# Multilevel Regression And Poststratification (A Primer)

Joseph T. Ornstein*

October 26, 2021

## 1 Running Example

To demonstrate how MRP works, we'll consider an example where we know the "real" answer, and can explore how various refinements to the model improve predictive accuracy. The approach we'll use mirrors that in Buttice and Highton (2013), taking responses from a large scale US survey of voters (schaffner ref).[1]

```
library(tidyverse)
library(ggrepel)

load('data/CES-2020.RData')
```

The data is available here, and we'll be using a tidied up version of the dataset created by `R/cleanup-ces-2020.R`. This tidied version of the data only includes the 33 states with at least 500 respondents.
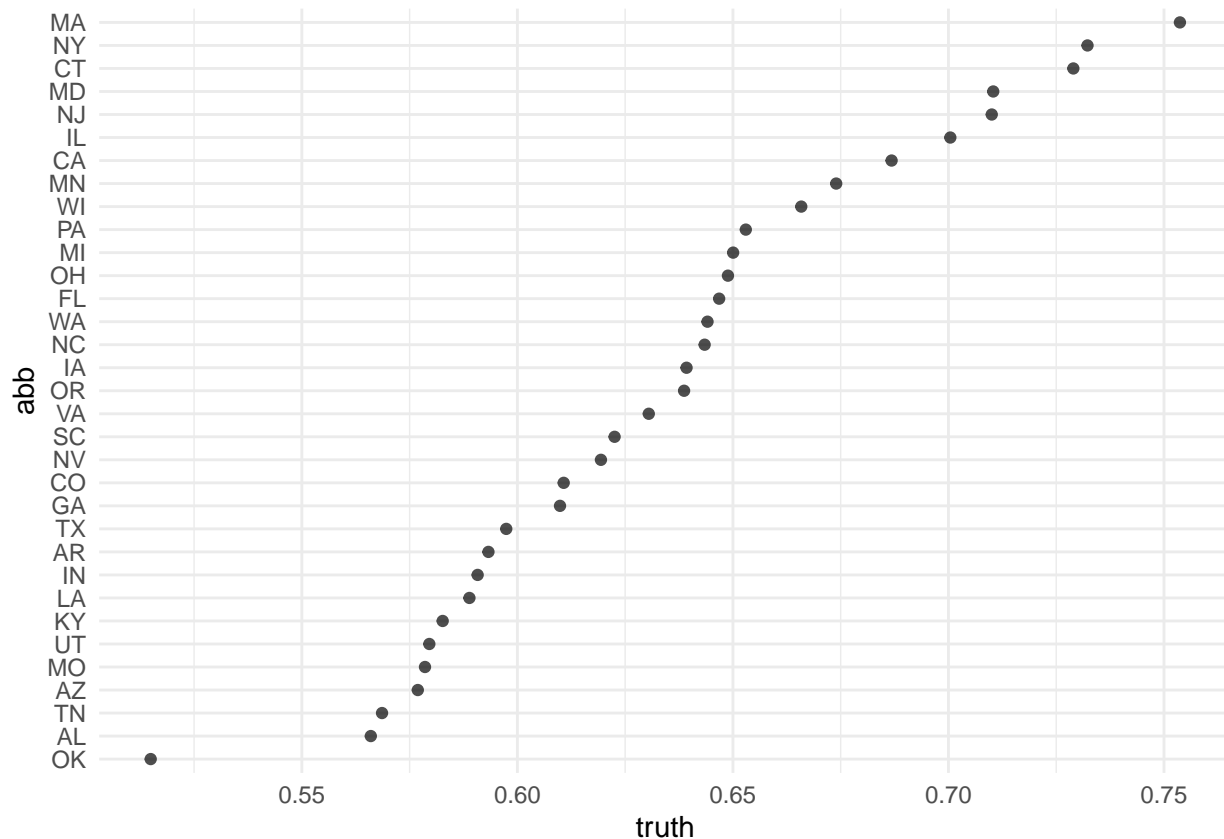
### 1.1 The Truth

```
truth <- ces %>%
  filter(!is.na(assault_rifle_ban)) %>%
  group_by(abb) %>%
  summarize(truth = sum(assault_rifle_ban == 'Support') / n())

# plot
truth %>%
  # reorder abb so the chart is organized by percent who support
  mutate(abb = fct_reorder(abb, truth)) %>%
  ggplot(mapping = aes(x=truth, y=abb)) +
  geom_point(alpha = 0.7) +
  theme_minimal()
```

