

# E s t a d í s t i c a   l l l   p a r a   l n

J o s e   D a n i e l   R a m i r e z   S o t o   2 0 2 3  
j d r 2 1 6 2 @ c o l u m b i a . e d u

# A g e n d a

a n u n c i o s   v a r i o s

Tarea 2 entrega lunes 8 de Mayo 2023 ( Preguntas  
modelos de analitica ( machine learning - ML ) Super

Á r b o l e s

Á r b o l e s   s i m p l e s

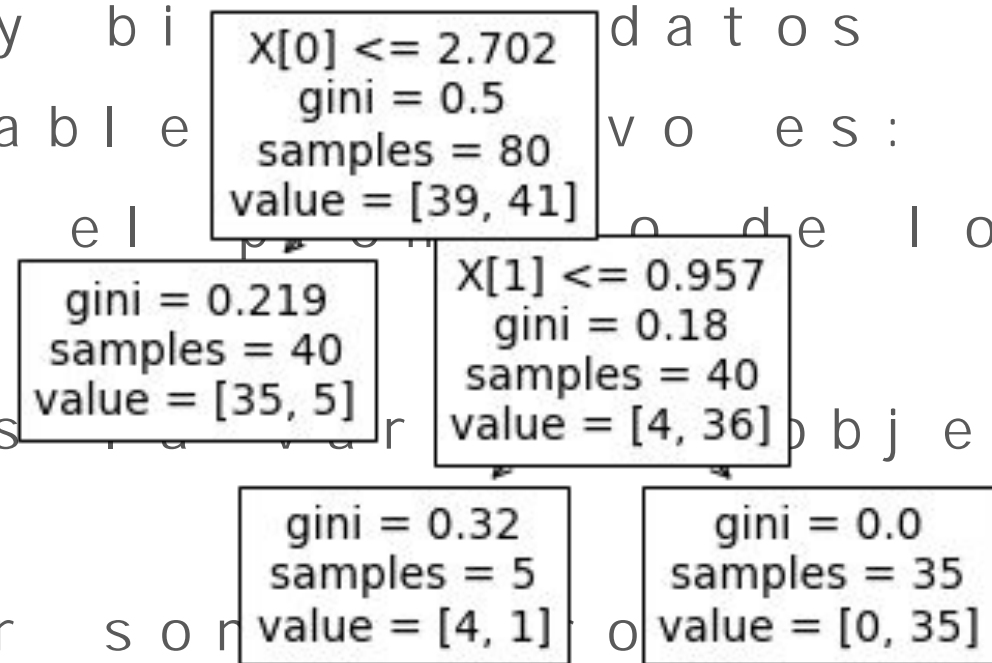
Ma t e m á t i c a   d e   l o s   á r b o l e s

R a n d o m F o r e s t

G B M

# Supervizado, Árboles

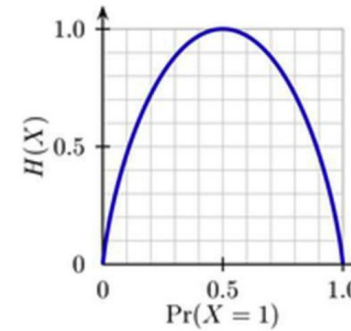
Árboles son modelos que se utilizan para clasificar los datos basados en la capacidad de una variable de dividir los datos. Funcionan muy bien con datos categóricos. Si la variable es continua, la predicción es el promedio de los vecinos. Si la variable es categórica, la predicción es la clase más común.



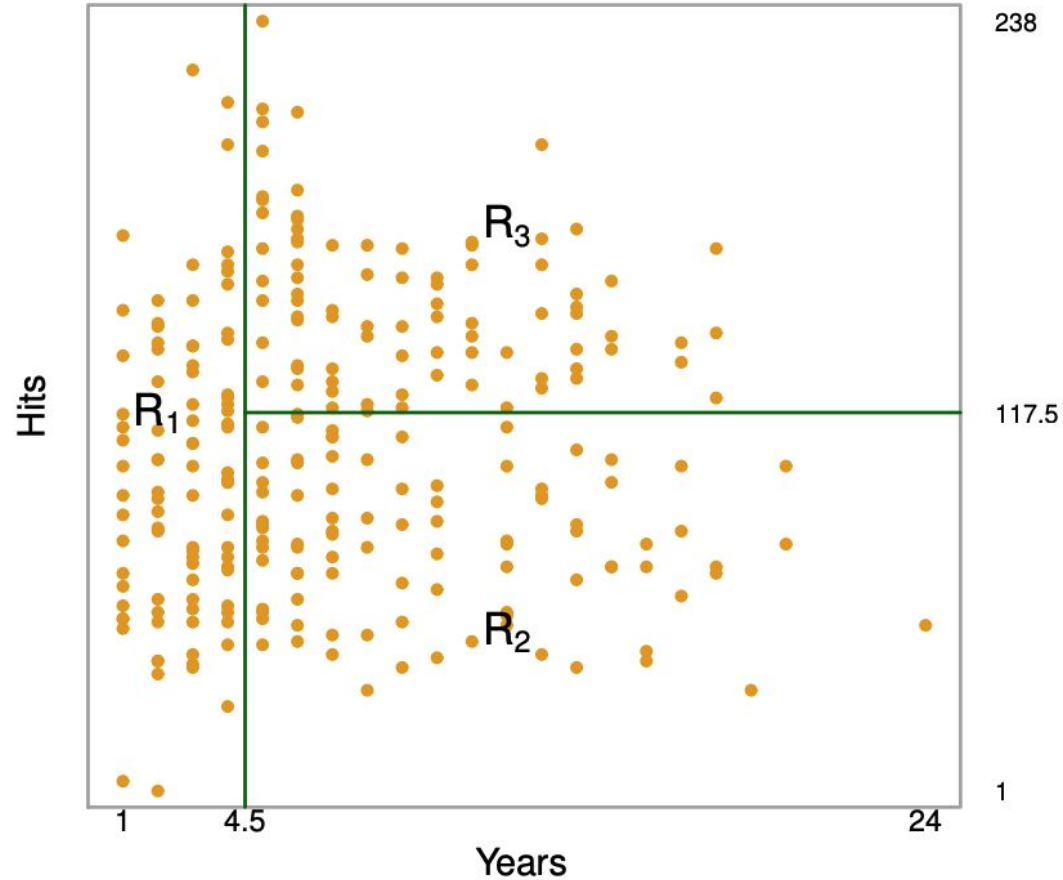
Los parámetros a definir son el tipo de árbol, número mínimo de samples entre otros.

