

# **HistGBODT: Gradient Boosting for Ordinal Classification**

by

Kemeng Zhang

A dissertation submitted to The Johns Hopkins University in conformity with the  
requirements for the degree of Bachelor of Science.

Baltimore, Maryland

May 14, 2020

© Kemeng Zhang 2020

All rights reserved

# Abstract

Tree-based models, including Gradient Boosting Trees and Random Forests, have recently demonstrated state-of-the-art performance in a variety of machine learning tasks. Most notably, popular implementations of Gradient Boosting such as **XGBoost**, **LightGBM**, and most recently, **CatBoost** have become the go-to machine learning solutions widely adopted in industry and academia by researchers. Gradient Boosting is compatible with a wide range of loss functions, thus can perform tasks such as: regression, binary/multi-classification, and learning-to-rank. Unfortunately, in the existing gradient boosting framework, there has not been any implementation that aims to address the situation where the target variable is discrete, but ordered. This class of problem is called “Ordinal Classification”. Traditionally, researchers ignore the ordered structure of the target variable and model it using multi-classification, or ignore the discrete nature of the target variable and model it using learning-to-rank. We introduce yet another extension to Gradient Boosting, called “Histogram-based Gradient Boosting Ordinal Decision Trees” (**HistGBODT**). We provide the theoretical foundation for learning discrete ordinal target variables and the algorithmic details

## ABSTRACT

for achieving such goals. To illustrate how **HistGBODT** addresses ordinal classification, we conduct extensive experiments, both on simulated and real-world datasets. **HistGBODT** typically yields improved performance over Gradient Boosting with multi-classification (cross-entropy) loss, while maintaining computational efficiency and interpretability. **HistGBODT** can be adopted when it is unclear when a learning task at hand should be approached by classification or learning-to-rank.

Primary Reader: Dr. John C. Miller

# Acknowledgments

I would like to thank my advisor, Dr. Miller, for his guidance during my undergraduate years at Johns Hopkins. I will never forget the exciting discussions we had both in academics and in life. I want to thank Dr. Fishkind and Dr. Athreya not only for their insightful comments and encouragement, but also for the thought-provoking questions which motivated me to broaden my research from various perspectives, in statistics, optimization, and matrix analysis. Last but not least I want to thank my family for supporting me obtaining my undergraduate degree, both financially and spiritually.

# Dedication

This thesis is dedicated to my parents.

# Contents

<b>Abstract</b>	<b>ii</b>
<b>Acknowledgments</b>	<b>iv</b>
<b>List of Tables</b>	<b>ix</b>
<b>List of Figures</b>	<b>x</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	2
1.2 Literature . . . . .	3
1.3 Main Contribution . . . . .	4
1.4 Thesis Structure . . . . .	4
<b>2 Ordinal Classification: Statistical Framework</b>	<b>6</b>
2.1 Problem Setting . . . . .	6
2.2 Ordinal Classification . . . . .	7

## CONTENTS

2.3	Loss Function Construction . . . . .	7
2.3.1	Binary Classification . . . . .	8
2.3.2	Extending Binary Classification - Threshold-Based Constructions	9
2.3.3	All-Threshold Loss . . . . .	10
<b>3</b>	<b>Gradient Boosted Trees</b>	<b>12</b>
3.1	Building Block: Decision Trees . . . . .	12
3.2	Regularized Learning Objective . . . . .	14
3.3	Gradient Tree Boosting . . . . .	15
3.3.1	Analysis of Hessian Matrix . . . . .	18
3.3.1.1	Inversion of Block Matrices . . . . .	18
3.3.1.2	Majorization Minimization: The Construction of In- vertible Hessian Matrix . . . . .	19
<b>4</b>	<b>Faster Training: Algorithmic Speed-up</b>	<b>27</b>
4.1	Weighted Quantile Sketch and Feature Quantization . . . . .	28
4.1.1	Approximate splitting with Quantile Sketch . . . . .	28
4.1.2	Feature Quantization . . . . .	29
4.1.3	Weighted Quantile . . . . .	30
4.2	Histogram Aggregation of Gradient Statistics . . . . .	31
4.3	Tree-growing Strategies . . . . .	31
4.4	Gradient-Aware Pruning and Adaptive Learning Rate . . . . .	35

## CONTENTS

4.4.1	Gradient-Aware Pruning Using Out-of-Bag Samples . . . . .	35
4.4.2	Adaptive Learning Rate Using Out-of-Bag Samples . . . . .	36
<b>5</b>	<b>Experiments</b>	<b>39</b>
5.1	Artificial Data . . . . .	39
5.2	Benchmark Data . . . . .	40
<b>6</b>	<b>Conclusion and Future Directions</b>	<b>45</b>
6.1	Conclusion . . . . .	45
6.2	Future Directions . . . . .	46
	<b>Bibliography</b>	<b>47</b>
	<b>Appendix A</b>	<b>53</b>
	<b>Appendix B</b>	<b>61</b>



# List of Tables

5.1	Datasets and their characteristics. “Attributes” state the number of numerical and nominal attributes. “Training Instances” and “Instances for Test” specify the size of training/test partitions. . . . .	42
5.2	Testing set results: the bold font indicates the cases in which the average value is the lowest in the results of minimizing the three loss functions. . . . .	43
5.3	Nonparametric Wilcoxon signed-rank test results . . . . .	44

# List of Figures

2.1	Binary Classification Loss Functions. . . . .	9
2.2	Ordinal Classification Loss Functions when $y_i = 3$ . . . . .	10
3.1	A visualization of decision tree. [1]. . . . .	13
3.2	Sharp Majorizing function at $x = 0$ . . . . .	23
4.1	We scan the training data from left to right, considering each value as a threshold candidate. [1] . . . . .	28
4.2	We restrict splitting candidates to only quantile points [1] . . . . .	29
4.3	Quantize numerical feature values into quantiles. [1] . . . . .	29
4.4	Aggregation of gradient statistics along one feature. Only one bar is displayed for simplicity. In reality, there are three bars, storing sum of gradients and Hessians w.r.t. $w$ (scalars) and second-order mixed gradient $\nabla_{w,\delta}\tilde{\mathcal{L}}$ (a vector) at each bar. . . . .	32
4.5	Splitting gradient statistics to left and right child. The sum of histograms of left and right child equals the histogram of parent node. . . . .	33
4.6	Loss-optimal tree growth strategy . . . . .	33
4.7	When a data point has missing values and they are compared against tree node thresholds, it will be split to the direction that gives rise to greater loss reduction. [2] . . . . .	34

# Chapter 1

## Introduction

Gradient Boosting Decision Trees (GBDT) is a highly effective machine learning method used in a variety of tasks. It has become the de-facto choice of method for problems with noisy and high-dimensional data such as weather forecasting [3, 4], recommendation systems [5], and ad click-through rate prediction [6]. A vast majority of winning solutions to machine learning challenges hosted by Kaggle incorporated Gradient Boosting in some fashion.

The reason why GBDT is gaining popularity in predictive modeling is that it can perform some of the most important tasks such as classification, regression, and ranking. GBDT is able to capture the complex data dependencies by natively incorporating feature interactions in the process of constructing decision trees. It produces an ensemble of weak learners using decision trees, which has many desirable properties as a statistical model. Currently, there are many fast and scalable implementations

of gradient boosting: `XGBoost`, `LightGBM`, and `CatBoost` [2, 7–9] that have found a wide following among practitioners.

## 1.1 Motivation

Gradient Boosting is compatible with a wide range of problems and objectives. In this paper, we consider using Gradient Boosting for the objective of Ordinal Classification. Our aim is to learn a rule to predict labels from a discrete ordered set. For example, users generate ratings for products on Amazon from several rating “levels”, e.g. one through five stars. Similar problems often arises when the target variable is human-generated ratings, severity of symptoms of diseases, or performance metrics.

Ordinal Classification is fundamentally different from traditional machine learning tasks such as classification, regression, and learning-to-rank. Similar to classification, our target variable consist of discrete values; as in regression, there is meaningful order between values; as in learning-to-rank, we aim to predict relative order between each instances. Different from classification, our label is discrete but ordered. We call this type of label “ordinal” as opposed to “nominal”. Ordinal classification formalizes the notion of order by incurring a greater penalty for predictions farther from the true label than those closer. As in regression, our target variable is not on a continuous scale. As in learning-to-rank, which also outputs on a continuous scale, fails to produce a discrete label for each instance. In this sense, it is possible for a ranking

## CHAPTER 1. INTRODUCTION

model (but not for ordinal classification) to predict the wrong label and incur no loss at all, as long as the relative order of the true labels and predicted labels are correct. Although ordinal classification and learning-to-rank are two different problems, the distinction is not always clear.

In the past, ordinal classification could be approached by all of the three tasks mentioned previously: multi-classification, regression, or learning-to-rank. However, none of these approaches reflect the specific structure of discrete ordered labels.

## 1.2 Literature

To date, the most popular approach to model discrete ordinal data is by Ordinal Logistic Regression. McCullagh (1980) introduced the Proportional Odds Model, which uses logits of cumulative probabilities, often called cumulative logits [10–12]. Recently, many algorithms for ordinal classification have been proposed from a machine learning perspective. Freund et al. [13] studied the ranking objective in collaborative filtering and proposes an AdaBoost based solution called **RankBoost**. Dembczyński et al. [14] further the boosting framework by providing a gradient descent algorithm for ordinal classification. A crucial aspect of the learning objective is the loss function. Rennie and Srebro [15] surveyed strategies for learning discrete ordered labels. They view ordinal classification as a generalization of binary classification. Thus, commonly used loss functions for binary classification can be easily

## CHAPTER 1. INTRODUCTION

extended to multi-category. Most notably, we adopt the “All-Threshold Loss” from threshold-based approach in Rennie and Srebro. Crammer and Singer [16] suggest a generalization of the Perceptron algorithm for discrete ordinal labels. Shashua and Levin [17] suggest a similar generalization to Support Vector Machine (SVM). Chu and Keerthi [18] proposes support vector approaches for ordinal classification, which optimize multiple thresholds to define parallel discriminant hyperplanes for the ordinal scales. This approach guarantees that the thresholds are properly ordered at the optimal solution. Finally, the theoretical properties of threshold-based loss function are discussed in Lin and Li [19,20], and Pedregosa et al. [21] characterizes the Fisher consistency of threshold-based loss functions in the context of ordinal classification.

### 1.3 Main Contribution

Our main contribution is to extend the powerful gradient boosting algorithm to learn discrete ordinal response variables using the All-Threshold Loss function. We provide an efficient algorithm for optimizing our objective using Majorization-Minimization (MM) with provable convergence guarantees.

### 1.4 Thesis Structure

In Chapter 2, I will introduce the statistical framework of Ordinal Classification and the All-Threshold loss function. In Chapter 3, I will review the Gradient Boosting

## CHAPTER 1. INTRODUCTION

learning objective in the traditional setting, and extend the objective to All-Threshold loss. In Chapter 4, I will describe the algorithms that speed up the learning speed during the growth of each decision tree in detail. Finally, Chapter 5 will provide numerical experimentation on simulated and real-world datasets, comparing the performance of multi-classification and ordinal classification.

## Chapter 2

# Ordinal Classification: Statistical Framework

### 2.1 Problem Setting

We present the general framework learning discrete ordinal responses. For a given data set with  $n$  examples and  $m$  features  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  ( $|\mathcal{D}| = n, \mathbf{x}_i \in \mathbb{R}^m, y_i \in \{r_1, \dots, r_K\}$ ), it is assumed our response variable is from a set of  $K$  discrete ordered class labels  $r_1 \prec \dots \prec r_{K-1} \prec r_K$ , where  $\prec$  denotes the ordering relation between the labels. Let us assume, without loss of generality, that  $r_k = k - 1$ .



## 2.2 Ordinal Classification

In ordinal classification, label can take one of several discrete, but ordered labels. For instance, movie ratings can go from zero to ten stars. Unlike multi-class classification or learning-to-rank, not all errors are equally bad, e.g., it is worse to mis-classify a 1-star with a 10-star than a 4-star movie with a 5-star one. Hence, we construct a loss matrix  $L \in \mathbb{R}^{K \times K}$  for ordinal classification to satisfy the following properties:

1.  $L(y, \hat{y}) = L_{ij}$ , where  $i = y, j = \hat{y} = F(\mathbf{x}) \in \{0, \dots, K-1\}$ .
2.  $L_{ii} = 0, \forall i$
3.  $L_{ij} \leq L_{ik}, \forall i < j < k$ , or  $k < j < i$ .

In other words, we incur no loss when the predicted label is the same as ground truth, and a larger loss when the absolute distance between the predicted label and the ground truth is larger.

## 2.3 Loss Function Construction

Note that in traditional multi-class classification with  $K$  classes, the problem is often solved as a collection of two class sub-problems. Friedman [22], and Hastie & Tibshirani [23] considers  $\binom{K}{2}$  pairwise subproblems, often referred to as the “One-vs-One” (OvO) strategy. Another strategy, often referred to as “One-vs-Rest” (OvR), considers building one classifier per class. For each classifier, the class is fitted against

all the other classes, hence  $K$  classifiers in total [24]. In Rennie and Srebro [15], the authors consider the strategy for building classifiers for  $K$  discrete ordered labels. They also view ordinal classification as a generalization of binary classification, and adopted the “One-vs-Rest” strategy by exploiting the ordering structure in the label when constructing the loss function. They extend the binary loss functions to the case of ordinal classification by introducing  $K - 1$  thresholds. We first review common loss functions used with binary responses, where  $y \in \pm 1$ . Then we extend this to more general loss functions for discrete ordinal labels.

