

1 Appendix A

2

3 In this appendix, GTB algorithm is briefly reviewed according to Chen and
4 Guestrin [16]. Let $x \in \mathbb{R}^m$ be m dimensional input variables, $y \in \mathbb{R}$ be an
5 output variable, and $D = \{(x_i, y_i)_{i=1}^n | x_i \in \mathbb{R}^m, y_i \in \mathbb{R}\}$ be n number of data. A
6 regression tree $f(x)$ having T leaves and leaf weight vector $w \in \mathbb{R}^T$ is defined as

$$7 \quad f(x) = w_{q(x)},$$
$$8 \quad q: \mathbb{R}^m \rightarrow \{1, \dots, T\}, w \in \mathbb{R}^T,$$

9 where q is a function that maps input variables $x \in \mathbb{R}^m$ to leaf index $i \in$
10 $\{1, \dots, T\}$. The function q specifies the tree structure of a regression tree. $w \in \mathbb{R}^T$
11 specifies the leaf weights of the regression tree. Let \mathcal{F} be the whole set of the
12 regression trees having T leaves.

13 In ensemble learning such as Bagging and Boosting, outcome y is predicted
14 using an ensemble $\phi(x)$ of K number of weak-learners $\{f_k(x)\}$:

$$15 \quad \hat{y}_i = \phi(x_i) = \sum_{k=1}^K f_k(x_i), f_k \in \mathcal{F}$$

16 For each regression tree $f_k(x)$, we need to learn tree structure q and leaf weights
17 $w \in \mathbb{R}^T$ from training data. Tree structure q and leaf weights $w \in \mathbb{R}^T$ are
18 learned by minimizing the following regularized objective function:

$$19 \quad L(\phi) = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k),$$
$$20 \quad \Omega(f) = \gamma T + \frac{1}{2} \lambda \|w\|^2,$$

21 Here l is a differentiable convex function that measures difference between
22 predicted values \hat{y}_i and the true value y_i (fitting loss). $\Omega(f)$ is a
23 regularization term introduced in order to prevent over-fitting. The penalty term
24 γT is added so that leaf number T is small and resulting regression tree $f_k(x)$ is

25 simple. The penalty term $\frac{1}{2} \lambda \|w\|^2$ is added so that l_2 -norm of the leaf weights w
26 is small. γ, λ are hyperparameters of the XGBoost.

27 In the Gradient Boosting, a weak-learner f_k is trained in additive manner. Let
28 $\hat{y}_i^{(t)}$ be the predicted value for the i -th data at iteration step t . Using a new
29 function f_t , objective function $L(\phi)$ is rewritten as follows:

1
$$L^{(t)}(f_t) = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t).$$

2 Given $\hat{y}_i^{(t-1)}$, regression tree f_t is greedily optimized to minimize the objective
 3 function $L^{(t)}(f_t)$. XGBoost exploits the second order approximation of $L^{(t)}(f_t)$
 4 using Taylor expansion, i.e.,

5
$$L^{(t)}(f_t) \approx \sum_{i=1}^n \left[l(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \Omega(f_t).$$

6 Here $g_i = \partial_{\hat{y}_i^{(t-1)}} l(y_i, \hat{y}_i^{(t-1)})$ and $h_i = \partial_{\hat{y}_i^{(t-1)}}^2 l(y_i, \hat{y}_i^{(t-1)})$ are gradient and
 7 hessian statistics computed from first and second order derivatives of the loss
 8 function l , respectively.

9 Objective function $L^{(t)}(f_t)$ as a function of f_t is simplified as

10
$$\tilde{L}^{(t)} = \sum_{i=1}^n \left[g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \gamma T + \frac{1}{2} \lambda \|w\|^2.$$

11 For each leaf $j \in \{1, \dots, T\}$ of a regression tree f_t with a tree structure q , we
 12 define an index set $I_j = \{i | q(x_i) = j\}$. We define $G_j = \sum_{i \in I_j} g_i$, $H_j = \sum_{i \in I_j} h_i$.

13 Objective function $\tilde{L}^{(t)}$ is rewritten as

17
$$\tilde{L}^{(t)} = \sum_{j=1}^T \left[G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2 \right] + \gamma T.$$

14 Note that the summation is only taken through leaves $j \in \{1, \dots, T\}$ and not
 15 through data. Since the objective function $\tilde{L}^{(t)}$ is quadratic with respect to leaf
 16 weight vector w , the optimal leaf weights can be solved as

18
$$w_j^* = -\frac{G_j}{H_j + \lambda} \quad j \in \{1, \dots, T\}.$$

19 By substituting the optimal weight vector w^* into the objective function $\tilde{L}^{(t)}$, we
 20 obtain

21
$$\tilde{L}^{(t)}(q) = -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{H_j + \lambda} + \gamma T$$

22 Note that once a tree structure q of a regression tree f_t is given, the optimal
 23 weight vector w^* and $\tilde{L}^{(t)}(q)$ are computed. Objective function $\tilde{L}^{(t)}(q)$ is used
 24 as a score of the tree structure q .

25 It is computationally severe to enumerate all candidates of tree structures q . To

1 obtain an efficient and approximate algorithm, a greedily optimal splitting of a leaf
2 node is explored using the score function $\tilde{L}^{(t)}(q)$. Let q_B be a tree structure before
3 splitting of a leaf node and q_A be a tree structure after splitting of the leaf node.
4 Then the gain of the splitting is measured by a score:

5
$$L_{split} = \tilde{L}^{(t)}(q_B) - \tilde{L}^{(t)}(q_A).$$

6 A greedily optimal splitting is determined by that maximizing the gain. A greedily
7 optimal tree structure q^* is explored by repeating the above procedure. XGBoost
8 is elaborated on speeding up the above algorithm and increasing the scalability,
9 with taking into account sparsity-aware algorithm, out-of-core computing, and
10 parallel computing. See Chen and Guestrin [16] for technical details.