

# Efficient Experimental Design for the Behrens-Fisher Problem With Application to Bioassay

Holger DETTE and Timothy E. O'BRIEN

---

A common experimental design for the problem of comparing two means from a normal distribution assumes knowledge of the ratio of the population variances. The optimal sampling ratio is proportional to the square root of this quantity. This article demonstrates that a misspecification of the ratio of the population variances can cause a substantial loss in power of the corresponding tests. As a robust alternative, a maximin approach is used to construct designs, which are efficient, whenever the experimenter is able to specify a specific region for the ratio of the population variances. The advantages of the robust designs for inference in the Behrens-Fisher problem are illustrated in a simulation study and an application to the design of experiment for bioassay is presented.

KEY WORDS: Behrens-Fisher problem; Bioassay; Design of experiment; Local optimal design; Robust designs.

---

## 1. INTRODUCTION

The problem of comparing the means of two populations based on a sample of observations is of fundamental importance in applied statistics. Let  $\mu_i, \sigma_i^2$  denote the population mean and variance of the  $i$ th population for  $i = 1, 2$ , then the parameter of interest is typically the difference of the means  $\mu = \mu_1 - \mu_2$  or the ratio  $\rho = \mu_2/\mu_1$ . If the ratio  $\kappa = (\sigma_2^2/\sigma_1^2)$  of the population variances is unknown and the underlying populations are assumed normally distributed, the scenario is the well known Behrens-Fisher problem (see Scheffé 1970). There is a large number of articles in which various tests are suggested concerning the hypothesis regarding the difference of the means  $\mu$ . In the case of testing simple hypotheses, Welch's approximate  $t$ -solution (see Welch 1936, 1938) appears to be a good compromise between a test that is unbiased and that is appealing to practitioners because of its simplicity; see, for example, Wang (1971) and Best and Rayner (1987). This approach was further extended by Dannenberg, Dette, and Munk (1994) for testing interval hypotheses.

In contrast to the goal of constructing useful tests for the Behrens-Fisher problem, the problem of allocating observations

from both populations if the total sample size has been fixed has not found much attention in the literature. It is well known (see, e.g., Staudte and Sheater 1990) that if  $n_1$  and  $n_2$  denote the sample sizes from both populations, the power of Welch's test is maximized if  $n_1/n_2 \approx \kappa^{-1/2} = \sigma_1/\sigma_2$ . A similar observation was made by Dannenberg, Dette, and Munk (1994) in the context of testing interval hypotheses of the form  $H_0 : \mu \notin [-\Delta, \Delta], H_1 : \mu \in [-\Delta, \Delta]$ . However, these results are "local" in the sense of Chernoff (1953) as they require knowledge of the population variances in order to determine the sample sizes  $n_1$  and  $n_2$ .

Section 2 demonstrates by means of a simulation study that the loss of power caused by such a misspecification can be substantial. Consequently, a misspecification of  $\kappa$  can yield a substantial loss in power if the sample sizes are chosen according to the rule  $n_1/n_2 \approx \kappa^{-1/2}$ . To obtain designs that are less sensitive with respect to such misspecifications, we propose the maximization of the minimum of an appropriately standardized power function taken over a certain range for the parameter  $\kappa$  with respect to the proportion of total observations in the first sample. We also give an explicit formula for the relative proportions for both samples with respect to the new criterion, and we demonstrate the ease with which this technique can be applied in practical settings. It is demonstrated by means of a simulation study that the new designs are robust and efficient whenever a range for the unknown ratio of the population variances can be specified.

Our new methodology is applied to the classical problem of testing the difference of two normal means and to the important problem of inference about the ratio of these means useful in direct bioassays.

## 2. LOCAL OPTIMAL ALLOCATION OF SAMPLE SIZES

Let  $X_1, \dots, X_{n_1}$  and  $Y_1, \dots, Y_{n_2}$  denote two independent samples of independent identically distributed observations such that  $X_i \sim \mathcal{N}(\mu_1, \sigma_1^2)$ ,  $i = 1, \dots, n_1$ ;  $Y_j \sim \mathcal{N}(\mu_2, \sigma_2^2)$ ,  $j = 1, \dots, n_2$ , and consider the one-sided problem of testing the hypotheses

$$H_0 : \mu = \mu_1 - \mu_2 \leq 0 \quad \text{versus} \quad H_1 : \mu > 0. \quad (1)$$

