

# Estadística parte 2

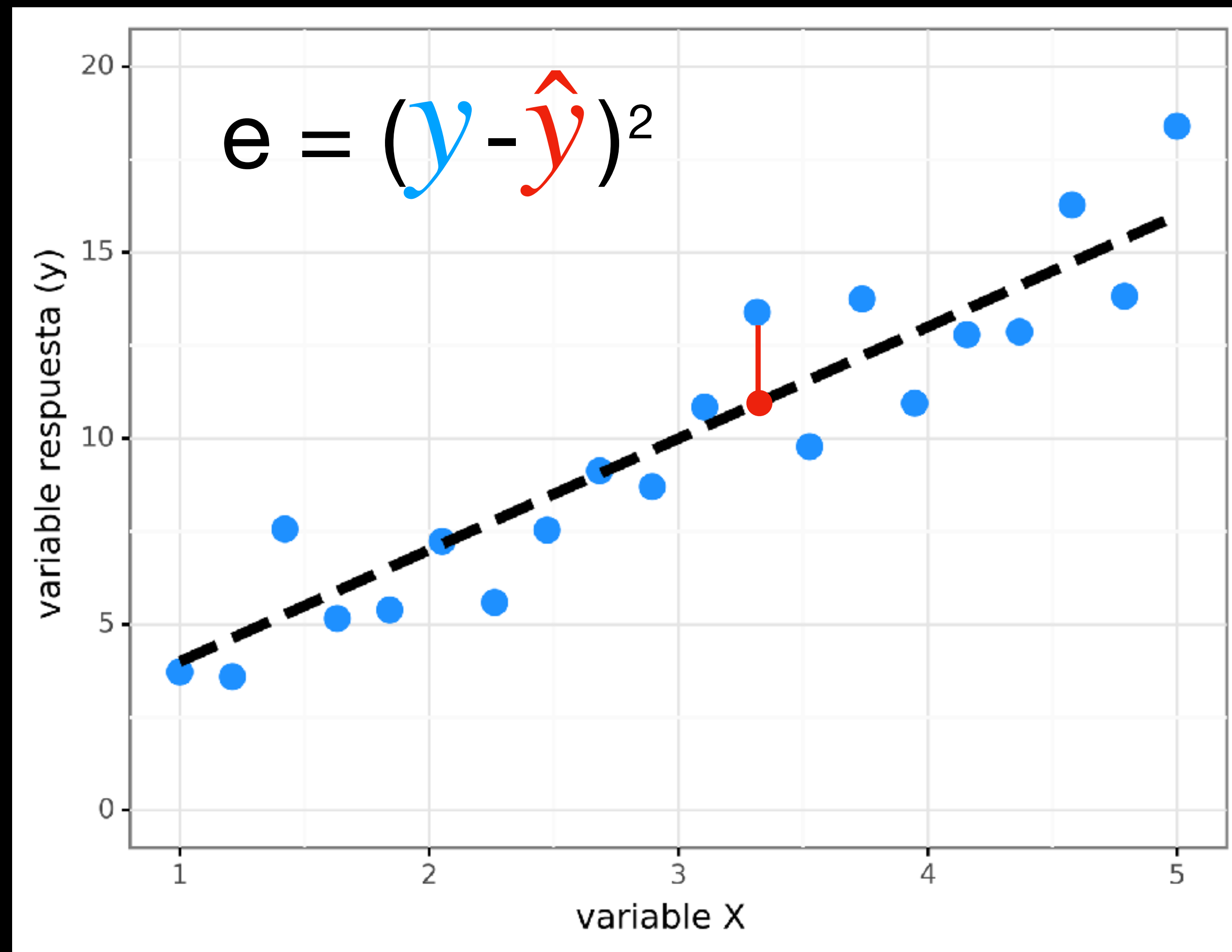
## Curso de Investigación en Radiología

# Contenidos

- Correlación entre variables
- Modelos de regresión lineal.
- Métricas de regresión lineal.
- Modelos de regresión logística.
- Métricas de clasificación.

