# Simplified power analyses

# for clustered sampling designs

# with compound symmetric covariance structure of *x* & y:

# A survey of sample size ratios (SSR)

UCSF CAPS Methods Core

May 11$^{th}$, 2021

Steve Gregorich

# Overview

Sample size ratios (SSR) provide convenient short-cuts
  for sample size calculations

## Assumptions

. 2- and 3-level clustered sampling designs

. Limited coverage of 3-level designs in this talk

. Compound symmetric correlation structure of both $x$ and $y$

## Regression modeling contexts

. GLMM

. GEE

. Survey Sampling (SS)

# Sample Size Ratios (SSR): Introduction

AKA design effect, misspecification effect, variance inflation/deflation factor.
I chose the SSR label because it is broadly applicable

Assume a simple random sample (SRS) of size $N$ drawn from a population with population mean $\mu$ and variance $\sigma^2$

We choose the usual estimator $\hat{\mu}$ of the sample mean of $x$

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} x_i$$

The variance of the estimator $\hat{\mu}$ is

$$\sigma_{\hat{\mu}}^2 = \sigma^2 / N$$

The *precision* of the estimator is the inverse of the above quantity, i.e., $1/\sigma_{\hat{\mu}}^2 = N/\sigma^2$, i.e., larger $N$ obtains higher precision.

# Sample Size Ratios (SSR): Introduction

Say we have an alternative estimator $\hat{\mu}_a$ with variance equal to

$$\sigma^2_{\hat{\mu}_a} = \sigma^2/N_a, \quad \text{and rearranging}$$

$$N_a = \sigma^2/\sigma^2_{\hat{\mu}_a},$$

i.e., $N_a$ equals population variance $(\sigma^2) \times$ estimator precision $(1/\sigma^2_{\hat{\mu}_a})$

Similarly, for estimator $\hat{\mu}$

$$N = \sigma^2/\sigma^2_{\hat{\mu}}$$

SSR represents relative (effective) sample size and relative precision, i.e.,

$$\text{SSR} = \frac{N}{N_a} = \frac{\frac{\sigma^2}{\sigma^2_{\hat{\mu}}}}{\frac{\sigma^2}{\sigma^2_{\hat{\mu}_a}}} = \frac{\sigma^2_{\hat{\mu}_a}}{\sigma^2_{\hat{\mu}}}$$

# Sample Size Ratios (SSR): Introduction

$$\mathrm{SSR} = \frac{N}{N_a} = \frac{\sigma^2_{\hat{\mu}_a}}{\sigma^2_{\hat{\mu}}}$$

Assume $N$=1000, estimator $\hat{\mu}_a$ has $\sigma^2_{\hat{\mu}_a}$=2, and estimator $\hat{\mu}$ has $\sigma^2_{\hat{\mu}}$=1
   . SSR, as defined above, equals 2÷1=2
   . I.e., $\hat{\mu}_a$ has larger variance and lower precision than $\hat{\mu}$

Knowing $N$ and SSR, we can calculate the effective sample size, $N$eff,
      for an application of $\hat{\mu}_a$

For $N$=1000 and SSR=2, when applying $\hat{\mu}_a$, $N$eff=$N$÷SSR=  500.
   . I.e., when applying $\hat{\mu}_a$ with $N$=1000, the $N$eff=500
   . Or, wrt precision, $\hat{\mu}_a$ with $N$=1000 is equivalent to $\hat{\mu}$  with $N$=500

At the same time, SSR=2 indicates the expectations that
   . the variance of $\hat{\mu}_a$ (i.e., $\sigma^2_{\hat{\mu}_a}$) will equal 2× the variance of $\hat{\mu}$ (i.e., $\sigma^2_{\hat{\mu}}$)
   . the std err of $\hat{\mu}_a$ (i.e., $\sigma_{\hat{\mu}_a}$) will equal $\sqrt{2}$ times the std err of $\hat{\mu}$ (i.e., $\sigma_{\hat{\mu}}$)

# Ex #1a: Planning for a Cluster-Randomized Trial (CRT)

Context
. Clustered sampling: :Level1 participants nested w/in Level2 clusters
. Level2 clusters are randomized w/ 1:1 allocation to experimental groups
. $N=1000$: $n_2=100$ clusters, each of size $n_1=10$
. $y \sim N(0,1)$, $x \sim B(0.50)$, where $x$ is the experimental group indicator
. Linear regression model
. Intra-cluster correlation (ICC) of $y$ ($\rho_y$) equals 0.05
. 80% power with two-tailed $\alpha = .05$

Goal
. Solve for minimum detectable effect size, $b_x$

In this context, the familiar Design Effect (Deff) is a useful SSR.

. $\mathrm{SSR} = \mathrm{Deff} = 1 + r\rho_y$, where $r = n_1 - 1$,

Deff was described by Kish

Kish (1965). *Survey Sampling.* New York: John Wiley & Sons, Inc

# Ex #1a: Planning for a Cluster-Randomized Trial (CRT)

Application of SSR (Deff) to solve for $b_x$

Step 1. Calculate SSR & $N_{\text{eff}}$ given $r=9$, $\rho_y=.05$, and $N=1000$

$\quad$ . $\text{SSR} = \text{Deff} = 1 + r\rho_y = 1 + (10 - 1) \times .05 = 1.45$

$\quad$ . $N_{\text{eff}} = N/\text{SSR} = 1000/1.45 = 689.7$

$\quad$ . Note. $N_{\text{eff}} < N$

Step 2. Calculate minimum detectable effect specifying $N=689.7$

$\quad$ Result: $b_x = .21244$ (from PASS Linear Regression routine, $\sigma_y^2=1$)

In this case, a GLMM/GEE/SS model fit to $N=1000$ clustered observations obtains the same power as plain linear regression model fit to $N \cong 690$ independent observations (SRS)

# Ex #1b: Planning for a Cluster-Randomized Trial (CRT)

If we instead began w/ values of $b_x$ and $n_1$, we could solve for $n_2$ and $N$
. $b_x$ = 0.21244
. $n_1$ = 10

---

Step 1. Calculate $N$ assuming $b_x$=.21244 & independent obs. ($\rho_y$=0)

    . from PASS Linear Regression routine, $N = 690$, if $\rho_y$=0
      (PASS only returns integer $N$ values)

    . In this case, $N$ from PASS is our target effective sample size, $N_{eff}$

---

Step 2. Calculate $N$ assuming $\rho_y$=.05 and $n_{L1}$=10

    . $\text{SSR} = 1 + r\rho_y = 1 + 9 \times .05 = 1.45$

    . $N = N_{\text{eff}} \times \text{SSR} = 690 \times 1.45 = 1000.5$

    . $n_2 = N/n_1 = 1000.5/10 \cong 100$

---

# Ex #2: Planning for a RCT Randomizing Level1 Units

## Context
. Clustered sampling: Level1 participants nested w/in Level2 clusters

. Level1 units are randomized with 1:1 allocation to experimental groups

. $N=1000$: $n_2=100$ clusters, each of size $n_1=10$

. $y \sim N(0,1),\ x \sim B(0.50)$, where $x$ is the experimental group indicator

. Linear regression model

. Intra-cluster correlation (ICC) of $y$ ($\rho_y$) equals 0.05

. 80% power with two-tailed $\alpha = .05$

## Goal
. Solve for minimum detectable effect size, $b_x$

## In this context, a familiar sample size ratio is

. $\mathrm{SSR} = 1 - \rho_y$

This is the same SSR that applies to a paired $t$-test

# Ex #2: Planning for a RCT Randomizing Level1 Units

Application of the SSR: solve for $b_x$

---

Step 1. Calculate SSR and Effective Sample Size ($N_{eff}$)

. $\text{SSR} = 1 - \rho_y = 1 - .05 = 0.95$

. $N_{eff} = N/\text{SSR} = 1000/0.95 = 1052.6$

. Note. $N_{eff} > N$

---

Step 2. Calculate minimum detectable effect specifying $N$=1052.6

Result: $b_x$ = .17222 (from PASS)

---

Here, a GLMM/GEE/SS model fit to $N$=1000 clustered observations obtains the same precision as a plain model fit to $N$=1053 independent observations

# Sample Size Ratios (SSR): 2-Level Sampling Designs

So far, we've discussed two sample size ratios

$$\text{SSR} = \text{Deff} = 1 + r\rho_y$$

and

$$\text{SSR} = 1 - \rho_y$$

When does each apply?

