

DATA SCIENCE II:
Machine Learning
MTH 9899
Baruch College

Lecture 4: Boosting, Neural Network Gradient Descent
Learning

Adrian Sisser

April 17, 2019

Outline

- 1 Boosting
 - Overview & AdaBoost
 - Gradient Boosting
 - XGBoost and lightGBM
- 2 Neural Networks
 - Optimizers
 - Activation Functions

Boosting

- Boosting is the idea of combining a set of weak learners to form a stronger learner.
- It's a type of ensemble, since it's a collection of learners, but is different in many important aspects.
- Boosting differs because it trains the learners in sequence, with a focus on errors made by the prior learners.

AdaBoost

AdaBoost was the first successful implementation of boosting. It applies to a binary classification problem:

$$Y \in \{-1, +1\}$$

We will learn weak learners, with a goal of combining them all into a final learner:

$$\begin{aligned} f_j &: X \rightarrow \{-1, +1\} \\ F_i(x) &= \sum_{j=1}^i \beta_j f_j(x) \end{aligned}$$

AdaBoost

function TRAINADABOOST(X, Y)

$W_1(i) \leftarrow \frac{1}{N}$ $\triangleright W$ are weights for the points

$F_0(X) = \sum_i W_0(i)y_i$

for $j \in 1..T$ **do** $\triangleright i$ indicates the Boosting Round

$f_j \leftarrow \text{TRAINWEAKLEARNER}(X, Y, W_j)$

$\epsilon_j \leftarrow \sum_i W_j(i)I(Y_i \neq f_j(x_i))$

$\alpha_j \leftarrow \frac{1}{2} \ln \frac{1-\epsilon_j}{\epsilon_j}$

$W_{j+1}(i) \leftarrow W_j(i) \exp(-\alpha_j y_i f_j(x_i))$

$W_{j+1}(i) \leftarrow \frac{W_j(i)}{\sum_k W_j(k)}$ \triangleright Normalize the weights to sum to 1

$F_j(X) \leftarrow F_{j-1}(X) + \alpha_j f_j(X)$

end for

return F_T

end function

So now we ask, where did those update rules come from? Let's define an error function for F_j :

$$\begin{aligned}
 E &= \sum_{i=1}^N e^{-y_i F_j(x_i)} && \text{Exponential Loss} \\
 &= \sum_{i=1}^N e^{-y_i (F_{j-1}(x_i) + \alpha_j f_j(x_i))} && \text{Using the def of } F_j \\
 &= \sum_{i=1}^N e^{-y_i F_{j-1}(x_i)} e^{-y_i \alpha_j f_j(x_i)} \\
 &= \sum_{f_j(x_i)=y_i} e^{-y_i F_{j-1}(x_i)} e^{-\alpha_j} + \sum_{f_j(x_i) \neq y_i} e^{-y_i F_{j-1}(x_i)} e^{\alpha_j} && \text{Split up into right/wrong examples} \\
 &= \sum_{i=1}^N e^{-y_i F_{j-1}(x_i)} e^{-\alpha_j} + \sum_{f_j(x_i) \neq y_i} e^{-y_i F_{j-1}(x_i)} (e^{\alpha_j} - e^{-\alpha_j})
 \end{aligned}$$

Now, we want to minimize the expression from before:

$$E = \sum_{i=1}^N e^{-y_i F_{j-1}(x_i)} e^{-\alpha_j} + \sum_{f_j(x_i) \neq y_i} e^{-y_i F_{j-1}(x_i)} (e^{\alpha_j} - e^{-\alpha_j})$$

In the algo, we said:

$$\begin{aligned} W_j(i) &= W_{j-1}(i) e^{-\alpha_{j-1} y_i f_{j-1}(x_i)} \\ &= [W_{j-2}(i) e^{-\alpha_{j-2} y_i f_{j-2}(x_i)}] e^{-\alpha_{j-1} y_i f_{j-1}(x_i)} \\ &= W_0 e^{\sum_{k=0}^{j-1} -\alpha_k y_i f_k} \\ &= e^{-y_i F_{j-1}(x_i)} \end{aligned}$$

By telescoping the prior statement ...

