

[Open in app](#)

Following ▾

529K Followers



This is your **last** free member-only story this month. [Upgrade for unlimited access.](#)

Ensemble Learning: Bagging & Boosting

How to combine weak learners to build a stronger learner to reduce bias and variance in your ML model



Fernando López · 2 days ago · 7 min read ★

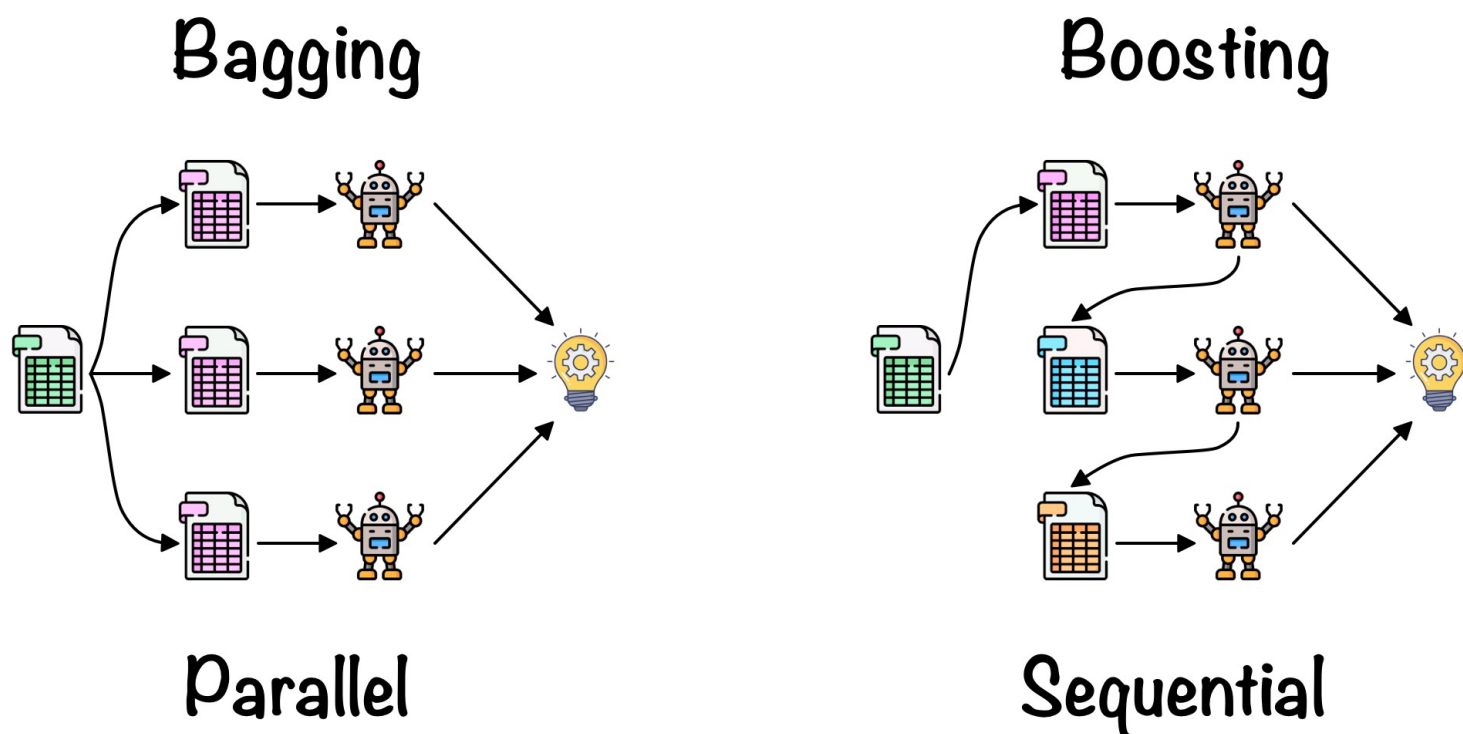


Figure 1. Bagging and Boosting | Spreadsheet, Robot and Idea icons by [Freepik](#) on [Flaticon](#)

