

# Análisis de regresión

24 de agosto  
(semana 3)

# Plan de trabajo

1. Introducción a la regresión lineal simple
2. Inferencia sobre los coeficientes de regresión

## **Nota:**

- No olviden aceptar la invitación al CLASSROOM y revisar acceso al DRIVE. Pendiente de mi parte crear MOODLE.
- La primera entrega del trabajo final será para el domingo 4 de septiembre por CLASSROOM.
  - Llenen la lista que les compartí con los integrantes.
  - Si van a recoger sus datos, NO empiecen hasta que yo les dé el aval.
  - En el DRIVE, videos de repaso de estadística descriptiva.
- En la semana del 29 de agosto al 2 de septiembre no habrá clase. Usen ese tiempo para la 1era entrega y el taller 1.
- Impriman por favor material módulo 2.

# Repaso de probabilidad

# Repaso: operador covarianza

## *Teorema: propiedades de la covarianza*

- (1)  $\text{cov}(X, Y) = E[XY] - E[X]E[Y];$
- (2) (*Simetría*)  $\text{cov}(X, Y) = \text{cov}(Y, X);$
- (3) (*Positividad semidefinida*)  $\text{cov}(X, X) = \text{Var}(X);$
- (4) (*Linealidad 1*)  $\text{cov}(X, d) = 0, \forall d \in \mathbb{R};$
- (5) (*Linealidad 2*)  $\text{cov}(aX + b, Y) = \text{cov}(aX, Y) + \text{cov}(b, Y)$   
 $= a \text{cov}(X, Y), \forall a, b \in \mathbb{R};$
- (6) (*Linealidad 3*)  $\text{cov}(aX + b, cY + d) = ac \text{cov}(X, Y),$   
 $\forall a, b, c, d \in \mathbb{R}.$
- (7) (*Independencia  $\Rightarrow \text{cov} = 0$* ) Si  $X, Y$  son independientes,  
 $\text{cov}(X, Y) = 0.$

# Repaso: operador covarianza

## *Teorema: propiedades de la covarianza (II)*

Para  $X, Y, X_1, X_2, \dots, X_n \in L^2(\Omega, \mathfrak{F}, p)$ , se tiene

(1) Desigualdad de Cauchy-Schwarz:

$$|\text{cov}(X, Y)| \leq \sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}.$$

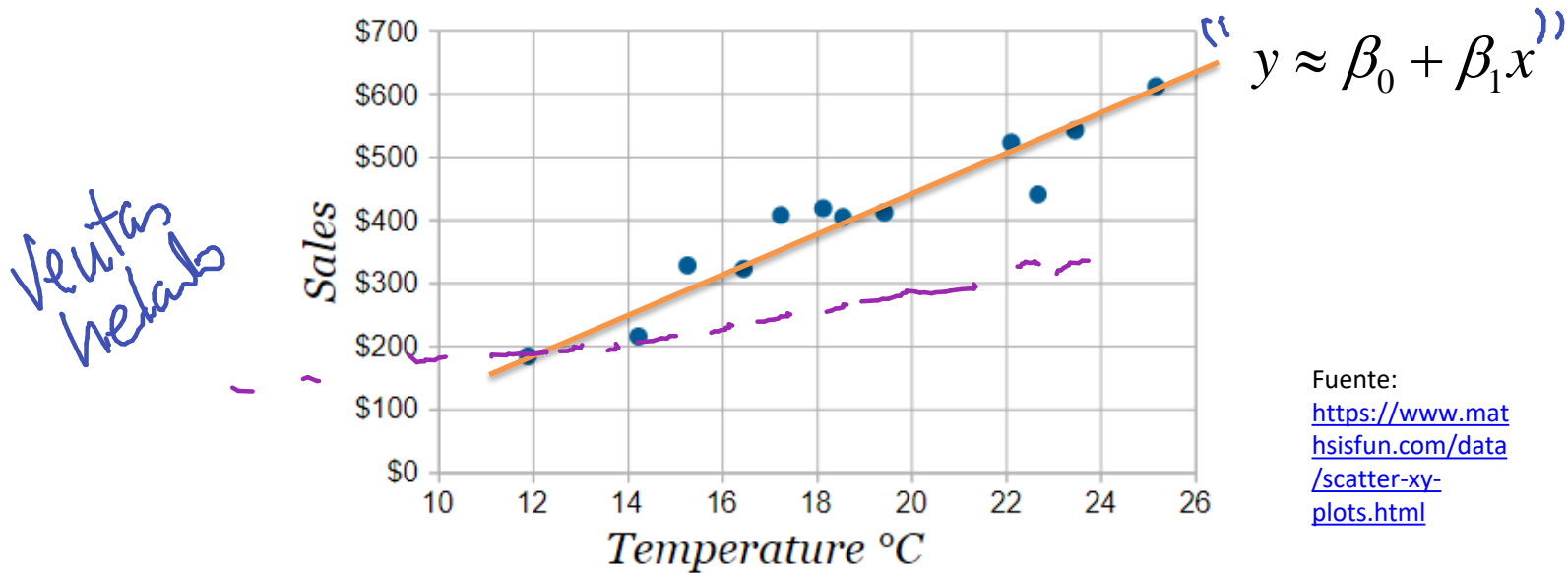
$$(2) \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{cov}(X, Y).$$

