# Boosting

Jinwon Sohn

# Reference

- Greedy Function approximation, Friedman, 2001.

- Elementary of Statistical Learning, Springer.

# Numerical optimization
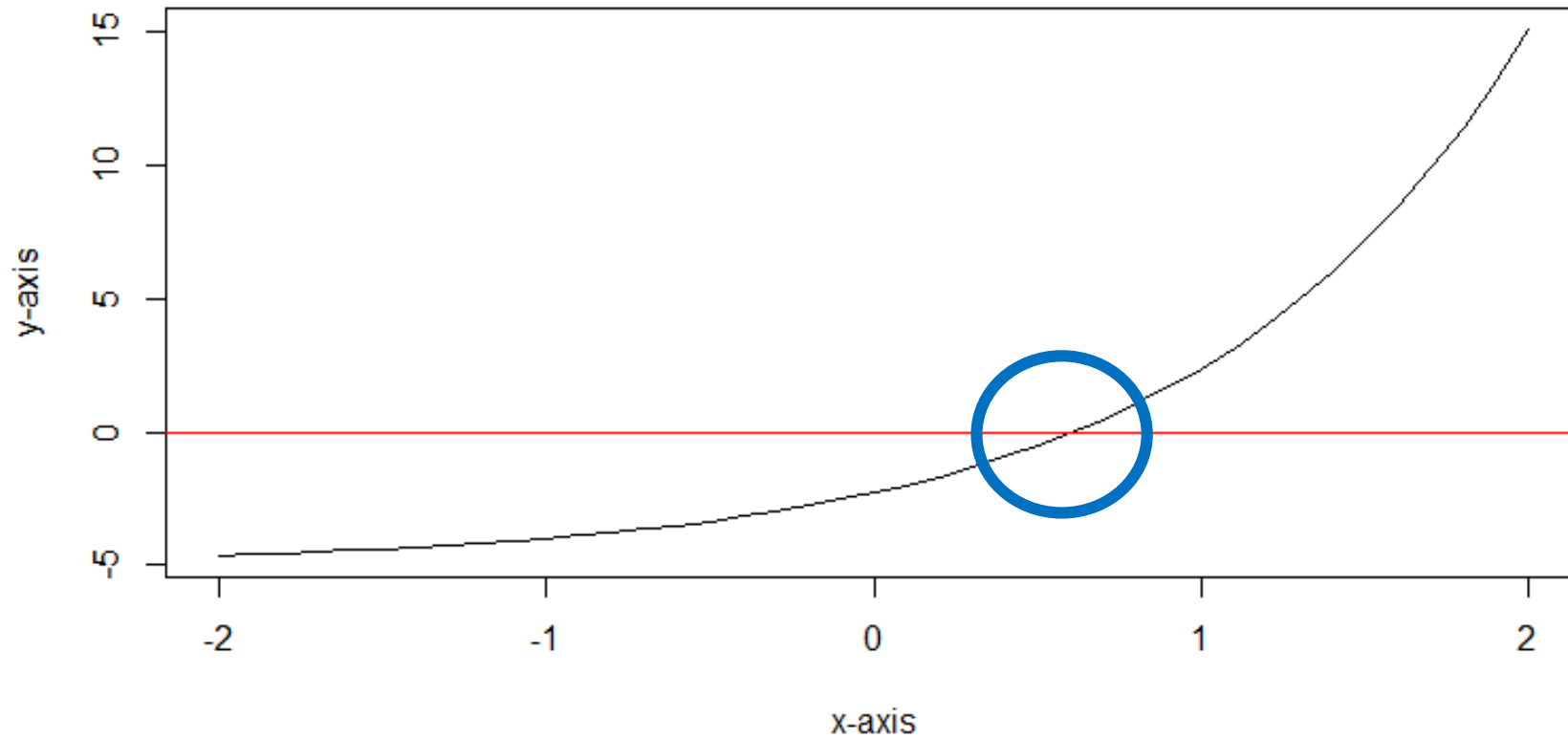
- In many problems, we can not obtain the closed solution for given equations.

- For example, the coefficients of logistic regression.