In a famous article, Welch (1938) suggested the rejection of the null hypothesis if

$$\frac{\bar{X}_{n_1} - \bar{Y}_{n_2}}{\sqrt{\frac{1}{n_1} \hat{S}_1^2 + \frac{1}{n_2} \hat{S}_2^2}} > t_{1-\alpha, \hat{f}} \quad (2)$$

where  $\bar{X}_{n_1}, \bar{Y}_{n_2}$  denote the sample means,  $\hat{S}_1^2 = \frac{1}{n_1-1} \sum_{i=1}^{n_1} (X_i - \bar{X}_{n_1})^2$ ,  $\hat{S}_2^2 = \frac{1}{n_2-1} \sum_{i=1}^{n_2} (Y_i - \bar{Y}_{n_2})^2$  are the common estimators of the population variances  $\sigma_1^2, \sigma_2^2$ , respectively, and  $t_{1-\alpha, \hat{f}}$

---

Holger Dette is Professor, Ruhr-Universität Bochum, Fakultät für Mathematik, 44780 Bochum, Germany (E-mail: holger.dette@ruhr-uni-bochum.de). Timothy E. O'Brien is Assistant Professor, Loyola University Chicago, Department of Mathematics and Statistics, 6525 N. Sheridan Road, Chicago, IL 60626. The authors are grateful to Isolde Gottschlich who typed parts of this article with considerable technical expertise. The financial support of the Deutsche Forschungsgemeinschaft (SFB 475, reduction for complexity in multivariate data structures) is gratefully acknowledged. The authors are also grateful to an unknown associate editor and to the editor for their helpful comments, which led to a substantial improvement of an earlier version of this article.

is the  $(1 - \alpha)$ -quantile of the  $t$  distribution with

$$\hat{f} = \frac{\left(\frac{\widehat{S}_1^2}{n_1} + \frac{\widehat{S}_2^2}{n_2}\right)^2}{\left(\frac{\widehat{S}_1^2}{n_2}\right)^2 / (n_1 - 1) + \left(\frac{\widehat{S}_2^2}{n_2}\right)^2 / (n_2 - 1)} \quad (3)$$

estimated degrees of freedom. It was pointed out by Scheffé (1970) and Wang (1971) that this test provides a good compromise between tests which should on the one hand be unbiased and on the other hand be easily implemented. The performance of a given test is measured by its power function under some alternative hypothesis. It is easy to see that for a fixed sample size the power of Welch's test must depend on the relative proportions  $n_1/(n_1 + n_2)$  and  $n_2/(n_1 + n_2)$ . For example, if  $n_1 = 0$  or  $n_2 = 0$  the power of the test is 0 and it is impossible to test hypotheses regarding the difference of the population means, because observations are available only from one population. However, what is a good choice of the relative sample sizes to obtain a most efficient inference?

Throughout this article we call any specification of the relative proportion  $n_1/(n_1 + n_2)$  of total observations for the first sample an experimental design. The optimal design problem is to maximize the power of the test with respect to the choice of  $n_1$  for a fixed sample size  $n_1 + n_2$ . Because it is not clear which alternative should be used for this calculation one usually considers "local" alternatives, very close to the null hypothesis if the total sample size is large, and for this reason, particularly difficult to detect. It is well known (see Staudte and Sheater 1990, p. 180) that for local alternatives of the form

$$\mu = \frac{\sigma_1}{\sqrt{n_1 + n_2}} \quad (4)$$

the asymptotic power function of this test is given by

$$\pi(\kappa) = \Phi \left( \left\{ \frac{1}{w} + \frac{\kappa}{1-w} \right\}^{-1/2} - u_{1-\alpha} \right), \quad (5)$$

where  $\kappa = \sigma_2^2/\sigma_1^2$  is the ratio of the population variances,  $u_{1-\alpha} = \Phi^{-1}(1 - \alpha)$  is the quantile of the standard normal distribution and

$$\lim_{\substack{n_1 \rightarrow \infty \\ n_2 \rightarrow \infty}} \frac{n_1}{n_1 + n_2} = w \in (0, 1) \quad (6)$$

is asymptotically the relative proportion of total observations in the first sample. It was pointed out by Dette and Munk (1997) that  $\pi(\kappa)$  also coincides with the asymptotic power function of the extension of Welch's test to the problem of testing the equivalence hypotheses

$$H_0 : \mu \notin [-\Delta, \Delta]; H_1 : \mu \in [-\Delta, \Delta] \quad (7)$$

under contiguous alternatives  $\mu = \Delta + \sigma_1(n_1 + n_2)^{-1/2}$ . A simple calculation shows that the power  $\pi(\kappa)$  is maximal if

