

# Super Learner for Prediction (Part 1)

David Whitney (based on materials by Karla Diaz-Ordaz and Jamie Burns)

8 Feb 2022

## Introduction

---

In this practical, you will apply the SuperLearner. We will use part of the National Health and Nutrition Examination Survey (NHANES) dataset. We want to write the raw dataset to a data frame in our environment, and can do this with the following code:

```
nhanes_installed <- require(NHANES)
if(!nhanes_installed) install.packages("NHANES", repos = "http://cran.us.r-project.org")
library(NHANES)

data("NHANES")
df <- NHANESraw
```

## Data Processing

---

Inspect the data

```
dim(df)
```

```
## [1] 20293    78
```

```
names(df)
```

```
## [1] "ID"           "SurveyYr"      "Gender"
## [4] "Age"          "AgeMonths"     "Race1"
## [7] "Race3"        "Education"     "MaritalStatus"
## [10] "HHIncome"     "HHIncomeMid"   "Poverty"
## [13] "HomeRooms"    "HomeOwn"       "Work"
## [16] "Weight"       "Length"        "HeadCirc"
## [19] "Height"       "BMI"           "BMICatUnder20yrs"
## [22] "BMI_WHO"      "Pulse"         "BPSysAve"
## [25] "BPDiaAve"     "BPSys1"        "BPDia1"
## [28] "BPSys2"       "BPDia2"        "BPSys3"
## [31] "BPDia3"       "Testosterone"  "DirectChol"
## [34] "TotChol"      "UrineVol1"     "UrineFlow1"
## [37] "UrineVol2"    "UrineFlow2"    "Diabetes"
## [40] "DiabetesAge"  "HealthGen"     "DaysPhysHlthBad"
## [43] "DaysMentHlthBad" "LittleInterest" "Depressed"
## [46] "nPregnancies" "nBabies"       "Age1stBaby"
## [49] "SleepHrsNight" "SleepTrouble"  "PhysActive"
```

```
## [52] "PhysActiveDays"    "TVHrsDay"         "CompHrsDay"
## [55] "TVHrsDayChild"    "CompHrsDayChild"  "Alcohol12PlusYr"
## [58] "AlcoholDay"       "AlcoholYear"      "SmokeNow"
## [61] "Smoke100"         "SmokeAge"         "Marijuana"
## [64] "AgeFirstMarij"    "RegularMarij"     "AgeRegMarij"
## [67] "HardDrugs"        "SexEver"          "SexAge"
## [70] "SexNumPartnLife"  "SexNumPartYear"   "SameSex"
## [73] "SexOrientation"   "WTINT2YR"         "WTMEC2YR"
## [76] "SDMVPSU"          "SDMVSTRA"         "PregnantNow"
```

```
table(df$SmokeNow, df$Smoke100, useNA='always')
```

```
##
##           No  Yes <NA>
## No         0 2779    0
## Yes        0 2454    0
## <NA> 6536     2 8522
```

We will be interested in smoking as an exposure. The dataset contains two smoking variables `Smoke100` which is a binary indicator of whether a person has smoked at least 100 cigarettes in their lifetime (but is not a current smoker), and `SmokeNow` which indicates if the person is a current smoker.

Combine these into a single factor variable that indicates whether the individual is an ex, current or never smoker.

```
df <- df %>% mutate(Smoke = ifelse(SmokeNow == "Yes", "Current",
                                   ifelse(Smoke100 == "Yes", "Ex", "Never"))) %>%
  mutate(Smoke = ifelse(is.na(Smoke), "Never", Smoke)) %>%
  mutate(Smoke = factor(Smoke))
```

A systolic blood pressure reading (`BPSysAve`), a continuous outcome, will be a primary outcome. If there is time, you can also consider diabetes status (`Diabetes`), a binary outcome.

Now we trim the variables in the data set, keeping only the ones relevant to the practical:

```
df <- df %>% dplyr::select(one_of(
  "BPSysAve", "BMI", "Age", "SleepHrsNight", "PhysActive", "Smoke",
  "Gender", "Race1", "Poverty", "Diabetes", "TotChol"))
```

Inspect our outcome variables, and do the same for key covariates BMI and Age.

```
table(df$Diabetes, useNA = 'always')
```

```
##
##    No   Yes  <NA>
## 17754 1706   833
```

```
summary(df$BPSysAve)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.     NA's
##      74.0   105.0   115.0   118.1   127.0   233.0     5426
```

```
table(df$Smoke, useNA = 'always')
```

```
##
## Current      Ex   Never   <NA>
##    2454     2779   15060      0
```

```
summary(df$BMI)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.     NA's
```

```
##    12.40    19.79    24.92    25.65    30.10    84.87    2279
```

```
summary(df$Age)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.00   10.00   28.00   32.02   53.00   80.00
```

