



# IIC 2433 Minería de Datos

<https://github.com/marcelomendoza/IIC2433>

# Ensembles

# Ensembles

En ensembles asumimos que al usar varios modelos podemos obtener mejores resultados que usando un modelo.

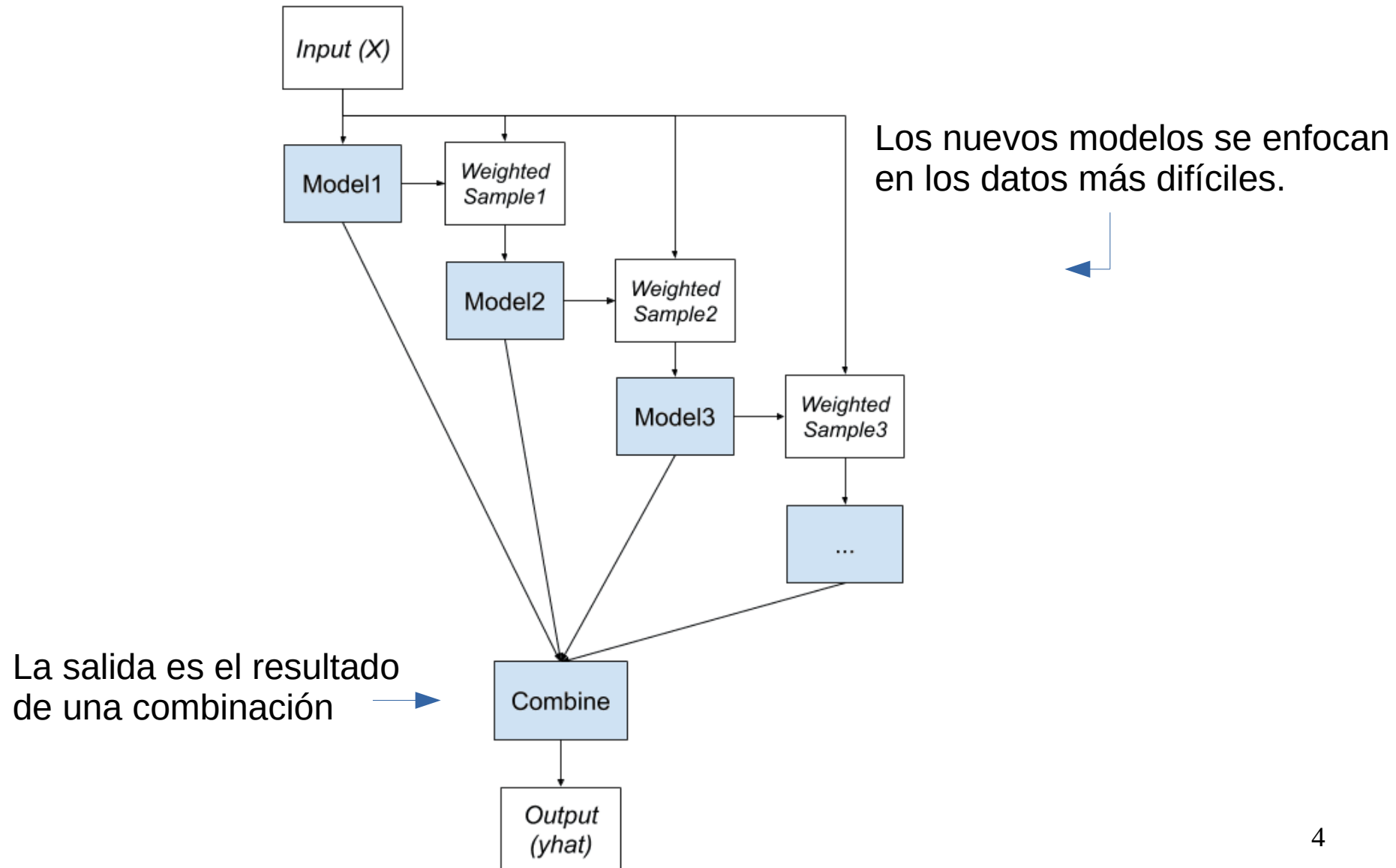
Los ensembles denominan a los modelos base *weak learners*.

Existen tres estrategias para combinar modelos:

- Boosting (AdaBoost, XGBoost, ...)
- Bagging (random forests, extra trees, ...)
- Stacking (canonical stacking, super ensemble, ...)