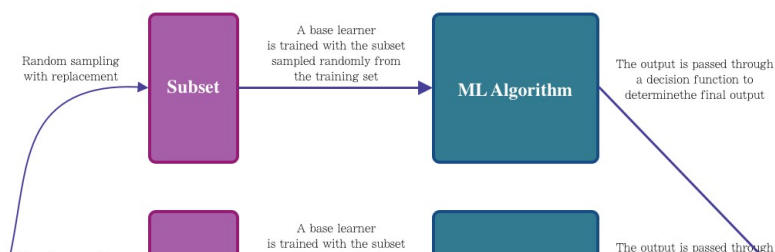
The bias and variance tradeoff is one of the key concerns when working with machine learning algorithms. Fortunately there are some **Ensemble Learning** based techniques that machine learning practitioners can take advantage of in order to tackle the bias and variance tradeoff, these techniques are **bagging** and **boosting**. So, in this blog we are going to explain how **bagging** and **boosting** works, what their components are and how you can implement them in your ML problem, thus this blog will be divided in the following sections:

- What is Bagging?
- What is Boosting?
- AdaBoost
- Gradient Boosting

What is Bagging?

Bagging or **Bootstrap Aggregation** was formally introduced by Leo Breiman in 1996 [3]. **Bagging** is an **Ensemble Learning** technique which aims to reduce the error learning through the implementation of a set of homogeneous machine learning algorithms. The key idea of **bagging** is the use of multiple base learners which are trained separately with a random sample from the training set, which through a voting or averaging approach, produce a more stable and accurate model.

The main two components of **bagging** technique are: the *random sampling with replacement* (**bootstrapping**) and the *set of homogeneous* machine learning algorithms (**ensemble learning**). The **bagging** process is quite easy to understand, first it is extracted “*n*” subsets from the training set, then these subsets are used to train “*n*” base learners of the same type. For making a prediction, each one of the “*n*” learners are feed with the test sample, the output of each learner is averaged (in case of regression) or voted (in case of classification). Figure 2 shows an overview of the **bagging** architecture.



