



MACHINE LEARNING

All You Need to Know about Gradient Boosting Algorithm

2. Classification

LATEST

EDITOR'S PICKS

DEEP DIVES

Algorithm explained with an example, math, and code

Tomonori Masui

Feb 7, 2022 14 min read

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

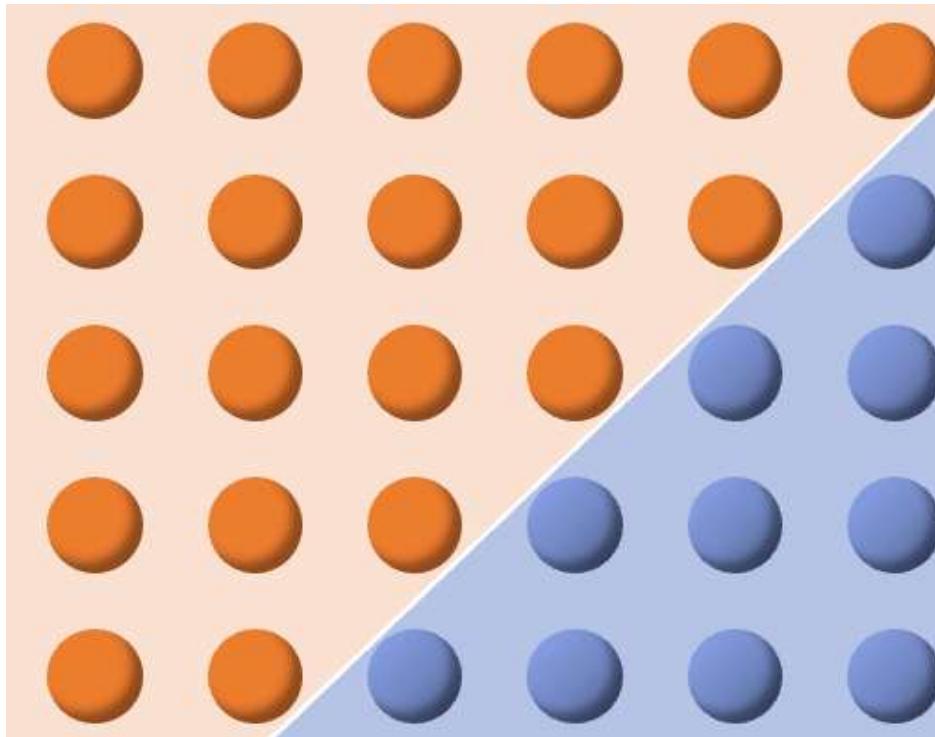


Image by author

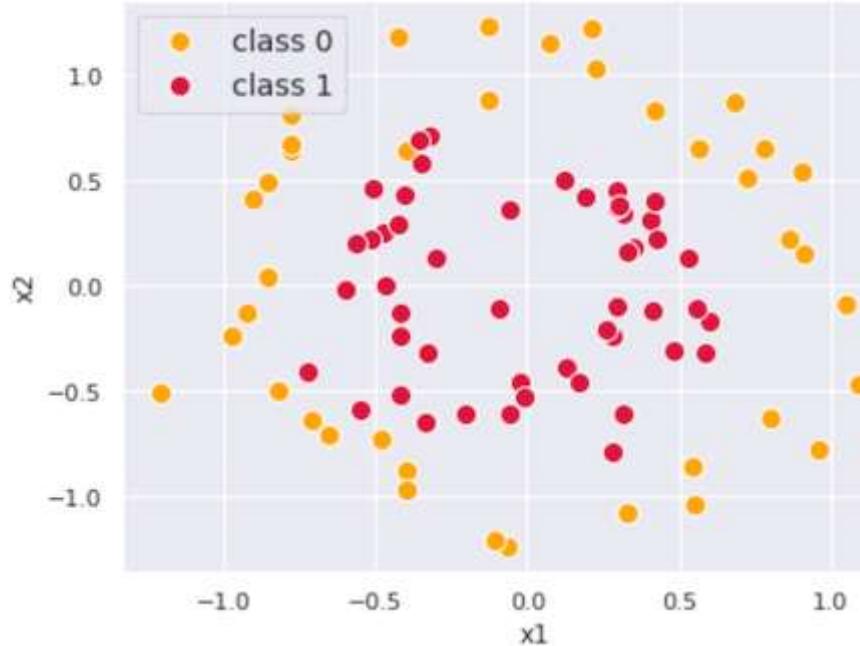
In [the Part 1 article](#), we learned the gradient boosting algorithm in its detail. As we reviewed in that post, it is flexible enough to deal with any loss function.

differentiable. That means if we just replace the loss function used for regression, specifically mean squared loss, with a loss function that deals with classification problems, we can perform classification without changing the algorithm itself. Even though the base algorithm is the same, there are still some differences that we want to know. In this post, we will dive into all the details of the classification algorithm.

Algorithm with an Example

Gradient boosting is one of the variants of ensemble learning where you create multiple weak models (they are decision trees) and combine them to get better performance. In this section, we are building a gradient boosting model using very simple example data to intuitively understand how it works.

The picture below shows the sample data. It has one target variable y (0 and 1) and two features x_1 and x_2 .

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[NEWSLETTER](#)[WRITE FOR TDS](#)[Sign in](#)[Submit an Article](#)

Sample for the classification problem (Image by author)

Our goal is to build a gradient boosting model that classifies those two classes. The first step is making a uniform prediction on a probability of class 1 (we will call it p) for all the data points. The most reasonable value for the uniform prediction might be the proportion of class 1 which is just a mean of y .

$$p = P(y = 1) = \bar{y}$$

LATEST

EDITOR'S PICKS

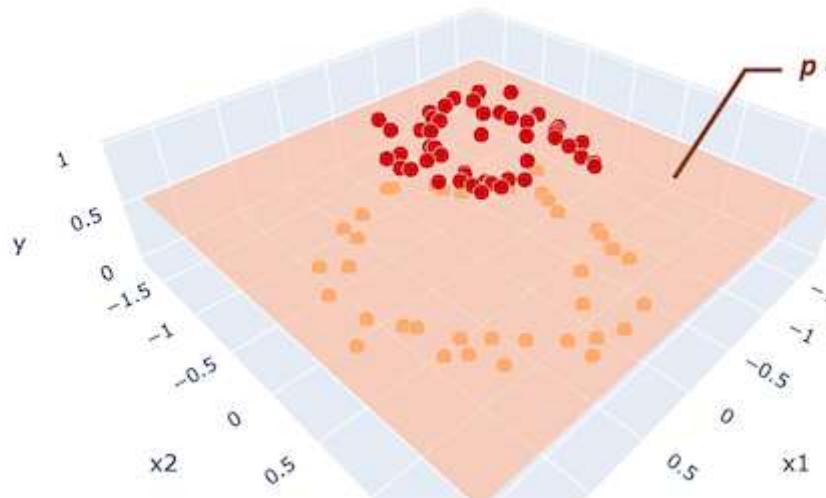
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

