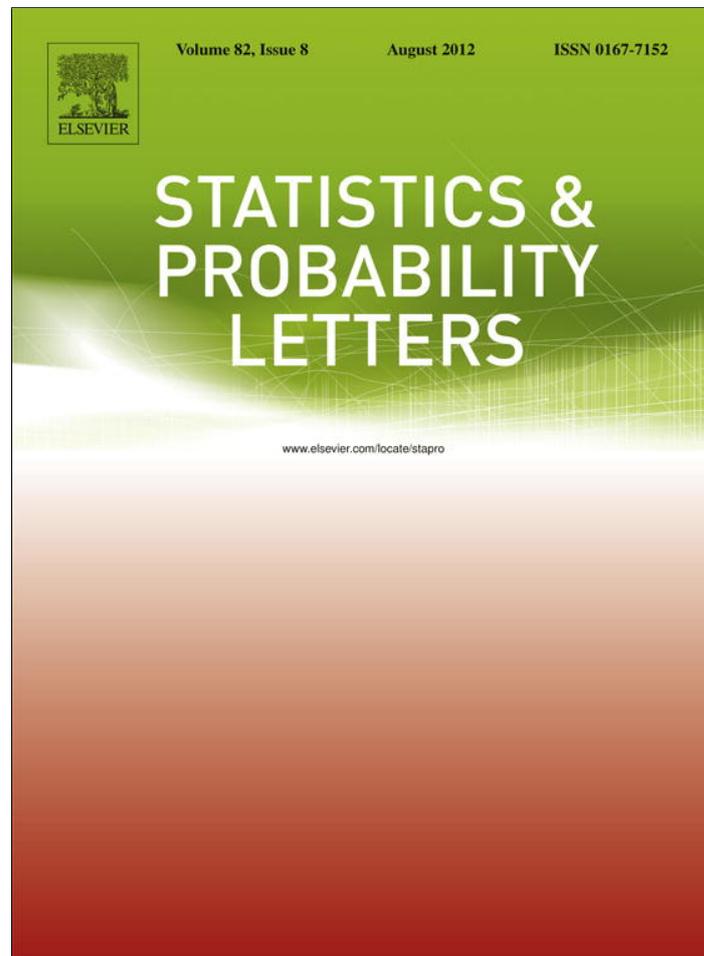


Provided for non-commercial research and education use.  
Not for reproduction, distribution or commercial use.



This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

<http://www.elsevier.com/copyright>



Contents lists available at SciVerse ScienceDirect

# Statistics and Probability Letters

journal homepage: [www.elsevier.com/locate/stapro](http://www.elsevier.com/locate/stapro)

## An inductive order construction for the difference of two dependent proportions

Weizhen Wang

Department of Mathematics and Statistics, Wright State University, Dayton, OH 45435, United States

### ARTICLE INFO

#### Article history:

Received 28 November 2011  
 Received in revised form 7 March 2012  
 Accepted 26 March 2012  
 Available online 9 May 2012

#### Keywords:

Confidence interval  
 Coverage probability  
 Multinomial distribution  
 Order  
 Set inclusion

### ABSTRACT

This paper concerns interval estimation for the difference of two dependent proportions. An order on the sample space is constructed using an inductive method; then the smallest one-sided  $1 - \alpha$  confidence interval under the order is derived. This interval is admissible under the set inclusion criterion, and is illustrated in two examples. An R-code is used to compute the proposed order and interval.

© 2012 Elsevier B.V. All rights reserved.

### 1. Introduction

The comparison of two dependent proportions,  $p_1$  and  $p_2$ , is useful in statistical practice. These proportions appear in a matched-pair experiment, which effectively removes the confounding between subject and treatment effects, and utilizes experimental units more efficiently compared with the two independent sample experiment since two measurements are collected from each unit. Researchers typically employ statistical tests to address the issue, and conduct [McNemar's test \(1947\)](#) for the hypothesis of equal proportions,  $H_0 : p_1 - p_2 = 0$ . However, if the goal of a study is to show a treatment better than the control and further find out how much better the treatment is, then it cannot be accomplished by the test, and a one-sided confidence interval is a must. The goal of this paper is to construct a smallest one-sided  $1 - \alpha$  confidence interval with a lower random limit for  $p_1 - p_2$  by carefully selecting an order on all sample points. The other one-sided interval with an upper random limit may be derived similarly.

