_j})} \omega(r_{l_i}, r_{l_j}, \mathbf{r}_z). \quad (23)$$

The deviation  $d(r_{l_i}, r_{l_j}, \mathbf{r}_l)$  is

$$d(r_{l_i}, r_{l_j}, \mathbf{r}_l) = |\omega(r_{l_i}, r_{l_j}, \mathbf{r}_l) - \bar{\omega}(r_{l_i}, r_{l_j})|.$$

From the definition, it can be known that smaller values of  $\gamma$  and  $\lambda$  reflect greater impacts of the deviations of item positions and position gaps in rankings on the consensus scores. It is worth noting that the consensus scores defined in the previous section are a special case of the weighted consensus scores with  $\gamma = 1, \lambda = 1$ . Here, we do not need to make any change to the overall consensus scores defined in Definition 3.

To calculate the weighted consensus scores with the matrix representation, equation (15) in Theorem 1 is changed to

$$A^{\mathbf{r}_l}[j, i] = \begin{cases} \gamma^{h(r_{l_i}, \mathbf{r}_l)}, & \text{if } i = j \text{ and } f(r_{l_i}, r_{l_j}) \geq q \\ \lambda^{d(r_{l_i}, r_{l_j}, \mathbf{r}_l)}, & \text{if } i < j \text{ and } f(r_{l_i}, r_{l_j}) \geq q \\ 0, & \text{otherwise.} \end{cases} \quad (24)$$

Small changes will be needed in Algorithm 1. We follow the steps of Algorithm 1 and change the way to calculate  $A^{\mathbf{r}_l}[j, i]$  on Line 8 to the following form

$$A^{\mathbf{r}_l}[j, i] = \begin{cases} H(A^{\mathbf{r}_x}[\pi(r_{l_j}, \mathbf{r}_x), \pi(r_{l_i}, \mathbf{r}_x)]) \gamma^{h(r_{l_i}, \mathbf{r}_l)}, & \text{if } i = j \\ H(A^{\mathbf{r}_x}[\pi(r_{l_j}, \mathbf{r}_x), \pi(r_{l_i}, \mathbf{r}_x)]) \lambda^{d(r_{l_i}, r_{l_j}, \mathbf{r}_l)}, & \text{if } i < j. \end{cases}$$

Line 19 is replaced by

$$A^{\mathbf{r}_l}[j, i] = \begin{cases} \gamma^{h(r_{l_i}, \mathbf{r}_l)}, & \text{if } i = j \\ \lambda^{d(r_{l_i}, r_{l_j}, \mathbf{r}_l)}, & \text{if } i < j, \end{cases}$$

and meanwhile the average position  $\bar{\pi}(r_{l_i})$  and the average position gap  $\bar{\omega}(r_{l_i}, r_{l_j})$  are recorded in here for further use on Line 8.

**Remark 1** (Rankings with ties). *Rankings with ties are used in the case that the preferences over some items are identical. Let  $\mathbf{r}_z = (\mathcal{T}_{z_1}, \mathcal{T}_{z_2}, \dots, \mathcal{T}_{z_n})$  be a ranking with ties, where  $\mathcal{T}_{z_i}, i \in [1, n]$  is a set of items with identical preference. For  $i < j$ , every item in  $\mathcal{T}_{z_i}$  is more preferred than all the items in  $\mathcal{T}_{z_j}$ . The proposed approach can be extended to rankings with ties by making a small change to the position function. Specifically, we can replace (12) with*

$$\pi(r_{l_i}, \mathbf{r}_z) = \begin{cases} p, & \text{if } r_{l_i} \in \mathcal{T}_{z_p} \\ 0, & \text{otherwise} \end{cases}$$

*to make the approach applicable to evaluate consensus of rankings with ties.*

## 5. Detecting outliers

The individual consensus scores  $\kappa_1^{\mathbf{r}_l}(q)$  and  $\kappa_2^{\mathbf{r}_l}(q)$  directly reflect the (weighted) numbers of  $q$ -support patterns that  $\mathbf{r}_l$  shares with the other rankings in  $\mathcal{R}$ . For instance, ranking  $\mathbf{r}_3$  in Example 1 shares less 3-support patterns with the others, thus it has much lower consensus scores. This can be used to detect outlier rankings, which have low consensus with most rankings. The following outlier detection method is naturally developed from the consensus quantifying approach.

Consider a ranking set  $\mathcal{R}$  with overall consensus scores  $\bar{\kappa}_1(q)$  and  $\bar{\kappa}_2(q)$  for a given  $q$ . Define the relative deviations of the individual consensus scores of ranking  $\mathbf{r}_l \in \mathcal{R}$  from the overall consensus scores as

$$v_1^{\mathbf{r}_l}(q) = \frac{\kappa_1^{\mathbf{r}_l}(q) - \bar{\kappa}_1(q)}{\bar{\kappa}_1(q)} \quad (25)$$

$$v_2^{\mathbf{r}_l}(q) = \frac{\kappa_2^{\mathbf{r}_l}(q) - \bar{\kappa}_2(q)}{\bar{\kappa}_2(q)}. \quad (26)$$

Note that  $v_1^{\mathbf{r}_l}(q) < 0$  and  $v_2^{\mathbf{r}_l}(q) < 0$  imply that the ranking  $\mathbf{r}_l$  has lower consensus scores than the overall averages. For given constants  $\epsilon_1 > 0$  and  $\epsilon_2 > 0$ , if  $v_1^{\mathbf{r}_l}(q) < -\epsilon_1$  or  $v_2^{\mathbf{r}_l}(q) < -\epsilon_2$ , we regards  $\mathbf{r}_l$  as an outlier of the ranking set. The values of  $\epsilon_1, \epsilon_2$  depend on the specific need for a system.

This outlier detection method can be used to figure out irrelevant rankings in the ranking set and consequently identify the majority of rankings with higher consensus.



