

BACI2: Before-After-Control-Impact (BACI) Power Analysis For Several Related Populations (With Unknown Variance Matrix)

Richard A. Hinrichsen

June 25, 2011

Caveat: This study design tool is for an idealized power analysis built upon several simplifying assumptions (Table 1). For a specific design, a more accurate portrayal of power may require changing these assumptions and the underlying equations. Therefore, this analysis should be treated as a rough guide to power.

Introduction

Currently there are many watershed projects underway in the Columbia Basin to determine the effects of various management actions on salmon survival rate. For example, there are a series of intensively monitored watersheds (IMWs) being established for the purpose of better understanding how salmon respond to approaches to restore habitat. When these projects are developed with a rigorous study design, it may be possible to identify the effectiveness of restoration and other management actions. This analysis was motivated by the need to design these studies so that they have a good chance of detecting significant survival rate changes when they occur. Using this tool can give an investigator a rough idea of the number of years to run an study and determine what statistical power can be expected based on what is known about the variance-covariance structures for survival rate among the salmon populations studied. The framework for the analysis given here, although developed with salmon in mind, fits into the framework of the Before-After-Control-Impact (BACI) design. Such BACI-type designs find application beyond Columbia River salmon survival rate (Osenberg and Schmitt 1996).

An *a priori* power analysis is developed for a BACI-type design aimed at estimating a change in survival rate for several populations. The design includes a Before period where all populations receive no treatment followed by an After period where only the treatment populations receive treatment. This is a generalization of the BACI-type design where the control population and impact population are sampled one time before and one time after the treatment (Green 1979, Osenberg and Schmitt 1996). The assumptions for the analysis are given in Table 1. It is assumed that in the absence of treatment, all populations have a common mean log survival rate. Because this assumption and others may not hold in practice, this analysis should be treated as a rough guide.

The main goal of this work is to demonstrate the probability of detecting an effect on survival rate when several related populations with a common mean survival rate are used in a BACI-type design. This goal is accomplished by describing the study design in a statistically rigorous way, setting up the likelihood function, and then using maximum likelihood theory to estimate power. Power is the probability of rejecting the null hypothesis of “no treatment effect” when it is false. Because theoretical formulas were not available for the estimate of the treatment effect when variance was unknown, Monte Carlo simulations were used to estimate power instead of a formula.

Table 1. —Assumptions used in the power analysis¹.

A1	The observations of log(survival rate) follow a multivariate normal distribution.
A2	There is no serial dependence in log(survival rate).
A3	All populations have a common mean log(survival rate) before treatment.
A4	After treatment, the control populations continue to have the same common mean as exhibited in the Before years, and the treatment populations also have a common mean, but shifted by a constant amount (treatment effect) that is the same for all treatment populations.
A5	The measurement errors in log(survival rate) follow a multivariate normal distribution and the errors are independent of the error due to actual year-to-year environmental variability.
A6	The estimator of the treatment effect is a maximum likelihood estimate.
A7	The variance-covariance matrix representing the error in log(survival rate) is not known and must be estimated.
A8	The variance-covariance matrix takes the form of an intraclass covariance matrix with equal variances and equal covariances.

¹ These are assumptions for an idealized design. For a specific application, a more accurate study design may require changing these assumptions and the underlying equations. Therefore this analysis should be treated as a rough guide to power.

The website www.onefishtwofish.net contains a web-based tool that implements this power analysis with the added assumptions that, in the data generating model, variances in log(survival rate) are equal for all populations and the correlations in log(survival rate) are equal for each population pair. This is the intraclass covariance structure studied by R.A. Fisher (1925). The code for implementing this power analysis,

which may be found in Appendix A, was implemented in R, a system for statistical computation and graphics (Venerables et al. 2010).

Methods

To conduct the power analysis, a model was formulated and maximum likelihood estimators (MLEs) were derived (Mood et al. 1974). These estimators were then used as the basis for testing the null hypothesis of “no treatment” effect using Monte Carlo simulation. Power was then calculated as the probability that the null hypothesis is rejected when it is false.

