

.....

# Power Analysis for Logistic Regression Models Fit to Clustered Data: Choosing the Right Rho

.....

Steve Gregorich  
May 16, 2014

# Abstract

## Context

Power analyses for logistic regression models fit to clustered data

## Approach

- . estimate *effective* sample size ( $N_{\text{eff}}$ : cluster-adjusted total sample size)
- . input  $N_{\text{eff}}$  into standard power analysis routines for independent obs.

## Wrinkle

- . in the context of logistic regression there are two general approaches to estimating the intra-cluster correlation of  $Y$ :
  - . phi-type coefficient and
  - . tetrachoric-type coefficient.

## Resolution

- . The phi-type coefficient should be used when calculating  $N_{\text{eff}}$

I will present background on this topic as well as some simulation results

# Simple random sampling (SRS)

- . Fully random selection of participants  
e.g., start with a list, select  $N$  units at random
- . Some key features wrt statistical inference:
  - representativeness
  - all units have equal probability of selection
  - all sampled units can be considered to be independent of one another
- . SRS with replacement versus without replacement

# Clustered sampling

- . Rnd sample of  $m$  clusters; rnd sample of  $n$  units w/in each cluster  
multi-stage area sampling  
patients within clinics
- . Repeated measures  
Random sample of  $m$  respondents;  $n$  repeated measures are taken  
repeated measures are clustered within respondents
- . Typically, elements within the same cluster are more similar to each other than elements from different clusters
- . The  $n$  units w/in a cluster usually do not contain the same amount of info wrt some parameter,  $\theta$ , as the same number of units in an SRS sample  
...the concept of effective sample size,  $N_{\text{eff}}$ ...

Therefore, it is usually true that  $\sigma_{\text{clus}}^2(\hat{\theta}) \neq \sigma_{\text{srs}}^2(\hat{\theta})$

# Two-stage clustered sampling design

*Unless otherwise noted, I assume*

- . Clustered sampling of  $m$  clusters, each with  $n$  units:

$$N = m \times n$$

- . Normally distributed unit-standardized  $x$ , binary  $y$   
exchangeable / compound symmetric correlation structure
  - $\rho_y > 0$ : intra-cluster correlation of  $y$  (outcome) response
  - $\rho_x = 0$  or  $1$ : intra-cluster correlation of  $x$  (explanatory var) response
- . Regression of  $y$  onto  $x$  via
  - . a mixed logistic model with random cluster intercepts or
  - . a GEE logistic model
- . Common effects of  $x$  across clusters, i.e., no random slopes for  $x$
- . Common between- and within-cluster effects of  $x$

## The design effect, $deff$

- $deff$  can be thought of as a design-attributable multiplicative change in variation that results from choice of a clustered sampling versus an SRS design

$$\widehat{deff} = \frac{\sigma_{\text{clus}}^2(\hat{\theta})}{\sigma_{\text{srs}}^2(\hat{\theta})} \quad \text{and} \quad \widehat{N}_{\text{eff}} = \frac{N}{\widehat{deff}}, \quad \text{where}$$

$\sigma_{\text{clus}}^2(\hat{\theta})$  is the estimated parameter variation given a clustered sampling design;

$\sigma_{\text{srs}}^2(\hat{\theta})$  is the estimated parameter variation given a SRS design;

$N$  is the common size of the SRS and clustered ( $N=m \times n$ ) samples;

$\widehat{N}_{\text{eff}}$  estimated effective size of the clustered sample wrt information about  $\hat{\theta}$ , relative to what would have been obtained with a SRS of size  $N$

*Assumes compound symmetric covariance structure of the response*

# The misspecification effect, $m_{eff}$

Conceptually similar to  $d_{eff}$  except that the multiplicative change corresponds to the effect of correctly modeling the clustering of observations versus ignoring the cluster structure

$$\widehat{m_{eff}} = \frac{\sigma_{clus}^2(\hat{\theta})}{\sigma_{\cancel{clus}}^2(\hat{\theta})} \quad \text{and} \quad \widehat{N}_{eff} = \frac{N}{\widehat{m_{eff}}}, \quad \text{where}$$

$\sigma_{clus}^2(\hat{\theta})$  is the estimated parameter variation given clustered responses;

$\sigma_{\cancel{clus}}^2(\hat{\theta})$  is the estimated parameter variation ignoring clustering of responses;

$N$  is the total size of the clustered sample;

$\widehat{N}_{eff}$  is the effective size of the clustered sample wrt information about  $\hat{\theta}$ , relative to what would have been obtained with a SRS of the same size