$$(3) \text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{Var}(X_i) + \sum_{i \neq j} \text{cov}(X_i, X_j) \\ = \sum_{i=1}^n \text{Var}(X_i) + 2 \sum_{i < j} \text{cov}(X_i, X_j).$$

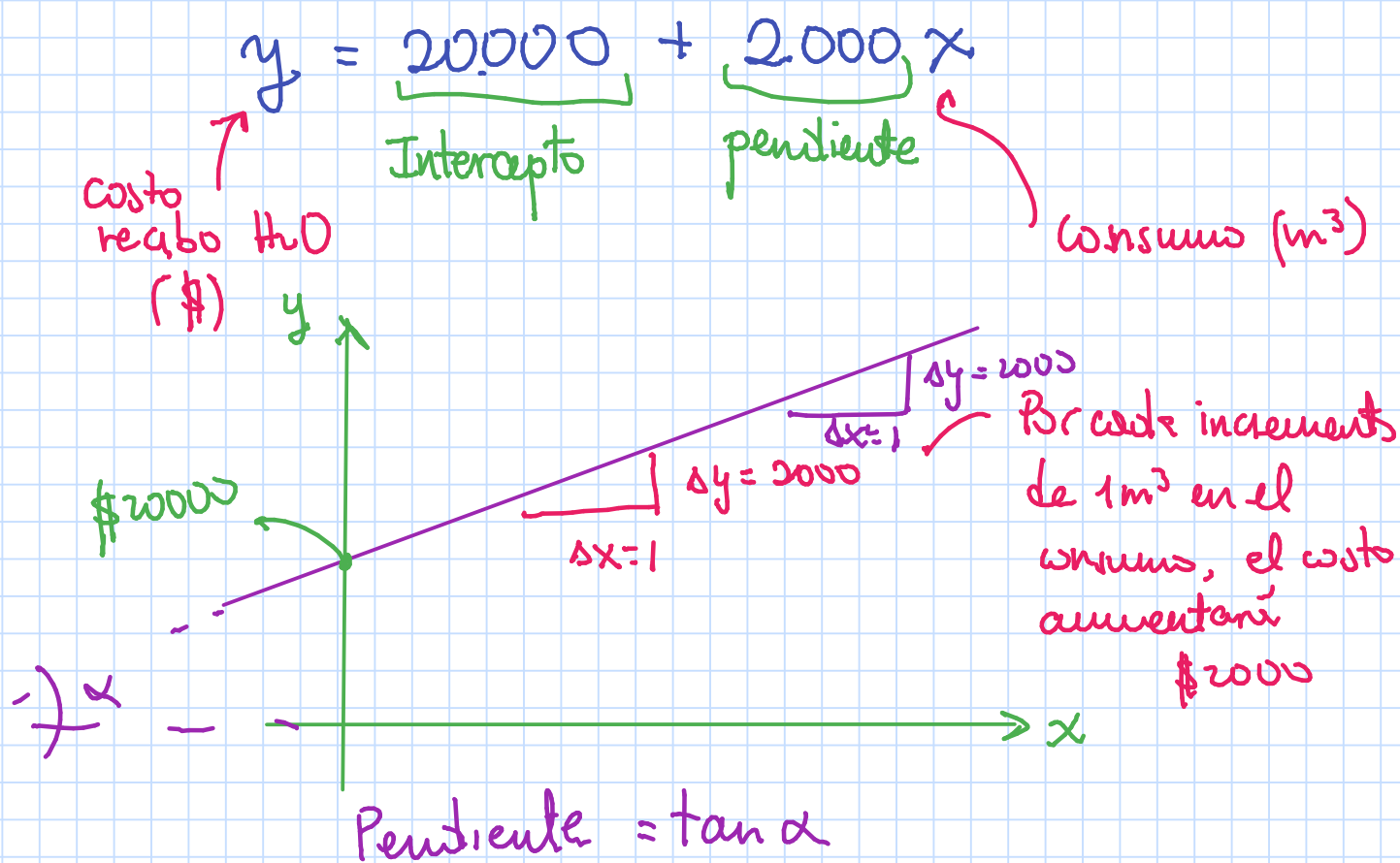
# Regresión lineal simple

# Coeficientes de regresión

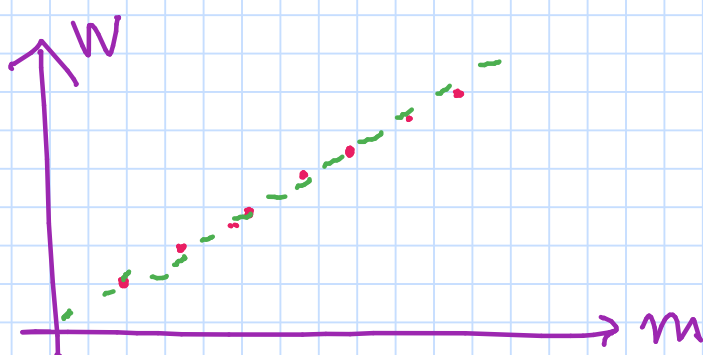
- ❖ Estos coeficientes determinan la “mejor” relación lineal que puede ser descrita entre las dos variables



- ❖ Mucho más útil porque permite hacer predicciones de una variable en función de la otra (una vez conozcamos estimaciones de los parámetros).
- ❖ ¿Cómo encontrar los valores de los coeficientes?



$W = mg$



Modelo de los helados

$Y = \boxed{\beta_0 + \beta_1 X} + \boxed{\epsilon}$

Ventas por helados en un día  $Y$

Temp ( $^{\circ}C$ )  $X$

Comp Aleatoria  $\epsilon$

Comp Aleatoria  $\epsilon$

Brwn no explicada por la componente sistemática

Componente sistemática ( $\mu$ )

Modelo de regresión lineal simple

$Y_k = \mu_k + \epsilon_k$

Comp. System  $\mu_k$

Comp. aleatoria  $\epsilon_k$

$\mu_k = \beta_0 + \beta_1 X_k$

$\epsilon_1, \epsilon_2, \dots, \epsilon_k$  m.a.