The choice depends on whether the $x$ has
. a between-cluster or
. a within-cluster effect

The intra-cluster correlation of $x$ ($\rho_x$) is important

# Sample Size Ratios (SSR): 2-Level Sampling Designs

We often think of ICC in terms of a variance component decomposition.

$$\rho = \frac{\sigma^2_{\text{between\_cluster}}}{\sigma^2_{\text{between\_cluster}} + \sigma^2_{\text{w/in\_cluster}}}$$

However, that formulation has positive bias.

When thinking about $\rho_x$, it is helpful to consider the unbiased formula (Harris 1913; Kish 1965; Wikipedia ICC page)

$$\rho_x = \frac{\sigma^2_{x.\text{between\_cluster}} - \sigma^2_{x.\text{w/in\_cluster}}/r}{\sigma^2_{x.\text{between\_cluster}} + \sigma^2_{x.\text{w/in\_cluster}}}$$

Harris JA (1913). On the calculation of intra-class and inter-class coefficients of correlation from class moments when the number of possible combinations is large. *Biometrika*, 9, 446–472.

Kish, Leslie (1965). *Survey Sampling.* New York: John Wiley & Sons, Inc (p. 170)

# Sample Size Ratios (SSR): 2-Level Sampling Designs

$$\text{SSR} = \text{Deff} = 1 + r\rho_y$$

In a regression context,
    this SSR applies when $x$ has a fully <u>between-cluster</u> effect on $y$

---

When will $x$ have a fully <u>between-cluster</u> effect on $y$?
    . When $x$ is a Level2 variable; it will have
        positive between-cluster & <u>zero</u> (**0**) within-cluster variance

    . In this case, the intra-cluster correlation of $x$ ($\rho_x$) equals 1.0.

Consider the unbiased formula for intra-cluster correlation as applied to $x$

$$\rho_x = \frac{\sigma^2_{x.\text{between\_cluster}} - \mathbf{0}/r}{\sigma^2_{x.\text{between\_cluster}} + \mathbf{0}} = 1$$

---

Use of this SSR (i.e.,. $\text{SSR} = \text{Deff} = 1 + r\rho_y$) assumes
    a fully between-cluster $x$ effect, $\rho_x = 1$

Therefore, I label this SSR as SSR$_b$  (i.e., sub-'b' for 'between')

# Summary Implications: SSR$_b$

. SSR$_b$ for a fully between-cluster effect, i.e., when $\rho_x = 1$

$$\text{SSR}_b = 1 + r\rho_y,$$

Basic results for SSR$_b$

| $\rho_y$ | SSR$_b$ result | $N_{\text{eff}}$ vs $N$ |
|---|---|---|
| $\rho_y = 0$ | $\text{SSR}_b = 1$ | $N_{\text{eff}} = N/1 \quad = N$ |
| $\rho_y > 0$ | $\text{SSR}_b > 1$ | $N_{\text{eff}} = N/\text{SSR}_b < N$ |

. I.e., when $\rho_x = 1$ and $\rho_y > 0$,

  $x$ has a between-cluster effect that will have <u>lower</u> precision
   versus within an alternative SRS design, all else being equal

Note. $\rho_y < 0$ is rare and not considered in this talk

# Sample Size Ratios (SSR): 2-Level Sampling Designs

$$\text{SSR} = 1 - \rho_y$$

In a regression context…
   . This SSR applies when $x$ has a fully <u>within-cluster</u> effect on $y$.
   . I.e., when $x$ has exactly <u>zero</u> (**0**) between-cluster variation.
   . In this circumstance, $\rho_x$ takes its minimum value.

$$\rho_x = \frac{\mathbf{0} - \sigma^2_{\text{w/in\_cluster}}/r}{\mathbf{0} + \sigma^2_{\text{w/in\_cluster}}} = \frac{-1}{r}$$

When will $x$ have a fully <u>within-cluster</u> effect?
   . Often, a Level1 <u>design variable</u> w/ zero between-cluster variation, e.g.,
   . RCT randomizing Level1 units w/ identical proportionate allocation
      across clusters
   . $x$ indicates scheduled assessment times (base, 6m, 12m)
   . $x$ indicates intra-cluster role/position, e.g., doctor versus patient
   . Not always design vars. e.g., $x$ holds deviations from cluster means

I label this SSR (i.e., $\text{SSR} = 1 - \rho_y$) as $\text{SSR}_\text{w}$     (i.e., sub-'w' for 'within')

# Summary Implications: SSR$_w$

. SSR$_w$ for a fully within-cluster effect, i.e., when $\rho_x = -1/r$

$$\text{SSR}_w = 1 - \rho_y$$

Basic results for SSR$_w$

| $\rho_y$ | SSR$_w$ result | $N_{eff}$ vs $N$ |
|----------|----------------|------------------|
| $\rho_y = 0$ | $\text{SSR}_w = 1$ | $N_{eff} = n_{L2} \times n_{L1} = N$ |
| $\rho_y > 0$ | $\text{SSR}_w < 1$ | $N_{eff} = N/\text{SSR}_w \; > N$ |

. I.e., when $\rho_x = -1/r$ and $\rho_y > 0$,
    x has a within-cluster effect that will have <u>higher</u> precision vs SRS

SSR$_w$ applies whenever $\rho_x = -1/r$

$\rho_x = -1/r$ when between-cluster x variation equals zero (exactly).
    . Often, but not always, by design (e.g., assessment times) or
    analysis (e.g., deviation scores)

# Sample Size Ratios (SSR): 2-Level Sampling Designs

So far, we've covered SSRs for
. fully between-cluster effects, i.e., when $\sigma^2_{x.\text{within}} = 0$, $\rho_x = 1$ and
. fully within-cluster effects, i.e., when $\sigma^2_{x.\text{between}} = 0$, $\rho_x = -1/r$

However, $\rho_x$ values are not limited to $-1/r$ and 1

When $-1/r < \rho_x < 1$,
    $x$ can have both between- and within-cluster effects

You might expect $-1/r < \rho_x < 1$ when $x$ is a(n)...
    . design var. w/ some between-cluster variance (often $-1/r < \rho_x < 0$)
    . observed random, eg, participant-reported, variable (often $0 \leq \rho_x < 1$)

Which SSR should be used when $-1/r < \rho_x < 1$?

The answer depends upon the regression modeling framework, i.e.,
    . Survey Sampling (SS) versus
    . Generalized Linear Mixed Models (GLMM) or GEE

Why is this the case?

# Sample Size Ratios (SSR): 2-Level Sampling Designs

First, a side note. The formulation…

$$\rho_x = \frac{\sigma^2_{x.\text{between\_cluster}} - \sigma^2_{x.\text{w/in\_cluster}}/r}{\sigma^2_{x.\text{between\_cluster}} + \sigma^2_{x.\text{w/in\_cluster}}}$$

makes it clear that $\rho_x = 0$ when $\sigma^2_{x.\text{between\_cluster}} = \sigma^2_{x.\text{w/in\_cluster}}/r$

Thus, when $\rho_x = 0$, positive between-cluster variation is expected

---

Reminder…

In this talk, when $-1/r < \rho_x < 1$,
   I assume equivalent between- and within-cluster effects of $x$

For reference, John Neuhaus describes modeling options
   that decompose between- and within-cluster effects.

Neuhaus, JM and Kalbfleisch, JD (1988).  Between- and within-cluster covariate effects
   in the analysis of clustered data.  *Biometrics*, 54, 638-645.

Neuhaus, JM (2001). Assessing change with longitudinal and clustered binary data.
   *Annual Review of Public Health*, 22, 115-118.

