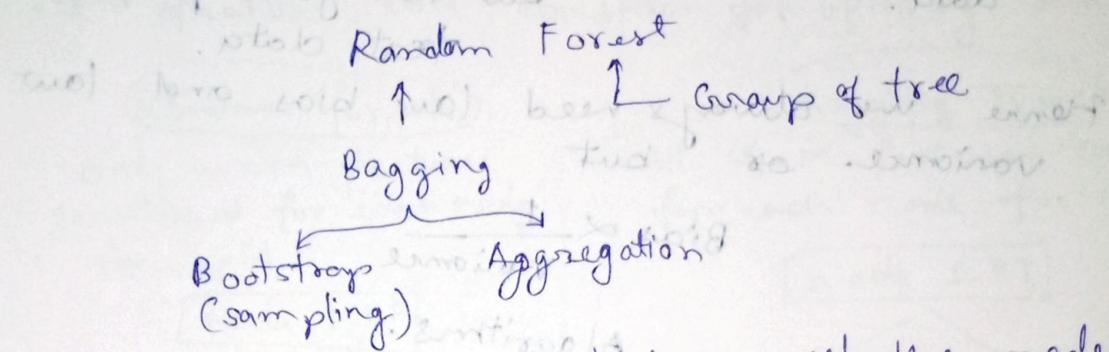


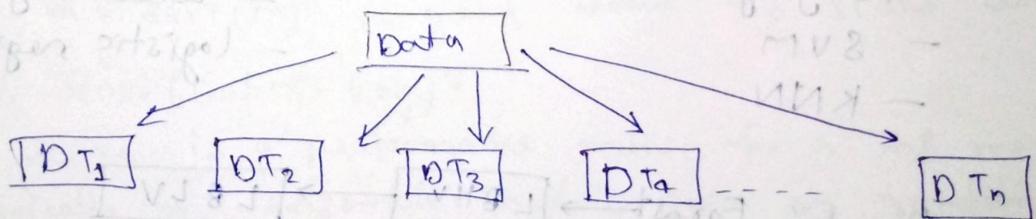
19-Oct-2023

RANDOM FOREST

- Random Forest is the ~~apply~~ for both classification and regression.
- Random Forest is the ensemble technique.
- it is a part of bagging techniques



- If in the Bagging Technique all the model are decision tree then this technique is called Random Forest.



where $DT_n = \text{Decision Trees (model)}$

- for each model we provide a subset of data to of data (original)
- if assume if original data has 1000 row then we provide each decision tree to 500 row.
- these 500 row is decided by these method

- (i) Row Sampling (Pasting)
- (ii) column Sampling (Random subspaces)
- (iii) row and column sampling (Random patches)
- (iv)

Bias Variance Trade-Off:

The error comes in machine learning due to two reasons. (i) Bias

(ii) Variance

when ML model not provide better result in training data.

when ML model is perfectly work on the training data (fit data रखते हैं) but not good performance on test data.

Hence we always need low bias and low variance. ~~or~~ but

$$\text{Bias} \propto \frac{1}{\text{variance}}$$

Algorithms types

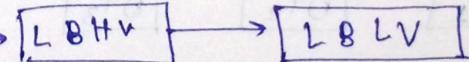
Low bias high variance

- Fully grown DT
- SVM
- KNN

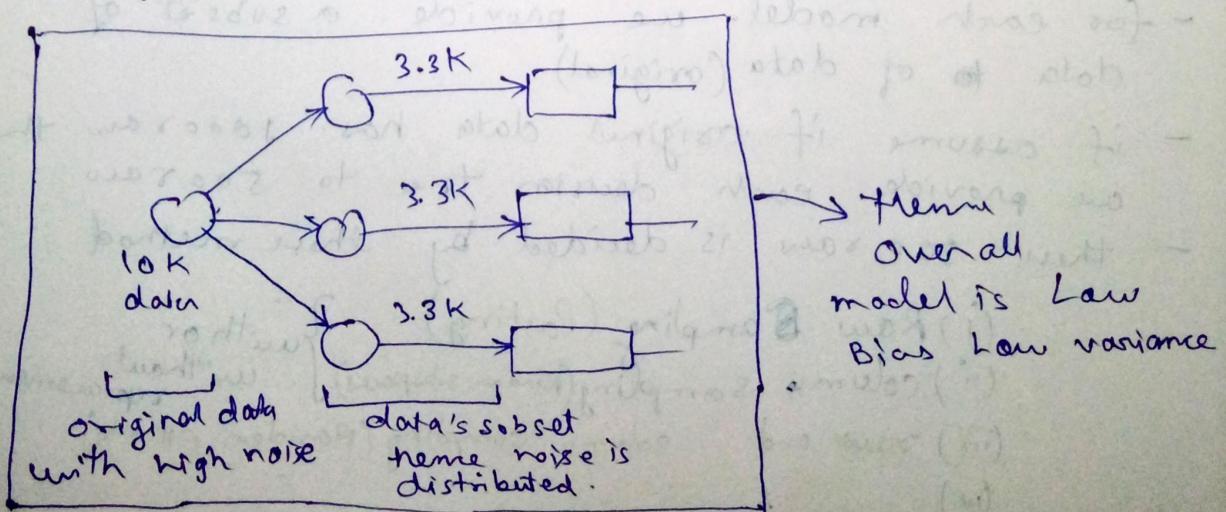
High bias low variance

- linear regression
- logistic regression

Random Forest



example How?

