

11-1-2002

# Adaptive Tests for Ordered Categorical Data

Vance W. Berger

*Biometry Research Group, National Cancer Institute*

Anastasia Ivanova

*University of North Carolina*

 Part of the [Applied Statistics Commons](#), [Social and Behavioral Sciences Commons](#), and the [Statistical Theory Commons](#)

## Recommended Citation

Berger, Vance W. and Ivanova, Anastasia (2002) "Adaptive Tests for Ordered Categorical Data," *Journal of Modern Applied Statistical Methods*: Vol. 1 : Iss. 2 , Article 36.

DOI: 10.22237/jmasm/1036108980

## Adaptive Tests for Ordered Categorical Data

Vance W. Berger  
Biometry Research Group  
National Cancer Institute

Anastasia Ivanova  
Department of Biostatistics  
University of North Carolina

---

Consider testing for independence against stochastic order in an ordered  $2 \times J$  contingency table, under product multinomial sampling. In applications one may wish to exploit prior information concerning the direction of the treatment effect, yet ultimately end up with a testing procedure with good frequentist properties. As such, a reasonable objective may be to simultaneously maximize power at a specified alternative and ensure reasonable power for all other alternatives of interest. For this objective, none of the available testing approaches are completely satisfactory. A new class of admissible adaptive tests is derived. Each test in this class strictly preserves the Type I error rate and strikes a balance between good global power and nearly optimal (envelope) power to detect a specific alternative of most interest. Prior knowledge of the direction of the treatment effect, the level of confidence in this prior information, and possibly the marginal totals might be used to select a specific test from this class.

Key words: Contingency table; exact conditional test; linear rank test; omnibus test; permutation test.

---

### Introduction

When comparing two treatments on the basis of an ordinal endpoint, the data can be summarized as a  $2 \times J$  contingency table. The objective tumor response data, e.g., from 35 ovarian cancer patients treated with cisplatin-based combination chemotherapy and salvage platinum-based therapy (Chiara et al., 1993) are (4,7,2,2) and (1,6,7,6) for patients with treatment-free intervals  $\leq 12$  months and  $> 12$  months, respectively, with categories for 'progressive disease', 'stable disease', 'partial response', and 'complete response'. Combining the two 'non-response' categories, as is common, yields counts  $C_1 = (11,2,2)$  and  $C_2 = (7,7,6)$  in the two groups. For simplicity, the case  $J = 3$  is treated, but with modification the results apply more generally. It is common in practice to dispense with the specification of the alternative hypothesis, and proceed directly to the analysis.

This failure to make the specific alternative hypothesis explicit is unfortunate, because it should serve as the basis for selecting and evaluating the analysis. Linear rank tests, based on assigning numerical scores to the categories, are the most powerful tests to detect point alternatives. If one wishes to test for the superiority of one treatment to another, then stochastic order serves as a reasonable (composite) alternative hypothesis (Cohen and Sackrowitz, 1998). Unless the margins satisfy pathological conditions, there is no uniformly most powerful test or monotone likelihood ratio. When testing for stochastic order, nonlinear rank tests, including the Smirnov, improved (Berger and Sackrowitz, 1997), convex hull (Berger, Permutt, and Ivanova, 1998; henceforth BPI), and  $COM(L)$  Fisher tests, tend to have better overall power profiles than linear rank tests do.

Berger's (1998) adaptive nonlinear rank test can be generalized to provide an entire class of exact, admissible, adaptive nonlinear rank tests, each of which balances omnibus power for any stochastically ordered alternative against optimal power to detect a specific alternative of greatest interest. The margins may be used to suggest the selection of one particular test from this novel class of tests. The exact conditional powers of some of the aforementioned tests are compared.

---

Vance W. Berger is Mathematical Statistician at the NCI and Adjunct Professor at University of Maryland Baltimore County. E-mail: [vb78c@nih.gov](mailto:vb78c@nih.gov). Anastasia Ivanova is Assistant Professor, Dept. of Biostatistics, School of Public Health., University of North Carolina – Chapel Hill. E-mail: [aivanova@bios.unc.edu](mailto:aivanova@bios.unc.edu).

Notation and Formulation

Consider product multinomial sampling, with  $n_1$  and  $n_2$  (each fixed by the design) patients treated with the control and active treatments, respectively. The vectors of cell probabilities (each summing to one) are  $\pi_1=(\pi_{11},\pi_{12},\pi_{13})$  and  $\pi_2=(\pi_{21},\pi_{22},\pi_{23})$ , respectively, and the corresponding trinomial random vectors are  $C_1=(C_{11},C_{12},C_{13})$  and  $C_2=(C_{21},C_{22},C_{23})$ , with  $n_i=C_{i1}+C_{i2}+C_{i3}$ ,  $i=1,2$ . The log odds ratios,  $\theta_1$  and  $\theta_2$ , are calculated from  $\pi_1$  and  $\pi_2$  as

$$\theta_1 = \log\{(\pi_{11}\pi_{23})/(\pi_{21}\pi_{13})\} \text{ and}$$

$$\theta_2 = \log\{(\pi_{12}\pi_{23})/(\pi_{22}\pi_{13})\}.$$

Let  $T_j = C_{1j} + C_{2j}$ ,  $j = 1,2,3$ . Conditional on  $T=(T_1,T_2,T_3)$ , the sample space  $\Gamma$  is the set of  $2 \times 3$  contingency tables with nonnegative integer cell counts, and row and column totals  $n=(n_1,n_2)$  and  $T$ , respectively. Given  $T$ ,  $n$ , and  $c=(C_{11},C_{12})$ , the entire  $2 \times 3$  contingency table can be reconstructed as  $C_{13} = n_1 - C_{11} - C_{12}$  and  $C_2 = T - C_1$ . Thus,  $c$  suffices to denote a point of  $\Gamma$ .

Figure 1. The permutation sample space  $\Gamma$  for the data set  $\{(11,2,2);(7,7,6)\}$ , with  $n=(15,20)$  and  $T=(18,9,8)$ .

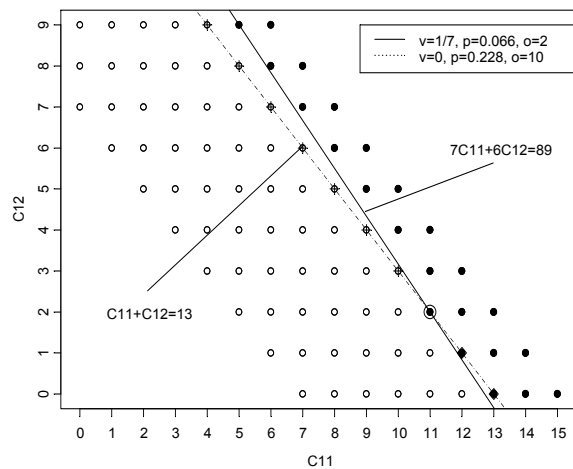


Figure 1 displays  $C_{12}$  plotted against  $C_{11}$  for all 87 tables of  $\Gamma$  for the example,  $\{(11,2,2);(7,7,6)\}$ , with observed table (11,2)

circled. With  $K(T;\theta)=1/\sum_{c \in \Gamma} H(c)\exp[\theta'c]$ ,  $\theta=(\theta_1,\theta_2)$ ,  $\pi=(\pi_1,\pi_2)$ , and  $H(c)=n_1!n_2!/\prod_{i=1}^2 \prod_{j=1}^3 C_{ij}!$ , the density follows the exponential family:

$$P_{\pi}\{c|T\} = P_{\theta}\{c|T\} = K(T;\theta)H(c)\exp[\theta'c]. \quad (2.1)$$

Let  $\Delta_1 = \pi_{11} - \pi_{21}$ , and  $\Delta_2 = (\pi_{11} + \pi_{12}) - (\pi_{21} + \pi_{22}) = \pi_{23} - \pi_{13}$ . If  $\Delta_1 \geq 0$ , and  $\Delta_2 \geq 0$ , at least one strictly, then the active treatment is objectively superior to the control. One may wish to test  $H: \pi_1 = \pi_2$  against the one-sided alternative hypothesis that the active response distribution is stochastically larger than the control response distribution,  $H_A': \Delta_1 \geq 0, \Delta_2 \geq 0, \pi_1 \neq \pi_2$ . As will be explained, this is not actually possible with a conditional test. By (2.1),  $P_{\pi}\{c|T\}$  depends on  $\pi$  only through  $\theta(\pi)$ , so if  $\theta(\pi) = \theta(\pi^*)$ , then  $c$  offers no information with which to distinguish  $\pi$  from  $\pi^*$ . To be identifiable, then, the hypotheses must be formulated in terms of  $\theta$  (Berger, 1998).

