

# Simpson's Paradox in Nonparametric Statistical Analysis: Theory, Computation, & Susceptibility in Public Health Data

James Boudreau<sup>1</sup>, Justin Ehrlich<sup>2</sup>, and Shane Sanders<sup>3</sup>

<sup>1</sup>Department of Economics, Kennesaw State University

<sup>2,3</sup>Falk College of Sport & Human Dynamics, Syracuse University

## Abstract

This study establishes sufficient conditions for observing instances of Simpson's (data aggregation) Paradox under rank sum scoring (RSS), as used, e.g., in the Wilcoxon-Mann-Whitney (WMW) rank sum test. The WMW test is a primary nonparametric statistical test in FDA drug product evaluation and other prominent medical settings. Using computational nonparametric statistical methods, we also establish the relative frequency with which paradox-generating Simpson Reversals occur under RSS when an initial data sequence is pooled with its ordinal replicate. For each 2-sample, n-element per sample or  $2 \times n$  case of RSS considered, strict Reversals

al. (2016) find evidence of aggregation paradox instances in randomized clinical trial data. Evidence of aggregation paradox has also been found in the settings of large-scale registry data (Gron et al. 2016), meta-analyses of an academic literature (Kuss 2016), and clinical risk reclassification (Cook et al. 2017).

In one respect, the Paradox can be viewed as a robustness check on a given statistical result. When the Paradox occurs, it follows that a given result is at least partly a function of data scale or sample size. As noted, the Paradox has been shown to occur for the Wilcoxon-Mann-Whitney (WMW) Rank Sum Test. However, there exists no computational or empirical evidence as to the frequency with which instances of the Paradox occur for the WMW Rank Sum Test and little such evidence for non-parametric statistical tests overall. Are Simpson Reversals pervasive or only a marginal concern for the WMW Test? Even previous research as to the incidence of the Paradox for parametric statistical tests is scarce and provides somewhat contrasting conclusions. There are two studies that directly estimate the incidence of Simpson's Paradox for parametric tests: one pertaining to contingency tables and the other pertaining to path models. Specifically, Pavlides and Perlman (2009) find that a Simpson Reversal occurs for one-sixtieth (1.67%) of all  $2 \times 2 \times 2$  contingency tables. Kock (2015) estimates the likelihood of a Simpson Reversal in path models as approximately 12.8%.

## 2 Material and Methods

### 2.1 Rank Sum Scoring and Simpson's Aggregation Paradox: Definitions and a Theorem

Let us formally define 2-group rank sum scoring. Consider two groups, A and B. Each group is defined as a rank-ordered sequence of  $n$  individual elements, where  $n$  is some integer greater than 1 ( $n \in \mathbb{Z}^+$ ). For example, A is defined as  $A = (a_1; a_2; a_3; \dots; a_n)$ , where the element  $a_i$  represents the  $i^{\text{th}}$  ranked element in A. We define an event as an objective process of comparison that generates a complete rank-order sequence of individuals across more than one group (i.e., both within and between groups). An event might be defined as a competition or as a statistical test. Consider an event in which elements of A and B are compared. If A and B are each composed of  $n$  elements, for example, then the event generates a rank-ordered outcome sequence of  $2n$  elements. One possible outcome sequence for the case in which  $n = 3$  is  $F_{AB} = (a_1; b_1; b_2; a_2; b_3; a_3)$ . If  $a_i$  precedes  $b_j$  in the outcome sequence, we say  $a_i \succ b_j$  ( $a_i$  ranks higher than  $b_j$ ). For simplicity, we assume that rank-order equality between two elements is not possible, an outcome that would obtain given continuous measurement of underlying parameter values. For any  $a_i \in A$  and  $b_j \in B$ , that is, we have that  $a_i \succ b_j \vee b_j \succ a_i$  is a tautology.

Formally, we represent the rank of an element  $a_i \in A$  in the outcome sequence  $F_{AB}$  as  $r(a_i \mid F_{AB})$ . Let  $X_i^+(F_{AB}) = \{x \in F_{AB} : x \succ a_i\}$  be the set of elements in  $F_{AB}$  that rank better than  $a_i$ . Then,  $r(a_i \mid F_{AB}) = |X_i^+(F_{AB})| + 1$ . From elemental rankings, we generate a rank sum score for each group as follows. The respective scores for A and B for outcome sequence  $F_{AB}$  are  $S(A \mid F_{AB}) = \sum_{a_i \in A} r(a_i \mid F_{AB})$  and  $S(B \mid F_{AB}) = \sum_{b_j \in B} r(b_j \mid F_{AB})$ , where it must be that  $S(A \mid F_{AB}) + S(B \mid F_{AB}) = \frac{2n(2n+1)}{2}$ . That is, the sum of ranks for a  $2n$  element sequence simply equals the sum of integers from 1 to  $2n$ . Merits 494 Tf

$B \succ_{FF^0} A$  (i.e., that A ranks strictly higher than B in F and  $F^0$ , but B ranks strictly higher than A in  $FF^0$ ) or that  $B \succ_F A$  and  $B \succ_{F^0} A$ , but  $A \succ_{FF^0} B$  (i.e., that B ranks strictly higher than A in F and  $F^0$ , but A ranks strictly higher than B in  $FF^0$ ).

