# Midterm Election Simulation Project

Group 2

Erica Nam, Xiao Xu, Xinyu Zhang, Yitao Li

For each of *nrep* replications:

- We first generate *N_contam* the size of the contaminating group using a bin($N$, $q$). $q$ is our assumption for the percentage of the contaminating group out of the entire sample which has a sample size $N$. $N$ - *N_contam* is the size of the sample that is not contaminated. Note that when $q$ = 0, there is no contamination in the sample.

- Second, we generate the number of people in the contaminating group who prefer candidate 1 ($N\_true_1$), candidate 2 ($N\_true_2$) and candidate 3 ($N\_true_3$) respectively using a multinomial ($N\_contam$, $p_1$, $p_2$, $p_3$) since the proportion of people in this group favors candidate 1, candidate 2 and candidate 3 is $p_1$, $p_2$, $p_3$.

- We also generate the number of people in the not contaminated group who prefer candidate 1 ($N\_contam_1$), candidate 2 ($N\_contam_2$) and candidate 3 ($N\_contam_3$) respectively using a multinomial ($N$ - $N\_contam$, $r_1$, $r_2$, $r_3$) since the proportion of people in this group favors candidate 1, candidate 2 and candidate 3 is $r_1$, $r_2$, $r_3$.

- Then we have the simulated number of people who prefer candidate 1 $N_1$ = $N\_true_1$ + $N\_contam_1$. People who prefer candidate 2 $N_2$ = $N\_true_2$ + $N\_contam_2$. People who prefer candidate 3 $N_3$ = $N\_true_3$ + $N\_contam_3$.

Continue for the same replication $t$:

- We first calculate $p_1\_hat = \frac{N_1}{N}$. $p_2\_hat = \frac{N_2}{N}$. $p_3\_hat = \frac{N_3}{N}$.

- Then we calculate $CI\_m1_{i,j}$ the 1- $\alpha$ confidence interval for $p_i$ - $p_j$ using the following formula from method 1:

$$p_i\_hat - p_j\_hat \pm \sqrt{\frac{Ad_{i,j}}{N}} \text{ where } A = \chi^2_{M-1}\left(\frac{\alpha}{M}\right), \; M = m\frac{(m-1)}{2}, \; m = 3, d_{i,j} = p_i\_hat + p_j\_hat - (p_i\_hat - p_j\_hat)^2$$

Calculate $CI\_m2_{i,j}$ the confidence interval for $p_i$ - $p_j$ using the following formula from method 2:

$$p_i\_hat - p_j\_hat \pm \frac{a}{\sqrt{N}} \text{ where} 1 - 2[1 - z_a] - 4[m-2][1 - z_{a\sqrt{2}}] = 1 - \alpha, \; m = 3$$

Note that if $p_1 - p_2 \in CI\_m1_{1,2}$ and $p_1 - p_3 \in CI\_m1_{1,3}$ and $p_2 - p_3 \in CI\_m1_{2,3}$, then we say the simultaneous confidence interval consist of $CI\_m1_{1,2}$, $CI\_m1_{1,3}$ and $CI\_m1_{2,3}$ covers the true value $p_1 - p_2$, $p_1 - p_3$ and $p_2 - p_3$ simultaneously for this replication. Same goes for method 2.

- The average confidence interval width is $avg\_width\_m1_t = \dfrac{2\sqrt{\frac{Ad_{1,2}}{N'}} + 2\sqrt{\frac{Ad_{1,3}}{N'}} + 2\sqrt{\frac{Ad_{2,3}}{N'}}}{3}$ for method 1 and $avg\_width\_m2_t = 2\frac{a}{\sqrt{N}}$ for method 2.

- The maximum confidence interval width is $max\_width\_m1_t = \max[\, 2\sqrt{\frac{Ad_{1,2}}{N'}} , 2\sqrt{\frac{Ad_{1,3}}{N'}} , 2\sqrt{\frac{Ad_{2,3}}{N'}} \,]$ for method 1 and $max\_width\_m2_t = 2\frac{a}{\sqrt{N}}$ for method 2.

