
Orthogonal Gradient Boosting for Interpretable Additive Rule Ensembles

Supplementary Information

Anonymous Author(s)

Affiliation

Address

email

1 A Proofs

2 A.1 Proof of Lemma 4.1

3 *Proof.* Let $\mathbf{f} = [\mathbf{Q}; \mathbf{g}] \alpha$ and $\tilde{\mathbf{f}} = [\mathbf{Q}; \mathbf{q}] \beta$. We can decompose the squared error

$$\begin{aligned} \|\mathbf{f} - \tilde{\mathbf{f}}\|^2 &= \|[\mathbf{Q}; \mathbf{g}] \alpha - [\mathbf{Q}; \mathbf{q}] \beta\|^2 \\ &= \|[\mathbf{Q}; \mathbf{g}_{\parallel} + \mathbf{g}_{\perp}] \alpha - [\mathbf{Q}; \mathbf{q}_{\parallel} + \mathbf{q}_{\perp}] \beta\|^2 \\ &= \|[\mathbf{Q}; \mathbf{g}_{\parallel}] \alpha + \alpha_t \mathbf{g}_{\perp} - [\mathbf{Q}; \mathbf{q}_{\parallel}] \beta + \beta_t \mathbf{q}_{\perp}\|^2 \\ &= \|[\mathbf{Q}; \mathbf{g}_{\parallel}] \alpha - [\mathbf{Q}; \mathbf{q}_{\parallel}] \beta\|^2 + \|\alpha_t \mathbf{g}_{\perp} - \beta_t \mathbf{q}_{\perp}\|^2 \end{aligned}$$

4 where the last step follows from the Pythagorean theorem and the fact that $\alpha_t \mathbf{g}_{\perp} - \beta_t \mathbf{q}_{\perp}$ is an element
5 from the orthogonal complement of $\text{range}[\mathbf{Q}; \mathbf{g}_{\parallel}] = \text{range}[\mathbf{Q}; \mathbf{g}_{\parallel}] = \text{range } \mathbf{Q}$. The equality of these
6 ranges also implies that $\beta_1, \dots, \beta_{t-1}$ can, for all choices of β_t , be chosen such that the left term of
7 the error decomposition is 0. Setting $\gamma = \beta_t / \alpha_t$, it follows that

$$\begin{aligned} \min_{\beta \in \mathbb{R}^t} \|\mathbf{f} - \tilde{\mathbf{f}}\|^2 &= \min_{\beta \in \mathbb{R}^t} \|\alpha_t \mathbf{g}_{\perp} - \beta_t \mathbf{h}_{\perp}\|^2 \\ &= \min_{\gamma \in \mathbb{R}^t} \alpha_t^2 \|\mathbf{g}_{\perp} - \gamma \mathbf{q}_{\perp}\|^2 \\ &= \min_{\gamma \in \mathbb{R}^t} \alpha_t^2 (\|\mathbf{g}_{\perp}\|^2 - 2\gamma \mathbf{q}_{\perp}^T \mathbf{g}_{\perp} + \gamma^2 \|\mathbf{q}_{\perp}\|^2) \end{aligned}$$

8 and plugging in the minimizing $\gamma = \mathbf{q}_{\perp}^T \mathbf{g}_{\perp} / \|\mathbf{q}_{\perp}\|^2$

$$= \alpha^2 (\|\mathbf{g}_{\perp}\|^2 - (\mathbf{g}_{\perp}^T \mathbf{q}_{\perp})^2 / \|\mathbf{q}_{\perp}\|^2) ,$$

9 from which, noting that $\mathbf{g}_{\perp}^T \mathbf{q}_{\perp} = \mathbf{g}_{\perp}^T \mathbf{q}$, the claim follows. □

10 A.2 Proof of Lemma 4.2

11 *Proof.* After the weight correction step β is a stationary point of $R(\mathbf{Q}(\cdot))$, i.e., we have for all $j \in [t]$

$$0 = \frac{\partial R(\mathbf{Q}\beta)}{\partial \beta_j} = \sum_{i=1}^n \frac{\partial l(\tilde{\mathbf{q}}_i^T \beta, y_i)}{\partial \beta_j} = \sum_{i=1}^n q_{ij} \underbrace{\frac{\partial l(\tilde{\mathbf{q}}_i^T \beta, y_i)}{\partial \tilde{\mathbf{q}}_i^T \beta}}_{g_i} = \mathbf{q}_j^T \mathbf{g} .$$