# Sample Size Ratios (SSR): 2-Level Sampling Designs

SSR$_{SS}$: SSR for the Survey Sampling regression modeling framework

$$\text{SSR}_{SS} = 1 + r\rho_x\rho_y$$

Essentially, SSR$_{SS}$ adds $\rho_x$ to the product term of SSR$_b$

SSR$_{SS}$ versus SSR$_b$ and SSR$_w$

| $\rho_x$ | Result |
|---|---|
| 1.0 | $\text{SSR}_{SS} = \text{SSR}_b = 1 + r\rho_y$ |
| $-1/r$ | $1 - \rho_y = \text{SSR}_w = \text{SSR}_{SS}$ |
| $-1/r < \rho_x < 1$ | $1 - \rho_y < \quad \text{SSR}_{SS} \quad < 1 + r\rho_y$ |

SSR$_{SS}$ was described by

Scott, AJ and Holt, D (1982). The Effect of Two-Stage Sampling on Ordinary Least Squares Methods. *Journal of the American Statistical Association*, 77(380), 848- 854.

(the SSR$_{SS}$ label is mine)

# Sample Size Ratios (SSR): 2-Level Sampling Designs

$$\text{SSR}_\text{SS} = 1 + r\rho_x\rho_y \qquad \text{vs} \qquad \text{SSR}_\text{b} = 1 + r\rho_y \qquad \text{vs} \qquad \text{SSR}_\text{w} = 1 - \rho_y$$

Example SSR$_\text{SS}$ results assuming SS model, $\rho_y = 0.10$, $N$=1000, $r$ =10

| $\rho_x$ | $\rho_y$ | $x$ effect type | Neff$_\text{SS}$= 1000/SSR$_\text{SS}$ | Neff$_\text{b}$= 1000/SSR$_\text{b}$ | Neff$_\text{w}$= 1000/SSR$_\text{w}$ |
|---|---|---|---|---|---|
| 1.00 | 0.10 | btw-cluster | 500 | 500 | -- |
| 0.10 | 0.10 | btw- & w/in- | 909 | -- | -- |
| 0 | 0.10 | btw- & w/in- | 1000 | -- | -- |
| -0.05 | 0.10 | btw- & w/in- | 1053 | -- | -- |
| -0.10† | 0.10 | w/in-cluster | 1111 | -- | 1111 |

† $\rho_x = -1/r = -0.10$. '--': Inappropriate applications of SSR$_\text{b}$ & SSR$_\text{w}$

When applying the SS modeling framework

    . Do not use SSR$_\text{b}$ unless $\rho_x = 1$ (or, trivially, $\rho_y = 0$)
      If $\rho_x < 1$, then use of SSR$_\text{b}$ can <u>underestimate</u> $N_\text{eff}$, power

    . Do not use SSR$_\text{w}$ unless $\rho_x = -1/r$
      if $-1/r < \rho_x$, then use of SSR$_\text{w}$ can <u>overestimate</u> $N_\text{eff}$, power

# Sample Size Ratios (SSR): 2-Level Sampling Designs

SSR$_{GE}$ for the GEE and GLMM modeling frameworks

$$\text{SSR}_{\text{GE}} = \frac{(1+r\rho_y)(1-\rho_y)}{1-\rho_y+r\rho_y(1-\rho_x)} = \frac{1-\rho_y+r\rho_y(1-\rho_y)}{1-\rho_y+r\rho_y(1-\rho_x)}$$

SSR$_{GE}$ is based upon the seminal (but underutilized) work of
  Basagaña, Liao, and Spiegelman (2011; BLS).

BLS were focused on power of longitudinal studies estimating the effects
  of a time-varying binary *x* variable.

BLS reported a SSR assuming compound symmetric (CS) $\rho_x$ and $\rho_y$
  versus assuming $\rho_x = 1$ and CS $\rho_y$.  See their Eq. 3.5.

SSR$_{GE}$ manipulates BLS Eq. 3.5 to reflect comparison of
  . a clustered sampling design with CS $\rho_x$ and $\rho_y$ versus
    *N* independently sampled units (SRS).

  . See Appendix A

Basagaña, X., Liao, X., and Spiegelman. D. (2011). Power and sample size calculations for longitudinal studies estimating a main effect of a time-varying exposure.  *Statistical Methods in Medical Research*, 29, 181-192.

# Sample Size Ratios (SSR): 2-Level Sampling Designs

SSR$_{GE}$ for the GEE and GLMM modeling frameworks

$$\text{SSR}_{GE} = \frac{1 - \rho_y + r\rho_y(1 - \rho_y)}{1 - \rho_y + r\rho_y(1 - \rho_x)}$$

SSR$_{GE}$ versus SSR$_b$, SSR$_w$, and SSR$_{SS}$

| $\rho_x$ | Result | | |
|---|---|---|---|
| 1.0 | $\text{SSR}_{GE} = \text{SSR}_{SS} = \text{SSR}_b = 1 + r\rho_y$ | | |
| $-1/r$ | $1 - \rho_y = \text{SSR}_w = \text{SSR}_{GE} = \text{SSR}_{SS}$ | | |
| $-1/r < \rho_x < 1$ | $1 - \rho_y <$ | $\text{SSR}_{GE} < \text{SSR}_{SS}$ | $< 1 + r\rho_y$ |

Notes.

. SSR$_{GE}$ ≤ SSR$_{SS}$

. When $-1/r < \rho_x < 1$, $N$eff$_{GE}$ > $N$eff$_{SS}$, i.e.,

   superior power via the GEE/GLMM vs SS modeling framework

# Sample Size Ratios (SSR): 2-Level Sampling Designs

$$\text{SSR}_{\text{GE}} = \frac{1-\rho_y+r\rho_y(1-\rho_y)}{1-\rho_y+r\rho_y(1-\rho_x)} \qquad \text{vs} \qquad \text{SSR}_{\text{SS}} = 1 + r\rho_x\rho_y$$

Results: $Neff_{\text{GE}}$ versus $Neff_{\text{SS}}$: $N$=1000, $\rho_y = 0.10$, $r$=10

| $\rho_x$ | $\rho_y$ | $x$ effect type | $Neff_{\text{GE}}=$ 1000/SSR$_{\text{GE}}$ | $Neff_{\text{SS}}=$ 1000/SSR$_{\text{SS}}$ | $\dfrac{Neff_{\text{GE}}}{Neff_{\text{SS}}}$ |
|---|---|---|---|---|---|
| 1.00 | 0.10 | btw-cluster | 500 | 500 | = |
| 0.10 | 0.10 | btw- & w/in- | 1000 | 909 | +10% ‡ |
| 0 | 0.10 | btw- & w/in- | 1056 | 1000 | +6% ‡ |
| -0.05 | 0.10 | btw- & w/in- | 1083 | 1053 | +3% ‡ |
| -0.10† | 0.10 | w/in-cluster | 1111 | 1111 | = |

note. † $\rho_x = -1/r = -0.10$.  ‡ GEE/GLMM has a power advantage over SS

The tabled results are not dramatic, but the GEE/GLMM advantage can be stark for some combinations of $\rho_x$, $\rho_y$, and $r$

When you have a choice, GEE/GLMM can be more efficient than SS

# Sample Size Ratios (SSR): 2-Level Sampling Designs

$$\text{SSR}_{\text{GE}} = \frac{1-\rho_y+r\rho_y(1-\rho_y)}{1-\rho_y+r\rho_y(1-\rho_x)} \qquad \text{vs} \qquad \text{SSR}_{\text{SS}} = 1 + r\rho_x\rho_y$$

Results: $N\text{eff}_{\text{GE}}$ vs $N\text{eff}_{\text{SS}}$: $N=1000$, $\rho_y=.50$, $r=1$ (e.g., repeated measures)

| $\rho_x$ | $\rho_y$ | $x$ effect type | $N\text{eff}_{\text{GE}}=$ <br> $1000/\text{SSR}_{\text{GE}}$ | $N\text{eff}_{\text{SS}}=$ <br> $1000/\text{SSR}_{\text{SS}}$ | $\dfrac{N\text{eff}_{\text{GE}}}{N\text{eff}_{\text{SS}}}$ |
|---|---|---|---|---|---|
| 1.00 | 0.50 | btw-cluster | 667 | 667 | = |
| 0.10 | 0.50 | btw- & w/in- | 1267 | 952 | +33% ‡ |
| 0 | 0.50 | btw- & w/in- | 1333 | 1000 | +33% ‡ |
| -0.50 | 0.50 | btw- & w/in- | 1667 | 1333 | +25% ‡ |
| -1.00[†] | 0.50 | w/in-cluster | 2000 | 2000 | = |

note. [†] $\rho_x = -1/r = -1.0$.   ‡ GEE/GLMM has a power advantage over SS

The following slide compares $\text{SSR}_{\text{GE}}$ and $\text{SSR}_{\text{SS}}$ values
    for a range of $\rho_x$ and $\rho_y$ values and $r=1$

