



Step-up and step-down methods for testing multiple hypotheses in sequential experiments

Shyamal K. De, Michael Baron *

Department of Mathematical Sciences, The University of Texas at Dallas, Richardson, TX 75080, United States

ARTICLE INFO

Article history:

Received 18 July 2011

Received in revised form

2 February 2012

Accepted 6 February 2012

Available online 14 February 2012

Keywords:

Bonferroni methods

Familywise error rate

Holm procedure

Sequential probability ratio test

Stopping boundaries

Wald approximation

ABSTRACT

Sequential methods are developed for testing multiple hypotheses, resulting in a statistical decision for each individual test and controlling the familywise error rate and the familywise power in the strong sense. Extending the ideas of step-up and step-down methods for multiple comparisons to sequential designs, the new techniques improve over the Bonferroni and closed testing methods proposed earlier by a substantial reduction of the expected sample size.

© 2012 Elsevier B.V. All rights reserved.

1. Introduction

1.1. Motivation

The problem of multiple inferences in sequential experiments arises in many fields. Typical applications are in sequential clinical trials with both efficacy and safety endpoints (Jennison and Turnbull, 1993) or several outcome measures of efficacy (O'Brien, 1984; Pocock et al., 1987), acceptance sampling with several different criteria of acceptance (Baillie, 1987; Hamilton and Lesperance, 1991), multichannel change-point detection (Tartakovsky and Veeravalli, 2004; Tartakovsky et al., 2003) and in microarray experiments (Dudoit et al., 2003). It is often necessary to find the statistical answer to each posed question by testing each individual hypothesis rather than giving one global answer by combining all the tests into one and testing a composite hypothesis.

Methods developed in this paper aim to test *multiple hypotheses* based on sequentially collected data, resulting in *individual decisions* for each individual test. They control the familywise error rate and the familywise power in the strong sense. That is, both probabilities of rejecting at least one true null hypothesis and accepting at least one false null hypothesis are kept within the chosen levels α and β under any set of true hypotheses. This condition is a multi-testing analogue of controlling both probabilities of Type I and Type II errors in sequential experiments. As a result, the *familywise power*, defined as the probability of detecting *all* significant differences at the specified alternative parameter values, is controlled at the level $(1 - \beta)$ (see Shaffer, 1995, for three alternative definitions of familywise power).

* Corresponding author.

E-mail addresses: shyamal@utdallas.edu (S.K. De), mbaron@utdallas.edu (M. Baron).

Under these conditions, proposed stopping rules and decision rules achieve substantial reduction of the expected sample size over all the existing (to the best of our knowledge) sequential multiple testing procedures.

1.2. Sequential multiple comparisons in the literature

The concept of multiple comparisons is not new in sequential analysis. Sequential methods exist for inferences about multivariate parameters (Ghosh et al., 1997, Sections 6.8 and 7.5). They are widely used in studies where inferences about individual parameters are not required.

Most of the research in sequential multiple testing is limited to two types of problems.

One type is the study of several ($k > 2$) treatments comparing their effects. Sampled units are randomized to k groups where treatments are administered. Based on the observed responses, one typically tests a composite null hypothesis $H_0 : \theta_1 = \dots = \theta_k$ against H_A : not H_0 , where θ_j is the effect of treatment j for $j = 1, \dots, k$ (Betensky, 1996; Edwards, 1987; Edwards and Hsu, 1983; Hughes, 1993; Jennison and Turnbull, 2000, Chapter 16; O'Brien and Fleming, 1979; Siegmund, 1993; Wilcox, 2004; Zacks, 2009, Chapter 8). Sometimes each treatment is compared to the accepted standard (e.g., Paulson, 1962), and often the ultimate goal is selection of the best treatment (Jennison et al., 1982; Paulson, 1964).

The other type of studies involves a sequentially observed sequence of data that needs to be classified into one of the several available sets of models. In a parametric setting, a null hypothesis $H_0 : \theta \in \Theta_0$ is tested against several alternatives, $H_1 : \theta \in \Theta_1$ vs ... vs $H_k : \theta \in \Theta_k$, where θ is the common parameter of the observed sequence (Armitage, 1950; Baum and Veeravalli, 1994; Novikov, 2009; Simons, 1967).

The optimal stopping rules for such tests are (naturally!) extensions of the classical Wald's sequential probability ratio tests (Govindarajulu, 2004; Wald, 1947; Wald and Wolfowitz, 1948; Siegmund, 1985). For the case of three alternative hypotheses, Sobel and Wald (1949) obtained a set of four stopping boundaries for the likelihood-ratio statistic. Their solution was generalized to a larger number of alternatives resulting in the *multi-hypothesis sequential probability ratio tests* (Dragalin et al., 1999; Lai, 2000).

1.3. Our goal—simultaneous testing of individual hypotheses

The focus of this paper is different and more general. We assume that the sequence of sampled units is observed to answer several questions about its parameters. Indeed, once the sampling cost is already spent on each sampled unit, it is natural to use it to answer more than just one question! Therefore, there are d individual hypotheses about parameters $\theta_1, \dots, \theta_d$ of sequentially observed vectors $\mathbf{X}_1, \mathbf{X}_2, \dots$

$$H_0^{(1)} : \theta_1 \in \Theta_{01} \text{ vs } H_A^{(1)} : \theta_1 \in \Theta_{11},$$

$$H_0^{(2)} : \theta_2 \in \Theta_{02} \text{ vs } H_A^{(2)} : \theta_2 \in \Theta_{12},$$

⋮

$$H_0^{(d)} : \theta_d \in \Theta_{0d} \text{ vs } H_A^{(d)} : \theta_d \in \Theta_{1d}.$$

(1)

A few sequential procedures have been proposed for multiple tests similar to (1). One can conduct individual sequential tests of $H_0^{(1)}, \dots, H_0^{(d)}$ and stop after the first rejection or acceptance, as in Jennison and Turnbull (2000, Chapter 15). Hypotheses that are not rejected at this moment will be accepted, conservatively protecting the familywise Type I error rate (FWER-I).

Alternatively, one can assign level α_j and the corresponding Pocock or O'Brien–Fleming rejection boundary to the j th hypothesis. Then one conducts sequential or group sequential tests in a hierarchical manner, as proposed in Glimm et al. (2010), Tamhane et al. (2010), and Maurer et al. (2011) for testing primary, secondary, and possibly tertiary endpoints of a clinical trial. This procedure controls FWER-I at the level $\alpha = \sum \alpha_j$.

A different approach proposed in Tang and Geller (1999) and further developed in Bartroff and Lai (2010) allows to control FWER-I by testing a *closed set* of hypotheses. Along with the individual hypotheses $H_0^{(1)}, \dots, H_0^{(d)}$, this method requires to test all the composite hypotheses consisting of intersections $\cap H_0^{(j_k)}$, $1 \leq j_k \leq d$, $1 \leq k \leq d$. This results in mandatory testing of $(2^d - 1)$ instead of d hypotheses. As shown in Section 4, controlling the overall familywise Type I error rate will then require a rather large expected sample size.

While focusing on the Type I FWER, these procedures do not control the *familywise Type II error rate* and the familywise power. On the other hand, a Type II error, for example, on one of the tests of a safety clinical trial implies a failure to notice a side effect of a treatment, which is important to control.

Notice that sequential tests of single hypotheses are able to control probabilities of both the Type I and Type II errors. Extending this to multiple testing, our goal is to control *both familywise error rates* I and II and to do so at a *low sampling cost* by computing the optimal stopping boundaries and the optimal stopping rule followed by the optimal terminal decisions.

