

Gradient Boost

→ Gradient Boost builds the model in a stage-wise fashion like other boosting methods do, and it generalizes them by adding the feature of optimizing an arbitrary differentiable loss function.

REGRESSION

→ Input: $\left\{ (x_i, y_i) \right\}_{i=1}^n \rightarrow \text{Data}$
 $L(y_i, F(x)) \rightarrow \text{Loss function.}$

Most ~~common~~ common loss function = $\frac{1}{2} (\text{observed} - \text{Predicted})^2$

→ Step 1: Initialize model with a const value

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

for, $\gamma = \text{avg}(\text{all elements})$, Loss function is least.

Thus, Prediction = average of all values for each

→ Step 2:

for $m = 1$ to M [M : NO of trees (8-32) or more]

(A) $r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial (F(x_i))} \right] \quad \text{for } i=1 \dots n.$

$F(x) = F_{m-1}(x)$

$$= - \frac{d}{d \text{ Pred.}^2} \frac{1}{2} (\text{observed} - \text{Pred}^2)^2 = \text{observed} - \text{Predicted}$$

Thus this calculates residuals.

but, γ_m is called pseudo-residuals.

(B) Greedy tree to fit residuals NOT outputs

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

↳ Argmin of loss function takes average of values

↳ this means convert all leaves of tree to average of all values in that leaf.

(D) $F_m(x) = F_{m-1}(x) + \gamma \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

New pred = old pred + learning rate \times value of residual from tree traversal

NOTE! Step 2 is done M times

→ Step 3! Output $F_m(x)$

NOTE! Gradient boost have full sized trees not just stumps.

CLASSIFICATION

Input: Data: $\{(x_i, y_i)\}_{i=1}^n$

Loss fn: $- [\text{Observed} \cdot \log(p) + (1 - \text{Observed}) \cdot \log(1-p)]$

$= - [\text{Observed} \times \log(\text{odds})] + \log(1 + e^{\log(\text{odds})})$

$\frac{d}{d(\log \text{odds})} (\text{Loss function}) = -\text{Observed} + \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}}$

$= -\text{Observed} + p$

Step 1:

Initialize model with constant value

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

$[\gamma = \log(\text{odds})]$

$p =$ This is again average of all values, $1 \rightarrow \text{yes}$
 $0 \rightarrow \text{no}$

$\log(\text{odds}) = \log\left(\frac{p}{1-p}\right)$

$F_0(x) = \log(\text{odds})$

Step 2:

for $m = 1$ to M .

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

pseudo residual $= \text{Observed} - \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}}$

$$= \text{Observed} - p \quad [\text{Any one}]$$

Calculate for entire dataset.

(B)

→ Build regression tree to residuals not output

(C) for $j=1 \dots J_m$ compute,

$$\gamma_{jm} = \underset{\gamma}{\operatorname{argmin}} \sum_{x_i \in R_j} L(y_i, F_{m-1}(x_i) + \gamma)$$

This reduces to,

$$\frac{\sum_{j=1}^{J_m} \text{Residuals}_j}{\sum_{j=1}^{J_m} P_j(1-P_j)}$$

[In a particular leaf]

[$P \rightarrow$ last prediction]

(D)

$$F_m(x) = F_{m-1}(x) + \gamma \sum_{j=1}^{J_m} \gamma_{jm} \mathbb{I}(x \in R_{jm})$$

New pred = Last pred + Learning Rate \times Output from tree

This is in log(odds)

Do everything until $m=M$

→ Step 3:

Output

$F_m(x)$ if $P > 0.5 \rightarrow 1$ else $\rightarrow 0$