

Introduction to XGBoost

March 13, 2019

- 1 Preliminary: Objective Functions
- 2 Preliminary: Regression Trees (CARTs)
 - A Brief Introduction
 - Generate Regression Tree
 - Pre-Stopping and Pruning
- 3 The Model of XGBoost: Tree Ensembles
 - Tree Ensembles
- 4 Tree Boosting
 - Preparation
 - Objective Function
 - Train Our Model

Objective Functions

Definition (Objective Functions)

$\theta \in \Theta$ - the parameter space.

$$Obj(\theta) = L(\theta) + \Omega(\theta)$$

L : **Loss function** measures how well model fit on training data

Ω : **Regularization** measures complexity of model

Example (Objective Function for Ridge Regression)

$\beta \in \mathbb{R}^n$. λ is a constant number.

$$Obj(\beta) = \|y - X\beta\|_2^2 + \lambda\|\beta\|_2^2$$

- 1 Preliminary: Objective Functions
- 2 Preliminary: Regression Trees (CARTs)
 - A Brief Introduction
 - Generate Regression Tree
 - Pre-Stopping and Pruning
- 3 The Model of XGBoost: Tree Ensembles
 - Tree Ensembles
- 4 Tree Boosting
 - Preparation
 - Objective Function
 - Train Our Model

Regression Trees (CARTs)

How do we use **Classification And Regression Tree (CART)** to regress or classify?

Example (Regression Tree)

Every regression tree T has the following form:

$$\hat{y}_i = \sum_{m=1}^M \omega_m 1_{x \in R_m}.$$

Two Elements completely determine a regression tree:

- 1 Partition of Sample Space $\{R_1, \dots, R_M\}$,
- 2 and Weights Vector $\omega = (\omega_1, \dots, \omega_M)$.

Example (Regression Tree: Revisit)

Given a partition of sample space $\{R_1, \dots, R_M\}$.

For every sample point x , there must exist $n \in \{1, 2, \dots, M\}$ such that $x \in R_n$. Define

$$q : x \mapsto n,$$

which describes the partition $\{R_1, \dots, R_M\}$.

\implies re-write regression tree as

$$T = \omega_q$$

- 1 Preliminary: Objective Functions
- 2 Preliminary: Regression Trees (CARTs)
 - A Brief Introduction
 - **Generate Regression Tree**
 - Pre-Stopping and Pruning
- 3 The Model of XGBoost: Tree Ensembles
 - Tree Ensembles
- 4 Tree Boosting
 - Preparation
 - Objective Function
 - Train Our Model

Generate Regression Tree

Problem: Given N observations (x_i, y_i) , $i = 1, 2, \dots, N$. We want to build a regression tree such that

$$L = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

is minimized.

Solution: Notice the following **fact**:

Given q , there is unique $\hat{\omega}$ such that L is minimized.

$$\begin{aligned}\hat{\omega}_m &= \frac{1}{|R_m|} \sum_{x_i \in R_m} y_i \\ &:= \text{avg} R_m \quad m = 1, \dots, M.\end{aligned}$$

Generate Regression Tree

Because of the fact, our problem become:

How to find the best tree structure q ?

Answer: We generate the tree from top to bottom.

Example (Generate a regression tree)

Define $R_{j,s} = \{X \mid X_j \leq s\}$.

Solve the following optimization problem:

$$\min_{j,s} \left(\sum_{x_i \in R_{j,s}} (y_i - \text{avg} R_{j,s})^2 + \sum_{x_i \notin R_{j,s}} (y_i - \text{avg} R_{j,s}^c)^2 \right)$$

How to solve it: s is chosen from $\{x_{ij}\}_{i=1,\dots,N}$. Enumerate every pair (j, s) .

Then we will add a split at each subnode of the first node. But **when should we stop adding a new split?**

- 1 Preliminary: Objective Functions
- 2 Preliminary: Regression Trees (CARTs)
 - A Brief Introduction
 - Generate Regression Tree
 - Pre-Stopping and Pruning
- 3 The Model of XGBoost: Tree Ensembles
 - Tree Ensembles
- 4 Tree Boosting
 - Preparation
 - Objective Function
 - Train Our Model

Pre-Stopping and Pruning

Pre-Stopping: We can stop adding a new spit when ...

- *max_depth*
- *min_child_weight*
- *gamma*
- ...

Pruning

Suppose we have a partition into M regions R_1, \dots, R_M . Tree T is built on the partition.

