

Rapporto n. 195

***Power Estimation for Multiple Co-Primary Endpoints:
a comparison among conservative solutions***

A. Lucadamo , N. Accoto, D. De Martini

Ottobre 2010

Dipartimento di Metodi Quantitativi per le Scienze Economiche ed Aziendali

Università degli Studi di Milano Bicocca

Via Bicocca degli Arcimboldi 8 - 20126 Milano - Italia

Tel +39/02/64483102/3 - Fax +39/2/64483105

Segreteria di redazione: Andrea Bertolini

Power Estimation for Multiple Co-Primary Endpoints: a comparison among conservative solutions

A. Lucadamo¹, N. Accoto², D. De Martini³

¹: TEDASS Dep.- Università del Sannio - alucadam@unina.it

²: Dep. of Decision Sciences - Università Comm. L.Bocconi - nadia.accoto@unibocconi.it

³: DIMEQUANT Dep.- Università Milano-Bicocca - daniele.demartini@unimib.it

SUMMARY. The problem of estimating the power of the multivariate Intersection Union test (IUT) is studied. Four classical parametric solutions and a bootstrap non-parametric one, providing statistical lower bounds (i.e. one directional confidence intervals) for the power, are considered. The performances of these techniques in several bivariate IUT settings are compared through a simulation study. All solutions are biased, since their actual coverage probabilities are higher than the nominal one. The bootstrap solution shows the smallest bias, and the variability of its estimates is the lowest. Moreover, the bias of the bootstrap solution reduces faster than those of the other techniques when the pilot sample size, or the correlation, or the rate between the two noncentrality parameters increases. Also, the nonparametric bootstrap solution can be improved by calibration, with a considerable bias reduction.

KEYWORDS. Sample size estimation; conservative approach; bootstrap solution.

1 Introduction

In the last decade multiple endpoint statistical problems have received increasing attention. On many occasions, indeed, clinical research aims to demonstrate the efficacy of a new drug on more than one endpoint. Related techniques for statistical analysis (Sankoh et al., 2003, Dmitrienko et al., 2003) and for sample size computation (Dilba et al. 2006, Senn and Bretz, 2007) have then been developed.

For certain disorders, indeed quite a few, a new treatment is required by regulatory agencies to demonstrate efficacy on multiple co-primary endpoints, all significant at the one-sided 2.5% level. The adequate statistical test for treating this latter problem is the so-called (Gleser, 1973, Lehmann, 1952) Intersection Union Test (IUT). Berger (1982) has proposed the use of IUT for acceptance sampling problems. Recently, many authors provided interesting contributes on IUT and its applications to biomedical statistics. Among others, Chuang-Stein et al. (2007) proposed an approach based on the notion of controlling the maximum false positive rate over the restricted null space in order to use a higher significance level to test individual endpoints; Offen et al. (2007), who formed a team of experts from the Pharmaceutical Research and Manufacturers of America, provided medical and statistical solutions for multiple co-primary endpoints.

As regards power and sample size computation for the IUT, a conservative approach based on a mathematical sharp lower bound for the power function can be found in Eaton and Muirhead (2007); these authors obtained the lower bound induced under no correlation and showed other interesting results. Song (2009) provided sample size formulas for IU testing of rate differences in non-inferiority trials. Considering a mathematical lower bound for the power, Yeo and Qu (2009) adopted the plug-in pointwise sample size estimation for the IUT, but they did not take the variability of pilot data into account.

In biomedical statistics further experiments are usually planned on the basis of the results of previous studies available in literature. In particular, phase III clinical trials are planned referring to phase I and II results. This kind of sample size computation technique falls under Sample Size Estimation methodology (SSE), which is often adopted in many different applied research contexts (see, for example Johnston et al., 2009, Eng, 2003, Devane et al., 2004). In these situations, if one forgets

to take the variability of pilot data into account, wrong experimental planning may occur. Conservative approaches to sample size estimation (CSSE) have, therefore, been proposed (Chuang-Stein, 2006, Wang et al., 2006, De Martini, 2010).

It is worth noting that the core of CSSE is the estimation of the unknown true power of the test, i.e. the power function computed in correspondence to the unknown true value(s) of the parameter(s). One sided confidence intervals, i.e. statistical lower bounds, for the true power are then needed.

In this paper the problem of estimating the true power (or, simply, the power) of the IUT, and consequently its sample size, on the basis of a set of pilot data is considered, accounting for the variability of these latter. The aim is, therefore, to provide tools to compute statistical lower bounds for the power of the IUT. It is anticipated that technical difficulties will be due to the multidimensional nature of the parameter, which is the argument of the power function, and to the bias of IUT.

In Section 2, the theoretical framework of IUT, together with its power, are recalled. In Section 3, some different approaches for estimating the power of IUT are introduced. A comparison of the performances of the different techniques is shown in Section 4 and in Section 5, where the best estimation technique is refined through calibration. In Section 6 a numerical example of CSSE for IUT is presented and in Section 7 the conclusions are reported. Computational details follow in appendix (i.e. Section 8).

2 Theoretical framework of IU test and power

Let $\mathbf{X} = (X_1, \dots, X_\ell)$ be the observations of the ℓ endpoints for an individual receiving the new drug and $\mathbf{Y} = (Y_1, \dots, Y_\ell)$ be those for an individual who received the control drug. Furthermore, assume that $\mathbf{X} \sim N_\ell(\mu_X, \Sigma_X)$ and that $\mathbf{Y} \sim N_\ell(\mu_Y, \Sigma_Y)$, where Σ_\bullet are the covariance matrices.

The correlation coefficients are the off diagonal elements of the matrices Σ_\bullet , i.e. $\Sigma_{\bullet,ij} = \rho_{\bullet,ij}$ with $i \neq j$. Without loss of generality we can assume that the diagonal elements of Σ_\bullet are all equal to 1. Being $\delta = (\delta_1, \dots, \delta_\ell) = \mu_X - \mu_Y$ the vector of the effect sizes, the statistical hypotheses for the non-inferiority multiple test are:

$$\begin{cases} H_0 : \delta_j \leq 0 & \text{for at least one } j \\ H_1 : \delta_j > 0 & \text{for all } j \end{cases} \quad (1)$$

In practice, the multivariate null hypothesis is rejected if, and only if, all univariate null hypotheses are rejected.

Now, consider drawing a sample of m individuals from each group, that is \mathbf{X}_i and \mathbf{Y}_i , $i = 1, \dots, m$, are i.i.d. with common distribution function $N_\ell(\mu_X, \Sigma_X)$ and $N_\ell(\mu_Y, \Sigma_Y)$, respectively. Then, compute the vector of the test statistics $\mathbf{T}_m = \sqrt{\frac{m}{2}}(\bar{\mathbf{X}}_m - \bar{\mathbf{Y}}_m) = (T_{1,m}, \dots, T_{\ell,m})$, where $\bar{\mathbf{X}}_m = \sum_{i=1,m} \mathbf{X}_i/m$, and $\bar{\mathbf{Y}}_m$ is analogous. Moreover, $\alpha \in (0, 1)$ is the type I error probability, $z_{1-\alpha} = \Phi^{-1}(1 - \alpha)$ and Φ is the cumulative distribution function of the standard Normal law.

