

Tree-Based Models

Big Data y Machine Learning para Economía Aplicada

Ignacio Sarmiento-Barbieri

Universidad de los Andes

Agenda

- 1 Recap
- 2 Tree-Based Models
 - Motivation
 - Árboles
 - Bagging
 - Random Forests
 - Boosting
 - AdaBoost
 - Boosting Trees
 - XGBoost

Agenda

1 Recap

2 Tree-Based Models

- Motivation
- Árboles
- Bagging
 - Random Forests
- Boosting
 - AdaBoost
 - Boosting Trees
 - XGBoost

Recap

- Queremos predecir y en función de observables (\mathbf{x}_i)

$$y = f(\mathbf{x}_i) + u \quad (1)$$

- donde la estimación de f implica la que minimize el riesgo empírico (prediga mejor fuera de muestra):

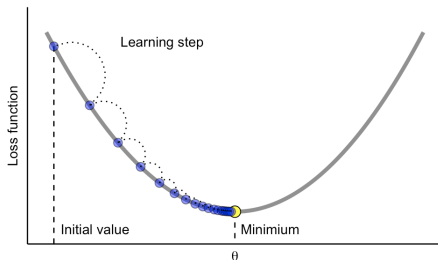
$$\hat{f} = \underset{f}{\operatorname{argmin}} \left\{ \sum_{i=1}^n L(y_i, f(\mathbf{x}_i; \Theta)) \right\} \quad (2)$$

Recap

- ▶ Hasta ahora vimos modelos lineales o linealizables.
- ▶ Por ejemplo si L es la perdida cuadrática y $f(x_i) = x_i\beta$ el problema se resume a encontrar los β .

$$\hat{\beta} = \underset{\tilde{\beta}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - \mathbf{x}_i\beta)^2 \right\} \quad (3)$$

- ▶ para estimar los β usamos en ML: gradiente descendiente



Agenda

1 Recap

2 Tree-Based Models

- Motivation
- Árboles
- Bagging
 - Random Forests
- Boosting
 - AdaBoost
 - Boosting Trees
 - XGBoost

Agenda

1 Recap

2 Tree-Based Models

- Motivation
- Árboles
- Bagging
 - Random Forests
- Boosting
 - AdaBoost
 - Boosting Trees
 - XGBoost

Motivación

► Inspirados por Leo Breiman:

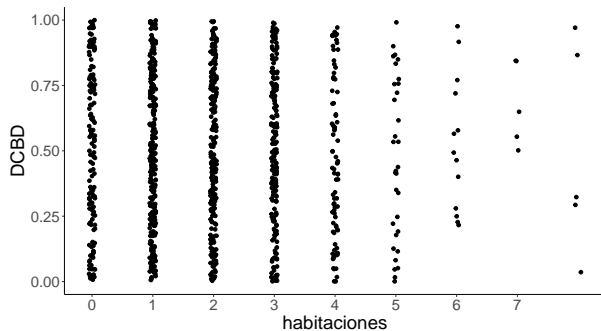
“There are two cultures in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a given stochastic data model. The other uses algorithmic models and treats the data mechanism as unknown.” Breiman [2001b], p199.

“The statistical community has been committed to the almost exclusive use of data models. This commitment has led to irrelevant theory, questionable conclusions, and has kept statisticians from working on a large range of interesting current problems. Algorithmic modeling, both in theory and practice, has developed rapidly in fields outside statistics. It can be used both on large complex data sets and as a more accurate and informative alternative to data modeling on smaller data sets. If our goal as a field is to use data to solve problems, then we need to move away from exclusive dependence on data models and adopt a more diverse set of tools.” Breiman [2001b], p199.