1  
2  
3  
4  
5  
6  
7  
8  
9 It is of great importance in many scenarios, e.g., design of auto-suggestion queries in  
10 search engine. It is worth noting that one potential application of the obtained detection  
11 method is to improve rank aggregation. Rank aggregation is the task of aggregating  
12 the preferences of different agents to generate a final ranking. The outliers of rank-  
13 ings/agents play a negative role in drawing a consensus ranking. Even though many  
14 285 existing studies have been carried out on rank aggregation [40, 9, 5], there is still room  
15 to improve aggregated rankings so that the aggregated result is as close to the ground  
16 truth as possible. This will be studied in a separate paper.  
17  
18  
19  
20  
21  
22

## 23 6. Experimental studies

24  
25 290 This section shows how the proposed approach can be used to evaluate consen-  
26 sus for a set of rankings. The source code is available at [https://github.com/  
27 zhiweiuu/secs](https://github.com/zhiweiuu/secs).  
28  
29  
30

### 31 6.1. Analysis of the Mechanical Turk Dots datasets

32  
33 The Mechanical Turk Dots datasets [30] include four publicly available datasets  
34 295 obtained for four dots tasks. These datasets each contain rankings obtained by 794  
35 to 800 voters over four candidates. Each candidate corresponds to a certain number  
36 of random dots. The voters were asked to rank the candidates from those with the  
37 least dots to the most. Each task contains candidates with 200,  $200 + i$ ,  $200 + 2i$ , and  
38  $200 + 3i$  dots, where  $i = 3, 5, 7, 9$  respectively for the four tasks. Figure 1 shows the  
39 proportions of rankings in each dataset with different Spearman's  $\rho$  to the ground truth  
40 ranking. The values of different Spearman's  $\rho$  are distinguished by colors. It can be seen  
41 that the proportions of rankings with high Spearman coefficients 0.8 and 1.0 increase  
42 from Dataset 1 to Dataset 4, while that with coefficient 0.4 decreases significantly. The  
43 300 ranking consensus degrees seem increasing from Dataset 1 to Dataset 4. We apply the  
44 proposed approach to accurately compare these datasets.  
45  
46  
47  
48  
49  
50  
51 305

52 The overall consensus scores without weighting are first considered. Since the  
53 datasets have complete rankings, i.e., all the candidates under consideration are ranked  
54 in the rankings, the consensus scores of the single items satisfy  $\bar{\kappa}_1(q) = 4$  for all  $q$  and  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

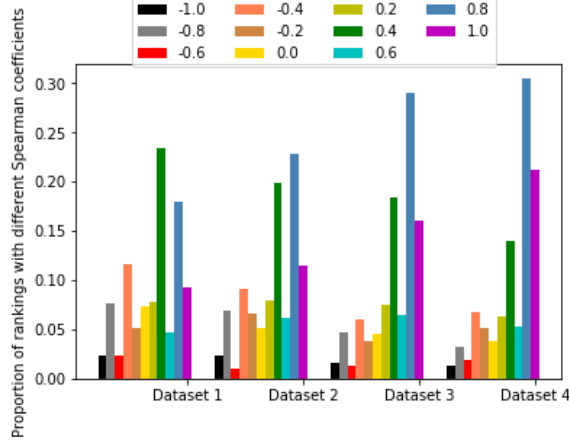


Figure 1: Spearman's  $\rho$  between the rankings and the ground truth ranking

all the datasets. Figure 2 gives the overall consensus scores  $\bar{\kappa}_2(q)$  with respect to  $\frac{q}{N}$ , where  $\frac{q}{N} \geq 0.5$  indicates that the commonality embedded in half or more than half of the rankings is evaluated. The trend of the overall consensus scores for the four datasets is clear. Dataset 4 has the largest overall consensus score, which indicates that Dataset 4 has the most  $q$ -support common patterns. Specifically, it can be seen from the figure that, when  $\frac{q}{N}$  is 0.5, the consensus score  $\bar{\kappa}_2(q)$  is 0.59, 0.62, 0.68, and 0.71 respectively for Dataset 1, 2, 3 and 4. This means that on average, 59.00%, 62.00%, 68.00%, and 71.00% of the pairwise patterns of a ranking are  $\lceil \frac{N}{2} \rceil$ -support patterns in Dataset 1, 2, 3 and 4, respectively. As the value of  $q$  increases, the consensus scores decrease. When  $\frac{q}{N}$  reaches 0.67, the consensus score is zero for Dataset 1, which means that arbitrary  $q \geq 0.67N$  rankings in the dataset have no common pattern. On the other hand, the consensus scores are 0.12, 0.37, 0.38 for Dataset 2, 3, 4. In other words, on average, 12.00%, 37.00%, 38.00% of the patterns of a ranking are supported by at least  $\lceil 0.67N \rceil$  rankings in the corresponding dataset.

The overall consensus scores with weightings are then evaluated. Figure 3 shows the consensus scores with respect to the weights  $\gamma$  and  $\lambda$  for a fixed  $q = \lceil \frac{N}{2} \rceil$ . As shown,  $\bar{\kappa}_1(\lceil \frac{N}{2} \rceil)$  and  $\bar{\kappa}_2(\lceil \frac{N}{2} \rceil)$  decrease with the increase of weightings on the deviations of positions and position gaps. Dataset 1 has the lowest overall consensus scores and

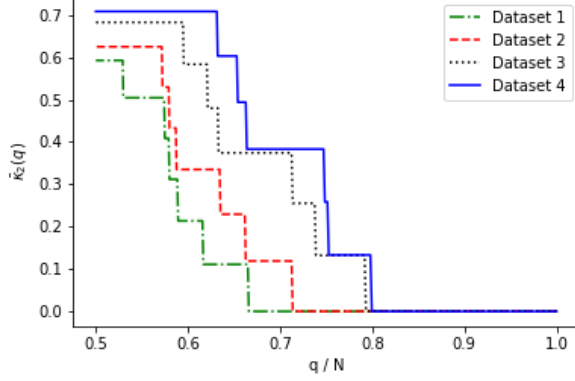


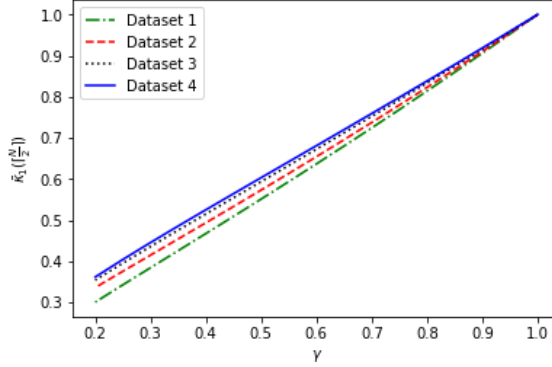
Figure 2: Consensus scores  $\bar{\kappa}_2(q)$  of Dots datasets without weighting