1.4. Approach—extension of non-sequential ideas

To approach this problem, we use the step-up and step-down ideas developed for *non-sequential* multiple comparisons. Detailed overviews of non-sequential methods were given at the NSF-CBMS Conference “New Horizons in Multiple Comparison Procedures” in 2001 (Benjamini et al., 2004), in a 2008 monograph (DudoitLaan and van der Laan, 2008), and at the 2009 Workshop on Modern Multiple Testing by Prof. S. Sarkar.

It was noted that the elementary Bonferroni adjustment for multiple comparisons takes care of the familywise error rate at the expense of power. However, some power can be regained by wisely designed *step-up and step-down methods*, ordering p -values of individual tests, choosing one of the ordered p -values, and proceeding from it into one or the other direction making decisions on the acceptance or rejection of individual null hypotheses (Benjamini and Hochberg, 1995; Hochberg and Tamhane, 1987; Holm, 1979; Sarkar, 1998, 2002; Sidak, 1967; Simes, 1986).

Lehmann and Romano (2005) introduced the *generalized error rate*, which is the probability of making at least r incorrect inferences instead of at least one. This weaker requirement on the error control allows to regain power in studies with a large number of simultaneous inferences, where Bonferroni-type adjustments result in very low significance levels and a great loss of power. The new concept was quickly developed, and multiple comparison methods controlling the generalized error rate were proposed (Romano and Wolf, 2007; Sarkar, 2007; Sarkar and Guo, 2009).

Fixed-sample studies are able to control either the Type I or the Type II error probabilities, but in general, not both. Conversely, Wald’s sequential probability test and subsequent sequential procedures for testing a single hypothesis can be designed to satisfy both the given significance level and the given power. Similarly, sequential testing of multiple hypotheses, considered here, can be set to guarantee a strong control of both familywise Type I and Type II error rates

$$FWER_I = \max_{\mathcal{T} \neq \emptyset} \mathbf{P} \{ \text{reject at least one true null hypothesis} | \mathcal{T} \} = \max_{\mathcal{T} \neq \emptyset} \mathbf{P} \left\{ \bigcup_{j \in \mathcal{T}} \text{reject } H_0^{(j)} | \mathcal{T} \right\}, \quad (2)$$

$$FWER_{II} = \max_{\mathcal{F} \neq \emptyset} \mathbf{P} \{ \text{accept at least one false null hypothesis} | \mathcal{T} \} = \max_{\mathcal{F} \neq \emptyset} \mathbf{P} \left\{ \bigcup_{j \in \mathcal{F}} \text{accept } H_0^{(j)} | \mathcal{T} \right\}, \quad (3)$$

where $\mathcal{T} \subset \{1, \dots, d\}$ is the index set of the true null hypotheses and $\mathcal{F} = \overline{\mathcal{T}}$ is the index set of the false null hypotheses.

In the sequel, the non-sequential Bonferroni, step-down Holm, and step-up Benjamini–Hochberg methods for multiple comparisons are generalized to the sequential setting. Essentially, at every step of the multiple testing scheme, the continue-sampling region is inserted between the acceptance and rejection boundaries in such a way that the resulting sequential procedure controls the intended error rate. Then, the so-called *intersection stopping rule* is introduced that controls both familywise error rates simultaneously.

2. Sequential Bonferroni method

Let us begin with the rigorous formulation of the problem. Suppose a sequence of independent and identically distributed vectors $\mathbf{X}_1, \mathbf{X}_2, \dots \in \mathbf{R}^d$ that are observed as a result of purely sequential or group sequential sampling. Components (X_{i1}, \dots, X_{id}) of the i th random vector may be dependent, and the j th component has a marginal density $f_j(\cdot | \theta_j)$ with respect to a reference measure $\mu_j(\cdot)$. For every $j=1, \dots, d$, measures $\{f_j(\cdot | \theta_j), \theta_j \in \Theta_j\}$ are assumed to be mutually absolutely continuous, and the Kullback–Leibler information numbers

$$K(\theta_j, \theta'_j) = \mathbf{E}_{\theta_j} \log \{f_j(X_j | \theta_j) / f_j(X_j | \theta'_j)\}$$

are strictly positive and finite for $\theta_j \neq \theta'_j$.

Consider a battery of one-sided (right-tail, with a suitable parameterization) tests about parameters $\theta_1, \dots, \theta_d$

$$H_0^{(j)} : \theta_j \leq \theta_0^{(j)} \text{ vs. } H_A^{(j)} : \theta_j \geq \theta_1^{(j)}, \quad j = 1, \dots, d. \quad (4)$$

A stopping rule T is to be found, accompanied with decision rules $\delta_j = \delta_j(X_{1j}, \dots, X_{Tj})$, $j=1, \dots, d$, on the acceptance or rejection of each of the null hypotheses $H_0^{(1)}, \dots, H_0^{(d)}$. This procedure has to control both familywise error rates (2) and (3), i.e., guarantee that

$$FWER_I \leq \alpha \quad \text{and} \quad FWER_{II} \leq \beta, \quad (5)$$

for pre-assigned $\alpha, \beta \in (0, 1)$.

Technically, it is not difficult to satisfy condition (5). Wald’s sequential probability ratio test (SPRT) for the j th hypothesis controls the probabilities of Type I and Type II errors at the given levels α_j and β_j . Choosing $\alpha_j = \alpha/d$ and $\beta_j = \beta/d$, we immediately obtain (5) by the Bonferroni inequality.

This testing procedure is based on log-likelihood ratios

$$A_n^{(j)} = \sum_{i=1}^n \log \frac{f_j(X_{ij} | \theta_1^{(j)})}{f_j(X_{ij} | \theta_0^{(j)})}, \quad j = 1, \dots, d, \quad n = 1, 2, \dots$$

Wald's classical stopping boundaries are

$$a_j = \log \frac{1-\beta_j}{\alpha_j} \quad \text{and} \quad b_j = \log \frac{\beta_j}{1-\alpha_j}. \quad (6)$$

Wald's SPRT for the single j th hypothesis $H_0^{(j)}$ rejects it (i.e., chooses $H_A^{(j)}$) after n observations if $A_n^{(j)} \geq a_j$, accepts it (i.e., chooses $H_0^{(j)}$) if $A_n^{(j)} \leq b_j$, and continues sampling if $A_n^{(j)} \in (b_j, a_j)$.

Assuming that the marginal distributions of $A_1^{(j)}$ have the monotone likelihood ratio property (e.g. Casella and Berger, 2002), the error probabilities are maximized when $\theta_j = \theta_0^{(j)}$ and when $\theta_j = \theta_1^{(j)}$, respectively, for all $j=1, \dots, d$. Then, separately performed SPRT for the j th hypothesis with stopping boundaries (6) controls the probabilities of Type I and Type II errors approximately

$$\mathbf{P}\{A_T^{(j)} \geq a_j | \theta_j = \theta_0^{(j)}\} \approx \alpha_j, \quad \mathbf{P}\{A_T^{(j)} \leq b_j | \theta_j = \theta_1^{(j)}\} \approx \beta_j,$$

where $T_j = \inf\{n : A_n^{(j)} \notin (b_j, a_j)\}$ (Govindarajulu, 1987, 2004; Wald, 1947). This Wald's approximation (Basseville and Nikiforov, 1993) results from ignoring the overshoot over the stopping boundary and assuming that having just crossed the stopping boundary for the first time, the log-likelihood ratio $A_n^{(j)}$ approximately equals to that boundary.

Extending SPRT to the case of multiple hypothesis, $d \geq 2$, continue sampling until all the d tests reach decisions. Define the stopping rule

$$T = \inf \left\{ n : \bigcap_{j=1}^d \{A_n^{(j)} \notin (b_j, a_j)\} \right\}. \quad (8)$$

Lemma 1. For any pairs (a_j, b_j) , the stopping rule defined by (8) is proper.