# Boosting

## Boosting Ensemble



# Boosting

---

**Input:** Sample distribution  $\mathcal{D}$ ;  
Base learning algorithm  $\mathcal{L}$ ;  
Number of learning rounds  $T$ .

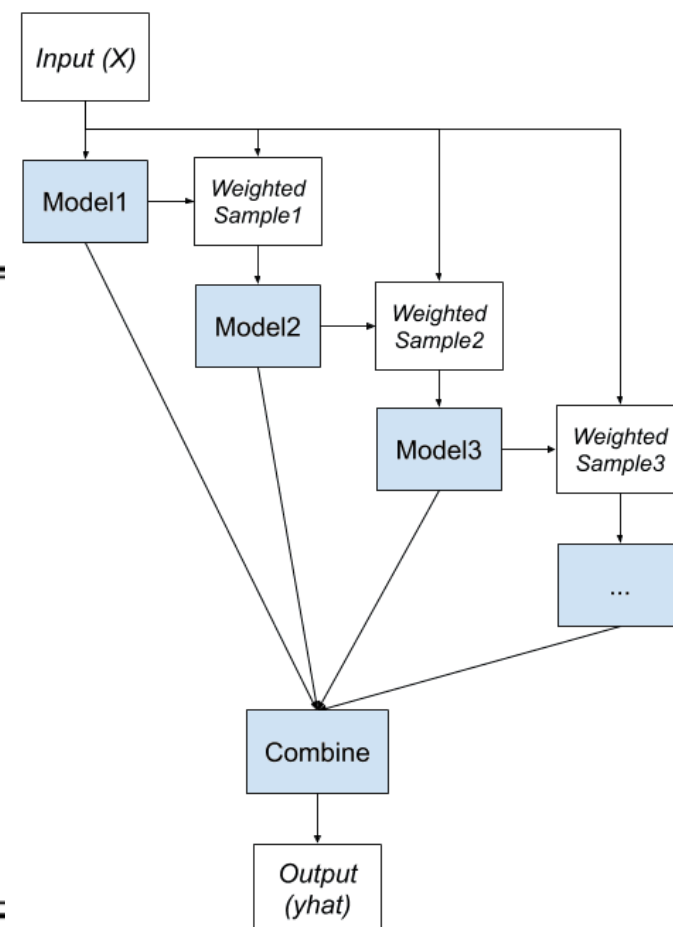
**Process:**

1.  $\mathcal{D}_1 = \mathcal{D}$ .     % Initialize distribution
2. **for**  $t = 1, \dots, T$ :
3.      $h_t = \mathcal{L}(\mathcal{D}_t)$ ;   % Train a weak learner from  $\mathcal{D}_t$
4.      $\epsilon_t = P_{\mathbf{x} \sim \mathcal{D}_t}(h_t(\mathbf{x}) \neq f(\mathbf{x}))$ ;   % Evaluate the error
5.      $\mathcal{D}_{t+1} = \text{Adjust\_Distribution}(\mathcal{D}_t, \epsilon_t)$
6. **end**

**Output:**  $H(\mathbf{x}) = \text{Combine\_Outputs}(\{h_1(\mathbf{x}), \dots, h_t(\mathbf{x})\})$

---

Boosting Ensemble



# AdaBoost

AdaBoost es un algoritmo basado en Boosting muy usado.

AdaBoost tiene por objetivo minimizar la pérdida según:

$$\ell_{\text{exp}}(h \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})h(\mathbf{x})}] \quad \text{Exponential loss}$$

└─ Clases en  $\{-1, +1\}$

# AdaBoost

AdaBoost es un algoritmo basado en Boosting muy usado.

AdaBoost tiene por objetivo minimizar la pérdida según:

$$\ell_{\text{exp}}(h \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})h(\mathbf{x})}] \quad \text{Exponential loss}$$

└─▶ Clases en  $\{-1, +1\}$

Los *weak learners* se combinan de forma lineal:

$$H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) .$$

└─▶ Hay que ajustarlos

## AdaBoost

$$H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) .$$

Minimizando:  $\ell_{\text{exp}}(h \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})h(\mathbf{x})}]$



$$\begin{aligned} \frac{\partial e^{-f(\mathbf{x})H(\mathbf{x})}}{\partial H(\mathbf{x})} &= -f(\mathbf{x})e^{-f(\mathbf{x})H(\mathbf{x})} \\ &= -e^{-H(\mathbf{x})}P(f(\mathbf{x}) = 1 \mid \mathbf{x}) + e^{H(\mathbf{x})}P(f(\mathbf{x}) = -1 \mid \mathbf{x}) \\ &= 0 . \end{aligned}$$