- The CI and its width is calculated *nrep* times for *nrep* replications.

- The estimated coverage probability is $\frac{nrep\_cov\_m1}{nrep}$ if among *nrep* replications *nrep_cov_m1* of them cover the true value simultaneously with the CI constructed with method 1. Same foes for method 2.

- The estimated average simultaneous CI width is $\frac{\sum_t avg\_width\_m1_t}{nrep}$ for method 1. Same goes for method 2.
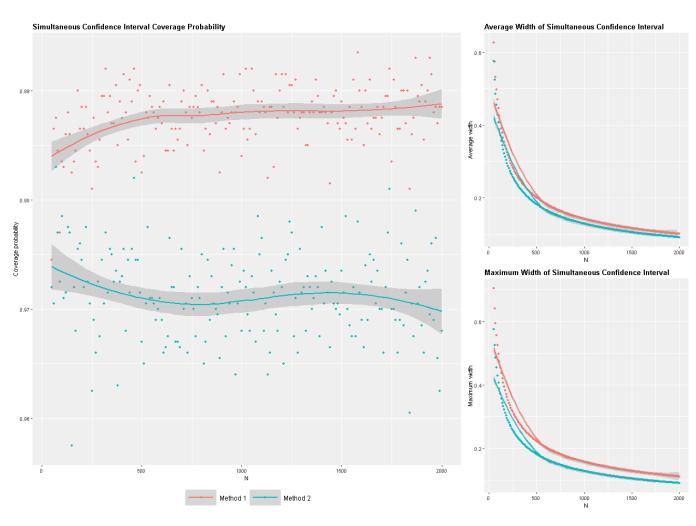
- The estimated maximum simultaneous CI width is $\frac{\sum_t max\_width\_m1_t}{nrep}$ for method 1. Same goes for method 2.

Shiny app:

https://zxynj.shinyapps.io/Midterm_election_simulation_project_Stat_6341/

- When there is no contamination, method 1 is better than method 2 but with wider CI.
- CI width decreases as $N$ increase.

# Shiny app interactive plot

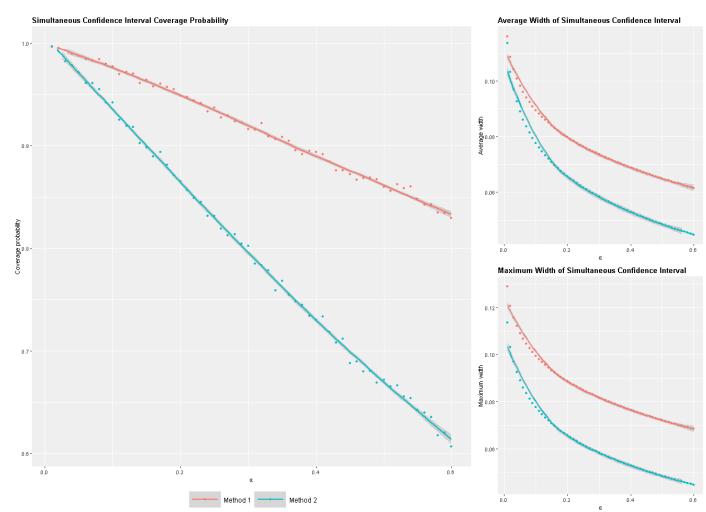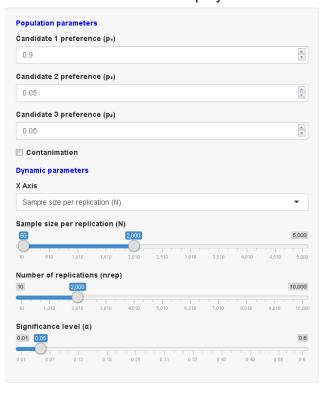- The variance of the coverage probability decreases as *nrep* increase.