To be precise, suppose there are  $n$  independent and identical trials in an experiment, and each trial is inspected by two criteria 1 and 2. By criterion  $i$ , each trial is classified as  $S_i$  or  $F_i$  for  $i = 1, 2$ . The numbers of trials with outcomes  $(S_1, S_2)$ ,  $(S_1, F_2)$ ,  $(F_1, S_2)$  and  $(F_1, F_2)$  are the observations, and are denoted by  $N_{11}$ ,  $N_{12}$ ,  $N_{21}$  and  $N_{22}$ , respectively. It is clear that  $\underline{X} = (N_{11}, N_{12}, N_{21})$  follows a multinomial distribution with probabilities  $p_{11}$ ,  $p_{12}$ ,  $p_{21}$ , respectively. Let  $p_i = P(S_i)$  be the two dependent proportions. We are interested in estimating a difference

$$\theta_D \stackrel{\text{def}}{=} p_1 - p_2 = (p_{11} + p_{12}) - (p_{11} + p_{21}) = p_{12} - p_{21} \quad (1)$$

with one-sided  $1 - \alpha$  confidence interval of form  $[L, 1]$  since only the lower bound of  $p_1 - p_2$  should be controlled. An example of this setting is to compare a new eye medicine with a control, where the new medicine and the control are applied to each eye for each of  $n$  patients, and  $N_{12}$ , for example, is the number of patients who show improvement when using the new medicine but show no improvement when using the control.

E-mail address: [weizhen.wang@wright.edu](mailto:weizhen.wang@wright.edu).

To make the order and the interval construction easier, we now reduce the number of parameters and simplify the sample space. Let  $T = N_{11} + N_{22}$  and  $p_T = p_{11} + p_{22}$ , and take a look at the conditional probability mass function of  $(N_{11}, N_{12}, N_{21})$  for given  $(N_{12}, T)$

$$p(n_{11}, n_{12}, n_{21} | n_{12}, t) = \frac{t!}{n_{11}!n_{22}!} \left(\frac{p_{11}}{p_T}\right)^{n_{11}} \left(\frac{p_{22}}{p_T}\right)^{n_{22}}, \tag{2}$$

where  $p_{22} = 1 - p_{11} - p_{12} - p_{21}$ . This function does not involve  $p_{12}$  and  $p_{21}$ . Therefore, following a similar logic of sufficient statistic, we consider intervals for  $\theta_D$  of form  $[L(\underline{Z}), 1]$ , where  $\underline{Z} = (N_{12}, T)$  also follows a multinomial distribution with probabilities  $p_{12}$  and  $p_T$ . The simplified sample space is

$$S_D = \{(n_{12}, t) : 0 \leq n_{12} + t \leq n\} \tag{3}$$

with a reduced parameter space

$$H_D = \{(p_{12}, p_T) : 0 \leq p_{12}, p_T \leq 1; p_{12} + p_T \leq 1\}.$$

Note  $\theta_D = 2p_{12} + p_T - 1$ . Rewrite  $H_D$  as

$$H_D = \{(\theta_D, p_T) : p_T \in D(\theta_D), -1 \leq \theta_D \leq 1\},$$

where

$$D(\theta_D) = \{p_T : 0 \leq p_T \leq \min\{1 + \theta_D, 1 - \theta_D\}\}. \tag{4}$$

The probability mass function of  $(N_{12}, T)$  in terms of  $\theta_D$  and  $p_T$  is

$$p_D(n_{12}, t; \theta_D, p_T) = \frac{n!}{n_{12}!t!n_{21}!} \left(\frac{1 + \theta_D - p_T}{2}\right)^{n_{12}} p_T^t \left(\frac{1 - \theta_D - p_T}{2}\right)^{n_{21}}. \tag{5}$$

