

governing governing

Spencer Braun

11/12/2020

```
library(DOS2)
library(optmatch)
library(RIttools)
library(rcbalance)

library(readstata13)
library(tidyverse)
library(haven)

source("utility.R")
```

Main Analysis

Here I do two things: define an governing dataframe with the covariates used in the paper and define a larger governing dataframe including more covariates. I will compare these approaches in generating propensity scores.

```
main <- read.csv("processed_data/main.csv")
# main %>% head()
# main %>% names()

analysis.cols <- main %>% select(
  p2p,
  current_vote,
  transfersshare,
  donationssshare,
  gov,
  y2006,
  female,
  exminister,
  previous_vote,
  media_mentions,
  quality,
  years_served
)

covariates <- c(
  "transfersshare",
  "donationssshare",
  "y2006",
  "female",
  "exminister",
```

```

"previous_vote",
"media_mentions",
"quality",
"years_served"
)

sapply(analysis.cols, function(x) sum(is.na(x)))

##           p2p    current_vote transfersshare donationsshare          gov
##           0             0             29             29             0
##        y2006        female    exminister  previous_vote media_mentions
##           0             0             0             0             2
##      quality    years_served
##           0             0

# for now, drop NAs

analysis <- analysis.cols %>%
  drop_na() %>%
  rename(zb = p2p) %>%
  rename(y = current_vote)

analysis.cols.large <- main %>% select(
  p2p,
  current_vote,
  transfersshare,
  donationsshare,
  gov,
  y2006,
  female,
  exminister,
  previous_vote,
  media_mentions,
  quality,
  years_served,
  # province, #prov if need numeric
  election,
  winner,
  pop_per_km2,
  immigrants,
  citizens,
  unemployment_rate,
  median_family_income
)

analysis.large <- analysis.cols.large %>%
  drop_na() %>%
  rename(zb = p2p) %>%
  rename(y = current_vote)

governing <- analysis %>% filter(gov == 1) %>% select(-gov)
governing.large <- analysis.large %>% filter(gov == 1) %>% select(-gov)

```

```
governing %>%
  group_by(zb) %>%
  summarise(count=n())

## `summarise()` ungrouping output (override with `.groups` argument)
## # A tibble: 2 x 2
##       zb count
##   <int> <int>
## 1     0    63
## 2     1    35
```

Propensity Scores

```
prop.scores <- glm(zb ~ . - y, family=binomial, data=governing)
governing$prop <- prop.scores$fitted.values

prop.scores.large <- glm(zb ~ . - y, family=binomial, data=governing.large)
governing.large$prop.large <- prop.scores.large$fitted.values

governing$prop.large <- prop.scores.large$fitted.values
```

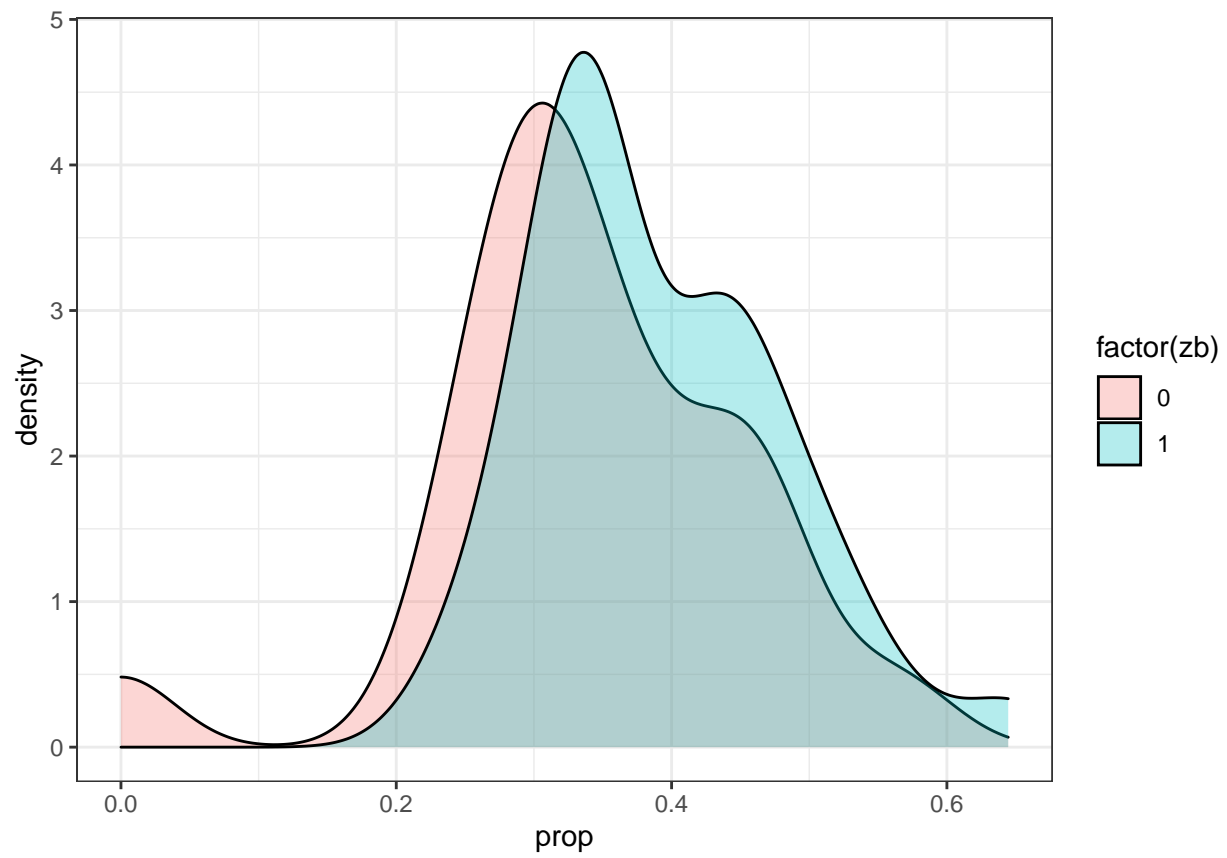
Imbens / Rubin page 282:

First, it is important to note, however, that the goal is not simply to get the best estimate of the propensity score in terms of mean-integrated-squared-error, or a similar criterion based on minimizing the difference between the estimated and true propensity score. Such a criterion would always suggest that using the true propensity score is preferable to using an estimated propensity score. In contrast, for our purposes, it is often preferable to use the estimated propensity score. The reason is that using the estimated score may lead to superior covariate balance in the sample compared to that achieved when using the true super-population propensity score.

