

# Práctica 5

## Árboles de decisión

Mario Emilio Jiménez Vizcaíno  
A01173359@itesm.mx  
Tecnológico de Monterrey  
Ingeniería en Tecnologías Computacionales  
Monterrey, N.L., México

Jesus Abraham Haros Madrid  
A01252642@itesm.mx  
Tecnológico de Monterrey  
Ingeniería en Tecnologías Computacionales  
Monterrey, N.L., México

### ABSTRACT

En la actualidad, uno de los problemas más importantes para los científicos es la clasificación: el etiquetado de elementos, basándose en las características de estos, seleccionando de un conjunto de clases, una que mejor lo represente. Los árboles de decisión, el objeto de estudio en esta práctica, son clasificadores que predicen las clases de los elementos usando algoritmos simples, lo que facilita su uso e implementación.

## 1 INTRODUCCIÓN

Un árbol de decisión clasifica instancias de datos planteando una serie de preguntas sobre las características de estos elementos. Cada pregunta se representa con un nodo, y cada nodo apunta a un nodo hijo, que puede ser un nodo terminal (que presenta el resultado del árbol: una clase o etiqueta), u otro nodo de decisión. Las preguntas forman así una jerarquía de decisiones capturada en una estructura de árbol.

Para clasificar un elemento se sigue el camino desde el nodo superior o raíz, hasta un nodo terminal, dependiendo las características del nodo y las preguntas que cada hoja del camino presenten.

Una ventaja de los árboles de decisión es que muchas veces son más interpretables que otros clasificadores, como las redes neuronales y las máquinas de vectores de soporte[1], porque combinan preguntas sencillas sobre los datos de forma comprensible. Por desgracia, pequeños cambios en los datos de entrada pueden provocar a veces grandes cambios en el árbol construido. Los árboles de decisión son lo suficientemente flexibles como para manejar elementos con una mezcla de características de valor real y categóricas, así como elementos con algunas características ausentes.

## 2 CONCEPTOS PREVIOS

- Programación básica en Python
- Conocimiento de las librerías *scikit-learn*, *matplotlib* y *numpy*
- Conocimientos básicos de estadística

## 3 METODOLOGÍA

### 3.1 Dataset Iris

TODO

### 3.2 Dataset Wine

TODO

### 3.3 Dataset Breast Cancer

TODO

### 3.4 Modelo de árbol de decisión

Para la generación del modelo del árbol de decisión se utilizó la implementación de la librería *sklearn*, específicamente la clase *sklearn.tree.DecisionTreeClassifier*[2], que, aunque provee la función de elegir qué algoritmo utilizar para medir la calidad de las preguntas dentro del árbol, utiliza por defecto el algoritmo de impureza de Gini.

El código que ejecutamos para realizar el análisis del dataset se encuentra en el apéndice A.

## 4 RESULTADOS

TODO

### 4.1 Dataset Iris

TODO

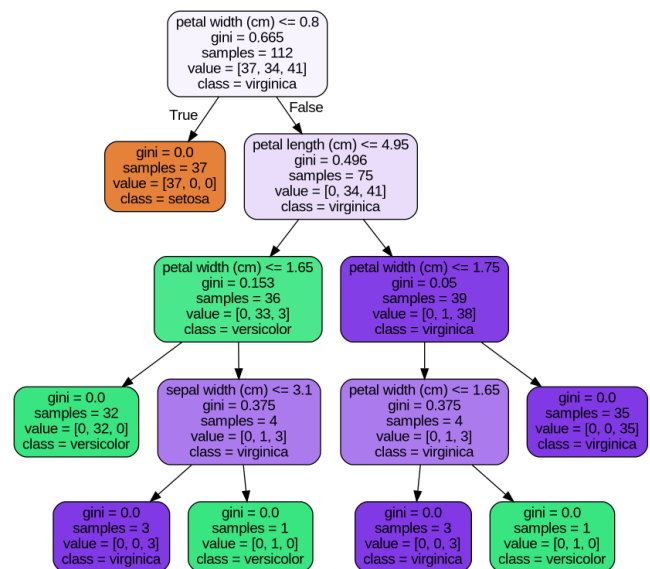


Figure 1: Modelo de árbol de decisión generado para el dataset Iris

### 4.2 Dataset Wine

TODO

### 4.3 Dataset Breast Cancer

TODO

## 5 CONCLUSIONES Y REFLEXIONES

### TODO

#### 5.1 Refrexi3n de Abraham

### TODO

#### 5.2 Reflexi3n de Mario

Por mi parte, considero que esta pr3ctica me ayud3 a comprender por qu3 se utilizan los 3rboles de decisi3n a pesar de que toman decisiones "codiciosas" y con son muy inestables cuando se seleccionan los datos utilizados para entrenar el 3rbol. Las caracteristicas que yo pienso hacen del 3rbol de decisi3n un clasificador f3cil de aprender y ense1ar son que se puede representar gr3ficamente como lo hicimos en la secci3n resultados, adem3s de que el proceso de clasificaci3n de una instancia o elemento es una serie de preguntas simples.

## REFERENCES

- [1] Carl Kingsford and Steven L Salzberg. 2008. What are decision trees? *Nature biotechnology* 26, 9 (2008), 1011–1013.
- [2] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.

## A CÓDIGO PARA LA GENERACIÓN DEL ÁRBOL DE DECISIÓN

```
1 import os
2
3 from graphviz import Source
4 import numpy as np
5 import matplotlib as mpl
6 import matplotlib.pyplot as plt
7 from sklearn.datasets import load_iris
8 from sklearn.model_selection import train_test_split
9 from sklearn.tree import DecisionTreeClassifier, export_graphviz
10
11 mpl.rc('axes', labelsizes=14)
12 mpl.rc('xtick', labelsizes=12)
13 mpl.rc('ytick', labelsizes=12)
14
15 # Setup folders
16 PROJECT_ROOT_DIR = "."
17 CHAPTER_ID = "decision_trees"
18 IMAGES_PATH = os.path.join(PROJECT_ROOT_DIR, "images", CHAPTER_ID)
19 os.makedirs(IMAGES_PATH, exist_ok=True)
20
21
22 def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=300):
23     path = os.path.join(IMAGES_PATH, fig_id + "." + fig_extension)
24     print("Saving figure", fig_id)
25     if tight_layout:
26         plt.tight_layout()
27     plt.savefig(path, format=fig_extension, dpi=resolution)
28
29
30 # Load and split Iris dataset
31 iris = load_iris()
32 xTrain, xTest, yTrain, yTest = train_test_split(
33     iris.data, iris.target, random_state=0)
34
35 # Create and train the decision tree model
36 tree_clf = DecisionTreeClassifier(random_state=0)
37 tree_clf.fit(xTrain, yTrain)
38
39 # Export an image of the tree
40 dot_src = export_graphviz(tree_clf, feature_names=iris.feature_names,
41                             class_names=iris.target_names, rounded=True,
42                             filled=True)
43 image_filename = os.path.join(IMAGES_PATH, "iris_tree")
44 Source(dot_src).render(image_filename, format="png", cleanup=True)
45
46 # Test the decision tree
47 accuracy = tree_clf.score(xTest, yTest)
48 print("Accuracy:", accuracy)
```