# Coding Tests

Wenzhi Wang

October 26, 2024

A Github repository (link here) is created to store codes and summary report for the coding tasks.

## 1. Question 1: Simulations

To present the results in a more simplified and unified way, from questions 1.1. to questions 1.4., I create a Python class that is able to conduct all the procedures required in the task. The complete codes are presented in Appendix Figure A.1, and are stored in the "codes/Question1.py" script file.

The `ExpDistSimulation` class has three key attributes: `lambda_value`, `n_sample`, and `n_simulation`, which represents the distribution parameter of the exponential distribution, the sample size in each simulation, and the number of simulations, respectively. It has the following key methods and perform procedures required in the coding tasks:

- `gen_exp_sample` generates a random draw from an exponential distribution (with default randomization seed 1234).

- `cal_sample_mean` calculates the (unadjusted) sample mean from any random draw.

- `adj_sample_mean` calculates the adjusted sample mean from any random draw.

- `simulation` conducts simulation, and for each simulated sample, it calculates the adjusted sample mean and stores it.

- `simulation_plot` plots the empirical distribution of the simulated adjusted sample mean.

## 1.1. Question 1.1.

First, I initialize different object instances with different sample size parameters. Then, I generate a random sample for each sample size. The codes are presented in Figure 1, with one example output shown in Figure 2 (*I don't show all the random sample, as they actually don't matter so much for subsequent key results but can greatly reduce the readability of this report.*).

Figure 1: Codes for Question 1.1.

```
1   # ??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??
2   # ?? question 1.1.
3   # ??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??
4
5   ExpDistSimulation_5 = ExpDistSimulation(n_sample=5, lambda_value=0.125)
6   ExpDistSimulation_20 = ExpDistSimulation(n_sample=20, lambda_value=0.125)
7   ExpDistSimulation_50 = ExpDistSimulation(n_sample=50, lambda_value=0.125)
8   ExpDistSimulation_100 = ExpDistSimulation(n_sample=100, lambda_value=0.125)
9   ExpDistSimulation_500 = ExpDistSimulation(n_sample=500, lambda_value=0.125)
10
11  sample_5 = ExpDistSimulation_5.gen_exp_sample()
12  sample_20 = ExpDistSimulation_20.gen_exp_sample()
13  sample_50 = ExpDistSimulation_50.gen_exp_sample()
14  sample_100 = ExpDistSimulation_100.gen_exp_sample()
15  sample_500 = ExpDistSimulation_500.gen_exp_sample()
```
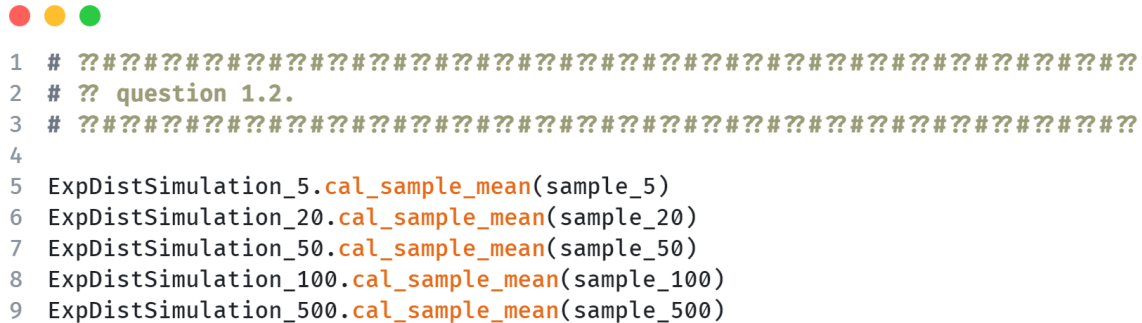
Figure 2: Output for Question 1.1.

```
1  sample_5 = ExpDistSimulation_5.gen_exp_sample()
2  sample_5
✓ 0.0s
array([12.18984195,  5.71720643, 14.50340573,  6.50334787,  7.85841657])
```