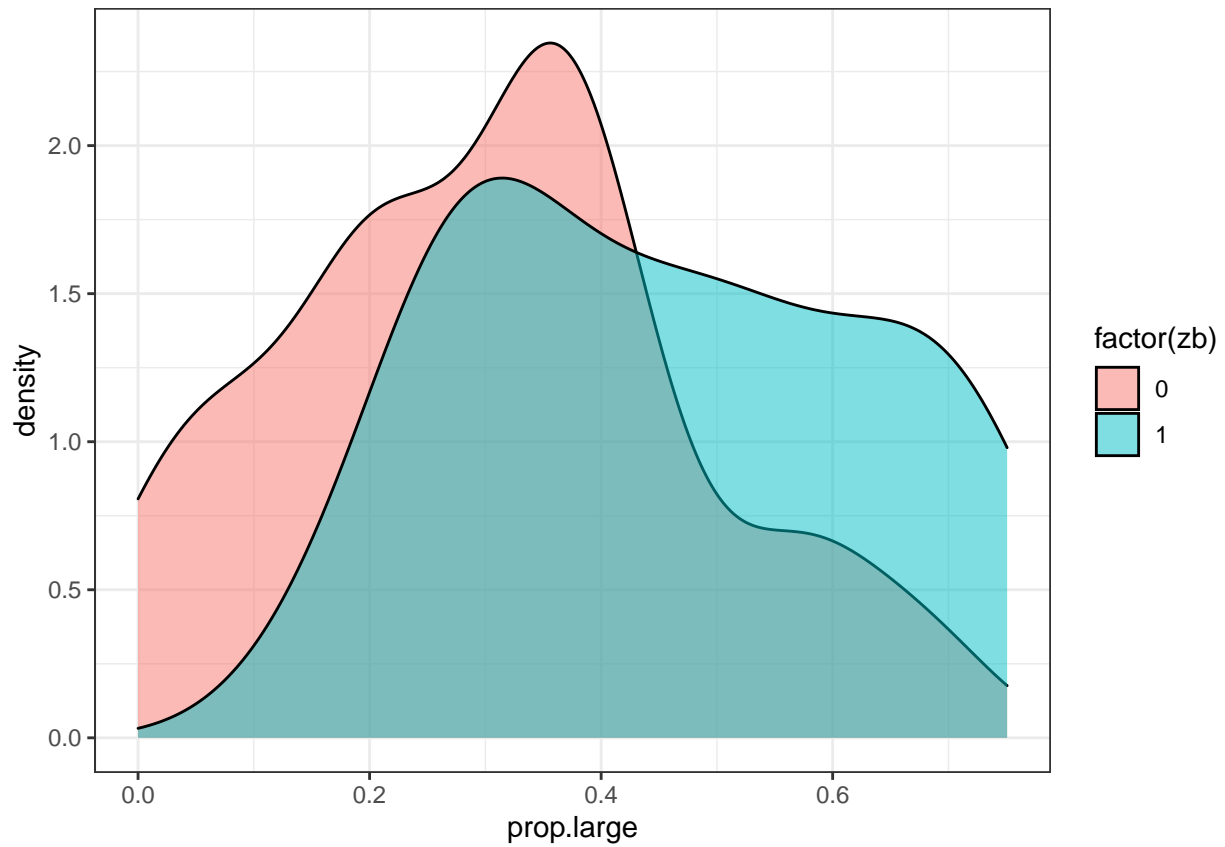
This makes the simpler prop score look better since they are more balanced between treatment and control.

Overlap much better for simple prop score

```
governing %>% ggplot() +
  geom_density(aes(x=prop, fill=factor(zb)), alpha=0.3) +
  theme_bw()
```

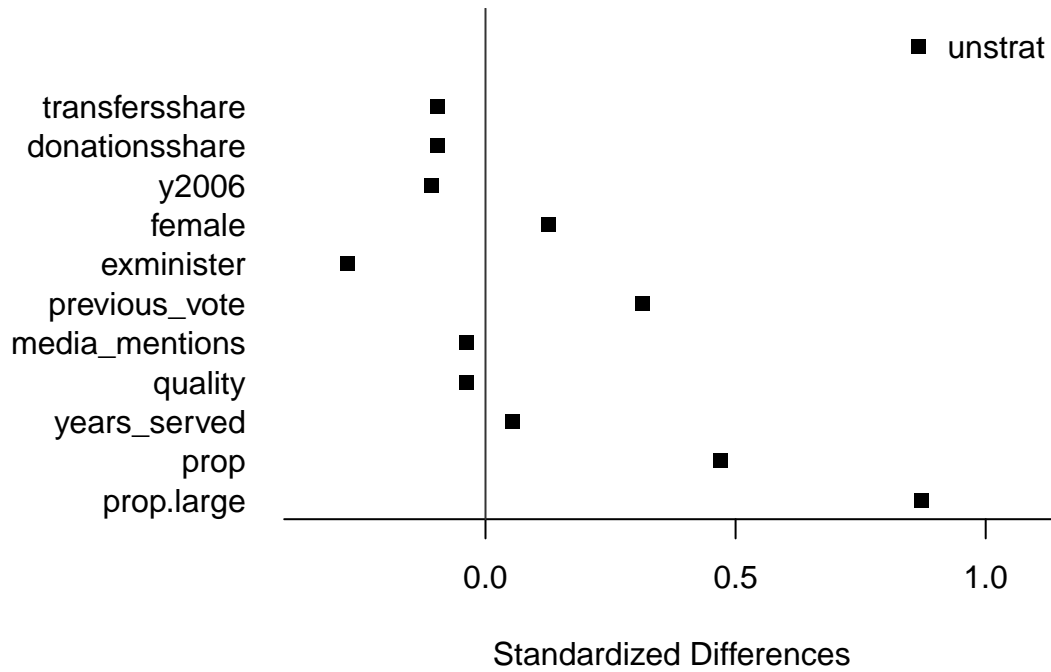


```
governing %>% ggplot() +  
  geom_density(aes(x=prop.large, fill=factor(zb)), alpha=0.5) +  
  theme_bw()
```



