

Gradient Boosted Regression Trees

A brief introduction

This document provides a brief mathematical introduction to the idea of gradient boosting, specifically in the context of gradient boosted regression trees (GBRT). For a more fleshed out and properly sourced version, see the supplemental methods of the GRL publication doi 10.1002/2017GL075661.

A.K. April 2017

Contents

1 Problem Setting	1
2 Regression Trees	2
3 Boosted Regression	4
3.1 Boosted Regression as Gradient Descent	5
4 Gradient Boosted Regression Tree	7
4.1 Regularization	7
4.2 Stability	8

1 Problem Setting

Consider a typical regression problem in which one seeks to predict a continuous variable from a vector of continuous or categorical *features*. The *training data* is a collection of observations of the form (\mathbf{x}, y) where $\mathbf{x} = (x_1, \dots, x_p)$ is a feature vector belonging to a p -dimensional feature space and y is the quantity of interest for prediction. Given n such observations:

$$\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\},$$

we wish to find a *predictor* $f : \mathbb{R}^p \rightarrow \mathbb{R}$ relating feature vectors to the quantity of interest y . This is achieved by formulating the problem of finding f as an optimization problem over a pre-specified family \mathbb{F} of possible predictors. In other words, we seek an optimal (in a sense to be determined) function f in \mathbb{F} .

For a given predictor f we denote its predictions over the training data by:

$$\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n) := (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)),$$

which we will refer to as the *prediction vector*. For convenience, we define the evaluation mapping $E : \mathbb{F} \rightarrow \mathbb{R}^n$ where n is the number of training samples, which maps each predictor f to its prediction vector:

$$f \mapsto \hat{\mathbf{y}} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$$

Goodness of fit is formulated in terms of a *loss function* $L : \mathbb{R}^n \rightarrow \mathbb{R}^+$ which assigns a measure of error $L(\hat{\mathbf{y}})$ to the vector $\hat{\mathbf{y}}$ of n predictions with respect to observed values $\mathbf{y} = (y_1, \dots, y_n)$. We will be using the common *sum of squares* loss function, also commonly used in linear regression, which is given by:

$$L(\hat{\mathbf{y}}) = \sum_{k=1}^n |\hat{y}_k - y_k|^2$$

Finally, our task for fitting a predictor f to the given samples is to solve the optimization problem:

$$\min_{f \in \mathbb{F}} L(E(f))$$

The description of a statistical model involves the specification of the family \mathbb{F} of predictors together with an algorithm that finds the optimal predictor f in \mathbb{F} with respect to L . In the remainder of this section we first describe ordinary regression trees as such a model. We then specify GBRT as an additive (i.e boosted) extension of regression trees or, in a more general sense, as a form of gradient descent over the predictor family \mathbb{F} . For this discussion, we will assume a fixed set of n training samples $\{(\mathbf{x}_k, y_k)\}$ is given.

2 Regression Trees

A regression tree partitions the feature space into a collection of box-like regions D_i and assigns a constant value y to each region. In other words, each such region is the subset of the feature space constrained by conditions of the form $x_k > T$ or $x_k < T$ where each x_k is an individual feature. Alternatively, such predictors can be viewed as a decision trees in which each node is a comparison of the form $x_k > T$ on some feature x_k . Given a feature vector \mathbf{x} , a predictor descends down the decision tree by comparing feature values with the conditions at each node. The bottom-most nodes in the tree, namely the *leaves*, are then assigned constant values which are returned by the predictor (1)

The predictor family \mathbb{F} is the collection of all such piecewise constant functions. An individual predictor $f \in \mathbb{F}$ is identified by the specification of its box-like regions as well as the value it assigns to each region. In general, minimizing the loss function over the entire family \mathbb{F} is intractable and instead, an approximating algorithm is typically employed which motivates its interpretation as a decision tree (Algorithm 1).

Since individual features are treated separately in Algorithm 1, it is clear that regression trees pose no difficulty when the feature set contains both continuous and categorical variables.

