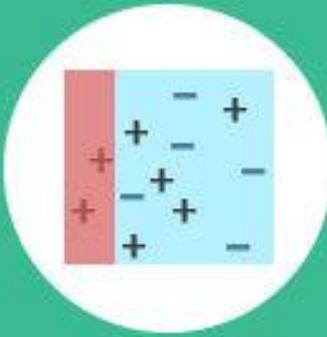


**Gradient  
boosting**



**AdaBoost**

[www.educba.com](http://www.educba.com)

## Table of Contents

1. Understanding Boosting
2. Understanding AdaBoost
3. Solving and Example on AdaBoost
4. Understanding Gradient Boosting
5. Solving an Example on Gradient Boosting
6. AdaBoost Vs Gradient Boosting

## Boosting

→ Boosting is an ensemble technique (involves several trees) that starts from a weaker decision & keeps on building the model such that the final prediction is the weighted sum of all the weaker decisions made.

the weights are assigned based on performance of individual trees

→ in boosting, ensemble parameters are calculated in "stagewise way" which means that while calculating

~~the subsequent weight, the learning from~~

~~previous tree is considered as well. But the~~

~~learning of next tree is done by ignoring the~~

\* i.e., the learning of previous tree boosts the learning of next tree and goes on -

→ we mostly used decision tree as weak classifier.

Any other algorithm can be used as well;

but reasons for choosing tree are -

~~Decision tree gives more interpretable results~~

### Pros

1) Robust to outliers

2) Feature scaling not required

3) Handles missing values

4) can deal with irrelevant inputs

5) computational scalability

~~only Adaboost~~

1) ability to extract a useful combination of features

2) high variance leading to small computation power

→ Boosting minimizes the variance by taking into consideration the result from various trees.

## ⑧ general understanding eg as boost

① you wanna travel to different place.

→ And you know a friend who is a traveller, so you ask him about the place (give your info, he will tell what he know & he will ask his friends).

→ He asks his friend who visited there to give some phone no or any agency.

→ And that friend calls another friend who lives there.

First, we have to boost other person with info we have, and he will boost next person with info he have & finally get result

## ② dad wants to buy car

→ He asks you & without any knowledge you say yes. (Bad weightage)

→ Grandpa thinks and says what model should be bought consider grand children

→ also wants car. (Good weightage)

→ consider son & grandpa, mom says what to buy, first new or second hand, based on financial status (Better weightage)

## (6) AdaBoost (Adaptive Boosting)

$$\begin{matrix} f_1 \\ f_2 \\ \vdots \\ f_N \end{matrix}$$

→ Assume that the no. of training samples are "N", and the no. of iterations (tree created) is "M". Note that possible class outputs are  $Y = \{-1, 1\}$ .

① Initially, total weightage should be equal "1".

so. Initially consider observation weight as  $w_i = \frac{1}{N}$

② for  $m = 1$  to  $M$ ,

→ fit a classifier  $g_m(m)$  to the training data with  $w_i$

→ compute  $\text{err}_m = \sum_{i=1}^N w_i \cdot I(g_i \neq g_m(m))$  Prediction

→ compute  $\alpha_m = \frac{1}{2} \log \frac{(1-\text{err}_m)}{\text{err}_m}$  This is the contribution of that tree to final result

→ calculate new weights using the formula,

$$w_i \leftarrow w_i \cdot \exp[\alpha_m \cdot I(g_i \neq g_m(m))]$$

③ normalize the new sample weights, so that their sum is 1.

④ construct the new tree using the new weights,

i.e; consider weights obtained in 1st tree as the weights for the second tree and, weights of 2nd tree for the 3rd and so on.

⑤ at end, compare summation of results from all trees and final result of either the one with highest sum (for Regr) or the class with most weightage (for classification),

Eg

- ① understand ~~read~~ about Adaboost by  
Simple dataset for heart patient prediction

Step 1

	<u>is chest pain</u>	<u>Arteries blocked</u>	<u>weight or Person</u>	<u>is Heart Patient</u>	<u>new weight</u>
0	Yes	Yes	205	Yes	0.05
1	No	Yes	180	Yes	0.05
2	Yes	No	210	Yes	0.05
3	Yes	Yes	162	Yes	0.33
4	No	Yes	156	No	0.05
5	No	Yes	125	No	0.05
6	Yes	No	160	No	0.05
7	Yes	Yes	172	No	0.05