# Entropía, Como funcionan los árboles

## Information Entropy



- For a Bernoulli trial ( $X = \{0, 1\}$ ) the graph of entropy vs.  $\text{Pr}(X=1)$ . The highest  $H(X) = 1 = \log(2)$



$$H_{(S)} = \sum_{i=1}^C -p_i \log_2 p_i$$

$$H_{(T,X)} = \sum_{c \in X} p(c) H_{(c)}$$

$$\text{Gain}_{(T,X)} = H_{(T)} - H_{(T,X)}$$

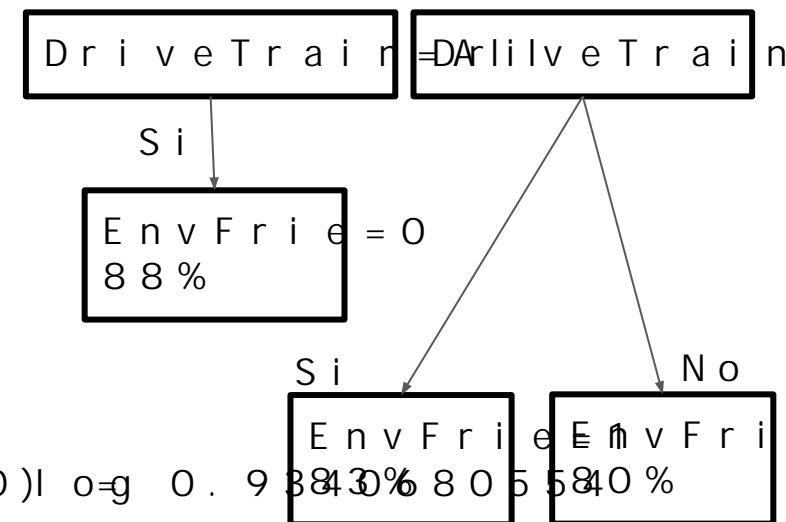
# Calcular el árbol a mano para

Type	DriveTrain	Cylinders	EnvFriendly
Sedan	All	6.0	0
SUV	All	6.0	0
Wagon	Front	4.0	1
Sedan	Rear	8.0	0
Sedan	Front	4.0	1
SUV	Front	4.0	1
SUV	All	6.0	0
Sedan	All	6.0	0
Sedan	All	6.0	0
Sedan	Front	4.0	1
Sports	Rear	6.0	0
Wagon	Rear	6.0	0
Wagon	All	6.0	0
Sedan	Front	6.0	1
Sedan	Rear	8.0	0
SUV	All	6.0	0
Wagon	All	4.0	1
Sports	Rear	4.0	1
SUV	All	8.0	0
Sedan	Front	6.0	0

$$H(S) = \sum_{i=1}^C -p_i \log_2 p_i$$

$$H_{(T,X)} = \sum_{c \in X} p(c) H_{(c)}$$

$$Gain_{(T,X)} = H_{(T)} - H_{(T,X)}$$



$$H(\text{EnvFriendly}) = -\left(\frac{12}{20} \log_2 \frac{12}{20} + \frac{8}{20} \log_2 \frac{8}{20}\right) = 0.9183$$

$$H(\text{EnvFriendly}, \text{DriveTrain}) = p(\text{Train}=\text{All}) H(\text{Train}=\text{All}) + p(\text{Train}=\text{Front}) H(\text{Train}=\text{Front}) + p(\text{Train}=\text{Rear}) H(\text{Train}=\text{Rear}) = 0.9183$$

$$Gain(\text{EnvFriendly}) = H(S) - H(\text{EnvFriendly}) = 1.0 - 0.9183 = 0.0817$$

$$Gain_{\text{DriveTrain}} \text{ es mayor que la ganancia de la variable Type}$$