$E[\epsilon_i] = 0$  (s1)

$Var[\epsilon_i] = \sigma^2 \forall i$  (s2. Homoscedasticidad)

3 parámetros

$E[Y_k] = E[\mu_k + \epsilon_k] = E[\mu_k] + E[\epsilon_k]$

$= E[\beta_0 + \beta_1 X_k]$

$= \beta_0 + \beta_1 X_k$

$\beta_0$ : Ventas promedio o esperadas cuando  $x = 0^{\circ}C$

$\beta_1$ : Por cada  $1^{\circ}C$  que aumente la temperatura, las ventas esperadas o promedio cambiarán  $\beta_1$  unidades



# Supuestos del modelo de regresión lineal simple

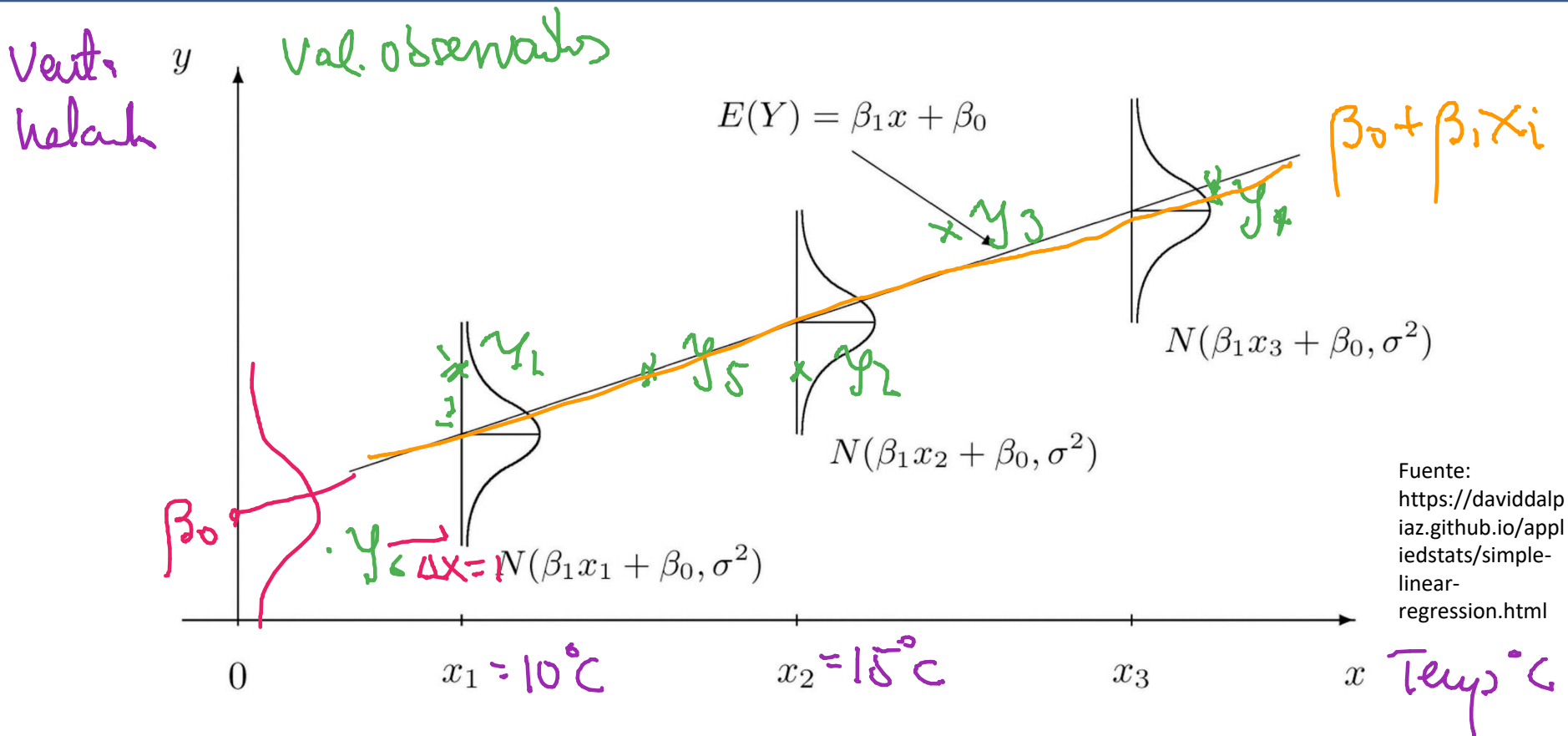
## *Supuestos del modelo de regresión lineal simple*