The model. —It was assumed that mean log(survival rate) before treatment was the same for each population and equal to μ_1 . After treatment, the mean log(survival rate) of the treatment populations shifts by the amount δ for the treatment populations while the control populations continue to have a mean log(survival rate) of μ_1 . It was also assumed that year-to-year variability in log(survival rate) and measurement error followed a multivariate normal distribution with variance $\Sigma = \Sigma_y + \Sigma_m$, where Σ_y is the variance-covariance matrix that describes year-to-year variability in the absence of measurement error, and Σ_m represents the variance-covariance matrix of the measurement error. In this implementation of the BACI design, the variance-covariance matrix, Σ , was assumed to be unknown so that it must be estimated along with other two model parameters (μ_1 and δ).

Maximum likelihood estimators (MLEs). —To derive MLEs, an expression for the likelihood function is needed. For the model described above, the log-likelihood function is

$$l(\theta, \Sigma) = C + (n/2) \ln |\Sigma^{-1}| - (1/2) \sum_{t=1}^{n_1} (\mathbf{x}_t - [\mathbf{e} \quad \mathbf{0}] \theta)' \Sigma^{-1} (\mathbf{x}_t - [\mathbf{e} \quad \mathbf{0}] \theta) \quad (1)$$

$$- (1/2) \sum_{t=n_1+1}^n (\mathbf{x}_t - [\mathbf{e} \quad \mathbf{e}_2] \theta)' \Sigma^{-1} (\mathbf{x}_t - [\mathbf{e} \quad \mathbf{e}_2] \theta);$$

where $\theta = [\mu \quad \delta]'$; Σ is the unknown variance-covariance matrix with equal variances and equal covariances; C is a constant that does not depend on the parameters; n_1 is the number of years prior to treatment; n is the total number of years of the study; \mathbf{x}_t is a k -vector of observed survival rates in year t ; k is the number of populations (treatment + control) used in the study; \mathbf{e} is a k -vector of 1s; \mathbf{e}_2 is a k -vector of k_1 0s followed by k_2

1s, where k_1 represents the number of control populations and k_2 represents the number of treatment populations. The vector \mathbf{x}_i is arranged so that the k_1 control populations precede the k_2 treatment populations.

MLEs for μ_1, δ and Σ are sought. When Σ is estimated all populations have the same variance and all pairs of populations have the same covariance. This gives rise to the intraclass covariance matrix structure studied by Fisher (1925) where all diagonal entries are equal and all off-diagonal entries are equal.

Using maximum likelihood theory, estimating equations for the Before mean, treatment effect, and the covariance matrix are developed. In the case of the before mean and treatment effect, maximizing the likelihood function is equivalent to solving the generalized least squares problem of minimizing

$$SS = \left(\begin{bmatrix} \bar{\mathbf{x}}_1 \\ \bar{\mathbf{x}}_2 \end{bmatrix} - \begin{bmatrix} \mathbf{e} & 0 \\ \mathbf{e} & \mathbf{e}_2 \end{bmatrix} \begin{bmatrix} \mu \\ \delta \end{bmatrix} \right)' \begin{bmatrix} \Sigma^{-1}n_1 & 0 \\ 0 & \Sigma^{-1}n_2 \end{bmatrix} \left(\begin{bmatrix} \bar{\mathbf{x}}_1 \\ \bar{\mathbf{x}}_2 \end{bmatrix} - \begin{bmatrix} \mathbf{e} & 0 \\ \mathbf{e} & \mathbf{e}_2 \end{bmatrix} \begin{bmatrix} \mu \\ \delta \end{bmatrix} \right); \quad (2)$$

where $\bar{\mathbf{x}}_1$ represents the k -vector of sample means of log(survival rate) in the Before period, and $\bar{\mathbf{x}}_2$ represents the k -vector of sample means of log(survival rate) in the After period.

