

# REGRESIÓN LINEAL: UN PROBLEMA DE APRENDIZAJE ESTADÍSTICO

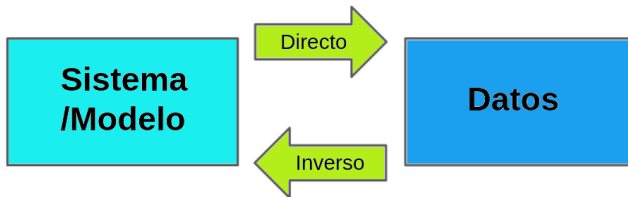
Mario I. Caicedo

17 de mayo de 2022



# APRENDIZAJE ESTADÍSTICO

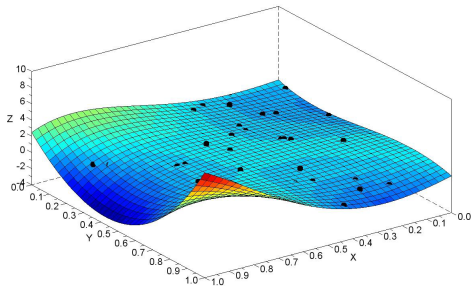
$$Y = H(\Theta|X)$$



$$Y^{(obs)} = H(\Theta|X)$$

# REGRESIÓN LINEAL

$$z_i = z_0 + a_i x_i + b_i y_i + c_i x_i^2 + d_i y_i^2 + e_{ij} x_i y_j + [\dots] + \epsilon_i$$



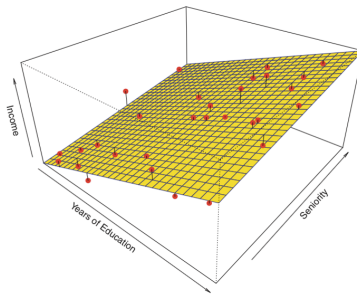
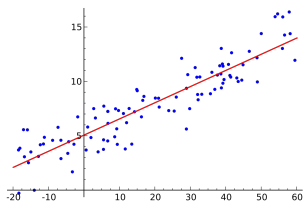
- $\mathbf{y}^{(obs)} = \mathbf{H}(\boldsymbol{\Theta}|\mathbf{X})$
- $\mathbf{y}^{(obs)} = \boldsymbol{\Theta}^T \mathbf{X}$



# CASOS MÁS SENCILLOS

$$y_i = b + m x_i + \epsilon_i$$

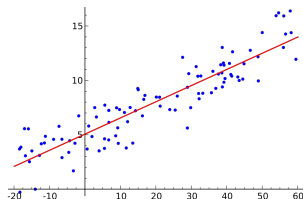
$$z_i = z_0 + a_i x_i + b_i y_i + \epsilon_i$$



# PRIMER CONTACTO

La **Regresión Lineal Simple** es el caso más elemental de entre los problemas de regresión lineal.

DATOS:  $N$  pares  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$



Los valores  $y_i$  tienen ciertos errores de medición ( $\epsilon_i$ ) y en un gráfico cruzado (crossplot)  $y$  vs  $x$  encontramos que parece haber una relación de la forma

$$y_i = \theta_1 x_i + \theta_0 + \epsilon_i$$

# NOMENCLATURA

En la notación estándar de **ML**, la **hipótesis/modelo** se resume en las  $N$  fórmulas

$$y_i = h_{\Theta}(x_i) = \theta_0 + \theta_1 x_i + \varepsilon_i$$

donde los valores  $y_i$  y  $x_i$  son conocidos como **valores objetivo** y **predictores** respectivamente y las cantidades  $\theta_0$  y  $\theta_1$  **pesos**.



