

# Guía de Laboratorio - Semana 5: Árboles de Decisión

Inteligencia Artificial (100000S14F)

Ingeniería de Software

Ciclo 2 Agosto 2025

13 de septiembre de 2025

## 1 Objetivos de la Sesión

Al finalizar esta sesión de laboratorio, el estudiante será capaz de:

- Comprender los fundamentos teóricos y matemáticos de los árboles de decisión
- Implementar un árbol de decisión utilizando scikit-learn con diversas configuraciones
- Realizar preprocesamiento de datos adecuado para modelos de árboles de decisión
- Evaluar el rendimiento del modelo utilizando múltiples métricas
- Visualizar e interpretar la estructura del árbol de decisión
- Optimizar hiperparámetros mediante validación cruzada
- Aplicar técnicas para evitar sobreajuste (overfitting)

## 2 Marco Teórico

Los árboles de decisión son estructuras jerárquicas utilizadas para clasificación y regresión. Constan de:

- **Nodos raíz e internos:** Representan pruebas sobre atributos
- **Ramas:** Resultados posibles de las pruebas
- **Nodos hoja:** Decisiones finales (clases o valores)

### 2.1 Algoritmos de Construcción

El algoritmo ID3 (Iterative Dichotomiser 3) utiliza la **ganancia de información** basada en la **entropía**:

$$\text{Entropía: } H(S) = - \sum_{i=1}^c p_i \log_2 p_i \quad (1)$$

$$\text{Ganancia de informaci3n: } IG(S, A) = H(S) - \sum_{t \in T} p(t)H(t) \quad (2)$$

El algoritmo C4.5 (extensi3n de ID3) utiliza la **ganancia de raz3n**:

$$\text{Ganancia de raz3n: } GR(S, A) = \frac{IG(S, A)}{H(A)} \quad (3)$$

El algoritmo CART utiliza el **3ndice Gini**:

$$\text{3ndice Gini: } Gini(S) = 1 - \sum_{i=1}^c p_i^2 \quad (4)$$

## 3 Configuraci3n del Entorno

### 3.1 Requisitos Previos

```
1 # Instalar las bibliotecas necesarias
2 pip install scikit-learn==1.2.2 pandas==2.0.1 matplotlib==3.7.1 numpy
   ==1.24.3 graphviz==0.20.1
```

### 3.2 Configuraci3n de Entorno de Desarrollo

Recomendamos utilizar Jupyter Notebook o Google Colab para este laboratorio. Para configurar el entorno:

```
1 # Crear un entorno virtual (opcional pero recomendado)
2 python -m venv arboles_decision
3 source arboles_decision/bin/activate # Linux/Mac
4 # o
5 .\arboles_decision\Scripts\activate # Windows
6
7 # Instalar paquetes
8 pip install -r requirements.txt
```

## 4 Implementaci3n Paso a Paso

### 4.1 Paso 1: Importar Bibliotecas

```
1 # Bibliotecas fundamentales
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from IPython.display import Image, display
7
8 # Scikit-learn para rboles de decisi n
9 from sklearn.datasets import load_iris
10 from sklearn.tree import DecisionTreeClassifier, export_graphviz,
   plot_tree
```

```

11 from sklearn.model_selection import train_test_split, cross_val_score,
    GridSearchCV
12 from sklearn.metrics import (accuracy_score, confusion_matrix,
    classification_report,
13                               precision_score, recall_score, f1_score,
                                roc_curve, auc)
14 from sklearn.preprocessing import label_binarize
15 from sklearn.utils import resample
16
17 # Graphviz para visualizaci n
18 import graphviz
19
20 # Configuraci n de estilo
21 plt.style.use('seaborn-v0_8-whitegrid')
22 sns.set_palette("husl")
23 %matplotlib inline