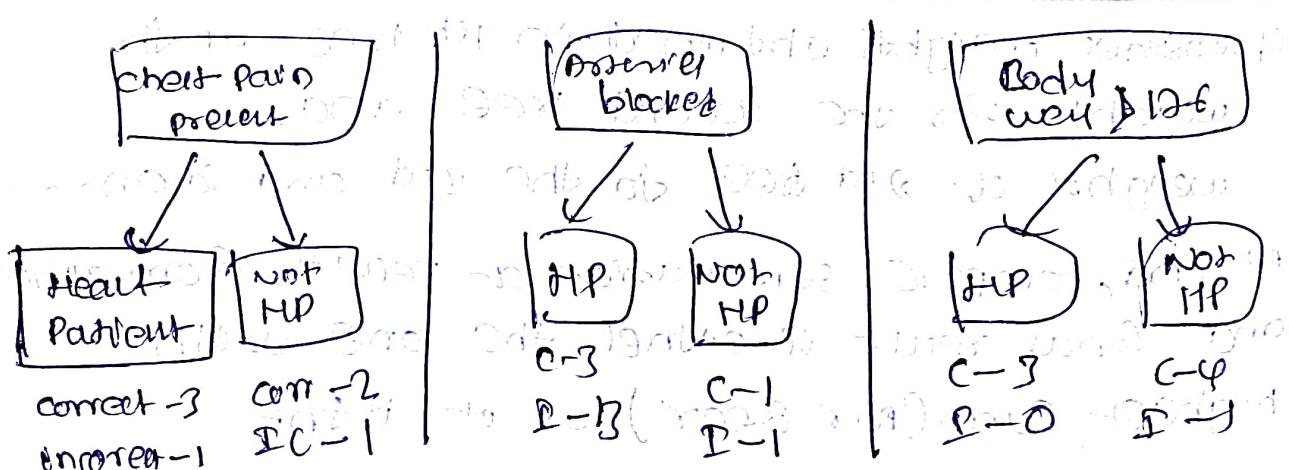
Step ② → Generalized weights

→ there are 8 rows in our dataset.

$$\therefore w = \frac{1}{N} = \frac{1}{8} \quad (\text{samples are equally important})$$

Step ③

consider columns to create a weak decision rule and then try to figure out what all correct & incorrect predictions based on that column



→ we will calculate "alpha index" of each column.

And we select tree with the lowest AI.

& this will be the first decision maker for our model.

now we

Step ③ Now, we calculate the contribution of tree to our final decision

→ we got weight of first decision maker & we have one incorrectly predicted by it.

$$\text{error} = 1/p$$

$$\text{contribution} = \frac{1}{8} \log\left(\frac{1 - \text{err}}{\text{err}}\right) \approx 0.97$$

→ Step ④

modify the weights

that's how AdaBoost works in each iteration

→ increase the sample weight for

incorrectly predicted, new weight = old weight  $e^{0.97}$

$$w_{i+1} = w_i \cdot e^{0.97} = 0.33$$

→ decrease the sample weight for correctly predicted

$$\text{new weight} = \text{old weight} \cdot e^{-0.97}$$

$$w_{i+1} = \frac{1}{8} \cdot e^{-0.97} = 0.08$$

and add weights of all trees by summation to

normalize it to 1.0 for further iterations

then, new normalized weights will act as sample weight for the next iteration,

### Step ①

- And we create new tree which consider dataset prepared with new sample weight.
- Suppose, m tree classify a person as HP & n tree classify a person as Not HP.

then contribution of all m tree & all n-tree  
 } all m-tree  
 } all n-tree  
 } are added.

Higher  
is  
taken

### \* Gradient Boosting Tree

- Gradient Boosted tree use decision tree at each node.  
 It can work with different loss function, evaluate its gradient & approximate with simple tree.

Adaboost is a special case of Gradient boosting where it use exponential loss function.

#### Algorithm

- ① calculate the avg. over the label column & initially this avg shall minimize the total error.
- ② calculate pseudo residual.

$$\text{pseudo residual} = \text{actual label} - \text{predicted result}$$

On this this negative deviation corrects

prediction

mathematically,

$$\textcircled{4} \text{ derivative of pseudo residual} = \frac{\delta L(y_i, p_{mi})}{\delta (f_m(i))}$$

→ here, gradient of error term is getting calculated at the goal is to minimize the error. Hence, name is gradient boosted tree.

\textcircled{3} create a tree to predict the pseudo residuals instead of a tree to predict for actual val. of

$$\text{new result} = (\text{previous result}) + (\text{learning rate} \times \text{residual})$$

$$f_i(n) = f_0(n) + v \cdot \epsilon_n$$

v - learning rate  
ε - residual.