## AdaBoost

$$H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) .$$

Minimizando:  $\ell_{\text{exp}}(h \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})h(\mathbf{x})}]$



$$\begin{aligned} \frac{\partial e^{-f(\mathbf{x})H(\mathbf{x})}}{\partial H(\mathbf{x})} &= -f(\mathbf{x})e^{-f(\mathbf{x})H(\mathbf{x})} \\ &= -e^{-H(\mathbf{x})}P(f(\mathbf{x}) = 1 \mid \mathbf{x}) + e^{H(\mathbf{x})}P(f(\mathbf{x}) = -1 \mid \mathbf{x}) \\ &= 0 . \end{aligned}$$





$$H(\mathbf{x}) = \frac{1}{2} \ln \frac{P(f(x) = 1 \mid \mathbf{x})}{P(f(x) = -1 \mid \mathbf{x})}$$


## AdaBoost

$$H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) .$$

Minimizando:  $\ell_{\text{exp}}(h \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})h(\mathbf{x})}]$


$$\begin{aligned} \frac{\partial e^{-f(\mathbf{x})H(\mathbf{x})}}{\partial H(\mathbf{x})} &= -f(\mathbf{x})e^{-f(\mathbf{x})H(\mathbf{x})} \\ &= -e^{-H(\mathbf{x})}P(f(\mathbf{x}) = 1 \mid \mathbf{x}) + e^{H(\mathbf{x})}P(f(\mathbf{x}) = -1 \mid \mathbf{x}) \\ &= 0 . \end{aligned}$$


$$H(\mathbf{x}) = \frac{1}{2} \ln \frac{P(f(x) = 1 \mid \mathbf{x})}{P(f(x) = -1 \mid \mathbf{x})}$$


$$\begin{aligned} \text{sign}(H(\mathbf{x})) &= \text{sign} \left( \frac{1}{2} \ln \frac{P(f(x) = 1 \mid \mathbf{x})}{P(f(x) = -1 \mid \mathbf{x})} \right) \\ &= \begin{cases} 1, & P(f(x) = 1 \mid \mathbf{x}) > P(f(x) = -1 \mid \mathbf{x}) \\ -1, & P(f(x) = 1 \mid \mathbf{x}) < P(f(x) = -1 \mid \mathbf{x}) \end{cases} \\ &= \arg \max_{y \in \{-1, 1\}} P(f(x) = y \mid \mathbf{x}), \end{aligned}$$

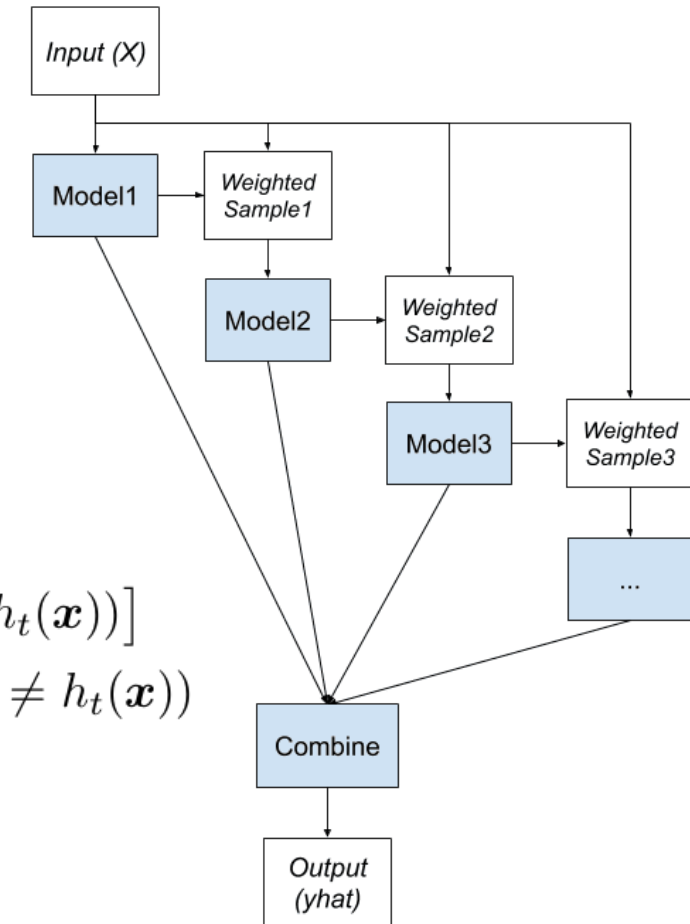
# AdaBoost

$$H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \text{ se produce iterativamente.}$$

