



The AI-ML and Data Science Club of IITH

Bagging and Boosting techniques

Contents

1. What is Bootstrapping?
2. What is Bagging
3. How Bagging Works?
 - 3.1 Bias and variance
 - 3.2 How Bagging reduces variance?
 - 3.3 Effect of Correlation
 - 3.4 Random Forests
4. What is boosting?
 - 4.1 Weighted dataset
 - 4.2 Different type of boosting algorithm
 - 4.3 AdaBoost
 - 4.4 GBM
 - 4.5 XGBM
 - 4.6 LightGBM
 - 4.7 CatBoost
 - 4.8 Comparison of different boosting algorithms
5. Limitations and Assumptions
6. Questions
7. Implementing AdaBoosting

Compiled by : Anshul Sangrame

1. WHAT IS BOOTSTRAPPING?

Let's first start by understanding what is Bootstrapping. This statistical technique consists in generating samples of size B (called bootstrap samples) from an initial dataset of size N by randomly drawing with replacement B observations.

2. WHAT IS BAGGING?

The algorithm is explained in the following steps:

- 1) Given a dataset D with N training points and a training model
- 2) Create M bootstrap samples of D i.e. $\{\tilde{D}_i\}_{i=1}^M$ with same number of training points i.e. N .
- 3) Create M copies of untrained model $\{h_i\}_{i=1}^M$ and train each h_i on \tilde{D}_i
- 4) Let y_i be predicted value of model h_i . Then the final predicted value will be $y = \frac{1}{M} \sum_{i=1}^M y_i$

3. HOW BAGGING WORKS?

Now that we know the algorithms, we will try to understand how does it work. But for that, we need some knowledge of bias and variance.

3.1 Bias and variance

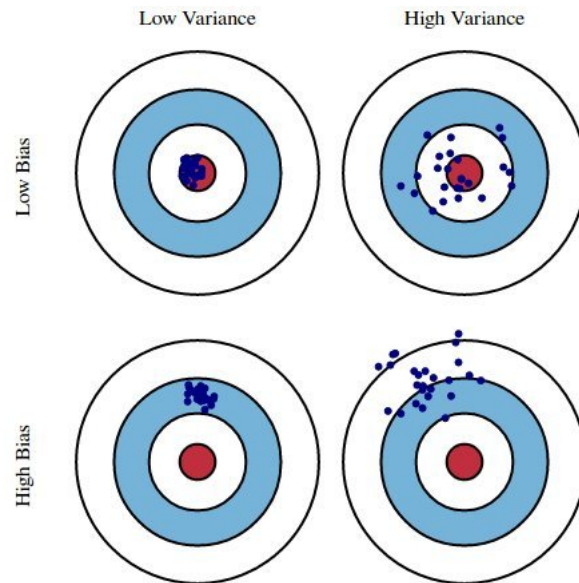
Let \hat{y} be the predicted target value of a random input X and y be the actual target value. y_* be the best-predicted target value of input X that can be made. We can show that,

$$E[(\hat{y} - y)^2] = (y_* - E[\hat{y}])^2 + Var[\hat{y}] + Var[y]$$

We can interpret the following terms in the above equation:

- 1) $E[(\hat{y} - y)^2]$ is the expected loss of the predicted value. We need to minimize this.
- 2) $(y_* - E[\hat{y}])^2$ is called the bias. It indicates how close is the predicted value to the best prediction that can be achieved. Higher bias corresponds to underfitting.
- 3) $Var[\hat{y}]$ is called the variance. It indicates the amount of variability in the predictions (a higher value corresponds to overfitting).
- 4) $Var[y]$ is called the Bayes error and is the inherent unpredictability of the targets. We cannot reduce this.

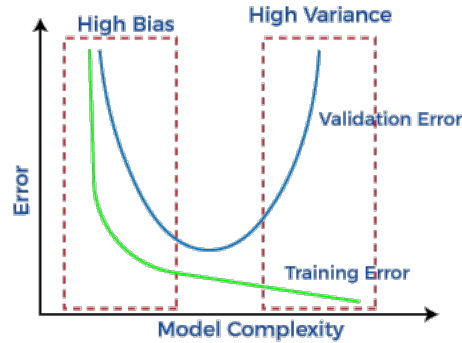
For a better understanding refer to the image below:



Other useful observations related to bias and variance are:

- 1) High bias indicates that the model is less complex and is not trained sufficiently (i.e. training error is high) which indicates underfitting.
- 2) High variance indicates that the model is more complex and is over-trained (i.e. training error is low) which indicates overfitting.