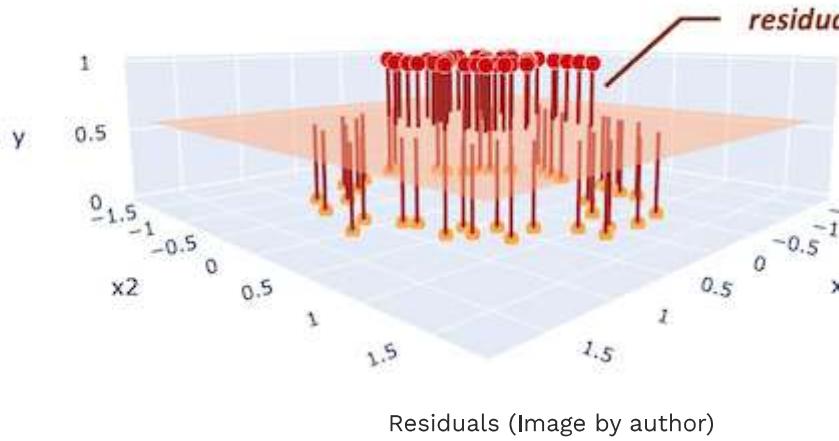
Submit an Article



Prediction plane (Image by author)

In our data, the mean of y is 0.56. As it is bigger than 0.5, everything is classified into class 1 with this initial prediction. Some of you might feel that this uniform value $p = 0.56$ does not make sense, but don't worry. We will improve this prediction as we add more weak models to it.

To improve our prediction quality, we might want to focus on the residuals (i.e. prediction error) from our initial prediction as that is what we want to minimize. The residuals are defined as $r_i = y_i - p$ (i represents the index of each data point). In the figure below, the residuals are shown as the brown lines that are the perpendicular lines from each data point to the prediction plane.



LATEST
EDITOR'S PICKS
DEEP DIVES
NEWSLETTER

WRITE FOR TDS

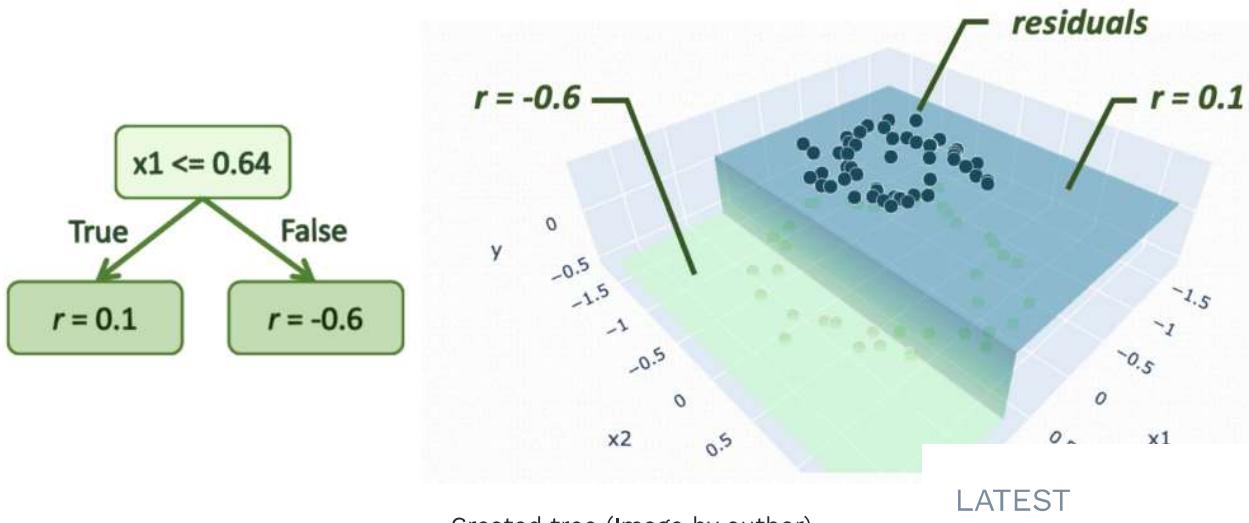
Sign in

Submit an Article

To minimize these residuals, we are building a regression model with both x_1 and x_2 as its features and r as its target. If we can build a tree that finds some patterns between x and r , we can reduce the residuals of the prediction p by utilizing those found patterns.

To simplify the demonstration, we are building very shallow trees, each of which only has one split and two terminal nodes, called "stump". Please note that gradient boosting can have a little deeper trees such as ones with 8 to 10 nodes.

Here we are creating the first tree predicting the two different values $r = \{0.1, -0.6\}$.



LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

$$\gamma_j = \frac{\sum_{x_i \in R_j} (y_i - p)}{\sum_{x_i \in R_j} p(1-p)}$$

y is computed for each terminal node j

Aggregating for all the samples that belongs to terminal node j

$\sum_{x_i \in R_j}$ means we are aggregating the values in the sample x_i s that belong to the terminal node j . The index of each terminal node. You might notice that the numerator of the fraction is the sum of the residuals for all the samples that belong to terminal node j . We will go through all the calculations of this formula in the next section, but let's just focus on the intuition for now. Below is the computed values of γ_1 and γ_2 .

$$\gamma_1 = \frac{\sum_{x_i \in R_1} (y_i - 0.56)}{\sum_{x_i \in R_1} 0.56 \cdot (1 - 0.56)} = 0.3$$

$$\gamma_2 = \frac{\sum_{x_i \in R_2} (y_i - 0.56)}{\sum_{x_i \in R_2} 0.56 \cdot (1 - 0.56)} = -2.2$$

This γ is not simply added to our initial prediction. Instead, we are converting p into log-odds (we will call this converted value $F(x)$), then adding γ to it. For those familiar with log-odds, it is defined below. You may have used in [logistic regression](#).