## 1.2. Question 1.2.

The codes for answering this question are in Figure 3. Again, for the $n = 5$ case, the raw sample mean is $9.354443709511397$.
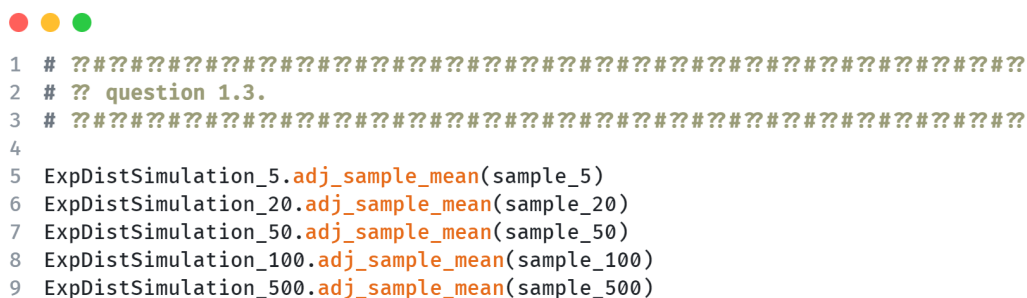
Figure 3: Codes for Question 1.2.

```
1  # ??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??
2  # ?? question 1.2.
3  # ??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??
4
5  ExpDistSimulation_5.cal_sample_mean(sample_5)
6  ExpDistSimulation_20.cal_sample_mean(sample_20)
7  ExpDistSimulation_50.cal_sample_mean(sample_50)
8  ExpDistSimulation_100.cal_sample_mean(sample_100)
9  ExpDistSimulation_500.cal_sample_mean(sample_500)
```

## 1.3. Question 1.3.

The codes for answering this question are in Figure 4. Again, for the $n = 5$ case, the adjusted sample mean is $3.0286282061644623$.

Figure 4: Codes for Question 1.3.

```
1  # ??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??
2  # ?? question 1.3.
3  # ??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??#??
4
5  ExpDistSimulation_5.adj_sample_mean(sample_5)
6  ExpDistSimulation_20.adj_sample_mean(sample_20)
7  ExpDistSimulation_50.adj_sample_mean(sample_50)
8  ExpDistSimulation_100.adj_sample_mean(sample_100)
9  ExpDistSimulation_500.adj_sample_mean(sample_500)
```
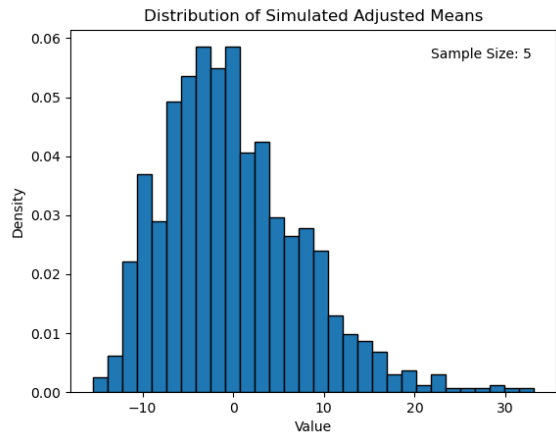
## 1.4. Question 1.4.

Next, I repeat the procedures in question 1.1.-1.3. for 1000 times, and plots the empirical distribution of the simulated adjusted means.
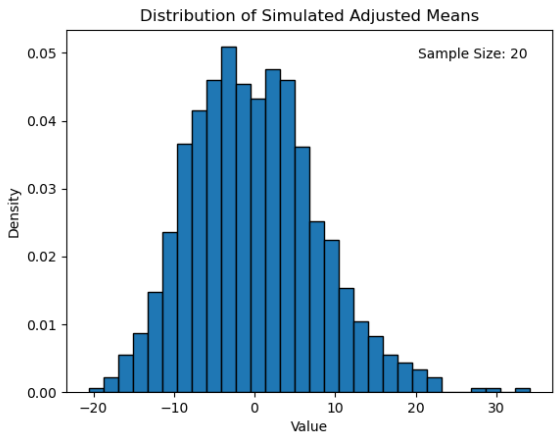
The results are in Figure 5.