The null hypothesis  $\pi_1 = \pi_2$  is equivalent to  $H: \theta(\pi) = \mathbf{0}$ , but unless  $0 \leq \theta_2 \leq \theta_1$ ,  $\theta(\pi)$  provides insufficient information with which to determine if  $\pi$  satisfies  $H_A'$  because no conditional alternative hypothesis is equivalent to  $H_A$ . Note, e.g., that  $\{(3,3,4)/10;(2,4,4)/10\}$  satisfies  $H_A'$  and  $\{(21,51,328)/400; (7,34,164)/205\}$  does not, yet  $\theta = (\log(3/2), \log(3/4))$  for both. The conditional power to detect  $\pi$  depends on  $\theta(\pi)$  only, so no conditional test that preserves the  $\alpha$ -level whenever  $H_A$  does not hold can be globally powerful whenever it does hold.

However, if  $\pi$  satisfies  $H_A$ , then  $\theta_1(\pi) > 0$ ; and if  $\theta_1 > 0$ , then for any  $\theta_2$  there exists (Berger and Sackowitz, 1997)  $\pi$  satisfying  $H_A$  such that  $\theta(\pi) = (\theta_1, \theta_2)$ . As such,  $\theta_1$  is the key parameter; the active treatment is superior on  $\Omega_A = \{\theta | \theta_1 >$

0}, no different on  $\Omega_0 = \{\theta | \theta_1 = 0\}$ , and inferior on  $\Omega_C = \{\theta | \theta_1 < 0\}$ . It is reasonable, then, to test  $H$  against  $H_A : \theta_1 > 0$ . The large unconditional indifference region, where neither group stochastically dominates the other, has, by conditioning, been absorbed into  $\Omega_0 \cup \Omega_A \cup \Omega_C$ .

Let  $\delta(\theta) = 1 - \theta_2/\theta_1$  be the *direction* of the effect. As  $\theta_1$  increases in both  $\Delta_1$  and  $\Delta_2$ , while  $\theta_2$  ( $\theta_1 - \theta_2$ ) increases in  $\Delta_2$  ( $\Delta_1$ ), and decreases in  $\Delta_1$  ( $\Delta_2$ ), the superiority of the active treatment to the control is due primarily to a shift from the middle to the best outcome ( $\Delta_2 > \Delta_1$ ) if  $\delta(\theta)$  is small, or from the worst to the middle outcome ( $\Delta_1 > \Delta_2$ ) if  $\delta(\theta)$  is large. Let  $\Omega_v = \{\theta | \theta_1 > 0, \delta(\theta) = v\}$ . As  $\delta(\theta)$  is generally unknown *a priori*, omnibus tests that are sensitive to departures from  $H_0$  in each direction of  $\Omega_A = \cup_{v \in \mathfrak{R}^1} \Omega_v$  are preferred to tests that lack this desirable property.

If the  $\phi$  rejection region  $R_\alpha(\phi)$  contains  $D[\Gamma]$ , the set of directed extreme points of  $\Gamma$  (BPI, 1998), then  $\phi$  is omnibus. The challenge is to exploit prior information about  $\delta(\theta)$  to construct omnibus tests with especially good power in one preferred direction,  $\Omega_v$ . For reasons articulated by Berger (2000) and Berger et al. (2002), we consider only exact conditional tests in this formulation.

A New Look at Linear Rank Tests

Linear rank tests are based on numerical scores  $(v_1, v_2, v_3)$ ,  $v_1 < v_3$ , assigned to the three outcome levels. With  $v = (v_2 - v_1)/(v_3 - v_1)$ ,  $\phi_v$  uses test statistic  $z_v(\mathbf{c}) = C_{11} + (1 - v)C_{12}$ . New notation allows for greater insight into linear rank tests. Let  $M_v(\mathbf{c}) = \{\mathbf{c}^* \in \Gamma | z_v(\mathbf{c}^*) \geq z_v(\mathbf{c})\}$  be the  $\phi_v$  extreme region of  $\mathbf{c}$ , with boundary  $B_v(\mathbf{c})$  and p-value  $p_v(\mathbf{c}) = P_{\mathbf{0}}\{M_v(\mathbf{c})|T\}$ . The level set (Frick, 2000, p. 719) of  $z_v(\mathbf{c})$  is  $B_v(\mathbf{c}) \cap \Gamma$ , with  $o_v(\mathbf{c})$  its order, or the number of points of  $B_v(\mathbf{c}) \cap \Gamma$ . If  $\mathbf{c} =$

$(C_{11}, C_{12}) \in \Gamma$  and  $\mathbf{c}^* = (C_{11}^*, C_{12}^*) \in \Gamma - \mathbf{c}$ , then  $z_v(\mathbf{c}^*) = z_v(\mathbf{c})$  if and only if  $v = 1 - (C_{11} - C_{11}^*)/(C_{12}^* - C_{12})$ , say  $v = v_{\mathbf{c}, \mathbf{c}^*}$  (vector valued for  $J > 3$ ). Let  $V(\mathbf{c}) = \{v_1(\mathbf{c}), v_2(\mathbf{c}), \dots, v_{K_c}(\mathbf{c})\}$  be the ordered set  $\{v_{\mathbf{c}, \mathbf{c}^*} | |v_{\mathbf{c}, \mathbf{c}^*}| < \infty, \mathbf{c}^* \in \Gamma - \mathbf{c}\}$ , and let  $v_0(\mathbf{c}) = -\infty$  and  $v_{K_c+1}(\mathbf{c}) = \infty$ . For finite  $v$ ,  $o_v(\mathbf{c}) > 1$  if and only if  $v \in V(\mathbf{c})$ .

Let  $\varepsilon(\mathbf{c}) = \min_k [v_{k+1}(\mathbf{c}) - v_k(\mathbf{c})]/2$ ,  $z_v^\perp(\mathbf{c}) = C_{12} + (v - 1)C_{11}$ ,  $B_v^+(\mathbf{c}) = \{\mathbf{c}^* \in B_v(\mathbf{c}) \cap \Gamma | z_v^\perp(\mathbf{c}^*) > z_v^\perp(\mathbf{c})\}$ ,  $B_v^-(\mathbf{c}) = \{\mathbf{c}^* \in B_v(\mathbf{c}) \cap \Gamma | z_v^\perp(\mathbf{c}^*) < z_v^\perp(\mathbf{c})\}$ ,  $v^*(\mathbf{c}) = \{v^* | p_{v^*}(\mathbf{c}) \leq p_v(\mathbf{c}) \text{ for all } v^*\}$ .

By Lemma 1 (in the Appendix),  $v^*(\mathbf{c})$  consists of the scores that minimize not just  $p_v(\mathbf{c})$  but also  $p_{\min(v)}(\mathbf{c}) = \min(\lim_{u \downarrow v} p_u(\mathbf{c}), \lim_{u \uparrow v} p_u(\mathbf{c})) = p_v(\mathbf{c}) - \max(P_0\{B_v^-(\mathbf{c})\}, P_0\{B_v^+(\mathbf{c})\})$ . Hence,  $p_{\min(v)}(\mathbf{c})$ , which also equals  $\min\{p_{v-\varepsilon}(\mathbf{c}), p_{v+\varepsilon}(\mathbf{c})\}$ , is a true p-value. As  $\Gamma$  has finitely many subsets, there can be only a finite number of values for  $p_v(\mathbf{c})$ , so the minimum p-value is attained, and  $v^*(\mathbf{c}) \neq \emptyset$ . If  $v \in V(\mathbf{c})$ , then  $o_v(\mathbf{c}) > 1$ ,  $B_v^-(\mathbf{c}) \cup B_v^+(\mathbf{c}) \neq \emptyset$ ,  $p_{\min(v)}(\mathbf{c}) < p_v(\mathbf{c})$ , and  $v \notin v^*(\mathbf{c})$ . Hence,  $v^*(\mathbf{c}) \cap V(\mathbf{c}) = \emptyset$ , and, by Lemma 1,  $v^*(\mathbf{c})$  consists of one or more open intervals of the form  $(v_k(\mathbf{c}), v_{k+1}(\mathbf{c}))$ . For  $\{(11,2,2);(7,7,6)\}$ ,  $\mathbf{c} = (11,2)$ ,  $K_c = 42$ ,  $\varepsilon(11,2) = 1/84$ , and  $V(\mathbf{c}) = \{-6, -5, -4, -3, -5/2, -2, -5/3, -3/2, -4/3, -5/4, -6/5, -1, -5/6, -4/5, -3/4, -2/3, -3/5, -4/7, -1/2, -3/7, -2/5, -1/3, -2/7, -1/4, -1/5, -1/6, -1/7, 0, 1/7, 1/6, 1/5, 1/4, 2/7, 1/3, 2/5, 1/2, 2/3, 1, 3/2, 2, 5/2, 3, 4, 5, 6\}$ .