Matching

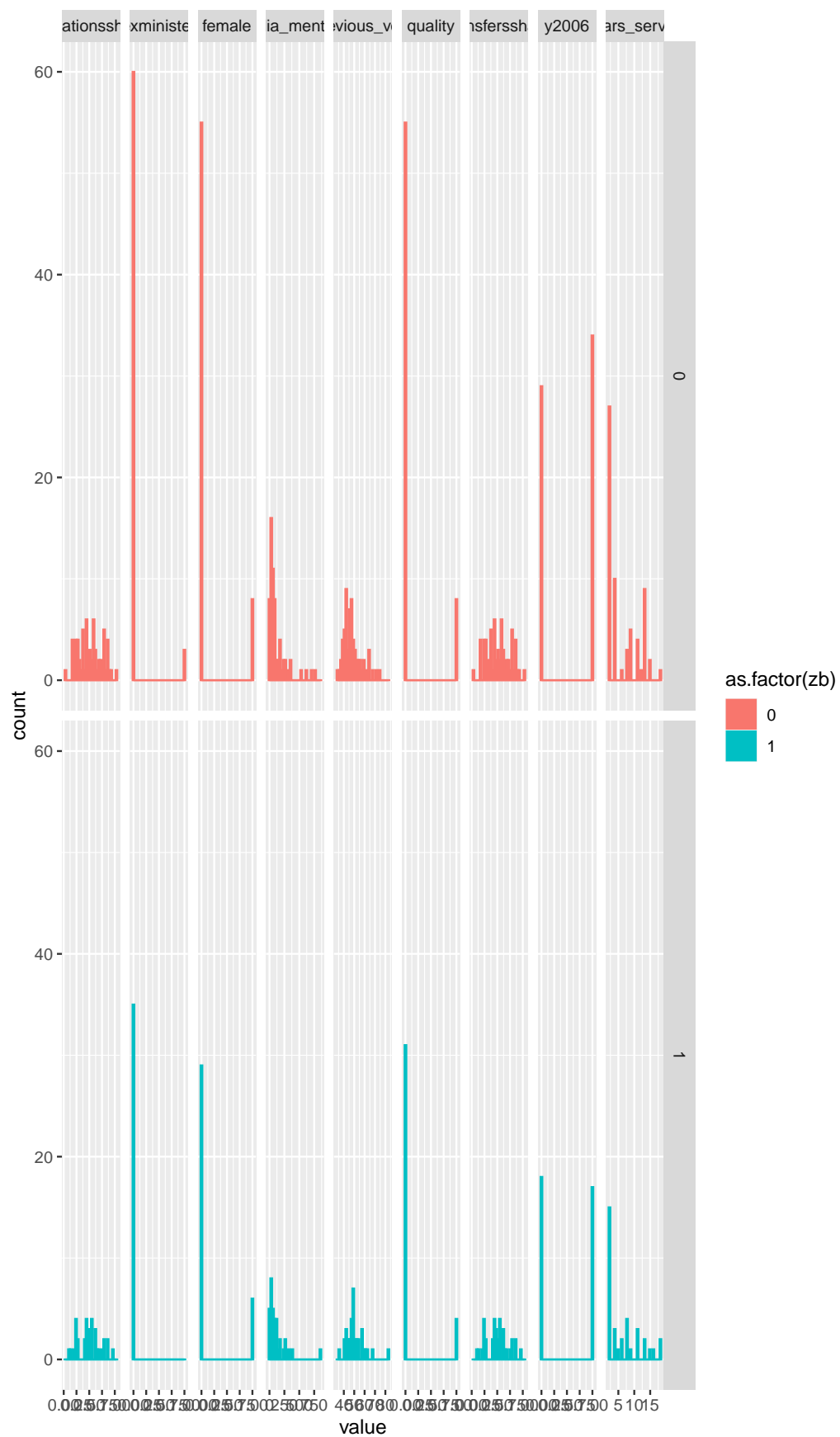
```
plot(xBalance(zb ~ . - 1 - y, data=governing))
```



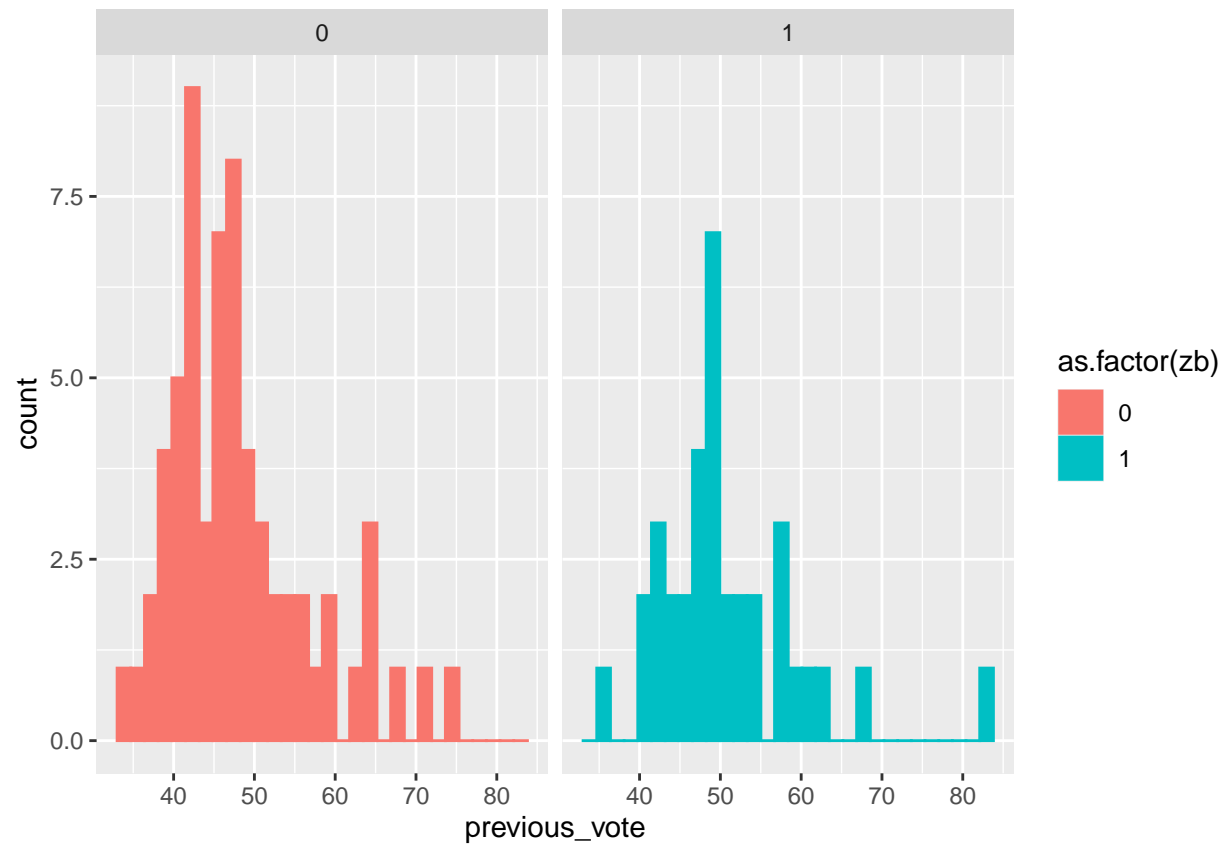
```
# plot(xBalance(zb ~ . - 1 -y, data=governing.large))
```

```
governing %>%  
  dplyr::select(zb, covariates) %>%  
  pivot_longer(-zb, names_to="covariate", values_to="value") %>%  
  ggplot(aes(x = value, color = as.factor(zb), fill = as.factor(zb))) +  
  geom_histogram(bins=30) +  
  facet_grid(cols=vars(covariate), rows=vars(zb), scales='free_x')
```