Comments: First, we should expect these adjusted means across different simulated samples should have a normal distribution if the sample size in each simulation sample goes to infinity, as the Central Limit Theorem suggests. It is clear that the empirical distributions become more similar to a normal distribution, as $n$, the sample size in each simulation increases. For example, the $n = 5$ case is obviously left-skewed and have a fat right tail, but these properties almost vanish for the $n = 500$ case. This is because the Central Limit Theorem states the asymptotic properties of our statistic of interest, and the asymptotic distribution holds only if we draw enough sample points from the distribution.

Figure 5: Distribution of the simulated adjusted means with different sample sizes
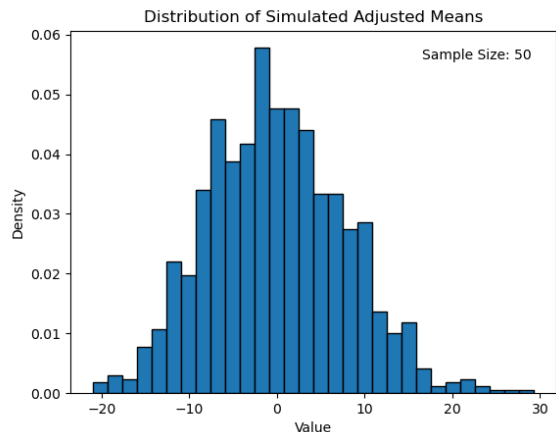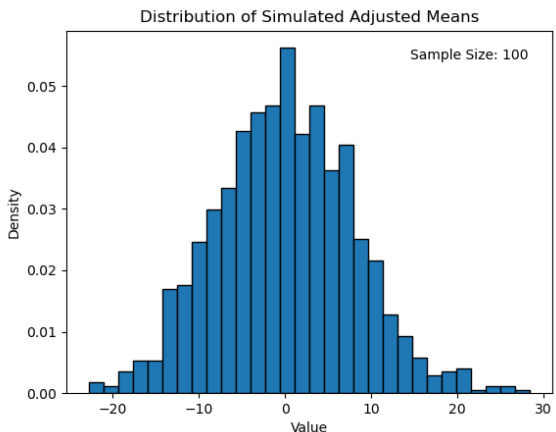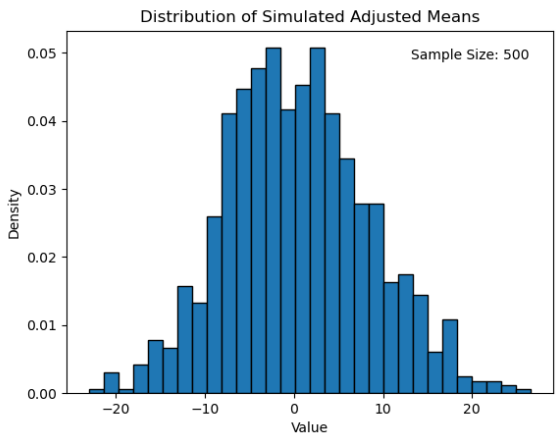
(a) $n = 5$

(b) $n = 20$



(c) $n = 50$

(d) $n = 100$



(e) $n = 500$

## 1.5. Question 1.5.

I guess this question is about MLE, thus the moment condition shouldn't be

$$5000 \cdot \log\left(\lambda\right) - \lambda \sum_{i=1}^{5000} x_i = 0.$$

*I use the following optimization problem for estimating the $\lambda$:*

$$\max_{\lambda} \left[ 5000 \cdot \log\left(\lambda\right) - \lambda \sum_{i=1}^{5000} x_i \right].$$

In particular, I first create a sample with 5000 random draws, then define the objective function (which is the negative log-likelihood, since I use a minimization routine), as shown in Figure 6.

