

# Multiple Endpoints: An Overview and New Developments

Ajit C. Tamhane                      and                      Brent R. Logan  
Department of Statistics                      Division of Biostatistics  
Northwestern University                      Medical College of Wisconsin

Division of Biostatistics  
Medical College of Wisconsin

Technical Report 43

September 2003

Division of Biostatistics  
Medical College of Wisconsin  
8701 Watertown Plank Road  
Milwaukee, WI 53226  
Phone:(414) 456-8280



# Multiple Endpoints: An Overview and New Developments

Ajit C. Tamhane

Department of Statistics

Northwestern University

2006 ShereAn7Mo10u

## Summary

In the last two decades a large number of papers have been published on the topic of analysis of multiple endpoints in clinical trials. We provide a comprehensive review of this vast literature focusing on the statistical aspects. We make comparisons between competing procedures, present some new developments and extensions/modifications of existing procedures, make recommendations for use and note some open problems for research.

**Keywords:** Multiple comparisons; multiple tests; one-sided multivariate tests; Bonferroni test; chi-bar squared distribution; multivariate normal distribution; clinical decision rules; global tests; endpoint specific tests; closure method; resampling; adjusted  $p$ -values; family-wise error rate.

## 1. Introduction

Most clinical trials are conducted to compare a treatment group with a control group on multiple endpoints. Often, the treatment is expected to have a positive effect on all endpoints. Depending on the nature of the disease the endpoints may be grouped into primary and secondary types. We mainly focus on the case where all endpoints are primary and provide a comprehensive review of the vast literature and some new results focusing on the statistical aspects. Shorter review articles by Chi (1998), Huque and Sankoh (1997), Sankoh, Huque and Dubey (1997), Sankoh, Huque, Russell and D'Agostino (1999) and Zhang, Quan, Ng and Stepanavage (1997) also discuss some clinical aspects with examples.

The corresponding correlation matrix will be denoted by  $\mathbf{R}$  with elements

$$r_{jk} = \text{Corr}(x_{ij^k}, x_{ij}) = \frac{c_{jk}}{\sqrt{c_{jj}c_{kk}}} \quad (1 \leq j, k \leq m).$$

In the heteroscedastic case,  $c_{11}$  and  $c_{22}$  are not assumed to be equal. The elements of

### 3.1 Homoscedastic Case

#### 3.1.1 Exact Likelihood Ratio (LR) Tests

It is well-known that because Hotelling's  $T^2$  test is designed for the omnibus (two-sided) alternative  $H_2 : \boldsymbol{\mu} = \mathbf{0}$ , it lacks power for the one-sided alternative  $H_1$  of (2.1) (Meier 1975, O'Brien 1984). Kudô (1963) derived an exact LR test when  $\boldsymbol{\Sigma}$  is *known* for the one-sample problem which can be easily extended to the two-sample problem as follows. Let  $\boldsymbol{\mu}^*$  be the projection of  $\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2$  in the positive orthant with respect to the distance function  $d(\mathbf{u}, \mathbf{v}) = (\mathbf{u} - \mathbf{v})' \boldsymbol{\Sigma}^{-1}(\mathbf{u} - \mathbf{v})$

Stein-type two-stage test that is free of . It also has other desirable properties such as

The null distribution of  $g(\mathbf{u})$  is the  $\chi^2$  distribution with symmetric binomial probability weights given by

$$\Pr_{H_0}\{g(\mathbf{u}) > c\} = \sum_{k=0}^m \binom{m}{k} 2^{-m} \Pr\left\{\frac{z}{k} > c\right\}$$





We see that both the OLS and GLS statistics are standardized weighted sums of the individual  $t$ -statistics for the  $m$  endpoints. The OLS statistic uses equal weights, while the GLS statistic uses unequal weights determined by the sample correlation matrix  $\mathbf{R}$ . If some endpoint is highly correlated with the others then the GLS statistic gives a correspondingly lower weight to its  $t$ -statistic.