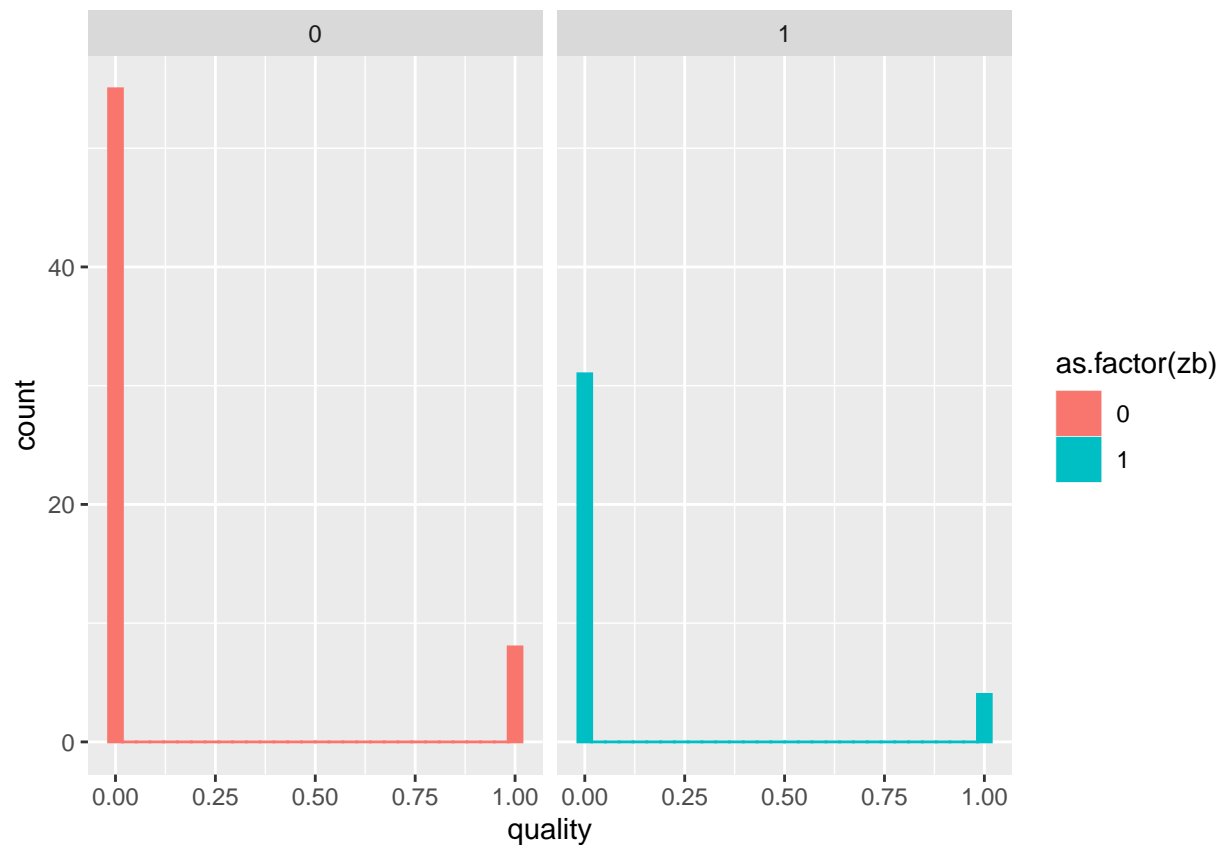
```
## Note: Using an external vector in selections is ambiguous.  
## i Use `all_of(covariates)` instead of `covariates` to silence this message.  
## i See <https://tidyselect.r-lib.org/reference/faq-external-vector.html>.  
## This message is displayed once per session.
```



```
governing.large %>%
  ggplot(aes(x = previous_vote, color = as.factor(zb), fill = as.factor(zb))) +
  geom_histogram(bins=30) +
  facet_wrap(~zb)
```



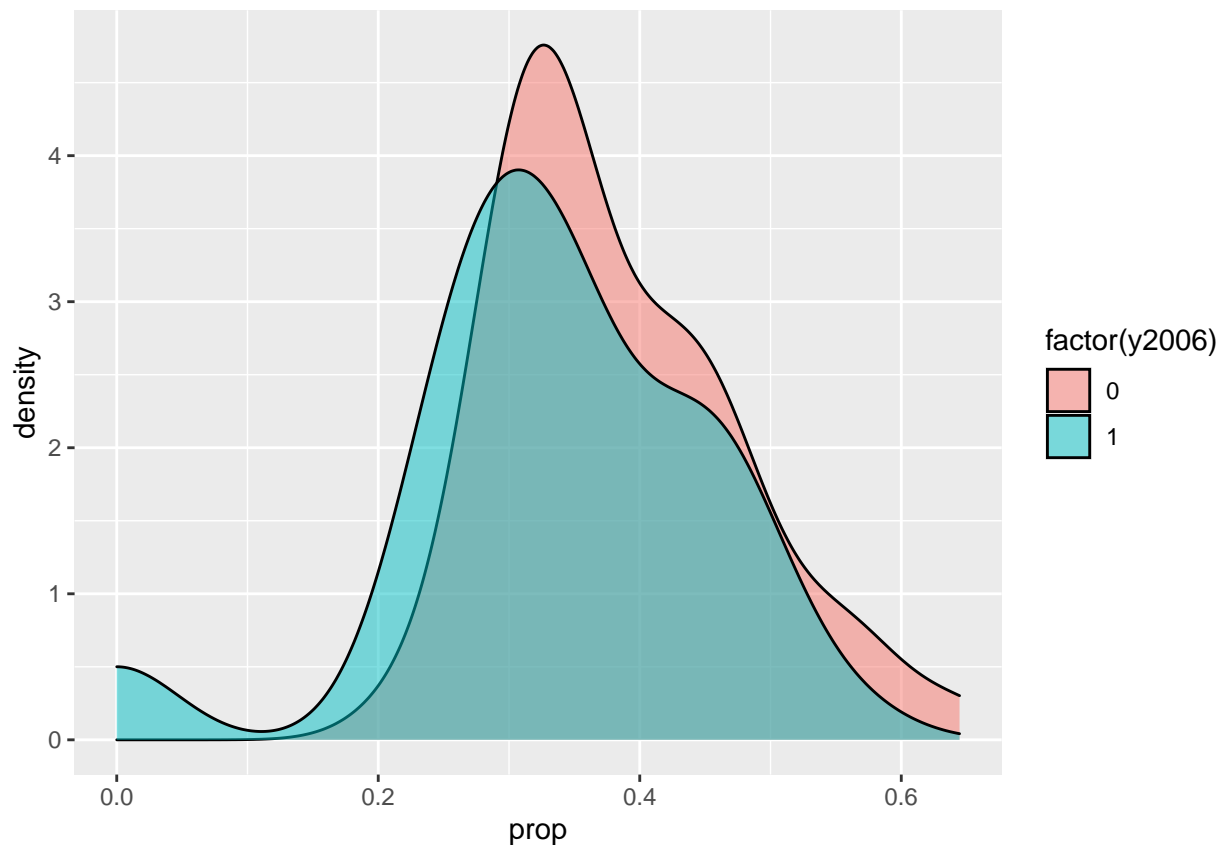
```
governing %>%
  ggplot(aes(x = quality, color = as.factor(zb), fill = as.factor(zb))) +
  geom_histogram(bins=30) +
  facet_wrap(~zb)
```

```
names(governing)
```

```
## [1] "zb"          "y"           "transfersshare" "donationssshare"
## [5] "y2006"       "female"      "exminister"    "previous_vote"
## [9] "media_mentions" "quality"     "years_served"  "prop"
## [13] "prop.large"
```

```
governing %>%
  ggplot() +
  geom_density(aes(x=prop, fill=factor(y2006)), alpha=0.5)
```