Figure 6: Codes for Question 1.5.

```python
# ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ??
# ?? question 1.5.
# ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ?? # ??

sample_for_est = ExpDistSimulation(n_sample=5000).gen_exp_sample()


def obj_func(lambda_value: float, sample: NDArray) → float:
    moment = 5000 * np.log(lambda_value) - lambda_value * np.sum(sample)
    obj = -moment
    return obj


bounds = [(1e-6, 100)]
initial_guess = 0.2
result = minimize(
    obj_func, x0=initial_guess, args=(sample_for_est,), bounds=bounds
)
print(
    f"The lambda value that maximizes the objective function is: {result.x}."
)
```

The results (togher with the codes screenshot) are in Figure 7. The estimated $\lambda$, $\widehat{\lambda}$ is 0.123.

Figure 7: Answer for Question 1.5.

```python
print(
    f"The lambda value that maximizes the objective function is: {result.x[0]}."
)
```
✓ 0.0s

The lambda value that maximizes the objective function is: 0.12307785112006046.

## 1.6. Question 1.6.

When restricting the search over intervals $[0.5, 1]$ and $[0.2, 1]$, the optimizer value is $0.5$, and $0.2$, respectively. These boundary values suggest that the maximization solution in these two cases fail to reach the true maximization (which can be calculated by the first-order condition). Therefore, we need to extend these narrow bound conditions when conducting the numerical maximization.

But when the restriction is $[0.1, 1]$, we can reach the true optimizer: $0.123$, since we have extended our bound condition to a broader interval.
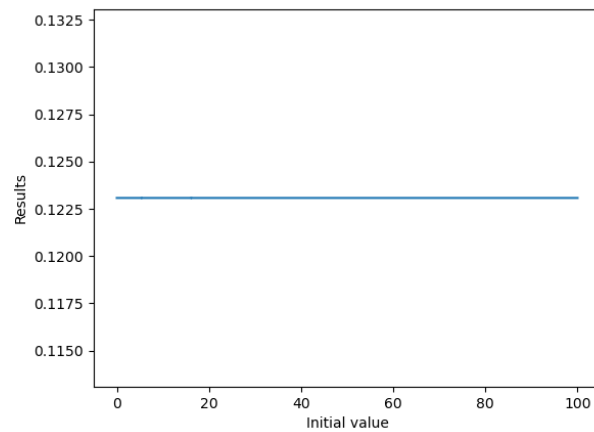
## 1.7. Question 1.7.

To answer this question, I create a equi-distanced grid over $[10^{-6}, 100]$ as different initial values to start the search, and then for each initial value, I conduct the optimization routine, and stores the results in an array, and plots the results against different starting values. The codes are presented in Figure 8, and the results are in Figure 9.

Figure 8: Codes for Question 1.7.

```python
bounds = [(1e-6, 100)]

initial_guess_array = np.linspace(1e-6, 100, 10000)
results = np.zeros(shape=(10000))

i = 0
for guess in initial_guess_array:
    result_guess = minimize(
        obj_func, x0=guess, args=(sample_for_est,), bounds=bounds
    )
    lambda_estimated = result_guess.x[0]
    # print(result_guess.x[0])
    results[i] = lambda_estimated
    i = i + 1

plt.plot(initial_guess_array, results)
mean_result = np.mean(results)
plt.ylim(mean_result - 0.01, mean_result + 0.01)
plt.xlabel("Initial value")
plt.ylabel("Results")

final_file_path = os.path.join(results_path, "initial_guess_plot.png")
plt.savefig(final_file_path)
```
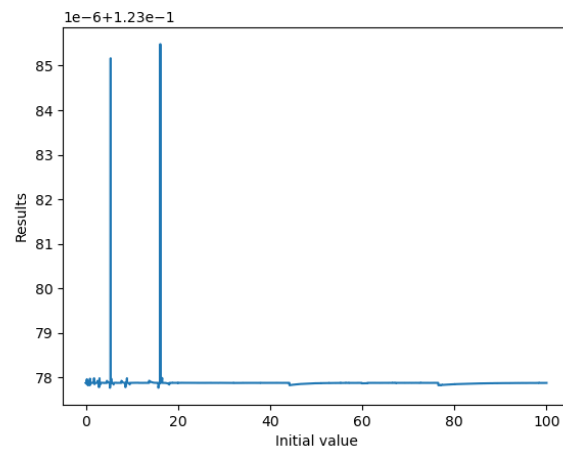
Figure 9: Answer for Question 1.7.