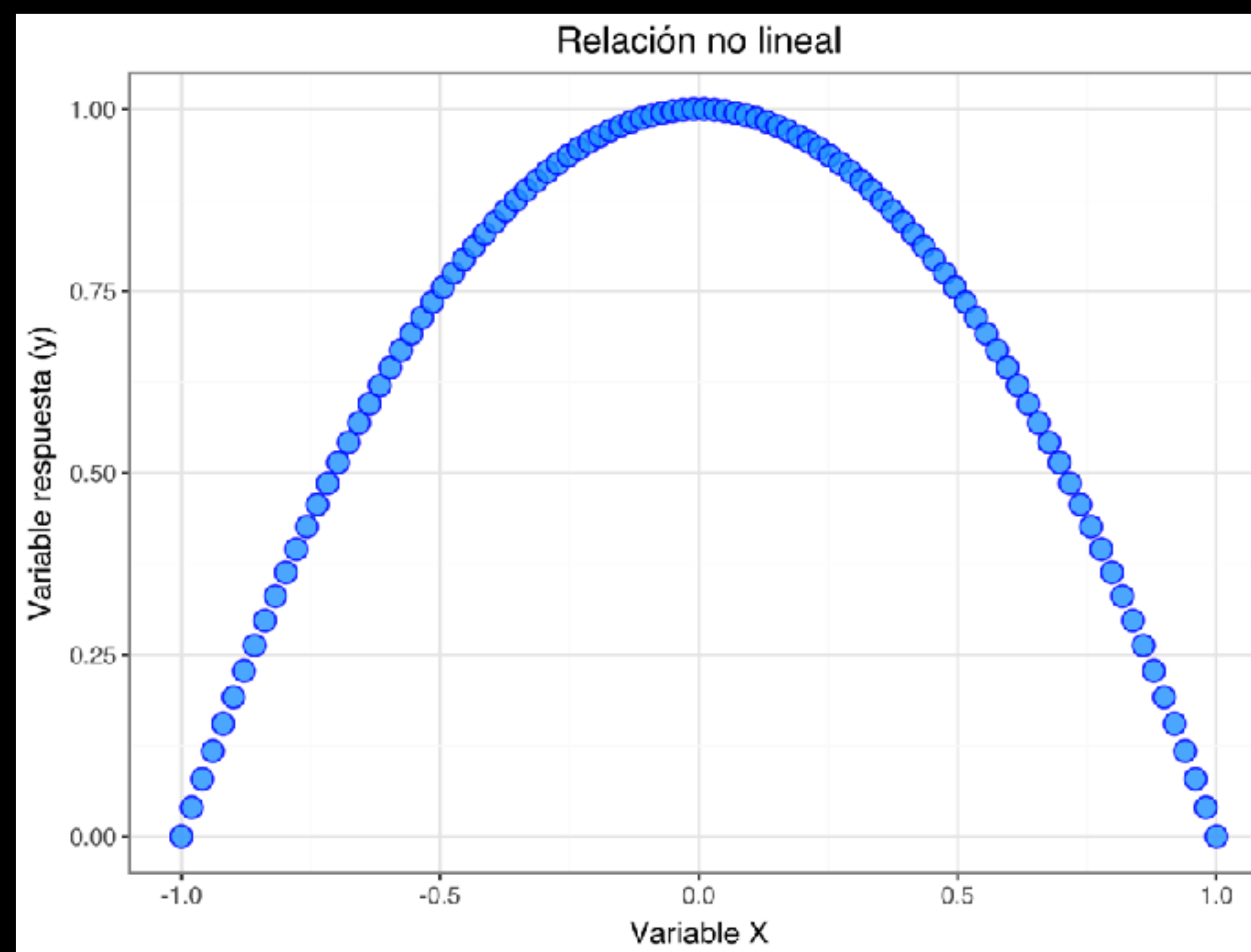
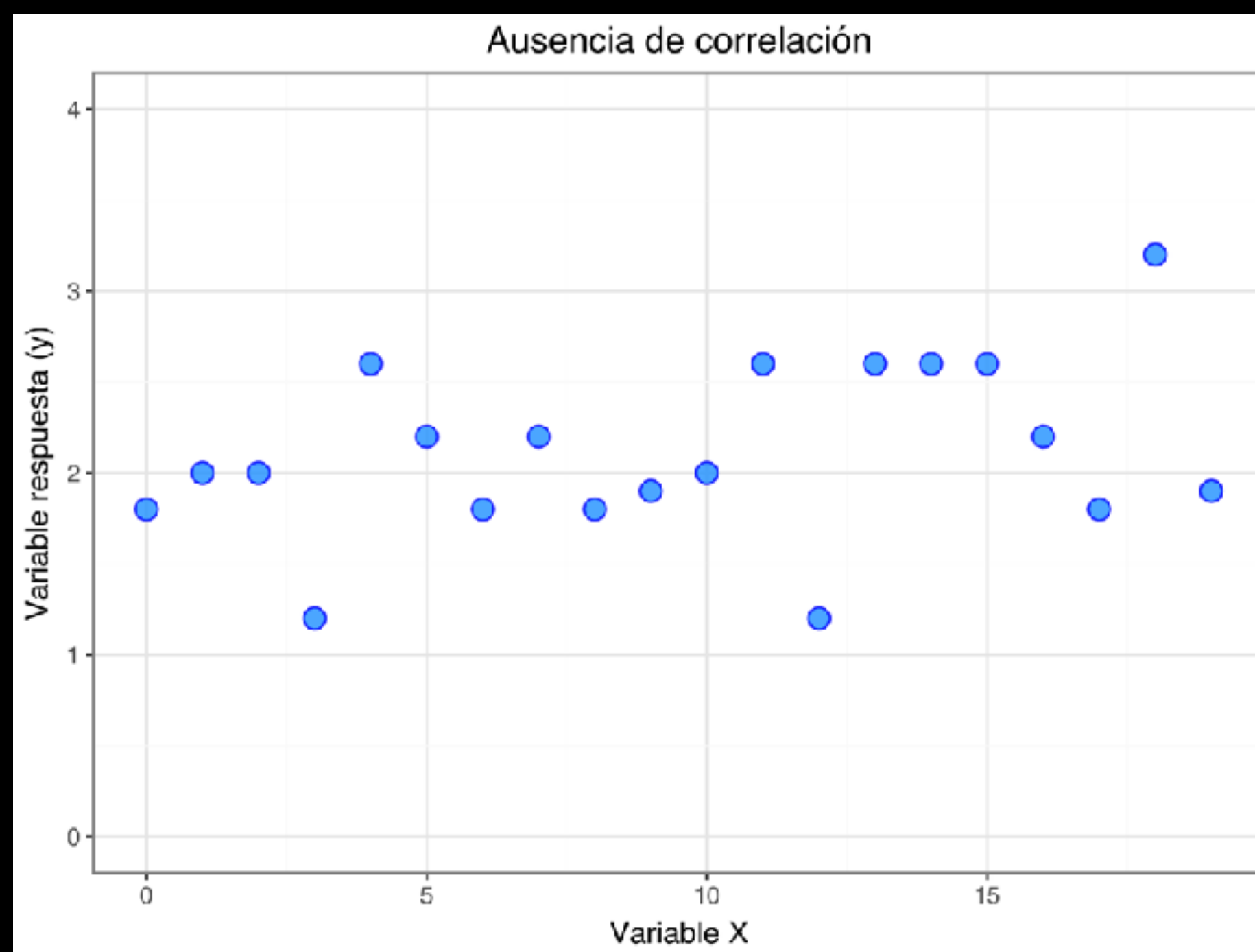
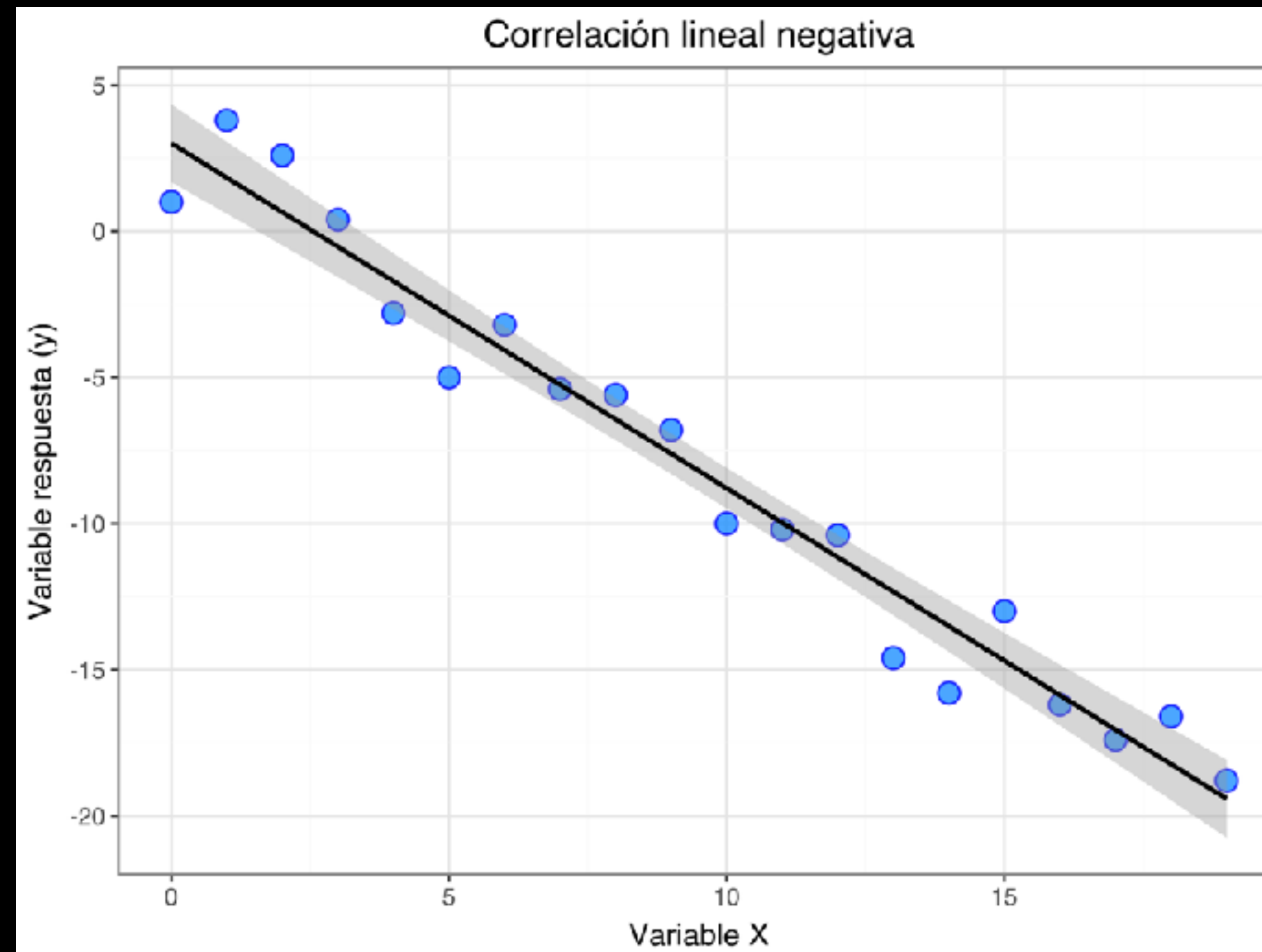
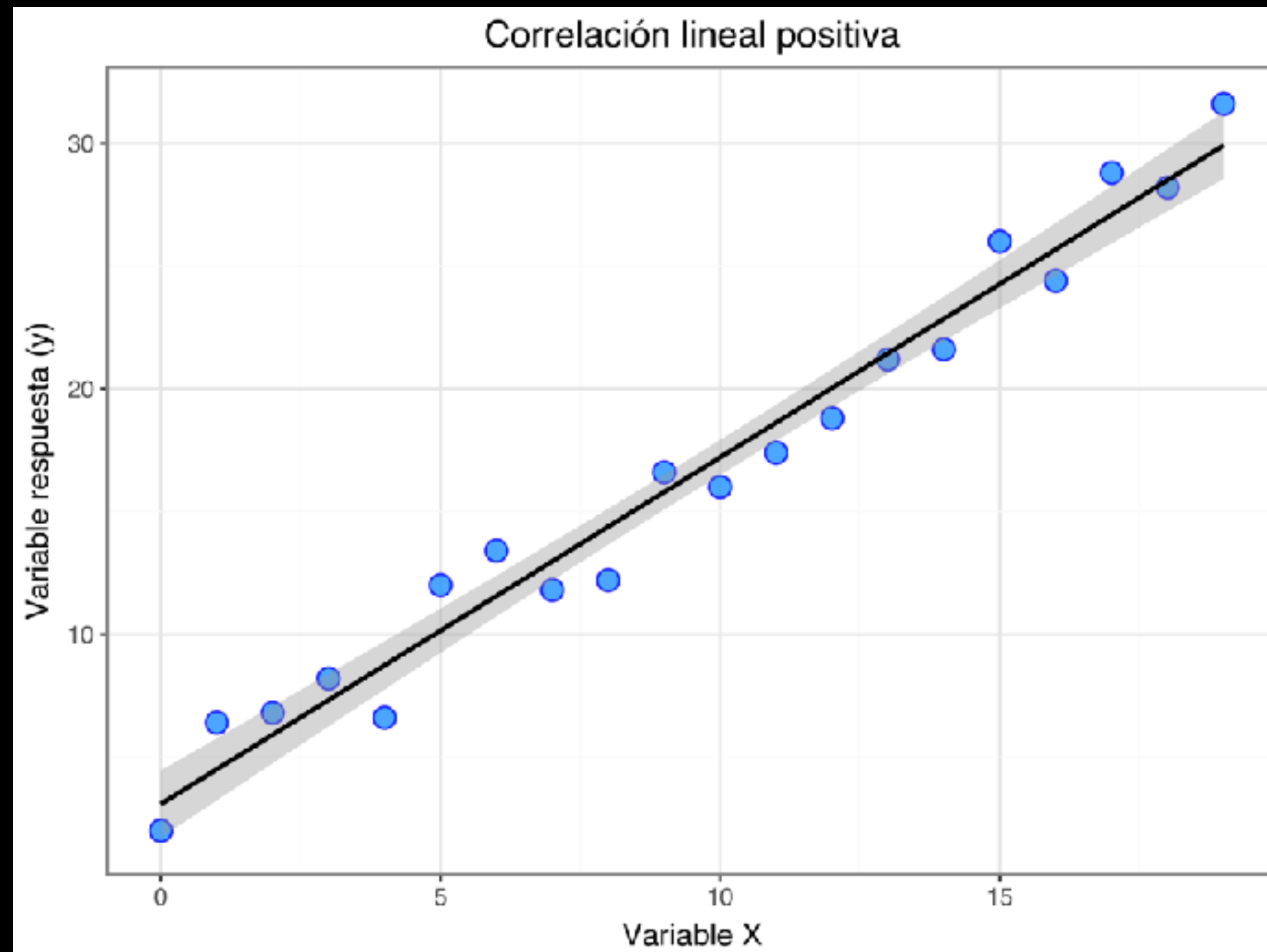
# ABC de modelos predictivos

$$f(x) = y \quad \longrightarrow \quad \hat{y} = \beta_0 + \beta_1 X + \epsilon$$



- *Regresión*
- *Clasificación*





# Medidas de dependencia lineal

- Covarianza:
  - Si hay relación lineal positiva, la covarianza es positiva y grande.
  - Si hay relación lineal negativa, la covarianza es negativa y valor absoluto grande.
  - Si no hay relación entre las variables o marcadamente no lineal, la covarianza es próxima a cero.
  - El problema: covarianza depende de las unidades de medida de las variables.