1:1 Exact Matching Note: we have the potential to match members to themselves, since ~150 had power to propose in one of the two years. However `y2006` is a covariate in the distance formula and imbalance is low, meaning we are mostly matching within the same year. We could add the specification that it has to match exactly on `y2006` so no member matches to themselves.

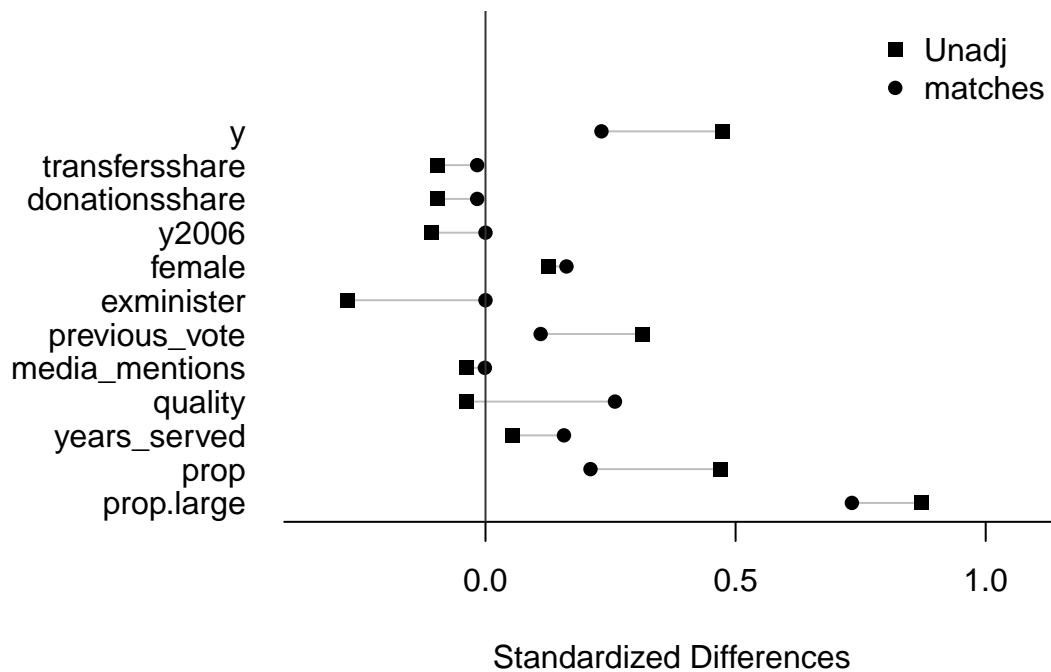
Here I do matching with an without a propensity caliper, but I end up using the propensity caliper one since the prop scores are still one of the worse balanced measures between groups.

```
z <- governing$zb
X <- governing %>% dplyr::select(-c(zb, prop, y, prop.large))

distance <- smahal(z, X)
distance.cal <- addcaliper(distance, z=governing$zb, p=governing$prop, caliper=0.1)

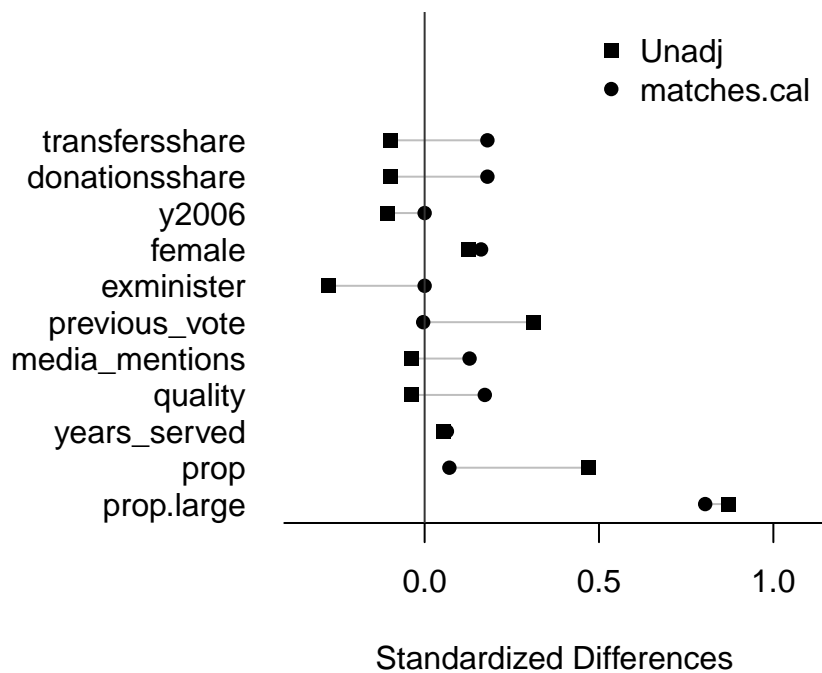
# governing %>%
#   group_by(zb) %>%
#   summarise(count=n())

matches <- pairmatch(distance, data=governing)
plot(xBalance(zb ~ . - 1 + strata(matches) , data=governing))
```



```
matches.cal <- pairmatch(distance.cal, data=governing)
matches.df <- summarize.match(governing, matches.cal)

plot(xBalance(zb ~ . + strata(matches.cal) - 1 - y, data=governing))
```