Figure 1 shows  $M_{1/7}(11,2)$  by dark dots and  $M_0(11,2) - M_{1/7}(11,2)$  by crosses. Because  $(11,2)$  minimizes  $z_{1/7}^\perp(11,2) = 7C_{12} - 6C_{11}$  over  $B_{1/7}(11,2) \cap \Gamma$  (Table 1),  $B_{1/7}^-(11,2) = \emptyset$  and  $p_{1/7}(11,2) = \lim_{u \uparrow 1/7} p_u(11,2) = 0.066$ . Also  $p_v(11,2) = 0.020$  for  $v \in (1.0, 1.5) = v^*(11,2)$ . If  $v$

$\in V(11,2)$ , then  $P_0\{B_v^-\} \leq P_0\{B_v^+\}$  for  $v > 1.5$ , and  $P_0\{B_v^-\} \geq P_0\{B_v^+\}$  for  $v < 1.0$ . The optimality of most powerful (MP) test  $\phi_{\delta}(\theta)$  to detect  $I\theta$ , for  $I > 0$  (BPI, 1998), is offset by its

potentially poor power on  $\Omega_A - \Omega_{\delta}(\theta)$ . In fact,  $D[\Gamma]$  may not be contained in the  $\phi_v$  critical region  $R_{\alpha}(\phi_v)$  for any  $v$ , so for

Table 1. All possible linear rank tests with scores  $(0,v,1)$ , with middle score  $v \in [0,2]$ , for the data set  $\{(11,2,2);(7,7,6)\}$ , along with the number of points in its level set, the endpoints and null probabilities of each segment of its level set, and various p-values. (null probabilities of various extreme regions).

$v$	$\alpha_v(11,2)$	Endpoints of: $B_v^+$ $B_v^-$	$p_v$ (minimum is underlined)	$p_v^-$	$p_v^+$	$P_0\{B_v^+\}$	$P_0\{B_v^-\}$	$p_{v,\infty}$	$M_v - M_{v,\infty}$
$v \in (-1/7,0)$	1		0.2262	0.2262	0.2262			0.2262	
$v = 0$	10	(4,9)   (12,1) -(10,3)   -(13,0)	0.2277	0.2262	<u>0.0661</u>	0.1615	0.0015	0.0726	(7,6)- (10,3)
$v \in (0,1/7)$	1		0.0661	0.0661	0.0661			0.0661	
$v = 1/7$	2	(5,9)	0.0661	0.0661	<u>0.0661</u>	$2.1 \cdot 10^{-5}$		0.0661	
$v \in (1/7,1/6)$	1		0.0661	0.0661	0.0661			0.0661	
$v = 1/6$	2	(6,8)	0.0661	0.0661	<u>0.0657</u>	0.0004		0.0661	
$v \in (1/6,1/5)$	1		0.0657	0.0657	0.0657			0.0657	
$v = 1/5$	2	(7,7)	0.0657	0.0657	<u>0.0629</u>	0.0028		0.0657	
$v \in (1/5,1/4)$	1		0.0629	0.0629	0.0629			0.0629	
$v = 1/4$	2	(8,6)	0.0629	0.0629	<u>0.0538</u>	0.0091		0.0629	
$v \in (1/4,2/7)$	1		0.0538	0.0538	0.0538			0.0538	
$v = 2/7$	2	(6,9)	0.0538	0.0538	<u>0.0538</u>	$5.7 \cdot 10^{-6}$		0.0538	
$v \in (2/7,1/3)$	1		0.0538	0.0538	0.0538			0.0538	
$v = 1/3$	3	(7,8) -(9,5)	0.0538	0.0538	<u>0.0387</u>	0.0152		0.0387	(9,5)
$v \in (1/3,2/5)$	1		0.0387	0.0387	0.0387			0.0387	
$v = 2/5$	2	(8,7)	0.0387	0.0387	<u>0.0382</u>	0.0005		0.0387	
$v \in (2/5,1/2)$	1		0.0382	0.0382	0.0382			0.0382	
$v = 1/2$	4	(9,6)   (12,0) -(10,4)	0.0385	0.0382	<u>0.0237</u>	0.0148	0.0003	0.0249	(10,4)
$v \in (1/2,2/3)$	1		0.0237	0.0237	0.0237			0.0237	
$v = 2/3$	2	(10,5)	0.0237	0.0237	<u>0.0220</u>	0.0017		0.0237	
$v \in (2/3,1)$	1		0.0220	0.0220	0.0220			0.0220	
$v = 1$	5	(11,4)   (11,1) -(11,3)   -(11,0)	0.0276	0.0220	<b><u>0.0198</u></b>	0.0078	0.0056	0.0276	
$v \in (1,3/2)$	1		<b><u>0.0198</u></b>	<b><u>0.0198</u></b>	<b><u>0.0198</u></b>			0.0198	
$v = 3/2$	2	(10,0)	0.0205	<b><u>0.0198</u></b>	0.0205		0.0008	0.0205	
$v \in (3/2,2)$	1		0.0205	0.0205	0.0205			0.0205	
$v = 2$	4	(12,3)   (10,1) -(9,0)	0.0294	<u>0.0205</u>	0.0289	0.0005	0.0089	0.0294	
$v \in (2,5/2)$	1		0.0289	0.0289	0.0289			0.0289	

Note that all the values are calculated at the outcome  $(11,2)$ ;  $p_{v,\infty}$  and  $M_{v,\infty}$  are the p-value and extreme region, respectively, of the adaptive test based on  $v$  and  $\tau = \infty$ .

each  $v$  there will exist  $\theta \in \Omega_A$  for which the power of  $\varphi_v$  to detect  $l\theta$  tends to zero as  $l$  gets large (BPI, 1998). Podgor, Gastwirth, and Mehta (1996) proposed the maximin efficiency robust test (MERT) in hopes of providing better power than linear rank tests. Ironically, the MERT is itself a linear rank test; its rejection region may also fail to contain  $D[\Gamma]$ , leading to poor power on parts of  $\Omega_A$  and no power in the limit in some directions. Berger and Ivanova (2002) showed that at certain  $\alpha$ -levels the most stringent linear rank test is  $\varphi_{v_S}$ , where  $v_S$  is such that the two points of  $D[\Gamma]$  that are furthest (in Euclidean distance) from each other are equated by  $z_{v_S}(c)$ . For  $\{(11,2,2),(7,7,6)\}$ , this gives  $v_S = 0$ , because  $\Gamma$  has two directed extreme points,  $D[\Gamma]=\{(15,0);(6,9)\}$ , and  $z_0(15,0) = 15+(1-0)(0)=15=6+(1-0)(9)=z_0(6,9)$ .

#### Nonlinear Rank Tests

By allowing the boundary of  $R_\alpha(\varphi)$  to curve, nonlinear rank tests often require smaller  $\alpha$ -levels to ensure that  $D[\Gamma] \subset R_\alpha(\varphi)$  than linear rank tests would. However, this is not always the case. Berger and Ivanova (2002) provide an example in which the proportional odds and proportional hazards tests (McCullagh, 1980) are not nonlinear enough to be omnibus at reasonable  $\alpha$ -levels. The Smirnov test,  $\varphi_S$ , uses as the test statistic the largest of three quantities, 0,  $D_1 = C_{11}/n_1 - C_{21}/n_2$ , and  $D_2 = (C_{11} + C_{12})/n_1 - (C_{21} + C_{22})/n_2$ . Among tests routinely available in standard statistical software packages ( $\varphi_S$  is a standard feature of StatXact),  $\varphi_S$  minimizes the  $\alpha$ -level required for its rejection region to contain  $D[\Gamma]$ . However,  $\varphi_S$  is not generally admissible (Berger, 1998).

