

# Regresión lineal por gradiente descendente

Gamaliel Moreno Chávez

MCPI

Evaluación  
2021

- Aprendizaje supervisado
- Regresión lineal
- Descenso por gradiente

# Aprendizaje supervisado

Aprendizaje supervisado. métodos entrenados con:

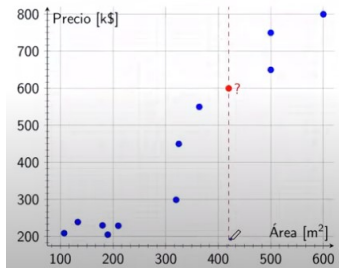
- Conjunto de entrenamiento con pares ordenados  $(\mathbf{x}^{(i)}, y^{(i)})$
- $\mathbf{x}^{(i)}$  es el i-ésimo vector de entrada
- $y^{(i)}$  es la correspondiente etiqueta (label) correcta que se desea predecir posteriormente
- Descenso por gradiente

Es el tipo de aprendizaje más común.

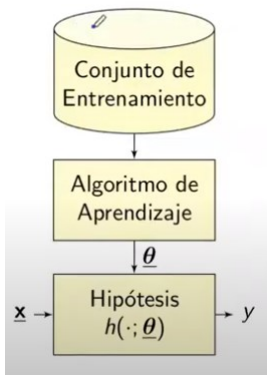
# Problema

Dado un conjunto de entrenamiento como el que se muestra en la figura, ¿Cómo se puede encontrar la relación de salida en términos de la entrada?

Área [m <sup>2</sup> ]	Plantas	Hab.	Precio [k\$]
600	3	5	800
190	2	2	205
210	2	2	229
364	2	2	550
325	2	4	450
180	2	2	230
133	2	2	239
500	2	3	650
107	1	2	209
320	2	3	299
500	2	4	750



# Notación y modelo supervisado



- $m$ : números de muestras de entrenamiento
- $\mathbf{x}$ : datos de entrada
- $n$ : dimensión de la entrada  $\mathbf{x}$  (número de características)
- $y$ : variable de salida u objetivo (target)
  - Clasificación:  $y \in \{C_1, \dots, C_k\} k \in \mathbb{N}$
  - Regresión:  $y \in \mathbb{R}$
- $(\mathbf{x}^{(i)}, y^{(i)})$ :  $i$ -ésima muestra de entrenamiento
- $\theta$ : parámetros
- $y = h(\mathbf{x}; \theta)$ : hipótesis

# Hipótesis para regresión lineal

- Ejemplo de hipótesis: regresión lineal

$$y = h(\mathbf{x}; \boldsymbol{\theta}) = h_{\boldsymbol{\theta}}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \cdots + \theta_n x_n$$

- Ejemplo con  $n=3$ 
  - $x_1$ : área de casa
  - $x_2$ : # de habitaciones
  - $x_3$ : # de pisos
- Convención para simplificar notación:  $x_0 = 1$

$$y = h(\mathbf{x}; \boldsymbol{\theta}) = h_{\boldsymbol{\theta}}(\mathbf{x}) = \theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n = \sum_{i=0}^n \theta_i x_i$$

$$= \boldsymbol{\theta}^T \mathbf{x} = \langle \boldsymbol{\theta}, \mathbf{x} \rangle = \boldsymbol{\theta} \cdot \mathbf{x}$$

$$\boldsymbol{\theta} = [\theta_0, \theta_1, \dots, \theta_n]^T$$

$$\mathbf{x} = [x_0, x_1, \dots, x_n]^T$$

# Función objetivo y minimización de cuadrados

Para encontrar  $\theta$  minimizamos la función de error  $J(\theta)$  con

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

- $r^i = h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}$  se le llama residuo
- El factor  $\frac{1}{2}$  se coloca por conveniencia
- Planteamos problemas de optimización de mínimos cuadrados ordinarios

$$\theta^* = \arg \min_{\theta}$$

Objetivo: Se buscan parámetros  $\theta$  que producen el menor valor de  $\theta$

# Ejemplo de regresión de precios de casa

El caso general de regresión **lineal** minimiza entonces a

$$\begin{aligned} J(\theta) &= \frac{1}{2} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2} \sum_{i=1}^m (\theta^T \mathbf{x} - y^{(i)})^2 \end{aligned}$$

y para el caso de precio =  $f(\text{area})$

$$J(\theta_0, \theta_1) = \frac{1}{2} \sum_{i=1}^m (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)})^2$$



# Minimización de la función objetivo

- Hay varias posibilidades para minimizar  $J(\theta)$
- En general, las técnicas de aprendizaje
  - Toma un valor inicial de  $\theta$
  - Modifican iterativamente  $\theta$  para reducir  $J(\theta)$

## Caso particular descenso por gradiente

- 1 Tome un valor  $\theta^{(0)}$  inicial,  $t=0$
- 2 Calcule en  $\theta^{(t)}$  el gradiente (máxima dirección de cambio)

$$\nabla_{\theta} J(\theta^{(t)}) = \left[ \frac{\partial J}{\partial \theta_0} \frac{\partial J}{\partial \theta_1} \cdots \frac{\partial J}{\partial \theta_n} \right]^T$$