Now we can combine these:

$$E = \sum_{i=1}^N e^{-y_i F_{j-1}(x_i)} e^{-\alpha_j} + \sum_{f_j(x_i) \neq y_i} W_j(i) (e^{\alpha_j} - e^{-\alpha_j})$$

$$E = \sum_{i=1}^N e^{-y_i F_{j-1}(x_i)} e^{-\alpha_j} + \sum_{f_j(x_i) \neq y_i} W_j(i) (e^{\alpha_j} - e^{-\alpha_j})$$

The left part of the error above has no dependence on f_j . The term on the right is just the weighted sum of the incorrect examples according to f_j . Given that, we try to minimize the weighted sum of the incorrect points.

Let's look at how we choose α , the weight we use when we add the new learner f_j to the ensemble. We'll start with an earlier version of the equation, since it will make solving it easier.

$$\begin{aligned}
 E &= \sum_{f_j(x_i)=y_i} e^{-y_i F_{j-1}(x_i)} e^{-\alpha_j} + \sum_{f_j(x_i) \neq y_i} e^{-y_i F_{j-1}(x_i)} e^{\alpha_j} \\
 \frac{\partial E}{\partial \alpha_j} &= - \sum_{f_j(x_i)=y_i} e^{-y_i F_{j-1}(x_i)} e^{-\alpha_j} + \sum_{f_j(x_i) \neq y_i} e^{-y_i F_{j-1}(x_i)} e^{\alpha_j} \\
 0 &= - \sum_{f_j(x_i)=y_i} e^{-y_i F_{j-1}(x_i)} + e^{2\alpha_j} \sum_{f_j(x_i) \neq y_i} e^{-y_i F_{j-1}(x_i)} \\
 \alpha_j &= \frac{1}{2} \ln \frac{\sum_{f_j(x_i)=y_i} e^{-y_i F_{j-1}(x_i)}}{\sum_{f_j(x_i) \neq y_i} e^{-y_i F_{j-1}(x_i)}} \\
 \alpha_j &= \frac{1}{2} \ln \frac{1 - \epsilon_j}{\epsilon_j}
 \end{aligned}$$

We've now seen what AdaBoost is really doing. It's a greedy optimization at each iteration - where we train a learner to best fit the weighted points, and the weight of points increases as they are misclassified.

Gradient Boosting is another boosting technique. It's applicable to problems like regression. We'll take the same approach as we did with AdaBoost. At each step, we'll add a new weak learner into the overall predictor, and find the optimal weight for it.

- You can use any weak learner, typically, we use Regression Trees
- Like other boosting algorithms, it can be slower than a normal Random Forest, because it can't be parallelized.

function TRAINGB(X, Y)

$$F_0(X) = \bar{y}$$

$$E \leftarrow \mathcal{L}(F_0(X), Y)$$

for $j \in 1..T$ **do**

$$r_j \leftarrow -\frac{\partial E}{\partial F_{j-1}}$$

$$f_j \leftarrow \text{TRAINWEAKLEARNER}(X, r_j)$$

$$\alpha_j \leftarrow \arg \min_{\alpha} \mathcal{L}(F_{j-1}(X) + \alpha f_j(X), Y)$$

$$F_j(X) = F_{j-1}(X) + \alpha_j f_j(X)$$

$$E \leftarrow \mathcal{L}(F_j(X), Y)$$

end for

return F_T

end function

Let's look in more detail how this is the same as AdaBoost. In AdaBoost, our loss function was $e^{-y_i F(x_i)}$ and $E = \sum_{i=1}^N e^{-y_i F_j(x_i)}$. Since

$$r_j = -\frac{\partial E}{\partial F_{j-1}}$$

$$r_j = \sum_{i=1}^N y_i e^{-y_i F_{j-1}(x_i)}$$

$$r_j = \sum_{i=1}^N y_i W_j(i)$$

From earlier slide

So, we can see that r_j is just the weighted inputs from earlier.

What happens in the case of regression? Let's define an error function:

$$E = \sum_{i=1}^N \frac{1}{2} (y_i - F_j(i))^2$$

Now let's calculate r_{ij} (point i in the j th iteration):

$$\begin{aligned} r_{ij} &= -\frac{\partial E}{\partial F_{j-1}} && \text{For GB} \\ &= (y_i - F_{j-1}(i)) \end{aligned}$$

So for regression, r_{ij} is just the residual from the previous iteration!! Now, we just have to fit learners at each iteration to our residual, and add them in.