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[NEWSLETTER](#)[WRITE FOR TDS](#)[Sign in](#)[Submit an Article](#)

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

log-odds

One more tweak on the prediction update is that we scale down by **learning rate** ν , which ranges between 0 and 1. This is added to the log-odds-converted prediction $F(x)$ to prevent the model not to overfit the training data.

$$F_1(x) = F_0(x) + \nu \cdot \gamma$$

In this example, we use a relatively big learning rate. This makes the optimization process easier to understand but is usually supposed to be much smaller values such that the model converges quickly.

By substituting actual values for the variables in the right side of the above equation, we get our updated prediction $F_1(x)$.

$$F_1(x) = \begin{cases} \log\left(\frac{0.56}{1 - 0.56}\right) + 0.9 \cdot 0.3 = 0.5 & \text{if } x_1 \leq 0.64 \\ \log\left(\frac{0.56}{1 - 0.56}\right) - 0.9 \cdot 2.2 = -1.7 & \text{otherwise} \end{cases}$$

If we convert log-odds $F(x)$ back into the predicted probability $p(x)$ (we will cover how we can convert it in the next section), it looks like a stair-like object below.

LATEST

EDITOR'S PICKS

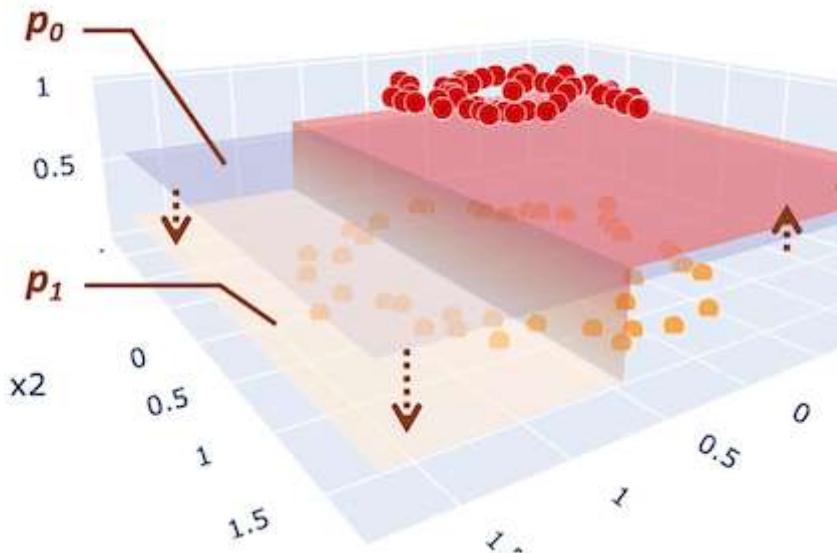
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

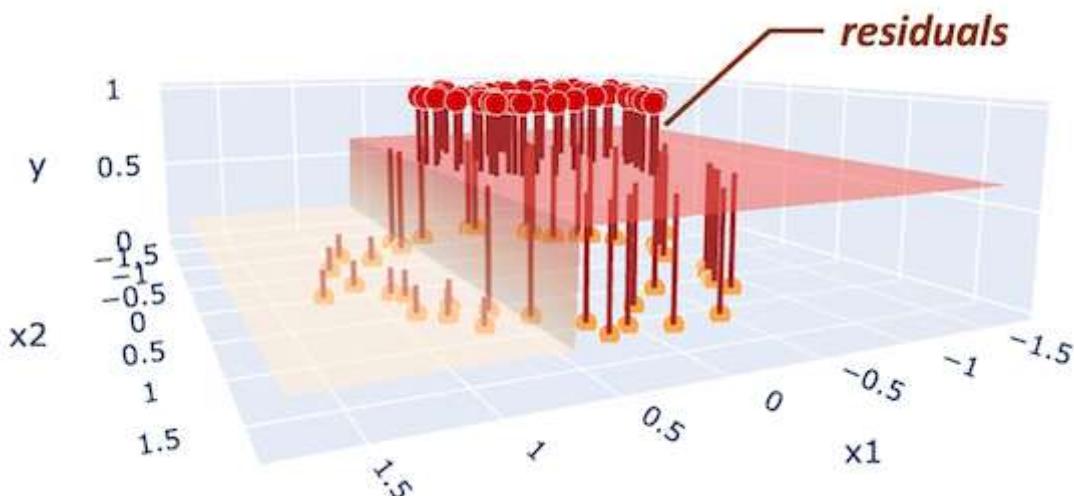
Submit an Article



Updated prediction plane (Image by author)

The purple-colored plane is the initial prediction, and the red and yellow plane p_1 is the updated prediction.

Now, the updated residuals r looks like this:



Updated residuals (Image by author)

LATEST

EDITOR'S PICKS

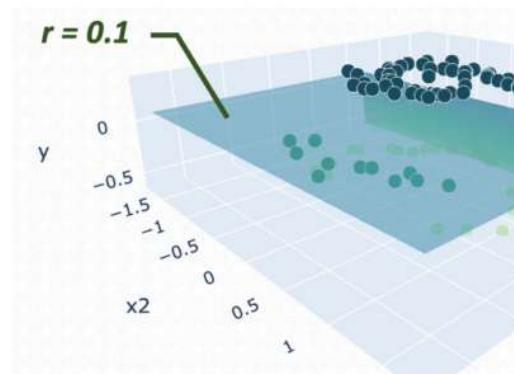
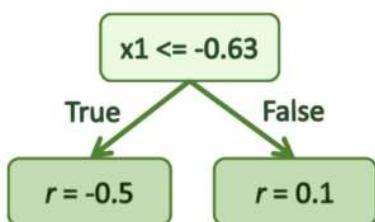
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article



Created tree (Image by author)

We apply the same formula to compute γ . The along with the updated prediction $F_2(x)$ are as follows:

$$F_2(x) = \begin{cases} F_1(x) - \nu \cdot 2.3 = 0.5 - 0.9 \cdot 2.3 = -1.6 & \text{if } x_1 \leq -0.63 \\ F_1(x) + \nu \cdot 0.4 = 0.5 + 0.9 \cdot 0.4 = 0.9 & \text{else if } -0.63 < x_1 \leq 0.64 \\ F_1(x) + \nu \cdot 0.4 = -1.7 + 0.9 \cdot 0.4 = -1.3 & \text{otherwise} \end{cases}$$

These are y computed with this formula:

$$\gamma_j = \frac{\sum_{x_i \in R_j} (y_i - p)}{\sum_{x_i \in R_j} p(1-p)}$$

Again, if we convert log-odds $F_2(x)$ back into the probability $p_2(x)$, it looks like something below.

