# Assignment 2 – Solutions: Part 2 (STAR Dataset)

Applied Quantitative Methods II, UC3M

## 1. Data preparation

**a)** Load the data:

```
library(dplyr)
library(broom)
library(ggplot2)
library(modelsummary)


star = read.csv("https://raw.githubusercontent.com/franvillamil/AQM2/master/datasets/star.csv")
```

**b)** Create factor for class type:

```
star$classtype = factor(star$classtype,
  levels = 1:3,
  labels = c("Small", "Regular", "Regular+Aide"))
```

**c)** Create factor for race:

```
star$race = factor(star$race,
  levels = 1:6,
  labels = c("White", "Black", "Asian", "Hispanic",
             "Native American", "Other"))
```

**d)** Binary indicator for small class:

```
star$small = ifelse(star$classtype == "Small", 1, 0)
```

**e)** Report observations:

```
nrow(star)
```

```
## [1] 6325
```

```
sum(!is.na(star$g4reading))
```

```
## [1] 2353
```

```
sum(!is.na(star$g4math))
```

```
## [1] 2395
```

Many students lack 4th grade scores because they were not followed up to that point.

# 2. Comparing groups

**a)** Mean reading scores by class type:

```
star %>%
  group_by(classtype) %>%
  summarise(mean_reading = mean(g4reading, na.rm = TRUE))
```

```
## # A tibble: 3 x 2
##   classtype    mean_reading
##   <fct>               <dbl>
## 1 Small                723.
## 2 Regular              720.
## 3 Regular+Aide         721.
```

Students in small classes score highest on average.

**b)** Bivariate regression:

```
m1 = lm(g4reading ~ small, data = star)
tidy(m1)
```

```
## # A tibble: 2 x 5
##   term         estimate std.error statistic p.value
##   <chr>           <dbl>     <dbl>     <dbl>   <dbl>
## 1 (Intercept)    720.       1.30     554.       0
## 2 small            3.10      2.34       1.32    0.185
```

The coefficient on `small` is the average difference in 4th grade reading scores between students assigned to small classes vs. all others. Because this is an experiment, it has a causal interpretation.

**c)** The regression coefficient equals the difference in group means:

```
star %>%
  group_by(small) %>%
  summarise(mean_reading = mean(g4reading, na.rm = TRUE))
```

```
## # A tibble: 2 x 2
##   small mean_reading
##   <dbl>        <dbl>
## 1     0         720.
## 2     1         723.
```

**d)** Repeat for math:

```
m1_math = lm(g4math ~ small, data = star)
tidy(m1_math)
```

```
## # A tibble: 2 x 5
##   term         estimate std.error statistic p.value
##   <chr>           <dbl>     <dbl>     <dbl>   <dbl>
## 1 (Intercept)    709.       1.06     669.       0
## 2 small            0.591     1.91       0.310   0.756
```

Similar pattern: small-class students score higher on average.

# 3. Adding controls

**a)** Multiple regression:

```
m2 = lm(g4reading ~ small + race + yearssmall, data = star)
tidy(m2)
```

```
## # A tibble: 7 x 5
##    term          estimate std.error statistic  p.value
##    <chr>            <dbl>     <dbl>     <dbl>     <dbl>
## 1 (Intercept)      724.       1.40   517.      0
## 2 small             -4.00     4.98    -0.804  4.22e- 1
## 3 raceBlack        -33.8      2.95   -11.4     1.80e-29
## 4 raceAsian         14.8     19.3      0.767  4.43e- 1
## 5 raceHispanic       8.43    36.1      0.234  8.15e- 1
## 6 raceOther         80.3     36.1      2.23   2.62e- 2
## 7 yearssmall         2.17     1.29     1.68   9.33e- 2
```

**b)** The coefficient on `small` should remain similar to the bivariate model. This is expected because treatment was randomly assigned, so covariates should be balanced across groups. Stability of the coefficient confirms successful randomization.

**c)** The coefficient on `yearssmall` captures the cumulative benefit of spending additional years in a small class. A positive coefficient suggests each extra year is associated with higher scores.

# 4. Interactions

**a)** Interaction model:

```
m3 = lm(g4reading ~ small * race + yearssmall, data = star)
```

**b)** Print results:

```
tidy(m3)
```

```
## # A tibble: 11 x 5
##    term                 estimate std.error statistic  p.value
##    <chr>                   <dbl>     <dbl>     <dbl>     <dbl>
##  1 (Intercept)            725.       1.43   507.      0
##  2 small                   -5.32     5.12    -1.04    2.99e- 1
##  3 raceBlack              -36.0      3.59   -10.0     3.40e-23
##  4 raceAsian               21.3     20.9      1.02    3.07e- 1
##  5 raceHispanic             9.14    36.1      0.253   8.00e- 1
##  6 raceOther               53.3     51.0      1.05    2.96e- 1
##  7 yearssmall               2.25     1.29     1.74    8.25e- 2
##  8 small:raceBlack          6.97     6.33     1.10    2.71e- 1
##  9 small:raceAsian        -46.7     55.1     -0.847   3.97e- 1
## 10 small:raceHispanic      NA       NA       NA       NA
## 11 small:raceOther         54.3     72.2      0.753   4.52e- 1
```

**c)** Effect of small class for White students (reference category) = coefficient on `small`. For Black students = `small` + `small:raceBlack`:

```r
# White students
coef(m3)["small"]
```

```
##     small
## -5.317517
```

```r
# Black students
coef(m3)["small"] + coef(m3)["small:raceBlack"]
```

```
##     small
## 1.656508
```

