

Ensemble Learning

Ensemble learning is a machine learning technique where we combine predictions from multiple models (called *base learners* or *weak learners*) to produce a more accurate and robust final model.

Why Ensemble Learning?

1. Reduces Overfitting – combining models smooths out individual biases.
2. Improves Accuracy – multiple models correct each other's mistakes.
3. More Robust – works well even if individual models are weak.

Sources of Error in Machine Learning

1. Bias

- Error due to wrong assumptions in the learning algorithm.
- Model is too simple to capture the underlying patterns.
- Leads to underfitting (poor performance on training & test data).

☞ Example:

- Trying to fit a linear model to data that's clearly non-linear.
-

2. Variance

- Error due to too much sensitivity to training data.
- Model is too complex, learns noise along with signal.
- Leads to overfitting (great on training, poor on test).

☞ Example:

- Deep decision tree that perfectly memorizes training data but fails on unseen data.
-

3. Irreducible Error (Noise)

- Error caused by factors outside the model's control (randomness in data).
- Even a perfect model cannot eliminate this.
- Comes from measurement errors, missing features, randomness in real-world systems.

☞ Example:

- Two students with same study hours may get different exam scores due to stress, luck, or health.

Pros of Ensemble Learning

✓ 1. Higher Accuracy

- Combines multiple models → usually outperforms individual models.
- Example: Random Forest often beats a single Decision Tree.

✓ 2. Reduces Overfitting (Variance)

- Bagging methods like Random Forest average out noise → more stable.

✓ 3. Reduces Underfitting (Bias)

- Boosting methods (AdaBoost, Gradient Boosting) improve weak learners step by step.

✓ 4. Works for Both Classification & Regression

- Can be applied to almost all supervised learning tasks.

✓ 5. Robustness

- Less sensitive to noise/outliers than a single model (especially bagging).

Cons of Ensemble Learning

✗ 1. Computationally Expensive

- Training many models = more CPU/GPU time, memory.
- Boosting is sequential → slower training.

✗ 2. Less Interpretable

- A single decision tree is easy to explain.
- A forest of 500 trees or stacked models = hard to interpret (black-box).

✗ 3. Risk of Overfitting (Boosting)

- If not tuned properly, boosting methods can overfit noisy data.

✗ 4. Longer Inference Time

- Making predictions requires running multiple models, which can be slow in real-time applications.

✕ 5. Complexity in Implementation

- Simple models (linear regression, single tree) are easy.
- Ensembles need careful design (number of learners, depth, learning rate, etc.).

Use ensembles when you need high accuracy and robustness.

Avoid ensembles if you need a very fast, simple, or interpretable model.

What is a Bootstrap Sample?

A bootstrap sample is a dataset created by:

1. Taking the original dataset of size N .
2. Randomly selecting N samples with replacement.

Ensemble Learning Architectures

Ensemble learning can be broadly grouped into two strategies:

◆ 1. Parallel Ensemble Learning

- Definition: Base learners are trained independently at the same time (parallel).
- Key Idea: Errors of one learner do not affect the training of others.
- Goal: Reduce variance by averaging out errors.
- How it works:
 - Divide the data (using bootstrapping or random subsets).
 - Train multiple learners in parallel.
 - Aggregate their predictions (majority vote / averaging).

✓ Examples:

- Bagging (Bootstrap Aggregating)
 - Random Forest
 - Voting Classifier / Averaging Ensembles
-

◆ 2. Sequential Ensemble Learning

- Definition: Base learners are trained one after another (sequentially).

- Key Idea: Each new learner depends on the performance of the previous ones.
- Goal: Reduce bias by focusing more on misclassified or poorly predicted samples.
- How it works:
 - Train the first learner.
 - Give more weight (importance) to the mistakes of that learner.
 - Train the next learner on the updated distribution.
 - Combine learners in a weighted manner.