To add in a new residual, we need to find α , the weight.

$$E = \sum_{i=1}^N \frac{1}{2} (y_i - F_j)^2$$

$$E = \sum_{i=1}^N \frac{1}{2} (y_i - (F_{j-1}(i) + \alpha_j f_j(i)))^2$$

$$\frac{\partial E}{\partial \alpha_j} = \sum -f_j(i) (y_i - F_{j-1}(i) - \alpha_j f_j(i))$$

There is no easy way to solve this. Instead, in a general Gradient Boosting algorithm, the α value is effectively a learning rate - another hyperparameter that we have to tune.

Let's summarize what we just did for Gradient Boosting.

- After each iteration, we calculated the residual of each point.
- We fit a function to these residuals
- We added it into a composite learner with a small weight to improve the overall result

XGBoost

Last week, we talked about Extra Trees as one of the best 'off the shelf' modelling options. **XGBoost** is an open source package that is also, arguably, the best stock ML algo. XGBoost is short for Extreme Gradient Boosting and has a few interesting features:

- Very well implemented - FAST
- Distributed - It can be split up across a cluster.
- Custom Loss Functions

XGBoost

XGBoost takes care to reduce overfitting by regularizing - they add a penalty to avoid growing overly complex trees.

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

Here γ and λ are hyper parameters, T is the number of leaves, and w_j is the value in each leaf of the tree. The first component simple penalizes tree size. The second part is more interesting: We can think of it as penalizing us for vastly different values in the leaves - so we need a significant improvement in loss to justify it. *This is very similar to Ridge Regression or any other regularization.*

XGBoost

The next interesting technique in XGBoost is how they model the loss given a new function f_j to add in. We can take a Taylor expansion of the loss:

$$\begin{aligned}\mathcal{L}(Y, Y_j) &\approx \mathcal{L}(Y, Y_{j-1}) + \frac{\partial \mathcal{L}}{\partial Y_{j-1}} + \frac{1}{2} \frac{\partial^2 \mathcal{L}}{\partial Y_{j-1}^2}^2 \\ \mathcal{O} &\approx \mathcal{L}(Y, Y_{j-1}) + \frac{\partial \mathcal{L}}{\partial Y_{j-1}} + \frac{1}{2} \frac{\partial^2 \mathcal{L}}{\partial Y_{j-1}^2}^2 + \Omega(f) \\ \min &\quad \mathcal{L}(Y, Y_{j-1}) + \frac{\partial \mathcal{L}}{\partial Y_{j-1}} + \frac{1}{2} \frac{\partial^2 \mathcal{L}}{\partial Y_{j-1}^2}^2 + \Omega(f)\end{aligned}$$

Now, when we want to determine if we should add a new layer to the tree, we can see if the new split, will increase or decrease the overall objective based on just the first and

LightGBM

LightGBM is a newer package than XGBoost.

- It is typically significantly faster than XGBoost
- LightGBM introduced a histogram based approach to splitting which led to the speed boost. XGBoost now offers a similar splitting approximation.
- LightGBM uses a “leaf-wise” splitting approach, as opposed to “level-wise”. This can lead to deeper trees so overfitting is an issue.

Overall, lightGBM seems to offer slightly better average performance and run slightly faster.

AdaGrad

AdaGrad is another optimization technique, based on the idea that we don't update all parameters at the same rate. It also accounts for the fact that we want the learning rate to decrease over time.

First, let's define the sum of the squared gradients over time, G_i will be a diagonal matrix:

$$G_T = \text{diag} \sum_{j=1}^T \left(\frac{\partial E}{\partial \theta_j} \right)^2$$

Now, we learn with this, where ϵ is a small constant to avoid dividing by zero:

$$\theta_{i+1} = \theta_i - \frac{\eta}{\sqrt{G_T + \epsilon}} \odot \frac{\partial E}{\partial \theta_i}$$

What are the drawbacks?