- La variable  $X$  (cualitativa o cuantitativa) se asume siempre condicionante, es decir, dado que  $X = x$ .
- La relación entre ambas variables es lineal, es decir:

$$Y_i | \{X = x_i\} = \beta_0 + \beta_1 x_i + \varepsilon_i$$

- $\{\varepsilon_i\}$  son una m.a. **generalmente** de una distr.  $N(0, \sigma^2)$ 
  - × Las  $\{Y_i\}$  no son una m.a.
  - ×  $Y_i | \{X = x_i\} \sim N(\underbrace{\beta_0 + \beta_1 x_i}, \sigma^2), \forall i.$

# Supuestos del modelo de regresión simple (II)

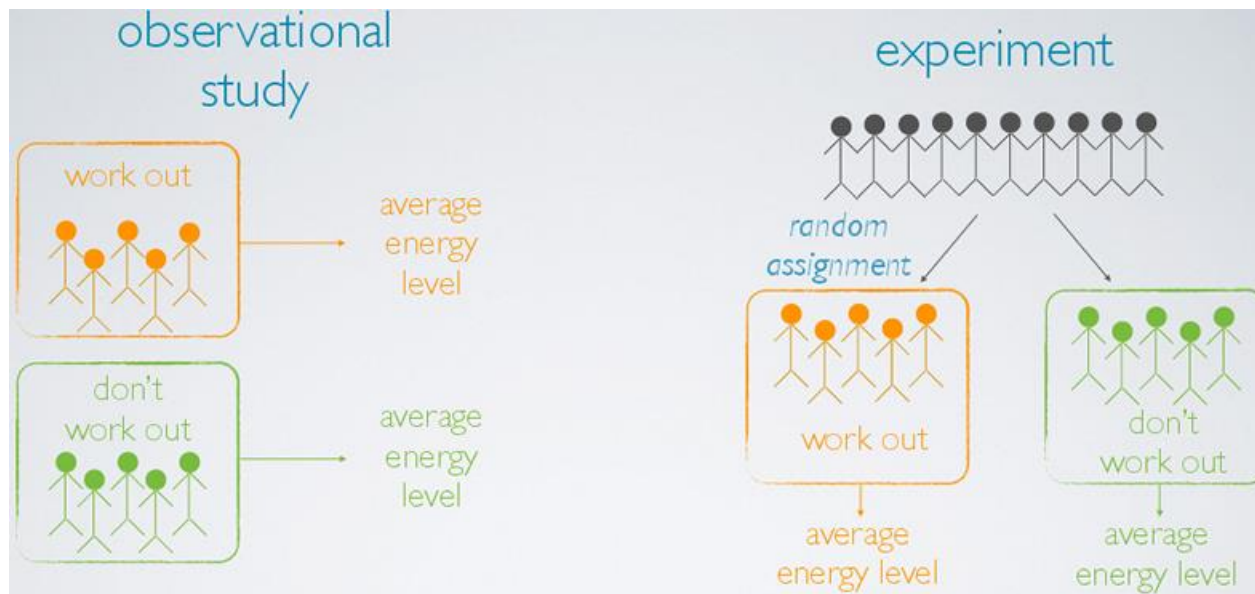


Fuente:  
<https://davidalpiaz.github.io/appliedstats/simple-linear-regression.html>

- ❖ ¿Por qué no usar los puntos correspondientes a un valor de  $x$  para estimar los parámetros de cada curva?

# Supuestos del modelo de regresión lineal simple

- ❖ La lógica detrás de que la variable  $X$  no se asuma como aleatoria viene de asumir que esa variable es un factor controlado en un **estudio experimental**.
- ❖ Sin embargo, también puede asumirse como aleatoria y por ende, se habla de la distribución condicional una vez que se conoce  $X$ . Según esta lógica, aunque  $X$  no es controlable, sí es más fácil de conocer en un **estudio observacional**.
- ❖ En ambos casos, ojo con la **confusión de efectos**.



