

Missingness Mechanisms

Comparing multiple imputation for a dataset under MCAR, MAR, MNAR mechanisms with varying degrees of missingness.

- Dataset is a linear model $y_i = \alpha + x_i\beta + \epsilon_i$ where parameter beta (true value = 1) is estimated
- Will observe if imputation strength decreases going from MCAR -> MAR -> MNAR w/ low missingness -> MNAR w/ high missingness as expected.

MCAR

```
MCAR.create.data <- function(beta = 1, sigma2 = 1, n = 200,
                             run = 1) {
  set.seed(seed = run)
  x <- rnorm(n)
  y <- beta * x + rnorm(n, sd = sqrt(sigma2))
  cbind(x = x, y = y)
}
```

```
MCAR.make.missing <- function(data, p = 0.5){
  rx <- rbinom(nrow(data), 1, p)
  data[rx == 0, "x"] <- NA
  data
}
```

```
MCAR.test.impute <- function(data) {
  imp <- mice(data, print = FALSE)
  fit <- with(imp, lm(y ~ x))
  tab <- summary(pool(fit), "all", conf.int = TRUE)
  as.numeric(tab[2, c("estimate", "2.5 %", "97.5 %")])
}
```

```
MCAR.simulate <- function(runs = 10) {
  res <- array(NA, dim = c(1, runs, 3))
  dimnames(res) <- list(c("MCAR"),
                        as.character(1:runs),
                        c("estimate", "2.5 %", "97.5 %"))
  for(run in 1:runs) {
    data <- MCAR.create.data(run = run)
    data <- MCAR.make.missing(data)
    res[1, run, ] <- MCAR.test.impute(data)
  }
  res
}
```

```
MCAR.res <- MCAR.simulate(10)
```

```
apply(MCAR.res, c(1, 3), mean, na.rm = TRUE)
```

MCAR code inspiration courtesy of)

```
##      estimate      2.5 %    97.5 %
```

```
## MCAR 0.9933456 0.8157855 1.170906
true <- 1
RB <- mean(MCAR.res[, "estimate"]) - true
PB <- 100 * abs((mean(MCAR.res[, "estimate"]) - true) / true)
CR <- mean(MCAR.res[, "2.5 %"] < true & true < MCAR.res[, "97.5 %"])
AW <- mean(MCAR.res[, "97.5 %"] - MCAR.res[, "2.5 %"])
RMSE <- sqrt(mean(MCAR.res[, "estimate"] - true)^2)
data.frame(RB, PB, CR, AW, RMSE)

##          RB          PB CR          AW          RMSE
## 1 -0.006654436 0.6654436 1 0.3551201 0.006654436
```

MAR

- MAR mechanisms involve different degrees of missingness for groups dependent on an observed variable in the dataset.
- I will run simulations for varied differences in missingness between the groups ($|p_2 - p_1|$) and see how multiple imputation performs.

```
# Creates function that create a dataset with observations for a
# MULTIPLE linear regression model with a categorical variable.
MAR.create.data <- function(beta = 1, sigma2 = 1, n = 20000,
                             run = 1, beta2 = 10, categorical_var = "gender",
                             category_one = "M", category_zero = "F") {
  set.seed(seed = run)
  x <- rnorm(n)
  # p represents probability of being in some group (E.g. Male = 1 vs Female = 0)
  z <- rbinom(n, 1, p = 0.5)
  temp <- cbind(x = x, z = z)
  # Note the beta2 term for categorical variable effect.
  y <- beta * x + rnorm(n, sd = sqrt(sigma2)) + beta2 * z
  z[z == 0] <- category_zero
  z[z == 1] <- category_one
  # cat_var = categorical_var
  data <- cbind(x = x, y = y, categorical_var = z)
  # colnames(data) <- c("x", "y", categorical_var)
  # sort df by categorical variable.
  df <- as.data.frame(data) %>%
    arrange(categorical_var)
}
```