Proof. Section A.1. \square

Accepting or rejecting the j th null hypothesis at time T depending on whether $A_T^{(j)} \leq b_j$ or $A_T^{(j)} \geq a_j$, we could obtain (approximately, subject to Wald's approximation) a strong control of probabilities of Type I and Type II errors by the Bonferroni inequality

$$\begin{aligned} \mathbf{P}\{\text{at least one Type I error}\} &\leq \sum_{j=1}^d \mathbf{P}\{A_T^{(j)} \geq a_j\} \leq \sum_{j=1}^d \alpha_j = \alpha, \\ \mathbf{P}\{\text{at least one Type II error}\} &\leq \sum_{j=1}^d \mathbf{P}\{A_T^{(j)} \leq b_j\} \leq \sum_{j=1}^d \beta_j = \beta. \end{aligned} \quad (9)$$

However, Wald's approximation is only accurate when the overshoot of A_T over the stopping boundary is negligible. When testing d hypotheses, the corresponding d log-likelihood ratios may cross their respective boundaries at different times. Then, at the stopping time T , when sampling is halted, a number of log-likelihood ratios may be deep inside the stopping region, creating a considerable overshoot. Wald's approximation is no longer accurate for the stopping time T ! It has to be replaced by rigorous statements.

Lemma 2. Let T be a proper stopping time with respect to the vector sequence $(\mathbf{X}_1, \mathbf{X}_2, \dots)$, such that

$$\mathbf{P}\{A_T^{(j)} \in (b, a) | \mathcal{T}\} = 0,$$

for some $j \in \{1, \dots, d\}$, $b < 0 < a$, and any combination of true null hypotheses \mathcal{T} . Consider a test that rejects $H_0^{(j)}$ at time T if and only if $A_T^{(j)} \geq a$. For this test

$$\mathbf{P}\{\text{Type I error on } H_0^{(j)}\} \leq \mathbf{P}\{A_T^{(j)} \geq a | \theta_0^{(j)}\} \leq e^{-a}, \quad (10)$$

$$\mathbf{P}\{\text{Type II error on } H_0^{(j)}\} \leq \mathbf{P}\{A_T^{(j)} \leq b | \theta_1^{(j)}\} \leq e^b. \quad (11)$$

Proof. Section A.2. \square

Avoiding the use of Wald's (and any other) approximation, replace Wald's stopping boundaries (6) for $A_n^{(j)}$ by

$$a_j = -\log \alpha_j \quad \text{and} \quad b_j = \log \beta_j, \quad (12)$$

and use the stopping rule (8). Then, according to Lemma 2, the corresponding test of $H_0^{(j)}$ controls the Type I and Type II error probabilities rigorously at levels $e^{-a_j} = \alpha_j$ and $e^{b_j} = \beta_j$. Therefore, by the Bonferroni arguments in (9), the described multiple testing procedure controls both error rates. The following theorem is then proved.

Theorem 1. Sequential Bonferroni procedure for testing multiple hypotheses (4) with the stopping rule (8), rejection regions $A_T^{(j)} \geq -\log(\alpha/d)$, and acceptance regions $A_T^{(j)} \leq \log(\beta/d)$ controls both error rates at levels $\text{FWER}_I \leq \alpha$ and $\text{FWER}_{II} \leq \beta$.

Further development of the Bonferroni methods and comparison of the associated stopping rules can be found in De and Baron (in press).

3. Step-down and step-up methods

Since Bonferroni methods are based on an inequality that appears rather crude for moderate to large d , controlling both familywise rates I and II can only be done at the expense of a large sample size. For non-sequential statistical inferences, a number of elegant stepwise (step-up and step-down) methods have been proposed, attaining the desired FWER-I and improving over the Bonferroni methods in terms of the required sample size. In this Section, we develop a similar approach for sequential experiments.

Following the Holm method for multiple comparisons (e.g., Neter et al., 1996), we order the tested hypotheses $H_0^{(1)}, \dots, H_0^{(d)}$ according to the significance of the collected evidence against them and set the significance levels for individual tests to be

$$\alpha_1 = \frac{\alpha}{d}, \quad \alpha_2 = \frac{\alpha}{d-1}, \quad \alpha_3 = \frac{\alpha}{d-2}, \dots, \alpha_j = \frac{\alpha}{d+1-j}, \dots, \alpha_d = \alpha.$$

Similarly, in order to control the familywise Type II error rate, choose

$$\beta_j = \frac{\beta}{d+1-j} \quad \text{for } j = 1, \dots, d.$$

Comparing with Bonferroni method, increasing the individual Type I and Type II error probabilities has to cause the familywise error rates to increase. On the other hand, since all the stopping boundaries become *tighter*, this will necessarily reduce the expected sample size $E(T)$ under any combination of true and false hypotheses, \mathcal{T} and \mathcal{F} . Therefore, if both rates can still be controlled at the pre-specified levels α and β , then the resulting stepwise procedure is an improvement of Bonferroni schemes under the given constraints on FWE rates.

In the following algorithm, we combine the stepwise idea for efficient multiple testing with Wald's sequential probability ratio testing of individual hypotheses. The first scheme controls FWER-I, the second scheme controls FWER-II, and the third "intersection" scheme based on them controls *both* familywise error rates.

Scheme 1: Reject one by one, accept all (step-down scheme). Choose the stopping boundaries $a_j = -\log \alpha_j = -\log(\alpha/(d-j+1))$ and arbitrary $b_j < 0$ for $j = 1, \dots, d$.

After observing vectors $\mathbf{X}_1, \dots, \mathbf{X}_n$, order the log-likelihood ratio statistics in their *non-increasing* order

$$A_n^{[1]} \geq A_n^{[2]} \geq \dots \geq A_n^{[d]},$$

and let $H_0^{[j]}$ for $j = 1, \dots, d$ be the corresponding tested null hypotheses arranged in the same order. Proceed according to a *step-down* scheme, from the most significant log-likelihood ratio statistic $A_n^{[1]}$ down to the second most significant statistic, etc.

- Step1.* If $A_n^{[1]} \leq b_1$ then accept all $H_0^{[1]}, \dots, H_0^{[d]}$
 If $A_n^{[1]} \in (b_1, a_1)$ then continue sampling; collect \mathbf{X}_{n+1}
 If $A_n^{[1]} \geq a_1$ then reject $H_0^{[1]}$ and go to $A_n^{[2]}$
Step2. If $A_n^{[2]} \leq b_2$ then accept all $H_0^{[2]}, \dots, H_0^{[d]}$
 If $A_n^{[2]} \in (b_2, a_2)$ then continue sampling; collect \mathbf{X}_{n+1}
 If $A_n^{[2]} \geq a_2$ then reject $H_0^{[2]}$ and go to $A_n^{[3]}$

etc. (Fig. 1).