Dataset 4 has the highest. The ratios of the consensus scores between Dataset 4 and Dataset 3, Dataset 3 and Dataset 2, and Dataset 2 and Dataset 1 are shown in Table 1 for the cases without weighting and with weighting parameters  $\gamma = 0.5, \lambda = 0.5$ . By comparing the two cases, it can be found that the ratios with weightings on the deviations of the position and position gaps are higher than those without weightings. This reveals that the differences of the positions of the single  $q$ -support items and the position gaps of the  $q$ -support patterns decrease from Dataset 1 to Dataset 4.

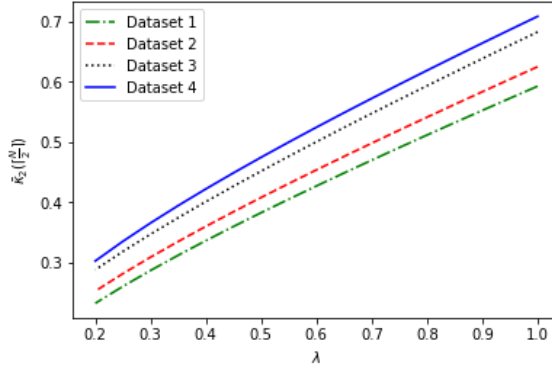
Table 1: Ratios of the consensus scores between datasets

|  | Dataset 4/Dataset 3 | Dataset 3/Dataset 2 | Dataset 2/Dataset 1 |
|--|---------------------|---------------------|---------------------|
| $\bar{\kappa}_1(\lceil \frac{N}{2} \rceil), \gamma = 1$    | 1.00                | 1.00                | 1.00                |
| $\bar{\kappa}_1(\lceil \frac{N}{2} \rceil), \gamma = 0.5$  | 1.02                | 1.04                | 1.04                |
| $\bar{\kappa}_2(\lceil \frac{N}{2} \rceil), \lambda = 1$   | 1.04                | 1.09                | 1.05                |
| $\bar{\kappa}_2(\lceil \frac{N}{2} \rceil), \lambda = 0.5$ | 1.05                | 1.11                | 1.07                |

The relative deviations of  $\kappa_2^{r_i}(\lceil \frac{N}{2} \rceil)$  from the overall consensus score  $\bar{\kappa}_2(\lceil \frac{N}{2} \rceil)$  is also studied to verify the effectiveness of the proposed outlier detection method. By choosing  $\lambda = 0.5$ , the result in Table 2 can be obtained. The deviations are very high for  $\mathbf{r}_{21}, \mathbf{r}_{24}, \mathbf{r}_{22}, \mathbf{r}_{13}$  of Dataset 1,  $\mathbf{r}_{20}, \mathbf{r}_{22}, \mathbf{r}_{19}, \mathbf{r}_{17}$  of Dataset 2,  $\mathbf{r}_{19}, \mathbf{r}_{20}, \mathbf{r}_{22}, \mathbf{r}_{15}$  of Dataset 3, and  $\mathbf{r}_{21}, \mathbf{r}_{22}, \mathbf{r}_{24}, \mathbf{r}_{20}$  of Dataset 4. These rankings are regarded as out-



(a)  $\bar{\kappa}_1(\lceil \frac{N}{2} \rceil)$  with respect to  $\gamma$



(b)  $\bar{\kappa}_2(\lceil \frac{N}{2} \rceil)$  with respect to  $\lambda$

Figure 3: Weighted consensus scores of Dots datasets

liers of the datasets. They are  $(4, 3, 2, 1)$ ,  $(4, 3, 1, 2)$ ,  $(4, 2, 3, 1)$ ,  $(3, 4, 2, 1)$  respectively in each dataset. Note that Spearman's  $\rho$  between  $(4, 3, 2, 1)$  and the ground truth  $(1, 2, 3, 4)$  are  $-1$ , and all the Spearman coefficients of the rest three to the ground truth are  $-0.8$ . After deleting these outlier rankings, the consensus scores  $\bar{\kappa}_1(\lceil \frac{N}{2} \rceil)$  with  $\gamma = 0.5$  increase from 0.55, 0.57, 0.59, 0.60 to 0.58, 0.59, 0.61, 0.62 for Dataset 1, 2, 3, 4, respectively. The consensus scores  $\bar{\kappa}_2(\lceil \frac{N}{2} \rceil)$  change from 0.38, 0.41, 0.45, 0.47 to 0.42, 0.44, 0.48, 0.49 for the four datasets. This confirms the effectiveness of the proposed outlier detection method.

It is further found that the four datasets have the same set of the  $\lceil \frac{N}{2} \rceil$ -support pat-

