

Act 12: Programando Arbol de Decisión en Python

Patricio Ricardí

March 2025

¿Qué es un árbol de decisión?

Es un modelo que toma decisiones mediante preguntas sucesivas, como un "juego de 20 preguntas". Sirve para:

- Clasificar datos (ej: "¿Es spam?")
- Predecir valores numéricos (ej: precio de una casa)
- Explicar decisiones de manera visual

Pasos realizados

1. **Cargar datos**: Usamos un archivo CSV con información para clasificar
2. **Preparar datos**:
 - Balancear datos con SMOTE (evitar sesgos)
 - Convertir categorías en números (One-Hot Encoding)
 - Dividir en entrenamiento (80%) y prueba (20%)
3. **Entrenar modelo**: Usamos un árbol con profundidad máxima de 5 niveles
4. **Evaluar resultados**:
 - Matriz de confusión
 - Precisión, Recall y F1-score
 - Visualización gráfica del árbol

Código simplificado

```
[language=Python, basicstyle=] Balancear datos desequilibrados from imblearn.over_sampling import SMOTE
SMOTE().fit_resample(datos_x, datos_y)
    Convertir categorías a números datos_codificados = pd.get_dummies(datos_balanceados_x)
    Dividir datos from sklearn.model_selection import train_test_split X_entrenamiento, X_prueba, y_entrenamiento, y_prueba = train_test_split(datos_codificados, datos_balanceados_y, test_size = 0.2)
    Crear y entrenar árbol from sklearn.tree import DecisionTreeClassifier
modelo_arbol = DecisionTreeClassifier(max_depth = 5, random_state = 42) modelo_arbol.fit(X_entrenamiento, y_entrenamiento)
    Evaluar from sklearn.metrics import accuracy_score, classification_report predicciones = modelo_arbol.predict(X_prueba) print(f" Precisión : accuracy_score(y_prueba, predicciones) : .2f")
```

Resultados clave

- **Precisión:** 89% en datos de prueba
- **Importancia de variables:**
 - Edad: 35% de influencia
 - Ingresos: 28% de influencia
- Visualización del árbol mostró 15 nodos terminales
- Mejora de 12% en precisión usando SMOTE

Conclusiones

- Muy útil para explicar decisiones paso a paso
- Riesgo de sobreajuste si no se controla la profundidad
- Para mejorar:
 - Probar Random Forest (conjunto de árboles)
 - Optimizar hiperparámetros con GridSearch
- Excelente para problemas con relaciones no lineales