I checked on the NAs in the matches - turns out these are just the extra control units. All treatment units are 1:1 matched with a control unit, so we should be 100% good here.

```
sum(is.na(matches)) #NAs are just extra controls !
```

```
## [1] 28
```

```
sum(!is.na(matches))/2 == governing %>% summarise(sum(zb))
```

```
##      sum(zb)
## [1,]    TRUE
```

```
# governing[which(is.na(matches)),]
```

```
sum(is.na(matches.cal))
```

```
## [1] 28
```

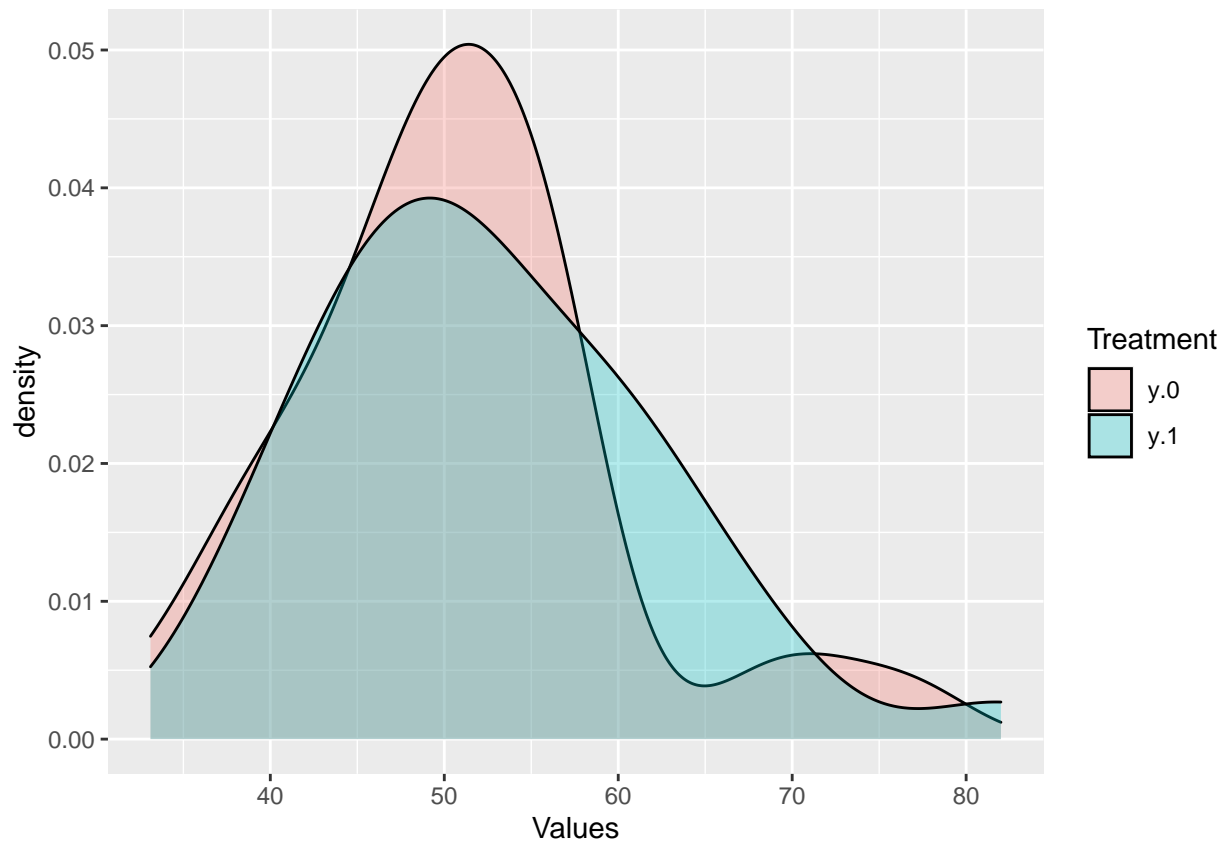
```
sum(!is.na(matches.cal))
```

```
## [1] 70
```

FRT

In matched pairs, the vote share received looks nearly identical.

```
matches.df %>%
  select(y.1, y.0) %>%
  pivot_longer(c(y.1, y.0), names_to='Treatment', values_to='Values') %>%
  ggplot() +
  geom_density(aes(x=Values, fill=Treatment), alpha=0.3)
```



Run the FRT and our p-value is large, and tau (T.obs) small and opposite the expected sign.

```
# DIM current assignment
T.obs1 <- matches.df %>%
  summarise(mean(y.1) - mean(y.0)) %>%
```

```

pull(.)

T.obs1 #T obs small and negative!

## [1] 1.882858

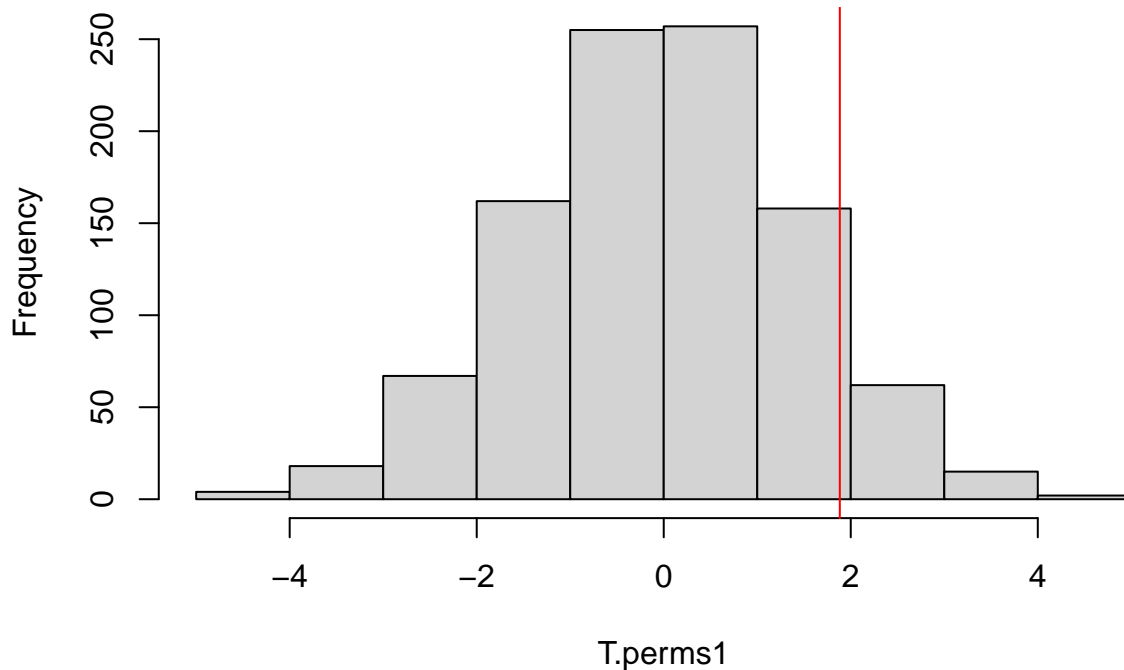
genPermute <- function(x, matches) {
  treated_unit <- sample(c(0,1), nrow(matches), replace=TRUE)
  matches %>%
    select(y.0, y.1) %>%
    mutate(treated=treated_unit) %>%
    mutate(treated_y = ifelse(treated == 1, y.1, y.0),
           control_y = ifelse(treated == 1, y.0, y.1)) %>%
    summarise(mean(treated_y) - mean(control_y)) %>%
    pull(.)
}

set.seed(123)
iters <- 1000
reps <- rep(NA, iters)

# Generate vector of test statistics under permutations
T.perms1 <- sapply(reps, function(x) genPermute(x, matches.df))
hist(T.perms1)
abline(v=T.obs1, col='red')