Definition (Loss function)

Let $N_m = |R_m|$, $\hat{w}_m = \text{avg} R_m$, and $Q_m(T) = \frac{1}{|R_m|} \sum_{x_i \in R_m} (y_i - \hat{w}_m)^2$.

$$L(T) := \sum_{m=1}^{|T|} N_m Q_m(T).$$

Definition (Complexity of our model)

Given a constant α . Define

$$\Omega(T) = \alpha |T|.$$

Problem: We have built a large tree T_0 . Find a subtree T_α such that

$$\begin{aligned} C_\alpha(T) &= L(T) + \Omega(T) \\ &= \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T| \end{aligned}$$

is minimized.

Solution: Let $\mathfrak{N} \in T^\circ := \{\text{All internal nodes of } T\}$. The operator

$$\sigma_{\mathfrak{N}} : T \mapsto \sigma_{\mathfrak{N}} T$$

collapses the node \mathfrak{N} to get the subtree $\sigma_{\mathfrak{N}} T$.

Solving the following optimization problem

$$\min_{\mathfrak{N} \in T^\circ} (L(T) - L(\sigma_{\mathfrak{N}} T))$$

we get $\tilde{\mathfrak{N}}$.

Now let's begin from T_0 . We successively collapse node and finally get a sequence of subtrees

$$\{\sigma_{\tilde{\mathfrak{N}}_0} T_0, \sigma_{\tilde{\mathfrak{N}}_1} \sigma_{\tilde{\mathfrak{N}}_0} T_0, \dots, \{\mathfrak{R}\}\},$$

where \mathfrak{R} is the root of T_0 .

Find T_α such that the value of $(L(T) - L(\sigma_{\mathfrak{N}} T))$ is the smallest in this set.

- 1 Preliminary: Objective Functions
- 2 Preliminary: Regression Trees (CARTs)
 - A Brief Introduction
 - Generate Regression Tree
 - Pre-Stopping and Pruning
- 3 The Model of XGBoost: Tree Ensembles**
 - Tree Ensembles
- 4 Tree Boosting
 - Preparation
 - Objective Function
 - Train Our Model

Definition (Tree Ensemble Model)

Given a fixed sample point $x_i = (x_{i1}, \dots, x_{in})$, our model will predict:

$$\hat{y}_i = \sum_{k=1}^t f_k(x_i),$$

where f_k is a CART, for each k .

Question: How to find parameters $\{f_1, \dots, f_t\}$ of this model?

Answer: Define an objective function, and optimize it!.

- 1 Preliminary: Objective Functions
- 2 Preliminary: Regression Trees (CARTs)
 - A Brief Introduction
 - Generate Regression Tree
 - Pre-Stopping and Pruning
- 3 The Model of XGBoost: Tree Ensembles
 - Tree Ensembles
- 4 **Tree Boosting**
 - **Preparation**
 - Objective Function
 - Train Our Model

Tree Boosting

Aim: Find all parameters $\{f_1, \dots, f_S\}$ of our model

$$\hat{y} = \sum_{k=1}^S f_k(x).$$

Method: We'll find f_k one by one. Assume we have a part of our model:

$$\hat{y}_i^{(t-1)} = \sum_{j=1}^{t-1} f_j(x_i).$$

We want to find the t-th parameter in the space of **all CARTs**

$$\Theta = \{f_1^{(t)}, \dots, f_K^{(t)}, \dots\}.$$

We construct the following optimization problem to find f_t :

$$\min_{f_t \in \Theta} \text{Obj}^{(t)}(f_t) = L^{(t)}(f_t) + \Omega(f_t)$$

- 1 Preliminary: Objective Functions
- 2 Preliminary: Regression Trees (CARTs)
 - A Brief Introduction
 - Generate Regression Tree
 - Pre-Stopping and Pruning
- 3 The Model of XGBoost: Tree Ensembles
 - Tree Ensembles
- 4 **Tree Boosting**
 - Preparation
 - **Objective Function**
 - Train Our Model

$L^{(t)}$ - Loss Function

Re-Write Loss Function $L^{(t)}$:

$$\begin{aligned} L^{(t)}(f_t) &= \sum_{j=1}^N l(y_j, \hat{y}_j^{(t)}) \\ &= \sum_{j=1}^N l(y_j, \hat{y}_j^{(t-1)} + f_t(x_j)). \end{aligned}$$

Definition (Objective Function $Obj^{(t)}$)

$$Obj^{(t)}(f_t) = \sum_{j=1}^N l(y_j, \hat{y}_j^{(t-1)} + f_t(x_j)) + \Omega(f_t) + \sum_{k=1}^{t-1} \Omega(f_k)$$

Note: It is hard to minimize it. We need a simplified version of loss.

