

HW2

Sam Olson

Outline

- Q1: SOLID
- Q2: SOLID
- Q3: SOLID
- Q4: Skeleton
- Q5: Skeleton
- Q6: Skeleton
- Q7: Skeleton

Problem 1 (20 pt)

A city has a total of 100,000 dwelling units, of which 35,000 are houses, 45,000 are apartments, and 20,000 are condominiums. A stratified sample of size $n = 1000$ is selected using proportional allocation (and rounding the sample sizes to the nearest integer). The three strata are houses ($h = 1$), apartments ($h = 2$), and condominiums ($h = 3$). The table below gives the estimates of the mean energy consumption per dwelling unit for the three strata and the corresponding standard errors.

Stratum (h)	Estimated Mean Energy Consumption (\bar{y}_h) (kWh per dwelling unit)	Estimated Standard Error ($\hat{SE}(\bar{y}_h)$)
House ($h = 1$)	915	4.84
Apartments ($h = 2$)	641	2.98
Condominium ($h = 3$)	712	7.00

1.

Estimate the total energy consumption for the full population of 100,000 dwelling units.

Answer

```
nh <- c(35000, 45000, 20000)
barY <- c(915, 641, 712)
tStr <- sum(nh * barY)
tStr
```

```
## [1] 75110000
```

$$\hat{T}_{str} = \sum_{h \in H} N_h \bar{y}_h = 915(35,000) + 641(45,000) + 712(20,000) = 75,110,000$$

Units of kWh per dwelling unit

2.

Estimate the standard error of the estimator used in (1).

Answer

```
nh2 <- nh^2
seBarY <- c(4.84, 2.98, 7.00)
seHatT <- sqrt(sum(nh2 * (seBarY^2)))
seHatT
```

```
## [1] 257447.4
```

$$SE(\hat{T}_{str}) = \sqrt{\text{Var} \left(\sum_{h \in H} N_h \bar{y}_h \right)} = \sqrt{\left(\sum_{h \in H} \text{Var} (N_h \bar{y}_h) \right)} = \sqrt{\left(\sum_{h \in H} N_h^2 \text{Var} (\bar{y}_h) \right)}$$

$$SE(\hat{T}_{str}) = \sqrt{(35,000^2)(4.84^2) + (45,000^2)(2.98^2) + (20,000^2)(7.00^2)} = 257,447.4$$

Units of kWh per dwelling unit

3.

What would be the sample size if the optimal allocation is to be used (under $n = 1000$) for this population? Assume that the survey costs are the same for each stratum.

Hint: Use the following steps:

a)

What is the sample size n_h for each stratum under proportional allocation?

Answer

$$\frac{n}{N} = \frac{1,000}{100,000} = 0.01$$

Under proportional allocation, we just take the above sampling rate multiplied by the population of the strata, i.e.:

$$n_h = N_h \times 0.01$$

Corresponding to

$$n_1 = 350$$

$$n_2 = 450$$

$$n_3 = 200$$

b)

Note that:

$$\hat{SE}(\bar{y}_h) = \sqrt{\frac{1}{n_h} \left(1 - \frac{n_h}{N_h}\right) s_h^2}$$

Thus, you can obtain s_h^2 .

Answer By definition:

$$SE(\bar{y}_h) = \sqrt{\left(\frac{1}{n_h} - \frac{1}{N_h}\right) S_h^2}$$

```
littlenH <- c(350, 450, 200)

# should be better with notation, but I love camel case
bigSh <- (seBarY^2)/(littlenH^-1 - nh^-1)
sqrt(bigSh)
```

```
## [1] 91.00427 63.53381 99.49367
```

We want the population-level standard errors by strata. To that end:

$$\sqrt{\left(\frac{1}{350} - \frac{1}{35000}\right)} S_1^2 = 4.84$$

$$\sqrt{\left(\frac{1}{450} - \frac{1}{45000}\right)} S_2^2 = 2.98$$

$$\sqrt{\left(\frac{1}{200} - \frac{1}{20000}\right)} S_3^2 = 7.00$$

Solving for S_1, S_2, S_3 respectively:

$$S_1 \approx 91.00, \quad S_2 \approx 63.53, \quad S_3 \approx 99.49$$

c)

Apply Neyman allocation (optimal allocation) using s_h in place of S_h .

