

Bagging and Random Forests

Nate Wells

Math 243: Stat Learning

November 15th, 2021

Outline

In today's class, we will . . .

- Implement random forests in R
- Investigate boosting as an **learning** method for improving decision trees

Section 1

Bagging and Random Forests in R

Random Forests

To create a random forest:

- 1 Select the number of models m to build and a number of predictors k to use at each step t
- 2 Generate a bootstrap sample for each model
- 3 Build a tree on the bootstrap sample where at each step, a random selection of k of the p predictors can be used (independent of prior predictors selected)
- 4 Aggregate the models to create an ensemble model.

Random Forests

To create a random forest:

- 1 Select the number of models m to build and a number of predictors k to use at each step t
- 2 Generate a bootstrap sample for each model
- 3 Build a tree on the bootstrap sample where at each step, a random selection of k of the p predictors can be used (independent of prior predictors selected)
- 4 Aggregate the models to create an ensemble model.

Random Forest in R

- To create both bagged trees and random forests, we use the `randomForest` function in the `randomForest` package in R:

```
library(randomForest)
rfmodel <- randomForest(Carbon_Sequestration_lb ~ ., data = my_pdxTrees_train)
rfmodel

##
## Call:
## randomForest(formula = Carbon_Sequestration_lb ~ ., data = my_pdxTrees_train)
##           Type of random forest: regression
##           Number of trees: 500
## No. of variables tried at each split: 2
##
##           Mean of squared residuals: 123.925
##           % Var explained: 84.38
```

Modifications

We can control how many trees are generated with `ntree` and the number of predictors at each split with `mtry`

Modifications

We can control how many trees are generated with `ntree` and the number of predictors at each split with `mtry`

- By default, `randomForest` uses $p/3$ predictors for regression and \sqrt{p} predictors for classification

Modifications

We can control how many trees are generated with `ntree` and the number of predictors at each split with `mtry`

- By default, `randomForest` uses $p/3$ predictors for regression and \sqrt{p} predictors for classification

```
set.seed(1)
rfmodel2 <- randomForest(Carbon_Sequestration_lb ~ ., data = my_pdxTrees_train,
                          ntree = 10, mtry = 5)
rfmodel2
```

```
##
## Call:
## randomForest(formula = Carbon_Sequestration_lb ~ ., data = my_pdxTrees_train, ntree = 10,
##              Type of random forest: regression
##              Number of trees: 10
## No. of variables tried at each split: 5
##
##              Mean of squared residuals: 142.6305
##              % Var explained: 82.02
```

Modifications

We can control how many trees are generated with `ntree` and the number of predictors at each split with `mtry`

- By default, `randomForest` uses $p/3$ predictors for regression and \sqrt{p} predictors for classification

```
set.seed(1)
rfmodel2 <- randomForest(Carbon_Sequestration_lb ~ ., data = my_pdxTrees_train,
                          ntree = 10, mtry = 5)
rfmodel2
```

```
##
## Call:
## randomForest(formula = Carbon_Sequestration_lb ~ ., data = my_pdxTrees_train, ntree = 10,
##              Type of random forest: regression
##              Number of trees: 10
## No. of variables tried at each split: 5
##
##              Mean of squared residuals: 142.6305
##              % Var explained: 82.02
```

How can we create a bagged model using the `randomForest` function?

Modifications

We can control how many trees are generated with `ntree` and the number of predictors at each split with `mtry`

- By default, `randomForest` uses $p/3$ predictors for regression and \sqrt{p} predictors for classification

```
set.seed(1)
rfmodel2 <- randomForest(Carbon_Sequestration_lb ~ ., data = my_pdxTrees_train,
                          ntree = 10, mtry = 5)
rfmodel2
```

```
##
## Call:
## randomForest(formula = Carbon_Sequestration_lb ~ ., data = my_pdxTrees_train, ntree = 10,
##              Type of random forest: regression
##              Number of trees: 10
## No. of variables tried at each split: 5
##
##              Mean of squared residuals: 142.6305
##              % Var explained: 82.02
```

How can we create a bagged model using the `randomForest` function?