Sampling continues while at least one ordered log-likelihood ratio $A_n^{[j]}$ belongs to its continue-sampling region (b_j, a_j) . The stopping rule corresponding to this scheme is

$$T_1 = \inf \left\{ n : \bigcap_{j=1}^d A_n^{[j]} \notin (b_j, a_j) \right\}. \quad (13)$$

Theorem 2. The stopping rule T_1 is proper; Scheme1 strongly controls the familywise Type I error rate. That is, for any set \mathcal{T} of true hypotheses

$$P_{\mathcal{T}}\{1 < \infty\} = 1$$

and

$$P_{\mathcal{T}}\{\text{at least one Type I error}\} \leq \alpha. \quad (14)$$

Proof. Section A.3. \square

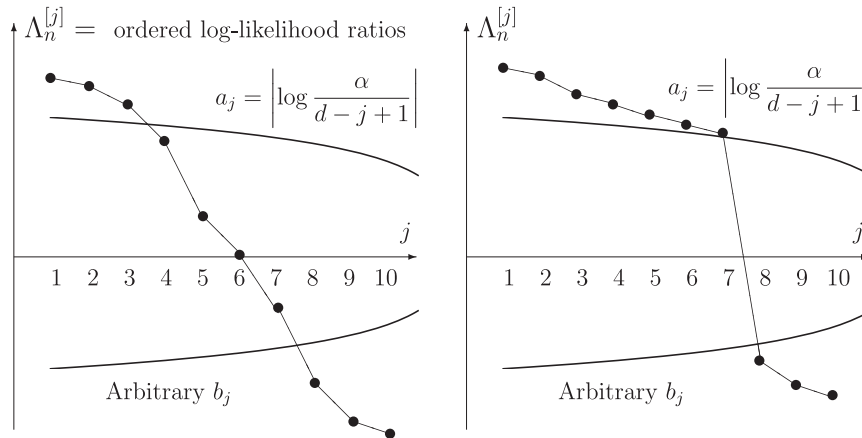


Fig. 1. Example of Scheme 1. On the left, sampling continues. On the right, stop sampling; reject $H_0^{(1)}, \dots, H_0^{(7)}$; accept $H_0^{(8)}, \dots, H_0^{(10)}$.

Theorem 2 holds for arbitrary $b_j < 0$, thus it also shows that the rejection boundary a_j alone controls FWER-I.

Symmetrically to Scheme 1, introduce the following Scheme 2 which controls FWER-II through the choice of acceptance boundary b_j alone.

Scheme 2: *Accept one by one, reject all (step-up scheme).* For this scheme, choose the stopping boundaries $b_j = \log \beta_j = \log(\beta/(d-j+1))$ and arbitrary $a_j > 0$ for $j = 1, \dots, d$.

After observing $\mathbf{X}_1, \dots, \mathbf{X}_n$, order the log-likelihood ratio statistics in their *non-decreasing* order

$$\Lambda_n^{(1)} \leq \Lambda_n^{(2)} \leq \dots \leq \Lambda_n^{(d)},$$

and let $H_0^{(j)}$ for $j = 1, \dots, d$ be the corresponding tested null hypotheses arranged in the same order. Proceed according to the *step-up scheme*, from the least significant log-likelihood ratio to the second least significant, etc.

- Step1.** If $\Lambda_n^{(1)} \geq a_1$ then reject all $H_0^{(1)}, \dots, H_0^{(d)}$
 If $\Lambda_n^{(1)} \in (b_1, a_1)$ then continue sampling; collect \mathbf{X}_{n+1}
 If $\Lambda_n^{(1)} \leq b_1$ then accept $H_0^{(1)}$ and go to $\Lambda_n^{(2)}$
- Step2.** If $\Lambda_n^{(2)} \geq a_2$ then reject all $H_0^{(2)}, \dots, H_0^{(d)}$
 If $\Lambda_n^{(2)} \in (b_2, a_2)$ then continue sampling; collect \mathbf{X}_{n+1}
 If $\Lambda_n^{(2)} \leq b_2$ then accept $H_0^{(2)}$ and go to $\Lambda_n^{(3)}$

etc.

According to this scheme, the stopping rule is, similarly to (13)

$$T_2 = \inf \left\{ n : \bigcap_{j=1}^d \Lambda_n^{(j)} \notin (b_j, a_j) \right\}. \quad (15)$$

Theorem 3. The stopping rule T_2 is proper; Scheme 2 strongly controls the familywise Type II error rate. That is, for set \mathcal{T} of true hypotheses

$$\mathbf{P}_{\mathcal{T}}\{T_2 < \infty\} = 1$$

and

$$\mathbf{P}_{\mathcal{T}}\{\text{at least one Type II error}\} \leq \beta. \quad (16)$$

Proof. Section A.4. \square

Notice that Scheme 2 can be equivalently formulated in terms of the non-increasing statistics $\Lambda_n^{[j]}$ instead of non-decreasing $\Lambda_n^{(j)}$. Indeed, one rearranges the log-likelihood ratios from their non-increasing ordering to their nonincreasing ordering by simply reverting the order, i.e.,

$$\Lambda_n^{[j]} = \Lambda_n^{[d-j+1]}.$$

Rearranging the boundary values accordingly, i.e., replacing $b_j = \log(\beta/d-j+1)$ with $b_j = \log(\beta/j)$, one obtains Scheme 2 in terms of $\Lambda_n^{[j]}$ instead of $\Lambda_n^{(j)}$. In short

$$\Lambda_n^{[j]} \leq \log \frac{\beta}{d-j+1} \Leftrightarrow \Lambda_n^{(d-j+1)} \leq \log \frac{\beta}{j} \Leftrightarrow \Lambda_n^{[j]} \leq \log \frac{\beta}{j}.$$

This is illustrated in Fig. 2.

Comparing the logic of Schemes 1 and 2, we see that the step-down Scheme 1 starts with the most significant log-likelihood ratio statistic $\Lambda_n^{[1]}$ and carefully/conservatively rejects one null hypothesis at a time. It focuses on controlling Type I errors of wrong rejection and results in controlling the overall FWER-I.

On the contrary, the step-up Scheme 2 starts with the least significant statistic and conservatively accepts one null at a time, controlling FWER-II.

It is actually possible to combine both schemes and to develop a sequential procedure controlling both familywise error rates as follows.

The Intersection Scheme: Combining Schemes 1 and 2, define the *intersection stopping rule*

$$T^* = \inf \left\{ n : \bigcap_{j=1}^d \Lambda_n^{[j]} \notin (b_j^*, a_j^*) \right\}, \quad a_j^* = -\log \frac{\alpha}{d-j+1}, \quad b_j^* = \log \frac{\beta}{j}. \quad (17)$$

Acceptance and rejection regions in this case are simply *intersections* of rejection and acceptance regions for $\Lambda_n^{[j]}$ for Schemes 1 and 2, as long as the lower boundary in Scheme 1 and the upper boundary in Scheme 2 are inside the interval $[\log \beta, -\log \alpha]$.

According to the Intersection Scheme, reject the hypothesis $H_0^{[j]}$ corresponding to the j th ordered log-likelihood ratio $\Lambda_n^{[j]}$ if $\Lambda_n^{[j]} \geq a_j^*$ and accept it if $\Lambda_n^{[j]} \leq b_j^*$ (Fig. 3).

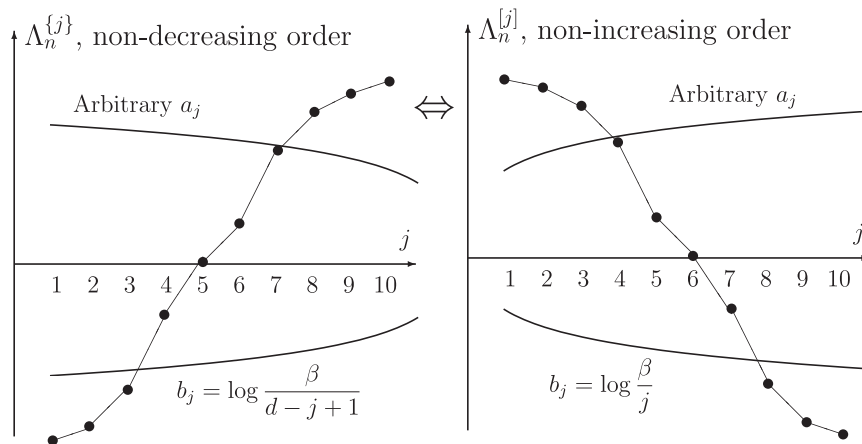


Fig. 2. Scheme 2 in terms of log-likelihood ratios in the non-decreasing order (left), and equivalently, in the non-increasing order (right).

