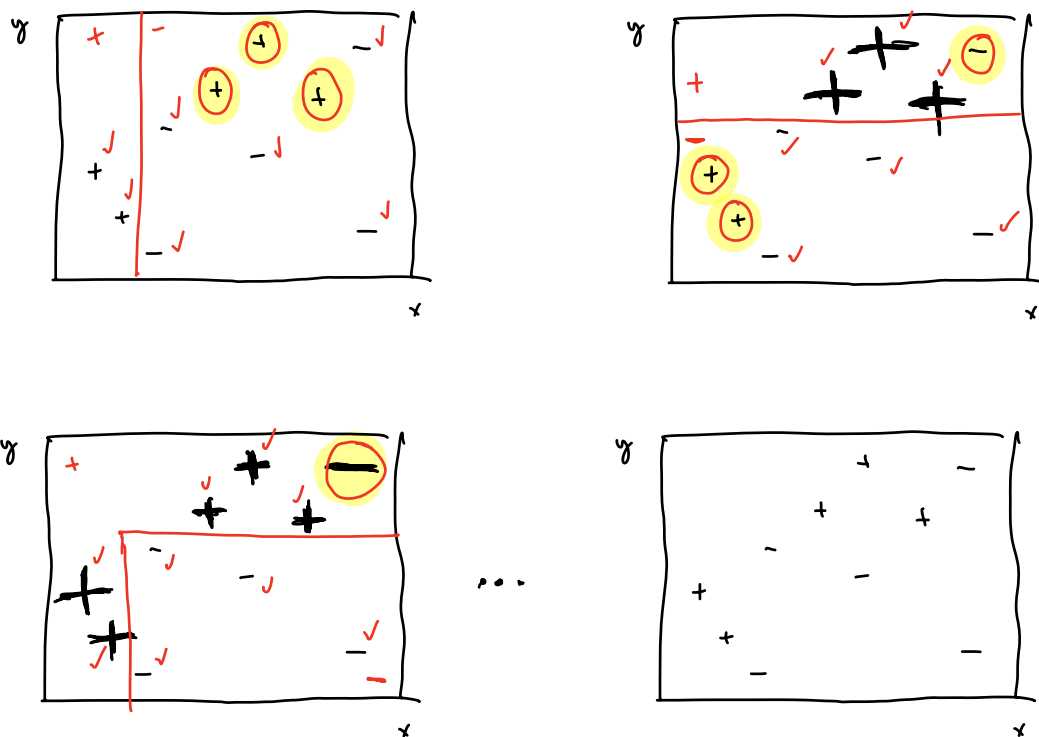


Adaboost: en cada iteración construimos un árbol de clasificación y le damos más peso a las observaciones en las que el algoritmo se equivocó en la iteración pasada:



Tenemos:  $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$  ,  $x_i$  vector de  $p$  dimensiones

↳ construir modelo  $g^{[1]}(\cdot)$  para predecir  $y_i$ :

$(x_1, y_1), \dots (x_n, y_n) \longrightarrow g^{[1]}(\cdot)$  árbol

⇓

actualizamos los pesos de cada observación

$(x_1, y_1), \dots (x_n, y_n) \longrightarrow g^{[2]}(\cdot)$

⋮

$$\hat{f}_{\text{adaboost}}(x) = \sum_{m=1}^M \alpha_m g^{[m]}(x)$$

→ predicción final del adaboost es votación ponderada de cada uno de los árboles

$\alpha_m$ : peso de cada árbol

Inicializamos los pesos:  $w_i^{[1]} = \frac{1}{n}$ ,  $i = 1, \dots, n$ .

Para  $m = 1, \dots, M$ :

• construir un árbol de decisión utilizando los pesos  $w_i^{[m-1]}$  para obtener un modelo  $g^{[m]}$ .

$$\text{error}^{[m]} = \frac{\sum_{i=1}^n w_i^{[m-1]} \mathbb{I}\{y_i \neq g^{[m]}(x_i)\}}{\sum_{i=1}^n w_i^{[m-1]}}$$

$$\alpha^{[m]} = \log \left( \frac{1 - \text{error}^{[m]}}{\text{error}^{[m]}} \right)$$

$$\tilde{w}_i^{[m]} = w_i^{[m-1]} \exp \left[ \alpha^{[m]} \mathbb{I}\{y_i \neq g^{[m]}(x_i)\} \right]$$

$$w_i^{[m]} = \frac{\tilde{w}_i^{[m]}}{\sum_{j=1}^n \tilde{w}_j^{[m]}}$$

Xgboost:  $(x_1, y_1) \dots (x_n, y_n) \rightarrow$  construye un modelo  $f(x)$  para predecir.

$$\begin{aligned} y_1 &= 0.9 \\ y_2 &= 1.4 \\ &\vdots \end{aligned}$$

$$\begin{aligned} f(x_1) &= 0.8 \\ f(x_2) &= 1.3 \\ &\vdots \end{aligned}$$

equivocó en 0.1

equivocó en 0.1

Quisiera construir un segundo modelo que me ayude a "corregir" esos errores:

$$\begin{aligned} f(x_1) + h(x_1) &= y_1 \\ f(x_2) + h(x_2) &= y_2 \\ &\vdots \end{aligned}$$

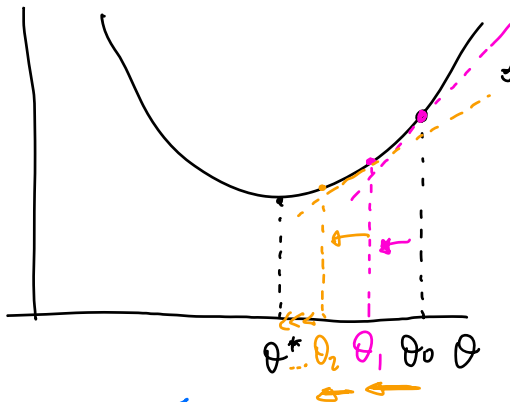
"h" completara a "f"

$$h(x_1) = y_1 - f(x_1)$$

$$h(x_2) = y_2 - f(x_2)$$

→ quiero que  $h$  me ayude a predecir los errores del modelo  $f$ .

Algoritmos de gradient descent:



$$\theta_{i+1} = \theta_i - \rho \frac{\partial J}{\partial \theta}$$

learning rate.  
derivada de la función de pérdida con respecto al parámetro.

Si  $\rho$  es baja, necesito muchas iteraciones para que el algoritmo converja.

Si  $\rho$  es más alto, necesito menos iteraciones para converger.

Si  $\rho$  es muy alto, puede que no converja.

Función de pérdida cuadrática:

$$L(y, f(x)) = \sum_{i=1}^n \frac{(y_i - f(x_i))^2}{2}$$

$$\frac{\partial L(y, f(x))}{\partial f(x_i)} = (y_i - f(x_i)) \rightarrow \text{error del modelo en la observación } i.$$

$$f_2(x_i) = f(x_i) + h(x_i) \quad \sim \quad h \text{ aproxima } (y_i - f(x_i)) \text{ error del modelo } f.$$

$$f_2(x_i) = f(x_i) + (y_i - f(x_i))$$

$$= f(x_i) - \underbrace{(f(x_i) - y_i)}_{\frac{\partial L(y, f(x))}{\partial f(x_i)}}$$

$$f_2(x_i) = \underline{f(x_i)} - \underbrace{\underline{1}}_{\underline{\frac{\partial L(y, f(x))}{\partial f(x_i)}}} \quad \text{learning rate} = 1.$$