Permutt and Berger (2000) and Ivanova and Berger (2001) each proposed refinements of  $\varphi_S$  that break its ties. Although such refinements are necessarily uniformly more powerful than  $\varphi_S$  (Rohmel and Mansmann, 1999, p. 158), the term

“improvement of  $\varphi$ ” is reserved for a test whose exact (possibly randomized) version is uniformly more powerful than the exact (possibly randomized) version of  $\varphi$ . By this definition, refinements are rarely improvements. Berger and Sackrowitz (1997) developed methodology for constructing improvements of a given inadmissible test. In fact, by improving the “ignore-the-data” test,  $\varphi_{ITD}(c) = \alpha$  for all  $c \in \Gamma$ , Berger and Sackrowitz (1997) constructed the first known test for this problem that is simultaneously admissible and unbiased. However, rejection regions at different  $\alpha$ -levels need not be nested, so these improved tests may not yield unambiguous p-values, and thus are of somewhat limited value.

Berger (1998) established the one-to-one correspondence between the class of convex hull type tests and the minimal complete class of admissible tests. The convex hull test (BPI, 1998),  $\varphi_{CH}$ , is the simplest member of this convex hull class, and is qualitatively similar to the improvements of both  $\varphi_S$  and  $\varphi_{ITD}$ , while minimizing, among all families of tests, the  $\alpha$ -level required for its rejection region to contain  $D[\Gamma]$ .

In addition,  $\varphi_{CH}$  is based on a test statistic, so rejection regions at different  $\alpha$ -levels are nested, and p-values are provided. As such,  $\varphi_{CH}$  is about as good a test as there is for testing  $H$  against  $H_A$ , which is about as close as one can get to testing  $H$  against  $H'_A$  when dealing with  $\theta$  instead of  $\pi$ . Specifically, admissible (unbiased) tests of  $H$  against  $H_A$  are conditionally admissible (unbiased) as tests of  $H$  against  $H'_A$  (Berger and Sackrowitz, 1997). However,  $\theta(\pi)$  is a nonlinear function, and maps small corners of  $\pi$ -space (neighborhoods of structural zeros) into large regions of  $\theta$ -space. By giving each direction  $\delta(\theta)$  equal consideration,  $\varphi_{CH}$  accommodates these small corners as much as it does the large regions of  $\pi$ -space that are of greatest unconditional interest. As such,  $\varphi_{CH}$  may not be ideal when viewed unconditionally. Cohen and Sackrowitz (1998) proposed another member of the convex hull class, called the *COM(L)* Fisher test, or

$\varphi_{COM(L)}$ , based on repeatedly adding to the critical region those directed extreme points of the current acceptance region that are least likely under  $H_0$ . Because the test statistics of  $\varphi_{COM(L)}$  and  $\varphi_{CH}$  are defined not algebraically but relationally, by the relative position of  $c$  within  $\Gamma$ , the rejection regions need to be constructed recursively. This feature is a barrier to their use.

Adaptive Tests

Gross (1981, Section 5) suggested that an "analysis based on ... data-dependent scores may yield procedures that compare favorably to fixed-score procedures ...". Distinct from another definition used, e.g., by Rukhin and Mak (1992), Hogg (1974, p. 917) and Edgington (1995, pp. 371-373) defined adaptive tests as tests with data-based test statistics. This allows  $\Gamma$  to be partitioned into regions sharing a common test statistic. Because the region need not be even nearly ancillary, conditioning on the region (as suggested by Donegani, 1991, and Good, 1994, p. 122) may entail a loss of power. Comparing the value of the test statistics across regions avoids this loss of power. The intuitive objection to "comparing apples to oranges" notwithstanding, such an approach is "good" or "bad" only to the extent to which it produces a "good" or "bad" test. This approach results in tests with excellent power properties. In fact, Gastwirth (1985) stated that "when the MERT for a particular problem has a low  $r^2$ , adaptive procedures are needed".

Without knowing  $\theta$  *a priori*, it is unclear where to maximize the power. One could estimate  $\delta(\theta)$  from  $c$ , say as  $\hat{\delta}_p(c)$ , perhaps using maximum likelihood, and use the MP test  $\varphi_{\hat{\delta}_p}$ . The p-value of  $\varphi_{\hat{\delta}_p}$  evaluated at observed outcome  $c$ ,  $p_{\hat{\delta}_p}(c)$ , is stochastically too small to serve as a valid p-value, but  $p_{\hat{\delta}_p}(c)$  can be used as a *test statistic*, to be compared to its null distribution (Rohmel and Mansmann, 1999, p. 165). Variation in  $c$  is reflected in  $p_{\hat{\delta}_p}(c)$  through *both* the argument and the subscript. Using either  $p_{\hat{\delta}_p}(c)$  or  $z_{\hat{\delta}_p}(c)$ , suitably normalized, as a test statistic, any

estimator  $\hat{\delta}_p(c)$  of  $\delta(\theta)$  induces an adaptive test, with regions  $\Gamma_v = \hat{\delta}^{-1}(v) = \{c \in \Gamma \mid \hat{\delta}_p(c) = v\}$ . If the regions are  $\Gamma_0 = \{c \in \Gamma \mid C_{12} > n_1 T_2 / (n_1 + n_2)\}$ ,  $\Gamma_1 = \Gamma - \Gamma_0$ , and  $\Gamma_v = \emptyset$  for  $v \notin \{0,1\}$ , and the  $\varphi_v$  test statistic  $z_v(c)$  is used on  $\Gamma_v$ , with  $C_{11} + C_{12}$  ( $v = 0$ ) and  $C_{11}$  ( $v = 1$ ) normalized to  $D_2$  and  $D_1$ , respectively, to facilitate the comparison of points from  $\Gamma_1$  ( $D_1 > D_2$ ) to those from  $\Gamma_0$  ( $D_2 \geq D_1$ ), then  $\varphi_S$  results. Similar binary adaptive tests might define  $\Gamma_0$  and  $\Gamma_1$  by whichever of  $\varphi_0$  and  $\varphi_1$  yields a smaller p-value or a larger  $\chi^2$ .

Berger (1998) proposed judging outcome  $c$  by how small a p-value it can yield with an MP test; that is,  $\varphi_A$  uses  $p_{v^*}(c) = \min_{-\infty \leq v \leq \infty} p_v(c)$  as the test statistic. This is a continuous version of the adaptive test based on  $\min(p_0(c), p_1(c))$ , and estimates  $\delta(\theta)$  non-uniquely as  $\hat{\delta}_c = v$  for any value  $v \in v^*(c)$ . The induced regions are  $\Gamma_v = \{c \in \Gamma \mid v \in v^*(c)\}$ . The  $\varphi_A$  critical region is  $R_\alpha(\varphi_A) = \cup_{v \in R^1} R_{\alpha^*}(v)(\varphi_v)$  for some set of  $\alpha^*(v) < \alpha$ , so  $\varphi_A$  is intuitively similar to union-intersection tests (Roy, 1953; Marden, 1991). Despite being constructed non-recursively,  $\varphi_A$  is a convex hull type test (Berger, 1998); hence  $\varphi_A$  is always admissible. Also,  $\varphi_A$  tends to be omnibus, as  $D[\Gamma] \subset R_\alpha(\varphi_A)$  for reasonable  $\alpha$ -levels.