This generalized sum of squares may be written in the familiar form

$$SS = (\mathbf{y} - \mathbf{B}\boldsymbol{\theta})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \mathbf{B}\boldsymbol{\theta}); \quad (3)$$

where $\mathbf{y}' = [\bar{\mathbf{x}}_1' \quad \bar{\mathbf{x}}_2']$, $\mathbf{B} = \begin{bmatrix} \mathbf{e} & 0 \\ \mathbf{e} & \mathbf{e}_2 \end{bmatrix}$, and $\boldsymbol{\Omega}^{-1} = \begin{bmatrix} \Sigma^{-1}n_1 & 0 \\ 0 & \Sigma^{-1}n_2 \end{bmatrix}$. In this form, the generalized least squares solution, the called the Aitken estimator (Press 2005), is known to be

$$\hat{\boldsymbol{\theta}} = (\mathbf{B}^T \boldsymbol{\Omega}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \boldsymbol{\Omega}^{-1} \mathbf{y}. \quad (4)$$

After considerable matrix algebra, we may write

$$\hat{\boldsymbol{\theta}} = \begin{bmatrix} \hat{\mu} \\ \hat{\delta} \end{bmatrix} = \frac{\begin{bmatrix} (\mathbf{e}_2' \Sigma^{-1} \mathbf{e}_2)(\mathbf{e}' \Sigma^{-1} \bar{\mathbf{x}}) - (\frac{n_2}{n})(\mathbf{e}' \Sigma^{-1} \mathbf{e}_2)(\mathbf{e}_2' \Sigma^{-1} \bar{\mathbf{x}}_2) \\ (\mathbf{e}' \Sigma^{-1} \mathbf{e})(\mathbf{e}_2' \Sigma^{-1} \bar{\mathbf{x}}_2) - (\mathbf{e}_2' \Sigma^{-1} \mathbf{e})(\mathbf{e}' \Sigma^{-1} \bar{\mathbf{x}}) \end{bmatrix}}{(\mathbf{e}' \Sigma^{-1} \mathbf{e})(\mathbf{e}_2' \Sigma^{-1} \mathbf{e}_2) - (\frac{n_2}{n})(\mathbf{e}_2' \Sigma^{-1} \mathbf{e})^2}; \quad (5)$$

where $\bar{\mathbf{x}}$ is a k -vector representing population-specific sample means over the entire duration of the study. Also well known is the conditional variance of the estimate of $\boldsymbol{\theta}$ (given the variance-covariance matrix):

$$\text{var} \hat{\boldsymbol{\theta}} | \Sigma = (\mathbf{B}^T \boldsymbol{\Omega}^{-1} \mathbf{B})^{-1} = \frac{\begin{bmatrix} n_2 \mathbf{e}_2' \Sigma^{-1} \mathbf{e}_2 & -n_2 \mathbf{e}_2' \Sigma^{-1} \mathbf{e} \\ -n_2 \mathbf{e}_2' \Sigma^{-1} \mathbf{e} & n \mathbf{e}' \Sigma^{-1} \mathbf{e} \end{bmatrix}}{nn_2(\mathbf{e}' \Sigma^{-1} \mathbf{e})(\mathbf{e}_2' \Sigma^{-1} \mathbf{e}_2) - (n_2 \mathbf{e}_2' \Sigma^{-1} \mathbf{e})^2}. \quad (6)$$

Next, the estimating equation for the variance covariance matrix is derived. To do this, the partial derivatives of the likelihood function are calculated with respect to the inverse of the covariance matrix (Σ^{-1}), and set to zero. The variance matrix has the form of an intraclass covariance matrix. In this case, the inverse of the variance-covariance matrix also has an intraclass covariance matrix structure and may be written as

$$\Sigma^{-1} = (a - b)\mathbf{I} + b\mathbf{e}\mathbf{e}'. \quad (7)$$

Note that all of the diagonal entries of the inverse covariance matrix are equal to the scalar quantity a , and all of the off-diagonal entries are equal to the scalar quantity b . In this special case, the log-likelihood function may be written as

$$l(\mathbf{\Sigma}) = C + (n/2) \ln |\mathbf{\Sigma}^{-1}| - (1/2) \sum_{t=1}^n \mathbf{z}_t' \mathbf{\Sigma}^{-1} \mathbf{z}_t. \quad (8)$$

where $\mathbf{z}_t = \mathbf{x}_t - [\mathbf{e} \quad \mathbf{0}]' \boldsymbol{\theta}$ when $t \leq n_1$ and $\mathbf{z}_t = \mathbf{x}_t - [\mathbf{e} \quad \mathbf{e}_2]' \boldsymbol{\theta}$ when $t > n_1$.

