

Week 5: Surveys and Weighting

PLSC 40502 - Statistical Models

Review

Previously

- Stan
 - **Hamiltonian Monte Carlo** - use the gradient of the likelihood \times prior to get good M-H proposals
 - Avoids proposals that are either too auto-correlated and too likely to be rejected
 - Common structure to statistical modeling: **data, parameters, model**
 - Frequently used models have canned routines in `rstanarm`
- Diagnosing model fit
 - Posterior predictive checks
 - **Information criteria** vs. **Cross-validation**
 - log-predictive density as a measure of prediction quality (how much probability mass do we place on the "correct" outcome)

This week

- **Survey weighting and post-stratification**
 - The problem of non-response
 - Using auxiliary variables to construct weights
 - Post-stratification vs. raking
- **Calibration weighting**
 - Generalizing post-stratification-style adjustments to allow for a variety of balance conditions
- **Combining multilevel regression and post-stratification**
 - Using multi-level models and population-level information to estimate quantities for **small areas**

Survey sampling

Surveys

- A common task in survey research is estimating an unknown population parameter θ from a sample of respondents $i \in \{1, 2, 3, \dots, N\}$
 - What share of voters in Michigan plan to vote for Joe Biden in the 2020 election?
 - What is the share of adults who are vaccinated in each U.S. county?
 - What share of residents in Nevada plan to vote in the 2024 general election?
- Historically, two approaches to survey sampling design
 - **Quota sampling** - Define a set of known demographic targets and recruit respondents to match.
 - **Probability sampling** - Select respondents from a **sampling frame** at random with a known probability.
- Dominance of **probability sampling** in the late 20th century
 - Random Digit Dialing allowed for (near)-simple random samples from the U.S. adult population
 - High response rates

Decline of pure probability samples

- Two big factors have lead to the decline of exclusively probability-based sampling approaches
 - Decline in population coverage - fewer individuals have landlines!
 - Extremely high non-response rates.
- Non-random non-response can bias our estimates
 - Non-responders have different characteristics to the responders that may be correlated with the target quantity of interest.
 - Huge concern in recent efforts to poll elections (e.g. "shy tories")
- Modern polling
 - Combine probability and quota approaches
 - Weighting ex-post to match population targets.

Survey inference

- Our goal is to estimate some parameter θ_y related to a population outcome variable y (e.g. a proportion, mean, median, etc...)
- We construct an estimator $\hat{\theta}_y$
 - For population means, we'll use a sample mean, etc...
- In addition to the outcome of interest y , we observe **auxiliary** variables \mathbf{x}
 - Units are sampled from some unknown target population $f_U(y, \mathbf{x})$
 - The in-sample distribution of the **responders** is denoted $g_R(y, \mathbf{x})$
- Our auxiliary information involves some **known features** of the target distribution: $\check{I}_{\mathbf{x}}$
 - We'll use this to construct **targets** for the population distribution $\tilde{T}_{\mathbf{x}} = \{\tilde{T}_{\mathbf{x}1}, \dots, \tilde{T}_{\mathbf{x}M}\}$
 - For example, suppose we know the full population distribution of age, gender, education, income, and party ID
- In some settings, the auxiliary distribution maps easily to the target - but in other settings we have to use the auxiliary information to **estimate** our target
 - (e.g) we might have features of our target distribution but we don't know who will turn out to **vote** before the election
 - "Likely voter models" are a **target estimation** problem

Inference with known sampling weights

- If we know the probability that a unit is selected into the sample from the sample frame π_i , it is straightforward to construct an estimator $\hat{\theta}_y$ of the population mean θ_y .
- The **Horvitz-Thompson** estimator weights each unit by $d_i = \frac{1}{\pi_i}$, the inverse probability of being selected into the sample

$$\hat{\theta}_y^{\text{HT}} = \frac{\sum_{i=1}^N d_i Y_i}{\mathbb{E}[\sum_{i=1}^n d_i]}$$

- When sampling probabilities are equivalent, this reduces to the sample mean.
- More commonly, rather than using the expectation of the weights in the denominator, we'll use the actual observed sum of the weights, giving the Hajek estimator

$$\hat{\theta}_y^{\text{H}} = \frac{\sum_{i=1}^N d_i Y_i}{\sum_{i=1}^n d_i}$$

Unknown sampling weights

- When d_i is not known, we will need to estimate **adjustment weights** using a combination of modeling assumptions and auxiliary data
 - Even when d_i is known, if non-response is high, we still don't know the probability of selection into the **observed** data ρ_i
- With adjustment weights \tilde{w}_i , our Hajek estimator becomes

$$\hat{\theta}_y^w = \frac{\sum_{i=1}^N \tilde{w}_i Y_i}{\sum_{i=1}^n \tilde{w}_i}$$

- Now the weights are not necessarily known but must be obtained from our population targets \tilde{T}_x

Post-stratification

- The easiest approach to adjusting a non-representative survey is to weight to match the **known joint** distribution of x in the population $f_U(\mathbf{x})$
 - This requires a lot of auxiliary information about the target $f_U(\mathbf{x})$ - typically obtained from high-quality census data
 - (e.g.) U.S. Census Public-Use Microdata: What is the share of Black, college educated, 30-45 year olds in Massachusetts?
- In post-stratification, we divide our sample up into C **cells** c that are mutually exclusive and exhaustive.
 - $\tilde{T}_{\mathbf{x}} = \{\tilde{P}_1, \tilde{P}_2, \dots, \tilde{P}_C\}$ is our population distribution of these cells
- Our **sampling/response model** assumes that **within** each of these cells we have a simple-random sample from that particular stratum of the population
 - There may be variation in non-response or over/under-sampling, but only **across** cells
 - "post"-stratification because the intuition is akin to a design where we actually *did* stratify ex-ante
- Our **measurement model** assumes we observe the joint distribution of auxiliary variables \mathbf{x} in the target population
 - Difficult in many cases!

Post-stratification

- Two ways to think of post-stratification:
- **First** - Let $\hat{\theta}_y^c$ denote our estimator for the population mean **within** cell c (possibly using design weights d_i)
 - Then, our post-stratification estimator $\hat{\theta}_y^{PS}$ is:

$$\hat{\theta}_y^{PS} = \sum_{c=1}^C \tilde{P}_c \hat{\theta}_y^c$$

- Alternatively, we'll sometimes write $\tilde{P}_c = \frac{\tilde{N}_c}{\tilde{N}}$ where \tilde{N} is the size of the population and \tilde{N}_c is the number of units in that cell.

$$\hat{\theta}_y^{PS} = \sum_{c=1}^C \frac{\tilde{N}_c}{\tilde{N}} \hat{\theta}_y^c$$

- Fixed "constant" weights on the within-cell estimators.