Motivación

- Queremos predecir:

$$\text{Precio} = f(\text{habitaciones}, \text{ubicacion}, \dots) \quad (4)$$



Agenda

1 Recap

2 Tree-Based Models

- Motivation
- Árboles
- Bagging
 - Random Forests
- Boosting
 - AdaBoost
 - Boosting Trees
 - XGBoost

Árboles

- Podemos asumir f es lineal

$$f(\mathbf{x}_i) = \beta_0 + \beta_1 \text{Habitaciones}_i + \beta_2 \text{DCBD}_i + u_i \quad (5)$$

- o ajustar un modelo flexible e interpretable como son los arboles

$$f(x_i; \{R_j, \gamma_j\}_1^J) = T(x_i; \{R_j, \gamma_j\}_1^J) \quad (6)$$

- donde

$$T(x_i; \{R_j, \gamma_j\}_1^J) = \sum_{j=1}^J \gamma_j I(x \in R_j) \quad (7)$$

Árboles

- El problema de optimización

$$\hat{f} = \operatorname{argmin}_f \left\{ \sum_{i=1}^n L(y_i, f(\mathbf{x}_i; \Theta)) \right\} \quad (8)$$

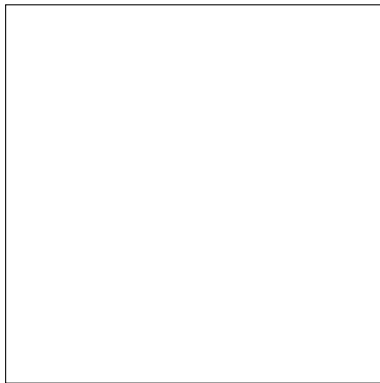
- es ahora

$$\{\hat{R}_j, \hat{\gamma}_j\}_1^J = \operatorname{argmin}_{\{R_j, \gamma_j\}_1^J} \left\{ \sum_{i=1}^n L(y_i, T(\mathbf{x}_i; \{R_j, \gamma_j\}_1^J)) \right\} \quad (9)$$

Árboles

“Recursive binary splitting”

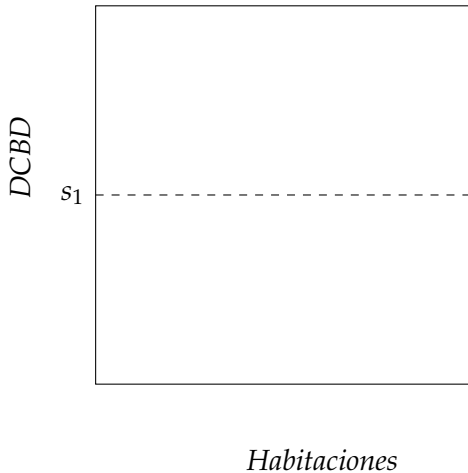
DCBD



Habitaciones

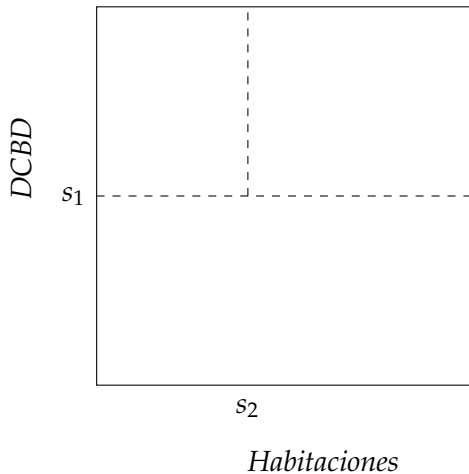
Árboles

“Recursive binary splitting”

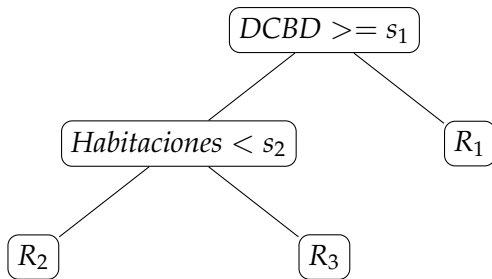


Árboles

“Recursive binary splitting”



Árboles



Ejemplo



photo from <https://www.dailydot.com/parsec/batman-1966-labels-tumblr-twitter-vine/>

Agenda

1 Recap

2 Tree-Based Models

- Motivation
- Árboles
- **Bagging**
 - Random Forests
- Boosting
 - AdaBoost
 - Boosting Trees
 - XGBoost

Bagging

- ▶ Problema con CART: pocos robustos.
- ▶ Podemos mejorar mucho el rendimiento mediante la agregación
- ▶ Idea: la varianza del promedio es menor que la de una sola predicción.