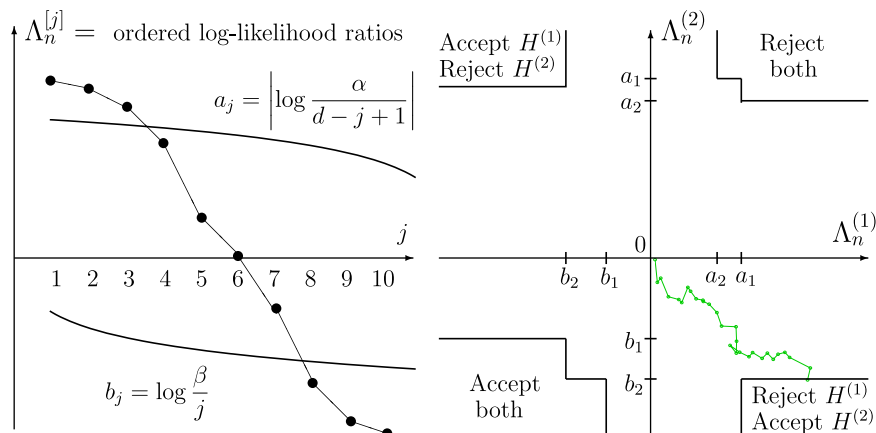


Fig. 3. The Intersection Scheme. Stopping boundaries for the ordered log-likelihood ratios for $d=10$ tests, snap shot after n observations (left). Stopping boundaries for the unordered log-likelihood ratios for $d=2$ tests, the entire path until the stopping time T^* (right).

Being a special case of Scheme 1 and Scheme 2 at the same time, the Intersection Scheme controls both FWER-I and FWER-II.

Theorem 4. *The stopping rule T^* is proper; the Intersection Scheme strongly controls both familywise Type I and Type II error rates. That is, for any set \mathcal{T} of true hypotheses*

$$P_{\mathcal{T}}\{T^* < \infty\} = 1,$$

$$P_{\mathcal{T}}\{\text{at least one Type I error}\} \leq \alpha,$$

$$P_{\mathcal{T}}\{\text{at least one Type II error}\} \leq \beta.$$

Proof. Section A.5. \square

4. Comparison of multiple testing schemes

Performance of the proposed schemes is evaluated and compared with the Holm–Bonferroni non-sequential multiple testing procedure and the multistage step-down procedure of Bartroff and Lai (2010).

First, consider testing three null hypotheses

$$H_0^{(1)} : \theta_1 = 0 \text{ vs } H_A^{(1)} : \theta_1 = 0.5,$$

$$H_0^{(2)} : \theta_2 = 0 \text{ vs } H_A^{(2)} : \theta_2 = 0.5,$$

$$H_0^{(3)} : \theta_3 = 0.5 \text{ vs } H_A^{(3)} : \theta_3 = 0.75,$$

based on a sequence of random vectors $\mathbf{X}_1, \mathbf{X}_2, \dots$, where $\mathbf{X}_i = (X_{i1}, X_{i2}, X_{i3})$, $X_{i1} \sim \text{Normal}(\theta_1, 1)$, $X_{i2} \sim \text{Normal}(\theta_2, 1)$, and $X_{i3} \sim \text{Bernoulli}(\theta_3)$, which is the scenario considered in Bartroff and Lai (2010).

For each combination of null hypotheses and alternatives, $N=55,000$ random sequences are simulated (omitting the combinations $H_A^{(1)} \cap H_0^{(2)} \cap H_0^{(3)}$ and $H_A^{(1)} \cap H_0^{(2)} \cap H_A^{(3)}$, because of their interchangeability with $H_0^{(1)} \cap H_A^{(2)} \cap H_0^{(3)}$ and $H_0^{(1)} \cap H_A^{(2)} \cap H_A^{(3)}$, respectively, therefore yielding exactly the same performance of the considered stopping and decision rules). Based on them, the familywise Type I and Type II error rates are estimated along with the expected stopping time and the standard error of our estimator of the expected stopping time. Scheme 1 and the Intersection Scheme are set to control $\text{FWER}_I \leq 0.05$. Also, Scheme 2 and the Intersection Scheme are set to control $\text{FWER}_{II} \leq 0.10$.

Results are shown in Tables 1 and 2. It can be seen that under each combination of true null hypotheses, the expected sample sizes of sequential Bonferroni, step-down Scheme 1, step-up Scheme 2, and the Intersection Scheme are of the same order; however, they are about a half of the expected size of the closed testing procedure or the fixed-sample size required by the non-sequential Holm–Bonferroni procedure.

Table 1

Comparison of the proposed sequential schemes with nonsequential Holm and sequential closed testing procedures.

Procedure	Expected st. time	Standard error	FWER-I	FWER-II
Under $H_0^{(1)}, H_0^{(2)}, H_0^{(3)}; \theta_1 = 0, \theta_2 = 0, \theta = 0.5$				
Bonferroni	46.8	0.10	0.022	–
Scheme 1	35.4	0.09	0.026	–
Scheme 2	35.7	0.08	0.055	–
Intersection Scheme	37.0	0.09	0.023	–
Nonsequential Holm	105.0	–	0.040	–
Multistage step-down	104.6	–	0.048	–
Under $H_0^{(1)}, H_0^{(2)}, H_A^{(3)}; \theta_1 = 0, \theta_2 = 0, \theta = 0.75$				
Bonferroni	49.3	0.10	0.013	0.010
Scheme 1	41.2	0.09	0.021	0.034
Scheme 2	41.3	0.09	0.036	0.027
Intersection Scheme	45.7	0.09	0.019	0.028
Nonsequential Holm	105.0	–	0.044	–
Multistage step-down	98.3	–	0.042	–
Under $H_0^{(1)}, H_A^{(2)}, H_0^{(3)}; \theta_1 = 0, \theta_2 = 0.65, \theta = 0.5$				
Bonferroni	41.8	0.09	0.017	0.001
Scheme 1	33.9	0.07	0.026	0.005
Scheme 2	36.1	0.08	0.043	0.004
Intersection Scheme	38.7	0.08	0.024	0.004
Nonsequential Holm	105.0	–	0.044	–
Multistage step-down	96.9	–	0.049	–

Table 2

Comparison of the proposed sequential schemes with nonsequential Holm and sequential closed testing procedures (continued).

Procedure	Expected st. time	Standard error	FWER-I	FWER-II
Under $H_0^{(1)}, H_A^{(2)}, H_A^{(3)}; \theta_1 = 0, \theta_2 = 0.5, \theta = 0.75$				
Bonferroni	51.6	0.10	0.006	0.021
Scheme 1	43.8	0.09	0.016	0.055
Scheme 2	43.1	0.09	0.019	0.034
Intersection Scheme	48.0	0.10	0.016	0.030
Nonsequential Holm	105.0	–	0.043	–
Multistage step-down	92.3	–	0.032	–
Under $H_A^{(1)}, H_A^{(2)}, H_0^{(3)}; \theta_1 = 0.5, \theta_2 = 0.5, \theta = 0.5$				
Bonferroni	51.6	0.10	0.007	0.022
Scheme 1	44.3	0.09	0.016	0.059
Scheme 2	42.3	0.09	0.020	0.037
Intersection Scheme	47.8	0.10	0.016	0.032
Nonsequential Holm	105.0	–	0.038	–
Multistage step-down	93.1	–	0.027	–
Under $H_A^{(1)}, H_A^{(2)}, H_A^{(3)}; \theta_1 = 0.5, \theta_2 = 0.5, \theta = 0.75$				
Bonferroni	53.7	0.10	–	0.029
Scheme 1	42.1	0.08	–	0.073
Scheme 2	42.2	0.09	–	0.038
Intersection Scheme	43.9	0.09	–	0.031
Nonsequential Holm	105.0	–	–	–
Multistage step-down	86.1	–	–	–

Table 3

Sequential Bonferroni scheme versus sequential stepwise procedures.