*Assumes compound symmetric covariance structure of the response*

## *deff*, *meff*, and the sample size ratio

A 'context free' label for *deff* and *meff* is the sample size ratio, SSR

$$SSR = \frac{N}{\hat{N}_{\text{eff}}}$$

- . *deff*, *meff*, and SSR have equivalent meaning wrt power analysis, but *deff* and *meff* are conceptually distinct
- . *deff* assumes that you are considering SRS versus clustered sampling
- . *meff* assumes that you have chosen a clustered sampling design and want to make adjustments to an analysis that assumed SRS
- . I will use *meff* for this talk

# Estimating $m_{eff}$ via the intra-cluster correlation

- . Given positive intra-cluster correlation of  $y$ :  $\rho_y > 0$ ,  
the  $m_{eff}$  estimator depends on  $\rho_x$

#1. Level-2 (cluster-level)  $x$  variables will have zero *within*-cluster variation and  $\rho_x = 1$

$$\rho = \frac{\sigma_{btw}^2}{(\sigma_{btw}^2 + \sigma_{w/in}^2)}$$

- . In this case

$$\widehat{m_{eff}} = \frac{\sigma_{clus}^2(\hat{\theta})}{\sigma_{clus}^2(\hat{\theta})} = \frac{N}{N_{eff}} = 1 + (n - 1)\rho_y,$$

- . note: when estimating  $\bar{y}$ , assume  $\rho_x = 1$

# Estimating $m_{eff}$ via the intra-cluster correlation

#2. Consider a level-1 stochastic  $x$  variable with positive within-cluster variation and zero between-cluster variation:  $\rho_x = 0$ :

$$\rho = \frac{\sigma_{btw}^2}{(\sigma_{btw}^2 + \sigma_{w/in}^2)}$$

. In this case

$$\widehat{m_{eff}} = \frac{\sigma_{clus}^2(\hat{\theta})}{\sigma_{clus}^2(\hat{\theta})} = \frac{N}{N_{eff}} \approx 1 - \rho_y^{(n/(n-1))}$$

note:  $n/(n-1) \rightarrow 1$  as  $n \rightarrow \infty$

(In this talk, I do not cover Level-1  $x$  variables with  $0 < \rho_x < 1$ )

# Power analysis for clustered sampling designs using $m_{eff}$ :

## Option 1

*Option 1. Given a chosen model, power, and alpha level, plus a proposed clustered sample of size  $N=m \times n$ , and a  $m_{eff}$  estimate*

$$\cdot \widehat{N}_{eff} = \frac{N}{\widehat{m_{eff}}}$$

. Use standard power analysis software, plug in  $\widehat{N}_{eff}$  (instead of  $N$ ), and estimate

# Power analysis for clustered sampling designs using *meff*: Option 1 Example

## Estimate Power by Simulation

- . Simulate data from a CRT with 100 clusters ( $j$ ) and 30 individuals/cluster ( $i$ )

$$y_{ij} = \text{group}_j \mathbf{0.5} + u_j + e_{ij}$$

where,  $\text{VAR}(u_j) = \text{VAR}(e_{ij}) = 1$ ,  
 $\text{VAR}(u_j) + \text{VAR}(e_{ij}) = \mathbf{2}$ , and  
 $\rho_y = \sigma_u^2 / (\sigma_u^2 + \sigma_e^2) = \mathbf{0.50}$

needed later for PASS

- . Linear mixed model results from analysis of 2000 replicate samples

- .  $\rho_y = \mathbf{0.501}$
- . residual std dev =  $\mathbf{1.416} \approx \sqrt{2}$
- .  $\hat{b}_{\text{group}} = \mathbf{0.495}$
- . simulated power for group effect: 67.7%

all relatively  
unbiased

# Power analysis for clustered sampling designs using *meff*: Option 1 Example

- . Simulation result: power = 67.7%
- . Use PASS Linear Regression routine to solve for power
  - .  $\widehat{meff} = 1 + (30 - 1) \times \mathbf{0.501} = 15.529$
  - .  $\widehat{N}_{eff} = 100 \times 30 \div 15.529 \approx 193$
  - . specify 193 as *N* in PASS
  - . specify H<sub>1</sub> slope = **0.495**
  - . specify Residual Std Dev = **1.416** (resid. @ level-1 plus level-2)
- . PASS result: power = 67.6%

## Summary

- . choose *meff* estimator and estimate *meff*
- . estimate  $N_{eff}$
- . plug  $N_{eff}$  into power analysis software (w/ other parameters)
- . estimate power