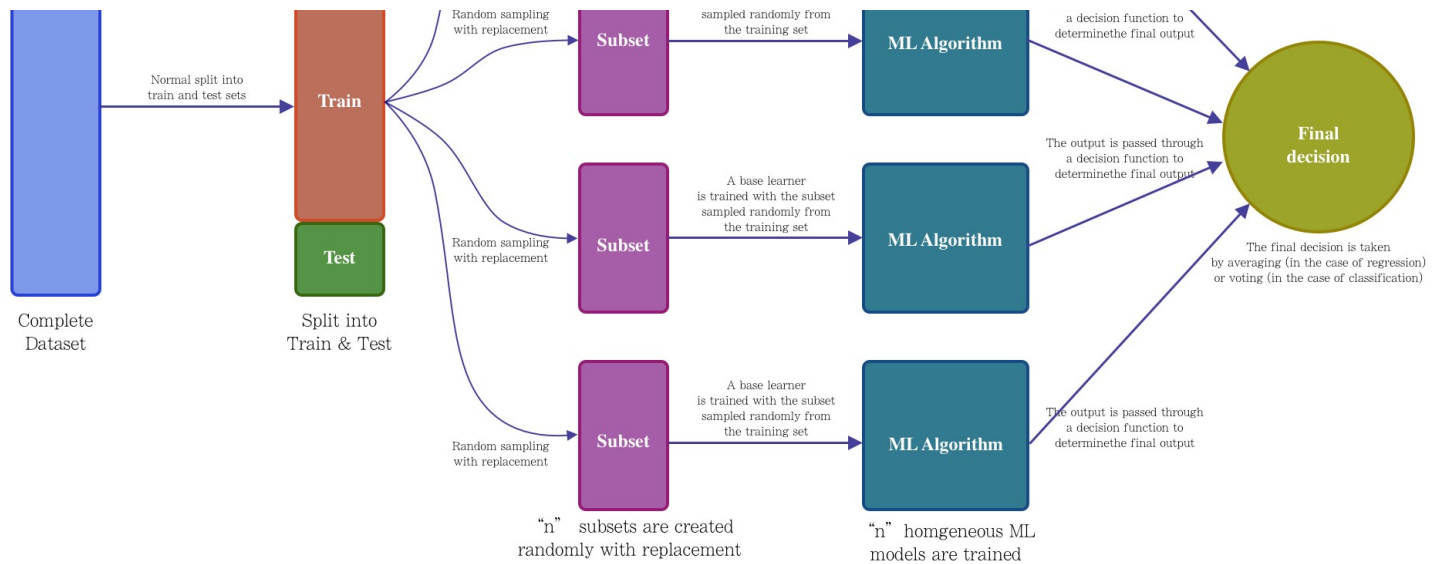


Figure 1. Bagging | Image by Author

It is important to notice that the number of *subsets* as well as the number of items per *subset* will be determined by the nature of your ML problem, the same for the type of ML algorithm to be used. In addition, Leo Breiman mention in his paper that he noticed that for classification problems are required more *subsets* in comparison with regression problems.

For implementing **bagging**, scikit-learn provides a function to do it easily. For a basic execution we only need to provide some parameters such as the *base learner*, the *number of estimators* and the *maximum number of samples* per subset.

```

1  # For this basic implementation, we only need these modules
2  from sklearn.datasets import load_breast_cancer
3  from sklearn.model_selection import train_test_split
4  from sklearn.tree import DecisionTreeClassifier
5  from sklearn.ensemble import BaggingClassifier
6
7  # Load the well-known Breast Cancer dataset
8  # Split into train and test sets
9  x, y = load_breast_cancer(return_X_y=True)
10 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=23)
11
12 # For simplicity, we are going to use as base estimator a Decision Tree with fixed parameters
13 tree = DecisionTreeClassifier(max_depth=3, random_state=23)
14
15 # The bagging ensemble classifier is initialized with:
16 # base_estimator = DecisionTree

```

```
17 # n_estimators = 5 : it's gonna be created 5 subsets to train 5 Decision Tree models
18 # max_samples = 50 : it's gonna be taken randomly 50 items with replacement
19 # bootstrap = True : means that the sampling is gonna be with replacement
20 bagging = BaggingClassifier(base_estimator=tree, n_estimators=5, max_samples=50, bootstrap=True)
21
22 # Training
23 bagging.fit(x_train, y_train)
24
25 # Evaluating
26 print(f"Train score: {bagging.score(x_train, y_train)}")
27 print(f"Test score: {bagging.score(x_test, y_test)}")
```

bagging.py hosted with  by GitHub

[view raw](#)

Code snippet 1. Bagging implementation

In the previous code snippet was created a *bagging based model* for the well know *breast cancer dataset*. As base learner was implemented a Decision Tree, 5 subsets were created randomly with replacement from the training set (to train 5 decision tree models). The number of items per subset were 50. By running it we will get:

```
Train score: 0.9583568075117371
Test score: 0.941048951048951
```

One of the key advantages of **bagging** is that it can be executed in parallel since there is no dependency between estimators. For *small datasets*, a few estimators will be enough (such as the example above), *larger dataset* may require more estimators.

Great, so far we've already seen what **bagging** is and how it works. Let's see what **boosting** is, its components and why it is related to **bagging**, let's go for it!

What is Boosting?

Boosting is an **Ensemble Learning** technique that, like **bagging**, makes use of a set of *base learners* to improve the stability and effectiveness of a ML model. The idea behind a **boosting** architecture is the generation of sequential hypotheses, where each hypothesis tries to improve or correct the mistakes made in the previous one [4]. The central idea of **boosting** is the implementation of *homogeneous ML algorithms* in a **sequential way**, where each of these ML algorithms tries to improve the stability of the model by

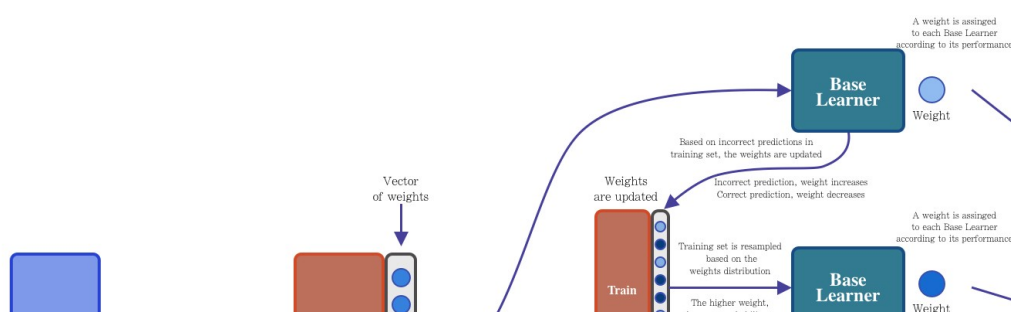
focusing on the errors made by the previous ML algorithm. The way in which the errors of each *base learner* is considered to be improved with the next *base learner* in the sequence, is the key differentiator between all variations of the **boosting** technique.