Procedure	Expected st. time	Standard error	FWER-I	FWER-II
Scenario: μ_1, \dots, μ_{10} are parameters of Normal($\cdot, 1$) distributions of $X_{i,1}, \dots, X_{i,10}$; p_1, \dots, p_{10} are parameters of Bernoulli(\cdot) distributions of $X_{i,11}, \dots, X_{i,20}$; test $\mu_j = 0$ vs 0.5 and $p_j = 0.5$ vs 0.75; odd-numbered null hypotheses are true				
Bonferroni	118.4	0.13	0.002	0.003
Intersection Scheme	111.1	0.13	0.004	0.006
Scenario: μ_1, \dots, μ_{10} are parameters of Normal($\cdot, 1$) distributions of $X_{i,1}, \dots, X_{i,10}$; p_1, \dots, p_{10} are parameters of Bernoulli(\cdot) distributions of $X_{i,11}, \dots, X_{i,20}$; test $\mu_j = 0$ vs 0.5 and $p_j = 0.5$ vs 0.75; null hypotheses $H_0^{(1-9, 11-19)}$ are true				
Bonferroni	114.6	0.13	0.004	0.001
Intersection Scheme	96.1	0.12	0.007	0.006
Scenario: μ_1, \dots, μ_{10} are parameters of Normal($\cdot, 1$) distributions of $X_{i,1}, \dots, X_{i,10}$; p_1, \dots, p_{10} are parameters of Bernoulli(\cdot) distributions of $X_{i,11}, \dots, X_{i,20}$; test $\mu_j = 0$ vs 0.5 and $p_j = 0.5$ vs 0.75; null hypotheses $H_0^{(10,20)}$ are true				
Bonferroni	121.6	0.13	0.001	0.005
Intersection Scheme	101.5	0.12	0.004	0.008
Scenario: μ_1, \dots, μ_{10} are parameters of Poisson(\cdot) distributions of $X_{i,1}, \dots, X_{i,10}$; p_1, \dots, p_{10} are parameters of Bernoulli(\cdot) distributions of $X_{i,11}, \dots, X_{i,20}$; test $\mu_j = 5$ vs 6 and $p_j = 0.5$ vs 0.75; odd-numbered null hypotheses are true				
Bonferroni	148.4	0.18	0.002	0.003
Intersection Scheme	138.7	0.17	0.003	0.006
Scenario: μ_1, \dots, μ_5 are parameters of Normal($\cdot, 1$) distributions of $X_{i,1}, \dots, X_{i,5}$; with $\text{Cov}(X_{i,j}, X'_{i,j}) = 0.5$ for $\forall j \neq j'$; p_1, \dots, p_5 are parameters of Bernoulli(\cdot) distributions of $X_{i,6}, \dots, X_{i,10}$; test $\mu_j = 0$ vs 0.5 and $p_j = 0.5$ vs 0.75; odd-numbered null hypotheses are true				
Bonferroni	91.2	0.12	0.002	0.004
Intersection Scheme	85.3	0.12	0.004	0.007

Advantage of stepwise schemes versus the sequential Bonferroni scheme is more significant for a larger number of tests. This is seen in Table 3 for different multiple testing problems. Reduction in the expected sample size ranges from 6% when 50% of null hypotheses are true to 16.5% when most null hypotheses are either true or false. Results are also based on $N=55,000$ simulated sequences for each considered scenario.

The last example in Table 3 deals with correlated components of the observed random vectors. Indeed, all the results in this paper make no assumption about the joint distribution of $(X_{i,1}, \dots, X_{i,d})$ for each $i = 1, 2, \dots$. When components are correlated, positively or negatively, the expected sample size of each procedure should reduce.

We also notice that results of Theorems 1–4 are based on Bonferroni-type inequalities and corollaries from them. For a large number of tests, Bonferroni inequality tends to be rather crude, and therefore, the familywise error rates, guaranteed

by Theorems 1–4, are often satisfied with a rather wide margin. This certainly leaves room for the improvement of the proposed sequential testing schemes!

Acknowledgments

The authors are grateful to the Editor Professor N. Balakrishnan, the Associate Editor, and to the anonymous referee for deep, insightful, and encouraging comments that helped us tremendously. Research of both authors is funded by the National Science Foundation Grant DMS 1007775. Research of the second author is partially supported by the National Security Agency Grant H98230-11-1-0147. This funding is greatly appreciated.

Appendix A. Proofs

A.1. Proof of Lemma 1

By the weak law of large numbers

$$\mathbf{P}\{A_n^{(j)} \in (b_j, a_j) | H_0^{(j)}\} \rightarrow 0 \quad \text{and} \quad \mathbf{P}\{A_n^{(j)} \in (b_j, a_j) | H_A^{(j)}\} \rightarrow 0 \text{ as } n \rightarrow \infty,$$

because the non-zero expected values of $A_n^{(j)}$ are guaranteed by the assumptions of Section 2 on Kullback–Leibler information numbers. Then

$$\mathbf{P}_T\{T > n\} \leq \sum_{j=1}^d \mathbf{P}_T\{A_n^{(j)} \in (b_j, a_j)\} \rightarrow 0,$$

and $\mathbf{P}_T\{T = \infty\} \leq \mathbf{P}_T\{T > n\} \rightarrow 0$, therefore, $\mathbf{P}_T\{T = \infty\} = 0$.

A.2. Proof of Lemma 2

The proof is based on Doob's maximal inequality for submartingales (e.g., Williams, 1991, Section 14.6; Kuo, 2006, Section 4.5).

For any $j=1, \dots, d$ and $n=1, 2, \dots$, the likelihood ratio

$$\lambda_n^{(j)} = \exp\{A_n^{(j)}\} = \prod_{i=1}^n \frac{f_j(X_{ij} | \theta_1^{(j)})}{f_j(X_{ij} | \theta_0^{(j)})}$$

is a non-negative martingale under $\theta_0^{(j)}$ with respect to the filtration generated by (X_{1j}, X_{2j}, \dots) . Then, by Doob's inequality, for all $N \geq 1$

$$\mathbf{P}\left\{\max_{1 \leq n \leq N} \lambda_n^{(j)} \geq e^a \middle| \theta_0^{(j)}\right\} \leq e^{-a} \mathbf{E}\{\lambda_N^{(j)} | \theta_0^{(j)}\} = e^{-a}$$

and

$$\mathbf{P}(\text{Type I error on } H_0^{(j)}) \leq \mathbf{P}\{A_T^{(j)} \geq a | \theta_0^{(j)}\} \leq \mathbf{P}\left\{\max_{1 \leq n \leq N} A_n^{(j)} \geq a \middle| \theta_0^{(j)}\right\} + \mathbf{P}\{T > N | \theta_0^{(j)}\} \leq e^{-a} + \mathbf{P}\{T > N | \theta_0^{(j)}\}. \quad (18)$$

Taking the limit as $N \rightarrow \infty$ proves inequality (10) because $\mathbf{P}\{T > N | \theta_0^{(j)}\} \rightarrow 0$ since the stopping time T is proper.

To prove inequality (11), we notice that $1/\lambda_n^{(j)} = \exp\{-A_n^{(j)}\}$ is a non-negative martingale under $\theta_1^{(j)}$. Applying Doob's inequality, we obtain

$$\mathbf{P}\left\{\min_{1 \leq n \leq N} A_n^{(j)} \leq b \middle| \theta_1^{(j)}\right\} = \mathbf{P}\left\{\max_{1 \leq n \leq N} 1/\lambda_n^{(j)} \geq e^{-b} \middle| \theta_1^{(j)}\right\} \leq e^b,$$

and the arguments similar to (18) conclude the proof.

