

Pset1

Gabriel

October 9, 2025

1 Simulation

```
1 # Step 1: Create a population
2 # Set seed for reproducibility
3 set.seed(123)
4
5 # Define population size
6 N <- 100000
7 # Define population with some traits
8 prop_party_Pop <- c(0.4, 0.4, 0.2) # define party distribution in population
9 prop_gender_Pop <- c(0.5, 0.5) # define gender distribution in population
10 prop_income_Pop <- c(0.4, 0.5, 0.1) # define income distribution in population
11
12 population <- data.frame(
13   # income status: Low, Medium, High
14   income = sample(c("Low", "Medium", "High"), N, replace = TRUE, prob = prop_party_
15     Pop),
16   # Gender: Male, Female
17   gender = sample(c("Male", "Female"), N, replace = TRUE, prob = prop_gender_Pop),
18   # Party: Democrat, Republican, Independent
19   party = sample(c("Democrat", "Republican", "Independent"), N, replace = TRUE, prob
20     = prop_income_Pop)
21 )
22
23 head(population)
24
25 # sample sizes to demonstrate "as n increases"
26 n_vals <- c(100, 200, 500, 1000, 5000, 8000)
27
28 # ---- simulate once for each n ----
29 library(tidyverse)
30 party_proportions <- data.frame() # define once outside the loop
31
32 for (n in n_vals){
33   # Step 1: sample n observations
34   sample_data <- sample_n(population, size = n, replace = TRUE)
35
36   # Step 2: random assignment; 1 = treatment, 0 = control
37   treatment <- rbinom(n, size = 1, prob = 0.5)
38
39   # Step 3: make dataframe
40   df <- data.frame(sample_data, treat = treatment)
41
42   # Step 4: Calculate proportions
43   prop_party_treat <- prop.table(table(df$party[df$treat == 1]))
44   prop_party_control <- prop.table(table(df$party[df$treat == 0]))
45   prop_gender_treat <- prop.table(table(df$gender[df$treat == 1]))
46   prop_gender_control <- prop.table(table(df$gender[df$treat == 0]))
47   prop_income_treat <- prop.table(table(df$income[df$treat == 1]))
48   prop_income_control <- prop.table(table(df$income[df$treat == 0]))
49
50   # proportions for all
```

```

51 prop_party <- prop.table(table(df$party))
52 prop_gender <- prop.table(table(df$gender))
53 prop_income <- prop.table(table(df$income))
54
55 party_proportions <- bind_rows(
56   party_proportions,
57   data.frame(n = n, group = "All", party = names(prop_party), prop = as.numeric
58     (prop_party)),
59   data.frame(n = n, group = "Treat", party = names(prop_party), prop = as.numeric
60     (prop_party_treat)),
61   data.frame(n = n, group = "Control", party = names(prop_party), prop = as.numeric
62     (prop_party_control))
63 )
64
65 print(party_proportions)
66
67 # step 5: join sample proportions with population proportions
68 party_prop_full <- party_proportions %>%
69   left_join(
70     data.frame(
71       party = names(prop_party),
72       prop_party_Pop = as.numeric(prop_party_Pop)
73     ),
74     by = "party"
75   ) %>%
76   mutate(imbalance = prop - prop_party_Pop)
77
78 head(party_prop_full)
79
80 # simple wide table
81 party_wide <- party_prop_full %>%
82   select(-imbalance) %>%
83   pivot_wider(names_from = group, values_from = prop) %>%
84   rename(prop_party_Pop = prop_party_Pop) %>%
85   arrange(n, party)
86
87 head(party_wide)
88
89 # plot 1: groups approach population line as n grows
90 party_proportions_baseline <- tibble(
91   Party=c("Democrat", "Republican", "Independent"),
92   pop_prop=c(0.4, 0.4, 0.2)
93 )
94
95 ggplot(party_prop_full, aes(x = n, y = prop, color = group)) +
96   geom_point() + geom_line() +
97   geom_hline(data = party_proportions_baseline, aes(yintercept = pop_prop), linetype
98     = 2) +
99   facet_wrap(~ party, nrow = 2) +
100   labs(x = "Sample size (n)", y = "Proportion",
101     title = "Treatment & Control proportions approach population as n increases",
102     subtitle = "Dashed line = population proportion per category") +
103   theme_minimal()
104
105 # plot 2: imbalance across parties shrinks with n
106 imbalance <- party_prop_full %>%
107   filter(group %in% c("Treat", "Control")) %>%
108   group_by(n, party, group) %>%
109   summarise(prop = mean(prop), .groups = "drop") %>%
110   pivot_wider(names_from = group, values_from = prop) %>%
111   mutate(abs_diff = abs(Treat - Control)) %>%
112   group_by(n) %>%
113   summarise(
114     max_abs_diff = max(abs_diff),
115     ll_sum_diff = sum(abs_diff),
116     .groups = "drop"
117   )