RMSProp

RMSProp is a variant of AdaGrad (similar to AdaDelta, which we will omit for time considerations) developed by Hinton. We basically replace the sum of squared gradients from AdaGrad with an EMA (Exponentially Moving Average) version.

$$G_i = \lambda G_{i-1} + (1 - \lambda) \left(\frac{\partial E}{\partial \theta_i} \right)^2$$
$$\theta_{i+1} = \theta_i - \frac{\eta}{\sqrt{G_i + \epsilon}} \circ \frac{\partial E}{\partial \theta_i}$$

Adam

Adam is short for Adaptive Momentum Estimation that extends RMSProp and momentum. Remember that momentum kept a simple sum of squared gradient while RMS Prop kept an EMA of the squared gradient.

$$m_T = \beta_1 m_{T-1} + (1 - \beta_1) \frac{\partial E}{\partial \theta_T}$$
$$G_T = \beta_2 G_{T-1} + (1 - \beta_2) \left(\frac{\partial E}{\partial \theta_T} \right)^2$$

These exponential moving averages are very biased initially, so Adam corrects them:

$$m_T^* = \frac{m_T}{1 - \beta_1^T}$$
$$G_T^* = \frac{G_T}{1 - \beta_2^T}$$

Adam takes these two updates, and combines them similarly to AdaGrad:

$$\theta_{i+1} = \theta_i - \frac{\eta}{\sqrt{G_i^*} + \epsilon} m_i^*$$

And many more

- AdaDelta - A precursor to AdaGrad, based on a moving window of gradients ...
- Nadam - A blend of Nesterov Momentum and Adam
- AMSGrad is a variant of Adam that adapts to the fact that some minibatches are more useful than others, particularly in image recognition.

So far, we've mostly spoken about the sigmoid activation function, as it's the 'original' and has nice derivative properties:

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

Output Activation Functions For Classification

For classification outputs, we need a special activation function. The most popular choice is the Softmax activation:

$$\mathcal{A}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}}$$

This converts the output to a PDF across the different class of outputs.

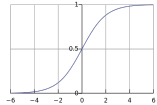
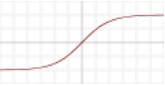
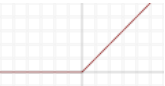
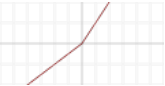
What are the requirements for a good activation function:

- Non-Linear - otherwise we're just fitting a linear model.
- Differentiable? - Not exactly - it just needs to be subdifferentiable
- Not easily saturated..

Does a step function work?

$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

No! It is NOT subdifferentiable, and the gradient is 0 for most of it, so we can't make any progress with Gradient Descent.

Name	$\mathcal{A}(x)$	$\frac{d\mathcal{A}}{dx}$	
Sigmoid	$\sigma(x)$	$\sigma(x)(1 - \sigma(x))$	
Tanh	$\tanh(x)$	$1 - \tanh^2(x)$	
Rectified Linear Unit (ReLU)	$\max(x, 0)$	$\begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	
Parameterized ReLU (aka 'Leaky' ReLU)	$\begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$\begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	

So why do we have 2 versions of ReLU?

ReLU

ReLU (and its variants) have become extremely popular:

- PReLU helps the fact that we have a large area with 0 gradient in 'vanilla' ReLU - this means that we can get 'stuck' in this region and the neuron becomes useless.
- ReLU has a gradient of 1 in its 'active' area. This means that if we build very deep networks, we don't have to worry about the gradient vanishing.

Vanishing Gradient

The “Vanishing Gradient Problem” is a well known one in training deep neural networks. The basic idea is that when you use a sigmoid activation function, for a very large area, the gradient is nearly 0. When we have a deep network, with many layers and we apply back propagation, by the time we get back to the early layers, the gradient is extremely small.

ReLU fixes this because the derivative when you are in the ‘active’ region of the neuron is 1.

This has been a fundamental change in NN and enabled much deeper networks.

Vanishing Gradient

There are other ways to fix the Vanishing Gradient Problem.

- Other activation functions that don't become saturated.
- Training methods - we can train a network layer by layer in an unsupervised manner.
- Residual Network (ResNet) - We feed the input from previous layer(s) plus the output of the previous hidden layer to each layer.