A.3. Proof of Theorem 2

1. The stopping time T_1 is almost surely bounded by

$$T' = \inf \left\{ n : \bigcap_{j=1}^d \{A_n^{(j)} \notin (\min b_j, \max a_j)\} \right\}. \quad (19)$$

Since T' is proper by Lemma 1, so is T_1 .

2. The proof of control of FWER-I borrows ideas from the classical derivation of the experimentwise error rate of the non-sequential Holm procedure, extending the arguments to sequential tests.

Let $\mathcal{T} \subset \{1, \dots, d\}$ be the index set of true null hypotheses, and $\mathcal{F} = \overline{\mathcal{T}}$ be the index set of the false null hypotheses, with cardinalities $|\mathcal{T}|$ and $|\mathcal{F}|$. Then, arrange the log-likelihood ratios at the stopping time T_1 in their non-increasing order, $A_{T_1}^{(1)} \geq \dots \geq A_{T_1}^{(d)}$ and let m be the smallest index of the ordered log-likelihood ratio that corresponds to a true hypothesis. In other words, if $H_0^{(j)}$ denotes the null hypothesis that is being tested by the log-likelihood ratio $A_{T_1}^{(j)}$ for $j=1, \dots, d$, then m is such that all $H_0^{(j)}$ are false for $j < m$ whereas $H_0^{(m)}$ is true. Thus, there are at least $(m-1)$ false hypotheses, so that $m-1 \leq |\mathcal{F}| = d - |\mathcal{T}|$.

No Type I error can be made on false hypotheses $H_0^{(1)}, \dots, H_0^{(m-1)}$. If the Type I error is not made on $H_0^{(m)}$ either, then there is no Type I error at all because according to Scheme 1, acceptance of $H_0^{(m)}$ implies automatic acceptance of the remaining hypotheses $H_0^{(m+1)}, \dots, H_0^{(d)}$.

Therefore

$$\mathbf{P}_{\mathcal{T}}\{\text{at least one Type I error}\} = \mathbf{P}_{\mathcal{T}}\{A_{T_1}^{(m)} \geq a_m\} \leq \mathbf{P}_{\mathcal{T}}\{A_{T_1}^{(m)} \geq a_{d-|\mathcal{T}|+1}\} = \mathbf{P}_{\mathcal{T}}\left\{\max_{j \in \mathcal{T}} A_{T_1}^{(j)} \geq a_{d-|\mathcal{T}|+1}\right\} \leq \sum_{j \in \mathcal{T}} \mathbf{P}_{\mathcal{T}}\{A_{T_1}^{(j)} \geq a_{d-|\mathcal{T}|+1}\}.$$

Recall that rejection boundaries for Scheme 1 are chosen as $a_j = -\log \alpha_j$, where $\alpha_j = \alpha/(d+1-j)$. Therefore, $a_{d-|\mathcal{T}|+1} = -\log(\alpha/|\mathcal{T}|)$, and by Lemma 2

$$\mathbf{P}_{H_0^{(j)}}\{A_{T_1}^{(j)} \geq a_{d-|\mathcal{T}|+1}\} \leq \exp\{-a_{d-|\mathcal{T}|+1}\} = \alpha/|\mathcal{T}|.$$

Finally, we have

$$\mathbf{P}_{\mathcal{T}}\{\text{at least one Type I error}\} \leq \sum_{j \in \mathcal{T}} \alpha/|\mathcal{T}| = \alpha.$$

A.4. Main steps of the Proof of Theorem 3

Ideas of Section A.3 are now translated to Scheme 2 and control of FWER-II. Let us outline the main steps of the proof, especially because control of FWER-II has not been studied in sequential multiple testing, to the best of our knowledge.

1. Similarly to T_1 , the stopping time T_2 is also bounded by the proper stopping rule (19), and therefore, it is also proper.

2. Following Scheme 2, arrange $A_{T_2}^{(j)}$ in their non-decreasing order, $A_{T_2}^{(1)} \leq \dots \leq A_{T_2}^{(d)}$. Then let ℓ be the smallest index of the ordered log-likelihood ratio that corresponds to a false null hypothesis, so that all $H_0^{(1)}, \dots, H_0^{(\ell-1)}$, corresponding to $A_{T_2}^{(1)}, \dots, A_{T_2}^{(\ell-1)}$, are true but $H_0^{(\ell)}$ is false. The number of true hypotheses is then at least $(\ell-1)$, so that $\ell \leq |\mathcal{T}| + 1 = d - |\mathcal{F}| + 1$, where $|\mathcal{T}|$ and $|\mathcal{F}|$ are the numbers of true and false null hypotheses.

If any Type II error is made during Scheme 2, then it has to occur on $H_0^{(\ell)}$, because its (correct) rejection leads to the automatic (correct) rejection of the remaining hypotheses $H_0^{(\ell+1)}, \dots, H_0^{(d)}$, according to the scheme.

Therefore, applying (11) to $A_{T_2}^{(j)}$, we obtain

$$FWER_{II} = \mathbf{P}_{\mathcal{T}}\{A_{T_2}^{(\ell)} \leq b_{\ell}\} \leq \mathbf{P}_{\mathcal{T}}\{A_{T_2}^{(\ell)} \leq b_{d-|\mathcal{F}|+1}\} = \mathbf{P}_{\mathcal{T}}\left\{\min_{j \in \mathcal{F}} A_{T_2}^{(j)} \leq b_{d-|\mathcal{F}|+1}\right\} \leq \sum_{j \in \mathcal{F}} \mathbf{P}_{\mathcal{T}}\{A_{T_2}^{(j)} \leq \log(\beta/|\mathcal{F}|)\} \leq \sum_{j \in \mathcal{F}} \beta/|\mathcal{F}| = \beta.$$

A.5. Proof of Theorem 4

The Intersection Scheme satisfies all the conditions of Theorem 2, therefore, the stopping time T^* is proper, and the scheme controls $FWER_I \leq \alpha$. At the same time, it satisfies Theorem 3, and therefore, it controls $FWER_{II} \leq \beta$.

References

- Armitage, P., 1950. Sequential analysis with more than two alternative hypotheses, and its relation to discriminant function analysis. *Journal of the Royal Statistical Society Series B* 12, 137–144.
- Baillie, D.H., 1987. Multivariate acceptance sampling—some applications to defence procurement. *The Statistician* 36, 465–478.
- Bartroff, J., Lai, T.-L., 2010. Multistage tests of multiple hypotheses. *Communications in Statistics—Theory and Methods* 39, 1597–1607.
- Basseville, M., Nikiforov, I.V., 1993. *Detection of Abrupt Changes: Theory and Application*. PTR, Prentice-Hall, Inc.
- Baum, C.W., Veeravalli, V.V., 1994. A sequential procedure for multihypothesis testing. *IEEE Transactions on Information Theory* 40, 1994–2007.
- Benjamini, Y., Hochberg, Y., 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society* 57 (1), 289–300.
- Benjamini, Y., Bretz, F., Sarkar, S. (Eds.), 2004. *Recent Developments in Multiple Comparison Procedures*, IMS Lecture Notes—Monograph Series, Beachwood, OH.
- Betensky, R.A., 1996. An O'Brien–Fleming sequential trial for comparing three treatments. *Annals of Statistics* 24 (4), 1765–1791.
- Casella, G., Berger, R.L., 2002. *Statistical Inference*. Duxbury Press, Belmont, CA.
- De, S., Baron, M. Sequential Bonferroni methods for multiple hypothesis testing with strong control of familywise error rates I and II. *Sequential Analysis*, in press.
- Dragalin, V.P., Tartakovsky, A.G., Veeravalli, V.V., 1999. Multihypothesis sequential probability ratio tests. Part I: asymptotic optimality. *IEEE Transactions on Information Theory* 45 (7), 2448–2461.
- Dudoit, S., Shaffer, J.P., Boldrick, J.C., 2003. Multiple hypothesis testing in microarray experiment. *Statistical Science* 18, 71–103.