---

*Department of Political Science, University of Georgia
[1]Throughout, I will use R functions from the "tidyverse" to make the code more human-readable.

Note what I mean by the "truth" here is the true percentage of CES respondents who supported the assault rifle ban. That's our target. This overstates the percent of the total population that support such a ban, since the CES sample is not a simple random sample.

## 1.2  Draw a Sample

Step 1: draw a sample.

```
sample_data <- ces %>%
  slice_sample(n = 500)
```

```
sample_summary <- sample_data %>%
  filter(!is.na(assault_rifle_ban)) %>%
  group_by(abb) %>%
  summarize(pct_support = sum(assault_rifle_ban == 'Support') / n(),
            num = n())
```

```
sample_summary
```

```
## # A tibble: 33 x 3
##     abb   pct_support   num
##    <chr>        <dbl> <int>
## 1 AL          0.778       9
## 2 AR          0.4         5
## 3 AZ          0.706      17
## 4 CA          0.725      40
## 5 CO          0.714       7
```

```
##  6 CT            0.5         2
##  7 FL            0.510      49
##  8 GA            0.667      18
##  9 IA            1           5
## 10 IL            0.667      21
## # ... with 23 more rows
```

For readers who are less familiar with American politics, rest assured that this is an unrepresentative draw from the state of Iowa. So simply disaggregating and taking sample means will not yield good estimates, as you can see by comparing the percent of respondents from the sample who supported the ban against the percent of CES respondents.

```
sample_summary %>%
  left_join(truth, by = 'abb') %>%
  ggplot(mapping = aes(x=pct_support,
                       y=truth,
                       label=abb)) +
  geom_point(alpha = 0.5) +
  geom_text_repel() +
  theme_minimal() +
  geom_abline(intercept = 0, slope = 1, linetype = 'dashed') +
  labs(x = 'Sample Average',
       y = 'Truth',
       title = 'Disaggregated Estimates')
```
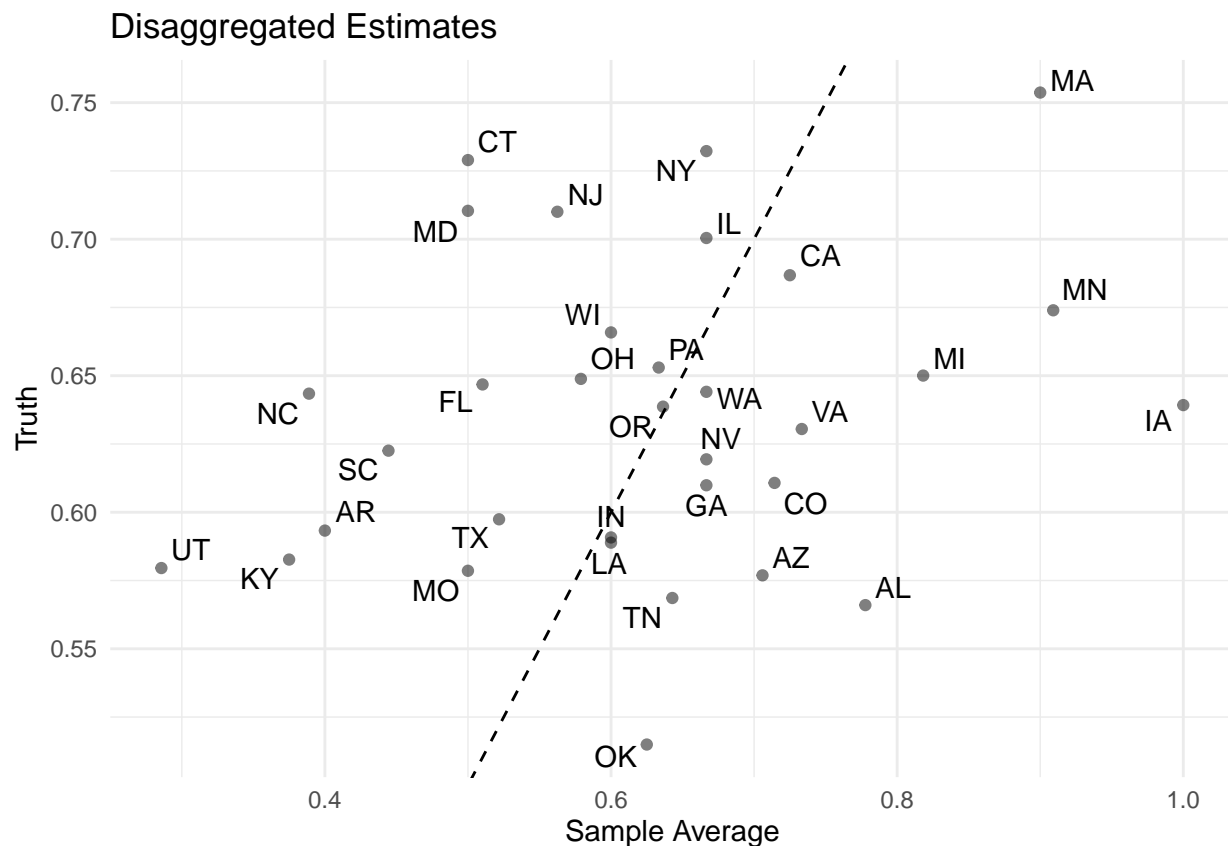