The above two points are portrayed in the following image:



3.2 How Bagging reduces variance?

In the previous section, we saw different terms in the expected loss function. Now we will try to understand how each term is affected in bagging.

- 1) Bayes error: Unchanged, since we have no control over it.
- 2) Bias: Unchanged

$$\begin{aligned}
 E[\hat{y}] &= E\left[\frac{1}{M} \sum_{i=1}^M \hat{y}_i\right] \\
 &= \frac{1}{M} \sum_{i=1}^M E[\hat{y}_i] \\
 &= E[\hat{y}_1]
 \end{aligned}$$

- 3) Variance: Reduced, since we're averaging over independent samples

$$\begin{aligned}
 Var[\hat{y}] &= Var\left[\frac{1}{M} \sum_{i=1}^M \hat{y}_i\right] \\
 &= \frac{1}{M^2} \sum_{i=1}^M Var[\hat{y}_i] \\
 &= \frac{1}{M} Var[\hat{y}_1]
 \end{aligned}$$

3.3 Effect of Correlation

Till now we assumed that the data are independent and then we saw the variance decreases by a factor of M . However, if the data are dependent on each other, the prediction made by each model is also not independent. Hence, we cannot claim the same. In fact, Variance is given by the following formula:

$$Var\left[\frac{1}{M} \sum_{i=1}^M \hat{y}_i\right] = \frac{1}{M} (1 - \rho) \sigma^2 + \rho \sigma^2$$

Where ρ is the correlation factor and σ is the standard deviation.

3.4 Random Forests

Random Forest is a bagging method with a learning model as decision trees. In addition, while selecting a bootstrap dataset for each decision tree, it chooses a random set of features on which the decision tree will split. This additional trick will ensure all predictions made by decision trees are not dependent.

Random forests are probably the best black-box machine learning algorithm. They often work well with no tuning whatsoever and are the most widely used algorithm in Kaggle competition.

4. WHAT IS BOOSTING?

Boosting methods work in the same spirit as bagging methods: we build a family of models that are aggregated to obtain a strong learner that performs better. However, unlike bagging which mainly aims at reducing variance, boosting is a technique that consists of fitting sequentially multiple weak learners in a very adaptative way: each model in the sequence is fitted giving more importance to observations in the dataset that were badly handled by the previous models in the sequence. In this way, the bias is lowered.

The base model in boosting is the starting weak learning model. The base model is often a high-bias and low-variance model because boosting mainly focuses on reducing bias (and maybe increasing the variance as boosting can cause overfitting). Boosting algorithms also use weighted dataset which is discussed below.

4.1 Weighted dataset

Till now the loss function was giving equal treatment to all data points i.e

$$Loss = \frac{1}{M} \sum_{i=1}^M L_i(\hat{y}_i, y_i)$$

The key idea of having a weighted dataset is that the learning algorithm can give more focus on data points with higher weights. In this way, we can control which data points should the learning algorithm put focus on. This is used in boosting since we can increase the weights of misclassified data points so that the next weak learner focuses on it.

4.2 Different type of boosting algorithm

There are mainly 5 types of boosting algorithms:

- 1) AdaBoost (Adaptive Boosting)
- 2) GBM (Gradient Boosting Machine)
- 3) XGBM (Extreme Gradient Boosting Machine)
- 4) LightGBM
- 5) CatBoost

4.3 AdaBoost

An adaptative boosting (often called 'Adaboost'), we try to define our ensemble model as a weighted sum of L weak learners:

$$H_T(X) = \sum_{t=1}^T \alpha_t h_t(x)$$

α_t and h_t are chosen such that we minimize the fitting loss which is given as follows:

$$\alpha_t, h_t = \arg \min_{\alpha_t, h_t} \frac{1}{N} \sum_{n=1}^N \text{loss}(y^{(n)}, H_{T-1}(x^{(n)}) + \alpha_t h_T(x^{(n)}))$$