For this illustration, we want to keep only the complete cases in the dataset. We do this using a single line of code.

```
df <- df[complete.cases(df),]
```

Now split the data into training and test sets. The split should be 50/50. Call the sets `df_train` and `df_holdout`.

```
set.seed(101)
train_obs <- sample(nrow(df), size = nrow(df)*0.5)

df_train <- df[train_obs, ]

# Create a holdout set for evaluating model performance.
# Note: cross-validation is even better than a single holdout sample.
df_holdout <- df[-train_obs, ]
```

We want the outcome (BPSysAve) in vector form (also split into train and holdout) and kept separate from the main data frame. Do this now.

```
# The continuous outcome will be the average of the 3 systolic blood pressure measurement
# BPSysAve
y_train <- df$BPSysAve[train_obs]
y_holdout <- df$BPSysAve[-train_obs]
```

## Fitting Individual Algorithms

---

### Linear Regression

---

We begin by trying different learners individually. For example: a linear regression, lasso (glmnet), random-Forest, XGBoost. These should ideally be tested with multiple hyperparameter settings for each algorithm. (You will get to do this later!)

Carry out a simple linear regression for the systolic blood pressure measurement, using BMI, Age, SleepHrsNight, PhysActive, Smoke, Gender, Race1, Poverty, Diabetes and TotChol as covariates.

Using the generic `predict()` function to make predictions, manually calculate the sum of square errors for the test set.

```
form <- "BPSysAve ~ BMI + Age + SleepHrsNight + PhysActive +
        Smoke + Gender + Race1 + Poverty + Diabetes +TotChol"
mod.reg <- glm(form, data=df_train, family=gaussian)
Yhat.reg <- predict(mod.reg, newdata=df_holdout, type='response')

# sum squares error
SSE.reg <- sum((y_holdout - Yhat.reg)^2)
```

## Boosting

---

Now we run a gradient boosting algorithm on the data. We specify the formula with only the main terms, but recall that trees automatically include interactions up to the specified depth.

Use the `gbm()` function to fit a gradient boosting model to the data.

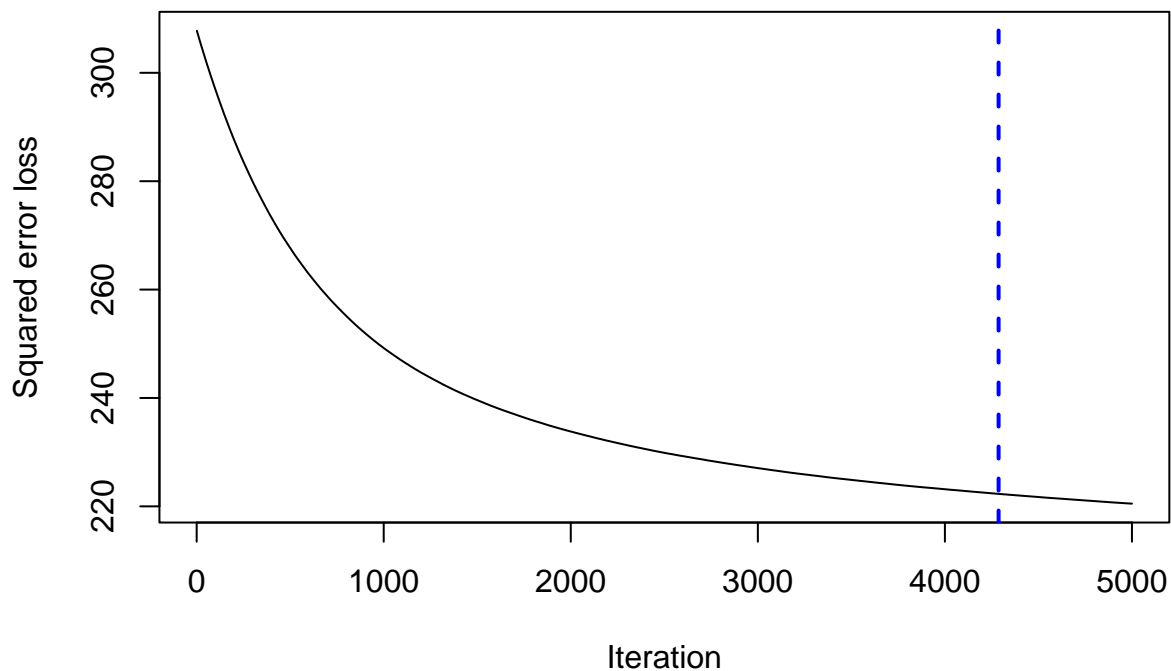
```
#' Set the seed for reproducibility.
set.seed(1)

boost <- gbm( formula = as.formula(form),
              data = df_train,
              distribution = "gaussian",
              shrinkage = 0.001,
              n.trees = 5000,
              cv.folds = 5,
              interaction.depth=3 )
```

