

hw2

Johnstone Tcheou

2025-03-04

Contents

Fitting smoothing spline	5
Fitting multivariate adaptive regression spline	9
Fitting generalized additive model (GAM)	12

```
library(caret)
library(tidymodels)
library(splines)
library(mgcv)
library(pdp)
library(earth)
library(ggplot2)
```

First, we need to set a seed for reproducibility. Then we will import the data and conduct a training/testing split, using 80% for training and 20% for testing.

```
set.seed(81062)
college <- read.csv("College.csv")

# college <- college /> na.omit()

data_split <- initial_split(college, prop = 0.8)
training <- training(data_split)
testing <- testing(data_split)

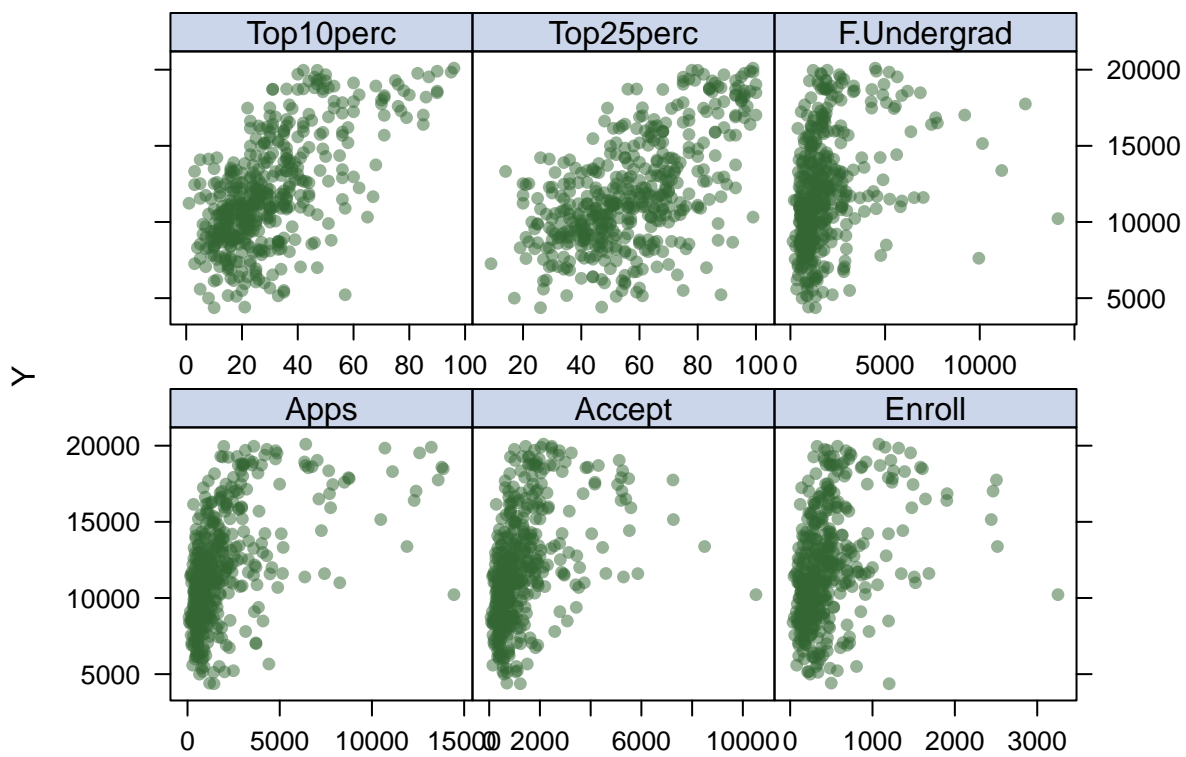
training_y <- training$Outstate
training_x <- model.matrix(Outstate ~ ., training)[,-1]

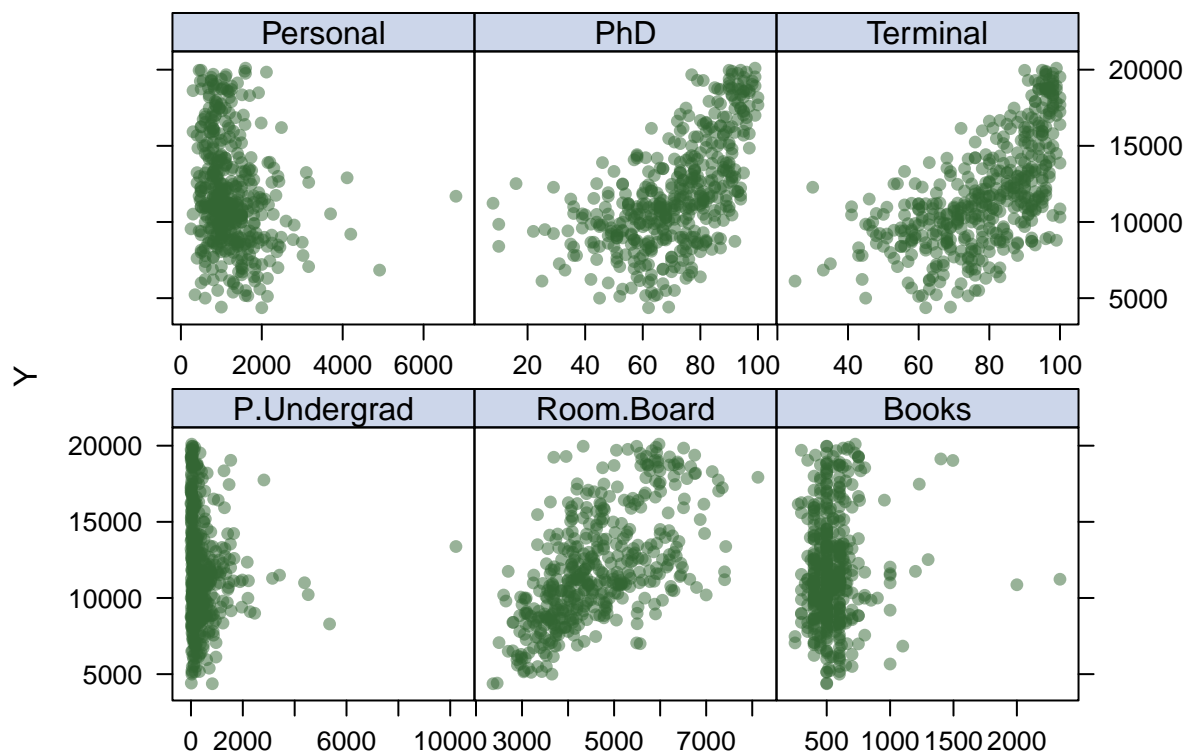
testing_y <- testing$Outstate
testing_x <- model.matrix(Outstate ~ ., testing)[,-1]
```