- Dudoit, S., van der Laan, M.J., 2008. Multiple Testing Procedures with Applications to Genomics. Springer, New York.
- Edwards, D., 1987. Extended-Paulson sequential selection. *Annals of Statistics* 15 (1), 449–455.
- Edwards, D.G., Hsu, J.C., 1983. Multiple comparisons with the best treatment. *Journal of the American Statistical Association* 78, 965–971.
- Ghosh, M., Mukhopadhyay, N., Sen, P.K., 1997. Sequential Estimation. Wiley, New York.
- Glimm, E., Maurer, W., Bretz, F., 2010. Hierarchical testing of multiple endpoints in group-sequential trials. *Statistics in Medicine* 29, 219–228.
- Govindarajulu, Z., 1987. The Sequential Statistical Analysis of Hypothesis Testing, Point and Interval Estimation, and Decision Theory. American Sciences Press, Columbus, OH.
- Govindarajulu, Z., 2004. Sequential Statistics. World Scientific Publishing, Co., Singapore.
- Hamilton, D.C., Lesperance, M.L., 1991. A consulting problem involving bivariate acceptance sampling by variables. *Canadian Journal of Statistics* 19, 109–117.
- Hochberg, Y., Tamhane, A.C., 1987. Multiple Comparison Procedures. Wiley, New York.
- Holm, S., 1979. A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics* 6, 65–70.
- Hughes, M.D., 1993. Stopping guidelines for clinical trials with multiple treatments. *Statistics in Medicine* 12, 901–915.
- Jennison, C., Turnbull, B.W., 1993. Group sequential tests for bivariate response: interim analyses of clinical trials with both efficacy and safety endpoints. *Biometrics* 49, 741–752.
- Jennison, C., Turnbull, B.W., 2000. Group Sequential Methods with Applications to Clinical Trials. Chapman & Hall, Boca Raton, FL.
- Jennison, C., Johnstone, I.M., Turnbull, B.W., 1982. Asymptotically optimal procedures for sequential adaptive selection of the best of several normal means. In: Gupta, S.S., Berger, J.O. (Eds.), *Statistical Decision Theory and Related Topics III*, vol. 2. , Academic Press, New York, pp. 55–86.
- Kuo, H.H., 2006. Introduction to Stochastic Integration. Springer, New York.
- Lai, T.L., 2000. Sequential multiple hypothesis testing and efficient fault detection-isolation in stochastic systems. *IEEE Transactions on Information Theory* 46 (2), 595–608.
- Lehmann, E.L., Romano, J.P., 2005. Generalizations of the familywise error rate. *Annals of Statistics* 33, 1138–1154.
- Maurer, W., Glimm, E., Bretz, F., 2011. Multiple and repeated testing of primary, coprimary, and secondary hypotheses. *Statistics in Biopharmaceutical Research* 3 (2), 336–352.
- Neter, J., Kutner, M., Nachtsheim, C., Wasserman, W., 1996. Applied Linear Statistical Models, fourth ed. McGraw-Hill.
- Novikov, A., 2009. Optimal sequential multiple hypothesis tests. *Kybernetika* 45 (2), 309–330.
- O'Brien, P.C., 1984. Procedures for comparing samples with multiple endpoints. *Biometrics* 40, 1079–1087.
- O'Brien, P.C., Fleming, T.R., 1979. A multiple testing procedure for clinical trials. *Biometrika* 35, 549–556.
- Paulson, E., 1962. A sequential procedure for comparing several experimental categories with a standard or control. *Annals of Mathematical Statistics* 33, 438–443.
- Paulson, E., 1964. A sequential procedure for selecting the population with the largest mean from k normal populations. *Annals of Mathematical Statistics* 35, 174–180.
- Pocock, S.J., Geller, N.L., Tsiatis, A.A., 1987. The analysis of multiple endpoints in clinical trials. *Biometrics* 43, 487–498.
- Romano, J.P., Wolf, M., 2007. Control of generalized error rates in multiple testing. *Annals of Statistics* 35, 1378–1408.
- Sarkar, S.K., 1998. Some probability inequalities for ordered mtp_2 random variables: a proof of the simes conjecture. *Annals of Statistics* 26 (2), 494–504.
- Sarkar, S.K., 2002. Some results on false discovery rate in stepwise multiple testing procedures. *Annals of Statistics* 30 (1), 239–257.
- Sarkar, S.K., 2007. Step-up procedures controlling generalized FWER and generalized FDR. *Annals of Statistics* 35, 2405–2420.
- Sarkar, S.K., Guo, W., 2009. On a generalized false discovery rate. *Annals of Statistics* 37 (3), 1545–1565.
- Shaffer, J.P., 1995. Multiple hypothesis testing. *Annual Review of Psychology* 46, 561–584.
- Sidak, Z., 1967. Rectangular confidence regions for the means of multivariate normal distributions. *Journal of the American Statistical Association* 62, 626–633.
- Siegmund, D., 1985. Sequential Analysis: Tests and Confidence Intervals. Springer-Verlag, New York.
- Siegmund, D., 1993. A sequential clinical trial for comparing three treatments. *Annals of Statistics* 21, 464–483.
- Simes, R.J., 1986. An improved Bonferroni procedure for multiple tests of significance. *Biometrika* 73, 751–754.
- Simons, G., 1967. Lower bounds for the average sample number of sequential multihypothesis tests. *Annals of Mathematical Statistics* 38 (5), 1343–1364.
- Sobel, M., Wald, A., 1949. A sequential decision procedure for choosing one of three hypotheses concerning the unknown mean of a normal distribution. *Annals of Mathematical Statistics* 20 (4), 502–522.
- Tamhane, A.C., Mehta, C.R., Liu, L., 2010. Testing a primary and a secondary endpoint in a group sequential design. *Biometrics* 66, 1174–1184.
- Tang, D.-I., Geller, N.L., 1999. Closed testing procedures for group sequential clinical trials with multiple endpoints. *Biometrics* 55, 1188–1192.
- Tartakovsky, A.G., Veeravalli, V.V., 2004. Change-point detection in multichannel and distributed systems with applications. In: Mukhopadhyay, N., Datta, S., Chattopadhyay, S. (Eds.), *Applications of Sequential Methodologies*, Marcel Dekker, Inc., New York, pp. 339–370.
- Tartakovsky, A.G., Li, X.R., Yaralov, G., 2003. Sequential detection of targets in multichannel systems. *IEEE Transactions on Information Theory* 49 (2), 425–445.
- Wald, A., 1947. Sequential Analysis. Wiley, New York.
- Wald, A., Wolfowitz, J., 1948. Optimal character of the sequential probability ratio test. *Annals of Mathematical Statistics* 19, 326–339.
- Wilcox, R.R., 2004. Extension of Hochberg's two-stage multiple comparison method. In: Mukhopadhyay, N., Datta, S., Chattopadhyay, S. (Eds.), *Applications of Sequential Methodologies*, Marcel Dekker, Inc., New York, pp. 371–380.
- Williams, D., 1991. Probability with Martingales. Cambridge University Press, Cambridge, UK.
- Zacks, S., 2009. Stage-Wise Adaptive Designs. Wiley, Hoboken, NJ.