We can use the `gbm.perf()` function to check the performance of this model, and by specifying the option `method`, we can do this in multiple ways. Check the performance using out-of-bag validation (`method = "OOB"`) and cross-validation (`method = "CV"`).

How do you interpret the plot that results from running `gbm.perf()`?

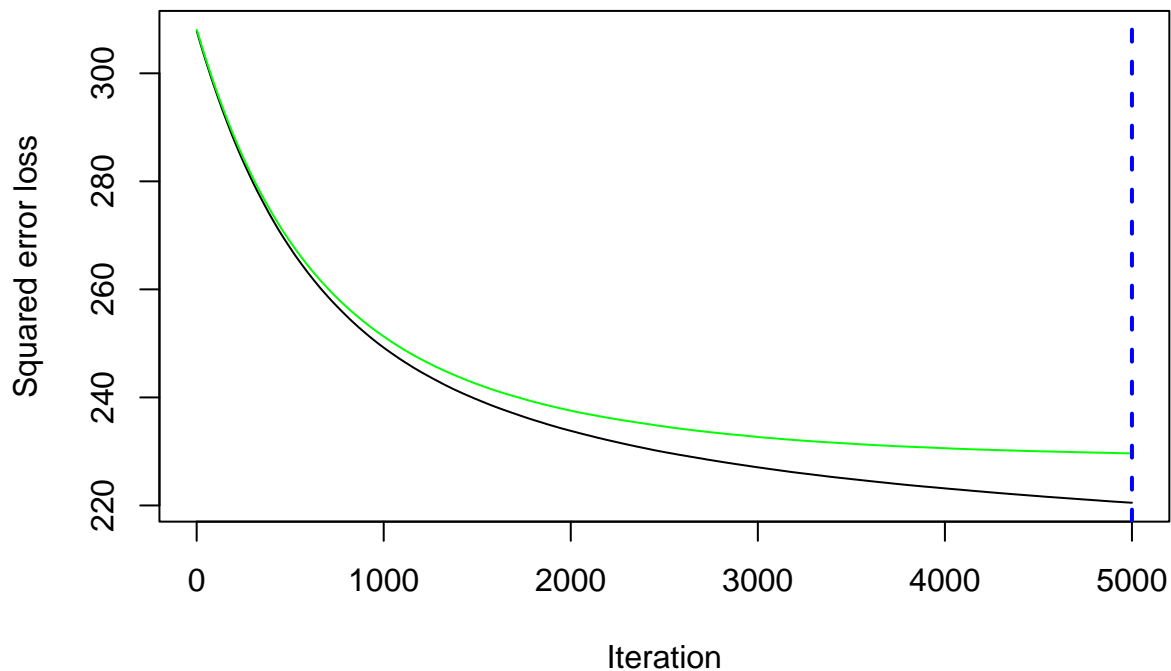
```
#' Checks performance using the out-of-bag (OOB) error
best.iter <- gbm.perf(boost, method = "OOB")
```



```
print(best.iter)

## [1] 4287
## attr(,"smoother")
## Call:
## loess(formula = object$oobag.improve ~ x, enp.target = min(max(4,
##     length(x)/10), 50))
##
## Number of Observations: 5000
## Equivalent Number of Parameters: 39.99
## Residual Standard Error: 0.001578

# Checks performance using 5-fold cross-validation
best.iter2 <- gbm.perf(boost, method = "cv")
```



```
print(best.iter2)

## [1] 5000

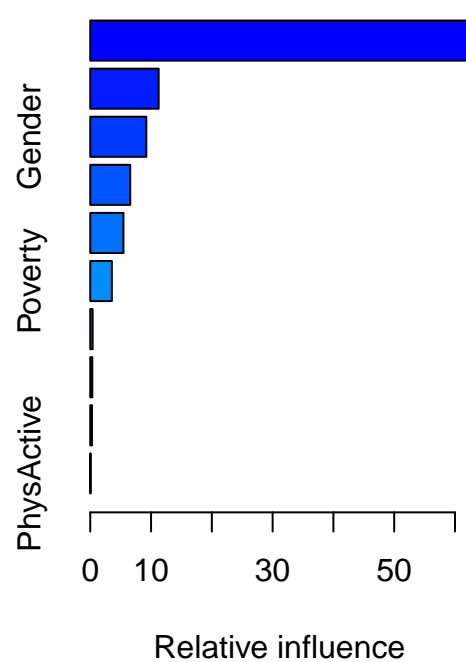
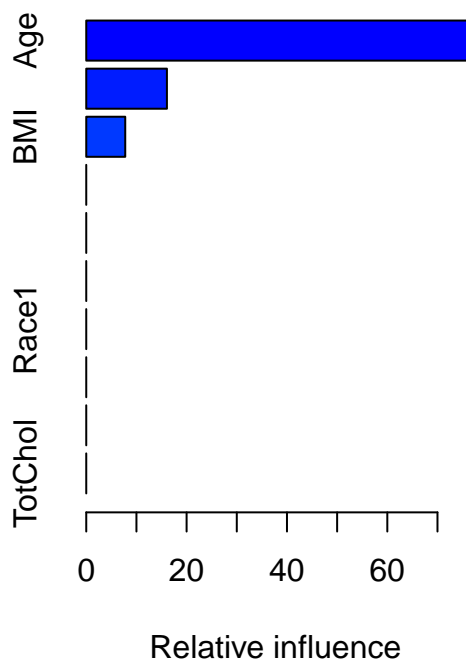
The OOB error typically underestimates the optimal number of iterations.

Using the summary() function with the n.trees options, generate plots that show the relative influence of
each variable after (a) a single tree and (b) with the optimal number of iterations as determined in the last
question.

# Plots relative influence of each variable
par(mfrow = c(1, 2))
summary(boost, n.trees = 1)           # using first tree
```

```
##           var    rel.inf
## Age           Age 76.152577
## Gender        Gender 16.076925
## BMI           BMI  7.770499
## SleepHrsNight SleepHrsNight 0.000000
## PhysActive    PhysActive 0.000000
## Smoke         Smoke 0.000000
## Race1         Race1 0.000000
## Poverty       Poverty 0.000000
## Diabetes      Diabetes 0.000000
## TotChol       TotChol 0.000000
```

```
summary(boost, n.trees = best.iter) # using estimated best number of trees
```



```
##           var    rel.inf
## Age           Age 62.85775470
## BMI           BMI 11.25197645
## Gender        Gender 9.22223588
## TotChol       TotChol 6.57923313
## Race1         Race1 5.46377909
## Poverty       Poverty 3.55175178
## Diabetes      Diabetes 0.38909641
## SleepHrsNight SleepHrsNight 0.33934926
## Smoke         Smoke 0.28974810
## PhysActive    PhysActive 0.05507517
```