The **boosting** technique has been studied and improved over the years, several variations have been added to the core idea of boosting, some of the most popular are: **AdaBoost** (Adaptive Boosting), **Gradient Boosting** and **XGBoost** (Extreme Gradient Boosting). As mentioned above, the key differentiator between *boosting-based techniques* is the way in which errors are penalized (by modifying **weights** or minimizing a **loss function**) as well as how the data is sampled.

For a better understanding of the differences between some of the **boosting** techniques, let's see in a general way how **AdaBoost** and **Gradient Boosting** work, two of the most common variations of the boosting technique, let's go for it!

AdaBoost

AdaBoost is an algorithm based on the **boosting** technique, it was introduced in 1995 by Freund and Schapire [5]. **AdaBoost** implements a *vector of weights* to penalize those samples that were incorrectly inferred (by increasing the weight) and reward those that were correctly inferred (by decreasing the weight). Updating this *weight vector* will generate a distribution where it will be more likely to extract those samples with higher weight (that is, those that were incorrectly inferred), this sample will be introduced to the next *base learner* in the sequence. This will be repeated until a stop criterion is met. Likewise, each base learner in the sequence will have assigned a weight, the higher the performance, the higher the weight and the greater the impact of this base learner for the final decision. Finally, to make a prediction, each base learner in the sequence will be fed with the test data, each of the predictions of each model will be voted (for the classification case) or averaged (for the regression case). In Figure 3 we observe the descriptive architecture of the **AdaBoost** operation.



