

Regresión lineal I

introducción y conceptos base

Profesor: Raimundo Undurraga

Auxiliares: Brandon Galarza, Camila Jáuregui Leonardo Meneses, Francisca Monetta,
Matias Reyes, Bastian Urzua, Antonia Villegas

P1.- Demostraciones básicas

Es sabido que la Regresión Lineal Simple supone que el fenómeno en estudio sigue una tendencia lineal de la forma:

$$y = \hat{\beta}_0 + \hat{\beta}_1 X$$

donde "y" es la variable dependiente, "X" la variable independiente y $\hat{\beta}_1, \hat{\beta}_0$ son parámetros a determinar. Considerando que $y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + \epsilon_i$ (la demanda observada $y_i, i = 1, \dots, n$ comete un error ϵ_i), estos parámetros son escogidos de forma que la suma del cuadrado de los errores ϵ_i , provenientes de la observación i sea mínima. Demuestre que bajo lo anterior se tiene que:

$$\hat{\beta}_1 = \frac{\frac{1}{n-1} \sum_{i=1}^n X_i y_i - \bar{y} \bar{X} N}{\frac{1}{n-1} \sum_{i=1}^n X_i^2 - \bar{X}^2 N} = \frac{Cov(X_i, y_i)}{Var(X_i)} \quad y \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{X}$$

Con \bar{X} e \bar{y} promedios muestrales de X e y respectivamente.

P2.- Aplicación

El ministerio de salud desea realizar un análisis de productividad de los médicos que atienden en consultorios. Se desea estudiar como la experiencia de los doctores (medida en años transcurridos desde que comenzaron a trabajar en los consultorios), X - afecta la productividad, medido en número de pacientes atendidos por hora, y. Para esto, se recopiló información de 200 médicos de distintas especialidades de distintas zonas del país. Para cada médico, se calculó su experiencia (exp) y la productividad durante el último año (prod). La siguiente tabla muestra algunos estadísticos obtenidos para los 200 doctores. La correlación entre las dos variables es 0.64. Asuma que la varianza de los residuos es 1,65.

Variable	Obs	Mean	Std. Dev.	Min	Max
exp(x)	200	6,81	4,27	0,00	14,00
prod(y)	200	2,69	1,39	0,45	10,91

1.- Se desea estimar una regresión $y = \beta_0 + \beta_1 X + \epsilon$. Con los datos proporcionados, estime β_0 y β_1 mediante Mínimos Cuadrados Ordinarios

2.-¿Qué se puede decir sobre la significancia de β_1 ? Realice un Test de Hipótesis.

P3.- Interpretación de resultados

Se realizó una regresión lineal en R, para estudiar la incidencia del PIB per cápita en el índice de obesidad mundial, por cada país. Donde *obesidad_pais* es la variable dependiente y *pib_per_capita* la independiente.

1.- ¿Qué otras variables podrían ser útiles para mejorar la estimación de la variable obesidad, además de pib per cápita? Proponga otros 2 regresores y justifique.

2.- ¿Qué efecto se espera de pib per cápita y las otras dos variables propuestas, sobre la variable dependiente?

Luego se obtienen los resultados de la regresión, siendo los siguientes:

Regresión p3	
Obesidad por país	
pib_per_capita	0.0002*** (0.00004)
Constant	17.1735*** (0.9750)
<i>N</i>	185
R ²	0.0912
Adjusted R ²	0.0863
Residual Std. Error	10.7483 (df = 183)
F Statistic	18.3684*** (df = 1; 183)
<i>Notes:</i>	*** Significant at the 1 percent level. ** Significant at the 5 percent level. * Significant at the 10 percent level.

Figura 1: Resultado de la regresión

3.- ¿Cuáles son los valores de beta 1 y del intercepto? ¿Cuál es la desviación estándar del valor de beta 1?

4.- ¿Se obtuvo el efecto esperado que se mencionó en 2? ¿Es significativo el resultado obtenido?

5.- Con la información de la tabla responda: ¿basta con el regresor pib per cápita elegido para explicar el índice de obesidad?

Resumen

Regresión lineal simple: Busca aproximar el comportamiento de una variable de forma lineal a través de otra variable, de la forma más exacta posible, es decir, minimizando los residuos. Se escribe de la forma:

$$y_i = \beta_0 + \beta_{1i}X_i + \epsilon_i$$

Donde y es la variable a explicar, o **variable dependiente**, X_i es la variable explicativa o **regresor** y ϵ es el **error**, que corresponde a un factor aleatorio no observable.