12 □

13 A.3 Proof of Theorem 4.3

14 The condition of Theorem 4.3 states that:

15 Let $\mathbf{Q} \in \mathbb{R}^{n \times (t-1)}$ be the selected query matrix, \mathbf{g} the corresponding gradient vector after a full
16 weight correction, and \mathbf{q}^* be a maximizer of the **orthogonal gradient boosting objective** function
17 defined by

$$\text{obj}_{\text{ogb}}(q) = \frac{|\mathbf{g}^T \mathbf{q}|}{\|\mathbf{q}_{\perp}\| + \epsilon}$$

18 where \mathbf{q}_{\perp} is the projection of q onto the orthogonal complement of range \mathbf{Q} .

19 A.3.1 Property a

20 **Proposition A.1.** For $\epsilon \rightarrow 0$, $[\mathbf{q}_1, \dots, \mathbf{q}_{t-1}, \mathbf{q}^*]$ is the best approximation to $[\mathbf{q}_1, \dots, \mathbf{q}_{t-1}, \mathbf{g}]$.

Proof. If $\epsilon \rightarrow 0$, then $\text{obj}_{\text{ogb}}(q) \rightarrow \frac{|g^T q|}{\|q_{\perp}\|}$. If \mathbf{q}^* is a maximizer of obj_{ogb} , then as shown in Lemma
4.1, \mathbf{q}^* minimises the minimum distance from all

$$\mathbf{f} \in \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_{t-1}, \mathbf{g}\}$$

to the subspace of

$$\text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_{t-1}, \mathbf{q}^*\}.$$

21 Therefore, the subspace spanned by $[\mathbf{q}_1, \dots, \mathbf{q}_{t-1}, \mathbf{q}^*]$ is the best approximation to the subspace
22 spanned by $[\mathbf{q}_1, \dots, \mathbf{q}_{t-1}, \mathbf{g}]$. \square

23 A.3.2 Property b

24 **Proposition A.2.** For $\epsilon \rightarrow \infty$, \mathbf{q}^* is also a maximizer of obj_{gs} and any maximizer of obj_{gs} is also a
25 maximizer of obj_{ogb} .

26 *Proof.* Let q_1 and q_2 be any two queries and denote by $\text{obj}_{\text{ogb}}^{(\epsilon)}(q)$ the obj_{ogb} -value of q for a specific
27 ϵ . Then

$$\begin{aligned} & \lim_{\epsilon \rightarrow \infty} \epsilon \left(\text{obj}_{\text{ogb}}^{(\epsilon)}(q_1) - \text{obj}_{\text{ogb}}^{(\epsilon)}(q_2) \right) \\ &= \lim_{\epsilon \rightarrow \infty} \epsilon \left(\frac{|g^T \mathbf{q}_1|}{\|\mathbf{q}_1^{\perp}\| + \epsilon} - \frac{|g^T \mathbf{q}_2|}{\|\mathbf{q}_2^{\perp}\| + \epsilon} \right) \\ &= \lim_{\epsilon \rightarrow \infty} \left(\frac{|g^T \mathbf{q}_1|}{\|\mathbf{q}_1^{\perp}\|/\epsilon + 1} - \frac{|g^T \mathbf{q}_2|}{\|\mathbf{q}_2^{\perp}\|/\epsilon + 1} \right) \\ &= |g^T \mathbf{q}_1| - |g^T \mathbf{q}_2| \\ &= \text{obj}_{\text{gs}}(q_1) - \text{obj}_{\text{gs}}(q_2) \end{aligned}$$

28 Thus for large enough ϵ , the signs of $\text{obj}_{\text{ogb}}^{(\epsilon)}(q_1) - \text{obj}_{\text{ogb}}^{(\epsilon)}(q_2)$ and $\text{obj}_{\text{gs}}(q_1) - \text{obj}_{\text{gs}}(q_2)$ agree.
29 Therefore, a query q is a obj_{gs} -maximizer, i.e., $\text{obj}_{\text{gs}}(q) \geq \text{obj}_{\text{gs}}(q')$ for all $q' \in \mathcal{Q}$, if and only if q
30 is a obj_{ogb} -maximizer, i.e., $\text{obj}_{\text{ogb}}(q) \geq \text{obj}_{\text{ogb}}(q')$ for all $q' \in \mathcal{Q}$. \square