The following formulas are used for calculating the partial derivatives of this log-likelihood function

$$\frac{\partial \ln |\mathbf{\Sigma}^{-1}|}{\partial a} = k \sigma_{11} \text{ and } \frac{\partial \ln |\mathbf{\Sigma}^{-1}|}{\partial b} = k(k-1) \sigma_{12}, \quad (9)$$

where σ_{11} is the common variance term in the variance matrix and σ_{12} is the common covariance value. Also used are the formulas

$$\frac{\partial \mathbf{z}' \mathbf{\Sigma}^{-1} \mathbf{z}}{\partial a} = \mathbf{z}' \mathbf{z} \text{ and } \frac{\partial \mathbf{z}' \mathbf{\Sigma}^{-1} \mathbf{z}}{\partial b} = \mathbf{e}' \mathbf{z} \mathbf{z}' \mathbf{e} - \mathbf{z}' \mathbf{z} \quad (10)$$

Armed with these equations, it may be shown that

$$\frac{\partial l(\mathbf{\Sigma}^{-1})}{\partial a} = (n/2) k \sigma_{11} - (1/2) \sum_{t=1}^n \mathbf{z}_t' \mathbf{z}_t, \quad (11)$$

and

$$\frac{\partial l(\Sigma^{-1})}{\partial b} = (n/2)k(k-1)\sigma_{12} - (1/2)\sum_{t=1}^n \mathbf{e}'\mathbf{z}_t\mathbf{z}_t'\mathbf{e} - \mathbf{z}_t'\mathbf{z}_t. \quad (12)$$

Setting these two partial derivatives equal to zero yields the estimating equations

$$\hat{\sigma}_{11} = \frac{\sum_{t=1}^n \mathbf{z}_t'\mathbf{z}_t}{nk} \quad (13)$$

and

$$\hat{\sigma}_{12} = \frac{\sum_{t=1}^n \mathbf{e}'\mathbf{z}_t\mathbf{z}_t'\mathbf{e} - \mathbf{z}_t'\mathbf{z}_t}{nk(k-1)}. \quad (14)$$

The MLE for the intraclass covariance matrix is therefore equal to

$$\hat{\Sigma} = (\hat{\sigma}_{11} - \hat{\sigma}_{12})\mathbf{I} + \hat{\sigma}_{12}\mathbf{e}\mathbf{e}'. \quad (15)$$

MLE numerical algorithm.—Armed with the estimating equations an algorithm to solve them for the MLEs is now derived. An iterative procedure is used because the MLE of $\boldsymbol{\theta}$ depends on the MLE of Σ in equation (5). The algorithm is based on a procedure called iteratively reweighted least squares or IRLS, which is a special case of iterative estimating equations (IEE), with known convergence properties (Jiang et al. 2007). In practice, the method will converge quickly if the number of observations is sufficiently greater than the number of estimated parameters. The total number of estimated model parameters is always 4: a control mean, a treatment effect, a common variance, and a common covariance parameter.

The iterative procedure begins by setting the initial estimate of the variance matrix, call it $\hat{\Sigma}^{(0)}$, equal to the identity matrix. The next step is to make an initial

estimate of the θ vector. This is accomplished by using $\hat{\Sigma}^{(0)}$ in place of Σ in equation (6) and solving for $\hat{\theta}^{(0)}$. The estimate $\hat{\theta}^{(0)}$ is then used in equations (13)-(15) to get an updated estimate of the variance matrix, $\hat{\Sigma}^{(1)}$. This entire procedure is repeated with the most recent updates of the parameters until the likelihood function fails to decrease by some specified tolerance.

Statistical power calculations.—Statistical power is estimated with a Monte Carlo procedure. Nominal values of the model parameters are specified, and then Monte Carlo replications of the MLEs are constructed. The test statistic used in this power analysis is

$$\hat{T} = \hat{\delta} / se(\hat{\delta}), \quad (16)$$

where $se(\hat{\delta})$ is the square root of the variance of $\hat{\delta}$ obtained by substituting $\hat{\Sigma}$ for Σ in equation (6). This statistic is known to have a student's t -distribution in the case of a linear regression with uniform variance and uncorrelated errors, which is a special case of the model considered here. A theoretical distribution for the test statistic was not assumed. Instead, Monte Carlo simulation was used to derive an estimate of its distribution under the null hypothesis of “no treatment effect.” This was then used to estimate power. An alternative to this test statistic would be based on the likelihood ratio (Mood et al. 1974), which would be compared to a chi-square distribution with one degree of freedom.