Accommodating a Favored Alternative

Suppose that one believes *a priori* that  $\delta(\theta) = \delta_p$ . Let  $\tau \geq 0$  be a measure of the strength in the belief that  $\delta(\theta) = \delta_p$ . The dual objectives are ensuring nearly MP power on  $\Omega_{\delta_p}$  and reasonable power on  $\Omega_A - \Omega_{\delta_p}$ , with relative importance dictated by  $\tau$ . One might use  $\varphi_{\delta_p}$  (which is MP on  $\Omega_{\delta_p}$ ) for large  $\tau$ , or  $\varphi_A$  (which is a good omnibus test) for small  $\tau$ , but none of the aforementioned test suffices for intermediate values of  $\tau$ . Linear

combinations such as  $(\tau \varphi_{\delta_p} + \varphi_A)/(\tau + 1)$  would not suffice either, because they have large randomization regions and small critical regions, consisting only of the intersection  $R_\alpha(\varphi_{\delta_p}) \cap R_\alpha(\varphi_A)$ . Of course, these inadmissible tests could be improved to admissibility, but then the procedure would be complicated, and p-values may not be defined. There is another approach to bridge the gap between  $\varphi_{\delta_p}$  and  $\varphi_A$ . Specifically, start with  $\varphi_A$ , but penalize those  $\mathbf{c}$  whose minimizing MP p-value is obtained by  $\nu$  far from  $\delta_p$ . To this end, let  $\varphi_{\delta_p, \tau, \alpha}$  (or  $\varphi_{\delta_p, \tau}$ ) be the level- $\alpha$  adaptive test based on the test statistic

$$A(\delta_p, \tau, \mathbf{c}) = \min_{-\infty \leq \nu \leq \infty} [\rho_{\min(\nu)}(\mathbf{c})(1 + |\delta_p - \nu|)^\tau].$$

Let  $\nu_{[\delta_p, \tau]}(\mathbf{c}) = \{\nu \mid p_{\min(\nu)}(\mathbf{c})(1 + |\delta_p - \nu|)^\tau = A(\delta_p, \tau, \mathbf{c})\}$ . Clearly,  $\varphi_{\delta_p, 0} = \varphi_A$  for any  $\delta_p$  and  $p_{\min(\nu)}(\mathbf{c})(1 + |\delta_p - \nu|)^\tau \leq 1$  if  $\nu \in \nu_{[\delta_p, \tau]}(\mathbf{c})$ . Lemmas 2-4 confine  $\nu_{[\delta_p, \tau]}(\mathbf{c})$  to a finite subset of an interval that shrinks, as  $\tau$  gets large, to  $\{\delta_p\}$ . By Lemma 4,  $\varphi_{\delta_p, \infty}$  induces the same ordering on  $\Gamma$  as  $\varphi_{\delta_p}$  does, thereby optimizing power on  $\Omega_{\delta_p}$ . Yet because the  $\varphi_{\delta_p, \infty}$  test statistic is  $p_{\min(\delta_p)}(\mathbf{c})$ , and not necessarily  $p_{\delta_p}(\mathbf{c})$ ,  $\varphi_{\delta_p, \infty}$  is a refinement of  $\varphi_{\delta_p}$ , and  $p_{\min(\nu)}(\mathbf{c}) \leq p_{\nu, \infty}(\mathbf{c}) \leq p_\nu(\mathbf{c})$  for all  $\nu$  and  $\mathbf{c}$ . From Table 1, e.g.,  $p_{0.5}(11,2) = 0.0385$ , but  $p_{0.5, \infty}(11,2) = 0.0385 - P_0\{(10,4)|T\} = 0.0249$ . Each test in the class of adaptive tests is admissible.

*Theorem 1.* For any triple  $\delta_p \in \mathfrak{R}^1$ ,  $\tau \geq 0$ , and  $\alpha \in [0,1]$ ,  $\varphi_{\delta_p, \tau, \alpha}$  is admissible. Graubard and Korn (1987) suggested that without a reason to use a different  $\delta_p$ ,  $\varphi_{0.5}$  should be used. The desire to focus power on the "central" direction,  $\Omega_{0.5}$ , is understandable, but the use of linear rank tests in general (BPI, 1998; Berger and Ivanova, 2002), and  $\varphi_{0.5}$  in particular (Ivanova and Berger, 2001), have been criticized. Now  $\varphi_{0.5, \tau}$  offers good central power without sacrificing global power

(unless  $\tau = \infty$ ). But even if  $\tau = \infty$ ,  $\varphi_{0.5, \infty}$  is still more powerful than, and hence preferable to  $\varphi_{0.5}$ .

#### Margin-Based Selection of $\delta_p$ and $\tau$

Recall that  $\nu_S$  can be determined from the margins ( $\mathbf{n}$  and  $\mathbf{T}$ , summarized by  $\Gamma$ ). In some cases, it may be reasonable to use  $\nu_S$  as  $\delta_p$ . In others, it may be reasonable to use the margins to find the largest  $\tau$  that allows  $R_\alpha(\varphi_{\delta_p, \tau, \alpha})$  to contain  $D[\Gamma]$ . Unless  $|\delta_p - \nu_S|$  is small, the larger  $\tau$  is, the less  $\varphi_{\delta_p, \tau}$  focuses on omnibus power. Hence, the  $\alpha$ -level required for  $R_\alpha(\varphi_{\delta_p, \tau, \alpha})$  to contain  $D[\Gamma]$  tends to increase in  $\tau$ . If a range of  $\alpha$ -levels would be considered, say  $0.01 \leq \alpha \leq 0.1$ , then use the smallest  $\alpha$ -level in selecting  $\tau$ . Restricting attention to the integer values of  $\tau$ , and using  $\delta_p = 0.5$ , note that for  $\{(11,2,2), (7,7,6)\}$ ,  $D[\Gamma] = \{(6,9); (15,0)\}$  is contained by  $R_{0.01}(\varphi_{0.5, 18})$ ,  $R_{0.025}(\varphi_{0.5, 20})$ ,  $R_{0.05}(\varphi_{0.5, 22})$ , and  $R_{0.1}(\varphi_{0.5, 24})$ ; but none of  $R_{0.01}(\varphi_{0.5, 19})$ ,  $R_{0.025}(\varphi_{0.5, 21})$ ,  $R_{0.05}(\varphi_{0.5, 23})$ , or  $R_{0.1}(\varphi_{0.5, 25})$  contain  $(6,9)$ . Consequently,  $\varphi_{0.5, 18}$  would be used by this approach.

#### Comparisons of Tests

The exact conditional power of the one-sided nonrandomized versions of  $\varphi_{0.0}$ ,  $\varphi_{0.5}$ ,  $\varphi_{1.0}$ ,  $\varphi_S$ ,  $\varphi_{CH}$ ,  $\varphi_{COM(L)}$ , and some adaptive tests, at  $\alpha \leq 0.05$ , are compared considering all 87  $2 \times 3$  tables with row and column margins as in the example,  $T = (18,9,8)$ ,  $\mathbf{n} = (15,20)$ . Figure 2 illustrates extreme regions. The exact conditional power of  $\varphi$  to detect  $\theta$  is calculated as  $P_\theta\{R_{0.05}(\varphi)|T\}$ . Here  $4 \times 7 = 28$  alternatives, with  $\theta_1 \in \{0.5, 1.0, 1.5, 2.0\}$  and  $\theta_2 = \{-1.5, -1.0, -0.5, 0.0, 0.5, 1.0, 1.5\}$ , are considered, along with the null case,  $\theta_1 = \theta_2 = 0$ . Bold entries represent the best power, for given  $\theta$ , among the six targeted tests in columns 4-9 and among five omnibus tests in columns 10-13. Because the linear rank tests  $\varphi_{0.0}$  ( $\alpha = 0.005$ ),  $\varphi_{0.5}$  ( $\alpha = 0.038$ ), and  $\varphi_{1.0}$  ( $\alpha = 0.028$ ) are excessively conservative, per the top