Answer Combining everything from the above together:

$$n_h = \frac{N_h S_h}{\sum_{h=1}^H N_h S_h} n$$

where $n = 1000$.

```
sh <- c(91, 63.53, 99.49)

neyman <- round(((nh * sh) / sum(nh * sh)) * 1000,
               digits = 0)
neyman
```

```
## [1] 396 356 248
```

```
sum(neyman)
```

```
## [1] 1000
```

Explicitly,

$$n_1 = \frac{35,000(91.00)}{35,000(91.00) + 45,000(63.53) + 20,000(99.49)} \cdot 1000 \approx 396$$

$$n_2 = \frac{45,000(63.53)}{35,000(91.00) + 45,000(63.53) + 20,000(99.49)} \cdot 1000 \approx 356$$

$$n_3 = \frac{10,000(99.49)}{35,000(91.00) + 45,000(63.53) + 20,000(99.49)} \cdot 1000 \approx 248$$

$$n_1 \approx 396 \quad n_2 \approx 356, \quad n_3 \approx 248$$

Where we round to the nearest integer, and validate by checking the sum of samples is 1,000 as expected.

4.

What would be the estimated standard error of the total estimator under the optimal allocation in (3)? Compare it with the answer in (2). Which one is smaller?

Answer

$$SE(\hat{T}_{str}) = \sqrt{\sum_{h=1}^H \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) S_h^2}$$

```
firstTerm <- nh^2 / littlenH
secondTerm <- (1 - (littlenH/nh))
thirdTerm <- sh^2
sqrt(sum(firstTerm * secondTerm * thirdTerm))
```

```
## [1] 257435.2
```

Given:

$$N_1 = 35,000, \quad N_2 = 45,000, \quad N_3 = 20,000$$

$$n_1 = 396, \quad n_2 = 356, \quad n_3 = 248$$

$$S_1 \approx 91.00, \quad S_2 \approx 63.53, \quad S_3 \approx 99.49$$

We compute:

$$\frac{N_1^2}{n_1} \left(1 - \frac{n_1}{N_1}\right) S_1^2$$

$$\frac{(35,000)^2}{396} \left(1 - \frac{396}{35,000}\right) (91.00)^2$$

$$\frac{N_2^2}{n_2} \left(1 - \frac{n_2}{N_2}\right) S_2^2$$

$$\frac{(45,000)^2}{356} \left(1 - \frac{356}{45,000}\right) (63.53)^2$$

$$\frac{N_3^2}{n_3} \left(1 - \frac{n_3}{N_3}\right) S_3^2$$

$$\frac{(20,000)^2}{248} \left(1 - \frac{248}{20,000}\right) (99.49)^2$$

Summing these three terms and taking the square root:

Calculating:

$$T_1 = \frac{(35,000)^2}{396} \left(1 - \frac{396}{35,000}\right) (91.00)^2 = 25,326,894,797.98$$

$$T_2 = \frac{(45,000)^2}{356} \left(1 - \frac{356}{45,000}\right) (63.53)^2 = 22,776,307,940.68$$

$$T_3 = \frac{(20,000)^2}{248} \left(1 - \frac{248}{20,000}\right) (99.49)^2 = 15,766,970,443.16$$

Summing these:

$$T_1 + T_2 + T_3 = 25,326,894,797.98 + 22,776,307,940.68 + 15,766,970,443.16 = 63,870,173,181.82$$

Taking the square root:

$$SE_{Prop}(\hat{T}_{str}) = \sqrt{63,870,173,181.82} \approx 257,435.2 < 257,447.4 = SE_{Srs}(\hat{T}_{str})$$

SE under proportional allocation is smaller than the previously calculated SE, which is expected.

Problem 2 (10 pt)

Consider a simple random sample of size $n = 200$ from a finite population with size $N = 10,000$, measuring (X, Y) , taking values on $\{(0, 0), (0, 1), (1, 0), (1, 1)\}$. The finite population has the following distribution.

	$X = 1$	$X = 0$	
$Y = 1$	N_{11}	N_{10}	N_{1+}
$Y = 0$	N_{01}	N_{00}	N_{0+}
	N_{+1}	N_{+0}	N

The population count N_{ij} are unknown.

Suppose that the realized sample has the following sample counts:

	$X = 1$	$X = 0$	
$Y = 1$	70	30	100
$Y = 0$	50	50	100
	120	80	200

1.

If it is known that $N_{+1} = N_{+0} = 5000$, how can you make use of this information to obtain a post-stratified estimator of $\theta = E(Y)$, using X as the post-stratification variable?

Answer

We can obtain a post-stratified estimator of $\theta = E(Y)$ via:

$$\hat{\theta} = W_1 \bar{y}_1 + W_2 \bar{y}_2 = 0.5 \cdot \left(\frac{70}{120} \right) + 0.5 \cdot \left(\frac{30}{80} \right) = 0.479$$

2.

If we are interested in estimating $\theta = P(Y = 1 | X = 1)$, discuss how to estimate θ from the above sample and how to estimate its variance (Hint: Use Taylor expansion of ratio estimator to obtain the sampling variance).