We see that in both (first and best tree) Age has the largest influence. Gender and BMI have second and third

largest influence, using the first tree, while the best tree has BMI as higher than Gender.

Now using the generic `predict()` function, make predictions for the holdout set. Don't forget to stipulate how many iterations of the gradient boosting model should be used. This can be done using the `n.trees` option.

Also generate the sum of square errors that results from applying this model to the holdout set.

```
#' predictions will be on the link scale
Yhat.boost <- predict(boost, newdata = df_holdout, n.trees = best.iter, type = "link")

#' Sum of Squared Errors SSE
SSE.boost <- sum((y_holdout - Yhat.boost)^2)
```

## LASSO via SuperLearner

We will continue to fit single learners, but in order to get familiar it we do so using the `SuperLearner` package

First, check which learners have been integrated into the `SuperLearner` package. We can use any of these when we run the `SuperLearner`:

```
listWrappers(what = "SL")

## [1] "SL.bartMachine"      "SL.bayesglm"        "SL.biglasso"
## [4] "SL.caret"           "SL.caret.rpart"     "SL.cforest"
## [7] "SL.earth"           "SL.extraTrees"      "SL.gam"
## [10] "SL.gbm"             "SL.glm"             "SL.glm.interaction"
## [13] "SL.glmnet"          "SL.ipredbag"        "SL.kernelKnn"
## [16] "SL.knn"             "SL.ksvm"            "SL.lda"
## [19] "SL.leekasso"        "SL.lm"              "SL.loess"
## [22] "SL.logreg"          "SL.mean"            "SL.nnet"
## [25] "SL.nnls"            "SL.polymars"        "SL.qda"
## [28] "SL.randomForest"    "SL.ranger"          "SL.ridge"
## [31] "SL.rpart"           "SL.rpartPrune"      "SL.speedglm"
## [34] "SL.speedlm"         "SL.step"            "SL.step.forward"
## [37] "SL.step.interaction" "SL.stepAIC"         "SL.svm"
## [40] "SL.template"        "SL.xgboost"
```

`SuperLearner` (SL) likes the outcome and matrix of covariates to be kept separate. Do this now.

```
X <- df %>% dplyr::select(-one_of("BPSysAve"))

# Also divide our design matrix into training and testing sets
x_train <- X[train_obs, ]
x_holdout <- X[-train_obs, ]
```

Now let's fit penalised regression LASSO, but using SL. For now using all the defaults which we will explain later:

```
# Fit lasso model
sl_lasso <- SuperLearner(Y = y_train,
                        X = x_train,
                        family = gaussian(),
                        SL.library = "SL.glmnet",
                        cvControl = list(V=5L))
```

```
# Review the elements in the SuperLearner object.
names(sl_lasso)
```

```
## [1] "call"           "libraryNames"   "SL.library"
## [4] "SL.predict"     "coef"           "library.predict"
## [7] "Z"             "cvRisk"         "family"
## [10] "fitLibrary"     "cvFitLibrary"   "varNames"
## [13] "validRows"      "method"         "whichScreen"
## [16] "control"        "cvControl"      "errorsInCVLibrary"
## [19] "errorsInLibrary" "metaOptimizer"  "env"
## [22] "times"
```

Again using the generic `predict()` function, calculate the sum of square errors on the test set.

```
# Predict outcome in the holdout using lasso
Yhat.lasso <- predict(sl_lasso, x_holdout, onlySL = TRUE)$pred

# calculate sum of square errors
SSE.lasso <- sum((y_holdout - Yhat.lasso)^2)
```

The SSE corresponding to LASSO is larger than the one obtained from boosting.

## Random Forest via SuperLearner

Fit a random forest model using the wrapper `SL.ranger`. All other options as before. Also calculate the sum of square errors on the test set.

```
# Fit random forest using the wrapper function SL.ranger, with all the defaults
# You can find these by typing ?SL.ranger

sl_rf <- SuperLearner(Y = y_train,
                     X = x_train,
                     family = gaussian(),
                     SL.library = "SL.ranger",
                     cvControl = list(V=5L))

# predict Y in the holdout
Yhat.rf <- predict(sl_rf, x_holdout, onlySL = TRUE)$pred

# root least squares error
SSE.ranger <- sum((y_holdout - Yhat.rf)^2)
```

The `SL.ranger` wrapper fits a random forest with 500 trees, minimum node size 5, and the number of variables used to split the squared root of the number of available predictors.

With these defaults, the SSE is a bit larger than the one we obtained with boosting.

## Multiple Learners Stacked in SuperLearner

Instead of fitting the models separately and looking at the performance using sum of square errors as we have been doing, we now fit them simultaneously by including them all in the SuperLearner library.

For now, we include those algorithms we tried up to now:



```

# Select candidate algorithms
my.library <-c("SL.glm","SL.glmnet", "SL.ranger")

# Set seed
set.seed(101)

# Execute the call to SuperLearner
sl0 <- SuperLearner(Y = y_train, # Y is the outcome variable
                   X = x_train, # X is a dataframe of predictor variables, in this case
                                # everything except for outcome
                   family = gaussian(), # family will be discussed in more detail when
                                         # we see how wrappers are written.
                                         # for now gaussian (outcome is continuous)
                                         # binomial() for 0/1 outcome
                   method = "method.NNLS",
                   # method specifies how the ensembling is done (i.e. how the optimal
                   # combination is chose)
                   # for now we will use the  $\sum_{k=1}^K \alpha_k f_{k,n}$  method by default
                   SL.library = my.library,
                   cvControl = list(shuffle = F, V = 5)
                   # cvControl specifies parameters related to cross validation,
                   # used to estimate the risk on future data
                   # the default is for V = 10-fold
)