Fuente:  
<https://researchhubs.com/post/ai/data-analysis-and-statistical-inference/observational-studies-and-experiments-sampling-and-source-bias.html>

# Estimación de los coeficientes de regresión

❖ **Parámetros:**  $\beta_0$  : intercepto,  $\beta_1$  : pendiente

**Método de mínimos cuadrados ordinarios:** Suponga que se tiene el siguiente modelo:

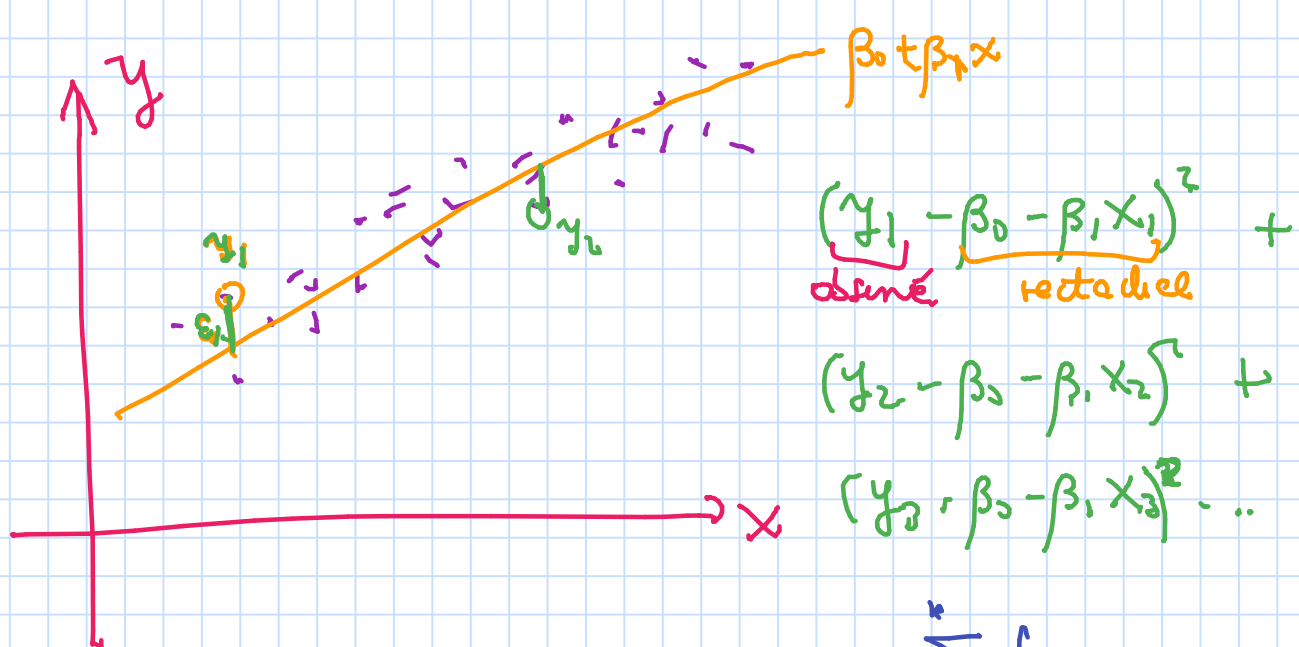
$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i; \quad 1 \leq i \leq n$$

Donde:

1.  $E(\varepsilon_i) = 0 \quad \forall i$
2.  $Var(\varepsilon_i) = \sigma^2 \quad \forall i$
3.  $Cov(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i \neq j$

La estimación de mínimos cuadrados (dados los datos) es:

$$(\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{\beta_0, \beta_1} Q(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$



$$(\hat{\beta}_0, \hat{\beta}_1) = \operatorname{argmin} Q(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial Q}{\partial \beta_0} = \frac{\partial}{\partial \beta_0} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right\}$$

$$= \sum_{i=1}^n \frac{\partial}{\partial \beta_0} \{ (y_i - \beta_0 - \beta_1 x_i)^2 \}$$