Bagging

- ▶ Bagging:
 - ▶ Obtenga repetidamente muestras aleatorias $(X_i^b, Y_i^b)_{i=1}^N$ de la muestra observada (bootstrap).
 - ▶ Para cada muestra, ajuste un árbol de regresión $\hat{f}^b(x)$
 - ▶ Promedie las muestras de bootstrap

$$\hat{f}_{bag} = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(x) \quad (10)$$

- ▶ Básicamente estamos suavizando las predicciones.

Random Forests

- ▶ Problema con el bagging: si hay un predictor fuerte, diferentes árboles son muy similares entre sí.
- ▶ Bosques (forests): reduce la correlación entre los árboles en el bootstrap.
- ▶ Si hay p predictores, en cada partición use solo $m < p$ predictores, elegidos al azar.
- ▶ Bagging es forests con $m = p$ (usando todo los predictores en cada partición).
- ▶ m es un hiper-parámetro, $m = \sqrt{p}$ es un benchmark

Ejemplo



photo from <https://www.dailydot.com/parsec/batman-1966-labels-tumblr-twitter-vine/>

Agenda

1 Recap

2 Tree-Based Models

- Motivation
- Árboles
- Bagging
 - Random Forests
- **Boosting**
 - AdaBoost
 - Boosting Trees
 - XGBoost

Boosting: Motivation

- ▶ Problema con CART: varianza alta.
- ▶ Podemos mejorar mucho el rendimiento mediante la agregación
- ▶ El boosting toma esta idea pero lo “encara” de una manera diferente → viene de la computación
- ▶ Va a usar arboles “pequeños” y a aprender secuencialmente

Boosting: Motivation

- ▶ Boosting es una de las ideas de aprendizaje más poderosas introducidas en los últimos veinte años.
- ▶ Originalmente fue diseñado para problemas de clasificación,
- ▶ Pero también puede extenderse provechosamente a la regresión.

Boosting: Motivation

- ▶ Intuitivamente lo que hacen es aprender de los errores lentamente.
- ▶ Esto lo hace "ajustando" árboles pequeños.
- ▶ Esto permite mejorar lentamente aprendiendo $f(\cdot)$ en áreas donde no funciona bien.
- ▶ OJO: a diferencia de *bagging*, la construcción de cada árbol depende en gran medida de los árboles anteriores

AdaBoost.M1

- ▶ El AdaBoost.M1 es una implementación para clasificación

AdaBoost.M1

- ▶ El AdaBoost.M1 es una implementación para clasificación
- ▶ Como viene de la computación tiene un lenguaje ligeramente diferente.
- ▶ Vocabulario:
 - ▶ $y \in -1, 1$, X vector de predictores.
 - ▶ $\hat{y} = G(X)$ (clasificador)
 - ▶ $err = \frac{1}{N} \sum_i^N I(y_i \neq G(x_i))$

AdaBoost.M1

- ▶ Vocabulario:
 - ▶ $y \in -1, 1$, X vector de predictores.
 - ▶ $\hat{y} = G(X)$ (clasificador)
 - ▶ $err = \frac{1}{N} \sum_i^N I(y_i \neq G(x_i))$

AdaBoost.M1

- 1 Comenzamos con ponderadores $w_i = 1/N$
- 2 Para $m = 1$ hasta M :
 - 1 Estimar $G_m(x)$ usando ponderadores w_i .
 - 2 Computar el error de predicción