# NOTACIÓN

En notación matricial, la hipótesis se escribe como

$$\mathbf{Y} = \mathbf{H}(\boldsymbol{\Theta}|\mathbf{X}) = \mathbf{X} \boldsymbol{\Theta}$$

Donde:

$$\mathbf{Y} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{pmatrix} \quad \text{y} \quad \boldsymbol{\Theta} = \begin{pmatrix} \theta_0 \\ \theta_1 \end{pmatrix}$$



# PLANTEAMIENTO DEL PROBLEMA

- El problema de regresión (ajuste) ó **aprendizaje** consiste en encontrar (calcular/estimar) el vector de pesos  $\Theta$
- Para atacar el problema de **aprendizaje** se introduce un problema de minimización en el espacio de los pesos
- Una vez que se alcanza el **aprendizaje**, es decir, se encontraron los pesos  $\theta_0$  y  $\theta_1$ , estos valores se pueden utilizar junto con la hipótesis para llevar a cabo predicciones/generalizaciones (extrapolaciones)





# EL PROBLEMA DE MINIMIZACIÓN

- La función de pérdida/costo (error cuadrático medio) se define por:

$$\begin{aligned} L(\Theta) &= \left[ \mathbf{Y}^{(obs)} - \mathbf{H}(\Theta|\mathbf{X}) \right]^T \left[ \mathbf{Y}^{(obs)} - \mathbf{H}(\Theta|\mathbf{X}) \right] = \\ &= \left[ \mathbf{Y}^{(obs)} - \mathbf{X}\Theta \right]^T \left[ \mathbf{Y}^{(obs)} - \mathbf{X}\Theta \right] \end{aligned}$$

- La solución al problema de aprendizaje se reduce a minimizar la función de costo con respecto al vector de pesos, esto es, resolviendo la ecuación

$$\nabla_{\Theta^T} L(\Theta) = 0.$$



# SOLUCIÓN ANALÍTICA

- El gradiente de la función de pérdida es

$$\nabla_{\Theta^T} L(\Theta) = -\mathbf{X}^T (\mathbf{Y}^{(obs)} - \mathbf{X}\Theta)$$

- Al resolver para los puntos críticos de  $L$ , se obtiene de inmediato

$$\hat{\Theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}^{(obs)}$$



# SOLUCIÓN ANALÍTICA

## COMENTARIOS

- 1 La solución analítica requiere el cálculo de la inversa de la matriz  $\mathbf{X}^T \mathbf{X}$ , una matriz  $2 \times 2$ , lo que convierte al problema en **casi trivial** [la existencia de la inversa no está garantizada]
- 2 El problema de aprendizaje también puede resolverse minimizando la función de costo a través de la técnica de **Descenso por Gradiente**.



# FÓRMULAS EXPLÍCITAS PARA LA SOLUCIÓN ANALÍTICA

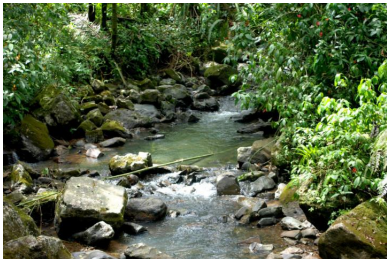
$$\Theta = \begin{pmatrix} N & \sum x_i \\ \sum x_i & \sum (x_i)^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum y_i \\ \sum y_i x_i \end{pmatrix}$$

$$\theta_0 = \textit{intercepto} = \frac{N \sum y_i - m \sum x_i}{N}$$

$$\theta_1 = \textit{pendiente} = \frac{N \sum x_i y_i - \sum x_i \sum y_j}{N \sum x_i^2 - (\sum x_i)^2}$$

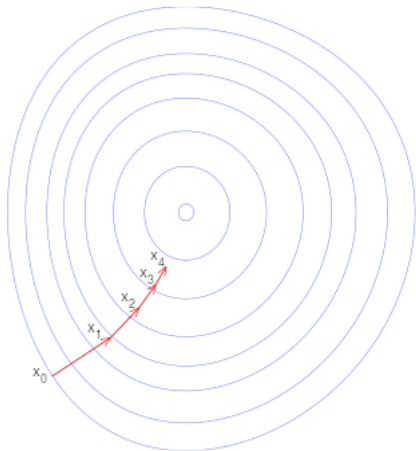


# DESCENSO POR GRADIENTE: INTUICIÓN



El método de descenso por gradiente intenta emular el comportamiento de un arroyo de montaña, que siempre fluye “hacia abajo” siguiendo la ruta de mayor pendiente posible.