Post-stratification

- **Second** - We can think of it as weighting individual observations using our adjustment weights \tilde{w}_i^{PS}
- Let $c(i)$ denote the class to which unit i belongs. Then the post-stratification weights are

$$\tilde{w}_i^{PS} = (\tilde{P}_{c(i)} / \hat{P}_{c(i)}^S) \times d_i$$

where $\hat{P}_{c(i)}^S$ is the estimated **in-sample** proportion of observations in class c

- When design weights are constant, $\hat{P}_{c(i)}^S$ is just the sample mean of the indicator of class membership $\mathbf{1}_{i \in c}$
 - More generally, you can write it as $\hat{P}_{c(i)}^S = (\sum_{i=1}^N d_i \mathbf{1}_{i \in c}) / (\sum_{i=1}^N d_i)$
- **Intuition**
 - The weights **up-weight** observations that are under-represented in the sample relative to the target population
 - The weights **down-weight** observations that are over-represented in the sample relative to the target population.

Raking

- Often post-stratification with many covariates is challenging!
 - The number of cells grows rapidly as we add more covariates
 - e.g. gender x party x state = $2 \times 3 \times 50 = 300$ cells!
- Sometimes our population data only give us the marginal distributions but not the joint distributions
- **Raking** weights are designed to match the marginal distribution of the auxiliary covariate **in-sample**
 - No closed-form expression, but iterative algorithms exist to compute raking weights.
- **Intuition**
 - Raking works well when things are additive
 - Doesn't work as well when things are **interactive**

Example: CCES 2020

- Let's dive in to the 2020 CCES.

```
library(survey)
cces <- read_csv("data/CCES_subset.csv") %>% filter(!is.na(trumpApprove))
cces$educ_bin[cces$educ_bin == "College degree"] <- "College graduate" # Fix a naming inconsistency
```

- We'll be using a subset of the outcome data, looking at the share of respondents who state that they strongly or somewhat approve of then-President Donald Trump.
- Start by making a **svydesign** object - we'll pretend that the sampling weights don't exist for now

```
cces_surv_unwt <- svydesign(~1, weights = ~1, data=cces)
```

- In-sample, what's the proportion of respondents who approve of Trump?

```
svymean(~trumpApprove, design=cces_surv_unwt)
```

```
##           mean SE
## trumpApprove 0.383  0
```

Example: CCES 2020

- How does this compare to the properly weighted mean?

```
cces_surv_wt <- svydesign(~1, weights = ~commonweight, data=cces)
svymean(~trumpApprove, design=cces_surv_wt)
```

```
##              mean SE
## trumpApprove 0.444  0
```

- Trump's approval is a few points higher after the weighting adjustment.
- From the CCES Guide:

the completed cases were weighted to the sampling frame using entropy balancing. The 2019 ACS was used as the frame for weighting the common content and the team samples. The CES sample was weighted to match the distributions of the 2019 ACS on gender, age, race, Hispanic origin, and education level. The moment conditions included age, gender, education, race, plus their interactions. The resultant weights were then post-stratified by age, gender, education, race, "born again" status, voter registration status, 2016 Presidential vote choice, and 2020 Presidential vote choice as needed.

Population targets

- We'll be using the 2020 ACS 5-year Public Use Microdata Sample
 - Obtain the complete **joint** distribution of age, gender and education in the U.S.

```
acs_targets <- read_csv("data/ACS_2020_microdata_3cat.csv")
```

- What's the **marginal** distribution of education in the target population?

```
acs_targets %>% group_by(educ_bin) %>% summarize(n = sum(Count)) %>% ungroup() %>% mutate(prop
```

```
## # A tibble: 4 × 3
##   educ_bin      n  prop
##   <chr>      <dbl> <dbl>
## 1 College graduate 48194135 0.190
## 2 H.S. or less    99042042 0.391
## 3 Postgraduate   28425804 0.112
## 4 Some college   77634550 0.306
```

Population targets

- How does it compare to the **unweighted** marginal distribution in the sample?

```
cces %>% group_by(educ_bin) %>% summarize(n = n()) %>% ungroup() %>% mutate(prop = n/sum(n))
```

```
## # A tibble: 4 × 3
##   educ_bin      n  prop
##   <chr>    <int> <dbl>
## 1 College graduate 14146 0.232
## 2 H.S. or less    18592 0.305
## 3 Postgraduate    8373 0.137
## 4 Some college   19857 0.326
```

Population targets

- What about the **joint** distribution of gender and education?
- In the **population**

```
acs_targets %>% group_by(educ_bin, gender_bin) %>% summarize(n = sum(Count)) %>% ungroup() %>%
```

```
## # A tibble: 8 × 4
##   educ_bin      gender_bin      n    prop
##   <chr>         <chr>    <dbl> <dbl>
## 1 College graduate Female  25521371 0.101
## 2 College graduate Male    22672764 0.0895
## 3 H.S. or less   Female  48061006 0.190
## 4 H.S. or less   Male    50981036 0.201
## 5 Postgraduate   Female  15043514 0.0594
## 6 Postgraduate   Male    13382290 0.0528
## 7 Some college   Female  41305976 0.163
## 8 Some college   Male    36328574 0.143
```

Population targets

- In the **sample**

```
cces %>% group_by(educ_bin, gender_bin) %>% summarize(n = n()) %>% ungroup() %>% mutate(prop =
```

```
## # A tibble: 8 × 4
##   educ_bin      gender_bin      n    prop
##   <chr>         <chr>    <int> <dbl>
## 1 College graduate Female    7478 0.123
## 2 College graduate Male      6668 0.109
## 3 H.S. or less   Female   12021 0.197
## 4 H.S. or less   Male     6571 0.108
## 5 Postgraduate   Female    4153 0.0681
## 6 Postgraduate   Male     4220 0.0692
## 7 Some college   Female   11536 0.189
## 8 Some college   Male     8321 0.136
```

Post-stratification

- We can construct the post-stratification weights manually
 - Start by calculating the population proportions in each bin

```
acs_targets <- acs_targets %>% mutate(strata = str_c(gender_bin, educ_bin, age_bin, sep="-"),  
                                     pop_proportion = Count/sum(Count))
```

- Do the same for the sample

```
cces <- cces %>% mutate(strata = str_c(gender_bin, educ_bin, age_bin, sep="-"))  
cces_strat <- cces %>% group_by(strata) %>% summarize(n=n()) %>% ungroup() %>% mutate(samp_prop =
```

- Join the datasets

```
cces <- cces %>% left_join(acs_targets%>% select(strata, pop_proportion), by="strata")  
cces <- cces %>% left_join(cces_strat %>% select(strata, samp_proportion), by="strata")
```