- Ex)Solving, $\arg\max\limits_{\beta_0,\beta_1} \prod_{i=1}^{N} \left(\frac{1}{1+e^{-(\beta_0+\beta_1 x_i)}}\right)^{r_i} \left(1 - \frac{1}{1+e^{-(\beta_0+\beta_1 x_i)}}\right)^{1-r_i}$

# No closed solutions exist!

# Newton's method

- $x_{n+1} = x_n - \dfrac{f(x_n)}{f'(x_n)}$

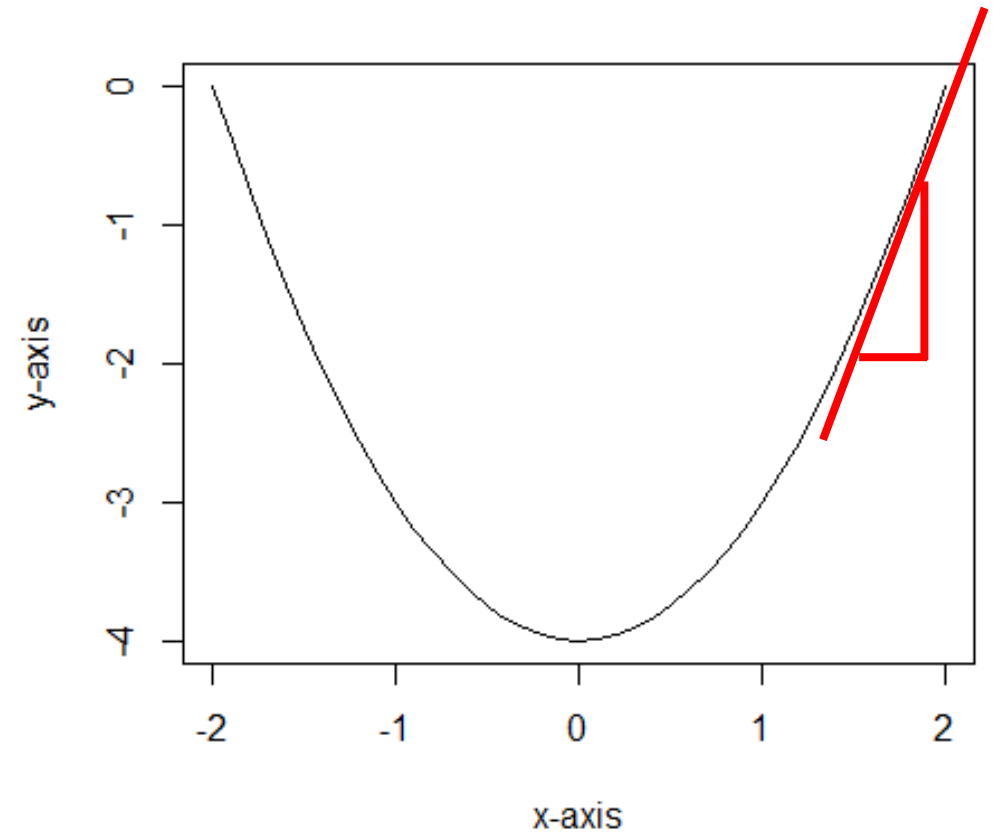- $x_{n+1} = x_n - \dfrac{f(x_n)(x_n - x_{n-1})}{f(x_n) - f(x_{n-1})}$

# Newton's method

# Newton's method

- The solution might be local optimum which is the usual problem in numerical optimization.
  > By using multiple initial points, we can bypass the problem


- Other mathematical properties..

# Gradient descent (or ascent)

- $x_{n+1} = x_n - \alpha \dfrac{df}{dx_n}$

- $\theta_{n+1} = \theta_n - \alpha \dfrac{dL}{d\theta_n}$

- $x_{n+1} = x_n - \alpha f'(x_n)$

- $x_{n+1} = x_n - \alpha \dfrac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$

# Line search gradient descent

- $x_{n+1} = x_n - \alpha_n \frac{df}{dx_n}$

- $\theta_{n+1} = \theta_n - \alpha_n \frac{dL}{d\theta_n}$

where $\alpha_n$ is $\alpha$ minimize a function $L\left(x_n - \alpha \frac{df}{dx_n}\right)$.

# Line search gradient descent

- $\theta_{n+1} = \theta_n - \alpha_n \frac{dL}{d\theta_n}$

- $n \longrightarrow \infty,\ \theta_n \longrightarrow \theta^*$ (Solution!)

- $\theta^* = \sum_{n=0}^{\infty} \left( \theta_0 - \alpha_n \frac{dL}{d\theta_n} \right)$ where $\theta_0$ is initial value set normally as 0.

- $\theta^* \approx \sum_{n=0}^{N} \left( -\alpha_n \frac{dL}{d\theta_n} \right)$ for some large $N$.

- The term, $-\alpha_n \frac{dL}{d\theta_n}$, is called 'boost' or 'step'. We will see this later.