```

Histogram of T.perms1



```

pval <- (1/iters) * (sum(ifelse(T.perms1 >= T.obs1, 1, 0)) ) #calculate one sided p-value
pval

```

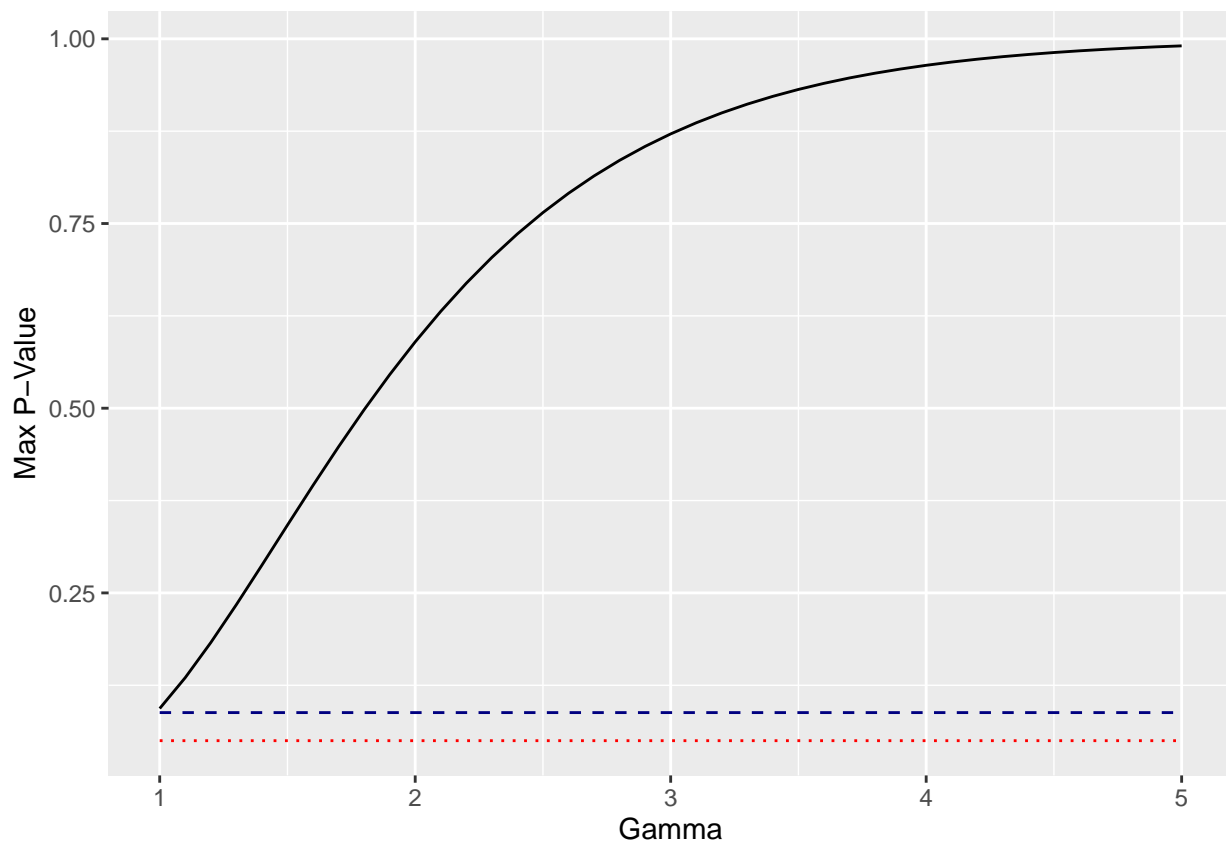
```
## [1] 0.088
```

Run sensitivity governing over all values of gamma needed. Uninterestingly since it starts insignificant it only gets more so.

```
senm.data1 <- cast.senm(governing, matches.cal)
sen.out <- sensitivitymult::senm(senm.data1$y, senm.data1$z, senm.data1$mset, gamma=1.2, trim = Inf, in

gamma.compute <- function(gamma) sensitivitymult::senm(senm.data1$y, senm.data1$z, senm.data1$mset, gam
grange = seq(1,5,by=0.1)
sensitivity <- sapply(grange, gamma.compute)

data.frame(gamma=grange, pvalue=sensitivity) %>%
  ggplot(aes(x=gamma, y=sensitivity)) +
  geom_line() +
  geom_line(aes(x=seq(1,5,length=length(grange)), y=pval), col='navy', linetype='dashed') +
  geom_line(aes(x=seq(1,5,length=length(grange)), y=0.05), col='red', linetype='dotted') +
  labs(x="Gamma", y="Max P-Value")
```



IPW Estimators

Horvitz Thompson I fixed the HT estimator - the denominator was wrong. The nice thing about these IPW estimators is they do not use the matches, so kind of a separately analytical check that there is an issue with the raw difference in means.

```
horvitz.thompson <- governing %>%
  mutate(ZY1 = (zb * y)/prop) %>%
  mutate(ZY0 = ((1-zb) * y)/(1-prop)) %>%
  summarise(
    Units = n(),
```

```

YBar1 = sum(ZY1) / nrow(governing),
YBar0 = sum(ZY0) / nrow(governing),
.groups='drop'
) %>%
mutate(DIM=(YBar1 - YBar0)) %>%
summarise(sum(DIM)) %>%
pull(.)