Post-stratification

- Construct the post-stratification weights

```
cces <- cces %>% mutate(postStratWt = pop_proportion/samp_proportion)
```

- Take the weighted average to estimate Trump Approval

```
weighted.mean(cces$trumpApprove, cces$postStratWt)
```

```
## [1] 0.397
```

Population targets

- Did the weights equalize the distributions? Let's look again at the joint distribution of gender and education
- In the **population**

```
acs_targets %>% group_by(educ_bin, gender_bin) %>% summarize(n = sum(Count)) %>% ungroup() %>%
```

```
## # A tibble: 8 × 4
##   educ_bin      gender_bin      n    prop
##   <chr>        <chr>    <dbl> <dbl>
## 1 College graduate Female  25521371 0.101
## 2 College graduate Male    22672764 0.0895
## 3 H.S. or less Female  48061006 0.190
## 4 H.S. or less Male    50981036 0.201
## 5 Postgraduate Female  15043514 0.0594
## 6 Postgraduate Male    13382290 0.0528
## 7 Some college Female  41305976 0.163
## 8 Some college Male    36328574 0.143
```

Population targets

- In the re-weighted **sample**

```
cces %>% group_by(educ_bin, gender_bin) %>% summarize(n = sum(postStratWt)) %>% ungroup() %>% n
```

```
## # A tibble: 8 × 4
##   educ_bin      gender_bin      n    prop
##   <chr>        <chr>    <dbl> <dbl>
## 1 College graduate Female   6143. 0.101
## 2 College graduate Male     5457. 0.0895
## 3 H.S. or less  Female  11568. 0.190
## 4 H.S. or less  Male   12271. 0.201
## 5 Postgraduate  Female   3621. 0.0594
## 6 Postgraduate  Male    3221. 0.0528
## 7 Some college  Female   9942. 0.163
## 8 Some college  Male    8744. 0.143
```


Post-stratification.

- We can also use the **survey** package - create the post-stratification weights using **postStratify** in **survey**

```
cces_postStrat <- postStratify(cces_surv_unwt, strata=~gender_bin + age_bin + educ_bin,  
                             population = acs_targets %>% select(gender_bin, age_bin, educ_bi
```

- Then use **svymean**

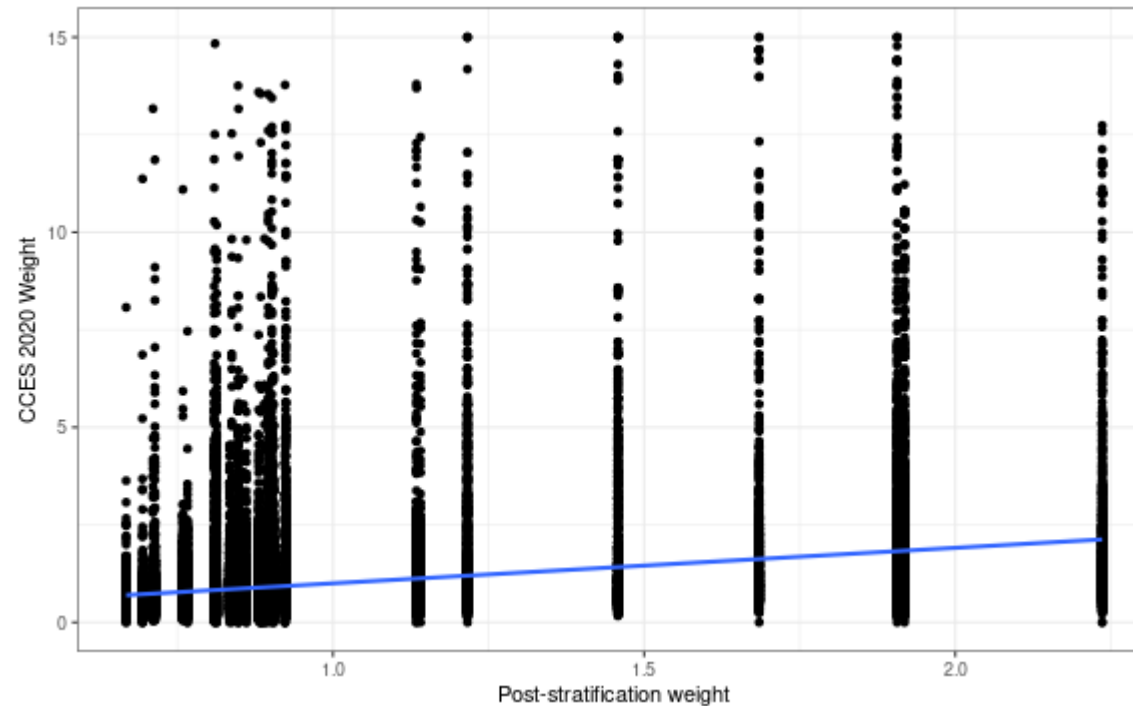
```
svymean(~trumpApprove, cces_postStrat)
```

```
##              mean SE  
## trumpApprove 0.397  0
```

Post-stratification

- Plotting our post-stratification weights against the actual CCES weights

```
cces %>% ggplot(aes(x=postStratWt, y=commonweight)) + geom_point() + geom_smooth(method="lm") +  
  xlab("Post-stratification weight") + ylab("CCES 2020 Weight")
```



Post-stratification

- Note that we can combine survey design weights with additional post-stratification weights.
 - For example, if we want to take a national survey and re-weight it to different demographic targets
- What happens if we combine our post-stratification weights with the CCES design weights

```
cces_postStrat_wt <- postStratify(cces_surv_wt, strata=~gender_bin + age_bin + educ_bin,  
                                population = acs_targets %>% select(gender_bin, age_bin, educ_bin))
```

- Our estimated Trump approval is closer to that original 44 percent (as expected, since our post-stratification variables are a subset of all the covariates that go into the CCES weights)

```
svymean(~trumpApprove, cces_postStrat_wt)
```

```
##               mean SE  
## trumpApprove 0.443  0
```

Raking

- Sometimes we only know the marginals and not the joint distributions. We can still construct weights that get us balance on the *marginals*.
- Start by generating the marginal counts (in the real world, these would be **all** that we have).

```
acs_marginals <- list()
acs_marginals[["gender_bin"]] <- acs_targets %>% group_by(gender_bin) %>% summarize(Freq = sum(Count))
acs_marginals[["age_bin"]] <- acs_targets %>% group_by(age_bin) %>% summarize(Freq = sum(Count))
acs_marginals[["educ_bin"]] <- acs_targets %>% group_by(educ_bin) %>% summarize(Freq = sum(Count))
```

