

# Ensemble Learning in Machine Learning: A Comparative Study of Bagging, Boosting, and Stacking

Name: Prudhvi Nath Thota

Student Id: 23096383

GitHUB LINK : <https://github.com/thotaprudhvinath/Ensemble-Learning-in-Machine-Learning-A-Comparative-Study-of-Bagging-Boosting-and-Stacking>

## Introduction:

Ensemble learning is a powerful machine learning technique that combines multiple models to improve predictive performance and robustness. Instead of relying on a single weak model, ensemble methods aggregate the predictions of multiple models, leading to better accuracy and stability. The three primary ensemble learning techniques are Bagging, Boosting, and Stacking.

## 1. Bagging

Bagging is an ensemble learning technique that involves training multiple models in parallel on different subsets of the dataset. The final prediction is made by averaging the predictions (for regression) or using majority voting (for classification).

### Key Characteristics:

- Reduces variance and prevents overfitting.
- Works well with high-variance models like decision trees.
- Random Forest is a popular example of Bagging.

## 2. Boosting

Boosting builds models sequentially, with each new model correcting the mistakes of the previous ones. The method gives more importance to misclassified instances, improving overall performance.

### Key Characteristics:

- Reduces bias and variance.
- Models are trained in sequence, unlike Bagging.
- Gradient Boosting, AdaBoost, and XGBoost are popular Boosting algorithms.

### 3. Stacking

Stacking is a more complex ensemble technique that combines multiple models by training a meta-learner to aggregate their predictions. It leverages the strengths of different models to enhance performance.

#### Key Characteristics:

- Uses diverse base models (e.g., decision trees, SVM, neural networks).
- Employs a meta-learner (e.g., logistic regression) to optimize predictions.
- Can outperform individual models but requires careful tuning.

#### Comparative Analysis of Ensemble Techniques

Feature	Bagging	Boosting	Stacking
Training Approach	Parallel	Sequential	Layered
Overfitting Risk	Low	Moderate	Moderate
Computational Complexity	Moderate	High	High
Robustness	High	Medium	High
Example Algorithms	Random Forest	XGBoost, AdaBoost	Blending, Meta-Learner Models

#### Applications of Ensemble Learning

1. **Finance:** Credit risk assessment, fraud detection.
2. **Healthcare:** Disease diagnosis, medical imaging analysis.
3. **Marketing:** Customer segmentation, churn prediction.
4. **E-commerce:** Product recommendation systems.

The dataset "processed.cleveland.data" comes from the Cleveland Heart Disease dataset, which is part of the UCI Machine Learning Repository. Here's a brief description

## Dataset Description

- Domain: Medical/Healthcare
- Purpose: Predicting the presence of heart disease in a patient based on various medical attributes.
- Number of Instances: 303
- Number of Features: 14 (including the target variable)

## Feature Overview

Feature Name	Description	Type
age	Age of the patient	Numeric
sex	Gender (1 = Male, 0 = Female)	Categorical
cp	Chest pain type (1-4)	Categorical
trestbps	Resting blood pressure (mm Hg)	Numeric
chol	Serum cholesterol (mg/dL)	Numeric
fbs	Fasting blood sugar > 120 mg/dL (1 = Yes, 0 = No)	Categorical
restecg	Resting electrocardiographic results (0-2)	Categorical
thalach	Maximum heart rate achieved	Numeric
exang	Exercise-induced angina (1 = Yes, 0 = No)	Categorical
oldpeak	ST depression induced by exercise	Numeric
slope	Slope of peak exercise ST segment (0-2)	Categorical
ca	Number of major vessels (0-3) colored by fluoroscopy	Numeric
thal	Thalassemia (3 = Normal, 6 = Fixed defect, 7 = Reversible defect)	Categorical
target	Presence of heart disease (0 = No, 1-4 = Yes)	Categorical

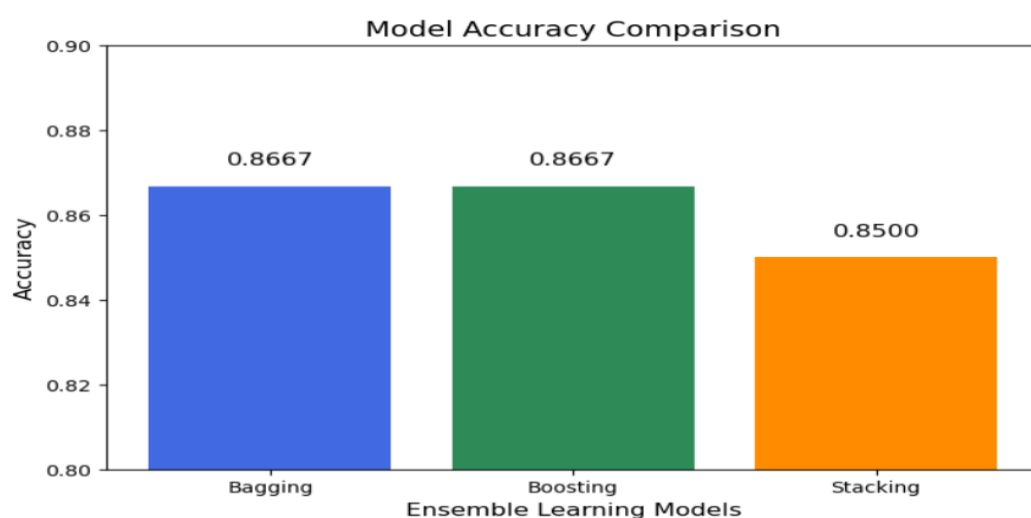
## Data Cleaning Notes

- **Missing Values:** The dataset contains missing values represented by "?", which should be handled before training. For Numerical data, Missing values are replaced by Mean and For Categorical data, Most Frequently repeated value which is mode are replaced.
- **Feature Engineering:** The target variable is originally multi-class (0-4) but is often converted into a binary classification problem (0 = No disease, 1 = Disease).

This dataset is widely used in machine learning and medical analytics to develop predictive models for heart disease detection.

### 1.Creating Model and Comparing

Three machine learning models—Bagging, Boosting, and Stacking—were trained using the dataset. The trained models were then applied to the test set ( $y_{test}$ ), generating predictions ( $y_{test\_predict}$ ). The accuracy of each model was subsequently computed and compared. Below is a graphical representation illustrating the variations in accuracy among the three models.



## 2. Evaluation of Model Performance

Below is the table that compares three ML Models, Bagging, Boosting, and Stacking ensemble methods based on training speed, accuracy, and overfitting.

ML Model	Training Speed	Accuracy Score	Overfitting
Bagging	Fast (Parallel Training)	0.867	Low
Boosting	Slow (Sequential Training)	0.867	Medium
Stacking	Moderate	0.850	Medium

## Conclusion

Ensemble learning methods significantly enhance model performance by leveraging multiple models' strengths. Bagging is ideal for reducing variance, Boosting improves weak learners by addressing errors sequentially, and Stacking provides an optimized final prediction by combining diverse models. The choice of technique depends on the dataset characteristics, computational resources, and the problem at hand.

## References

- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189-1232.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning*. Springer Science & Business Media.
- Wolpert, D. H. (1992). Stacked generalization. *Neural Networks*, 5(2), 241-259.