In accordance with (1), the so-called Intersection-Union Test (IUT) introduced by Berger (1982) is the following:

$$\Psi_m(\mathbf{T}_m) = \begin{cases} 0 & \text{if } T_{j,m} \leq z_{1-\alpha} \text{ for at least one } j \\ 1 & \text{if } T_{j,m} > z_{1-\alpha} \text{ for all } j = 1, \dots, \ell \end{cases} \quad (2)$$

Berger (1982) showed that if all the ℓ univariate tests are α -level tests, then the global test too is α -level. In fact, under the null hypothesis the *Sup* of the power is α and is achieved when $\delta_{\bar{i}} = 0$ and δ_j tends to $+\infty$ with $j = 1, \dots, \ell$, $j \neq \bar{i}$.

Nevertheless, it is very important to note that the IUT is biased, because under the alternative hypothesis the power can be lower than α . Indeed, when $\delta_j = \epsilon > 0 \forall j$, we fall under H_1 and if ϵ is “very small” then the power of each univariate test is $\approx \alpha$; consequently, when $\rho_{\bullet,ij} = 0$, the power of the IUT turns out to be $\approx \alpha^\ell$, which is lower than α since $\alpha < 1$.

Through simple algebra we obtain that $\mathbf{T}_m = (T_{1,m}, \dots, T_{\ell,m}) \sim N_\ell(\sqrt{m/2}\delta, \Sigma_T)$, where $\Sigma_{T,ii} = 1$ and $\Sigma_{T,ij} = (\rho_{x,ij} + \rho_{y,ij})/2$, $1 \leq i, j \leq \ell$. Then, the power of (2) is

$$E[\Psi_m(\mathbf{T}_m)] = P(T_{1,m} > z_{1-\alpha}, \dots, T_{\ell,m} > z_{1-\alpha}) \quad (3)$$

and this can be computed as a function of δ , Σ_X , Σ_Y , m and α .

We begin studying power estimation techniques for IUT under the simplest non-trivial situation, that is the bivariate case (i.e. $\ell = 2$) with equal dependence structure in the treatment and control groups. This implies $\rho_{x,12} = \rho_{y,12} = \rho$ so that $\Sigma_X = \Sigma_Y = \Sigma$. Hence, the power function in (3) simplifies to:

$$\pi(\delta_1, \delta_2, \rho, m, \alpha) = P(T_{1,m} > z_{1-\alpha}, T_{2,m} > z_{1-\alpha}) \quad (4)$$

and this can be computed as a function of the effect sizes δ_1 and δ_2 , of the correlation

ρ , of m and α , becoming:

$$\pi(\delta_1, \delta_2, \rho, m, \alpha) = 1 - \Phi_{\delta_1 \sqrt{m/2}, 1}(z_{1-\alpha}) - \Phi_{\delta_2 \sqrt{m/2}, 1}(z_{1-\alpha}) + \Phi_{\delta \sqrt{m/2}, \Sigma}(z_{1-\alpha}, z_{1-\alpha}) \quad (5)$$

Given α , and being $1 - \beta$ the power to achieve, the ideal sample size is:

$$M_I = \min\{m | \pi(\delta_1, \delta_2, \rho, m, \alpha) > 1 - \beta\} \quad (6)$$

Note that ρ plays the role of a nuisance parameter in the power function and, consequently, in sample size computation. In some papers tables showing how M_I varies with different ρ s are presented (Chuang-Stein et al., 2007, Yeo and Qu, 2009), and differences are not negligible.

In practice, δ_1 , δ_2 and ρ are unknown and so is M_I . If pilot samples are available, then M_I can be estimated and the conservative approach is suggested. So, let us suppose two samples of size n are drawn from the treatment and the control group, respectively, i.e. \mathbf{X}_i , $i = 1, \dots, n$, i.i.d., $\mathbf{X}_i \sim N_2(\mu_X, \Sigma)$, and \mathbf{Y}_i , $i = 1, \dots, n$, i.i.d., $\mathbf{Y}_i \sim N_2(\mu_Y, \Sigma)$. The challenge is now to estimate $\pi(\delta_1, \delta_2, \rho, m, \alpha)$ given m and α and, so, indirectly to estimate M_I given $1 - \beta$.

3 Some different approaches for estimating IUT power

3.1 The parametric approach and related techniques

Being ρ a nuisance parameter in this testing context, the power does not depend primarily on ρ . Consequently, the conservative estimation approach is here applied to the vector (δ_1, δ_2) and its lower bounds are plugged-into the power function for obtaining conservative estimates of the true power. As regards the correlation coefficient, two solutions are considered: the first one consists in plugging-into the power function the pointwise estimate of ρ , say r_n , using the pooled estimator proposed by Donner and Rosner (1980); the second one adopts the mathematical lower bound for the power proposed by Eaton and Muirhead (2007), which considers $\rho = 0$.

Specifically, considering a confidence region D_n^γ for (δ_1, δ_2) where γ is the amount of conservativeness, i.e. $P((\delta_1, \delta_2) \in D_n^\gamma) = \gamma$, the lower bound of the power is given by $\min_{D_n^\gamma} \{\pi(\delta_1, \delta_2, \bullet, m, \alpha)\}$, where \bullet stands for the generic solution adopted for substituting the unknown value of ρ in the power function.

Remark 1. Unusefulness of IUT inversion. The logical direct way for conservatively estimating (δ_1, δ_2) is by inverting the IUT at a level γ (see also Wilson, 1927). In practice, when the point estimate $\mathbf{d}_n = \bar{\mathbf{X}}_n - \bar{\mathbf{Y}}_n = (d_{n,1}, d_{n,2})$ is observed, the confidence region is given by the points $(\bar{\delta}_1, \bar{\delta}_2)$ for which the IUT with null hypothesis $H_0 : \{\delta_1 \leq \bar{\delta}_1 \text{ or } \delta_2 \leq \bar{\delta}_2\}$ is non significant. This region turns out to be $D_n^\gamma = R^2 - \{(\delta_1, \delta_2) < (d_{n,1} - z_\gamma \sqrt{2/n}, d_{n,2} - z_\gamma \sqrt{2/n})\}$, i.e. the entire plane without an open square in the low-left part. Consequently, we have that $\min_{D_n^\gamma} \{\pi(\delta_1, \delta_2, \bullet, m, \alpha)\} = 0$, so that this region is not useful for conservatively estimating the power.

Two different approaches for D_n^γ are here adopted: the first one consists in the classical elliptical confidence region for (δ_1, δ_2) (see for example Morrison, 2005); the second is based on two simultaneous lower bounds for δ_1 and δ_2 , according to Anderson (1958) and Roy & Bose (1953). In the following of this section both are briefly recalled.