**Table 1: Sufficient Condition for Presence of Reversal and Observation of at Least One**

Sufficient Condition for Presence of Reversal	Observation of at least one Reversal		
		F	T
	F	162	12
T	0	78	

**Table 2: Sufficient Condition for Absence of Reversal and Observation of at Least One**

Sufficient Condition for Absence of Reversal	Observation of at least one Reversal		
		F	T
	F	84	90
T	78	0	

**Table 3: Sufficient Condition for Absence of Reversal and Observation of at Least One**

Sufficient Condition for Absence of Reversal	Observation of at least one Reversal		
		F	T
	F	364	520
T	220	0	

Of the 252 initial sequences,  $F_{AB}$ , Table 1 tells us that the sufficient condition for presence of at least one Reversal across all poolings of  $F_{AB}$  and  $F_{AB}^0$  holds for 78 of those sequences. Empirically, we observe at least one Reversal for each of those sequences. The second table shows that for a distinct 78 of the 252 initial  $2 \times 5$  sequences,  $F_{AB}$ , the sufficient condition for absence of Reversals across all poolings of  $F_{AB}$  and  $F_{AB}^0$  holds. Empirically, we do not observe a Reversal in any of those sequences. For the  $2 \times 5$  case, then, the sufficient conditions from Theorems 1 and 2 assure us whether or not Reversal is possible for 156 of the 252 initial sequences (61.9%).

Tables 3 and 4 deal with sufficient conditions for the  $2 \times 6$  case. Of the 924 initial sequences,  $F_{AB}$ , for the  $2 \times 6$  case, Table 3 shows that the sufficient condition for presence of at least one Reversal across all possible poolings of  $F_{AB}$  and  $F_{AB}^0$  holds for 220 of those sequences. Empirically, we observe at least one Reversal for each of those sequences. The fourth table shows that for a distinct 364 of the 924 initial  $2 \times 6$  sequences,  $F_{AB}$ , the sufficient condition for absence of Reversals across all poolings of  $F_{AB}$  and  $F_{AB}^0$  holds. Empirically, we do not observe a Reversals for any of those sequences. For the  $2 \times 6$  case, then, the sufficient conditions from Theorems 1 and 2 assure us whether or not Reversal is possible for 584 of the 924 initial sequences (63.2%). In each observed case, the sufficient conditions determine unambiguously whether an initial sequence is susceptible to Reversal in more than three-fths of cases. Therefore, we can usually assess the general robustness of a rank sum result in terms of

susceptibility to Simpson Reversals. As such an assessment can determine whether a given result is scale-variant, we conclude that Theorems 1 and 2 can usually combine to offer a "quick and dirty" robustness check on a rank sum result.

## 2.2 The Sample Space: A Combinatorial Description

For the  $2 \times n$  case, there are  $\frac{(2n)!}{(n!)^2}$  initial sequences,  $F$ . We are arranging  $2n$  elements |  $n$  elements from each of 2 groups | where we do not distinguish between respective objects of a given group. For each initial sequence, we then ask in how many ways  $F$  can be pooled with its ordinal replicate,  $F^0$ . This is equivalent to a "stars and bars" combinatorial problem, in which we are placing  $2n$  "stars" or elements from  $F^0$  into  $2n$  "bars" or potential pooling positions amongst the elements of  $F$ . From this characterization, there are  $\frac{(4n)!}{((2n)!)^2}$  poolings for each initial sequence and  $\frac{(2n)!}{(n!)^2}$  initial sequences. The number of poolings for a given  $2 \times n$  case equals the product of the number of initial sequences and the number of poolings per initial sequence, or  $\frac{(2n)!}{(n!)^2} \cdot \frac{(4n)!}{((2n)!)^2}$ , for each case,  $2 \times n$ . For example, in the  $2 \times 7$  case, there are  $\frac{(2 \cdot 7)!}{(7!)^2} = 3,432$  initial sequences,  $F$ . Moreover, there are  $\frac{(4 \cdot 7)!}{((2 \cdot 7)!)^2} = 40,116,600$  poolings per initial sequence. As such, there are  $3,432 \cdot 40,116,600$  or approximately 137.68 billion possible poolings for the  $2 \times 7$  case. We provide the sample space for each  $2 \times n$  case in Table 5 of the subsequent section.