Figure 1: Esimates from disaggregated sample data

## 1.3 Multilevel Regression

```
# TODO: multilevel model; show how partial pooling fixes a bunch here.

# logistic model
model <- glm(as.numeric(assault_rifle_ban == 'Support') ~
               gender + educ + race + age,
           data = sample_data,
           family = 'binomial')

summary(model)
```

```
##
## Call:
## glm(formula = as.numeric(assault_rifle_ban == "Support") ~ gender +
##     educ + race + age, family = "binomial", data = sample_data)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -2.2034  -1.1676   0.6614   0.9510   1.6246
##
## Coefficients:
##                           Estimate Std. Error z value Pr(>|z|)
## (Intercept)               1.035670   0.897739   1.154   0.2486
## genderMale               -0.956595   0.204229  -4.684 2.81e-06 ***
## educHigh school graduate -1.187415   0.615309  -1.930   0.0536 .
## educSome college         -1.373742   0.616990  -2.227   0.0260 *
## educ2-year               -0.613440   0.639581  -0.959   0.3375
## educ4-year               -0.488137   0.621276  -0.786   0.4320
## educPost-grad            -0.395022   0.653821  -0.604   0.5457
## raceBlack                 0.936116   0.724075   1.293   0.1961
## raceHispanic              0.001496   0.703360   0.002   0.9983
## raceOther                -0.774865   0.778793  -0.995   0.3198
## raceWhite                 0.022585   0.648816   0.035   0.9722
## age                       0.013543   0.005617   2.411   0.0159 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 665.03  on 499  degrees of freedom
## Residual deviance: 607.47  on 488  degrees of freedom
## AIC: 631.47
##
## Number of Fisher Scoring iterations: 4
```