3.1.1 The elliptical confidence region

This region is based on the joint distribution of the sample difference mean vector with the sample covariance matrix of the bivariate normal distribution, given by the pooled within-groups covariance estimators, i.e. the distribution of (d_n, \mathbf{S}_n) , where $d_n = \bar{\mathbf{X}}_n - \bar{\mathbf{Y}}_n \sim N_2(\delta, \frac{2}{n}\Sigma)$ and $\mathbf{S}_n = [(n-1)S_{n,1} + (n-1)S_{n,2}]/2(n-1)$. Being $\tau^2 = \frac{n}{2}(\bar{\mathbf{X}}_n - \bar{\mathbf{Y}}_n)' \mathbf{S}_n^{-1} (\bar{\mathbf{X}}_n - \bar{\mathbf{Y}}_n)$, in consequence of the Hotelling and Wishart distribution properties, it is obtained that $\tau^2 \sim T_2^2(2(n-1); \zeta^2)$, where $\zeta^2 = \frac{n}{2} \delta' \Sigma^{-1} \delta$ is the noncentrality parameter. If $\delta = 0$, then $\zeta^2 = 0$ and $\tau^2 \sim T_2^2(2(n-1))$, that is $(2n-3)\tau^2/4(n-1) \sim F(2, 2n-3)$. Hence, the boundary of the elliptical region for δ , centered at d_n , with $100(\gamma)$ per cent (approximated) confidence, is given by the following equation: $\frac{n}{2}(\delta - d_n)' \mathbf{S}_n^{-1} (\delta - d_n) = 4(n-1)F_{2, 2n-3}^{-1}(1-\gamma)/(2n-3)$. Note that the confidence region has, in this case, only approximated confidence level, since the diagonal elements of Σ (i.e. the variances σ_i^2) are supposed to be known and equal to 1.

We will refer to the techniques stemming from the elliptical shape of D_n^γ as ELLP and ELLM, when ρ is either estimated pointwise or set, according to the mathematical *Minoration*, equal to zero, respectively.

Remark 2. Inversion of UIT. It can be noted that the elliptical region corresponds to the inversion of the Union Intersection Test (UIT), not the IUT in study. Indeed, this confidence region is given by the points $(\bar{\delta}_1, \bar{\delta}_2)$ for which the UIT with null hypothesis $H_0 : (\delta_1, \delta_2) = (\bar{\delta}_1, \bar{\delta}_2)$ (versus $H_1 : (\delta_1, \delta_2) \neq (\bar{\delta}_1, \bar{\delta}_2)$) is non significant.

3.1.2 The simultaneous bounds region

The simultaneous bounds are obtained through simultaneous one-directional confidence intervals. Recall first that bi-directional simultaneous confidence intervals for the mean of a bivariate Normal distribution are: $(\bar{X}_i - K_{\frac{1-\gamma}{2}} \sqrt{\frac{s_i^2}{n}}; \bar{X}_i + K_{\frac{1-\gamma}{2}} \sqrt{\frac{s_i^2}{n}})$, $i = 1, 2$, where s_i^2 are the estimated variance and $K_{\frac{1-\gamma}{2}} = \sqrt{2(n-1)F_{2,n-2}^{-1}(1-\gamma)/(n-2)}$. Then, supposing $\sigma_i^2 = 1$ to be known, the pivotal distribution simplifies to a χ_2^2 , and the $100(\gamma)$ per cent conservative estimate for the effect size δ is $(d_{n,1} - \sqrt{c_{2,2(1-\gamma)}/n}, d_{n,2} - \sqrt{c_{2,2(1-\gamma)}/n})$, where $c_{2,2(1-\gamma)}$ is such that $P(\chi_2^2 \leq c_{2,2(1-\gamma)}) = 2\gamma - 1$.

We will refer to these techniques as SIMP and SIMM, when ρ is either estimated pointwise or set equal to zero, respectively.

3.1.3 Parametric computational algorithm

As regards the calculations, $\min_{D_n^\gamma} \{\pi(\delta_1, \delta_2, r_n, m, \alpha)\}$ is obtained through an algorithm detecting the level curve of the power (say the iso-power curve), which is tangent to the low-left part of the (elliptical or rectangular) confidence region. Note that iso-power curves behave almost like hyperboles (see also Section 8 for related computational details). In practice, the problem consists in detecting the curve tangent to the ellipse/open-rectangle centered in $(d_{n,1}, d_{n,2})$ among the family of “hyperboles”. Since we use the elliptical confidence region for computing just a lower bound for (δ_1, δ_2) , note also that the real coverage probability of the ellipse is $\gamma_r = \frac{\gamma+1}{2}$, with $\gamma \in [0, 1]$.

3.2 The nonparametric approach and the bootstrap technique

Sometimes, quite severe technical difficulties arise within parametric frameworks. On these occasions, nonparametric methods might be useful to solve parametric problems. As has just been shown in Subsection 3.1, the task of providing bounds

for the power of the IUT is, actually, problematic. It is widely known that Efron's bootstrap is a highly versatile nonparametric method. Recently, a nonparametric bootstrap technique for estimating the power of statistical tests (even conservatively) has been presented by De Martini (2011), where applications to sample size estimation for the Wilcoxon rank-sum test are shown. We, therefore, adopt this bootstrap technique for conservatively estimating the power of the IUT and we here provide a brief reminder.

Let us denote the bivariate empirical distribution functions of the treatment group and of the control group $F_{T,n}$ and $F_{C,n}$, respectively. Note that these functions contain information both on the shifts (δ_1, δ_2) and on the correlation ρ . Then, when the power is viewed as a functional of the distributions (i.e., $\pi = \pi(N_2(\mu_X, \Sigma), N_2(\mu_Y, \Sigma), m, \alpha)$), the simple bootstrap plug-in estimate of the true power is $\pi(F_{T,n}, F_{C,n}, m, \alpha)$.

Now, in order to provide a lower bound for the true power, draw two samples of size n from $F_{T,n}$ and $F_{C,n}$ respectively, and let $F_{T,n}^*$ and $F_{C,n}^*$ be the empirical distribution functions so obtained. Hence, $\pi(F_{T,n}^*, F_{C,n}^*, m, \alpha)$ is the bootstrap estimator of the true power. Finally, denote with $\pi^\gamma(F_{T,n}, F_{C,n}, m, \alpha)$ the $1 - \gamma$ p-tile of this latter estimator: this is the (approximated) γ -lower bound for the true power: $P(\pi^\gamma(F_{T,n}, F_{C,n}, m, \alpha) \leq \pi(N_2(\mu_X, \Sigma), N_2(\mu_Y, \Sigma), m, \alpha)) \simeq \gamma$. A theoretical justification of this bootstrap solution, which will be denoted by BO, can be found in De Martini (2011).

4 Simulation Study

In this Section we evaluate and compare the performances of the five different estimation techniques for the power of IUT introduced in the previous Section (viz. ELLP, ELLM, SIMP, SIMM and BO).

4.1 Design of the Study

In order to evaluate the performances of different techniques we vary δ_1 , $k = \delta_1/\delta_2$, ρ and also the ideal sample size M_I in such a way that the power is always 90%. We consider a small sample situation, with $M_I = 60$, and a larger one, with $M_I = 180$. For each M_I we consider $k = 1, 1.5, 2$ and $\rho = 0.2, 0.4, 0.5, 0.6, 0.8$. However, we do not

evaluate estimation performances in all the 15 possible cases (i.e. 3 k s \times 5 ρ s): we just consider the 8 couples (k, ρ) that are reported in Table 1, together with the corresponding values of δ_1 for each M_I (see Appendix A for computational details).