### 2.3.1 Binary Classification

Consider the loss function over training observations:

$$L(\mathbf{y}, \tilde{\mathbf{y}}) = \frac{1}{n} \sum_{i=1}^n l(y_i, \tilde{y}_i) \quad (2.1)$$

where  $\tilde{y}_i = f(\mathbf{x}_i) \in \mathbb{R}$ ;  $\hat{y}_i = \text{sgn}(\tilde{y}_i) \in \pm 1$ . The loss function  $l$  should have several desired properties. It should be close to zero when  $\tilde{y}_i$  and  $y_i$  agree in sign and have non-negative value when they have opposite signs. These common choices are:

- Zero-one loss,  $I(\text{sgn}(y_i \tilde{y}_i) \neq 1)$ , where  $I$  is the indicator function;
- Hinge loss,  $\max(0, 1 - \tilde{y}_i y_i)$
- Logistic loss,  $\log(1 + \exp(-\tilde{y}_i y_i))$

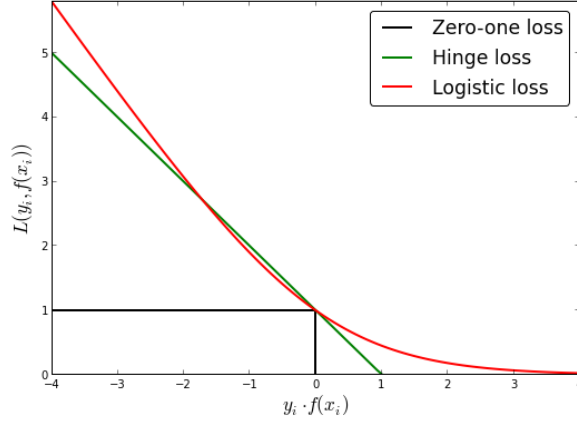


Figure 2.1: Binary Classification Loss Functions.

### 2.3.2 Extending Binary Classification - Threshold-Based Constructions

In the binary case, we use a single threshold, namely 0, to separate the real line into two distinct segments:  $\tilde{y}_i = f(\mathbf{x}_i) > 0$  corresponds to  $\hat{y}_i = +1$  predicted label, while  $\hat{y}_i = -1$  otherwise. To extend binary loss to the case of discrete ordinal regression, [15] introduce  $K - 1$  thresholds  $\theta_0 < \theta_1 < \dots < \theta_{K-2}$  partitioning the real line into  $K$  segments. We denote this vector of threshold as  $\boldsymbol{\theta} \in \mathbb{R}^{K-1}$ . Each of the  $K$  segment corresponds to one of the  $K$  labels. The predicted class label  $\hat{y}_i$  can thus

be derived from the predicted score  $\tilde{y}_i$  and thresholds vector  $\boldsymbol{\theta}$ :

$$\hat{y}_i = \begin{cases} 0 & \text{if } \tilde{y}_i \in (-\infty, \theta_0) \\ k & \text{if } \tilde{y}_i \in [\theta_{k-1}, \theta_k) \\ K-1 & \text{if } \tilde{y}_i \in [\theta_{K-2}, \infty) \end{cases} \quad (2.2)$$

### 2.3.3 All-Threshold Loss

The idea behind [15] all-threshold loss is to sum up all threshold violation penalties. Here, not only the function increases when thresholds are violated, but the slope of the loss function increases each time a threshold is crossed.

$$l_{ord}(y_i, \tilde{y}_i, \boldsymbol{\theta}) = \sum_{k=0}^{y_i-1} l(1, \tilde{y}_i - \theta_k) + \sum_{k=y_i}^{K-2} l(-1, \tilde{y}_i - \theta_k) \quad (2.3)$$

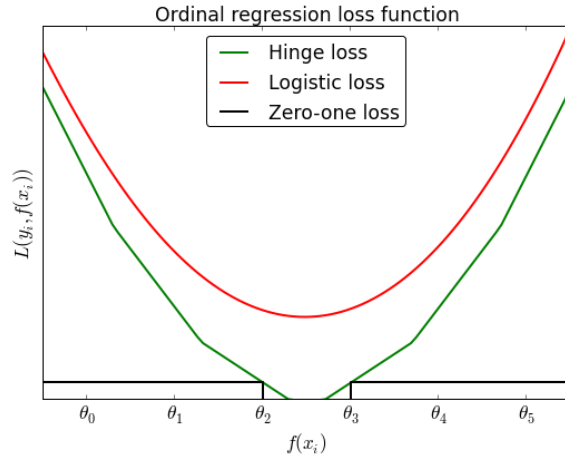


Figure 2.2: Ordinal Classification Loss Functions when  $y_i = 3$ .

## CHAPTER 2. ORDINAL CLASSIFICATION: STATISTICAL FRAMEWORK

In Figure 2.2, consider the case  $K = 7, y_i = 3$ . The x-axis denotes  $\tilde{y}_i = f(\mathbf{x}_i)$ , and the y-axis denotes the loss incurred by  $\tilde{y}_i$ . All-threshold loss checks if  $\tilde{y}_i > \theta_k$  for  $k = 0, 1, 2$ ; if any threshold is violated, a loss is incurred  $l(1, \tilde{y}_i - \theta_k)$ , where  $l$  is the binary classification loss mentioned previously. Similarly, it also checks if  $\tilde{y}_i < \theta_k$  for  $k = 3, 4, 5$ ; a loss  $l(-1, \tilde{y}_i - \theta_k)$  is incurred if any threshold value is violated. We can make our expression more compact by defining the following binary function:

$$s(x) = \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{if } x \leq 0 \end{cases} \quad (2.4)$$

$$l_{ord}(y_i, \tilde{y}_i, \boldsymbol{\theta}) = \sum_{k=0}^{K-2} l(s(y_i - k), \tilde{y}_i - \theta_k) \quad (2.5)$$

A problem arises if in the learned threshold  $\boldsymbol{\theta}$ , there exists  $k$  such that  $\theta_k > \theta_{k+1}$ . In Li et. al [20] Theorem 2, the authors proved that under this framework as long as  $l$  is convex there exists optimal solution  $(\tilde{\mathbf{y}}^*, \boldsymbol{\theta}^*)$  such that  $\boldsymbol{\theta}^*$  is ordered.

### Probabilistic Interpretation of the All-Threshold Loss

Define

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\psi(x) = \log(1 + e^{-x}) = -\log(\sigma(x)) \text{ (Negative log-likelihood)}$$

$$\mathbb{P}(y_i \leq k) = \sigma(s(y_i - k)(\tilde{y}_i - \theta_k)) = \sigma(\theta_k - \tilde{y}_i)$$

# Chapter 3

## Gradient Boosted Trees

We review the gradient tree boosting algorithms [2,25]. Furthermore, we introduce learning rules under the context of the statistical framework of threshold-based ordinal classification.

### 3.1 Building Block: Decision Trees

The most popular base learner in gradient boosting is the decision tree, also known as CART [26]. It is represented as a directed acyclic graph data structure (See Fig. 3.1). Given a training instance  $\mathbf{x}_i \in \mathbb{R}^m$  to be predicted, we traverse the tree from top to bottom, performing a series of boolean operations stored in each intermediate node, comparing a threshold value against a feature in  $\mathbf{x}_i$  in the process. We proceed to either the left child or the right child depending on the result of the comparison.

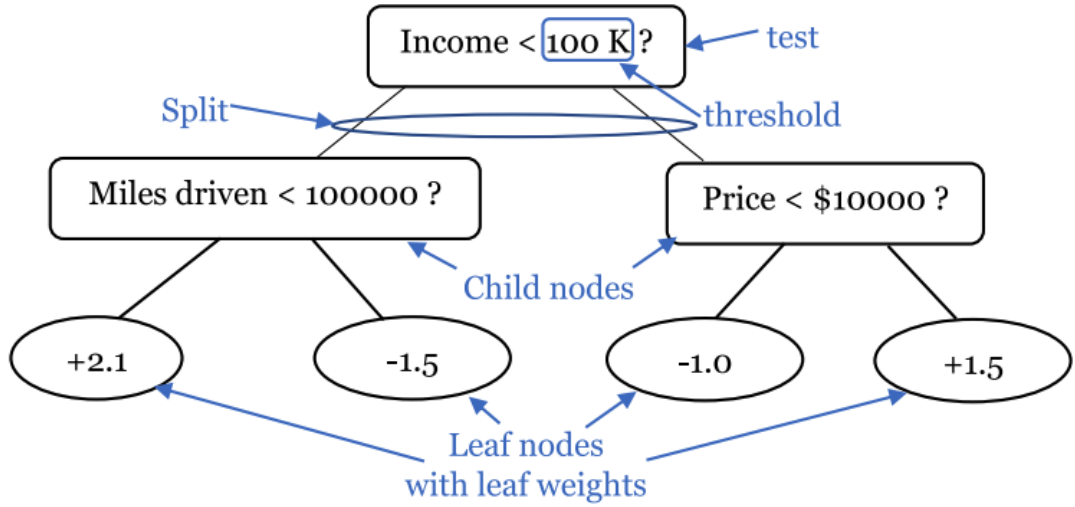


Figure 3.1: A visualization of decision tree. [1].

The intermediate node with a feature and a threshold value represents a “split” on the dataset. We repeat the splitting process until we reach the terminal leaf region, where the node has no left or right child, and stores an output value called the “leaf weight”. Each leaf node gives rise to an instance set  $I_j$ , the set of all data points for which traversal ended at that leaf node [1].

Formally, the predicted output of a decision tree can be expressed as:

$$\tilde{y}_i = \sum_{j=1}^J w_j \mathbf{1}_{\{\mathbf{x}_i \in I_j\}} \quad (3.1)$$

The strategy for growing a decision tree will be discussed in Chapter 4 where the algorithmic aspect of HistGBODT will be explained in detail.

## 3.2 Regularized Learning Objective

A tree ensemble model uses  $T$  additive functions to predict output.

$$\tilde{y}_i = \sum_{t=1}^T f_t(\mathbf{x}_i), f_t \in \mathcal{F}, \quad (3.2)$$

where  $\mathcal{F} = \{f(\mathbf{x}) = w_{q(\mathbf{x})}\}$  is the space of regression trees (also known as CART). Here  $q : \mathbb{R}^m \rightarrow \{1, 2, \dots, J\}$  is the mapping from the feature space to the terminal leaf assignment, i.e. the tree structure; define instance set  $I_j = \{i | q(\mathbf{x}_i) = j\}$  as the set of examples that get assigned to leaf  $j$ ;  $f_t(\mathbf{x}_i) = w_j$  if  $q(\mathbf{x}_i) = j$ , i.e.  $i \in I_j$ ;  $J$  is the number of leaves in a tree;  $w_j \in \mathbb{R}$  is the predicted output on the  $j$ -th leaf. For a given example, we will use decision rules  $q$  in each tree to assign each example into the leaves and calculate the final prediction by summing up the scores (given by  $w$ ) in the corresponding leaves. To learn the set of functions in the ensemble model, we minimize the following *regularized* objective:

$$\mathcal{L} = \sum_{i=1}^n l_{ord}(y_i, \tilde{y}_i, \boldsymbol{\theta}) + \sum_{t=1}^T \Omega(f_t) \quad (3.3)$$

$$\text{where } \Omega(f) = \gamma J + \frac{1}{2} \lambda \|w\|^2$$

Here the second term  $\Omega$  penalizes the complexity of the regression trees. The additional term regularizes the model to smooth out the final prediction in order to avoid over-fitting. Intuitively, the regularized objective tends to select simpler and predictive regression trees.



### 3.3 Gradient Tree Boosting

The tree ensemble model in Eq. (3.3) includes functions as parameters and cannot be optimized directly. Instead, the model is trained in an additive manner. Formally, let  $\tilde{y}_i^{(t)}$  be the prediction of the  $i$ -th example at the  $t$ -th iteration, we need to add  $f_t$ , and  $\boldsymbol{\delta}^{(t)}$  to minimize the following objective:

$$\mathcal{L}^{(t)} = \sum_{i=1}^n l_{ord}(y_i, \tilde{y}_i^{(t-1)} + f_t(\mathbf{x}_i), \boldsymbol{\theta}^{(t-1)} + \boldsymbol{\delta}^{(t)}) + \Omega(f_t) \quad (3.4)$$

where  $\tilde{y}_i^{(t-1)}$  is the combined output of the previous  $t - 1$  trees with the  $i$ -th data point as input.

We greedily add new regression tree  $f_t$  and threshold improvement  $\boldsymbol{\delta}^{(t)}$  that most improve our model according to Eq. (3.3). To quickly optimize the objective using second-order approximation, we choose the logistic loss in Sec. 2.3.1 that is convex and twice differentiable.