# Boosting

- An ensemble model composed of the sum of week models.

- A stump or simple linear regression model are kinds of week model.

- In other word, the week model has the high bias.

- Let's start!

# Boosting

- So,

$$G_H(x) = \sum_{h=1}^{H} w_h g_h(x; \theta_h)$$

The learner, $G_H(x)$, averages the **week learners**, $g_h(x)$, with the weights $w_h$. In order to improve the performance of $F(x)$, $\Theta = \{\theta_1, \dots, \theta_H\}$ and $\omega = \{w_1, \dots, w_H\}$ have to be optimized.

# Boosting

- Namely, the following loss function should be minimized with respect to $\Theta, \omega$.

$$\underset{\Theta, \omega}{\arg\min} \sum_{i=1}^{N} L(y_i, G_H(x_i))$$

$$= \underset{\Theta, \omega}{\arg\min} \sum_{i=1}^{N} L\left(y_i, \sum_{h=1}^{H} w_h g_h(x; \theta_h)\right)$$

# Boosting

- However, it requires intensive computation.

- Imagine what if we have to find the optimum values in $|\Theta| \times |\omega|$-spaces.

- A simple alternative can approximate the loss function. We call "*Forward Stagewise Additive Modeling*".

# Forward Stagewise Additive Modeling

- Optimize the parameters one by one by moving in the forward direction.

- Let $f_h(x) = f_{h-1}(x) + wg(x; \theta)$.

- Then, $f_h(x)$ will be decided by optimizing $w$ and $\theta$ in terms of

$$\arg\min_{w,\theta} \sum_{i=1}^{N} L(y_i, f_h(x_i)), \qquad h = 1,2,\dots,H.$$

$$= \arg\min_{w,\theta} \sum_{i=1}^{N} L(y_i, f_{h-1}(x_i) + wg(x_i; \theta)), \quad h = 1,2,\dots,H.$$

# Forward Stagewise Additive Modeling

- Let $L(y, f(x)) = (y - f(x))^2$, squared loss.
- Let $f_0(x) = 0$, then
- $f_1(x)$ can be obtained by solving

$$\arg\min_{w,\theta} \sum_{i=1}^{N} L(y_i, f_0(x_i) + wg(x_i; \theta))$$

$$= \arg\min_{w,\theta} \sum_{i=1}^{N} (y_i - wg(x_i; \theta))^2$$

# Forward Stagewise Additive Modeling

- $f_2(x)$ can be obtained by solving

$$\underset{w,\theta}{\arg\min} \sum_{i=1}^{N} L(y_i, f_1(x_i) + wg(x_i; \theta))$$

$$\underset{w,\theta}{\arg\min} \sum_{i=1}^{N} (y_i - f_1(x_i) - wg(x_i; \theta))^2$$

$$\underset{w,\theta}{\arg\min} \sum_{i=1}^{N} \left(r_{1,i} - wg(x_i; \theta)\right)^2 \qquad \text{residual}$$

- Denote these optimized parameters as $w_1$ and $\theta_1$.

# Forward Stagewise Additive Modeling

- $f_3(x)$ can be obtained by solving

$$\arg\min_{w,\theta} \sum_{i=1}^{N} L(y_i, f_2(x_i) + wg(x_i; \theta))$$

$$\arg\min_{w,\theta} \sum_{i=1}^{N} (y_i - f_2(x_i) - wg(x_i; \theta))^2$$

$$\arg\min_{w,\theta} \sum_{i=1}^{N} (y_i - f_1(x_i) - w_1 g(x_i; \theta_1) - wg(x_i; \theta))^2$$

# Forward Stagewise Additive Modeling

$$\underset{w,\theta}{\arg\min} \sum_{i=1}^{N} \left( r_{1,i} - w_1 g(x_i; \theta_1) - w g(x_i; \theta) \right)^2$$

$$\underset{w,\theta}{\arg\min} \sum_{i=1}^{N} \left( r_{2,i} - w g(x_i; \theta) \right)^2 \qquad \textcolor{red}{\text{residual}}$$

- Denote these optimized parameters as $w_2$ and $\theta_2$.