- Make the raking design

```
cces_rake_unwt <- rake(cces_surv_unwt, sample.margins = list(~gender_bin, ~age_bin, ~educ_bin),
                      population.margins=acs_marginals)
cces$rakeWt <- (weights(cces_rake_unwt)/sum(weights(cces_rake_unwt)))*nrow(cces)
```

- In this case, we actually do about as well as when we have the full joint distribution!

```
svymean(~trumpApprove, cces_rake_unwt)
```

```
##               mean SE
## trumpApprove 0.399  0
```

Raking

- Note that raking will **not** guarantee balance on the full joint distribution
- In the **population**

```
acs_targets %>% group_by(educ_bin, gender_bin) %>% summarize(n = sum(Count)) %>% ungroup() %>%
```

```
## # A tibble: 8 × 4
##   educ_bin      gender_bin      n    prop
##   <chr>        <chr>    <dbl> <dbl>
## 1 College graduate Female  25521371 0.101
## 2 College graduate Male    22672764 0.0895
## 3 H.S. or less   Female  48061006 0.190
## 4 H.S. or less   Male    50981036 0.201
## 5 Postgraduate   Female  15043514 0.0594
## 6 Postgraduate   Male    13382290 0.0528
## 7 Some college   Female  41305976 0.163
## 8 Some college   Male    36328574 0.143
```

Raking

- In the re-weighted **sample**

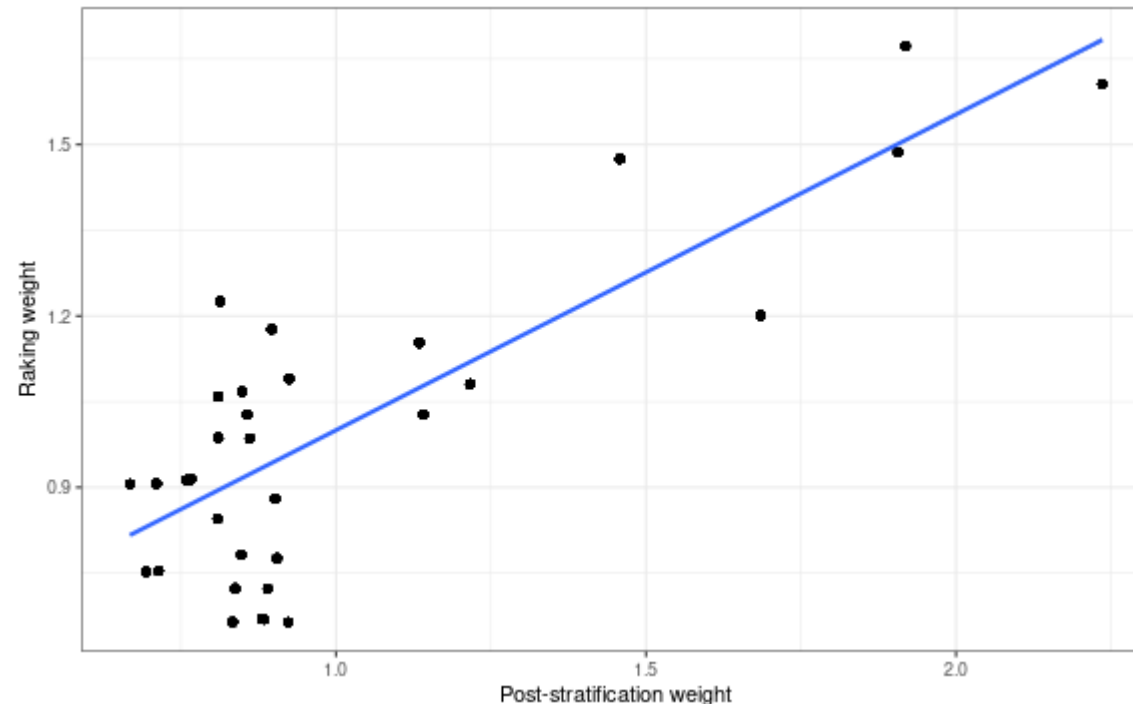
```
cces %>% group_by(educ_bin, gender_bin) %>% summarize(n = sum(rakeWt)) %>% ungroup() %>% mutate
```

```
## # A tibble: 8 × 4
##   educ_bin      gender_bin      n    prop
##   <chr>         <chr>    <dbl> <dbl>
## 1 College graduate Female   5268. 0.0864
## 2 College graduate Male     6332. 0.104
## 3 H.S. or less   Female  13650. 0.224
## 4 H.S. or less   Male   10189. 0.167
## 5 Postgraduate   Female   2890. 0.0474
## 6 Postgraduate   Male    3952. 0.0648
## 7 Some college   Female   9466. 0.155
## 8 Some college   Male    9220. 0.151
```

Raking vs. Post-stratification

- How do the raking and post-stratification weights compare?

```
cces %>% ggplot(aes(x=postStratWt, y=rakeWt)) + geom_point() + geom_smooth(method="lm") + theme_minimal() +  
  xlab("Post-stratification weight") + ylab("Raking weight")
```



Calibration Weighting

Calibration

- The choice of weights in a survey can be framed more generally as a **calibration** problem
- Find the set of weights \tilde{w}_i that minimize some distance measure

$$L(\tilde{w}, d) = \sum_{i=1}^N D(\tilde{w}_i, d_i)$$

subject to K constraints

$$\tilde{T}_{\mathbf{x}k} = \sum_{i=1}^N \tilde{w}_i z_{ik}$$

and

$$\sum_{i=1}^N \tilde{w}_i = 1$$

$$\tilde{w}_i \geq 0, \quad i = 1, \dots, n$$

Calibration

- Our constraints are defined in terms of population targets $\tilde{T}_{\mathbf{x}k}$ and sample quantities z_{ik}
 - z_{ik} is some transformation of the auxiliary variable x_{ik}
 - Typically just the identity, but we could consider higher-order moments with continuous covariates.
- **Intuition:** Generalizing raking to also include continuous variables
 - Match **sample moments** to **population moments**

Calibration

- What is this distance metric and how do we choose it?
- Lots of different ways to weight the sample to match a set of moment constraints.
 - But some weights are better than others!
- One popular distance metric is the **entropy loss**:

$$L(\tilde{w}, d) = \sum_{i=1}^N \tilde{w}_i \log \left(\frac{\tilde{w}_i}{d_i} \right)$$

- Equivalent to the Kullback-Leibler divergence between the distribution of the new weights and the original design weights (often uniform)
 - Without any constraints, this is minimized at the design weights
- **Intuition**: If the moment conditions are met, we prefer weights that are closer to uniform (or the original design weights) to minimize the loss of power to weighting
- Alternative metrics: squared distance

$$L(\tilde{w}, d) = \frac{1}{2} \sum_{i=1}^N (\tilde{w}_i - d_i)^2$$

Effective sample size

- Kish (1965) proposed a formula for the "effective sample size" of a weighted mean:

$$ESS = \frac{(\sum_{i=1}^N \tilde{w}_i)^2}{\sum_{i=1}^n \tilde{w}_i^2}$$

- Intuitively, this is maximized at uniform weights and becomes smaller as the weights become more imbalanced
 - Extremely large weights on a single observation lead to a **huge** loss of efficiency
 - Typically we see weights "winsorized" or cut-off to avoid these losses