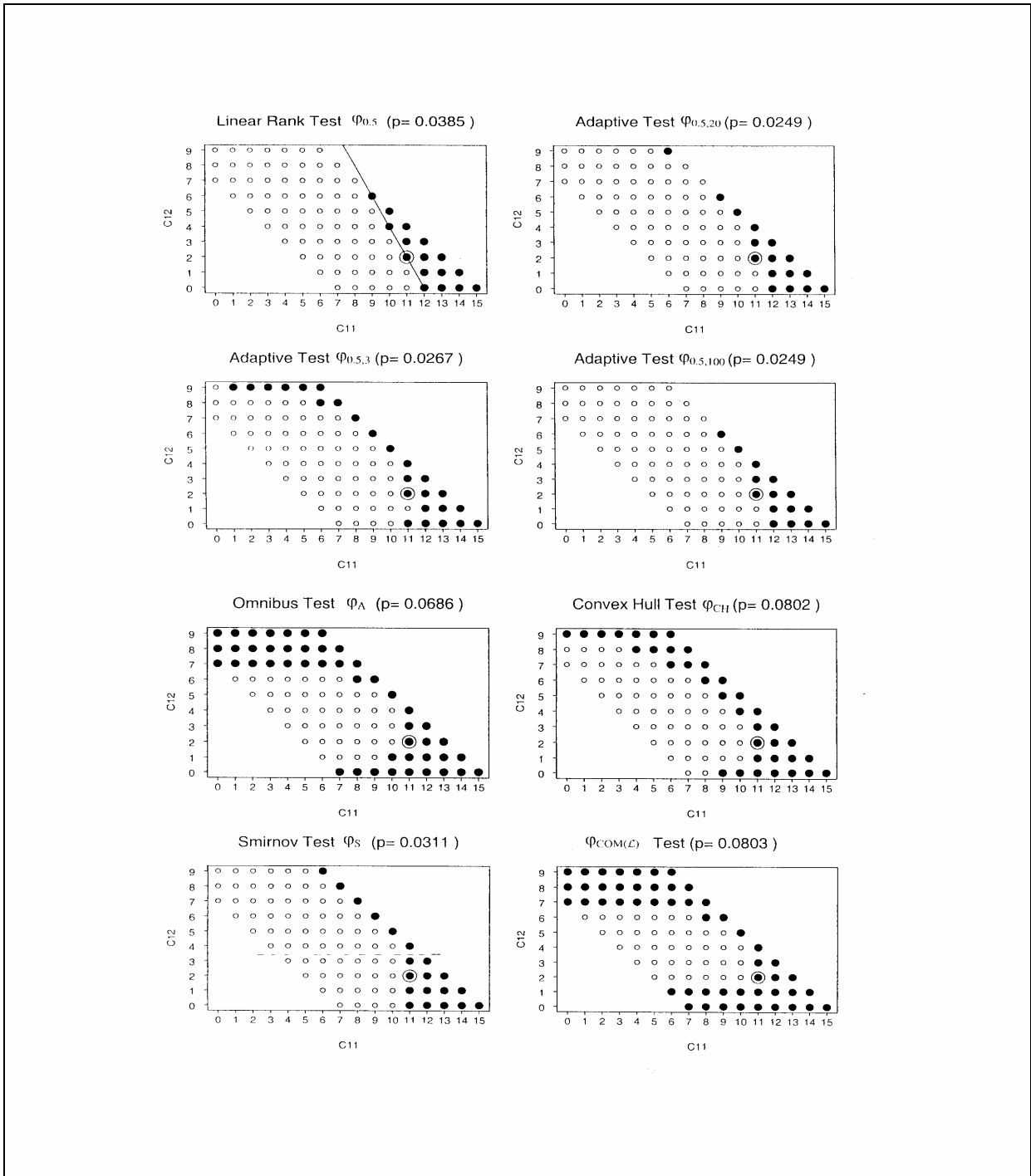


Figure 2. Extreme regions and p-values for  $\{(11,2,2);(7,7,6)\}$  and several tests including the linear rank test with equally-spaced scores  $\varphi_{0,5}$ , the adaptive tests with similar direction but varying second parameter  $\varphi_{0.5,3}$ ,  $\varphi_{0.5,20}$ ,  $\varphi_{0.5,100}$ , the omnibus adaptive test  $\varphi_A$ , the Smirnov test  $\varphi_S$ , the convex hull test  $\varphi_{CH}$ , and the  $\varphi_{COM(L)}$  test.

Table 2. Exact conditional power of the conservative (nonrandomized) versions of linear rank tests ( $\varphi_0, \varphi_1, \varphi_{0.5}$ ), adaptive tests ( $\varphi_{0,100}, \varphi_{1,100}, \varphi_{0.5,100}, \varphi_{0.5,1}$ ), omnibus adaptive test  $\varphi_A$ , the  $\varphi_{COM(L)}$  test, Smirnov test  $\varphi_S$ , and convex hull test  $\varphi_{CH}$ , with  $\alpha \leq 0.05$ , and table margins  $T=(18,9,8)$ ,  $n=(15,20)$ . Bold entries represent the best power among the tests in each block (narrow and omnibus) for each given  $\theta$ .

$\delta(\theta)$	$\theta$	$\varphi_0$	$\varphi_{0,100}$	$\varphi_{0.5}$	$\varphi_{0.5,100}$	$\varphi_1$	$\varphi_{1,100}$	$\varphi_{0.5,1}$	$\varphi_A$	$\varphi_{COM(L)}$	$\varphi_S$	$\varphi_{CH}$
	0.0 0.0	0.005	0.040	0.038	0.044	0.028	0.039	0.046	0.047	0.050	0.031	0.035
-2.000	0.5 1.5	0.054	<b>0.232</b>	0.046	0.063	0.006	0.015	0.258	<b>0.375</b>	0.316	0.058	0.255
-1.000	0.5 1.0	0.038	<b>0.163</b>	0.071	0.080	0.021	0.032	0.150	<b>0.198</b>	0.152	0.053	0.145
-0.500	1.0 1.5	0.107	<b>0.325</b>	0.151	0.174	0.039	0.067	0.290	<b>0.332</b>	0.244	0.131	0.285
0.000	0.5 0.5	0.025	<b>0.120</b>	0.103	0.110	0.057	0.070	<b>0.109</b>	0.108	0.090	0.073	0.093
0.000	1.0 1.0	0.079	<b>0.264</b>	0.212	0.223	0.099	0.126	<b>0.219</b>	0.215	0.169	0.151	0.200
0.000	1.5 1.5	0.184	<b>0.447</b>	0.352	0.371	0.149	0.208	<b>0.366</b>	0.361	0.292	0.270	0.349
0.250	2.0 1.5	0.280	0.606	0.603	<b>0.615</b>	0.370	0.445	<b>0.524</b>	0.491	0.460	0.485	0.489
0.333	1.5 1.0	0.143	0.417	0.442	<b>0.455</b>	0.288	0.328	<b>0.379</b>	0.333	0.310	0.346	0.330
0.500	1.0 0.5	0.055	0.231	0.274	<b>0.291</b>	0.200	0.225	<b>0.244</b>	0.196	0.189	0.223	0.193
0.500	2.0 1.0	0.231	0.593	0.689	<b>0.704</b>	0.560	0.597	<b>0.615</b>	0.543	0.537	0.605	0.542
0.667	1.5 0.5	0.109	0.390	0.521	<b>0.550</b>	0.454	0.481	<b>0.483</b>	0.391	0.395	0.475	0.390
0.750	2.0 0.5	0.188	0.560	0.754	<b>0.785</b>	0.723	0.741	<b>0.738</b>	0.634	0.640	0.735	0.634
1.000	0.5 0.0	0.015	0.096	0.137	<b>0.157</b>	0.121	0.147	<b>0.140</b>	0.104	0.116	0.127	0.100
1.000	1.0 0.0	0.038	0.201	0.333	<b>0.378</b>	0.332	0.368	<b>0.347</b>	0.258	0.283	0.339	0.257
1.000	1.5 0.0	0.082	0.349	0.585	<b>0.646</b>	0.612	0.642	<b>0.621</b>	0.499	0.521	0.617	0.499
1.000	2.0 0.0	0.153	0.514	0.799	0.851	0.836	<b>0.852</b>	<b>0.841</b>	0.736	0.748	0.839	0.736
1.250	2.0 -0.5	0.126	0.467	0.830	0.896	0.906	<b>0.924</b>	<b>0.908</b>	0.828	0.844	0.906	0.828
1.333	1.5 -0.5	0.062	0.302	0.634	0.729	0.736	<b>0.779</b>	<b>0.744</b>	0.628	0.665	0.737	0.628
1.500	1.0 -0.5	0.026	0.167	0.384	0.472	0.471	<b>0.536</b>	<b>0.483</b>	0.377	0.432	0.472	0.375
1.500	2.0 -1.0	0.106	0.429	0.854	0.920	0.944	<b>0.963</b>	<b>0.948</b>	0.899	0.915	0.944	0.899
1.667	1.5 -1.0	0.048	0.262	0.671	0.784	0.822	<b>0.876</b>	<b>0.834</b>	0.752	0.795	0.823	0.750
1.750	2.0 -1.5	0.093	0.401	0.874	0.930	0.965	<b>0.983</b>	<b>0.970</b>	0.945	0.958	0.965	0.945
2.000	0.5 -0.5	0.010	0.077	0.171	0.221	0.212	<b>0.273</b>	<b>0.227</b>	0.167	0.217	0.214	0.163
2.000	1.0 -1.0	0.018	0.136	0.426	0.552	0.593	<b>0.692</b>	<b>0.617</b>	0.524	0.604	0.593	0.520
2.000	1.5 -1.5	0.038	0.231	0.703	0.811	0.877	<b>0.934</b>	<b>0.895</b>	0.848	0.889	0.877	0.845
2.500	1.0 -1.5	0.013	0.111	0.463	0.602	0.687	<b>0.810</b>	0.730	0.668	<b>0.756</b>	0.687	0.660
3.000	0.5 -1.0	0.006	0.059	0.203	0.292	0.318	<b>0.433</b>	0.350	0.284	<b>0.377</b>	0.318	0.275
4.000	0.5 -1.5	0.004	0.044	0.231	0.348	0.419	<b>0.591</b>	0.487	0.435	<b>0.558</b>	0.419	0.416
Mean power		0.083	0.293	0.447	0.500	0.458	<b>0.505</b>	<b>0.519</b>	0.469	0.481	0.482	0.457
p-value for (11,2,2;7,7,6)		0.228	0.073	0.038	0.025	0.028	0.028	0.037	0.069	0.080	0.031	0.080

Table 3. Pairwise comparisons of 11 tests for  $4 \times 7 = 28$  values of  $\theta$ , where each entry is the number of parameter values (out of 28 considered in the power calculations) for which the test to the left (defining the row) had greater power than the test above (defining the column).

