

TAREA 1

CIENCIA DE DATOS

Brain de Jesús Salazar, César Ávila, Iván García

12 de septiembre de 2025

Solución del problema 1.

Solución del problema 2.

Solución del problema 3.

Solución del problema 4.

Solución del problema 5. Consideremos una muestra Y_1, \dots, Y_n de variables aleatorias independientes e idénticamente distribuidas con media μ y varianza σ^2 , con indicadores $R_i \in \{0, 1\}$, de tal manera que $R_i \perp Y_i$ para cada $i \in \{1, \dots, n\}$. Dichas R_i existen pues estamos bajo el modelo MCAR, y se interpretan como $R_i = 1$ si y solo si el dato Y_i fue observado.

Notemos que si n_{obs} representa el número de datos observados, entonces $n_{obs} = \sum_{i=1}^n R_i$. Además, por la definición de los R_i ,

$$\bar{Y}_{obs} = \frac{1}{n_{obs}} \sum_{i=1}^n R_i Y_i = \frac{1}{n_{obs}} \sum_{i: R_i=1} Y_i.$$

Así pues, si $\mathbf{R} = (R_1, \dots, R_n)$, entonces

$$\begin{aligned} \mathbb{E} [\bar{Y}_{obs} | \mathbf{R}] &= \mathbb{E} \left[\frac{1}{\sum_{i=1}^n R_i} \sum_{i=1}^n R_i Y_i \middle| \mathbf{R} \right] = \frac{1}{\sum_{i=1}^n R_i} \sum_{i=1}^n \mathbb{E} [R_i Y_i | \mathbf{R}] \\ &= \frac{1}{\sum_{i=1}^n R_i} \sum_{i=1}^n R_i \mathbb{E} [Y_i | \mathbf{R}] \\ &= \frac{1}{\sum_{i=1}^n R_i} \sum_{i=1}^n R_i \mathbb{E} [Y_i] \\ &= \mu, \end{aligned}$$

en donde hemos usado que las Y_i son iid con media μ y son independientes de \mathbf{R} (por las hipótesis del modelo MCAR). Por consiguiente,

$$\mathbb{E} [\bar{Y}_{obs}] = \mathbb{E} [\mathbb{E} [\bar{Y}_{obs} | \mathbf{R}]] = \mu.$$

Por otra parte,

$$\bar{Y}_{obs}^2 = \frac{1}{(\sum_{i=1}^n R_i)^2} \left[\sum_{i=1}^n R_i^2 Y_i^2 + \sum_{i \neq j} R_i R_j Y_i Y_j \right],$$

por lo que

$$\begin{aligned} \mathbb{E} [\bar{Y}_{obs}^2 | \mathbf{R}] &= \frac{1}{(\sum_{i=1}^n R_i)^2} \left[\sum_{i=1}^n R_i^2 \mathbb{E} [Y_i^2 | \mathbf{R}] + \sum_{i \neq j} R_i R_j \mathbb{E} [Y_i Y_j | \mathbf{R}] \right] \\ &= \frac{1}{(\sum_{i=1}^n R_i)^2} \left[\sum_{i=1}^n R_i^2 \mathbb{E} [Y_i^2] + \sum_{i \neq j} R_i R_j \mathbb{E} [Y_i Y_j] \right] \\ &= \frac{1}{(\sum_{i=1}^n R_i)^2} \left[(\sigma^2 + \mu^2) \sum_{i=1}^n R_i^2 + \mu^2 \sum_{i \neq j} R_i R_j \right] \\ &= \mu^2 + \sigma^2 \frac{\sum_{i=1}^n R_i^2}{(\sum_{i=1}^n R_i)^2}. \end{aligned}$$

De lo anterior se sigue que

$$\begin{aligned}
\text{Var} [\bar{Y}_{obs}] &= \mathbb{E} [\bar{Y}_{obs}^2] - (\mathbb{E} [\bar{Y}_{obs}])^2 = \mathbb{E} [\mathbb{E} [\bar{Y}_{obs}^2 \mid \mathbf{R}]] - \mu^2 = \sigma^2 \mathbb{E} \left[\frac{\sum_{i=1}^n R_i^2}{(\sum_{i=1}^n R_i)^2} \right] \\
&= \sigma^2 \mathbb{E} \left[\frac{\sum_{i=1}^n R_i}{(\sum_{i=1}^n R_i)^2} \right] \\
&= \sigma^2 \mathbb{E} \left[\frac{1}{\sum_{i=1}^n R_i} \right] \\
&= \sigma^2 \mathbb{E} \left[\frac{1}{n_{obs}} \right] \\
&\geq \frac{\sigma^2}{n} = \text{Var} [\bar{Y}] ,
\end{aligned}$$

