

Random Forests and Boosting

Travers Parsons-Grayson

April 7, 2019

Random Forests vs. Bagging (Motivation)

The Difference

The process of creating a random forest is very similar to the process of bagging with one small caveat. In a random forest every time a split is considered a random sample of m predictors from the total p predictors are chosen as candidates for the split. Bagging is a special case of Random Forests when $m = p$.

Why?

Decorrelation: The weakness of bagging is that many trees end up looking the same because they will almost always use the strongest predictors in the same order. Random Forests *decorrelate* trees and thus reduce the variance of the prediction.

Bagging vs Random Forests

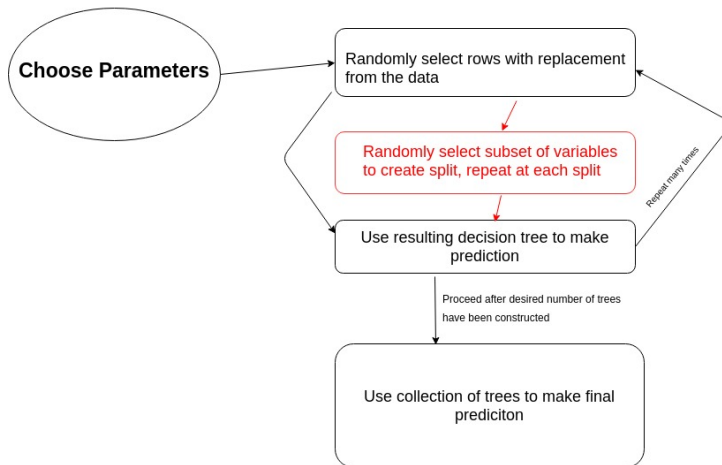


Figure 1:

Construction of a Random Forest

Parameters

- ① Number of trees, k
- ② Number of variables to select randomly at each split, m
- ③ (optional) Size of training set, the rows that we sample from without replacement
- ④ (optional) Maximum size of the trees grown, by number of nodes j

Process

- ① Randomly select rows with replacement from data (typically use 2/3's of rows)

Construction of a Random Forest

Parameters

- ① Number of trees, k
- ② Number of variables to select randomly at each split, m
- ③ (optional) Size of training set, the rows that we sample from without replacement
- ④ (optional) Maximum size of the trees grown, by number of nodes j

Process

- ① Randomly select rows with replacement from data (typically use 2/3's of rows)
- ② **Randomly select m variables to create split** (typically $m \equiv \sqrt{p}$)

Construction of a Random Forest

Parameters

- ① Number of trees, k
- ② Number of variables to select randomly at each split, m
- ③ (optional) Size of training set, the rows that we sample from without replacement
- ④ (optional) Maximum size of the trees grown, by number of nodes j

Process

- ① Randomly select rows with replacement from data (typically use 2/3's of rows)
- ② **Randomly select m variables to create split** (typically $m \equiv \sqrt{p}$)
- ③ Repeat step 2 at each split until decision tree is built

Construction of a Random Forest

Parameters

- ① Number of trees, k
- ② Number of variables to select randomly at each split, m
- ③ (optional) Size of training set, the rows that we sample from without replacement
- ④ (optional) Maximum size of the trees grown, by number of nodes j

Process

- ① Randomly select rows with replacement from data (typically use 2/3's of rows)
- ② **Randomly select m variables to create split** (typically $m \equiv \sqrt{p}$)
- ③ Repeat step 2 at each split until decision tree is built
- ④ Use resulting decision tree to make prediction

Construction of a Random Forest

Parameters

- ① Number of trees, k
- ② Number of variables to select randomly at each split, m
- ③ (optional) Size of training set, the rows that we sample from without replacement
- ④ (optional) Maximum size of the trees grown, by number of nodes j

Process

- ① Randomly select rows with replacement from data (typically use 2/3's of rows)
- ② **Randomly select m variables to create split** (typically $m \equiv \sqrt{p}$)
- ③ Repeat step 2 at each split until decision tree is built
- ④ Use resulting decision tree to make prediction
- ⑤ Repeat steps 1-3 k times

Example in R (Ames Housing)