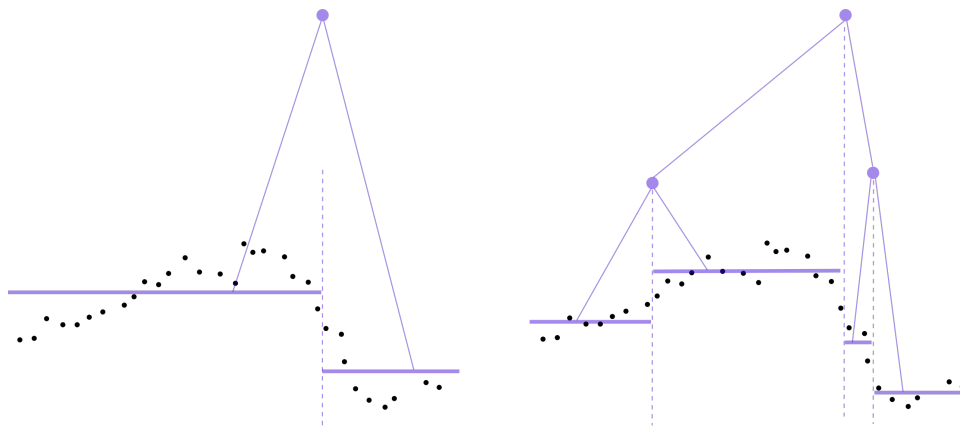


Figure 1: Example of a one dimensional regression tree after one (*left*) and two (*right*) iterations.

Algorithm 1 Ordinary Regression Tree

$D = \mathbb{R}^p$	▷ The starting region is the entire feature space
$f(D) \leftarrow \text{avg}(y)$	▷ Initialize with constant function
$j \leftarrow 0$	▷ current tree depth
while $j < J$ do	▷ J is the maximum tree depth
for region R do	▷ each region is a leaf in the tree
$k \leftarrow$ "best" feature $k = 1, \dots, p$ to split R on	
$T \leftarrow$ "best" threshold for feature x_k in R	
$D_1, D_2 \leftarrow$ partition D along the x_k axis with threshold T	
$f(D_i) \leftarrow \text{avg}(y)$ in D_i	▷ assign constants to each new region (leaf)
end for	
end while	
return f	

3 Boosted Regression

Boosting is a general heuristic in statistical learning which proposes combining multiple predictors, referred to as *weak learners*, into a stronger predictor. Given a family \mathbb{F} , a boosted predictor F based on \mathbb{F} is a linear combination of predictors in \mathbb{F} , that is:

$$F = \alpha_1 f_1 + \dots + \alpha_m f_m$$

for some $f_i \in \mathbb{F}$ and $\alpha_i \in \mathbb{R}$. When the family \mathbb{F} is closed under scalar multiplication (as is the case with regression trees), we can simplify the general form of a boosted predictor to:

$$F = f_1 + \dots + f_m$$

for some $f_i \in \mathbb{F}$. One can describe the problem of training a boosted model in similar terms as before: we seek the optimal predictor F in the boosted predictor family

$$\text{span}(\mathbb{F}) = \left\{ f_1 + \dots + f_m ; m \in \mathbb{N}, f_i \in \mathbb{F} \right\}$$

such that the training loss

$$\sum_{k=1}^n L(F(\mathbf{x}_k))$$

is minimized. In a boosted regression tree model the underlying predictor family \mathbb{F} consists of ordinary regression trees. However, there are many algorithms for approximating the optimal boosted predictor. Gradient boosting (Algorithm 2) is one such algorithm which is shown to contain many earlier boosting algorithms as a special case. In gradient boosting, at each stage a weak learner f_i is fit to the current *residual* vector which is the difference between known values \mathbf{y} and the current prediction vector $\hat{\mathbf{y}}$ produced by $f_1 + \dots + f_{i-1}$.

Our statistical model can thus be summarized as follows. Gradient boosting is used to approximate the optimal boosted regression tree $f_1 + \dots + f_m$ that minimizes the sum of squares loss.