$$\frac{n_1}{n_1 + n_2} \approx w_\kappa^* = \frac{1}{1 + \sqrt{\kappa}} = \frac{1}{1 + \sigma_2/\sigma_1}, \quad (8)$$

and we will call  $w_\kappa^*$  the local optimal design for testing the hypotheses (1) or (7). The phrase "local" is due to Chernoff (1953) and used because the optimal allocation to both samples depends

on the unknown parameter  $\kappa = \sigma_2^2/\sigma_1^2$ . If some information regarding the ratio of population variances is available, the power of Welch's test can be increased substantially by using the rule (8). However, the following example shows that in general the local optimal design is indeed sensitive with respect to misspecification of the parameter  $\kappa$ .

*Example 1.* We have conducted a small simulation study, where  $\mu = 1, \sigma_1^2 + \sigma_2^2 = 5$ , and the "true" ratio  $\kappa_t^{1/2} = \sigma_2/\sigma_1$  varies between 1 and 1/5. We have calculated the power of Welch's test (2) with nominal level 5% for the hypotheses (1) for various designs, which are calculated under the respective assumptions that the ratio is given by  $\kappa_a^{1/2} = 1, 1/3, 1/5$ . In other words, if  $\kappa_t \neq \kappa_a$  the design was calculated under a misspecification for the ratio of the population variances. The local optimal designs are obtained by a simple rounding procedure from the values  $(n_1 + n_2) \cdot w_\kappa^* = (n_1 + n_2) \cdot (1 + \sqrt{\kappa})^{-1/2}$ , which gives the sample size for the first sample. The rejection probabilities of the test (2) are calculated by 10,000 simulation runs, while the total sample sizes satisfy  $n_1 + n_2 = 25$  or  $n_1 + n_2 = 50$ .

Table 1 shows the loss of efficiency if a design has been calculated by a misspecification of the parameter  $\kappa$ . The efficiency losses are believed accurate to the reported precision. The loss of efficiency is remarkably large. For example, if the "true" ratio of the population variances is given by  $\kappa_t^{1/2} = 1$ , but the local optimal design is found under the assumption that  $\kappa_a^{1/2} = 1/3$ , then we obtain for the sample size  $n_1 + n_2 = 50$  the power 0.581, while the best design yields power 0.715. This corresponds to a loss of power of approximately 19%  $\approx (0.715 - 0.581)/0.715$ , which is the value listed in Table 1. The results indicate that the optimal allocation rule (8) is rather sensitive with respect to a misspecification of the unknown ratio of the population variances. For example, the allocation rule  $n_1 = 19, n_2 = 6$  (corresponding to the assumption  $\kappa_a^{1/2} = 1/3$ ) yields a loss of efficiency of 21% ( $\kappa_t^{1/2} = 1$ ) and 1% ( $\kappa_t^{1/2} = 1/5$ ) while it is the best for  $\kappa_t = 1/3$ . Similarly, the loss of efficiency of the allocation rule  $n_1 = 21, n_2 = 4$  (corresponding to the assumption  $\kappa_a^{1/2} = 1/5$ ) is approximately 40% ( $\kappa_t^{1/2} = 1$ ) and 7% ( $\kappa_t^{1/2} = 1/3$ ).

In the following section robust designs will be calculated by a maximin approach, which uses only the information that the ratio of the population standard deviations lies in the interval  $[1/5, 1]$ . We feel this is the more realistic setting because practitioners will rarely be able to give an accurate point estimate for the ratio of the variances, whereas an accurate interval estimate can usually be given.

### 3. ROBUST DESIGNS FOR THE BEHRENS-FISHER PROBLEM

Note from (5) that the power function of the test (2) increases with the expression

$$f(w, \kappa) = \left\{ \frac{1}{w} + \frac{\kappa}{1-w} \right\}^{-1}, \quad (9)$$

and that the local optimal design  $w_\kappa^* = 1/(1 + \kappa^{1/2})$  is found by maximizing  $f(w, \kappa)$  with respect to  $w$  for given  $\kappa$  [see the

Table 1. Loss of Efficiency of Welch's Test (2) for the Hypotheses (1) for Various Designs and Ratios  $\kappa_t = \sigma_2^2/\sigma_1^2$  of Population Variances. The results are based on 10,000 simulation runs.