- Set `mtry = p`, where p is the total number predictors available

Making predictions

- So you have your `randomForest` model. How do you make predictions?

```
my_preds <- predict(rfmodel, my_pdxTrees_test)
results <- data.frame(obs = my_pdxTrees_test$Carbon_Sequestration_lb, preds = my_preds)

results %>% head()
```

```
##      obs      preds
## 1  70.7  71.56668
## 2  38.6  45.41645
## 3  39.5  41.68216
## 4  60.9  53.92886
## 5  79.9  95.81464
## 6  77.0  79.18872
```

Variable Importance

Bagging and Random Forests increase prediction accuracy by reducing variance of the model.

Variable Importance

Bagging and Random Forests increase prediction accuracy by reducing variance of the model.

- But the cost comes in interpretability. We no longer have a single decision tree to follow to reach our prediction.

Variable Importance

Bagging and Random Forests increase prediction accuracy by reducing variance of the model.

- But the cost comes in interpretability. We no longer have a single decision tree to follow to reach our prediction.
- How can we determine which predictors are most influential?

Variable Importance

Bagging and Random Forests increase prediction accuracy by reducing variance of the model.

- But the cost comes in interpretability. We no longer have a single decision tree to follow to reach our prediction.
- How can we determine which predictors are most influential?

One possibility is to record the total amount of RSS/Purity that is decreased due to splits of the given predictor, averaged across all trees in the random forest.

Importance in R

```
importance(rfmodel)
```

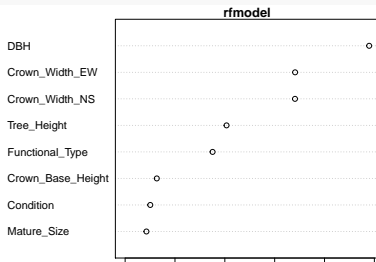
```
##              IncNodePurity
## DBH              489591.55
## Condition         50434.99
## Tree_Height      203426.25
## Crown_Width_NS   340775.81
## Crown_Width_EW   340999.72
## Crown_Base_Height 63524.57
## Functional_Type  175428.72
## Mature_Size      42685.72
```

Importance in R

```
importance(rfmodel)
```

```
##              IncNodePurity
## DBH              489591.55
## Condition        50434.99
## Tree_Height     203426.25
## Crown_Width_NS  340775.81
## Crown_Width_EW  340999.72
## Crown_Base_Height 63524.57
## Functional_Type 175428.72
## Mature_Size     42685.72
```

```
varImpPlot(rfmodel)
```

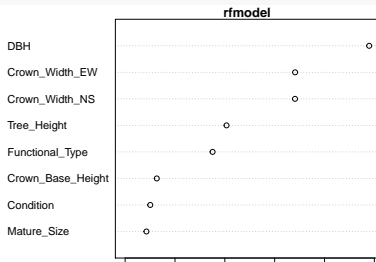


Importance in R

```
importance(rfmodel)
```

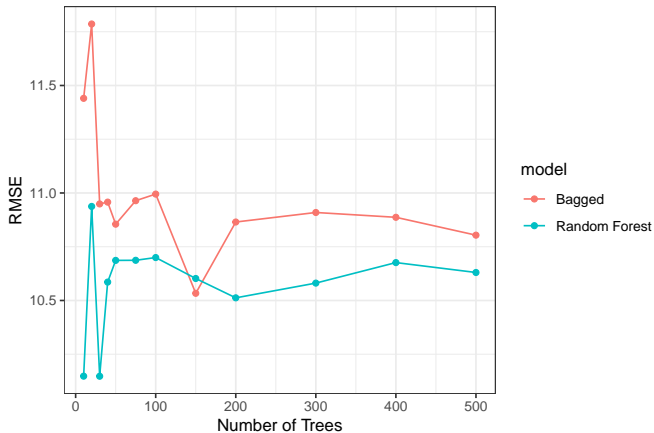
```
##              IncNodePurity
## DBH              489591.55
## Condition        50434.99
## Tree_Height     203426.25
## Crown_Width_NS  340775.81
## Crown_Width_EW  340999.72
## Crown_Base_Height 63524.57
## Functional_Type 175428.72
## Mature_Size     42685.72
```

```
varImpPlot(rfmodel)
```



- For regression trees, node impurity is calculated using RSS.
- For classification trees, node impurity is calculated using Gini Index.

Comparison of Bagged Trees versus Random Forests



Section 2

Boosting

Motivation

Suppose you have a model which, given a binary classification dataset, always returned a classifier with training error strictly lower than 50%.

Motivation

Suppose you have a model which, given a binary classification dataset, always returned a classifier with training error strictly lower than 50%.