The exact small sample null distributions of  $t_{OLS}$  and  $t_{GLS}$  are not known. O'Brien (1984) proposed a  $t$ -distribution with  $n_1 + n_2 - 2m$  d.f. as an approximation. For large sample sizes the standard normal ( $z$ ) distribution may be used as an approximation. The  $t$ -approximation is exact for  $m = 1$ , but is conservative for  $m > 1$ ; on the other hand, the  $z$ -approximation is liberal. The convergence of  $t_{GLS}$  to the standard normal distribution is slower than that of  $t_{OLS}$  because of the use of the estimated correlation matrix  $\mathbf{R}$  both in the calculation of  $t_{GLS}$  and in the estimate of  $SD(t_{GLS})$ . Also, the simulation study by Reitmeir and Wassmer (1996) has shown that the powers of the OLS and GLS tests are comparable when used to test subset hypotheses in closed testing procedures (see Section 4.1). Finally, the linear combination  $\mathbf{j} \mathbf{R}^{-1}$  used by the GLS test can have some negative weights, which can lead to anomalous results; this problem does not occur with the OLS test. For all these reasons, the OLS test is recommended.

Finally we note that Tang, Gnecco and Pocock (1993) have generalized the GLS test statistic for an arbitrary ray alternative  $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2 = (\alpha_1, \dots, \alpha_m)$ , where the vector  $(\alpha_1, \dots, \alpha_m)$  with all positive elements is specified. However, if the observed mean difference  $\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2$  is not close to this ray then the power of the test may be adversely affected. Since the vector  $(\alpha_1, \dots, \alpha_m)$  is in general difficult to specify, Tang, Gnecco and Pocock suggest following the maxmin approach (maximize the minimum power over all ray alternatives) of Abelson and Tukey (1963).

### 3.1.4 Lauter's Exact Tests

Lauter (1996) proposed a class of test statistics for the hypotheses (2.1) having the property that they are exactly  $t$ -distributed with  $n_1 + n_2 - 2$  d.f. under  $H_0$ . Recall that  $\bar{\mathbf{x}}_i = (\bar{x}_{i,1}, \bar{x}_{i,2}, \dots, \bar{x}_{i,m})$  denotes the vector of sample means for the  $i$ th group ( $i = 1, 2$ ) and

let

$$\bar{\mathbf{x}}_{..} = \frac{n_1 \bar{\mathbf{x}}_{1.} + n_2 \bar{\mathbf{x}}_{2.}}{n_1 + n_2} = (\bar{x}_{..1}, \bar{x}_{..2}, \dots, \bar{x}_{..m})$$

denote the vector of overall sample means. Define the total cross-products matrix by

$$\mathbf{V} = \sum_{i=1}^2 \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_{..})(\mathbf{x}_{ij} - \bar{\mathbf{x}}_{..}) = (n_1 + n_2 - 2) \sum_{i=1}^2 n_i (\bar{\mathbf{x}}_{i.} - \bar{\mathbf{x}}_{..})(\bar{\mathbf{x}}_{i.} - \bar{\mathbf{x}}_{..})$$

Because the total pooled standard deviation overestimates the true standard deviation since



of Hotelling's  $T^2$  test by eliminating outcomes with negative differences on all endpoints, but its rejection region is not monotone.

### 3.2 Heteroscedastic Case

#### 3.2.1 Approximate Likelihood Ratio (ALR) Test

In Tamhane and Logan (2001) we proposed to extend the ALR test to the heteroscedastic case as follows. Let

$$\mathbf{V}_i = \frac{1}{n_i} \mathbf{V}_i \quad (i = 1, 2) \quad = \quad \mathbf{V}_1 + \mathbf{V}_2 \quad \text{and} \quad = \frac{n_1 n_2}{n_1 + n_2} .$$

The sample estimates of these matrices are denoted by putting carets over them; thus  $\hat{\mathbf{V}}_i = (1/n_i) \hat{\mathbf{V}}_i$ ,  $\hat{\mathbf{V}} = \hat{\mathbf{V}}_1 + \hat{\mathbf{V}}_2$  and

$$= \frac{n_1 n_2}{n_1 + n_2} .$$

The transformation matrix  $\mathbf{A}$  in (3.1) is chosen such that  $\mathbf{A} \mathbf{A} = \mathbf{I}^{-1}$  and  $\mathbf{A} \mathbf{A} = \mathbf{I}$ .