Example: Calibrating with Party ID

- We have data on party ID in the population as well.
 - From **Pew Research** we have **gender** x **party** for registered voters in 2020
- But we don't have the full **party** x **age** x **gender** x **education** interaction (at least without additional research)
- So how do we weight to two partially overlapping datasets?
 1. *Census*: **Gender** x **Age** x **Education**
 2. *Pew*: **Party** x **Gender**
- Let's use **entropy balancing** to construct weights that match both sets of joint distributions

Example: Calibrating with Party ID

- From our Pew Party ID data

```
pew_party <- read_csv("data/pew_2020_party_gender.csv") %>% mutate(share = share/100) %>% rename(pew_party)
```

```
## # A tibble: 8 × 3
##   gender_bin party_bin   share
##   <chr>      <chr>     <dbl>
## 1 Male      Republican 0.31
## 2 Male      Democratic 0.26
## 3 Male      Independent 0.39
## 4 Male      Other      0.04
## 5 Female    Republican 0.28
## 6 Female    Democratic 0.39
## 7 Female    Independent 0.3
## 8 Female    Other      0.03
```

- Let's convert this into the **joint** distribution using the known age distribution

```
pew_party$Freq <- pew_party$share*(acs_marginals[["gender_bin"]]$Freq/sum(acs_marginals[["gender_bin"]]))
pew_party$Freq[pew_party$gender_bin == "Male"] <- pew_party$share[pew_party$gender_bin == "Male"]
pew_party <- pew_party %>% select(-share)
```

Example: Calibrating with Party ID

- Next, we'll recode our Party ID variable in the CCES to match

```
cces <- cces %>% mutate(party_bin = case_when(party_id == 1 ~ "Democratic",  
                                              party_id == 2 ~ "Republican",  
                                              party_id == 3 ~ "Independent",  
                                              party_id == 4 ~ "Other",  
                                              party_id == 5 ~ "Other",  
                                              is.na(party_id) ~ "Other"))
```

Example: Calibrating with Party ID

- Our **population target** distribution of gender and party ID

```
pew_party %>% group_by(gender_bin, party_bin) %>% summarize(n=Freq[1])
```

```
## # A tibble: 8 × 3
## # Groups:   gender_bin [2]
##   gender_bin party_bin      n
##   <chr>      <chr>    <dbl>
## 1 Female    Democratic 0.200
## 2 Female    Independent 0.154
## 3 Female    Other      0.0154
## 4 Female    Republican 0.144
## 5 Male      Democratic 0.127
## 6 Male      Independent 0.190
## 7 Male      Other      0.0195
## 8 Male      Republican 0.151
```


Example: Calibrating with Party ID

- Our **in-sample** distribution of gender and party ID

```
cces %>% group_by(gender_bin, party_bin) %>% summarize(n = n()) %>% ungroup() %>% mutate(prop =
```

```
## # A tibble: 8 × 3
##   gender_bin party_bin      prop
##   <chr>      <chr>      <dbl>
## 1 Female    Democratic  0.238
## 2 Female    Independent 0.142
## 3 Female    Other       0.0589
## 4 Female    Republican  0.138
## 5 Male      Democratic  0.135
## 6 Male      Independent 0.139
## 7 Male      Other       0.0382
## 8 Male      Republican  0.111
```

Example: Calibrating with Party ID

- Post-stratification on age x gender x education wasn't quite enough (in fact, we did *worse* for some bins)

```
cces %>% group_by(gender_bin, party_bin) %>% summarize(n = sum(postStratWt)) %>% ungroup() %>%
```

```
## # A tibble: 8 × 3
##   gender_bin party_bin    prop
##   <chr>      <chr>      <dbl>
## 1 Female    Democratic  0.210
## 2 Female    Independent 0.125
## 3 Female    Other       0.0513
## 4 Female    Republican  0.127
## 5 Male     Democratic  0.150
## 6 Male     Independent 0.156
## 7 Male     Other       0.0506
## 8 Male     Republican  0.131
```

Example: Calibrating with Party ID

- Setting up our entropy-weighting targets

```
calibration_targets <- acs_targets %>% arrange(strata) %>% pull(pop_proportion)
names(calibration_targets) <- colnames(model.matrix(~strata, acs_targets))
calibration_targets[1] <- 1 # Weights sum to 1 ("intercept")
pew_targets <- pew_party %>% mutate(strata2 = str_c(gender_bin, party_bin, sep="-")) %>% select(
calibration_targets_2 <- pew_targets %>% arrange(strata2) %>% pull(Freq)
names(calibration_targets_2) <- colnames(model.matrix(~strata2, pew_targets))
calibration_targets <- c(calibration_targets, calibration_targets_2[-1])
```

- Make the same strata in our dataset

```
cces$strata2 <- str_c(cces$gender_bin, cces$party_bin, sep="-")
```

- Find the entropy weights

```
cces$strata <- as.factor(cces$strata)
cces$strata2 <- as.factor(cces$strata2)
cces_surv_unwt <- svydesign(~1, weights = ~1, data=cces)
cces_entropy <- calibrate(cces_surv_unwt, formula = ~strata + strata2,
                          population=calibration_targets*nrow(cces), calfun="raking")
cces$entropyWt <- weights(cces_entropy)
```

Entropy weighting

- Our **population target** distribution of gender and party ID

```
pew_party %>% group_by(gender_bin, party_bin) %>% summarize(n=Freq[1])
```

```
## # A tibble: 8 × 3
## # Groups:   gender_bin [2]
##   gender_bin party_bin      n
##   <chr>      <chr>    <dbl>
## 1 Female    Democratic 0.200
## 2 Female    Independent 0.154
## 3 Female    Other      0.0154
## 4 Female    Republican 0.144
## 5 Male      Democratic 0.127
## 6 Male      Independent 0.190
## 7 Male      Other      0.0195
## 8 Male      Republican 0.151
```