It is easy to see that the starting guess values doesn't matter for this MLE problem.

To see the tiny difference between different starting points, we can rescale the y-axis and have the following graph - Figure 10. Notice that the scale of the y-axis is $10^{-6}$!

Figure 10: Answer for Question 1.7. - Rescaling the y-axis

# 2. Estimation

Appendix Figure A.2 stores some Stata commands for initial setup, such as codes for defining several path macros, and specifying the version. The corresponding do file is "codes/Question2.do".

## 2.2. Dataset Construction

The Stata codes for this question are in Figure 11. The required moments are

- Mean: 311065.

- Standard deviation: 3193502.

- P10: 2.597.

- P50: 1161.257.

- P90: 2346.26.

- Number of observations: 28803.

Figure 11: Codes for Question 2.2.

```
1  *-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?
2  *-? question 2.2. - i.
3  *-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?
4
5  import delimited "${data}/trade.csv", clear
6
7  collapse (sum) trade, by(origin destination year)
8  sort origin destination year
9
10 *-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?
11 *-? question 2.2. - ii.
12 *-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?
13
14 preserve
15     import delimited "${data}/gravity.csv", clear
16     save "${data}/temp_gravity.dta", replace
17 restore
18
19 merge 1:1 origin destination year using "${data}/temp_gravity.dta", keep(match) nogenerate
20
21 *-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?
22 *-? question 2.2. - iii.
23 *-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?*-?
24
25 summarize trade if year==2015, detail
```
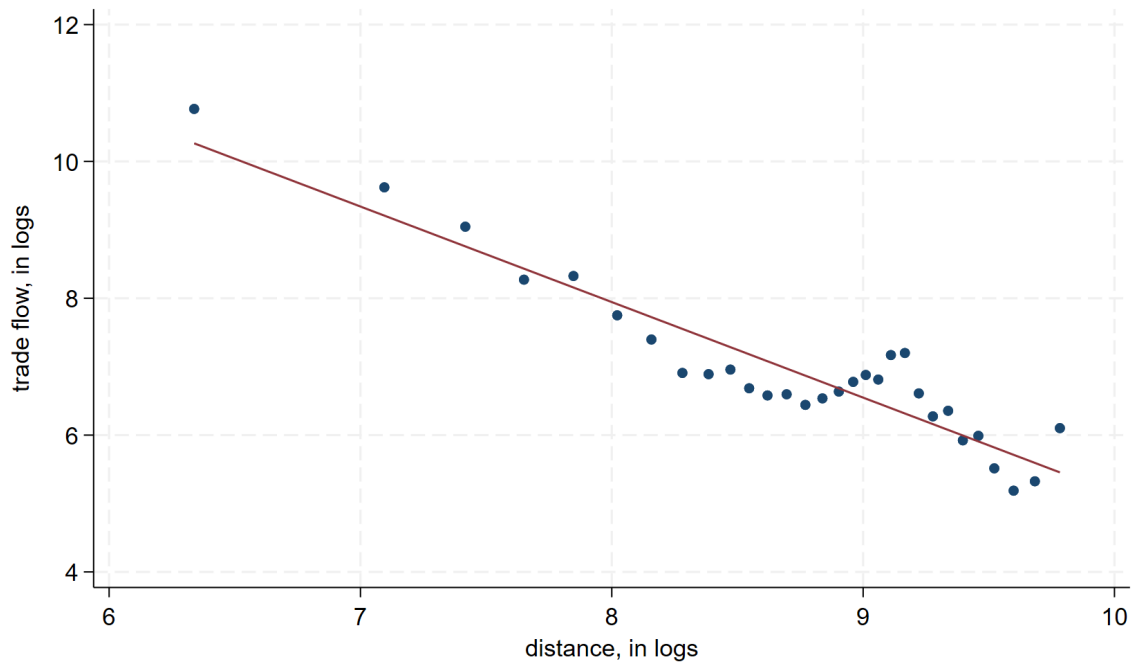
## 2.3. Estimation