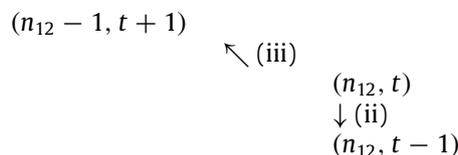
Buehler (1957) first proposed an idea to construct a one-sided confidence interval based on an order  $\leq$  on a sample space  $S$ . Order  $\leq$  is a binary relation among sample points  $x'$ 's that satisfies antisymmetry, transitivity, totality and sets  $U_x = \{x' \in S : x' \leq x\}$  and  $V_x = \{x' \in S : x' \equiv x\}$  both  $\sigma$ -measurable. See [http://en.wikipedia.org/wiki/Total\\_order](http://en.wikipedia.org/wiki/Total_order) for the first three properties. Chen (1993) and Lloyd and Kabaila (2003) extended the result to a general case. Wang (2010) rediscovered this technique and derived a smallest confidence interval for the difference of two independent proportions. However, this general approach requires a predetermined order on  $S$  to be valid. In the next section, we first generate an order using an inductive construction method, then obtain the corresponding smallest one-sided interval for  $p_1 - p_2$ . Two examples are given to illustrate the construction of order and interval.

## 2. An inductive order construction

Let  $\leq_D$  denote the order to be constructed on  $S_D$ . The purpose of such an order is to rank the lower confidence limit  $L$  over  $S_D$ . Here are several natural rules for  $\leq_D$ :

- (i)  $(n_{12}, t) = (n, 0)$  is the largest point on  $S_D$  under  $\leq_D$ ;
- (ii)  $(n_{12}, t) \leq_D (n_{12}, t')$ , i.e.,  $(n_{12}, t)$  is not larger than  $(n_{12}, t')$ , if  $t \leq t'$ ;
- (iii)  $(n_{12}, t) \leq_D (n'_{12}, t')$  if  $n_{12} + t = n'_{12} + t'$  and  $n_{12} \leq n'_{12}$ .

Rule (i) is obvious since  $(n, 0)$  tends to yield the largest estimate of  $\theta_D$ . If  $T$  increases but  $N_{12}$  does not change, then  $N_{21}$  decreases. So rule (ii) should be true. If  $n_{12} + t = n'_{12} + t'$ , then  $N_{21}$  remains unchanged. Therefore, rule (iii) is true due to  $n_{12} \leq n'_{12}$ . The following diagram shows the relation of three points:



where “ $\leftarrow$ ” means “ $\leq_D$ ”. These rules simplify the construction of  $\leq_D$ . However, the determination of a larger point between  $(n_{12} - 1, t + 1)$  and  $(n_{12}, t - 1)$  requires a numerical evaluation, described in step  $k + 1$ -c) below. Let  $R_k$  be a subset of  $S_D$  that contains the  $k$ th largest point(s) under  $\leq_D$ . We now introduce an inductive method to determine all  $R_k$ 's, which is equivalent to construct order  $\leq_D$ .

Step 1: For  $k = 1$ ,  $(n, 0)$  is the largest under  $\leq_D$  due to rule (i). So  $R_1 = \{(n, 0)\}$ . Let  $S_k$  be the subset of  $S_D$  on which order  $\leq_D$  is defined up to step  $k$ . So  $S_1 = R_1$ .

...

Step  $k$ : For  $k \geq 1$ , suppose order  $\leq_D$  is defined on  $S_k = \cup_{i=1}^k R_i$ . Thus,  $S_k$  contains the largest through  $k$ th largest points in  $S_D$ . Note that the construction of  $\leq_D$  is complete if  $S_{k_0} = S_D$  for some positive integer  $k_0$ .

Step  $k + 1$ : Now we determine  $R_{k+1}$  that contains the  $(k + 1)$ th largest point(s) under  $\leq_D$  in  $S_D$ .