## 2.3 Computational Methods Materials

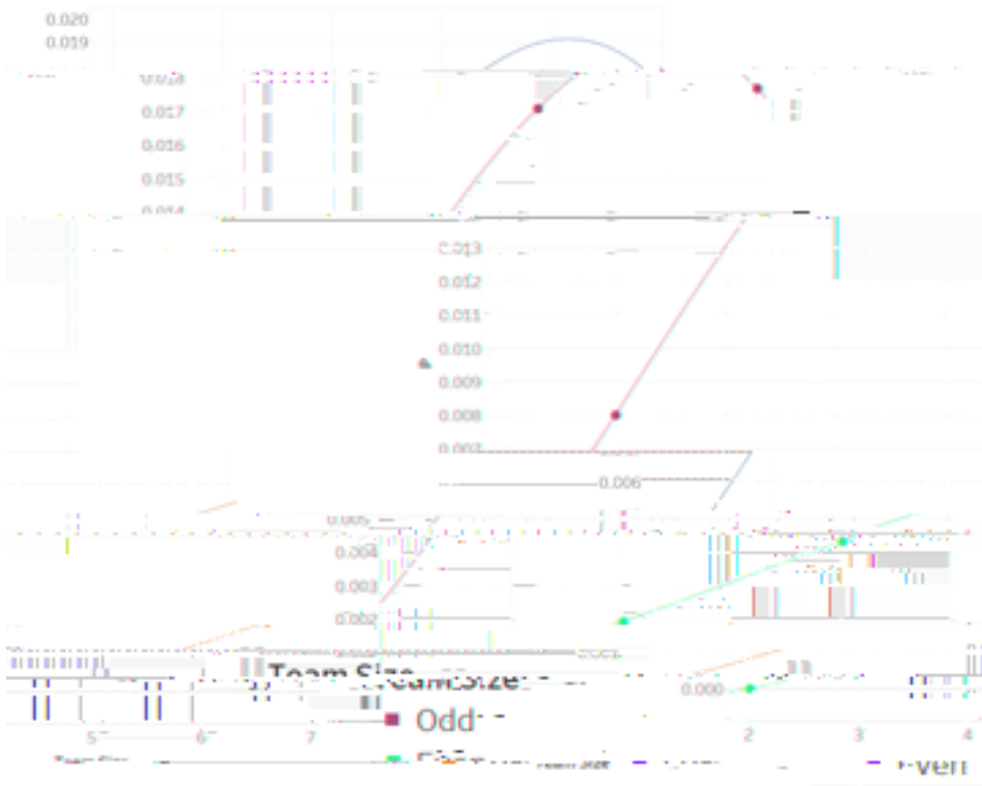
We wrote a computational algorithm in Java by which to search the sample space of each case where  $0 < n(2 \cdot Z^+) < 7$ : The algorithm is shown in Appendix 1. It programmatically generates all possible initial sequences,  $F_{AB}$  ( $F_{AB}^0$ ), for a case, then creates all possible pooled sequences,  $FF_{AB}^0$ ; for each initial sequence. For each initial sequence, rank sum scores for A and B are computed. This scoring task is then repeated for each pooling  $FF_{AB}^0$  of  $F_{AB}$  and  $F_{AB}^0$  and iteratively for each pooling of each initial sequence. Then, instances of Simpson Reversal are checked using the

**Table 5: Relative Frequency of Simpson Reversal by Case**

Case	Score	Reversal	Frequency	Score	Reversal	Frequency
2x1	0.000	0	0.000	0.000	0	0.000
2x2	0.000	0	0.000	0.000	0	0.000
2x3	0.000	0	0.000	0.000	0	0.000
2x4	0.000	0	0.000	0.000	0	0.000
2x5	0.000	0	0.000	0.000	0	0.000
2x6	0.000	0	0.000	0.000	0	0.000
2x7	0.000	0	0.000	0.000	0	0.000
2x8	0.000	0	0.000	0.000	0	0.000