We need to solve the above optimization problem. Here in this example, we will consider an exponential loss function given by:

$$Loss(y, \hat{y}) = \exp(-y\hat{y})$$

Let $\mathbb{L}(h(X^{(n)}) \neq y^{(n)}) = \frac{1}{2}(1 - h(X^{(n)}) \cdot y^{(n)})$. Using this we can simplify the optimization problem as follow:

$$\begin{aligned} \alpha_t, h_t &= \arg \min_{\alpha_t, h_t} \sum_{n=1}^N \exp(-y^{(n)}(H_{T-1}(x^{(n)}) + \alpha_t h_T(x))) \\ &= \arg \min_{\alpha_t, h_t} \sum_{n=1}^N \exp(-y^{(n)}(H_{T-1}(x^{(n)}))) \exp(-y^{(n)} \alpha_t h_T(x^{(n)})) \\ &= \arg \min_{\alpha_t, h_t} w_T^{(n)} \exp(-y^{(n)} \alpha_t h_T(x^{(n)})) \end{aligned}$$

Solving the optimization problem we get

$$\alpha_t = \frac{1}{2} \log \frac{1 - err_t}{err_t}$$

Where, $err_t = \frac{\sum_{n=1}^N w_t^{(n)} \mathbb{L}(h_t(X^{(n)}) \neq y^{(n)})}{\sum_{n=1}^N w^{(n)}}$

With that, we also find that h_t minimizes the weighted 0/1-loss i.e.

$$h_t = \arg \min_h \sum_{i=0}^N w_t^{(i)} \mathbb{L}(h(X^{(i)}) \neq y^{(i)})$$

We also find the relationship between weights as well:

$$w_{t+1}^{(n)} = w_t^{(n)} \exp(-y^{(n)} \alpha_t h_t(x^{(n)}))$$

The algorithm is given as follows:

- 1) Input: Data D with N data points, weak classifier WeakLearn (a classification procedure that returns a classifier h , e.g. best decision stump, from a set of classifiers H , e.g. all possible decision stumps), number of iterations T
- 2) Output: Classifier $H(x)$
- 3) Initialize all weights ($w \in \mathbb{R}^N$) to $\frac{1}{N}$
- 4) For $t = 1 \dots T$ do the following:
 - (a) Train the classifier h_t on the weighted dataset.
 - (b) Compute weighted error as follows:

$$err_t = \frac{\sum_{n=1}^N w^{(n)} \mathbb{L}(h_t(X^{(n)}) \neq y^{(n)})}{\sum_{n=1}^N w^{(n)}}$$

- (c) Compute classifier coefficient as follows:

$$\alpha_t = \frac{1}{2} \log \frac{1 - err_t}{err_t}$$

- (d) Update data weights as follow

$$w^{(n)} \leftarrow w^{(n)} \exp(-\alpha_t y^{(n)} h_t(X^{(n)}))$$

- (e) Normalize all weights

- 5) Return $H(x) = \text{sign}(\sum_{t=1}^T \alpha_t h_t(x))$

4.4 GBM

Just like AdaBoost, we try to define our ensemble model as a weighted sum of L weak learners. The main difference with adaptive boosting is in the definition of the sequential optimization process. Indeed, gradient boosting casts the problem into a gradient descent one: at each iteration, we fit a weak learner to the opposite of the gradient of the current fitting error concerning the current ensemble model. To put this mathematical perspective, the ensemble can be written as:

$$H_T(x) = H_{T-1}(x) - \alpha_T \nabla H_{T-1} E(H_{T-1})(x)$$

Where $E(.)$ is the fitting error. The coefficient α_t is computed following a one-dimensional optimization process (line-search to obtain the best step size α_t). Moreover, Adaptive boosting tries to solve at each iteration exactly the “local” optimization problem (find the best weak learner and its coefficient to add to the strong model), gradient boosting uses instead a gradient descent approach and can more easily be adapted to a large number of loss functions. Thus, gradient boosting can be considered as a generalization of AdaBoost to arbitrary differentiable loss functions.

4.5 XGBM

XGBoost is an enhanced version of the gradient boosting method. Firstly, it improves overfitting by using regularisation. Secondly, it improves the runtime speed by optimizing sorting using parallel running. Lastly, it uses the maximum depth of the decision tree as the parameter to prune the tree which reduces runtime significantly.

4.6 LightGBM

As the name suggests, Light Gbm further improves the runtime of the program by making the computing workload ‘light’. However, it can still maintain the same or higher level of model performance compared to other algorithms.

Light Gbm optimizes runtime speed and accuracy in mainly two ways:

- 1) It adopts the histogram-based algorithm, splitting the continuous variables into different buckets (rather than sorting them individually). This improves the runtime a lot.
- 2) It uses the leaf-wise tree growth method instead of the level-wise tree growth method

4.7 CatBoost

CatBoost stands for Categorical Boosting. It has the great feature of automatically handling categorical variables without the need to convert them into numerics.

