

EST-25134: Aprendizaje Estadístico

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

Objetivo: En esta sesión estudiaremos una extensión de los árboles de decisión como mecanismos iterativos de partición del espacio de atributos. Estudiaremos un método reciente, que resulta más robusto a selecciones de hiper-parámetros. Además conectaremos esto con una técnica estadística de mucha utilidad en la era del *Big Data*, pruebas de hipótesis múltiples..

Lectura recomendada: Capítulo 3 de [2]. Pruebas de hipótesis múltiples pueden ser consultadas en el libro de Efron and Hastie [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 = 3)
```

```
1 # A tibble: 100 × 8
2   y      ch2 m2      m4      m10      m20      nor      uni
3   <dbl> <dbl> <fct> <fct> <fct> <fct> <dbl> <dbl>
4 1 2.08   1.86 1      2      10      16      0.107 0.00758
5 2 0.804  2.40 2      3      9       4      -0.309 0.804
6 3 0.313  1.29 2      2      4      20      1.18  0.658
7 # ... with 97 more rows
8 # Use 'print(n = ...)' to see more rows
```

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.

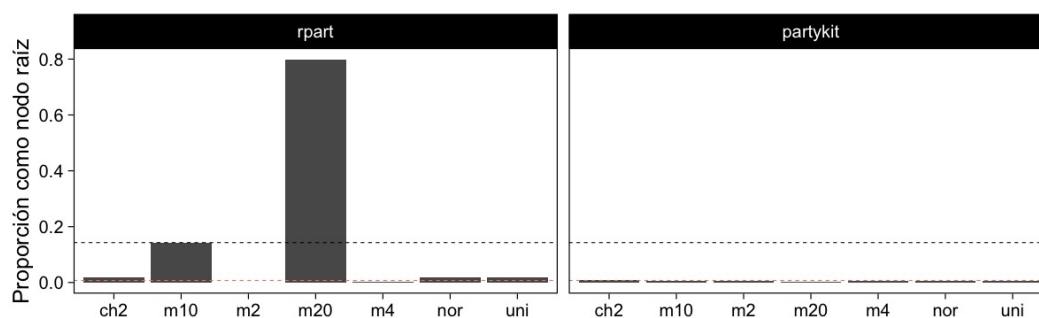


FIGURA 1. La línea negra representa la probabilidad de seleccionar una variable al azar como nodo raíz ($p = 1/7$) y la curva salmón representa la tolerancia de error (tasa de falsos positivos) ajustada ($\alpha = 0.05/7$).

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 χ^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 μ_h el valor esperado de la función de influencia condicional en la permutación σ y por Σ_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 μ y Σ 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

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

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

Lo ideal es poder tener conocimiento de la distribución de muestreo de c_{quad} y c_{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. Y evaluar si se satisface un umbral, definido por el analista, para controlar el nivel de significancia estadística con

$$\alpha = \mathbb{P}\{\text{rechazar } H_0 \text{ cuando es verdadera}\}. \quad (8)$$

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

$$c_{\text{quad}} \sim \chi_1^2, \quad c_{\text{max}} \sim N(0, 1). \quad (9)$$

2.1. Ejemplos de pruebas de hipótesis de independencia

Bajo el caso de dos variables continuas podemos utilizar el mecanismo de inferencia condicional para probar independencia estadística. Por ejemplo, regresemos a nuestro ejemplo de jugadores de *baseball*.

```
1 library(coin)
2 independence_test(Salary ~ Years, data = hitters, teststat = "quadratic")
```

```
1
2      Asymptotic General Independence Test
3
4 data:  Salary by Years
5 chi-squared = 42, df = 1, p-value = 9e-11
```

Incluso también podríamos aplicar lo mismo para nuestro conjunto de datos ficticio.