# DESCENSO POR GRADIENTE I



- El método de descenso por gradiente es un **algoritmo** de minimización **iterativo**
- El algoritmo utiliza la propiedad geométrica fundamental del gradiente.



# DESCENSO POR GRADIENTE II

- 1 El algoritmo comienza dando una semilla (valor inicial)  $\Theta^{(0)}$
- 2 A partir de la semilla se itera (el superíndice  $k$  indica la iteración) para conseguir nuevos valores de los parámetros

$$\Theta^{(k+1)} = \Theta^{(k)} - \mu \nabla_{\Theta^T} L(\Theta^{(k)})$$

- 3 La iteración se detiene con algún criterio. [Un criterio típico consiste en la estabilización de la función de costo].



# PSEUDOCÓDIGO

---

---

**procedure** GRADIENT DESCENT

**Input**  $\Theta^{(0)}$ : Semilla

$\Theta^{(k)} = \Theta^{(0)}$

**While**  $\epsilon \geq \text{criterio}$

$L^{(k)} = L(\Theta^{(k)})$

**Paso de Gradiente:**

$$\Theta^{(k+1)} = \Theta^{(k)} - \mu \nabla_{\Theta^T} L(\Theta^{(k)})$$

$L^{(k+1)} = L(\Theta^{(k+1)})$

$\epsilon = |L^{(k+1)} - L^{(k)}|$

**Output:** return  $\Theta^{(k+1)}$

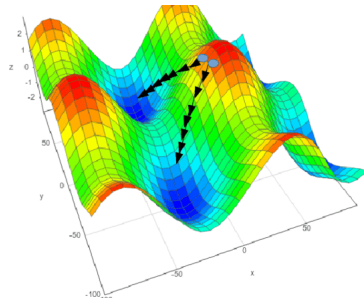
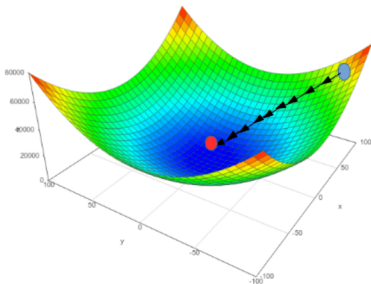
---





# DESCENSO POR GRADIENTE

Hay que tener cuidado con la posibilidad de multimodalidad



# REGRESIÓN LINEAL

## DEFINICIÓN

*Un problema de Regresión Lineal es un problema de Aprendizaje Automático en que la Hipótesis  $\mathbf{Y} = \mathbf{H}(\boldsymbol{\Theta}|\mathbf{X})$  tiene la forma matricial*

$$\mathbf{Y} = \mathbf{F}[\mathbf{X}]\boldsymbol{\Theta},$$

*donde  $\mathbf{F}[\mathbf{X}]$  es una matriz que depende de los atributos predictores*

## OBSERVACIÓN

*La solución a un problema de regresión lineal se consigue siguiendo exactamente los mismos pasos utilizados en el caso de la Regresión Lineal Simple, es decir, utilizando el método de mínimos cuadrados*



# REGRESIÓN LINEAL: OBSERVACIONES I

- La solución analítica del problema ya es conocida:

$$\hat{\Theta} = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Y}^{(obs)}$$

- Desafortunadamente, no hay garantías acerca de la inversibilidad o buen condicionamiento de  $\mathbf{F}^T \mathbf{F}$
- Para  $\dim(\Theta) = N$ ,  $\dim(\mathbf{F}^T \mathbf{F}) = N \times N$ , en consecuencia, y a pesar de las lindas propiedades de  $\mathbf{F}^T \mathbf{F}$  e incluso si la matriz es inversible y bien condicionada, la inversión puede ser un problema extremadamente pesado de álgebra lineal numérica
- Para las regresiones lineales, la función de pérdida de mínimo cuadrados es convexa lo que garantiza que el método de descenso encontrará un mínimo (se alcanza el aprendizaje)