	$\varphi_0$	$\varphi_{0,100}$	$\varphi_{0.5}$	$\varphi_{0.5,100}$	$\varphi_1$	$\varphi_{1,100}$	$\varphi_{0.5,1}$	$\varphi_A$	$\varphi_{COM(L)}$	$\varphi_S$	$\varphi_{CH}$	Total
$\varphi_0$	-	<b>0</b>	1	<b>0</b>	4	3	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	8
$\varphi_{0,100}$	<b>28</b>	-	7	6	10	9	7	7	9	9	9	101
$\varphi_{0.5}$	27	21	-	<b>0</b>	14	12	6	14	13	12	17	136
$\varphi_{0.5,100}$	<b>28</b>	<b>22</b>	<b>28</b>	-	18	15	13	21	18	17	21	201
$\varphi_1$	24	18	14	10	-	<b>0</b>	<b>0</b>	19	14	<b>0</b>	20	119
$\varphi_{1,100}$	25	19	16	13	<b>28</b>	-	17	20	21	19	20	198
$\varphi_{0.5,1}$	<b>28</b>	21	22	15	<b>28</b>	11	-	25	23	<b>28</b>	<b>28</b>	<b>229</b>
$\varphi_A$	<b>28</b>	21	14	7	9	8	3	-	10	8	<b>28</b>	136
$\varphi_{COM(L)}$	<b>28</b>	19	15	10	14	7	5	18	-	12	20	148
$\varphi_S$	<b>28</b>	19	16	11	<b>28</b>	9	<b>0</b>	20	16	-	21	168
$\varphi_{CH}$	<b>28</b>	19	11	7	8	8	<b>0</b>	<b>0</b>	8	7	-	96
Total	272	179	144	79	161	82	<b>51</b>	144	132	112	184	

row of Table 2, they are dominated at  $\alpha = 0.05$  by their corresponding adaptive tests  $\varphi_{0,0,100}$  ( $\alpha = 0.040$ ),  $\varphi_{0.5,100}$  ( $\alpha = 0.044$ ), and  $\varphi_{1,0,100}$  ( $\alpha = 0.039$ ). This is not surprising, and will be the case quite generally. Note that  $\varphi_{0.5,1}$  maximizes the average power, at 0.519, or the area under the power curve. The non-adaptive tests did not fare as well. Among the omnibus tests ( $\varphi_A$ ,  $\varphi_{COM(L)}$ ,  $\varphi_S$ , and  $\varphi_{CH}$ ),  $\varphi_{0.5,1}$  maximizes the power for 22 of the 28  $\theta$  values ( $\varphi_A$  and  $\varphi_{COM(L)}$  each maximize the power for three  $\theta$  values). Also,  $\varphi_{0.5,1}$  ( $p = 0.037$ ) and  $\varphi_S$  ( $p = 0.031$ ) are the only omnibus tests to yield statistical significance at  $\alpha = 0.05$  for  $\{(11,2,2);(7,7,6)\}$ . Table 3, above, shows that  $\varphi_{0.5,1}$  dominates both  $\varphi_S$  and  $\varphi_{CH}$ , and almost dominates  $\varphi_A$  and  $\varphi_{COM(L)}$  too, and does dominate them when  $\delta(\theta)$  is near the  $\delta_P$  value of 0.5 used by  $\varphi_{0.5,1}$ . In fact, only where  $\delta(\theta) \leq -0.5$  or  $\delta(\theta) \geq 2.5$  is  $\varphi_A$  or  $\varphi_{COM(L)}$  more powerful than  $\varphi_{0.5,1}$ . Among pairwise comparisons,  $\varphi_{0.5,1}$  has larger power than its competitor (each of the other ten tests are considered for each of 28 alternatives) for 229 out

of 280 comparisons, and 104 of the 112 comparisons to omnibus tests. The non-adaptive tests did not fare as well, but  $\varphi_S$  attained 168/280 or 57/112, respectively, which is quite respectable.

### Conclusion

In an effort to improve the comparison of two treatments on the basis of ordinal data, a new class of adaptive tests was defined, and shown to be admissible, while providing unambiguous p-values and a non-iterative construction. If one is interested in testing for  $\theta_1 > 0$ , and has no particular preference for any subset of  $\Omega_A$  relative to any other, then  $\varphi_{CH}$  would be a fine test to use.

However,  $\varphi_A$  and  $\varphi_{0.5,1}$  are also excellent omnibus tests, and are easier to compute than  $\varphi_{CH}$ . If one is interested in testing for stochastic order, and uses  $\theta_1 > 0$  only as a surrogate, then  $\varphi_A$  and  $\varphi_{0.5,1}$  are probably better tests than  $\varphi_{CH}$ . Certainly if one is in the situation treated in this article, with a preferred direction, then an appropriate adaptive test would be the test of choice. There is nothing particular about ordered trinomial distributions that makes this problem especially amenable to treatment with the adaptive

approach. For any hypothesis testing problem with a composite alternative hypothesis, one can enumerate the alternatives and the corresponding MP test for each. One can then apply each of these MP tests to a given outcome, and find the smallest of the resulting p-values. Using this minimized MP p-value as a test statistic produces a test analogous to  $\phi_A$ , and reduces to the uniformly most powerful test if one exists. If not, then the adaptive tests that bridge the gap between  $\phi_A$  and the MP tests to detect a favored direction should have good properties in a variety of contexts.

#### References

- Berger, V. W. (1998). Admissibility of exact conditional tests of stochastic order. *Journal of Statistical Planning and Inference*, 66, 39-50.
- Berger, V. W. (2000). Pros and cons of permutation tests. *Statistics in Medicine*, 19, 1319-1328.
- Berger, V. W., & Ivanova, A. (2002). The bias of linear rank tests when testing for stochastic order in ordered categorical data. *Journal of Statistical Planning and Inference*, 107, 237-247.
- Berger, V. W., Lunneborg, C., Ernst, M. D., Levine, J. G. (2002), Parametric analyses in randomized clinical trials, *Journal of Modern Applied Statistical Methods*, 1, 74-82.
- Berger, V. W., Permutt, T., & Ivanova A. (1998). The Convex hull test for ordered categorical data. *Biometrics*, 54, 1541-1550.
- Berger, V. W., & Sackrowitz, H. (1997). Improving tests for superior treatments in contingency tables. *Journal of American Statistical Association*, 92, 700-705.
- Chiara, S., Compora, E., Merlini, L., et al. (1993). Recurrent ovarian carcinoma: salvage treatment with platinum in patients responding to first-line platinum-based regimens. *European Journal of Cancer*, 29A, 652.
- Cohen, A., & Sackrowitz, H. (1998). Directional tests for one-sided alternatives in multivariate models. *Annals of Statistics*, 26, 2321-2338.
- Donegani, M. (1991). An adaptive and powerful randomization test. *Biometrika*, 78, 930-933.
- Edgington, E. S. (1995). *Randomization tests*. (3<sup>rd</sup> ed.). New York: Marcel Dekker.
- Frick, H. (2000). Undominated p-values and property  $c$  for unconditional one-sided two-sample binomial tests. *Biometrical Journal*, 42, 715-728.
- Gastwirth, J. L. (1985). The use of maximum efficiency robust tests in combining contingency tables and survival analysis. *The Journal of the American Statistical Association*, 80, 380-384.
- Good, P. (1994). *Permutation tests: a practical guide to resampling methods for testing hypotheses*. New York: Springer-Verlag.
- Gross, S. T. (1981). On asymptotic power and efficiency of tests of independence in contingency tables with ordered classifications. *Journal of American Statistical Association*, 76, 935-941.
- Graubard, B. I., & Korn, E. L. (1987). Choice of column scores for testing independence in ordered 2xk contingency tables. *Biometrics*, 43, 471-476.
- Hogg, R. V. (1974). Adaptive robust procedures: a partial review and some suggestions for future applications and theory. *Journal of American Statistical Association*, 69, 909-923.
- Ivanova, A., & Berger, V. W. (2001). Drawbacks of integer scoring of ordered categorical data. *Biometrics*, 57, 567-570.
- Marden, J. I. (1991). Sensitive and sturdy p-values. *Annals of statistics*, 19, 918-934.
- McCullagh, P. (1980). Regression methods for ordinal data. *Journal of Royal Statistical Society, B* 42, 109-142.
- Permutt, T., & Berger, V. W. (2000). Rank tests in ordered 2xk contingency tables. *Communications in Statistics, Theory and Methods*, 29, 989-1003.
- Podgor, M. J., Gastwirth, J. L., Mehta, C. R. (1996). Efficiency robust tests of independence in contingency tables with ordered categories. *Statistics in Medicine*, 15, 2095-2105.
- Rohmel, J. & Mansmann, U. (1999). Unconditional non-asymptotic one-sided tests for independent binomial proportions when the interest lies in showing non-inferiority and/or superiority. *Biometrical Journal*, 41, 149-170.
- Roy, S. N. (1953). On a heuristic method of test construction and its use in multivariate analysis. *Annals of Statistics*, 24, 220-238.

