

GBDT Paper Summary

1. Problem Setting \longrightarrow Function Estimation

$Y \rightarrow$ response variable
 $X \rightarrow$ explanatory variable

① \swarrow

\longrightarrow We hope to do: $F^* = \underset{F}{\operatorname{argmin}} \mathbb{E}_{x,y} [L(y, F(x))]$

$$= \underset{F}{\operatorname{argmin}} \mathbb{E}_x [\mathbb{E}_y [L(y, F(x)) | X]]$$

intractable !!!

\longrightarrow We do not the exact distribution
 $(X, Y) \sim P$

\longrightarrow it is also intractable to estimate
the conditional expectation $\mathbb{E}_y [\cdot | X]$

Pf Sketch $\mathbb{E}_{x,y} [L] = \int_x \int_y L(y, F(x)) p(x,y) dx dy$
 $= \int_x p(x) \cdot \int_y L(y, F(x)) \cdot p(y|x) dy dx$
 $= \int_x \mathbb{E}_y [L | X=x] p(x) dx$
 $= \mathbb{E}_x [\mathbb{E}_y [L(y, F(x)) | X]]$

Rmk1: 1. MSE Optimality $\Rightarrow F(x) = \mathbb{E}[Y | X=x]$
 \hookrightarrow if $L(y, F(x)) = (y - F(x))^2$
 \longleftrightarrow Linear Regression (Wasserman Lec 26)

2. the previous result is unconstrained result

actually what we do is CONSTRAINED (under some specific parametric model) i.e., tree model etc.

2. Optimization Framework (Infinite Data) \longrightarrow expectation case

a) Parametric Perspective

① Parametrized Model $F(x; P) = F(x, \{\beta_m, \alpha_m\}_{m=1}^M)$
 $= \sum_{m=1}^M \beta_m h(x; \alpha_m)$

② $\Phi(P) := \mathbb{E}_{(x,y)} [L(y, F(x; P))]$

$P^* = \underset{P}{\operatorname{argmin}} \Phi(P) \Rightarrow \underline{F^*(x) = F(x; P^*)}$

③ From the optimization perspective, we always have:

$$p^* = \sum_{m=0}^M p_m \Rightarrow \text{Gradient Descent Framework etc.}$$

optimization perspective

$$g_m = \nabla_p \Phi(p) \Big|_{p=p_{m-1}}$$

$$p_{m-1} = \sum_{i=0}^{m-1} d_i$$

multi-dim

$$d_m = -p_m g_m \quad \text{and} \quad p_m = \underset{p}{\operatorname{argmin}} \Phi(p_{m-1} - p g_m)$$

Rmk: Here, we try to characterize the function family through parametric form, i.e.,

$$f(x) = \sum_{i=1}^K \alpha_i \mathbb{1}\{x \in R_i\} \rightarrow \text{Decision Tree Model}$$

Then formulate the problem as optimization problem:

$$\rightarrow \hat{\omega} = \underset{\omega \in \Omega}{\operatorname{argmin}} R_{\text{pop}}(\omega)$$

b) Non-parametric Perspective

$$\begin{aligned} \textcircled{1} \quad \Phi(F) &= \mathbb{E}_{(x,y)} [L(y, F(x))] \\ &= \mathbb{E}_x [\mathbb{E}_y [L(y, F(x)) | x]] \\ &:= \mathbb{E}_x [\phi(F(x))] \end{aligned}$$

$$\phi(F(x)) = \mathbb{E}_y [L(y, F(x)) | x]$$

$$F^* = \underset{F}{\operatorname{argmin}} \Phi(F) \Leftrightarrow F^*(x) = \underset{F(x)}{\operatorname{argmin}} \phi(F(x)) \quad \text{for each } x$$

$$\Leftrightarrow F^*(x) = \operatorname{argmin} \mathbb{E}_y [L(y, F(x)) | X=x]$$

② Note:

→ for the above problem, there are infinitely many parameters

→ But for data sets, there are only a finite number

$\{F(x_i)\}_{i=1}^N$ involved!

→ Discussed Below!