Algorithm 2 Gradient Boosted Regression Tree

```

i ← 0                                ▷ number of weak learners trained so far
r ← y                                ▷ n-dimensional vector of residuals
while i < M do                    ▷ M is the number of weak learners to be trained
    i ← i + 1
    fi ← best weak learner fit to r    ▷ i.e with training data {(x1, r1), ..., (xn, rn)}
    fi = νf                          ▷ regularization via shrinkage (see discussion below)
    r̂ ← (fi(x1), ..., fi(xn))
    r ← r − r̂                        ▷ updated residuals
end while
return f1 + f2 + ... + fM

```

Boosting has proved to be a fruitful strategy in various settings on multiple fronts. First, by extending the family of predictors from \mathbb{F} to $\text{span}(\mathbb{F})$, it allows for lower bias and variance. For instance, if \mathbb{F} is the very limited family of step functions (i.e regression trees of depth 1), the extended class $\text{span}(\mathbb{F})$ can approximate any function arbitrarily closely. This, however

comes with the risk of overfitting which is mitigated through regularization (discussed below). Furthermore, through the statistical power of averaging, boosting has been shown to increase robustness to noise (see discussion of stability below).

3.1 Boosted Regression as Gradient Descent

In this section we describe an alternative, and more general, description of boosted regression in terms of a form of gradient descent in a function space. This description is the motivation for the name “gradient boosting”. However, we will see that under sum of squares loss the two are exactly identical.

First recall the familiar gradient descent framework: given a scalar field \mathcal{L} over a vector space V we wish to find the point $\mathbf{p} \in V$ that minimizes \mathcal{L} . We use the fact that at any point p the gradient of \mathcal{L} at \mathbf{p} , here denoted by $\nabla_{\mathbf{p}} \mathcal{L}$, points in the direction of steepest ascent. Therefore, an iterative improvement approach along the direction of steepest descent, is guaranteed to find the global minimizer of \mathcal{L} as long as \mathcal{L} is convex (which is the case here and in most statistical models):

Algorithm 3 Standard Gradient Descent

```

 $\mathbf{p} \leftarrow \mathbf{v}_1$  ▷ arbitrary starting point in  $V$ 
 $i \leftarrow 0$  ▷ the number of steps so far
while  $|\nabla_{\mathbf{p}} \mathcal{L}| > \epsilon$  ▷ alternatively, enforce an upper bound on  $i$ 
   $i \leftarrow i + 1$ 
   $\mathbf{u} = -\nabla_{\mathbf{p}} \mathcal{L}$  ▷ direction of steepest descent
   $\mathbf{v}_i = \operatorname{argmin}_{\mathbf{v}=\alpha\mathbf{u}} \mathcal{L}(\mathbf{v})$  ▷ line search in one variable  $\alpha$  (along  $\hat{u}$ )
   $\mathbf{p} \leftarrow \mathbf{p} + \mathbf{v}_i$ 
end while
return  $\mathbf{p} = \mathbf{v}_1 + \mathbf{v}_2 + \dots + \mathbf{v}_i$ 

```

Now instead of the usual finite dimensional vector space of calculus, let V be the vector space of functions $\operatorname{span}(\mathbb{F})$ which is the domain of optimization in boosted regression. In boosted regression trees \mathbb{F} is the family of piecewise constant functions of a specified maximum tree depth J ¹. The function to be minimized over V is $\mathcal{L} = L \circ E$ where E evaluates a predictor f at the training samples $\mathbf{x}_1, \dots, \mathbf{x}_n$ and the loss function L assigns an error to the prediction vector $\hat{\mathbf{y}} = E(f)$.

Recall that the gradient of a scalar field $\mathcal{L} : V \rightarrow \mathbb{R}$ at a point \mathbf{p} is the linear map

$$\nabla_{\mathbf{p}} \mathcal{L} : V \rightarrow \mathbb{R}$$

that assigns to any vector \mathbf{v} the rate of change of \mathcal{L} at \mathbf{p} along \mathbf{v} :

$$\mathbf{v} \mapsto \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{v}) := \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [\mathcal{L}(\mathbf{p} + \epsilon \mathbf{v}) - \mathcal{L}(\mathbf{p})]$$

In boosting, we have $V = \operatorname{span}(\mathbb{F})$ and the “gradient”² at any “point” $f \in V$ is a linear map $\nabla_f \mathcal{L} : V \rightarrow \mathbb{R}$ assigning to each “direction” ϕ (which is now a function in V) the

¹However, the discussion in this section is independent of the choice of \mathbb{F} as long as it is closed under scalar multiplication. In fact, *any* boosting additive model can be viewed as gradient descent in a function space as described in this section.

²technically known as the *Fréchet derivative* over a function space.

