

Biostatistics 140.656, 2017-18
Lab 5 Solution

Topics:

- Simulating multi-level data
- Estimating power to detect an effect of interest given fixed sample size

Learning Objectives:

Students who successfully complete this lab will be able to:

- Simulate 2-level continuous multilevel data with specific between, within and contextual effects
- Estimate the power to detect a significant contextual effect given fixed sample size

If time permits, we will consider a 2nd example:

- Simulate 2-level binary multilevel data where the goal is to compare the odds of the binary outcome across two exposure groups, assuming no contextual effect.
- Estimate the power to detect a significant odds ratio given fixed sample size.

Scientific Background:

Consider the data setting from Lab 2: You are an obstetrician interested in women's satisfaction with their labor and childbirth experiences.

For each woman, you will have the following:

- Maternal case-mix: score for how complicated the birth is likely to be; higher scores mean the birth is expected to be more complicated based on pre-existing conditions, prenatal care, and socioeconomic status. Assume this variable is a z-score with mean 0 and variance 1.
- Patient satisfaction score (Y): scores for patient satisfaction with labor and childbirth experiences; higher scores indicate greater satisfaction with the birth experience. Assume this variable is a z-score with mean 0 and variance 1.

Your goal is to determine if the context of the hospital matters; i.e. do women of the same maternal case-mix benefit from attending a hospital with lower than average maternal case-mix.

You are going to design an observational study where you will recruit women from M hospitals from a national network of hospitals with 500 hospitals. After identifying the M hospitals, you have the resources and time to interview 30 mothers. You need to figure out how many hospitals you need to recruit to the study to achieve roughly 80% power to detect a significant contextual effect of maternal case-mix.

Once you conduct your study, you will fit the following model: Let i denote the hospital ($i = 1, \dots, M$), j denote the mother within hospital i ($j = 1, \dots, 30$), Y_{ij} is the patient satisfaction score and X_{ij} is the casemix for woman j from hospital i .

$$Y_{ij} = \beta_0 + b_{0i} + \beta_1(X_{ij} - \bar{X}_i) + \beta_2\bar{X}_i + \varepsilon_{ij}, b_{0i} \sim N(0, \tau^2), \varepsilon_{ij} \sim N(0, \sigma^2)$$

The contextual effect is measured as: $\beta_2 - \beta_1$.

Note that from this model, you can derive the following:

$$\begin{aligned} \text{Var}(Y_{ij}) &= \text{Var}(b_{0i} + \varepsilon_{ij}) = \tau^2 + \sigma^2 \\ \text{ICC}(Y) &= \frac{\tau^2}{\tau^2 + \sigma^2} \end{aligned}$$

We can also think about decomposing the variation in the casemix variable in the same way:

$$\begin{aligned} \text{Var}(X_{ij}) &= 1 = \text{Var}[(X_{ij} - \bar{X}_{i.}) + \bar{X}_{i.}] = \text{Var}(X_{ij} - \bar{X}_{i.}) + \text{Var}(\bar{X}_{i.}) \\ \text{ICC}(X) &= \frac{\text{Var}(\bar{X}_{i.})}{\text{Var}(X_{ij} - \bar{X}_{i.}) + \text{Var}(\bar{X}_{i.})} \end{aligned}$$

Here is some additional relevant information that you know from women giving birth within the last month at a random sample of 20 hospitals within the national network of hospitals.

