

Regresión lineal

Laboratorio de Datos, IC - FCEN - UBA - 1er. Cuatrimestre 2024

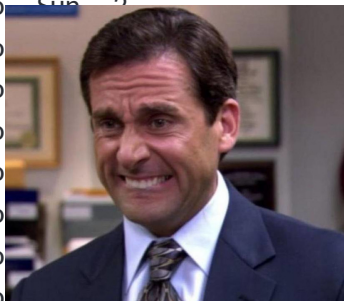
¿Qué modelo uso? 🤔

¿Cómo analizarías estos datos?

0	16.99	1.01	Female	No	Sun	2
1	10.34	1.66	Male	No	Sun	3
2	21.01	3.50	Male	No	Sun	3
3	23.68	3.31	Male	No	Sun	2
4	24.59	3.61	Female	No	Sun	4
5	25.29	4.71	Male	No	Sun	4
6	8.77	2.00	Male	No	Sun	2
7	26.88	3.12	Male	No	Sun	4
8	15.04	1.96	Male	No	Sun	2
:	:	:	:	:	:	:

¿Cómo analizarías estos datos?

0	16.99	1.01	Female	No	Sun	2
1	10.34	1.66	Male	No		
2	21.01	3.50	Male	No		
3	23.68	3.31	Male	No		
4	24.59	3.61	Female	No		
5	25.29	4.71	Male	No		
6	8.77	2.00	Male	No		
7	26.88	3.12	Male	No	Sun	1
8	15.04	1.96	Male	No	Sun	2
:	:	:	:	:	:	:



Primero, **¿qué es un modelo?**

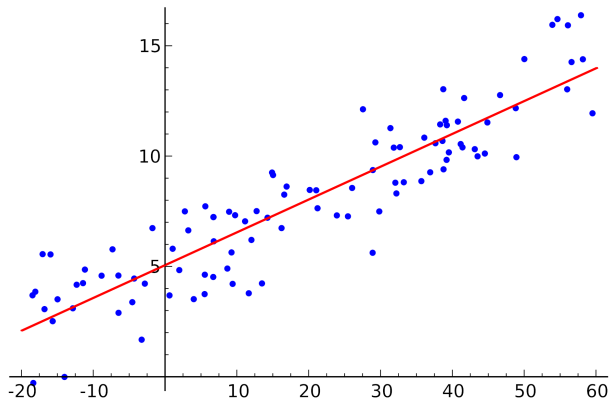
Un modelo es **una representación de fenómenos o procesos del mundo real**. En el contexto de Ciencias de Datos, los modelos son representaciones matemático-computacionales utilizadas para explicar relaciones potencialmente existentes entre las variables de los datos disponibles.

Primero, **¿qué es un modelo?**

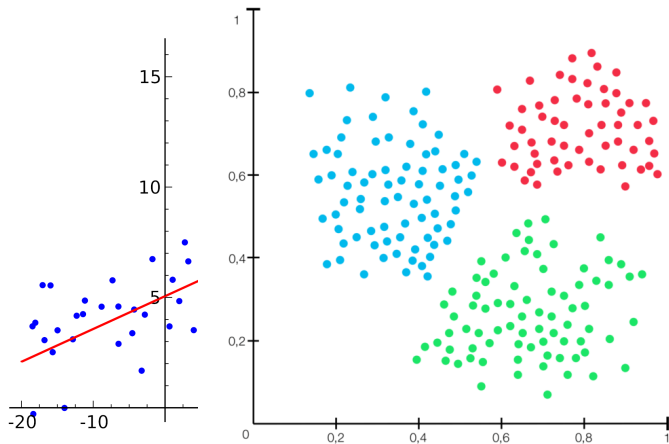
Un modelo es **una representación de fenómenos o procesos del mundo real**. En el contexto de Ciencias de Datos, los modelos son representaciones matemático-computacionales utilizadas para explicar relaciones potencialmente existentes entre las variables de los datos disponibles.

