



Bagging v/s Boosting



Roshmita Dey

Follow

6 min read · Jan 30, 2024

197

1



Bagging and boosting are both ensemble learning techniques that aim to improve the performance of machine learning models by combining the predictions of multiple base learners. These approaches differ in their methodologies, strategies, and goals. In this detailed explanation, we'll delve into bagging and boosting, exploring their key concepts, algorithms, advantages, and potential challenges.

Bagging (Bootstrap Aggregating):

Bagging is a popular ensemble learning technique that focuses on reducing variance and improving the stability of machine learning models. The term "bagging" is derived from the idea of creating multiple subsets or bags of the training data through a process known as bootstrapping. Bootstrapping involves randomly sampling the dataset with replacement to generate multiple subsets of the same size as the original data. Each of these subsets is then used to train a base learner independently.

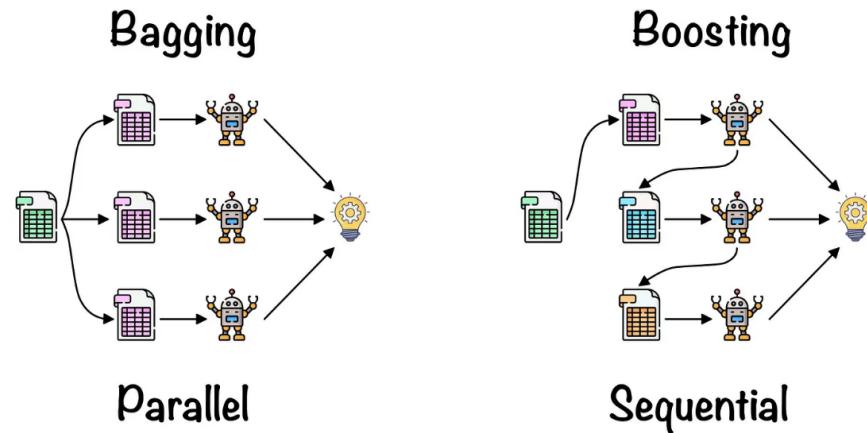
One of the primary goals of bagging is to reduce overfitting by exposing each base learner to slightly different variations of the training data. Since each subset is created by sampling with replacement, some instances may be duplicated while others may be omitted. This diversity helps the ensemble model generalize well to unseen data.

The most well-known algorithm for bagging is the Random Forest. In a Random Forest, a collection of decision trees is built, each trained on a different subset of the data. During the training process, each tree is constructed by selecting a random subset of features for each split, adding an extra layer of randomness and diversity to the ensemble. The final prediction is then made by averaging or taking a vote among the predictions of individual trees.

One key advantage of bagging is its ability to handle noisy data and outliers effectively. Since the ensemble model aggregates predictions from multiple base learners, the impact of outliers on the overall performance is reduced. Additionally, bagging is parallelizable, as each base learner can be trained independently, leading to efficient and scalable implementations.

Despite its strengths, bagging might not significantly improve the

performance of an already stable model or one that is not prone to overfitting. It is particularly useful when dealing with complex models that have high variance, such as deep decision trees or neural networks.



Boosting:

Boosting, like bagging, is an ensemble learning technique, but it aims to improve the performance of weak learners by combining them in a sequential manner. The core idea behind boosting is to give more weight to misclassified instances during the training process, enabling subsequent learners to focus on the mistakes made by their predecessors.

Unlike bagging, boosting does not rely on bootstrapped subsets of the data. Instead, it assigns weights to each instance in the training set and adjusts these weights throughout the boosting iterations. In each iteration, a new weak learner is trained on the data, and the weights of misclassified instances are increased. This allows the subsequent learner to pay more attention to the previously misclassified examples.

The most well-known boosting algorithm is AdaBoost (Adaptive Boosting). In AdaBoost, the weak learners are usually simple models with low predictive power, such as shallow decision trees or stumps (trees with a single split). Each weak learner is trained sequentially, and at each iteration, the weights of misclassified instances are increased, forcing the model to focus on the difficult-to-classify examples.

AdaBoost assigns a weight to each weak learner based on its performance, and the final prediction is made by combining the weighted predictions of all weak learners. Instances that are consistently misclassified by the ensemble receive higher weights, allowing subsequent weak learners to give more emphasis to these challenging cases.