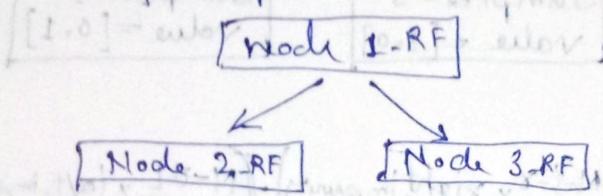


Difference between the Bagging and ~~Boos~~ Random Forests.

Fundamental difference between bagging and the random forest is that in Random forests, only a subset of features is selected randomly out of the entire and therefore the best split feature from the subset is employed to separate each node during a tree, unlike in bagging where all features are considered for splitting a node.

Random forests.)

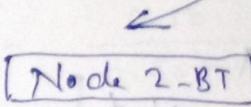
only m_{try} features
considered for each node
for a split



Bagging Trees

All of M features considered
for each node for a split

Node 1-BT



Node 2-BT

where m can be selected via out-of-bag error,
but $m = \sqrt{M}$ is good value to start with

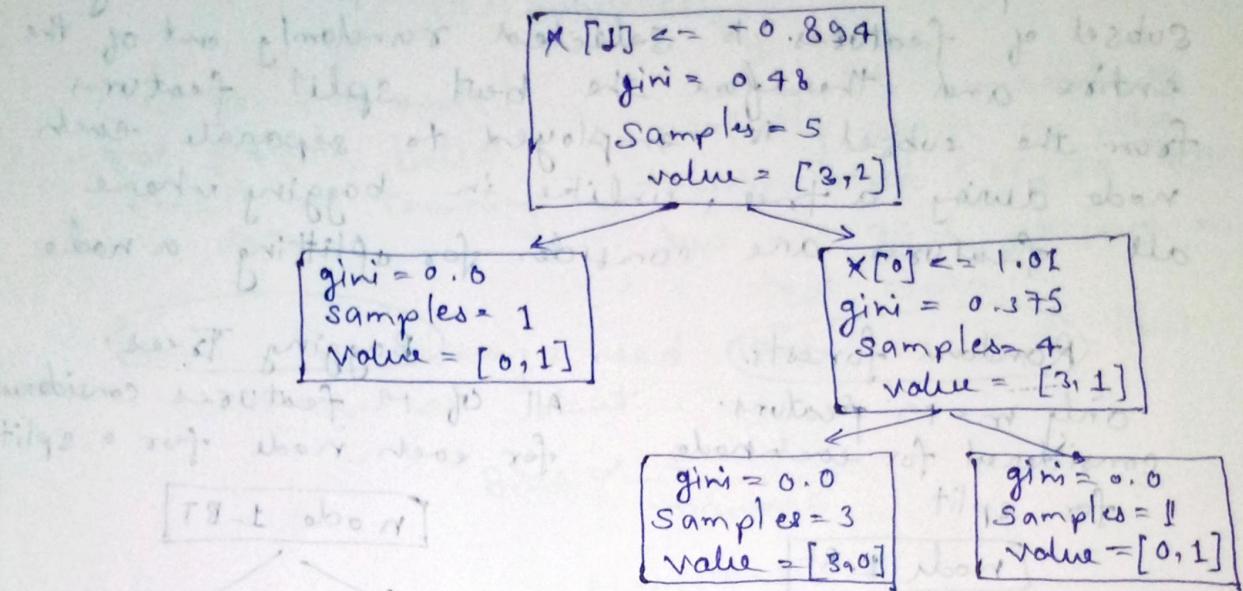
OOB Score (out-of-bag):

OOB score is a performance metric for a ml model, specifically for ensemble models such as random forest, it is calculated using the samples that are not used in the training of the model, which is called out-of-bag samples. These samples are used to provide an unbiased estimate of the model's performance, which is known as the OOB score.

Feature Importance Calculation.