$$\text{Cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

# Medidas de dependencia lineal

- Coeficiente de correlación lineal:
  - Medida de dependencia lineal que no depende de las unidades de medida de las variables.
  - Valores entre -1 y 1

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

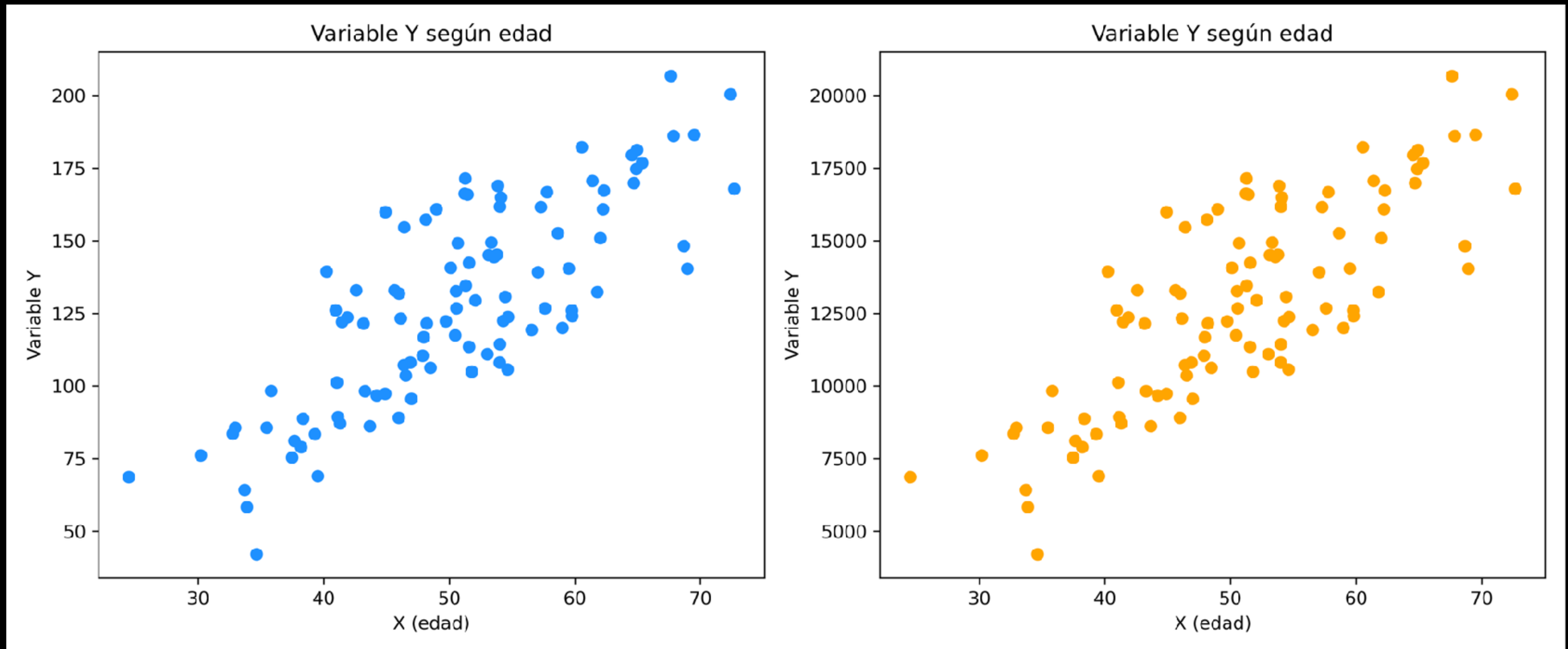
# Medidas de dependencia lineal

Covarianza: 280.05

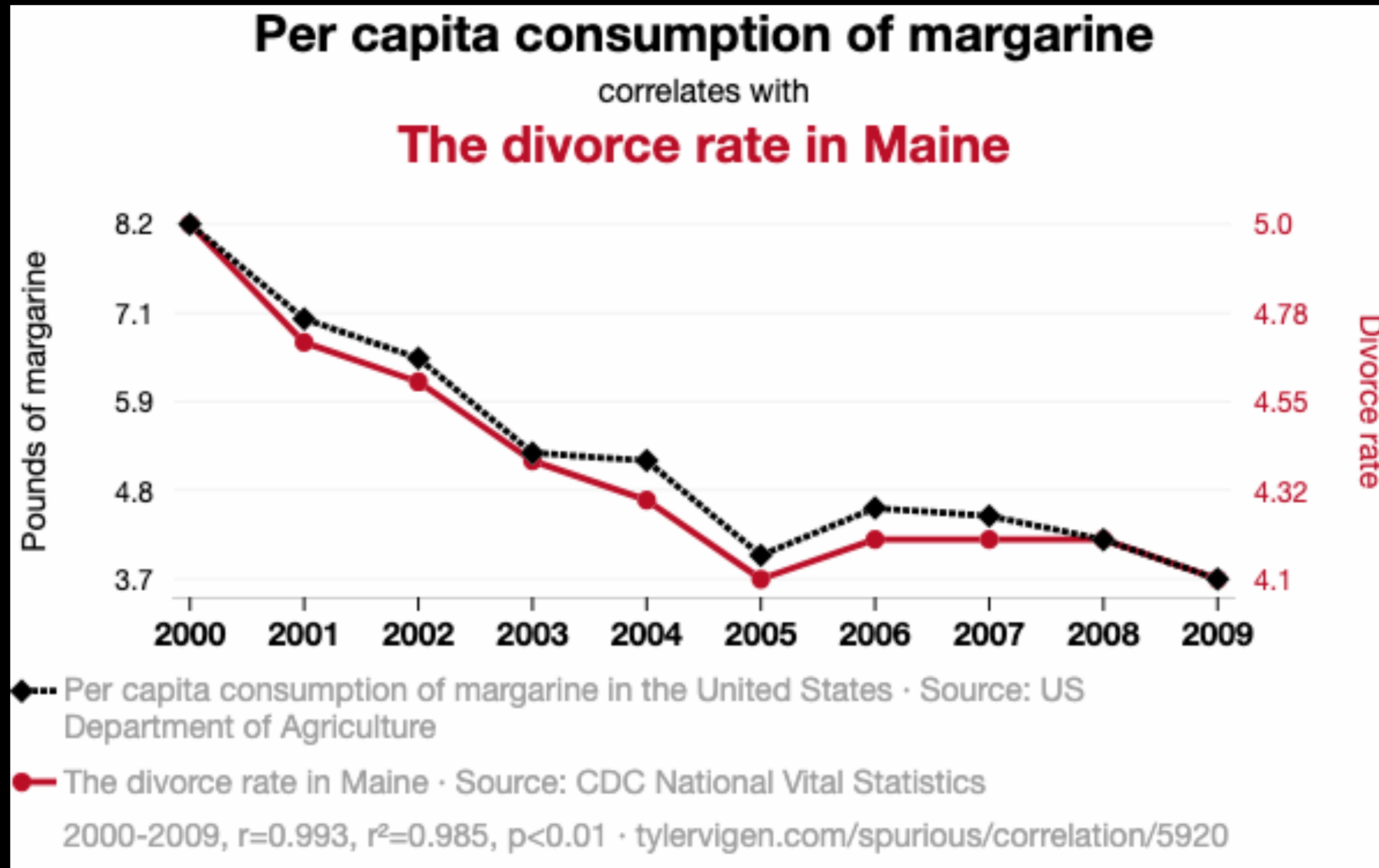
Correlación lineal: 0.80

Covarianza (y escalado): 28005.99