As regards the size n of pilot samples, recent works on CSSE (Wang et al., 2006, De Martini, 2010) indicate that, in order to obtain sufficiently accurate power estimates, n should be of the same order of magnitude as M_I . So, for every setting here considered, we evaluate the performances of our techniques with pilot samples of size around M_I . Specifically, we set $n = 2M_I/3, 4M_I/3$. Hence, the total number of experimental points is 32 ($2 M_{IS} \times 8 \delta_{1S} \times 2 ns$). In Table 2 the 32 Scenarios so obtained are defined in detail.

For every Scenario we simulate the behavior of our techniques by generating $B_0 = 5000$ samples from the bivariate normal distributions of the treatment and the control groups. We, thus, obtain 5000 conservative estimates of the power for each one of the 5 techniques, for each conservative level $\gamma = 0.5, \dots, 0.99$, with step 0.01. The resulting estimates are evaluated by considering the correctness of conservative levels and the variability of the estimates.

As regards the former point, for every γ the bias is intended to be the difference between the actual coverage probability (ACP) and the nominal one (NCP, viz. γ); hence, the average bias is computed.

As regards the variability of the estimates, the weighted average of the means of the absolute standardized differences between the γ -conservative estimators π^γ and the true power π , namely $I2$, is adopted in accordance with De Martini (2011). In practice, being $D_\gamma = (\pi^\gamma - \pi)/(1 - \pi)$ if $\pi^\gamma > \pi$, and $D_\gamma = (\pi - \pi^\gamma)/\pi$ otherwise, we have that $I2$ is the weighted average of $E[D_\gamma]$ over the set of γ s considered, that is:

$$I2 = \sum_{i=0}^4 w_{.5+.1i} E[D_{.5+.1i}]$$

Since $\pi^{50\%}$ is merely pointwise (and not conservative) and $\gamma = 90\%$ may be too severe a conservativeness, in practice the most used conservative levels are around 60–80%. Consequently, we used $w_{50\%} = w_{90\%} = 0.125$ and $w_{60\%} = w_{70\%} = w_{80\%} = 0.25$.

4.2 Results

Our five techniques present a clear bias, since they all are, in most Scenarios, too conservative.

As regards parametric techniques, those based on the so-called mathematical lower bounds (viz. ELLM and SIMM) provide results similar to the respective ones obtained with the pointwise estimate of ρ (the former approach is a little more conservative than the latter). Since all parametric techniques were too conservative, estimating ρ pointwise provides less biased (i.e. better) results. Moreover, SIMP lower bounds are less biased than ELLP ones in all Scenarios, leaving SIMP the best parametric performer.

In general, the bias of the nonparametric BO is somewhat lower than that given by SIMP in all settings. The nonparametric BO should, therefore, be preferred.

The bias of all techniques decreases as n , ρ , and k increase. To show these behaviors, we focus on small sample settings (i.e. $M_I = 60$); moreover, Scenario #1 represents the basic setting, whereas Scenarios #2, #5 and #21 show the behavior of power estimation techniques as n , ρ and k increase, respectively. In particular, for the above settings the ACPs of the different techniques against NCP are plotted in Figures 1-4. The average bias of the five techniques under these Scenarios, together with the rate of improvement, and the I_2 index are reported in Table 3.

In Scenario #1 the biases of the different techniques are quite similar, around 15%, and the lowest is BO (13.84%). As n increases passing from $2M_I/3 = 40$ to $4M_I/3 = 80$, all techniques improve, albeit marginally (see Figures 1-2); nevertheless, the improvement of BO is the highest (12.3%). As ρ passes from 0.2 to 0.8, the bias of all techniques clearly decreases, and the average bias of BO remains somewhat lower than the others (see Figures 1 and 3); once again the improvement of BO is the highest (59.3%). Finally, the highest bias reductions of all five techniques can be observed comparing Figures 1 and 4, i.e. when k passes from 1 to 2: in practice, the bias of BO disappears (average bias 0.34%, i.e. 97.5% improvement).

As far as the variability of the estimates is concerned, under Scenarios 1, 2, 5 and 21 BO presented lower I_2 s than those of the parametric techniques in all cases but one. Moreover, the values shown by BO are very similar to those observed in estimating the power of the widely used Wilcoxon rank sum test, with the same ns and M_I (see De Martini, 2011, Table 2 and 3).

Focusing now on BO (i.e. the best conservative power estimator among those here considered), the most interesting results from the practical point of view are those with $n = 2M_I/3$. In fact, BO performances under the eight Scenarios with

$M_I = 60$ are reported in Table 4 (2nd and 6th columns). It can be noted that the highest biases are observed with $k = 1$ and small values of ρ (viz. Scenarios 1 and 3), whereas bias is very small when $k = 2$ (viz. Scenarios 21, 23 and 25). On the contrary, the variability index $I2$ shows small differences among these 8 settings.

When n passes to 80, the average bias decreases a little, where $I2$ decreases significantly. For example, in Scenarios 4 and 14 (to be compared with Scenarios 3 and 13) the average bias is 8.62% and 2.72%, where the values of $I2$ are 0.3771 and 0.3862, respectively.

Finally, biases and variabilities observed under large sample settings (i.e. $M_I = 180$) are similar to the corresponding ones with small samples: the bias is a little larger, where $I2$ is a little smaller. For example, Scenarios 7 and 27 (to be compared with Scenarios 1 and 21) provide average biases of 14.13% and 1.43%, and $I2$ values of 0.5081, 0.5006, respectively.

4.3 Discussion

The clear improvement shown by the techniques increasing ρ or k is due to the fact that the IUT becomes univariate when ρ tends to 1 or when k diverges. In these cases, indeed, the test is unbiased, as are power estimation techniques (BO only approximately). When n diverges power estimators theoretically converge, and in fact the variability index $I2$ decreases; the bias also decreases, just a little.

In parametric techniques, the bias is mainly due to the different shapes of iso-power and confidence region curves. The former are quite similar to hyperboles, whereas the latter (i.e. D_γ) are either elliptical or rectangular. This implies that a certain amount of probability mass lies between D_γ and the tangent iso-power curve. Consequently, the resulting γ -conservative power estimate, i.e. the level of the iso-power curve tangent to D_γ , is more conservative than its actual nominal coverage (i.e. γ).

The cause of BO bias is mainly the biasedness of IUT. In particular, even if the two n sized samples fit $N_2(\mu_X, \Sigma)$ and $N_2(\mu_Y, \Sigma)$, respectively, well, many re-sampled n sized samples can fall close to each other, generating a couple of empirical distribution still under H_1 , but in reality close to H_0 . These re-samples inherit the bias of IUT and carry it into bootstrap estimation. For large M_I s analogous estimation problems do exist.