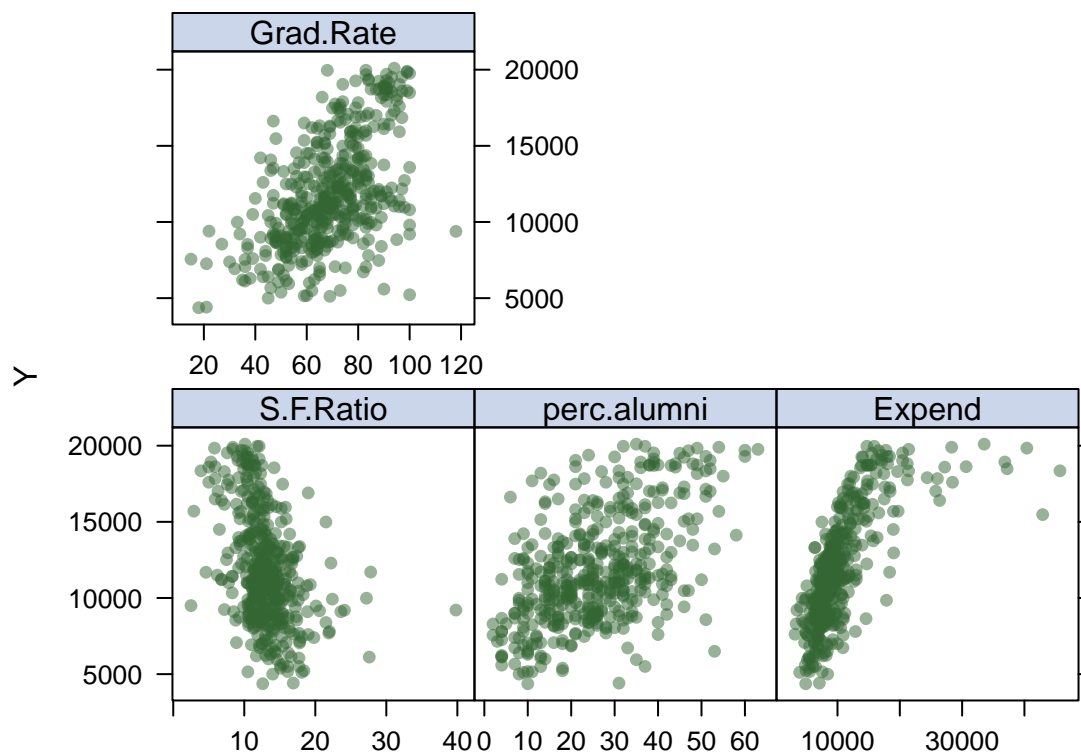
Except for the dummy variables created for each college, the other predictors are continuous. We can examine their relationships with our outcome variable, `Outstate` - out of state tuition.

```
theme1 <- trellis.par.get()
theme1$plot.symbol$col <- rgb(.2, .4, .2, .5)
theme1$plot.symbol$pch <- 16
theme1$plot.line$col <- rgb(.8, .1, .1, 1)
theme1$plot.line$lwd <- 2
theme1$strip.background$col <- rgb(.0, .2, .6, .2)
trellis.par.set(theme1)

featurePlot(training_x[, -c(1:451)], training_y, plot = "scatter", labels = c("", "Y"),
            type = c("p"), layout = c(3, 2))
```







Potential variables with nonlinear relationships with Outstate include Apps, Accept, Enroll, F.Undergrad, Expend, S.F.Ratio, PhD, and Terminal (both of which appear to have some exponential aspect).

