

GD Regression:

Step 0: The data and the loss function

Input Data  $\{(x_i, y_i)\}_{i=1}^n$  and a differentiable Loss Function  $L(y_i, F(x))$

- $x$  = Features
- $y$  = Target
- $n$  = Number of rows
- $y_i$  = Observed Value (Target)
- $F(x)$  = Predicted Value (Target)

$$\text{Loss Function} = \frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

$$\frac{d}{d \text{ Predicted}} \frac{1}{2} (\text{Observed} - \text{Predicted})^2 = -(\text{Observed} - \text{Predicted})$$

Step 1: Initialize model with a constant value:

$$F_0(x) = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, \gamma)$$

- $\gamma$  = Predicted Value (Target)
- argmin over gamma means we need to find a Predicted Value that minimizes this sum
- $F_0(x)$  = Initial Predicted Value that predicts that all sample will equal  $F_0(x)$
- Given this Loss Function,  $F_0(x)$  = Average of all the Observed Values

Step 2: Built M trees

For  $m = 1$  to  $M$

- $m$  = index tree
- $M$  = number of trees

### (A) Calculate Residuals

Compute  $r_{im} = -\left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}\right]_{F(x)=F_{m-1}(x)}$  For  $i = 1, \dots, n$

- $\left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}\right]_{F(x)=F_{m-1}(x)}$  = the derivative of the Loss Function
- $r$  = residual
- $i$  = sample number
- $r_{im}$  = residual for each sample in a specific tree (pseudo residual)

$r_{im} = (\text{Observed} - \text{Predicted})$

### (B) Fit a regression tree to the residuals

Fit a regression tree to the  $r_{im}$  values and create terminal region  $R_{jm}$  for  $j = 1 \dots J_m$

- $j$  = index leaf
- $J_m$  = Total number of leaves
- $R_{jm}$  = Leaves

### (C) Optimize leaf output values

For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \underset{\gamma}{\operatorname{argmin}} \sum_{x \in R_{jm}} L(y_i, F_{m-1}(x) + \gamma)$

$$\gamma_{jm} = \underset{\gamma}{\operatorname{argmin}} \sum_{x \in R_{jm}} \frac{1}{2} (y_i - (F_{m-1}(x) + \gamma))^2$$

The Output Values for  $\gamma_{jm}$  is always the average of the Residuals that end up in the same leaf

(D) Update predictions with the new tree

$$\text{Update } F_{(m)}(x) = F_{m-1}(x) + \eta \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$$

Step 3: Output the final prediction

$$\text{Output } F_M(x) = F_0(x) + \eta \sum_{j=1}^{J_1} \gamma_{j1} I(x \in R_{j1}) + \eta \sum_{j=1}^{J_2} \gamma_{j2} I(x \in R_{j2}) + \eta \sum_{j=1}^{J_M} \gamma_{jM} I(x \in R_{jM})$$

GB Classification:

Step 0: The data and the loss function

Input Data  $\{(x_i, y_i)\}_{i=1}^n$  and a differentiable Loss Function  $L(y_i, F(x))$

- $x$  = Features
- $y$  = Target
- $n$  = Number of rows
- $y_i$  = Observed Value (Target)
- $F(x)$  = Predicted Value (Target)

$$\text{Log(likelihood)} = \sum_{i=1}^N y_i \log(p) + (1 - y_i) \log(1 - p)$$

$$\text{Log Loss} = -\text{Log(likelihood)} \text{ or negative log(likelihood)} = -\sum_{i=1}^N y_i \log(p) + (1 - y_i) \log(1 - p)$$

$$\text{odds} = \frac{p}{1-p}$$

$$p = \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}}$$

$$\text{Loss Function} = -y_i \log(\text{odds}) + \log(1 + e^{\log(\text{odds})})$$

$$\frac{d}{d \log(\text{odds})} - y_i \log(\text{odds}) + \log(1 + e^{\log(\text{odds})})$$

$$= -y_i + \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}} \text{ or } -y_i + p$$

Step 1: Initialize model with a constant value:

$$F_0(x) = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, \gamma)$$

- $\gamma = \log(\text{odds})$  value
- argmin over gamma means we need to find a  $\log(\text{odds})$  that minimizes this sum
- $F_0(x)$  = Initial Predicted Value that predicts that all sample will equal  $F_0(x)$
- Given this Loss Function,  $F_0(x) = \log \frac{\text{Number of Yes}}{\text{Number of No}}$

Step 2: Built M trees

For  $m = 1$  to  $M$

- $m$  = index tree
- $M$  = number of trees

(A) Calculate Residuals

Compute  $r_{im} = -\left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}\right]_{F(x)=F_{m-1}(x)}$  For  $i = 1, \dots, n$

- $\left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}\right]_{F(x)=F_{m-1}(x)}$  = the derivative of the Loss Function
- $r$  = residual
- $i$  = sample number
- $r_{im}$  = residual for each sample in a specific tree (pseudo residual)

$$r_{im} = (\text{Observed} - p)$$

(B) Fit a regression tree to the residuals

Fit a regression tree to the  $r_{im}$  values and create terminal region  $R_{jm}$  for  $j = 1 \dots J_m$

- $j$  = index leaf
- $J_m$  = Total number of leaves
- $R_{jm}$  = Leaves

(C) Optimize leaf output values

For  $j = 1 \dots J_m$  compute  $\gamma_{jm} = \underset{\gamma}{\operatorname{argmin}} \sum_{x \in R_{ij}} L(y_i, F_{m-1}(x) + \gamma)$

$$\gamma_{jm} = \underset{\gamma}{\operatorname{argmin}} \sum_{x \in R_{ij}} -y_i [F_{m-1}(x_i) + \gamma] + \log(1 + e^{F_{m-1}(x_i) + \gamma})$$

$$\gamma_{jm} = \frac{\sum_{j=1}^{J_m} \text{Residual}}{\sum_{j=1}^{J_m} p(1-p)} I(j \in R_{jm})$$

(D) Update predictions with the new tree

$$\text{Update } F_{(m)}(x) = F_{m-1}(x) + \eta \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$$

Step 3: Output the final prediction

$$\text{Output } F_M(x) = F_0(x) + \eta \sum_{j=1}^{J_m} \gamma_{j1} I(x \in R_{j1}) + \eta \sum_{j=1}^{J_m} \gamma_{j2} I(x \in R_{j2}) + \eta \sum_{j=1}^{J_m} \gamma_{jM} I(x \in R_{jM})$$

The process:

1. Initial predicted value
2. Loss function
3. Residual (negative gradient of loss function) for each observation
4. Build a first weak learner
5. Added into a model:  $F(x) = F(x_0) + \eta$  (first weak learner)
6. Calculate new prediction from the model for each observation
7. Calculate new residual for each observation
8. Build a second weak learner
9. Added into a model:  $F(x) = F(x_0) + \eta$  (first weak learner) +  $\eta$  (second weak learner)