# 2042 MLM Mini Project (Spring 2020)
## Group 1

May 13 2020

## Team members and division of work

Frank Jiang, Lisa Song, Yuyue Hua, Seeun Jang, Tong Jin

```
set.seed(2042001)
```

### Question 1

You will generate simulated data for a single school with 100 classrooms, each of which has 200 students.

a. Outcome for student $i$ in classroom $j$: $Y_{ij}$.

b. There is a single predictor, $X_{ij} \sim U(0,1)$ (uniform on $[0,1]$)

c. There is a classroom random effect, $\eta_j \sim N(0, \sigma_\eta^2)$, where $\sigma_\eta^2 = 2$.

d. Subject level error, $\varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2)$, where $\sigma_\varepsilon^2 = 2$.

e. `set.seed(2042001)` once at the beginning of your code.

f. Generate the random quantities in this order to ensure the same solution for everyone: $X$, $\eta_j$, $\varepsilon_{ij}$

g. The outcome has the following form (DGP, given the modeling parameters above):

$$ Y\_{ij} = 0 + 1X\_{ij} + \eta j + \varepsilon\{ij\} , \setminus $$
$$ \eta j \ N(0, \hat{} 2\{ \eta \}), \ \varepsilon\{ij\} \ N(0, \hat{} 2\{ \varepsilon \}), \text{indep.} $$

h. Generate a single simulated dataset (you will need a "classid" variable to track classrooms); you can optionally assign a "studentid")

i. **Important:** construct classid such that classrooms appear consecutively within the dataframe. As per: `rep(1:J,each=n_j)`

```
# Compute variables related to the equation
N_j <- 100 # Total number of classrooms
n_i <- 200 # Number of students in each classroom
N_i <- N_j * 200 # Total number of students

X_ij <- runif( # Single predictor: Uniform on [0, 1]
  N_i,
  min = 0,
  max = 1
```

```
)
eta_j <- rnorm( # Classroom random effect: Normal on (0, 2)
  N_j,
  mean = 0,
  sd = sqrt(2)
)
epsilon_ij <- rnorm( # Subject level error: Normal on (0, 2)
  N_i,
  mean = 0,
  sd = sqrt(2)
)

# Create equation elements
eta_j <- rep( # Assign classroom random effect to students in each classroom
  eta_j,
  each = n_i
)
Y_ij <- 0 + 1 * X_ij + eta_j + epsilon_ij

classid <- rep(
  1:N_j,
  each = n_i
)
studentid <- seq(
  1:N_i
)

# Create a dataframe to store all elements
classroom_sim <- data.frame(
  outcome = Y_ij,
  predictor = X_ij,
  cls_raneff = eta_j,
  subject_error = epsilon_ij,
  classid = classid,
  studentid = studentid,
  row.names = studentid
)
```

**Question 2**

Fit the model corresponding to the DGP on your simulated data.

```
# Fit the model in Q1
fit_q2 <- lmer(
  outcome ~ predictor + (1 | classid),
  data = classroom_sim
)
# Report the model fit
# summary(fit_q2)

# Calculate the coefficient estimate of slope on X
slope_X_q2 <- round(
  summary(fit_q2)$coefficients['predictor', 'Estimate'],
```

```
  digits = 4
)
```

a. Report coefficient estimate for slope on X.

**Response:** 0.9864.

b. Does a 95% confidence band for this coefficient estimate cover the "truth" that you used to generate the data?

```
CI_Q2_b <- confint(
  fit_q2,
  parm = "predictor",
  level = 0.95
)
lower <- round(
  CI_Q2_b[1],
  digits = 4
)
upper <- round(
  CI_Q2_b[2],
  digits = 4
)
```

**Response:** Yes. The confidence band (95% of confidence) for this coefficient estimate is between 0.9179 and 1.0549. This confirms that a 95% confidence band covers the "truth" of the slope, which is 1.