We observe that Simpson Reversals are not possible for sufficiently small  $n$  (i.e.,  $n < 3$ ). In the context of Theorem 1, the largest possible  $n$  is not sufficiently large to motivate a strict Simpson Reversal in these cases. For the  $2 \times 1$  and  $2 \times 2$  cases, a group that is strictly outranked in  $F_{AB}$  cannot have a positive  $\alpha$ , and therefore a strict Simpson Reversal is not possible for these cases. We can also consider computed cases where  $n > 2$ . From even to odd case, the results suggest a wavelike movement in the likelihood of a Simpson Reversal. In general, there is a lower likelihood of strict Simpson Reversal in even cases than in neighboring odd cases due to the possibility of ties for  $n$ -even cases of pairwise rank sum scoring (but not for  $n$ -odd cases). With some probability mass allowing for a pairwise tie in the  $n$ -even case, strict Simpson Reversals are less likely. This result also holds for other social choice violations (e.g., violations of Transitivity and of IIA; see Boudreau et al. 2014). To evaluate the marginal effect of increases in  $n$ , as distinct from the effect of changes from even to odd case, one should compare the iterative trend between  $n$  and  $n + 2$  rather than that between  $n$  and  $n + 1$ . We do this for the even and odd cases, respectively, in Figure 1.

**Figure 1: Frequency of Simpson Reversal by Case**

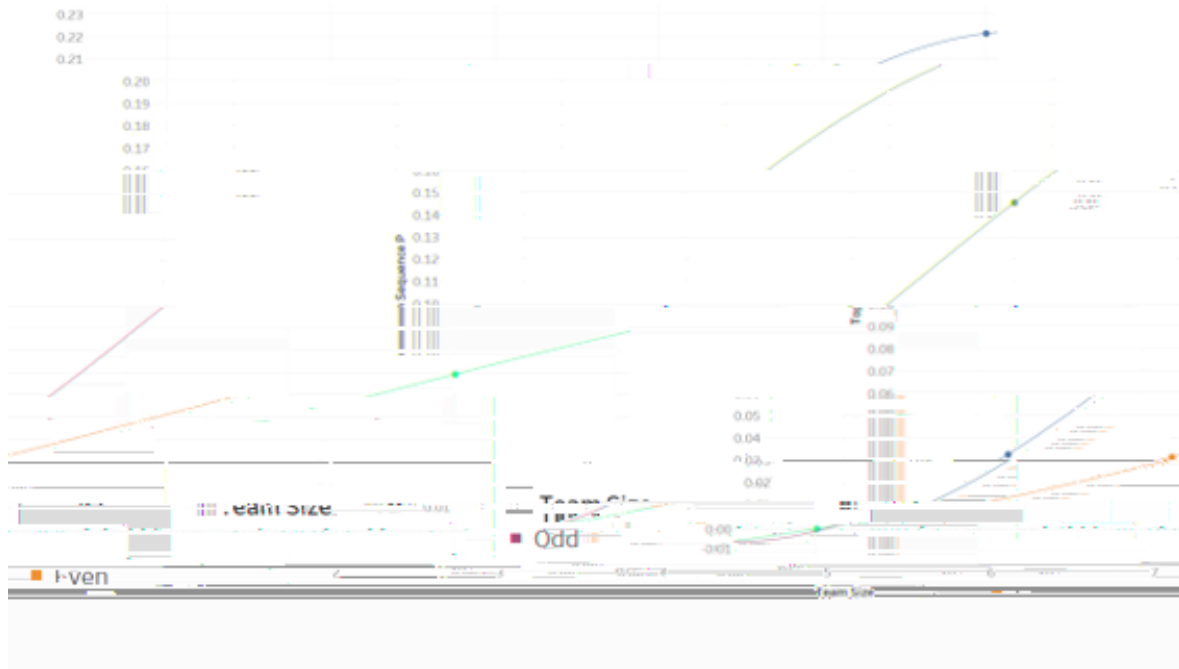


Over the set of cases computed, the relative frequency of Reversal rises for both the even and odd sets of cases. For the  $2 \times 8$  case, we run a simulation to estimate whether this trend might continue. Specifically, we randomly

select and generate one-quarter of all possible initial sequences,  $F_{AB}$ , (without replacement) for this case and then replicate each selected initial sequence. For each selected initial sequence and its replicate, we then randomly select approximately 0.1% of all possible poolings, or a little more than 600,000 poolings per sampled initial sequence. For each pooling, we check for Reversals



**Figure 3: Highest Initial Sequence Level Reversal Likelihood by Case**



Team Size	Team Class	Reversal Likelihood	Sequence	Closest Margin of Victory	
2	1	0/6	NA	NA	
2	2	0/10	NA	NA	
2	3	30/924 = 3.25%	abbaab	10 - 11	
2	4	402/12376 = 3.12%	abbbabab	19 - 17	
2	2	5	26,872/184,756 = 14.54%	azbbbbazab	27 - 25