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

Table 2: Deviations of the consensus scores

| $q = \lceil \frac{N}{2} \rceil$ | Dataset 1    | Dataset 2    | Dataset 3    | Dataset 4    |
|---------------------------------|--------------|--------------|--------------|--------------|
| $v_2^{\mathbf{r}^1}(q)$         | 0.72         | 0.68         | 0.56         | 0.55         |
| $v_2^{\mathbf{r}^2}(q)$         | 0.14         | 0.38         | 0.29         | 0.20         |
| $v_2^{\mathbf{r}^3}(q)$         | 0.44         | 0.29         | 0.28         | 0.18         |
| $v_2^{\mathbf{r}^4}(q)$         | 0.44         | 0.36         | 0.15         | 0.08         |
| $v_2^{\mathbf{r}^5}(q)$         | 0.09         | 0.08         | -0.03        | -0.09        |
| $v_2^{\mathbf{r}^6}(q)$         | 0.06         | 0.06         | -0.03        | -0.21        |
| $v_2^{\mathbf{r}^7}(q)$         | 0.47         | -0.01        | -0.15        | -0.10        |
| $v_2^{\mathbf{r}^8}(q)$         | 0.11         | -0.27        | -0.04        | -0.11        |
| $v_2^{\mathbf{r}^9}(q)$         | 0.02         | -0.04        | -0.35        | -0.41        |
| $v_2^{\mathbf{r}^{10}}(q)$      | -0.22        | 0.02         | -0.37        | -0.21        |
| $v_2^{\mathbf{r}^{11}}(q)$      | -0.29        | -0.28        | -0.17        | -0.20        |
| $v_2^{\mathbf{r}^{12}}(q)$      | -0.38        | -0.11        | -0.42        | -0.41        |
| $v_2^{\mathbf{r}^{13}}(q)$      | <b>-0.71</b> | -0.13        | -0.20        | -0.56        |
| $v_2^{\mathbf{r}^{14}}(q)$      | -0.25        | -0.15        | -0.49        | -0.35        |
| $v_2^{\mathbf{r}^{15}}(q)$      | -0.06        | -0.46        | <b>-0.74</b> | -0.21        |
| $v_2^{\mathbf{r}^{16}}(q)$      | -0.09        | -0.45        | -0.20        | -0.54        |
| $v_2^{\mathbf{r}^{17}}(q)$      | -0.36        | <b>-0.72</b> | -0.28        | -0.49        |
| $v_2^{\mathbf{r}^{18}}(q)$      | -0.39        | -0.37        | -0.50        | -0.47        |
| $v_2^{\mathbf{r}^{19}}(q)$      | -0.40        | <b>-0.74</b> | <b>-1.00</b> | -0.55        |
| $v_2^{\mathbf{r}^{20}}(q)$      | -0.46        | <b>-1.00</b> | <b>-0.75</b> | <b>-0.74</b> |
| $v_2^{\mathbf{r}^{21}}(q)$      | <b>-1.00</b> | -0.44        | -0.50        | <b>-1.00</b> |
| $v_2^{\mathbf{r}^{22}}(q)$      | <b>-0.73</b> | <b>-0.75</b> | <b>-0.77</b> | <b>-0.75</b> |
| $v_2^{\mathbf{r}^{23}}(q)$      | -0.12        | -0.45        | -0.52        | -0.54        |
| $v_2^{\mathbf{r}^{24}}(q)$      | <b>-0.74</b> | -0.47        | -0.52        | <b>-0.77</b> |

1  
2  
3  
4  
5  
6  
7  
8  
9 terms  $\mathbf{S}_2(\lceil \frac{N}{2} \rceil) = \{12, 13, 14, 23, 24, 34\}$ . By aggregating these  $\lceil \frac{N}{2} \rceil$ -support patterns,  
10 we can obtain the ranking (1, 2, 3, 4), i.e., the ground truth ranking. This enhances  
11 the advantage of the proposed consensus quantifying approach over the rank correla-  
12 350 tion and distance functions approaches, where no common patterns of the rankings are  
13 specified.  
14  
15  
16

## 17 6.2. Evaluation of the information retrieval results of the 2015 CLEFeHealth Lab Task

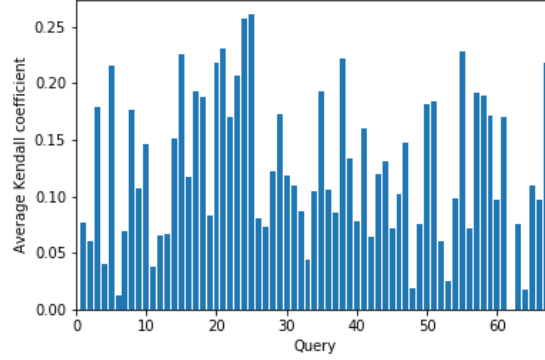
### 18 2

19  
20  
21 355 This experiment focuses on top- $k$  rankings using the dataset of the CLEF 2015  
22 eHealth Evaluation Lab Task 2 [32], instead of the complete rankings as in the previ-  
23 ous section. The CLEF 2015 eHealth Evaluation Lab Task 2 aimed to foster the design  
24 of web search engines in providing access to medical information especially for self-  
25 diagnosis information, since commercial search engines were far from being effective  
26 in the field. The problem considered in the task was to retrieve web pages for queries  
27 360 related to different medical conditions. The queries were pre-generated by showing im-  
28 ages and videos of medical conditions to potential users. There were 67 queries selected  
29 to be used in the task for 23 medical conditions, among which 22 conditions had three  
30 queries and one condition had one query. The queries were first created in English and  
31 then translated into several other languages. The document collection made available  
32 365 to the participants for information retrieval contains approximately one million web  
33 pages on a broad range of health topics. The participants were asked to submit up to  
34 ten runs for the English queries. The first run of each team was with the highest priority  
35 for selection of documents to contribute to the final assessment. Twelve participating  
36 370 teams submitted their English information retrieval results.  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46