```

## 4.2 Paso 2: Cargar y Explorar Datos

```

1 # Cargar dataset
2 iris = load_iris()
3 X = iris.data
4 y = iris.target
5
6 # Crear DataFrame para mejor visualizaci n
7 df = pd.DataFrame(X, columns=iris.feature_names)
8 df['target'] = y
9 df['species'] = df['target'].apply(lambda x: iris.target_names[x])
10
11 print("=== INFORMACI N DEL DATASET IRIS ===")
12 print(f"Dimensiones: {X.shape}")
13 print(f"Caracter sticas: {iris.feature_names}")
14 print(f"Clases: {iris.target_names}")
15
16 print("\n=== PRIMERAS FILAS ===")
17 print(df.head())
18
19 print("\n=== ESTAD STICAS DESCRIPTIVAS ===")
20 print(df.describe())
21
22 print("\n=== DISTRIBUCI N DE CLASES ===")
23 print(df['species'].value_counts())
24
25 # Visualizaci n de distribuciones
26 fig, axes = plt.subplots(2, 2, figsize=(12, 10))
27 for i, feature in enumerate(iris.feature_names):
28     row, col = i // 2, i % 2
29     for species in iris.target_names:
30         species_data = df[df['species'] == species][feature]
31         axes[row, col].hist(species_data, alpha=0.7, label=species)
32         axes[row, col].set_title(f'Distribuci n de {feature}')
33         axes[row, col].set_xlabel(feature)
34         axes[row, col].set_ylabel('Frecuencia')
35         axes[row, col].legend()
36
37 plt.tight_layout()

```

```

38 plt.savefig('distribuciones_caracteristicas.png', dpi=300, bbox_inches=
    'tight')
39 plt.show()
40
41 # Matriz de correlaci n
42 plt.figure(figsize=(10, 8))
43 correlation_matrix = df[iris.feature_names].corr()
44 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', center=0)
45 plt.title('Matriz de Correlaci n de Caracter sticas')
46 plt.savefig('matriz_correlacion.png', dpi=300, bbox_inches='tight')
47 plt.show()

```

### 4.3 Paso 3: Preprocesamiento de Datos

```

1 # Dividir en conjunto de entrenamiento y prueba
2 X_train, X_test, y_train, y_test = train_test_split(
3     X, y,
4     test_size=0.3,
5     random_state=42,
6     stratify=y, # Mantener proporci n de clases
7     shuffle=True
8 )
9
10 print(f"Tama o conjunto de entrenamiento: {X_train.shape}")
11 print(f"Tama o conjunto de prueba: {X_test.shape}")
12
13 # Verificar distribuci n de clases en ambos conjuntos
14 print("\nDistribuci n en entrenamiento:")
15 unique_train, counts_train = np.unique(y_train, return_counts=True)
16 for cls, count in zip(unique_train, counts_train):
17     print(f"Clase {iris.target_names[cls]}: {count} muestras ({count/
18         len(y_train)*100:.2f}%)")
19
20 print("\nDistribuci n en prueba:")
21 unique_test, counts_test = np.unique(y_test, return_counts=True)
22 for cls, count in zip(unique_test, counts_test):
23     print(f"Clase {iris.target_names[cls]}: {count} muestras ({count/
24         len(y_test)*100:.2f}%)")

```

### 4.4 Paso 4: Crear y Entrenar el Modelo

```

1 # Configuraci n de hiperpar metros
2 parametros_arbol = {
3     'criterion': 'entropy', # Criterio de divisi n: 'gini' o '
    'entropy'
4     'max_depth': 3, # Profundidad m xima del rbol
5     'min_samples_split': 2, # M nimo de muestras para dividir un
    nodo
6     'min_samples_leaf': 1, # M nimo de muestras en un nodo hoja
7     'max_features': None, # N mero de caracter sticas a
    considerar
8     'random_state': 42, # Semilla para reproducibilidad
9     'ccp_alpha': 0.0, # Par metro de podado cost-
    complexity

```

```

10     'class_weight': None                # Pesos de clases ( til para
      datasets desbalanceados)
11 }
12
13 # Crear el clasificador de rbol de decisi n
14 clf = DecisionTreeClassifier(**parametros_arbol)
15
16 # Entrenar el modelo
17 clf.fit(X_train, y_train)
18
19 # Hacer predicciones
20 y_pred = clf.predict(X_test)
21 y_pred_proba = clf.predict_proba(X_test) # Probabilidades para cada
      clase
22
23 # Mostrar informaci n del rbol resultante
24 print("=== INFORMACI N DEL RBOL ===")
25 print(f"N mero de nodos: {clf.tree_.node_count}")
26 print(f"Profundidad del rbol : {clf.get_depth()}")
27 print(f"N mero de hojas: {clf.get_n_leaves()}")

