

Weekly Report 2 - Random Forest

Ganji Varshitha
AI20BTECH11009

Introduction

Random Forest is an ensemble classifier which combines multiple classifiers to achieve better accuracy. It trains several models using bootstrapped dataset and selects the majority vote for classification problems and average for regression problems.

Algorithm

Algorithm 1 Random Forest Algorithm

```
Given a training set S
for i = 1 to k do
    Build subset  $S_i$  by sampling with replacement from S
    Learn tree  $T_i$  from  $S_i$ 
    for each node do
        Choose best split from random subset of F features
        Each tree grows to the largest extent, and no pruning
    end for
end for
Make predictions according to majority vote of the set of k trees.
```

The value of F needs to be constant during the algorithm and it should be very less compared to total number of features M.

Possible values of F are $\frac{1}{2}\sqrt{M}$, \sqrt{M} , $2\sqrt{M}$.

Why does bagging work?

Decision trees are prone to overfit which results in high variance of the model. Bagging reduces the variance of the model.

Let S be the training dataset.

Let S_k be a sequence of training sets containing a sub-set of S.

Let P be the underlying distribution of S.

Bagging replaces the prediction of the model with the majority of the predictions given by the classifiers S.

$$\phi(x, P) = \mathbb{E}_s(\phi(x, S_k)) \quad (1)$$