Table 6 shows that reversals are more likely given sequences that feature both close rank sum scores and uninterrupted clusters of one group and then of the other within the rank sequence. Note that the maximum Reversal likelihood generating sequence for each case is not unique. In each case, one could transpose the elements 'a' and the elements 'b' to obtain the same Reversal likelihood. We find that the maximum Reversal likelihood generating sequence also generates the closest margin of victory in each case (i.e., 1 rank sum unit for n-odd cases and 2 rank sum units for n-even cases). While the overall likelihood of Reversal

Rank in Set (Lower SARs rank higher)	Phone Type	SAR
1	Apple iPhone (4GB)	0.974
2	Apple iPhone (8GB)	0.974
3	Apple iPhone 3G (16GB)	1.38
4	Apple iPhone 3G (8GB)	1.38
5	Apple iPhone 3G-S (16GB)	1.45
6	Apple iPhone 3G-S (32GB)	1.45
7	Apple iPhone 4 (32GB)	1.17
8	Nokia E61i	0.83
9	Nokia E63	1.24
10	Nokia E65	0.74
11	Nokia E70	0.9
12	Nokia E71	1.53
13	Nokia E72	0.99
14	Nokia E75	0.99
15	Nokia E90 Communicator	0.59

From this data, we find that Nokia E Series phones from this time period rank higher than Apple iPhones in terms of emitting lower levels of radiation. The rank sum score for the 8 Nokia (Apple) phones is 62 (74). We also compare subsets of these two mobile phone series. For example, we compare the 7 (6; 5; 4; 3; 2; 1) most recently released Nokia E phones in the dataset with the 7 (6; 5; 4; 3; 2; 1) most recently released Apple iPhones. For each of these subsets, Nokia E Series phones also rank better than Apple iPhones under rank sum scoring. Given these subset results, we might expect Simpson Reversals to not occur in this application data.

In this application setting, there are two main ways in which to think of Simpson Reversals. One can think of them in the specific: Is there an alternative set of data comparing the two phone series such that, when pooled with the original data, yields a Reversal? Alternatively, one can think of them generally: For what proportion of poolings of this data and its ordinal replicate does a strict Reversal arise? Though the specific question dominates applications in the previous literature on Simpson Reversals, the general question has certain conceptual advantages. Under the general question, one can determine how globally robust a given data is against Reversal when pooled with an ordinal data that individually generates an identical test result. When one ordinally replicates a data set, no new information is introduced by which to evaluate the two groups. By definition, the original data and its ordinal replicate yield the very same rank sum test result. By considering incidence of Reversal under pooling of the two data sets, one can determine the general robustness of the original result by considering to what extent that result relies upon the interaction of the test itself with scale-variant features of the data. In the present application, therefore, we consider the general question as a means to determine the general robustness of the data against (susceptibility to) Reversal. In so doing, one can characterize the strength of the original result in terms of data scale invariance.

In the empirical exercise, we first consider the 2 groups and 8 phone types per group case (i.e., the 2 x 8 case). We sort the data from lowest to highest SAR level to obtain SAR rankings for each of the 16 phones. We then add the 8 rank positions of Apple iPhones and the 8 rank positions of Nokia E phones, respectively, to obtain each brand's empirically-observed rank sum score. We then consider each "most-recent sub-sample" of the data. That is, the 2 x 7 case is developed by rank sum scoring the 7 most recently marketed Apple iPhones in the sample against the 7 most recently marketed Nokia E phones. The same procedure was followed inductively to obtain the 2 x n case 8 n 2 f1; 2; 3; :::; 6g. For each case, rank sum scores are shown in Table 8. In Table 9, incidence of empirically observed Reversal is reported for each case.

	$N > 1$	$2 \times 1$	$\langle a_i \rangle$	2	1
$\langle a_i, i, n \rangle$		$2+3=5$	$1+4=5$	$N > 1$	$2 \times 2$
$\langle n, i, i \rangle$	$4+5+7+8=24$	$1+2+3=6=12$		$2 \times 4$	$\langle n, n, i, i \rangle$
$\langle 7+10, 1+4+5+8+9+12 \rangle$	$39$	$N > 1$	$2 \times 6$	$\langle n, i, i, n, i, i, n, i, i, n \rangle$	$2+3+6=11=39$
$\langle 9+11+14+15 \rangle$	$53$	$N > 1$	$2 \times 8$	$\langle n, i, i, n, i, i, n, i, i, n \rangle$	$5+10+13=28$
$\langle 14+15+17 \rangle$	$46$			$2 \times 8$	$\langle n, i, i, n, i, i, n, i, i, n \rangle$

Table 1: Results of the simulation for the distribution of the number of Reversals

Reversals	Strict Reversals	(if specified)	Res
0	0	$2 \times 1$	6
0	0	$2 \times 2$	70
0	0	$2 \times 6$	2,704,156
0	0	$2 \times 7$	40,116,400
$6.01 \times 10^8$	0	$2 \times 8$	