$$= \sum_{i=1}^n -2 (y_i - \beta_0 - \beta_1 x_i)$$

$$= -2 \left( \sum_{i=1}^n y_i - \sum_{i=1}^n \beta_0 - \sum_{i=1}^n \beta_1 x_i \right)$$

$$= -2 [n\bar{y} - n\beta_0 - n\beta_1 \bar{x}] = -2n [\bar{y} - \beta_0 - \beta_1 \bar{x}] \quad (1)$$

$$\frac{\partial Q}{\partial \beta_1} = \frac{\partial}{\partial \beta_1} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right\}$$

$$= \sum_{i=1}^n \frac{\partial}{\partial \beta_1} \{ (y_i - \beta_0 - \beta_1 x_i)^2 \}$$

$$= \sum_{i=1}^n 2 (y_i - \beta_0 - \beta_1 x_i) (-x_i) \quad (2)$$

$$= -2 \left( \sum_{i=1}^n y_i x_i - \sum_{i=1}^n \beta_0 x_i - \sum_{i=1}^n \beta_1 x_i^2 \right)$$

$$= -2 \left( \sum_{i=1}^n y_i x_i - \beta_0 n \bar{x} - \beta_1 \sum_{i=1}^n x_i^2 \right) \quad (2')$$

Ignorando a 0 (1) y (2):

$$-2n [\bar{y} - \hat{\beta}_0 - \hat{\beta}_1 \bar{x}] = 0 \quad (1')$$

$$-2 \left( \sum_{i=1}^n y_i x_i - \hat{\beta}_0 n \bar{x} - \hat{\beta}_1 \sum_{i=1}^n x_i^2 \right) = 0 \quad (2')$$

⋮

# Solución del problema de minimización (I)

## 1. Puntos críticos.

$$\left\{ \begin{array}{l} \frac{\partial Q}{\partial \beta_0}(\hat{\beta}_0, \hat{\beta}_1) = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \quad (1) \end{array} \right.$$

$$\left\{ \begin{array}{l} \frac{\partial Q}{\partial \beta_1}(\hat{\beta}_0, \hat{\beta}_1) = -2 \sum_{i=1}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \quad (2) \end{array} \right.$$

$$\left\{ \begin{array}{l} \bar{y} - \hat{\beta}_0 - \hat{\beta}_1 \bar{x} = 0 \therefore \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (1') \end{array} \right.$$

$$\left\{ \begin{array}{l} \sum_{i=1}^n x_i y_i - \hat{\beta}_0 \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0 \quad (2') \end{array} \right.$$

Reemplazando  $\hat{\beta}_0$  de (1') en (2') y despejando, se tiene que

$$\sum_{i=1}^n x_i y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) \sum_{i=1}^n x_i - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0$$

$$\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} - n\hat{\beta}_1 \bar{x}^2 - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0 \therefore \hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

# Solución del problema de minimización (II)

## 2. Mínimo local (y global).

$$\begin{cases} \frac{\partial^2 Q}{\partial \beta_0^2}(\hat{\beta}_0, \hat{\beta}_1) = -2 \sum_{i=1}^n (-1) = 2n \\ \frac{\partial^2 Q}{\partial \beta_1^2}(\hat{\beta}_0, \hat{\beta}_1) = -2 \sum_{i=1}^n x_i (-x_i) = 2 \sum_{i=1}^n x_i^2 \\ \frac{\partial^2 Q}{\partial \beta_0 \partial \beta_1}(\hat{\beta}_0, \hat{\beta}_1) = 2 \sum_{i=1}^n x_i \end{cases} \quad \mathbf{J} = \begin{pmatrix} 2n & 2 \sum_{i=1}^n x_i \\ 2 \sum_{i=1}^n x_i & 2 \sum_{i=1}^n x_i^2 \end{pmatrix}$$