③ Infinitely many parameters Scenario

→ optimization

Recap: $F^*(x) = \operatorname{argmin} \mathbb{E}_y [L(y, F(x)) | X=x]$
 $= \operatorname{argmin} \phi(F(x))$ for individual x

GD Framework → intractable

$$\Rightarrow F^*(x) = \sum_{m=0}^M f_m(x)$$

$$f_m(x) = -\rho_m g_m(x)$$

where $g_m(x) = \frac{\partial \phi(F(x))}{\partial F(x)} \Big|_{F(x)=F_{m-1}(x)}$ $F_{m-1}(x) = \sum_{i=0}^{m-1} f_i(x)$

\downarrow
1-dim
 $= \frac{\partial \mathbb{E}_y [L(y, F(x)) | x]}{\partial F(x)} \Big|_{F(x)=F_{m-1}(x)}$

under
regularity cond.
 $\mathbb{E}_y \left[\frac{\partial L(y, F(x))}{\partial F(x)} \Big| x \right]_{F(x)=F_{m-1}(x)}$

$$\rho_m = \operatorname{argmin}_{\rho} \mathbb{E}_{(x,y)} [L(y, F_{m-1}(x) - \rho g_m(x))]$$

3. Finite Data Case \rightarrow Tractable Scenario

① Greedy strategy + empirical loss

$$\left\{ \begin{array}{l} (\beta_m, \alpha_m) = \underset{\beta, \alpha}{\operatorname{argmin}} \sum_{i=1}^N L(y_i, F_{m-1}(x_i) + \beta h(x_i, \alpha)) \\ F_m(x) = F_{m-1}(x) + \beta_m h(x; \alpha_m) \end{array} \right.$$

\rightarrow when $\left\{ \begin{array}{l} y_i \in \{-1, +1\}, h(x; \alpha) \in \{-1, +1\} \\ L(y, F) = \exp\{-yF\} \end{array} \right.$

then this is exactly AdaBoost!

② Stage-wise Strategy to find (β_m, α_m)

Motivation \rightarrow in previous (non-parametric) discussion, we finally express $F^*(x)$ as $F^*(x) = \sum_{m=0}^M f_m(x)$

$$= \sum_{m=0}^M -\rho_m g_m(x)$$

Issue Here, we do not have complete expected gradient

$$g_m(x) = \frac{\partial (\phi(F(x)))}{\partial F(x)} \Bigg|_{F(x) = F_{m-1}(x)}$$

$$= \frac{\partial (\mathbb{E}_y[L(y, F(x)) | x])}{\partial F(x)} \Bigg|_{F(x) = F_{m-1}(x)}$$

What we have is:

$$g_m(x_i) = \frac{\partial (L(y, F(x)))}{\partial F(x)} \quad \Bigg| \quad F(x) = F_{m-1}(x_i)$$

Solution $F_m(x) = F_{m-1}(x) + \beta_m \cdot h(x; d_m)$

$$\begin{cases} 1. \alpha_m = \underset{\alpha, \beta}{\operatorname{argmin}} \sum_{i=1}^N (-g_m(x_i) - \beta h(x_i; d_m))^2 \\ 2. \rho_m = \underset{\rho}{\operatorname{argmin}} \sum_{i=1}^N L(y_i, F_{m-1}(x_i) + \rho h(x_i, d_m)) \end{cases}$$

GBDT Algorithm

4. Application

Here $L(y, F) = (y - F)^2/2$. The pseudo-response in line 3 of Algorithm 1 is $\tilde{y}_i = y_i - F_{m-1}(x_i)$. Thus, line 4 simply fits the current residuals and the line search (line 5) produces the result $\rho_m = \beta_m$, where β_m is the minimizing β of line 4. Therefore, gradient boosting on squared-error loss produces the usual *stagewise* approach of iteratively fitting the current residuals:

Algorithm 2: LS_Boost

$F_0(\mathbf{x}) = \bar{y}$

For $m = 1$ to M do:

$\tilde{y}_i = y_i - F_{m-1}(\mathbf{x}_i), \quad i = 1, N$

$(\rho_m, \mathbf{a}_m) = \underset{\mathbf{a}, \rho}{\operatorname{argmin}} \sum_{i=1}^N [\tilde{y}_i - \rho h(\mathbf{x}_i; \mathbf{a})]^2$

$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m h(\mathbf{x}; \mathbf{a}_m)$

endFor

end Algorithm

12.05

→ idea diagram (non-trivial)