We suggested the same  $\bar{F}$  approximation (3.4) to the null distribution of  $g(\mathbf{u})$  in the heteroscedastic case, but with the following Welch-Satterthwaite estimated d.f. derived by Yao (1965) for the multivariate Behrens-Fisher problem:

$$\frac{1}{\nu} = \frac{1}{(\mathbf{d}^{-1} \mathbf{d})^2} \left[ \frac{(\mathbf{d}^{-1} \mathbf{d}_1^{-1} \mathbf{d})^2}{n_1 - 1} + \frac{(\mathbf{d}^{-1} \mathbf{d}_2^{-1} \mathbf{d})^2}{n_2 - 1} \right] ,$$

where  $\mathbf{d} = (\mathbf{x}_1, -\mathbf{x}_2)$ . Note that Yao derived this formula (also using the moment matching method) to approximate the distribution of

$$\mathbf{u} \mathbf{u} = \frac{n_1 n_2}{n_1 + n_2} (\mathbf{x}_1, -\mathbf{x}_2)^{-1} (\mathbf{x}_1, -\mathbf{x}_2)$$

by Hotelling's  $T_{m, -m+1}^2 = \frac{m}{-m+1} F_{m, -m+1}$  distribution with an estimated  $\nu$ . We simply extended Yao's approximation to the  $\bar{F}$  distribution. Our simulations for selected values of  $m, n_1 = n_2 = n, \nu_1$  and  $\nu_2$  showed that this approximation is quite accurate for controlling the type I error probability at the nominal level  $\alpha = 0.05$  for  $m = 4$  if  $n \geq 20$  and for  $m = 8$  if  $n \geq 30$ .

### 3.2.2 Ordinary Least Squares (OLS) and Generalized Least Squares (GLS) Tests

Pocock, Geller and Tsiatis (1987) extended O'Brien's GLS test to the heteroscedastic case as follows. Initially assume that  $\sigma_1^2$  and  $\sigma_2^2$  are known. Then the statistic for comparing the treatment with the control on the  $k$ th endpoint is

$$Z_k = \frac{\bar{X}_{1\cdot k} - \bar{X}_{2\cdot k}}{\sqrt{\sigma_{1,kk}/n_1 + \sigma_{2,kk}/n_2}} \quad (1 \leq k \leq m). \quad (3.11)$$

Let  $\mathbf{z} = (z_1, z_2, \dots, z_m)$  and  $\bar{\mathbf{R}} = (n_1 \mathbf{R}_1 + n_2 \mathbf{R}_2)/(n_1 + n_2)$ . In analogy with (3.9), Pocock et al. proposed the statistic

$$Z_{\text{GLS}} = \frac{\mathbf{j} \bar{\mathbf{R}}^{-1} \mathbf{z}}{\sqrt{\mathbf{j} \bar{\mathbf{R}}^{-1} \mathbf{j}}}.$$

However, this is just an ad-hoc extension. Furthermore, the covariance (correlation) matrix of  $\mathbf{z}$  is not  $\bar{\mathbf{R}}$ , but  $\mathbf{C} = \{c_k\}$  with elements

$$c_k = \sigma_{1,k}^2/n_1 + \sigma_{2,k}^2/n_2$$

for  $i = 1, 2$ . Note that  $\rho_{1k} - \rho_{2k} = \rho_k$  for all  $k$ . Also note that  $\Sigma_1$  and  $\Sigma_2$  are not correlation matrices, and  $\Sigma = \Sigma_1 + \Sigma_2$  if  $n_1 = n_2$ .