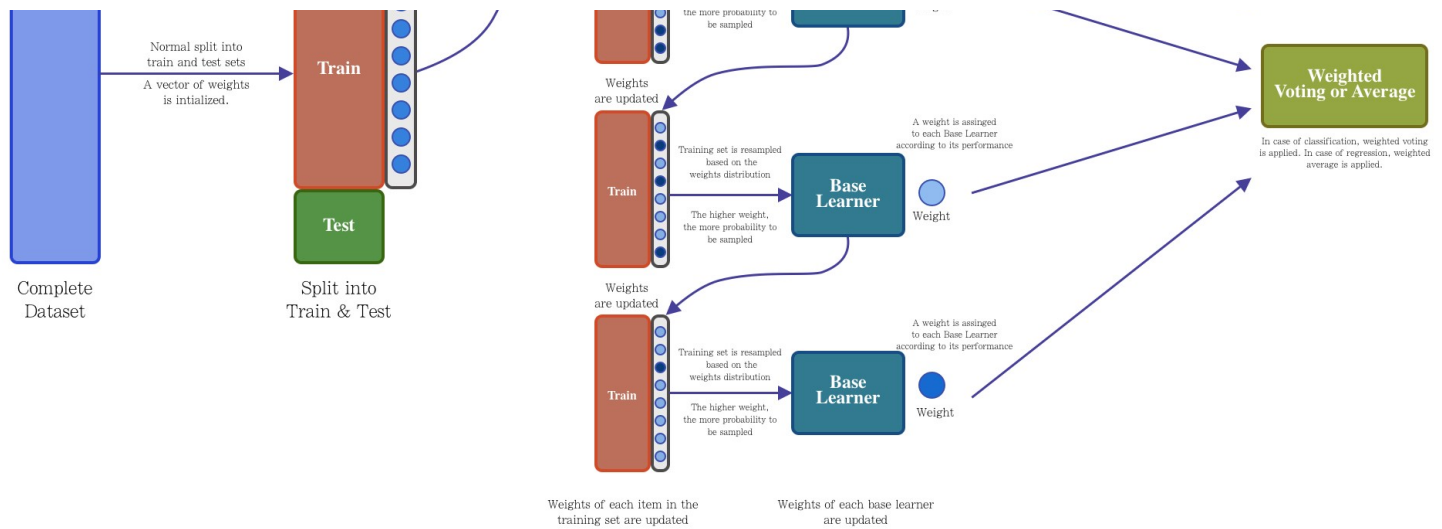


Figure 3. AdaBoost: a descriptive architecture | Image by Author

Scikit-learn provides the function to implement the **AdaBoost** technique, let's see how to perform a basic implementation:

```

1  # For this basic implementation, we only need these modules
2  from sklearn.datasets import load_breast_cancer
3  from sklearn.model_selection import train_test_split
4  from sklearn.tree import DecisionTreeClassifier
5  from sklearn.ensemble import AdaBoostClassifier
6
7  # Load the well-known Breast Cancer dataset
8  # Split into train and test sets
9  x, y = load_breast_cancer(return_X_y=True)
10 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=23)
11
12 # The base learner will be a decision tree with depth = 2
13 tree = DecisionTreeClassifier(max_depth=2, random_state=23)
14
15 # AdaBoost initialization
16 # It's defined the decision tree as the base learner
17 # The number of estimators will be 5
18 # The penalizer for the weights of each estimator is 0.1
19 adaboost = AdaBoostClassifier(base_estimator=tree, n_estimators=5, learning_rate=0.1, random_state=23)
20
21 # Train!
22 adaboost.fit(x_train, y_train)
23
24 # Evaluation
25 print(f"Train score: {adaboost.score(x_train, y_train)}")
26 print(f"Test score: {adaboost.score(x_test, y_test)}")

```

adaboost.py hosted with  by GitHub[view raw](#)

Code snippet 2. AdaBoost implementation

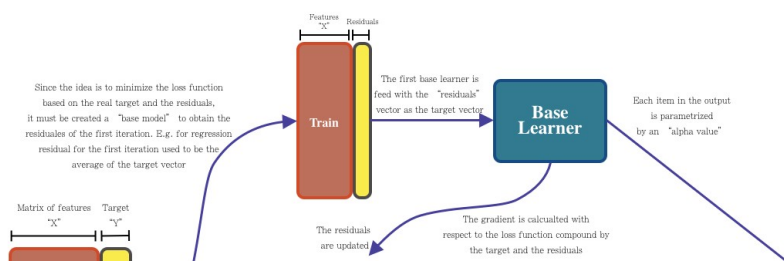
As we can see, the *base learner* that we are using is a *decision tree* (it is suggested that it be a decision tree, however, you can try some other ML algorithm), we are also defining only 5 estimators for the base learner sequence (this is enough for the toy dataset that we are trying to predict), running this we would obtain the following results:

```
Train score: 0.9694835680751174
Test score: 0.958041958041958
```

Great, we already saw in a general way how **AdaBoost** works, now let's see what about **Gradient Boosting** and how we can implement it.

Gradient Boosting

The **Gradient Boosting** method does not implement a *vector of weights* like **AdaBoost** does. As its name implies, it implements the calculation of the *gradient* for the optimization of a given loss function. The core idea of **Gradient Boosting** is based on minimizing the residuals of each learner base in a sequential way, this minimization is carried out through the calculation of the gradient applied to a specific loss function (either for classification or regression). Then each *base learner* added to the sequence will minimize the residuals determined by the previous *base learner*. This will be repeated until the error function is as close to zero or until a specified number of *base learners* is completed. Finally, to make a prediction, each of the *base learners* are fed with the test data whose outputs are parameterized and subsequently added to generate the final prediction.