(i) For Decision Tree

the addition of all features is always one



Formula

$$M_i = \frac{N-t}{N} \left[\text{impurity} - \left(\frac{N-t_L}{N-t} \times \text{right impurity} \right) - \left(\frac{N-t_R}{N-t} \times \text{left impurity} \right) \right]$$

$$f_{ik} = \frac{\sum n_i}{\text{je node split feature}} \quad \text{je all nodes}$$

In the above tree

the column 0 and 1

$X[0] <= 1.01$
gini = 0.375
samples = 4
value = [3, 1]

splitting based on zero

$X[1] <= -0.897$
gini = 0.48
samples = 5
value = [3, 2]

splitting based on 1

$$n_i(0^m) = \frac{X}{X+Y}$$

$$n_i(1^m) = \frac{Y}{X+Y}$$

$$h_i = \frac{N-t}{N} \left[\text{impurity} - \left(\frac{N-t}{N-t} \times \text{right-impurity} \right) - \left(\frac{N-t-L}{N-t} \times \text{left-impurity} \right) \right]$$

$$\gamma = \frac{5}{5} \left[0.48 - \left(\frac{4}{5} \times 0.375 \right) - \frac{1}{5} \times 0 \right] = 0.18$$

$$x = \frac{4}{5} \left[0.375 - (0) - 0 \right] = 0.30$$

feature¹
importance of 0th feature [0th] = $\frac{x}{x+y}$

$$= \frac{0.30}{0.30 + 0.18} = 0.625$$

feature²
importance of 1st feature [1st] = $\frac{y}{y+x}$

$$= \frac{0.18}{0.30 + 0.18} = 0.375$$

(ii) For Random Forest.

print(rf.estimators_[0].feature_importance_)

print(rf.estimators_[1].feature_importance_)

output. [1. 0.]

[0.5555 0.4444]

feature importance of 1st column

print((1+0.5555)/2) → 0.777500.

feature importance of 2nd column

print((0+0.4444)/2)

warning: impurity-based feature importances can be misleading if for high cardinality features (many unique values) See for alternative we use
sklearn.inspection.permutation_importance

$$\text{GDP} = \left[\text{GDP} - (\text{GDP})_{\text{mean}} \right] \frac{1}{2} \approx 0$$

$$X_{11} = [x_{11}]_{\text{mean}} \rightarrow \text{geography}$$

$$\text{GDP} = \frac{\text{GDP}}{\text{GDP}_{\text{mean}}}$$

$$\frac{x_{11}}{x_{11}} = [x_{11}]_{\text{mean}} \rightarrow \text{geography}$$

$$\text{GDP} = \frac{\text{GDP}}{\text{GDP}_{\text{mean}}}$$

(geography, gdp, oil) - (1 - correlation) being
(geography, oil) - (1 - correlation) being

↳ 0.27 significance

[PPT p. 222/20]

↳ 0.27 significance

↳ 0.27 significance

↳ 0.27 significance

↳ 0.27 significance

1. What is an OOB error (Out-Of-Bag error) and how is it useful?

- OOB error is a metric used in bagging algorithm like RF. It calculates the prediction error for each data point based on the votes of the decision trees in the ensemble that did not include that data point in their bootstrap sample during training.
- The OOB error serves as a reliable estimate of the model's performance without the need for a separate validation set, making it useful for assessing model accuracy and preventing overfitting.
- For example in RF with 100 DT, the OOB error is calculated as the average prediction error across all data points. While OOB error simplifies the modeling process and provides a robust estimate, it can be computationally intensive with a large number of trees and may have some variability due to its estimation nature.

2. In what scenario DT should be preferred over random forest?

(i) When to Prefer:

DT: When interpretability is crucial, and you need a single, understandable tree.

RF: When you seek higher predictive accuracy and robustness to outliers or noisy data.

(ii) Example use case:

DT (Medical Diagnosis): In a scenario where medical practitioners need clear, explainable rule for diagnosing a common illness based on a small dataset of patient symptoms.

RF(churn prediction): when predicting customer churn in a telecom company with a large dataset of diverse features, aiming for improved accuracy.

(iii) Description:

DT: it provides a clear, understandable decision path, crucial in scenarios where interpretability is more important than marginal gains in accuracy.

RF: Combines multiple decision trees, reducing overfitting and improving generalization performance, making it suitable for complex, high-dimensional data.

(iv) Advantages:

DT: Easy to visualize and explain; work well with small to medium-sized datasets.

RF: Reduces overfitting through ensemble learning; captures complex relationship in data

(v) Disadvantages:

DT: Prone to overfitting on large datasets or complex data.

RF: May not provide a transparent, interpretable model.

(vi) Numerical example:

DT (Medical Diagnosis): consider a small dataset of patient symptoms for diagnosing a common illness. A single decision tree can provide a clear set of rules that a medical practitioner can follow for diagnosis.

RF (customer churn Prediction):

In a large telecom dataset with hundreds of features, a random forest can combine multiple DT to predict customer churn accurately, considering various factors like call duration, contract length and customer demographics.