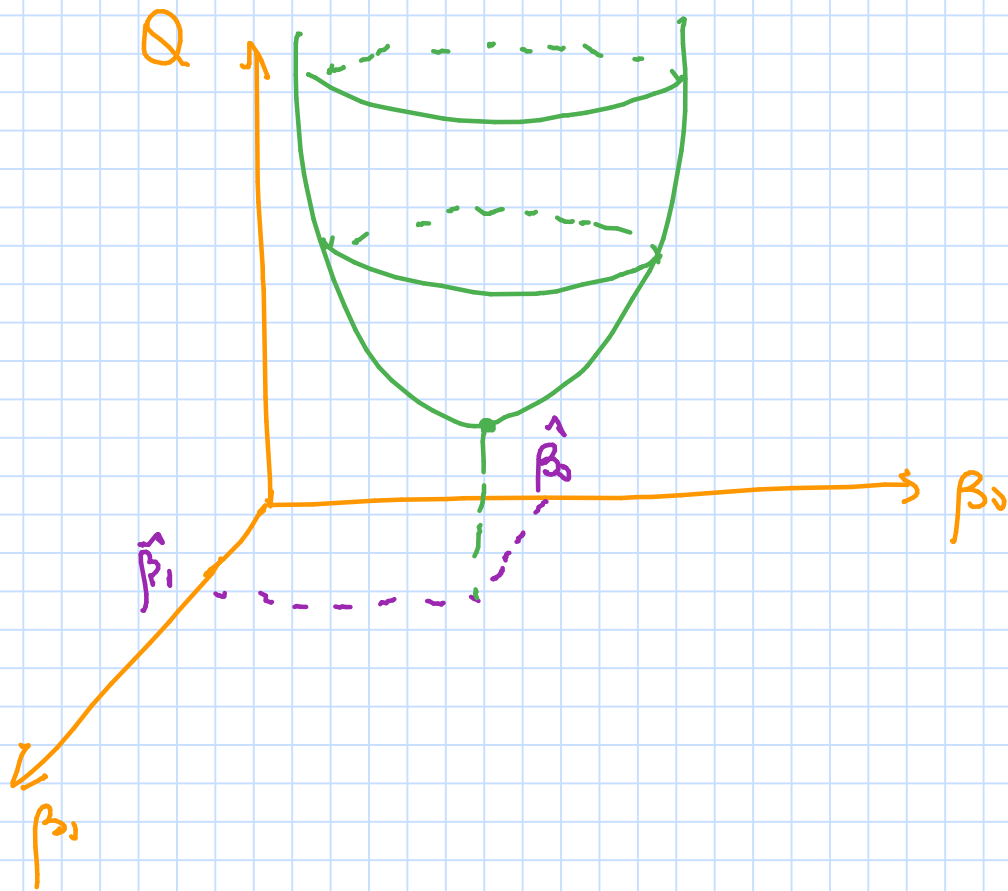
Luego, por el criterio de Sylvester, como

$$2n > 0 \text{ y } \det(\mathbf{J}) = 4n \sum_{i=1}^n x_i^2 - 4n^2 \bar{x}^2 = 4n \sum_{i=1}^n (x_i - \bar{x})^2 > 0;$$

por ende,  $\mathbf{J}$  es positiva definida, lo que implica que la solución es un mínimo local.

Como  $\mathbf{J}$  no depende de  $\hat{\beta}_0, \hat{\beta}_1$ ; quiere decir que  $Q$  es convexa y eso hace que la solución sea global.

No es difícil ver que  $\lim_{\|\beta\| \rightarrow \infty} Q(\beta_0, \beta_1) = \infty$ ; lo que completa la prueba.





# Una definición importante y un lema

## *Mejor estimador linealmente insesgado (BLUE)*

Sean  $Y_1, Y_2, \dots, Y_n$  variables aleatorias involucradas en un proceso de estimación. Se dice que un estimador de la forma  $\sum_{i=1}^n \phi_i Y_i$  con constantes  $\phi_i$  no aleatorias y conocidas es un **BLUE** si

- 1) es un estimador insesgado,
- 2) es el estimador con menor varianza dentro de todos los estimadores lineales insesgados.

❖ Como siempre la definición no es constructiva. Usaremos el lema del máximo (Lema 11.2.7 de Casella & Berger) más adelante para verificar que un estimador es BLUE.