

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.

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

```

1 class Random_Forest():
2     def __init__(self,n_trees, bootstrap_samples , n_features):
3         '''Constructor'''
4         self.n_trees = n_trees
5         self.bootstrap_samples = bootstrap_samples
6         self.n_features = n_features
7
8
9     # Function to get a subset of data with replacement
10    def Subset(self,data):
11        indices = np.random.choice(data.shape[0], size=self.
12        bootstrap_samples ,replace = True)
13        train_data = data[indices]
14        OOB_data = np.delete(data,indices,0)
15        return train_data,OOB_data
16
17    def get_label(self,data):
18        label_column = data[:, -1]
19        unique_classes, counts_unique_classes = np.unique(label_column,
20        return_counts=True)
21
22        index = counts_unique_classes.argmax()
23        classification = unique_classes[index]
24
25        return classification
26
27    def test_error(self,y_true, y_pred):
28        misclassified = 0
29        for i in range(len(y_true)):
30            if y_pred[i] != y_true[i]:
31                misclassified += 1
32        return misclassified / len(y_true)
33
34    def learn_trees(self,data):
35        # Creating an empty list for storing trees
36        trees_list = []
37        Decision_tree = Decision(self.n_features)
38        for i in range(self.n_trees):
39            train,OOB_data = self.Subset(data)
40            tree_learnt = Decision_tree.learn(train,{},0)
41            trees_list.append(tree_learnt)
42        return trees_list
43
44    def OOB_score(self,data):
45        OOB_error= []
46        Decision_tree = Decision(self.n_features)
47        for i in range(self.n_trees):
48            train,OOB_data = self.Subset(data)
49            tree_learnt = Decision_tree.learn(train,{},0)
50
51            Y_oob = Decision_tree.predict_test(OOB_data,tree_learnt)
52            OOB_error.append(self.test_error(OOB_data[:,-1],Y_oob))
53        OOB_score = np.mean(OOB_error)

```

```

52     return OOB_score
53
54
55     # Bagging - most important part
56     def predict(self, test, trees_list):
57         Decision_tree = Decision(self.n_features)
58         len_samples = len(test)
59         Preds_all = np.empty((len(trees_list), len_samples))
60         Preds = []
61         for i in range(len(trees_list)):
62             predict = Decision_tree.predict_test(test, trees_list[i])
63             Preds_all[i]=predict
64
65
66
67         for p in range(len_samples):
68             list1 = list(Preds_all[:,p])
69
70             Preds.append(max(set(list1), key=list1.count))
71
72
73     return Preds

```

Listing 1: Random Forest Code