Cuando entrenamos un *weak learner*, calculamos su coeficiente minimizando:

$$\begin{aligned} \ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-f(\mathbf{x})\alpha_t h_t(\mathbf{x})}] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-\alpha_t} \mathbb{I}(f(\mathbf{x}) = h_t(\mathbf{x})) + e^{\alpha_t} \mathbb{I}(f(\mathbf{x}) \neq h_t(\mathbf{x}))] \\ &= e^{-\alpha_t} P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) = h_t(\mathbf{x})) + e^{\alpha_t} P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) \neq h_t(\mathbf{x})) \\ &= e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t \end{aligned}$$

## Boosting Ensemble



# AdaBoost

$$H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \text{ se produce iterativamente.}$$

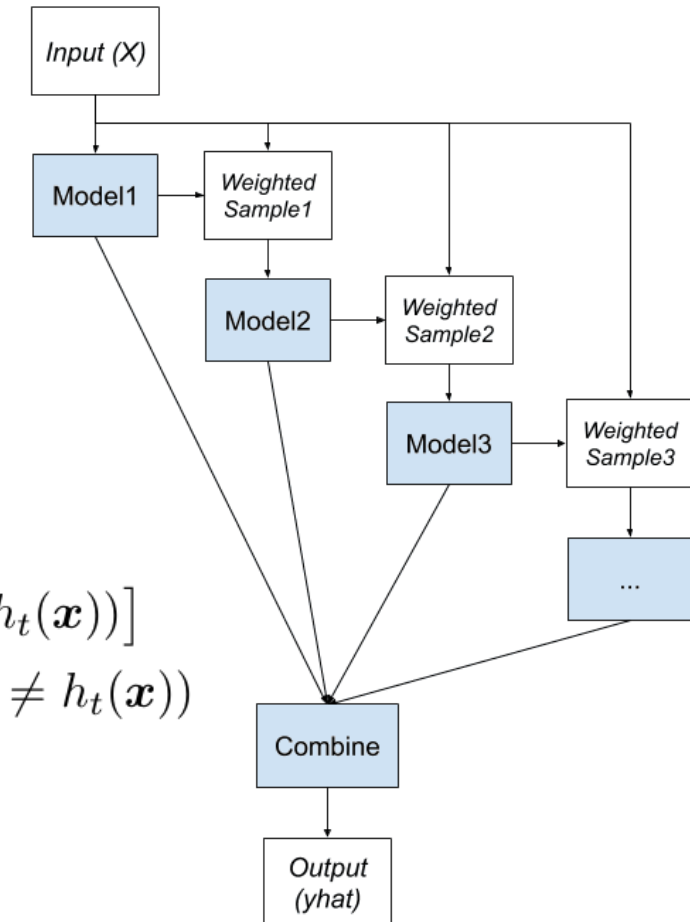
Cuando entrenamos un *weak learner*, calculamos su coeficiente minimizando:

$$\begin{aligned} \ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-f(\mathbf{x})\alpha_t h_t(\mathbf{x})}] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-\alpha_t} \mathbb{I}(f(\mathbf{x}) = h_t(\mathbf{x})) + e^{\alpha_t} \mathbb{I}(f(\mathbf{x}) \neq h_t(\mathbf{x}))] \\ &= e^{-\alpha_t} P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) = h_t(\mathbf{x})) + e^{\alpha_t} P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) \neq h_t(\mathbf{x})) \\ &= e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t \end{aligned}$$

por lo que:

$$\frac{\partial \ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t)}{\partial \alpha_t} = -e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t = 0$$

## Boosting Ensemble



# AdaBoost

$$H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \text{ se produce iterativamente.}$$

Cuando entrenamos un *weak learner*, calculamos su coeficiente minimizando:

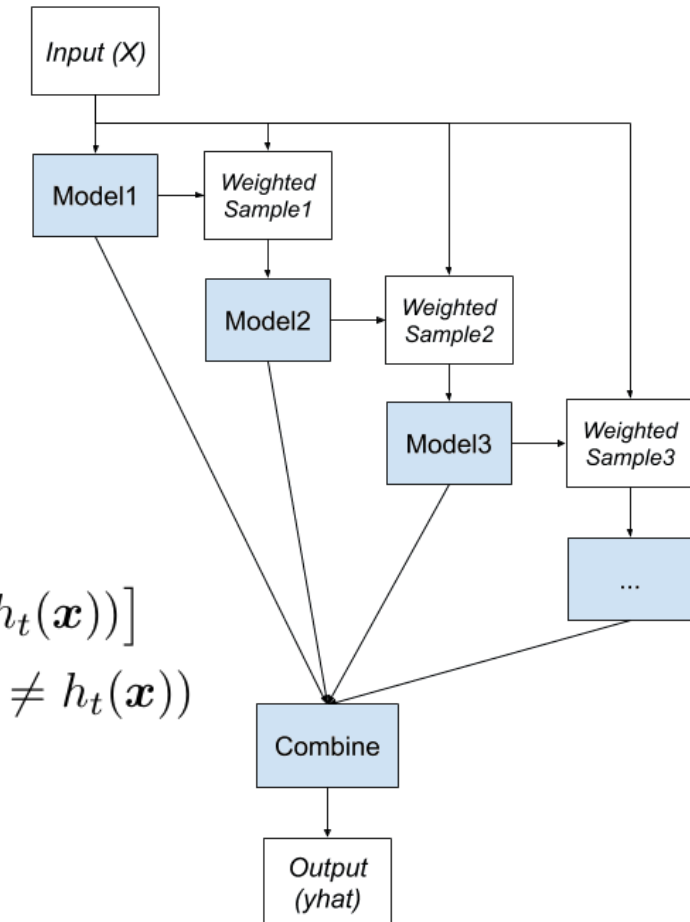
$$\begin{aligned} \ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-f(\mathbf{x}) \alpha_t h_t(\mathbf{x})}] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-\alpha_t} \mathbb{I}(f(\mathbf{x}) = h_t(\mathbf{x})) + e^{\alpha_t} \mathbb{I}(f(\mathbf{x}) \neq h_t(\mathbf{x}))] \\ &= e^{-\alpha_t} P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) = h_t(\mathbf{x})) + e^{\alpha_t} P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) \neq h_t(\mathbf{x})) \\ &= e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t \end{aligned}$$

por lo que:

$$\frac{\partial \ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t)}{\partial \alpha_t} = -e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t = 0$$

$$\rightarrow \alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$

## Boosting Ensemble



## AdaBoost

Ahora vamos a obtener una expresión útil para la distribución:

$$\mathcal{D}_{t+1}(\mathbf{x}) = \frac{\mathcal{D}(\mathbf{x}) e^{-f(\mathbf{x}) H_t(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x}) H_t(\mathbf{x})}]} \longrightarrow \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x}) H_{t-1}(\mathbf{x})} e^{-f(\mathbf{x}) h_t(\mathbf{x})}]$$

## AdaBoost

Ahora vamos a obtener una expresión útil para la distribución:

$$\mathcal{D}_{t+1}(\mathbf{x}) = \frac{\mathcal{D}(\mathbf{x}) e^{-f(\mathbf{x}) H_t(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x}) H_t(\mathbf{x})}]} \longrightarrow \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x}) H_{t-1}(\mathbf{x})} e^{-f(\mathbf{x}) h_t(\mathbf{x})}]$$

Es decir:

$$\mathcal{D}_{t+1}(\mathbf{x}) = \mathcal{D}_t(\mathbf{x}) \cdot e^{-f(\mathbf{x}) \alpha_t h_t(\mathbf{x})} \frac{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x}) H_{t-1}(\mathbf{x})}]}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x}) H_t(\mathbf{x})}]}$$

# AdaBoost

Unamos las piezas para obtener el algoritmo.

---

---

**Input:** Data set  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ ;  
Base learning algorithm  $\mathfrak{L}$ ;  
Number of learning rounds  $T$ .

