

# EST-25134: Aprendizaje Estadístico

**Profesor:** Alfredo Garbuno Iñigo — Primavera, 2023 — Extensiones Arboles.

**Objetivo:** Que veremos.

**Lectura recomendada:** Capítulo 3 de [1].

## 1. INTRODUCCIÓN

Los árboles de decisión por construcción son insensibles a:

- Variables continuas con *outliers*.
- Variables continuas con diferentes escalas.
- Variables categóricas (no se necesitan transformar a *dummies*).

El algoritmo de ajuste de un árbol de decisión, en regresión por ejemplo, busca minimizar la función de pérdida

$$R_{\alpha}(T) = \sum_{m=1}^{|T|} \sum_{i: x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|, \quad (1)$$

donde  $\hat{y}_{R_m} = \sum_{i: x_i \in R_m} y_i / n_m$  es el promedio de las respuestas en la región  $m$ -ésima, a través de ir buscando secuencialmente las regiones  $R_m$  que logren la máxima disminución de RSS.

La búsqueda se puede lograr por medio de un algoritmo voraz (*greedy*) que escoge variables y puntos de corte.

*1.0.1. Para pensar:* En nuestra última discusión sobre árboles de decisión es que tienen un sesgo al utilizar variables categóricas con muchas etiquetas. ¿Por qué?

### 1.1. Motivación

Consideremos el experimento siguiente. Tenemos una colección de 7 atributos y una respuesta. Los datos son generados de manera aleatoria.

```
1 set.seed(108727)
2 generate_data(nsamples = 100) >
3 print(n = 5)
```

De manera artificial consideraremos que la respuesta  $y$  está relacionada con los demás atributos (sabemos que no es cierto).

```
1 fit_tree <- function(engine){
2   data_train <- generate_data()
3
4   tree_spec <- decision_tree(tree_depth = 2) >
5   set_engine(engine) >
6   set_mode("regression")
7
8   tree_spec >
9   fit(y ~ ., data = data_train) >
10  extract_fit_engine() >
11  vi() >
12  filter(Importance > 0) >
13  mutate(rank = min_rank(desc(Importance)))
14 }
```

Repetimos este proceso (generar datos aleatorios y ajustar un árbol) un número determinado de veces y registramos cuántas veces cada atributo fue utilizado en el nodo raíz.

En contraste, con un modelo denominado **conditional trees** somos capaces de evitar ese sesgo. Incluso podemos evitar escoger variables de corte que no tienen asociación con la respuesta.

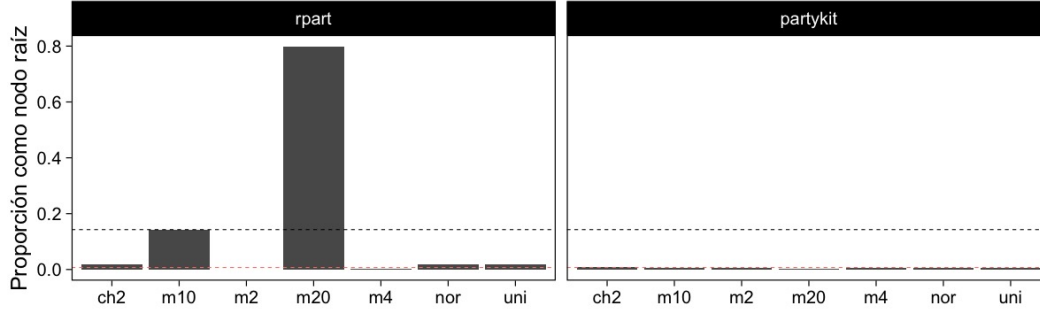


FIGURA 1. La línea negra representa la probabilidad de seleccionar una variable al azar como nodo raíz ( $p = 1/7$ ) y la crua salmón representa la tolerancia de error (tasa de falsos positivos) ajustada ( $\alpha = 0.05/7$ ).

## 2. INFERENCIA CONDICIONAL

Recordemos que lo que necesitamos para construir un árbol de decisión es iterar:

1. Seleccionar una variable;
2. Seleccionar un punto de corte.

La variable se puede escoger primero por medio de comparaciones estadísticas entre un atributo  $X_j$  y la respuesta  $Y$ :

1. Hacer una prueba de  $\chi^2$  si ambas son variables nominales.
2. Hacer una prueba de correlación si ambas son continuas.

El punto de corte se puede escoger con una métrica de impureza.

Consideremos una colección de datos  $(x_i, y_i) \stackrel{\text{iid}}{\sim} \mathbb{P}_{\mathcal{X}, \mathcal{Y}}$  con  $n = 1, \dots, n$ . Lo que queremos contrastar es

$$H_0 : \quad \mathbb{P}(Y|X) = \mathbb{P}(Y). \quad (2)$$

La prueba adecuada se puede escoger de acuerdo a las características de  $X$  y  $Y$ .