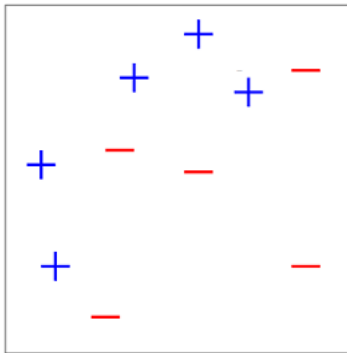
$$err_m = \frac{\sum_{i=1}^N w_i I(y_i \neq G_m(x_i))}{\sum_{i=1}^N w_i} \quad (11)$$

- 3 Obtener $\alpha_m = \ln \left[\frac{(1-err_m)}{err_m} \right]$
- 4 Actualizar los ponderadores : $w_i \leftarrow w_i c_i$

$$c_i = \exp [\alpha_m I(y_i \neq G_m(x_i))] \quad (12)$$

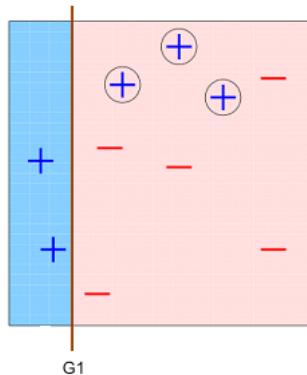
- 3 Resultado: $G(x) = \text{sign}[\sum_{m=1}^M \alpha_m G_m(x)]$

AdaBoost.M1



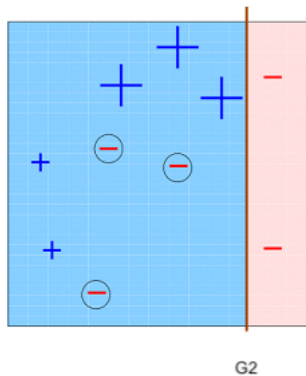
AdaBoost.M1

► Iteración 1



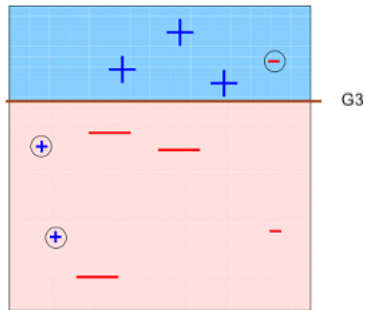
AdaBoost.M1

► Iteración 2



AdaBoost.M1

► Iteración 3



AdaBoost.M1

► Resultado Final

$$G_{\text{final}} = \text{sign} \left(\alpha_1 \begin{array}{|c|} \hline \text{blue} \\ \hline \text{red} \end{array} + \alpha_2 \begin{array}{|c|} \hline \text{blue} \\ \hline \text{blue} \end{array} + \alpha_3 \begin{array}{|c|} \hline \text{blue} \\ \hline \text{red} \end{array} \right) = \begin{array}{|c|c|c|} \hline \text{blue} & + & + \\ \hline \text{blue} & + & + \\ \hline \text{red} & - & - \\ \hline \end{array}$$

The diagram illustrates the final result of the AdaBoost.M1 algorithm. It shows the combination of three weak classifiers, each represented by a square divided into blue and red regions, with weights α_1 , α_2 , and α_3 . The final result is a square divided into blue and red regions, with blue regions containing '+' signs and red regions containing '-' signs.