5 Improving bootstrap performances through calibration

Calibration is usually adopted for correcting the bias of asymptotic confidence intervals. The NCP (viz. γ) of asymptotic confidence intervals is achieved when the sample size n tends to ∞ . In practice, with finite ns the ACP is different from γ . Nevertheless, there exists a correct coverage, say γ_c , which provides the confidence interval with the desired NCP of γ .

Practical calibration at first makes use of the available sample to estimate γ_c . Once the estimate $\hat{\gamma}_c$ is calculated, it is adopted to compute a confidence interval with nominal coverage $\hat{\gamma}_c$. The ACP of the confidence interval so obtained is, then, closer to γ than that of the simple confidence interval with NCP = γ .

For an introduction to calibration see Efron and Tibshirani (1993). Here, we adopt calibration in the context of IUT power estimation. Moreover, since the bias of BO is lower than those of parametric techniques, we apply calibration to BO.

From the results of the simulation study in the previous Section, it is worth noting that the ACP of BO presents a parabolic shape. So, we evaluate here if a parabolic model for ACP fits the bias well. It is natural to assume that there is no bias with extreme coverage probabilities, i.e. setting the bias at zero when $NCP = 0$ or 1 . The following model for the ACP of BO is, hence, derived:

$$ACP = aNCP^2 + (1 - a)NCP \quad (7)$$

The values of a are, then, computed with the classic least squares method for all the 32 scenarios of Section 4 (a subset of a values is reported in Table 4). The model (7) fits the ACP data very well: the correlation between the observed ACP and those provided by the model is, approximately, 100% in all 32 scenarios.

The parabolic behavior of ACP can be exploited in our calibration: γ_c should not be computed separately for every single γ , but it can be provided in a general way by inverting the parabola in (7). In practice, the bias is at first estimated on the basis of the pilot sample by estimating the parameter a ; once \hat{a} is obtained, $\hat{\gamma}_c$ is computed through the inversion of (7) at the given γ (i.e. $\hat{\gamma}_c = (\hat{a} - 1 + \sqrt{(1 - \hat{a})^2 + 4\hat{a}\gamma})/2\hat{a}$); finally, the γ -conservative estimate of the power is computed by adopting $\hat{\gamma}_c$.

Remark 3. Bias smoothing. It is worth noting that the use of (7) for modelling the bias can be viewed as a kind of bias smoothing, which also allows considerable

saving of computational time.

5.1 Simulation Study

The aim is to evaluate the performances of this BO calibrated technique (namely BOC). Although computing power estimates with calibration for a single practical case can be completed in a few minutes, to perform a simulation study becomes, computationally, quite heavy. In order to evaluate the improvement given by calibration we, therefore, consider just 8 scenarios among the 32 of the previous Section, i.e. only those with $M_T = 60$ and $n = 40$. The number of power estimates is also decreased to $B_0 = 1000$, and γ varies from 0.5 to 0.98 with a step of 0.02.

The results are most favorable: the global average of the absolute values of mean biases provided by BO (i.e. 4.86%) is reduced by calibration to 2.11%, with an improvement rate higher than 50%. The improvement rate is around 70% for the two highest average biases in particular (viz. those under Scenarios 1 and 3). Detailed results are reported in Table 4. The bias of BOC can also be observed in Figures 1, 3 and 4, where the ACP of BOC is shown.

It should be noted that calibration can sometimes invert the sign of the bias, since it tends to balance the bias itself. Moreover, when BO bias is small that of BOC can be a little higher (viz. Scenarios 21 e 23), but on these occasions calibration is not needed. These biases of BOC may be reduced by increasing B_0 and the Monte Carlo parameters of bootstrap calibration. It should also be remembered that calibration can be iterated to obtain further reductions of the bias (Hall and Martin, 1988). Finally, the variability I_2 index of BOC is substantially equal to that of BO (only slightly smaller).

Hence, calibration improves BO bias significantly, but not estimation variability.

6 An example of conservative sample size estimation

The problem of estimating the sample size for a clinical study on sleep disorders is studied. In the phase II trial two groups of $n = 48$ patients are recruited; these same undergo the drug and placebo treatment, respectively. Two clinical parameters concerning the quality of sleeping are recorded before and after the treatment period,

and the post-pre differences represent the clinical variables of statistical interest. In practice, $\ell = 2$ and $\mathbf{X}_i, \mathbf{Y}_i, i = 1, \dots, 48$ are observed.

In order to show the efficacy of the treatment drug the statistical significance of the differences between groups of both variables should be obtained. Consequently, the IUT should be used.

The values of the standardized differences between the means are $d_{48,1} = 0.827$ and $d_{48,2} = 0.553$, and those of the correlation coefficients result $\rho_X = 0.358$ and $\rho_Y = 0.396$. Since the research team considers these results to be scientifically relevant, the phase III trial is launched. The sample size can, then, be computed on the basis of preliminary data and the conservative approach is adopted.

In the light of the above results, the bootstrap calibrated technique (BOC) is used. The conservative estimated power curves are shown in Figure 5, where four conservative γ levels are considered, i.e. $\gamma = 50\%$ (viz. pointwise approach), 60%, 70% and 80%. In order to achieve a power of 90% the conservative estimates of the sample size are: 65, 78, 98 and 129.

Moreover, the simple bootstrap conservative power estimates (BO) are computed, and are also reported in Figure 5. With the same conservative γ levels and power to be achieved, the resulting sample size estimates are: 76, 95, 117 and 161. It can be noted that these estimates are higher than BOC ones, in keeping with simulation results of Sections 4 and 5.

7 Conclusions

The parametric techniques we considered have provided poor performances, both when ρ is estimated pointwise and when it is set, conservatively, equal to 0. In our opinion, instead of studying other estimation or bounding solutions for ρ , it would be better to focus on confidence regions of the same (or similar) shape as the iso-power IUT curves. Although such curves are defined by complicated equations, providing analytical solutions of this kind is an interesting challenge for the future.

Conversely, a general nonparametric solution for power estimation was available, and it can be applied to univariate or multivariate tests. This technique is based on bootstrap, it has already provided satisfactory results when applied to the Wilcoxon rank-sum test, and it can also be useful when applied to complex

parametric situations. Here, this bootstrap technique has been applied to the IU test and yielded favorable performances; it presented a certain amount of coverage probability bias merely in some circumstances. Nevertheless, its performances can be improved through calibration, obtaining a considerable bias reduction. Finally, bootstrap power estimation of IUT can be applied in the same way even when different correlations within groups or deviations from normality of data distributions arise.

8 Appendix A: computational details

8.1 Computation of effect sizes and of iso-power curves

As shown in section 2, the power of the test depends on δ_1 , δ_2 , ρ , m and α . For every value of ρ , M_I and α there are infinite couples (δ_1, δ_2) providing a given power $1 - \beta$. For this reason, we built the iso-power curves starting from the couples where $\delta_1 = \delta_2$ and used an algorithm with subsequent approximations which recalls the multivariate normal distribution. The curves so obtained turn out to behave almost like hyperboles and the maximum error with respect to the chosen value of $1 - \beta$ is 0.000001. Finally, a bisection method is applied to compute the values of δ_1 and δ_2 for which the ratio $k = \delta_1/\delta_2=1, 1.5, 2$.