- The coverage probability decreases as $\alpha$ increase.
- The coverage probability for Method 2 drops faster than Method 1. This could be explained by Method 2's narrower CI width and faster deacreasing CI width.
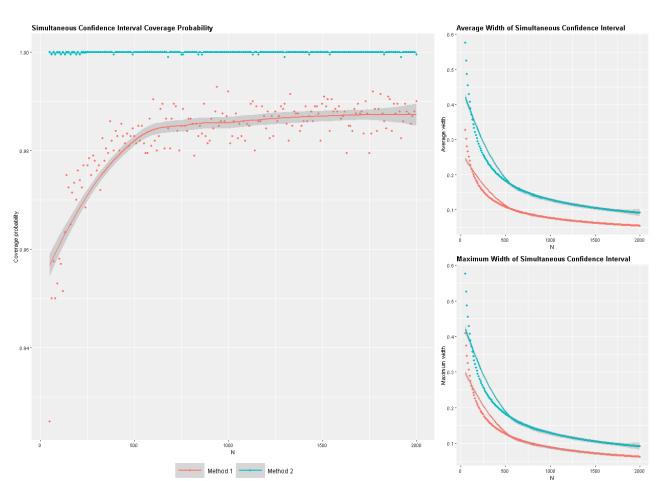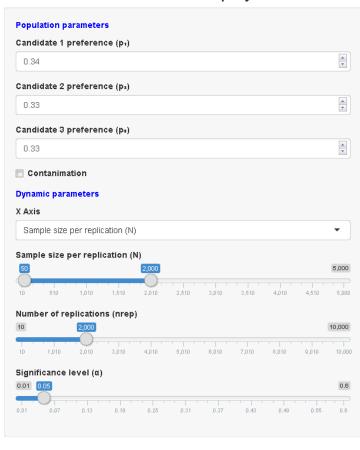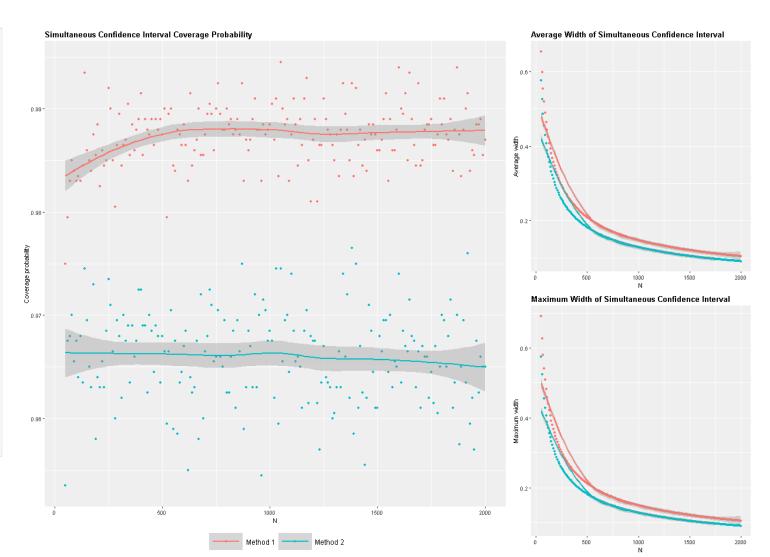
