

# PRONOSTICAR EL RESULTADO DEL ICFES UTILIZANDO ARBOLES DE DECISIÓN

Andrés Múnera  
Universidad EAFIT  
Colombia  
ajmunerau@eafit.edu.co

Jorge Daniel Ruiz  
Universidad EAFIT  
Colombia  
jdruizl@eafit.edu.co

Mauricio Toro  
Universidad EAFIT  
Colombia  
mtorobe@eafit.edu.co

## RESUMEN

Determinar los factores que determinan los resultados de los estudiantes colombianos en los exámenes institucionales podría ser la clave para la mejoría de la educación en nuestro país. Debido a lo anterior, el objetivo del presente trabajo es presentar un algoritmo que permita ubicar, clasificar y prever los resultados de este grupo de personas en las pruebas, de modo que a partir de los resultados que arroje dicho ejercicio puedan tomarse medidas al respecto.

## PALABRAS CLAVES

Estructura de datos; árboles de decisión; complejidad; tiempo de ejecución; ArrayList; consumo de memoria.

## PALABRAS CLAVE DE LA CLASIFICACIÓN DE LA ACM

Theory of computation → Computational complexity and cryptography → Problems, reduction and completeness  
Theory of computation → Design and analysis of algorithms → Graph algorithms analysis → Shortest path

## 1. INTRODUCCIÓN

El estudiante promedio en Colombia no sabe leer ni entiende de matemáticas o de ciencias; prueba de ello, los resultados de estos en las pruebas Pisa 2019 [1]. ¿Qué factores determinan estos resultados y cómo trabajar en ellos para mejorar la educación en el país? Conocer las variables que influyen en un buen o un mal resultado del estudiante podría dar lugar a una mejor educación y, por lo tanto, a mejores resultados en pruebas internacionales y aumento de posibles becas y oportunidades para el estudiante.

## 2. PROBLEMA

¿Cómo anticipar el éxito o el fracaso del estudiante (entendiendo por éxito el haber alcanzado el promedio) en las PRUEBAS SABER PRO, a partir de factores que influyen en su rendimiento (nivel socioeconómico, condiciones familiares, etc.)?

## 3. TRABAJOS RELACIONADOS

### 3.1 ID3

Es un algoritmo utilizado en el contexto de la inteligencia artificial, es muy rápido y construye un árbol pequeño.

La estructura de los árboles de decisión está formada por:

1. **Nodos:** nombres de los atributos.
2. **Ramas:** posibles valores del atributo asociado al nodo.

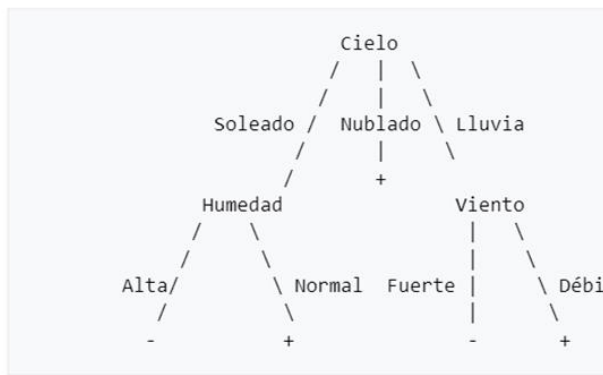
3. **Hojas:** conjuntos clasificados con el nombre de una clase.

La entropía permite calcular el porcentaje de incertidumbre de una muestra y se calcula así:

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

Ej.	Cielo	Temperatura	Humedad	Viento	Jugar tenis
D1	Sol	Alta	Alta	Débil	-
D2	Sol	Alta	Alta	Fuerte	-
D3	Nubes	Alta	Alta	Débil	+
D4	Lluvia	Suave	Alta	Débil	+
D5	Lluvia	Baja	Normal	Débil	+
D6	Lluvia	Baja	Normal	Fuerte	-
D7	Nubes	Baja	Normal	Fuerte	+
D8	Sol	Suave	Alta	Débil	-
D9	Sol	Baja	Normal	Débil	+
D10	Lluvia	Suave	Normal	Débil	+
D11	Sol	Suave	Normal	Fuerte	+
D12	Nubes	Suave	Alta	Fuerte	+
D13	Nubes	Alta	Normal	Débil	+
D14	Lluvia	Suave	Alta	Fuerte	-

En ese caso el árbol finalmente obtenido sería así:



### 3.2 C4.5

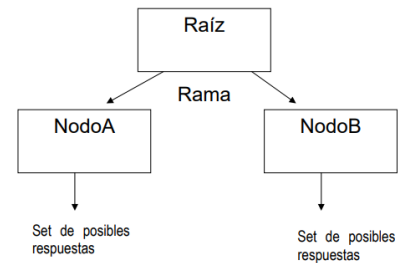
Es uno de los algoritmos más populares de *clasificación basada en reglas*. Es utilizado para resolver *minería de datos*, es decir, puede encontrar patrones repetitivos, tendencias y predecir resultados. El C4.5 se considera un avance del algoritmo ID3 ya que fue desarrollado por el mismo creador años después.