Get Roshmita Dey's stories in your inbox

Join Medium for free to get updates from this writer.

Enter your email

Subscribe

One of the significant advantages of boosting is its ability to handle complex relationships in the data and improve the performance of weak learners significantly. Boosting often outperforms bagging when it comes to reducing both bias and variance. However, boosting is more sensitive to noisy data and outliers compared to bagging.

Top highlight

Differences Between Bagging and Boosting:

Sequential vs. Parallel:

- Bagging: The base learners are trained independently in parallel, as each learner works on a different subset of the data. The final prediction is typically an average or vote of all base learners.
- Boosting: The base learners are trained sequentially, and each learner focuses on correcting the mistakes of its predecessors. The final prediction is a weighted sum of the individual learner predictions.

Data Sampling:

- Bagging: Utilizes bootstrapping to create multiple subsets of the training data, allowing for variations in the training sets for each base learner.
- Boosting: Assigns weights to instances in the training set, with higher weights given to misclassified instances to guide subsequent learners.

Weighting of Base Learners:

- Bagging: All base learners typically have equal weight when making the final prediction.
- Boosting: Assigns different weights to each base learner based on its performance, giving more influence to learners that perform well on challenging instances.

Handling Noisy Data and Outliers:

- Bagging: Robust to noisy data and outliers due to the averaging or voting mechanism, which reduces the impact of individual errors.
- Boosting: More sensitive to noisy data and outliers, as the focus on misclassified instances might lead to overfitting on these instances.

Model Diversity:

- Bagging: Aims to create diverse base learners through random subsets of the data and, in the case of Random Forests, random feature selection for each tree.
- Boosting: Focuses on improving the performance of weak learners sequentially, with each learner addressing the weaknesses of its predecessors.

Bias and Variance:

- Bagging: Primarily reduces variance by averaging predictions from multiple models, making it effective for models with high variance.

multiple models, making it effective for models with high variance.

- Boosting: Addresses both bias and variance, with a focus on reducing bias by sequentially correcting mistakes made by weak learners.

Advantages of Bagging:

1. Variance Reduction: Bagging is effective in reducing variance, making it particularly useful for unstable models or models prone to overfitting.
2. Robustness to Noisy Data: The ensemble nature of bagging makes it robust to noisy data and outliers, as the impact of individual errors is mitigated by the aggregation of predictions.
3. Parallelization: Bagging algorithms, such as Random Forests, can be parallelized, leading to efficient implementations and faster training times on distributed computing systems.
4. Versatility: Bagging can be applied to a wide range of base learners, making it a versatile technique applicable to different types of models.

Advantages of Boosting:

1. Improved Model Accuracy: Boosting often leads to improved model accuracy compared to individual weak learners, as it focuses on correcting errors made by previous models.
2. Handling Complex Relationships: Boosting is effective in capturing complex relationships in the data, making it suitable for tasks where the underlying patterns are intricate.
3. Bias and Variance Reduction: Boosting addresses both bias and variance, making it suitable for models with high bias or models that struggle with both underfitting and overfitting.
4. Adaptability to Weak Learners: Boosting can boost the performance of weak learners, allowing for the creation of strong predictive models from simple base learners.

Challenges and Considerations:

1. Overfitting in Boosting: Boosting may be susceptible to overfitting, especially when the focus on correcting misclassifications is too aggressive. Tuning parameters like learning rate can mitigate this risk.
2. Sensitivity to Noisy Data: Boosting can be sensitive to noisy data, as it assigns higher weights to misclassified instances, potentially leading to overemphasis on noisy patterns.
3. Computational Complexity: Some boosting algorithms, especially when using complex base learners, can be computationally expensive and may require more training time compared to bagging.
4. Interpretability: Ensembles generated by boosting, particularly with a large number of weak learners, might become complex and challenging to interpret compared to simpler bagging ensembles.

5. Data Requirements: Boosting may require more substantial amounts of data compared to bagging, as it focuses on iteratively improving the model's performance, which requires diverse and informative instances.

Conclusion:

In summary, bagging and boosting are both powerful ensemble learning techniques that address different aspects of model performance. Bagging, with its emphasis on reducing variance and providing robustness to noisy data, is well-suited for unstable models. Random Forests, a popular bagging algorithm, have been widely used for tasks like classification and regression.