## 1.4 Poststratification

```
psframe <- ces %>%
  count(abb, gender, educ, race, age)

head(psframe)
```

```
## # A tibble: 6 x 6
```

```
##    abb   gender educ  race     age      n
##   <chr> <chr>  <fct> <chr> <dbl> <int>
## 1 AL      Female No HS Black    28     1
## 2 AL      Female No HS Black    29     1
## 3 AL      Female No HS Black    34     1
## 4 AL      Female No HS Black    54     1
## 5 AL      Female No HS Black    64     1
## 6 AL      Female No HS Other    36     1
```

Append predicted probabilities to the poststratification frame.

```
psframe <- psframe %>%
  mutate(predicted_probability = predict(model, psframe, type = 'response'))
```

Poststratified estimates are the population-weighted predictions

```
poststratified_estimates <- psframe %>%
  group_by(abb) %>%
  summarize(estimate = weighted.mean(predicted_probability, n))
```

Merge and compare:

```
d <- left_join(poststratified_estimates,
               truth,
               by = 'abb')

library(ggrepel)

ggplot(data = d,
       mapping = aes(x = estimate,
                     y = truth,
                     label = abb)) +
  geom_point(alpha = 0.7) +
  geom_text_repel() +
  theme_minimal() +
  geom_abline(intercept = 0, slope = 1, linetype = 'dashed') +
  coord_equal() +
  labs(caption = paste0('Correlation = ', round(cor(d$estimate, d$truth), 2)))
```

This highlights one of the dangers of producing MRP estimates from an underspecified model. When the first-stage model is underfit, poststratified estimates tend to collapse towards the global mean [CITATION NEEDED + BETTER TERMINOLOGY].

This compression means we should be wary of MRP studies that show policy outcomes "leapfrogging" estimated public opinion (Simonovits and Payson 2020). It could be that policymakers are more extreme than their constituents, or that MRP produces estimates of constituency preferences that are too moderate.

## 1.5   The Other Extreme: Overfitting

To illustrate the other extreme, let's estimate a model with a separate intercept term for each state – a "fixed effects" model. Because our sample contains several states with very few observations, these state-specific intercepts will likely overfit to sampling variability.

```
# fit the model
model2 <- glm(as.numeric(assault_rifle_ban == 'Support') ~
                gender + educ + race + age +
                  abb,
              data = sample_data,
```
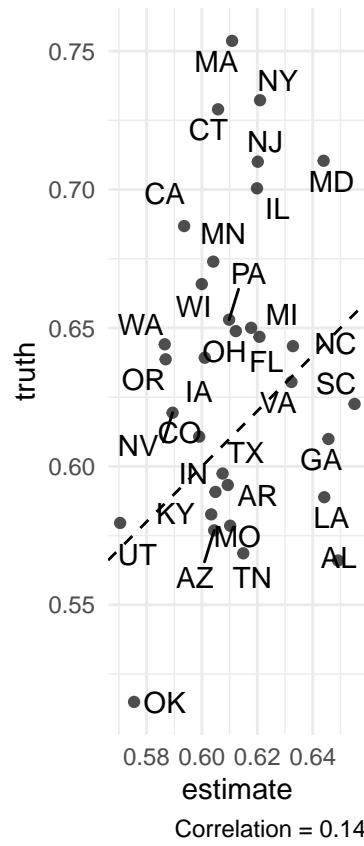
Figure 2: Underfit MRP estimates from complete pooling model

```r
                 family = 'binomial')

# make predictions
psframe <- psframe %>%
  mutate(predicted_probability = predict(model2, psframe, type = 'response'))

# poststratify
poststratified_estimates <- psframe %>%
  group_by(abb) %>%
  summarize(estimate = weighted.mean(predicted_probability, n))


d <- left_join(poststratified_estimates,
               truth,
               by = 'abb')

ggplot(data = d,
       mapping = aes(x = estimate,
                     y = truth,
                     label = abb)) +
  geom_point(alpha = 0.7) +
  geom_text_repel() +
  theme_minimal() +
  geom_abline(intercept = 0, slope = 1, linetype = 'dashed') +
  coord_equal() +
  labs(caption = paste0('Correlation = ', round(cor(d$estimate, d$truth), 2)))
```

Compare this to Figure 1 – these estimates are similarly overfit. Iowa is predicted to have roughly 100% support due to an idiosyncratic sample, while Maryland has the opposite problem.