Taylor Expansion Approximation of Loss

Notice that for the j -th sample point

$$l(y_j, \hat{y}_j^{(t-1)} + f_t(x_j)) \approx l(y_j, \hat{y}_j^{(t-1)}) + g_j f_t(x_j) + \frac{1}{2} h_j f_t^2(x_j),$$

where $g_j = \partial_{\hat{y}^{(t-1)}} l(y_j, \hat{y}^{(t-1)})$ and $h_j = \partial_{\hat{y}^{(t-1)}}^2 l(y_j, \hat{y}^{(t-1)})$.

Definition (Loss Function)

$$L^{(t)}(f_t) = \underbrace{\sum_{j=1}^N (l(y_j, \hat{y}_j^{(t-1)}))}_{\text{doesn't depend on } f_t} + \sum_{j=1}^N (g_j f_t(x_j) + \frac{1}{2} h_j f_t^2(x_j))$$

Note: $\sum_{j=1}^N (l(y_j, \hat{y}_j^{(t-1)}))$ doesn't depend on f_t .

Ω - Complexity of CART

We re-write a given CART T as follow

$$T(x) = \omega_{q(x)}.$$

Note: Every CART T is determined by two element:

- q - the structure of T
- ω - the weights of each leaf of T

Definition

Given constants γ and λ , we define

$$\Omega(T) = \gamma|T| + \frac{1}{2}\lambda \sum_{j=1}^{|T|} w_j^2$$

Revisit: Objective Function

Definition (Final Objective Function)

Set ...

$$I_j = \{ \}$$

$$\begin{aligned} \text{Obj}^{(t)}(f_t) &\approx \sum_{j=1}^N (g_j f_t(x_j) + \frac{1}{2} h_j f_t^2(x_j)) + \gamma |T| + \frac{1}{2} \lambda \sum_{j=1}^{|T|} w_j^2 \\ &= \sum_{j=1}^{|f_t|} \underbrace{\left[\left(\sum_{i \in I_j} g_i \right) \omega_j + \frac{1}{2} \left(\sum_{i \in I_j} h_i + \lambda \right) \omega_j^2 \right]}_{\text{red line}} + \gamma |f_t| \\ &:= \sum_{j=1}^{|f_t|} \underbrace{\left[G_j \omega_j + \frac{1}{2} (H_j + \gamma) \omega_j^2 \right]}_{\text{red line}} \end{aligned}$$

Revisit: Objective Function

Definition (Final Objective Function)

$$Obj^{(t)}(f_t) = \sum_{j=1}^{|f_t|} [G_j \omega_j + \frac{1}{2} (H_j + \gamma) \omega_j^2] + \gamma |f_t|$$

Note: Given a tree structure q , we can immediately get a CART minimizing $Obj^{(t)}$, by setting $\omega_j = -\frac{G_j}{H_j + \lambda}$.

And with the same q ,

$$Obj^{(t)}(f_t) = -\frac{1}{2} \sum_{j=1}^{|f_t|} \frac{G_j^2}{H_j + \lambda} + \gamma |f_t|.$$

- 1 Preliminary: Objective Functions
- 2 Preliminary: Regression Trees (CARTs)
 - A Brief Introduction
 - Generate Regression Tree
 - Pre-Stopping and Pruning
- 3 The Model of XGBoost: Tree Ensembles
 - Tree Ensembles
- 4 Tree Boosting
 - Preparation
 - Objective Function
 - Train Our Model

Reminder of Our Aim: Find all parameters $\{f_1, \dots, f_S\}$ of our model

$$\hat{y} = \sum_{k=1}^S f_k(x).$$

Method: We construct the following optimization problem to find t-th parameter f_t :

$$\min_{f_t \in \Theta} \text{Obj}^{(t)}(f_t) = \sum_{j=1}^{|f_t|} [G_j \omega_j + \frac{1}{2} (H_j + \gamma) \omega_j^2] + \gamma |f_t|.$$

Generate f_t

Similar to generate a CART, it suffices to find **the best tree structure** q .

Note: To get q_{new} , we add two nodes \mathfrak{N}_L and \mathfrak{N}_R at \mathfrak{N} of q_{old} .

Define

$$\begin{aligned}\text{Gain} &= \text{Obj}^{(t)}(q_{\text{old}}) - \text{Obj}^{(t)}(q_{\text{new}}) \\ &= \frac{1}{2} \left[\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma\end{aligned}$$

We add a new split with the highest value of Gain.

Pre-Stopping and Post-Pruning

Pre-Stopping We can stop adding a new split when

- Gain is negative.
- Reach the maximum depth.
- ...

Post-Pruning

- **CARTs part refers to this classical textbook.** *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. 2nd Edition. Chapter 9.
- **XGBoost refers to Tianqi Chen's slide, *Introduction to Boosted Trees*.**
`https://homes.cs.washington.edu/~tqchen/pdf/BoostedTree.pdf`
- **and to this document,**
`http://xgboost.readthedocs.io/en/latest/model.html`