

# Análisis CART – CART analysis

Authored by  
**memjavad**

November 12, 2025

## RECOMMENDED CITATION

memjavad (2025). *Análisis CART – CART analysis*. Spanish Psychological Databases.  
Retrieved from <https://spanish.arabpsychology.com/?p=3986>

## Análisis CART (Árboles de Clasificación y Regresión)

**Primary Disciplinary Field(s):** Estadística, Aprendizaje Automático (Machine Learning), Minería de Datos (Data Mining).

**Proponents:** [Leo Breiman](#), Jerome H. Friedman, Richard A. Olshen, Charles J. Stone.

### 1. Principios Fundamentales de CART

El análisis de Árboles de Clasificación y Regresión (**CART**, por sus siglas en inglés) constituye una metodología fundamental dentro de la estadística no paramétrica y el aprendizaje automático, introducida formalmente en 1984 por Breiman, Friedman, Olshen y Stone. Su objetivo primordial es la creación de un modelo predictivo que se representa visualmente como una estructura de árbol binario. Este modelo se construye mediante la partición recursiva del espacio de las variables predictoras en regiones mutuamente excluyentes y exhaustivas, buscando que cada región terminal, conocida como nodo hoja, sea lo más homogénea posible respecto al valor de la variable de respuesta. La gran fortaleza de CART reside en su capacidad para manejar datos complejos sin requerir suposiciones rígidas sobre la distribución subyacente de los datos, a diferencia de los métodos estadísticos paramétricos tradicionales.

La filosofía central de CART se basa en el principio de la **división binaria**. En cada paso del proceso de construcción del árbol, el algoritmo selecciona una única variable predictora y un punto de corte (umbral) que divide el conjunto de datos actual en exactamente dos subconjuntos o nodos hijos. Esta división se elige de manera que maximice la ganancia de información o, equivalentemente, minimice la impureza o la varianza dentro de los nodos resultantes. Este enfoque codicioso (greedy) y descendente continúa hasta que se alcanza un criterio de parada predefinido, lo que resulta en un árbol potencialmente muy grande y complejo. El resultado final es una serie de reglas condicionales simples que, en conjunto, definen la predicción para cualquier nueva observación de datos, facilitando enormemente la interpretación del modelo.

A diferencia de otros algoritmos de agrupamiento o clasificación lineal, CART se distingue por su enfoque en la interpretabilidad y la facilidad de uso. El proceso de segmentación busca aislar grupos de observaciones que comparten características similares y que exhiben un comportamiento uniforme en la variable objetivo. Por ejemplo, en un contexto de clasificación, si el objetivo es predecir la probabilidad de que un cliente abandone un servicio (churn), CART creará nodos hoja donde la mayoría de los miembros han abandonado o la mayoría no lo han hecho. Esta transparencia, a menudo referida como un modelo de "caja blanca", contrasta marcadamente con la opacidad inherente a modelos más complejos como las redes neuronales profundas, haciendo de CART una herramienta preferida en entornos donde la explicación de la decisión es tan crucial como la precisión de la predicción.

## 2. Distinción entre Árboles de Clasificación y Regresión

Aunque el algoritmo de construcción fundamental es idéntico, la aplicación de CART se bifurca en dos categorías principales dependiendo de la naturaleza de la variable de respuesta: Árboles de Clasificación y Árboles de Regresión. La distinción crucial radica en la métrica utilizada para medir la impureza de un nodo y, por ende, la eficacia de una división. Cuando la variable de respuesta es categórica (por ejemplo, sí/no, A/B/C), se emplea un Árbol de Clasificación. En este caso, el algoritmo busca la división que resulte en la mayor pureza de clase dentro de los nodos hijos. Las métricas más comunes para evaluar esta pureza son el índice de **Gini** y la ganancia de información basada en la entropía.

El **índice de Gini** es la medida de impureza preferida por el algoritmo CART original. Este índice cuantifica la probabilidad de clasificar incorrectamente una observación seleccionada al azar dentro del nodo si su clasificación se basa en la distribución de las clases en ese nodo. Una impureza Gini de cero indica una pureza perfecta (es decir, todas las observaciones en el nodo pertenecen a la misma clase). El algoritmo selecciona el punto de corte que produce la mayor reducción en el índice de Gini ponderado entre los nodos hijos, asegurando que la nueva partición sea lo más discriminatoria posible. Otros enfoques, como el uso de la entropía, también buscan maximizar la homogeneidad de la clase, pero el índice de Gini ofrece ventajas computacionales en la implementación estándar de CART.