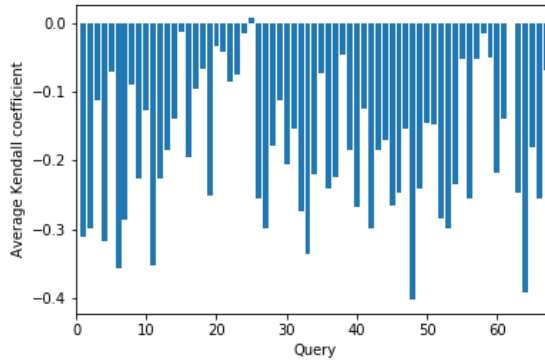
47 This section evaluates the information retrieval results of the first English runs.  
48 Given that the first two pages of a user’s search result probably draw the most atten-  
49 tion in practice, the top-20 retrieved documents for each query are considered in the  
50 evaluation. The conventional Spearman’s  $\rho$  and Kendall’s  $\tau$  measure the correlation  
51 375 of two complete rankings, as they compare the positions of same items in two rank-  
52 ings. For this dataset with incomplete rankings, the Spearman’s  $\rho$  and Kendall’s  $\tau$  for  
53 top- $k$  rankings proposed in [15] are employed to measure the correlations for the 67  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
 2  
 3  
 4  
 5  
 6  
 7  
 8  
 9 queries. Because a typo exists in the  $62^{nd}$  query, there is no record of some teams  
 10 for this query in the dataset. This query is not considered in the following analysis.  
 11  
 12 Since there is no ground truth ranking available, we pairwise compare the ranking  
 13 380 obtained by each team for a specific query with the rankings of the other teams and  
 14 take the average as the comparison result of the team. The obtained comparison results  
 15 of all the teams for the query are further aggregated by taking their average, and the  
 16 aggregated result measures the correlations of rankings obtained by all the teams for  
 17 the query. Figure 4 gives the results of Kendall’s  $\tau$ . Note that a key parameter  $p$  is  
 18 385 introduced in the calculation of Kendall’s  $\tau$  for top- $k$  rankings in [15]. This parameter  
 19 corresponds to the penalty for the case that two items  $\sigma_i$  and  $\sigma_j$  appear in one ranking  
 20  $\mathbf{r}_l$  and none of them are considered in the other compared ranking  $\mathbf{r}_z$ . In this case, the  
 21 term  $\text{sgn}(\pi(\sigma_i, \mathbf{r}_l) - \pi(\sigma_j, \mathbf{r}_l))\text{sgn}(\pi(\sigma_i, \mathbf{r}_z) - \pi(\sigma_j, \mathbf{r}_z))$  is set to be  $p$ . We normalize  
 22 Kendall’s  $\tau$  to the domain of  $[-1, 1]$ . The parameter  $p = 1$  gives an optimistic ap-  
 23 390 proach. It implies that  $\sigma_i$  and  $\sigma_j$  in  $\mathbf{r}_z$  are regarded as in the same order as in  $\mathbf{r}_l$  when  
 24 there is no enough information about them. When  $p = 0$ , it gives a neutral approach.  
 25 It can be found in Figure 4(a) and Figure 4(b) that Kendall’s coefficients are highly  
 26 depends on the value of  $p$ . The result of Spearman’s  $\rho$  is shown in Figure 5. If an item  
 27 395  $\sigma_i$  in one top- $k$  ranking  $\mathbf{r}_l$  does not appear in the other compared top- $k$  ranking  $\mathbf{r}_z$ , then  
 28 the position  $\pi(\sigma_i, \mathbf{r}_z)$  is set to  $\ell$ . In Figure 5,  $\ell$  is chosen to be  $k + 1$ . Spearman’s  $\rho$  also  
 29 depends on the value of  $\ell$ .

30  
 31 Unlike the Spearman’s  $\rho$  and Kendall’s  $\tau$  for top- $k$  rankings, where assumptions  
 32 about unknown factors are made without sufficient information and may consequently  
 33 400 lead to bias in the measurement results, the proposed approach has no such problem and  
 34 the consensus of a ranking set is measured more intuitively based on  $q$ -support patterns.  
 35 It provides a clear understanding about the commonality embedded in the rankings  
 36 obtained with different information retrieval approaches, and it can help to find hard  
 37 topics in the information retrieval task. Figure 6 shows the 6-support (i.e.,  $\frac{N}{2}$ -support)  
 38 405 consensus scores without weightings for the ranking sets of the 66 queries obtained  
 39 by the 12 teams. The relative values of the consensus scores are generally consistent  
 40 with the results in Figures 4 and 5. However, our results based on  $q$ -support patterns,  
 41 especially the pairwise patterns, reveal more obvious and detailed information. It can



(a) Optimistic approach



(b) Neutral approach

Figure 4: Average Kendall's  $\tau$  for the ranking sets of the 66 queries obtained by the 12 teams

be seen from Figure 6(a) that the consensus score  $\bar{\kappa}_1(6)$  is greater than 0.5 for queries  
 410 10, 13, 15, 20, 24, 25, 31, 38, 57, 58, 59, 67. This means that, on average, more than  
 47 50% of the ranked items in a ranking for these queries are embedded in at least half  
 48 of the ranking set. When the orders of these ranked items are further considered, Figure  
 49 6(b) shows that, on average, more than 15% of the pairwise patterns of a ranking are  
 50 supported by at least half of the rankings for queries 20, 24, 25, 38, 57, 58, 59, 67.  
 51  
 52  
 53  
 54  
 415 Figure 7 shows the 6-support consensus scores with the weighting parameters on the  
 55 deviations of positions and position gaps being  $\gamma = 0.9, \lambda = 0.9$ . It can be noticed  
 56  
 57  
 58  
 59  
 60  
 61  
 62  
 63  
 64  
 65



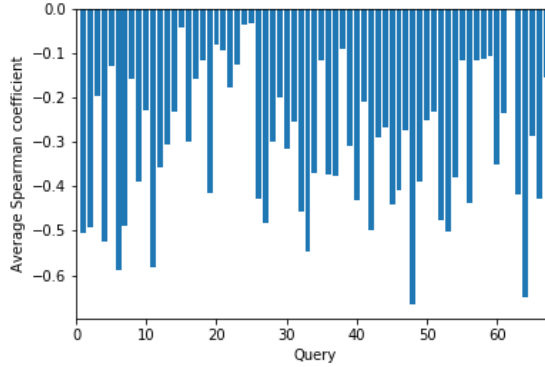


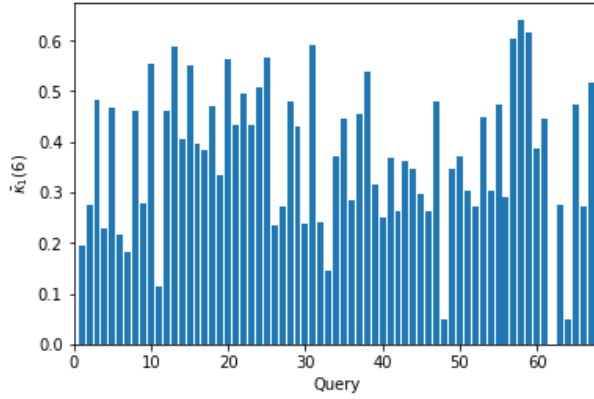
Figure 5: Average Spearman’s  $\rho$  for the ranking sets of the 66 queries obtained by the 12 teams

that queries 58, 25, 24, 55 have higher consensus scores  $\bar{\kappa}_2(\delta)$ , which indicates that the rankings of these queries share more weighted pairwise  $q$ -support patterns. Moreover, the consensus scores  $\bar{\kappa}_1(\delta)$  for these queries are also high. In contrast, the consensus scores of queries 64, 48, 11, 33 are much lower. The detailed information of these queries is given in Table 3 and Table 4. By comparing the two tables, it can be found that the queries with clear descriptions or for typical symptoms tend to have higher consensus scores, while vague descriptions or uncommon symptoms lead to retrieval results with lower consensus scores.