- Can one use it to build a strong classifier that has error close to 0?

Motivation

Suppose you have a model which, given a binary classification dataset, always returned a classifier with training error strictly lower than 50%.

- Can one use it to build a strong classifier that has error close to 0?



AdaBoost

In the 1990s, Shapire and Freund developed algorithms to do just that.

AdaBoost

In the 1990s, Shapire and Freund developed algorithms to do just that.

- Their algorithm (AdaBoost) generates a sequence of weak classifiers, where at each iteration the algorithm finds the best classifier based on the current sample weights.

AdaBoost

In the 1990s, Shapire and Freund developed algorithms to do just that.

- Their algorithm (AdaBoost) generates a sequence of weak classifiers, where at each iteration the algorithm finds the best classifier based on the current sample weights.
 - Observations that are incorrectly classified in the k th iteration receive more weight in the $(k + 1)$ th iteration.

AdaBoost

In the 1990s, Shapire and Freund developed algorithms to do just that.

- Their algorithm (AdaBoost) generates a sequence of weak classifiers, where at each iteration the algorithm finds the best classifier based on the current sample weights.
 - Observations that are incorrectly classified in the k th iteration receive more weight in the $(k + 1)$ th iteration.
- The overall sequence of classifiers are combined into an ensemble which has high chance of classifying more accurately than any individual model in the list.

AdaBoost

In the 1990s, Shapire and Freund developed algorithms to do just that.

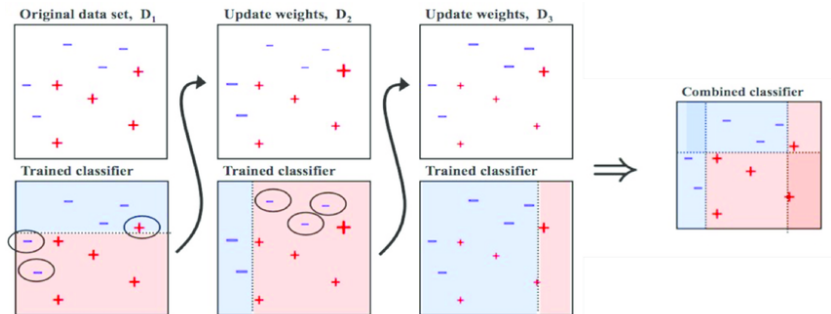
- Their algorithm (AdaBoost) generates a sequence of weak classifiers, where at each iteration the algorithm finds the best classifier based on the current sample weights.
 - Observations that are incorrectly classified in the k th iteration receive more weight in the $(k + 1)$ th iteration.
- The overall sequence of classifiers are combined into an ensemble which has high chance of classifying more accurately than any individual model in the list.
- The algorithm relies on using a sequence of **weak** learners (low variance, high bias)

AdaBoost

In the 1990s, Shapire and Freund developed algorithms to do just that.

- Their algorithm (AdaBoost) generates a sequence of weak classifiers, where at each iteration the algorithm finds the best classifier based on the current sample weights.
 - Observations that are incorrectly classified in the k th iteration receive more weight in the $(k + 1)$ th iteration.
- The overall sequence of classifiers are combined into an ensemble which has high chance of classifying more accurately than any individual model in the list.
- The algorithm relies on using a sequence of **weak** learners (low variance, high bias)
 - In the tree setting, we can create weak learners by restricting the depth of the tree.

AdaBoost Graphic



Boosting for regression

Boosting also works in the regression setting. The **gradient boosting machine** is a boosting algorithm that works as follows:

- 1 Select tree depth D and number of iterations K .
- 2 Compute the average response \hat{y} and use this as the initial predicted value for each observation
- 3 Compute the residual for each observation.
- 4 Fit a regression tree of depth D , using the **residuals** as the response.
- 5 Predict each observation using the regression tree from the previous step.
- 6 Update the predicted value of each observation by adding the previous iteration's predicted value to the predicted value generated in the previous step.
- 7 Repeat at total of K times.