- When one candidate dominates the race, the CI width for Method 1 → 0 which cause its coverage probability to perform poorly comparing to Method 1 which has a fixed CI width.
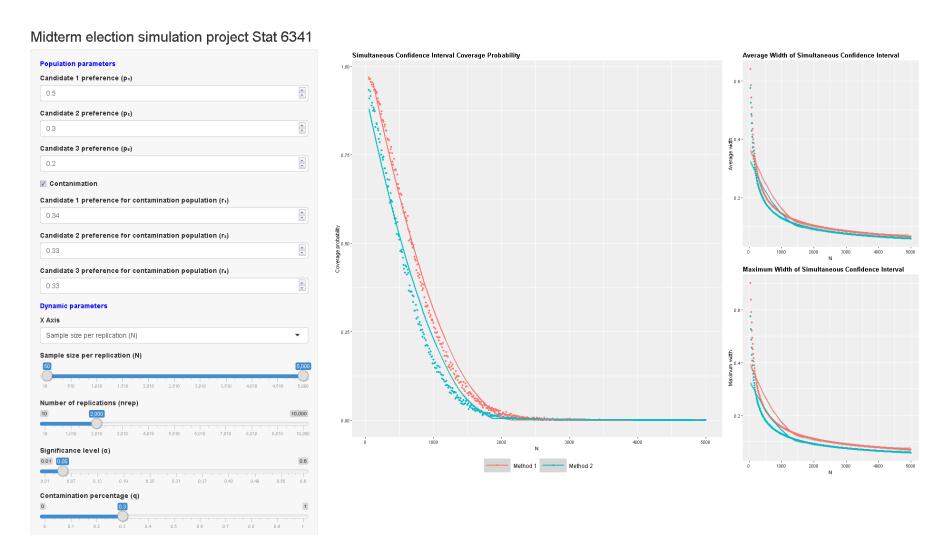- $d_{i,j} = p_i\_hat + p_j\_hat - (p_i\_hat - p_j\_hat)^2$
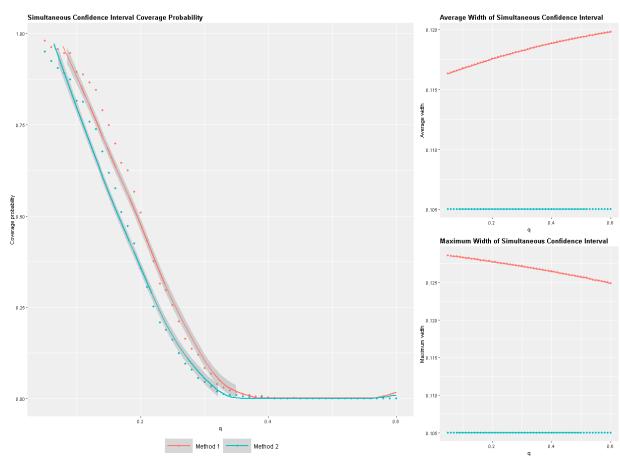
- When all three candidates are very close.