The hypotheses (3.5) can be tested by using a univariate regression framework as in (3.6):

$$y_{ijk} = \mu_k + \frac{1}{2} I_{ijk} + \epsilon_{ijk} \quad (i = 1, 2; 1 \leq j \leq n_i; 1 \leq k \leq m), \quad (3.12)$$

where  $\mu_k = (\rho_{1k} + \rho_{2k})/2$ ,  $I_{ijk} = +1$  if  $i = 1$  and  $-1$  if  $i = 2$ , and  $\epsilon_{ij} = (\epsilon_{ij1}, \epsilon_{ij2}, \dots, \epsilon_{ijm})$  are independently distributed as  $N(\mathbf{0}, \Sigma_i)$ . Using the same methods as those used in the homoscedastic case, the OLS and GLS statistics are as given below; for derivations, see Logan (2001).

Assuming that  $\Sigma_1$  and  $\Sigma_2$  are known, it is straightforward to show that

$$t_{\text{OLS}} = \frac{\mathbf{j}'(\mathbf{y}_1 - \mathbf{y}_2)}{m} = \bar{y}_{1..} - \bar{y}_{2..} \quad \text{and} \quad \text{SD}(t_{\text{OLS}}) = \sqrt{\mathbf{j}' \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \mathbf{j}}.$$

Hence the OLS statistic with the  $\Sigma_i$  replaced by their sample estimates  $\hat{\Sigma}_i$  equals

$$t_{\text{OLS}} = \frac{\mathbf{j}'(\mathbf{y}_1 - \mathbf{y}_2)}{\mathbf{j}' \left( \hat{\Sigma}_1/n_1 + \hat{\Sigma}_2/n_2 \right) \mathbf{j}}, \quad (3.13)$$

where the elements of  $\hat{\Sigma}_i$  are given by

$$\hat{\Sigma}_{i,k} = \frac{s_{i,k}^2}{(s_{i,kk} + s_{2,kk})(s_{i,1} + s_{2,1})}.$$

For  $n_1 = n_2 = n$ , the above OLS statistic reduces to

$$t_{\text{OLS}} = \frac{\mathbf{j}' \mathbf{t}}{\mathbf{j}' \hat{\Sigma} \mathbf{j}},$$

where  $\mathbf{t}$  is a vector of



Next we derive the GLS test. Assuming that  $\beta$  is known, it can be shown

$$b_{GLS} = \frac{4j \left( \frac{1}{n_1} + \frac{1}{n_2} \right)^{-1} (y_1 - y_2)}{j [(I - B)^{-1} \frac{1}{n_1} + (I + B)^{-1} \frac{1}{n_2}] j}$$

and

$$SD(b_{GLS}) = \frac{4 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)^{-1} j}{j [(I - B)^{-1} \frac{1}{n_1} + (I + B)^{-1} \frac{1}{n_2}] j} \cdot \frac{1}{1 - \frac{1}{2} \frac{1}{n_1} \frac{1}{n_2} \frac{1}{n_1 + n_2}}$$

The problems associated with the Bonferroni test have been well-documented; see, e.g., O'Brien (1984), Pocock, Geller and Tsiatis (1987): (i) it is overly conservative especially if  $m$  is large or the endpoints are highly correlated, and (ii) it is powerful if only one endpoint has a large treatment effect, but not if most or all endpoints have moderate treatment effects.

It should be noted that the Bonferroni test is a union-intersection (UI) test when  $H_0$  is viewed as  $H_0 = \bigcap_{k=1}^m H_{0k}$ . Therefore rejection of  $H_0$  implies rejection of any  $H_{0k}$  with  $p_k < \alpha/m$ ; this implied multiple test procedure for testing null hypotheses on the individual endpoints controls the FWE at level  $\alpha$  (Hochberg and Tamhane 1987, pp. 28-29).