Answer

For $\theta = P(Y = 1 | X = 1)$, we can make use of conditional probability, i.e.,

$$P(Y = 1 | X = 1) = \frac{P(Y = 1 \cap X = 1)}{P(X = 1)}$$

Substituting values known, and using hats for observed sample counts, we obtain:

$$\hat{\theta} = \frac{\hat{P}(X = 1, Y = 1)}{\hat{P}(X = 1)} = \frac{n_{11}}{n_{1+}} = \frac{70}{120} = 0.583$$

where $n_{1+} = n_{11} + n_{10}$.

Next, for variance estimation, we use the available hint and consider the Taylor expansion of the ratio estimator:

$$\hat{\theta} \approx \theta + \frac{1}{E(n_{1+})}(n_{11} - \theta n_{1+}) = \theta + \frac{1}{E(n_{1+})/n} \left(\frac{n_{11}}{n} - \theta \frac{n_{1+}}{n} \right) = \theta + \frac{1}{E(\hat{P}_{1+})} (\hat{P}_{11} - \theta \hat{P}_{1+})$$

This gives us the variance of $\hat{\theta}$,

$$V(\hat{\theta}) \approx \frac{1}{P_{1+}^2} \left\{ V(\hat{P}_{11}) + \theta^2 V(\hat{P}_{1+}) - 2\theta \text{Cov}(\hat{P}_{11}, \hat{P}_{1+}) \right\}$$

Under SRS then, we have:

$$V(\hat{P}_{11}) = \frac{1}{n}(1-f)P_{11}(1-P_{11})$$

$$V(\hat{P}_{1+}) = \frac{1}{n}(1-f)P_{1+}(1-P_{1+})$$

$$\text{Cov}(\hat{P}_{11}, \hat{P}_{1+}) = \frac{1}{n}(1-f)P_{11}(1-P_{1+})$$

Where f is defined as usual, i.e. $f = \frac{n}{N} = \frac{200}{10,000} = 0.02$.

As given, we know, $\hat{\theta} = \hat{P}_{11}/\hat{P}_{1+} \rightarrow \theta = P_{11}/P_{1+}$.

The above allows us to further simplify the variance formula, specifically:

$$V(\hat{\theta}) = \frac{1}{n}(1-f) \frac{1}{P_{1+}^2} \left\{ P_{11}(1-P_{11}) + \frac{P_{11}^2}{P_{1+}^2} P_{1+}(1-P_{1+}) - 2 \frac{P_{11}}{P_{1+}} P_{11}(1-P_{1+}) \right\}$$

$$V(\hat{\theta}) = \frac{1}{n}(1-f) \frac{1}{P_{1+}} \left\{ P_{11} - \frac{P_{11}^2}{P_{1+}} \right\} = \frac{1}{n}(1-f) \frac{1}{P_{1+}} \theta(1-\theta)$$

Using the above simplification, we can turn to the estimated variance formula:

$$\hat{V}(\hat{\theta}) = (1-f) \frac{1}{n_{1+}} \hat{\theta}(1-\hat{\theta})$$

Using $f = 0.02$, $n_{1+} = 120$, and $\hat{\theta} = 70/120$, we then have:

```
f <- 0.02
n1Plus <- 120
hatTheta <- 70/120

(1 - f) * (1 / n1Plus) * hatTheta * (1 - hatTheta)

## [1] 0.001984954
```

$$\hat{V}(\hat{\theta}) = 0.001985$$

Problem 3 (10 pt)

Suppose that we have a finite population of $(Y_{hi}(1), Y_{hi}(0))$ generated from the following superpopulation model:

$$\begin{pmatrix} Y_{hi}(0) \\ Y_{hi}(1) \end{pmatrix} \sim \left[\begin{pmatrix} \mu_{h0} \\ \mu_{h1} \end{pmatrix}, \begin{pmatrix} \sigma_{h0}^2 & \sigma_{h01} \\ \sigma_{h01} & \sigma_{h1}^2 \end{pmatrix} \right] \quad (1)$$

for $i = 1, \dots, N_h$ and $h = 1, \dots, H$. Instead of observing $(Y_{hi}(0), Y_{hi}(1))$, we observe $T_{hi} \in \{0, 1\}$ and

$$Y_{hi} = T_{hi}Y_{hi}(1) + (1 - T_{hi})Y_{hi}(0)$$

The parameter of interest is the average treatment effect:

$$\tau = \sum_{h=1}^H W_h (\mu_{h1} - \mu_{h0}),$$

where $W_h = N_h/N$. The estimator is:

$$\hat{\tau}_{\text{sre}} = \sum_{h=1}^H W_h \hat{\tau}_h$$

where

$$\hat{\tau}_h = \frac{1}{N_{h1}} \sum_{i=1}^{N_h} T_{hi} Y_{hi} - \frac{1}{N_{h0}} \sum_{i=1}^{N_h} (1 - T_{hi}) Y_{hi}$$

1.