31 A.3.3 Property c

32 **Proposition A.3.** For $\epsilon = 0$ and $\|\mathbf{q}_{\perp}\| > 0$, the ratio $(\frac{\text{obj}_{\text{ogb}}(q)}{\text{obj}_{\text{gb}}(q)})^2$ is equal to $1 + (\frac{\|\mathbf{q}_{\parallel}\|}{\|\mathbf{q}_{\perp}\|})^2$.

33 *Proof.* If $\epsilon = 0$ and $\|\mathbf{q}_\perp\| > 0$, then

$$\begin{aligned} \left(\frac{\text{obj}_{\text{ogb}}(q)}{\text{obj}_{\text{gb}}(q)} \right)^2 &= \frac{\frac{|\mathbf{g}^T \mathbf{q}|^2}{\|\mathbf{q}_\perp\|^2}}{\frac{|\mathbf{g}^T \mathbf{q}|^2}{\|\mathbf{q}\|^2}} = \frac{\|\mathbf{q}\|^2}{\|\mathbf{q}_\perp\|^2} \\ &= \frac{\|\mathbf{q}_\parallel\|^2 + \|\mathbf{q}_\perp\|^2}{\|\mathbf{q}_\perp\|^2} \\ &= 1 + \left(\frac{\|\mathbf{q}_\parallel\|}{\|\mathbf{q}_\perp\|} \right)^2 \end{aligned}$$

34

□

35 **A.3.4 Property d**

36 **Proposition A.4.** The objective value $\text{obj}_{\text{ogb}}(q)$ is upper bounded by $\|\mathbf{g}\|$.

37 *Proof.* If we divide the numerator and denominator of $\text{obj}_{\text{ogb}}(\mathbf{q})$ with $\|\mathbf{q}_\perp\|$, then we can get

$$\begin{aligned} \text{obj}_{\text{ogb}}(\mathbf{q}) &= \frac{|\mathbf{g}^T \mathbf{q}|}{\|\mathbf{q}_\perp\| + \epsilon} \\ &= \frac{\frac{|\mathbf{g}^T \mathbf{q}_\perp|}{\|\mathbf{q}_\perp\|}}{1 + \frac{\epsilon}{\|\mathbf{q}_\perp\|}} \end{aligned}$$

38 according to the Cauchy–Schwarz inequality, $\frac{|\mathbf{g}^T \mathbf{q}|}{\|\mathbf{q}_\perp\|} \leq \frac{\|\mathbf{g}\| \|\mathbf{q}_\perp\|}{\|\mathbf{q}_\perp\|} = \|\mathbf{g}\|$, so,

$$\text{obj}_{\text{ogb}}(\mathbf{q}) \leq \frac{\|\mathbf{g}\|}{1 + \frac{\epsilon}{\|\mathbf{q}_\perp\|}}$$

39 as $\|\mathbf{q}_\perp\|$ is upper bounded by the number of data points n ,

$$\begin{aligned} \text{obj}_{\text{ogb}}(\mathbf{q}) &\leq \frac{\|\mathbf{g}\|}{1 + \frac{\epsilon}{n}} \\ \text{obj}_{\text{ogb}}(\mathbf{q}) &\leq \|\mathbf{g}\|. \end{aligned}$$

40

□

41 **A.4 Proof of Theorem 4.4**

42 *Proof.* To see the claim, we first rewrite the objective value for the i -th prefix as

$$\frac{\mathbf{g}^T \mathbf{q}^{(i)}}{\|\mathbf{q}_\perp^{(i)}\| + \epsilon} = \frac{\mathbf{g}^T \mathbf{q}^{(i)}}{\|\mathbf{q}^{(i)}\| - \|\mathbf{q}_\parallel^{(i)}\| + \epsilon}.$$

