# Gradient-boosting (Analytics Vidhya blog) + statquest

Adaboost VS Gradient-boosting

| ↓ | ↘ |
|---|---|
| we build "stumps" | first we make a single leaf |
| ↓ | (represents initial guess) |
| amount of say of new "stump" or how well it compensates prev-errors | ↓ |
| | The gradient-boost builds a "tree" (new) |
| | ↓ |
| ↓ | Tree is built on errors made by previous trees. Tree is typically larger than a stump. |
| builds next "stump" and so on. | ↓ |
| | Gradient boost also scales the trees but it scales them by a same amount. |

## Gradient-boosting

| GB regressor | GB classifier |
|---|---|
| ↓ | ↓ |
| when target variable is continuous | when target variable is discrete. |

"Gradient-boosting" → As objective is to minimise loss function by adding weak learners using "gradient-descent"

as scale is the same. ←

<u>steps in GB algo :→</u>

step 1) First we build a base model. (single leaf). For simplicity we take the avg. of target column.

why avg.?

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

Base
model
$\gamma$ ↓ we pick such a value that min. loss function.

step 2) Then we calculate the pseudo residuals ⇒ (observed val - pred. val)

why (observed - predicted) ?

$$r_{im} = -\left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x) = F_{m-1}(x)} \quad \text{for } i=1,2,\ldots,n$$

n datapoints

⤷ residual for ith datapoint for mth model

↓ "gradient descent type algo."

the previous model

$$\frac{\partial L}{\partial \gamma} = -(y_i - \gamma) \quad \text{hence} \quad r_{im} = (y_i - \gamma) = (obs. - pred.)$$

we use above residuals to make the ~~next~~ mth model.

↓

we use residuals as our goal is to decrease these residuals.

**step 3)** Let's say $h_m(x)$ is our decision tree made on ~~poor~~ residuals.

In this step we find output value for each leaf of ~~our~~ our dec. tree. If there is a case ~~way~~ where a leaf gets more than one ~~du~~ ~~trul~~ residual → we only take the final output of that leaf

↓

which is the avg. of all residuals in that leaf

↓

why avg?
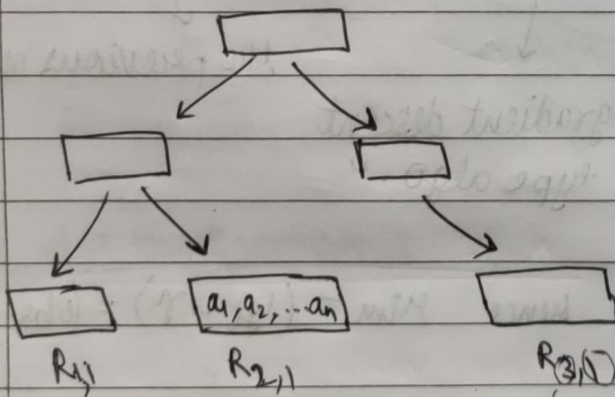
↓

for $j = 1, 2, \ldots, J_m$    $\boxed{\gamma_{jm}} = \underset{\gamma}{\text{argmin}} \sum_{i=1}^{n} L\left(y_i, F_{m-1}(x_i) + \gamma \boxed{h_m(x_i)}\right)$

↓                                                ↓

output of a particular jth leaf.
                         $h_m(x_i)$ is the dec. tree on passing $x_i$ datapoint

suppose
                                 ↙          ↓

                            $1$          $0$

                        for the leaf    for all other it ends up in     leaves.



$R_{1,1}$       $a_1, a_2, \ldots a_n$                   ▢

      $R_{2,1}$                       $R_{3,1}$ → 1st model

            $\boxed{\gamma_{2,1}} = \dfrac{\sum a_i}{n}$     3rd leaf

                          $h_m(x_i) = 0$ for all datapoints who don't reach $R_{2,1}$

output of $R_{2,1}$     proof using above formula → $h_m(x_i) = 1$ for all which reach $R_{2,1}$

                                ↳ derivative says

                                   $\gamma_{2,1} = $ avg. of all residuals in $R_{2,1}$

step 4) We update the predictions of the prev. model.

$$\boxed{F_m(x)} = \boxed{F_{m-1}(x)} + \boxed{V_m}\,\boxed{h_m(x)}$$

↓ new prediction    ↓ prev. prediction    ↓ learning rate    → output of tree made on residuals.

if learning rate = 1 → variance is high

↓

"overfitting of the data."

step 5) Then we recalculate the residuals = (observed − predicted) and repeat from step 3. We do until we reach the no. of trees - limit or exactly fit the data.

<u>more mathematically</u>

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

step 1) Initialise model with constant value :

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

→ no. of trees in our GB algo

step 2) for m=1 to $\boxed{M}$:

A) compute $r_{im} = -\left[\dfrac{\partial L(y_i, F(x_i))}{\partial F(x_i)}\right]$   for $i = 1, 2, \cdots, n$

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

B) Fit a regression tree to the $r_{im}$ values to get dec.
tree $h_m(x)$ and terminal region $R_{jm}$ for $j = 1, \ldots, J_m$

↓

no. of leaves
in mth tree.

C) For $j = 1, \ldots, J_m$ compute

$$r_{jm} = \underset{r}{\text{argmin}} \sum_{x_i \in R_{ij}} L(y_i, F_{m-1}(x_i) + r)$$

↳ only the datapoints that
reach $R_{ij}$

D) update $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} r_{jm} \mathbb{1}(x \in R_{jm})$

the dec. tree
we made

summation is just
in case a datapoint ends
up in multiple leaves

⟶ then we sum up the associated $r$'s (outputs)
of such leaves.

## Gradient boosting classifier (from analytics vidhya)

This is used when target column is binary.
Same as GB regressor except log likelihood is the loss function.

$$L = -\left[ \sum_{i=1}^{n} y_i \log(p) + (1 - y_i) \log(1 - p_i) \right]$$

$$L = -\left[ \sum_{i=1}^{n} \left( y_i \log(odds) - \log(1 + e^{\log(odds)}) \right) \right]$$

we minm. wrt log(odds) ⇒ makes calculations
simpler.

$$\text{odds}_i = \frac{P_i}{1 - P_i}$$

$$\frac{\partial L}{\partial \log(\text{odds}_i)} = -y_i + P_i$$

$$r_{i,m} = -\left[\frac{\partial L}{\partial \log(\text{odds}_i)}\right] = y_i - P_i = \text{observed} - \text{predicted}$$

Now we build a decision tree. If a leaf has more than one
residual we use the following formula

$$\widehat{r} = \frac{\sum\limits_{i=1}^{n} \text{Residual}_i}{\sum\limits_{i=1}^{n} \left[\text{previous probability}_i \times (1 - \text{previous probability}_i)\right]}$$

output of leaf
with multiple
residuals.

### GB classifier (from statquest)

Follow exact same mathematical steps as GB regressor with different
function.