Muchos factores influyen en la elección del modelo

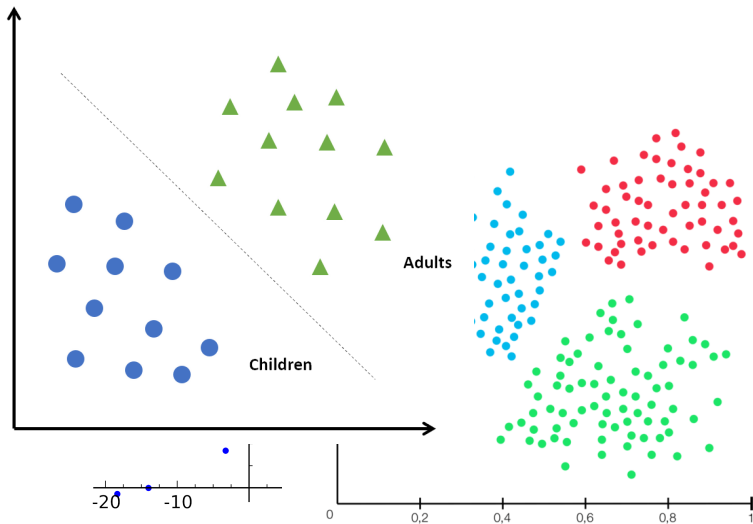
¿Cuál es el problema? ¿Cuál es el objetivo del análisis?



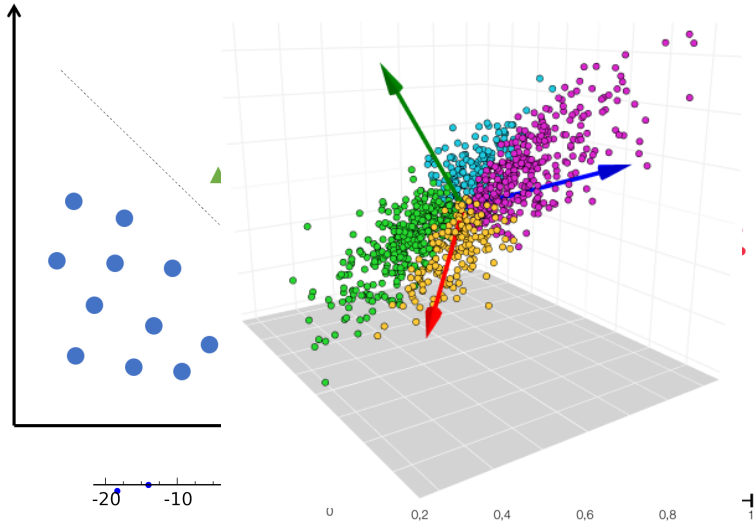
¿Cuál es el problema? ¿Cuál es el objetivo del análisis?



¿Cuál es el problema? ¿Cuál es el objetivo del análisis?



¿Cuál es el problema? ¿Cuál es el objetivo del análisis?



- ¿Cuál es el problema? ¿Cuál es el objetivo del análisis?

- ¿Cuál es el problema? ¿Cuál es el objetivo del análisis?
- ¿Qué tipos de variables tengo?

- ¿Cuál es el problema? ¿Cuál es el objetivo del análisis?
- ¿Qué tipos de variables tengo?
- ¿Cuántos datos tengo?

- ¿Cuál es el problema? ¿Cuál es el objetivo del análisis?
- ¿Qué tipos de variables tengo?
- ¿Cuántos datos tengo?
- ¿Tengo muchos outliers? ¿Qué tan robusto debe ser el modelo?

- ¿Cuál es el problema? ¿Cuál es el objetivo del análisis?
- ¿Qué tipos de variables tengo?
- ¿Cuántos datos tengo?
- ¿Tengo muchos outliers? ¿Qué tan robusto debe ser el modelo?
- ¿Con cuántos recursos computacionales cuento?

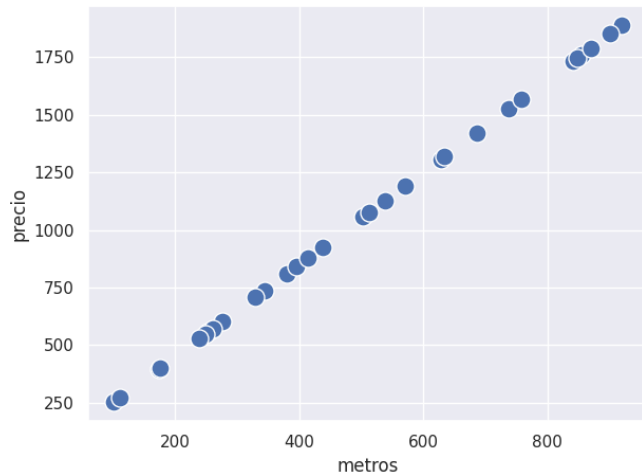
- ¿Cuál es el problema? ¿Cuál es el objetivo del análisis?
- ¿Qué tipos de variables tengo?
- ¿Cuántos datos tengo?
- ¿Tengo muchos outliers? ¿Qué tan robusto debe ser el modelo?
- ¿Con cuántos recursos computacionales cuento?
- ¿Es importante poder entender cómo el modelo toma decisiones?