$\kappa_a^{1/2}$	$n_1 + n_2 = 25$				$n_1 + n_2 = 50$			
	1	1/3	1/5	robust	1	1/3	1/5	robust
	$n_1 = 13$	$n_1 = 19$	$n_1 = 21$	$n_1 = 17$	$n_1 = 25$	$n_1 = 37$	$n_1 = 41$	$n_1 = 33$
$\kappa_t^{1/2}$	$n_2 = 12$	$n_2 = 6$	$n_2 = 4$	$n_2 = 8$	$n_2 = 25$	$n_2 = 13$	$n_2 = 9$	$n_2 = 17$
1	0%	21%	40%	9%	0.0%	19%	37%	6%
1/3	14%	0%	7%	0%	12%	0.0%	2%	0%
1/5	22%	1%	0%	7%	17%	1.0%	0%	5%

derivation of (8)]. The performance of a particular given design can be measured by its efficiency

$$\text{eff}(w, \kappa) = \frac{f(w, \kappa)}{\max_v f(v, \kappa)} = \frac{(1 + \sqrt{\kappa})^2}{\frac{1}{w} + \frac{\kappa}{1-w}}. \quad (10)$$

Roughly speaking  $1 - \text{eff}(w, \kappa)$  measures the loss in power if  $\kappa$  is the “true” ratio of population variances and the design  $w$  is used instead of the local optimal design  $w_\kappa^*$ , which requires the knowledge of  $\kappa$ . Note that the efficiency varies between 0 and 1 and that a design with efficiency close to 1 yields the best power. For example, if  $\kappa = 1$  the local optimal design advises the experimenter to take equal sample sizes in both samples (i.e.,  $w_1^* = 0.5$ ) and this design has efficiency  $\text{eff}(0.5, 1) = 1$ . On the other hand the design, which takes 82% of the observations in the first sample, has efficiency  $\text{eff}(0.82, 1) = 0.59$ .

In Example 1, we showed that local optimal designs are not necessarily robust with respect to a misspecification of the unknown ratio of the population variances. For the construction of a more robust design, we assume that an interval, say  $[\kappa_L, \kappa_U]$ , for the unknown population variance can be specified by the experimenter and determine a design that maximizes the worst efficiency over this interval. It follows that the resulting design will have reasonable efficiencies over the full interval  $[\kappa_L, \kappa_U]$ . We call a design  $w^*$  standardized maximin optimal if it maximizes the minimum efficiency

$$g(w) = \min_{\kappa \in [\kappa_L, \kappa_U]} \text{eff}(w, \kappa) \quad (11)$$

over the interval  $[\kappa_L, \kappa_U]$ . This design criterion is similar to the standardized optimality criteria used by Dette (1997) and Imhof (2001). Further, the Appendix establishes that for fixed  $w$  the function  $\kappa \rightarrow \text{eff}(w, \kappa)$  is unimodal with at most one maximum in the interval  $[\kappa_L, \kappa_U]$  (see Lemma A.1). It therefore follows that

$$g(w) = \min\{\text{eff}(w, \kappa_L), \text{eff}(w, \kappa_U)\}. \quad (12)$$

Moreover, Lemma A.2 (see the Appendix) shows that for the standardized maximin optimal design

$$w^* = \arg \max_{w \in [0,1]} g(w)$$

it follows that  $\text{eff}(w^*, \kappa_L) = \text{eff}(w^*, \kappa_U)$ . This equality determines the optimal design as

$$w^* = \frac{2 + \kappa_L^{1/2} + \kappa_U^{1/2}}{2(1 + \kappa_L^{1/2})(1 + \kappa_U^{1/2})} \quad (13)$$

for which the minimal efficiency is

$$g(w^*) = \frac{(2 + \kappa_L^{1/2} + \kappa_U^{1/2})\{\kappa_L^{1/2}(1 + \kappa_U^{1/2}) + \kappa_U^{1/2}(1 + \kappa_L^{1/2})\}}{2(1 + \kappa_L^{1/2})(1 + \kappa_U^{1/2})(\kappa_L^{1/2} + \kappa_U^{1/2})}. \quad (14)$$

*Example 2.* If the experimenter specifies the interval  $(\kappa_L^{1/2}, \kappa_U^{1/2}) = (1/5, 1)$  for the ratio of the standard deviations the standardized maximin optimal design weight is  $w^* = 2/3$  and the minimal efficiency is  $g(w^*) = 8/9$ . This high value of the minimal value of the design efficiency underscores the remarkable robustness of our robust design. Incidentally, the corresponding weight is translated into a practical design allocation for the first sample by rounding

$$(n_1 + n_2) \cdot w^* = (n_1 + n_2)2/3$$

to the nearest integer (as in Table 1). In the fourth column (labeled “robust”) Table 1 also contains the loss of efficiency of this robust design for all situations under consideration. For example, if  $n_1 + n_2 = 25$  the loss of efficiency of the allocation rule  $n_1 = 17, n_2 = 8$  compared to the best design is only approximately 9% ( $\kappa_t^{1/2} = 1$ ), 0% ( $\kappa_t^{1/2} = 1/3$ ), and 7% ( $\kappa_r^{1/2} = 1/5$ ). Thus, the new design constructed by the maximin approach is quite robust and efficient. The results of Table 1 along with additional simulations (not shown for the sake of brevity) indicate that robust and efficient designs are available if an interval for the unknown ratio of the population variances can be specified by the experimenter.