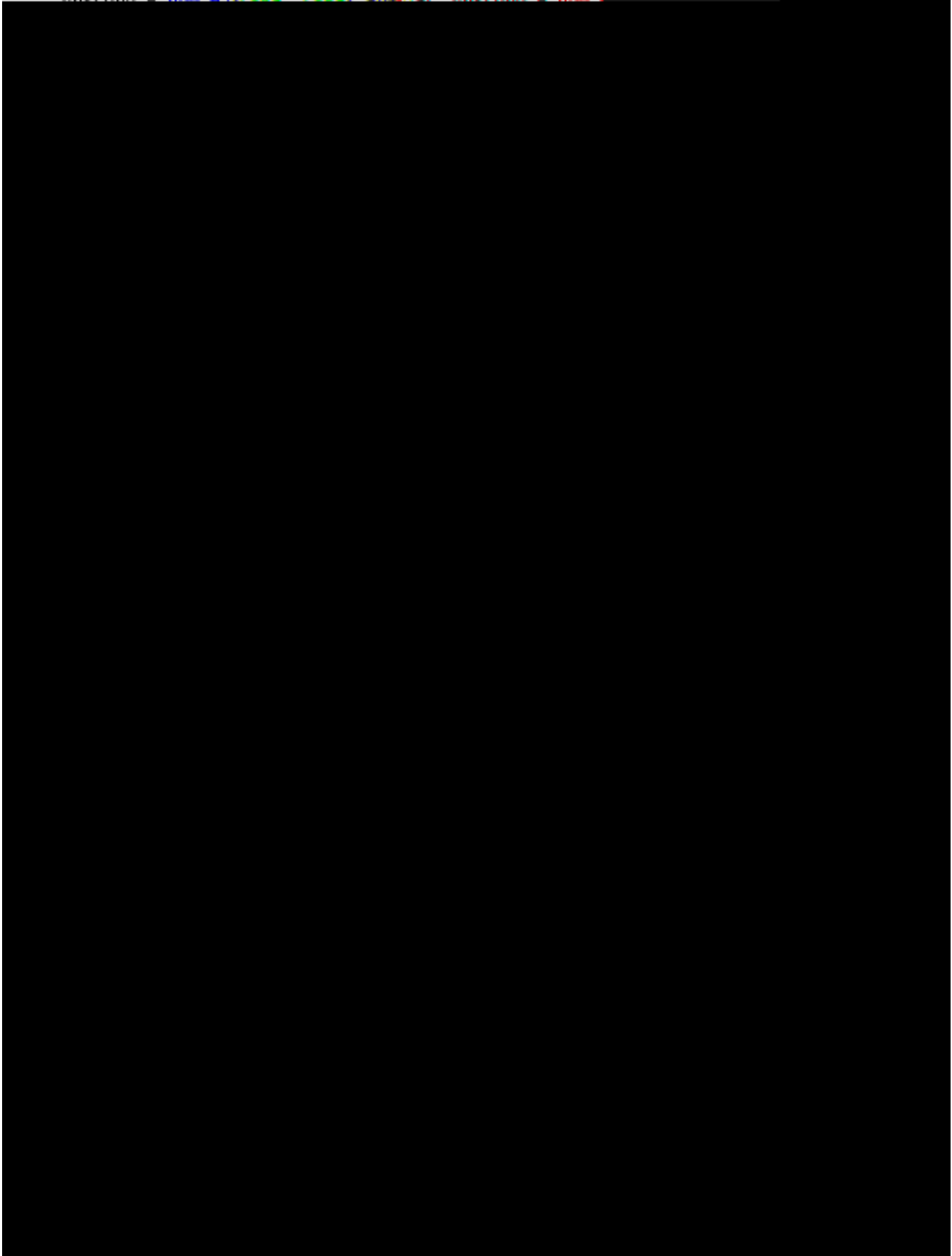
Unlike in our computational treatment, note that a single outcome for F is given (observed) in the empirical treatment. For the empirical application, then, we need only consider all possible poolings of the specified sequence, F, and its ordinal replicate, F<sup>0</sup>. In the computational section, we observed that the likelihood of a strict Reversal has a high degree of variability across initial sequences. As this application selects a single sequence F based solely on market characteristics of two cellular phone product series (e.g., similar market time period, status as a popular line of phones during that time period, and number of models in series) and not on parametric properties of the underlying data, there was no a priori reason to believe that instances of strict Reversal would occur at all in the application. For two of the 2 x n cases considered, the 2 x 1 and 2 x 2 cases, we have established that Reversals are not possible for any pooling of the data. However, we observe a cluster of Reversals