Step  $k + 1$ -a): For each point  $\underline{z} = (n_{12}, t) \in S_D$ , let

$$N_{\underline{z}} = \{(n_{12}, t - 1), (n_{12} - 1, t + 1)\} \cap S_D$$

be the neighbor set of  $\underline{z}$  that contains points next to but smaller than  $\underline{z}$  under order  $\preceq_D$  as shown in the diagram. Let

$$N_k = [\cup_{\underline{z} \in S_k} N_{\underline{z}}] \cap S_k^c = \{(n_{12}, t) : (n_{12}, t + 1) \in S_k \text{ or } (n_{12} + 1, t - 1) \in S_k\} \cap S_k^c \tag{6}$$

be the neighbor set of  $S_k$ .

Step  $k + 1$ -b): Some points in  $N_k$  should not be in  $R_{k+1}$  due to rules (ii) and (iii). To eliminate these points, consider the following candidate set

$$C_k = \{\underline{z} = (n_{12}, t) \in N_k : (n_{12}, t + 1) \notin N_k, (n_{12} + 1, t - 1) \notin N_k\}, \tag{7}$$

a subset of  $N_k$ , and from which  $R_{k+1}$  is to be selected. In many cases,  $C_k$  is much smaller than  $N_k$ . Therefore, selecting  $R_{k+1}$  from  $C_k$  rather than  $N_k$  is simpler. The three sets have a relation:  $R_{k+1} \subseteq C_k \subseteq N_k$ .

Step  $k + 1$ -c): For each point  $\underline{z}' \in C_k$ , solve the following equation:

$$f_{\underline{z}'}(\theta_D) = \inf_{p_T \in D(\theta_D)} \left[ 1 - \sum_{(n_{12}, t) \in S_k \cup \underline{z}'} p_D(n_{12}, t; \theta_D, p_T) \right] = 1 - \alpha. \tag{8}$$

Define

$$L(\underline{z}') = \begin{cases} \text{the smallest solution of (8),} & \text{if there exists a solution for (8);} \\ -1, & \text{otherwise.} \end{cases} \tag{9}$$

Then

$$R_{k+1} = \{\underline{z} \in C_k : L_D^*(\underline{z}) = \max\{L(\underline{z}') : \underline{z}' \in C_k\}\} \quad \text{and} \quad S_{k+1} = \cup_{i=1}^{k+1} R_i. \tag{10}$$

Since  $S_D$  is a finite set and  $S_k$  is strictly increasing in  $k$ , eventually,  $S_{k_0} = S_D$  for some positive integer  $k_0 (\leq (n + 1)(n + 2)/2)$ . The construction of  $\preceq_D$  is complete.

Once an ordered partition  $\cup_{k=1}^{k_0} R_k$  (equivalent to order  $\preceq_D$ ) on  $S_D$  is constructed, the smallest  $1 - \alpha$  confidence interval  $[L_S(N_{12}, T), 1]$  for  $\theta_D$  can be automatically computed in an interval class  $\mathcal{B}$  following Lemma 1. Class  $\mathcal{B}$  contains all  $1 - \alpha$  confidence intervals of form  $[L(N_{12}, T), 1]$ , where  $L$  is a constant on each  $R_k$  and  $L$  at  $R_k$  is not smaller than  $L$  at  $R_{k+1}$  for any  $k < k_0$ .

**Lemma 1.** For  $\alpha \in (0, 1)$ ,  $k \leq k_0$  and  $\underline{z} = (n_{12}, t) \in R_k$ , let

$$f_k(\theta_D) = \inf_{p_T \in D(\theta_D)} \left[ 1 - \sum_{\underline{z}' \in S_k} p_D(\underline{z}'; \theta_D, p_T) \right], \tag{11}$$

where  $D(\theta_D)$  is given in (4),  $p_D$  is given in (5) and  $S_k = \cup_{i=1}^k R_i$ , and let  $G_k = \{\theta_D \in [-1, 1] : f_k(\theta_D) = 1 - \alpha\}$ . Define

$$L_D(\underline{z}) = \begin{cases} \inf G_k, & \text{if } G_k \neq \emptyset; \\ -1, & \text{otherwise.} \end{cases} \tag{12}$$

Then, (a) interval  $[L_D(N_{12}, T), 1]$  belongs to  $\mathcal{B}$ ; (b)  $[L_D(N_{12}, T), 1]$  is the smallest in  $\mathcal{B}$ , i.e., for any interval  $[L(N_{12}, T), 1] \in \mathcal{B}$ ,  $L(N_{12}, T) \leq L_D(N_{12}, T)$ .