LATEST

EDITOR'S PICKS

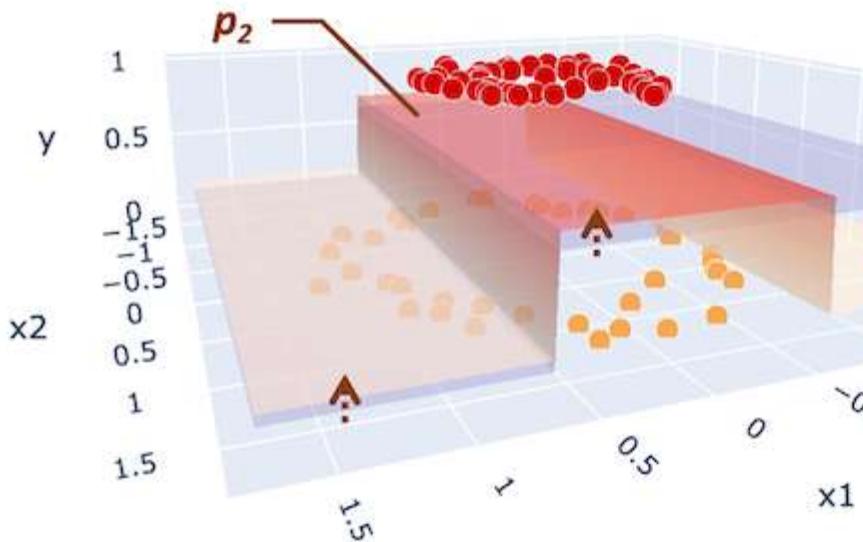
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

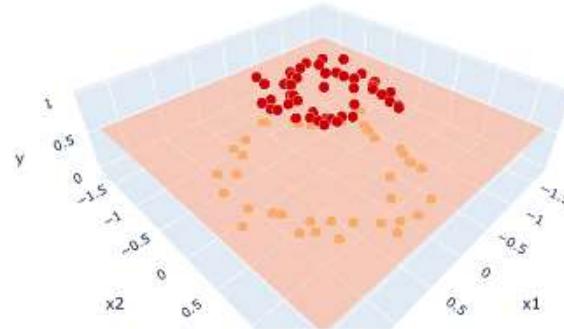
Submit an Article



Updated prediction plane (Image by author)

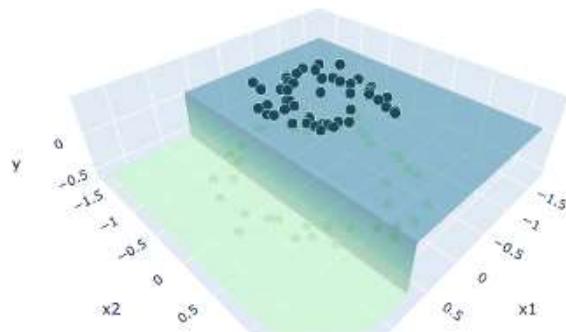
We iterate these steps until the model predictive improving. The figures below show the optimization 0 to 4 iterations.

Predictions of iteration 0



Residuals of iteration 1

Prediction LATEST



EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

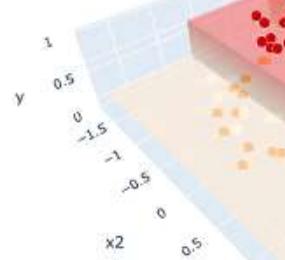
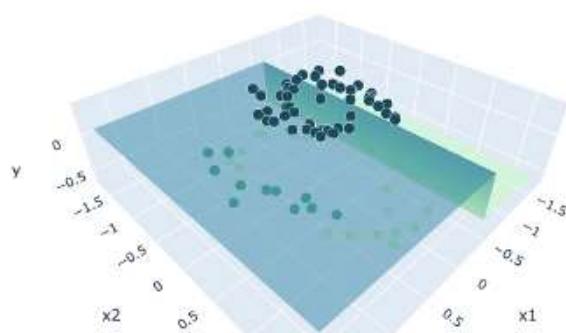
WRITE FOR TDS

Sign in

Submit an Article

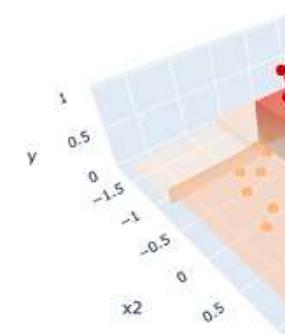
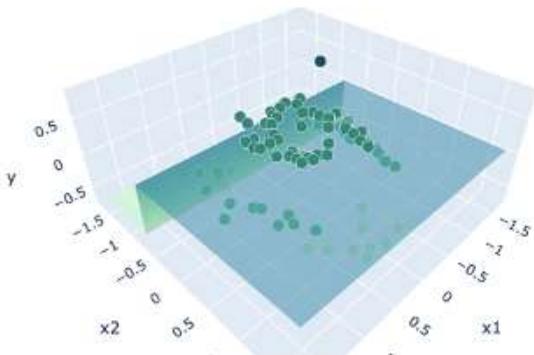
Residuals of iteration 2