An improvement on the Bonferroni test was proposed by Simes (1986). To apply the Simes test first order the  $p$ -values:  $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(m)}$  and denote the corresponding hypotheses by  $H_{0(1)}, H_{0(2)}, \dots, H_{0(m)}$ . Then reject  $H_0$  if

$$p_{(k)} < \frac{(m - k + 1)\alpha}{m} \text{ for some } k = 1, 2, \dots, m. \quad (3.17)$$

Simes proved that this is an  $\alpha$ -level test under the assumption that the  $p$ -values are independent.

global test of each null hypothesis  $H_{0k}$  for  $K = M$ . Any of the global tests discussed in the previous section can be used for this purpose.

## 4.2 Normal Theory Based Tests

It is conceivable to test the null hypotheses  $H_{0k}$  of (2.2) on the individual endpoints using the test statistics  $t_k$  from (3.8) in the homoscedastic case and from (3.14) in the heteroscedastic case. To control the FWE at level  $\alpha$ , we would need the upper  $\alpha$  critical point of  $\max_{1 \leq k \leq m} t_k$  under the overall null hypothesis  $H_0$  in each case. However, the joint distribution of  $(t_1, t_2, \dots, t_m)$  is not multivariate  $t$  even in the homoscedastic case because the standard deviations  $\sqrt{s_{kk}}$  used to standardize the  $t_k$  statistics are different though correlated for  $k = 1, 2, \dots, m$ . Furthermore these correlations (as well as those between the numerators of the  $t_k$  statistics) are unknown being the correlations between the corresponding endpoints. Therefore the standard Dunnett-type (1955) test or its stepwise versions (Dunnett and Tamhane 1991, 1992) cannot be applied to test the hypotheses  $H_{0k}$ .

## 4.3 Procedures Based on Adjusted $p$ -Values

Let  $p_k$  be the  $p$ -value for testing  $H_{0k}$  as discussed in Section 3.3 and let  $P_k$  be the corresponding r.v. This  $p$ -value is not adjusted for multiplicity of tests on all  $H_{0k}$ . A way to control the FWE at level  $\alpha$  is to find multiplicity adjusted  $p$ -values (see Dunnett and Tamhane 1991, 1992 and Wright 1992), denoted by  $\rho_k$ , and reject  $H_{0k}$  if  $\rho_k < \alpha(1 - k/m)$ .

The adjusted  $p$ -values corresponding to a single-step test procedure (see Hochberg and Tamhane 1987, Ch. 2) are given by

$$\rho_k = \Pr_{H_0} \left( \min_{1 \leq m} P \leq p_k \mid 1 \leq k \leq m \right). \quad (4.1)$$

The joint distribution of  $(P_1, P_2, \dots, P_m)$  is unknown because of the unknown correlations among the endpoints. Therefore an approximation is often needed. The simplest such approximation is the Bonferroni adjustment (corresponding to the Bonferroni test) given by

$$\rho_k = mp_k \alpha(1 - k/m).$$

Various sharpened versions of the Bonferroni adjusted  $p$ -values are available based on the Šidák (1968) inequality and its modifications. The Šidák adjustment assumes that the  $P$ 's

are independent and is given by

$$p_k = 1 - (1 - p_k)^m \quad (1 \leq k \leq m).$$

If the  $P$ 's are positively dependent then this adjustment is conservative. Armitage and Parmar (1986) gave the following ad-hoc approximation to the adjusted  $p$ -values that takes into account the correlations between the endpoints:

$$p_k = 1 - (1 - p_k)^{m^f} \quad (1 \leq k \leq m),$$

where  $f$  is an empirically determined function of the  $\rho_k$ 's. Dubey (1985) suggested using a different function  $\bar{f}_k = 1 - \bar{\rho}_k$  for each  $k$ , where  $\bar{\rho}_k$  is the average of the correlations of the  $k$ th endpoint with the others. However, it is readily seen from the definition (4.1) of the adjusted  $p$ -value that  $f$  must be a symmetric function of all correlations. Therefore  $\bar{\rho}_k$  in Dubey's formula should be replaced by  $\bar{\rho}$ , namely, the average of all  $\rho_k$ 's. Notice that if all  $\rho_k = 0$  then we get the Šidák adjustment and if all  $\rho_k = 1$  then  $p_k = p_k$ , i.e., there is no adjustment. Tukey, Ciminera and Heyse (1985) suggested using  $f = 1/2$ , i.e.,  $p_k = 1 - (1 - p_k)^{\bar{m}}$ , which assumes that the average correlation is  $1/2$ . An analytic approximation to  $p_k$  for jointly normally distributed endpoints was proposed by James (1991). Finally, Westfall and Young's (1989,1993) resampling method, which is distribution-free and implicitly takes the correlations between the endpoints into account can always be applied to estimate the  $p_k$ .

For multivariate binary endpoints, a bootstrap method was given by Westfall and Young (1989) which was further extended to many other multiple testing problem in their 1993 book. Chen (1998) proposed using the generalized estimating equation (GEE) approach to estimate the unknown correlations of binary endpoints to find the adjusted  $p$ -values.

Another approach to sharpen the Bonferroni adjustment is to use a stepwise procedure for testing. The adjusted  $p$ -values for a step-down test procedure are given by

$$\begin{aligned} p_{(m)} &= \Pr_{H_0} \left[ \min_m P \leq p_{(m)} \right] \quad \text{and} \\ p_{(k)} &= \max \{ p_{(k+1)}, \Pr_{H_0} \left[ \min_k P \leq p_{(k)} \right] \} \quad \text{for } k = 1, \dots, m-1. \end{aligned} \quad (4.2)$$

Conservative approximations to the above adjusted  $p$ -values can be obtained by using the Bonferroni inequality and are given by

$$p_{(m)} = mp_{(m)} \quad \text{and} \quad p_{(k)} = \max \{ p_{(k+1)}, kp_{(k)} \} \quad (1 \leq k \leq m-1).$$

These approximations correspond to Holm's (1979) step-down test procedure, which rejects  $H_{0(k)}$  if  $p_{(k)} < \alpha/k$  for  $k = 1, 2, \dots, m$ . This procedure can be derived by using the Bonferroni test (3.16) to test subset null hypotheses in the closure method.

Hommel (1988) derived a stepwise procedure by using the Simes test (3.17) to test subset null hypotheses in the closure method. Hochberg (1988) offered a slightly conservative but a much simpler procedure. It is of step-up type in that it is the exact opposite of Holm's step-down procedure in terms of sequence of testing. The adjusted  $p$ -values for the Hochberg procedure are given by

$$p_{(1)} = p_{(1)} \text{ and } p_{(k)} = \min(p_{(k-1)}, kp_{(k)}) \quad (2 \leq k \leq m).$$

Hochberg's procedure accepts  $H_{0(k)}$  if  $p_{(k)} < \alpha/k$  for  $k = 1, 2, \dots, m$ . Troendle (1996) gave a bootstrap-based permutational step-up procedure.

#### 4.4 A Hybrid Method Combining Global and Endpoint-Specific Tests

As we have seen, there are two main approaches to identify the significant endpoints: (i) adjusting the  $p$ -values of individual endpoints, and (ii) using the closure method that employs one of the global tests to test subset null hypotheses. The first approach is more powerful when only a few endpoints have positive treatment effects, while the second approach is more powerful when all or most of the endpoints have an effect. A test procedure with a more uniform power performance can be obtained by combining these two approaches along the lines of Hothorn's (1999)  $T_{\max}$  testing principle.