This lemma extends Theorem 4 in Wang (2010), but it is easier to implement in terms of computation. The proof is similar and is then omitted.

**Remark 1.** The idea of selecting  $R_{k+1}$  given that  $R_1$  up to  $R_k$  are already selected is to pick a sample point  $\underline{z}$  in  $C_k$  that yields the largest possible lower confidence limit, i.e., at step  $k + 1$ , we choose a point with the largest confidence limit as the  $(k + 1)$ th largest point under order  $\preceq_D$ . This choice makes  $[L_D, 1]$  admissible under the set inclusion criterion (see Wang, 2006), i.e., for any  $1 - \alpha$  interval  $[L, 1]$ ,  $L$  is identical to  $L_D$  if  $L \geq L_D$ .  $\square$

**Remark 2.** If  $R_k$  contains only one point, then the smallest interval under  $\preceq_D$   $[L_D(\underline{z}), 1]$  is equal to  $[L_D^*(\underline{z}), 1]$  on  $R_k$ , where  $L_D^*(\underline{z})$  is defined in (10). If  $R_k$  contains one point for all  $k$ , then  $L_D(\underline{z}) = L_D^*(\underline{z})$  for all  $\underline{z} \neq (n, 0)$ . Note that  $L_D^*(n, 0)$  is not introduced in the construction, but it is easy to find  $L_D(n, 0)$  following Lemma 1.  $\square$

**Remark 3.** Order  $\preceq_D$  is level dependent, i.e., the order generated for a confidence level  $1 - \alpha_1$  may not be identical to that for  $1 - \alpha_2$  when  $\alpha_1 > \alpha_2$ . This may happen for any distribution. For example, it is easy to see that order  $\preceq_{L_\alpha}$ , that corresponds to the well known one-sided t-interval  $[L_\alpha(\bar{X}, S), +\infty) = [\bar{X} - t_\alpha \frac{S}{\sqrt{n}}, +\infty)$ , is also level dependent, where order  $\preceq_{L_\alpha}$  is

defined as:

$$(\bar{x}_1, s_1) \preceq_{L_\alpha} (\bar{x}_2, s_2) \text{ if and only if } L_\alpha(\bar{x}_1, s_1) \leq L_\alpha(\bar{x}_2, s_2). \tag{13}$$

In consequence, when a discrete distribution is involved, as in our case, and two  $\alpha$ 's are close, the  $1 - \alpha_1$  interval might not be included in the  $1 - \alpha_2$  interval at certain sample points. If the inclusion is violated, then we say that the  $1 - \alpha_1$  interval is not nested in the  $1 - \alpha_2$  interval. It seems no guarantee that the proposed interval is always nested between any two levels. As regards practice, however, the three commonly used 90%, 95% and 99% intervals are still nested evidenced in a limited numerical study for any  $n \leq 32$  even though different orders are detected at three levels. In short, the confidence level plays a more important role than the order to determine the interval. See more details in [Examples 1 and 2](#).  $\square$

**Example 1.** For illustration purpose, we show the details of the construction of  $\preceq_D$  (i.e., an ordered partition  $\cup_{i=1}^{k_0} R_k$  of  $S_D$ ) and the corresponding smallest 95% confidence interval  $[L_D(N_{12}, T), 1]$  for  $\theta_D$  when  $n = 4$ .

Step 1: Let  $R_1 = \{(4, 0)\}$  be the set containing the largest point under  $\preceq_D$ . We determine  $L_D(4, 0)$  by finding the smallest solution of equation

$$f_1(\theta_D) = \inf_{p_T \in D(\theta_D)} (1 - p_D(4, 0; \theta_D, p_T)) = 0.95$$

and obtain  $L_D(4, 0) = -0.05427$  due to (11) and (12). This is done in two mini-steps. First, for each fixed value of  $\theta_D$ , find the infimum by numerically evaluating  $1 - p_D(4, 0; \theta_D, p_T)$  for each  $p_T$  from 0 to  $\min\{1 + \theta_D, 1 - \theta_D\}$  with a small increment. Second, solve  $L_D(4, 0)$  by evaluating each infimum as a function of  $\theta_D$  from  $\theta_D = -1$  through  $\theta_D = 1$  with a small increment as well.