Compute the variance of $\hat{\tau}_{\text{sre}}$ using the model parameters in (1).

Answer

Under the given assumptions, by definition:

$$E(\hat{\tau}_{\text{sre}} \mid \mathcal{F}_N) = \sum_{h=1}^H W_h \bar{\tau}_h$$

where

$$\bar{\tau}_h = N_h^{-1} \sum_{i=1}^{N_h} \{Y_{hi}(1) - Y_{hi}(0)\}$$

Via EVVE, we can express the total variance as:

$$V(\hat{\tau}_{\text{sre}}) = V\{E(\hat{\tau}_{\text{sre}} \mid \mathcal{F}_N)\} + E\{V(\hat{\tau}_{\text{sre}} \mid \mathcal{F}_N)\}$$

For the first term, we note that:

$$V\{E(\hat{\tau}_{sre} \mid \mathcal{F}_N)\} = V\left\{\sum_h^H W_h \bar{\tau}_h\right\}$$

Also, noting the lecture slides, we know the second term of this expression too!

$$V(\hat{\tau}_{sre} \mid \mathcal{F}_N) = \sum_{h=1}^H W_h^2 \frac{1}{N_h} \left(\frac{N_{h0}}{N_{h1}} S_{h1}^2 + \frac{N_{h1}}{N_{h0}} S_{h0}^2 + 2S_{h01} \right)$$

Substituting these into our initial total variance formula:

$$V(\hat{\tau}_{sre}) = V\left\{\sum_{h=1}^H W_h \bar{\tau}_h\right\} + E\left\{\sum_{h=1}^H W_h^2 \frac{1}{N_h} \left(\frac{N_{h0}}{N_{h1}} S_{h1}^2 + \frac{N_{h1}}{N_{h0}} S_{h0}^2 + 2S_{h01} \right)\right\}$$

As given from our model, we know our variance-covariance matrix, such that we have a known formula for the firm term in our total variance equation:

$$V\left\{\sum_{h=1}^H W_h \bar{\tau}_h\right\} = \sum_{h=1}^H W_h^2 \frac{1}{N_h} (\sigma_{h1}^2 + \sigma_{h0}^2 - 2\sigma_{h01})$$

Taking expectation, note that only the assignment of treatments is random (all else are considered fixed), giving us:

$$E\left\{\sum_{h=1}^H W_h^2 \frac{1}{N_h} \left(\frac{N_{h0}}{N_{h1}} S_{h1}^2 + \frac{N_{h1}}{N_{h0}} S_{h0}^2 + 2S_{h01} \right)\right\} = \sum_{h=1}^H W_h^2 \frac{1}{N_h} \left(\frac{N_{h0}}{N_{h1}} \sigma_{h1}^2 + \frac{N_{h1}}{N_{h0}} \sigma_{h0}^2 + 2\sigma_{h01} \right)$$

So we have effectively simplified both parts of our total variance formula. We will find that a number of terms cancel when we bring them all back together, specifically:

$$V(\hat{\tau}_{sre}) = \sum_{h=1}^H W_h^2 \frac{1}{N_h} (\sigma_{h1}^2 + \sigma_{h0}^2 - 2\sigma_{h01}) + \sum_{h=1}^H W_h^2 \frac{1}{N_h} \left(\frac{N_{h0}}{N_{h1}} \sigma_{h1}^2 + \frac{N_{h1}}{N_{h0}} \sigma_{h0}^2 + 2\sigma_{h01} \right) = \sum_{h=1}^H W_h^2 \left(\frac{\sigma_{h1}^2}{N_{h1}} + \frac{\sigma_{h0}^2}{N_{h0}} \right)$$

Noting that:

$$\frac{1}{N_h} \sigma_{h1}^2 + \frac{1}{N_h} \frac{N_{h0}}{N_{h1}} \sigma_{h1}^2 = \frac{\sigma_{h1}^2}{N_{h1}}$$

As $N_h = N_{h1} + N_{h0}$.

2.

Assuming the model parameters in (1) are known, what is the optimal sample allocation such that $\text{Var}(\hat{t}a_{u_{Se}})$ is minimized subject to $N_h = N_{h1} + N_{h0}$ for $h = 1, \dots, H$ are fixed? That is, how to choose N_{h1} and N_{h0} for a given N_h ?