Correlación lineal(y escalado): 0.80

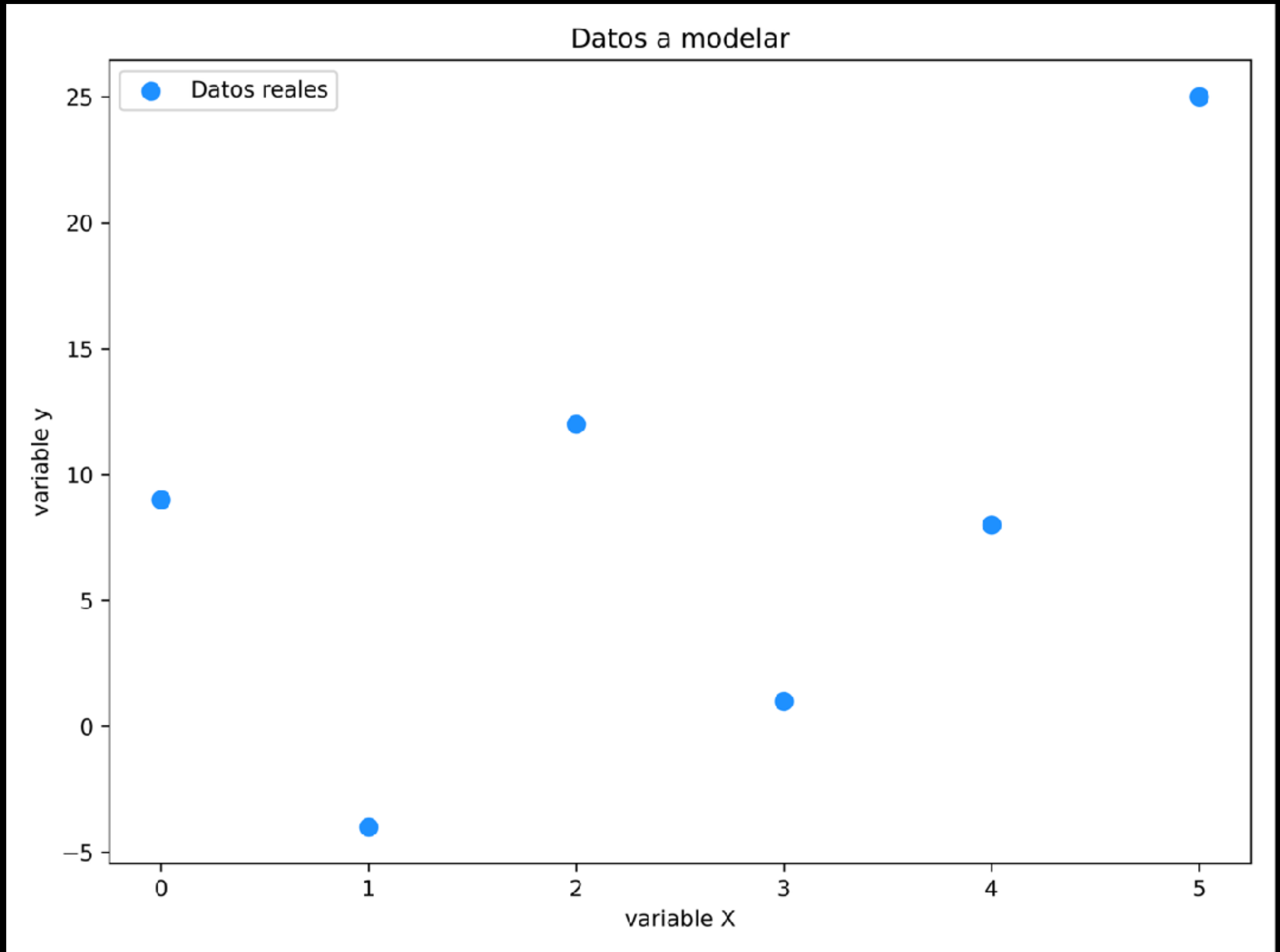


# Correlación no es causalidad





# Regresión lineal simple



# Regresión lineal simple

$$\hat{y} = \beta_0 + \beta_1 X_1 + \epsilon$$

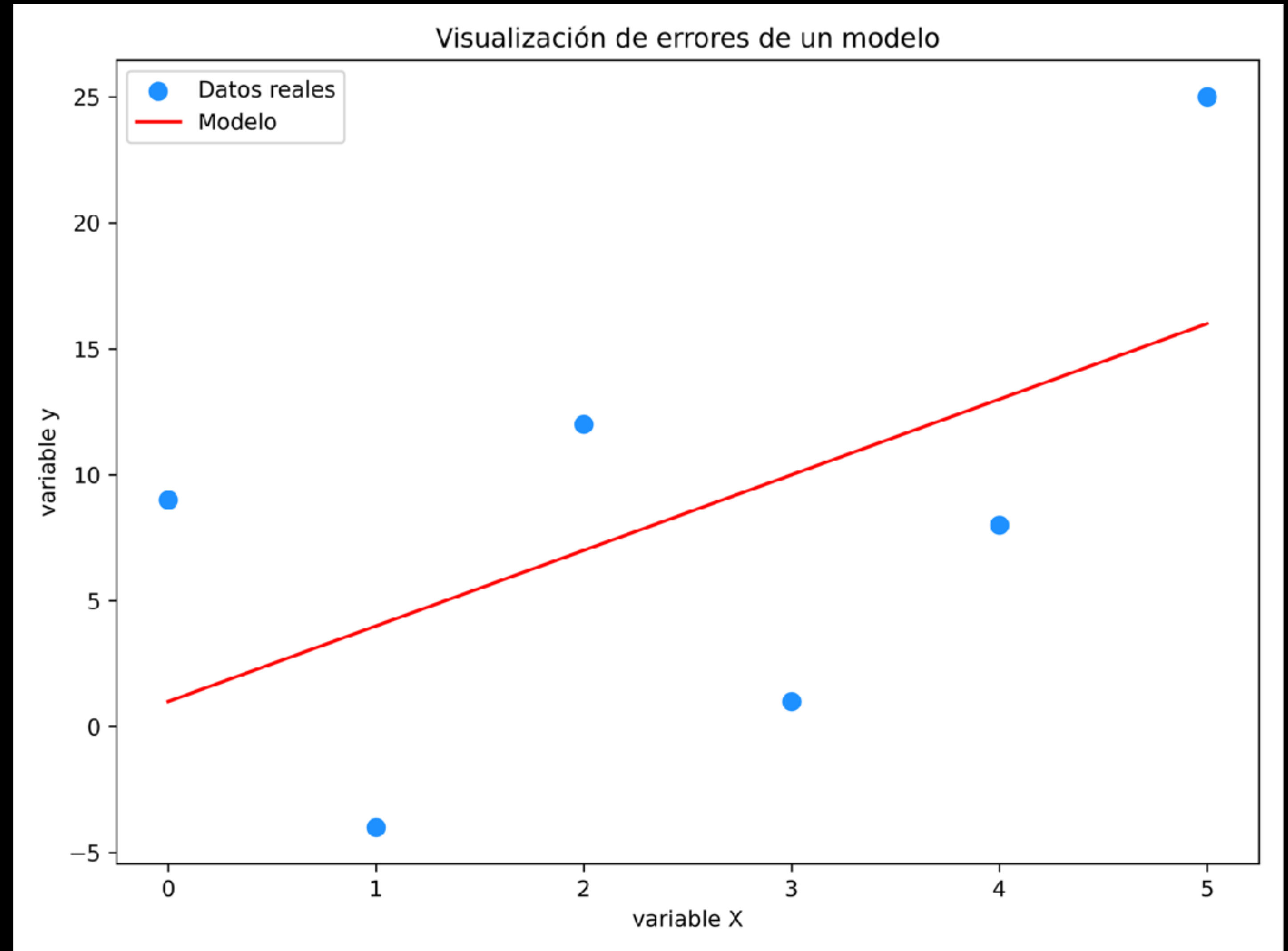
$\hat{y}$  : *Predicción*

$\beta_0$  : *Intercepto*

$\beta_1$  : *Pendiente*

$X_1$  : *Variable independiente*

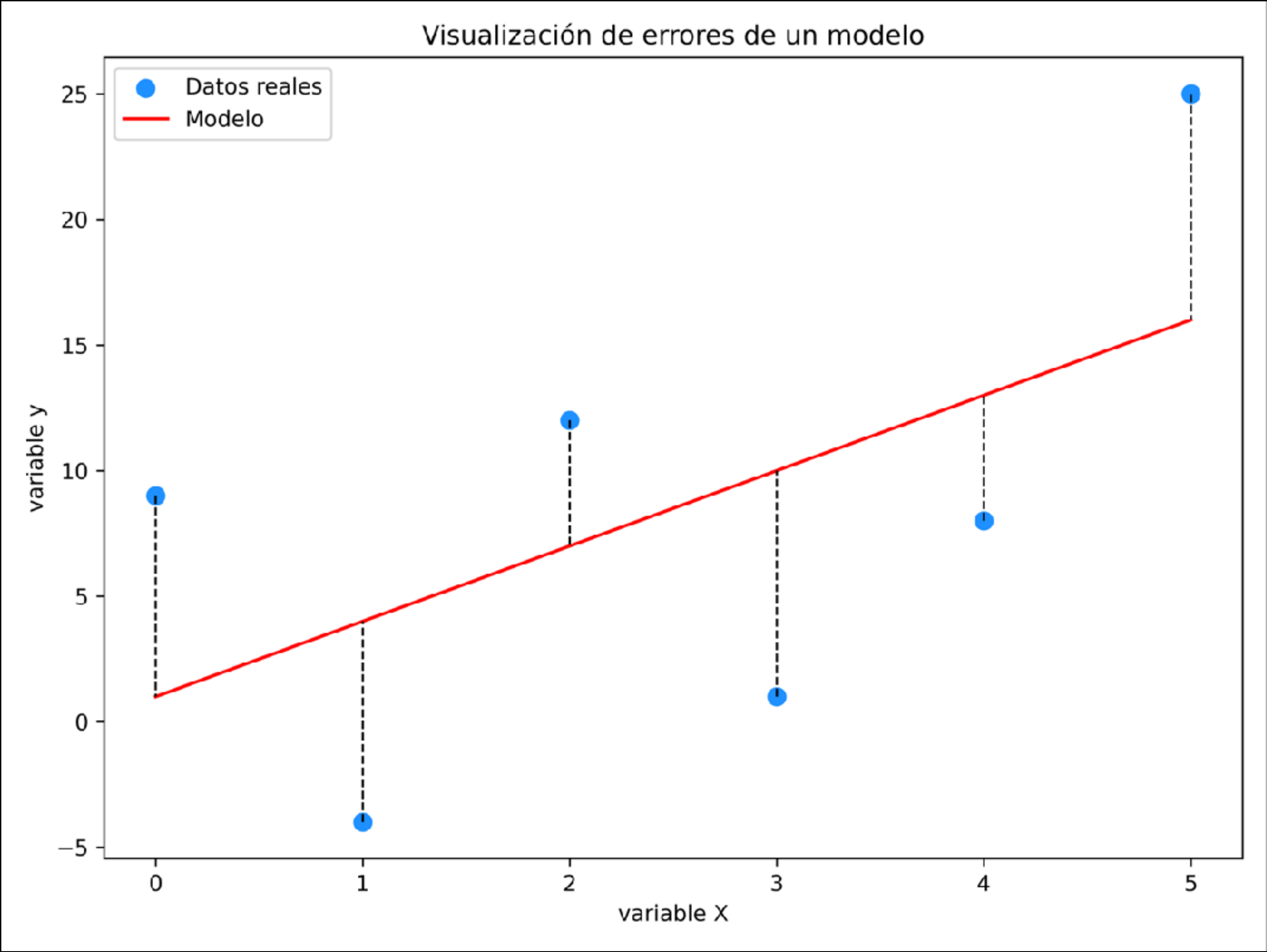
$\epsilon$  : *Ruido*



# Regresión lineal simple

$$\hat{y} = 1 + 3X_1$$