# Power analysis for clustered sampling designs using *meff*: Option 1 Example

**Linear Regression**

File View Run Procedures Tools Window Help

Reset Open Save As

**Calculate**

**Design**

Solve For: Power

**Test**

Alternative Hypothesis: Two-Sided

**Alpha**

Alpha: 0.05

**Sample Size**

N (Sample Size): 193

**Effect Size**

Slope

B0 (Slope|H0): 0.0

B1 (Slope|H1): 0.495

Standard Deviation of X's

SX (Standard Deviation of X's): xs 0, 1

Residual Variance Calculation

Residual Variance Method: S (Std. Dev. of Residuals)

S (Standard Deviation of Residuals): 1.416

# Power analysis for clustered sampling designs using *meff*: Option 1 Example

**PASS Output**

File View Edit Window Help

Save As Print Copy Find Add Output to Gallery Auto Add

Navigation Pane

- Linear Regression Power Analysis
  - Numeric Results
  - References
  - Report Definitions
  - Summary Statements
  - Chart Section

5/12/2014 9:30:13 PM 1

**Linear Regression Power Analysis**  
**Numeric Results for Two-Sided Testing of B = B0 where B0 = 0.000**

	Sample Size (N)	Slope (B)	Standard Deviation of X (SX)	Standard Deviation of Residuals (S)	Alpha	Beta
Power	193	0.495	0.500	1.416	0.05000	0.32417

**References**  
 Neter, J., Wasserman, W., and Kutner, M. 1983. Applied Linear Regression Models. Richard D. Irwin, Inc. Chicago, Illinois.

**Report Definitions**  
 Power is the probability of rejecting a false null hypothesis. It should be close to one.  
 N is the size of the sample drawn from the population. To conserve resources, it should be small.  
 B0 is the slope under the null hypothesis.  
 B is the slope at which the power is calculated.  
 SX is the standard deviation of the X values.  
 S is the standard deviation of the residuals.  
 Alpha is the probability of rejecting a true null hypothesis. It should be small.  
 Beta is the probability of accepting a false null hypothesis. It should be small.

**Summary Statements**  
 A sample size of 193 achieves 68% power to detect a change in slope from 0.000 under the null hypothesis to 0.495 under the alternative hypothesis when the standard deviation of the X's is 0.500, the standard deviation of the residuals is 1.416, and the two-sided significance level is 0.05000.

PASS: power = 67.6%

Simulation: power = 67.7%

# Power analysis for clustered sampling designs using *meff*:

## Option 2 example

*Option 2. Given a clustered sample design, chosen model, power, and alpha level, plus an effect size estimate and a meff estimate*

. Use standard power analysis software to estimate required sample size assuming independent observations, i.e.,  $N_{\text{eff}}$ . Then estimate required  $N$

$$. \hat{N} = \hat{N}_{\text{eff}} \times \widehat{meff}$$

### Option 2: Step 1

Start with...

- . the group effect (b=**0.495**),
- . a residual standard deviation of **1.416**,
- . and power equal to 67.6%,

. Use PASS to estimate the required effective sample size,  $\hat{N}_{\text{eff}} = 193$

# Power analysis for clustered sampling designs using *meff*: Option 2 example

**Linear Regression**

File View Run Procedures Tools Window Help

Reset Open Save As

**Calculate**

**Design**

Design

Reports

Plots

**Solve For:** Sample Size

**Test**

Alternative Hypothesis: Two-Sided

**Power and Alpha**

Power: .676

Alpha: 0.05

**Effect Size**

**Slope**

B0 (Slope|H0): 0.0

B1 (Slope|H1): 0.495

**Standard Deviation of X's**

SX (Standard Deviation of X's): xs 0, 1

**Residual Variance Calculation**

Residual Variance Method: S (Std. Dev. of Residuals)

S (Standard Deviation of Residuals): 1.414

Result:  $\widehat{N}_{\text{eff}} = 193$

# Power analysis for clustered sampling designs using *meff*:

## Option 2 example

### Option 2: Step 2

- . Given  $\widehat{N}_{\text{eff}} = 193$ , clusters of size  $n=30$ , and  $\rho_y = 0.501$ ,  
adjust  $\widehat{N}_{\text{eff}} = 193$  to obtain the required needed sample size
  - . for a CRT,  $\rho_x = 1$  and  $\widehat{meff} = 1 + (n - 1)\rho_y$
  - .  $\widehat{N} = 193 \times [1 + (30 - 1) \times 0.501] \approx 3000$
- . Given clusters of size  $n=30$ ,  $\widehat{N}=3000$  suggests that  
100 clusters need to be sampled and randomized (i.e.,  $3000 \div 30$ )