Answer

We want to minimize:

$$Q(N_{h1}, N_{h0}) = \frac{\sigma_{h1}^2}{N_{h1}} + \frac{\sigma_{h0}^2}{N_{h0}}$$

subject to the constraint:

$$N_h = N_{h1} + N_{h0}$$

Given this is a constrained optimization problem, we consider using Lagrange:

$$\mathcal{L}(N_{h1}, N_{h0}, \lambda) = \frac{\sigma_{h1}^2}{N_{h1}} + \frac{\sigma_{h0}^2}{N_{h0}} + \lambda(N_{h1} + N_{h0} - N_h)$$

Taking partial derivatives, setting both zero:

$$\frac{\partial \mathcal{L}}{\partial N_{h1}} = -\frac{\sigma_{h1}^2}{N_{h1}^2} + \lambda = 0$$

$$\frac{\partial \mathcal{L}}{\partial N_{h0}} = -\frac{\sigma_{h0}^2}{N_{h0}^2} + \lambda = 0$$

From the first equation:

$$\lambda = \frac{\sigma_{h1}^2}{N_{h1}^2}$$

From the second equation:

$$\lambda = \frac{\sigma_{h0}^2}{N_{h0}^2}$$

Setting these two expressions equal to one another gives us:

$$\frac{\sigma_{h1}^2}{N_{h1}^2} = \frac{\sigma_{h0}^2}{N_{h0}^2}$$

Taking the square root on both sides to isolate σ_{h1}, σ_{h0} :

$$\frac{\sigma_{h1}}{N_{h1}} = \frac{\sigma_{h0}}{N_{h0}}$$

Now we want to isolate N_{h1} and utilize our constraint:

$$N_{h1} = N_{h0} \frac{\sigma_{h1}}{\sigma_{h0}}$$

As we are constrained by $N_h = N_{h1} + N_{h0}$, we have: $N_{h0} = N_h - N_{h1}$:

$$N_{h1} = (N_h - N_{h1}) \frac{\sigma_{h1}}{\sigma_{h0}}$$

Solving for N_{h1} :

$$N_{h1} \left(1 + \frac{\sigma_{h1}}{\sigma_{h0}} \right) = N_h \frac{\sigma_{h1}}{\sigma_{h0}}$$

$$N_{h1}^* = N_h \frac{\sigma_{h1}}{\sigma_{h1} + \sigma_{h0}}$$

Noting:

$$\frac{\sigma_{h1}}{\sigma_{h1}} = 1 \text{ as by definition } \sigma_{h1}, \sigma_{h0} > 0$$

With N_{h1}^* and our constraint $N_h = N_{h1} + N_{h0}$, we can now easily find N_{h0}^* :

$$N_{h0}^* = N_h - N_{h1}^* = N_h \frac{\sigma_{h0}}{\sigma_{h1} + \sigma_{h0}}$$

Problem 4 (10 pt)

Assume that a simple random sample of size n is selected from a population of size N and (x_i, y_i) are observed in the sample. In addition, we assume that the population mean of x , denoted by \bar{X} , is known.

1.

Use a Taylor linearization method to find the variance of the product estimator $\frac{\bar{x}\bar{y}}{\bar{X}}$, where (\bar{x}, \bar{y}) is the sample mean of (x_i, y_i) .

Answer

The product estimator is:

$$\hat{\theta} = \frac{\bar{x}\bar{y}}{\bar{X}}.$$

Using Taylor linearization, we approximate $\hat{\theta}$ using a first-order expansion around the true means:

$$\hat{\theta} \approx \frac{\bar{Y}\bar{X} + (\bar{x} - \bar{X})\bar{Y} + (\bar{y} - \bar{Y})\bar{X}}{\bar{X}}.$$

Simplifying,

$$\hat{\theta} \approx \bar{Y} + (\bar{x} - \bar{X})\frac{\bar{Y}}{\bar{X}} + (\bar{y} - \bar{Y}).$$

Taking variances,

$$V(\hat{\theta}) \approx V(\bar{y}) + \frac{\bar{Y}^2}{\bar{X}^2} V(\bar{x}) + 2\frac{\bar{Y}}{\bar{X}} \text{Cov}(\bar{x}, \bar{y}).$$

Since in simple random sampling:

$$V(\bar{x}) = \frac{S_x^2}{n} \left(1 - \frac{n}{N}\right), \quad V(\bar{y}) = \frac{S_y^2}{n} \left(1 - \frac{n}{N}\right),$$

$$\text{Cov}(\bar{x}, \bar{y}) = \frac{S_{xy}}{n} \left(1 - \frac{n}{N}\right).$$

Thus,

$$V(\hat{\theta}) \approx \left(\frac{S_y^2}{n} + \frac{\bar{Y}^2}{\bar{X}^2} \frac{S_x^2}{n} + 2\frac{\bar{Y}}{\bar{X}} \frac{S_{xy}}{n} \right) \left(1 - \frac{n}{N}\right).$$

2.

Find the condition that this product estimator has a smaller variance than the sample mean \bar{y} .