8.2 Computation of the bounds for the parametric approaches

In order to build elliptical confidence regions some functions have been implemented in R package (R Development Core Team, 2005), but many of them do not correspond to the inversion of the UIT. For this reason, we preferred to build our elliptical confidence region ex-novo. As stated in section 3.1.1, the boundary of the elliptical region for $\delta=(\delta_1, \delta_2)$ is given by the equation:

$$\frac{n}{2}(\delta - d_n)'S^{-1}(\delta - d_n) = 4(n - 1)F_{2,2n-3}^{-1}(1 - \gamma)/(2n - 3)$$

which can be written as:

$$F_{2,2n-3}^{-1}(1 - \gamma) = \frac{n(2n - 3)}{8(n - 1)}[\delta - d_n]'S_n^{-1}[\delta - d_n]$$

where $d_n = (d_{n,1}, d_{n,2})$. If we consider $\frac{n(2n-3)}{8(n-1)} = \lambda$ and $[\delta - d_n] = \psi$ we have:

$$F_{2,2n-3}^{-1}(1-\gamma) = \lambda \begin{pmatrix} \psi_1 & \psi_2 \end{pmatrix} \begin{bmatrix} S_{11}^{-1} & S_{12}^{-1} \\ S_{21}^{-1} & S_{22}^{-1} \end{bmatrix} \begin{pmatrix} \psi_1 \\ \psi_2 \end{pmatrix}$$

where $S_{..}^{-1}$ are the elements of the inverse of the covariance matrix. After appropriate algebraic and matricial operations, we obtained this second degree equation:

$$\lambda S_{22}^{-1} \psi_2^2 + 2\lambda S_{12}^{-1} \psi_1 \psi_2 + \lambda S_{11}^{-1} \psi_1^2 - F_{2,2n-3}^{-1}(1-\gamma) = 0$$

whose solution provides the boundary of the ellipse centered in $(d_{n,1}, d_{n,2})$. The power lower bound can be found by looking for the iso-power curve that is tangent to the low-left part of the elliptical confidence region. Then, we considered the point of the ellipse that has the minimum value on the x-axis, we calculated the power at this point and considered the corresponding iso-power curve. If the curve did not intersect the ellipse at other points, then this point represented the lower bound for the power, otherwise we moved on the ellipse and we repeated the operation until we found the tangent curve.

For the simultaneous confidence intervals the computation of the lower bound of the power was obtained considering simply the power estimated through (2) at the point $(d_{n,1} - \sqrt{c_{2,2(1-\gamma)}/n}, d_{n,2} - \sqrt{c_{2,2(1-\gamma)}/n})$.

8.3 Bootstrap computational details

The distribution of $\pi(F_{T,n}^*, F_{C,n}^*, m, \alpha)$ was approximated with $B_1 = 500$ points generated with the Monte Carlo technique. Each point was computed on the basis of a couple of samples of size n drawn from $F_{T,n}$ and $F_{C,n}$, whose empirical distribution functions, namely $F_{T,n}^*, F_{C,n}^*$, provided $\pi(F_{T,n}^*, F_{C,n}^*, m, \alpha)$. This latter power value was computed by generating $B_2 = 500$ couples of samples of size m from $F_{T,n}^*$ and $F_{C,n}^*$ which underwent the IUT, and by considering the rate of statistically significant tests.

To implement calibration, we first computed the simple plug-in estimate of the power, i.e. $\pi(F_{T,n}, F_{C,n}, m, \alpha)$. Then, assuming this latter value to be the true power, we generated $B_c = 500$ estimates of the power for each conservative level $\gamma \in (0, 1)$ and we computed the ACP. Finally, through least squares formulas, we computed the parameter a in (7) and the related corrected level γ_c for each γ -level of interest was obtained by inverting (7). Hence, $\pi^{\gamma_c}(F_{T,n}, F_{C,n}, m, \alpha)$ was computed.

Acknowledgements. Partial support was provided by the Italian Ministry of Research (MIUR) (protocol 2007 AYHZWC Statistical methods for learning in clinical research).

References

- Anderson, T.W. (1958). *An introduction to multivariate statistical analysis*. New York: John Wiley & Sons.
- Berger, R.L. (1982). Multiparameter hypothesis testing and acceptance sampling. *Technometrics* 24: 295–300.
- Chuang-Stein, C. (2006). Sample size and the probability of a successful trial. *Pharmaceutical Statistics* 5: 305–309.
- Chuang-Stein, C., Stryszak, P., Dmitrienko, A., and Offen, W. (2007). Challenge of multiple co-primary endpoints: A new approach. *Statistics in Medicine* 26: 1181–1192.
- De Martini, D. (2010). Adapting by calibration the sample size of a phase III trial on the basis of phase II data. *Pharmaceutical Statistics*.(Article first published online: 8 FEB 2010).
- De Martini, D. (2011). Conservative Sample Size Estimation in Nonparametrics. *Journal of Biopharmaceutical Statistics*.Vol. 21, n.1, in press.
- Devane, D., Begley, C.M., Clarke, M. (2004). How many do I need? Basic principles of sample size estimation. *Journal of Advanced Nursing* 47:297–302.
- Dilba, G., Bretz, F., Hothorn, L.A., Guiard, V. (2006). Power and sample size computations in simultaneous tests for non-inferiority based on relative margins. *Statistics in Medicine* 25: 1131–1147.
- Dmitrienko, A., Offen, W.W., Westfall, P.H. (2003). Gatekeeping strategies for clinical trials that do not require all primary effects to be significant. *Statistics in Medicine* 22: 2387–2400.
- Donner, A., Rosner, B. (1980). On inferences concerning a common correlation coefficient. *Journal of the Royal Statistical Society, Series C* 29: 69–76.

- Eaton, M.L., Muirhead, R.J. (2007). On a multiple endpoints testing problem. *JSPI* 137, 11: 3416–3429.
- Efron, B., Tibshirani, R.J. (1993). *An introduction to the Bootstrap*. New York: Chapman & Hall.
- Eng, J. (2003). Sample size estimation: how many individuals should be studied? *Radiology* 227:309–313.
- Gleser, L.J. (1973). On a theory of Intersection-Union Tests (Preliminary Report), *IMS Bulletin* 2, 233.
- Johnston, M.F., Hays, R.D., Hui, K.K. (2009). Evidence-based effect size estimation: An illustration using the case of acupuncture for cancer-related fatigue *BMC Complementary and Alternative Medicine* 9: 1–9.
- Lehmann, E.L. (1952). Testing multiple hypotheses. *Annals of Mathematical Statistics* 23: 541–552.
- Morrison, D.F. (2005). *Multivariate Statistical Methods*. New York: Thomson/Brooks/Cole.
- Offen W. et al. (2007). Multiple co-primary endpoints: Medical and statistical solutions - A report from the Multiple Endpoints Expert Team of the Pharmaceutical Research and Manufacturers of America. *Drug Information Journal* 41, 1: 31–46.
- R Development Core Team (2005). *R: A language and environment for statistical computing*. Vienna: R Foundation for Statistical Computing. URL <http://www.R-project.org>.
- Roy, S.N., Bose, R.C. (1953). Simultaneous confidence interval estimation. *Annals of Mathematical Statistics* 24, 4: 513–536.
- Sankoh, A.J., D’Agostino, R.B., Huque, M.F. (2003). Efficacy endpoint selection and multiplicity adjustment methods in clinical trials with inherent multiple endpoint issues. *Statistics in Medicine* 22: 3133–3150.
- Senn, S., Bretz, F. (2007). Power and sample size when multiple endpoints are considered. *Pharmaceutical Statistics* 6: 161–170.