Statistical power was estimated by the following 4-step procedure:

Step 1. First obtain $Nsim$ replications of the test statistic (defined in equation 16) when the true treatment effect is δ . These replications are expressed as

$\hat{T}_1^*, T_2^*, \dots, T_{Nsim}^*$, where $\hat{T}_i^* = \hat{\delta}_i^* / se(\hat{\delta}_i^*)$ where $\hat{\delta}_i^*$ is the i th replication of the treatment effect estimate, and $se(\hat{\delta}_i^*)$ is calculated by substituting $\hat{\Sigma}_i^*$ for Σ in equation (6).

Step 2. Use these replications of the test statistic to build a set of replications of the test statistic under the null hypothesis. To do this, define the replications under the null hypothesis as $\hat{T}_i^* - \delta / se(\hat{\delta}_i^*)$

Step 3. Assuming that the probability of a Type I error is set to the value α , the critical value is then estimated as the average of the values $|q_{\alpha/2}^*|$ and $q_{1-\alpha/2}^*$,

where $q_{\alpha/2}^*$ represents the $\alpha/2$ quantile of the replications of the treatment effect estimate in Step 2 and $q_{1-\alpha/2}^*$ represents the $1-\alpha/2$ quantile.

Step 4. Power is estimated as the fraction of the Monte Carlo replications of the test statistic (generated in Step 1) whose absolute value exceeds the critical value calculated in Step 3.

In Step 2, a short cut was used: $Nsim$ replications from the null distribution of the test statistic were generated by subtracting $\delta / se(\hat{\delta}_i^*)$ from replications generated in Step 1. This shortcut was possible because if \mathbf{x} is a random sample of log survival rates during a single year from the After period in Step 1 with true treatment effect of δ , then $\mathbf{x} - \mathbf{e}_2\delta$ is a random sample during a single year of the After period when the null hypothesis is assumed (true treatment effect of zero). In other words, the statistic $\hat{\delta}$ is a location statistic, for which increasing log(survival rate) of treated populations during the After period by δ increases the statistic $\hat{\delta}$ itself by δ . When a location statistic is used, the t-like statistic in equation (16) is known to be an appropriate test statistic (Efron and Tibshirani 1993).

Investigators often choose a design such that power of 0.8 is achieved (e.g, Peterman and Bradford 1987, Liermann and Roni 2008).

Standard error and coefficient of variation. — The standard error of $\hat{\delta}$ was estimated as the square root of the sample variance of the $Nsim$ replications of $\hat{\delta}$:

$$se^*(\hat{\delta}) = \sqrt{\sum_{i=1}^{Nsim} (\hat{\delta}_i^* - \bar{\hat{\delta}}^*)^2 / (Nsim - 1)} . \quad (17)$$

where $\bar{\hat{\delta}}^*$ was the sample mean of the $Nsim$ replications of $\hat{\delta}$. The Monte Carlo estimate of the coefficient of variation was estimated as

$$CV^*(\hat{\delta}) = se^*(\hat{\delta}) / \delta . \quad (18)$$

Example. — Consider the case where the number of treatment populations equals the number of control populations ($k_1 = k_2 = 2$); the number of before years equals the number of after years $n_1 = n_2 = 10$ or, alternatively, $n_1 = n_2 = 5$; the generating model

uses a common variance of 1.0 and a common covariance of 0.5; measurement error is zero; and the probability of a type I error is $\alpha = 0.05$. Let the treatment effect vary from 0 to $\log(2)$ and compare the power obtained when the variance matrix is known to the power obtained when the variance matrix is estimated. The results of this exercise may be found in Table 2. Notice that the power calculations in the case when variance is treated as known and in the case when it is treated as unknown are close, even when the total years of study decreases from $n=20$ to $n=10$. Notice also how power declines as the total number of years of study declines.

