

Let's start from "Adaboost" & move towards "XGBoost": this assumes that we are familiar with "Decision trees" & "Random Forest".

### 1. Adaboost (Additive boosting)

⇒ Three main ideas:

1. In adaboost each time we make a tree, we make full tree (no predetermined maximum depth).

In contrast trees made with Random Forest are usually just a node & two leaves.



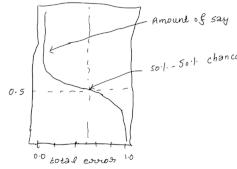
⇒ Stump can only use one variable to make a decision.  
was learned

2. In contrast to Random Forest, in forest of stumps made with Adaboost, some stumps get more "say" in the final classification than others.

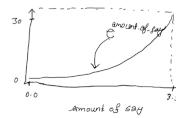
3. Similarly, order in which stumps are made is important.  
while random forest where the tree order is not important.

→ Reasons that the first stump makes influence the chance the second stump makes ... ...

$$\text{# amount of say for a stump} = \frac{1}{2} \log \left( \frac{1 - \text{total chance}}{\text{total chance}} \right)$$



$$\text{# new sample weight} = \text{sample weight} \times \exp^{\frac{t}{2}}$$



\* Strategy is to increase the sample weights of "incorrectly classified" samples and decrease the weight of "correctly classified" samples.

- \* If we have a "weighted count function" then we will use the sample weights, otherwise we use the sample weights to make a new dataset that reflects these weights.

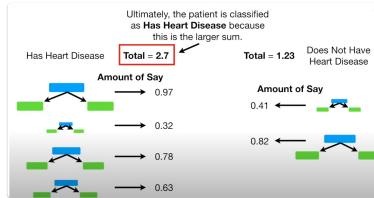


Fig: Example of how Adaboost makes predictions

### 2. Gradient Boosting (Regression):

1. Gradient Boost starts by making a single leaf, instead of a tree or stump, this leaf represents an initial guess for the weight of all of the samples.  
or learning values.

2. Then Gradient Boost builds a tree, but still maintains the size of the tree (to 32 leaves generally).

3. Like Adaboost, gradient boost scales the tree, however it scales all draws by same amount. Then gradient boost builds another tree based on the errors made by previous tree then it repeats the tree ... ... .

$$\text{Prediction} = \boxed{\text{initial value}} + \beta_1 T_1 + \beta_2 T_2 + \beta_3 T_3 + \dots \dots$$

Average of initially  
value is given

where,  $\beta_i$  = learning rate & taking small step in weight  
 $T_1, T_2, \dots, T_n$  = Tree

#### Algorithm

Input: Data  $\{(x_i, y_i)\}_{i=1}^n$  and a differentiable Loss function  $L(y_i; F(x_i))$ .

Step I: Initialize model with a constant value:  $F_0(x) = \arg\min_{f_0} \sum_{i=1}^n L(y_i; f_0)$

Step II: for m = 1 to M:

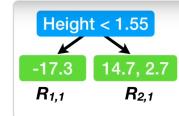
$$\textcircled{4} \text{ compute } g_{im} = \underbrace{\left[ \frac{\partial L(y_i; F(x_i))}{\partial F(x_i)} \right]}_{F(x) = f_{m-1}(x)} \text{ for } i = 1, \dots, n$$

(observed - predicted)  
or  
(Residual)

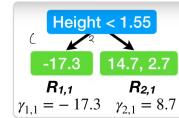
Height (m)	Favorite Color	Gender	Weight (kg)	$R_{1,1}$
1.6	Blue	Male	88	14.7
1.6	Green	Female	76	2.7
1.5	Blue	Female	56	-17.3

fig: Example of residual calculation  
for first tree ( $m=1$ ).

- ④ Fit a regression tree to the stem values & create terminal regions  $R_{jm}$ , for  $j=1, \dots, J_m$

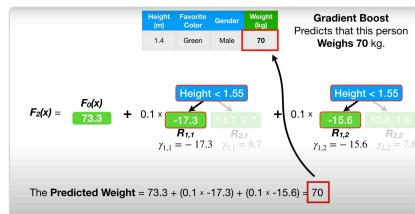


- ⑤ For  $j=1 \dots J_m$  compute  $\hat{y}_{jm} = \underset{\mathcal{D}}{\operatorname{argmin}} \sum_{x_i \in R_{ij}} L(y_i, f_{m-1}(x_i) + \hat{y})$



- ⑥ Update  $f_m(x) = f_{m-1}(x) + \nu \sum_{j=1}^{J_m} \hat{y}_{jm} I(x \in R_{jm})$   
↳ learning rate [0, 1]

Step 3: Output  $F_m(x)$



Figures Making prediction Example  
for M= 1 )

#### # classificaiton using Gradient Boosting:

Input: Data  $\{(x_i, y_i)\}_{i=1}^m$ , and a differentiable Loss function  $L(y_i, F(x))$ .