Rukhin, AL, Mak, KS (1992). Adaptive Test Statistics and Bahadur Efficiency. *Statistica Sinica*, 2, 541-552.

Appendix

Lemmas (with Proofs), and Proofs of Theorems

*Lemma 1.* Let  $c \in \Gamma$  and  $k \in \{0, 1, \dots, K_c\}$ . If  $|v_k(c) \pm \varepsilon(c)| < \infty$  then  $v_k(c) \pm \varepsilon(c) \notin V(c)$ . If  $v \in (v_k(c), v_{k+1}(c))$ , then  $M_v(c) = M_{v_{(k+1)}(c)}(c) - B_{v_{(k+1)}(c)}^-(c) = M_{v_{(k)}(c)}(c) - B_{v_{(k)}(c)}^+(c)$ .

*Proof.* Increasing (decreasing)  $v$  by  $\varepsilon(c)$  moves  $B_v^-(c)$  ( $B_v^+(c)$ ) into the interior of, and  $B_v^+(c)$  ( $B_v^-(c)$ ) completely out of, the new critical region, but if  $v \in V(c)$ , then no points of  $\Gamma - M_v(c)$  are moved into the new critical region (Table 1). Hence,  $o_{v-\varepsilon(c)}(c) = o_{v+\varepsilon(c)}(c) = 1$ , and neither  $v_k(c) - \varepsilon(c)$  nor  $v_k(c) + \varepsilon(c)$  is in  $V(c)$ . If  $v \notin V(c)$ , say  $v_k(c) < v < v_{k+1}(c)$ , then  $o_v(c) = 1$ , so  $B_v^+(c) = B_v^-(c) = \emptyset$  and  $M_v(c)$  will not change when  $v$  varies within  $(v_k(c), v_{k+1}(c))$ .

*Lemma 2.* If  $\delta_P \in \mathfrak{R}^1$ ,  $\tau > 0$ ,  $v_* \in v_{[\delta_P, \tau]}(c)$ , and  $v^* \in v^*(c)$ , then  $|\delta_P - v_*| \leq |\delta_P - v^*|$ .

*Proof.* If there exist  $v^* \in v^*(c)$  and  $v_* \in v_{[\delta_P, \tau]}(c)$  such that  $|\delta_P - v^*| < |\delta_P - v_*|$ , then  $p_{v^*}(c)(1 + |\delta_P - v^*|)^\tau < p_{\min(v_*)}(c)(1 + |\delta_P - v_*|)^\tau$ , and  $v_*$  cannot be in  $v_{[\delta_P, \tau]}(c)$ .

*Lemma 3.* For any  $\delta_P$ ,  $\tau > 0$ , and  $c \in \Gamma$ ,  $v_{[\delta_P, \tau]}(c) \subset V(c) \cup \delta_P$ .

*Proof.* Assume there exists  $v \neq \delta_P$  in  $v_{[\delta_P, \tau]}(c) - V(c)$ , say  $v_k(c) < v < v_{k+1}(c)$ . Let  $v^* = v_k(c)$  if  $\delta_P \leq v_k(c)$ ,  $v^* = \delta_P$  if  $v_k(c) < \delta_P < v_{k+1}(c)$ , or  $v^* = v_{k+1}(c)$  if  $v_{k+1}(c) \leq \delta_P$ . Now  $v^* \subset V(c) \cup \delta_P$  and

$$p_{\min(v)}(c)(1 + |\delta_P - v|)^\tau > p_{\min(v^*)}(c)(1 + |\delta_P - v^*|)^\tau.$$

*Lemma 4.* For any  $\delta_P$  and  $c \in \Gamma$ ,  $v_{[\delta_P, \tau]}(c) = \{\delta_P\}$  for sufficiently large  $\tau$ .

*Proof.* Let  $D_c(\delta_P) = \min_{v \in V(c) - \delta_P} |\delta_P - v| > 0$ . For  $\tau > 0$ , let  $v \in v_{[\delta_P, \tau]}(c) - \delta_P$ . By Lemma 3,  $v \in V(c) - \delta_P$ , so  $|\delta_P - v| \geq D_c(\delta_P)$ . If  $\tau > -\ln(p_{\min(\delta_P)}(c))/\ln(1 + D_c(\delta_P))$ , then  $p_{\min(v)}(c)(1 + |\delta_P - v|)^\tau \geq p_{\min(v)}(c)(1 + |D_c(\delta_P)|)^\tau > 1$ , contradicting  $v \in v_{[\delta_P, \tau]}(c)$ .

*Proof of Theorem 1.* By Theorem 3.3 of Berger (1998), it suffices to show that for any  $B \subset \Gamma$ , if  $c^*$  minimizes  $A(\delta_P, \tau, c)$  over  $B$ , then  $c^* \in D[B]$ . If  $c^* \notin D[B]$ , then  $c^*$  cannot, for any  $v$ , uniquely minimize  $p_v$  over  $B$ , and for every  $v$  there exists  $c \in B - c^*$  such that  $p_v(c) \leq p_v(c^*)$ . If  $v \notin V(c^*)$ , then  $o_v(c^*) = 1$ , so  $p_v(c) \neq p_v(c^*)$ , and  $p_v(c) \leq p_v(c^*) - \min_{c \in \Gamma} P_0\{c|\Gamma\}$ . Let  $v_1 \in v_{[\delta_P, \tau]}(c^*)$ . By the continuity in  $v$  of the function  $(1 + |\delta_P - v|)^\tau$ , one can, for any  $\varepsilon > 0$ , choose  $v_2 \notin V(c^*)$  suitably close to  $v_1$  to satisfy  $p_{v_2}(c^*) = p_{\min(v_1)}(c^*)$ , and, thus,

$$\begin{aligned} A(\delta_P, \tau, c) &= \min_{-\infty \leq v \leq \infty} [p_{\min(v)}(c)(1 + |(\delta_P - v)|)^\tau] \leq \\ & p_{v_2}(c)(1 + |(\delta_P - v_2)|)^\tau \\ & \leq [p_{v_2}(c^*) - \min_{c \in \Gamma} P_0\{c|\Gamma\}](1 + |\delta_P - v_2|)^\tau \\ & = [p_{\min(v_1)}(c^*) - \min_{c \in \Gamma} P_0\{c|\Gamma\}](1 + |\delta_P - v_2|)^\tau \\ & < A(\delta_P, \tau, c^*) - \min_{c \in \Gamma} P_0\{c|\Gamma\}(1 + |\delta_P - v_2|)^\tau + \varepsilon < A(\delta_P, \tau, c^*), \end{aligned}$$

the last inequality holding for  $\varepsilon < \min_{c \in \Gamma} P_0\{c|\Gamma\}$ . This is a contradiction.