We initialize at  $t = 0$ :

$$\tilde{y}^{(0)}, \boldsymbol{\theta}^{(0)} = \underset{\tilde{y}, \boldsymbol{\theta}}{\operatorname{argmin}} \sum_{i=1}^n l_{ord}(y_i, \tilde{y}, \boldsymbol{\theta}) \quad (3.5)$$

$$\tilde{y}_i^{(0)} = \tilde{y}^{(0)} \quad \forall i = \{1, \dots, n\}$$

$$\begin{aligned} \mathcal{L}^{(t)} \simeq & \sum_{i=1}^n \left[ l_{ord}(y_i, \tilde{y}_i^{(t-1)}, \boldsymbol{\theta}^{(t-1)}) + g_i f_t(\mathbf{x}_i) + \frac{1}{2} h_i f_t^2(\mathbf{x}_i) \right] + \\ & \left( \sum_{i=1}^n \frac{\partial l_i}{\partial \boldsymbol{\theta}} \right)^\top \boldsymbol{\delta}^{(t)} + \frac{1}{2} \boldsymbol{\delta}^{(t)\top} \left( \sum_{i=1}^n \frac{\partial^2 l_i}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right) \boldsymbol{\delta}^{(t)} + \left( \sum_{i=1}^n f_t(\mathbf{x}_i) \frac{\partial^2 l_i}{\partial \boldsymbol{\theta} \partial \tilde{y}} \right)^\top \boldsymbol{\delta}^{(t)} + \Omega(f_t) \end{aligned} \quad (3.6)$$

### CHAPTER 3. GRADIENT BOOSTED TREES

Exploiting the tree structure  $q_t$  of  $f_t$ , we can re-write  $\mathcal{L}^{(t)}$  in terms of instance sets  $I_j$ .

We drop all constant term to get a simplified objective:

$$\begin{aligned} \tilde{\mathcal{L}}^{(t)} = & \sum_{j=1}^J \left[ \left( \sum_{i \in I_j} \left( g_i + \frac{\partial^2 l_i}{\partial \boldsymbol{\theta} \partial \tilde{y}} \boldsymbol{\delta}^{(t)\top} \right) w_j + \frac{1}{2} \left( \sum_{i \in I_j} h_i + \lambda \right) w_j^2 \right) + \right. \\ & \left. \left( \sum_{i=1}^n \frac{\partial l_i}{\partial \boldsymbol{\theta}} \right)^\top \boldsymbol{\delta}^{(t)} + \frac{1}{2} \boldsymbol{\delta}^{(t)\top} \left( \sum_{i=1}^n \frac{\partial^2 l_i}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right) \boldsymbol{\delta}^{(t)} + \gamma J \right] \end{aligned} \quad (3.7)$$

where

$$\begin{aligned} l_i &= l_{ord}(y_i, \tilde{y}, \boldsymbol{\theta})|_{\tilde{y}=\tilde{y}_i^{(t-1)}, \boldsymbol{\theta}=\boldsymbol{\theta}^{(t-1)}} \\ g_i &= \frac{\partial}{\partial \tilde{y}} l_{ord}(y_i, \tilde{y}, \boldsymbol{\theta})|_{\tilde{y}=\tilde{y}_i^{(t-1)}, \boldsymbol{\theta}=\boldsymbol{\theta}^{(t-1)}} \\ h_i &= \frac{\partial^2}{\partial \tilde{y}^2} l_{ord}(y_i, \tilde{y}, \boldsymbol{\theta})|_{\tilde{y}=\tilde{y}_i^{(t-1)}, \boldsymbol{\theta}=\boldsymbol{\theta}^{(t-1)}} \end{aligned}$$

Note that our simplified loss function could also be rewritten in matrix form for simplicity. We define

Trainable parameters (descent direction in a Newton-Raphson step):

$$\boldsymbol{\beta} = \begin{bmatrix} w_1 & \dots & w_J & \delta_0 & \dots & \delta_{K-2} \end{bmatrix}^\top$$

Full gradients:

$$\boldsymbol{g} = \begin{bmatrix} G_1 & \dots & G_J & \boldsymbol{u}^\top \end{bmatrix}^\top$$

### CHAPTER 3. GRADIENT BOOSTED TREES

Hessian matrix:

$$\mathbf{H} = \begin{bmatrix} H_1 + \lambda & & 0 & -\mathbf{L}_1^\top & \\ & \ddots & & \vdots & \\ 0 & & H_J + \lambda & -\mathbf{L}_J^\top & \\ | & & | & & \\ \mathbf{L}_1 & \dots & \mathbf{L}_J & \mathbf{V} & \\ | & & | & & \end{bmatrix} \quad (3.8)$$

where  $\tilde{\mathcal{L}}^{(t)}$  can be re-expressed in matrix form as:

$$\tilde{\mathcal{L}}^{(t)} = \mathbf{g}^\top \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\beta}^\top \mathbf{H} \boldsymbol{\beta} + \gamma J \quad (3.9)$$

here

$$G_j = \sum_{i \in I_j} g_i \quad \mathbf{u} = \sum_{i=1}^n \frac{\partial l_i}{\partial \boldsymbol{\theta}} \quad H_j = \sum_{i \in I_j} h_i$$

$$\mathbf{L}_j = \sum_{i \in I_j} \frac{\partial^2 l_i}{\partial \boldsymbol{\theta} \partial \tilde{\mathbf{y}}} \quad \mathbf{V} = \sum_{i=1}^n \frac{\partial^2 l_i}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}$$

See Appendix B for full expressions of the gradient statistics above. We can easily solve for the optimal parameter from the quadratic program above:

$$\boldsymbol{\beta}^* = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \tilde{\mathcal{L}}^{(t)} = -\mathbf{H}^{-1} \mathbf{g} \quad (3.10)$$

## CHAPTER 3. GRADIENT BOOSTED TREES

The corresponding value for  $\tilde{\mathcal{L}}^{(t)}$  is

$$\tilde{\mathcal{L}}^{(t)*} = -\frac{1}{2}\mathbf{g}^\top \mathbf{H}^{-1} \mathbf{g} + \gamma J \quad (3.11)$$

### 3.3.1 Analysis of Hessian Matrix

Computing the optimal parameter  $\beta^*$  requires inverting the Hessian matrix  $\mathbf{H}$ , which runs in  $\mathcal{O}((J+K)^3)$  time. Moreover, note when the  $\ell_2$  regularization parameter  $\lambda = 0$ , the rows and columns of  $\mathbf{H}$  sum to 0, which means it has an eigenvalue 0 and associated eigenvector  $\mathbf{1}$ . The Hessian matrix is thus singular and we cannot apply Newton-Raphson method directly to find the minimum of our objective.

It turns out we could both get the inversion time complexity down to  $\mathcal{O}(K^3)$ , and construct a “majorizing” Hessian matrix that is positive definite regardless of the value of  $\ell_2$  regularization parameter  $\lambda$ .

#### 3.3.1.1 Inversion of Block Matrices

The special structure of the Hessian matrix can be exploited. Note that  $\mathbf{H}$  is a  $2 \times 2$  block matrix, with two diagonal matrices along its diagonal entries. Inverting diagonal matrices is trivial. Inverse formulae also exists for  $2 \times 2$  block matrices [27]:

$$\mathbf{H}^{-1} = \begin{bmatrix} \mathbf{H}_\lambda^{-1} + \mathbf{H}_\lambda^{-1} \mathbf{L}^\top (\mathbf{V} - \mathbf{L} \mathbf{H}_\lambda^{-1} \mathbf{L}^\top)^{-1} \mathbf{L} \mathbf{H}_\lambda^{-1} & -\mathbf{H}_\lambda^{-1} \mathbf{L}^\top (\mathbf{V} - \mathbf{L} \mathbf{H}_\lambda^{-1} \mathbf{L}^\top)^{-1} \\ -(\mathbf{V} - \mathbf{L} \mathbf{H}_\lambda^{-1} \mathbf{L}^\top)^{-1} \mathbf{L} \mathbf{H}_\lambda^{-1} & (\mathbf{V} - \mathbf{L} \mathbf{H}_\lambda^{-1} \mathbf{L}^\top)^{-1} \end{bmatrix} \quad (3.12)$$

where

$$\mathbf{H}_\lambda = \begin{bmatrix} H_1 + \lambda & & 0 \\ & \ddots & \\ 0 & & H_J + \lambda \end{bmatrix}$$

$\mathbf{H}_\lambda$  is nonsingular, and Hessian  $\mathbf{H}$  is invertible whenever the Schur complement  $\mathbf{V} - \mathbf{L}\mathbf{H}_\lambda^{-1}\mathbf{L}^\top$  of  $\mathbf{H}_\lambda$  is also invertible. In other words, in order to compute the inverse of the Hessian matrix we only need to invert  $\mathbf{V} - \mathbf{L}\mathbf{H}_\lambda^{-1}\mathbf{L}^\top$ , in addition to matrix multiplication. Hence, we are able to reduce the computational cost such that the inversion does not scale with decision tree growth, as the number of leaves grow larger and larger.

For clarity, we can now re-express the optimal value for each parameter:

$$\begin{aligned} \boldsymbol{\delta}^{(t)*} &= -\left(\mathbf{V} - \mathbf{L}\mathbf{H}_\lambda^{-1}\mathbf{L}^\top\right)^{-1}\left(\mathbf{u} - \mathbf{L}\mathbf{H}_\lambda^{-1}\mathbf{g}\right) \\ &= -\left(\mathbf{V} - \sum_{j=1}^J \frac{\mathbf{L}_j\mathbf{L}_j^\top}{H_j + \lambda}\right)^{-1}\left(\mathbf{u} - \sum_{j=1}^J \frac{G_j}{H_j + \lambda}\mathbf{L}_j\right) \\ w_j^* &= -\frac{G_j + \mathbf{L}_j^\top\boldsymbol{\delta}^{(t)*}}{H_j + \lambda} \end{aligned} \tag{3.13}$$

### 3.3.1.2 Majorization Minimization: The Construction of Invertible Hessian Matrix

We observe that gradient boosting is, in its essence, a Newton-Raphson optimization algorithm, where we leverage the second order Taylor expansion to the original objective function, and at each boosting iteration train a decision tree to minimize

## CHAPTER 3. GRADIENT BOOSTED TREES

the objective function in the direction of steepest descent. However, there are cases where it falls apart:

- The objective function is not twice differentiable;
- The objective function is twice differentiable, but the Hessian is singular;

We encountered the second case. The true condition behind gradient boosting to converge is to derive an upper bound of original function. This works exactly like Majorization-Minimization (Expectation-Maximization is its special case), where we try to find a quadratic upper bound of the objective function, and minimize the upper bound [28,29]. In order to find that, we just need to find a majorizing Hessian matrix  $\tilde{\mathbf{H}} \succ \mathbf{H}$ . MM is worth considering since we have access to an invertible surrogate hessian matrix, making the objective function much easier to optimize.

Formally, we are minimizing a complicated objective function  $\mathcal{L}$  iteratively. We construct a majorizing function by second-order Taylor expansion around the current best solution  $\mathbf{x}^{(t)}$ , replacing the original hessian by its upper bound. We then find a new solution  $\mathbf{x}^{(t+1)}$  by minimizing the majorizing function. Then we construct a new majorizing function around  $\mathbf{x}^{(t+1)}$ , and so on.

**Definition 1.** (*Majorization*) Suppose  $f$  and  $g$  are twice differentiable at  $\mathbf{x}^{(t)}$ . If  $g$  majorizes  $f$  at  $\mathbf{x}^{(t)}$  then

- $g(\mathbf{x}^{(t)}) = f(\mathbf{x}^{(t)})$ ,

### CHAPTER 3. GRADIENT BOOSTED TREES

- $\nabla_{\mathbf{x}}g(\mathbf{x}^{(t)}) = \nabla_{\mathbf{x}}f(\mathbf{x}^{(t)})$
- $\nabla_{\mathbf{x}}^2g(\mathbf{x}^{(t)}) \succ \nabla_{\mathbf{x}}^2f(\mathbf{x}^{(t)})$

*Proof.* If  $g$  majorizes  $f$  at  $\mathbf{x}^{(t)}$  then  $d = g - f$  has a minimum at  $\mathbf{x}^{(t)}$ . Now use the familiar necessary conditions for the minimum of a differentiable function, which says that the derivative at the minimum is zero and the hessian is symmetric positive definite. Then we have  $g(\mathbf{x}) > f(\mathbf{x}) \forall \mathbf{x} \neq \mathbf{x}^{(t)}$ .  $\square$

**Majorization Minimization (MM) Algorithm** Suppose our current best approximation to the minimum of  $f$  is  $\mathbf{x}^{(t)}$ , and we have a  $g$  that majorizes  $f$  at  $\mathbf{x}^{(t)}$ . If  $\mathbf{x}^{(t)}$  already minimizes  $g$  we stop, otherwise we update  $\mathbf{x}^{(t)}$  to  $\mathbf{x}^{(t+1)}$  by minimizing  $g$ . If we do not stop we have the sandwich inequality:

$$f(\mathbf{x}^{(t+1)}) < g(\mathbf{x}^{(t+1)}) < g(\mathbf{x}^{(t)}) = f(\mathbf{x}^{(t)}) \quad (3.14)$$

Repeating these steps produces a decreasing sequence of function values, and under appropriate additional compactness and continuity conditions this guarantees convergence of the algorithm.

Now we can consider the critical question: how can we construct  $\tilde{\mathbf{H}} \succ \mathbf{H}$ ? One way to achieve this is to find a sharp upper bound along the diagonal entries of  $\mathbf{H}$ , such that  $\tilde{\mathbf{H}} - \mathbf{H}$  is a diagonal matrix with positive entries. This operation shifts the eigenvalues of  $\mathbf{H}$  to the positive direction. As a result,  $\tilde{\mathbf{H}}$  becomes invertible.