GBDT

ideally, $F^* = \operatorname{argmin}_F \mathbb{E}_{(x,y)} [L(y, F(x))]$

break down to

→ parametric:

$$\Rightarrow p^* = \operatorname{argmin}_p \mathbb{E}_{(x,y)} [L(y, f(x; p))]$$

★ Non-parametric: → our interest

$$\begin{aligned} \Rightarrow F^* &= \operatorname{argmin}_F \mathbb{E}_{(x,y)} [L(y, F(x))] \\ &= \operatorname{argmin}_F \mathbb{E}_x [\mathbb{E}_y [L(y, F(x)) | x]] \end{aligned}$$

optimization framework

$$\begin{aligned} \Rightarrow F^*(x) &= \sum_{m=0}^M f_m(x) \\ &= F_{m-1}(x) - \rho_m g_m(x) \\ &\rightarrow g_m(x) = \frac{\partial (\mathbb{E}_y [L(y, F(x)) | x])}{\partial F(x)} \Big|_{F(x) = F_{m-1}(x)} \end{aligned}$$

Still infeasible



Surrogate: Finite Data Case

Motivation: we want to approximate $g_m(x) = \frac{\partial \phi(F(x))}{\partial F(x)} \Big|_{F(x) = F_{m-1}(x)}$

Solution: what we have $\Rightarrow \{g_m(x_i)\}_{i=1}^N \rightsquigarrow$ N-discrete point

$$g_m(x_i) = \frac{\partial (L(y_i, F(x_i)))}{\partial F(x_i)}$$

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

- \Rightarrow {

 ① train $h(x; \alpha_m)$ to fit $\{g_m(x_i)\}_{i=1}^N$

 ② then $h(x; \alpha_m) \approx g_m(x) \rightarrow$ our aim

Last Step: select optimal step length ρ_m through Loss Funcⁿ:

$$\rho_m = \arg \min_{\rho} \sum_{n=1}^N L(y_n, F_{m-1}(x_n) + \rho h(x_n; \alpha_m))$$

XGBoost (Boosting Algorithm)



- { An approximation on Empirical Loss (Taylor)

Tree Structure

Consider the model defined as follows :

$$\begin{aligned}
 \rightarrow F_t(x) &= \sum_{l=1}^t f_l(x) \\
 &= F_{t-1}(x) + f_t(x)
 \end{aligned}$$

$$\Rightarrow \hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + f_t(x_i)$$

$$\textcircled{2} g_i := \nabla_{\hat{y}_i^{(t-1)}} L(y_i, \hat{y}_i^{(t-1)})$$

$$\textcircled{3} h_i := \nabla_{\hat{y}_i^{(t-1)}}^2 L(y_i, \hat{y}_i^{(t-1)})$$

Then, $R_{\text{emp}}^{(t)} = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t)}) + \Omega(F_t)$

$$= \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(F_t)$$

the empirical loss for

first-t base learner $\approx \sum_{i=1}^n [L(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t(x_i)^2] + \Omega(f_t)$

$$F_t(x) = \sum_{l=1}^t f_l(x)$$

$$\propto \sum_{i=1}^n [g_i f_t(x_i) + \frac{1}{2} h_i f_t(x_i)^2] + \Omega(f_t)$$

$$= \sum_{i=1}^n [g_i f_t(x_i) + \frac{1}{2} h_i f_t(x_i)^2] + \gamma \cdot T + \frac{1}{2} \lambda \sum_{t=1}^T w_t^2$$

$$= \sum_{j=1}^T [(\sum_{i \in I_j} g_i) \cdot w_j + \frac{1}{2} (\sum_{i \in I_j} h_i + \lambda) w_j^2] + \gamma T$$

Assume we know the structure of the t-th base learner

i.e., $\{I_1, I_2, \dots, I_T\}$ \rightarrow the partition of $[n]$

Conclusion: for a fixed structure $\{I_1, I_2, \dots, I_T\}$ $\Leftrightarrow q(x)$

the optimal weight $w_j^* = - \frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$

Note:
 $I_j := \{i : q(x_i) = j\}$

then, $R_{\text{emp}}^{(t)}(q) = - \frac{1}{2} \sum_{j=1}^T \frac{(\sum_{i \in I_j} g_i)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma T$ for a given structure $q(\cdot)$



Rmk: Interpretation $\rightarrow R_{\text{emp}}^{(t)}(q)$ can be viewed as a score to measure

the quality of tree structure $q(\cdot)$ [like Gini impurity & entropy]

Combining with Recursive Binary Search (like DT)