- Creates function that removes p_0 % of data from category zero,
-

p_1 % of data from category_one. Note this is effectively MCAR within each group (definition of MAR). Note if $p_0 = p_1$, have a MAR mechanism.

```
MAR.make.missing <- function(data, p_0 = 0.4, p_1 = 0.6){
  # don't generate nrow(data) times, generate count of category times.
  counts <- count(as.data.frame(data), categorical_var)
  num_cat_zero <- counts[1,2]
```

```

num_cat_one <- counts[2,2]
r_zero <- rbinom(counts[1,2], 1, p_0)
r_one <- rbinom(counts[2,2], 1, p_1)
data[1:num_cat_zero,1][r_zero == 1] <- NA
data[(num_cat_zero + 1):nrow(data), 1][r_one == 1] <- NA
data
}

```

*# Function that calls mice (applying imputation) and applies Rubin's Rules,
and creates 95% confidence intervals for parameter*

```

MAR.test.impute <- function(data) {
  # Convert numerical vars to doubles (from character).
  data_num <- as.data.frame(apply(data[,c(1:2)], 2, as.numeric))
  # create copy
  data_num_complete <- data_num
  # add the categorical var
  data_num_complete$categorical_var <- data[,3]
  imp <- mice(data_num_complete, print = FALSE)
  fit <- with(imp, lm(y ~ x + categorical_var))
  tab <- summary(pool(fit), "all", conf.int = TRUE)
  as.numeric(tab[2, c("estimate", "2.5 %", "97.5 %")])
}

```

```

simulate <- function(runs = 10) {
  res <- array(NA, dim = c(5, runs, 3))
  dimnames(res) <- list(c("No_miss", "MCAR", "lightMAR", "moderateMAR", "extremeMAR"),
                        as.character(1:runs),
                        c("estimate", "2.5 %", "97.5 %"))
  for(run in 1:runs) {
    data <- MAR.create.data(run = run)
    none_data <- MAR.make.missing(data, p_0 = 0, p_1 = 0)
    MCAR_data <- MAR.make.missing(data, p_0 = 0.5, p_1 = 0.5)
    light_data <- MAR.make.missing(data, p_0 = 0.4, p_1 = 0.6)
    moderate_data <- MAR.make.missing(data, p_0 = 0.3, p_1 = 0.7)
    heavy_data <- MAR.make.missing(data, p_0 = 0.2, p_1 = 0.8)
    res[1, run, ] <- MAR.test.impute(none_data)
    res[2, run, ] <- MAR.test.impute(MCAR_data)
    res[3, run, ] <- MAR.test.impute(light_data)
    res[4, run, ] <- MAR.test.impute(moderate_data)
    res[5, run, ] <- MAR.test.impute(heavy_data)
  }
  res
}

```

```
res <- simulate(10)
```

```
apply(res, c(1, 3), mean, na.rm = TRUE)
```

```

##           estimate      2.5 %   97.5 %
## No_miss    0.9979153 0.9841060 1.011724
## MCAR        0.9826710 0.9562152 1.009127
## lightMAR    0.9863792 0.9614206 1.011338
## moderateMAR 0.9857348 0.9569590 1.014511
## extremeMAR  0.9841665 0.9556249 1.012708

```

```

true <- 1
RB <- rowMeans(res[, "estimate"]) - true
PB <- 100 * abs((rowMeans(res[, "estimate"]) - true) / true)
CR <- rowMeans(res[, "2.5 %"] < true & true < res[, "97.5 %"])
AW <- rowMeans(res[, "97.5 %"] - res[, "2.5 %"])
RMSE <- sqrt(rowMeans((res[, "estimate"] - true)^2))
data.frame(RB, PB, CR, AW, RMSE)

```

```

##           RB      PB  CR      AW      RMSE
## No_miss   -0.002084746 0.2084746 0.9 0.02761842 0.008541116
## MCAR       -0.017329035 1.7329035 0.6 0.05291144 0.020096735
## lightMAR   -0.013620804 1.3620804 0.9 0.04991714 0.016401763
## moderateMAR -0.014265179 1.4265179 0.9 0.05755165 0.017432401
## extremeMAR -0.015833519 1.5833519 0.7 0.05708324 0.019741332

