

Bagging + Random Forests

ways to improve existing tree-based models

bagging works for any type of model

regression tree \rightarrow low bias + high variance

decrease variance = better model

bagging \rightarrow average good regression trees, variance \downarrow

RF \rightarrow average "weak" regression trees, bias \uparrow variance $\downarrow\downarrow$

Bagging

take B bootstrap samples of the data

\hookrightarrow fit a good regression tree for each bootstrap sample

\hookrightarrow no pruning

\hookrightarrow final model an average of these B separate trees

\hookrightarrow each tree a "base learner"

\hookrightarrow no impact on bias, reduce variance

$$\hat{f}(x) = \left(\frac{1}{B}\right) \sum_{b=1}^B \hat{f}^{x_b}(x)$$

\nwarrow new bootstrap sample will be "similar"
= similar splits on each tree

Random Forest

- almost identical to bagging

- make each tree a bit more random

- for every split on every tree:

\hookrightarrow take a sample of $m \leq p$ predictors

\hookrightarrow use the best split from the m sampled

- "weak learners"

- each tree will now have some bias

- less correlation between the trees \rightarrow variance $\downarrow\downarrow$

variable importance

- decrease in RSS

- decrease in accuracy { If you remove the information in a predictor, how much worse does your model get?