Step I: Initialize model with a constant value:  $F_0(x) = \underset{\mathcal{D}}{\operatorname{argmin}} \sum_{i=1}^m L(y_i, s)$   
↳ log(odds)

note:  $\log(\text{odds}) = \log\left(\frac{\text{Probability of samples in 'yes' class}}{\text{Probability of samples in 'no' class}}\right)$

if "p" is the probability for "Yes" class then,

$$\log(\text{odds}) = \log(p/(1-p))$$

Step II: For  $m = 1$  to  $M$ :

$$④ \text{compute } g_{im} = \left[ \frac{\partial L(y_i, F(x))}{\partial F(x)} \right]_{F(x)=F_{m-1}(x)} \text{ for } i=1, \dots, n$$

- ⑤ Fit a regression tree to the stem values & create terminal regions  $R_{jm}$ , for  $j=1, \dots, J_m$ .

$$⑥ \text{For } j=1 \dots J_m \text{ compute } \hat{y}_{jm} = \underset{\mathcal{D}}{\operatorname{argmin}} \sum_{x_i \in R_{ij}} L(y_i, f_{m-1}(x_i) + \hat{y})$$

$$⑦ \text{Update } f_m(x) = f_{m-1}(x) + \nu \sum_{j=1}^{J_m} \hat{y}_{jm} I(x \in R_{jm})$$

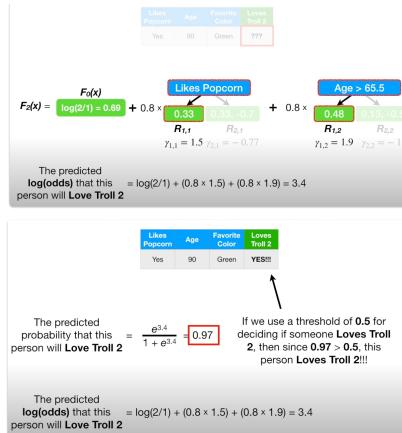


Fig. Example showing classification using gradient boosting

# Important considerations for classification using gradient boosting.

- Log likelihood of the observed data given the prediction

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

where,  $p$  = predicted probability

The goal is to maximize this log-likelihood function.  
or, minimize -log-likelihood function.

On simplification we get:

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

$y$  = observed  
 $\log(\text{odds}) = \log(\frac{p}{1-p})$

also,  $\frac{\partial L}{\partial \log(\text{odds})} = -\text{observed} + p = -y + p$

- In step ② where we need to determine

$$\delta_{j,m} = \underset{j}{\operatorname{argmin}} \sum_{x_i \in R_j} L(y_i, F_{m-1}(x_i) + \delta_j)$$

for single iteration  $R_{11}$ ,

$$\delta_{11} = \underset{j}{\operatorname{argmin}} -y_{11} \times [F_{m-1}(x_{11}) + \delta_j] + \log(1 + e^{F_{m-1}(x_{11}) + \delta_j})$$

Solving for optimal gamma:

Approximate the loss function: Taylor polynomial

$$L(y, F_{m-1}(x) + \delta) = L(y, F_{m-1}(x)) + \frac{\partial}{\partial \delta} (y, F_{m-1}(x)) \delta +$$

$$\frac{1}{2} \frac{\partial^2}{\partial \delta^2} (y, F_{m-1}(x)) \delta^2$$

$$\text{so, } \frac{\partial L(y, F_{m-1}(x) + \delta)}{\partial \delta} = \frac{\partial}{\partial \delta} (y, F_{m-1}(x)) + \frac{\partial^2}{\partial \delta^2} (y, F_{m-1}(x)) \delta$$

now, solving for  $\delta$  with  $\frac{\partial L}{\partial \delta} = 0$  we get,

$$\delta = \frac{-\frac{\partial}{\partial \delta} (y, F_{m-1}(x))}{\frac{\partial^2}{\partial \delta^2} (y, F_{m-1}(x))}$$

$$\text{or, } \delta = \frac{\text{observed} - p}{p \times (1-p)}$$

$$\text{or, } \delta = \frac{\text{Residual}}{p \times (1-p)}$$

For multiple samples in a leaf:

$$\left[ \frac{\text{observed} - p}{p \times (1-p)} \right] = \frac{\sum \text{Residuals}_i}{\sum \text{Residuals}_i}$$