- Song, J.X. (2009). Sample size for simultaneous testing of rate differences in non-inferiority trials with multiple endpoints. *Computational Statistics and Data Analysis* 53: 1201–1207.
- Wang, S.J., Hung, H.M.J., O’Neill, R.T. (2006). Adapting the sample size planning of a phase III trial based on phase II data. *Pharmaceutical Statistics* 5: 85–97.
- Wilson, E.B. (1927). Probable Inference, the Law of Succession, and Statistical Inference. *Journal of the American Statistical Association* 22, 158: 209–212.
- Yeo, A., Qu, Y. (2009). Evaluation of statistical power for multiple tests: a case study. *Pharmaceutical Statistics* 8: 5–11.

Table 1	Distributional parameters δ_{1s} providing 90% power				
$M_I = 60$	ρ				
	0.2	0.4	0.5	0.6	0.8
$k = 1$.65315		.64576		.63044
$k = 1.5$.88916		.88835	
$k = 2$	1.18366		1.18363		1.18363
$M_I = 180$	ρ				
	0.2	0.4	0.5	0.6	0.8
$k = 1$.37710		.37283		.36399
$k = 1.5$.51336		.51289	
$k = 2$.68337		.68337		.68339

Table 1. Design of the simulation study and shift parameters.

Table 2	All Scenario settings				
scenario #	n	m	δ_1	ρ	k
1	40	60	0.65315	0.2	1
2	80	60	0.65315	0.2	1
3	40	60	0.64576	0.5	1
4	80	60	0.64576	0.5	1
5	40	60	0.63044	0.8	1
6	80	60	0.63044	0.8	1
7	120	180	0.3771	0.2	1
8	240	180	0.3771	0.2	1
9	120	180	0.37283	0.5	1
10	240	180	0.37283	0.5	1
11	120	180	0.36399	0.8	1
12	240	180	0.36399	0.8	1
13	40	60	0.88916	0.4	1.5
14	80	60	0.88916	0.4	1.5
15	40	60	0.88835	0.6	1.5
16	80	60	0.88835	0.6	1.5
17	120	180	0.51336	0.4	1.5
18	240	180	0.51336	0.4	1.5
19	120	180	0.51289	0.6	1.5
20	240	180	0.51289	0.6	1.5
21	40	60	1.18366	0.2	2
22	80	60	1.18366	0.2	2
23	40	60	1.18363	0.5	2
24	80	60	1.18363	0.5	2
25	40	60	1.18363	0.8	2
26	80	60	1.18363	0.8	2
27	120	180	0.68339	0.2	2
28	240	180	0.68339	0.2	2
29	120	180	0.68339	0.5	2
30	240	180	0.68339	0.5	2
31	120	180	0.68339	0.8	2
32	240	180	0.68339	0.8	2

Table 2. All Scenario settings.

Table 3		Average bias			
Scenario	#1	#2	#5	#21	
ELLP	16.88%	15.86%	13.64%	12.06%	
ELLM	16.93%	16.02%	14.88%	12.07%	
SIMP	14.09%	13.40%	8.68%	5.21%	
SIMM	14.26%	13.66%	10.47%	5.25%	
BO	13.84%	12.13%	5.63%	0.34%	
Rate of improvement					
	w.r.t. n	w.r.t. ρ	w.r.t. k		
ELLP		6.1%	19.2%	28.5%	
ELLM		5.4%	12.1%	28.7%	
SIMP		4.9%	38.4%	63.0%	
SIMM		4.2%	26.6%	63.2%	
BO		12.3%	59.3%	97.5%	
Values of $I2$ index					
ELLP	0.5906	0.4145	0.5756	0.5743	
ELLM	0.5971	0.4224	0.6300	0.5747	
SIMP	0.5484	0.3923	0.5209	0.5171	
SIMM	0.5593	0.4035	0.5612	0.5178	
BO	0.5176	0.3706	0.5252	0.5132	

Table 3. Average biases of ACP of the five techniques with respect to NCP and $I2$ values under four Scenarios with $M_I = 60$.

Table 4 scenario	Average biases				I_2 values	
	BO	(a)	BOC	% of impr.	BO	BOC
1	13.84%	(-0.8690)	4.56%	67.0%	0.5176	0.4772
3	9.79%	(-0.6163)	2.67%	72.8%	0.5207	0.5169
5	5.63%	(-0.3576)	3.69%	34.5%	0.5252	0.5231
13	4.89%	(-0.2743)	-2.53%	48.2%	0.5030	0.5034
15	2.41%	(-0.1442)	-0.27%	88.9%	0.5132	0.5182
21	0.34%	(-0.0277)	-2.10%	-512.7%	0.5132	0.5028
23	-0.69%	(0.0476)	1.04%	-51.0%	0.5114	0.5157
25	-1.24%	(0.0663)	-0.04%	96.5%	0.5178	0.5189

Table 4. Average biases of ACP of BO with respect to NCP and I_2 values under the eight Scenarios with $n = 40$, in comparison with those of BOC.

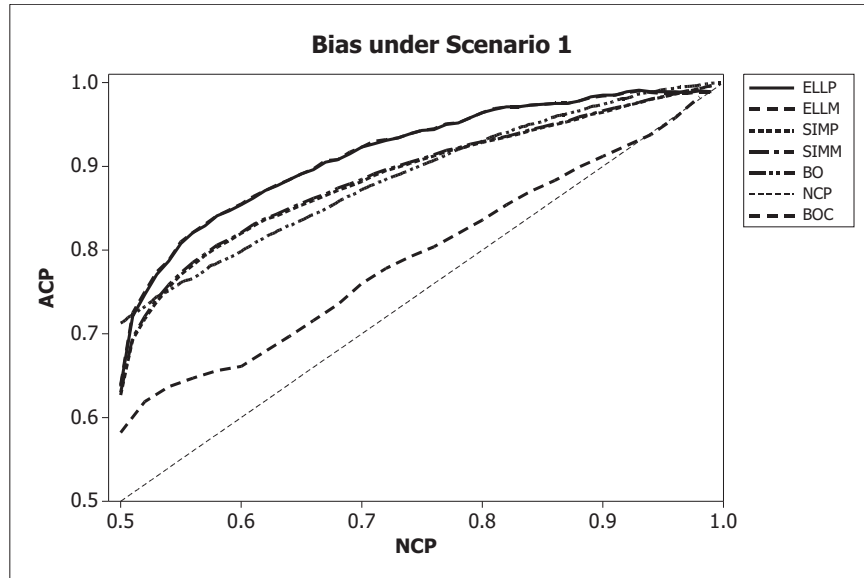


Figure 1. Bias of ACPs of the five techniques, together with that of BOC, under Scenario 1.