Answer

For $\hat{\theta}$ to be more efficient than \bar{y} , we require:

$$V(\hat{\theta}) < V(\bar{y}).$$

Substituting,

$$\frac{S_y^2}{n} + \frac{\bar{Y}^2}{\bar{X}^2} \frac{S_x^2}{n} + 2 \frac{\bar{Y}}{\bar{X}} \frac{S_{xy}}{n} < \frac{S_y^2}{n}.$$

Canceling S_y^2/n ,

$$\frac{\bar{Y}^2}{\bar{X}^2} S_x^2 + 2 \frac{\bar{Y}}{\bar{X}} S_{xy} < 0.$$

Rearranging,

$$2\bar{Y}S_{xy} < -\frac{\bar{Y}^2}{\bar{X}^2} S_x^2.$$

Dividing by \bar{Y} ,

$$2S_{xy} < -\frac{\bar{Y}}{\bar{X}^2} S_x^2.$$

Thus, the product estimator has lower variance when the covariance between x and y is sufficiently negative.

3.

Prove that if the population covariance of x and y is zero, then the product estimator is less efficient than \bar{y} .

Answer

If $S_{xy} = 0$, then the variance formula simplifies to:

$$V(\hat{\theta}) = V(\bar{y}) + \frac{\bar{Y}^2}{\bar{X}^2} V(\bar{x}).$$

Since $\frac{\bar{Y}^2}{\bar{X}^2} V(\bar{x}) > 0$, it follows that:

$$V(\hat{\theta}) > V(\bar{y}).$$

Thus, when x and y are uncorrelated, the product estimator is less efficient than \bar{y} .

Problem 5 (10 pt)

In a population of 10,000 businesses, we want to estimate the average sales \bar{Y} . For that, we sample $n = 100$ businesses using simple random sampling. Furthermore, we have at our disposal the auxiliary information “number of employees”, denoted by x , for each business. It is known that $\bar{X} = 50$ in the population. From the sample, we computed the following statistics:

- $\bar{y}_n = 5.2 \times 10^6$ (average sales in the sample)
- $\bar{x}_n = 45$ employees (sample mean)
- $s_y^2 = 25 \times 10^{10}$ (sample variance of y_k)
- $s_x^2 = 15$ (sample variance of x_k)
- $r = 0.8$ (sample correlation coefficient between x and y)

Answer the following questions:

1.

Compute a 95% confidence interval for \bar{Y} using the ratio estimator.

Answer

The **ratio estimator** for the population mean sales is:

$$\hat{Y}_R = \bar{y}_n \frac{\bar{X}}{\bar{x}_n}.$$

Substituting the given values:

$$\hat{Y}_R = (5.2 \times 10^6) \times \frac{50}{45} = 5.778 \times 10^6.$$

The variance of the ratio estimator is approximated by:

$$V(\hat{Y}_R) \approx \bar{Y}^2 \left(\frac{1}{n} \right) \left(\frac{s_y^2}{\bar{y}_n^2} + \frac{s_x^2}{\bar{x}_n^2} - 2r \frac{s_y}{\bar{y}_n} \frac{s_x}{\bar{x}_n} \right).$$

Substituting the values:

$$V(\hat{Y}_R) = (5.778 \times 10^6)^2 \times \frac{1}{100} \left(\frac{25 \times 10^{10}}{(5.2 \times 10^6)^2} + \frac{15}{45^2} - 2(0.8) \frac{5 \times 10^5}{5.2 \times 10^6} \frac{3.873}{45} \right).$$

Computing the standard error, the **95% confidence interval** is given by:

$$\hat{Y}_R \pm 1.96 \times \sqrt{V(\hat{Y}_R)}.$$

2.

Compute a 95% confidence interval for \bar{Y} using the regression estimator based on the simple linear regression of y on x (with intercept).

Answer

The **regression estimator** for the population mean is:

$$\hat{Y}_{reg} = \bar{y}_n + b(\bar{X} - \bar{x}_n),$$

where the estimated slope b is given by:

$$b = r \frac{s_y}{s_x}.$$

Substituting the values:

$$b = 0.8 \times \frac{5 \times 10^5}{3.873} = 1.033 \times 10^5.$$

Thus, the regression estimator is:

$$\hat{Y}_{reg} = (5.2 \times 10^6) + (1.033 \times 10^5)(50 - 45) = 5.7165 \times 10^6.$$

The variance of the regression estimator is:

$$V(\hat{Y}_{reg}) = \frac{s_y^2}{n}(1 - r^2).$$

Substituting the values:

$$V(\hat{Y}_{reg}) = \frac{25 \times 10^{10}}{100}(1 - 0.64) = 9 \times 10^9.$$

Computing the standard error, the **95% confidence interval** is given by:

$$\hat{Y}_{reg} \pm 1.96 \times \sqrt{V(\hat{Y}_{reg})}.$$

Problem 6 (10 pt)