```
1 fake_data <- generate_data()
2 independence_test(y ~ nor, data = fake_data, teststat = "quadratic")
```

```
1
2      Asymptotic General Independence Test
3
4 data:  y by nor
5 chi-squared = 0.3, df = 1, p-value = 0.6
```

Por supuesto, también podríamos hacer una comparación para atributos categóricos. Por ejemplo, regresando a nuestro ejemplo de *Scooby Doo*. El que el monstruo del capítulo sea real o no, ¿tiene relación con que *Scooby* haya sido capturado en el episodio?

```
1 independence_test(monster_real ~ caught_scooby,
2                   data = scooby_data, teststat = "quadratic")
```

```
1
2      Asymptotic General Independence Test
3
4 data:  monster_real by caught_scooby (caught, not caught)
5 chi-squared = 13, df = 1, p-value = 3e-04
```

Adicionalmente, podríamos contrastar con una aproximación por remuestreo a la distribución de permutaciones:

```
1 independence_test(monster_real ~ caught_scooby,
2                   data = scooby_data, teststat = "quadratic",
3                   distribution = approximate(nresample = 10000))
```

```
1
2      Approximative General Independence Test
3
4 data:  monster_real by caught_scooby (caught, not caught)
5 chi-squared = 13, p-value = 1e-04
```

3. ÁRBOLES CONDICIONALES

El objetivo es: construir árboles de decisión con una estrategia de particionado recursivo sin sesgo basada inferencia condicional.

En cada paso de selección de variables lo que hacemos es realizar una prueba de hipótesis por cada atributo. Es decir, buscamos

$$H_0^j : \quad \mathbb{P}(Y|X_j) = \mathbb{P}(Y), \quad j = 1, \dots, m. \quad (10)$$

Dado que cada estadístico asociado $c(T^j, \mu^j, \Sigma^j)$ puede tener un escala muy distinta necesitamos estandarizar todos a la misma. Es por esto, que contrastar los valores p asociados a cada prueba nos ayudan a tener control sobre la escala de las comparaciones.

Adicionalmente, por diseño de las pruebas de hipótesis podríamos aceptar algunos falsos positivos. Esto puede ocurrir con una tasa de α (el nivel de tolerancia de error de las pruebas). Así que lo que buscamos es buscar que la colección de pruebas que realicemos tenga una tasa de error controlada

$$\text{FWER} \leq \alpha. \quad (11)$$

Esto lo podemos lograr, por ejemplo, con la corrección de Bonferroni, al considerar cada umbral como α/m .

El argumento es

$$\text{FWER} = \mathbb{P}\{\text{reject any true } H_0^j\} \quad (12)$$

$$= \mathbb{P}\left\{\bigcup_{j \in I_0} \text{reject } H_0^j \text{ when it's true}\right\} \quad (13)$$

$$= \mathbb{P}\left\{\bigcup_{j \in I_0} \left(p_j \leq \frac{\alpha}{m}\right)\right\} \quad (14)$$

$$\leq \sum_{j=1}^m \mathbb{P}\left(p_j \leq \frac{\alpha}{m}\right) = \alpha, \quad (15)$$

donde I_0 denota el conjunto de hipótesis nulas verdaderas.

Decimos que la variable aleatoria X es una variable aleatoria **super-uniforme** cuando domina estocásticamente una variable aleatoria uniforme. Es decir, cuando

$$\mathbb{P}\{X \leq u\} \leq u \quad \text{para toda } u \in [0, 1]. \quad (16)$$

Con la tasa de error controlada, y una vez seleccionada la variable para realizar el corte entonces procedemos a proponer puntos de corte para dicho atributo. El cual nos lleva a contrastar hipótesis para las distribuciones condicionales

$$\{y_i | x_{ij} < s\} \quad \text{y} \quad \{y_i | x_{ij} \geq s\}. \quad (17)$$

4. CONCLUSIONES

El modelo de **CTree** (disponible a través de **partykit**) 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).
4. el nivel α es un hiper-parámetro para estos modelos.

De acuerdo a Greenwell [2] , 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 [3] 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. Efron and T. Hastie. *Computer Age Statistical Inference*. Cambridge University Press, 2016. [1](#)
- [2] 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](#), [6](#)
- [3] 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. [6](#)