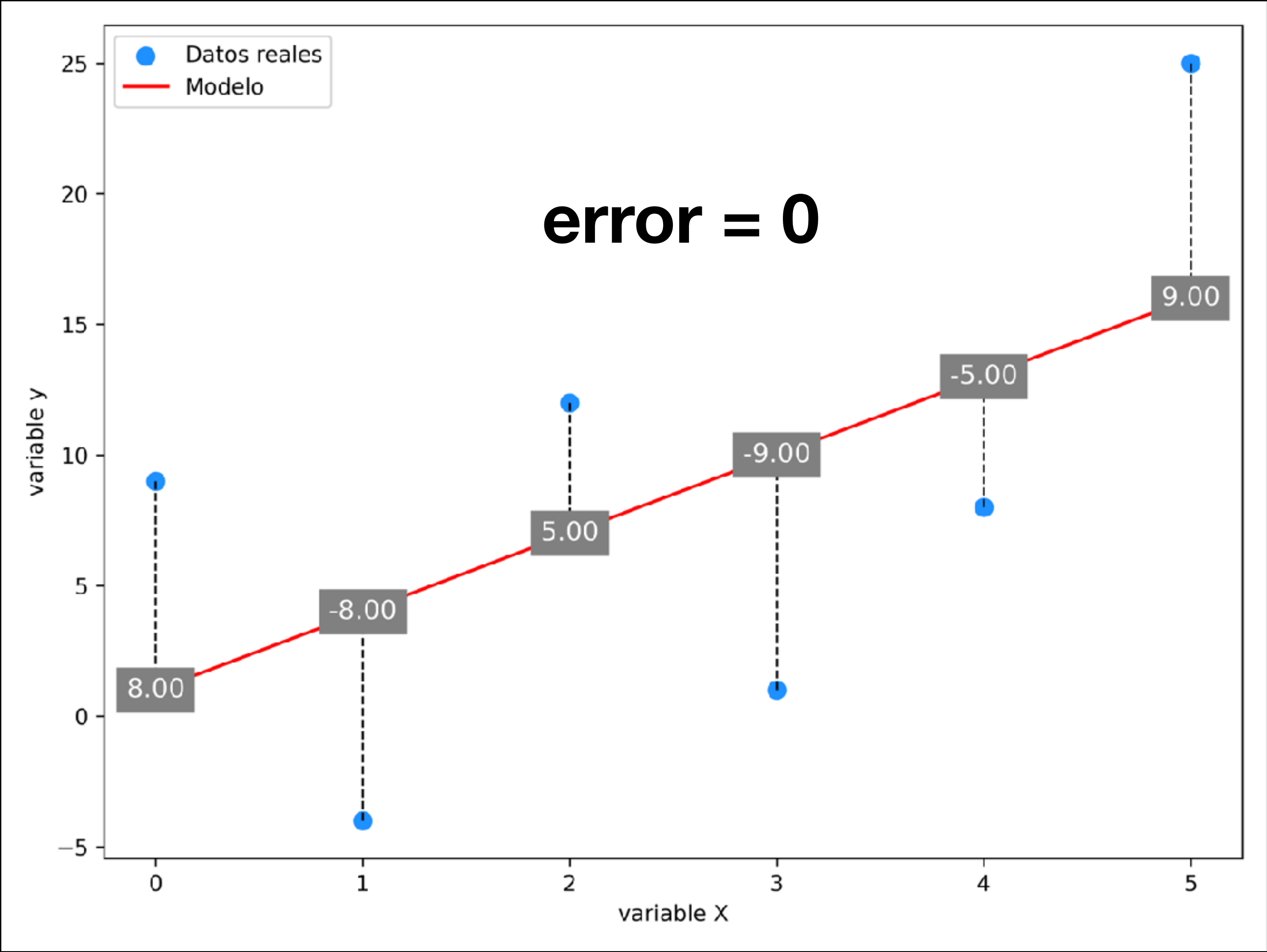
x	y	Pred
0	9	1
1	-4	4
2	12	7
3	1	10
4	8	13
5	25	16



# Regresión lineal simple

$$\hat{y} = 1 + 3X_1$$

x	y	Pred
0	9	1
1	-4	4
2	12	7
3	1	10
4	8	13
5	25	16





$$\hat{y} = 1 + 3X_1$$

Error Medio Absoluto

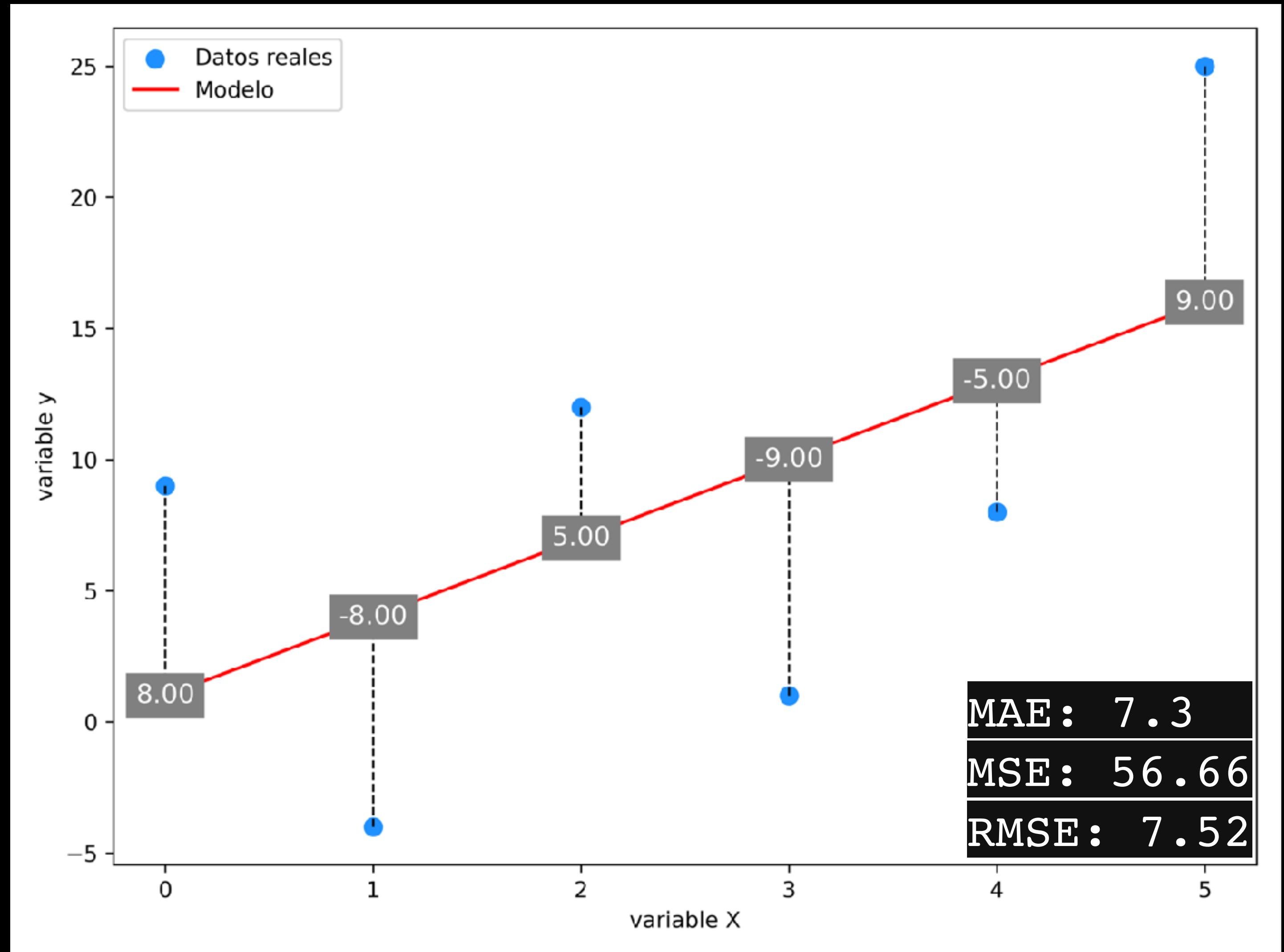
$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Error Cuadrático Medio

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

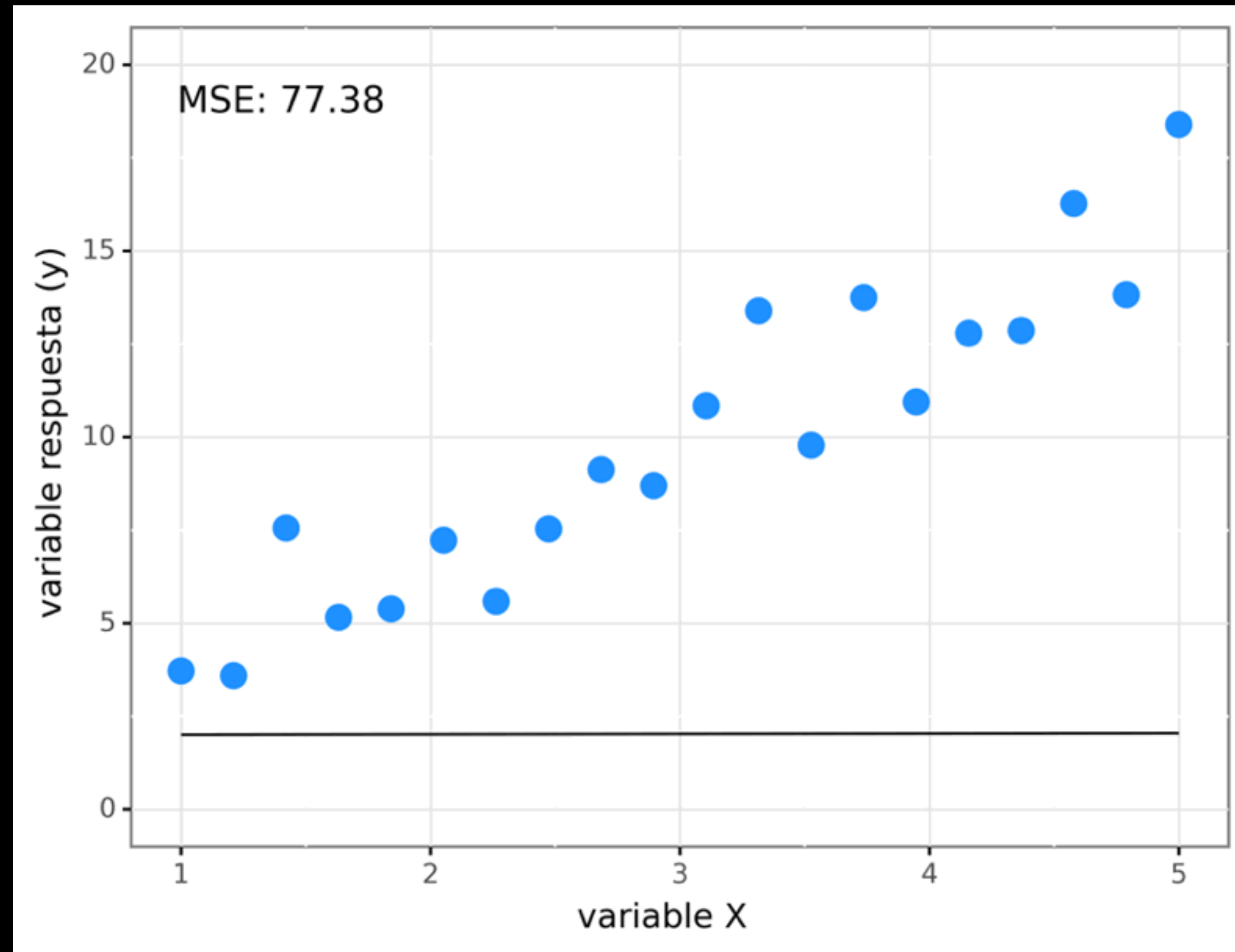
Raíz Cuadrada del Error Cuadrático Medio

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$



# Entrenamiento del modelo

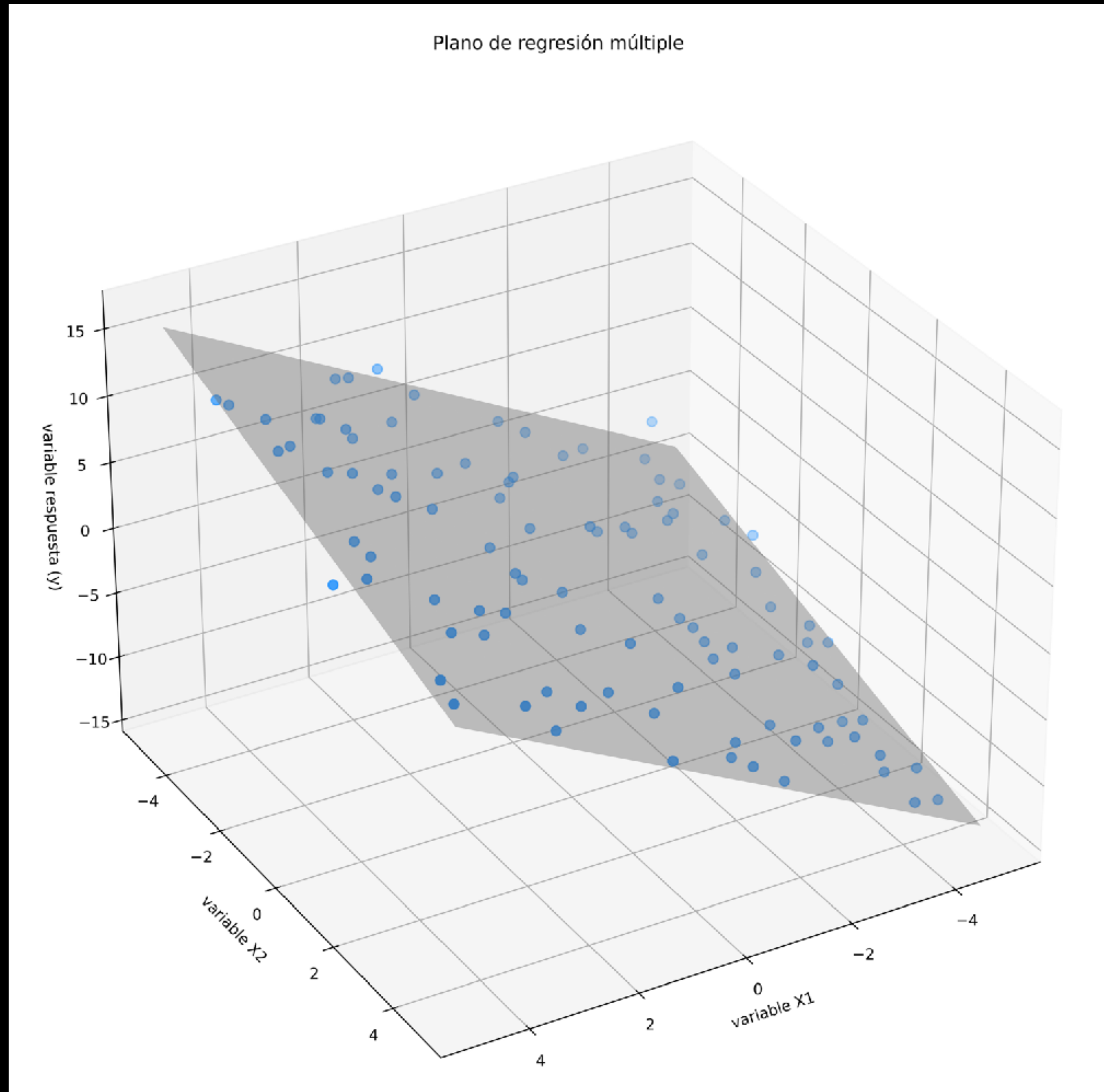
$$f(x) = y \quad \longrightarrow \quad \hat{y} = \beta_0 + \beta_1 X + \epsilon$$



Iteración	MSE
1	77,38
2	
3	
4	
5	
6	
7	
8	
9	
10	

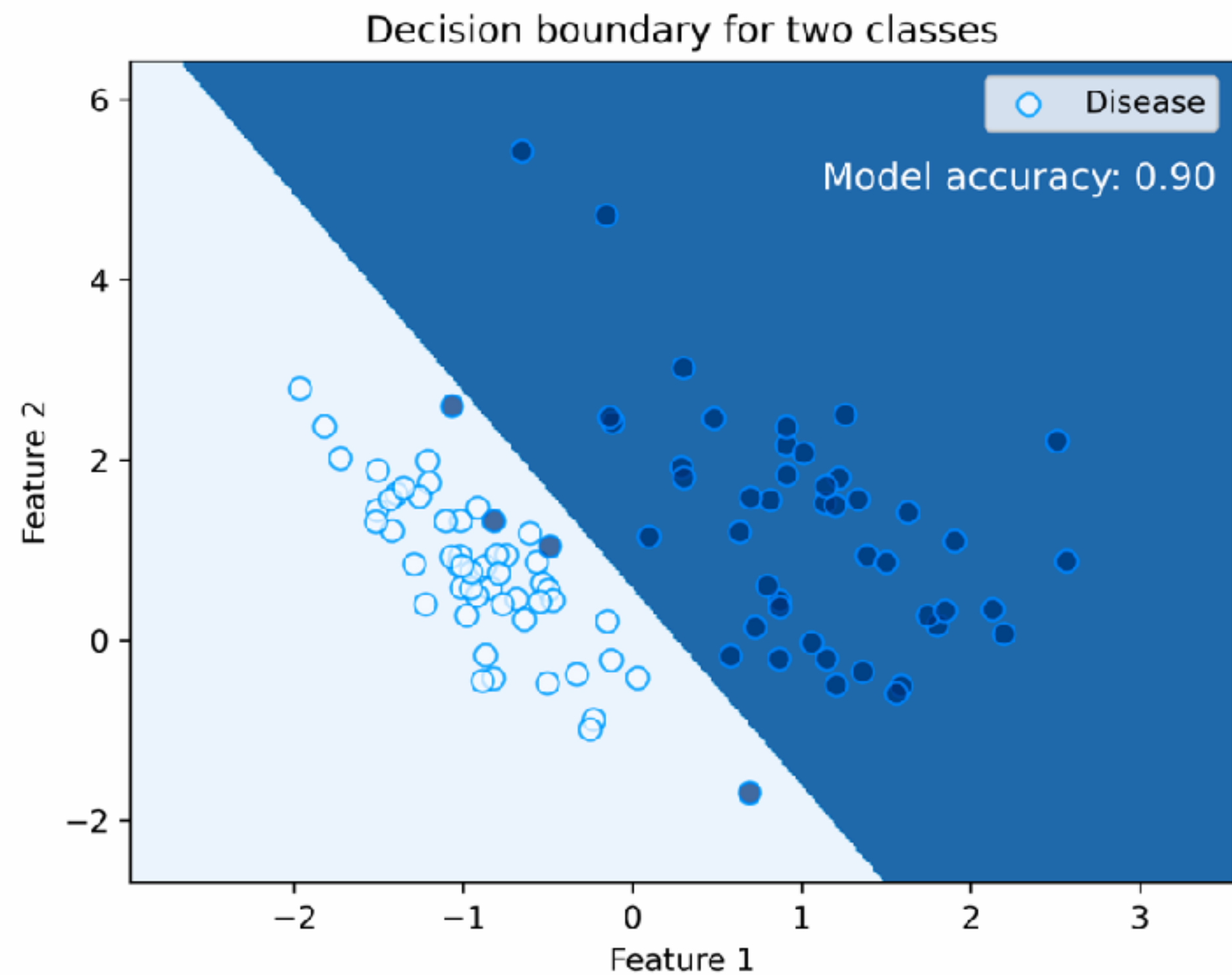
# Regresión lineal múltiple

$$\hat{y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$



# Modelos de clasificación

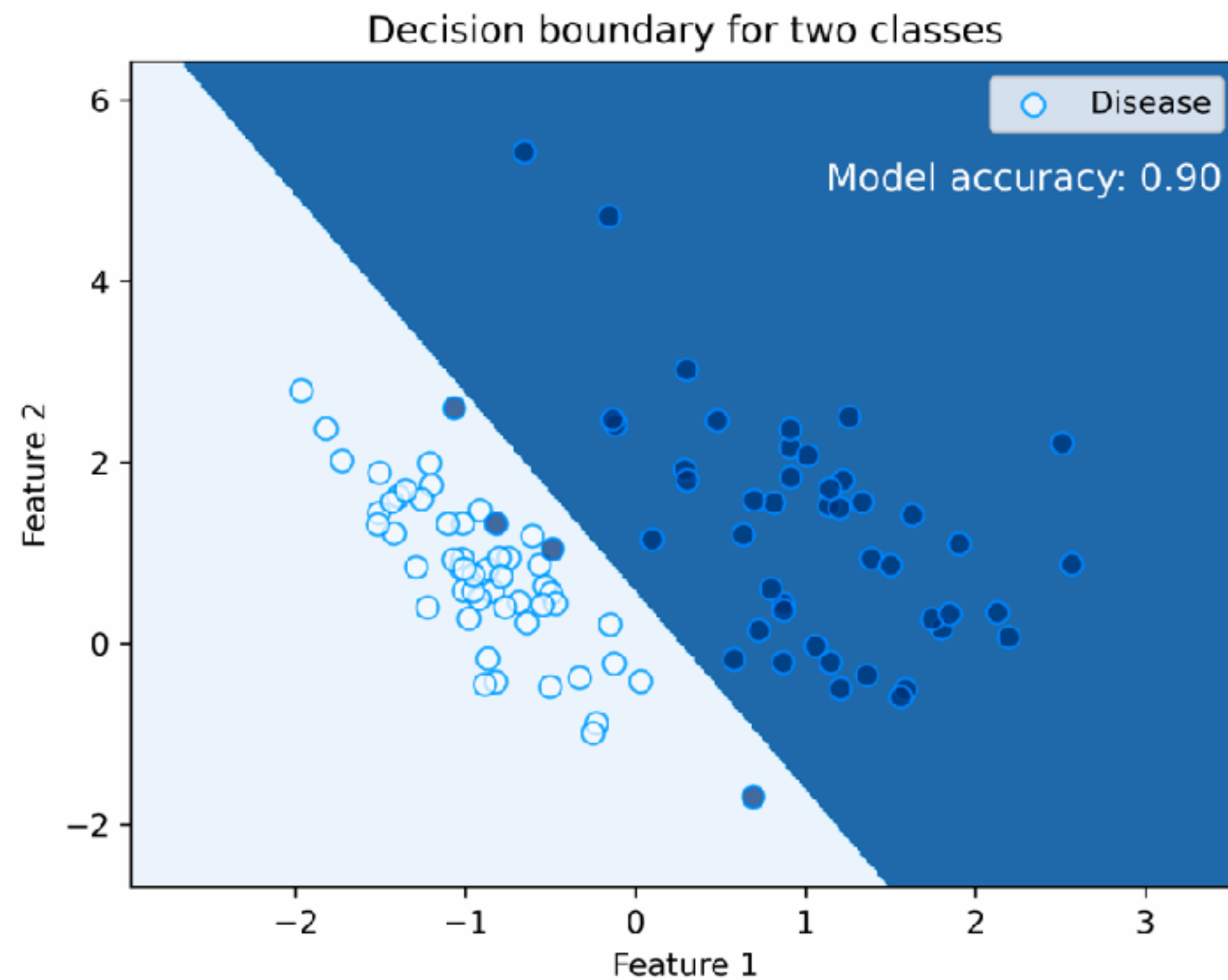
## *Modelo lineal*



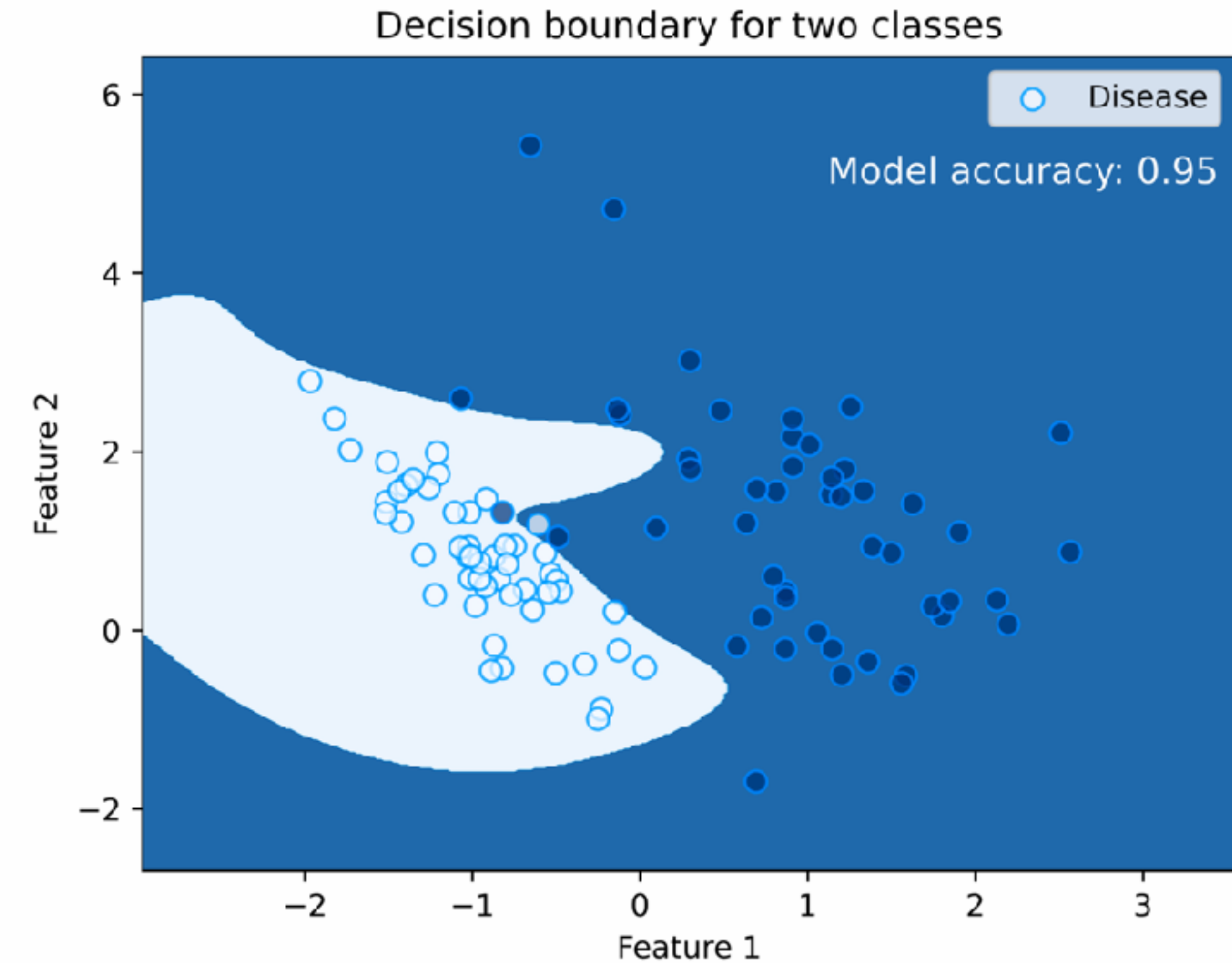


# Modelos de clasificación

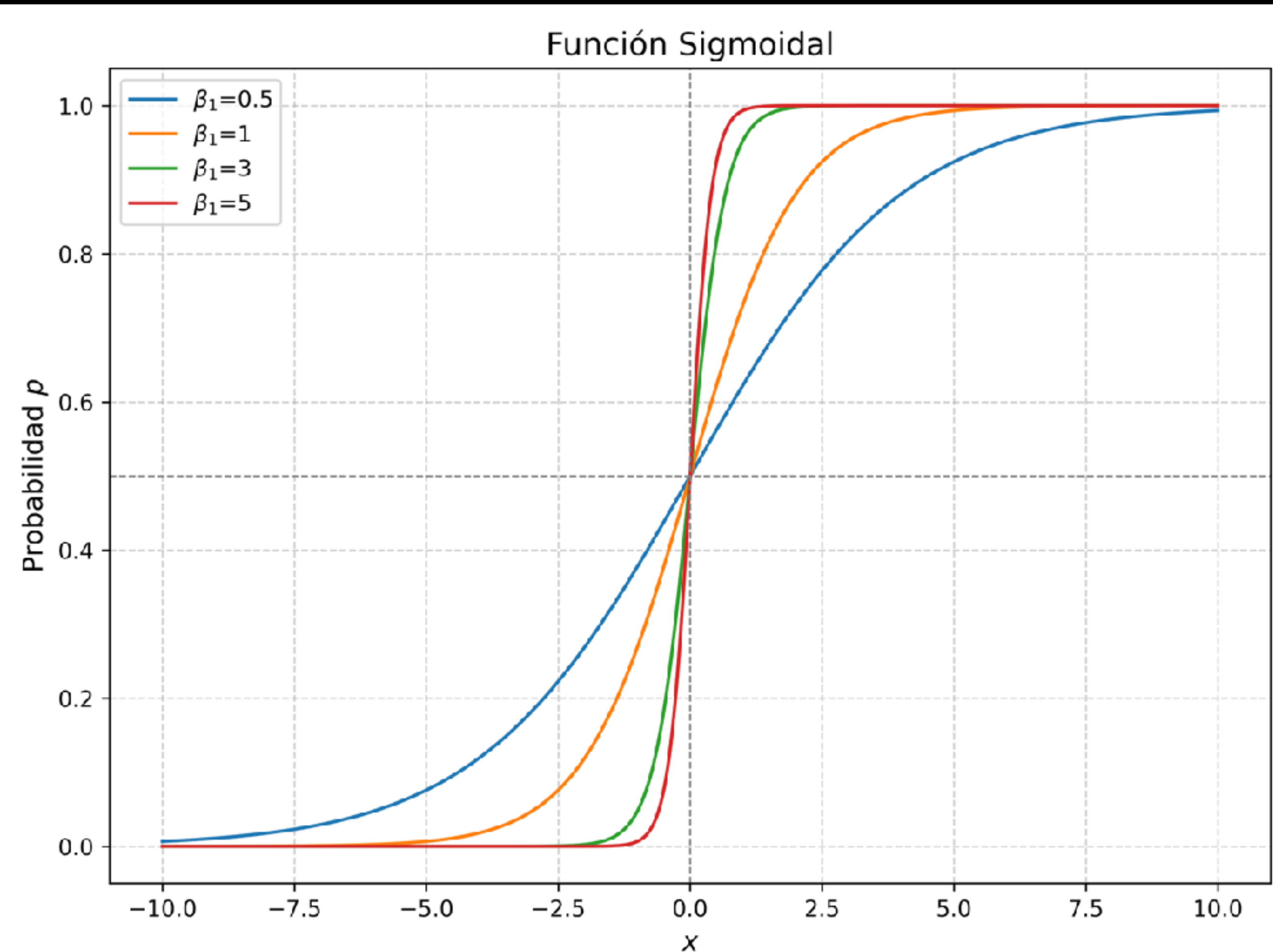
*Modelo lineal*



*Support Vector Machine*



# Regresión logística



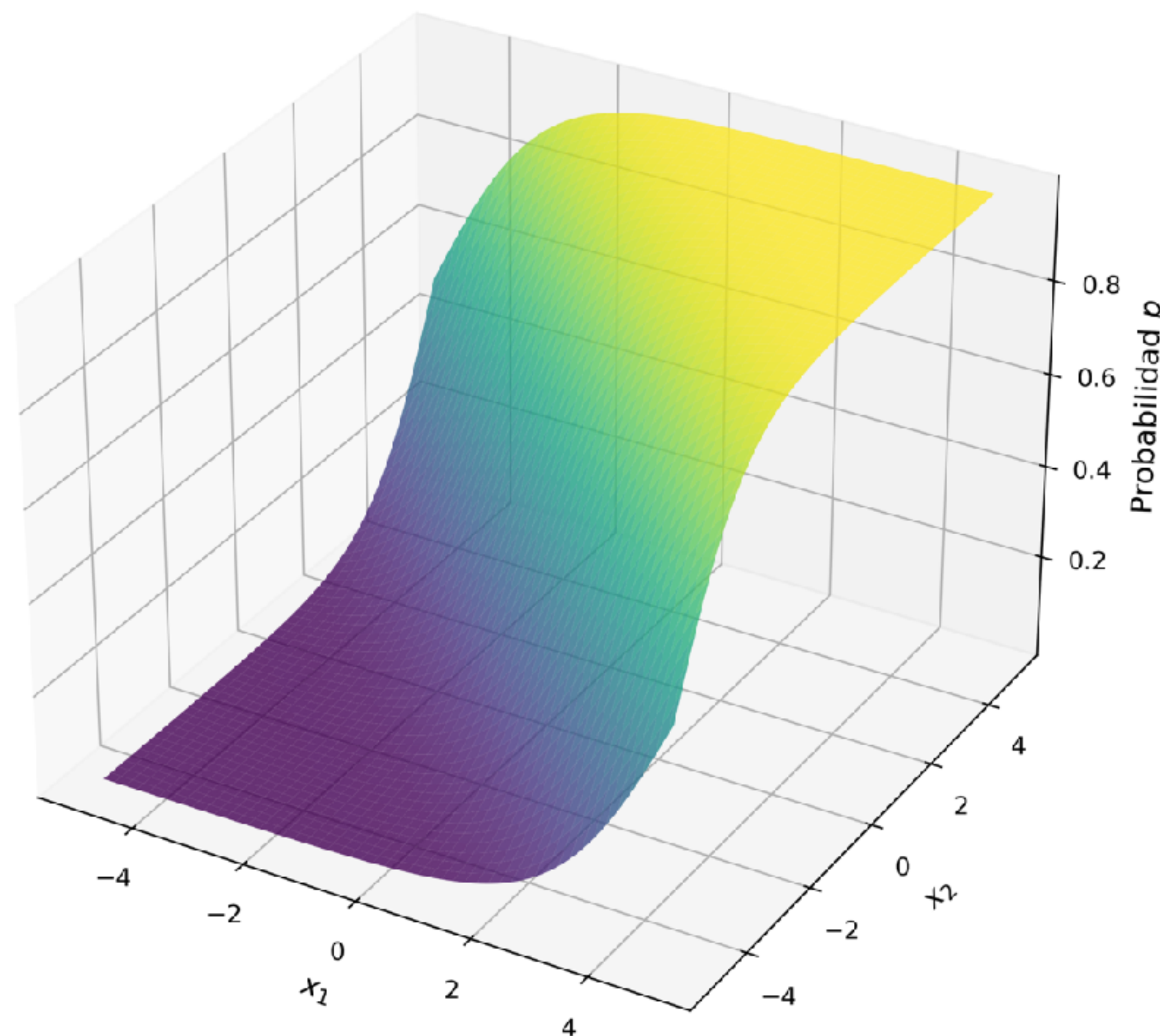
$$\text{logit}(p) = \log \left( \frac{p}{1-p} \right)$$

