

---

## EC 320: Introduction to Econometrics

### Problem Set 1

**Due: On canvas, January 17th at 5 pm**

**Learning Outcomes:**

- Understanding random variables and probabilities
- Understanding bias and variance of estimators
- Basic data summary and visualization

**Checklist Before Handing In:**

- Did you answer all questions?
- Did you answer all parts for each question?
- Were your answers too vague? If so, make them more precise to make sure they really answer the question being asked.

**Instructions:** You are encouraged to work with other students in the class, but you must provide original responses. To receive full credit, justify your answers and list your collaborators. For full credit on the computational exercises, include your code and output in addition to your answers. You will turn in digital copies of your responses on Canvas. Please note the list of acceptable file types on the submission page. Finally, you can earn up to 5 extra credit points by typing your responses, or having a clear, neat, well-formatted assignment.

Name:

Collaborator 1:

Collaborator 2:

Collaborator 3:

---

## Analytical Questions

1. Consider an experiment where you flip three fair coins, record whether each coin landed heads or tails, and assign heads a value of 2 and tails a value of -1. You then multiply the values from each flip and describe this product as the random variable  $X$ . For example, if the coins landed **heads, heads, tails**, then you would multiply  $2 \cdot 2 \cdot -1$  to find an outcome of  $X = -4$ .
  - (a) Is  $X$  a discrete or continuous random variable?
  - (b) Outline the probability distribution of  $X$  in a table.
  - (c) What is the population mean of  $X$ ? In other words, what is  $E(X)$ ?
  - (d) What is the population variance of  $X$ ?
2. Consider the distributions of three estimators of the population mean  $\mu$ :

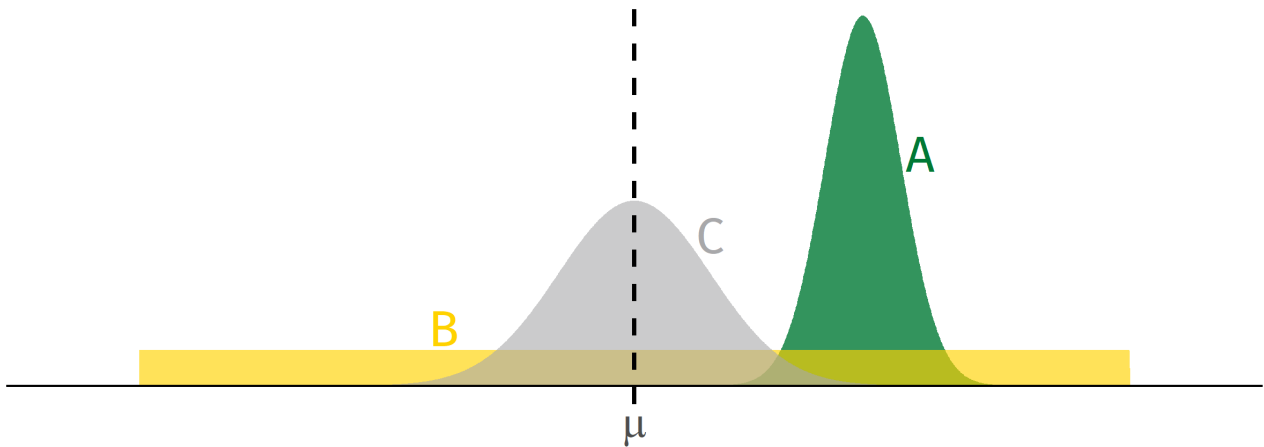


Figure 1: \*

**Note:** Shows the distributions of three estimators (A, B, and C) that provide estimates of the unknown population parameter  $\mu$ .  $E(A) = \mu + 3$ ,  $E(B) = \mu$ , and  $E(C) = \mu$ .

- (a) Which of the estimators above is unbiased?
  - (b) Which of the estimators above has the smallest variance?
  - (c) Which of the estimators above is the “best” unbiased estimator?
3. Use the rules of the expected values operator to **prove** that the population variance  $E[(X - E(X))^2]$  is equivalent to  $E(X^2) - [E(X)]^2$ .

---

<b>Student</b>	<b>Hours Slept</b>	<b>Exam Score</b>
Ahmed	8	95%
Chen	10	85%
Klara	7	70%
Peter	4	60%

4. Consider data gathered from students on the number of hours they slept the night before their midterm exam and their subsequent exam scores:
- (a) In the space provided below, draw a scatter plot with **Exam Score** on the  $Y$ -axis and **Hours Slept** on the  $X$ -axis. Do they appear positively or negatively correlated?
  - (b) What is the sample mean of **Exam Score**?
  - (c) What is the sample variance of **Hours Slept**?
  - (d) What is the sample correlation coefficient between **Exam Score** and **Hours Slept**?

---

## Computational Questions

For this portion of the problem set, you will use the file `election.2016.csv` in the `Problem_Set_1` folder on Canvas. The file contains county-level data on the 2016 presidential election from the MIT Election Data and Science Lab and the US Census Bureau. You will first need to download the data from Canvas and import it into R. Then you will need to load the `tidyverse` package.

Variable Name	Description
<code>stname</code>	Name of state
<code>fips</code>	County identifier (FIPS code)
<code>trump</code>	Number of votes cast for Donald Trump
<code>clinton</code>	Number of votes cast for Hillary Clinton
<code>totalvotes</code>	Total number of votes cast for any presidential candidate
<code>white</code>	Percentage of eligible voters who are non-Hispanic white
<code>poverty</code>	Percentage of residents below the poverty line
<code>population</code>	Number of residents

Use the data to complete the tasks and questions below.

1. Reduce your data set to include only those counties with a population of 50,000 or more.
2. Generate new variables that give the percentage of votes cast for each candidate using the `mutate` function. With these new variables, generate a third variable called `trump_margin` that gives Trump's "margin of victory" in each county. (For counties that Clinton won, the value of `trump_margin` will be negative.)
3. Produce summary statistics (min, max, mean, median, standard deviation, and the number of observations) for `trump_margin`, `white`, `poverty`, and `population` using the `summarize` function. Use your summary statistics to answer the following:
  - (a) What is Trump's average vote margin?
  - (b) What is the median poverty rate?
  - (c) How many people live in the least populous county? How many live in the most populous county?
4. Create a histogram of Trump's margin of victory. Which candidate won more counties? Is it necessarily the case that this candidate won more votes nationally? *Hint*: Look at a histogram of county population, too.
5. Create a scatter plot with `trump_margin` on the *Y*-axis and `poverty` on the *X*-axis. Does support for Trump appear positively or negatively correlated with poverty rates? Use the `cor` function to calculate the correlation coefficient. Interpret your result.

- 
6. Create a scatter plot with `trump_margin` on the *Y*-axis and `white` on the *X*-axis. Does support for Trump appear positively or negatively correlated with the percentage of white voters? Use the `cor` function to calculate the correlation coefficient. Interpret your result.
  7. Using `group_by` and `summarize`, aggregate the county-level data on `white`, `poverty`, and `population` to the state level. The first two variables are rates. To obtain accurate state-level rates, you need to account for differences in population across counties. You can do this by using the `weighted.mean` function inside `summarize`. For population, you can use the `sum` function. Your new `tibble` should have 50 rows.
  8. Using the state-level data, create a scatter plot with `poverty` on the *Y*-axis and `white` on the *X*-axis. To illustrate differences in population across states, vary the size of each point by population. You can do this by including `size = population` in the `aes` function of `ggplot`. Describe any patterns you see.