```

## 4.5 Paso 5: Evaluar el Modelo

```

1 # Calcular m tricas de evaluaci n
2 accuracy = accuracy_score(y_test, y_pred)
3 precision = precision_score(y_test, y_pred, average='weighted')
4 recall = recall_score(y_test, y_pred, average='weighted')
5 f1 = f1_score(y_test, y_pred, average='weighted')
6
7 print("=== M TRICAS DE EVALUACI N ===")
8 print(f"Precisi n (Accuracy): {accuracy:.4f}")
9 print(f"Precisi n (Precision): {precision:.4f}")
10 print(f"Sensibilidad (Recall): {recall:.4f}")
11 print(f"Puntuaci n F1: {f1:.4f}")
12
13 # Matriz de confusi n
14 plt.figure(figsize=(8, 6))
15 cm = confusion_matrix(y_test, y_pred)
16 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
17             xticklabels=iris.target_names,
18             yticklabels=iris.target_names)
19 plt.title('Matriz de Confusi n')
20 plt.ylabel('Etiqueta Real')
21 plt.xlabel('Etiqueta Predicha')
22 plt.savefig('matriz_confusion.png', dpi=300, bbox_inches='tight')
23 plt.show()
24
25 # Reporte de clasificaci n detallado
26 print("\n=== REPORTE DE CLASIFICACI N ===")
27 print(classification_report(y_test, y_pred, target_names=iris.
      target_names))
28
29 # Curvas ROC (para clasificaci n multiclase)
30 y_test_bin = label_binarize(y_test, classes=[0, 1, 2])
31 n_classes = y_test_bin.shape[1]
32

```

```

33 fpr = dict()
34 tpr = dict()
35 roc_auc = dict()
36
37 for i in range(n_classes):
38     fpr[i], tpr[i], _ = roc_curve(y_test_bin[:, i], y_pred_proba[:, i])
39     roc_auc[i] = auc(fpr[i], tpr[i])
40
41 # Plot all ROC curves
42 plt.figure(figsize=(10, 8))
43 colors = ['aqua', 'darkorange', 'cornflowerblue']
44 for i, color in zip(range(n_classes), colors):
45     plt.plot(fpr[i], tpr[i], color=color, lw=2,
46             label='ROC curve of class {0} (area = {1:0.2f})'
47             ''.format(iris.target_names[i], roc_auc[i]))
48
49 plt.plot([0, 1], [0, 1], 'k--', lw=2)
50 plt.xlim([0.0, 1.0])
51 plt.ylim([0.0, 1.05])
52 plt.xlabel('False Positive Rate')
53 plt.ylabel('True Positive Rate')
54 plt.title('Curvas ROC para Clasificaci n Multiclase')
55 plt.legend(loc="lower right")
56 plt.savefig('curvas_roc.png', dpi=300, bbox_inches='tight')
57 plt.show()

```

## 4.6 Paso 6: Visualizar el Árbol de Decisión

```

1 # Visualizaci n con Graphviz (alta calidad)
2 dot_data = export_graphviz(
3     clf,
4     out_file=None,
5     feature_names=iris.feature_names,
6     class_names=iris.target_names,
7     filled=True,          # Nodos coloreados
8     rounded=True,        # Bordes redondeados
9     special_characters=True,
10    proportion=True,      # Mostrar proporciones en lugar de counts
11    impurity=True,        # Mostrar impureza
12    node_ids=True         # Mostrar IDs de nodos
13 )
14
15 graph = graphviz.Source(dot_data)
16 graph.format = 'png'
17 graph.render('arbol_decision_iris', view=False, cleanup=True)
18
19 # Visualizaci n con matplotlib (m s simple pero interactiva)
20 plt.figure(figsize=(20, 12))
21 plot_tree(clf,
22           feature_names=iris.feature_names,
23           class_names=iris.target_names,
24           filled=True,
25           rounded=True,
26           fontsize=10)
27 plt.title(" rbol de Decisi n - Dataset Iris")