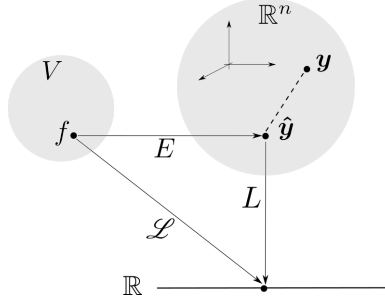


Figure 2: Schematic representation of the functional optimization problem to minimize $\mathcal{L} = L \circ E$ over the vector space $V = \text{span}(\mathbb{F})$.

directional rate of change:

$$\phi \mapsto \nabla_f \mathcal{L}(\phi) := \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [\mathcal{L}(f + \epsilon\phi) - \mathcal{L}(f)]$$

In vector calculus, V is a finite dimensional space, say \mathbb{R}^m and $\nabla_p \mathcal{L}$, being a linear map, can be represented as multiplication by a $m \times 1$ matrix. This matrix multiplication can be viewed as a dot product such that for any $v \in V$ the rate of change along v is $\nabla_p \mathcal{L} \cdot v$. Unfortunately, the dot product representation of the gradient does not generalize in a computationally tractable way to function spaces. However, in our case the structure of our scalar field $\mathcal{L} = L \circ E$ allows a similar formulation due to the simplicity of the evaluation function $E : V \rightarrow \mathbb{R}^n$. An application of the familiar chain rule for $\mathcal{L} = L \circ E$ (which is valid in infinite dimensional vector spaces) gives:

$$\nabla_f \mathcal{L} = \nabla_{E(f)} L \circ \nabla_f E$$

But it can easily be verified that E is linear and thus $\nabla_f E \equiv E$ for any f . Therefore, for any predictor $f \in V$ and any direction of change $\phi \in V$ we have:

$$\nabla_f \mathcal{L}(\phi) = \nabla_{E(f)} L \cdot E(\phi)$$

where $\nabla_{E(f)} L$ is the familiar n -dimensional gradient of $L : \mathbb{R}^n \rightarrow \mathbb{R}$ and the dot product is the familiar dot product in \mathbb{R}^n . Thus, at each step of gradient descent on $V = \text{span}(F)$ to minimize \mathcal{L} , the problem of finding the direction of steepest descent at is reduced to:

$$\max_{\phi \in \text{span}(F)} -\nabla_{E(f)} L \cdot E(\phi)$$

where $f = f_1 + f_2 + \dots + f_i$ is the current position (predictor) in V . However, since this is not in general tractable, the domain of search is typically restricted from $\text{span}(\mathbb{F})$ to \mathbb{F} . Furthermore, under regularity conditions, the problem can be equivalently formulated in terms of distance minimization:

$$\min_{\phi \in \mathbb{F}} \left\| -\nabla_{E(f)} L - E(\phi) \right\|$$

where $\|\mathbf{a} - \mathbf{b}\|$ is the usual Euclidean distance in \mathbb{R}^n . However, under sum of squares loss, the latter is precisely what weak learners optimize for. More specifically, a weak learner finds $\phi \in \mathbb{F}$ that minimizes the distance between $\hat{y} = E(\phi)$ and some known vector, which here is

$-\nabla_{E(f)} L$ (called *pseudo-residuals*). The “direction” of steepest descent within \mathbb{F} , therefore, can be obtained by fitting a weak learner to pseudo-residuals. Note that under the sum of squares L the pseudo-residuals are, up to a factor, identical to ordinary residuals of the vanilla boosting algorithm:

$$-\nabla_{E(f)} L = 2(\mathbf{y} - \hat{\mathbf{y}})$$

In other words, under sum of squares loss, the gradient descent generalization described here is exactly identical to the simple additive model described in the previous section.

Algorithm 4 Gradient Boosted Regression

```

 $f \leftarrow \text{avg}(y)$  ▷ arbitrary starting point in  $\mathbb{F}$ 
 $i \leftarrow 0$  ▷ the number of steps so far
while  $i < M$  do ▷  $M$  is the max number of steps
   $i \leftarrow i + 1$ 
   $\phi = \text{argmax}_{\phi \in \mathbb{F}} -\nabla_{E(f)} L \cdot E(\phi)$  ▷ direction of steepest descent
   $f_i = \text{argmin}_{f = \alpha \phi} \mathcal{L}(f)$  ▷ line search in one variable  $\alpha$  (along  $\phi$ )
   $f_i = \nu f_i$  ▷ shrinkage regularization, see discussion below
   $f \leftarrow f + f_i$ 
end while
return  $f = f_1 + \dots + f_M$ 