Brief Example

Compute the mean:

```
mu <- mean(my_pdxTrees_train$Carbon_Sequestration_lb)
mu
```

```
## [1] 35.47188
```

Brief Example

Compute the mean:

```
mu <- mean(my_pdxTrees_train$Carbon_Sequestration_lb)
mu
```

```
## [1] 35.47188
```

Compute residuals:

```
my_pdxTrees_train_boost <- my_pdxTrees_train %>%
  mutate(residuals1 = Carbon_Sequestration_lb - mu)
```

Brief Example

Compute the mean:

```
mu <- mean(my_pdxTrees_train$Carbon_Sequestration_lb)
mu
```

```
## [1] 35.47188
```

Compute residuals:

```
my_pdxTrees_train_boost <- my_pdxTrees_train %>%
  mutate(residuals1 = Carbon_Sequestration_lb - mu)
```

Fit a new tree

```
boost_tree_model <- rpart(residuals1 ~ Crown_Base_Height,
  data = my_pdxTrees_train_boost,
  control = rpart.control(maxdepth = 2))
```

Brief Example

Compute the mean:

```
mu <- mean(my_pdxTrees_train$Carbon_Sequestration_lb)
mu
```

```
## [1] 35.47188
```

Compute residuals:

```
my_pdxTrees_train_boost <- my_pdxTrees_train %>%
  mutate(residuals1 = Carbon_Sequestration_lb - mu)
```

Fit a new tree

```
boost_tree_model <- rpart(residuals1 ~ Crown_Base_Height,
  data = my_pdxTrees_train_boost,
  control = rpart.control(maxdepth = 2))
```

Predict

```
predictions <- predict(boost_tree_model, data = my_pdxTrees_test) + mu
```

Brief Example

Compute the mean:

```
mu <- mean(my_pdxTrees_train$Carbon_Sequestration_lb)
mu
```

```
## [1] 35.47188
```

Compute residuals:

```
my_pdxTrees_train_boost <- my_pdxTrees_train %>%
  mutate(residuals1 = Carbon_Sequestration_lb - mu)
```

Fit a new tree

```
boost_tree_model <- rpart(residuals1 ~ Crown_Base_Height,
  data = my_pdxTrees_train_boost,
  control = rpart.control(maxdepth = 2))
```

Predict

```
predictions <- predict(boost_tree_model, data = my_pdxTrees_test) + mu
```

And so on...

Boosting Properties

Boosting is similar to random forests: the final prediction is sum of predictions from an ensemble of models.

Boosting Properties

Boosting is similar to random forests: the final prediction is sum of predictions from an ensemble of models.

- But in Random Forests, all trees are created independently, are of maximum depth, and contribute equally to the final model.

Boosting Properties

Boosting is similar to random forests: the final prediction is sum of predictions from an ensemble of models.

- But in Random Forests, all trees are created independently, are of maximum depth, and contribute equally to the final model.
- In boosting, subsequent trees are highly dependent on past trees, have minimal depth, and contribute unequally.

Boosting Properties

Boosting is similar to random forests: the final prediction is sum of predictions from an ensemble of models.

- But in Random Forests, all trees are created independently, are of maximum depth, and contribute equally to the final model.
- In boosting, subsequent trees are highly dependent on past trees, have minimal depth, and contribute unequally.

Unlike random forests, boosting is susceptible to over-fitting (since it uses a greedy algorithm to maximize gradient at each step).

Boosting Properties

Boosting is similar to random forests: the final prediction is sum of predictions from an ensemble of models.

- But in Random Forests, all trees are created independently, are of maximum depth, and contribute equally to the final model.
- In boosting, subsequent trees are highly dependent on past trees, have minimal depth, and contribute unequally.

Unlike random forests, boosting is susceptible to over-fitting (since it uses a greedy algorithm to maximize gradient at each step).

- To remedy, we introduce a shrinkage penalty (like in Ridge Regression/LASSO)

Boosting Properties

Boosting is similar to random forests: the final prediction is sum of predictions from an ensemble of models.

- But in Random Forests, all trees are created independently, are of maximum depth, and contribute equally to the final model.
- In boosting, subsequent trees are highly dependent on past trees, have minimal depth, and contribute unequally.

Unlike random forests, boosting is susceptible to over-fitting (since it uses a greedy algorithm to maximize gradient at each step).

- To remedy, we introduce a shrinkage penalty (like in Ridge Regression/LASSO)
 - Instead of adding the full value for a sample to the previous iteration's predicted value, only a fraction of the current predicted value is added.

Boosting Properties

Boosting is similar to random forests: the final prediction is sum of predictions from an ensemble of models.

- But in Random Forests, all trees are created independently, are of maximum depth, and contribute equally to the final model.
- In boosting, subsequent trees are highly dependent on past trees, have minimal depth, and contribute unequally.