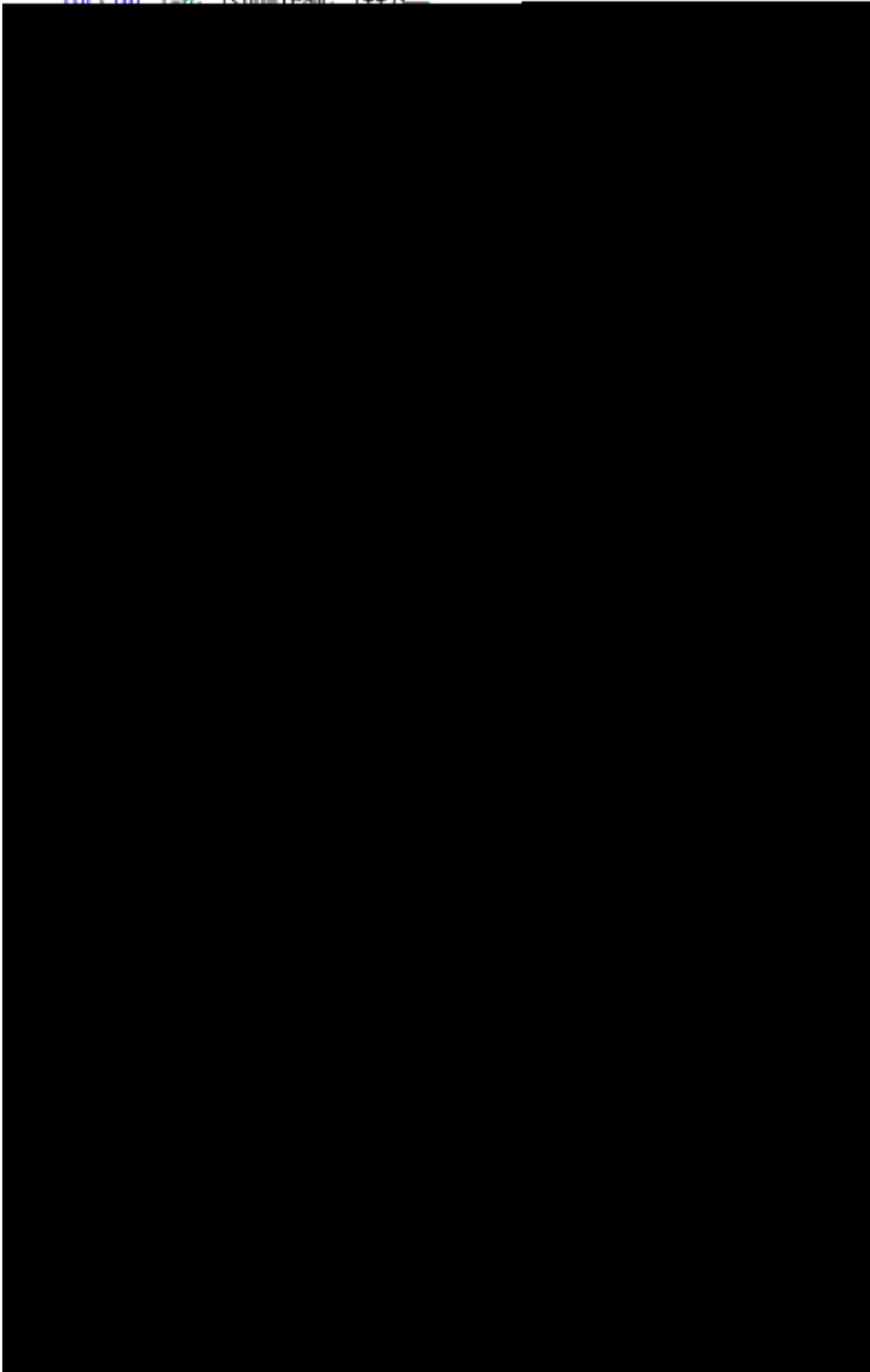
- [11] Haunsperger, D. B., & Saari, D. G. (1991). The Lack of Consistency for Statistical Decision Procedures. *The American Statistician*, 45(3), 252{255. <https://doi.org/10.1080/00031305.1991.10475814>
- [12] Huang, L., Zalkikar, J., Tiwari, R. (2019). Likelihood-ratio-test methods for drug safety signal detection from multiple clinical datasets. *Computational and mathematical methods in medicine*, 2019.
- [13] Kock, N. (2015). How Likely is Simpson's Paradox in Path Models?: *International Journal of E-Collaboration*, 11(1), 1{7. <https://doi.org/10.4018/ijec.2015010101>
- [14] Kuss, O. (2015). Statistical methods for meta-analyses including information from studies without any events-add nothing to nothing and succeed nevertheless. *Statistics in Medicine*, 34(7), 1097{1116. <https://doi.org/10.1002/sim.6383>
- [15] Lin, T., Chen, T., Liu, J., & Tu, X. M. (2021). Extending the Mann-Whitney-Wilcoxon rank sum test to survey data for comparing mean ranks. *Statistics in Medicine*, 40(7), 1705-1717.
- [16] Nagaraja, H. N., & Sanders, S. (2020). The Aggregation Paradox in Statistical Rankings. *PLOS ONE*, Forthcoming.
- [17] Pavlides, M. G., & Perlman, M. D. (2009). How Likely Is Simpson's Paradox? *The American Statistician*, 63(3), 226{233. <https://doi.org/10.1198/tast.2009.09007>
- [18] Pineiro, G., Paruelo, J. M., & Oesterheld, M. (2006). Potential long-term impacts of livestock introduction on carbon and nitrogen cycling in grasslands of Southern South America. *Global Change Biology*, 12(7), 1267{1284. <https://doi.org/10.1111/j.1365-2486.2006.01173.x>
- [19] Pordanjani, S. R., Kavousi, A., Mirbagheri, B., Shahsavani, A., & Etemad, K. (2021). Spatial analysis and geoclimatic factors associated with the incidence of acute lymphoblastic leukemia in Iran during 2006{2014: An environmental epidemiological study. *Environmental Research*, 202, 111662.
- [20] Simpson, E. H. (1951). The Interpretation of Interaction in Contingency Tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 13(2), 238{241. <https://doi.org/10.1111/j.2517-6161.1951.tb00088.x>
- [21] Tran, P., & Waller, L. (2015). Variability in results from negative binomial models for lyme disease measured at different spatial scales. *Environmental research*, 136, 373-380.
- [22] Winner, L. (2021). Cell Phone Radiation Ratings by Model/Brand: [users.stat.u.edu/~winner/datasets.html](https://users.stat.u.edu/~winner/datasets.html)
- [23] Yule, G. U. (1903). NOTES ON THE THEORY OF ASSOCIATION OF ATTRIBUTES IN STATISTICS. *Biometrika*, 2(2), 121{134. <https://doi.org/10.1093/biomet/2.2.121>