# Expected $\text{Neff}_{GE}$ (**bold**) & $\text{Neff}_{SS}$ (regular): $r=1$ & $N=1000$

| | | $\rho_y$ | | | |
|---|---|---|---|---|---|
| | | 0 | 0.3 | 0.6 | 0.9 |
| $\rho_x$ | $-1/r$ | **1000** / 1000 | **1429** / 1429 | **2500** / 2500 | **10,000** / 10,000 |
| | 0 | **1000** / 1000 | **1099** / 1000 | **1563** / 1000 | **5263** / 1000 |
| | 0.3 | **1000** / 1000 | **1000** / 917 | **1281** / 847 | **3842** / 787 |
| | 0.6 | **1000** / 1000 | **901** / 847 | **1000** / 735 | **2421** / 649 |
| | 0.9 | **1000** / 1000 | **802** / 787 | **718** / 649 | **1000** / 552 |
| | 1.0 | **1000** / 1000 | **769** / 769 | **625** / 625 | **526** / 526 |

. If $\rho_y = 0$ (green), $\rho_x = -1/r$ (purple), or $\rho_x = 1$ (pink), then $\text{Neff}_{GE} = \text{Neff}_{SS}$. Otherwise, $\text{Neff}_{GE} > \text{Neff}_{SS}$

. If $\rho_x = \rho_y$, then $\text{Neff}_{GE} = N$ (yellow)

. If $\rho_x < \rho_y$, then $\text{Neff}_{GE} > N$ (purple & blue)

. If $\rho_x > \rho_y$, then $\text{Neff}_{GE} < N$ (orange & pink)

# SSR inputs: linear versus logistic models

When planning for a <u>logistic</u> regression analysis there are some wrinkles

Population average (GEE, SS) versus unit-specific (GLMM)…
  . estimates of $x$ effects ($b_x$) as well as
      model-predicted vs observed outcome probabilities, and
  . estimates of $\rho_x$ and $\rho_y$

Power analyses described in this talk require population average inputs

If inputting an effect estimate into power analysis for logistic regression,
    choose a population average estimate
    (e.g., based upon observed means or a GEE, ALR, or SS analysis)

Obtain $\rho_x$ & $\rho_y$ estimates from a GEE logistic model,
    not a mixed logistic model

    . GEE ICCs reflect intra-cluster homogeneity of observed values

    . In contrast, mixed logistic model ICCs reflect underlying latent values

# Sample Size Ratios (SSR): 2-Level Sampling Designs

Simulation comparing calculated versus simulated SSRs

Simulate 2-level data for $c$ = 1 to 63 combinations of $\rho_x$ and $\rho_y$ values

. $\rho_x$ ranging from $-1/r$ to 1.0
. $\rho_y$ ranging from 0 to 0.9

Generate $i$ = 1 to 10K replicate samples from each of 63 combinations

Fit regression model to each replicate sample assuming independent obs.
Save standard error estimates for fixed effect of $x$, $\hat{\sigma}_{\text{ind.}ci}$

Fit regression model to each replicate sample assuming clustered obs.
Save standard error estimates for fixed effect of $x$, $\hat{\sigma}_{\text{clus.}ci}$

A simulated SSR for combination $c$ averages
$(\hat{\sigma}_{\text{clus.}ci}/\hat{\sigma}_{\text{ind.}ci})^2$ values across $i$=1 to 10K replicate samples

Compare simulated SSR to SSR$_{\text{GE}}$ or SSR$_{\text{SS}}$, as appropriate

Simulated (bold) and Expected (plain) Values of $SSR_{GE}$ assuming small clusters ($r$=1) and 63 Combinations of $\rho_x$ and $\rho_y$: GEE Linear Regression Modeling Framework with exchangeable working correlation structure.

| $\rho_x$ \ $\rho_y$ | 0 | .05 | .10 | .25 | .50 | .75 | .90 |
|---|---|---|---|---|---|---|---|
| **-1.0** | **1.000** | **0.950** | **0.900** | **0.749** | **0.499** | **0.249** | **0.099** |
| $(-1/r)$ | 1.000 | 0.950 | 0.900 | 0.750 | 0.500 | 0.250 | 0.100 |
| **0** | **0.998** | **0.996** | **0.988** | **0.936** | **0.749** | **0.437** | **0.188** |
| | 1.000 | 0.998 | 0.990 | 0.938 | 0.750 | 0.438 | 0.190 |
| **.05** | **0.998** | **0.998** | **0.993** | **0.948** | **0.769** | **0.454** | **0.196** |
| | 1.000 | 1.000 | 0.995 | 0.949 | 0.769 | 0.455 | 0.199 |
| **.10** | **0.998** | **1.000** | **0.998** | **0.960** | **0.789** | **0.471** | **0.206** |
| | 1.000 | 1.003 | 1.000 | 0.962 | 0.789 | 0.473 | 0.209 |
| **.25** | **0.998** | **1.008** | **1.013** | **0.998** | **0.855** | **0.537** | **0.243** |
| | 1.000 | 1.010 | 1.015 | 1.000 | 0.857 | 0.538 | 0.245 |
| **.50** | **0.998** | **1.021** | **1.040** | **1.069** | **0.998** | **0.699** | **0.341** |
| | 1.000 | 1.023 | 1.042 | 1.071 | 1.000 | 0.700 | 0.345 |
| **.75** | **0.997** | **1.034** | **1.068** | **1.152** | **1.198** | **0.996** | **0.578** |
| | 1.000 | 1.036 | 1.070 | 1.154 | 1.200 | 1.000 | 0.585 |
| **.90** | **0.997** | **1.043** | **1.086** | **1.207** | **1.362** | **1.343** | **0.991** |
| | 1.000 | 1.045 | 1.088 | 1.210 | 1.364 | 1.346 | 1.000 |
| **1.00** | **0.997** | **1.048** | **1.098** | **1.249** | **1.499** | **1.750** | **1.901** |
| | 1.000 | 1.050 | 1.110 | 1.250 | 1.500 | 1.750 | 1.900 |

Note. Simulated data for each combination of $\rho_x$ and $\rho_y$ included $m$=500 clusters, $n$=2 units per cluster, $N$=1000, and 10K replicate samples. Simulated $SSR_{GE}$ estimated from comparison of std errs estimated from GEE model versus cluster-naïve model. Cell shading codes for combinations of $\rho_x$ and $\rho_y$ values: Grey: E[$SSR_{GE}$]=1.0. Green: E[$SSR_{GE}$]<1.0; Orange: E[$SSR_{GE}$]>1.0

Simulated (bold) and Expected (plain) Values of SSR$_{GE}$ assuming small clusters ($r$=1) and 63 Combinations of $\rho_x$ and $\rho_y$: GLMM Linear Regression Modeling Framework with exchangeable working correlation structure.