```

```

28 plt.savefig('arbol_decision_matplotlib.png', dpi=300, bbox_inches='
    tight')
29 plt.show()
30
31 # Mostrar importancia de caracter sticas
32 importances = clf.feature_importances_
33 indices = np.argsort(importances)[::-1]
34
35 plt.figure(figsize=(10, 6))
36 plt.title("Importancia de las Caracter sticas")
37 plt.bar(range(X.shape[1]), importances[indices], align="center")
38 plt.xticks(range(X.shape[1]), [iris.feature_names[i] for i in indices])
39 plt.xlim([-1, X.shape[1]])
40 plt.tight_layout()
41 plt.savefig('importancia_caracteristicas.png', dpi=300, bbox_inches='
    tight')
42 plt.show()
43
44 print("Importancia de caracter sticas:")
45 for i, idx in enumerate(indices):
46     print(f"{i+1}. {iris.feature_names[idx]}: {importances[idx]:.4f}")

```

## 5 Experimentación con Hiperparámetros

### 5.1 Análisis de Profundidad del Árbol

```

1 # Probar diferentes profundidades m ximas
2 profundidades = range(1, 21)
3 train_scores = []
4 test_scores = []
5 cv_scores = []
6
7 for profundidad in profundidades:
8     # Crear y entrenar modelo
9     clf_temp = DecisionTreeClassifier(max_depth=profundidad,
        random_state=42)
10    clf_temp.fit(X_train, y_train)
11
12    # Calcular precisi n en entrenamiento y prueba
13    train_score = clf_temp.score(X_train, y_train)
14    test_score = clf_temp.score(X_test, y_test)
15
16    # Calcular validaci n cruzada
17    cv_score = cross_val_score(clf_temp, X, y, cv=5).mean()
18
19    train_scores.append(train_score)
20    test_scores.append(test_score)
21    cv_scores.append(cv_score)
22
23    print(f"Profundidad: {profundidad:2d} | Train: {train_score:.4f} |
        Test: {test_score:.4f} | CV: {cv_score:.4f}")
24
25 # Graficar resultados
26 plt.figure(figsize=(12, 6))

```

```

27 plt.plot(profundidades, train_scores, 'o-', label='Precisi n en
    Entrenamiento')
28 plt.plot(profundidades, test_scores, 'o-', label='Precisi n en Prueba'
    )
29 plt.plot(profundidades, cv_scores, 'o-', label='Validaci n Cruzada (5-
    fold)')
30 plt.xlabel('Profundidad M xima del rbol ')
31 plt.ylabel('Precisi n')
32 plt.title('Efecto de la Profundidad del rbol en el Rendimiento')
33 plt.legend()
34 plt.grid(True)
35 plt.savefig('profundidad_vs_rendimiento.png', dpi=300, bbox_inches='
    tight')
36 plt.show()

```

## 5.2 Optimizaci3n de Hiperpar3metros con Grid Search

```

1 # Definir espacio de b squeda de hiperpar metros
2 param_grid = {
3     'criterion': ['gini', 'entropy'],
4     'max_depth': [3, 5, 7, 10, None],
5     'min_samples_split': [2, 5, 10],
6     'min_samples_leaf': [1, 2, 5],
7     'max_features': [None, 'sqrt', 'log2'],
8     'ccp_alpha': [0.0, 0.01, 0.1] # Par metro de podado
9 }
10
11 # Crear el modelo base
12 dt = DecisionTreeClassifier(random_state=42)
13
14 # Configurar la b squeda en grid
15 grid_search = GridSearchCV(
16     estimator=dt,
17     param_grid=param_grid,
18     scoring='accuracy',
19     cv=5, # Validaci n cruzada de 5 folds
20     n_jobs=-1, # Usar todos los cores disponibles
21     verbose=1 # Mostrar progreso
22 )
23
24 # Ejecutar la b squeda en grid
25 grid_search.fit(X_train, y_train)
26
27 # Mejores par metros y resultados
28 print("Mejores par metros encontrados:")
29 print(grid_search.best_params_)
30 print(f"\nMejor precisi n en validaci n cruzada: {grid_search.
    best_score_:.4f}")
31
32 # Evaluar el mejor modelo en el conjunto de prueba
33 best_clf = grid_search.best_estimator_
34 y_pred_best = best_clf.predict(X_test)
35 best_accuracy = accuracy_score(y_test, y_pred_best)
36 print(f"Precisi n del mejor modelo en prueba: {best_accuracy:.4f}")
37
38 # Comparar con el modelo inicial