Entropy weighting

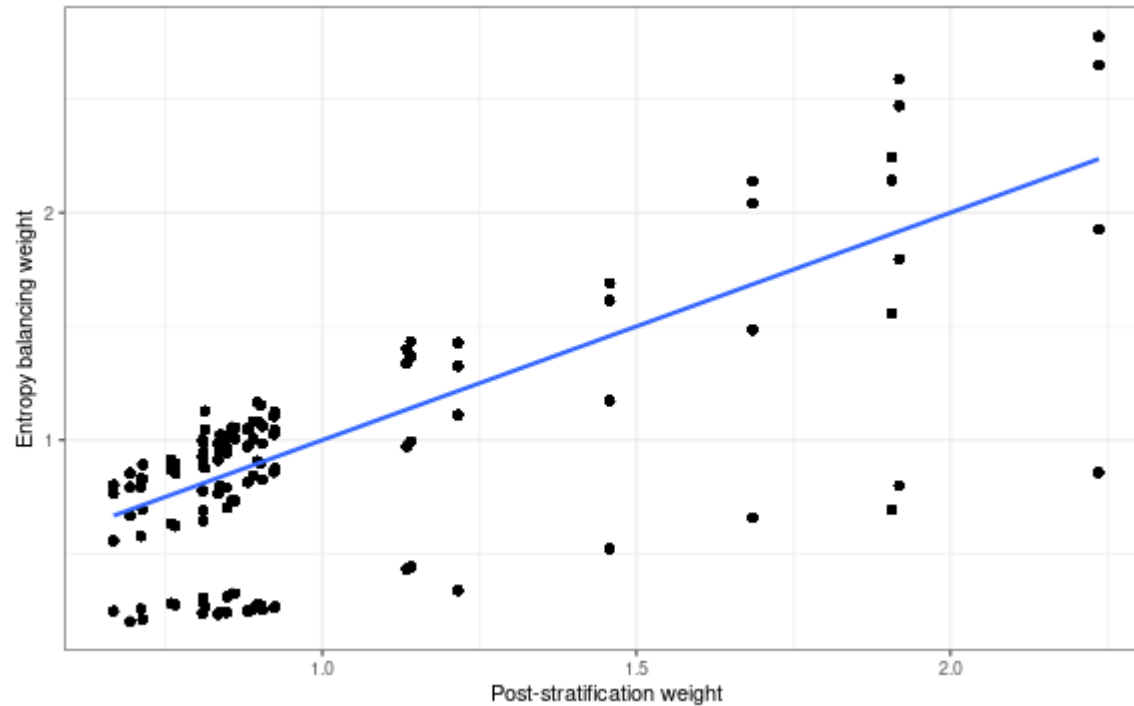
- Our **post-weighting** distribution of gender and party ID

```
cces %>% group_by(gender_bin, party_bin) %>% summarize(n = sum(entropyWt)) %>% ungroup() %>% mu
```

```
## # A tibble: 8 × 3
##   gender_bin party_bin      prop
##   <chr>      <chr>      <dbl>
## 1 Female    Democratic  0.200
## 2 Female    Independent 0.154
## 3 Female    Other       0.0154
## 4 Female    Republican  0.144
## 5 Male      Democratic  0.127
## 6 Male      Independent 0.190
## 7 Male      Other       0.0195
## 8 Male      Republican  0.151
```

Comparing the weights

```
cces %>% ggplot(aes(x=postStratWt, y=entropyWt)) + geom_point() + geom_smooth(method="lm") + theme_minimal() +  
  xlab("Post-stratification weight") + ylab("Entropy balancing weight")
```



Results

- Now that we've adjusted for party ID, we actually get closer to the right answer!

```
svymean(~trumpApprove, cces_entropy)
```

```
##               mean SE  
## trumpApprove 0.429  0
```

- And what's our effective sample size relative to our actual sample size?

```
ess_entropy <- (sum(cces$entropyWt)^2/sum(cces$entropyWt^2))  
ess_entropy
```

```
## [1] 50459
```

- How does this compare to our actual sample size (what's the efficiency loss?)

```
ess_entropy/sum(cces$entropyWt)
```

```
## [1] 0.828
```

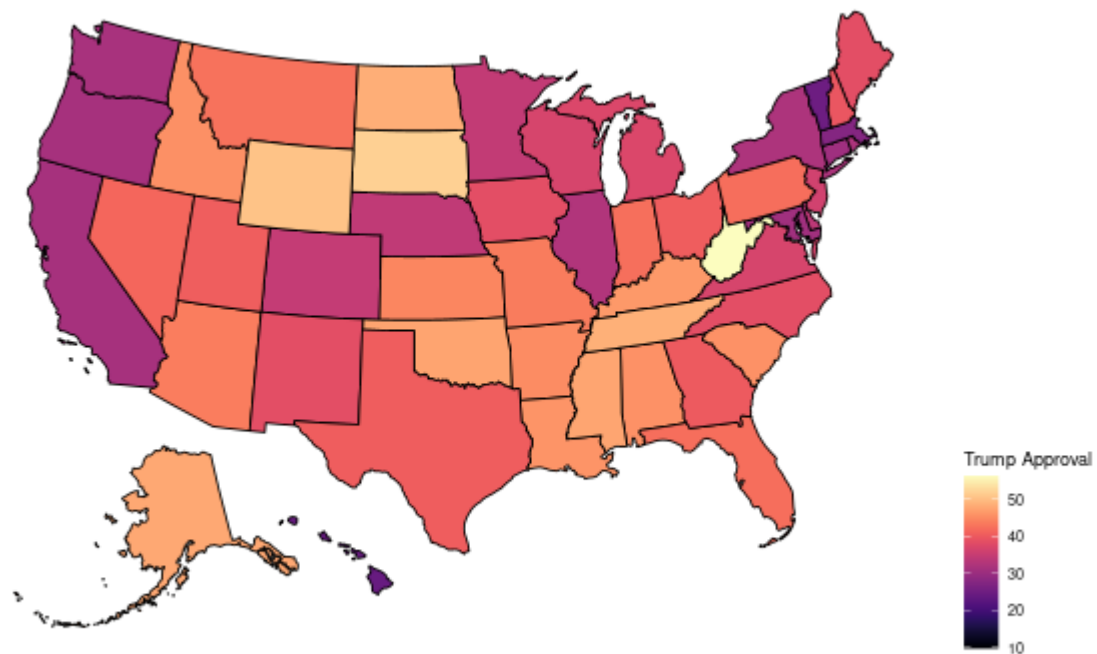
MrP

Multilevel Regression and Post-stratification

- Often we're interested in estimating political attitudes at very small areas
 - Within a nationally-representative sample, we may only have a minor fraction of respondents who are even in a particular region.
- Beyond that, our sample may not be demographically representative for a given region
 - Can we do better than just using the sub-group means for our outcome of interest?