Regresión

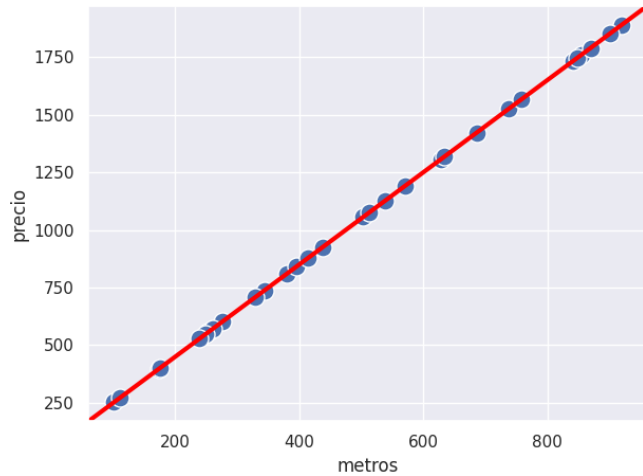
Regresión

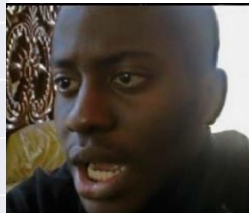
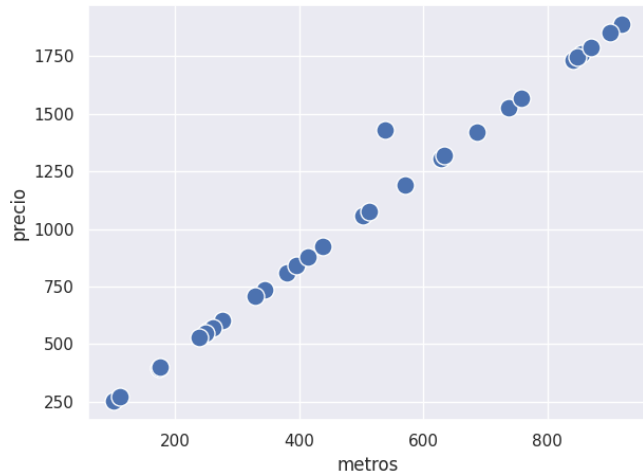
Queremos utilizar los datos que tenemos para poder **estimar datos que no conocemos** o **predecir observaciones futuras**. Los valores a predecir son valores **numéricos**, más precisamente, continuos.

Ejemplo: predecir el valor de un inmueble a partir de su tamaño (m^2)

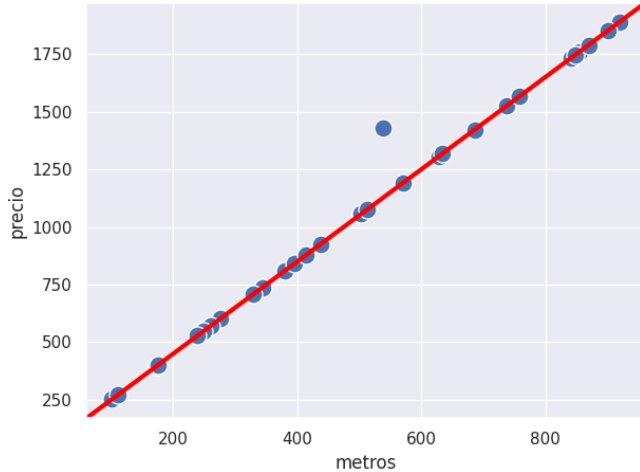


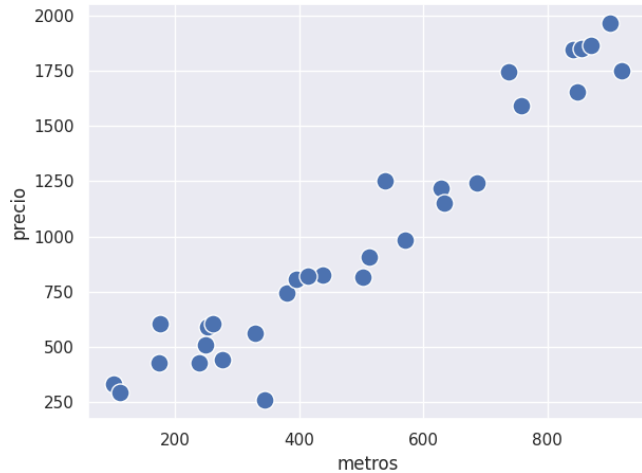
$$y = 2x + 50$$





$$y = 2x + 50$$





¿Qué hacemos ahora?