In Logan and Tamhane (2001) we gave a closed testing procedure by combining two tests for testing each intersection hypothesis: (i) the Bonferroni  $p_{\min}$  test and (ii) O'Brien's OLS test. According to this latter hybrid method, the adjusted  $p$ -value for any intersection hypothesis  $H_{0K} = \bigcap_{k \in K} H_{0k}$  is defined as

$$p_K = \Pr_{H_0} \left( \min_{k \in K} P_k, P_{K,OLS} \leq \min_{k \in K} p_k, p_{K,OLS} \right), \quad (4.3)$$

where, as before, the lower case  $p$ 's denote the unadjusted observed  $p$ -values (e.g.,  $p_k$  is the  $p$ -value for  $H_{0k}$  and  $p_{K,OLS}$  is the  $p$ -value for  $H_{0K}$  using the OLS test) and the upper case  $P$ 's denote the corresponding r.v.'s. In Logan (2001) a third test was added, namely the ALR

test. In a closed testing procedure, a hypothesis  $H_{0K}$  is rejected at level  $\alpha$  if all hypotheses  $H_{0L}$  for  $L \subseteq K$  are rejected at level  $\alpha$  and  $p_K < \alpha$ . In practice, the  $p_K$  defined in (4.3) need to be estimated by bootstrap resampling. C language programs for this purpose for both the homoscedastic as well as the heteroscedastic case are posted on the first author's home page (<http://users.iems.northwestern.edu/~ajit>).

configurations at which the type I error probability is maximized ( $= \alpha$ ) can be shown to be of the type  $\mu_k = 0$  for some  $k$  and  $\mu_j = \alpha$  for  $j = k$ . Cappizi and Zhang (1996) argued that the resulting MIN test is overly conservative. If the null hypothesis is restricted to  $H_0: \prod_{k=1}^m (\mu_k = 0)$  as in (2.1) then a much less conservative test is obtained. Snappin (1987)

for assessing the efficacy of a treatment under an intersection null hypothesis framework. The MIN test is useful for dealing with a union null hypothesis. Often, protocols for drug approval specify decision rules based on a combination of union and intersection null hypotheses. Many examples of such decision rules are given in Chi (1998, 2000). In this section we present two common types of clinical decision rules, give some examples, and discuss how formulating these decision rules as a combination of union and intersection null hypotheses can lead to FWE controlling procedures.

A typical decision rule leads to several paths for finding a significant treatment effect. For example, given three endpoints (e.g., one primary and two secondary), one might conclude effectiveness if either  $\tau_1 > 0$  or  $(\tau_2 > 0 \text{ and } \tau_3 > 0)$ , i.e., if the primary endpoint shows an effect or both secondary endpoints show an effect. As another example, given four endpoints, two primary and two secondary, a possible decision rule might be to conclude effectiveness if at least one primary endpoint and at least one secondary endpoint is significant, i.e., if  $(\tau_1 > 0 \text{ or } \tau_2 > 0)$  and  $(\tau_3 > 0 \text{ or } \tau_4 > 0)$ .

In each of the above cases, the decision rule corresponds to an alternative hypothesis, from which an appropriate null hypothesis can be constructed by taking the complement. Let  $H_{0i} : \tau_i = 0$  and  $H_{1i} : \tau_i > 0$  for each endpoint  $i$ . Then the alternative hypothesis for the first example is

$$H_1 : H_{11} \cap (H_{12} \cap H_{13}),$$

and the null hypothesis is

$$H_0 : H_{01} \cap (H_{02} \cap H_{03}) = (H_{01} \cap H_{02}) \cup (H_{01} \cap H_{03}).$$

Then applying the IU principle, we can test each intersection null hypothesis at level  $\alpha$  and conclude that the treatment is effective if both intersection null hypotheses are rejected. Similarly for the second case, the alternative hypothesis is

$$H_1 : (H_{11} \cap H_{12}) \cap (H_{13} \cap H_{14}),$$

and the null hypothesis is

$$H_0 : (H_{01} \cap H_{02}) \cup (H_{03} \cap H_{04}).$$



Again applying the IU principle, we can test each intersection null hypothesis at level  $\alpha$  and conclude effectiveness of the treatment if both intersection null hypotheses are rejected. Neuhäuser, Steinijans and Bretz (1999) gave an example of this method using the Simes test

Again using the IU principle, the resulting method is to test each equivalency hypothesis at level  $\alpha$  and to test the intersection hypothesis at level  $\alpha$  as well. If all hypotheses are rejected then conclude that the treatment is effective at level  $\alpha$ .