- a. The maternal case-mix varies between and within hospitals. You have data to suggest that ICC for maternal case-mix is 0.3; that is 30% of the variance in the maternal case-mix is attributable to differences between hospitals.
 - b. Patient satisfaction scores also cluster across hospitals. You have data to suggest that the ICC for patient satisfaction is 0.4; that is 40% of the variance in the maternal case-mix is attributable to differences between hospitals.
 - c. The within hospital relationship between patient satisfaction and maternal case-mix is -1.
 - d. The between hospital relationship between patient satisfaction and maternal case-mix is -1.25.
 - e. Therefore, the contextual effect is -0.25; that is, for two women with the same maternal case-mix, the woman delivering at the hospital with lower average maternal case-mix has an average patient satisfaction score that is 0.25 standard deviations greater than the women who delivers at the hospital with higher average material case-mix.
1. Based on this available information, walk through the “example1.do” file and identify the key steps to simulate the data.
 - Initialize M
 - Generate the number of subjects within each cluster, this can be fixed or random (we generate on average n subjects per cluster)
 - Generate the cluster specific features of the model:
 - Based on the ICC and total variance in Y, determine an estimate of the random intercept variance and sample a random intercept variance for each cluster.
 - Baseline the ICC and total variance in X, determine an estimate of the variance in the cluster mean X and generate a random sample of cluster means.
 - Generate the subject specific features of the model:
 - Within each cluster generate X based on the ICC and total variance in X and the cluster mean X
 - Within each cluster generate Y based on the specified between and within cluster association plus a randomly generated within cluster residual (the variance depends on the ICC and variance of Y)

Once you can generate data for a potential study, generate many hypothetical studies. For each study, fit the model and store the estimate of the contextual effect and p-value for the test that the contextual effect is 0. Determine for what proportion you would reject the null hypothesis and this is the estimated power.

2. Change the number of hospital included in the sample and identify the number of hospitals required to achieve 80% power to detect the contextual effect of -0.25.

If you set the number of hospitals to 140, then you have roughly 80% power to detect the contextual effect of -0.25.

```
. simulate contextual = r(estimate) se = r(se), reps(100) seed(734): ///
> example1 140 30 0 1 0.3 0 1 0.4 -1 -1.25
```

```
      command:  example1 140 30 0 1 0.3 0 1 0.4 -1 -1.25
contextual:    r(estimate)
              se:  r(se)
```

```
Simulations (100)
```

```
-----+----- 1 -----+----- 2 -----+----- 3 -----+----- 4 -----+----- 5
..... 50
..... 100
```

```
. gen ts = - abs(contextual) / se
```

```
. gen p = normal(ts)
```

```
. gen reject = p < 0.05
```

```
. summ reject
```

Variable	Obs	Mean	Std. Dev.	Min	Max
reject	100	.8	.4020151	0	1

3. What impact does reducing the number of women interviewed at each hospital have on the power to detect the contextual effect?

Here I generated 500 hypothetical studies each with 20 hospitals and 30 or 15 women selected within each hospital. Below are the results:

```
. simulate contextual = r(estimate) se = r(se), reps(500) seed(734):
///  
> example1 20 30 0 1 0.3 0 1 0.4 -1 -1.25
```

```
      command:  example1 20 30 0 1 0.3 0 1 0.4 -1 -1.25
contextual:    r(estimate)
      se:       r(se)
```

Simulations (500)

```
-----+--- 1 ----+--- 2 ----+--- 3 ----+--- 4 ----+--- 5
..... 50
..... 100
..... 150
..... 200
..... 250
..... 300
..... 350
..... 400
..... 450
..... 500
```

```
. gen ts = - abs(contextual) / se
```

```
. gen p = normal(ts)
```

```
. gen reject = p < 0.05
```

```
. summ reject
```

Variable	Obs	Mean	Std. Dev.	Min	Max
reject	500	.288	.4532846	0	1

```
.
. simulate contextual = r(estimate) se = r(se), reps(500) seed(734):
///  
> example1 20 15 0 1 0.3 0 1 0.4 -1 -1.25
```

```
      command:  example1 20 15 0 1 0.3 0 1 0.4 -1 -1.25
contextual:    r(estimate)
      se:       r(se)
```

Simulations (500)

```
-----+--- 1 ----+--- 2 ----+--- 3 ----+--- 4 ----+--- 5
```

```

..... 50
..... 100
..... 150
..... 200
..... 250
..... 300
..... 350
..... 400
..... 450
..... 500

```

```
. gen ts = - abs(contextual) / se
```

```
. gen p = normal(ts)
```

```
. gen reject = p < 0.05
```

```
. summ reject
```

Variable	Obs	Mean	Std. Dev.	Min	Max
reject	500	.292	.4551377	0	1

The impact on reducing the number of women within hospital has little to no effect on the power to detect a contextual effect. The driver for power in this case is the number of clusters.