**Process:**

1.  $\mathcal{D}_1(\mathbf{x}) = 1/m$ .    % Initialize the weight distribution
2. **for**  $t = 1, \dots, T$ :
3.     $h_t = \mathfrak{L}(D, \mathcal{D}_t)$ ; % Train a classifier  $h_t$  from  $D$  under distribution  $\mathcal{D}_t$
4.     $\epsilon_t = P_{\mathbf{x} \sim \mathcal{D}_t}(h_t(\mathbf{x}) \neq f(\mathbf{x}))$ ; % Evaluate the error of  $h_t$
5.    **if**  $\epsilon_t > 0.5$  **then break**
6.     $\alpha_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$ ; % Determine the weight of  $h_t$
7.     $\mathcal{D}_{t+1}(\mathbf{x}) = \mathcal{D}_t(\mathbf{x}) \cdot e^{-f(\mathbf{x})\alpha_t h_t(\mathbf{x})} \frac{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_t(\mathbf{x})}]}$
8. **end**

**Output:**  $H(\mathbf{x}) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right)$

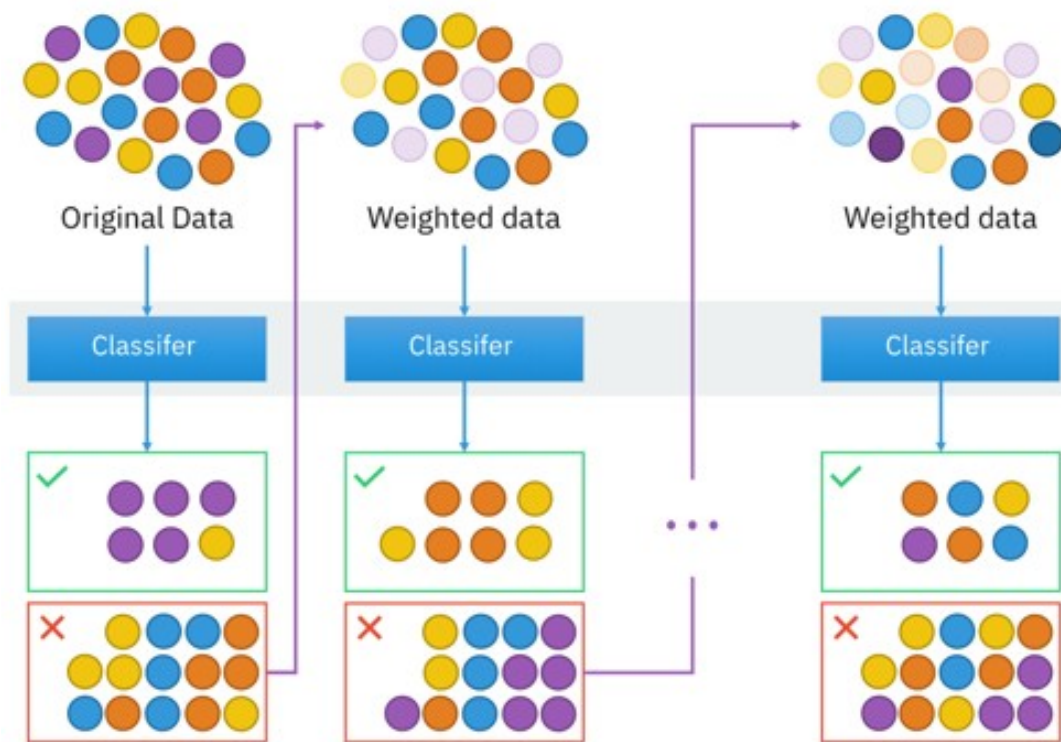
---

---



# AdaBoost

## BOOSTING LEARNING PROCEDURE



Strong Learner      Weak Learners

$$f(x) = \sum_t \alpha_t h_t(x)$$

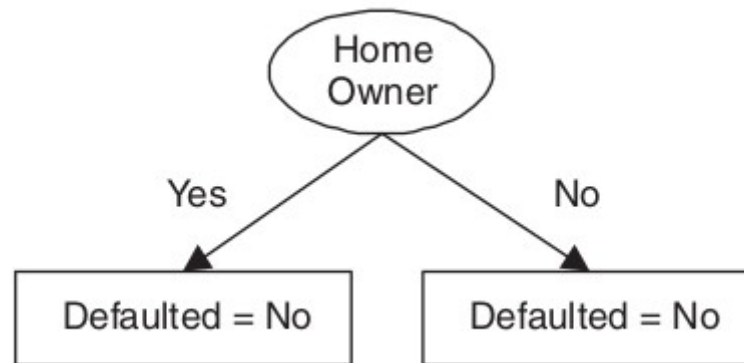
Ensemble Classifier

Weight calculated by considering the last iteration's error