Table 2. —Statistical power under two different estimation assumptions: known variance, estimated variance. In both cases, the variance matrix is assumed to have the intraclass covariance matrix structure. In this exercise, $k_1 = k_2 = 2$ and simulations are run with variance of 1.0 and a correlation of 0.5 and a measurement error of zero. The probability of a type I error was set to 0.05. The true treatment effect (δ) varies from 0 to $\log(2)$. A total of 10,000 Monte Carlo simulations were used when Σ was estimated.

delta	$n_1=n_2=10$		$n_1=n_2=5$	
	known Σ	estimated Σ	known Σ	estimated Σ
0.0000	0.05	0.05	0.05	0.05
0.0365	0.05	0.06	0.05	0.05
0.0730	0.06	0.06	0.06	0.06
0.1094	0.08	0.09	0.07	0.07
0.1459	0.11	0.13	0.08	0.08
0.1824	0.14	0.14	0.09	0.09
0.2189	0.18	0.16	0.11	0.11
0.2554	0.22	0.23	0.14	0.13
0.2919	0.28	0.28	0.16	0.16
0.3283	0.34	0.34	0.19	0.18
0.3648	0.40	0.37	0.23	0.22
0.4013	0.47	0.42	0.27	0.25
0.4378	0.54	0.50	0.31	0.29
0.4743	0.60	0.62	0.35	0.33
0.5107	0.67	0.68	0.40	0.37
0.5472	0.73	0.72	0.44	0.41
0.5837	0.78	0.76	0.49	0.45
0.6202	0.83	0.84	0.54	0.50
0.6567	0.87	0.83	0.59	0.55
0.6931	0.90	0.89	0.63	0.60

Acknowledgements

This work was supported by Bonneville Power Administration. Thanks to Charlie Paulsen, Rishi Sharma, Tracy Hillman, and Amber Parsons for their valuable reviews. Thanks to Brian Maschhoff for implementing this analysis as a web tool at www.onefishtwofish.net.

References

- Dwyer, P.S. 1967. Applications of matrix derivatives in multivariate analysis. *Journal of the American Statistical Association* 62: 607-625.
- Efron, B. and R.J. Tibshirani. 1993. An introduction to the bootstrap. *Monographs on Statistics and Applied Probability* 57. Chapman & Hall/CRC, New York, New York.
- Fisher, R.A. 1925. *Statistical Methods for Research Workers*. Oliver and Boyd, Edinburgh, Scotland.
- Green, R.H. 1979. *Sampling design and statistical methods for environmental biologists*. Wiley and Sons, New York, New York.
- Jiang, J. Luan, Y., and Y. Wang. 2007. Iterative estimating equations: linear convergence and asymptotic properties. *The Annals of Statistics* 35:2233-2260.
- Liermann, M., and P. Roni. 2008. More sites or more years? Optimal study design for monitoring fish response to watershed restoration. *North American Journal of Fisheries Management*. 28:935-943.
- Mood, A.M, Graybill, F.A., and D.C. Boes. 1974. *Introduction to the theory of statistics*, Third Edition. McGraw-Hill, New York, New York.
- Osenberg, C.W. and R.J. Schmidt. 1996. Detecting ecological impacts caused by human activities. In *Detecting Ecological Impacts: Concepts and Applications in Coastal Habitats*, R.J Schmitt and C.W. Osenberg, Editors. Academic Press, New York, New York.
- Peterman, W. M., and M. J. Bradford. 1987. Statistical power of trends in fish abundance. *Canadian Journal of Fisheries and Aquatic Sciences* 44: 1879:1889.
- Press, S.J. 2005. *Applied multivariate analysis: using Bayesian and frequentist methods of inference*. Dover Publications, Inc., Mineola, New York.
- Venerables, W.N., Smith, D.M., and R Development Core Team. 2010. *An Introduction to R. Notes on R: A Programming Environment for Data Analysis and Graphics Version 2.11.1* (2010-05-31). <http://www.r-project.org/>