Prediction



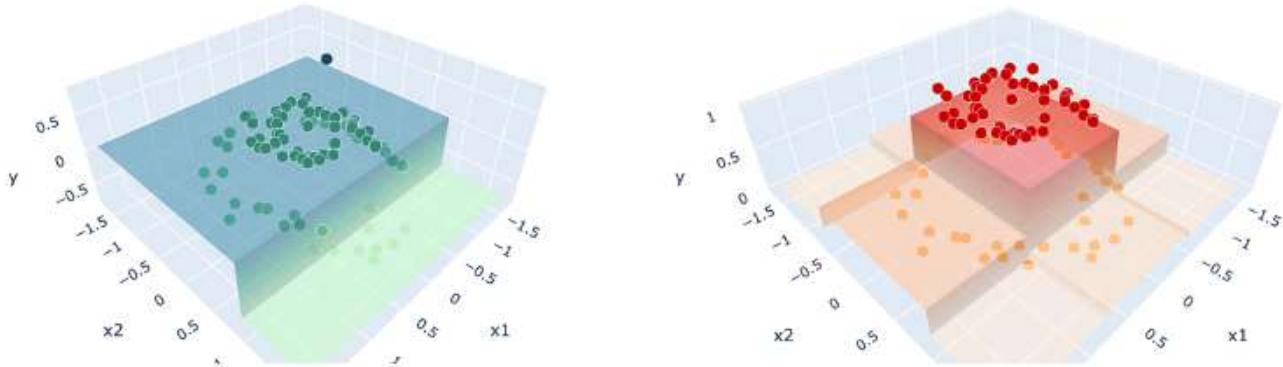
Residuals of iteration 3

Prediction



Residuals of iteration 4

Prediction



Prediction update through iterations (Image by auth~)

You can see the combined prediction $p(x)$ (red) is getting closer to our target y as we add more combined model. This is how gradient boosting complex targets by combining multiple weak models.

LATEST

EDITOR'S PICKS

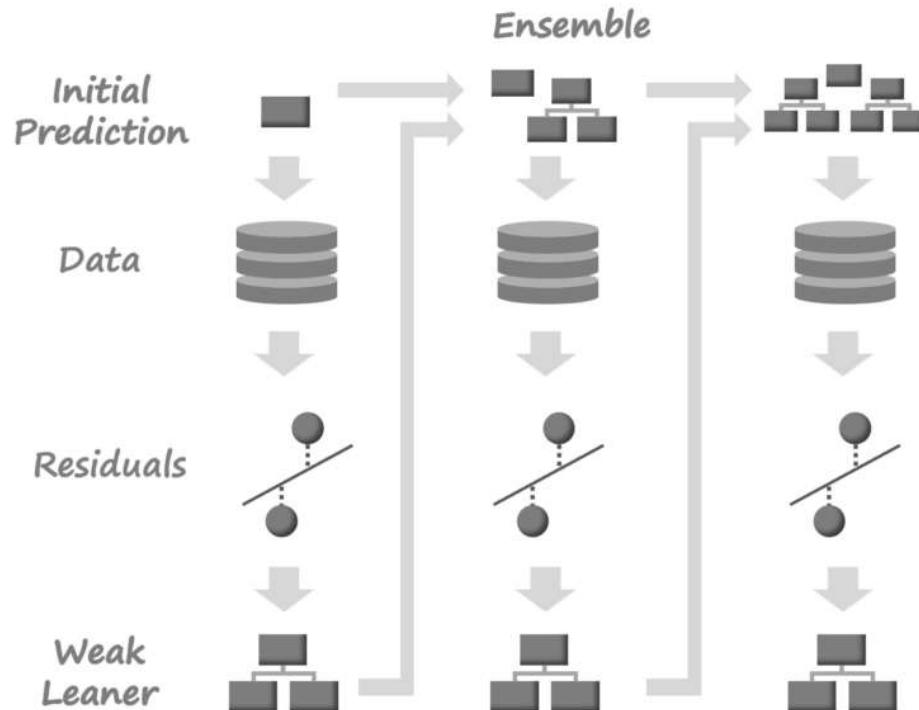
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article



Process of the algorithm (Image by Author)

Math

In this section, we are learning a more generalized algorithm by looking at its math details. Here is the whole algorithm in math formulas.

Gradient Boosting Algorithm

1. Initialize model with a constant value:

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

2. for $m = 1$ to M :

- 2-1. Compute residuals $r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}}$

- 2-2. Train regression tree with features x against r and c

- reasons R_{jm} for $j = 1, \dots, J_m$

- 2-3. Compute $\gamma_{jm} = \underset{\gamma}{\operatorname{argmin}} \sum_{x_i \in R_{jm}} L(y_i, F_{m-1}(x_i) + \gamma)$ for

- 2-4. Update the model:

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

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

Source: adapted from [Wikipedia](#) and [Friedman's paper](#)

Let's take a close look at it line by line.

Step 1

1. Initialize model with a constant value:

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

The first step is creating an initial constant prediction. L is the loss function and we are using log loss.

generally called [cross-entropy loss](#)) for it.

$$L = - (y_i \cdot \log(p) + (1 - y_i) \cdot \log(1 - p))$$

Log Loss

y_i is our classification target and it is either 0 or 1. p is the predicted probability of class 1. You might see L taking different values depending on the target class y_i .

[LATEST](#)

[EDITOR'S PICKS](#)

[DEEP DIVES](#)

[NEWSLETTER](#)

[WRITE FOR TDS](#)

[Sign in](#)

[Submit an Article](#)

*argmin means we are searching for the value γ (minimizes $\sum L(y_{i+}, \gamma)$. While it is more straightforward to think of γ as the predicted probability p , we assume γ is the true probability. *** it makes all the following computations easier. who forgot the log-odds definition which we reviewed in the previous section, it is defined as $\log(\text{odds}) = \log(p/(1-p))$.*

To be able to solve the argmin problem in terms of p we are transforming the loss function into the function

$$\begin{aligned}
 L &= -(y_i \cdot \log(p) + (1 - y_i) \cdot \log(1 - p)) \\
 &= -(y_i \cdot (\log(p) - \log(1 - p)) + \log(1 - p)) \\
 &= -\left(y_i \cdot \log\left(\frac{p}{1-p}\right) + \log(1 - p)\right) \\
 &= -(y_i \cdot \log(odds) + \log(1 - p))
 \end{aligned}$$

LATEST

Now, we might want to replace p in the above equation with something that is expressed in terms of log-odds. Transforming the log-odds expression shown earlier into p represented by log-odds:

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

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

$$\frac{p}{1-p} = e^{\log(odds)}$$

$$p = (1 - p)e^{\log(odds)}$$

$$(1 + e^{\log(odds)})p = e^{\log(odds)}$$

$$p = \frac{e^{\log(odds)}}{1 + e^{\log(odds)}}$$

Then, we are substituting this value for p in the equation and simplifying it.