---

**Algorithm 8.1** *Building a Regression Tree*

---

1. Use recursive binary splitting to grow a large tree on the training data, stopping only when each terminal node has fewer than some minimum number of observations.
  2. Apply cost complexity pruning to the large tree in order to obtain a sequence of best subtrees, as a function of  $\alpha$ .
  3. Use K-fold cross-validation to choose  $\alpha$ . That is, divide the training observations into  $K$  folds. For each  $k = 1, \dots, K$ :
    - (a) Repeat Steps 1 and 2 on all but the  $k$ th fold of the training data.
    - (b) Evaluate the mean squared prediction error on the data in the left-out  $k$ th fold, as a function of  $\alpha$ .Average the results for each value of  $\alpha$ , and pick  $\alpha$  to minimize the average error.
  4. Return the subtree from Step 2 that corresponds to the chosen value of  $\alpha$ .
-

# Bootstrap y Árboles

"The bootstrap is a widely applicable and extremely powerful method for estimating the uncertainty associated with a given estimate of a model parameter."  
Learning

Consiste en crear muchos experimentos tomando pequeñas muestras de los datos para reducir la varianza.

# Bagging y Árboles

“Bootstrap aggregation, or bagging, is a general-purpose learning method; we introduce it here because it is particularly effective.”  
Taken from An Introduction to Statistical Learning

Uso: Entrenar muchos árboles y crear un promedio ponderado.

$$\hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x).$$



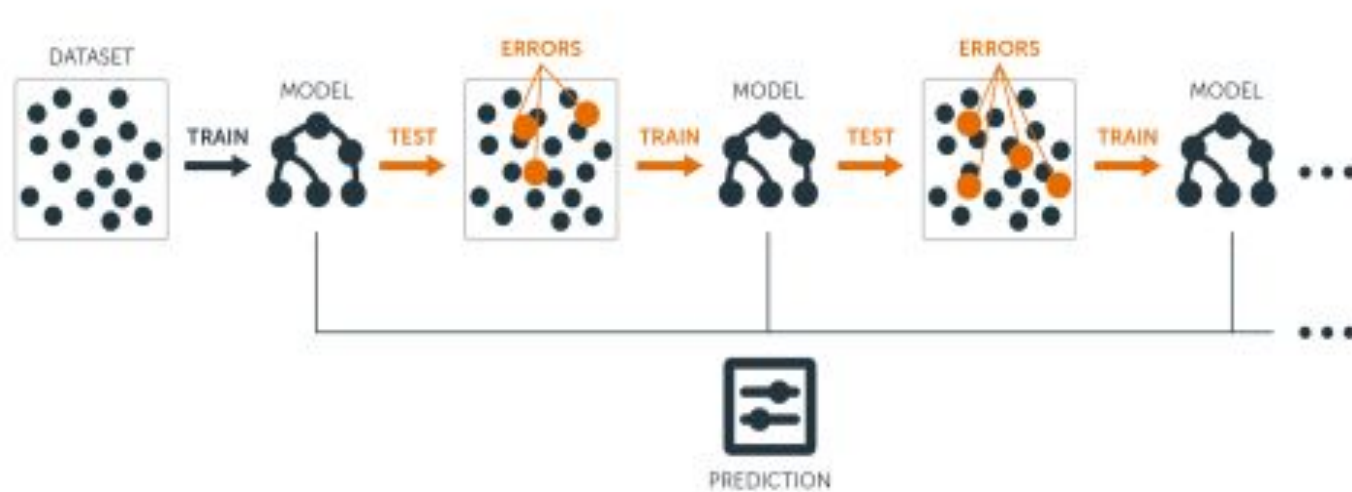
# Random Forest (Árboles)

“ Random forests provide an improvement over bagged  
decorrelates the trees. As in bagging, we build a nu  
But when building these decision trees, each time a  
predictors is chosen as split candidate from the  
Learning

Entrenar muchos árboles en paralelo utilizando diferentes muestras

# Boosting GBM( Á r b o l e s )

“ Boosting works in a similar way, except that the information from previously grown trees. Boosting does on a modified version of the original data that is statistically



# Boosting GBM( Á r b o l e s )

---

**Algorithm 8.2** *Boosting for Regression Trees*

---

1. Set  $\hat{f}(x) = 0$  and  $r_i = y_i$  for all  $i$  in the training set.
2. For  $b = 1, 2, \dots, B$ , repeat:
  - (a) Fit a tree  $\hat{f}^b$  with  $d$  splits ( $d + 1$  terminal nodes) to the training data  $(X, r)$ .
  - (b) Update  $\hat{f}$  by adding in a shrunk version of the new tree:

$$\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^b(x). \quad (8.10)$$

- (c) Update the residuals,

$$r_i \leftarrow r_i - \lambda \hat{f}^b(x_i). \quad (8.11)$$

3. Output the boosted model,

$$\hat{f}(x) = \sum_{b=1}^B \lambda \hat{f}^b(x). \quad (8.12)$$