Appendix A. R code used to calculate statistical power for the BACI-type design when the variance-covariance matrix is estimated

```

#Program to estimate power of a baci design
#when the variance-covariance matrix is unknown. Variance is estimated
#along with the treatment effect and the control population mean
#The estimated variance-covariance matrix has the
#form of an intraclass covariance matrix (a common variance and a common
#covariance)

#Baci code using Monte Carlo estimates of power
#Nsim number of Monte Carlo simulations
#s2 is variances (assumed equal for all populations)
#rho is correlation between each pair of populations
#n1 number of before years
#n2 number of after years
#k1 number of control populations
#k2 number of treatment populations
#me measurement error
#alpha probability of a type I error (rejecting null hypothesis when true)
#delta true treatment effect representing difference in natural log survival rate,
#ln(Streatment/Scontrol)
#the intraclass covariance matrix structure is assumed.
library(MASS)

baci2<-
function(Nsim=1000,s2=1,rho=.9,n1=5,n2=5,k1=1,k2=1,me=log(1.10),alpha=0.05,
delta=log(1.50)){
  k<-k1+k2
  SIG2<-matrix(s2*rho,ncol=k,nrow=k)
  diag(SIG2)<-s2+me*me
  INVSIG2<-solve(SIG2)

  deltas<-rep(NA,Nsim)
  ses<-rep(NA,Nsim)
  #Do Monte Carlo simulations to get replications of delta and se
  for(ii in 1:Nsim){
    bres<-baci.estimates(s2=s2,rho=rho,n1=n1,n2=n2,k1=k1,k2=k2,me=me,
alpha=alpha,delta=delta)
    if(!is.null(bres)){
      deltas[ii]<-bres$par[2]
      ses[ii]<-get.se(bres$SIG2,n1=n1,n2=n2,k1=k1)
    }
  }
}

```

```

    }
  }
  se<-sqrt(var(deltas,na.rm=T))
  #get critical value of distribution under null hypothesis
  crit2<-quantile(x=(deltas-delta)/ses,probs=c(alpha/2,1-alpha/2),na.rm=T)
  crit<-mean(abs(crit2))
  powerx<-abs(deltas/ses)>crit
  power<-mean(powerx,na.rm=T)
  ngood<-sum(!is.na(deltas/ses))
  return(list(Nsim=Nsim,ngood=ngood,s2=s2,rho=rho,n1=n1,n2=n2,k1=k1,k2=k2,me=me,
    alpha=alpha,delta=delta,se=se,cv=se/delta,power=power))
}
#outputs
#ngood is the number of simulations that result in a valid estimate
#se standard error
#cv coefficient of variation
#power is the probability of rejecting the null hypothesis of no effect when it is false

baci.estimates<-
function(s2=1,rho=.9,n1=5,n2=5,k1=1,k2=1,me=log(1.10),alpha=0.05,delta=log(1.50)){
  n<-n1+n2
  k<-k1+k2
  par<-c(log(.10),delta)
  SIG2<-matrix(s2*rho,ncol=k,nrow=k)
  diag(SIG2)<- s2+me*me
  xmat1<-mvrnorm(n=n1,mu=rep(par[1],k),Sigma=SIG2)
  xmat2<-mvrnorm(n=n2,mu=c(rep(par[1],k1),rep(par[1]+par[2],k2)),Sigma=SIG2)
  xmat<-cbind(t(xmat1),t(xmat2))
  res2<-myoptim2(xmat=xmat,s2=s2,rho=rho,me=me,n1=n1,k1=k1)
  return(res2)
}

#Iterate until maximum likelihood estimates are obtained
#solving the estimating equations which were
#determined by setting the partial derivatives of the
#likelihood function to zero.
myoptim2<-function(xmat,s2,rho,me,n1,k1){
  k<-dim(xmat)[1]
  n<-dim(xmat)[2]
  #begin with OLS regression estimates of theta parameters
  SIG2<-diag(1,k)
  par<-get.pars(xmat,SIG2,n1,k1)
  SIG2<-get.SIG2(par,xmat,n1,k1)
  #check condition number of SIG2

