

$$P \rightarrow \boxed{f_A(P)} \rightarrow a$$

## Aprendizaje Supervisado

$f_A$  es desconocida

$$f_A: X \rightarrow Y$$

$x$  es un vector  
típicamente  $x \in \mathbb{R}^n$   
 $x = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$

Si  $Y = \mathbb{R}$  Regresión

Si  $Y = \{-1, 1\}$  Clasificación binaria

Si  $Y = \{C_1, C_2, \dots, C_k\}$  Clasificación en varias clases

$\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$  es una muestra de  $X$

$$D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$$

donde  $y^{(i)} = f_A(x^{(i)}) + \epsilon$  V.A. dist. desconocida

El Problema encontrar una

$$\text{función } h^*: X \rightarrow Y$$

tal que  $h^* \approx f_A$

$h^* \in H$  hipótesis posibles

$$H = \{h_\alpha \mid h_\alpha: \mathbb{R} \rightarrow \mathbb{R} \text{ donde } h_\alpha(x) = \alpha x, \forall x \in \mathbb{R}\}$$

$$h: X \times \Theta \rightarrow Y, \Theta = \mathbb{R}^p$$

$$h(x^{(i)}, \theta) = \hat{y}^{(i)}$$

Si  $\theta$  es un vector de parámetros físicos

$$h_\theta: X \rightarrow Y$$

$$h_w(x) = \sum_{j=1}^T w_j x^j + w_0, x^j \in \mathbb{R}$$

$$h_w: \mathbb{R} \rightarrow \mathbb{R}$$

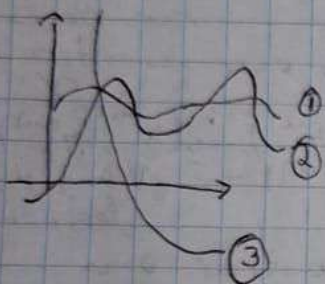
$$w = (w_0, w_1, \dots, w_T) \in \mathbb{R}^{T+1}$$

$$\hat{y} = w_1 x + w_0, \hat{y} = w_0 x^2 + w_1 x + b$$



Hipótesis:

1.  $f: X \rightarrow y$  existe y es desconocida
2. Tengo  $\{x^{(1)}, \dots, x^{(m)}\} \subseteq X$  una muestra de  $m$  datos con distribución desconocida.
3.  $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ , donde  $y^{(i)} = f(x^{(i)}) + e^{(i)}$  donde  $e^{(i)}$  son VA vector de variables
4. Tenemos una función "Parametrizada"  $h: X \times \Theta \rightarrow y$ , típicamente  $\Theta = \mathbb{R}^p$
5. Para un valor específico de  $\theta$ ,  $h_\theta: X \rightarrow y$
6.  $h \in H$  Conjunto de hipótesis
7. El aprendizaje supervisado consiste en encontrar  $h^* \in H$  tal que  $h^* \approx f$



$L: y \times y \rightarrow \mathbb{R}$  función de Pérdida

$$L(y, \hat{y}) = L(f(x), h^*(x))$$

$$L(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$$

$$= |y - \hat{y}| \quad y \in \mathbb{R}$$

$$L(y - \hat{y}) = \begin{cases} 0 & \text{si } y = \hat{y} \\ 1 & \text{si } y \neq \hat{y} \end{cases} = \mathbb{I}\{y \neq \hat{y}\} \quad y \in \{-1, 1\}$$

Error fuera de muestra

$$E_{out}(h^*) = \mathbb{E}_{x \in X} [L(f(x), h^*(x))]$$

$$f \approx h^* \text{ si y sólo si } E_{out}(h^*) \approx 0$$

$$E_{in} = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, h(x^{(i)}))$$

Error en muestra

$$f \approx h^* \text{ si } \begin{cases} E_{in}(h^*) \approx 0 \\ E_{in}(h^*) \approx E_{out}(h^*) \end{cases}$$