Step 2: Note that  $N_1 = C_1 = \{(3, 1)\}$  only contains a single point. Then  $R_2 = \{(3, 1)\}$  and  $L_D(3, 1) = -0.1073$ , which is the smallest solution of equation

$$f_2(\theta_D) = \inf_{p_T \in D(\theta_D)} (1 - p_D(4, 0; \theta_D, p_T) - p_D(3, 1; \theta_D, p_T)) = 0.95.$$

Step 3:  $N_2 = C_2 = \{(3, 0), (2, 2)\}$  by steps  $k + 1$ -a) and  $k + 1$ -b). It contains multiple points. Following step  $k + 1$ -c), for each point  $z' \in C_2$ , let  $L_D^*(z')$  be the smallest solution of equation

$$\inf_{p_T \in D(\theta_D)} (1 - p_D(4, 0; \theta_D, p_T) - p_D(3, 1; \theta_D, p_T) - p_D(z'; \theta_D, p_T)) = 0.95.$$

We obtain  $L_D^*(3, 0) = -0.5029$  and  $L_D^*(2, 2) = -0.2240$ , respectively. Since  $L_D^*(2, 2)$  is larger,  $(2, 2)$  is the third largest point (i.e.,  $R_3 = \{(2, 2)\}$  contains a single point), and  $L_D(2, 2) = L_D^*(2, 2) = -0.2240$  by [Remark 2](#).

The details for the entire order construction are given in [Table 1](#). In many cases,  $C_k$  is a much smaller subset of  $N_k$ . For example, in Step 5 ( $k = 4$ ),  $N_4$  contains 4 points, but  $C_4$  has only two points due to rules (ii) and (iii), and  $R_5 = \{(3, 0)\}$ . The three sets are shown below. This results in a simpler determination of  $R_5$ . Order  $\preceq_D$

$t$	$(n_{12}, t) : n_{12} + t \leq 4$	$n_{12}$			
		0	1	2	3
4	$N_4, C_4$	*	*	*	*
3	-	$R_4$	*	*	*
2	-	$N_4$	$R_3$	*	*
1	-	-	$N_4$	$R_2$	*
0	-	-	-	$N_4, C_4, R_5$	$R_1$

and the corresponding 95% smallest interval are reported in the second and fifth columns of [Table 1](#), respectively.

We also report  $\preceq_D$  and its smallest interval at the level of 90% in the last two columns of [Table 1](#). The ranks of point  $(2, 1)$  are different under the two orders, indicating the dependence of order and level. However, the 90% interval is still nested in the 95% interval on all sample points as desired because the two levels are not close and the level effect on interval overwrites the order effect.  $\square$

**Example 2.** For the sake of practice, we apply the proposed interval  $[L_D, 1]$  to a real data set presented in [Karacan et al. \(1976\)](#), where 32 marijuana users are compared with 32 matched controls with respect to their sleeping difficulties. So  $n = 32$ . The original data set is displayed in [Table 1](#) of their paper and is summarized below.

Sleeping difficulties using marijuana	Sleeping difficulties using controls	
	No	Yes
No	16(= $n_{11}$ )	9(= $n_{12}$ )
Yes	3(= $n_{21}$ )	4(= $n_{22}$ )

The researchers wish to see how much more help the marijuana use provides for sleeping by using a 95% confidence interval  $[L_D(n_{12}, t), 1]$  for  $\theta_D = p_1 - p_2$  at  $(n_{12}, t) = (9, 20)$ , where  $p_1$  is the proportion of marijuana users who have sleeping

**Table 1**  
The details to construct  $\leq_D$  and  $L_D(z)$  with three restrictions when  $n = 4$  and  $1 - \alpha = 0.95$  and  $0.9$  in Example 1.