AdaBoost.M1

- 1 Comenzamos con ponderadores $w_i = 1/N$
- 2 Para $m = 1$ hasta M :
 - 1 Estimar $G_m(x)$ usando ponderadores w_i .
 - 2 Computar el error de predicción

$$err_m = \frac{\sum_{i=1}^N w_i I(y_i \neq G_m(x_i))}{\sum_{i=1}^N w_i} \quad (13)$$

- 3 Obtener $\alpha_m = \ln \left[\frac{(1-err_m)}{err_m} \right]$
- 4 Actualizar los ponderadores : $w_i \leftarrow w_i c_i$

$$c_i = \exp [\alpha_m I(y_i \neq G_m(x_i))] \quad (14)$$

- 3 Resultado: $G(x) = \text{sign}[\sum_{m=1}^M \alpha_m G_m(x)]$

AdaBoost.M1

- ▶ $c_i = \exp[\alpha_m I(y_i \neq G_m(x_i))]$
- ▶ Si fue correctamente predicho, $c_i = 1$.
- ▶ En caso contrario, $c_i = \exp(\alpha_m) = \frac{(1 - \text{err}_m)}{\text{err}_m} > 1$
- ▶ En cada paso el algoritmo da mas importancia relativa a las predicciones incorrectas.
- ▶ Paso final: promedio ponderado de estos pasos

$$G(x) = \text{sign}\left[\sum_{m=1}^M \alpha_m G_m(x)\right] \quad (15)$$

Boosting Trees

- El objetivo es

$$\hat{f} = \operatorname{argmin}_f \left\{ \sum_{i=1}^n L(y_i, f(\mathbf{x}_i; \Theta)) \right\} \quad (16)$$

- Boosting

$$f(x_i; \{R_j, \gamma_j\}_1^J, M) = \sum_{m=1}^M T(x_i; \Theta_m) \quad (17)$$

- el problema se vuelve

$$\hat{f} = \operatorname{argmin}_f \left\{ \sum_{i=1}^n L(y_i, \sum_{m=1}^M T(x_i; \Theta_m)) \right\} \quad (18)$$

Boosting Trees

- El problema

$$\hat{f} = \operatorname{argmin}_f \left\{ \sum_{i=1}^n L(y_i, \sum_{m=1}^M T(x_i; \Theta_m)) \right\} \quad (19)$$

- Se puede aproximar por

$$\hat{\Theta}_m = \operatorname{argmin}_{\Theta_m} \left\{ \sum_{i=1}^n L(y_i, f_{m-1} + T(x_i; \Theta_m)) \right\} \quad (20)$$

- Para una función de pérdida cuadrática

$$\hat{\Theta}_m = \operatorname{argmin}_{\Theta_m} \left\{ \sum_{i=1}^n (y_i - f_{m-1} - T(x_i; \Theta_m))^2 \right\} \quad (21)$$

Boosting Trees: Algoritmo

1 Iniciamos fijando $\hat{f}(x) = 0$ y $r_i = y_i$ para todos los i del training set

2 Para $m = 1, 2, \dots, M$

1 Ajustamos un árbol $\hat{f}^m(x) = 0$ con J bifurcaciones

2 Actualizamos $\hat{f}(x)$ con una versión "shrunk" del nuevo árbol

$$\hat{f}(x) \leftarrow \hat{f}(x) + \hat{f}^m(x) \quad (22)$$

3 Actualizamos los residuales

$$r_i \leftarrow r_i - \hat{f}^m(x) \quad (23)$$

3 El modelo final es

$$\hat{f}_{boost} = \sum_{m=1}^M \hat{f}^m(x) \quad (24)$$

Boosting Trees: Hiper-parámetros

- ▶ Cuántas iteraciones (M) usar?
 - ▶ Cada iteración generalmente reduce el error de ajuste, de modo que para M lo suficientemente grande este error puede hacerse arbitrariamente pequeño (sesgo se va a cero).
 - ▶ Sin embargo, ajustar demasiado bien los datos de entrenamiento puede llevar a overfit (sobreajuste)
 - ▶ Por lo tanto, hay un número óptimo M^* que minimiza el error fuera de muestra
 - ▶ Una forma conveniente de encontrar M^* con validación cruzada

Boosting Trees: Hiper-parámetros

- Podemos utilizar “shrinkage”? \Rightarrow Yes!
- Se escala la contribución de cada árbol por un factor $\lambda \in (0, 1)$