43 The value of $\|\mathbf{q}^{(i)}\|$ is trivially given as \sqrt{i} , and $\mathbf{g}^T \mathbf{q}^{(i)}$ can be easily computed for all $i \in [l]$ in time
44 $O(n)$ via cumulative summation. Finally we can reduce the problem of computing the (squared)
45 norms of the l projected prefixes to computing the t (squared) norms of the prefixes on the subspaces
46 given by the individual orthonormal basis vectors via

$$\|\mathbf{q}_\parallel^{(i)}\|^2 = \left\| \sum_{k=1}^t \mathbf{o}_k \mathbf{o}_k^T \mathbf{q}^{(i)} \right\|^2 = \sum_{k=1}^t \|\mathbf{o}_k \mathbf{o}_k^T \mathbf{q}^{(i)}\|^2.$$

Each of these t sequences of (squared) norms can be computed in time $O(n)$ by rewriting

$$\begin{aligned}\|\mathbf{o}_k \mathbf{o}_k^T \mathbf{q}^{(i)}\| &= \left\| \mathbf{o}_k \mathbf{o}_k^T \left(\sum_{j=1}^i \mathbf{e}_{\sigma(j)} \right) \right\| \\ &= \|\mathbf{o}_k\| \sum_{j=1}^i \mathbf{o}_k^T \mathbf{e}_{\sigma(j)} \\ &= \sum_{j=1}^i o_{k, \sigma(j)}\end{aligned}$$

where the last equality shows how an $O(n)$ -computation is achieved via cumulative summation of the k -th basis vector elements in the order given by σ . \square

B Greedy approximation to bounding function

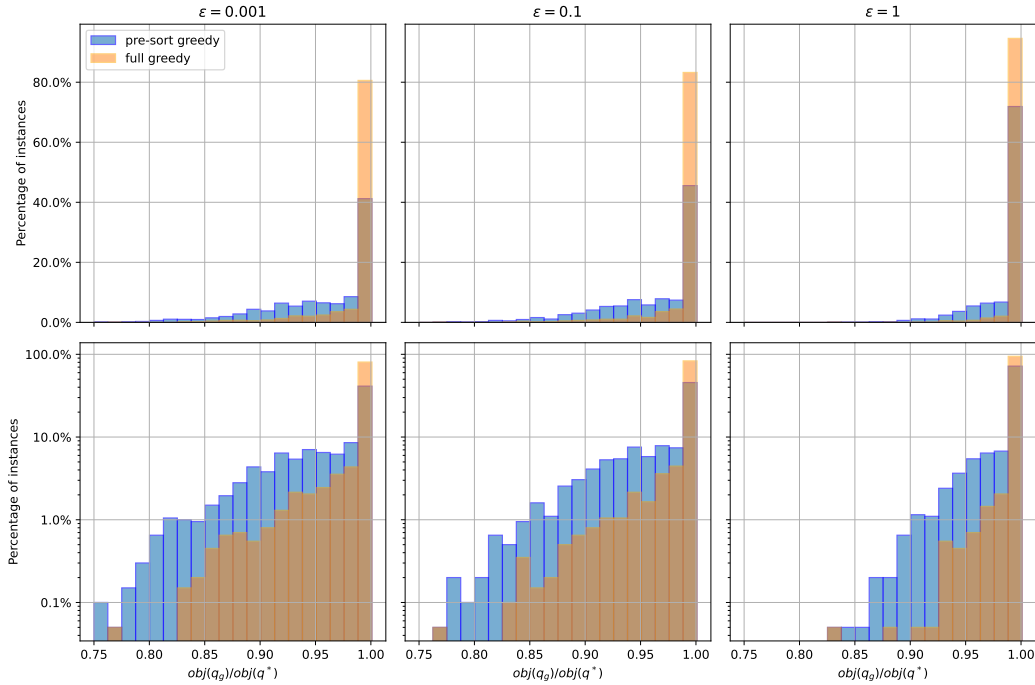


Figure 6: The number of instances of ratios between the best objective values obtained from the greedy search and the true optimal objective value. The upper figures are in linear scales and the lower figures are in log scales. The total variation distances for these three values of ϵ are 0.394, 0.377 and 0.227.