# Forward Stagewise Additive Modeling

- Thus, $f_H(x)$ can be obtained by solving

$$\arg\min_{w,\theta} \sum_{i=1}^{N} L(y_i, f_{H-1}(x_i) + wg(x_i; \theta))$$

$$\arg\min_{w,\theta} \sum_{i=1}^{N} \left(r_{H-1,i} - wg(x_i; \theta)\right)^2 \qquad \text{residual}$$

- $f_H(x) = f_{H-1}(x) + w_H g(x_i; \theta_H).$

- Thus, the final model is $\sum_{h=1}^{H} f_h(x) = \sum_{h=1}^{H} w_H g(x_i; \theta_H).$

# AdaBoost

- Set $L\big(y, f(x)\big) = e^{(-yf(x))}$. 'exponential loss'.
- Set the base estimator $g_h = g(x; \theta_h)$ be a 'stump', decision tree with one depth. 'Boost'
- Weights of observations are considered. 'Adaptive!'
- If applying *FSAM* to above setting,
- then you can obtain following algorithm.
- Please, refer to the page 344 in ESL for the proof.

# AdaBoost

1. Initialize the observation weight $w_i = \frac{1}{N}, \quad i = 1, \dots, N.$

2. For $h = 1$ to $H$:
   (a) Fit a classifier $g_h$ to the training data using weights $w_i$.
   (b) Compute

$$err_h = \frac{\sum_i w_i I(y_i \neq g(x_i; \theta_h))}{\sum_i w_i} .$$

   (c) Compute $\alpha_h = \log\big((1 - err_h)/err_h\big).$

   (d) Set $w_i \leftarrow w_i e^{\alpha_h I(y_i \neq g(x_i; \theta_h))}, \quad i = 1, \dots, N.$
3. Output $G(x) = sign[\sum_h \alpha_h g(x; \theta_h)]$

# Gradient Boosting

- Greedy Function Approximation, Friedman, 2001.

- History...

# Gradient Boosting

- Consider previous steepest descent with line search algorithm.

- $\theta_{n+1} = \theta_n - \alpha_n \frac{dL}{d\theta_n}$

- $\theta^* \approx \sum_{n=0}^{N} \left( -\alpha_n \frac{dL}{d\theta_n} \right) = \sum_{n=0}^{N} p_n$ for some large $N$.

- The increment, $p_n = -\alpha_n \frac{dL}{d\theta_n}$ , is called 'boost' or 'step'.

# Gradient Boosting

- We can regard **a function or classifier $F(x)$ as a parameter**, and optimize it numerically. This means that numerical optimization is used to estimate nonparametric function.

- $F_{h+1} = F_h - \alpha_h \dfrac{dL}{dF_h}$

- $F^* \approx \sum_{h=0}^{H} \left( -\alpha_h \dfrac{dL}{dF_h} \right) = \sum_{h=0}^{H} f_h$  for some large $H$.

- The classifier is composed of many increment functions!

# Gradient Boosting

- Let $F_{m-1} = \sum_{h=0}^{m-1} f_h$ .

- Then $F_m = F_{m-1} + f_m$
  The increment $f_m$ consists of $-\alpha_m \dfrac{dL}{dF_{m-1}}$ where $-\dfrac{dL}{dF_{m-1}}$ is the steepest gradient and $\alpha_m$ is found via the line search algorithm.

$$\alpha_m = arg \min_{\alpha} L\left(y, F_{m-1}(x) - \alpha \frac{dL}{dF_{m-1}}\right)$$

# Gradient Boosting

- The convergence steps or sequences of GB implicitly have the concept of Forward Stagewise Additive modeling.

- Since the negative gradients for each step are defined only at the specific data points, we have to construct models to generate the negative gradients.

- For $m$-step,

$$g_m(x; \theta_m) \approx -\frac{dL}{dF_{m-1}}$$

# Gradient Boosting

- In addition, the $\alpha_m$ come to be a weight parameter for the $m$-th model.

- So, the $F^*(x)$ can be approximated as,

$$F^* \approx \sum_{h=1}^{H} \alpha_h g_h(x; \theta_h)$$

# Gradient Boosting

- For example,
  set $L_2$ loss function for a GB model.
  set initial guess $f_0(x) = 0$ or $f_0(x) = \bar{y}$.