## 1.6 The Sweet Spot: Partial Pooling

Gelman and Little (1997) solution: multilevel models that partially pool across regions. (Explanation of partial pooling goes...here)

```r
library(lme4)

model3 <- glmer(as.numeric(assault_rifle_ban == 'Support') ~  gender + educ + race + age +
                (1|abb),
                data = sample_data,
                family = 'binomial')

# make predictions
psframe <- psframe %>%
  mutate(predicted_probability = predict(model3, psframe, type = 'response'))

# poststratify
poststratified_estimates <- psframe %>%
  group_by(abb) %>%
  summarize(estimate = weighted.mean(predicted_probability, n))


d <- left_join(poststratified_estimates,
               truth,
```

Figure 3: Overfit MRP estimates from fixed effects model

```
                by = 'abb')

ggplot(data = d,
       mapping = aes(x = estimate,
                     y = truth,
                     label = abb)) +
  geom_point(alpha = 0.7) +
  geom_text_repel() +
  theme_minimal() +
  geom_abline(intercept = 0, slope = 1, linetype = 'dashed') +
  coord_equal() +
  labs(caption = paste0('Correlation = ', round(cor(d$estimate, d$truth), 2)))
```

Figure 4: MRP estimates from model with partial pooling

TODO: Well that's not the approach we should take, then. Partial pooling isn't magic. It just undoes the damage that fixed effects does. The magic is in good geographic predictors.

## 1.7 Stacking

```
sample_data <- mutate(sample_data,
                      Y = as.numeric(assault_rifle_ban == 'Support'))

library(SuperLearner)
```

9

```r
# fit Super Learner
SL.library <- c("SL.ranger", "SL.gam", "SL.xgboost", "SL.glm")

X <- sample_data %>%
  select(gender, educ, race, age, abb)

newX <- psframe %>%
  select(gender, educ, race, age, abb)

sl.out <- SuperLearner(Y = sample_data$Y,
                       X = X,
                       newX = newX,
                       family = binomial(),
                       SL.library = SL.library, verbose = FALSE)
```

```
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
```
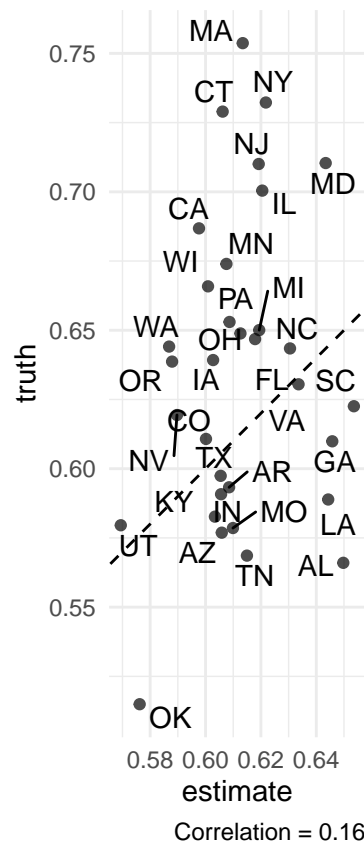
```
## Error in predict.xgb.Booster(model, newdata = newX) :
##   Feature names stored in `object` and `newdata` are different!
## Error in s(educ, 2) :
##   unordered factors cannot be used as smoothing variables
```

```r
#sl.out$SL.predict
sl.out$coef
```

```
##  SL.ranger_All      SL.gam_All SL.xgboost_All     SL.glm_All
##      0.4979383       0.0000000      0.0000000      0.5020617
```

```r
# make predictions
psframe <- psframe %>%
  mutate(predicted_probability = sl.out$SL.predict)

# poststratify
poststratified_estimates <- psframe %>%
  group_by(abb) %>%
  summarize(estimate = weighted.mean(predicted_probability, n))

d <- left_join(poststratified_estimates,
               truth,
               by = 'abb')

ggplot(data = d,
       mapping = aes(x = estimate,
                     y = truth,
                     label = abb)) +
  geom_point(alpha = 0.7) +
  geom_text_repel() +
  theme_minimal() +
  geom_abline(intercept = 0, slope = 1, linetype = 'dashed') +
  coord_equal() +
  labs(caption = paste0('Correlation = ', round(cor(d$estimate, d$truth), 2)))
```