**Question 3**

3. Next, we simulate missing data in several ways. This is the first:

a. Make a copy of the data, then modify the copy following these instructions:

```
classroom_miss <- classroom_sim
```

b. Generate $Z_{ij} \sim \text{Bernoulli}(p)$, with $p = 0.5$.

```
Z_ij <- rbinom(
  N_i,
  size = c(0, 1),
  prob = 0.5
)
```

c. Set $Y_{ij}$ to NA when $Z_{ij} == 1$. This should look a lot like "MCAR" missingness.

```
classroom_miss$missing <- Z_ij
classroom_miss$outcome <- ifelse(
  classroom_miss$missing == 1,
  yes = NA,
  no = classroom_miss$outcome
)
```

d. Refit the model on the new data and report the coefficient estimate for slope on X. Look at the other parameter estimates as well.

```
fit_q3 <- lmer(
  outcome ~ predictor + (1 | classid),
  data = classroom_miss
)
summary(fit_q3)
```

```
## Linear mixed model fit by REML. t-tests use Satterthwaite's method [
## lmerModLmerTest]
## Formula: outcome ~ predictor + (1 | classid)
##    Data: classroom_miss
##
## REML criterion at convergence: 53333.6
##
## Scaled residuals:
##     Min      1Q  Median      3Q     Max
## -4.0171 -0.6742  0.0023  0.6631  3.7668
##
## Random effects:
##  Groups   Name        Variance Std.Dev.
##  classid  (Intercept) 1.878    1.371
##  Residual             2.000    1.414
## Number of obs: 14963, groups:  classid, 100
##
## Fixed effects:
##              Estimate Std. Error        df t value Pr(>|t|)
## (Intercept) 1.270e-02  1.390e-01 1.033e+02   0.091    0.927
## predictor   9.612e-01  4.029e-02 1.486e+04  23.856   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##           (Intr)
## predictor -0.146
```

```
slope_X_q3 <- round(
  summary(fit_q3)$coefficients['predictor', 'Estimate'],
  digits = 4
)

# Calculate the conficience band of this model
CI_Q3_d <- confint(
  fit_q3,
  parm = "predictor",
  level = 0.95
)
```

```
## Computing profile confidence intervals ...
```

```
lower <- round(
  CI_Q3_d[1],
  digits = 4
)
upper <- round(
  CI_Q3_d[2],
  digits = 4
)
```

   e. Do you see any real change in the $\beta_X$ estimate?

     **Response:** No. The $\beta_X$ estimate of the model with missing data values is 0.9612, which is very close to the original estimate, 0.9864.

       i. Does a 95% confidence band for this coefficient estimate cover the "truth" that you used to generate the data?

     **Response:** Yes. The confidence band (95% of confidence) for this coefficient estimate is between 0.8822 and 1.0402. This confirms that a 95% confidence band covers the "truth" of the slope, which is 1.

   f. What is the total sample size $N$ used in the model fit?

     **Response:** The total sample size used in this model fit is 14963.

**Question 4:**

Missing Data II: Make another copy of the original data, then modify the copy as follows: a. Generate $Z_{ij} \sim$ Bernoulli($X_{ij}$), with $X_{ij}$ your predictor generated previously. b. Set Y to NA when $Z_{ij} == 1$. This should look a lot like "MAR" missingness.

```
# Insert code the generate your data
```

   c. Refit the model on the new data and report the coefficient estimate for slope on X. Look at the other parameter estimates as well. **comment**

```
# Insert code to fit model and compute confidence interval
```

Response:

   d. Do you see any real change in the $\beta_X$ estimate?

       i. Does a 95% confidence band for this coefficient estimate cover the "truth" that you used to generate the data? **comment**

     Response:

   e. What is the total sample size $N$ used in the model fit? **comment**

     Response:

**Question 5:**

Missing Data III: Make another copy of the original data, then modify the copy as follows:

```
# Insert code to make a copy of the original data
```

a. First, define the expit function: `expit <- function(x) exp(x)/(1+exp(x))`

```
# Insert code to define expit function
```

b. Generate $Z_{ij} \sim$ Bernoulli($expit(Y_{ij})$), with $Y_{ij}$ your *outcome* generated previously.

c. Set Y to NA when $Z_{ij} == 1$. This should look like a violation of "MAR" missingness (missingness depedents on outcome and cannot be *simply* predicted with the predictor set – Y should be correlated with X, though, so it might not be too bad a violation).

```
# Insert code the generate your data
```

d. Refit the model on the new data and report the coefficient estimate for slope on X. Look at the other parameter estimates as well. **comment**

```
# Insert code to fit model and compute confidence interval
```

`Response:`

e. Do you see any real change in the $\beta_X$ estimate? **comment**

    i. Does a 95% confidence band for this coefficient estimate cover the "truth" that you used to generate the data? **comment**

Response:

f. What is the total sample size $N$ used in the model fit? **comment**

Response: