

Mathematical Foundations of Biostatistics

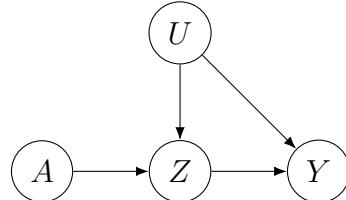
Assignment 2

The assignment is due Wednesday October 16 by 11:59pm. Questions Q1, Q2 and Q3 have the weights of 25%, 25% and 50%, respectively.

1. (a) Let $\mu_n = 1/n$ for $n = 1, 2, \dots$ and let $X_n \sim \text{Poisson}(\mu_n)$.
- Show that $X_n \xrightarrow{p} 0$.
 - Let $Y_n = nX_n$. Show that $Y_n \xrightarrow{p} 0$.
- (b) Suppose we have a computer program consisting of $n = 100$ pages of code. Let X_i be the number of errors on the i th page of code. Suppose that the X_i 's are Poisson with mean 1 and that they are independent. Let $Y = \sum_{i=1}^n X_i$ be the total number of errors. Use the central limit theorem to approximate $P(Y < 90)$.
- (c) Let X_1, \dots, X_n be IID Poisson random variables with mean μ .
- Find the limiting distribution of $\sqrt{n}(\log(\bar{X}_n) - \log(\mu))$.
 - Find a function g such that $\sqrt{n}(g(\bar{X}_n) - g(\mu)) \xrightarrow{d} N(0, 1)$ (known as a variance stabilizing transformation).
2. Consider the data-generating mechanism

$$Y_i = \theta Z_i + U_i + \varepsilon_i$$
$$Z_i = 0.5A_i + U_i + \tau_i,$$

where ε_i , τ_i , A_i and U_i are independent and standard normal distributed. We are interested in estimating the parameter θ , and consider the situation where the observed data are (a_i, y_i, z_i) , $i = 1, \dots, n$, but where U_i is unobserved, so we cannot adjust for it in the analysis. The situation can be illustrated through the below DAG.



Olli Saarela, Dalla Lana School of Public Health, 155 College Street, 6th floor, Toronto, Ontario M5T 3M7, Canada. E-mail: olli.saarela@utoronto.ca

Suppose that, because U_i is unobserved, we fit the model

$$Y_i = \theta_0^* + \theta_1^* Z_i + \varepsilon_i^*,$$

to use the ordinary least squares (OLS) estimator

$$\hat{\theta}_1^* = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(Z_i - \bar{Z})}{\sum_{i=1}^n (Z_i - \bar{Z})}$$

as an estimator of θ .

- (a) Find $\text{plim}_{n \rightarrow \infty} \hat{\theta}_1^*$. Under which circumstances $\hat{\theta}_1^*$ is a consistent estimator of θ ? Is this true in the present setting? How do you interpret the result? What is this phenomenon called? (Hint: you can assume that sample variance and covariance are consistent estimators of their theoretical counterparts, and make use the properties of theoretical covariance.)
- (b) An alternative estimator is obtained by taking $\tilde{\theta} = \hat{\phi}_1/\hat{\psi}_1$, where $\hat{\phi}_1$ and $\hat{\psi}_1$ are the OLS estimators of coefficients ϕ_1 and ψ_1 in the models

$$\begin{aligned} Y_i &= \phi_0 + \phi_1 A_i + \kappa_i \\ Z_i &= \psi_0 + \psi_1 A_i + \xi_i. \end{aligned}$$

Find $\text{plim}_{n \rightarrow \infty} \tilde{\theta}$. Under which circumstances $\tilde{\theta}$ is a consistent estimator of θ ? Is this true in the present setting? What is this estimator/method called?

3. Carry out and report a simulation study of the properties of (the sampling distributions of) $\hat{\theta}_1^*$ and $\tilde{\theta}$. The data generating mechanism is specified as above. Choose $\theta = 0.5$, $m = 1000$ and sample size scenarios $n = 250, 500$. Report mean, bias and standard deviation of the estimator, Monte Carlo error of the mean, mean standard error and 95% normal approximation confidence interval coverage probability. The standard error of $\hat{\theta}_1^*$ may be obtained from linear regression. For $\tilde{\theta}$ you can use the variance estimator

$$\hat{V}[\tilde{\theta}] = \frac{\hat{V}[Y_i - \tilde{\theta}Z_i]}{n\hat{V}[Z_i]\hat{\rho}_{A,Z}^2},$$

where $\hat{\rho}$ is the sample correlation coefficient. Report also the bootstrap standard error, and coverage probability of 95% bootstrap normal approximation and bootstrap percentile confidence intervals. Which one of the two point estimators is better and why? How do the results compare to the theoretical ones in Q2?

- (a) Let $\mu_n = 1/n$ for $n = 1, 2, \dots$ and let $X_n \sim \text{Poisson}(\mu_n)$.
 - Show that $X_n \xrightarrow{p} 0$.
 - Let $Y_n = nX_n$. Show that $Y_n \xrightarrow{p} 0$.
- (b) Suppose we have a computer program consisting of $n = 100$ pages of code. Let X_i be the number of errors on the i th page of code. Suppose that the X_i 's are Poisson with mean 1 and that they are independent. Let $Y = \sum_{i=1}^n X_i$ be the total number of errors. Use the central limit theorem to approximate $P(Y < 90)$.
- (c) Let X_1, \dots, X_n be IID Poisson random variables with mean μ .
 - Find the limiting distribution of $\sqrt{n}(\log(\bar{X}_n) - \log(\mu))$.
 - Find a function g such that $\sqrt{n}(g(\bar{X}_n) - g(\mu)) \xrightarrow{d} N(0, 1)$ (known as a variance stabilizing transformation).

a) i. Need to show $\lim_{n \rightarrow \infty} P(|X_n - \mu| > \varepsilon) = 0$

Since $X_n \sim \text{Poisson}(\mu_n)$, then $E(X_n) = \frac{1}{n} = V(X_n)$

From Chebyshev's Inequality,

$$P(|X_n - \mu| \geq \varepsilon) \leq \frac{V(X_n)}{\varepsilon^2}$$

$$= \frac{1/n}{\varepsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty$$

Thus, $X_n \xrightarrow{P} 0$ as needed. □

ii. We need to show $\lim_{n \rightarrow \infty} P(|Y_n - 0| > \varepsilon) = 0$.

$$\text{Note } P(|Y_n| > \varepsilon) = P(nX_n > \varepsilon)$$

$$= P(X_n > \varepsilon/n)$$

Since $X_i \sim \text{Poisson}(1/n)$, then $P(X_n = k) = \frac{(1/n)^k e^{-1/n}}{k!}$

$$\textcircled{1} P(X_n = 1) = \frac{1}{n} e^{-\frac{1}{n}} < \frac{1}{n} \text{ as } n \rightarrow \infty$$

Note $1/n \rightarrow 0$ as $n \rightarrow \infty$

$$\textcircled{2} P(X_n \geq 1) = 1 - P(X_n = 0) = 1 - e^{-\frac{1}{n}}$$

From Taylor expansion

$$e^x = 1 - \frac{1}{n} + \frac{(-\frac{1}{n})^2}{2!} - \frac{(-\frac{1}{n})^3}{3!} + \dots$$

$$\Rightarrow e^{-\frac{1}{n}} = 1 - \frac{1}{n} + O_n \quad \text{where } O_n \text{ is negligible}$$

$$\Rightarrow 1 - e^{-\frac{1}{n}} \approx \frac{1}{n}$$

$$\text{So } P(X_n \geq 1) = 1 - e^{-\frac{1}{n}} \approx \frac{1}{n} \rightarrow 0 \text{ as } n \rightarrow \infty$$

$$\textcircled{3} P(X_n = 0) = e^{-\frac{1}{n}} \approx 1 - \frac{1}{n} \rightarrow 1 \text{ as } n \rightarrow \infty$$

Then, $P(X_n = 0) = P(nX_n = 0) = P(X_n = 0) \rightarrow 1$
as $n \rightarrow \infty$.

Hence, combining ①, ②, and ③ we can say

$$P(|Y_n - \mu| > \varepsilon) = 1 - P(Y_n = \mu) = 1 - 1 = 0$$

$$\Rightarrow Y_n \xrightarrow{P} \mu$$

b) Let $Y = \sum X_i$. Since $X_i \sim \text{Poisson}(1)$, then

$Y \sim \text{Poisson}(100) \approx N(100, 100)$ by CLT. Then

$$P(Y < 90) = P\left(Z < \frac{90 - 100}{\sqrt{100}}\right)$$

$$= P(Z < -1)$$

$$= 0.1587.$$

(c) Let X_1, \dots, X_n be IID Poisson random variables with mean μ .

- Find the limiting distribution of $\sqrt{n}(\log(\bar{X}_n) - \log(\mu))$.
- Find a function g such that $\sqrt{n}(g(\bar{X}_n) - g(\mu)) \xrightarrow{d} N(0, 1)$ (known as a variance stabilizing transformation).

c) i. Recall $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$ by CLT

Then by Delta Method,

$$\nabla v(x_i) = \mu$$

$$\begin{aligned}\sqrt{n} [g(\bar{x}_n) - g(\mu)] &\xrightarrow{d} N(0, (g'(\mu))^2 \sigma^2) \\ \Rightarrow \sqrt{n} (\log(\bar{x}_n) - \log(\mu)) &\xrightarrow{d} N(0, ((\log \mu)')^2 \sigma^2) \\ &= N(0, 1/\mu)\end{aligned}$$

ii. Since $\sqrt{n}(\bar{x}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$ by CLT.

Then by Delta Method,

$$\sqrt{n} [g(\bar{x}_n) - g(\mu)] \xrightarrow{d} N(0, (g'(\mu))^2 \sigma^2)$$

$$\text{We want } (g'(\mu))^2 \sigma^2 = 1$$

Since $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}$, $E(X_i) = V(X_i) = \mu$

$$\text{So then } (g'(\mu))^2 \mu = 1$$

$$\Rightarrow g(x) = \int \frac{1}{\sqrt{x}} dx = 2\sqrt{x}$$

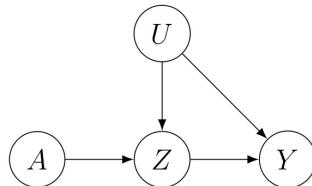
Hence $g(x) = 2\sqrt{x}$ is the variance stabilizing transformation

2. Consider the data-generating mechanism

$$Y_i = \theta Z_i + U_i + \varepsilon_i$$

$$Z_i = 0.5 A_i + U_i + \tau_i,$$

where ε_i , τ_i , A_i and U_i are independent and standard normal distributed. We are interested in estimating the parameter θ , and consider the situation where the observed data are (a_i, y_i, z_i) , $i = 1, \dots, n$, but where U_i is unobserved, so we cannot adjust for it in the analysis. The situation can be illustrated through the below DAG.



Suppose that, because U_i is unobserved, we fit the model

$$Y_i = \theta_0^* + \theta_1^* Z_i + \varepsilon_i^*,$$

to use the ordinary least squares (OLS) estimator

$$\hat{\theta}_1^* = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(Z_i - \bar{Z})}{\sum_{i=1}^n (Z_i - \bar{Z})}$$

as an estimator of θ .

- (a) Find $\text{plim}_{n \rightarrow \infty} \hat{\theta}_1^*$. Under which circumstances $\hat{\theta}_1^*$ is a consistent estimator of θ ? Is this true in the present setting? How do you interpret the result? What is this phenomenon called? (Hint: you can assume that sample variance and covariance are consistent estimators of their theoretical counterparts, and make use the properties of theoretical covariance.)

a) As $n \rightarrow \infty$, from WLLN

$$\begin{aligned} \sum (Y_i - \bar{Y})(Z_i - \bar{Z}) &\xrightarrow{P} \text{cov}(Y, Z) \quad \text{and} \\ \bar{Z} (Z_i - \bar{Z})^2 &\xrightarrow{P} V(Z) \end{aligned}$$

So, we have

$$\lim_{n \rightarrow \infty} \hat{\theta}_i^* = \lim_{n \rightarrow \infty} \frac{\sum (Y_i - \bar{Y})(Z_i - \bar{Z})}{\sum (Z_i - \bar{Z})^2} = \frac{\text{Cov}(Y, Z)}{V(Z)}$$

Since $Y_i = \theta Z_i + U_i + \varepsilon_i$ by the mechanism,
 $Z_i = 0.5A_i + U_i + T_i$

① Then

$$\begin{aligned}\text{Cov}(Y, Z) &= \text{Cov}(\theta Z_i + U_i + \varepsilon_i, Z_i) \\ &= \text{Cov}(\theta Z_i, Z_i) + \text{Cov}(U_i, Z_i) + \text{Cov}(\varepsilon_i, Z_i) \\ &= \theta V(Z) + \underline{\text{Cov}(U_i, Z_i)}\end{aligned}$$

$$\begin{aligned}\text{Cov}(U_i, Z_i) &= \text{Cov}(U_i, 0.5A_i + U_i + T_i) \\ &= \text{Cov}(U_i, 0.5A_i) \xrightarrow{\circ} + V(U_i) + \text{Cov}(U_i, T_i) \\ &= V(U_i) \\ &= 1 \quad \text{as } U_i \sim N(0, 1)\end{aligned}$$

$$\text{So } \text{Cov}(Y, Z) = \theta V(Z) + 1$$

$$(2) \text{ Now, } V(Z_i) = V(0.5A_i + U_i + T_i)$$

$$\stackrel{\text{ind}}{=} (0.5)^2 V(A_i) + V(V_i) + V(I_i)$$

$$= 0.25 + 1 + 1 \quad \text{as } A_i, V_i, I_i \sim N(0, 1)$$

$$= 2.25$$

$$\text{Hence, } \hat{\theta}_1^* = \frac{\sum (Y_i - \bar{Y})(Z_i - \bar{Z})}{\sum (Z_i - \bar{Z})^2}$$

$$\approx \frac{\text{cov}(Y, Z)}{V(Z)} = \frac{\theta V(Z)}{V(Z)} + \frac{V(U)}{V(Z)}$$

$$= \frac{2.25\theta + 1}{2.25}$$

$$= \theta + \frac{4}{9}$$

By defⁿ, $\hat{\theta}_n$ is consistent if $\hat{\theta}_n \xrightarrow{P} \theta$, which is not the case for $\hat{\theta}_1^*$. This occurs as both Z and Y (regressor and outcome) depends on the unobserved V_i , but V_i is not omitted from the model. It is consistent if $\text{Cov}(Z, V) = 0$,

V_i is a confounding variable. It violates the independence assumption in OLS models. This phenomenon is called Omitted-variable bias (Wikipedia).

- (b) An alternative estimator is obtained by taking $\tilde{\theta} = \hat{\phi}_1/\hat{\psi}_1$, where $\hat{\phi}_1$ and $\hat{\psi}_1$ are the OLS estimators of coefficients ϕ_1 and ψ_1 in the models

$$Y_i = \phi_0 + \phi_1 A_i + \kappa_i$$

$$Z_i = \psi_0 + \psi_1 A_i + \xi_i.$$

Find $\text{plim}_{n \rightarrow \infty} \tilde{\theta}$. Under which circumstances $\tilde{\theta}$ is a consistent estimator of θ ? Is this true in the present setting? What is this estimator/method called?

To use our new estimator $\tilde{\theta} = \hat{\phi}_1/\hat{\psi}_1$, recall our original data generating mechanism

$$\begin{cases} Y_i = \theta Z_i + U_i + \varepsilon_i \\ Z_i = 0.5 A_i + V_i + \tau_i \end{cases}$$

① Now for our model $Y_i = \phi_0 + \phi_1 A_i + \kappa_i$

Substitute the expression Y_i from the original data-generating mechanism,

$$\begin{aligned} Y_i &= \theta Z_i + U_i + \varepsilon_i \\ &= \theta(0.5 A_i + V_i + \tau_i) + U_i + \varepsilon_i \\ &= 0.5\theta A_i + \text{other terms} \end{aligned}$$

As $A_i \perp \text{other terms}$, we see $\hat{\phi}_1 = 0.5\theta$.
As $n \rightarrow \infty$, since OLS estimators are consistent

i.e. $\text{plim}_{n \rightarrow \infty} \hat{\phi}_1 = \phi_1 = 0.5\theta$

(2) Now for $Z_i = \psi_0 + \psi_1 A_i + \epsilon_i$, from the original data generating mechanism,

$$Z_i = 0.5A_i + V_i + \epsilon_i$$

We see $\hat{\psi}_1 = 0.5$ as $A_i \perp \text{other terms}$.