- When there is contamination, the CI width was wide enough at the beginning to achieve 90%+ coverage probability.
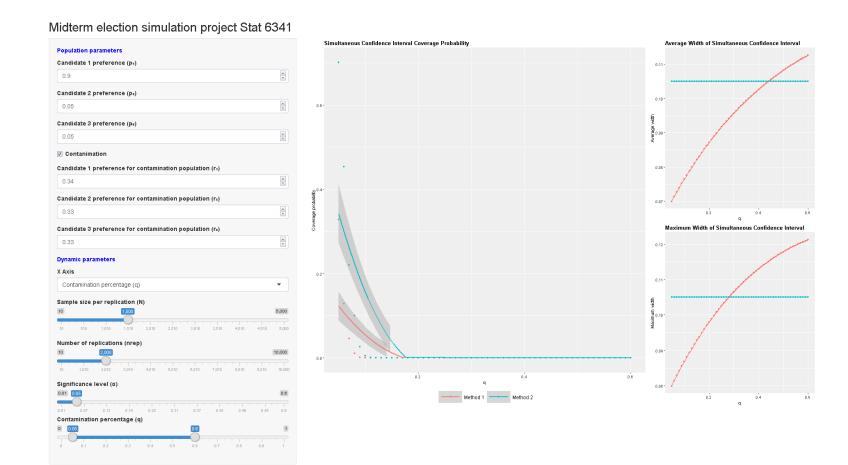
- When there is contamination, the coverage probability drops as contamination proportion increases.
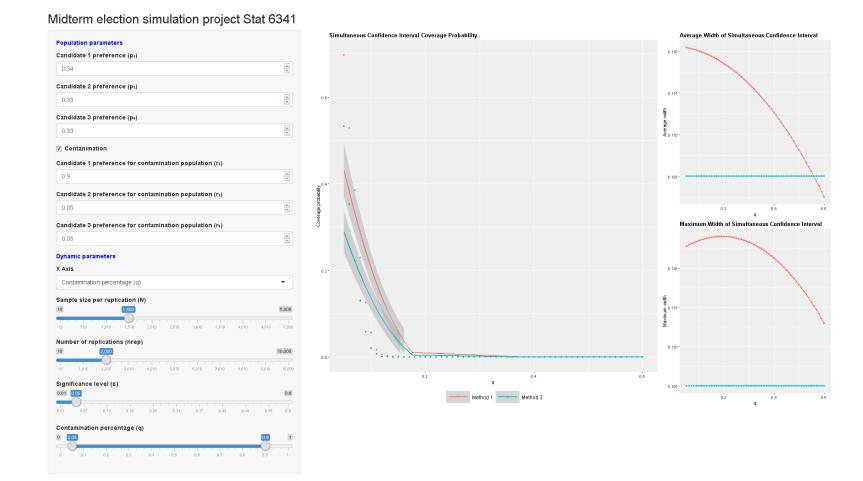
- When single domination is contaminated with close preference, the CI width for method 1 is smaller than Method 2 as expected.
- As the sample gets more contaminated, the CI center shifts towards the contaminating group's preference probability (close preference) which has a larger CI.
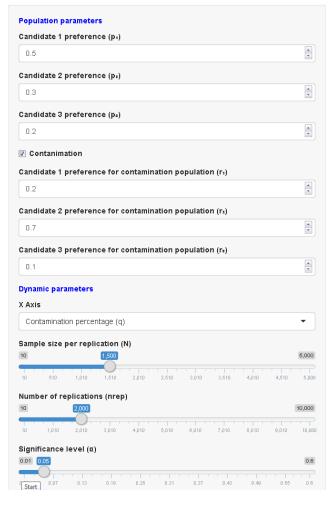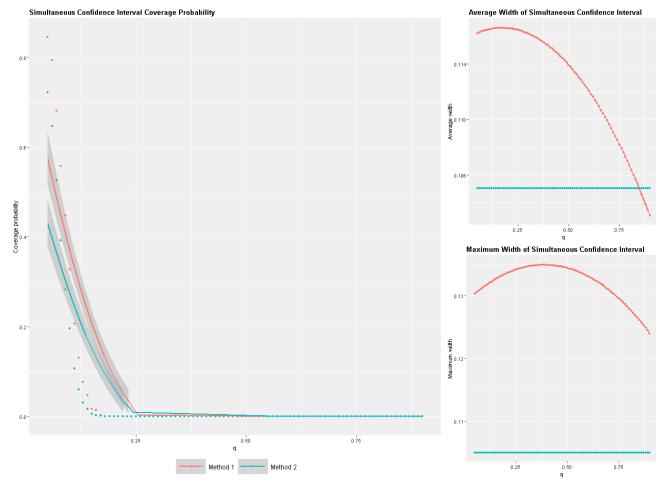
- When close preference is contaminated with single domination, as the sample gets more contaminated the CI center shifts towards the contaminating group's preference probability (single domination) which has a smaller CI.
- Even the overall trend of maximum CI with is decreasing, it actually goes up a little bit as the mixed $p_1$, $p_2$, $p_3$ becomes a little bit different.

- Compare $p_1$, $p_2$, $p_3$ mixed with $r_1$, $r_2$, $r_3$ and $p_1$, $p_2$, $p_3$ mixed with $r_1$, $r_3$, $r_2$.
- Coverage probability becomes 0 at q = 0.25 and average CI width has a small ascent right after mixing.

- Compare $p_1$, $p_2$, $p_3$ mixed with $r_1$, $r_2$, $r_3$ and $p_1$, $p_2$, $p_3$ mixed with $r_1$, $r_3$, $r_2$.
- Coverage probability becomes 0 at q =0.35.



Midterm election simulation project Stat 6341