$$L = - \left(y_i \cdot \log(\text{odds}) + \log \left(1 - \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}} \right) \right)$$

p is replaced with this

$$= - \left(y_i \cdot \log(\text{odds}) + \log \left(\frac{1}{1 + e^{\log(\text{odds})}} \right) \right)$$

$$= - \left(y_i \cdot \log(\text{odds}) + \underbrace{\log(1) - \log(1 + e^l)}_{= 0} \right)$$

$$= - \left(y_i \cdot \log(\text{odds}) - \log(1 + e^{\log(\text{odds})}) \right)$$

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

$$\begin{aligned} \frac{\partial}{\partial \log(\text{odds})} \sum_{i=1}^n L &= - \frac{\partial}{\partial \log(\text{odds})} \sum_{i=1}^n [y_i \cdot \log(\text{odds}) - \log(1 + e^{\log(\text{odds})})] \\ &= - \sum_{i=1}^n y_i + n \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}} \\ &= - \sum_{i=1}^n y_i + np \end{aligned}$$

In the equations above, we replaced the fraction odds with p to simplify the equation. Next, we set $\partial \log(\text{odds})$ equal to 0 and solving it for p .

$$-\sum_{i=1}^n y_i + np = 0$$

$$np = \sum_{i=1}^n y_i$$

$$p = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$$

In this binary classification problem, y is either mean of y is actually the proportion of class 1. see why we used $p = \text{mean}(y)$ for our initial predi

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

As γ is log-odds instead of probability p , we ar into log-odds.

WRITE FOR TDS

$$F_0(x) = \gamma^* = \log \left(\frac{\bar{y}}{1 - \bar{y}} \right)$$

Sign in

Submit an Article

Step2

2. for $m = 1$ to M :

The whole step2 processes from 2–1 to 2–4 are m denotes the number of trees we are creating represents the index of each tree.

Step 2–1

2-1. Compute residuals $r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)}$

We are calculating residuals r_{im} by taking a derivative of the loss function with respect to the previous prediction F_{m-1} and multiplying it by -1 . As you can see in the subscript index, r_{im} is computed for each single sample i . Some of you might be wondering why we are calling this r_{im} residuals. This value is actually **negative gradient** that gives us the directions (+/-) and the magnitude in which the loss function can be minimized. You will see why we are calling it residuals shortly. The technique where you use a gradient to minimize model is very similar to gradient descent technique typically used to optimize neural networks. (In fact, they are slightly different from each other. If you are interested, look at this article detailing that topic.)

Let's compute the residuals here. F_{m-1} in the equation is the prediction from the previous step. In this first iteration, F_0 . As in the previous step, we are taking a derivative with respect to log-odds instead of p since our prediction is expressed by log-odds. Below we are using L “expressed by log-odds” in the previous step.

$$\begin{aligned} r_{im} &= -\frac{\partial}{\partial \log(\text{odds})} L \\ &= \frac{\partial}{\partial \log(\text{odds})} [y_i \cdot \log(\text{odds}) - \log(1 + e^{\log(\text{odds})})] \\ &= y_i - \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}} \end{aligned}$$

In the previous step, we also got this equation:

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

$$p = \frac{e^{\log(odds)}}{1 + e^{\log(odds)}}$$

So, we can replace the second term in r_{im} equation with p .

$$r_{im} = y_i - p$$

You might now see why we call r residuals. This interesting insight that the negative gradient has direction and the magnitude to which the loss is actually just residuals.

LATEST
EDITOR'S PICKS
DEEP DIVES
NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

Step 2–2

2-2. Train regression tree with features x against r and create regions R_{jm} for $j = 1, \dots, J_m$

j represents a terminal node (i.e. leave) in the tree index, and capital J means the total number of terminal nodes.

Step 2–3

2-3. Compute $\gamma_{jm} = \operatorname{argmin}_\gamma \sum_{x_i \in R_{jm}} L(y_i, F_{m-1}(x_i) + \gamma)$ for each terminal node j . $\sum_{x_i \in R_{jm}} L$ means we are summing the loss on all the x_i s that belong to the terminal node j .

We are searching for γ_{jm} that minimizes the loss for each terminal node j . $\sum_{x_i \in R_{jm}} L$ means we are summing the loss on all the x_i s that belong to the terminal node j . We will plugin the loss function into the equation.

$$\begin{aligned}\gamma_{jm} &= \underset{\gamma}{\operatorname{argmin}} \sum_{x_i \in R_{jm}} L(y_i, F_{m-1}(x_i) + \gamma) \\ &= \underset{\gamma}{\operatorname{argmin}} \sum_{x_i \in R_{jm}} - \left(y(F_{m-1}(x_i) + \gamma) - \log(1 + e^{F_{m-1}(x_i) + \gamma}) \right)\end{aligned}$$

Solving this equation for γ_{jm} is going to be very hard. To make it more easily solvable, we are making the approximation of L “ by using the second-order Taylor polynomial. A Taylor polynomial is a way to approximate any function as a polynomial of an infinite/finite number of terms. While we are not going to go into the details here, you can look at this tutorial which explains it nicely if you are interested.

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

$$L(y, F_{m-1}(x_i) + \gamma) \approx L(y, F_{m-1}(x_i)) + \frac{\partial}{\partial F} L(y_i, F_{m-1}(x_i)) \gamma + \dots$$

Submit an Article

We are substituting this approximation for L in the equation for γ_{jm} , then finding the value of γ_{jm} that makes the sum $\Sigma(*)$ equal zero.

$$\frac{\partial}{\partial \gamma} \sum_{x_i \in R_{jm}} \left(L(y_i, F_{m-1}(x_i)) + \frac{\partial}{\partial F} L(y_i, F_{m-1}(x_i)) \gamma + \frac{1}{2} \frac{\partial^2 L(y_i, F_{m-1}(x_i))}{\partial F^2} \gamma^2 \right) = 0$$

Derivative of this term is zero

$$\sum_{x_i \in R_{jm}} \left(\frac{\partial}{\partial F} L(y_i, F_{m-1}(x_i)) + \frac{\partial^2 L(y_i, F_{m-1}(x_i))}{\partial F^2} \gamma \right) = 0$$

$$\sum_{x_i \in R_{jm}} \frac{\partial^2 L(y_i, F_{m-1}(x_i))}{\partial F^2} \gamma = - \sum_{x_i \in R_{jm}} \frac{\partial}{\partial F} L(y_i, F_{m-1}(x_i))$$