Regresión Lineal

Regresión

Lineal

Queremos predecir
un valor continuo **con una recta**

Modelo matemático: $Y = \beta_0 + \beta_1 X$

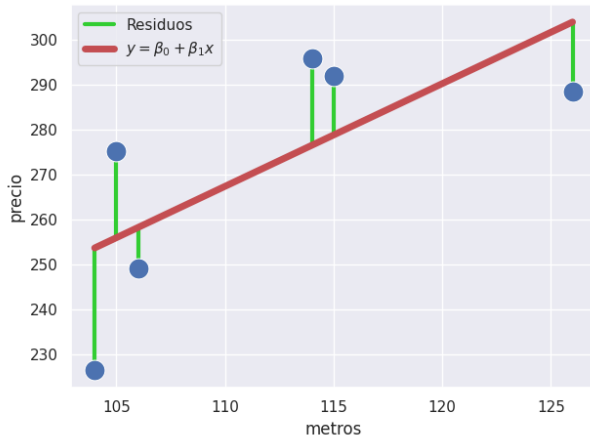
- β_0 es la ordenada al origen
- β_1 es la pendiente
- X es la variable predictora
- Y es la variable dependiente

Modelo de regresión lineal (simple): $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

- (x_i, y_i) son los datos observados
- ε_i es un error aleatorio (variación de Y no explicada por X)
- β_0 y β_1 son **los parámetros del modelo**

Residuo: dados β_0 y β_1 , definimos al residuo como la diferencia entre el valor observado (y_i) y el valor predicho (\hat{y}_i):

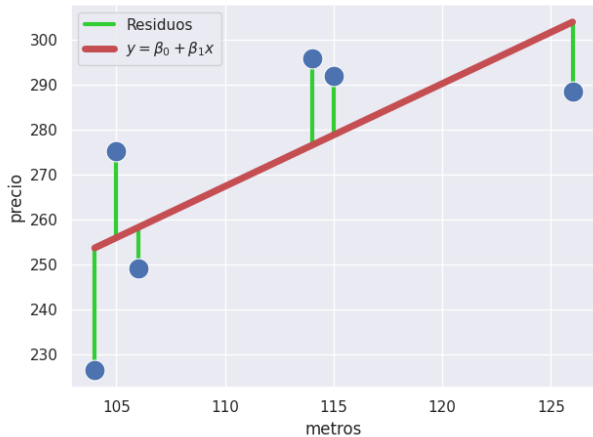
$$y_i - \underbrace{(\beta_0 + \beta_1 x_i)}_{\hat{y}_i}$$



Residuo: dados β_0 y β_1 , definimos al residuo como la diferencia entre el valor observado (y_i) y el valor predicho (\hat{y}_i):

$$y_i - \underbrace{(\beta_0 + \beta_1 x_i)}_{\hat{y}_i}$$

→ Queremos encontrar valores para β_0 y β_1 que minimicen los residuos



Cuadrados Mínimos

Cuadrados Mínimos

Minimizar la suma de los residuos al cuadrado:

$$\begin{aligned}RSS(\beta_0, \beta_1) &= (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \cdots + (y_n - \hat{y}_n)^2 = \\&= (y_1 - (\beta_0 + \beta_1 x_1))^2 + \cdots + (y_n - (\beta_0 + \beta_1 x_n))^2 = \\&= \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2\end{aligned}$$

Hallamos $\hat{\beta}_0$ y $\hat{\beta}_1$ tales que $\nabla RSS(\hat{\beta}_0, \hat{\beta}_1) = (0, 0)$:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

donde:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$
$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

¿Qué tan bien ajusta el modelo?

Error cuadrático medio (ECM): cuantifica qué tan cerca está un valor predicho del valor real:

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

¿Qué tan bien ajusta el modelo?

Variabilidad del modelo:

Variabilidad total:	$\sum_{i=1}^n (y_i - \bar{y})^2$	(\approx Varianza muestral)
Variabilidad no explicada:	$\sum_{i=1}^n (y_i - \hat{y}_i)^2$	(RSS)
Variabilidad explicada:	$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	

¿Qué tan bien ajusta el modelo?

La proporción de la variabilidad de Y explicada por X se puede explicar como:

$$R^2 = \frac{\text{Variabilidad explicada}}{\text{Variabilidad total}} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

¿Qué tan bien ajusta el modelo?

La proporción de la variabilidad de Y explicada por X se puede explicar como:

$$R^2 = \frac{\text{Variabilidad explicada}}{\text{Variabilidad total}} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Obs: $0 \leq R^2 \leq 1$

A mayor R^2 más cercanos están los puntos a la recta y, por lo tanto, tiene más poder de predicción.