

Ficha Técnica Avanzada: Árboles de Decisión

Jordi Pozo
I.E.S. Ribera de Castilla

Curso 2025/26

1. Concepto

Un **Árbol de Decisión** es un modelo de **aprendizaje supervisado** (utilizado tanto para clasificación como para regresión) que toma decisiones basándose en una estructura de grafo similar a un árbol.

El modelo aprende reglas de decisión simples (p. ej. “¿Es el ingreso > 50.000 €?”) inferidas a partir de las características de los datos. La estructura se compone de:

- **Nodo Raíz (Root Node):** El nodo inicial que representa a todo el conjunto de datos.
- **Nodos Internos (Internal Nodes):** Nodos donde se evalúa una característica para dividir los datos.
- **Ramas (Branches):** Representan el resultado de una decisión (p. ej. “Sí” o “No”).
- **Nodos Hoja (Leaf Nodes):** Nodos terminales que representan el resultado final.

2. ¿Cómo aprende el árbol? Métricas de División

El algoritmo (como CART o C4.5) construye el árbol de forma recursiva. En cada nodo, busca la “mejor” característica y el “mejor” umbral para dividir los datos: “mejor” significa producir nodos hijos lo más “puros” posible.

2.1. Impureza de Gini (Gini Impurity)

Es la métrica por defecto en `scikit-learn` para clasificación:

$$G = 1 - \sum_{i=1}^C p_i^2$$

- **Gini = 0 (Puro):** $1 - (1^2) = 0$.
- **Gini = 0.5 (2 clases, 50/50):** $1 - (0,5^2 + 0,5^2) = 0,5$.

2.2. Entropía y Ganancia de Información

$$H = - \sum_{i=1}^C p_i \log_2 p_i$$

La **ganancia de información** (IG) es la reducción de entropía tras dividir:

$$IG(\text{Padre}, \text{Split}) = H(\text{Padre}) - \sum_{\text{hijo}} \frac{N_{\text{hijo}}}{N_{\text{Padre}}} H(\text{hijo})$$

En `scikit-learn`: `criterion='entropy'`.

3. Ilustración del Modelo

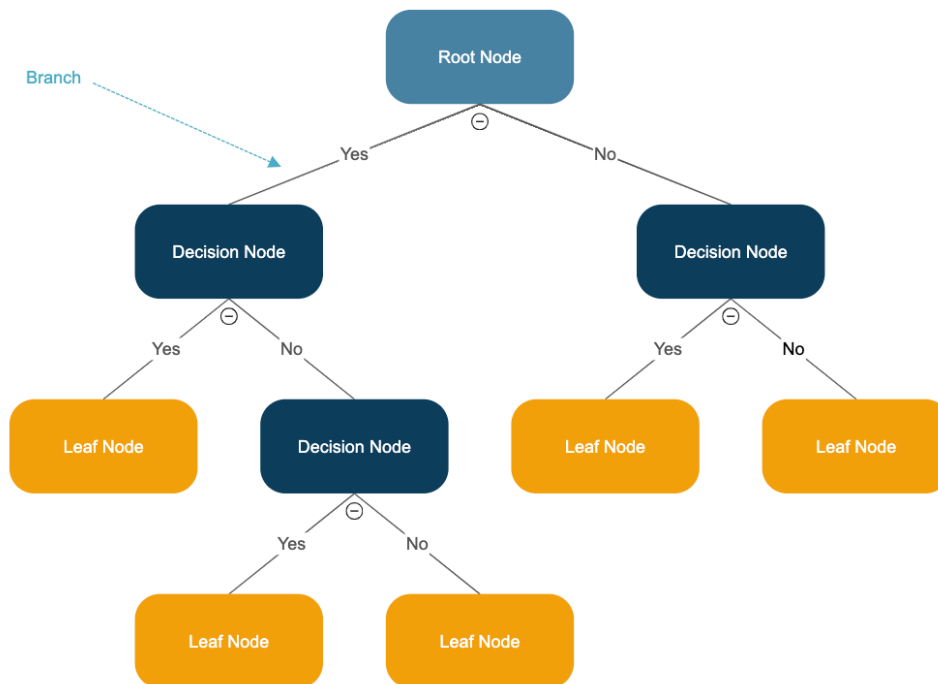


Figura 1: Representación gráfica de un árbol de decisión (placeholder).