CCES 2020

- Our **unweighted** average Trump support for each state:



MrP

- **Multilevel Regression** and **Post-stratification** combines two well-known tools to try to accomplish this task
 1. **Multilevel regression**: Using the entire dataset, specify a flexible model for the outcome of interest
 2. **Post-stratification**: Using known population characteristics, generate predictions for each covariate "cell" in each region
- The national-level regression allows us to improve **precision** in small-area estimation
 - e.g. if some of the variability between regions is driven by individual or group-level characteristics, we can model the relationship between those characteristics and the outcome using **all** the data.
- The population-level data allows us to address non-response/sampling bias
 - Ideally all of the variables that would go in the survey weighting model also go in the MrP model (so in principle we don't need to include survey weights in the analysis)

MrP Workflow

1. Get survey data that contains individual and group-level predictors (e.g. demographic characteristics + "small area" of residence)
2. Get information on the "small areas" themselves
3. Estimate a multilevel regression model using these two sets of data
4. Construct a post-stratification frame that has the joint distribution of our demographics for each "small area"
5. Predict using the model for each cell in the post-stratification frame
6. Aggregate to get predictions for each small area

Example: CCES 2020

- Let's fit a regression model
 - We'll use the **brms** package which has pre-existing implementations of popular models like the random slopes/random intercepts generalized linear mixed model
- We'll fit a simple logistic regression with coefficients on gender, age and education as well as a random intercept for state

```
library(brms)
```

```
state_mlm <- brm(trumpApprove ~ factor(gender_bin) + factor(age_bin) + factor(educ_bin) + (1 | state),  
                 family = bernoulli(), data=cces, cores=4, warmup=1000, iter=2000)
```

Model diagnostics

- Summarize the results

```
summary(state_mlm)
```

```
## Family: bernoulli
## Links: mu = logit
## Formula: trumpApprove ~ factor(gender_bin) + factor(age_bin) + factor(educ_bin) + (1 | inputstate)
## Data: cces (Number of observations: 60968)
## Draws: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
## total post-warmup draws = 4000
##
## Group-Level Effects:
## ~inputstate (Number of levels: 51)
##
```

	Estimate	Est.Error	l-95% CI	u-95% CI	Rhat	Bulk_ESS	Tail_ESS
sd(Intercept)	0.32	0.04	0.26	0.41	1.00	646	1187

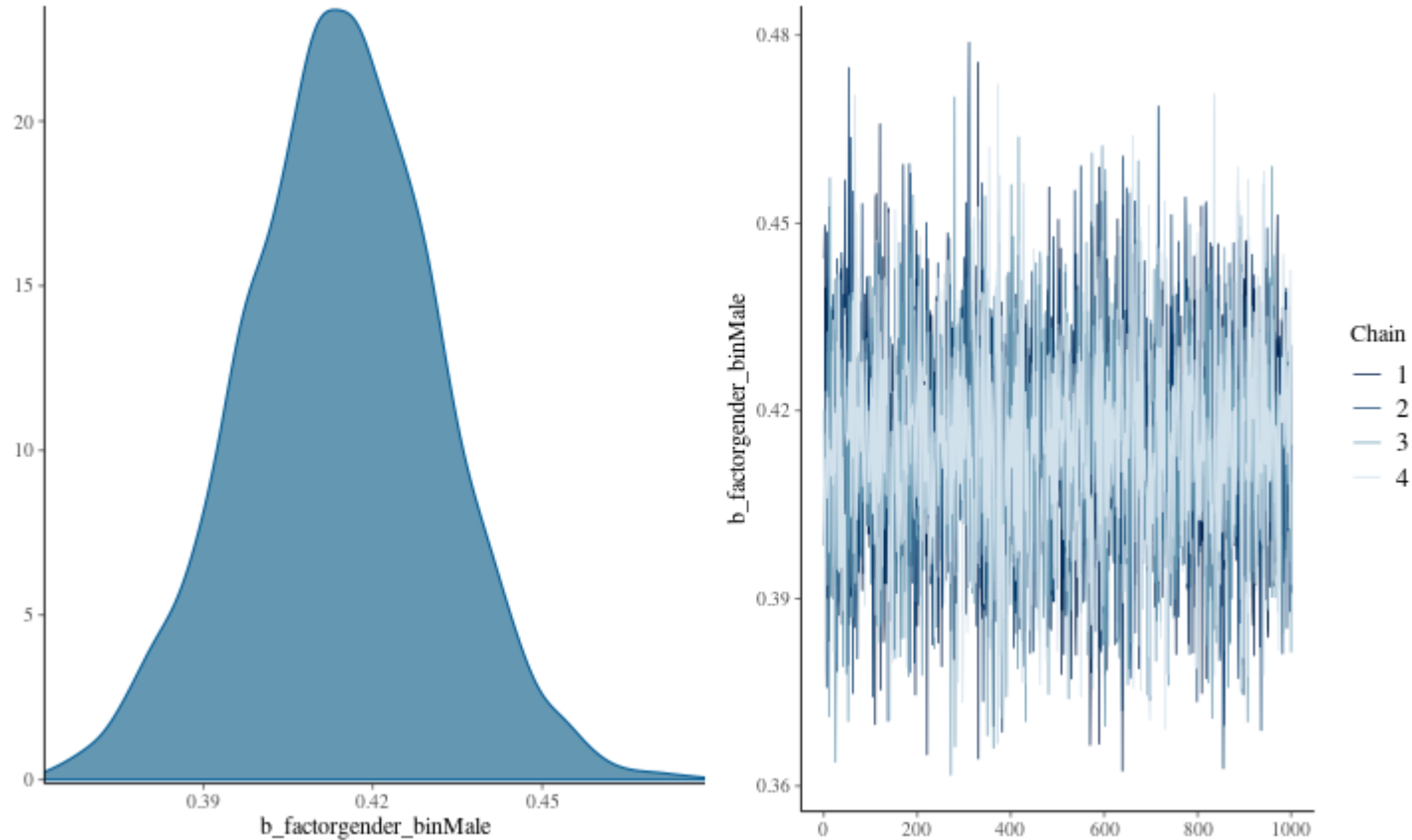
```
##
## Population-Level Effects:
##
```

	Estimate	Est.Error	l-95% CI	u-95% CI	Rhat	Bulk_ESS
Intercept	-1.31	0.05	-1.43	-1.21	1.00	460
factorgender_binMale	0.41	0.02	0.38	0.45	1.00	4202
factorage_bin30M44	0.23	0.03	0.17	0.28	1.00	2924
factorage_bin45M64	0.58	0.02	0.53	0.63	1.00	2947
factorage_bin65P	0.76	0.03	0.71	0.81	1.00	2874
factoreduc_binH.S.orless	0.61	0.02	0.56	0.66	1.00	2758
factoreduc_binPostgraduate	-0.37	0.03	-0.44	-0.31	1.00	3060
factoreduc_binSomecollege	0.30	0.02	0.25	0.34	1.00	2739