## CHAPTER 3. GRADIENT BOOSTED TREES

Let's analyze the diagonal entries of  $\mathbf{H}_\lambda$  and  $\mathbf{V}$

$$\begin{aligned} H_j + \lambda &= \sum_{i \in I_j} \sum_{k=0}^{K-2} \sigma(s_{ik}d_{ik})(1 - \sigma(s_{ik}d_{ik})) + \lambda \\ \mathbf{V}_{kk} &= \sum_{i=1}^n \sigma(s_{ik}d_{ik})(1 - \sigma(s_{ik}d_{ik})) \end{aligned} \quad (3.15)$$

where

$$\begin{aligned} \sigma(x) &= \frac{1}{1 + e^{-x}} \\ s_{ik} &= s(y_i - k) \\ d_{ik} &= \tilde{y}_i - \theta_k \end{aligned} \quad (3.16)$$

We observe it is sufficient to find a sharp majorizer to  $\sigma(x)(1 - \sigma(x))$ . Sec. 4.1 of [29] provides a procedure for constructing a sharp majorizer in the differentiable case. Sec 5.1 of [29] conveniently provided us with our desired majorizer  $\frac{2\sigma(x)-1}{2x}$  to  $\sigma(x)(1 - \sigma(x))$ .

**Sharp Quadratic Majorization in One Dimension** Since we would like  $\tilde{\mathbf{H}}$  to majorize  $\mathbf{H}$  at each of its diagonal entry, it is sufficient to consider the one-dimensional case. Consider twice differentiable  $f(x) \leq g(x) = f(x^{(t)}) + f'(x^{(t)})(x - x^{(t)}) + \frac{1}{2}h(x^{(t)})(x - x^{(t)})^2$ .  $f(x) = \log(1 + \exp(-x))$  is strongly convex, and is majorized by  $g$  at  $x^{(t)}$ . Hence by the inequality  $h$  must satisfy:

$$h(x^{(t)}) \geq \delta(x; x^{(t)}) = \frac{f(x) - f(x^{(t)}) - f'(x^{(t)})(x - x^{(t)})}{\frac{1}{2}(x - x^{(t)})^2} \quad \forall x \neq x^{(t)} \quad (3.17)$$



## CHAPTER 3. GRADIENT BOOSTED TREES

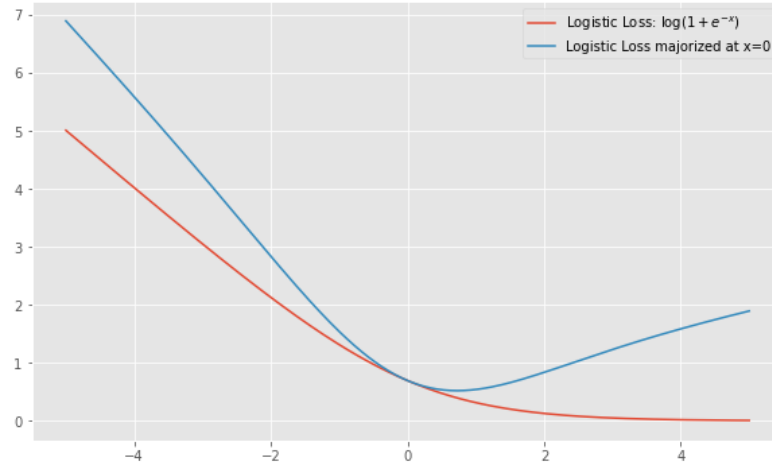


Figure 3.2: Sharp Majorizing function at  $x = 0$ .

We can find such  $h$  if and only if

$$\sup_x \delta(x; x^{(t)}) < \infty$$

And we achieve **sharp quadratic majorization** if we set  $h(x^{(t)}) = \sup_x \delta(x; x^{(t)})$  since we have found the tightest lower bound of  $h(x^{(t)})$  in order to achieve majorization conditions and the convergence rate will also be the fastest.

**Proposition:** We achieve sharp quadratic majorization when  $x = -x^{(t)}$

## CHAPTER 3. GRADIENT BOOSTED TREES

*Proof.*

$$\begin{aligned}
\frac{\partial}{\partial x} \delta(x; x^{(t)}) &= \frac{\frac{1}{2}(x - x^{(t)})^2 (f'(x) + f'(x^{(t)})) - (x - x^{(t)})(f(x) - f(x^{(t)}))}{\frac{1}{4}(x - x^{(t)})^4} \\
\frac{\partial}{\partial x} \delta(x; x^{(t)}) \Big|_{x=-x^{(t)}} &= \frac{-\frac{1}{2}(-2x^{(t)})^2 - (-2x^{(t)})x^{(t)}}{\frac{1}{4}(-2x^{(t)})^4} = 0 \\
\frac{\partial^2}{\partial x^2} \delta(x; x^{(t)}) \Big|_{x=-x^{(t)}} &= \frac{\frac{1}{2}((x - x^{(t)})^2 f''(x) - (x - x^{(t)})(f'(x) - f'(x^{(t)})))}{\frac{1}{4}(x - x^{(t)})^4} \\
&= \frac{2x^{(t)^2} \sigma(x^{(t)})(1 - \sigma(x^{(t)})) + x^{(t)}(\sigma(-x^{(t)}) - \sigma(x^{(t)}))}{4x^{(t)^4}} < 0 \quad \forall x^{(t)} \neq 0
\end{aligned} \tag{3.18}$$

□

Plug in  $x = -x^{(t)}$  and we get

$$h(x^{(t)}) = \delta(-x^{(t)}; x^{(t)}) = \frac{2\sigma(x^{(t)}) - 1}{2x^{(t)}} \tag{3.19}$$

In that case, we have found our modified hessian matrix  $\tilde{\mathbf{H}}$  with diagonal entries:

$$\begin{aligned}
\tilde{H}_j + \lambda &= \sum_{i \in I_j} \sum_{k=0}^{K-2} \frac{2\sigma(s_{ik}d_{ik}) - 1}{2s_{ik}d_{ik}} + \lambda \\
\tilde{V}_{kk} &= \sum_{i=1}^n \frac{2\sigma(s_{ik}d_{ik}) - 1}{2s_{ik}d_{ik}}
\end{aligned} \tag{3.20}$$

**Convergence Properties of Quadratic Majorization Minimization** Suppose

$\mathcal{L} : \mathcal{X} \rightarrow \mathbb{R}$  a convex objective function to be minimized, with gradient  $\mathbf{g} = \nabla_{\mathbf{x}} \mathcal{L}$  and

hessian matrix  $\mathbf{H} = \nabla_{\mathbf{x}}^2 \mathcal{L}$ ,

**Theorem 3.3.1.1.** (Böhning & Lindsay 1988 [30]) Let  $\mathbf{x}^{(0)} \in \mathcal{X}$  and suppose that

$(\mathbf{x}^{(t)})_{t \geq 1}$  is defined by the Majorization Minimization algorithm  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \tilde{\mathbf{H}}(\mathbf{x}^{(t)})^{-1} \mathbf{g}(\mathbf{x}^{(t)})$ ,

where  $\tilde{\mathbf{H}}(\mathbf{x}) \succeq \mathbf{H}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X}$ . The sequence  $(\mathbf{x}^{(t)})_{t \geq 1}$  has the following properties:

### CHAPTER 3. GRADIENT BOOSTED TREES

(a) *Monotonicity:*  $\mathcal{L}(\mathbf{x}^{(t+1)}) \geq \mathcal{L}(\mathbf{x}^{(t)})$ , with strict inequality if  $\mathbf{x}^{(t+1)} \neq \mathbf{x}^{(t)}$ .

(b) *Guaranteed convergence:* The sequence  $(\mathbf{g}(\mathbf{x}^{(t)}))_{t \geq 1}$  converges to 0 if  $\mathcal{L}$  is bounded below.

(c) *Rate of convergence:* The algorithm converges with rate  $\|\mathbb{I} - \tilde{\mathbf{H}}(\mathbf{x}^*)^{-1} \mathbf{H}(\mathbf{x}^*)\|_2$

*Proof.* To prove (a), let the descent direction  $\mathbf{d} = -\tilde{\mathbf{H}}(\mathbf{x}^{(t)})^{-1} \mathbf{g}(\mathbf{x}^{(t)})$ . At  $\mathbf{x}^{(t)}$  construct majorizing function  $\mathcal{M}$  of  $\mathcal{L}$ :

$$\mathcal{M}(\mathbf{x}) = \mathcal{L}(\mathbf{x}^{(t)}) + \mathbf{g}(\mathbf{x}^{(t)})^\top (\mathbf{x} - \mathbf{x}^{(t)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(t)})^\top \tilde{\mathbf{H}}(\mathbf{x}^{(t)}) (\mathbf{x} - \mathbf{x}^{(t)}) \quad (3.21)$$

By Definition 1,  $\mathcal{M}(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}) \forall \mathbf{x} \in \mathcal{X}$ ,  $\mathcal{M}(\mathbf{x}^{(t)}) = \mathcal{L}(\mathbf{x}^{(t)})$ . By Taylor expansion we have:

$$\begin{aligned} \mathcal{L}(\mathbf{x}^{(t+1)}) - \mathcal{L}(\mathbf{x}^{(t)}) &\leq \mathcal{M}(\mathbf{x}^{(t+1)}) - \mathcal{L}(\mathbf{x}^{(t)}) \\ &= \mathbf{g}(\mathbf{x}^{(t)})^\top \mathbf{d} + \frac{1}{2} \mathbf{d}^\top \tilde{\mathbf{H}}(\mathbf{x}^{(t)}) \mathbf{d} \\ &= -\frac{1}{2} \mathbf{g}(\mathbf{x}^{(t)})^\top \tilde{\mathbf{H}}(\mathbf{x}^{(t)})^{-1} \mathbf{g}(\mathbf{x}^{(t)}) \\ &\leq 0 \end{aligned}$$

To prove part (b), suppose for sake of contradiction  $\|\mathbf{g}(\mathbf{x}^{(t)})\|_2$  is bounded away from 0. But the increment above is clearly bounded above by 0, contradicting the boundedness of  $\mathcal{L}$ .

Part (c) was proved by Ostrowski's Theorem [Ortega & Rheinboldt, 1970, Theorem 10.1.3] [31]. □

At sharp quadratic majorization minimization, the convergence rate, similar to

## CHAPTER 3. GRADIENT BOOSTED TREES

that of Newton-Raphson, is still quadratic in our case. As a rough sketch, in the one-dimensional case, plug in  $\sigma(x)(1 - \sigma(x))$  and  $\frac{2\sigma(x)-1}{2x}$  we get  $1 - \frac{2xe^{-x}}{1-e^{-2x}} = \frac{x^2}{6} + \mathcal{O}(x^4)$ .

## Chapter 4

# Faster Training: Algorithmic Speed-up

In this section, we will discuss several algorithmic tricks that make Gradient Boosting fast to train. Also, techniques that enhance the generalizability of boosted trees will be introduced.

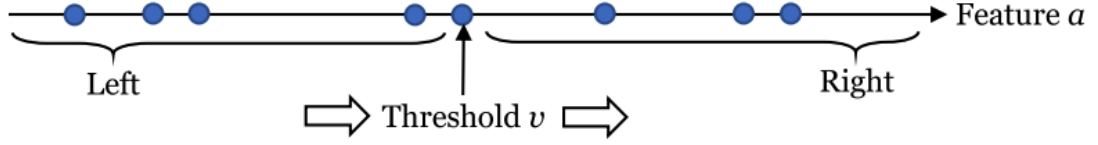


Figure 4.1: We scan the training data from left to right, considering each value as a threshold candidate. [1]

## 4.1 Weighted Quantile Sketch and Feature Quantization

### 4.1.1 Approximate splitting with Quantile Sketch

During one boosting step, at each node we will propose a split candidate by scanning through each feature and each threshold value. The reduction in loss function as a result of the split will be evaluated to assess whether the split candidate should be kept or discarded.

The naive approach is highly inefficient. We would first sort the feature and linear scan all possible unique values. It is unnecessary to do so, since there is not much difference splitting at  $x = 5$  and  $x = 5.00001$ .

To speed up the algorithm, instead of considering all possible unique values to split on, we only consider a fixed number of quantile points as our splitting candidates. We will build separate quantile sketches for each feature.

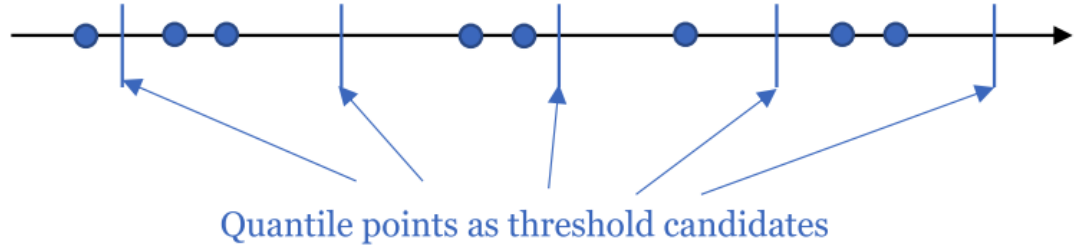


Figure 4.2: We restrict splitting candidates to only quantile points [1]

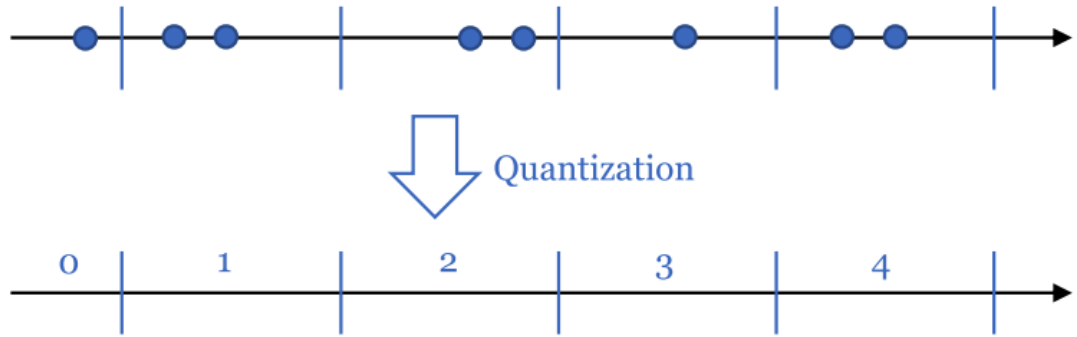


Figure 4.3: Quantize numerical feature values into quantiles. [1]

### 4.1.2 Feature Quantization

By leveraging quantile sketch, we can now place values of each feature to which quantile it belongs. All data points whose values for a particular feature fall between two consecutive quantile points will go to the same quantile index. As long as the value of the same feature of two data points are both placed between two consecutive quantile points, they are considered the same. As a result, we only need to remember the quantile indices instead of the actual feature value. Also, comparison at each split is faster, since it is faster to compare integers than to compare floats.