| $\rho_x$ | $\rho_y$ | | | | | | |
|---|---|---|---|---|---|---|---|
| | 0 | .05 | .10 | .25 | .50 | .75 | .90 |
| -1.0 | **1.000** | **0.947** | **0.900** | **0.750** | **0.500** | **0.251** | **0.100** |
| ($-1/r$) | 1.000 | 0.950 | 0.900 | 0.750 | 0.500 | 0.250 | 0.100 |
| | **0.998** | **0.996** | **0.988** | **0.936** | **0.750** | **0.439** | **0.191** |
| 0 | 1.000 | 0.998 | 0.990 | 0.938 | 0.750 | 0.438 | 0.190 |
| | **0.998** | **0.998** | **0.993** | **0.948** | **0.769** | **0.456** | **0.200** |
| .05 | 1.000 | 1.000 | 0.995 | 0.949 | 0.769 | 0.455 | 0.199 |
| | **0.998** | **1.001** | **0.998** | **0.960** | **0.790** | **0.473** | **0.210** |
| .10 | 1.000 | 1.003 | 1.000 | 0.962 | 0.789 | 0.473 | 0.209 |
| | **0.998** | **1.009** | **1.013** | **0.998** | **0.856** | **0.539** | **0.247** |
| .25 | 1.000 | 1.010 | 1.015 | 1.000 | 0.857 | 0.538 | 0.245 |
| | **0.999** | **1.023** | **1.040** | **1.069** | **0.999** | **0.701** | **0.346** |
| .50 | 1.000 | 1.023 | 1.042 | 1.071 | 1.000 | 0.700 | 0.345 |
| | **0.999** | **1.037** | **1.069** | **1.153** | **1.199** | **0.998** | **0.586** |
| .75 | 1.000 | 1.036 | 1.070 | 1.154 | 1.200 | 1.000 | 0.585 |
| | **0.999** | **1.047** | **1.088** | **1.209** | **1.363** | **1.345** | **1.000** |
| .90 | 1.000 | 1.045 | 1.088 | 1.210 | 1.364 | 1.346 | 1.000 |
| | **0.999** | **1.053** | **1.100** | **1.251** | **1.501** | **1.752** | **1.903** |
| 1.00 | 1.000 | 1.050 | 1.100 | 1.250 | 1.500 | 1.750 | 1.900 |

Note. Simulated data for each combination of $\rho_x$ and $\rho_y$ included $m$=500 clusters, $n$=2 units per cluster, $N$=1000, and 10K replicate samples. Simulated SSR$_{GE}$ estimated from comparison of std errs estimated from GLMM linear model versus cluster-naïve model. Cell shading codes for combinations of $\rho_x$ and $\rho_y$ values: Grey: E[SSR$_{GE}$]=1.0. Green: E[SSR$_{GE}$]<1.0; Orange: E[SSR$_{GE}$]>1.0

Simulated (bold) and Expected (plain) Values of SSR$_{SS}$ assuming small clusters ($r$=1) and 63 Combinations of $\rho_x$ and $\rho_y$: Survey Sampling Linear Regression Modeling Framework.

| $\rho_x$ | $\rho_y$ 0 | .05 | .10 | .25 | .50 | .75 | .90 |
|---|---|---|---|---|---|---|---|
| -1.0 | **1.001** | **0.951** | **0.901** | **0.750** | **0.501** | **0.251** | **0.100** |
| $(-1/r)$ | 1.000 | 0.950 | 0.900 | 0.750 | 0.500 | 0.250 | 0.100 |
| | **1.000** | **0.999** | **1.001** | **0.999** | **0.999** | **1.001** | **1.000** |
| 0 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 |
| | **0.999** | **1.003** | **1.004** | **1.012** | **1.025** | **1.037** | **1.043** |
| .05 | 1.000 | 1.003 | 1.005 | 1.013 | 1.025 | 1.038 | 1.045 |
| | **1.001** | **1.004** | **1.010** | **1.025** | **1.050** | **1.073** | **1.088** |
| .10 | 1.000 | 1.005 | 1.010 | 1.025 | 1.050 | 1.075 | 1.090 |
| | **0.999** | **1.012** | **1.024** | **1.062** | **1.123** | **1.185** | **1.223** |
| .25 | 1.000 | 1.013 | 1.025 | 1.063 | 1.125 | 1.188 | 1.225 |
| | **0.999** | **1.023** | **1.049** | **1.123** | **1.248** | **1.372** | **1.447** |
| .50 | 1.000 | 1.025 | 1.050 | 1.125 | 1.250 | 1.375 | 1.450 |
| | **0.997** | **1.035** | **1.073** | **1.186** | **1.372** | **1.560** | **1.673** |
| .75 | 1.000 | 1.038 | 1.075 | 1.188 | 1.375 | 1.563 | 1.675 |
| | **0.996** | **1.043** | **1.088** | **1.222** | **1.447** | **1.672** | **1.809** |
| .90 | 1.000 | 1.045 | 1.090 | 1.225 | 1.450 | 1.675 | 1.810 |
| | **0.997** | **1.047** | **1.096** | **1.246** | **1.496** | **1.748** | **1.900** |
| 1.00 | 1.000 | 1.050 | 1.100 | 1.250 | 1.500 | 1.750 | 1.900 |

Note. Simulated data for each combination of $\rho_x$ and $\rho_y$ included $m$=500 clusters, $n$=2 units per cluster, $N$=1000, and 10K replicate samples. Simulated SSR$_{SS}$ estimated from comparison of std errs estimated from SS linear model versus cluster-naïve model. Cell shading codes for combinations of $\rho_x$ and $\rho_y$ values: Grey: E[$SSR_{GE}$]=1.0. Green: E[$SSR_{GE}$]<1.0; Orange: E[$SSR_{GE}$]>1.0

Simulated (bold) and Expected (plain) Values of SSR$_{GE}$ assuming large clusters ($r=50$) and 63 Combinations of $\rho_x$ and $\rho_y$: GEE Linear Regression Modeling Framework with exchangeable working correlation structure.

| $\rho_x$ | $\rho_y$ | | | | | | |
|---|---|---|---|---|---|---|---|
| | 0 | .05 | .10 | .25 | .50 | .75 | .90 |
| -.02 | **1.000** | **0.951** | **0.902** | **0.754** | **0.507** | **0.257** | **0.104** |
| $(-1/r)$ | 1.000 | 0.950 | 0.900 | 0.750 | 0.500 | 0.250 | 0.100 |
| 0 | **0.999** | **0.964** | **0.917** | **0.768** | **0.517** | **0.262** | **0.106** |
| | 1.000 | 0.964 | 0.915 | 0.764 | 0.510 | 0.255 | 0.102 |
| .05 | **0.998** | **0.999** | **0.956** | **0.806** | **0.544** | **0.275** | **0.112** |
| | 1.000 | 1.000 | 0.956 | 0.802 | 0.536 | 0.268 | 0.107 |
| .10 | **0.995** | **1.038** | **1.000** | **0.847** | **0.572** | **0.291** | **0.118** |
| | 1.000 | 1.039 | 1.000 | 0.844 | 0.565 | 0.283 | 0.113 |
| .25 | **0.990** | **1.171** | **1.158** | **1.001** | **0.684** | **0.348** | **0.141** |
| | 1.000 | 1.177 | 1.161 | 1.000 | 0.675 | 0.339 | 0.136 |
| .50 | **0.984** | **1.491** | **1.573** | **1.444** | **1.010** | **0.519** | **0.212** |
| | 1.000 | 4.511 | 1.588 | 1.446 | 1.000 | 0.507 | 0.204 |
| .75 | **0.978** | **2.065** | **2.464** | **2.590** | **1.936** | **1.025** | **0.421** |
| | 1.000 | 2.111 | 2.512 | 2.613 | 1.926 | 1.000 | 0.405 |
| .90 | **0.976** | **2.693** | **3.754** | **4.965** | **4.330** | **2.457** | **1.038** |
| | 1.000 | 2.771 | 3.857 | 5.063 | 4.333 | 2.406 | 1.000 |
| 1.00 | **0.974** | **3.411** | **5.884** | **13.146** | **25.451** | **38.051** | **45.734** |
| | 1.000 | 3.500 | 6.000 | 13.500 | 26.000 | 38.500 | 46.000 |

Note. Simulated data for each combination of $\rho_x$ and $\rho_y$ included $m$=75 clusters, $n$=51 units per cluster, $N$=3825, and 10K replicate samples. Simulated SSR$_{GE}$ estimated from comparison of std errs estimated from GEE linear model versus cluster-naïve model. Cell shading codes for combinations of $\rho_x$ and $\rho_y$ values: Grey: E[$SSR_{GE}$]=1.0. Green: E[$SSR_{GE}$]<1.0; Orange: E[$SSR_{GE}$]>1.0

Simulation Results: Multivariate Logistic Models. Expected vs Simulated Sample Size Ratios and Statistical Power across the SS & GEE Modeling Frameworks: $r=2$, $m=500$, $\rho_y=0.349$ (population average), & 50K replicate samples.