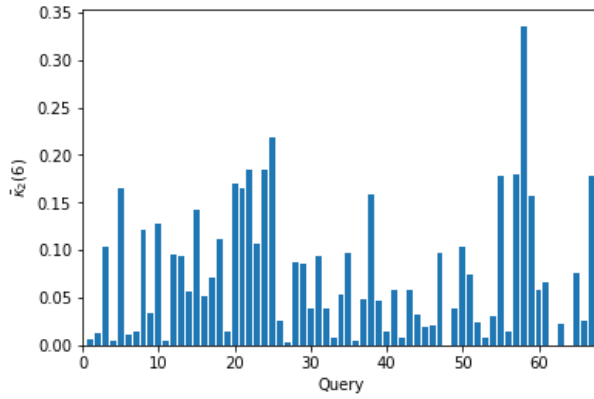
Table 3: Queries with higher consensus scores

| Query ID | Query                                 | $\bar{\kappa}_1(\delta)$ | $\bar{\kappa}_2(\delta)$ |
|----------|---------------------------------------|--------------------------|--------------------------|
| 58       | 39 degree and chicken pox             | 0.47                     | 0.27                     |
| 25       | red rash baby face                    | 0.45                     | 0.17                     |
| 24       | yellow gunk coming from one eye itchy | 0.42                     | 0.15                     |
| 55       | crate type mark in skin               | 0.40                     | 0.15                     |

The consensus of the information retrieval results for each topic is also evaluated with the proposed approach. The queries for each topic are supposed to link to an identical medical condition. The consensus based on 2-support patterns is studied for the 22 topics each with three queries. Topic 13 is not considered, since it associates with



(a) Consensus score  $\bar{\kappa}_1(6)$

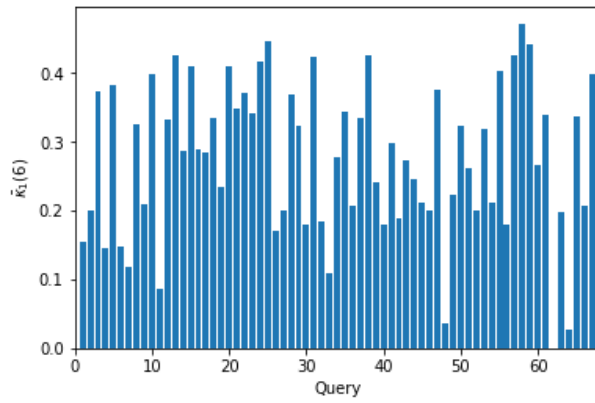


(b) Consensus score  $\bar{\kappa}_2(6)$

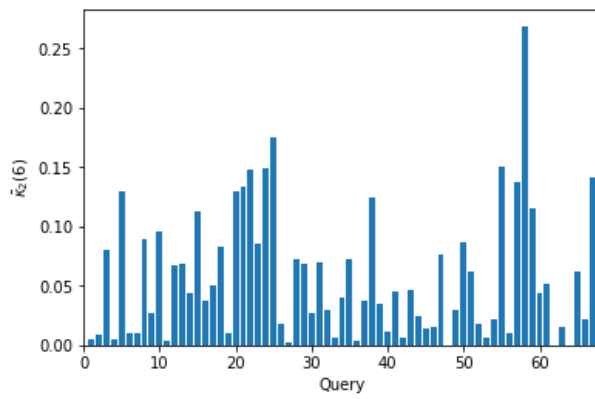
Figure 6: Consensus scores without weighting for the ranking sets of the 66 queries obtained by the 12 teams

query 62 having incomplete record in the dataset. We take the average of the consensus scores of the ranking sets of the 12 teams. The results are given in Figure 8. Specially, the rankings of topics 15 and 11 have the highest average consensus scores, and the average consensus scores for topic 21 and topic 18 are the lowest. By comparing the topics and the details of the related queries in Table 5 and Table 6, it can be found that the diseases of topics 15 and 11 are more common diseases to be easily self-diagnosed and the generated queries share more commonality. On the contrary, the topics with low consensus scores have more diverse queries, thus they can be regarded as hard

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



(a) Consensus score  $\bar{\kappa}_1(6)$



(b) Consensus score  $\bar{\kappa}_2(6)$

Figure 7: Weighted consensus scores for the ranking sets of the 66 queries obtained by the 12 teams

topics, which can be used in further tasks for the development of more advanced search engines.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

Table 4: Queries with lower consensus scores

| Query ID | Query  | $\bar{\kappa}_1(6)$ | $\bar{\kappa}_2(6)$ |
|----------|--|---------------------|---------------------|
| 64       | involuntary rapid left-right eye motion      | 0.03                | 0.00                |
| 48       | cannot stop moving my eyes medical condition | 0.04                | 0.00                |
| 11       | white patchiness in mouth                    | 0.09                | 0.00                |
| 33       | white infection in pharynx                   | 0.11                | 0.01                |