✓ Examples:

- AdaBoost (updates sample weights)
- Gradient Boosting (fits residual errors using gradients)
- XGBoost, LightGBM, CatBoost

Comparison: Parallel vs Sequential

Aspect	Parallel Ensemble	Sequential Ensemble
Training	Base learners trained independently	Base learners trained one after another
Dependency	No dependency between models	Each model depends on previous
Goal	Reduce variance (stabilize)	Reduce bias (improve accuracy)
Computation	Can be parallelized (faster)	Sequential → slower, harder to parallelize
Examples	Bagging, Random Forest, Voting	AdaBoost, Gradient Boosting, XGBoost

Types of Ensemble Methods

Ensemble methods can be broadly categorized into three main types:

1. Bagging (Bootstrap Aggregating)

- Idea: Train multiple models independently on random subsets of data (with replacement).
- Goal: Reduce variance (i.e., make model more stable).

- How it works:
 - Create several bootstrap datasets.
 - Train a base learner (e.g., decision tree) on each subset.
 - Aggregate predictions (majority vote for classification, average for regression).

◆ Examples:

- Bagging Classifier
 - Random Forest (Bagging + Random Feature Selection)
-

2. Boosting

- Idea: Train models sequentially, each one focusing on correcting the mistakes of the previous.
- Goal: Reduce bias (make weak learners stronger).
- How it works:
 - Start with a weak model.
 - Give more weight to misclassified samples.
 - Train the next model to fix those errors.
 - Combine models with weighted voting/averaging.

◆ Examples:

- AdaBoost (Adaptive Boosting)
 - Gradient Boosting Machines (GBM)
 - XGBoost
 - LightGBM
 - CatBoost
-

3. Stacking (Stacked Generalization)

- Idea: Use multiple models (called base learners) and combine their outputs using a meta-learner.
- Goal: Leverage different algorithms' strengths.
- How it works:

- Train multiple models (e.g., logistic regression, decision tree, SVM).
- Their predictions become new input features.
- A meta-model (often linear regression or another ML model) learns how to best combine them.

◆ Examples:

- Stacking Classifier
- Stacking Regressor

Bagging vs Boosting

Aspect	Bagging	Boosting
Meaning	<i>Bootstrap Aggregating</i> – trains multiple models independently on random subsets.	Sequential ensemble – trains models one after another, each correcting previous errors.
Training	Parallel (independent learners).	Sequential (each learner depends on previous).
Data Sampling	Uses bootstrap sampling (random sampling with replacement).	Uses the entire dataset, but re-weights samples (harder-to-predict samples get more focus).
Goal	Reduce variance (stability).	Reduce bias (accuracy).
Base Learner	Usually high-variance models like Decision Trees (deep trees).	Usually weak learners like shallow Decision Stumps (depth=1).
Combination	Simple averaging (regression) or majority vote (classification).	Weighted combination (stronger learners get higher weight).
Overfitting	Less prone to overfitting (averaging smooths results).	More prone to overfitting if not tuned (focuses strongly on errors).
Computation	Can be parallelized (faster training).	Sequential → cannot be parallelized easily (slower).
Examples	Random Forest, Bagging Classifier.	AdaBoost, Gradient Boosting, XGBoost, LightGBM, CatBoost.

What is Random Forest?

Random Forest is an **ensemble learning algorithm** that builds many **decision trees** and combines their outputs (by averaging or voting).

How Random Forest Works

1. Bootstrap Sampling (Bagging):

- For each tree, take a random sample (with replacement) of the dataset.

2. Random Feature Selection:

- At each node split, instead of using **all features**, choose a **random subset of features**.
- This increases diversity among trees.

3. Grow Decision Trees:

- Each tree is trained independently on its bootstrap sample with random features.

4. Aggregation:

- For **classification** → majority vote across trees.
- For **regression** → average predictions across trees.

Why Random Forest is Powerful

✓ Reduces Variance:

- Single decision tree = high variance (overfits easily).
- Random Forest = average of many trees → stable & robust.