| $x$ variables [a] | | Survey Sampling Modeling Framework [b] | | | | | GEE Modeling Framework [c] | | | | | $\dfrac{Neff_{GE}}{Neff_{SS}}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $SSR_{SS}$ | | $Neff_{SS}$ | Power | | $SSR_{GE}$ | | $Neff_{GE}$ | Power | | |
| | $\rho_x$ | expected [d] | simulated [e] | expected [f] | expected [g] | simulated [h] | expected [d] | simulated [e] | expected [f] | expected [g] | simulated [h] | |
| $x_1$ | -1.00 | 0.651 | 0.653 | 1535 | .715 | .701 | 0.651 | .647 | 1535 | .715 | .709 | 1.00 |
| $x_2$ | 0 | 1.000 | 1.000 | 1000 | .532 | .511 | 0.878 | .872 | 1138 | .586 | .572 | 1.14 |
| $x_3$ | 0.25 | 1.087 | 1.086 | 920 | .499 | .477 | 0.962 | .956 | 1039 | .547 | .530 | 1.13 |
| $x_4$ | 0.50 | 1.174 | 1.173 | 852 | .469 | .453 | 1.064 | 1.059 | 934 | .505 | .496 | 1.10 |
| $x_5$ | 0.75 | 1.261 | 1.259 | 793 | .443 | .427 | 1.189 | 1.187 | 841 | .464 | .452 | 1.06 |
| $x_6$ | 1.00 | 1.349 | 1.346 | 742 | .420 | .400 | 1.349 | 1.351 | 742 | .420 | .403 | 1.00 |

[a] $x_1$-$x_6$ jointly uncorrelated & each unit-standardized; column '$\rho_x$' reports the population intra-cluster correlation value of each $x$ variable.

[b] Survey sampling modeling framework: fixed effect parameters estimated assuming independent observations and standard errors estimated via the Taylor series method ($\hat{\sigma}_{Taylor}$) assuming a compound-symmetric covariance structure.

[c] GEE Modeling framework: fixed effect parameters and model-based standard errors ($\hat{\sigma}_{Mod}$) estimated by a GEE linear model with compound symmetric working correlation structure.

[d] $SSR_{SS} = 1 + r\rho_x\rho_y$ and $SSR_{GE} = \left[1 - \rho_y + r\rho_y(1 - \rho_y)\right]/\left[1 - \rho_y + r\rho_y(1 - \rho_x)\right]$, where $r=1$, $\rho_x$ values as tabled, and $\rho_y=0.349$.

[e] the quantity $\left(\hat{\sigma}_{Tay_i}/\hat{\sigma}_{ind_i}\right)^2$ or $\left(\hat{\sigma}_{Mod_i}/\hat{\sigma}_{ind_i}\right)^2$, as appropriate, averaged over $i$=1 to 50K replicate samples, where $\hat{\sigma}_{ind_i}$ denotes the corresponding standard error estimate assuming independent observations.

[f] $Neff_{SS}=N \div SSR_{SS}$ and $Neff_{GE}=N \div SSR_{GE}$, where $N=1000$.

[g] Statistical power calculated by PASS assuming $Neff_{SS}$ or $Neff_{GE}$, as appropriate, two-tailed $\alpha=.05$, fixed effect of $x$ equal $b_{pa}\approx0.130$ (population average), $P[y=1|x=0]\approx0.533$ (population average), $P[y=1|$any $x=1]\approx0.565$ (population average), and $\sigma_x=1.0$.

[h] Simulated statistical power represents the proportion of corresponding replicate-sample fixed effect parameter estimates with test $p$-value $<.05$.

# 3-Level Clustered Sampling Design: Simple Example

Example: Multisite RCT w/ sites (s), people (p), measures (m)

Levels and Sample sizes
. Sites        (s) @ Level3: $n3 = 30$ sites

. People      (p) @ Level2: $n2 = 10$ people per site. Units of randomization

. Measures (m) @ Level1: $n1 = 2$ assessment times per person

x variables
x3 is a site-level (L3) continuous covariate
   . x3 has positive between site variation and zero within site variation
   . $\rho_{x3} = 1.0.$                              The intra-site correlation of x3

x2 is the person-level (L2) binary experimental group indicator
   . Assume x2 has zero between-site variation
   . $\rho_{x2} = -1/(10-1) = -.\overline{111}.$        The intra-site correlation of x2

x1 is the binary assessment time indicator at L1
   . x1 has zero between-person variation
   . $\rho_{x1} = -1/(2-1) = -1.$              The intra-person correlation of x1

# 3-Level Clustered Sampling Design: Simple Example

Example: Multisite RCT w/ sites (s), people (p), measures (m)

Common types of $\rho_y$ estimates reported from a 3-level model

$$\rho_{y.s} = \frac{\sigma_{y.s}^2}{\sigma_{y.s}^2 + \sigma_{y.p}^2 + \sigma_{y.m}^2} \quad \text{Proportion of } y \text{ var. attributable to sites}$$

$$\rho_{y.s\&p} = \frac{\sigma_{y.s}^2 + \sigma_{y.p}^2}{\sigma_{y.s}^2 + \sigma_{y.p}^2 + \sigma_{y.m}^2} \quad \text{Prop. of } y \text{ var attributable to sites \& pts}$$

$$\rho_{y.p} = \frac{\sigma_{y.p}^2}{\sigma_{y.s}^2 + \sigma_{y.p}^2 + \sigma_{y.m}^2} \quad \text{Prop. of } y \text{ var. attributable to patients}$$

Both $\rho_{y.s\&p}$ and $\rho_{y.s}$ may be described as 'ICC at Level2'

. When reading the literature, be clear whether $\rho_{y.s\&p}$ or $\rho_{y.s}$ is reported

. When reporting, be clear whether you are reporting $\rho_{y.s\&p}$ or $\rho_{y.s}$

# 3-Level Clustered Sampling Design: Simple Example

Example: Multisite RCT w/ sites (s), people (p), measures (m)

1. SSR for $x3$. A site variable w/ $\rho_{x3}=1$ has a fully a between-site effect

My initial, incorrect conjecture

$$\text{SSR}_{x3} = \text{SSR}_{b(s)} = 1 + (10 \times 2 - 1)\rho_{y.s}$$

Correction

$$\text{SSR}_{x3} = \text{SSR}_{b(s)} = 1 + (10 \times 2 - 1)\rho_{y.s②},$$

where $\rho_{y.s②}$ is the intra-site correlation of $y$ estimated from a
2-level model that excludes Level2 cluster indicators (persons).
I.e., only top-level (site) clusters are modeled, i.e.,

$$\rho_{y.s②} = \frac{\sigma^2_{y.s②}}{\sigma^2_{y.s②}+\sigma^2_{y.m②}}$$    Proportion of $y$ var. attrib. to sites: from 2-level model

# 3-Level Clustered Sampling Design: Simple Example

Example: Multisite RCT w/ sites (s), people (p), measures (m)

1. SSR for $x3$, which has a fully a between-site effect

. Given a 3-level data structure, when a model ignores the 2nd level,
the Level2 variation is distributed to both Level3 & Level1 (Moerbeek).

| From a 3-level model: obtain prop. of variance explained at Levels 2 & 3 |
|---|
| $\rho_{y.s} = .05$        Proportion of $y$ variance attributable to sites |
| $\rho_{y.p} = .10$        Proportion of $y$ variance attributable to people |
| $\rho_{y.p} = .85$        Proportion of $y$ variance attributable to measures |

Given $\rho_{y.s} = .05$, $\rho_{y.p} = .10$, $n2=10$, and $n1=2$, estimate $\rho_{y.s②}\ldots$

$$\rho_{y.s②} = \rho_{y.s} + \rho_{y.p}(n1 - 1)/(n2 \times n1 - 1) = .05 + .10/19 = .055263$$

$$\text{SSR}_{x3} = 1 + (10 \times 2 - 1) \times 055263 = 2.05$$

Moerbeek, M (2004). The Consequence of Ignoring a Level of Nesting in Multilevel Analysis. *Multivariate Behavioral Research*, 39. 129-149.