Under the setup of Chapter 6, Part 1 lecture, prove the last two equalities on page 23:

$$\begin{aligned} \text{Cov} \left(\frac{1}{N_1} \sum_{i=1}^N T_i e_i(1), \frac{1}{N_0} \sum_{i=1}^N (1 - T_i) \mathbf{x}'_i \mathbf{B}_0 | \mathcal{F}_N \right) &= 0 \\ \text{Cov} \left(\frac{1}{N_0} \sum_{i=1}^N (1 - T_i) e_i(0), \frac{1}{N_0} \sum_{i=1}^N (1 - T_i) \mathbf{x}'_i \mathbf{B}_0 | \mathcal{F}_N \right) &= 0 \end{aligned}$$

Answer

Setup

- T_i is the treatment indicator (1 for treated, 0 for control).
- $N_1 = \sum_{i=1}^N T_i$ is the number of treated units.
- $N_0 = \sum_{i=1}^N (1 - T_i)$ is the number of control units.
- $e_i(1)$ and $e_i(0)$ are error terms under treatment and control.
- \mathbf{x}_i is the covariate vector.
- \mathbf{B}_0 is a fixed coefficient vector.

Proof

Since treatment assignments T_i are independent of errors and covariates, we have:

$$E[T_i e_i(1) | \mathcal{F}_N] = \pi_i e_i(1), \quad E[(1 - T_i) e_i(0) | \mathcal{F}_N] = (1 - \pi_i) e_i(0)$$

$$E[T_i \mathbf{x}'_i \mathbf{B}_0 | \mathcal{F}_N] = \pi_i \mathbf{x}'_i \mathbf{B}_0, \quad E[(1 - T_i) \mathbf{x}'_i \mathbf{B}_0 | \mathcal{F}_N] = (1 - \pi_i) \mathbf{x}'_i \mathbf{B}_0$$

where $\pi_i = P(T_i = 1)$.

Expanding the first covariance:

$$\text{Cov} \left(\frac{1}{N_1} \sum_{i=1}^N T_i e_i(1), \frac{1}{N_0} \sum_{i=1}^N (1 - T_i) \mathbf{x}'_i \mathbf{B}_0 | \mathcal{F}_N \right)$$

Expanding linearly:

$$\frac{1}{N_1 N_0} \sum_{i=1}^N \sum_{j=1}^N \text{Cov} (T_i e_i(1), (1 - T_j) \mathbf{x}'_j \mathbf{B}_0 | \mathcal{F}_N)$$

For $i \neq j$, T_i and $(1 - T_j)$ are independent, making cross terms vanish. For $i = j$:

$$\text{Cov}(T_i e_i(1), (1 - T_i) \mathbf{x}'_i \mathbf{B}_0 | \mathcal{F}_N) = E[T_i(1 - T_i) | \mathcal{F}_N] E[e_i(1) \mathbf{x}'_i \mathbf{B}_0 | \mathcal{F}_N]$$

Since $T_i(1 - T_i) = 0$, the covariance is zero.

Similarly, for the second covariance:

$$\text{Cov} \left(\frac{1}{N_0} \sum_{i=1}^N (1 - T_i) e_i(0), \frac{1}{N_0} \sum_{i=1}^N (1 - T_i) \mathbf{x}_i' \mathbf{B}_0 | \mathcal{F}_N \right)$$

Expanding:

$$\frac{1}{N_0^2} \sum_{i=1}^N \sum_{j=1}^N \text{Cov} \left((1 - T_i) e_i(0), (1 - T_j) \mathbf{x}_j' \mathbf{B}_0 | \mathcal{F}_N \right)$$

For $i \neq j$, the terms are independent and vanish. For $i = j$:

$$\text{Cov} \left((1 - T_i) e_i(0), (1 - T_i) \mathbf{x}_i' \mathbf{B}_0 | \mathcal{F}_N \right) = 0.$$

Thus, we have proven:

$$\begin{aligned} \text{Cov} \left(\frac{1}{N_1} \sum_{i=1}^N T_i e_i(1), \frac{1}{N_0} \sum_{i=1}^N (1 - T_i) \mathbf{x}_i' \mathbf{B}_0 | \mathcal{F}_N \right) &= 0 \\ \text{Cov} \left(\frac{1}{N_0} \sum_{i=1}^N (1 - T_i) e_i(0), \frac{1}{N_0} \sum_{i=1}^N (1 - T_i) \mathbf{x}_i' \mathbf{B}_0 | \mathcal{F}_N \right) &= 0 \end{aligned}$$

These results follow from the independence of treatment assignments and errors, ensuring that their covariances vanish.