```

horvitz.thompson

[1] 0.3821129

```

hayek <- governing %>%
  mutate(ZY1 = (zb * y)/prop) %>%
  mutate(ZY0 = ((1-zb) * y)/(1-prop)) %>%
  summarise(
    Units = n(),
    YBar1 = sum(ZY1) / sum(zb / prop),
    YBar0 = sum(ZY0) / sum((1-zb)/(1-prop)),
    .groups='drop'
  ) %>%
  mutate(DIM=(YBar1 - YBar0)) %>%
  summarise(sum(DIM)) %>%
  pull(.)

```

hayek

Hayek

[1] 2.328192

Subclassification with Neymanian CI's

```
quant.vec <- quantile(governing$prop, c(0.2, 0.4, 0.6, 0.8))
```

```

stratified <- governing %>%
  mutate(stratum = case_when(
    prop < quant.vec[1] ~ 1,
    (quant.vec[1] <= prop) & (prop < quant.vec[2]) ~ 2,
    (quant.vec[2] <= prop) & (prop < quant.vec[3]) ~ 3,
    (quant.vec[3] <= prop) & (prop < quant.vec[4]) ~ 4,
    prop >= quant.vec[4] ~ 5
  ))

```

Number of units by stratum and treatment status

```

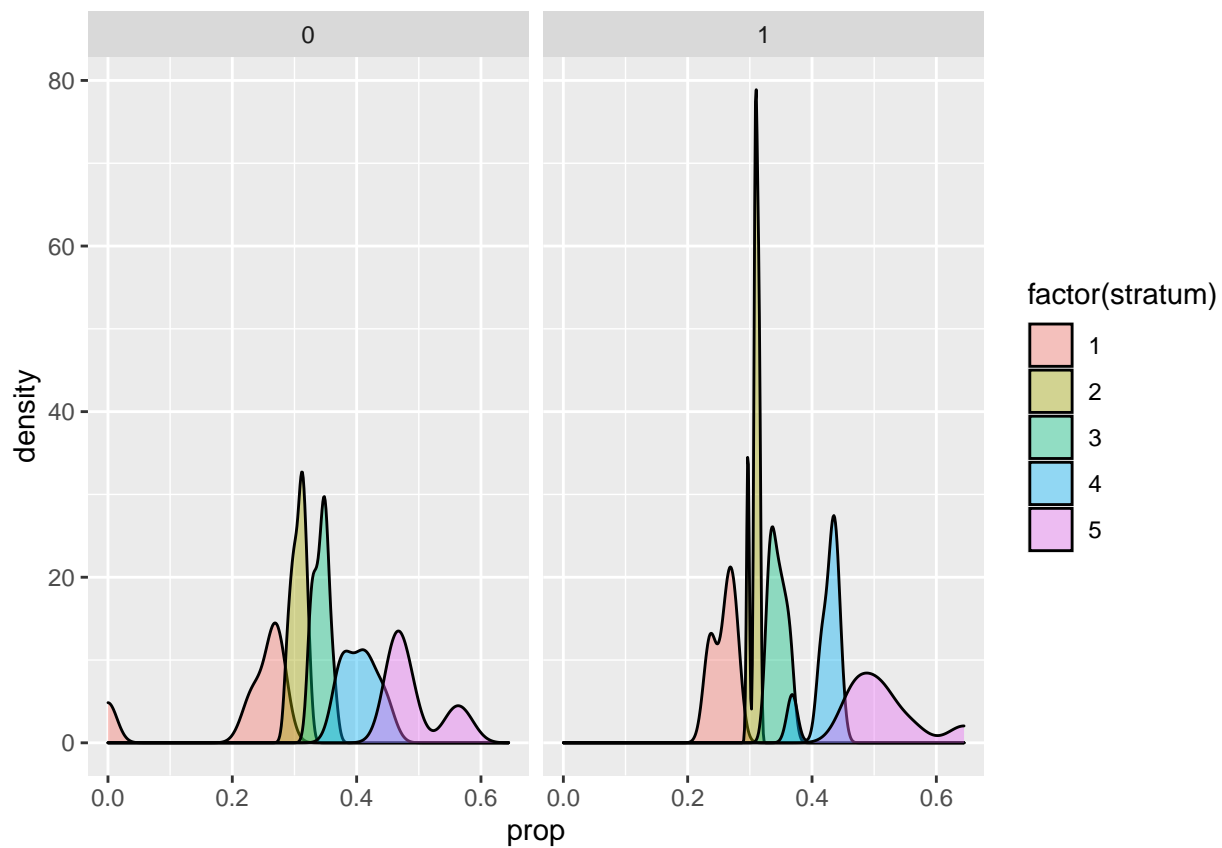
stratified %>%
  group_by(stratum) %>%
  summarise(
    treated = sum(zb),
    control = sum(1-zb),
    .groups='drop'
  )

```

A tibble: 5 x 3

```
##   stratum treated control
##   <dbl>   <int>   <dbl>
## 1      1      3     17
## 2      2      5     14
## 3      3     11      9
## 4      4      8     11
## 5      5      8     12
```

```
stratified %>%
  ggplot() +
  geom_density(aes(x=prop, fill=factor(stratum)), alpha=0.4) +
  facet_wrap(vars(zb))
```



```
tau_k <- stratified %>%
  mutate(ZY1 = zb * y) %>%
  mutate(ZY0 = (1-zb) * y) %>%
  group_by(stratum) %>%
  summarise(
    Units = n(),
    YBar1 = sum(ZY1) / sum(zb),
    YBar0 = sum(ZY0) / sum(1-zb),
    .groups='drop'
  ) %>%
  mutate(DIM=(YBar1 - YBar0) * Units/nrow(stratified)) %>%
  summarise(sum(DIM)) %>%
  pull(.)
```

```
tau_k
```



```
## [1] 2.646312
```

```
tau_k.var <- stratified %>%  
  group_by(stratum, zb) %>%  
  summarise(  
    N = n(),  
    Var = var(y),  
    .groups='drop'  
  ) %>%  
  mutate(weighted_V = Var / N) %>%  
  group_by(stratum) %>%  
  summarise(  
    stratum_var = sum(weighted_V) * (sum(N) / nrow(stratified))^2,  
    .groups='drop'  
  ) %>%  
  summarise(sum(stratum_var)) %>%  
  pull(.)
```

```
tau_k.var
```

```
## [1] 3.080705
```

Compute a 95% confidence interval

```
normalCI <- function(tau, variance) {  
  c(tau - sqrt(variance)*qnorm(.975), tau + sqrt(variance)*qnorm(.975))  
}
```

```
normalCI(tau_k, tau_k.var)
```

```
## [1] -0.7938045  6.0864285
```

Null result