And that's just "out-of-the-box!" What if we were more careful about it?

## 1.8  Synthetic Poststratification

Suppose that we did not have access to the entire joint distribution of individual-level covariates. Leemann and Wasserfallen (2017) suggest an extension of Mr. P, which they (delightfully) dub Multilevel Regression and Synthetic Poststratification (Mrs. P). Lacking the full joint distribution of covariates for poststratification, one can instead create a *synthetic* poststratification frame by assuming that additional covariates are statistically independent of one another. So long as your first stage model is linear-additive, this approach yields the same predictions as if you knew the true joint distribution![2] And even if the first-stage model is not linear-additive, simulations suggest that the improved performance from additional predictors tends to overcome the error introduced by synthetic poststratification.

To create a synthetic poststratification frame, convert frequencies to probabilities and multiply. For example, suppose we only had the joint distribution for gender, race, and age, and wanted to create a synthetic poststratification including education.

```r
# poststratification frame with 3 variables
psframe3 <- ces %>%
```

---
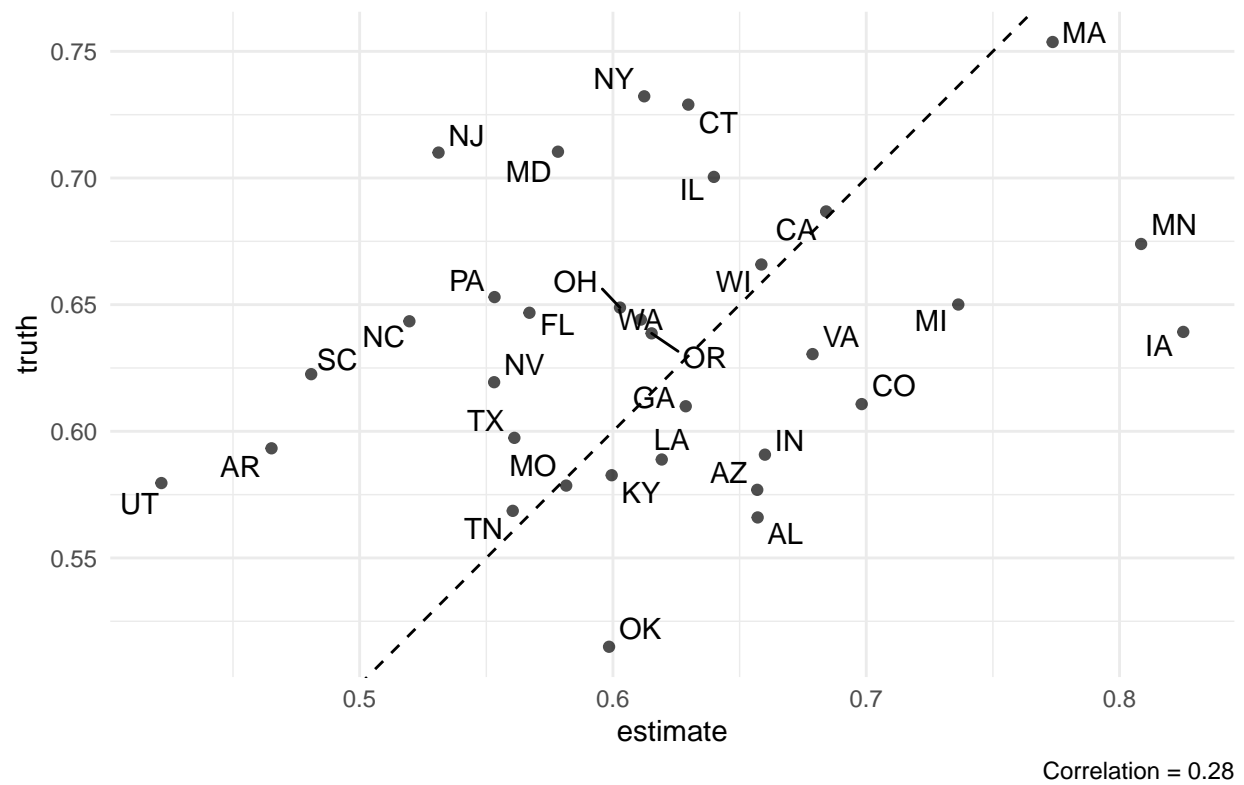
[2]See Ornstein (2020) appendix A for mathematical proof.