Por otro lado, cuando la variable de respuesta es continua o numérica (por ejemplo, precio, edad, ingreso), se utiliza un Árbol de Regresión. En este contexto, la noción de impureza se reemplaza por la varianza o, más comúnmente, el error cuadrático medio (**MSE**, Mean Squared Error). El objetivo del Árbol de Regresión no es clasificar, sino predecir un valor numérico. Por lo tanto, el algoritmo busca divisiones que minimicen la varianza de la variable de respuesta dentro de los nodos resultantes. La predicción final para cualquier nueva observación que caiga en un nodo hoja particular es típicamente el promedio de los valores de la variable de respuesta de todas las observaciones de entrenamiento que residen en ese nodo. Esta adaptación permite que la estructura de árbol sea utilizada eficazmente tanto para problemas de predicción continua como discreta.

## 3. El Algoritmo de Construcción del Árbol

La construcción de un árbol CART sigue un proceso iterativo y voraz. Se inicia con el nodo raíz, que contiene todas las observaciones del conjunto de datos de entrenamiento. En cada nodo, el algoritmo evalúa todas las posibles variables predictoras y todos los posibles puntos de corte para cada variable. El objetivo es identificar la combinación óptima de variable y umbral que, si se aplica, produce la mejor división binaria según la métrica de impureza o varianza seleccionada (Gini para clasificación, MSE para regresión). Este proceso se conoce como **partición recursiva**

## binaria.

El carácter "voraz" del algoritmo significa que en cada paso, la decisión de división se toma localmente para maximizar la mejora inmediata, sin considerar si esta decisión local conducirá al árbol globalmente óptimo a largo plazo. Una vez que se elige la mejor división, el nodo padre se divide en dos nodos hijos, y el proceso se repite de forma independiente en cada uno de estos nuevos nodos. Esta repetición continúa hasta que se cumplen ciertos criterios de parada, como alcanzar una profundidad máxima predefinida, tener un número mínimo de observaciones en un nodo para permitir una división, o cuando ninguna división posible mejora significativamente la pureza del nodo. Es crucial entender que, en la fase inicial de crecimiento, CART se permite crecer hasta ser un árbol "máximo" o "sobreajustado" (overfitted).

Un aspecto fundamental de la implementación de CART es su manejo de diferentes tipos de datos. Puede gestionar variables continuas, ordinales y nominales. Para variables continuas, el algoritmo busca puntos de corte discretos. Para variables categóricas, si son binarias, el corte es trivial. Si son nominales con múltiples categorías, CART evalúa todas las posibles agrupaciones binarias de esas categorías. Esta flexibilidad es una ventaja significativa, ya que elimina la necesidad de preprocesamiento complejo de variables, como la normalización o la creación de variables ficticias (dummy variables), que son esenciales para modelos como la regresión lineal o las máquinas de vectores de soporte (SVM).

## 4. Poda (Pruning) y Validación Cruzada

Uno de los principales desafíos en la construcción de árboles de decisión es el riesgo de **sobreajuste** (overfitting). Si se permite que el árbol crezca sin restricciones hasta que todos los nodos hoja sean perfectamente puros, el modelo resultante memorizará el ruido y las peculiaridades específicas del conjunto de entrenamiento, lo que resultará en un rendimiento pobre cuando se aplique a datos nuevos e invisibles. Para combatir este problema, el algoritmo CART emplea una técnica sofisticada conocida como poda de costo-complejidad (Cost-Complexity Pruning) o poda por el error más débil.

La poda de costo-complejidad opera en dos fases. Primero, se construye el árbol máximo (sobreajustado). Segundo, se utiliza un parámetro de complejidad, denotado como  $\alpha$  (alfa), que equilibra la complejidad del árbol (medida por el número de nodos terminales) con su error de entrenamiento. Para cada valor de  $\alpha$ , se encuentra el subárbol más pequeño que minimiza la suma del error de clasificación (o regresión) más una penalización proporcional a la complejidad del árbol. Al variar  $\alpha$  desde cero hasta el infinito, se genera una secuencia anidada de subárboles que van desde el árbol máximo hasta el nodo raíz único.

La selección del subárbol óptimo dentro de esta secuencia se realiza mediante la **validación cruzada** (cross-validation). Típicamente, se utiliza la validación cruzada de diez pliegues (10-fold

cross-validation). El conjunto de entrenamiento se divide en diez subconjuntos; el árbol se entrena en nueve y se prueba en el décimo, y este proceso se repite diez veces. La tasa de error promedio de estos diez ensayos se utiliza para estimar el rendimiento de generalización para cada subárbol en la secuencia de poda. El árbol final seleccionado es aquel que minimiza la tasa de error de validación cruzada. Esta metodología garantiza que el árbol final mantenga un equilibrio adecuado entre la complejidad y la capacidad de generalización, mitigando el riesgo de sobreajuste.