```



```

39 print(f"Mejora respecto al modelo inicial: {(best_accuracy - accuracy)
    :.4f}")
40
41 # Visualizar el mejor rbol
42 plt.figure(figsize=(20, 12))
43 plot_tree(best_clf,
44           feature_names=iris.feature_names,
45           class_names=iris.target_names,
46           filled=True,
47           rounded=True,
48           fontsize=10)
49 plt.title("Mejor rbol de Decisi n despu s de la Optimizaci n")
50 plt.savefig('mejor_arbol_decision.png', dpi=300, bbox_inches='tight')
51 plt.show()

```

## 6 Ejercicios Prácticos

### 6.1 Ejercicio 1: Análisis de Importancia de Características

1. Extrae la importancia de cada característica del modelo entrenado
2. Visualiza la importancia de características en un gráfico de barras horizontal
3. Interpreta qué características son más relevantes para la clasificación
4. Compara los resultados con la matriz de correlación generada anteriormente

**Solución:**

```

1 # Obtener importancia de caracter sticas
2 importances = best_clf.feature_importances_
3 feature_names = iris.feature_names
4
5 # Crear DataFrame para mejor visualizaci n
6 importance_df = pd.DataFrame({
7     'Caracter stica': feature_names,
8     'Importancia': importances
9 }).sort_values('Importancia', ascending=True)
10
11 # Gr fico de barras horizontal
12 plt.figure(figsize=(10, 6))
13 plt.barh(importance_df['Caracter stica'], importance_df['Importancia',
14 ])
15 plt.xlabel('Importancia')
16 plt.title('Importancia de Caracter sticas en el Mejor Modelo')
17 for i, v in enumerate(importance_df['Importancia']):
18     plt.text(v + 0.001, i, f'{v:.4f}', va='center')
19 plt.tight_layout()
20 plt.savefig('importancia_caracteristicas_mejor_modelo.png', dpi=300,
21           bbox_inches='tight')
22 plt.show()
23
24 # An lisis
25 print("An lisis de importancia de caracter sticas:")
26 for i, row in importance_df.iterrows():
27     print(f"- {row['Caracter stica']}: {row['Importancia']:.4f}")

```

## 6.2 Ejercicio 2: Validación Cruzada y Estabilidad del Modelo

1. Realiza validación cruzada con diferentes semillas aleatorias
2. Evalúa la estabilidad del modelo con diferentes divisiones de datos
3. Calcula intervalos de confianza para la precisión del modelo