Mínimos cuadrados ordinarios: (MCO u OLS) Minimizar el residuo significa hacer que sea lo más pequeño posible. Si se hiciera directamente, sería algo del estilo: $\sum_{i=1}^n \epsilon_i \rightarrow -\infty$, lo cual resultaría una estimación extremadamente distinta de y , por lo que no es correcto minimizarlo de esta manera. De aquí nace el concepto de **mínimos cuadrados ordinarios (MCO)**, con lo cual se busca minimizar la diferencia entre la estimación y la observación (residuo) al cuadrado, de forma que : $\sum_{i=1}^n \epsilon_i^2 \rightarrow 0$.

Reemplazando lo anterior de la fórmula de regresión lineal entonces se tendrá:

$$\begin{aligned} 0 &= \sum_{i=1}^n \epsilon_i^2 \\ 0 &= \sum_{i=1}^n (y_i - \beta_0 - \beta_{1i}X_i)^2 \end{aligned} \tag{1}$$

Resolviendo la ecuación anterior, se encuentran los óptimos $\hat{\beta}_0$ y $\hat{\beta}_1$ tales que finalmente el resultado de la regresión lineal simple será lo que se conoce como **valores predichos**:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

Visualmente lo que busca la regresión lineal es la recta que relacione ambas variables tal que minimice la distancia de cada observación a sí misma de la **mejor** forma:

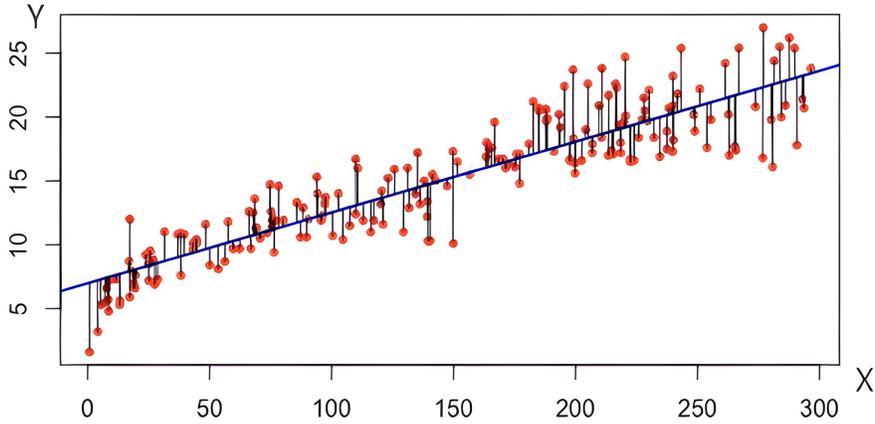


Figura 2: Regresión lineal en azul, observaciones en rojo

Importante destacar que el término residuo también se suele llamar más comunmente error, pero no confundir con error estándar pues no tienen nada que ver.

Resolucion del modelo con CPO: mediante el método de MCO se resuelve la ecuación (1), para lo cual se obtienen las ecuaciones:

- Para la pendiente: $\hat{\beta}_1 = \frac{Cov(X_i, y_i)}{Var(X_i)}$
- Para el intercepto: $\hat{\beta}_0 = \bar{y}_i - \hat{\beta}_1 X_i$

De lo anterior, $\hat{\beta}_1$ también puede escribirse como $\hat{\beta}_1 = \beta_1 + \frac{Cov(X, \varepsilon)}{Var(X)}$, entonces es directo obtener que si $\frac{Cov(X, \varepsilon)}{Var(X)} = 0$, entonces $\hat{\beta}_1$ es insesgado. En específico:

- Si $Cov(X, \varepsilon) = 0$, entonces $\hat{\beta}_1$ es insesgado con respecto a β_1 .
- Si $Cov(X, \varepsilon) \neq 0$, entonces $\hat{\beta}_1$ es sesgado con respecto a β_1 .

Resultados del MCO: se tienen los siguientes elementos como resultado de todo lo anterior:

- Valores predichos: $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$
- Residuo: $e_i = y_i - \hat{y}_i$
- Suma de cuadrados:
 - sum of squares total: $SST = \sum_{i=1}^n (y_i - \bar{y})^2$

- sum of squared regression: $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$
- sum of square errors: $SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$

Propiedades del MCO

- $\frac{1}{n} \sum_{i=1}^n \epsilon_i = 0$
- $Cov(\epsilon_i, X_i) = 0$, es decir: $\sum_{i=1}^n X_i \epsilon_i = 0$
- $SST = SSR + SSE$

R cuadrado: este coeficiente se obtiene de las regresiones lineales, y se calcula como:
 $R^2 = \frac{SSR}{SST} \Leftrightarrow R^2 = 1 - \frac{SSE}{SST}$

Es un valor entre 0 y 1 y representa el porcentaje de los datos correspondientes a y, que son explicados por el modelo.