

# **Lab 7: Decision Trees**

## **Bagging and Random Forests**

Here we apply bagging and random forests to the `Boston` data, using the `randomForest` package in `R`. The exact results obtained in this section may depend on the version of `R` and the version of the `randomForest` package installed on your computer. Recall that bagging is simply a special case of a random forest with  $m = p$ . Therefore, the `randomForest()` function can be used to perform both random forests and bagging. We perform bagging as follows:

`random  
Forest()`

```
> library(randomForest)
> set.seed(1)
> bag.boston=randomForest(medv~.,data=Boston,subset=train,
  mtry=13,importance=TRUE)
> bag.boston

Call:
randomForest(formula = medv ~ ., data = Boston, mtry = 13,
  importance = TRUE, subset = train)
      Type of random forest: regression
      Number of trees: 500
No. of variables tried at each split: 13

      Mean of squared residuals: 10.77
      % Var explained: 86.96
```

The argument `mtry=13` indicates that all 13 predictors should be considered for each split of the tree—in other words, that bagging should be done. How well does this bagged model perform on the test set?

```
> yhat.bag = predict(bag.boston,newdata=Boston[-train,])
> plot(yhat.bag, boston.test)
> abline(0,1)
> mean((yhat.bag-boston.test)^2)
[1] 13.16
```

The test set MSE associated with the bagged regression tree is 13.16, almost half that obtained using an optimally-pruned single tree. We could change the number of trees grown by `randomForest()` using the `ntree` argument:

```
> bag.boston=randomForest(medv~.,data=Boston,subset=train,
  mtry=13,ntree=25)
> yhat.bag = predict(bag.boston,newdata=Boston[-train,])
> mean((yhat.bag-boston.test)^2)
[1] 13.31
```

Growing a random forest proceeds in exactly the same way, except that we use a smaller value of the `mtry` argument. By default, `randomForest()` uses  $p/3$  variables when building a random forest of regression trees, and  $\sqrt{p}$  variables when building a random forest of classification trees. Here we use `mtry = 6`.

```
> set.seed(1)
> rf.boston=randomForest(medv~.,data=Boston,subset=train,
  mtry=6,importance=TRUE)
> yhat.rf = predict(rf.boston,newdata=Boston[-train,])
> mean((yhat.rf-boston.test)^2)
[1] 11.31
```

The test set MSE is 11.31; this indicates that random forests yielded an improvement over bagging in this case.

Using the `importance()` function, we can view the importance of each variable. `importance()`

```
> importance(rf.boston)
      %IncMSE  IncNodePurity
crim      12.384      1051.54
zn         2.103         50.31
indus      8.390      1017.64
chas        2.294         56.32
nox       12.791      1107.31
rm        30.754      5917.26
age       10.334         552.27
dis       14.641      1223.93
rad         3.583         84.30
tax         8.139         435.71
ptratio   11.274         817.33
black       8.097         367.00
lstat     30.962      7713.63
```

Two measures of variable importance are reported. The former is based upon the mean decrease of accuracy in predictions on the out of bag samples when a given variable is excluded from the model. The latter is a measure of the total decrease in node impurity that results from splits over that variable, averaged over all trees. In the case of regression trees, the node impurity is measured by the training RSS, and for classification trees by the deviance. Plots of these importance measures can be produced using the `varImpPlot()` function.

```
> varImpPlot(rf.boston)
```

`varImpPlot()`

The results indicate that across all of the trees considered in the random forest, the wealth level of the community (`lstat`) and the house size (`rm`) are by far the two most important variables.