$k$	95%				90%	
	$R_k$	$N_k$	$C_k$ $L_D^*(z')$	$L_D(z)$	$R_k$	$L_D(z)$
1	(4, 0)	(3, 1)	(3, 1) -0.1073,	-0.05427	(4, 0)	0.1246
2	(3, 1)	(3, 0), (2, 2)	(3, 0), (2, 2) -0.5029, -0.2240,	-0.1073	(3, 1)	0.06164
3	(2, 2)	(1, 3), (2, 1) (3, 0)	(1, 3), (3, 0) -0.3528, -0.5029	-0.2240	(2, 2)	-0.07709
4	(1, 3)	(0, 4), (1, 2) (2, 1), (3, 0)	(0, 4), (3, 0) -0.5272, -0.5029	-0.3528	(1, 3)	-0.2303
5	(3, 0)	(0, 4), (1, 2) (2, 1)	(0, 4), (2, 1) -0.5272, -0.5428	-0.5029	(3, 0)	-0.3592
6	(0, 4)*	(0, 3), (1, 2) (2, 1)	(2, 1) -0.5431	-0.5272	(2, 1)*	-0.4111
7	(2, 1)*	(0, 3), (1, 2) (2, 0)	(1, 2), (2, 0) -0.6302, -0.8049	-0.5431	(0, 4)*	-0.4377
8	(1, 2)	(0, 3), (1, 1) (2, 0)	(0, 3), (2, 0) -0.7515, -0.8049	-0.6302	(1, 2)	-0.5244
9	(0, 3)	(0, 2), (1, 1) (2, 0)	(2, 0) -0.8049	-0.7515	(0, 3)	-0.6796
10	(2, 0)	(0, 2), (1, 1)	(1, 1) -0.8315	-0.8049	(2, 0)	-0.7150
11	(1, 1)	(0, 2), (1, 0)	(0, 2), (1, 0) -0.9025, -0.9746	-0.8315	(1, 1)	-0.7543
12	(0, 2)	(0, 1), (1, 0)	(1, 0) -0.9746	-0.9025	(0, 2)	-0.8575
13	(1, 0)	(0, 1)	(0, 1) -0.9873	-0.9746	(1, 0)	-0.9481
14	(0, 1)	(0, 0)	(0, 0) -1	-0.9873	(0, 1)	-0.9741
15	(0, 0)			-1	(0, 0)	-1

\* indicates a different order.

improved, and  $p_2$  is the proportion in the controls. An R-code yields  $[L_D(9, 20), 1] = [0.00610, 1]$  following the steps described before. A small positive lower limit suggests that marijuana use do provide a little help to sleeping. This takes approximately 50 min on an HP desktop with Intel(R) Core(TM) 2 Quad CPU Q9300 @ 2.50 GHz and 8 GB RAM. The computing time severely depends on the rank of a sample point. If a point has a rank either small or large, then it is very quick to obtain the interval at the point. Point (9,20), however, has a rank of 172. We also compute the 90%, 95% and 99% intervals on all 561 sample points, 1683 intervals in total, for approximately 18 h. They are nested, and in particular the three intervals are equal to  $[0.03778, 1]$ ,  $[0.00610, 1]$  and  $[-0.06781, 1]$ , respectively, at (9, 20).

Due to  $n > 30$  one may consider using the following approximate 95% interval  $[L_A(n_{12}, t), 1]$ , where

$$L_A(n_{12}, t) = \frac{N_{12} - N_{21}}{n} - z_{0.05} \sqrt{\frac{1}{n} \left[ \frac{N_{12} + N_{21}}{n} - \left( \frac{N_{12} - N_{21}}{n} \right)^2 \right]}$$