Model diagnostics

- Check convergence

```
plot(state_mlm, variable = "b_factorgender_binMale")
```



Generating the post-stratification frame

- We'll use ACS 5-year microdata data at the state level to construct our post-stratification frame for each state

```
acs_state_targets <- read_csv("data/ACS_2020_microdata_by_state.csv") %>% filter(!is.na(Gender))
fips_codes <- read_csv("data/us-state-ansi-fips.csv")
acs_state_targets <- acs_state_targets %>% left_join(fips_codes, by=c(State="stname"))
acs_state_targets$inputstate <- as.numeric(acs_state_targets$st)
acs_state_targets <- acs_state_targets %>% pivot_longer(c(`18-29`, `30-44`, `45-64`, `65+`))
acs_state_targets <- acs_state_targets %>% rename(gender_bin = Gender, educ_bin = Education, ac
acs_state_targets <- acs_state_targets %>% mutate(prop_us = Count/sum(Count))
acs_state_targets <- acs_state_targets %>% group_by(inputstate) %>% mutate(prop_state = Count/s
```


Generating the post-stratification frame

- Each row contains a row for each **state** x **age** x **gender** x **education** combination

```
acs_state_targets
```

```
## # A tibble: 1,632 × 10
##   State  gender_...1 educ_...2 st   stusps input...3 age_bin  Count prop_us prop_...4
##   <chr>   <chr>      <chr>   <chr> <chr>   <dbl> <chr>   <dbl>   <dbl>   <dbl>
## 1 Alabama Male      H.S. o... 01    AL         1 18-29   195648 7.72e-4 0.0515
## 2 Alabama Male      H.S. o... 01    AL         1 30-44   198528 7.84e-4 0.0522
## 3 Alabama Male      H.S. o... 01    AL         1 45-64   294365 1.16e-3 0.0775
## 4 Alabama Male      H.S. o... 01    AL         1 65+     168153 6.64e-4 0.0442
## 5 Alabama Male      Some c... 01    AL         1 18-29   148968 5.88e-4 0.0392
## 6 Alabama Male      Some c... 01    AL         1 30-44   129718 5.12e-4 0.0341
## 7 Alabama Male      Some c... 01    AL         1 45-64   171277 6.76e-4 0.0451
## 8 Alabama Male      Some c... 01    AL         1 65+      91367 3.61e-4 0.0240
## 9 Alabama Male      Colleg... 01    AL         1 18-29    41573 1.64e-4 0.0109
## 10 Alabama Male      Colleg... 01    AL         1 30-44    73950 2.92e-4 0.0195
## # ... with 1,622 more rows, and abbreviated variable names 1gender_bin,
## # 2educ_bin, 3inputstate, 4prop_state
```

Predicting with the model

- Use the model to make predictions on each row of our post-stratification frame

```
acs_state_targets$predict_mlm <- predict(state_mlm, newdata=acs_state_targets)[,1]
```

- Aggregate the predictions by the stratum shares to get state-level estimates!

```
mrp_states <- acs_state_targets %>% group_by(st) %>% summarize(trumpApproveMLM = sum(predict_mlm
```

Compare

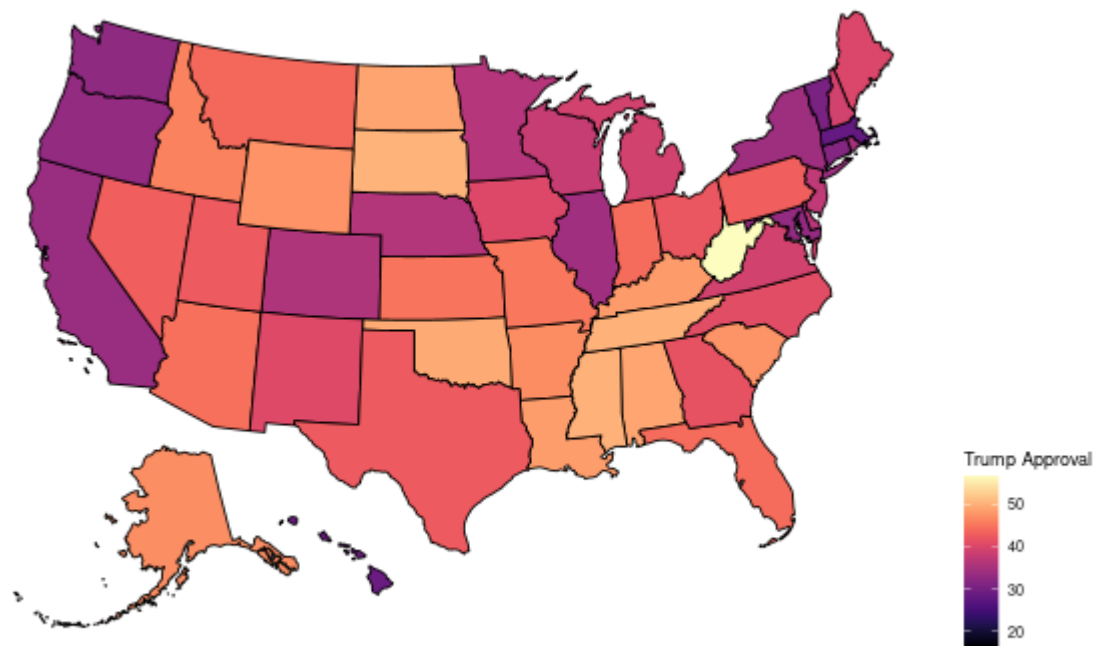
- Use the model to make predictions on each row of our post-stratification frame

```
state_summaries %>% left_join(mrp_states, by="fips")
```

```
## # A tibble: 51 × 3
##   fips values trumpApproveMLM
##   <chr>   <dbl>          <dbl>
## 1 01      45.6           49.0
## 2 02      47.8           47.1
## 3 04      43.7           44.6
## 4 05      45.1           47.0
## 5 06      31.2           33.5
## 6 08      34.5           35.7
## 7 09      30.4           31.4
## 8 10      33.8           36.4
## 9 11       9.64          15.9
## 10 12      42.2           44.1
## # ... with 41 more rows
```

Visualize

- Our **multilevel regression + post-stratification** estimates



How to improve this model

- We've provided a very simple example of an MrP model but there are a lot of ways it can be improved
 1. More group-level predictors (e.g. Trump state-level vote share)
 2. Individual-level interactions in the model + partial-pooling on the coefficients
 3. Actually knowing the joint distribution of the demographics in each state!
 4. Pool state-level means to regional rather than a common "grand" mean.
 5. Many more possibilities...
- More generally, we have a dataset where even within each state, we have many observations to work with, so the model isn't that necessary
 - MrP is most powerful when there are **very few** respondents from a particular small area in our data