```

```

cn<-kappa(SIG2)
if((1/cn)<=1.e-15){
  warning("SIG2 is computationally singular in myoptim2")
  return(NULL)}
lf1<-lf(par=par,x=xmat,n1=n1,k1=k1,SIG2)
etol<-1.e-5
err<-2.*etol*(abs(lf1)+etol)
iter<-0
#look for relative likelihood function convergence
while(err>etol*(abs(lf1)+etol)){
  par<-get.pars(xmat,SIG2,n1,k1)
  SIG2<-get.SIG2(par,xmat,n1,k1)
#check condition number of SIG2
  cn<-kappa(SIG2)
  if((1/cn)<=1.e-15){
    warning("SIG2 is computationally singular in myoptim2")
    return(NULL)}
  lf2<-lf(par=par,x=xmat,n1=n1,k1=k1,SIG2)
  err<-abs(lf2-lf1)
  lf1<-lf2
  iter<-iter+1
  if(iter>1000){
    warning("too many iterations in myoptim2")
    return(NULL)}
}

return(list(par=par,SIG2=SIG2))
}

#Use estimating equations to solve for parameter values
#Returns control mean (mu) and treatment effect (delta)
get.pars<-function(xmat,SIG2,n1,k1){
  n<-dim(xmat)[2]
  n2<-n-n1
  k<-dim(xmat)[1]
  E<-rep(1,k)
  E1<-c(rep(1,k1),rep(0,k-k1))
  E2<-c(rep(0,k1),rep(1,k-k1))
  xbar2<-apply(xmat[(n1+1):n],c(1),mean)
  xbar<-apply(xmat,c(1),mean)
  INVSIG2<-solve(SIG2)
  delta<-(t(E2)%*%INVSIG2%*%xbar2)*(t(E)%*%INVSIG2%*%E)-
t(E)%*%INVSIG2%*%xbar*(t(E2)%*%INVSIG2%*%E)
  den<-(t(E)%*%INVSIG2%*%E)*(t(E2)%*%INVSIG2%*%E2)-
(n2/n)*(t(E2)%*%INVSIG2%*%E)^2

```



```

delta<-delta/den
mu<-t(E2)%*%INVSIG2%*%xbar2-(t(E2)%*%INVSIG2%*%E2)*delta
den<-t(E2)%*%INVSIG2%*%E
mu<-mu/den
return(c(mu,delta))
}

```

```

#log-likelihood function
lf<-function(par,x,n1,k1,SIG2){
  INVSIG2<-solve(SIG2)
  n<-dim(x)[2]
  k<-dim(x)[1]
  z<-x
  like<-k*.5*n*log(2*pi)-.5*n*log(det(SIG2))
  for(ii in 1:n1){
    z[,ii]<-x[,ii]-rep(par[1],k)
    like<-like-.5*t(z[,ii])%*%INVSIG2%*%z[,ii]
  }
  for(ii in (n1+1):n){
    z[,ii]<-x[,ii]-c(rep(par[1],k1),rep(par[1]+par[2],k-k1))
    like<-like-.5*t(z[,ii])%*%INVSIG2%*%z[,ii]
  }
  return(like)
}

```

```

#Get the estimate variance-covariance matrix
#This is based on the estimating equations
#variance is unknown and has the
#form of an intraclass covariance matrix
get.SIG2<-function(par,x,n1,k1){
  iform<-1
  n<-dim(x)[2]
  k<-dim(x)[1]
  z<-x
  SIG2<-matrix(0,ncol=k,nrow=k)
  for(ii in 1:n1){
    z[,ii]<-x[,ii]-rep(par[1],k)
    SIG2<-SIG2+z[,ii]%*%t(z[,ii])/n
  }
  for(ii in (n1+1):n){
    z[,ii]<-x[,ii]-c(rep(par[1],k1),rep(par[1]+par[2],k-k1))
    SIG2<-SIG2+z[,ii]%*%t(z[,ii])/n
  }

  if(iform==1){

```

```

s2<-mean(diag(SIG2))
s12<-(sum(SIG2)-sum(diag(SIG2)))/(k*k-k)
SIG2<-matrix(s12,ncol=k,nrow=k)
diag(SIG2)<-s2
}
return(SIG2)
}

```

```

#Return se estimate based on SIG2
#This is a theoretical formula
#derived in the paper
get.se<-function(SIG2,n1=5,n2=5,k1=1){
  k<-dim(SIG2)[1]
  k2<-k-k1
  INVSIG2<-solve(SIG2)
  e<-rep(1,k)
  se<-(n1+n2)*t(e)%*%INVSIG2%*%e
  e1<-c(rep(1,k1),rep(0,k2))
  e2<-c(rep(0,k1),rep(1,k2))
  det<-n1*t(e)%*%INVSIG2%*%e+n2*t(e1)%*%INVSIG2%*%e1
  det<-det*n2*t(e2)%*%INVSIG2%*%e2-n2*n2*(t(e2)%*%INVSIG2%*%e1)^2
  se<-sqrt(se/det)
  return(se)
}

```