```

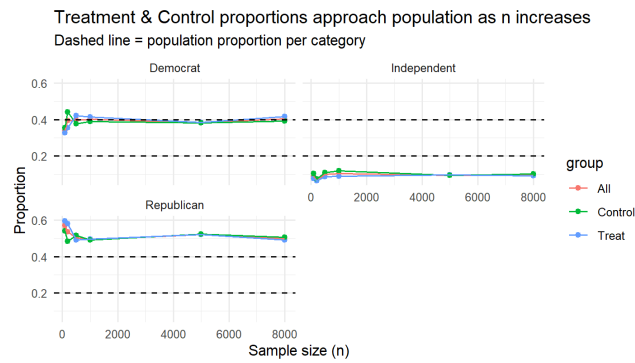


Figure 1: Treatment & Control proportions approach population as n increases

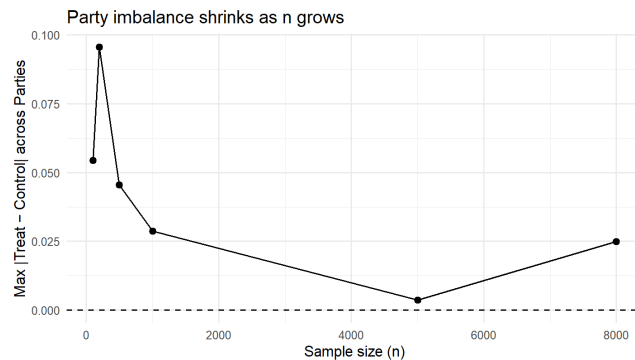


Figure 2: Party imbalance shrinks as n grows

```

114 )
115
116 ggplot(imbalance, aes(x = n, y = max_abs_diff)) +
117   geom_hline(yintercept = 0, linetype = 2) +
118   geom_point(size = 2) + geom_line() +
119   labs(x = "Sample size (n)", y = "Max |Treat - Control| across Parties",
120        title = "Party imbalance shrinks as n grows") +
121   theme_minimal()
122
123 table(party_prop_full$n)

```

2 Data Analysis

Question 1

```

1 # "message" is the treatment
2
3 voting <- read.csv("C:/WindowsD/TAM/25fall/Quant/Assignment/0930/voting.csv")
4 summary(voting$message)
5 # "message" is a continuous variable; its data type is "character"

```

Question 2

```

1 str(voting$message)
2 voting$treatment <- ifelse(voting$message == "no", 0, 1)

```

Question 3

```

1 library(dplyr)
2 mean(voting)

```

```

3 voting %>%
4   filter(treatment == 1) %>%
5   summarise(mean_voted1 = mean(voted, na.rm = TRUE)) %>%
6   print()
7   # mean_voted1: 0.3779482.

```

Among those who received the treatment, approximately 37.8% voted.

```

1 voting %>%
2   filter(treatment == 0) %>%
3   summarise(mean_voted0 = mean(voted, na.rm = TRUE)) %>%
4   print()
5   # mean_voted0: 0.2966383

```

Among those who received no treatment, approximately 29.7% voted.

Question 4

```

1 voting_treated <- voting %>%
2   filter(treatment == 1)
3 voting_nottreated <- voting %>%
4   filter(treatment == 0)

```

Question 5

```

1 voting %>%
2   filter(treatment == 0) %>%
3   summarise(mean_birth0 = mean(birth, na.rm = TRUE)) %>%
4   print()

```

The average birth year for the untreated group is 1956.186.

```

1 voting %>%
2   filter(treatment == 1) %>%
3   summarise(mean_birth1 = mean(birth, na.rm = TRUE)) %>%
4   print()

```

The average birth year for the treated group is 1956.147.

Question 6

```

1 mean_voted1 <- 0.3779482
2 mean_voted0 <- 0.2966383
3 average_causal_effect <- mean_voted1 - mean_voted0
4 print(average_causal_effect)
5 #average_causal_effect: 0.0813099

```

The average causal effect is 0.0813099. The average causal effect of **message** on **voted** is the average change (0.0813099) in **voted** caused by a one-unit increase in **message** for a group of individuals.

Question 7

The more messages that pressured people to vote by promising to tell their neighbors if they voted in the upcoming election, the more likely people were to vote.