\textcircled{4} Repeat these steps until residual stop decreasing.

EG to understand GBM.

Q.) we need to develop a model to predict an apartment's rent based on sq. footage data.

Sq feet	Rent
750	1160
800	1200
850	1200
900	1400
950	1600
1000	2000

Step① - set Empirical model

lets set avg rent prices in our model  $f_0(n) \geq \text{MP}$

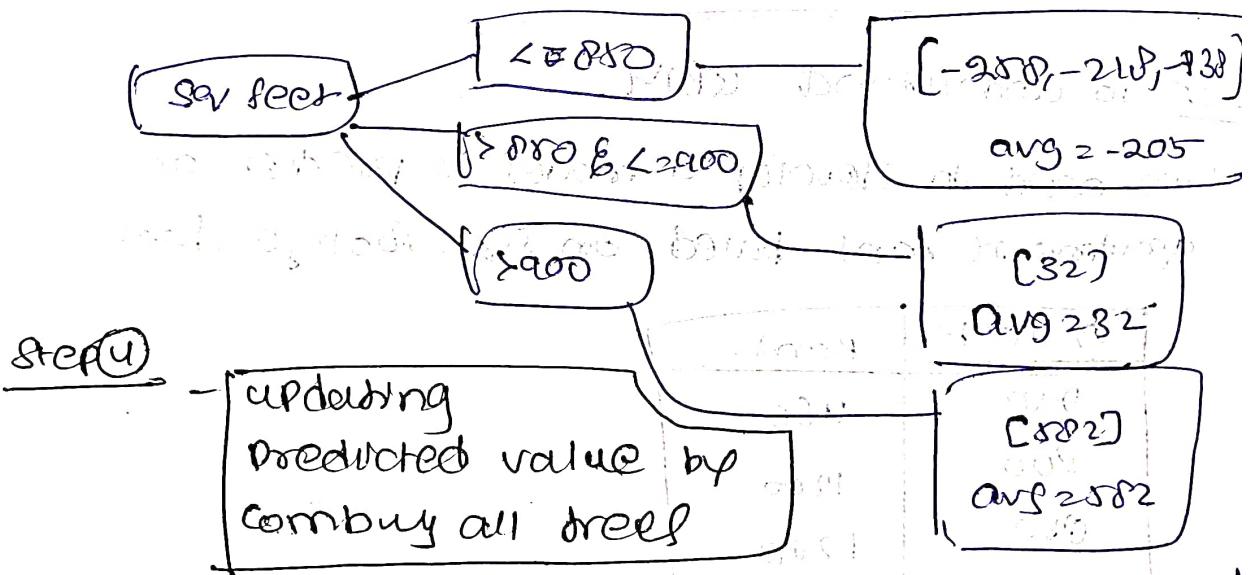
### Step ④ - calculate the residuals

→ For every instance, we compute the residual as diff b/w actual & predicted val (avg)

<u>sq feet</u>	<u>actual rent</u>	<u>Predicted rent (avg)</u>	<u>Residual</u>
280	1160	1018	-258
800	1200	1018	-218
880	1280	1018	-138
900	1300	1018	32
4500	2000	1018	582

### Step ⑤ → Build a model on residuals, instead on actual col. values

we build a decision tree to predict residual



### Step ⑥

- updating  
predicted value by  
combining all trees

from updated predicted val =  $P_0(n) + \text{learning rate} \times \frac{\text{Residual}}{\text{Predicted}}$

<u><math>P_0(n)</math></u>	<u>Residual</u>	<u>Learn rate</u>	<u><math>P_1(n)</math></u>
1018	-205	1	1213
1018	-205	1	1213
1018	-205	1	1213
1018	32	1	1050
1018	582	1	2000

step 6 — repeat step 2-4, for the predefined no. of iterations,

step 6 — finally, the composite model sums together all of the weak models into one strong prediction.

### GBM vs AdaBoost

→ essentially both AdaBoost & GBM have similar intuition, i.e., to convert a set of weak learners into a single strong learner.

④ however, the only point of difference is,

"How they create a weak learner during process?"

→ AdaBoost, changes the weights assigned to each of the instances, such that every subsequent learner focused more on the difficult ones and therefore, contributed towards creating a strong learner.

→ Gradient boosting, trains the weak learner on errors. Each iteration computes the residuals & fits the next learner. Finally, using an optimization process, the algorithm computes the contribution of each weak learner which minimized the overall error of the model.