Now

$$\operatorname{plim}_{n \rightarrow \infty} \hat{\psi}_1 = \psi_1 = 0.5$$

Then, $\tilde{\theta} = \hat{\beta}_1 / \hat{\psi}_1 = \frac{0.5\theta}{0.5} = \theta$. Both $\hat{\beta}_1$ and $\hat{\psi}_1$ converge to their true values due to OLS estimator consistency. Hence, $\operatorname{plim}_{n \rightarrow \infty} \tilde{\theta} = \theta$ is true.

This is called 2 stage least squares (2SLS Wikipedia). Another name that came up was instrumental variable analysis, which is the scenario for something like this.

3. Carry out and report a simulation study of the properties of (the sampling distributions of) $\hat{\theta}_1^*$ and $\tilde{\theta}$. The data generating mechanism is specified as above. Choose $\theta = 0.5$, $m = 1000$ and sample size scenarios $n = 250, 500$. Report mean, bias and standard deviation of the estimator, Monte Carlo error of the mean, mean standard error and 95% normal approximation confidence interval coverage probability. The standard error of $\hat{\theta}_1^*$ may be obtained from linear regression. For $\tilde{\theta}$ you can use the variance estimator

$$\widehat{V}[\tilde{\theta}] = \frac{\widehat{V}[Y_i - \tilde{\theta}Z_i]}{n\widehat{V}[Z_i]\widehat{\rho}_{A,Z}^2},$$

where $\widehat{\rho}$ is the sample correlation coefficient. Report also the bootstrap standard error, and coverage probability of 95% bootstrap normal approximation and bootstrap percentile confidence intervals. Which one of the two point estimators is better and why? How do the results compare to the theoretical ones in Q2?

Output:

1) $\hat{\theta}^*$ with $n=250$

```
[1] "theta_star"
  No. runs nobs  Mean Bias SD MCE Mean SE Normal CI coverage Mean BS SE
[1,] 1000 250 0.943 0.443 0.052 0.002 0.053 0.053 0 0.052
      Normal BS CI coverage BS Percent lower BS Percent upper
[1,] 0 0.845 1.043
```

2) $\tilde{\theta}$ with $n=250$

```
[1] "theta_tilde"
  #runs nobs  Mean Bias SD MCE Mean SE Normal CI coverage Mean BS SE Normal BS CI coverage BS Percent lower
[1,] 1000 250 0.481 -0.019 0.197 0.006 0.193 0.958 0.247 0.968 -0.032
      BS Percent upper
[1,] -0.032
```

3) $\hat{\theta}^*$ with $n=500$

```
[1] "theta_star"
No. runs nobs Mean Bias SD MCE Mean SE Normal CI coverage Mean BS SE Normal BS CI coverage BS Percent lower
[1,] 1000 500 0.943 0.443 0.038 0.001 0.037 0 0.037 0 0.873
BS Percent upper
[1,] 1.014
```

4) $\tilde{\theta}$ with $n=500$

```
[1] "theta_tilde"
#runs nobs Mean Bias SD MCE Mean SE Normal CI coverage Mean BS SE Normal BS CI coverage BS Percent lower
[1,] 1000 500 0.491 -0.009 0.132 0.004 0.13 0.954 0.136 0.962 0.201
BS Percent upper
[1,] 0.201
```

We can see $\tilde{\theta}$ is a better estimator than $\hat{\theta}^*$ from the Mean (around 0.5), Bias (close to 0). Also the 95% CI coverage probability is close to 1 for $\tilde{\theta}$ compared to 0 for $\hat{\theta}^*$ for both $n=250$ and $n=500$ which is good.

If matches with our theoretical analysis from Q2.

Note: Code included as QMD.