### 4.1.3 Weighted Quantile

We introduce sample weights to encourage hard-to-learn samples to be placed to the same quantile. As a result, a tree is more likely to make a split and isolate samples where the loss is high. In the original paper [2]:

$$\begin{aligned}
 \tilde{\mathcal{L}} + \text{constant} &= \sum_{i=1}^n [g_i f_t(\mathbf{x}_i) + \frac{1}{2} h_i f_t^2(\mathbf{x}_i) + \frac{g_i^2}{2h_i}] \\
 &= \sum_{i=1}^n \frac{1}{2} h_i [f_t^2(\mathbf{x}_i) + 2 \frac{f_t(\mathbf{x}_i) g_i}{h_i} + (g_i/h_i)^2] \\
 &= \sum_{i=1}^n \frac{1}{2} h_i [f_t(x_i) - (-g_i/h_i)]^2
 \end{aligned} \tag{4.1}$$

where  $-g_i/h_i$  is exactly the prediction to  $f_t(x_i)$  when  $\lambda = 0$ , hence the sum could be thought of as the weighted mean squared error (MSE), where the weight is determined by the hessian evaluated at each data point.

**Definition 2.** (*Weighted Quantile*) Let the multi-set  $\mathcal{D} = \{(x_{1k}, h_1), (x_{2k}, h_2), \dots, (x_{nk}, h_n)\}$  represent the  $k$ -th feature and hessian evaluated at each data point. We can define a quantile function  $r_k : \mathbb{R} \rightarrow [0, 1]$  as

$$r_k(\hat{x}) = \frac{1}{\sum_{(x,h) \in \mathcal{D}} h} \sum_{(x,h) \in \mathcal{D}, x < \hat{x}} h$$



## 4.2 Histogram Aggregation of Gradient Statistics

This is also where the `hist` in `histGBODT` comes from. One crucial step for determining leaf weight is the computation of sum of gradients and Hessians  $G_j, H_j$ , and  $L_j$ . Recall that, for the purpose of split finding, all data points in the same quantile are considered as identical. To speed up split finding process, we should sum up their first and second order gradients and store them into the histogram bin associated with the particular quantile. At each split, the left and right child will inherit parent node's histogram of gradient statistics for subsequent splitting. We can build two histogram with the effort of one, since we can get the other one by subtracting the histogram already built from parent node's histogram.

## 4.3 Tree-growing Strategies

We maintain a priority queue of nodes that will be split next. The node that provides the greatest global reduction in loss will be at the top of the priority queue. When the node at the top of the priority queue is popped, we propose the best splitting candidate and evaluate the reduction in loss after the split. If there is no positive reduction, the node will be finalized as a leaf node and no splitting will be made. Otherwise, calculate the loss reduction after splitting for each of the child

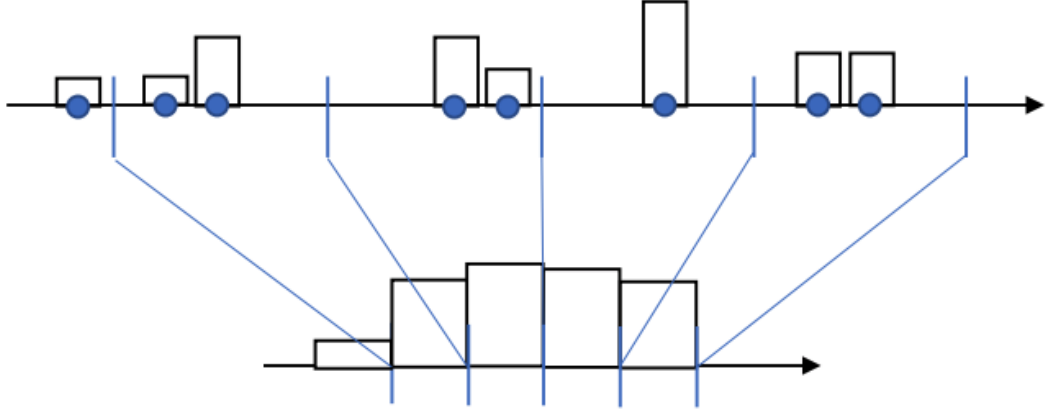


Figure 4.4: Aggregation of gradient statistics along one feature. Only one bar is displayed for simplicity. In reality, there are three bars, storing sum of gradients and Hessians w.r.t.  $w$  (scalars) and second-order mixed gradient  $\nabla_{w,\delta}\tilde{\mathcal{L}}$  (a vector) at each bar.

node, decide whether they will be added to the priority queue.

**Sparsity-Aware Split Finding:** When the training data contains missing values, we cannot compare the threshold hold value at each tree node to the missing value. When splitting the instance set into two, the data point is attempted to split to the left and split to the right. The direction that gives the greater amount of loss reduction is selected. See Appendix A for full split-finding algorithm detail.

We have found that, to compute the loss function evaluated at the optimal parameter, we need to invert the Hessian matrix, which implies that we need to have access to the inverse of the Schur complement  $\mathbf{V} - \mathbf{L}\tilde{\mathbf{H}}_{\lambda}\mathbf{L}^{\top}$ . Computing the inverse would be costly, since we will calculate the inverse at each candidate split point for each

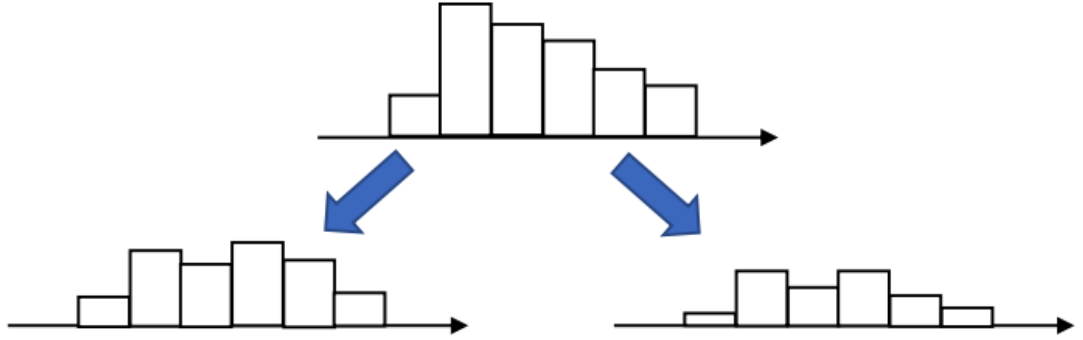


Figure 4.5: Splitting gradient statistics to left and right child. The sum of histograms of left and right child equals the histogram of parent node.

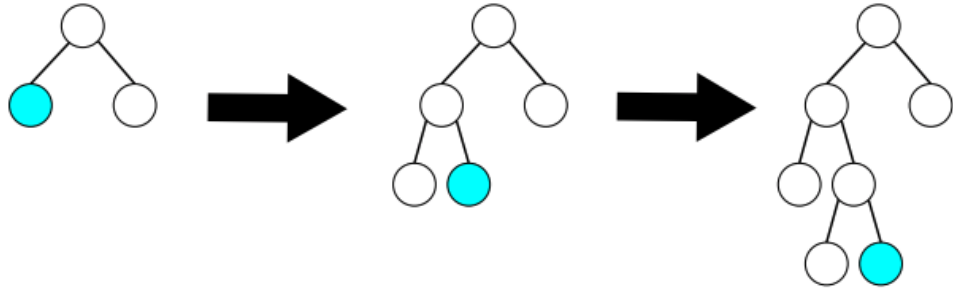


Figure 4.6: Loss-optimal tree growth strategy

feature. Instead, we could leverage the Sherman-Morrison Formula to compute the inverse, which only consists of inverting a diagonal matrix and matrix multiplication.

Suppose we are proposing a split candidate at leaf index  $k$ :

## CHAPTER 4. FASTER TRAINING: ALGORITHMIC SPEED-UP

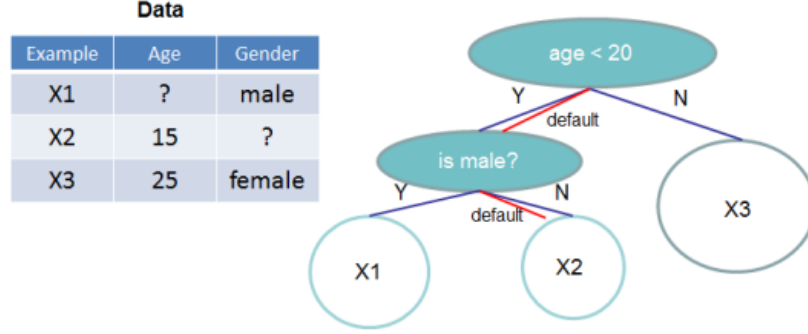


Figure 4.7: When a data point has missing values and they are compared against tree node thresholds, it will be split to the direction that gives rise to greater loss reduction. [2]

Before split:

$$\begin{aligned}
 (\mathbf{V} - \mathbf{L}\tilde{\mathbf{H}}_\lambda\mathbf{L}^\top)^{-1} &= (\mathbf{V} - \sum_{j=1}^J \frac{\mathbf{L}_j\mathbf{L}_j^\top}{H_j + \lambda})^{-1} \\
 &= (\mathbf{V} - \sum_{j \neq k}^J \frac{\mathbf{L}_j\mathbf{L}_j^\top}{H_j + \lambda} - \frac{\mathbf{L}_k\mathbf{L}_k^\top}{H_k + \lambda})^{-1} \\
 &= (\mathbf{A} - \frac{\mathbf{L}_k\mathbf{L}_k^\top}{H_k + \lambda})^{-1} \\
 &= \mathbf{A}^{-1} + \frac{\mathbf{A}^{-1}\mathbf{L}_k\mathbf{L}_k^\top\mathbf{A}^{-1}}{H_k + \lambda - \mathbf{L}_k^\top\mathbf{A}^{-1}\mathbf{L}_k}
 \end{aligned} \tag{4.2}$$

After split:

$$\begin{aligned}
 (\mathbf{V} - \mathbf{L}\tilde{\mathbf{H}}_\lambda\mathbf{L}^\top)^{-1} &= (\mathbf{V} - \sum_{j \neq k}^J \frac{\mathbf{L}_j\mathbf{L}_j^\top}{H_j + \lambda} - \frac{\mathbf{L}_{k_1}\mathbf{L}_{k_1}^\top}{H_{k_1} + \lambda} - \frac{\mathbf{L}_{k_2}\mathbf{L}_{k_2}^\top}{H_{k_2} + \lambda})^{-1} \\
 &= (\mathbf{A} - \frac{\mathbf{L}_{k_1}\mathbf{L}_{k_1}^\top}{H_{k_1} + \lambda} - \frac{\mathbf{L}_{k_2}\mathbf{L}_{k_2}^\top}{H_{k_2} + \lambda})^{-1}
 \end{aligned} \tag{4.3}$$

where

$$\mathbf{L}_k = \mathbf{L}_{k_1} + \mathbf{L}_{k_2}$$

$$\mathbf{H}_k = \mathbf{H}_{k_1} + \mathbf{H}_{k_2}$$

Note that at root node, the matrix  $\mathbf{A}$  above is just the diagonal matrix  $\tilde{\mathbf{V}}$  in our hessian matrix, which is trivial to invert. Hence we iteratively have access to the Schur complement by invoking Sherman-Morrison each time. To compute the inverse after the split, invoke Sherman-Morrison twice.

## 4.4 Gradient-Aware Pruning and Adaptive Learning Rate

We leverage out-of-bag samples to make our learned model more robust to overfitting. The idea is to reduce model complexity when the model does not generalize well to validation samples not used during training.

### 4.4.1 Gradient-Aware Pruning Using Out-of-Bag Samples

While growing a decision tree is a top-down procedure, pruning a tree is bottom-up. Without pruning, decision trees will be grown to the maximum height, unless there is not enough samples at each leaf, or splitting does not reduce loss evaluated

## CHAPTER 4. FASTER TRAINING: ALGORITHMIC SPEED-UP

in-sample. A loss-optimal greedy approach to tree growth is prone to over-fitting. Instead, after a decision tree has been built, we introduce validation samples to prune leaves that do not generalize well to new samples.

During pruning, the loss will be computed on the validation samples and compared between child nodes and parent node. The leaf will be pruned if keeping the left and right child does not decrease the out-of-bag loss. The deeper the decision tree grows, the fewer samples a leaf will contain, thus making the leaf weight estimates have higher variances. The pruning procedure robustify the training process by reducing the model complexity.

When implementing gradient-aware pruning, the histograms of gradient statistics are built on the validation samples top-down. Then, at each pair of children leaf nodes, from the out-of-bag gradient statistics we can evaluate loss before splitting, and after splitting respectively.