Unlike random forests, boosting is susceptible to over-fitting (since it uses a greedy algorithm to maximize gradient at each step).

- To remedy, we introduce a shrinkage penalty (like in Ridge Regression/LASSO)
 - Instead of adding the full value for a sample to the previous iteration's predicted value, only a fraction of the current predicted value is added.
 - This fraction is called the *learning rate* λ , with $0 < \lambda < 1$. (Typical values range from 0.001 to 0.01)

Boosting in R

We use the `gbm` function in the `gbm` package to create Boosted Trees

Boosting in R

We use the `gbm` function in the `gbm` package to create Boosted Trees

- For regression problems, we use the argument `distribution = "gaussian"` and for classification problems, we use `distribution = "bernoulli"`

Boosting in R

We use the `gbm` function in the `gmb` package to create Boosted Trees

- For regression problems, we use the argument `distribution = "gaussian"` and for classification problems, we use `distribution = "bernoulli"`
- The argument `n.trees` controls the number of iterations
- The argument `interaction.depth` controls the depth of each tree
- The argument `shrinkage` controls the learning rate λ

Boosting in R

We use the `gbm` function in the `gmb` package to create Boosted Trees

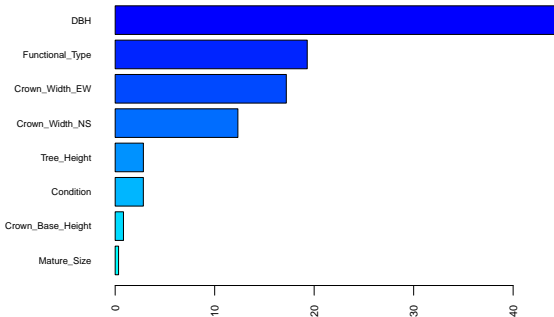
- For regression problems, we use the argument `distribution = "gaussian"` and for classification problems, we use `distribution = "bernoulli"`
- The argument `n.trees` controls the number of iterations
- The argument `interaction.depth` controls the depth of each tree
- The argument `shrinkage` controls the learning rate λ

```
library(gbm)
set.seed(10101)
boosted_tree<-gbm(Carbon_Sequestration_lb ~., my_pdxTrees_train,
  distribution = "gaussian",
  n.trees=1000,
  interaction.depth = 3,
  shrinkage = .02)
```


Summary Information

```
summary(boosted_tree )
```

```
##                var    rel.inf
## DBH                DBH 44.3189142
## Functional_Type    Functional_Type 19.3035576
## Crown_Width_EW      Crown_Width_EW 17.2013214
## Crown_Width_NS      Crown_Width_NS 12.3333653
## Tree_Height          Tree_Height  2.8359298
## Condition            Condition  2.8316083
## Crown_Base_Height    Crown_Base_Height 0.8345050
## Mature_Size          Mature_Size  0.3407984
```



Boosted Tree Performance

- How does the boosted tree do vs Random Forest? A pruned tree? A linear model?

Boosted Tree Performance

- How does the boosted tree do vs Random Forest? A pruned tree? A linear model?

```
results %>% group_by(model) %>% rmse(truth = obs, estimate = preds) %>% arrange(.estimate)
```

```
## # A tibble: 4 x 4
##   model      .metric .estimator .estimate
##   <chr>      <chr>    <chr>      <dbl>
## 1 boosted_tree rmse     standard    10.2
## 2 random_forest rmse     standard    10.6
## 3 pruned_tree  rmse     standard    13.3
## 4 linear_model  rmse     standard    17.0
```

Boosted Tree Performance

- How does the boosted tree do vs Random Forest? A pruned tree? A linear model?

```
results %>% group_by(model) %>% rmse(truth = obs, estimate = preds) %>% arrange(.estimate)
```

```
## # A tibble: 4 x 4
##   model      .metric .estimator .estimate
##   <chr>      <chr>   <chr>      <dbl>
## 1 boosted_tree rmse    standard    10.2
## 2 random_forest rmse    standard    10.6
## 3 pruned_tree  rmse    standard    13.3
## 4 linear_model  rmse    standard    17.0
```

- This behavior is typical. Boosted trees and Random Forests often have comparable performance, and both tend to be more accurate than other model types