*This example used the linear mixed models framework.*

*Now onto the models for clustered data with binary outcomes.*

# Logistic Regression Models Fit to Clustered Data

## misspecification effects

- . Consider a logistic model fit to 2-level clustered data
  - . e.g., primary care clinics, patients within clinics
  - . exchangeable correlation
- . Assume the GEE or GLMM (not the survey sampling) modeling framework
- . With binary outcomes, there is more than one type of  $\rho_y$  estimate
  - . a phi-type estimate
  - . a tetrachoric-type estimate
  - . note that for linear models, there is no corresponding distinction
- . Which estimate of  $\rho_y$  should be used when estimating  $m_{eff}$ ?
  - . Answer: the phi-type coefficient, whether modeling via GEE or GLMM
  - . Investigate via Monte Carlo simulation.

# Simulated data: Mixed Logistic Model

- .  $m=100$  clusters, each with  $n=50$  units:  $N = m \times n = 5000$  per replicate sample
- . Generate binary  $y$  values with exchangeable correlation structure via a mixed logistic model with random intercepts,

$$y_{ij}^* = 0.5 + 0.1x1_{ij} + 0.5x2_j + u_j + e_{ij}; \quad \text{if } y^* > 0 \text{ then } y = 1, \text{ else } y = 0$$

where

- .  $u_j \sim N(0, \pi^2/3)$ ; the level-2 residuals; between-cluster variation
- .  $e_{ij} \sim LOGISTIC(0, \pi^2/3)$ ; the level-1 residuals; within-cluster variation
- .  $\rho_y = 0.5$  and  $\hat{r}_{tet.y} = 0.54$
- .  $x1_{ij} \sim N(0,1)$ ; a stochastic level-1  $x$  variable with  $\rho_x=0$ ;  $meff_{x1} \approx 1 - \rho_y$
- .  $x2_j \sim N(0,1)$ ; a stochastic level-2  $x$  variable:  $\rho_x=1$ ;  $meff_{x2} = 1 + (n-1)\rho_y$
- .  $r_{x1,x2} = 0$

- . 500 replicate samples

# Simulation: Logistic Regression Models Fit to Clustered Data

Fit two models to each replicate sample:

GEE logistic and mixed logistic with random intercepts (Laplace)

. Save parameter and standard error estimates,  $\hat{\rho}_y$ , simulated power

# Simulation: Logistic Regression Models Fit to Clustered Data

## Results: Intra-cluster correlation of outcome response

	intra-cluster correlation
$\rho_{y(\text{GEE})}$	0.348
phi <sup>†</sup>	0.365
$\rho_{y(\text{GLMM})}$	0.493
tetrachoric <sup>†</sup>	0.543

<sup>†</sup> estimated from first two units of each cluster

As you would expect, GEE working correlations are phi-like,  
whereas mixed logistic model intra-cluster correlations are tetrachoric-like

# Simulation: Logistic Regression Models Fit to Clustered Data

## Results: Parameter and Standard error estimates

	GEE	GLMM
Intercept		
parameter (std dev)	0.330 (.123)	0.509 (.189)
standard error	.124	.186
<i>x1</i>		
parameter (std dev)	0.064 (.024)	0.099 (.036)
standard error	.024	.036
<i>x2</i>		
parameter (std dev)	0.327 (.128)	0.501 (.190)
standard error	.126	.187

### Summary

- . GLMM parameter estimates are relatively unbiased (green highlight)
- . GEE and GLMM standard error estimates relatively unbiased (yellow highlight)

# Simulation: Logistic Regression Models Fit to Clustered Data

## Results: GEE Parameter Estimates Relatively Unbiased

	GEE	GLMM	ratio
Intercept			
parameter est.	0.330	0.509	.648
<i>x</i> 1			
parameter est.	0.064	0.099	.651
<i>x</i> 2			
parameter est.	0.327	0.501	.652

GEE parameter estimates are relatively unbiased

- .  $\rho_{y(\text{GEE})} = 0.348$
- . Scaling factor:  $1 - \rho_{y(\text{GEE})} = .652$  (equal to  $meff_{x1(\text{GEE})}$  in this example)
- .  $b_{\text{GEE}} \approx b_{\text{GLMM}} \times (1 - \rho_{y(\text{GEE})})$