As demonstrated above, test procedures can be constructed for desired clinical decision rules which control the error rate at a pre-specified level  $\alpha$ .

## References

1. Abelson, R. P. and Tukey, J. W. (1963). Efficient utilization of non-numerical infor-

11. Dunnett, C. W. and Tamhane, A. C. (1991). Step-down multiple tests for comparing

23. Kieser, M., Reitmeir, P. and Wassmer, G. (1995). Test procedures for clinical trials with multiple endpoints. *Biometrie in der chemisch-pharma-zeitischen Industrie* (ed. J. Vollmar), **6**, Stuttgart: Gustav Fischer Verlag, 41 -60.
24. Kropf, ]TJ/Td[(,)-326(Stuttg)1(art955Tf222222228dures)-375(Appli22cuttg-330tionTf2326(f1.))-6

34. Neuhäuser, M., Steinijans, V. W. and Bretz, F. (1999). The evaluation of multiple clinical endpoints, with application to asthma. *Drug Information Journal*, **33**, 471-477.
35. O'Brien, P. C. (1984). Procedures for comparing samples with multiple endpoints. *Biometrics*, **40**, 1079-1087.
36. Perlman, M. D. (1969). One-sided testing problems in multivariate analysis. *Annals of Mathematical Statistics*, **40**, 549-567.

44. Sen, P. K. and Tsai, M-T (1999). Two-stage likelihood ratio and union-intersection tests for one-sided alternatives multivariate mean with nuisance dispersion matrix. *Journal of Multivariate Analysis*, **68**, 264 -282.
45. Šidák, Z. (1968). On multivariate normal probabilities of rectangles: Their dependence on correlations. *Annals of Mathematical Statistics*, **39**, 1425 -1434.
46. Simes, R. J. (1986). An improved Bonferroni procedure for multiple tests of significance. *Biometrika*, **73**, 751 - 754.
47. Silvapulle, M. J. (1997). A curious example involving the likelihood ratio test against one-sided hypotheses. *American Statistician*, **51**, 178 -180.
48. Snappin, S. M. (1987). Evaluating the efficacy of a combination therapy. *Statistics in Medicine*. **6**, 657–665.
49. Tamhane, A. C., Liu, W. and Dunnett, C. W. (1998). A generalized step-up-down multiple test procedure. *Canadian Journal of Statistics*, (1998), **26**, 353–363
50. Tamhane, A. C. and Logan, B. R. (2001). Accurate critical constants for the one-sided

55. Tukey, J. W., Ciminera, J. L. and Heyse, J. P. (1985). Testing the statistical certainty of a response with increasing doses of a compound. *Biometrics*, **41**, 295–301.
56. Wang, Y. and McDermott, M. P. (1998). Conditional likelihood ratio test for a non-negative normal mean vector. *Journal of the American Statistical Association*, **89**, 380–386.
57. Westfall, P. H. and Young, S. S. (1989).  $P$ -value adjustment for multiple tests in multivariate binary models. *Journal of the American Statistical Association*, **84**, 780–786.
58. Westfall, P. H. and Young, S. S. (1993) *Resampling Based Multiple Testing*. John Wiley: New York.
59. Wright, S. P. (1992). Adjusted  $p$ -value for simultaneous inference. *Biometrics*, 1005–1013.
60. Yao, Y. (1965). An approximate degrees of freedom solution to the multivariate Behrens-Fisher problem. *Biometrika*, **52**, 139–147.
61. Zhang, J., Quan, H., Ng J. and Stepanavage, M. E. (1997). Some statistical methods for multiple endpoints in clinical trials, *Controlled Clinical Trials*, **18**, 204–221.