```

```

MNAR.create.data <- function(beta = 1, sigma2 = 1, n = 200,
                             run = 1) {
  set.seed(seed = run)
  x <- rnorm(n)
  y <- beta * x + rnorm(n, sd = sqrt(sigma2))
  as.data.frame(cbind(x = x, y = y))
}

```

```

# Create missingness in x values greater than median with specified probability.
MNAR.make.missing <- function(data, prob_missing_higher = 0.2){
  higher <- data$x[data$x > median(data$x)]
  data$x[data$x > median(data$x)] = ifelse(sample(
    c(T, F), length(data$x[data$x > median(data$x)]), replace=T,
    prob=c(prob_missing_higher, 1 - prob_missing_higher)),
    NA,
    data$x[data$x > median(data$x)])
  data
}

```

```

MNAR.test.impute <- function(data, m = 5) {
  imp <- mice(data, m = m, print = FALSE)
  fit <- with(imp, lm(y ~ x))
  tab <- summary(pool(fit), "all", conf.int = TRUE)
  as.numeric(tab[2, c("estimate", "2.5 %", "97.5 %")])
}

```

```

simulate <- function(runs = 10) {
  res <- array(NA, dim = c(5, runs, 3))
  dimnames(res) <- list(c("lightest_MNAR", "light_MNAR", "moderate_MNAR",
    "heavy_MNAR", "heaviest_MNAR"),
    as.character(1:runs),
    c("estimate", "2.5 %", "97.5 %"))
  for(run in 1:runs) {
    data <- MNAR.create.data(run = run)
    lightest_data <- MNAR.make.missing(data, prob_missing_higher = 0.2)
    lighter_data <- MNAR.make.missing(data, prob_missing_higher = 0.4)
  }
}

```

```

moderate_data <- MNAR.make.missing(data, prob_missing_higher = 0.6)
heavier_data <- MNAR.make.missing(data, prob_missing_higher = 0.8)
heaviest_data <- MNAR.make.missing(data, prob_missing_higher = 1.0)
res[1, run, ] <- MNAR.test.impute(lightest_data)
res[2, run, ] <- MNAR.test.impute(lighter_data)
res[3, run, ] <- MNAR.test.impute(moderate_data)
res[4, run, ] <- MNAR.test.impute(heavier_data)
res[5, run, ] <- MNAR.test.impute(heaviest_data)
# print(paste("run", run, "completed"))
}
res
}

```

```
res <- simulate(10)
```

MNAR

- Obtain estimates for beta in regression model along with 95% CI for each of the given methods.

```
apply(res, c(1, 3), mean, na.rm = TRUE)
```

```

##           estimate      2.5 %   97.5 %
## lightest_MNAR 1.0170161 0.8712786 1.162754
## light_MNAR    0.9924430 0.8412146 1.143671
## moderate_MNAR 1.0049190 0.8460597 1.163778
## heavy_MNAR    0.9982643 0.8015096 1.195019
## heaviest_MNAR 1.1254224 0.6748302 1.576015

```

- Obtain performance measures

```

true <- 1
RB <- rowMeans(res[, , "estimate"]) - true
PB <- 100 * abs((rowMeans(res[, , "estimate"]) - true) / true)
CR <- rowMeans(res[, , "2.5 %"] < true & true < res[, , "97.5 %"])
AW <- rowMeans(res[, , "97.5 %"] - res[, , "2.5 %"])
RMSE <- sqrt(rowMeans((res[, , "estimate"] - true)^2))
data.frame(RB, PB, CR, AW, RMSE)

##           RB           PB  CR           AW           RMSE
## lightest_MNAR 0.017016074 1.7016074 1.0 0.2914750 0.04527563
## light_MNAR    -0.007556998 0.7556998 1.0 0.3024568 0.04241408
## moderate_MNAR 0.004918960 0.4918960 1.0 0.3177184 0.05984147
## heavy_MNAR    -0.001735672 0.1735672 1.0 0.3935094 0.10187087
## heaviest_MNAR 0.125422412 12.5422412 0.8 0.9011844 0.33298022

```