**d)** The interaction terms test whether the benefit of small classes differs by race. Some studies find that minority students benefit more, though interactions may not always be statistically significant given sample sizes.
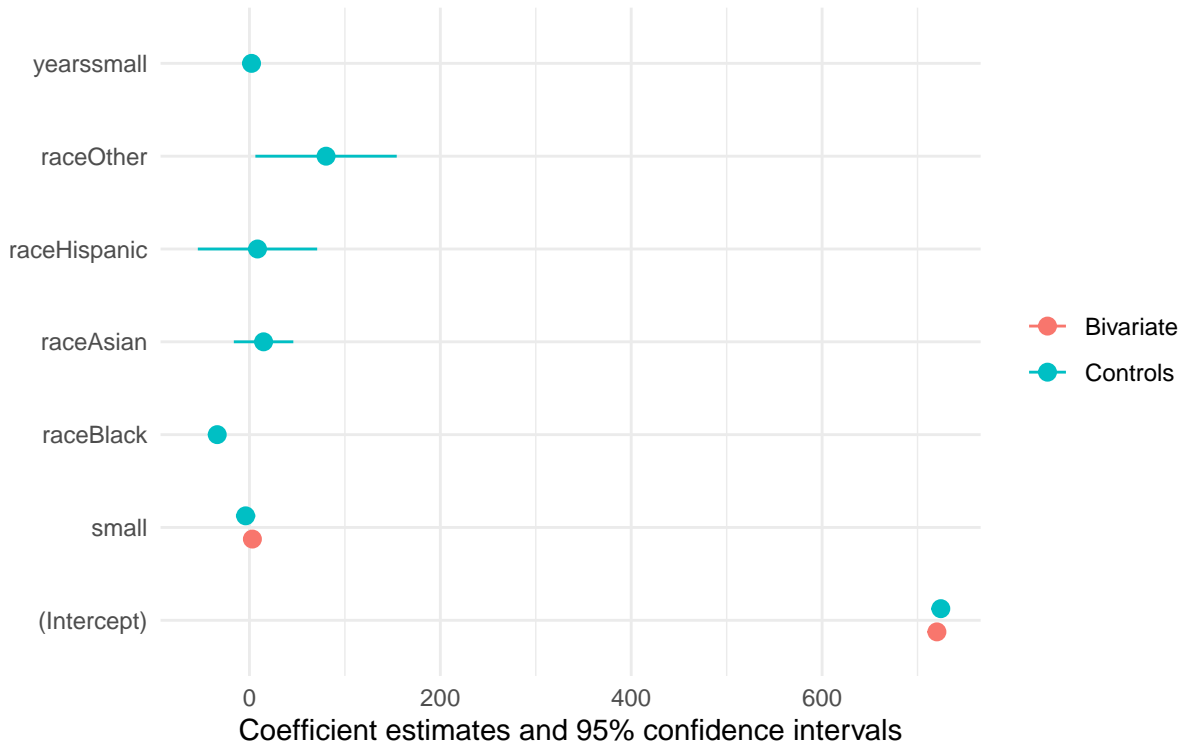
## 5. Presenting results

**a)** Comparison table:

```r
modelsummary(
  list("Bivariate" = m1, "Controls" = m2, "Interaction" = m3),
  vcov = "robust", output = "markdown")
```

|                    | Bivariate  | Controls   | Interaction |
|--------------------|------------|------------|-------------|
| (Intercept)        | 720.291    | 724.386    | 724.680     |
|                    | (1.309)    | (1.361)    |             |
| small              | 3.100      | -4.000     | -5.318      |
|                    | (2.319)    | (5.175)    |             |
| raceBlack          |            | -33.758    | -36.010     |
|                    |            | (3.077)    |             |
| raceAsian          |            | 14.803     | 21.320      |
|                    |            | (15.898)   |             |
| raceHispanic       |            | 8.433      | 9.140       |
|                    |            | (31.883)   |             |
| raceOther          |            | 80.274     | 53.320      |
|                    |            | (37.764)   |             |
| yearssmall         |            | 2.170      | 2.249       |
|                    |            | (1.359)    |             |
| small × raceBlack  |            |            | 6.974       |
| small × raceAsian  |            |            | -46.680     |
| small × raceOther  |            |            | 54.320      |
| Num.Obs.           | 2353       | 2353       | 2353        |
| R2                 | 0.001      | 0.057      | 0.058       |
| R2 Adj.            | 0.000      | 0.054      | 0.054       |
| AIC                | 25313.7    | 25188.4    | 25191.9     |
| BIC                | 25331.0    | 25234.5    | 25255.3     |
| Log.Lik.           | -12653.855 | -12586.194 | -12584.934  |
| F                  | 1.786      | 21.549     |             |
| RMSE               | 52.40      | 50.91      | 50.88       |

|  | Bivariate | Controls | Interaction |
|---|---|---|---|
| Std.Errors | Robust | Robust | Robust |

**b–c)** Coefficient plot:

```
modelplot(
  list("Bivariate" = m1, "Controls" = m2, "Interaction" = m3),
  vcov = "robust")
```



## 6. Brief discussion

**a)** The STAR data show a positive effect of small class sizes on student achievement: students randomly assigned to small classes score higher on both reading and math.

**b)** This evidence is more credible than typical observational studies because treatment was randomly assigned, eliminating confounders by design. In observational data, class size correlates with school resources, neighborhood characteristics, and student composition.

**c)** Limitations: substantial attrition (many students lack 4th grade scores, and attrition may differ by treatment group), potential compliance issues (some students may have switched classrooms), and the specific context (Tennessee, 1980s) may limit generalizability.