...MaxLoss, Loss, numSimonsBapadovViolationHighScore = pos.HackMaxLoss, Loss



```
for(int i=0; i<numTeam; i++){
```





```

int yScore = 0;
int counter = 0;
for (int i = 0; i < dataPoints.length(); i++) {
    counter++;
    for(int j=0; j<numTeam; j++){
        if(dataPoints.charAt(i) == groups[j]){

```

```

int minGroupsScores[0];
int minIndex = 0;
for (int i = 0; i < numTeam; i++){

```

```

//detect tie

```

```

return groupsScores[minIndex];
}
}

```

```

index(String dataPoints, int numGroups){

```

```

    numWinnerChanged;

```

```

    numWinnerChangedPossible;

```

```

    findWinner(dataPoints, numGroups);

```

```

private void findSimpsonsPe

```

```

    long prevNumWinnerChang

```

```

    long prevNumWinnerChang

```

```

    char winner = findDepas

```

```

    }

    long topScore = findDependentWinnerScore(dataPoints, numGroups);

    if (numSimponsParadoxViolationHighScore.containsKey(topScore)){
        numSimponsParadoxViolationHighScore.put(topScore, numSimponsParadoxVi
olationHighScore.get(topScore)+deltaNumWinnerChange);
    }
    else{
        numSimponsParadoxViolationHighScore.put(topScore, deltaNumWinnerChange);
    }
}

if (numSimponsParadoxViolationPossibleHighScore.containsKey(topScore)){
    numSimponsParadoxViolationPossibleHighScore.put(topScore, numSimponsP
ossibleHighScore.get(topScore)+deltaNumWinnerChange);
}
else{
    numSimponsParadoxViolationPossibleHighScore.put(topScore, deltaNumWinnerChange);
}
}

public static int countRepeats(String haystack, char needle)
{
    int count = 0;
    for (int i = 0; i < haystack.length(); i++)
    {
        if (haystack.charAt(i) == needle)
        {
            count++;
        }
    }
    return count;
}

public static int countRepeats(String haystack, char needle)
{
    int count = 0;
    for (int i = 0; i < haystack.length(); i++)
    {
        if (haystack.charAt(i) == needle)
        {
            count++;
        }
    }
    return count;
}
}

```

```
findSimpsonsParadox(dataPoints, numGroups):
```