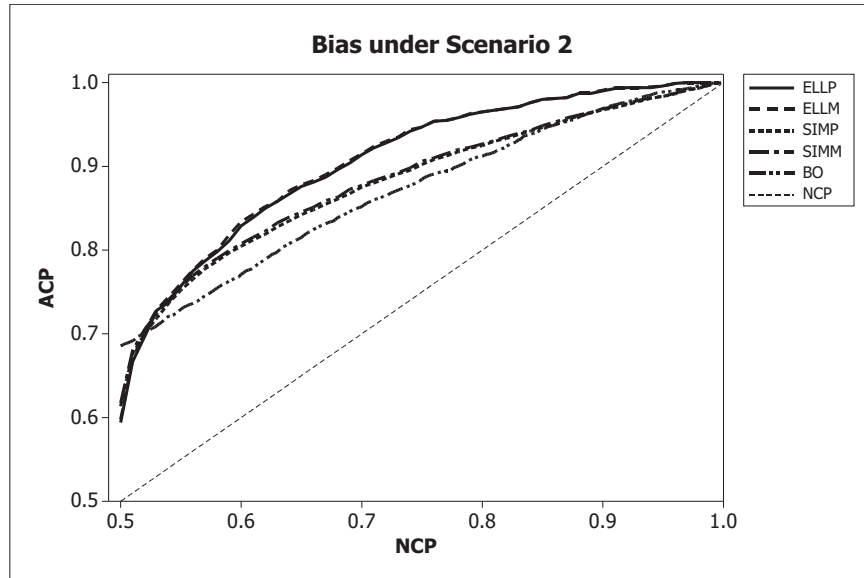


Figure 2. Bias of ACPs of the five techniques under Scenario 2.

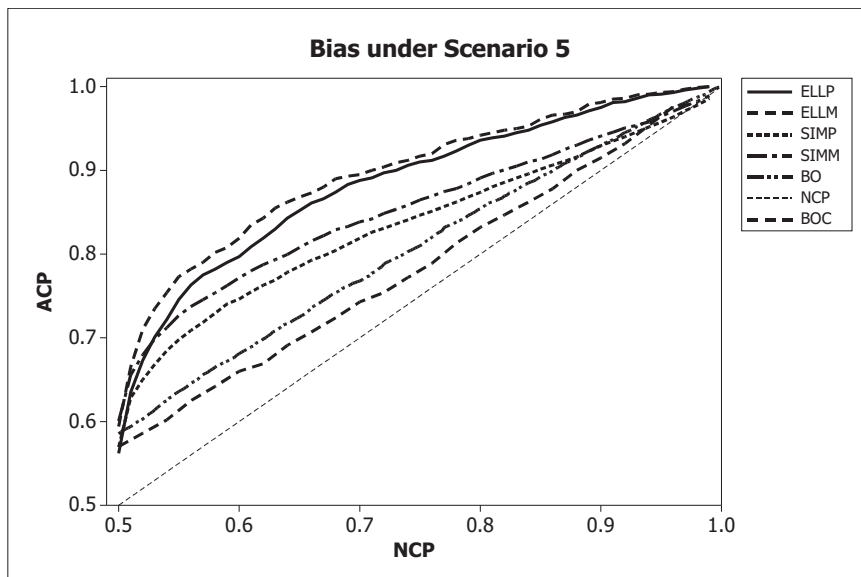


Figure 3. Bias of ACPs of the five techniques, together with that of BOC, under Scenario 5.

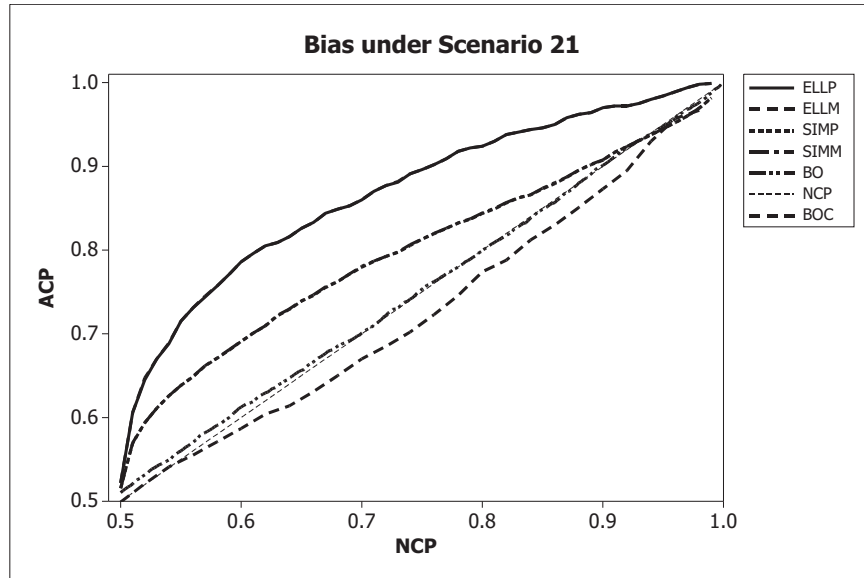


Figure 4. Bias of ACPs of the five techniques, together with that of BOC, under Scenario 21.

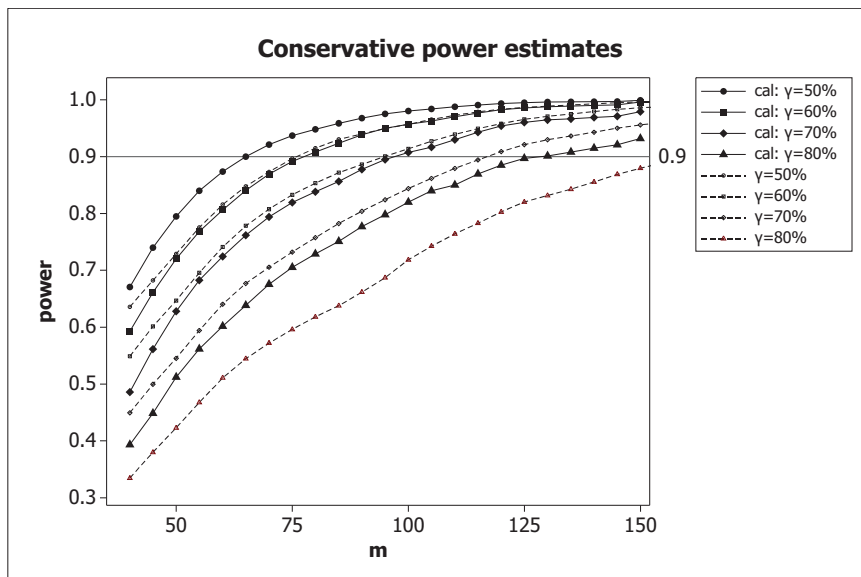


Figure 5. Conservative power estimates for the IU test obtained with BOC and BO, with $\gamma = 50\%$, 60% , 70% , and 80% , based on the $n = 48$ phase II data of the example in Section 6.