*Remark 1.* We also note that the design problem is symmetric in the following sense. If  $w_{\kappa_L, \kappa_U}^*$  denotes the standardized maximin optimal proportion for the first sample if the parameter  $\kappa$  is assumed to be in the interval  $[\kappa_L, \kappa_U]$ , then the corresponding quantity for the interval  $[1/\kappa_U, 1/\kappa_L]$  satisfies

$$w_{1/\kappa_U, 1/\kappa_L}^* = 1 - w_{\kappa_L, \kappa_U}^*.$$

It follows that the standardized maximin optimal design for the interval  $[1/\kappa_U, 1/\kappa_L]$  can be obtained from the corresponding design for the interval  $[\kappa_L, \kappa_U]$  by interchanging the role of the sample sizes  $n_1$  and  $n_2$ . For this reason the robust designs can easily be tabulated. Some designs for selected values of  $\kappa_L$  and  $\kappa_U$  are presented in Table 2. Finally, we note that this symmetry implies that the equal allocation rule  $w^* = 1/2$  is standardized maximin optimal for any interval of the form  $[1/\kappa_0, \kappa_0]$  where  $\kappa_0 > 1$ .

Table 2. Standardized Maximin Optimal Designs for Various Intervals  $[\kappa_L, \kappa_U]$  for the Unknown Ratio  $\kappa = \sigma_2^2/\sigma_1^2$  of the Population Variances. The value  $w^*$  in the table gives the relative proportion of total observations in the first sample.

$\kappa_L \backslash \kappa_U$	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
0.1	0.760	0.725	0.703	0.686	0.673	0.662	0.652	0.644	0.636	0.630
	0.2	0.691	0.668	0.652	0.638	0.627	0.618	0.609	0.602	0.595
		0.3	0.646	0.629	0.616	0.605	0.595	0.587	0.580	0.573
			0.4	0.613	0.599	0.588	0.579	0.570	0.563	0.556
				0.5	0.586	0.575	0.565	0.557	0.549	0.543
					0.6	0.563	0.554	0.546	0.538	0.532
						0.7	0.544	0.536	0.529	0.522
							0.8	0.528	0.521	0.514
								0.9	0.513	0.507
									1.0	0.5

Example 3. The results derived so far have been derived under the assumption that one-sided hypotheses are tested with Welch's approximate  $t$ -solution. It follows from Dette and Munk (1997) that these results are directly applicable to the problem of testing the equivalence hypotheses  $H_0 : \mu \notin [-\Delta, \Delta]$ ;  $H_1 : \mu \in [-\Delta, \Delta]$ , because the asymptotic power function coincides with that of the one-sided problem.

In principle, a similar analysis could be performed for cases where simple hypotheses  $H_0 : \mu = 0$ ;  $H_1 : \mu \neq 0$  or interval hypotheses  $H_0 : \mu \in [-\Delta, \Delta]$ ;  $H_1 : \mu \notin [-\Delta, \Delta]$  are of interest. However, our numerical results show that the designs derived for the one-sided problem are also very efficient for testing other hypotheses. By way of illustration, consider the situation of Example 1 where  $\sigma_1^2 + \sigma_2^2 = 5$  and a test with level 5% for the hypotheses  $H_0 : \mu = 0$ ;  $H_1 : \mu \neq 0$  has to be performed. To demonstrate the application of Remark 1, we consider the cases where  $\kappa_t^{1/2} = 1, 3, 5$  for the true value of the ratio of the variances, while we assumed  $\kappa_L^{1/2} = 1$  and  $\kappa_U^{1/2} = 5$  for the construction of the robust design. The optimal proportion for the first sample is now given by  $w^* = 1/3$  and the simulated loss of efficiency is given in Table 3 for sample sizes  $n_1 + n_2 = 25$  or 50. We observe a similar picture as for the one-sided case. The local optimal designs are sensitive with respect to misspecification of the unknown ratio of population variances, while the standard maximin optimal designs yield a reasonable power in all cases under consideration.

