

# Demystifying Bagging in Machine Learning :

Bagging, short for Bootstrap Aggregating, is a powerful ensemble learning technique that has become a cornerstone of many high-performing machine learning systems. As the name suggests, Bagging is composed of two core steps: Bootstrapping and Aggregation. Let's break it down and understand why it's so effective—especially for models prone to overfitting.

## Step 1: Bootstrapping – Sampling the Data

Bootstrapping is a statistical technique where we draw random samples with replacement from the original dataset. Each of these samples—often called bootstrap samples—acts as a training set for a separate model in the ensemble.

Let's consider a practical example. Suppose we have a classification dataset with 10,000 rows. Using Bagging, we decide to train 100 base models—say, 100 Decision Trees.

- For each model M1, M2, ... .., M100, we randomly draw (with replacement) a sample of, say, 1,000 rows from the original dataset.
- This number can be equal to or less than the total dataset size and is chosen based on the problem and resource constraints.
- Because the sampling is done with replacement, some data points may appear multiple times in the same sample, while others may not appear at all. This introduces variability across the models.
- Each base model is trained independently on its own bootstrap sample.

This process ensures diversity among the models, which is crucial for the success of any ensemble method.

## Step 2: Aggregation – Making a Collective Decision

Once all the base models are trained, we use them to make predictions on new data points.

- In the case of classification, each model votes for a class label.
- The final prediction is based on majority voting (i.e., the mode of all predictions).
- In the case of regression, we typically take the average of all model predictions.

This step—Aggregation—combines the outputs from all base models to produce a robust final prediction.

## Why Bagging Works: Reducing Variance Without Increasing Bias

One of the main challenges in machine learning is the bias-variance tradeoff. Most models fall into one of two categories:

- High bias, low variance (e.g., Linear Regression): Underfits the data
- Low bias, high variance (e.g., Decision Trees, KNN, SVM): Overfits the data

Bagging is particularly effective with models that have low bias but high variance, such as:

- Decision Trees (especially with `max_depth=None`)
- K-Nearest Neighbors (KNN)
- Support Vector Machines (SVM)

These models often perform very well on the training data but suffer from overfitting when tested on unseen data.

Here's where Bagging comes in:

- By training each model on a slightly different dataset (thanks to bootstrapping), we introduce variation.
- But when these diverse models vote together, the variance cancels out, and we get a more stable, generalizable prediction.
- If new data points are added or the dataset changes slightly, no single model is overly affected—because each one only sees a subset of the data. This reduces model sensitivity to noise and outliers.

Result: Lower Variance, Maintained Bias → Better Generalization

## Key Takeaways

- Bagging = Bootstrapping + Aggregation
- Helps reduce variance while maintaining low bias
- Best suited for high-variance models like Decision Trees
- Used in algorithms like Random Forest, which is essentially Bagging with randomized feature selection

## Final Thoughts

In a world where overfitting is a constant threat, Bagging provides a practical and effective defense. By embracing randomness in the training process and consensus in prediction, Bagging transforms unstable models into robust performers.

If you're building a machine learning model that's overfitting, consider trying Bagging—it might

just be the ensemble solution you need.