### 4.4.2 Adaptive Learning Rate Using Out-of-Bag Samples

Traditionally, in the most popular gradient boosting implementations, learning rate is set to a fix number and it is treated as a hyperparameter. Its optimal value is found via grid-search or random-search. Usually, smaller learning rate gives better performance. Inspired by adaptive learning rate commonly used in deep learning, the

## CHAPTER 4. FASTER TRAINING: ALGORITHMIC SPEED-UP

learning rate is reduced after every few iterations. During training, the benefit of each tree provides will diminish. Thus quite often the last few trees do not provide much predictive power.

Line search in optimization provides us with a solution to picking the optimal learning rate. The learning rate is computed using gradient statistics on out-of-bag samples to improve generalizability. As a result, the tree that does not have much predictive power will receive a learning rate close to 0, thus encouraging pruning at the whole tree-level.

Formally, we would like to pick the learning-rate  $\nu$  at each iteration right after the tree is built and pruned.  $\mathcal{D}_{\text{val}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n_{\text{val}}}$  is the out-of-bag (validation) samples.

$$\mathcal{L}^{\text{OOB}}(\nu, \mathbf{w}^{(t)}, \boldsymbol{\delta}^{(t)}) = \sum_{i=1}^{n_{\text{val}}} l_{\text{ord}}(y_i, F^{(t-1)}(\mathbf{x}_i) + \nu f_t(\mathbf{x}_i), \boldsymbol{\theta}^{(t-1)} + \nu \boldsymbol{\delta}^{(t)}) \quad (4.4)$$

For simplicity, we re-write this in terms of our previous best estimate,  $\mathbf{z}^{(t-1)}$ , consisting of the prediction of the target on out-of-bag samples, and  $\boldsymbol{\theta}^{(t-1)}$ ; a descent direction  $\boldsymbol{\beta}^{(t)}$ , consisting of leaf weights  $\mathbf{w}^{(t)}$  and  $\boldsymbol{\delta}^{(t)}$ ; and step size  $\nu$ .

$$\mathcal{L}^{\text{OOB}}(\mathbf{z}^{(t-1)} + \nu \boldsymbol{\beta}^{(t)}) \simeq \mathcal{L}^{\text{OOB}}(\mathbf{z}^{(t-1)}) + \nu \mathbf{g}^{\text{OOB}\top} \boldsymbol{\beta}^{(t)} + \frac{\nu^2}{2} \boldsymbol{\beta}^{(t)\top} \tilde{\mathbf{H}}^{\text{OOB}} \boldsymbol{\beta}^{(t)} \quad (4.5)$$

$$\nu^* = - \frac{\mathbf{g}^{\text{OOB}\top} \boldsymbol{\beta}^{(t)}}{\boldsymbol{\beta}^{(t)\top} \tilde{\mathbf{H}}^{\text{OOB}} \boldsymbol{\beta}^{(t)}} \quad (4.6)$$

Here the definition of  $\mathbf{g}^{\text{OOB}}$  and  $\tilde{\mathbf{H}}^{\text{OOB}}$  is the same to the counterpart during training in-sample. The only difference is the set of samples being evaluated on. In practice

## CHAPTER 4. FASTER TRAINING: ALGORITHMIC SPEED-UP

the over-fitting robust out-of-bag learning rate estimate is clipped to the range  $[0, 1]$ .



# Chapter 5

## Experiments

We conduct experiments on both simulated and real-world datasets. We start with simple synthetic datasets to understand the behavior of the algorithm. Then, we report our model’s performance on real-world data sets.

We employ two evaluation metrics as measures of predictive performance in ordinal classification tasks:

- *Mean Absolute Error* (MAE):  $\frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$
- *Mean Zero-One Error* (ZOE):  $\frac{1}{n} \sum_{i=1}^n I(\hat{y}_i \neq y_i)$

### 5.1 Artificial Data

Figure 5.1 presents the behavior of HistGBODT on 3 2-D synthetic datasets. Each dataset is generated from a Gaussian distribution divided into 3 ordinal ranks using

## CHAPTER 5. EXPERIMENTS

non-linear decision boundaries. The second dataset is labelled the same way as the first dataset, except the data points closest to the decision boundaries are randomized with probability  $\frac{1}{2}$ . The prediction accuracy on the three datasets are: 99.73%, 94.41%, and 98.01%, respectively. The simulation results show that the algorithm is working reasonably well on this task.

### 5.2 Benchmark Data

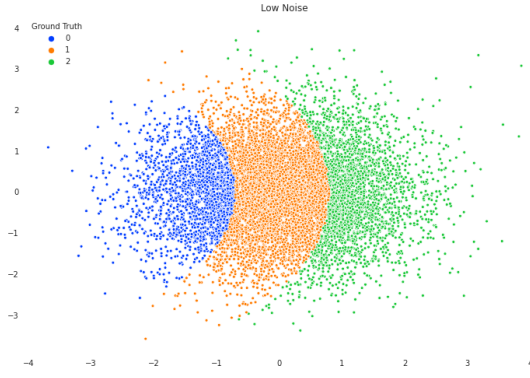
We performed experiments with 16 benchmark datasets that were used by Chu and Keerthi [32]<sup>1</sup>. The target values were discretized into 10 equally-sized quantiles. These bins thus divide the range of the target values into 10 ordinal ranks. The input vectors were normalized to zero mean and unit variance coordinate-wise. To show the advantage of explicitly modeling the ordinal nature of the response variables, we also employed the standard gradient boosting trained using the multi-class cross entropy loss, and using the mean-squared error with the predicted values rounded to the nearest ordinal scale. We randomly partitioned each dataset into training/testing splits as specified in table 5.1. The partition was repeated 20 times.

From the testing set results averaged over 20 trials, we could see a clear advantage of minimizing the All-Threshold loss over traditional methods such as multi-class cross entropy and mean-squared error. To determine statistical significance, we conduct

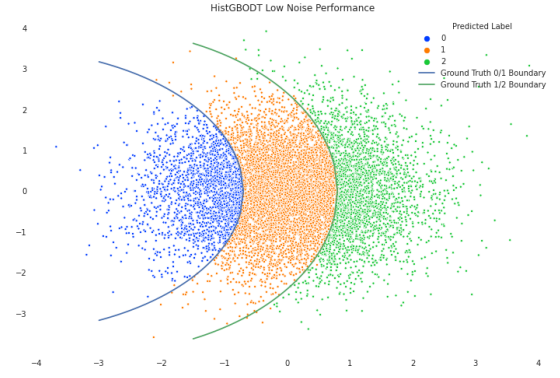
---

<sup>1</sup>Datasets available for download at: <http://www.gatsby.ucl.ac.uk/~chuwei/ordinalregression.html>

## CHAPTER 5. EXPERIMENTS



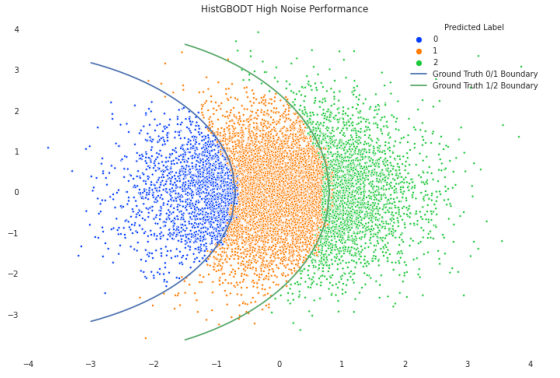
(a) Artificial dataset I with low noise



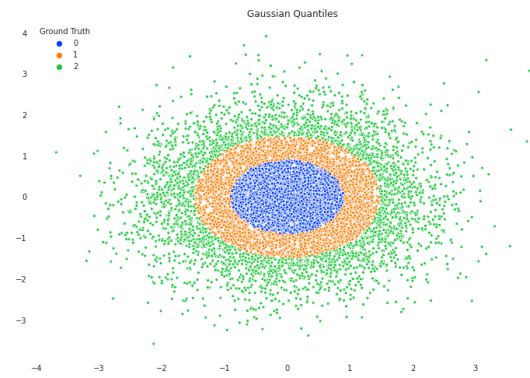
(b) Artificial dataset I prediction



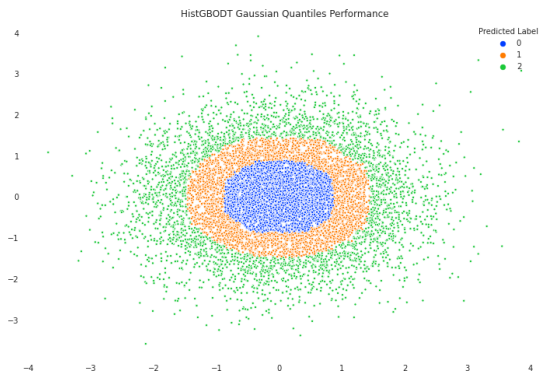
(c) Artificial dataset II with high noise



(d) Artificial dataset II prediction



(e) Artificial dataset III Gaussian Quantiles



(f) Artificial dataset III prediction

Figure 5.1: HistGBODT performance on synthetic datasets.

## CHAPTER 5. EXPERIMENTS

Datasets		Attributes(Numeric,Nominal)	Training Instances	Instances for Test
Set I	Diabetes	2(2,0)	30	13
	Pyrimidines	27(27,0)	50	24
	Triazines	60(60,0)	100	86
	Wisconsin Breast Cancer	32(32,0)	130	64
	Machine CPU	6(6,0)	150	59
	Auto MPG	7(4,3)	200	192
	Boston Housing	13(12,1)	300	206
	Stocks Domain	9(9,0)	600	350
	Abalone	8(7,1)	3177	1000
Set II	Bank Domains(1)	8(8,0)	5461	2731
	Bank Domains(2)	32(32,0)	5461	2731
	Computer Activity(1)	12(12,0)	5461	2731
	Computer Activity(2)	21(21,0)	5461	2731
	California Housing	8(8,0)	10427	5213
	Census Domains(1)	8(8,0)	11189	5595
	Census Domains(2)	16(16,0)	11189	5595

Table 5.1: Datasets and their characteristics. “Attributes” state the number of numerical and nominal attributes. “Training Instances” and “Instances for Test” specify the size of training/test partitions.

## CHAPTER 5. EXPERIMENTS

	Mean zero-one error			Mean absolute Error		
Data	AT	MultiClass	MSE	AT	MultiClass	MSE
Diabetes	69.23%	69.23%	69.23%	1.154	1.154	1.154
Pyrimidines	<b>70.83%</b>	75.62%	77.92%	<b>1.240</b>	1.577	1.383
Triazines	67.85%	<b>67.79%</b>	71.05%	<b>1.112</b>	1.319	1.153
Wisconsin Breast Cancer	86.33%	<b>85.78%</b>	88.83%	<b>2.123</b>	2.805	2.438
Machine CPU	<b>31.69%</b>	41.10%	37.46%	<b>0.4483</b>	0.939	0.5161
Auto MPG	<b>45.83%</b>	50.86%	47.24%	<b>0.5464</b>	0.6331	0.5576
Boston Housing	<b>43.91%</b>	45.02%	44.22%	<b>0.5204</b>	0.5680	0.5235
Stocks Domain	22.06%	<b>21.29%</b>	21.84%	0.2261	<b>0.2164</b>	0.2257
Abalone	<b>42.34%</b>	48.40%	44.68%	<b>0.5212</b>	0.7552	0.5384
Bank Domains(1)	72.28%	<b>69.17%</b>	78.4%	<b>1.456</b>	1.788	1.530
Bank Domains(2)	<b>38.81%</b>	44.75%	39.47%	<b>0.4094</b>	0.5145	0.4171
Computer Activity(1)	<b>51.45%</b>	52.28%	57.66%	<b>0.7039</b>	0.8177	0.7929
Computer Activity(2)	<b>41.60%</b>	47.05%	45.18%	<b>0.4687</b>	0.5881	0.5134
California Housing	<b>46.94%</b>	50.68%	49.18%	<b>0.5683</b>	0.6928	0.6016
Census Domains(1)	<b>60.58%</b>	62.59%	66.83%	<b>0.9296</b>	1.410	0.998
Census Domains(2)	62.25%	<b>60.25%</b>	66.93%	<b>0.9793</b>	1.130	1.022

Table 5.2: Testing set results: the bold font indicates the cases in which the average value is the lowest in the results of minimizing the three loss functions.

## CHAPTER 5. EXPERIMENTS

	Mean zero-one error		Mean absolute error	
Null hypothesis	AT < MultiClass	AT < MSE	AT < MultiClass	AT < MSE
Test Statistic	20	1	1	1
P-value	0.01155	0.0004026	0.0004026	0.0004026

Table 5.3: Nonparametric Wilcoxon signed-rank test results

Wilcoxon signed-rank tests with a p-value threshold of 0.025 (here, a Dunn–Šidák correction is applied for multiple comparison). We observe that the performance gains are statistically significant for all comparisons.

# Chapter 6

## Conclusion and Future Directions

### 6.1 Conclusion

Ordinal classification is an important supervised learning task but has not been widely implemented in popular machine learning frameworks such as `scikit-learn` or `XGBoost`, etc. In this paper, we propose minimizing the All-Threshold loss function combined with the powerful gradient boosting as a solution. To optimize the objective, we encounter the issue of non-invertibility of the Hessian matrix and adopted the Majorization-Minimization (MM) algorithm to perform Newton-Raphson updates. Experiments on benchmark datasets indicates the generalization performance is competitive and statistically significantly better than existing approaches such as minimizing the cross entropy and mean-squared error (MSE).