Desigualdad de Hoeffding (PAC - Learning)

$$P(|E_{out}(h^*) - E_{in}(h^*)| > \epsilon) \leq 2e^{-2\epsilon^2 m}$$

donde  $m$  es el tamaño de la muestra



dimension

$d_{vc}(H) \geq \#$  Parameters INDEPENDIENTES

EL APRENDIZAJE ES POSIBLE SI:

$$10^{-4} d_{vc}(H) \ll M \Leftrightarrow E_{in}(h^*) \approx E_{out}(h^*)$$

$$\begin{matrix} \begin{matrix} x^{(1)T} \\ x^{(2)T} \\ \vdots \\ x^{(m)T} \end{matrix} \rightarrow \begin{matrix} a_1 & a_2 & \dots & a_n \\ (x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)}) \\ (x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ (x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)}) \end{matrix} \begin{matrix} x \\ y \\ \vdots \\ y^{(m)} \end{matrix} \end{matrix}$$

instancias  $(m, n)$   $(m, 1)$

$x^{(i)} \in \mathbb{R}^n$   
 $y^{(i)} \in \mathbb{R}$

$$\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$$

Conjunto de aprendizaje

$$h(x) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b = \sum_{j=1}^n w_j x_j + b = w^T x + b$$

$$= x^T w + b$$

$$S: x = (x_1, \dots, x_n) \in \mathbb{R}^n$$

$$w = (w_1, \dots, w_n) \in \mathbb{R}^n$$

$$\theta = (w_1, \dots, w_n, b) \in \mathbb{R}^{n+1}$$

Aprendizaje:

$$E_{in}(h^*) \approx E_{out}(h^*)$$

$$E_{in}(h^*) \approx 0$$

$$h^* = \arg \min E_{in}(h) \Leftrightarrow \theta^* = w^*, b^* = \arg \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} E_{in}(h_w, b)$$

$$w^*, b^* = \arg \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (y^{(i)} - h_{w,b}(x^{(i)}))^2 \quad \text{MSE}$$

Mean Square Error

$$w^*, b^* = \arg \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (y^{(i)} - \underbrace{w^T x^{(i)} - b}_{= x^{(i)T} w - b})^2$$

$$= \arg \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{2m} [(y^{(1)} - x^{(1)T} w - b) \dots (y^{(m)} - x^{(m)T} w - b)]$$

$$= \arg \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{2m} [Y - [X, 1] \theta]^T [Y - [X, 1] \theta]$$

$\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} = \begin{bmatrix} x^{(1)T} \\ x^{(2)T} \\ \vdots \\ x^{(m)T} \end{bmatrix} w - \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} b$

$\begin{bmatrix} y^{(1)} - x^{(1)T} w - b \\ y^{(2)} - x^{(2)T} w - b \\ \vdots \\ y^{(m)} - x^{(m)T} w - b \end{bmatrix}$



Para encontrar un mínimo  
de una función: se deriva  
y se iguala a cero

En este caso como tenemos dos variables tenemos que derivar  
Parcialmente para encontrar el gradiente.

Ejemplo:  $f(x_1, \dots, x_n)$   
 $f(y)$   $x \in \mathbb{R}^n$ ,  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$$

Regresamos al Problema...

$$[Y - X\theta]^T [Y - X\theta]$$

$$\sum_{m=1}^M \sum_{i=1}^n [y - x_i \theta]^2$$

$$\frac{\partial (y - x_i)^2}{\partial x_i} = -2(y - x_i)x_i$$

$$\frac{1}{M} X^T [Y - X\theta] = \vec{0}$$

$$X^T Y - X^T X \theta = \vec{0}$$

$$X^T X \theta = X^T Y$$

$$\theta = [X^T X]^{-1} X^T Y \leftarrow \text{Este es el método mínimo de cuadrados}$$

$$P_{inv}(X)$$