

# Árboles de decisión

---

## Clasificación (i)

---

- La tarea de clasificación consiste en asignar objetos a una de las clases previamente definidas.
- Se trata de una tarea presente en multitud de aplicaciones:
  - detectar correo spam (spam, no spam)
  - clasificar células tumorales (benignas, malignas)
  - clasificar créditos bancarios (conceder, denegar)

## Clasificación (ii)

---

- Para llevar acabo la **tarea de clasificación** se dispone de un conjunto de objetos caracterizados por el par de atributos  $(x, y)$ , donde  $x$  es un **vector de características** e  $y$  es la conocida como **etiqueta de clase**.
- A la etiqueta de clase  $y$  se le conoce también como categoría.

# Algoritmo de Hunt

## Datos disponibles

		Atributos						Clase		
		$x_1$	$x_2$	$\cdot$	$\cdot$	$x_i$	$\cdot$	$\cdot$	$x_n$	$y$
				$\cdot$			$\cdot$			
Objeto	$id_1$	$x_1^1$	$x_2^1$	$\cdot$	$\cdot$	$x_i^1$	$\cdot$	$\cdot$	$x_n^1$	$y^1$
	$id_2$	$x_1^2$	$x_2^2$	$\cdot$	$\cdot$	$x_i^2$	$\cdot$	$\cdot$	$x_n^2$	$y^2$
	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$
	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$
	$id_j$	$x_1^j$	$x_2^j$	$\cdot$	$\cdot$	$x_{ji}$	$\cdot$	$\cdot$	$x_n^j$	$y^j$
	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$
	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$	$\cdot$
	$id_m$	$x_1^m$	$x_2^m$	$\cdot$	$\cdot$	$x_{mi}$	$\cdot$	$\cdot$	$x_n^m$	$y^m$

## Clasificación (ii)

---

- **Definición.-** La tarea de clasificación consiste en **aprender una función**  $f$  que asocia a cada vector  $x$  una de las clases predefinidas  $y$ .
- A la función  $f$  se la conoce también como **modelo de clasificación**.
- El modelo de clasificación puede adoptar diferentes formas: árbol de decisión, reglas, red neuronal ...
- El vector  $x$  puede tomar valores nominales o numéricos, pero la etiqueta  $y$  es siempre nominal. Cuando la variable  $y$  es numérica se construye un **árbol de regresión**.

## Utilidad de la clasificación

---

- **Modelo descriptivo.** Puede servir para distinguir entre objetos de diferentes clases identificando las características que las describen.
- **Modelo predictivo.** Puede usarse para predecir la clase a la que pertenece un objeto conociendo sus características. Es el uso que habitualmente se le da a la clasificación.

## Ejemplo 1

<i>ID</i>	PREVISIÓN	TEMPERATURA	HUMEDAD	VIENTO	JUGAR
<i>id</i> <sub>1</sub>	Soleado	Alta	Alta	Débil	No
<i>id</i> <sub>2</sub>	Soleado	Alta	Alta	Fuerte	No
<i>id</i> <sub>3</sub>	Nuboso	Alta	Alta	Débil	Sí
<i>id</i> <sub>4</sub>	Lluvioso	Media	Alta	Débil	Sí
<i>id</i> <sub>5</sub>	Lluvioso	Fría	Normal	Débil	Sí
<i>id</i> <sub>6</sub>	Lluvioso	Fría	Normal	Fuerte	No
<i>id</i> <sub>7</sub>	Nuboso	Fría	Normal	Fuerte	Sí
<i>id</i> <sub>8</sub>	Soleado	Media	Alta	Débil	No
<i>id</i> <sub>9</sub>	Soleado	Alta	Normal	Débil	Sí
<i>id</i> <sub>10</sub>	Lluvioso	Media	Normal	Débil	Sí
<i>id</i> <sub>11</sub>	Soleado	Media	Normal	Fuerte	Sí
<i>id</i> <sub>12</sub>	Nuboso	Media	Alta	Fuerte	Sí
<i>id</i> <sub>13</sub>	Nubos	Alta	Normal	Débil	Sí
<i>id</i> <sub>14</sub>	Lluvioso	Fría	Alta	Fuerte	No

## Ejemplo 2

---

<i>ID</i>	<b>CASA</b>	<b>ESTADO</b>	<b>INGRESOS</b>	<b>PRÉSTAMO</b>
<i>id<sub>1</sub></i>	Propiedad	Soltero	125000	Conceder
<i>id<sub>2</sub></i>	Alquiler	Casado	100000	Conceder
<i>id<sub>3</sub></i>	Alquiler	Soltero	70000	Conceder
<i>id<sub>4</sub></i>	Propiedad	Casado	12000	Conceder
<i>id<sub>5</sub></i>	Alquiler	Divorciado	95000	Denegar
<i>id<sub>6</sub></i>	Alquiler	Casado	60000	Conceder
<i>id<sub>7</sub></i>	Propiedad	Divorciado	220000	Conceder
<i>id<sub>8</sub></i>	Alquiler	Soltero	85000	Denegar
<i>id<sub>9</sub></i>	Alquiler	Casado	75000	Conceder
<i>id<sub>10</sub></i>	Alquiler	Soltero	90000	Conceder



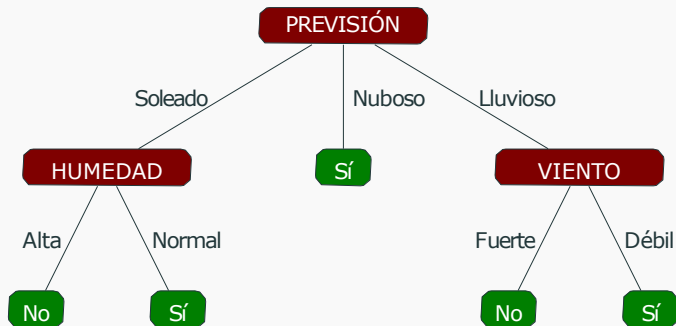
## ¿QUÉ SON Y CÓMO SE EMPLEAN?

---

# Árbol de decisión

## Ejemplo 1

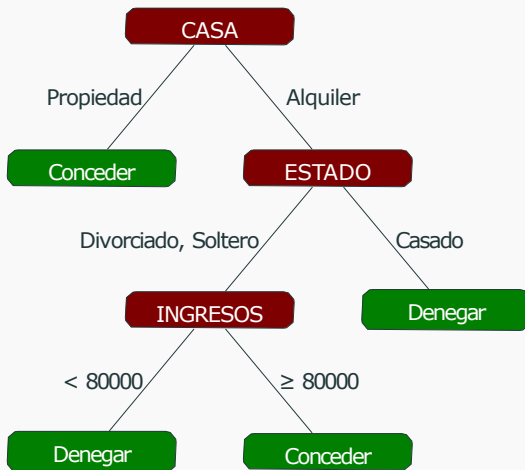
---



# Árbol de decisión

## Ejemplo 2

---



# Árbol de decisión

## Descripción y uso

---

- Cada nodo del árbol se corresponde con un atributo y de él parten tantas ramas como valores distintos tiene ese atributo.
- En las hojas del árbol se encuentran todos o algunos de los valores de la variable clase.
- Dada un árbol de decisión, para clasificar una nueva instancia se inspecciona el mismo desde la raíz hasta llegar a un nodo hoja.
- Cada nodo representa un test sobre un atributo y el valor correspondiente en la instancia indica la rama del árbol que debe recorrerse. El proceso se repite hasta alcanzar un nodo hoja. El valor de ese nodo suministra la clase a la que pertenece la instancia.

## Disyunciones de conjunciones (reglas)

---

- En general los árboles de decisión representan disyunciones de conjunciones de los valores de los atributos.
- Cada rama del árbol es una conjunción y el árbol en su conjunto una disyunción de esas conjunciones.

## Disyunciones de conjunciones (reglas)

---

• El árbol del ejemplo anterior se corresponde con el siguiente conjunto de reglas:

- **IF** (Previsión = Soleado) **and** (Humedad = alta) **THEN** (Jugar = No)
- **IF** (Previsión = Soleado) **and** (Humedad = normal) **THEN** (Jugar = Sí)
- **IF** (Previsión = Nuboso) **THEN** (Jugar = Sí)
- **IF** (Previsión = Lluvioso) **and** (Viento = Fuerte) **THEN** (Jugar = No)
- **IF** (Previsión = Lluvioso) **and** (Viento = Débil) **THEN** (Jugar = Sí)

## Problemas apropiados

---

- **Instancias representadas por pares atributo valor.** Los árboles son apropiados cuando cada atributo toma un número pequeño de valores.
- **La función objetivo toma valores discretos.** Aunque también existen algoritmos que permiten construir árboles de decisión cuando la variable de salida es continua.
- **Los datos de entrenamiento pueden contener errores.** Los algoritmos para construir árboles son robustos a errores de clasificación de los ejemplos de entrenamiento y a errores en los valores de los atributos.
- **Los datos de entrenamiento pueden contener valores desconocidos en algunos atributos.** Pueden construirse cuando algunos ejemplos de entrenamiento tienen valores desconocidos en algunos de los atributos.

## Problemas apropiados

---

- **Instancias representadas por pares atributo valor.** Los árboles son apropiados cuando cada atributo toma un número pequeño de valores.
- **La función objetivo toma valores discretos.** Aunque también existen algoritmos que permiten construir árboles de decisión cuando la variable de salida es continua.
- Los datos de entrenamiento pueden contener errores. Los algoritmos para construir árboles son robustos a errores de clasificación de los ejemplos de entrenamiento y a errores en los valores de los atributos.
- Los datos de entrenamiento pueden contener valores desconocidos en algunos atributos. Pueden construirse cuando algunos ejemplos de entrenamiento tienen valores desconocidos en algunos de los atributos.



## Problemas apropiados

---

- **Instancias representadas por pares atributo valor.** Los árboles son apropiados cuando cada atributo toma un número pequeño de valores.
- **La función objetivo toma valores discretos.** Aunque también existen algoritmos que permiten construir árboles de decisión cuando la variable de salida es continua.
- **Los datos de entrenamiento pueden contener errores.** Los algoritmos para construir árboles son robustos a errores de clasificación de los ejemplos de entrenamiento y a errores en los valores de los atributos.
- **Los datos de entrenamiento pueden contener valores desconocidos en algunos atributos.** Pueden construirse cuando algunos ejemplos de entrenamiento tienen valores desconocidos en algunos de los atributos.

## Problemas apropiados

---

- **Instancias representadas por pares atributo valor.** Los árboles son apropiados cuando cada atributo toma un número pequeño de valores.
- **La función objetivo toma valores discretos.** Aunque también existen algoritmos que permiten construir árboles de decisión cuando la variable de salida es continua.
- **Los datos de entrenamiento pueden contener errores.** Los algoritmos para construir árboles son robustos a errores de clasificación de los ejemplos de entrenamiento y a errores en los valores de los atributos.
- **Los datos de entrenamiento pueden contener valores desconocidos en algunos atributos.** Pueden construirse cuando algunos ejemplos de entrenamiento tienen valores desconocidos en algunos de los atributos.