On the other hand, boosting excels in improving the accuracy of weak learners, handling complex relationships in the data, and addressing both bias and variance. AdaBoost, a well-known boosting algorithm, has been successful in various applications, including face detection and object recognition.

The choice between bagging and boosting depends on the characteristics of the dataset, the nature of the problem, and the underlying model. In practice, it's not uncommon to experiment with both techniques and choose the one that yields the best results for a specific task. Additionally, variations and hybrid approaches, such as Gradient Boosting, have emerged to combine the strengths of both bagging and boosting, offering a more sophisticated and flexible ensemble learning framework. Understanding the nuances of bagging and boosting provides practitioners with valuable tools to enhance the performance of machine learning models across diverse applications.

Bagging Boosting Ensemble Ensemble Learning

197 1

Bookmark Share



Written by Roshmita Dey

425 followers · 32 following

Follow

Working as a Data Scientist in one of the leading Global banks, my expertise is in the field of Statistics and proficiency in Python, SQL and PySpark

Responses (1)



Write a response

What are your thoughts?

Great explanation!

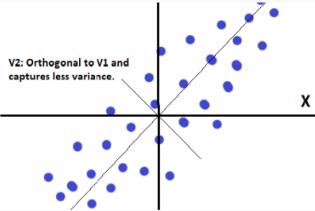
 4 [Reply](#)

More from Roshmita Dey

 Roshmita Dey

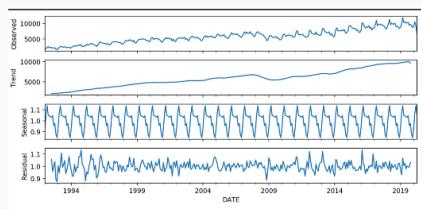
PySpark: Transformations v/s Actions

In PySpark, transformations and actions are two fundamental types of operations that yo...

Dec 9, 2023  90  2 Roshmita Dey

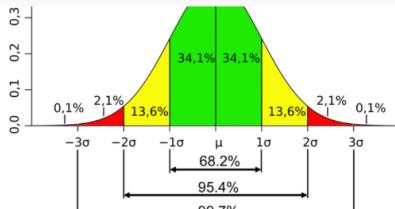
Understanding Principal Component Analysis (PCA)

Principal Component Analysis, or PCA, is a fundamental technique in the realm of data...

Oct 6, 2023  213  4 Roshmita Dey

Time Series Decomposition

Time series decomposition is a vital technique that helps in understanding the...

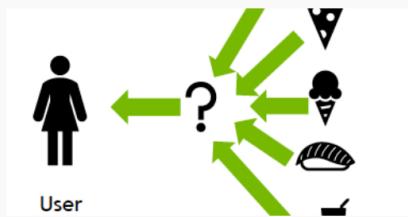
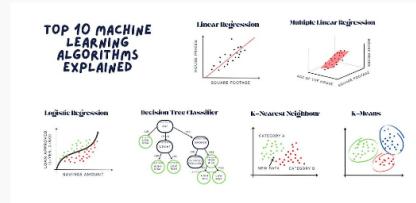
Jan 13, 2024  132  1 Roshmita Dey

Normal Distribution and its Significance

The normal distribution, also known as the Gaussian distribution or bell curve, is a...

Dec 8, 2023  37[See all from Roshmita Dey](#)

Recommended from Medium



In Learning Data by Rita Angelou

10 ML Algorithms Every Data Scientist Should Know—Part 1

I understand well that machine learning might sound intimidating. But once you break dow...

Jun 10 27



Sohail Saifi

The Mathematical Concept Behind Every Recommendation Algorithm

You know that moment when Netflix suggests exactly the show you didn't know you...

Jun 24 100 3



Anirban Mukherjee

50 Machine Learning Projects That Will Get You Hired in 2025

Not a member yet? Read for free here.

May 23 363 5



In The Algorithmic Minds by Rahul Agarwal

Crack ML System Design Interviews Like a Pro—Part II

Netflix Real-World Case Study

Feb 19 15



In Data Science Collective by James Wilkins

You're using ChatGPT wrong. Here's how to prompt like a pro

Smarter prompts lead to smarter responses.

Jun 5 2.6K 178



S. Moazeni, PhD

XGBoost: Optimized Gradient Boosting for Supervised Learning

Gradient Boosting is a method for regression and classification. It can handle regression...

Jul 14



[See more recommendations](#)