✓ Handles High-Dimensional Data:

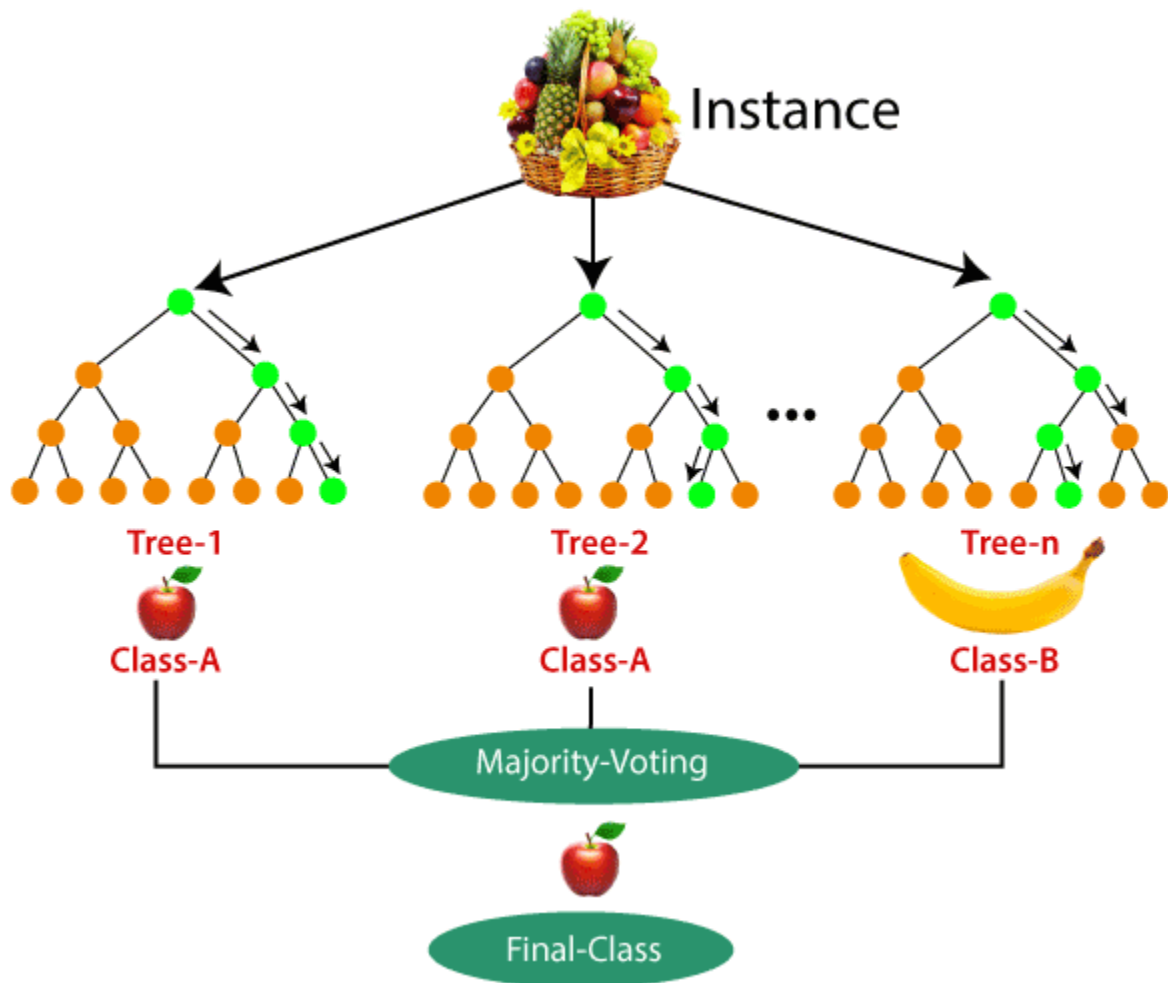
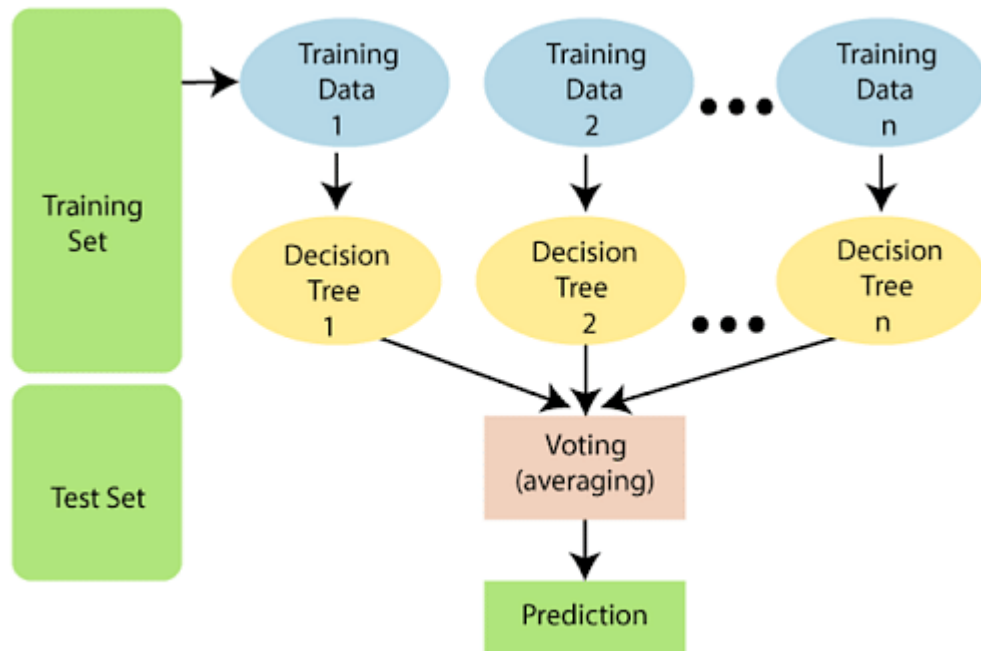
- Works well with lots of features.