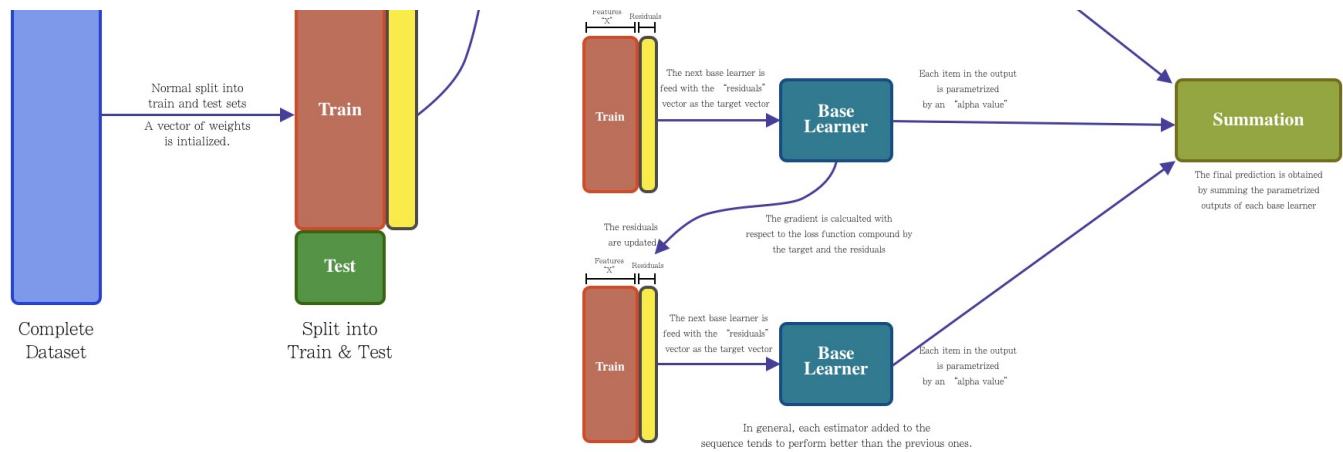


Figure 4. Gradient Boosting: a descriptive architecture | Image by author

Just like **Bagging** and **AdaBoost**, *scikit-learn* provides the function to implement **Gradient Boosting**, let's see how to make a basic implementation:

```

1  # For this basic implementation, we only need these modules
2  from sklearn.datasets import load_breast_cancer
3  from sklearn.model_selection import train_test_split
4  from sklearn.ensemble import GradientBoostingClassifier
5
6  # Load the well-known Breast Cancer dataset
7  # Split into train and test sets
8  x, y = load_breast_cancer(return_X_y=True)
9  x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=23)
10
11 # Gradient Boosting initialization
12 # The base learner is a decision tree as default
13 # The number of estimators is 5
14 # The depth for each decision tree is 2
15 # The learning rate for each estimator in the sequence is 1
16 gradientBoosting = GradientBoostingClassifier(n_estimators=5, learning_rate=1, max_depth=2, random_state=23)
17
18 # Train!
19 gradientBoosting.fit(x_train, y_train)
20
21 # Evaluation
22 print(f"Train score: {gradientBoosting.score(x_train, y_train)}")
23 print(f"Test score: {gradientBoosting.score(x_test, y_test)}")

```


Gradient Boosting works with *decision trees* by default, that is why in the implementation we do not define a specific base learner. We are defining that each tree in the sequence will have a maximum depth of 2, the number of trees will be 5 and the learning rate for each tree will be 0.1, running this we obtain:

```
Train score: 0.9906103286384976
Test score: 0.965034965034965
```

Fantastic, with this we finish this exploration on **bagging**, **boosting** and some **boosting** implementation, that's it!

Conclusion

In this blog we have seen two of the most widely implemented **Ensemble Learning** techniques.

As we have seen, **bagging** is a technique that performs random samples without replacement to train “n” *base learners*, this allows the model to be processed in parallel. It is because of this random sampling that **bagging** is a technique that mostly allows to reduce the variance. On the other hand, **boosting** is a sequentially constructed technique where each model in the sequence tries to focus on the error of the previous *base learner*. Although **boosting** is a technique that mainly allows to reduce the variance, it is highly prone to over-fitting the model.

References

[1] <https://quantdare.com/what-is-the-difference-between-bagging-and-boosting/>

[2] <https://arxiv.org/pdf/0804.2752.pdf>

[3] <https://link.springer.com/article/10.1023/A:1018054314350>

[4] <https://www.cs.princeton.edu/courses/archive/spr07/cos424/papers/boosting-survey.pdf>

[5] https://www.face-rec.org/algorithms/Boosting-Ensemble/decision-theoretic_generalization.pdf

Thanks to Ludovic Benistant.

[Bagging](#)

[Boosting](#)

[Ensemble Learning](#)

[Machine Learning](#)

[Editors Pick](#)

[About](#) [Help](#) [Legal](#)

Get the Medium app