- 3 Calcule la nueva posición

$$\theta^{(t+1)} := \theta^{(t)} - \alpha \nabla_{\theta} J(\theta^{(t)})$$

o de forma equivalente para cada  $\theta_j, j \in 1, \dots, n$

$$\theta_j^{(t+1)} := \theta_j^{(t)} - \alpha \frac{\partial J(\theta^{(t)})}{\partial \theta_j}$$

Para detener la búsqueda de mínimo:

- Usualmente se utilizan tasas de cambio
- Primera opción :  $J(\theta^t) - J(\theta^{t+1}) < \epsilon$
- segunda opción :  $\|\theta^t - \theta^{t+1}\| < \epsilon$
- tercera opción : Número máximo de iteraciones
- opción usual : combinación de anteriores

- Si el gradiente es fuertemente asimétrico (como en el caso actual), la tasa de aprendizaje  $\alpha$  debe elegirse muy pequeña y proceso necesitará demasiadas iteraciones para converger
- Datos deben normalizarse (preprocesamiento) para evitar estos problemas

# Descenso de gradiente para regresión lineal 1

Partiendo del caso concreto

$$J(\theta_0, \theta_1) = \frac{1}{2} \sum_{i=1}^m (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)})^2$$

podemos calcular el gradiente fácilmente

$$\begin{aligned} \nabla_{\theta} J(\theta_0, \theta_1) &= \begin{bmatrix} \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} \\ \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} \end{bmatrix}, \\ &= \begin{bmatrix} \sum_{i=1}^m (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)}) \cdot 1 \\ \sum_{i=1}^m (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)}) \cdot x_1^{(i)} \end{bmatrix} \end{aligned}$$

## Descenso de gradiente para regresión lineal 2

Observe que para el caso general de regresión lineal se tiene

$$\begin{aligned} J(\boldsymbol{\theta}) &= \frac{1}{2} \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2 \\ &= \frac{1}{2} \sum_{i=1}^m \left( (\theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n) - y^{(i)} \right)^2 \end{aligned}$$

La  $j$ -ésima componente del gradiente  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  es

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \theta_j} = \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \cdot x_j^{(i)}$$

# Descenso de gradiente para regresión lineal 3

Lo que finalmente implica que el gradiente es

$$\nabla J(\boldsymbol{\theta}) = \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \cdot \mathbf{x}^{(i)}$$

# Descenso de gradiente por lotes

Combinando todos los resultados anteriores tenemos el algoritmo de descenso de gradiente por lotes

- ➊ Inicialice  $t := 0$  y  $\theta^{(0)}$
  - ➋ **repeat**
  - ➌  $\theta^{(t+1)} := \theta^{(t)} - \alpha \sum_{i=1}^m (\theta^{(t)T} \mathbf{x}^{(i)} - y^{(i)}) \cdot \mathbf{x}^{(i)}$
  - ➍  $t := t + 1$
  - ➎ **go to 2 until converge**
- Lotes: cada paso usa todo el conjunto de entrenamiento
  - $\alpha$  tasa de aprendizaje, su ajuste es delicado
    - Si  $\alpha$  es muy grande oscila alrededor del mínimo
    - Si  $\alpha$  es muy pequeña, necesita muchos pasos para converger
  - El descenso de gradiente converge a extremos locales, que dependen del punto inicial



# Descenso por gradiente estocástico

Descenso por gradiente estocástico o incremental usa un ejemplo del conjunto de entrenamiento a la vez:

- ➊ Inicialice  $t := 0$  y  $\theta^{(0)}$
  - ➋ **repeat**
  - ➌ **for each**  $(\mathbf{x}^{(i)}, y^{(i)})$  in training set **do**
  - ➍  $\theta^{(t+1)} := \theta^{(t)} - \alpha(\theta^{(t)T} \mathbf{x}^{(i)} - y^{(i)}) \cdot \mathbf{x}^{(i)}$
  - ➎  $t := t + 1$
  - ➏ **end for**
  - ➐ **go to 2 until converge**
- No asegura convergencia
  - Trayectoria hacia el mínimo divaga pero en general se acerca al mínimo
  - Útil para conjuntos de entrenamiento gigantescos

# Descenso por gradiente estocástico

- Una época es la presentación de los  $m$  datos de entrenamiento
- Método estocásticos
  - Produce soluciones acertadas más pronto
  - Ruido de trayectorias puede sacar trayectoria de tascos
  - Una época requiere  $m$  pasos
- Métodos por lotes
  - Es trivialmente paralelizable
  - Tiende a estancarse más fácil
  - Si el conjunto de datos es muy grande, un paso es caro de calcular
  - Si la superficie de error es convexa converge rápido y suave
  - Cada paso es a la vez una época
- Pueden combinarse: mini-lotes
  - Compromiso entre los dos esquemas
  - Puede hacerse
    - Con reemplazo: Preferido si conjuntos de datos son gigantescos
    - Sin reemplazo (permutación): posible sólo con conjuntos pequeños
  - Si mini-lote tiene tamaño  $v$ , entonces una época tiene  $m/v$  pasos