1 Iniciamos fijando $\hat{f}(x) = 0$ y $r_i = y_i$ para todos los i del training set

2 Para $m = 1, 2, \dots, M$

1 Ajustamos un árbol $\hat{f}^m(x) = 0$ con J bifurcaciones

2 Actualizamos $\hat{f}(x)$

$$\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^m(x) \quad (25)$$

3 Actualizamos los residuales

$$r_i \leftarrow r_i - \lambda \hat{f}^m(x) \quad (26)$$

3 El modelo final es

$$\hat{f}_{boost} = \sum_{m=1}^M \lambda \hat{f}^m(x) \quad (27)$$

Boosting Trees: Iteraciones

- ▶ Los otros hiperparámetros a fijar son
 - ▶ λ la tasa a la que aprende, los valores típicos son 0.01 o 0.001
 - ▶ El tamaño del árbol.
 - ▶ $4 \leq J \leq 8$ funciona bien
 - ▶ $J = 2$ suele ser insuficiente
 - ▶ $J > 10$ rara vez se utiliza

Gradient Boosting Trees

- For **arbitrary** loss functions this “simplified” problem

$$\hat{\Theta}_m = \operatorname{argmin}_{\Theta_m} \left\{ \sum_{i=1}^n L(y_i, f_{m-1} + T(x_i; \Theta_m)) \right\} \quad (28)$$

- is still a tricky problem
- We can solve it by using an implementation of gradient descent

$$f_m = f_{m-1} - \epsilon \sum_{i=1}^n \nabla_{f_{m-1}} L(y_i, f_{m-1}(x)) \quad (29)$$

- How do we implement this search?

Gradient Boosting Trees: Algorithm

1 Initialize $f_0(x) = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$.

2 For $m=1$ to M :

1 For $i = 1, 2, \dots, N$ compute

$$r_{im} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}} \quad (30)$$

2 Fit a regression tree to the targets r_{im} giving terminal regions $R_{jm}, j = 1, 2, \dots, J_m$.

3 For $j = 1, 2, \dots, J_m$ compute

$$\gamma_{jm} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1} + \gamma) \quad (31)$$

4 Update

$$f_m = f_{m-1} + \lambda \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm}) \quad (32)$$

3 El modelo final es

$$\hat{f}(x) = \lambda \sum_{m=1}^M \hat{f}_m(x_i) \quad (33)$$

XGBoost: Motivation

- ▶ XGBoost is an implementation was specifically engineered for optimal performance and speed.
- ▶ Why talk about XGBoost?
 - ▶ Among the 29 challenge winning solutions published in Kaggle's blog during 2015, 17 solutions used XGBoost.
 - ▶ Among these solutions, eight solely used XGBoost to train the model, while most others combined XGBoost with neural nets in ensembles. (The second most popular method, deep neural nets, was used in 11 solutions)
 - ▶ The success of the system was also witnessed in 2015 Data Mining and Knowledge Discovery competition organized by ACM (KDD Cup) , where XGBoost was used by every winning team in the top-10.
 - ▶ Historically, XGBoost has performed quite well for structured, tabular data. But, if you are dealing with non-structured data such as images, neural networks are usually a better option (more on this later)

XGBoost

- Cual es la innovación? Penaliza la optimización, en cada paso la función objetivo para construir el árbol es

$$\mathcal{L} = \sum_{i=1}^N L(y_i, \sum_{m=1}^M T(x_i; \Theta_m)) + \sum_{m=1}^M \Omega(T(x_i; \Theta_m)) \quad (34)$$

- El segundo termino penaliza la complejidad del modelo,

$$\Omega(.) = \zeta J + \frac{1}{2} \phi \sum_{j=1}^J \gamma_j^2 \quad (35)$$

Example



photo from <https://www.dailydot.com/parsec/batman-1966-labels-tumblr-twitter-vine/>