en donde hemos usado que $R_i^2 = R_i$, pues $R_i \in \{0, 1\}$, para todo $i \in \{1, \dots, n\}$, y que $n_{obs} \leq n$. De hecho, siguiendo un procedimiento completamente análogo, pero ahora con varianzas condicionales, se sigue que

$$\text{Var} [\bar{Y}_{obs} \mid \mathbf{R}] = \frac{\sigma^2}{n_{obs}}.$$

De lo anterior podemos notar que \bar{Y}_{obs} es insesgado, pero que $\text{Var} [\bar{Y}_{obs}] \geq \text{Var} [\bar{Y}]$, por lo que \bar{Y}_{obs} tiene menor eficiencia (posee más varianza, pues la eliminación de datos hace que haya menos de ellos para poder estimar a la media de Y).

Solución del problema 6. Sean $\mathbf{Y} = (Y_{obs}, Y_{mis})$ y \mathbf{R} el patrón de datos faltantes. Bajo la definición del MAR,

$$\mathbb{P}[\mathbf{R} | Y_{obs}, Y_{mis}, \theta, \psi] = \mathbb{P}[\mathbf{R} | Y_{obs}, \psi].$$

Así pues, bajo este modelo,

$$\mathbb{P}[\mathbf{Y}, \mathbf{R} | \theta, \psi] = \mathbb{P}[\mathbf{Y} | \theta] \mathbb{P}[\mathbf{R} | \mathbf{Y}, \psi] = \mathbb{P}[\mathbf{Y} | \theta] \mathbb{P}[\mathbf{R} | Y_{obs}, \psi].$$

Por consiguiente, la verosimilitud de θ está dada por

$$\begin{aligned} L(\theta; Y_{obs}, \mathbf{R}) &= \int \mathbb{P}[\mathbf{Y}, \mathbf{R} | \theta, \psi] dY_{mis} = \int \mathbb{P}[\mathbf{Y} | \theta] \mathbb{P}[\mathbf{R} | Y_{obs}, \psi] dY_{mis} \\ &= \mathbb{P}[\mathbf{R} | Y_{obs}, \psi] \int \mathbb{P}[\mathbf{Y} | \theta] dY_{mis} \\ &= \mathbb{P}[\mathbf{R} | Y_{obs}, \psi] \mathbb{P}[Y_{obs} | \theta]. \end{aligned}$$

Ya que el factor $\mathbb{P}[\mathbf{R} | Y_{obs}, \psi]$ no depende de θ , se sigue que $L(\theta; Y_{obs}, \mathbf{R}) \propto \mathbb{P}[Y_{obs} | \theta]$. Para ver las condiciones *a priori* que garantizan ignorabilidad bajo el enfoque bayesiano, notemos que

$$\begin{aligned} \mathbb{P}[\theta | Y_{obs}, \mathbf{R}] &= \int \mathbb{P}[\theta, \psi | Y_{obs}, \mathbf{R}] d\psi \propto \int \mathbb{P}[Y_{obs}, \mathbf{R} | \theta, \psi] \pi(\theta, \psi) d\psi \\ &= \int \mathbb{P}[\mathbf{R} | Y_{obs}, \psi] \mathbb{P}[Y_{obs} | \theta] \pi(\theta, \psi) d\psi \\ &= \mathbb{P}[Y_{obs} | \theta] \int \mathbb{P}[\mathbf{R} | Y_{obs}, \psi] \pi(\theta, \psi) d\psi. \end{aligned}$$

Para concluir ignorabilidad buscamos que $L(\theta | Y_{obs}, \mathbf{R}) \propto \pi(\theta) \mathbb{P}[Y_{obs} | \theta]$, y ya que la integral anterior depende de θ solamente a través del factor $\pi(\theta, \psi)$, dicha ignorabilidad se logra cuando $\pi(\theta, \psi) = \pi(\theta)\pi(\psi)$. Es decir, si en la *a priori* se pide indistinguibilidad de los parámetros (i.e. que θ y ψ sean independientes), entonces el mecanismo es ignorable para inferir θ .

Solución del problema 7.

Solución del problema 8.

Solución del problema 9.

Solución del problema 10.

Solución del problema 11.

Solución del problema 12.

Solución del problema 13.