#### 4. APPLICATION TO BIOASSAY

One concern of bioassay, or biological assays, is the estimation of the potency of one drug (B) relative to another (A), typi-

cally involving comparing a new drug with a standard. Further, in contrast with indirect assays, direct assays hold that the necessary concentrations that produce the same therapeutic effect can be directly measured. In this setting, the relative potency ( $\rho$ ) of drug B to A is the ratio of the respective means, where the underlying respective distributions are assumed to be Gaussian  $A \sim \mathcal{N}(\mu_1, \sigma_1^2)$ ,  $B \sim \mathcal{N}(\mu_2, \sigma_2^2)$ ; thus,  $\rho = \mu_2/\mu_1$ . Further background of direct assays were given by Finney (1978, chap. 2) and Govindarajulu (2000, chap. 2).

Often practitioners are interested in a confidence interval for the relative potency, and experimental designs which produce shorter confidence intervals are therefore desired. In the case of independent populations, a standard calculation shows that the first-order approximation for the length of any reasonable confidence interval is proportional to the root of the function

$$g(w, \kappa, \rho) = \frac{1}{w} + \frac{\kappa/\rho^2}{1-w},$$

and all results of the previous sections are therefore applicable to this case but with  $\kappa$  replaced by  $\kappa/\rho^2$ . For example, the local optimal design uses

$$w_{\kappa/\rho}^* = \frac{1}{1 + \sqrt{\kappa/\rho}} \quad (15)$$

as the weight for the first sample. Similarly, if the experimenter is able to specify a region, say  $[\kappa_L, \kappa_U]$  for the quantity  $\kappa/\rho^2$  the optimal design is given by (13).

Consider for example the situation where the population variances are the same, that is,  $\kappa = 1$ , and a confidence interval is constructed using Fieller's theorem (Finney 1978). This interval

Table 3. Loss of Efficiency of Welch's Test of a Simple Hypothesis for Various Designs and Ratios  $\kappa_t = \sigma_2^2/\sigma_1^2$  of Population Variances. The results are based on 10,000 simulation runs.

$\kappa_a^{1/2}$	$n_1 + n_2 = 25$				$n_1 + n_2 = 50$			
	1	3	5	robust	1	3	5	robust
$\kappa_t^{1/2}$	$n_1 = 12$ $n_2 = 13$	$n_1 = 6$ $n_2 = 19$	$n_1 = 4$ $n_2 = 21$	$n_1 = 8$ $n_2 = 17$	$n_1 = 25$ $n_2 = 25$	$n_1 = 12$ $n_2 = 38$	$n_1 = 8$ $n_2 = 42$	$n_1 = 17$ $n_2 = 33$
1	0%	21%	36%	12%	0%	22%	40%	10%
3	24%	0%	3%	5%	15%	0%	4%	2%
5	31%	3%	0%	8%	24%	3%	0%	8%

Table 4. Loss of Efficiency in the Construction of Fieller's Confidence Interval for the Relative Potency for Various Designs and Different Values of  $\rho_t = \mu_2/\mu_1$ . The results are based on 10,000 simulation runs.

$\rho_a$	$n_1 + n_2 = 50$				
	1.0	2.25	4.0	6.25	robust
$\rho_t$	$n_1 = 25$ $n_2 = 25$	$n_1 = 35$ $n_2 = 15$	$n_1 = 40$ $n_2 = 10$	$n_1 = 43$ $n_2 = 7$	$n_1 = 30$ $n_2 = 20$
1.0	0%	7%	23%	41%	1%
2.25	3%	0%	10%	22%	2%
4.0	8%	0%	5%	13%	2%
6.25	11%	0%	2%	8%	4%

is of the form

$$(\rho_L, \rho_U) = \left( \frac{\hat{\rho} - \frac{ts}{\bar{X}_{n_1}} \left\{ \frac{1}{n_2} + \hat{\rho}^2 \frac{1}{n_1} - \frac{g}{n_2} \right\}^{1/2}}{1 - g}, \frac{\hat{\rho} + \frac{ts}{\bar{X}_{n_1}} \left\{ \frac{1}{n_2} + \hat{\rho}^2 \frac{1}{n_1} - \frac{g}{n_2} \right\}^{1/2}}{1 - g} \right),$$