## 6.2 Future Directions

Our implementation is under the `scikit-learn` framework, using `Cython` for native `C` speed performance and parallel computing. However, integration with highly optimized gradient boosting package such as `XGBoost` and `CatBoost` could be even more desirable due to access to GPU computing.

Traditionally, decision trees split along directions parallel to the coordinate axes. In high-dimensional space when classes seem to be inseparable along any single coordinate direction, or decision boundaries are highly irregular and non-linear, such splits require very deep trees with complicated structures, such increasing variance and overfitting. Several work has proposed to split along a linear combination of coordinates (Oblique Decision Trees, [33–37]). Integrating All-Threshold loss with boosted oblique decision trees could be an interesting future research direction.



# Bibliography

- [1] H. Cho, “Speeding up gradient boosting for training and prediction.”
- [2] T. Chen and C. Guestrin, “Xgboost: A scalable tree boosting system,” p. 785–794, 2016. [Online]. Available: <https://doi.org/10.1145/2939672.2939785>
- [3] T. Gneiting and A. E. Raftery, “Weather forecasting with ensemble methods,” *Science*, vol. 310, no. 5746, pp. 248–249, 2005. [Online]. Available: <https://science.sciencemag.org/content/310/5746/248>
- [4] T. Gneiting and M. Katzfuss, “Probabilistic forecasting,” *Annual Review of Statistics and Its Application*, vol. 1, no. 1, pp. 125–151, 2014. [Online]. Available: <https://doi.org/10.1146/annurev-statistics-062713-085831>
- [5] J. Bennett, C. Elkan, B. Liu, P. Smyth, and D. Tikk, “Kdd cup and workshop 2007,” *SIGKDD Explor. Newsl.*, vol. 9, no. 2, p. 51–52, Dec. 2007. [Online]. Available: <https://doi.org/10.1145/1345448.1345459>
- [6] X. He, J. Pan, O. Jin, T. Xu, B. Liu, T. Xu, Y. Shi, A. Atallah,

## BIBLIOGRAPHY

- R. Herbrich, S. Bowers, and et al., “Practical lessons from predicting clicks on ads at facebook,” p. 1–9, 2014. [Online]. Available: <https://doi.org/10.1145/2648584.2648589>
- [7] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu, “Lightgbm: A highly efficient gradient boosting decision tree,” pp. 3146–3154, 2017. [Online]. Available: <http://papers.nips.cc/paper/6907-lightgbm-a-highly-efficient-gradient-boosting-decision-tree.pdf>
- [8] A. V. Dorogush, V. Ershov, and A. Gulin, “Catboost: gradient boosting with categorical features support,” *CoRR*, vol. abs/1810.11363, 2018. [Online]. Available: <http://arxiv.org/abs/1810.11363>
- [9] A. V. Dorogush, A. Gulin, G. Gusev, N. Kazeev, L. O. Prokhorenkova, and A. Vorobev, “Catboost: unbiased boosting with categorical features,” *CoRR*, vol. abs/1706.09516, 2017. [Online]. Available: <http://arxiv.org/abs/1706.09516>
- [10] P. McCullagh, “Regression models for ordinal data,” *Journal of the Royal Statistical Society*, vol. 42, no. 2, pp. 109–142, 1980. [Online]. Available: <http://www.jstor.org/stable/2984952>
- [11] J. A. Anderson, “Regression and ordered categorical variables,” *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 46, no. 1, pp. 1–22, 1984. [Online]. Available: <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/j.2517-6161.1984.tb01270.x>

## BIBLIOGRAPHY

- [12] A. Agresti, “Modelling ordered categorical data: recent advances and future challenges,” *Statistics in Medicine*, vol. 18, no. 17-18, p. 2191–2207, 1999.
- [13] Y. Freund, R. Iyer, R. E. Schapire, and Y. Singer, “An efficient boosting algorithm for combining preferences,” *J. Mach. Learn. Res.*, vol. 4, p. 933–969, Dec. 2003.
- [14] K. Dembczynski, W. Kotlowski, and R. Slowinski, “Ordinal classification with decision rules,” 2007.
- [15] J. D. M. Rennie, “Loss functions for preference levels: Regression with discrete ordered labels,” pp. 180–186, 2005.
- [16] K. Crammer and Y. Singer, “Pranking with ranking,” pp. 641–647, 2002.  
[Online]. Available: <http://papers.nips.cc/paper/2023-pranking-with-ranking.pdf>
- [17] A. Shashua and A. Levin, “Ranking with large margin principle: Two approaches,” pp. 961–968, 2003. [Online]. Available: <http://papers.nips.cc/paper/2269-ranking-with-large-margin-principle-two-approaches.pdf>
- [18] W. Chu and S. S. Keerthi, “New approaches to support vector ordinal regression,” p. 145–152, 2005. [Online]. Available: <https://doi.org/10.1145/1102351.1102370>

## BIBLIOGRAPHY

- [19] H.-T. Lin and L. Li, “Large-margin thresholded ensembles for ordinal regression: Theory and practice,” pp. 319–333, 2006.
- [20] L. Li and H. tien Lin, “Ordinal regression by extended binary classification,” pp. 865–872, 2007. [Online]. Available: <http://papers.nips.cc/paper/3125-ordinal-regression-by-extended-binary-classification.pdf>
- [21] F. Pedregosa, F. Bach, and A. Gramfort, “On the consistency of ordinal regression methods,” *J. Mach. Learn. Res.*, vol. 18, no. 1, p. 1769–1803, Jan. 2017.
- [22] J. Friedman, “Another approach to polychotomous classification, technical report,” 1996.
- [23] T. H. . R. Tibshirani, “Classification by pairwise coupling,” *Annals of Statistics*, vol. 26, no. 2, pp. 451–471, 1998.
- [24] R. Rifkin and A. Klautau, “In defense of one-vs-all classification,” *J. Mach. Learn. Res.*, vol. 5, p. 101–141, Dec. 2004.
- [25] J. H. Friedman, “Greedy function approximation: A gradient boosting machine,” *Annals of Statistics*, vol. 29, pp. 1189–1232, 2000.
- [26] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, “Classification and regression trees,” 1984.
- [27] T.-T. Lu and S.-H. Shiou, “Inverses of 2 x 2 block matrices,” *Computers and Mathematics with Applications*, vol. 43, no. 1-2, pp. 119–129, 2002.

## BIBLIOGRAPHY

- [28] T. Chen, S. Singh, B. Taskar, and C. Guestrin, “Efficient Second-Order Gradient Boosting for Conditional Random Fields,” vol. 38, pp. 147–155, 09–12 May 2015. [Online]. Available: <http://proceedings.mlr.press/v38/chen15b.html>
- [29] J. de Leeuw and K. Lange, “Sharp quadratic majorization in one dimension.” 2011.
- [30] D. Böhning and B. Lindsay, “Monotonicity of quadratic-approximation algorithms,” *Annals of the Institute of Statistical Mathematics*, vol. 40, pp. 641–663, 02 1988.
- [31] J. M. Ortega and W. C. Rheinboldt, *Iterative Solution of Nonlinear Equations in Several Variables*. USA: Society for Industrial and Applied Mathematics, 2000.
- [32] W. Chu and Z. Ghahramani, “Gaussian processes for ordinal regression.” *Journal of Machine Learning Research*, vol. 6, pp. 1019–1041, 07 2005.
- [33] L. Breiman, “Random forests,” *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001. [Online]. Available: <http://dx.doi.org/10.1023/A%3A1010933404324>
- [34] J. J. Rodriguez, L. I. Kuncheva, and C. J. Alonso, “Rotation forest: A new classifier ensemble method,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, no. 10, pp. 1619–1630, 2006.

## BIBLIOGRAPHY

- [35] B. Menze, B. Kelm, D. Splitthoff, U. Köthe, and F. Hamprecht, “On oblique random forests,” pp. 453–469, 09 2011.
- [36] T. M. Tomita, M. Maggioni, and J. T. Vogelstein, “Randomer forests,” *ArXiv*, vol. abs/1506.03410, 2015.
- [37] T. M. Tomita, J. Browne, C. Shen, J. Chung, J. L. Patsolic, B. Falk, J. Yim, C. E. Priebe, R. Burns, M. Maggioni, and J. T. Vogelstein, “Sparse projection oblique randomer forests,” *ArXiv*, vol. abs/1506.03410, 2019.
- [38] L. Devroye, L. Györfi, and G. Lugosi, “A probabilistic theory of pattern recognition,” *New York: Springer Verlag*, vol. 31, 01 1996.
- [39] G. Biau, L. Devroye, and G. Lugosi, “Consistency of random forests and other averaging classifiers,” *Journal of Machine Learning Research*, vol. 9, pp. 2015–2033, 09 2008.
- [40] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, ser. Springer Series in Statistics. New York, NY, USA: Springer New York Inc., 2001.
- [41] K. P. Murphy, *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012.
- [42] P. Bickel and K. Doksum, *Mathematical Statistics: Basic Ideas and Selected Topics.*, 01 2007, vol. 56.

# Appendix A

---

**Algorithm 1** Histogram-based Gradient Boosting

---

**Input:** Training data  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$

$L_2$  regularization  $\lambda$

Learning rate  $\nu$

Leaf penalty  $\gamma$

Maximum iterations  $T$

Maximum depth of tree  $D$

Maximum number of nodes  $J$

Minimum number of training example in each leaf  $n_{\min}$

If pruning will be performed  $\mathbb{1}_{\text{prune}}$

If adaptive learning rate will be selected  $\mathbb{1}_{\text{adaptive}}$

Number of training observations to be set aside for pruning/adaptive

learning rate selection  $n_{\text{val}}$

**Output:** Collection of trees  $\{f_t\}_{t \geq 1, t \in \mathbb{N}}$ , Category threshold parameter  $\theta$

Shrinkage factor on each tree  $\{\nu_t\}_{t \geq 1, t \in \mathbb{N}}$

## APPENDIX A.

```

1: if  $\mathbb{1}_{\text{prune}} = 1$  or  $\mathbb{1}_{\text{adaptive}} = 1$  then

2:    $\mathcal{D}_{\text{train}}, \mathcal{D}_{\text{val}} = \text{TRAINVALIDATIONSPLIT}(\mathcal{D}, n_{\text{val}})$ 

3: else

4:    $\mathcal{D}_{\text{train}} = \mathcal{D}, \mathcal{D}_{\text{val}} = \emptyset$ 

5: end if

6:  $\dot{\mathbf{X}}_{\text{train}}, \dot{\mathbf{X}}_{\text{val}} = \text{FEATUREQUANTIZE}(\mathbf{X}_{\text{train}}, \mathbf{X}_{\text{val}})$ 

7:  $\tilde{y}^{(0)}, \boldsymbol{\theta}^{(0)} = \underset{\tilde{y}, \boldsymbol{\theta}}{\text{argmin}} \sum_{i \in \mathcal{D}_{\text{train}}} l_{\text{ord}}(y_i, \tilde{y}, \boldsymbol{\theta})$ 

8:  $\tilde{\mathbf{y}}_{\text{train}}^{(0)} = \tilde{\mathbf{y}}_{\text{val}}^{(0)} = \tilde{y}^{(0)} \mathbf{1}$ 

9: for  $t = 1 : T$  do

10:    $\{(g_i, h_i, (L_{ik})_{k=0}^{K-2})\}_{i \in \mathcal{D}_{\text{train}}} = \text{UPDATEGRADIENTSTATISTICS}(\mathbf{y}_{\text{train}}, \tilde{\mathbf{y}}_{\text{train}}^{(t-1)}, \boldsymbol{\theta}^{(t-1)})$ 

11:   if  $\mathbb{1}_{\text{prune}} = 1$  or  $\mathbb{1}_{\text{adaptive}} = 1$  then

12:      $\{(g_i, h_i, (L_{ik})_{k=0}^{K-2})\}_{i \in \mathcal{D}_{\text{val}}} = \text{UPDATEGRADIENTSTATISTICS}(\mathbf{y}_{\text{val}}, \tilde{\mathbf{y}}_{\text{val}}^{(t-1)}, \boldsymbol{\theta}^{(t-1)})$ 

13:   else

14:      $\{(g_i, h_i, (L_{ik})_{k=0}^{K-2})\}_{i \in \mathcal{D}_{\text{val}}} = \emptyset$ 

15:   end if

16:    $f_t, \boldsymbol{\delta}^{(t)} = \text{GROWTREE}(\dot{\mathbf{X}}_{\text{train}}, \{(g_i, h_i, (L_{ik})_{k=0}^{K-2})\}_{i \in \mathcal{D}_{\text{train}}}, J, D, \lambda, \gamma, n_{\min})$ 

17:   if  $\mathbb{1}_{\text{prune}} = 1$  then

18:      $f_t, \boldsymbol{\delta}^{(t)} = \text{DOPRUNE}(f_t, \{(g_i, h_i, (L_{ik})_{k=0}^{K-2})\}_{i \in \mathcal{D}_{\text{val}}})$ 

19:   end if

20:   if  $\mathbb{1}_{\text{adaptive}} = 1$  then

21:      $\nu_t = \text{COMPUTEADAPTIVELARNINGRATE}(f_t, \{(g_i, h_i, (L_{ik})_{k=0}^{K-2})\}_{i \in \mathcal{D}_{\text{val}}})$ 