```

Now check the output by simply calling the SuperLearner object.

```

sl0

##
## Call:
## SuperLearner(Y = y_train, X = x_train, family = gaussian(), SL.library = my.library,
##   method = "method.NNLS", cvControl = list(shuffle = F, V = 5))
##
##
##           Risk      Coef
## SL.glm_All    236.4878 0.4204597
## SL.glmnet_All 236.4112 0.0000000
## SL.ranger_All 233.7565 0.5795403

```

The output has two main components:

Firstly, the risk is a measure of model accuracy or performance, and we want our models to minimize the estimated risk (according to the specified loss function). Because we did not change it by an option, our SL has used the default mean square error, but this can be altered (see the help files for **SuperLearner**). In this case, the risks for each algorithm in the library are all broadly similar, with the random forest slightly outperforming the other two algorithms in the library.

The **coef** column tells us the importance of each algorithm in the final ensemble. By default (because we use NNLS) the weights are always greater than or equal to 0 and sum to 1 (a ‘convex combination’). If a coefficient is 0, it means that the algorithm is not being used in the SuperLearner ensemble. Here we see that **glm** has been given no weight in the final (ensemble) predictor and so is not used.

Now see which has the lowest risk, and how long it took to run, using the following code:

```

# Let's see which is the discrete SL (i.e min risk)
sl0$cvRisk[which.min(sl0$cvRisk)]

```

```
## SL.ranger_All
##      233.7565
# Review how long it took to run the SuperLearner:
s10$times$everything
```

```
##      user  system elapsed
## 11.284   0.481  10.929
```

Now that we have a SL ensemble predictor, make predictions on the holdout data set and review the results. This can be done using the generic `predict()` function, but we stipulate the option `onlySL = TRUE` so we do not fit algorithms in the library that had zero weight, saving computation.

```
pred.s10 <- predict(s10, x_holdout, onlySL = TRUE)
```

Check the structure of this prediction object using `str(pred.s10)`. You will see that the prediction object is a list with two objects. The first is a vector of SL predictions. The second is a matrix of predictions; these are the predictions for each of the individual learners in the SL library.

We want the first object in the list to serve as our predictions. Use this to calculate the sum of square errors resulting from applying the SL predictor to the test set.

```
# Pick out SL predictions
Yhat.s10 <- pred.s10[[1]]

# Calculate the sum of square errors
SSE.s10 <- sum((y_holdout - Yhat.s10)^2)
```

The SSE is the lower than those SSEs corresponding to LASSO and RF (which were included in the SL library).

## Extend the Ensemble

We will now add other base learners to the library. You can check which ones are implemented by running `listWrappers()`. (You can also write your own wrapper functions, we will see later how to do this).

For now, add `gbm()` and `glm.interaction`. [Note: `xgboost` has deprecated a function used in the SL wrapper so warnings “`reg:linear` is now deprecated in favor of `reg:squarederror`” will be displayed, so I won’t use for this exercise, though ordinarily I do!]

Run another SL predictor using the same three algorithms as before alongside the two new ones. Review the output and the time it took to run as before.

```
# Define a new, larger library
my.library.2 <- c("SL.glm", "SL.glm.interaction", "SL.glmnet", "SL.ranger", "SL.gbm")

# Set seed
set.seed(101)

# Run SL with new library
s11 <- SuperLearner( Y = y_train,
                    X = x_train,
                    family = gaussian(),
                    SL.library = my.library.2,
                    cvControl = list(shuffle = F, V = 5) )
```

```

# Review the SL object
sl1

##
## Call:
## SuperLearner(Y = y_train, X = x_train, family = gaussian(), SL.library = my.library.2,
##   cvControl = list(shuffle = F, V = 5))
##
##
##              Risk      Coef
## SL.glm_All      236.4878 0.0000000
## SL.glm.interaction_All 234.0326 0.2662155
## SL.glmnet_All    236.4585 0.0000000
## SL.ranger_All    233.6424 0.1555134
## SL.gbm_All       229.3069 0.5782711

# Review times
sl1$times$everything

##      user  system elapsed
## 81.447   2.045 187.981

```

You should see that the gradient boosted predictor (`gbm()`) has the most weight. The new `glm` with interactions is the second most important, and the others are either zero or make only small contributions to the ensemble.

As before, use the SL predictor to make predictions on the holdout set.

```

Yhat.sl1 <- predict(sl1, x_holdout, onlySL = TRUE)[[1]]

# SSE
SSE.sl1 <- sum((y_holdout - Yhat.sl1)^2)

```

We see that adding more learners slightly improves the SSE. Adding more (without removing) will result in further improvements.