Boosted Tree Performance

- How does the boosted tree do vs Random Forest? A pruned tree? A linear model?

```
results %>% group_by(model) %>% rmse(truth = obs, estimate = preds) %>% arrange(.estimate)
```

```
## # A tibble: 4 x 4
##   model      .metric .estimator .estimate
##   <chr>      <chr>   <chr>      <dbl>
## 1 boosted_tree rmse     standard    10.2
## 2 random_forest rmse     standard    10.6
## 3 pruned_tree  rmse     standard    13.3
## 4 linear_model  rmse     standard    17.0
```

- This behavior is typical. Boosted trees and Random Forests often have comparable performance, and both tend to be more accurate than other model types
- However, this performance comes at significant cost of interpretability.

Boosted Tree Performance

- How does the boosted tree do vs Random Forest? A pruned tree? A linear model?

```
results %>% group_by(model) %>% rmse(truth = obs, estimate = preds) %>% arrange(.estimate)
```

```
## # A tibble: 4 x 4
##   model      .metric .estimator .estimate
##   <chr>      <chr>   <chr>      <dbl>
## 1 boosted_tree rmse     standard    10.2
## 2 random_forest rmse     standard    10.6
## 3 pruned_tree  rmse     standard    13.3
## 4 linear_model  rmse     standard    17.0
```

- This behavior is typical. Boosted trees and Random Forests often have comparable performance, and both tend to be more accurate than other model types
- However, this performance comes at significant cost of interpretability.
- Note that boosted trees have a number of important hyperparameters: `n.trees`, `interaction.depth`, `shrinkage`.

Boosted Tree Performance

- How does the boosted tree do vs Random Forest? A pruned tree? A linear model?

```
results %>% group_by(model) %>% rmse(truth = obs, estimate = preds) %>% arrange(.estimate)
```

```
## # A tibble: 4 x 4
##   model      .metric .estimator .estimate
##   <chr>      <chr>   <chr>      <dbl>
## 1 boosted_tree rmse     standard    10.2
## 2 random_forest rmse     standard    10.6
## 3 pruned_tree  rmse     standard    13.3
## 4 linear_model  rmse     standard    17.0
```

- This behavior is typical. Boosted trees and Random Forests often have comparable performance, and both tend to be more accurate than other model types
- However, this performance comes at significant cost of interpretability.
- Note that boosted trees have a number of important hyperparameters: `n.trees`, `interaction.depth`, `shrinkage`.
 - How do we find the best values of these hyperparameters?

Boosted Tree Performance

- How does the boosted tree do vs Random Forest? A pruned tree? A linear model?

```
results %>% group_by(model) %>% rmse(truth = obs, estimate = preds) %>% arrange(.estimate)
```

```
## # A tibble: 4 x 4
##   model      .metric .estimator .estimate
##   <chr>      <chr>    <chr>      <dbl>
## 1 boosted_tree rmse     standard    10.2
## 2 random_forest rmse     standard    10.6
## 3 pruned_tree  rmse     standard    13.3
## 4 linear_model  rmse     standard    17.0
```

- This behavior is typical. Boosted trees and Random Forests often have comparable performance, and both tend to be more accurate than other model types
- However, this performance comes at significant cost of interpretability.
- Note that boosted trees have a number of important hyperparameters: `n.trees`, `interaction.depth`, `shrinkage`.
 - How do we find the best values of these hyperparameters?
 - Cross-validation!

Boosted Tree Performance

- How does the boosted tree do vs Random Forest? A pruned tree? A linear model?

```
results %>% group_by(model) %>% rmse(truth = obs, estimate = preds) %>% arrange(.estimate)
```

```
## # A tibble: 4 x 4
##   model      .metric .estimator .estimate
##   <chr>      <chr>    <chr>      <dbl>
## 1 boosted_tree rmse     standard    10.2
## 2 random_forest rmse     standard    10.6
## 3 pruned_tree  rmse     standard    13.3
## 4 linear_model  rmse     standard    17.0
```

- This behavior is typical. Boosted trees and Random Forests often have comparable performance, and both tend to be more accurate than other model types
- However, this performance comes at significant cost of interpretability.
- Note that boosted trees have a number of important hyperparameters: `n.trees`, `interaction.depth`, `shrinkage`.
 - How do we find the best values of these hyperparameters?
 - Cross-validation!
 - But tuning all three parameters by “hand” with `rsample` is tedious. We need a more powerful cv engine