# 3-Level Clustered Sampling Design: Simple Example

Example: Multisite RCT w/ sites (s), people (p), measures (m)

2. SSR for $x2$. A Level 2 variable w/ $\rho_{x2} = -1/r$
    has a fully within-site/fully between-people effect

My initial, incorrect conjecture

$$\text{SSR}_{x2} = \text{SSR}_{w(s)} \times \text{SSR}_{b(p)} = \left(1 - \rho_{y.s}\right) \times \left[1 + (2-1)\rho_{y.p}\right]$$

Correction

$$\text{SSR}_{x2} = \text{SSR}_{w(s)} \times \text{SSR}_{b(p)} = \left(1 - \rho_{y.s}\right) \times \left[1 + \left(\frac{2-1}{1-\rho_{y.s}}\right)\rho_{y.p}\right]$$

$$= \left(1 - \rho_{y.s}\right) \times \left[1 + (2-1)\right]\rho_{y.p\cancel{s}} ,$$

where $\rho_{y.p\cancel{s}}$ is estimated via var. components from a 3-level model, i.e.,

$\rho_{y.p\cancel{s}} = \dfrac{\sigma^2_{y.p}}{\sigma^2_{y.p} + \sigma^2_{y.m}}$   Prop. *y* var. attrib. to people, removing site variation

Given $\rho_{y.s} = .05$, $\rho_{y.p} = .10$, $\rho_{y.m} = .85$, & *n1*=2.

$\rho_{y.p\cancel{s}} = .10/(.10 + .85) = .10526$

$\text{SSR}_{x2} = (1 - .05) \times [1 + (2-1)].10526 = 1.05$

# 3-Level Clustered Sampling Design: Simple Example

Example: Multisite RCT w/ sites (s), people (p), measures (m)

3. SSR for $x1$, which has a fully within-site/fully within-people effect

My initial, incorrect conjecture
$$\text{SSR}_{x1} = \text{SSR}_{w(s)} \times \text{SSR}_{w(p)} = \left(1 - \rho_{y.s}\right)\left(1 - \rho_{y.p}\right)$$

Correction
$$\text{SSR}_{x1} = \text{SSR}_{w(s)} \times \text{SSR}_{w(p)} = \left(1 - \rho_{y.s}\right)\left(1 - \rho_{y.p\cancel{s}}\right),$$

where $\rho_{y.p\cancel{s}}$ is estimated as described above

---

Given $\rho_{y.s} = .05$ and $\rho_{y.p\cancel{s}} = .10526$

$$\text{SSR}_{x1} = (1 - .05)(1 - .10526) = 0.85$$

---

# 3-Level Clustered Sampling Design: Simple Example

Simulated data from 3-Level Linear Mixed Model w/ 2K replicate samples

. $n3$=30 sites (L3), $n2$=10 subjects/site (L2), $n1$=2 assessments/subject

. $x3$: a normal random variate at Level3

. $x2$: a binary randomized group indicator at Level2

$x1$: a binary assessment time indicator at Level1

| $x$(level) | $\rho_x$ | $x$ effect @ | | $\rho_y$ | $\rho_y$ adjustment | | SSR$_{GE}$ | |
| | | L3 | L2 | | $\rho_{y.s②}$ | $\rho_{y.p\bar{s}}$ | expected | simulated |
|---|---|---|---|---|---|---|---|---|
| $x3$ | 1.0 | btw | -- | 0.05 | .05526 | -- | 2.05 | 2.056 |
| $x2$ | $-1/r$ | w/in | btw | 0.10 | -- | .10526 | 1.05 | 1.053 |
| $x1$ | $-1/r$ | w/in | w/in | 0.85 | -- | -- | 0.85 | 0.854 |

SSR expected value calculations

$$\text{SSR}_{\text{GE}.x3} = \text{SSR}_{\text{b(s)}} \qquad\qquad = 1 + (10 \times 2 - 1) \times 0.05526 = 2.05$$

$$\text{SSR}_{\text{GE}.x2} = \text{SSR}_{\text{w(s)}} \cdot \text{SSR}_{\text{b(p)}} = (1 - 0.05) \times [1 + (2 - 1) \times 0.10526] = 1.50$$

$$\text{SSR}_{\text{GE}.x1} = \text{SSR}_{\text{w(s)}} \times \text{SSR}_{\text{w(p)}} \qquad = (1 - 0.05) \times (1 - 0.10526) = 0.85$$

SSR simulated values are relative size of std errs from LMM &

Independence models, averaged over $i$=1 to 2K replicates: $\left(\hat{\sigma}_{\text{LMM}_i} / \hat{\sigma}_{\text{Ind}_i}\right)^2$

# 3-Level Clustered Sampling Design: Alternative Design A

Simulated data from 3-Level Linear Mixed Model w/ 2K replicate samples

. $n3=30$ sites (L3), $n2=10$ subjects/site (L2), $n1=2$ assessments/subject

. $x3$: a normal random variate at Level3

. $x2$: a binary randomized group indicator at Level2

$x1$: a binary assessment time indicator at Level1

| $x$(level) | $\rho_x$ | $x$ effect @ | | $\rho_y$ | $\rho_y$ adjustment | | SSR$_{GE}$ | |
| | | L3 | L2 | | $\rho_{y.s②}$ | $\rho_{y.p\bar{s}}$ | expected | simulated |
|---|---|---|---|---|---|---|---|---|
| $x3$ | 1.0 | btw | -- | 0.2 | .2368 | -- | 5.50 | 5.564 |
| $x2$ | $-1/r$ | w/in | btw | 0.7 | -- | .875 | 1.50 | 1.523 |
| $x1$ | $-1/r$ | w/in | w/in | 0.1 | -- | -- | 0.10 | 0.102 |

SSR expected value calculations

$$\text{SSR}_{\text{GE}.x3} = \text{SSR}_{\text{b(s)}} \qquad = 1 + (10 \times 2 - 1) \times 0.2368 = 5.50$$

$$\text{SSR}_{\text{GE}.x2} = \text{SSR}_{\text{w(s)}} \cdot \text{SSR}_{\text{b(p)}} = (1 - 0.20) \times [1 + (2 - 1) \times 0.875] = 1.50$$

$$\text{SSR}_{\text{GE}.x1} = \text{SSR}_{\text{w(s)}} \cdot \text{SSR}_{\text{w(p)}} \qquad = (1 - 0.20) \times (1 - 0.875) = 0.10$$

SSR simulated values are relative size of std errs from LMM &

Independence models, averaged over $i$=1 to 2K replicates: $\left(\hat{\sigma}_{\text{LMM}_i}/\hat{\sigma}_{\text{Ind}_i}\right)^2$

# 3-Level Clustered Sampling Design: Alternative Design B

Simulated data from 3-Level Linear Mixed Model w/ 2K replicate samples

. $n3=30$ sites (L3), $n2=10$ subjects/site (L2), $n1=5$ assessments/subject

. $x3$: a normal random variate at Level3

. $x2$: a binary randomized group indicator at Level2

$x1$: a uniform categorical assessment time indicator at Level1

| $x$(level) | $\rho_x$ | $x$ effect @ | | $\rho_y$ | $\rho_y$ adjustment | | SSR$_{GE}$ | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | L3 | L2 | | $\rho_{y.s②}$ | $\rho_{y.p\cancel{s}}$ | expected | simulated |
| $x3$ | 1.0 | btw | -- | 0.05 | .0908 | -- | 5.45 | 5.479 |
| $x2$ | $-1/r$ | w/in | btw | 0.50 | -- | .5263 | 2.95 | 2.964 |
| $x1$ | $-1/r$ | w/in | w/in | 0.45 | -- | -- | 0.45 | 0.455 |

SSR expected value calculations

$$\text{SSR}_{\text{GE}.x3} = \text{SSR}_{\text{b(s)}} = 1 + (10 \times 5 - 1) \times 0.0908 = 5.45$$

$$\text{SSR}_{\text{GE}.x2} = \text{SSR}_{\text{w(s)}} \cdot \text{SSR}_{\text{b(p)}} = (1 - 0.05) \times [1 + (5 - 1) \times 0.5263] = 2.95$$

$$\text{SSR}_{\text{GE}.x1} = \text{SSR}_{\text{w(s)}} \cdot \text{SSR}_{\text{w(p)}} = (1 - 0.05) \times (1 - 0.5263) = 0.450$$

SSR simulated values are relative size of std errs from LMM &

Independence models, averaged over $i$=1 to 2K replicates: $\left(\hat{\sigma}_{\text{LMM}_i} / \hat{\sigma}_{\text{Ind}_i}\right)^2$

# Summary

Proper use of SSRs requires consideration of $\rho_x$

$\text{SSR}_b$ [$\text{SSR}_b = 1 + r\rho_y$] is well known, but perhaps over-applied.
Improper application of $\text{SSR}_b$ can lead to
    substantially under-estimated power, w/ cost and ethical implications

When $-1/r < p_x < 1$, choose GEE/GLMM over SS modeling framework

All results reported here assumed
    compound symmetric correlation structure of $x$ and $y$

Power analysis for 3-level logistic models entails a few more wrinkles,
    mostly regarding estimation of population average $\rho_y$ values
    (a future talk)

Some additional details and a quiz w/ answers included in the Appendix

## Thank you

# 2-Level Clustered Sampling Designs: $\rho_x$: Quiz Yourself

**Pre-post design**

    Goal: test pre-post mean $y$ difference in a one-arm longitudinal trial
    Will the pre-post comparison have a between- or within-cluster effect?
    What is the value of $\rho_x$?