```



## APPENDIX A.

```

22:          $\nu_t = \text{CLIP}(\nu_t, 0, \nu)$ 
23:     else
24:          $\nu_t = \nu$ 
25:     end if
26:      $\tilde{\mathbf{y}}_{\text{train}}^{(t)} = \tilde{\mathbf{y}}_{\text{train}}^{(t-1)} + \nu_t f_t(\dot{\mathbf{X}}_{\text{train}})$ 
27:      $\boldsymbol{\theta}^{(t)} = \boldsymbol{\theta}^{(t-1)} + \nu_t \boldsymbol{\delta}^{(t)}$ 
28:     if  $\mathbb{I}_{\text{prune}} = 1$  or  $\mathbb{I}_{\text{adaptive}} = 1$  then
29:          $\tilde{\mathbf{y}}_{\text{val}}^{(t)} = \tilde{\mathbf{y}}_{\text{val}}^{(t-1)} + \nu_t f_t(\dot{\mathbf{X}}_{\text{val}})$ 
30:     end if
31:      $\mathbb{I}_{\text{EarlyStop}} = \text{CHECKEARLYSTOPPING}(\mathbf{y}_{\text{train}}, \tilde{\mathbf{y}}_{\text{train}}, \mathbf{y}_{\text{val}}, \tilde{\mathbf{y}}_{\text{val}}, \nu_t)$ 
32:     if  $\mathbb{I}_{\text{EarlyStop}} = 1$  then
33:         BREAK
34:     end if
35: end for
36: end function

```

---

## APPENDIX A.

---

### Algorithm 2 Grow a tree at iteration $t$

---

**Input:** Quantized training observations  $\dot{\mathbf{X}}_{\text{train}}$

Gradient Statistics  $\{(g_i, h_i, (L_{ik})_{k=0}^{K-2})\}_{i \in \mathcal{D}_{\text{train}}}$

$L_2$  regularization  $\lambda$

Leaf penalty  $\gamma$

Maximum depth of tree  $D$

Maximum number of nodes  $J$

Minimum number of training example in each leaf  $n_{\min}$

**Output:** Tree  $f_t$  with loss optimal structure  $q_t$

Category threshold parameter Newton Raphson step  $\delta^{(t)}$

- 1:  $I_1 = \{1, \dots, n_{\text{train}}\}$  ▷ Initialize tree structure with a single node
- 2:  $\left(m_1, v_1, \mathbb{1}_{\text{MissingGoesLeft}}, \Delta \tilde{\mathcal{L}}^{(t)}(I_1)\right) = \text{EVALUATESPLIT}(I_1)$
- 3:  $d_1 = 1$  ▷ Depth level of node 1
- 4:  $S = \text{PriorityQueue}(\{[1 : \Delta \tilde{\mathcal{L}}^{(t)}(I_1)]\})$  ▷ Set of leaf nodes to be split
- 5:  $\text{LeafSet} = \text{HashSet}(\{\})$
- 6: **while**  $j = 1; j \leq J$  **do**
- 7:    $u = S.\text{pop}()$  ▷ Retrieve node with the largest loss reduction after split
- 8:   **if**  $d_u < D$  **then**
- 9:      $I_{j+1} = \{i \in I_u : \dot{\mathbf{x}}_{im_u} < \tau_u\}$  ▷ Split  $I_u$  using the  $m_u^{\text{th}}$  feature & threshold  $\tau_u$
- 10:     $I_{j+2} = \{i \in I_u : \dot{\mathbf{x}}_{im_u} \geq \tau_u\}$
- 11:     $d_{j+1}, d_{j+2} = d_u + 1, d_u + 1$  ▷ Assign depth level of new nodes

## APPENDIX A.

```

12:      if  $d_j + 1 = D$  then
13:          LeafSet.add( $I_{j+1}$ )
14:          LeafSet.add( $I_{j+2}$ )
15:      else
16:          if  $|I_{j+1}| \geq n_{\min}$  then
17:               $\left(m_{j+1}, \tau_{j+1}, \mathbb{1}_{\text{MissingGoesLeft}}, \Delta \tilde{\mathcal{L}}^{(t)}(I_{j+1})\right) = \text{EVALUATESPLIT}(I_{j+1})$ 
18:               $S.\text{add}(\{j + 1 : \Delta \tilde{\mathcal{L}}^{(t)}(I_{j+1})\})$ 
19:               $j = j + 1$ 
20:          else
21:              LeafSet.add( $I_{j+1}$ )
22:          end if
23:          if  $|I_{j+2}| \geq n_{\min}$  then
24:               $\left(m_{j+2}, \tau_{j+2}, \mathbb{1}_{\text{MissingGoesLeft}}, \Delta \tilde{\mathcal{L}}^{(t)}(I_{j+2})\right) = \text{EVALUATESPLIT}(I_{j+2})$ 
25:               $S.\text{add}(\{j + 2 : \Delta \tilde{\mathcal{L}}^{(t)}(I_{j+2})\})$ 
26:               $j = j + 1$ 
27:          else
28:              LeafSet.add( $I_{j+2}$ )
29:          end if
30:      end if
31:  end if
32:  if  $j == J$  then

```

## APPENDIX A.

```
33:      LeafSet.add( $I_{j+1}$ )
34:      LeafSet.add( $I_{j+2}$ )
35:      while  $|S| \neq 0$  do
36:          LeafSet.add( $S.pop()$ )
37:      end while
38:  end if
39: end while
40:  $(\boldsymbol{w}^{(t)}, \boldsymbol{\delta}^{(t)}) = \text{FINALIZELEAVES}(\text{LeafSet}) \triangleright$  Leaf weights & threshold Newton step
```

---

## APPENDIX A.

---

**Algorithm 3** EVALUATESPLIT( $I_j$ )     $\triangleright$  Compute greatest loss reduction for a node

---

```

1:  $(G_j, H_j, \mathbf{L}_j) = \sum_{i \in I_j} (g_i, h_i, (L_{ik})_{k=0}^{K-2})$ 

2: for feature  $m$  in  $\mathbf{X}_{\text{train}}$  do

3:    $M_m = \{i \in I_j : \mathbf{X}_{im} = \emptyset\}$   $\triangleright$  Provision for missing values

4:    $(G_m^{(\emptyset)}, H_m^{(\emptyset)}, \mathbf{L}_m^{(\emptyset)}) = \sum_{i \in M_m} (g_i, h_i, (L_{ik})_{k=0}^{K-2})$ 

5:    $\mathcal{T}_m = \{\text{set of quantile points for feature } m\}$ 

6:   for  $\tau \in \mathcal{T}_m$  do  $\triangleright$  Consider each threshold candidate

7:     LeftSet $_{m,\tau} = \{i \in I_j : \mathbf{X}_{im} < \tau\}$ 

8:     RightSet $_{m,\tau} = \{i \in I_j : \mathbf{X}_{im} \geq \tau\}$ 

9:      $(G_{m,\tau}^{(\text{left})}, H_{m,\tau}^{(\text{left})}, \mathbf{L}_{m,\tau}^{(\text{left})}) = \sum_{i \in \text{LeftSet}_{m,\tau}} (g_i, h_i, (L_{ik})_{k=0}^{K-2})$ 

10:     $(G_{m,\tau}^{(\text{right})}, H_{m,\tau}^{(\text{right})}, \mathbf{L}_{m,\tau}^{(\text{right})}) = \sum_{i \in \text{RightSet}_{m,\tau}} (g_i, h_i, (L_{ik})_{k=0}^{K-2})$ 

11:     $\Delta \tilde{\mathcal{L}}^{(t)}(m, \tau, \mathbb{1}_{\text{MissingGoesLeft}}) = \text{COMPUTELOSSREDUCTION} \left( \right.$ 

12:       $(G_{m,\tau}^{(\text{left})}, H_{m,\tau}^{(\text{left})}, \mathbf{L}_{m,\tau}^{(\text{left})}), (G_{m,\tau}^{(\text{right})}, H_{m,\tau}^{(\text{right})}, \mathbf{L}_{m,\tau}^{(\text{right})}), (G_m^{(\emptyset)}, H_m^{(\emptyset)}, \mathbf{L}_m^{(\emptyset)}), \lambda, \gamma \Big)$ 

13:    end for

14: end for

15:  $(m, \tau, \mathbb{1}_{\text{MissingGoesLeft}}) = \underset{m', \tau', \mathbb{1}'_{\text{MissingGoesLeft}}}{\text{argmax}} \quad \Delta \tilde{\mathcal{L}}^{(t)}(m', \tau', \mathbb{1}'_{\text{MissingGoesLeft}})$ 

16: return  $(m, \tau, \mathbb{1}_{\text{MissingGoesLeft}}, \Delta \tilde{\mathcal{L}}^{(t)}(m, \tau, \mathbb{1}_{\text{MissingGoesLeft}}))$ 

```

---

## APPENDIX A.

---

**Algorithm 4** DOPRUNE( $f_t, \{(g_i, h_i, (L_{ik})_{k=0}^{K-2})\}_{i \in \mathcal{D}_{\text{val}}}$ )

---

```

1:  $\{I_h^{(left)}, I_h^{(right)}\}_{h=1}^H = \text{COLLECTCHILDRENPairs}(f_t)$ 

2:  $S_{\text{merge}} = \text{PriorityQueue}([])$ 

3: for  $h = 1 : H$  do

4:    $(G_h^{(left)}, H_h^{(left)}, \mathbf{L}_h^{(left)}) = \sum_{i \in I_h^{(left)}} (g_i, h_i, (L_{ik})_{k=0}^{K-2})$ 

5:    $(G_h^{(right)}, H_h^{(right)}, \mathbf{L}_h^{(right)}) = \sum_{i \in I_h^{(right)}} (g_i, h_i, (L_{ik})_{k=0}^{K-2})$ 

6:    $\Delta_{\text{merge}} \tilde{\mathcal{L}}_h^{(t)} = \text{EVALUATEMERGE}\left((G_h^{(left)}, H_h^{(left)}, \mathbf{L}_h^{(left)}), (G_h^{(right)}, H_h^{(right)}, \mathbf{L}_h^{(right)})\right)$ 

7:   if  $\Delta_{\text{merge}} \tilde{\mathcal{L}}_h^{(t)} > \text{tolerance}$  then

8:      $S_{\text{merge}} \cdot \text{add}(\{(I_h^{(left)}, I_h^{(right)}) : \Delta_{\text{merge}} \tilde{\mathcal{L}}_h^{(t)}\})$ 

9:   end if

10: end for

11: while  $S_{\text{merge}} \neq \emptyset$  do

12:    $(I_h^{(left)}, I_h^{(right)}) = S_{\text{merge}} \cdot \text{pop}(); I_h = \text{MERGE}(I_h^{(left)}, I_h^{(right)})$ 

13:   if  $I_h \cdot \text{sibling} \cdot \text{isLeaf}() = \text{TRUE}$  then

14:      $(G_h^{(sibling)}, H_h^{(sibling)}, \mathbf{L}_h^{(sibling)}) = \sum_{i \in I_h^{(sibling)}} (g_i, h_i, (L_{ik})_{k=0}^{K-2})$ 

15:      $\Delta_{\text{merge}} \tilde{\mathcal{L}}_{h'}^{(t)} = \text{EVALUATEMERGE}\left((G_h, H_h, \mathbf{L}_h), (G_h^{(sibling)}, H_h^{(sibling)}, \mathbf{L}_h^{(sibling)})\right)$ 

16:     if  $\Delta_{\text{merge}} \tilde{\mathcal{L}}_{h'}^{(t)} > \text{tolerance}$  then

17:        $S_{\text{merge}} \cdot \text{add}(\{(I_h, I_h^{(sibling)}) : \Delta_{\text{merge}} \tilde{\mathcal{L}}_{h'}^{(t)}\})$ 

18:     end if

19:   end if

20: end while

```

---

## Appendix B

$$G_j = \sum_{i \in I_j} g_i = \sum_{i \in I_j} \sum_{k=0}^{K-2} -s_{ik} \sigma(-s_{ik} d_{ik})$$

$$H_j = \sum_{i \in I_j} h_i = \sum_{i \in I_j} \sum_{k=0}^{K-2} \frac{2\sigma(s_{ik} d_{ik}) - 1}{2s_{ik} d_{ik}}$$

$$\mathbf{L}_j = \sum_{i \in I_j} \frac{\partial^2 l_i}{\partial \boldsymbol{\theta} \partial \tilde{\mathbf{y}}} = \begin{bmatrix} \dots & \sum_{i \in I_j} -\sigma(s_{ik} d_{ik})(1 - \sigma(s_{ik} d_{ik})) & \dots \end{bmatrix}^\top \in \mathbb{R}^{K-1}$$

$$\mathbf{u} = \sum_{i=1}^n \frac{\partial l_i}{\partial \boldsymbol{\theta}} = \sum_{i=1}^n \begin{bmatrix} \dots & s_{ik} \sigma(-s_{ik} d_{ik}) & \dots \end{bmatrix}^\top \in \mathbb{R}^{K-1}$$

$$\mathbf{V} = \sum_{i=1}^n \frac{\partial^2 l_i}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} = \text{diag} \left( \begin{bmatrix} \dots & \sum_{i=1}^n \frac{2\sigma(s_{ik} d_{ik}) - 1}{2s_{ik} d_{ik}} & \dots \end{bmatrix} \right) \in \mathbb{R}^{(K-1) \times (K-1)}$$

## APPENDIX B.

$$\begin{aligned}\sigma(x) &= \frac{1}{1 + e^{-x}} \\ s_{ik} &= \mathbb{S}(y_i - k) \\ d_{ik} &= \tilde{y}_i - \theta_k\end{aligned}\tag{B.1}$$