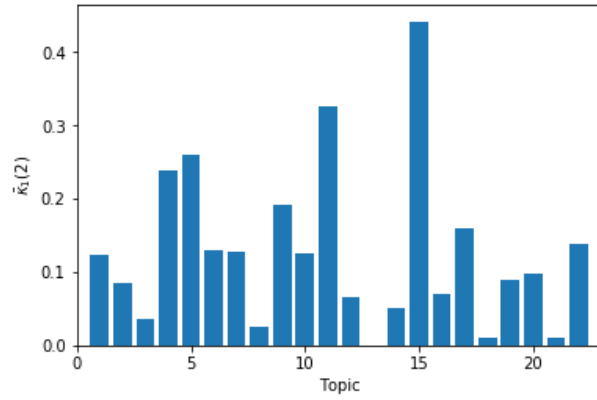
Table 5: Topics with higher average consensus scores

| Topic                                | Query  | $\bar{\kappa}_1(2)$ | $\bar{\kappa}_2(2)$ |
|--------------------------------------|--|---------------------|---------------------|
| 15: whooping cough<br>(pertussis)    | 12: baby has dry cough and has<br>problem to swallow saliva<br>46: baby cough<br>66: treatment of coughs in babies | 0.44                | 0.37                |
| 11: bronchiolitis<br>(caused by rsv) | 31: toddler having squeaky breath<br>49: baby always breathing with mouth closed<br>59: heavy and squeaky breath   | 0.32                | 0.17                |

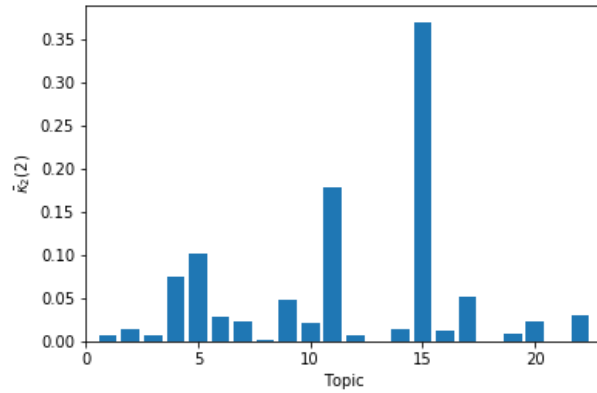
Table 6: Topics with lower average consensus scores

| Topic                  | Query  | $\bar{\kappa}_1(2)$ | $\bar{\kappa}_2(2)$ |
|------------------------|--|---------------------|---------------------|
| 21: nystagmus          | 36: eye are shaking<br>48: cannot stop moving my eyes medical condition<br>64: involuntary rapid left-right eye motion | 0.01                | 0.00                |
| 18: asthma<br>wheezing | 6: child make hissing sound when breathing<br>15: asthma attack<br>30: weird sounds when breathing                     | 0.01                | 0.00                |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



(a) Team average of the consensus score  $\bar{\kappa}_1(2)$



(b) Team average of the consensus score  $\bar{\kappa}_2(2)$

Figure 8: Average weighted consensus scores for the 22 topics of the 3 queries

1  
2  
3  
4  
5  
6  
7  
8  
9 **7. Conclusion**

10  
11 This paper presents a novel approach to quantifying the consensus degree of a rank-  
12 440 ing set. A new concept of  $q$ -support has been introduced to represent the common pat-  
13 terns embedded in rankings. A matrix representation has been developed to describe  
14 the commonality within a ranking set that is shared by an individual ranking, on the  
15 basis of which an algorithm has been developed to quantify the consensus efficiently.  
16  
17  
18  
19 445 Moreover, a scheme for detecting outliers in a ranking set is derived from the consen-  
20 sus quantifying approach. Consensus evaluation with weighting on item positions and  
21 position gaps has also been considered. Compared with the existing methods based  
22 on correlation and distance functions, our approach can characterize and quantify the  
23 group preferences more explicitly and it also lays the foundation for the effective de-  
24 tection of outliers and the development of rank aggregation algorithm, which have been  
25 450 illustrated in the experimental studies.  
26  
27  
28  
29  
30

31  
32 **Acknowledgements**

33  
34 This work is supported by the UK EPSRC under Grant No. EP/P031668/1.  
35  
36

37 **References**

- 38  
39 455 [1] Alcalde-Unzu J, Vorsatz M. Measuring consensus: Concepts, comparisons, and  
40 properties. In: Consensual Processes. Springer; 2011. p. 195–211.  
41  
42 [2] Alcalde-Unzu J, Vorsatz M. Measuring the cohesiveness of preferences: an ax-  
43 iomatic analysis. *Social Choice and Welfare* 2013;41(4):965–88.  
44  
45 [3] Baigent N. Preference proximity and anonymous social choice. *The Quarterly*  
46 460 *Journal of Economics* 1987;102(1):161–9.  
47  
48 [4] Bosch R. Characterizations of voting rules and consensus measures. Ph D Dis-  
49 sertation, Tilburg University 2005;.  
50  
51 [5] Caragiannis I, Chatzigeorgiou X, Krimpas GA, Voudouris AA. Optimizing posi-  
52 tional scoring rules for rank aggregation. *Artificial Intelligence* 2018;.  
53  
54  
55  
56  
57  
58

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

465 [6] Carpineto C, Romano G. A survey of automatic query expansion in information retrieval. *Acm Computing Surveys (CSUR)* 2012;44(1):1.

[7] Carron AV, Brawley LR. Cohesion: Conceptual and measurement issues. *Small group research* 2000;31(1):89–106.

[8] Carterette B. On rank correlation and the distance between rankings. In: *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*. ACM; 2009. p. 436–43. 470

[9] Chen Y, Suh C. Spectral mle: Top-k rank aggregation from pairwise comparisons. In: *International Conference on Machine Learning*. 2015. p. 371–80.

[10] Chiniara M, Bentein K. The servant leadership advantage: When perceiving 475 low differentiation in leader-member relationship quality influences team cohesion, team task performance and service ocb. *The Leadership Quarterly* 2018;29(2):333–45.

[11] Eckert D, Klamler C. Distance-based aggregation theory. In: *Consensual processes*. Springer; 2011. p. 3–22.

480 [12] Elzinga C, Wang H, Lin Z, Kumar Y. Concordance and consensus. *Information Sciences* 2011;181(12):2529–49.

[13] Erdamar B, García-Lapresta JL, Pérez-Román D, Sanver MR. Measuring consensus in a preference-approval context. *Information Fusion* 2014;17:14–21.