LATEST

$$\gamma = \frac{-\sum_{x_i \in R_{jm}} \frac{\partial}{\partial F} L(y_i, F_{m-1}(x_i))}{\sum_{x_i \in R_{jm}} \frac{\partial^2}{\partial F^2} L(y_i, F_{m-1}(x_i))}$$

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

As we already computed $\partial L / \partial F$ in the previous s

WRITE FOR TDS

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

Sign in

Submit an Article

We are substituting this for $\partial L / \partial F$ in the γ equa

$$\begin{aligned}
 \gamma &= \frac{\sum_{x_i \in R_{jm}} (y_i - p)}{\sum_{x_i \in R_{jm}} \frac{\partial}{\partial F} (-y_i + p)} \\
 &= \frac{\sum_{x_i \in R_{jm}} (y_i - p)}{\sum_{x_i \in R_{jm}} \frac{\partial}{\partial F} \left(-y_i + \frac{e^{\log(\text{odds})}}{1+e^{\log(\text{odds})}} \right)} \\
 &= \frac{\sum_{x_i \in R_{jm}} (y_i - p)}{\sum_{x_i \in R_{jm}} \frac{\partial}{\partial F} \left(-y_i + e^{\log(\text{odds})} (1 + e^{\log(\text{odds})})^{-1} \right)} \\
 &\quad \text{Derivative of } y \text{ is zero} \\
 &= \frac{\sum_{x_i \in R_{jm}} (y_i - p)}{\sum_{x_i \in R_{jm}} \frac{\partial}{\partial F} \left(e^{\log(\text{odds})} (1 + e^{\log(\text{odds})})^{-1} \right)} \\
 &= \frac{y_i - p}{\sum_{x_i \in R_{jm}} \left[e^{\log(\text{odds})} (1 + e^{\log(\text{odds})})^{-1} - e^{2\log(\text{odds})} (1 + e^{\log(\text{odds})})^{-2} \right]} \\
 &\quad \text{Apply product rule } d(u \cdot v)/dx \\
 &= \frac{\sum_{x_i \in R_{jm}} (y_i - p)}{\sum_{x_i \in R_{jm}} \left(\frac{e^{\log(\text{odds})}}{1+e^{\log(\text{odds})}} - \left(\frac{e^{\log(\text{odds})}}{1+e^{\log(\text{odds})}} \right)^2 \right)} \\
 &\quad = p \quad = p \\
 &= \frac{\sum_{x_i \in R_{jm}} (y_i - p)}{\sum_{x_i \in R_{jm}} (p - p^2)} \\
 &= \frac{\sum_{x_i \in R_{jm}} (y_i - p)}{\sum_{x_i \in R_{jm}} p(1 - p)}
 \end{aligned}$$

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

Finally, we got this simplified equation for the weight γ which we used in the previous section.

Step 2–4

2-4. Update the model:

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

In the final step, we are updating the prediction of the combined model F_m . $\gamma_{jm} 1(x \in R_{jm})$ means that we pick the value γ_{jm} if a given x falls in a terminal node R_{jm} . As all the terminal nodes are exclusive, any given single x falls into only a single terminal node and corresponding γ_{jm} is added to the previous LATEST and then it makes the updated prediction F_m .

As mentioned in the previous section, ν is learned between 0 and 1 which controls the degree of contribution of the additional tree prediction γ to the combined prediction. A smaller learning rate reduces the effect of the additional prediction, but it basically also reduces the chance of overfitting to the training data.

LATEST
EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

Now we have gone through the whole steps. To improve model performance, we want to iterate step 2 and 3. This means adding m trees to the combined model. I might often want to add more than 100 trees to improve model performance.

If you read [my article for a regression algorithm](#), the math computation of the classification algorithm is more complicated than the regression. While [an article](#) shows the derivative computation of the log-loss function, the underlining math algorithm which is shown in the first part of this section is exactly the same. That is actually

of the gradient boosting algorithm as exactly the same algorithm works for any loss function as long as it is differentiable. In fact, popular gradient boosting implementations such as [[XGBoost](#)] (<https://xgboost.readthedocs.io/en/stable/parameter.html#learning-task-parameters>) or [[LightGBM](#)] (<https://lightgbm.readthedocs.io/en/latest/Parameters.html#objective>) have a wide variety of loss functions, so you can choose whatever loss functions that suit your problem loss functions available in XGBoost or LightGBM

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

Code

In this section, we are translating the maths we into a viable python code to help us understand further. We are using `DecisionTreeRegressor` from `s` build trees which helps us just focus on the gra algorithm itself instead of the tree algorithm. W scikit-learn style implementation where you tra fit method and make predictions with `predict`

```

1  class CustomGradientBoostingClassifier:
2
3      def __init__(self, learning_rate, n_estimators, max_depth=1):
4          self.learning_rate = learning_rate
5          self.n_estimators = n_estimators
6          self.max_depth = max_depth
7          self.trees = []
8
9      def fit(self, X, y):
10
11         F0 = np.log(y.mean())/(1-y.mean()) # log-odds values
12         self.F0 = np.full(len(y), F0) # converting to array with
13         Fm = self.F0.copy()
14
15         for _ in range(self.n_estimators):

```

```

16     p = np.exp(Fm) / (1 + np.exp(Fm)) # converting back to probabilities
17     r = y - p # residuals
18     tree = DecisionTreeRegressor(max_depth=self.max_depth, random_state=0)
19     tree.fit(X, r)
20     ids = tree.apply(x) # getting the terminal node IDs
21
22     # looping through the terminal nodes
23     for j in np.unique(ids):
24         fltr = ids == j
25
26         # getting gamma using the formula ( $\sum \text{residuals} / \sum p(1-p)$ )
27         num = r[fltr].sum()                                LATEST
28         den = (p[fltr]*(1-p[fltr])).sum()
29         gamma = num / den                               EDITOR'S PICKS
30
31         # updating the prediction
32         Fm[fltr] += self.learning_rate * gamma          DEEP DIVES
33
34         # replacing the prediction value in the tree
35         tree.value[j, 0, 0] = gamma                     NEWSLETTER
36
37         self.trees.append(tree)
38
39     def predict_proba(self, X):
40
41         Fm = self.F0
42
43         for i in range(self.n_estimators):
44             Fm += self.learning_rate * self.trees[i].predict(X)
45
46         return np.exp(Fm) / (1 + np.exp(Fm)) # converting back to probabilities

```

[Sign in](#)[Submit an Article](#)

custom_gbdt.py hosted with ❤ by GitHub

Please note that all the trained trees are stored object and it is retrieved when we make predict predict_proba method.