Este genera un árbol mediante particiones de los datos realizadas recursivamente. El algoritmo C4.5 considera todas las posibles pruebas que se pueden dar, divide el conjunto de datos y selecciona la prueba que resulta en la mayor ganancia de información.

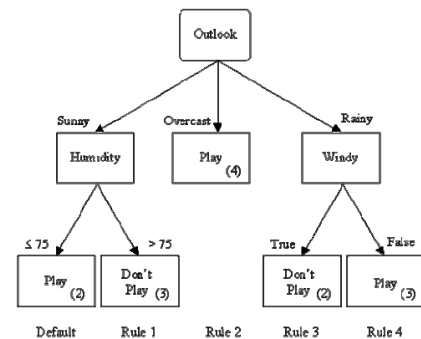
- Los nodos de las raíces son los nodos principales de los nodos del árbol. Se consideran todas las muestras y se seleccionan los atributos más importantes.
- La información de muestra es pasada a los nodos subsiguientes, llamados “nodos de la rama”, que

eventualmente terminan en los “nodos de las hojas”, que dan las decisiones.

- Las reglas son generadas mediante la ilustración del camino desde el “nodo de la raíz” hasta el “nodo de la hoja” [3].



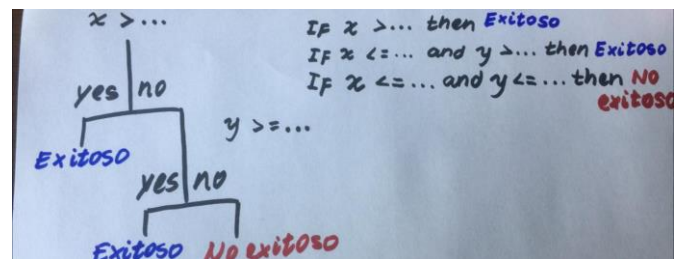
Ejemplo aplicado de Árbol de Decisión adaptado para C4.5



### 3.3 C5

### 3.4 CART

Los árboles de clasificación y regresión o también conocidos como la metodología CART, proceden del ámbito de la estadística.



#### 4. ESTRUCTURA DE DATOS

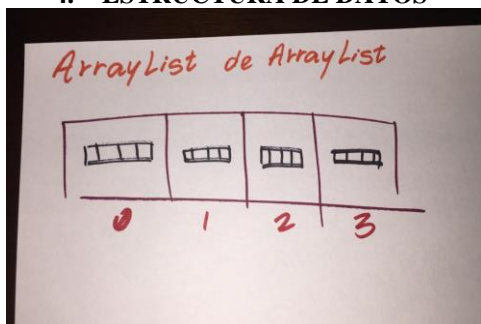
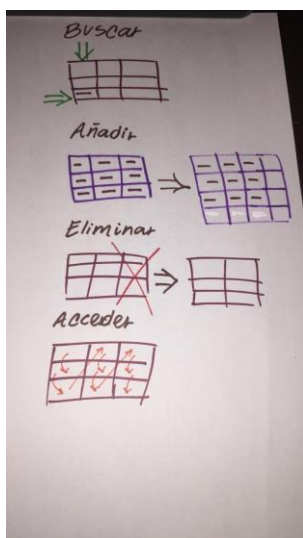


Figura 1. Un ArrayList permite crear una lista de elementos, pero las posiciones son dinámicas, es decir se pueden modificar.

El ArrayList externo es cada persona, y el ArrayList interno son los datos de la persona. En este caso va a ser un ArrayList de strings.

##### 4.1 OPERACIONES DE LAS ESTRUCTURAS DE DATOS



##### 4.2 CRITERIOS DE DISEÑO DE LA ESTRUCTURA DE DATOS

Hemos escogido esta estructura de datos, ya que es el mejor que se adapta al problema de las Pruebas Saber Pro, además, puede darnos una solución más efectiva.

Los ArrayList no son restringidos como los Arrays. En los ArrayList podemos hacer muchas cosas: insertar, eliminar, buscar, entre otros.

##### 4.3 ANÁLISIS DE COMPLEJIDAD

MÉTODO	COMPLEJIDAD
Añadir	$O(n)$
Eliminar	$O(1)$

Buscar	$O(n)$
Acceder	$O(1)$

Tabla 1. Es la complejidad de cada operador.

##### 4.4 TIEMPO DE EJECUCIÓN

##### 4.5 MEMORIA

##### 4.6 ANÁLISIS DE RESULTADOS

#### REFERENCIAS

1. Colprensa. Colombia, con la peor nota de la Oede en pruebas Pisa. *El Colombiano*.
2. J.R.Quinlan, Induction of Decision Trees, *Machine Learning*, 1986, pp81-106.
3. Mazid, M., Ali, S. and Tickle K. Improved C4.5 Algorithm for Rule Based Classification. *Central Queensland University*. 296-301
4. Díaz, J. and Correa, J. Comparación entre árboles de regresión CART y regresión lineal. *Universidad Santo Tomás*. Vol. 6, No. 2, 175-195