where  $g = t^2 s^2 / (n_1 \bar{X}_{n_1}^2)$ ,  $t$  is the  $(1 - \alpha)$ -quantile of the  $t$ -distribution with  $n_1 + n_2 - 2$  degrees of freedom,  $\hat{\rho} = \bar{Y}_{n_2} / \bar{X}_{n_1}$ , and  $\bar{S}^2$  is the pooled variance estimate. To highlight the benefits of our robust design strategy, we have performed a small simulation study to calculate the average length

$$\hat{L} = \rho_U - \rho_L$$

of this interval for different designs. For this simulation, the true relative potency  $\rho_t$  varies between 1, 2.25, 4, and 6.25, and for the construction of the locally optimal designs by formula (15) we again assume  $\sigma_1^2 = \sigma_2^2 = 0.25$ , (whence  $\kappa = 1$ ). The results are given in Table 4 and show the loss of efficiency, if the relative potency has been misspecified. We observe a strong dependence on the specification of the relative potency. Thus, a misspecification of this quantity can produce a substantially larger confidence interval. For example, if the true relative potency is  $\rho_t = 1$  but we use a design based on the assumption  $\rho_a = 4$ , the length of the resulting confidence interval is increased by  $23\% \approx (0.605 - 0.493)/0.493$ . On the other hand, the robust design given in the table is constructed under the assumption that the true  $\rho_t$  lies in the interval  $[1, 6.25]$ , and yields the optimal weight  $w^* = 0.607$  using formula (13). For the total sample size  $n_1 + n_2 = 50$ , this weight translates into the allocation  $n_1 = 30$  and  $n_2 = 20$ , for a total sample size of  $n_1 + n_2 = 50$ . From Equation (14), this robust design has an efficiency of at least 95.41%. This fundamental result is illustrated in our simulation study, which shows that the robust design is indeed both robust to the choice of  $\rho$  and very efficient with a loss of efficiency of at most 4% (see Table 4).

## 5. CONCLUDING REMARKS

This article determines efficient and robust designs for Welch's approximate  $t$  test for testing one-sided hypotheses. Our method is based on a maximin approach and we have shown their usefulness and superiority in the classical setting of inference for the difference of two means. An explicit formula for the

proportions of total observations for both samples is given and the designs can easily be implemented if the experimenter is able to specify a region  $[\kappa_L, \kappa_U]$  for the unknown ratio  $\kappa = \sigma_2^2/\sigma_1^2$  of the population variances. It is demonstrated by means of a simulation study that the derived designs yield to an efficient inference for all  $\kappa \in [\kappa_L, \kappa_U]$ , whenever  $0.2 \leq \kappa_L^{1/2} \leq \kappa_U^{1/2} \leq 1$  (equivalently  $1 \leq \kappa_L^{1/2} \leq \kappa_U^{1/2} \leq 5$ ). This should encompass most cases of practical interest. An experiment with a larger (smaller) ratio of standard deviations should never be performed because the power of the Welch test becomes very small.

We have concentrated on one-sided hypotheses of the form (1) for the sake of brevity. However, for the problem of testing the equivalence hypotheses  $H_0 : \mu \notin [-\Delta, \Delta]$ ;  $H_1 : \mu \in [-\Delta, \Delta]$  it was shown by Dette and Munk (1997) that the asymptotic power function of an extension of Welch's test coincides with the power function of the test for one-sided hypotheses. As a consequence the results obtained in this article are applicable for testing interval hypotheses by Welch's approximate  $t$  solution introduced by Dannenberg, Dette, and Munk (1994). Moreover, it is demonstrated that the designs derived in Section 3 also provide a robust and efficient allocation for the problem of testing simple hypotheses. For these reasons we recommend to use these designs for the Behrens-Fisher problem of testing the difference of two means whenever an interval for the ratio of the population variances can be specified.

The results are also applicable for the classical problem of bioassay where the goal of the experiment is the estimation of the potency of one drug relative to another. For this problem, robust and efficient designs can be obtained from the results of this article whenever the experimenter is able to specify an interval for the ratio  $\kappa/\rho^2$  where  $\rho$  is the unknown relative potency and  $\kappa$  the ratio of the population variances.

## APPENDIX

*Lemma A.1.* For fixed  $w$  the function  $\kappa \rightarrow \text{eff}(w, \kappa)$  defined in (10) is unimodal with at most one maximum in the interval  $[\kappa_L, \kappa_U]$ .