- Remember that the gradient of $L_2$ loss function is

$$\frac{dL_2}{dF(x)} = 2\big(y - F(x)\big).$$

# Gradient Boosting

- $F_1(x)$ can be obtained by solving

Step 1 :  $\theta_1 = arg \min_{\theta} \sum_{i=1}^{N} \left( -\dfrac{dL_2}{dF_0} - g(x_i; \theta) \right)^2$

$\theta_1 = arg \min_{\theta} \sum_{i=1}^{N} (-2(y_i - F_0(x_i)) - g(x_i; \theta))^2$

<span style="color:red">Negative gradient</span>

Step 2 :  $\alpha_1 = arg \min_{\alpha} \sum_{i=1}^{N} (y_i - F_0(x_i) - \alpha g(x_i; \theta_1))^2$

Step 3 :  $F_1(x) = F_0(x) + \alpha_1 g(x_i; \theta_1)$

# Gradient Boosting

- $F_2(x)$ can be obtained by solving

Step 1 :
$$\theta_2 = arg \min_\theta \sum_{i=1}^N \left( -\frac{dL_2}{dF_1} - g(x_i; \theta) \right)^2$$

$$\theta_2 = arg \min_\theta \sum_{i=1}^N (-2(y_i - F_1(x_i)) - g(x_i; \theta))^2$$

$$\theta_2 = arg \min_\theta \sum_{i=1}^N (-2(y_i - F_0(x_i) - \alpha_1 g(x_i; \theta_1)) - g(x_i; \theta))^2$$

<span style="color:red">Fitting on residuals!?</span>

Step 2 :
$$\alpha_2 = arg \min_\alpha \sum_{i=1}^N (y_i - F_1(x_i) - \alpha g(x_i; \theta_2))^2$$

Step 3 :
$$F_2(x) = F_1(x) + \alpha_2 g(x_i; \theta_2)$$

# Gradient Boosting

- $F_H(x)$ can be obtained by solving

Step 1 : $\theta_H = arg \min_{\theta} \sum_{i=1}^{N} \left( -\dfrac{dL_2}{dF_{H-1}} - g(x_i; \theta) \right)^2$

$\theta_H = arg \min_{\theta} \sum_{i=1}^{N} (-2(y_i - F_{H-1}(x_i)) - g(x_i; \theta))^2$

$\theta_H = arg \min_{\theta} \sum_{i=1}^{N} (-2(y_i - F_{H-2}(x_i) - \alpha_{H-1}g(x_i; \theta_{H-1})) - g(x_i; \theta))^2$

<span style="color:red">Fitting on residuals!?</span>

Step 2 : $\alpha_H = arg \min_{\alpha} \sum_{i=1}^{N} (y_i - F_{H-1}(x_i) - \alpha g(x_i; \theta_H))^2$

Step 3 : $F_H(x) = F_{H-1}(x) + \alpha_H g(x_i; \theta_H)$

# Gradient Boosting

- What is the significance of Gradient Boosting?

- Seemingly, it is equal to basic boosting in $L_2$ loss function.

- However, If we use $L_1$, or Huber loss function, GB has a more general applications.

- That's why we call the boost(or step) as the negative gradient, not residual even in $L_2$ loss function.

# Gradient Boosting

- $L_1$ loss function and its derivative.

$$L_1 = \sum_{i=1}^{N} |y_i - F(x_i)|, \qquad \frac{dL_1}{dF(x)} = sign(y - F(x))$$

This loss function is robust to outliers.

- Huber loss function and its derivative
  please, refer to https://en.wikipedia.org/wiki/Huber_loss

- You can customize your own loss function.

# Gradient Boosting

- You should specify the size of tree for each increment.

- Heuristically, $4 \leq d \leq 8$.

- The depth, $d$, reflect the order of an interaction!

- If a tree has 3 depth, then the tree includes not only main effect, but also up to third interactions.

# Gradient Boosting

- In classification problem, use the following setting.

$$L(\{y_k, F(x_k)\}_1^K) = -\sum_{k=1}^{K} y_k \log p_k(x)$$

$$F_k(x) = \log p_k(x) - \frac{1}{K} \sum_{k=1}^{K} y_k \log p_k(x)$$

# Extended Gradient Boosting

- Stochastic Gradient Booting
  > Subsampled data without replacement is used to fit a increment function. This do lighter computations.
  > A bit of resistance on overfitting

- Regularized Gradient Booting
  > Charge penalties onto the number of trees.

$$F_m(x) = F_{m-1}(x) + \nu_{penalty}\alpha_m g(x_i; \theta_m)$$

# XGBoost

- Gradient boosting needs burdensome computations which makes it difficult to apply GB models into Big data.

- T. Chen wrote in his paper how to improve computation ability by adjusting algorithms to build a tree.

- Refer to "https://gentlej90.tistory.com/87"

# LightGBM

- Structural difference.

- For more details, refer to https://towardsdatascience.com/catboost-vs-light-gbm-vs-xgboost-5f93620723db

# CatBoost

- Feature engineering on categorical predictors.

- For more details, refer to https://towardsdatascience.com/catboost-vs-light-gbm-vs-xgboost-5f93620723db