**Solución:**

```
1 # Validación cruzada con diferentes semillas
2 n_iterations = 10
3 cv_scores = []
4
5 for i in range(n_iterations):
6     # Crear modelo con diferentes semillas
7     clf_cv = DecisionTreeClassifier(**grid_search.best_params_)
8     clf_cv.random_state = i # Cambiar semilla
9
10    # Realizar validación cruzada
11    scores = cross_val_score(clf_cv, X, y, cv=5, scoring='accuracy')
12    cv_scores.append(scores.mean())
13    print(f"Iteración {i+1}: Precisión CV = {scores.mean():.4f}")
14
15 # Calcular estadísticas
16 mean_accuracy = np.mean(cv_scores)
17 std_accuracy = np.std(cv_scores)
18 confidence_interval = 1.96 * std_accuracy / np.sqrt(n_iterations) #
19    95% IC
20
21 print(f"\nPrecisión media: {mean_accuracy:.4f}")
22 print(f"Desviación estándar: {std_accuracy:.4f}")
23 print(f"Intervalo de confianza (95%): ({mean_accuracy -
24    confidence_interval:.4f}, {mean_accuracy + confidence_interval:.4f})")
25
26 # Visualizar distribución de precisiones
27 plt.figure(figsize=(10, 6))
28 plt.hist(cv_scores, bins=10, edgecolor='black', alpha=0.7)
29 plt.axvline(mean_accuracy, color='red', linestyle='--', label=f'Media:
30    {mean_accuracy:.4f}')
31 plt.xlabel('Precisión')
32 plt.ylabel('Frecuencia')
33 plt.title('Distribución de Precisiones en Validación Cruzada')
34 plt.legend()
35 plt.savefig('distribucion_precision_cv.png', dpi=300, bbox_inches='
36    tight')
37 plt.show()
```

## 6.3 Ejercicio 3: Aplicación a un Dataset Diferente

1. Carga el dataset de cáncer de mama de scikit-learn
2. Aplica el mismo proceso de preprocesamiento y modelado
3. Compara los resultados con los obtenidos con el dataset Iris

## Solución:

```
1 from sklearn.datasets import load_breast_cancer
2
3 # Cargar dataset de cáncer de mama
4 cancer = load_breast_cancer()
5 X_cancer = cancer.data
6 y_cancer = cancer.target
7
8 print(f"Dataset de cáncer de mama: {X_cancer.shape}")
9 print(f"Características: {cancer.feature_names}")
10 print(f"Clases: {cancer.target_names}")
11
12 # Preprocesamiento (escalado recomendado para este dataset)
13 from sklearn.preprocessing import StandardScaler
14
15 scaler = StandardScaler()
16 X_cancer_scaled = scaler.fit_transform(X_cancer)
17
18 # División en train/test
19 X_train_c, X_test_c, y_train_c, y_test_c = train_test_split(
20     X_cancer_scaled, y_cancer, test_size=0.3, random_state=42, stratify
21     =y_cancer
22 )
23
24 # Entrenar modelo con mejores parámetros del ejercicio anterior
25 clf_cancer = DecisionTreeClassifier(**grid_search.best_params_)
26 clf_cancer.fit(X_train_c, y_train_c)
27
28 # Evaluar
29 y_pred_c = clf_cancer.predict(X_test_c)
30 accuracy_c = accuracy_score(y_test_c, y_pred_c)
31 print(f"Precisión en dataset de cáncer de mama: {accuracy_c:.4f}")
32
33 # Comparar con modelo básico
34 clf_basic = DecisionTreeClassifier(random_state=42)
35 clf_basic.fit(X_train_c, y_train_c)
36 y_pred_basic = clf_basic.predict(X_test_c)
37 accuracy_basic = accuracy_score(y_test_c, y_pred_basic)
38 print(f"Precisión con modelo básico: {accuracy_basic:.4f}")
39 print(f"Mejora con parámetros optimizados: {accuracy_c -
40     accuracy_basic:.4f}")
```

## 7 Recursos Adicionales

### 7.1 Bibliografía Recomendada

- Mitchell, T. M. (1997). *Machine Learning*. McGraw-Hill. (Capítulo 3)
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning*. Springer. (Capítulo 9)
- Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly. (Capítulo 6)

## 7.2 Enlaces Útiles

- Documentación de scikit-learn: <https://scikit-learn.org/stable/modules/tree.html>
- Tutorial interactivo de árboles de decisión: <https://www.datacamp.com/community/tutorials/decision-tree-classification-python>
- Visualización interactiva de árboles de decisión: <http://www.r2d3.us/visual-intro-to-machine-learning-part-1.html>
- Repositorio de datasets para práctica: <https://archive.ics.uci.edu/ml/index.php>

## 7.3 Herramientas Alternativas

- WEKA: Herramienta gráfica para minería de datos
- RapidMiner: Plataforma de ciencia de datos con interfaz visual
- KNIME: Plataforma de analítica de datos de código abierto