### 2.3.1. The Binscatter Plot

The binscatter plot is in Figure 12. The raw correlation between log(trade flows) and log(distance) is $-0.2608$. The linear fit (with a constant) of the regression coefficient on log(distance) is $-1.396478$.
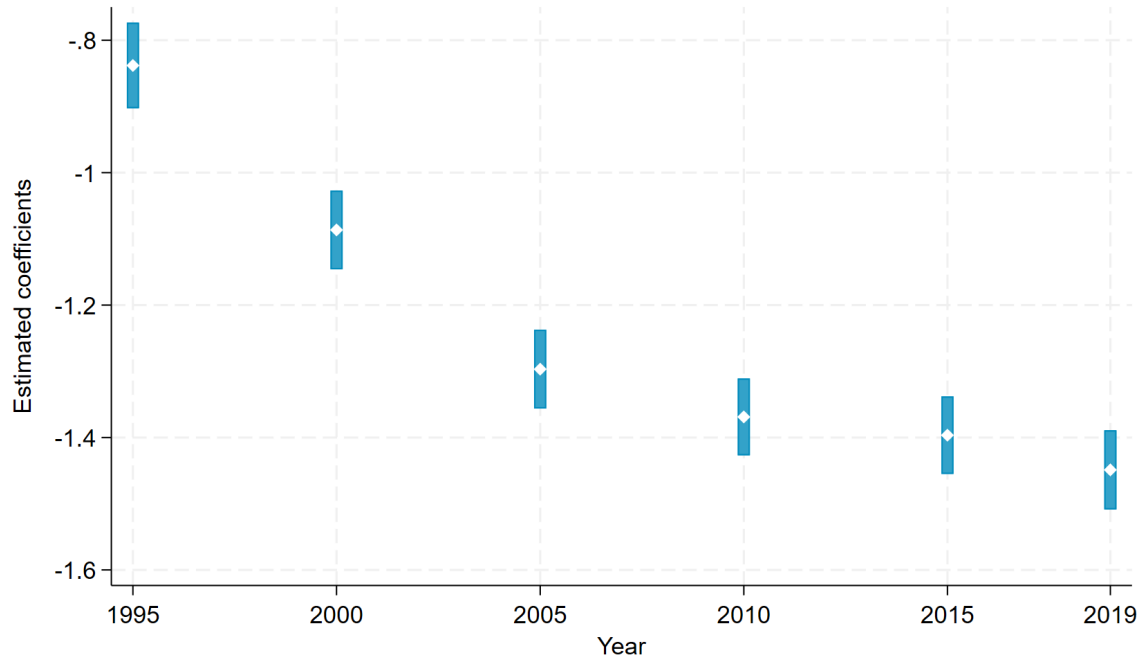
Figure 12: Binscatter plot of trade flows against distance

## 2.3.2. Distance Coefficients Across Years

The results are in Figure 13. In the graph, I report the estimated coefficients using sample from different years, and the 95% confidence intervals, which are calculated based on clustered standard errors (two-way clustering on both the origin and the destination countries).
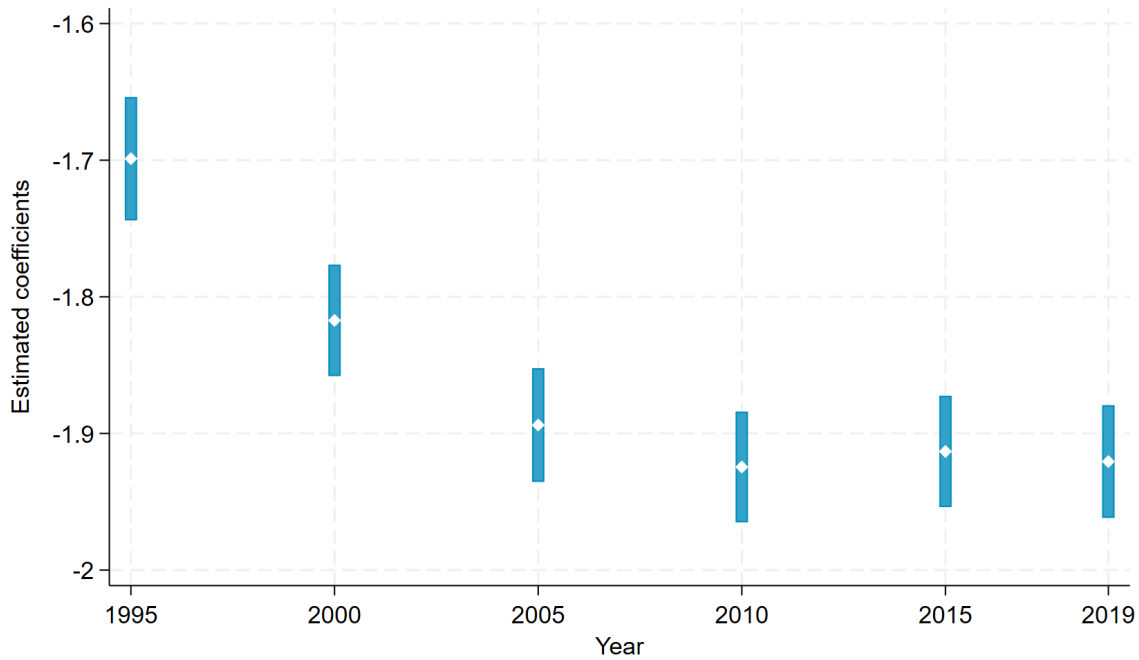
Figure 13: Coefficients on log(distance) for different years