$$p(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k)}}$$

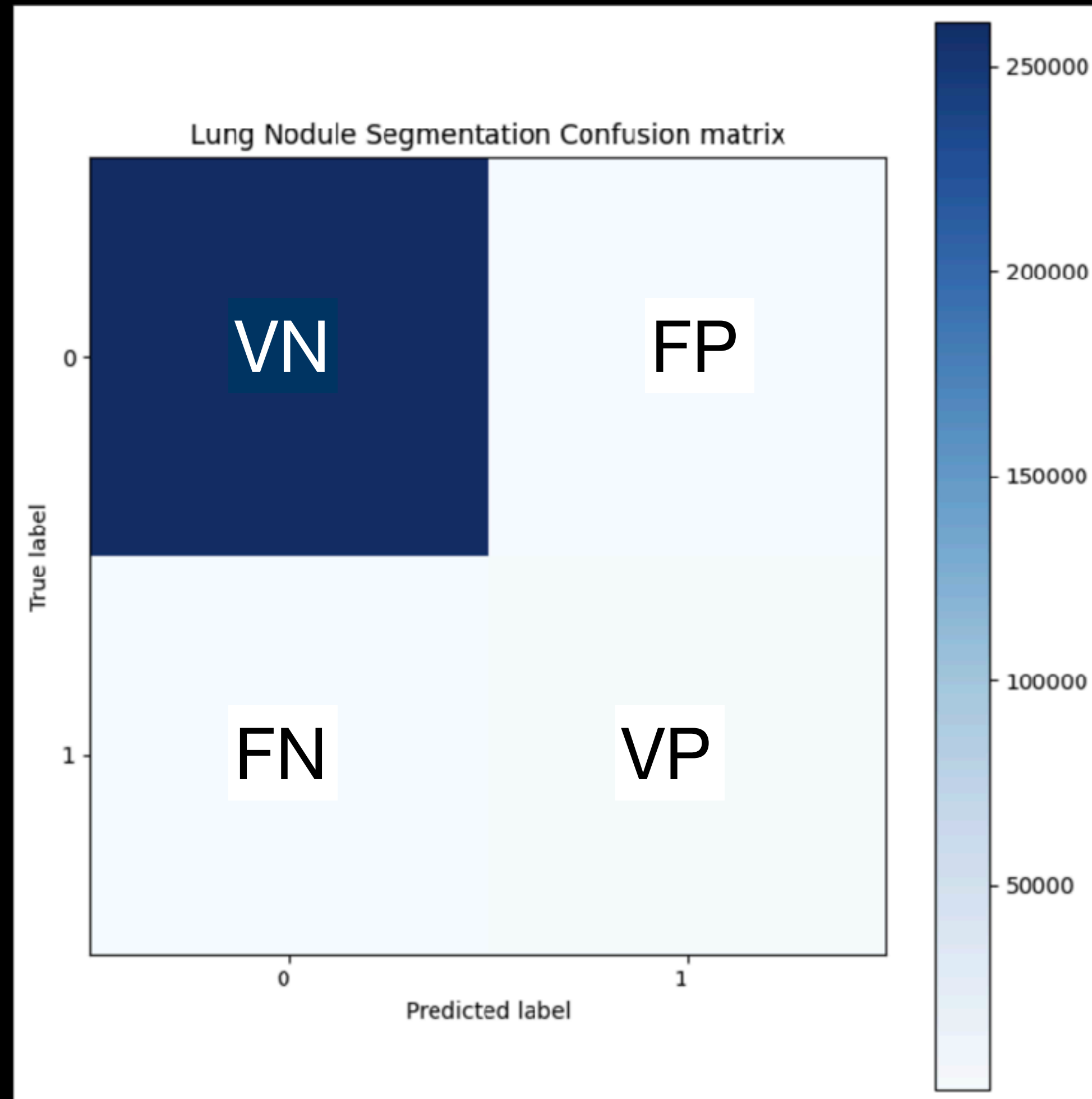
$$\log \left( \frac{p(x)}{1-p(x)} \right) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

# Regresión logística

Plano de la función logística (logit) para dos variables



# Métricas



$$exactitud = \frac{VP + VN}{VP + VN + FN + FP}$$

$$\text{Sensibilidad} = \frac{VP}{VP + FN}$$

(TPR)

$$especificidad = \frac{VN}{VN + FP}$$

(TNR)

$$precision = \frac{VP}{VP + FP}$$

(PPV)

$$F1score = 2 \cdot \frac{\text{Precisión} \cdot \text{Sensibilidad}}{\text{Precisión} + \text{Sensibilidad}}$$