# Weak learner

# Weak learner

## Árbol de decisión



Objetivo: el split produce nodos puros

Minimizar Gini index  $= 1 - \sum_{i=0}^{c-1} p_i(t)^2$

└─ Fracción de ejemplos de una clase en el nodo

# Weak learner

Árbol de decisión

Profundidad máxima

---

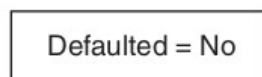
TreeGrowth ( $E, F$ )

```
1: if stopping_cond( $E, F$ ) = true then
2:   leaf = createNode().
3:   leaf.label = Classify( $E$ ).
4:   return leaf.
5: else
6:   root = createNode().
7:   root.test_cond = find_best_split( $E, F$ ).
8:   let  $V = \{v | v \text{ is a possible outcome of } root.test\_cond \}$ .
9:   for each  $v \in V$  do
10:     $E_v = \{e | root.test\_cond(e) = v \text{ and } e \in E\}$ .
11:    child = TreeGrowth( $E_v, F$ ).
12:    add child as descendent of root and label the edge ( $root \rightarrow child$ ) as  $v$ .
13:   end for
14: end if
15: return root.
```

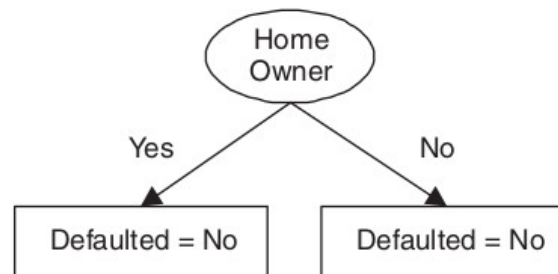
---

# Weak learner

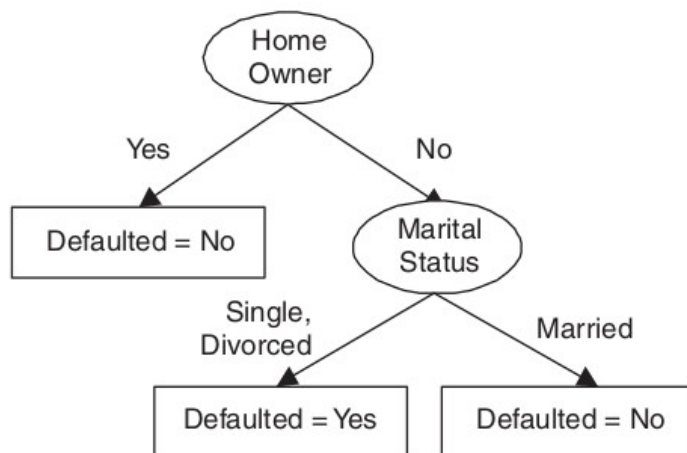
## Árbol de decisión



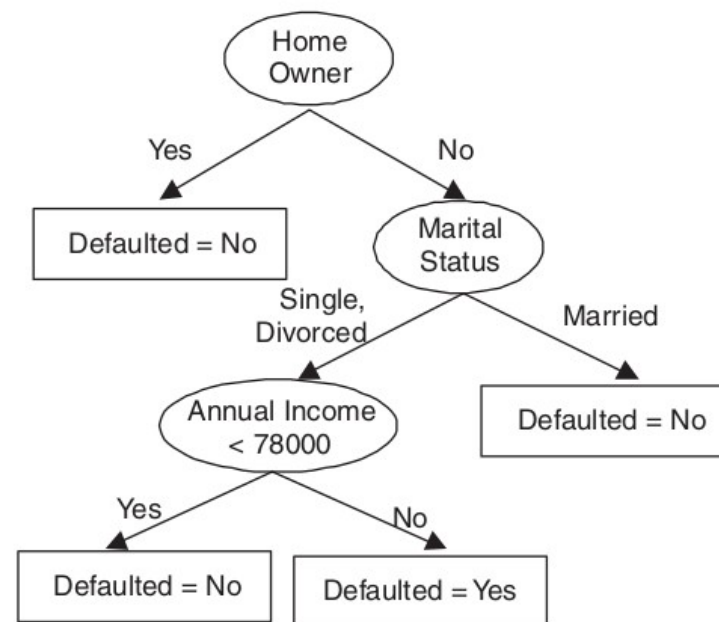
(a)



(b)



(c)



(d)