The interval at (9,20) equals  $[0.01799, 1]$ , which seems attractive. However, this interval is invalid, and has a zero minimum coverage probability (a 0% confidence level) due to two facts: (i) the interval at  $(n, 0)$  is  $[1, 1]$  and reduces to a point; (ii) the coverage probability at a sequence of parameter values of  $(p_{12}, p_T) = (1 - \frac{1}{m}, 0)$  approaches to zero when  $m$  goes to infinity. In fact  $[L_A, 1]$  has a zero confidence level for any  $n$  and  $1 - \alpha$ . It is unsafe to apply the interval even for a huge sample. See a similar result for the Wald interval in Agresti (2002, p. 32). A modified interval is obtained as follows: (a) define an order  $\leq_A$  by ranking all sample points with the values of  $L_A$  (for instance,  $L_A$  achieves its maximum at (32,0), so (32,0) is the largest under  $\leq_A$ ); (b) derive the smallest 95% interval, denoted by  $[L_A^*, 1]$ , under order  $\leq_A$  following Lemma 1. Then  $[L_A^*(9, 20), 1] = [-0.00842, 1]$ , wider than the proposed  $[L_D(9, 20), 1] = [0.00610, 1]$ , is not able to conclude significance for marijuana use. □

### 3. Discussion

Comparing two dependent proportions is a common problem in categorical data analysis. When a researcher wants to determine how much more effective for a treatment than a control, an optimal one-sided interval is needed for the difference of two proportions. These are the motivation of the paper. Lemma 1 indeed converts interval construction to order construction; also for any given interval  $[L, 1]$  an order  $\leq_L$  can be defined on the sample space  $S$  by

$$x \leq_L (\equiv) x' \text{ if and only if } L(x) \leq (=) L(x').$$

Two special cases are  $\preceq_{L_\alpha}$  in (13) and  $\preceq_A$  in Example 2. Therefore, there exists a one-to-one relation between orders and intervals. Order construction is relatively simple because only the rank is to be determined. Interval construction consists of two tasks: (A) order construction; (B) obtain the smallest interval under a given order. Lemma 1 only completes (B), and in this paper we accomplish (A) by proposing an inductive order. Such an order requires a starting point and a finite sample space, and these two requirements are met here. Although a related result on two independent proportions is discussed in Wang (2010), inference on two dependent proportions itself is important enough and a solution is expected for practice. So an optimal interval, as discussed in Remark 1 and Example 2, is derived under the inductive order. In the meantime, an R-code is available from the author upon request to compute the proposed order and interval efficiently, especially when the sample size is not large. This also makes the results more easily accessible. How to obtain a reasonable order on a general sample space is challenging and is of great interest.

### Acknowledgments

This work was partially supported by NSF(US) with grant number DMS-0906858. The author is grateful to an anonymous referee and an Associate Editor for their constructive suggestions. He also thanks Dr. Guogen Shan for an R-code template that improves the computation in examples.

### Appendix. Supplementary data

Supplementary material related to this article can be found online at <http://dx.doi.org/10.1016/j.spl.2012.03.035>.

### References

- Agresti, A., 2002. *Categorical Data Analysis*, second ed. Wiley, New York.
- Buehler, R.J., 1957. Confidence intervals for the product of two binomial Parameters. *Journal of the American Statistical Association* 52, 482–493.
- Chen, J., 1993. The Order relations in the sample spaces and the confidence limits for parameters. *Advances in Mathematics* 22, 542–552.
- Karacan, I., Fernandez-Salas, A., Coggins, W.J., Carter, W.E., Williams, R.L., Thornby, J.I., Salis, P.J., Okawa, M., Villaume, J.P., 1976. Sleep electroencephalographic-celeurooculographic characteristics of chronic marijuana users: Part 1. *Annals of the New York Academy of Sciences* 282, 348–374.
- Lloyd, C.J., Kabaila, P., 2003. On the optimality and limitations of buehler bounds. *Australian & New Zealand Journal of Statistics* 45, 167–174.
- McNemar, Q., 1947. Note on the sampling error of the difference between correlated proportions or percentages. *Psychometrika* 12, 153–157.
- Wang, W., 2006. Smallest confidence intervals for one binomial proportion. *Journal of Statistical Planning and Inference* 136, 4293–4306.
- Wang, W., 2010. On construction of the smallest one-sided confidence interval for the difference of two proportions. *The Annals of Statistics* 38, 1227–1243.