## 5. Ventajas Metodológicas y Aplicaciones

El análisis CART ofrece múltiples ventajas que explican su perdurable popularidad en diversos campos. Una de las más destacadas es la **interpretabilidad**. La estructura del árbol proporciona un mapa visual y lógico de las decisiones, permitiendo a los analistas rastrear exactamente cómo se llegó a una predicción. Esto es especialmente valioso en áreas reguladas como las finanzas (calificación crediticia) y la medicina (diagnóstico), donde la justificación de la decisión es un requisito legal o ético. Además, la sencillez de las reglas de decisión resultantes facilita su comunicación a audiencias no técnicas.

Otra ventaja significativa es la robustez de CART frente a la preparación de datos. Los árboles de decisión son inherentemente insensibles a las transformaciones monótonas de las variables predictoras (por ejemplo, logaritmos o raíces cuadradas), y pueden manejar datos faltantes de manera efectiva mediante técnicas como los **sustitutos** (surrogates) que utilizan otras variables para guiar la división cuando la variable principal está ausente. Además, CART es robusto a la presencia de valores atípicos (outliers) en las variables predictoras, ya que el proceso de división solo se preocupa por el ordenamiento de los datos y no por la magnitud exacta de los valores extremos.

Las aplicaciones de CART son extremadamente amplias. En el sector financiero, se utiliza para la modelización del riesgo de crédito y la detección de fraudes. En la investigación médica, ayuda a identificar grupos de pacientes con pronósticos similares o a determinar combinaciones de síntomas que predicen la aparición de una enfermedad. En marketing, se emplea para la segmentación de clientes y la predicción de la respuesta a campañas publicitarias. En la ingeniería y la fabricación, CART puede utilizarse para el diagnóstico de fallas y el control de calidad, identificando combinaciones de parámetros de entrada que conducen a productos defectuosos.

## 6. Limitaciones y Extensiones (Bagging, Random Forests y Boosting)

A pesar de sus muchas virtudes, el algoritmo CART presenta ciertas limitaciones inherentes. La principal es su **inestabilidad** o alta varianza. Un pequeño cambio en el conjunto de datos de entrenamiento, especialmente en las observaciones cerca de los puntos de corte óptimos, puede

llevar a una estructura de árbol completamente diferente en los niveles superiores, lo que resulta en predicciones inconsistentes. Además, aunque CART intenta encontrar la mejor división localmente, el enfoque voraz no garantiza que el árbol resultante sea el mejor modelo globalmente posible para el conjunto de datos. Los árboles de decisión únicos también tienden a ser menos precisos que los modelos modernos basados en ensambles.

Para superar la alta varianza y mejorar significativamente la precisión predictiva, las limitaciones de los árboles CART individuales llevaron al desarrollo de potentes **métodos de ensamble** que utilizan árboles de decisión como componentes básicos. El primero de ellos fue el **Bagging** (Bootstrap Aggregating), que entrena múltiples árboles CART en diferentes muestras bootstrap del conjunto de datos y promedia sus predicciones para reducir la varianza. Una extensión aún más popular es **Random Forests** (Bosques Aleatorios), propuesto también por Breiman. Este método introduce una aleatoriedad adicional al seleccionar solo un subconjunto aleatorio de variables predictoras en cada división de nodo, lo que asegura que los árboles en el ensamble sean descorrelacionados.

Otro avance crucial es el **Boosting**, específicamente el Gradient Boosting Machines (GBM) y su implementación moderna como [XGBoost](#). A diferencia de Bagging y Random Forests, que construyen árboles de forma independiente, el Boosting construye árboles secuencialmente. Cada nuevo árbol intenta corregir los errores residuales cometidos por el conjunto de árboles anteriores, centrándose en las observaciones que fueron mal clasificadas o mal predichas. Estos métodos de ensamble, al combinar la interpretabilidad de los árboles individuales con la potencia predictiva de la agregación, se han convertido en la columna vertebral de muchas soluciones de aprendizaje automático de vanguardia, demostrando que el legado de CART se extiende mucho más allá del árbol único original.

## 7. Lecturas Adicionales

[Decision tree learning \(Wikipedia\)](#)

[Leo Breiman \(Pionero de CART y Random Forests\)](#)

[Classification and Regression Trees \(Libro original de Breiman et al., 1984\)](#)

[Random Forests \(Sitio oficial de Breiman\)](#)

[Documentación de XGBoost \(Ejemplo moderno de Boosting basado en árboles\)](#)