4. Hiperparámetros Clave

4.1. Parámetros genéricos

- **criterion** (Gini/Entropía en clasificación; ECM en regresión).
- **max_depth**: tope de niveles del árbol.
- **min_samples_split**: mínimo para intentar dividir.
- **min_samples_leaf**: mínimo por hoja tras dividir.
- **ccp_alpha**: poda por complejidad.

4.2. Implementación en `scikit-learn`

- **criterion**: 'gini'/'entropy' ('squared_error' en regresión).
- **max_depth**: None por defecto; conviene fijarlo (p. ej. 3, 5, 8).
- **min_samples_split**: por defecto 2.
- **min_samples_leaf**: por defecto 1 (sube para suavizar).
- **ccp_alpha**: 0.0 por defecto (sin poda).

5. Características

5.1. Ventajas

- Interpretables ("white box").

- Poco preprocesado.
- Selección implícita de variables.

5.2. Desventajas

- Sobreajuste si no se limita.
- Inestables a pequeñas variaciones.
- Sesgo hacia variables con muchos niveles.

6. Usos

- **Clasificación:** diagnóstico, riesgo crediticio, spam.
- **Regresión:** precios, demanda.
- **Base de ensembles:** Random Forest, Gradient Boosting (XGBoost).

7. Ejemplos con Python y Scikit-learn

7.1. Ejemplo 1: Clasificación (Iris)

```

1  import pandas as pd
2  from sklearn.datasets import load_iris
3  from sklearn.model_selection import train_test_split
4  from sklearn.tree import DecisionTreeClassifier, plot_tree
5  from sklearn.metrics import accuracy_score, confusion_matrix
6  import matplotlib.pyplot as plt
7
8  # 1. Cargar datos
9  iris = load_iris()
10 X = pd.DataFrame(iris.data, columns=iris.feature_names)
11 y = iris.target
12
13 # 2. Dividir
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
15                                                    random_state=42)
16
17 # 3. Crear y entrenar (Ajustando max_depth)
18 clf = DecisionTreeClassifier(criterion='gini', max_depth=3, random_state=42)
19 clf.fit(X_train, y_train)
20
21 # 4. Evaluar
22 y_pred = clf.predict(X_test)
23 print(f"Accuracy (Precision): {accuracy_score(y_test, y_pred):.2f}")
24 print("Matriz de Confusion:\n", confusion_matrix(y_test, y_pred))
25
26 # 5. Visualizar el arbol
27 plt.figure(figsize=(15, 10))
28 plot_tree(clf,
29          filled=True,
30          rounded=True,
31          class_names=iris.target_names,
32          feature_names=iris.feature_names)
33 plt.title("Arbol de Decision - Clasificacion Iris (max_depth=3)")
34
35 # Guardar la imagen para usarla en LaTeX:

```

```
35 # plt.savefig("decision_tree_diagram.png")
```

Listing 1: Clasificación con Iris (criterion='gini', max_depth=3).

7.2. Ejemplo 2: Regresión (Diabetes)

```
1  from sklearn.datasets import load_diabetes
2  from sklearn.tree import DecisionTreeRegressor
3  from sklearn.metrics import mean_squared_error, r2_score
4  import numpy as np
5  import pandas as pd
6  from sklearn.model_selection import train_test_split
7  from sklearn.tree import plot_tree
8  import matplotlib.pyplot as plt
9
10 # 1. Cargar datos
11 diabetes = load_diabetes()
12 X = pd.DataFrame(diabetes.data, columns=diabetes.feature_names)
13 y = diabetes.target
14
15 # 2. Dividir
16 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
17                                                    random_state=42)
18
19 # 3. Crear y entrenar
20 reg = DecisionTreeRegressor(criterion='squared_error', max_depth=4, random_
21                             state=42)
22 reg.fit(X_train, y_train)
23
24 # 4. Evaluar
25 y_pred_reg = reg.predict(X_test)
26 mse = mean_squared_error(y_test, y_pred_reg)
27 r2 = r2_score(y_test, y_pred_reg)
28 print(f"MSE: {mse:.2f}")
29 print(f"R2: {r2:.2f}")
30
31 # 5. Visualizar
32 plt.figure(figsize=(20, 12))
33 plot_tree(reg, filled=True, rounded=True, feature_names=diabetes.feature_names
34           )
35 plt.title("Arbol de Decision - Regresion Diabetes (max_depth=4)")
36 # plt.savefig("diabetes_tree.png")
```

Listing 2: Regresión con Diabetes (criterion='squared_error', max_depth=4).