En general podemos utilizar un mecanismo general por medio de un vector de estadísticas

$$T = \text{vec} \left( \sum_{i=1}^n g(x_i) h(y_i)^\top \right) \in \mathbb{R}^{pq}, \quad (3)$$

donde el operador **vec** transforma matrices en vectores (por arreglo),  $g : \mathcal{X} \rightarrow \mathbb{R}^p$  es una **transformación de atributos** y  $h : \mathcal{Y} \rightarrow \mathbb{R}^q$  es un **función de influencia**.

Cómo escogemos  $g$  y  $h$  es lo que nos permite realizar las pruebas de hipótesis adecuadas:

- Pruebas de correlación;
- Pruebas de muestras pareadas;
- Pruebas de  $K$ -muestras similares a pruebas ANOVA;
- Pruebas de independencia en tablas de contingencia.

Independientemente de la prueba, para poder utilizar este mecanismo debemos de saber la **distribución de muestreo** de  $T$  bajo la hipótesis nula.

La ventaja que tenemos es que la estructura del problema de contraste de hipótesis nos permite dejar fijos los atributos y realizar permutaciones  $\sigma \in S$  sobre la respuesta. Es decir,

$$((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)) \mapsto ((x_1, y_{\sigma(1)}), (x_2, y_{\sigma(2)}), \dots, (x_n, y_{\sigma(n)})) , \quad (4)$$

Denotemos por  $\mu_h$  el valor esperado de la función de influencia condicional en la permutación  $\sigma$  y por  $\Sigma_h$  la matriz de varianzas-covarianzas asociada a ese valor esperado.

Esto nos permite calcular el valor esperado y matriz de varianzas covarianzas del estadístico  $T$  condicional en la permutación  $\sigma \in S$ , los cuales denotamos  $\mu$  y  $\Sigma$  respectivamente.

Finalmente con esto podemos calcular un estadístico de prueba para  $H_0$ . Esto se puede lograr a través de una forma cuadrática o un estadístico de orden

$$c_{\text{quad}} = (T - \mu)^\top \Sigma^{-1} (T - \mu) , \quad (5)$$

$$c_{\text{max}} = \text{máx} \left| \frac{T - \mu}{\text{diag}(\Sigma)^{1/2}} \right| . \quad (6)$$

Lo ideal es poder tener conocimiento de la distribución de muestreo de  $c_{\text{quad}}$  y  $c_{\text{max}}$ . Esto con la intención de poder construir valores  $p$  para contrastar  $H_0$ . Por ejemplo, podemos calcular

$$\mathbb{P}(c(T, \mu, \Sigma) \leq z | S) , \quad (7)$$

la cual se puede calcular como el número de permutaciones que tienen un estadístico menor que el nivel  $z$  dividido por el número total de permutaciones.

Utilizar remuestreo (*bootstrap*) nos puede ayudar a calcular aproximaciones hasta un nivel de tolerancia dado (mas detalles de esto en mi curso de simulación (EST-24107: Simulación)).

O podemos argumentar por algún resultado asintótico para determinar distribuciones de muestreo que permitan un cálculo mas eficiente. Por ejemplo, en el caso  $pq = 1$  podemos utilizar

$$c_{\text{quad}} \sim \chi_1^2 , \quad c_{\text{max}} \sim N(0, 1) . \quad (8)$$

## 2.1. Ejemplos de pruebas de hipótesis de independencia

Bajo el caso de dos variables continuas podemos utilizar el mecanismo

### 3. EJEMPLO DE PARTYKIT

### 4. CONCLUSIONES

El modelo de **CTree** es un modelo:

1. que utiliza pruebas de hipótesis para determinar variables y puntos de corte;
2. tiene un mecanismo de selección insesgado;
3. no requiere mucho post-procesamiento (poda);

De acuerdo a Greenwell [1] , aún cuando **CTree** tiene mejores propiedades estadística que **CART** hay un uso generalizado por el último debido a herramientas de código abierto.

Además, de acuerdo a Loh [2] mientras un árbol se escoja por cuestiones predictivas y no por inferencia tiene un riesgo bajo al utilizar procedimientos con sesgo. Por otro lado, validación cruzada puede ayudar a eliminar ramas redundantes durante el proceso de poda (mientras tengamos pocos atributos y un número suficiente de datos).

---

**REFERENCIAS**

- [1] B. M. Greenwell. *Tree-Based Methods for Statistical Learning in R*. Chapman and Hall/CRC, Boca Raton, first edition, may 2022. ISBN 978-1-00-308903-2. . [1](#), [3](#)
- [2] W.-Y. Loh. Fifty Years of Classification and Regression Trees. *International Statistical Review / Revue Internationale de Statistique*, 82(3):329–348, 2014. ISSN 0306-7734. [3](#)