```

4 Gradient Boosted Regression Tree

The GBRT model is exactly the model described in the previous section with \mathbb{F} being the family of regression trees of maximum depth J . Note that since we are allowing boosted (additive) regression trees we can use small values for J . Furthermore, note that in functional gradient descent, unlike the standard gradient descent, the stopping criteria cannot be, in general, formulated in terms of $|\nabla_f \mathcal{L}| < \epsilon$ since that requires imposing additional mathematical structure on the function space. The typical stoppage scheme is to terminate the search after M steps. This means the output of a GBRT is a sum of M regression trees each of which has depth at most J .

4.1 Regularization

A crucial difference between optimization problems and statistical modelling is the risk of overfitting. In fact, typically one can find predictors that drive the total training loss to zero (Fig. 3).

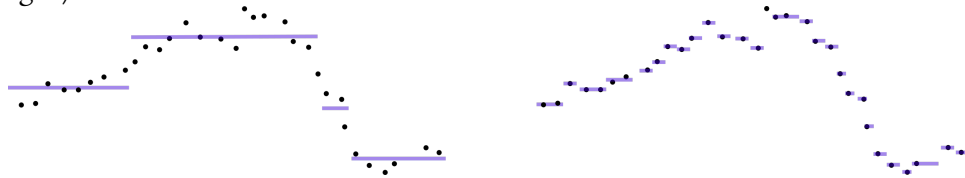


Figure 3: Arbitrarily small training loss with a regression tree (*right*) which is overfit to training data.

The standard scheme for containing the effects of overfitting (also known as *regularization*) for regression trees is to place an upper bound J on the decision tree depth. In gradient boosted trees, two additional regularization strategies are typically simultaneously used: 1) an upper bound M on the number of descent steps (discussed above) and 2) *shrinkage* which is used to scale down the contributions calculated in each gradient descent step by a parameter $0 < \nu \leq 1$. A choice of $\nu = 1$ implies no shrinkage, i.e non-regularized gradient descent, and smaller values of ν slow down the march towards arbitrarily small training losses, hence limiting the risk of overfitting.

4.2 Stability

An important desired characteristic of a statistical model is its *stability*, namely its robustness to noise. This is crucial in cases where one either has limited amounts of training data or data is known to be highly noisy. We wish our statistical model to be such that small perturbations in known labels y_1, \dots, y_n does not lead to large perturbations in the trained predictor. This is another domain where boosted regression trees are superior to ordinary regression trees: single trees are known to be unstable while boosted trees tend to be fairly stable.

In order to evaluate the stability of our model, we can add random noise to known y values of the training set and report the average difference between predictions under noise-free and noisy training data. This procedure is then repeated multiple times over random training and validation sets and with various noise amplitudes.

At various levels of noise we can repeat the standard cross-validation procedure with the addition of noise to training values.

The relative noise *amplitude* $0 \leq A \leq 1$ which is the ratio between the expected value of the added Gaussian noise η with respect to the average y value, namely $A = \mathbb{E}[\eta]/\bar{y}$ where \bar{y} is the average y over the training set. The extreme case of $A = 0$ corresponds to unperturbed training data.

For each partition of data to training and validation sets, multiple predictors should be considered. First, a baseline predictor f_0 is obtained from noise-free training data (i.e $A = 0$). Then, for each relative noise amplitude value A a predictor f_A is obtained from proportionally noisy data:

$$\left\{ \left(\mathbf{x}_k^{(t)}, y_k^{(t)} + \eta_k \right) \right\}_{k=1}^m$$

where

$$\frac{\eta_k}{\bar{y}} \sim \mathcal{N} \left(0, \frac{\pi A^2}{2} \right)$$

The perturbation in f_A is obtained by comparing its predictions over the validation set with those of the baseline predictor f_0 :

$$\|f_A - f_0\| = \left[\sum_{k=1}^{n-m} \left| f_A \left(\mathbf{x}_k^{(v)} \right) - f_0 \left(\mathbf{x}_k^{(v)} \right) \right|^2 \right]^{\frac{1}{2}}$$

This perturbation is then averaged over multiple rounds of partitioning the samples into training and validation sets.