Figure 5: Estimates from an ensemble first-stage model

Correlation = 0.28

```
  count(abb, gender, race, age) %>%
  group_by(abb) %>%
  mutate(prob1 = n / sum(n))

head(psframe3)
```

```
## # A tibble: 6 x 6
## # Groups:   abb [1]
##   abb   gender race   age     n  prob1
##   <chr> <chr> <chr> <dbl> <int>  <dbl>
## 1 AL    Female Asian    24     1 0.00106
## 2 AL    Female Asian    27     1 0.00106
## 3 AL    Female Asian    29     1 0.00106
## 4 AL    Female Asian    30     1 0.00106
## 5 AL    Female Asian    31     1 0.00106
## 6 AL    Female Asian    34     2 0.00211
```

```
# distribution of education variable by state
psframe_educ <- ces %>%
  count(abb, educ) %>%
  group_by(abb) %>%
  mutate(prob2 = n / sum(n))

head(psframe_educ)
```

```
## # A tibble: 6 x 4
## # Groups:   abb [1]
##   abb   educ                  n  prob2
##   <chr> <fct>             <int>  <dbl>
## 1 AL    No HS                49 0.0517
## 2 AL    High school graduate 287 0.303
## 3 AL    Some college        189 0.200
## 4 AL    2-year              122 0.129
## 5 AL    4-year              180 0.190
## 6 AL    Post-grad           120 0.127
```

```
synthetic_psframe <- left_join(psframe3, psframe_educ,
                               by = 'abb') %>%
  mutate(prob = prob1 * prob2)

head(synthetic_psframe)
```

```
## # A tibble: 6 x 10
## # Groups:   abb [1]
##   abb   gender race   age   n.x  prob1 educ              n.y prob2    prob
##   <chr> <chr> <chr> <dbl> <int>  <dbl> <fct>           <int> <dbl>   <dbl>
## 1 AL    Female Asian    24     1 0.00106 No HS             49 0.0517 5.46e-5
## 2 AL    Female Asian    24     1 0.00106 High school grad~ 287 0.303 3.20e-4
## 3 AL    Female Asian    24     1 0.00106 Some college      189 0.200 2.11e-4
## 4 AL    Female Asian    24     1 0.00106 2-year            122 0.129 1.36e-4
## 5 AL    Female Asian    24     1 0.00106 4-year            180 0.190 2.01e-4
## 6 AL    Female Asian    24     1 0.00106 Post-grad         120 0.127 1.34e-4
```