Next, we are checking if our CustomGradientBoosting performs as the same as GradientBoostingClassifier learn by looking at their log-loss on our data.

```

1  from sklearn.ensemble import GradientBoostingClassifier
2  from sklearn.metrics import log_loss
3
4  custom_gbm = CustomGradientBoostingClassifier(
5      n_estimators=20,
6      learning_rate=0.1,
7      max_depth=1
8  )
9  custom_gbm.fit(x, y)
10 custom_gbm_log_loss = log_loss(y, custom_gbm.predict_proba(x))
11 print(f"Custom GBM Log-Loss:{custom_gbm_log_loss:.15f}")
12
13 sklearn_gbm = GradientBoostingClassifier(
14     n_estimators=20,
15     learning_rate=0.1,
16     max_depth=1
17 )
18 sklearn_gbm.fit(x, y)
19 sklearn_gbm_log_loss = log_loss(y, sklearn_gbm.predict_proba(x))
20 print(f"Scikit-learn GBM Log-Loss:{sklearn_gbm_log_loss:.15f}")

```

compare_gbm_log_loss.py hosted with ❤ by GitHub

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

[Sign in](#)

[Submit an Article](#)

Output

Custom GBM Log-Loss: 0.46194306798
 Scikit-learn GBM Log-Loss: 0.46194306798

As you can see in the output above, both model the same log-loss.

Recommended Resources

In this blog post, we have reviewed all the details of the gradient boosting classification algorithm. If you are also interested in the gradient boosting regression algorithm, please look at the Part 1 article.

All You Need to Know about Gradient Boosting

Part 1. Regression

There are also some other great resources if you want further details of the algorithm:

- **StatQuest, Gradient Boost Part3 and Part 4** These are the YouTube videos explaining the gradient boosting classification algorithm with great visuals in a beginner-friendly way.
- **Terence Parr and Jeremy Howard, How to explain gradient boosting** While this article focuses on gradient regression instead of classification, it nicely details the detail of the algorithm.
- **Jerome Friedman, Greedy Function Approximation Gradient Boosting Machine** This is the original work by Jerome Friedman. While it is a little hard to understand, it shows the flexibility of the algorithm where it is a generalized algorithm that can deal with any function having a differentiable loss function.

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

You can also look at the full Python code in the [Google Colaboratory](#) or the Github link below.

Google Colaboratory

Gradient Boosting Algorithm – Classification masui/gradient-boosting

References

- Jerome Friedman, Greedy Function Approximation Gradient Boosting Machine

- Terence Parr and Jeremy Howard, [How to explain gradient boosting](#)
- Matt Bowers, [How to Build a Gradient Boosting Machine from Scratch](#)
- Wikipedia, [Gradient boosting](#)

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

WRITTEN BY

Tomonori Masui

[See all from Tomonori Masui](#)

Deep Dives

Gradient Boosting

Machine Learnir

Python

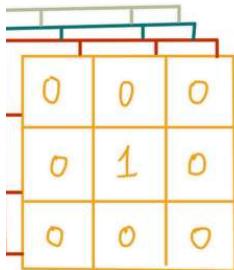
Share This Article



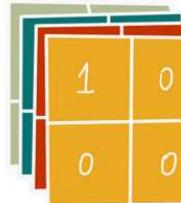
Towards Data Science is a community publishing platform where you can share your insights to reach our global audience and earn rewards through the TDS Author Payment Program.

[Write for TDS](#)

Related Articles



max
pooling



ARTIFICIAL INTELLIGENCE

Implementing Convolutional Neural Networks in TensorFlow

Step-by-step code guide to building a Convolutional Neural Network

Shreya Rao

August 20, 2024 6 min read



DATA SCIENCE

Hands-on Time Series Anomaly Detection using Autoencoders, with Python

Here's how to use Autoencoders to detect signals with anomalies in a few lines of...

Piero Paialunga

August 21, 2024 12 min read



LATEST
EDITOR'S PICKS
DEEP DIVES

How to Forecast Time Series

A beginner's guide to time series reconciliation

Dr. Robert Kübler

August 20, 2024

NEWSLETTER

WRITE FOR TDS

[Sign in](#)

[Submit an Article](#)



MACHINE LEARNING

3 AI Use Cases (Chatbot)

Feature engineer unstructured data

Shaw Talebi

August 21, 2024 7 min read



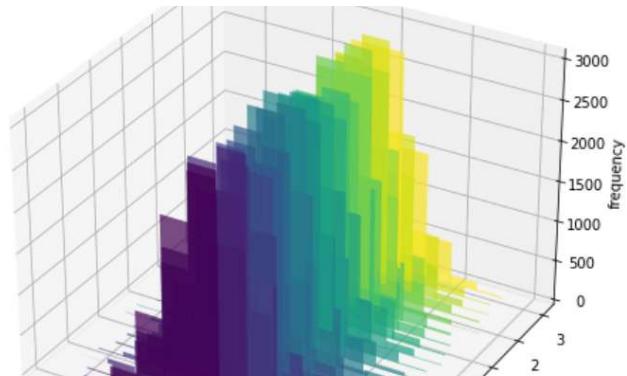
DATA SCIENCE

Back To Basics, Part Uno: Linear Regression and Cost Function

An illustrated guide on essential machine learning concepts

Shreya Rao

February 3, 2023 6 min read



DATA SCIENCE

Must-Know in Bivariate Normal Explained

Derivation and properties of the powerful concept

Luigi Battistoni

August 14, 2024 7

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER



DATA SCIENCE

Our Columns

Columns on TDS are carefully curated collections of posts on a particular idea or category...

TDS Editors

November 14, 2020 4 min read

WRITE FOR TDS

Sign in

Submit an Article



Your home for data science and AI. The world's leading publication for data science, data analytics, data engineering, machine learning, and artificial intelligence professionals.

© Insight Media Group, LLC 2025

[Subscribe to Our Newsletter](#)

[WRITE FOR TDS](#) · [ABOUT](#) · [ADVERTISE](#) · [PRIVACY POLICY](#)

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

[Sign in](#)

[Submit an Article](#)