# DESCENSO POR GRADIENTE

## APLICACIÓN A MÍNIMOS CUADRADOS

- Función de pérdida (loss). Necesaria para seguir la evolución del algoritmo y detener la corrida ( $\mathbf{F} = \mathbf{F}[\mathbf{X}]$ ).

$$L(\Theta) = \left[ \mathbf{Y}^{(obs)} - \mathbf{F}\Theta \right]^T \left[ \mathbf{Y}^{(obs)} - \mathbf{F}\Theta \right] \quad (1)$$

- Gradiente

$$\nabla_{\Theta^T} L(\Theta) = -\mathbf{F}^T (\mathbf{Y}^{(obs)} - \mathbf{F}\Theta) \quad (2)$$



## REGRESIÓN LINEAL: OBSERVACIONES II

- ¿De donde sale la función de costo? [Máxima Verosimilitud].
- Por diversas razones conviene “regularizar” la función de costo

$$L(\Theta) = \left[ \mathbf{Y}^{(obs)} - \mathbf{F}\Theta \right]^T \left[ \mathbf{Y}^{(obs)} - \mathbf{F}\Theta \right] + \lambda^2 \Theta^T \Theta$$

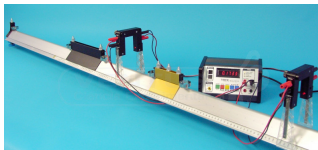
- La regularización implica modificaciones, por ejemplo, la solución exacta:

$$\hat{\Theta} = [\mathbf{F}^T \mathbf{F} + \lambda^2 \mathbf{I}]^{-1} \mathbf{F}^T \mathbf{Y}^{(obs)}$$

- ¿Como saber que tan “buena” es la regresión? [validación].



# EJEMPLO I



En un experimento de medición de posición en función del tiempo se encuentra que los datos parecen alinearse ( $x = x_0 + v_0 t$ )

$$\mathbf{X}^{(obs)} = \mathbf{F}\boldsymbol{\Theta} \iff \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{pmatrix} = \begin{pmatrix} 1 & t_1 \\ 1 & t_2 \\ \vdots & \vdots \\ 1 & t_N \end{pmatrix} \begin{pmatrix} x_0 \\ v_0 \end{pmatrix}$$



## EJEMPLO II

Para este problema, las fórmulas generales 1 y 2 quedan como sigue

$$\nabla_{\Theta^T} L(\Theta) = -\mathbf{F}^T (\mathbf{X}^{(obs)} - \mathbf{F}\Theta)$$

$$\nabla_{\Theta^T} L(\Theta) = - \begin{pmatrix} 1 & 1 & 1 \dots & 1 \\ t_1 & t_2 & t_3 \dots & t_N \end{pmatrix} \left[ \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{pmatrix} - \begin{pmatrix} 1 & t_1 \\ 1 & t_2 \\ 1 & t_3 \\ \vdots & \vdots \\ 1 & t_N \end{pmatrix} \begin{pmatrix} x_0 \\ v_0 \end{pmatrix} \right]$$



## EJEMPLO III

Y en definitiva, el gradiente tiene la forma

$$\nabla_{\Theta^T} L(\Theta^{(k)}) = - \begin{pmatrix} \sum x_i \\ \sum t_i x_i \end{pmatrix} + \begin{pmatrix} N & \sum t_i \\ \sum t_i & \sum (t_i)^2 \end{pmatrix} \begin{pmatrix} x_0^{(k)} \\ v_0^{(k)} \end{pmatrix}$$





# RESUMEN

- 1 Hemos explorado un problema (regresión lineal) muy general de aprendizaje automático
- 2 Mostramos el rango de aplicabilidad del problema y discutimos sus soluciones analítica e iterativa.



# RECURSOS

- **Lenguajes**

- ① [Python](#)

- **Bibliotecas**

- ① [Numpy, Python](#)

- ② [Sklern, Python](#)

- **Recurso Cloud**

- ① [Colab](#)