**Clustered sample of teachers and their current students**

    Goal: regress students' SAT ($y$) onto teacher's years of experience ($x$)

    Will teacher experience have a between- or within-cluster effect?

    What is the value of $\rho_x$?

**Multisite RCT. Randomization of patients within each site**

. Is the intervention group effect a between- or within-cluster effect?

. What can be said about the expected $\rho_x$ value?

**Observational study with geographic cluster sampling**

    Goal: regress smoking status ($y$) onto respondent income ($x$)

    Is income expected to have between- and/or within-cluster effects?

# Appendix A: Derivation of SSR$_{GE}$ from BLS (2011) Eq. 3.5.

For applications of the GLMM or GEE modeling frameworks,
  BLS Eq. 3.5 relates effective sample sizes under assumptions of
  . (i) CS correlation structures of $x$ and $y$ versus
  . (ii) CS correlation structure of $y$ with $\rho_x = 1$

$$\text{SSR}_{\text{BLS}} = \frac{\text{Neff}_{\rho_x, \rho_y}}{\text{Neff}_{\rho_x=1, \rho_y}} = \frac{1 - \rho_y + r\rho_y - r\rho_y\rho_x}{1 - \rho_y} \qquad \text{[BLS Eq. 3.5]}$$

A SSR that relates observed $N$ to $Neff$ assuming CS correlation structures
  of $x$ and $y$ can be derived from BLS Eq. 3.5, as follows.

$$\text{SSR}_{\text{GE}} = \frac{N}{\text{Neff}_{\rho x, \rho y}} = \frac{\dfrac{N}{\text{Neff}_{\rho x=1, \rho y}}}{\dfrac{\text{Neff}_{\rho x, \rho y}}{\text{Neff}_{\rho x=1, \rho y}}} = \frac{\text{SSR}_{\text{b}}}{\text{SSR}_{\text{BLS}}} =$$

$$= \frac{1 + r\rho_y}{\dfrac{1 - \rho_y + r\rho_y - r\rho_y\rho_x}{1 - \rho_y}} = \frac{(1 - \rho_y)(1 + r\rho_y)}{1 - \rho_y + r\rho_y - r\rho_y\rho_x}$$

$$= \frac{1 - \rho_y + r\rho_y(1 - \rho_y)}{1 - \rho_y + r\rho_y(1 - \rho_x)}$$

$$\text{SSR}_{\text{GE}} = \frac{(1-\rho_y)(1+r\rho_y)}{1-\rho_y[1-r(1-\rho_x)]} = \frac{1-\rho_y+r\rho_y(1-\rho_y)}{1-\rho_y+r\rho_y(1-\rho_x)}$$

if $\rho_x = -1/r$, then

$$\text{SSR}_{\text{GE}} = \frac{(1-\rho_y)(1+r\rho_y)}{1-\rho_y+r\rho_y(1-\rho_x)} = \frac{\text{SSR}_{\text{w}} \times \text{SSR}_{\text{b}}}{\text{SSR}_{\text{b}}} = \text{SSR}_{\text{w}}$$

if $\rho_x = 1$, then

$$\text{SSR}_{\text{GE}} = \frac{(1-\rho_y)(1+r\rho_y)}{1-\rho_y+r\rho_y(1-\rho_x)} = \frac{\text{SSR}_{\text{w}} \times \text{SSR}_{\text{b}}}{\text{SSR}_{\text{w}}} = \text{SSR}_{\text{b}}$$

if $\rho_x = 0$, then

$$\text{SSR}_{\text{GE}} = \frac{(1-\rho_y)(1+r\rho_y)}{1-\rho_y+r\rho_y(1-\rho_x)} = \frac{\text{SSR}_{\text{w}} \times \text{SSR}_{\text{b}}}{1+(r-1)\rho_y} \cong 1 - \rho_y^{(n_1/r)}, \quad \text{(where } n_1 = r+1\text{)}$$

I.e., as $r \to \infty$, $\dfrac{\text{SSR}_{\text{b}}}{1+(r-1)\rho_y} \to 1$, and $\text{SSR}_{\text{GE}} \to \text{SSR}_{\text{w}}$

Additionally, if $r$=1 then $\text{SSR}_{\text{GE}} = \dfrac{(1-\rho_y)(1+\rho_y)}{1} = 1 - \rho_y^2$

# 2-Level Clustered Sampling Designs: $\rho_x$: Quiz Answers

**Paired** t-test

   Goal: test pre-post mean *y* difference in a one-arm longitudinal trial

Here, respondents define the clusters and repeated measures (pre and post) are nested within respondents.

   Will the pre-post comparison have a between- or within-cluster effect?

Pre-post indicator (*x*) is defined at Level1. Each cluster (person) has 2 assessments: one pre (*x*=0) and one post (*x*=1). There is zero between-cluster variation of *x* and positive within-cluster variation of *x*. Therefore, the pre-post comparison is a fully within-cluster (within-person) effect.

   What is the value of $\rho_x$?

In this case $\rho_x = -1/(2-1) = -1.0$

# 2-Level Clustered Sampling Designs: $\rho_x$: Quiz Answers

Clustered sample of teachers and their current students

Goal: regress students' SAT (*y*) onto teacher's years of experience (*x*)

Teachers are clusters (Level2) and students (Level1) are nested within teachers

Will teacher experience have a between- or within-cluster effect?

Teacher experience is a Level2 variable. Therefore, teacher experience will have a fully between-cluster (between-teacher) effect.

What is the value of $\rho_x$?

Teacher experience will have positive between-cluster variation and zero within-cluster variation. Therefore, $\rho_x = 1$

# 2-Level Clustered Sampling Designs: $\rho_x$: Quiz Answers

Multisite RCT. Randomization of patients within each site

Sites are clusters (Level2) and patients (Level1) are nested in sites

. Is the intervention group effect a between- or within-cluster effect?

Intervention group indicator is a Level1 variable. Therefore, if the proportionate representation of Trt vs Ctrl assignment is identical across clusters, then the group effect will be a fully within-cluster effect. If the proportionate representation of group assignment varies slightly across site clusters, then a small of amount of between-cluster *x* variation will exist and the group comparison will be *mostly* a within-cluster effect.

. What can be said about the expected $\rho_x$ value?

$\rho_x = -1/r$ if the proportionate allocation to Trt v Ctrl is identical across clusters.

If proportionate treatment assignment varies across clusters, then $\rho_x > -1/r$ . Given sufficient cluster size and between-site variation, $\rho_x$ could become positive.

# 2-Level Clustered Sampling Designs: $\rho_x$: Quiz Answers

Observational study with geographic cluster sampling

   Goal: regress smoking status (*y*) onto respondent income (*x*)

Geographic areas are clusters (Level2) and respondents (Level1) are nested within clusters.

   Is income expected to have between- and/or within-cluster effects?

We expect that respondent income (*x*) will have both between- and within-cluster variation. Therefore, we expect $0 < \rho_x < 1$, which means that income (*x*) can have both between- and within-cluster effects on smoking status.