Problem 7 (10 pt)

Under the setup of Chapter 6, Part 2 lecture:

1.

Prove Lemma 3.

Lemma 3:

Let X be a $n \times p$ matrix such that

$$X = \begin{pmatrix} x'_1 \\ \vdots \\ x'_n \end{pmatrix}$$

and $\omega = (\omega_1, \dots, \omega_n)'$ be an n -dimensional weight vector ($n = N_1$). Given

$$\bar{x} = N^{-1} \sum_{i=1}^N x'_i,$$

and D ($p \times p$ symmetric, invertible matrix), the minimizer of

$$Q(\omega) = \gamma(\omega'X - \bar{x})'D(\omega'X - \bar{x}) + \omega'\omega$$

$$= \gamma(X'\omega - \bar{x})'D(X'\omega - \bar{x}) + \omega'\omega$$

is given by

$$\hat{\omega} = (\gamma XDX' + I_n)^{-1}\gamma XD\bar{x} \tag{10}$$

$$= X(X'X + \gamma^{-1}D^{-1})^{-1}\bar{x} \tag{11}$$

Answer

We seek to minimize the quadratic objective function:

$$Q(\omega) = \gamma(X'\omega - \bar{x})'D(X'\omega - \bar{x}) + \omega'\omega.$$

Taking the gradient with respect to ω :

$$\frac{dQ}{d\omega} = 2\gamma XD(X'\omega - \bar{x}) + 2\omega.$$

Setting this to zero:

$$\gamma XDX'\omega - \gamma XD\bar{x} + \omega = 0.$$

Rearranging:

$$(\gamma XDX' + I_n)\omega = \gamma XD\bar{x}.$$

Multiplying by $(\gamma XDX' + I_n)^{-1}$:

$$\hat{\omega} = (\gamma XDX' + I_n)^{-1}\gamma XD\bar{x}.$$

which proves equation (10), giving us Lemma 3.

2.

Show that the final weight in (13) satisfies a hard calibration for \mathbf{x}_1 :

$$\sum_{i \in A} \hat{\omega}_i \mathbf{x}_{1i} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_{1i}.$$

(9):

The implicit model is that

$$Y(1) = x'_1 \beta + x'_2 u + e(1) \quad (9)$$

where $u \sim (0, D_q \sigma_u^2)$ with known D_q and $e(1) \sim (0, \sigma_e^2)$.

(10-11): Given in Lemma 3

(12):

Using (11), the solution can be written as

$$\hat{\omega} = X (X' X + \Omega^{-1})^{-1} \bar{x} \quad (12)$$

where $\Omega^{-1} = \text{Diag}\{\gamma_1^{-1} D_p^{-1}, \gamma_2^{-1} D_q^{-1}\}$ and $\gamma_1 \rightarrow \infty$.

(13):

Under the mixed model setup in (9), the solution (12) can be written as

$$\hat{\omega}_i = \left(N^{-1} \sum_{i=1}^N x_i \right)' \left\{ \sum_{i=1}^N T_i x_i x'_i + \Omega^{-1} \right\}^{-1} x_i, \quad (13)$$

where $\Omega^{-1} = \text{Diag}\{0_p, \gamma_2^{-1} D_q^{-1}\}$ and $\gamma_2 = \sigma_u^2 / \sigma_e^2$.

Answer

We express the solution using the Woodbury identity:

$$(\gamma X D X' + I_n)^{-1} \gamma X D = X (X' X + \gamma^{-1} D^{-1})^{-1}.$$

Thus, we obtain equation (11):

$$\hat{\omega} = X (X' X + \gamma^{-1} D^{-1})^{-1} \bar{x}.$$

We need to show:

$$\sum_{i \in A} \hat{\omega}_i \mathbf{x}_{1i} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_{1i}.$$

From equation (13):

$$\hat{\omega}_i = \left(N^{-1} \sum_{i=1}^N x_i \right)' \left\{ \sum_{i=1}^N T_i x_i x_i' + \Omega^{-1} \right\}^{-1} x_i.$$

Summing over $i \in A$:

$$\sum_{i \in A} \hat{\omega}_i \mathbf{x}_{1i} = \sum_{i \in A} \left(N^{-1} \sum_{i=1}^N x_i \right)' \left\{ \sum_{i=1}^N T_i x_i x_i' + \Omega^{-1} \right\}^{-1} x_i \mathbf{x}_{1i}.$$

Since $\Omega^{-1} = \text{Diag}\{0_p, \gamma_2^{-1} D_q^{-1}\}$, for large γ_1 , the term simplifies to:

$$\frac{1}{N} \sum_{i=1}^N \mathbf{x}_{1i}.$$

Thus, the final weight satisfies hard calibration.