✓ Feature Importance:

- Can rank features by importance (based on how much they reduce impurity).

✓ Works with Missing Data:

- Can still work even if some data is missing.



Hyperparameter	Effect
<code>n_estimators</code>	Number of trees → higher = better but slower
<code>max_depth</code>	Max depth of each tree → controls overfitting
<code>min_samples_split</code>	Min samples to split a node
<code>min_samples_leaf</code>	Min samples at leaf → avoids tiny leaves
<code>max_features</code>	Number of features to consider at each split
<code>bootstrap</code>	Whether to sample with replacement
<code>max_samples</code>	Number of samples per tree (if <code>bootstrap=True</code>)
<code>class_weight</code>	Handle imbalanced classification
<code>n_jobs</code>	Use multiple CPU cores
<code>random_state</code>	Reproducibility

`from sklearn.ensemble import RandomForestClassifier`

Hyperparameters of Random Forest & When to Tune

1. `n_estimators` (Number of trees)

- **What it does:** More trees = more stable predictions (less variance), but slower training/inference.
- **When to tune:**
 - If the model seems unstable or highly variable → increase.
 - If training is too slow → reduce.
- **Rule of thumb:** Start with 100–300. Increase until performance stabilizes.

2. `max_depth` (Maximum depth of each tree)

- **What it does:** Controls how deep trees can grow. Deep trees = low bias but high variance.

- **When to tune:**
 - If **overfitting** → decrease depth.
 - If **underfitting** → allow deeper trees.
 - **Rule:** Tune when dataset is complex or noisy.
-

3. `min_samples_split` & `min_samples_leaf`

- **What they do:** Control how many samples are needed to split a node / remain in a leaf.
 - **When to tune:**
 - If overfitting → increase these (forces trees to be shallower & smoother).
 - If underfitting → decrease them (allows more detailed splits).
 - **Rule:** Small datasets → higher values; large datasets → lower values.
-

4. `max_features` (Number of features considered at each split)

- **What it does:** Adds randomness, prevents trees from being identical.
 - **When to tune:**
 - If features are **highly correlated** → reduce `max_features`.
 - If model is underfitting → increase `max_features`.
 - **Rule:**
 - Classification → default = $\sqrt{n_features}$.
 - Regression → default = $n_features/3$.
-

5. `bootstrap`

- **What it does:** Whether to use bootstrap samples (sampling with replacement).
 - **When to tune:**
 - Usually keep `True`.
 - If dataset is small → test `False` (each tree sees all data).
-

6. `class_weight` (for classification)

- **What it does:** Handles class imbalance.

- **When to tune:**
 - If dataset has skewed classes → use balanced or custom weights.
-