[14] Etesami O, Gohari A. Maximal rank correlation. *IEEE Communications Letters* 485 2016;20(1):117–20.

[15] Fagin R, Kumar R, Sivakumar D. Comparing top k lists. *SIAM Journal on discrete mathematics* 2003;17(1):134–60.

[16] García-Lapresta JL, Pérez-Román D. Consensus measures generated by weighted 490 kemeny distances on weak orders. In: *Intelligent Systems Design and Applications (ISDA), 2010 10th International Conference on*. IEEE; 2010. p. 463–8.

- 1  
2  
3  
4  
5  
6  
7  
8  
9 [17] García-Lapresta JL, Pérez-Román D. Measuring consensus in weak orders. In:  
10 Consensual processes. Springer; 2011. p. 213–34.  
11  
12 [18] Hashemi V, Endriss U. Measuring diversity of preferences in a group. In: Euro-  
13 pean Conference on Artificial Intelligence. 2014. p. 423–8.  
14  
15 [19] Hassanzadeh FF, Milenkovic O. An axiomatic approach to constructing distances  
495 for rank comparison and aggregation. IEEE Transactions on Information Theory  
16 2014;60(10):6417–39.  
17  
18 [20] Henzgen S, Hüllermeier E. Weighted rank correlation: a flexible approach based  
19 on fuzzy order relations. In: Joint European Conference on Machine Learning  
20 and Knowledge Discovery in Databases. Springer; 2015. p. 422–37.  
21  
22 500 [21] Hogg MA. Group cohesiveness: A critical review and some new directions. Eu-  
23 ropean review of social psychology 1993;4(1):85–111.  
24  
25 [22] Hotho A, Jäschke R, Schmitz C, Stumme G. Information retrieval in folk-  
26 sonomies: Search and ranking. In: European Semantic Web conference. Springer;  
27 2006. p. 411–26.  
28  
29 505 [23] Iman RL, Conover W. A measure of top–down correlation. Technometrics  
30 1987;29(3):351–7.  
31  
32 [24] Karpov A. Preference diversity orderings. Group Decision and Negotiation  
33 2017;26(4):753–74.  
34  
35 [25] Kemeny JG. Mathematics without numbers. Daedalus 1959;88(4):577–91.  
36  
37 510 [26] Kendall MG. A new measure of rank correlation. Biometrika 1938;30(1/2):81–  
38 93.  
39  
40 [27] Kim SH, Choi SH, Kim JK. An interactive procedure for multiple attribute group  
41 decision making with incomplete information: Range-based approach. European  
42 Journal of Operational Research 1999;118(1):139–52.  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54 515  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



- 1  
2  
3  
4  
5  
6  
7  
8  
9 [28] Kumar R, Vassilvitskii S. Generalized distances between rankings. In: Proceedings of the 19th international conference on World wide web. ACM; 2010. p. 571–80.
- 10  
11  
12  
13  
14 [29] Liu TY, et al. Learning to rank for information retrieval. Foundations and  
15 Trends® in Information Retrieval 2009;3(3):225–331.
- 16  
17  
18 [30] Mao A, Procaccia AD, Chen Y. Better human computation through principled  
19 voting. In: Conference on Artificial Intelligence. 2013. .
- 20  
21  
22 [31] Nurmi H. A comparison of some distance-based choice rules in ranking environ-  
23 ments. Theory and Decision 2004;57(1):5–24.
- 24  
25  
26 [32] Palotti JR, Zuccon G, Goeuriot L, Kelly L, Hanbury A, Jones GJ, Lupu M, Pecina  
27 P. Clef ehealth evaluation lab 2015, task 2: Retrieving information about medical  
28 symptoms. In: CLEF (Working Notes). 2015. p. 1–22.
- 29  
30  
31 [33] Poshyvanyk D, Gueheneuc YG, Marcus A, Antoniol G, Rajlich V. Feature lo-  
32 cation using probabilistic ranking of methods based on execution scenarios and  
33 information retrieval. IEEE Transactions on Software Engineering 2007;33(6).
- 34  
35  
36 [34] Qin J, Liu X. Multi-attribute group decision making using combined rank-  
37 ing value under interval type-2 fuzzy environment. Information Sciences  
38 2015;297:293–315.
- 39  
40  
41 [35] Salas E, Grossman R, Hughes AM, Coultas CW. Measuring team cohesion: Ob-  
42 servations from the science. Human factors 2015;57(3):365–74.
- 43  
44  
45 [36] Shieh GS. A weighted kendall’s tau statistic. Statistics & probability letters  
46 1998;39(1):17–24.
- 47  
48  
49 [37] Spearman C. The proof and measurement of association between two things. The  
50 American journal of psychology 1904;15(1):72–101.
- 51  
52  
53 [38] Tan L, Clarke CL. A family of rank similarity measures based on maximized  
54 effectiveness difference. IEEE Transactions on Knowledge and Data Engineering  
55 2015;27(11):2865–77.
- 56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4  
5  
6  
7  
8  
9 [39] Vigna S. A weighted correlation index for rankings with ties. In: Proceedings of  
10 the 24th international conference on World Wide Web. International World Wide  
11 Web Conferences Steering Committee; 2015. p. 1166–76.  
12 545
- 14 [40] Volkovs MN, Zemel RS. New learning methods for supervised and unsu-  
15 pervised preference aggregation. The Journal of Machine Learning Research  
16 2014;15(1):1135–76.  
17  
18
- 20 [41] Webber W, Moffat A, Zobel J. A similarity measure for indefinite rankings. ACM  
21 550 Transactions on Information Systems (TOIS) 2010;28(4):20.  
22
- 24 [42] Yilmaz E, Aslam JA, Robertson S. A new rank correlation coefficient for infor-  
25 mation retrieval. In: Proceedings of the 31st annual international ACM SIGIR  
26 conference on Research and development in information retrieval. ACM; 2008.  
27 p. 587–94.  
28  
29
- 31 555 [43] Zhu B, Xu Z, Xu J. Deriving a ranking from hesitant fuzzy preference  
32 relations under group decision making. IEEE transactions on cybernetics  
33 2014;44(8):1328–37.  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65