The SRP package contains a convenience function for this operation (see the vignette for more information).
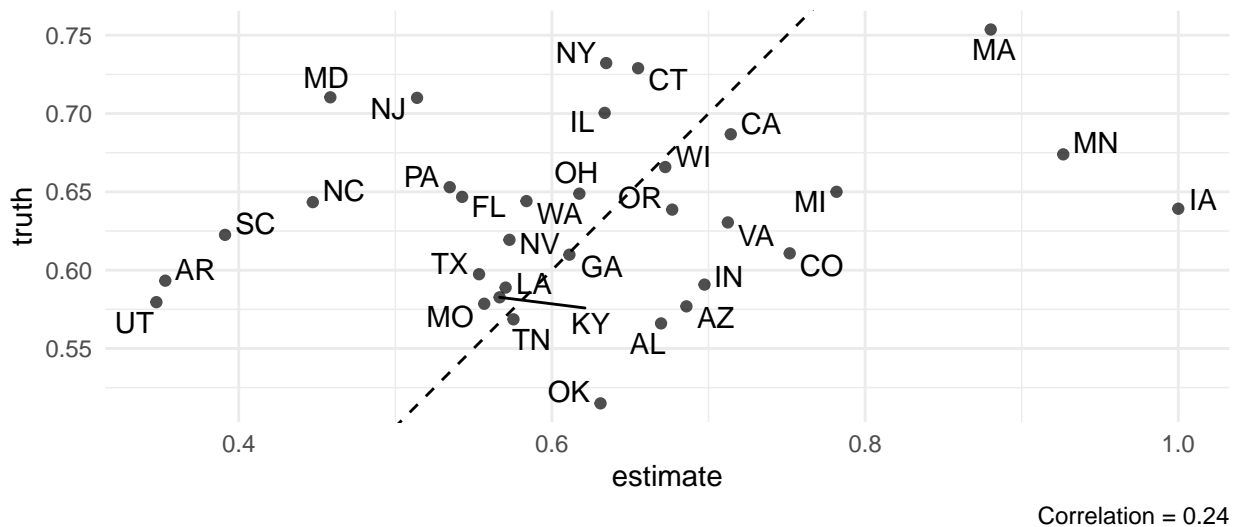
Then poststratify as normal.

TODO: Functions instead of comments.

```r
# make predictions
synthetic_psframe$predicted_probability <-  predict(model2, synthetic_psframe, type = 'response')

# poststratify
poststratified_estimates <- synthetic_psframe %>%
  group_by(abb) %>%
  summarize(estimate = weighted.mean(predicted_probability, prob))


d <- left_join(poststratified_estimates,
               truth,
               by = 'abb')

ggplot(data = d,
       mapping = aes(x = estimate,
                     y = truth,
                     label = abb)) +
  geom_point(alpha = 0.7) +
  geom_text_repel() +
  theme_minimal() +
  geom_abline(intercept = 0, slope = 1, linetype = 'dashed') +
  coord_equal() +
  labs(caption = paste0('Correlation = ', round(cor(d$estimate, d$truth), 2)))
```



Correlation = 0.24

Note that the performance is slightly worse than when we knew the true joint distribution. But is it worse

14

than omitting education entirely?

# References

Buttice, Matthew K., and Benjamin Highton. 2013. "How Does Multilevel Regression and Poststratification Perform with Conventional National Surveys?" *Political Analysis* 21 (4): 449–67. https://doi.org/10.1093/pan/mpt017.

Gelman, Andrew, and Thomas C Little. 1997. "Poststratification into Many Categories Using Hierachical Logistic Regression." *Survey Methodology* 23 (2): 127–35.

Leemann, Lucas, and Fabio Wasserfallen. 2017. "Extending the Use and Prediction Precision of Subnational Public Opinion Estimation." *American Journal of Political Science* 61 (4): 1003–22.

Ornstein, Joseph T. 2020. "Stacked Regression and Poststratification." *Political Analysis* 28 (2): 293–301. https://doi.org/10.1017/pan.2019.43.

Simonovits, Gabor, and Julia Payson. 2020. "Locally Controlled Minimum Wages Are No Closer to Public Preferences." *Working Paper*, 21.