CatBoost was developed most recently among the 5 boosting algorithms but very close to Light Gbm. It performs better when there are more categorical variables.

4.8 Comparison of different boosting algorithms

Points you can consider for choosing boosting algorithm

- 1) AdaBoost:
 - (a) Focuses on misclassified cases.
 - (b) It forms the foundation of boosting algorithm.
- 2) GBM:
 - (a) You can use any type of loss function.
 - (b) Gradient descent is used to minimize the loss function.
- 3) XGBoost
 - (a) Improves on overfitting
 - (b) Optimize running time by tree parallelism and tree pruning
- 4) LightGBM

- (a) further improves the speed of leaf-wise growth
- (b) Allow tuning of more parameter
- 5) CatBoost:
 - (a) Handles categorical features automatically.
 - (b) Works efficiently with categorical type data.

5. PROS AND CONS OF BOOSTING

- 1) Pros of boosting and cons of bagging:
 - (a) boosting reduces the bias whereas bagging doesn't decrease bias.
 - (b) In bagging we need to make sure predictions made by models are independent which is sometimes difficult in practice. Whereas, boosting doesn't need to worry about this.
- 2) Pros of bagging and cons of boosting:
 - (a) Since boosting is a sequential algorithm, it does not support parallel programming. whereas bagging can support parallel programming.
 - (b) Bagging decreases the variance whereas boosting can increase variance due to overfitting.

6. QUESTIONS

Subjective :

- 1) Can you justify briefly that the optimal prediction made for an input x is $y_* = E[y|x]$?
- 2) Can you prove $E[(\hat{y} - y)^2|x] = (E[t|x] - \hat{y})^2 + Var[t|x]$?
- 3) Using the equation proved in question 2 can you answer question 1 again?
- 4) Can you prove $E[(\hat{y} - y)^2] = (y_* - \hat{y})^2 + Var(\hat{y}) + Var(t)$?
- 5) Many machine libraries make use of parallel programming. Will libraries be more efficient in boosting algorithms or bagging algorithms? Justify
- 6) In Adaboost we constructed a optimization problem, can you show that $h_t(x)$ is the minimizer of the weighted 0/1-loss?
- 7) In bagging if the predictions of models are dependent then what will be the variance in terms of the Correlation factor (ρ) and standard deviation of one prediction(σ)?

Objective :

- 1) which of the following will be a good choice for the base model in boosting?
- a) SVM b) Neural network with 3 hidden layer c) decision stump d) decision tree

7. IMPLEMENTING ADABOOST

The full working code can be found on this [link](#). The AdaBoost model is made completely using numpy. We have used sklearn to import only the decision tree that can find the best split.

```

1 import numpy as np
2 from sklearn.tree import DecisionTreeClassifier
3
4 class AdaBoost(object):
5     def __init__(self) -> None:
6         self.hypothesis = None
7         self.hypothesis_weights = None
8
9     def train(self, X, y, num_iteration):
10         """
11         Input
12         -----
13         X: ndarray of shape (Num, features)
14         y: ndarray of shape (num,)
15         num_itteration: int
16         -----

```

```

17
18     Output:
19     -----
20     hist: list
21     -----
22     """
23     n = X.shape[0]
24     w = np.full(n, 1/n)
25     self.hypothesis = []
26     self.hypothesis_weights = []
27     hist = []
28     for t in range(num_iteration):
29         # training a model
30         stump = DecisionTreeClassifier(max_depth=1)
31         stump.fit(X,y,w)
32         self.hypothesis.append(stump)
33         # finding predicted value
34         y_pred = stump.predict(X)
35         # finding wieghted error
36         err = np.sum(w*(y_pred != y))/np.sum(w)
37         # finding alpha
38         alpha = (1/2) * np.log((1-err)/err)
39         self.hypothesis_weights.append(alpha)
40         #updating value of w
41         w = w*np.exp(-alpha*y*y_pred)
42         # computing the loss
43         loss = np.mean(w)
44         hist.append(loss)
45     return hist
46
47 def predict(self,X):
48     """
49     Input
50     -----
51     X: ndarray of shape (Num,features)
52     -----
53
54     Output:
55     -----
56     y: ndarray of shape (Num,)
57     -----
58     """
59     num = X.shape[0]
60     y = np.zeros(num)
61     for alpha,h in zip(self.hypothesis_weights,self.hypothesis):
62         y += alpha*(h.predict(X))
63     pos = (y >= 0).astype("int")
64     neg = (y < 0).astype("int")
65     y = pos - neg
66     return y

```