The branch-and-bound search described in Section 3.3 requires an efficient way of calculating the value of $\text{bnd}(\mathbf{q}) = \max\{\text{obj}(\mathbf{q}') : \mathbf{q}' \leq \mathbf{q}, \mathbf{q}' \in \{0, 1\}^n\}$. It is too expensive to enumerate all possible \mathbf{q}' s as there are 2^n cases in the worst case. One solution to this problem is that we can relax the constraint $\mathbf{q}' \in \{0, 1\}^n$ to $\mathbf{q}' \in [0, 1]^n$ and it can be solved by quadratic programming. However, this relaxation is too loose and inefficient. Instead, we consider relaxing the admission constraint and solve the problem using greedy algorithms.

A full greedy approach can be used to approximate the maximum objective value of the subset of data points selected by \mathbf{q} , which is the bounding value $\text{bnd}(\mathbf{q})$. Given a query $\mathbf{q}^{(t-1)} \leq \mathbf{q}$, we need to find the data point selected by \mathbf{q} which maximise the objective function, and use it with $\mathbf{q}^{(t-1)}$ to

60 form a $\mathbf{q}'(t)$.

$$i_*^{(t)} = \arg \max_{i \in I(\mathbf{q}) - I(\mathbf{q}'^{(t-1)})} \frac{\mathbf{g}^T(\mathbf{q}'^{(t-1)} + \mathbf{e}_i)}{\|(\mathbf{q}'^{(t-1)} + \mathbf{e}_i)_\perp\| + \epsilon}.$$

61 where $I(\mathbf{q}) = \{i : \mathbf{q}(x_i) = 1, 1 \leq i \leq n\}$, $0 \leq t \leq |I(\mathbf{q})|$, $\mathbf{q}'^{(0)} = \mathbf{0}$ and $\mathbf{q}'^{(t)} = \mathbf{q}'^{(t-1)} + \mathbf{e}_{i_*^{(t)}}$.

62 We use the maximum value of $\text{obj}(\mathbf{q}'^{(t)})$ as the bounding value for query \mathbf{q} . The computation time
63 complexity level of this approach is $O(n^2)$ for each query, which is not as efficient as the presorting
64 greedy approach described in Section 4.3.

65 The presorting greedy approach of solving the prefix optimization problem described in Section 4.3
66 leads to another approximation to the optimal objective function value for the queries which cover
67 subsets of data points covered by \mathbf{q} . As proved in Theorem 4.4, this approach has a time complexity
68 of $O(tn)$.

69 We test 2000 instances for different initial queries and initial gradient values to see the difference
70 between the approximation of $\text{bnd}(\mathbf{q})$ obtained by the full greedy approach, the pre-sorting greedy
71 approach, and the actual optimal objective values (obtained by a brute-force approach). We choose
72 three different values of ϵ : 0.001, 0.1 and 1.

73 Figure 6 compares the ratio between the approximations to $\text{bnd}(\mathbf{q})$ obtained by the two greedy
74 approaches and the true optimal objective value. The Y axis of Figure 6 represents the percentage of
75 instances of different ratios.

76 According to the comparison, the full greedy approach can approximate the true bounding function
77 better than the presorting greedy approach. For smaller ϵ values ($\epsilon = 0.001$), there are 90% instances
78 whose approximation values are more than 90% of the true bounding function values, while 96% of
79 instances approximate more than 90% of the value of $\text{bnd}(\mathbf{q})$ using the full greedy approach. For
80 $\epsilon = 0.1$, both algorithms have slight better (both 1% promotion) approximation than $\epsilon = 0.001$.
81 It can be observed that for $\epsilon = 1$, both algorithms have more instances where the approximations
82 are closed to the true bounding values. However, if the value of ϵ is too large, then the calculated
83 objective values cannot be accurate according to Theorem 4.3. Comparing the statistical distances of
84 these two greedy approaches, it is reasonable to use the presorting greedy approach to approximate
85 the bounding values.

86 To approximate the true bounding function more efficiently and more accurate, we adopt the presorting
87 greedy approach in this research.

88 C Experiments configurations

89 The experiments are conducted on a computer with CPU 'Intel(R) Core(TM) i5-10300H CPU @
90 2.50GHz' and memory of 24G.