### 2.3.3. Distance Coefficients Across Years With Importer and Exporter FEs

The results are in Figure 14, where the estimated coefficients using sample from different years, and the 95% confidence intervals, which are calculated based on robust standard errors are reported.

Figure 14: Coefficients on log(distance) for different years, with importer-year and exporter-year FEs



Discussion:

In Figure 14, where importer-year and exporter-year FEs are controlled, we are essentially reporting how much percentage of trade flows will decrease between two given countries in a given year, if their distance has increased by 1%. Controlling for importer and exported fixed effects allow us to partial out the effects of some country attributes, such as economy size (e.g., as measured by GDP) from the estimated coefficients. However, in Figure 13, where the FEs are not controlled, the variation used to identify the distance coefficient could also come from different importer-exporter pair with different distance.

### 2.3.4. Adding Other Controls

The required table is Table 1, where the robust standard errors are reported in parentheses.

Table 1

|  | Trade flows |
|---|---|
|  | (1) |
| distance, in logs | -1.588*** |
|  | (0.02) |
| sharing a border | 1.100*** |
|  | (0.10) |
| sharing a language | 1.008*** |
|  | (0.05) |
| having colonial ties | 0.558*** |
|  | (0.12) |
| having a trade agreement | 0.555*** |
|  | (0.04) |
| R-squared | 0.724 |
| Obs | 26938 |

Notes. Sample includes trade flows between different countries in 2015. Importer and exporter fixed effects are controlled. Robust standard error are reported.

Discussion:

In the table, the coefficient on distance is essentially reporting how much percentage of trade flows will decrease between two given countries who don't share a border, or a language, or colonial ties, and don't have a trade agreement in a given year, if their distance has increased by 1%. Compared with Figure 14, the effect of distance on the trade flows between two countries decreases when we add more controls that represent other notions of "distance" between two countries. The decrease is most likely due to the fact that these controls (especially "sharing a border") are negatively correlated with the geographical distance measure.

All other controls measure the easiness for trades to happen between different countries. In particular, "sharing a border" measures geographical distance, "sharing a language" and "sharing colonial

ties" measures some type of cultural linkages, and "having a trade agreement" measures some institutional costs for international trade. For instance, if two countries share the same language, then it is easier for importer firms to communicate with the exporter firms – they are more likely to match each, negotiate with each other and reach a business contract. Therefore, it is intuitive that these factors will contribute positively to trade flows between two countries.