*Proof.* Recall the definition of the efficiency in (10). A straightforward calculation shows that

$$\frac{\partial}{\partial \tilde{\kappa}} (\log(\text{eff}(w, \tilde{\kappa}^2))) = 2 \frac{(\tilde{\kappa} + 1)w - 1}{(1 + \tilde{\kappa})(w - 1 - w\tilde{\kappa})},$$

which vanishes only at the point  $\tilde{\kappa} = (1 - w)/w$ . A similar calculation of the second derivative yields

$$\frac{\partial^2}{\partial \tilde{\kappa}^2} \log(\text{eff}(w, \tilde{\kappa}^2)) \Big|_{\tilde{\kappa} = \frac{1-w}{w}} = \frac{2w^3}{(w-1)(w+(1-w))^2} < 0.$$

Consequently it follows that the function  $\text{eff}(w, \tilde{\kappa})$  has at most one extremum in the interval  $[\kappa_L, \kappa_U]$ , which is a maximum.

*Lemma A.2.* If  $w_{\kappa_L, \kappa_U}^*$  denotes the standardized maximin optimal design, then

$$\text{eff}(w_{\kappa_L, \kappa_U}^*, \kappa_L) = \text{eff}(w_{\kappa_L, \kappa_U}^*, \kappa_U).$$

*Proof.* We can split the maximization of the right-hand side of (12) in the maximization over the sets

$$\begin{aligned}\mathcal{M}_< &= \{w \in [0, 1] \mid \text{eff}(w, \kappa_L) < \text{eff}(w, \kappa_U)\}, \\ \mathcal{M}_> &= \{w \in [0, 1] \mid \text{eff}(w, \kappa_L) > \text{eff}(w, \kappa_U)\}, \\ \mathcal{M}_= &= \{w \in [0, 1] \mid \text{eff}(w, \kappa_L) = \text{eff}(w, \kappa_U)\}.\end{aligned}$$

Now assume that  $w_{\kappa_L, \kappa_U}^* \in \mathcal{M}_<$ . In this case we obtain  $w_{\kappa_L, \kappa_U}^* = 1/(1 + \sqrt{\kappa_L})$  and by the definition of  $\mathcal{M}_<$  the inequality

$$\text{eff}\left(\frac{1}{1 + \sqrt{\kappa_L}}, \kappa_L\right) < \text{eff}\left(\frac{1}{1 + \sqrt{\kappa_L}}, \kappa_U\right).$$

But this inequality is equivalent to

$$(\sqrt{\kappa_L} - \sqrt{\kappa_U})^2 < 0,$$

which yields a contradiction. A similar argument for the set  $\mathcal{M}_>$  shows that the maximum is attained in  $\mathcal{M}_=$ , which completes the proof.

[Received August 2003. Revised February 2004.]

## REFERENCES

- Best, D. J., and Rayner, J. C. W. (1987), "Welch's Approximate Solution for the

Behrens-Fisher Problem," *Technometrics*, 29, 205–210.

Chernoff, H. (1953), "Locally Optimal Designs for Estimating Parameters," *Annals of Mathematical Statistics*, 24, 586–602.

Dannenberg, O., Dette, H., and Munk, A. (1994), "An Extension of Welch's Approximate *t*-solution to Comparative Bioequivalence Trials," *Biometrika*, 81, 91–101.

Dette, H. (1997), "Designing Experiments With Respect to 'Standardized' Optimality Criteria," *Journal of the Royal Statistical Society*, Ser. B, 59, 97–110.

Dette, H., and Munk, A. (1997), "Optimum Allocation of Treatments for Welch's Test in Equivalence Assessment," *Biometrics*, 53, 1143–1150.

Finney, D. J. (1978), *Statistical Method in Biological Assay* (3rd ed.), London: Charles Griffin & Co.

Govindarajulu, Z. (2000), *Statistical Techniques in Bioassay* (2nd ed.), Basel: Karger.

Imhof, L. A. (2001), "Maximin Designs for Exponential Growth Models and Heteroscedastic Polynomial Models," *The Annals of Statistics*, 29, 561–576.

Scheffé, H. (1970), "Practical Solutions of the Behrens-Fisher Problem," *Journal of the American Statistical Association*, 332, 1501–1508.

Staudte, R. G., and Sheater, S. J. (1990), *Robust Estimation and Testing*, New York: Wiley.

Wang, Y. Y. (1971), "Probabilities of the Type I Errors of the Welch Tests for the Behrens-Fisher Problem," *Journal of the American Statistical Association*, 66, 605–608.

Welch, B. L. (1936), "Specification of Rules for Rejecting the Variable, a Product With Particular Reference to an Electric Lamp Problem," *Journal of the Royal Statistical Society*, Suppl. 3, 29–48.

——— (1938), "The Significance of the Difference Between Means When the Population Variances are Unequal," *Biometrika*, 29, 350–362.