The same scaling factor applies to standard error estimates

Neuhaus and Jewell (1990); Neuhaus, Kalbfleisch, and Hauck (1991); Neuhaus 1992 report #21, Eq. 14

# Using PASS to estimate power (compare to simulated power)

. For the GEE and GLMM results, calculate

$a. \Pr(y_{ij}=1 \mid x_1 = x_2 = 0)$ (intercept)
---

$b. \Pr(y_{ij}=1 \mid x_1 = 1)$
---------------------------------

$c. meff_{x_1} \approx 1 - \rho_y$ (because $\rho_x=0$ and $n$ is large)
--

$d. \Pr(y_{ij}=1 \mid x_2 = 1)$
---------------------------------

$e. meff_{x_2} = 1 + (n - 1)\rho_y$ (because $\rho_x=1$ )
---

. I estimated  $meff_{x_1}$  and  $meff_{x_2}$  using both  $\rho_{y(GEE)}$  and  $\rho_{y(GLMM)}$

. To solve for power for logistic regression, PASS requests

. specification of alpha: 0.05, two-tailed

. sample size: 5000  $\div$   $meff_{x_1}$  or 5000  $\div$   $meff_{x_2}$ , as appropriate

. baseline probability:  $a$

. alternative probability:  $b$  or  $d$ , as appropriate

. distribution of  $x$ : unit-standardized normal

PASS: estimate power for int.,  $x_1$ ,  $x_2$ , using both GEE- and GLMM-based  $meff$ s

# Simulation: Logistic Regression Models Fit to Clustered Data

## Results: Power

	GEE $\rho_{y(\text{GEE})} = 0.348$	GLMM $\rho_{y(\text{GLMM})} = 0.493$
<b>Intercept</b>		
power: simulated [PASS]	<b>.742 [.760]</b>	<b>.762 [.942]</b>
$m_{\text{eff}} = 1+(n-1)\rho_y$ ( $N_{\text{eff}}$ )	0.652 (277)	0.507 (199)
<b>x1</b>		
power: simulated [PASS]	<b>.788 [.787]</b>	<b>.778 [.997]</b>
$m_{\text{eff}} \approx 1-\rho_y$ ( $N_{\text{eff}}$ )	18.032 (7,664)	25.172 (9,868)
<b>x2</b>		
power: simulated [PASS]	<b>.726 [.734]</b>	<b>.756 [.942]</b>
$m_{\text{eff}} = 1+(n-1)\rho_y$ ( $N_{\text{eff}}$ )	0.652 (277)	0.507 (199)

- .  $m_{\text{eff}}$ -based estimates of  $N_{\text{eff}}$  in combination with PASS provided power estimates that were roughly equivalent to simulated values.
- . Clearly, when  $\rho_{y(\text{GLMM})}$  is used to estimate  $m_{\text{eff}}$ s, the result is not correct.

## Implications: Power for 2-level logistic models with exchangeable response correlation.

. If you have  $\hat{\rho}_{y(\text{GEE})}$  or  $\hat{\phi}$  as an estimate of intra-cluster correlation of binary response, then you can estimate power via *meffs* and standard software (PASS)

. When using *meff*-derived  $N_{\text{eff}}$  to help estimate power for logistic models, the regression parameters input into (or estimated by) the standard power analysis software will represent population average parameter estimates, i.e., the type of parameter estimates produced by GEE logistic regression

After completing a *meff*-driven power analysis, you can approximate the minimum detectable unit-specific parameter estimates from their population average counterparts using the scaling factor described by John Neuhaus

## Implications: Power for 2-level logistic models with exchangeable response correlation.

. If you only have  $\hat{\rho}_{y(\text{GLMM})}$  or  $\hat{r}_{tet}$  as an intra-cluster correlation estimate of binary response, then you should not use them to estimate power via *meffs*

Instead...

(i) estimate power by simulation using a GLMM data-generating model

When using a GLMM data-generating model, you subsequently can estimate power via GLMM or GEE logistic regression

It is your call, because given exchangeable response correlation GEE and GLMM models provide equivalent power

or

(ii) use the GLMM-generated data to estimate  $\hat{\rho}_{y(\text{GEE})}$  by simulation and then proceed with *meff*-based methods

# Limitations

Very limited simulation

- . 'large' number of clusters and 'large' clusters considered
  - . *meff*-based approximations may not work as well with smaller  $m$  or  $n$
- . simple two-level model
- . balanced cluster size
- . limited values of  $\rho_y$  and  $\rho_x$  considered.
- . limited replicate samples

*When in doubt, estimate power by simulation*

**Thank you**