# A. Appendix

Figure A.1

```python
1
2  class ExpDistSimulation:
3      """
4      This class conducts exercises required in question 1.1 - question 1.4.
5      """
6
7      def __init__(
8          self,
9          lambda_value: float = 1 / 8,
10         n_sample: int = 500,
11         n_simulation: int = 1000,
12     ):
13         """
14         Initialize an object instance with distributional and simulation pars.
15         """
16         self.n_sample = n_sample
17         self.n_simulation = n_simulation
18         self.lambda_value = lambda_value
19
20     def __repr__(self) -> str:
21         """
22         Print the distributional and simulation parameters.
23         """
24         return (
25             f"An exponential distribution simulation exercise with lambda = {self.lambda_value:.3f}, "
26             f"sample size = {self.n_sample:.0f}, "
27             f"and number of simulations = {self.n_simulation:.0f}.\n"
28         )
29
30     def gen_exp_sample(self, seed: int = 1234) -> NDArray:
31         """
32         This method generates a random draw from an exponential distribution.
33
34         Parameters:
35         self.lambda_value: the inverse of the expectation of the distribution
36         self.n_sample: sample size
37         seed: randomization seed
38
39         Returns:
40         sample: np.ndarray with shape (self.n_sample,)
41         """
42         scale_value = 1 / self.lambda_value
43         rng = np.random.default_rng(seed=seed)
44         sample = rng.exponential(scale=scale_value, size=self.n_sample)
45         return sample
46
47     def cal_sample_mean(self, array: NDArray) -> float:
48         """
49         This method calculates the sample mean from any sample array.
50
51         Parameters:
52         array: a np.ndarray with shape (self.n_sample,)
53         """
54         return array.mean()
55
56     def adj_sample_mean(self, array: NDArray) -> float:
57         """
58         This method calculates the adjusted sample mean from any sample array
59         drawn from an exponential distribution.
60
61         Parameters:
62         self.lambda_value: the inverse of the expectation of the distribution
63         self.n_sample: sample size
64         array: a np.ndarray with shape (self.n_sample,)
65         """
66         raw_sample_mean = self.cal_sample_mean(array)
67         expectation = 1 / self.lambda_value
68         res = np.sqrt(self.n_sample) * (raw_sample_mean - expectation)
69         return res
70
71     def simulation(self) -> NDArray:
72         """
73         This method conducts simulations. In particular, self.n_simulation
74         random samples will be generated using the self.gen_exp_sample()
75         method. Then all adjusted sample means calculated by the
76         self.adj_sample_mean() method will stored in a np.ndarray.
77
78         Parameters:
79         self.lambda_value: the inverse of the expectation of the distribution
80         self.n_sample: sample size
81         self.n_simulation: number of simulations
82
83         Returns:
84         res: a np.ndarray with shape (self.n_simulation,)
85         """
86         res = np.zeros(shape=(self.n_simulation,))
87         for i in range(self.n_simulation):
88             sample = self.gen_exp_sample(seed=i)
89             res_i = self.adj_sample_mean(array=sample)
90             res[i] = res_i
91         return res
92
93     def simulation_plot(self, file_name, **kwargs):
94         """
95         This method plots the empirical distribution of the simulated adjusted
96         sample mean (calculated using different self.n_sample), and saves the
97         figure into a png file.
98         """
99         simulated_means = self.simulation()
100        plt.hist(simulated_means, density=True, **kwargs)
101        plt.title("Distribution of Simulated Adjusted Means")
102        plt.xlabel("Value")
103        plt.ylabel("Density")
104        plt.text(
105            0.95,
106            0.95,
107            f"Sample Size: {self.n_sample}",
108            horizontalalignment="right",
109            verticalalignment="top",
110            transform=plt.gca().transAxes,
111        )
112        final_file_path = os.path.join(results_path, f"{file_name}.png")
113        plt.savefig(final_file_path)
114        plt.show()
115        plt.close()
116
```

Figure A.2

```stata
 1  *??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??
 2  *?? step 0. initial setup
 3  *??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??*??
 4
 5  clear all
 6  set more off
 7
 8  set maxvar 32767
 9  set varabbrev off
10
11  if  "`c(username)'" == "wang" {
12      global user "E:/RA/BoothTests"
13  }
14
15  cd "${user}"
16
17  global codes   "${user}/codes"
18  global data    "${user}/data"
19  global results "${user}/results"
20
21  version 17.0
```