Fitting smoothing spline

The range of degrees of freedom is from 1 to the number of unique predictor values, which is 57. Therefore, a range of degrees of freedom is tested, ranging from 10 to 50 by an increment of 10.

```
ss.10 <-
  smooth.spline(
    training_x[, "perc.alumni"],
    training_y,
    df = 10
  )

ss.20 <-
  smooth.spline(
    training_x[, "perc.alumni"],
    training_y,
    df = 20
  )

ss.30 <-
  smooth.spline(
    training_x[, "perc.alumni"],
```

```

    training_y,
    df = 30
  )

ss.40 <-
  smooth.spline(
    training_x[, "perc.alumni"],
    training_y,
    df = 40
  )

ss.50 <-
  smooth.spline(
    training_x[, "perc.alumni"],
    training_y,
    df = 50
  )

pred.ss.10 <-
  predict(
    ss.10,
    x = testing_x[, "perc.alumni"],
  )

pred.ss.20 <-
  predict(
    ss.20,
    x = testing_x[, "perc.alumni"]
  )

pred.ss.30 <-
  predict(
    ss.30,
    x = testing_x[, "perc.alumni"]
  )

pred.ss.40 <-
  predict(
    ss.40,
    x = testing_x[, "perc.alumni"]
  )

pred.ss.50 <-
  predict(
    ss.50,
    x = testing_x[, "perc.alumni"]
  )

combined <-
  cbind(
    x = testing_x[, "perc.alumni"],
    df10 = as.data.frame(pred.ss.10)[, "y"],
    df20 = as.data.frame(pred.ss.20)[, "y"],

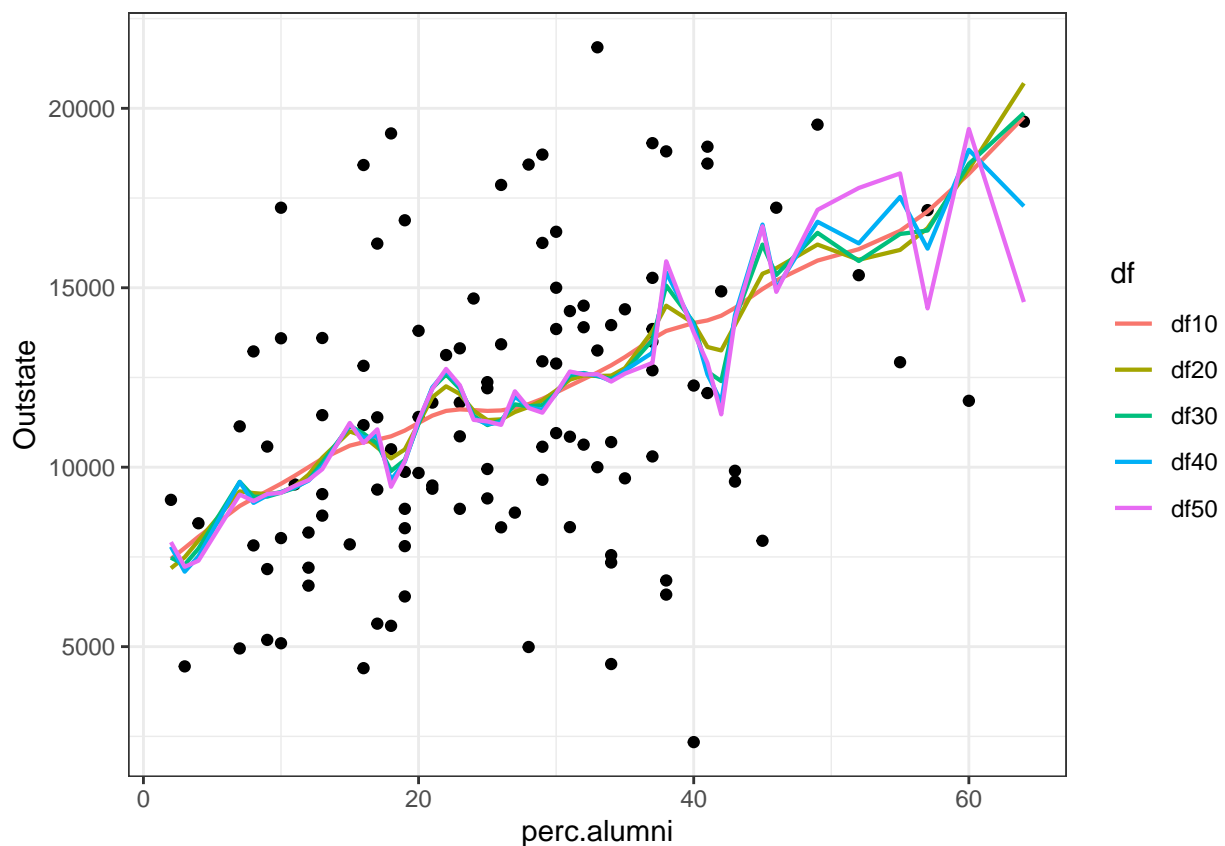
```

```

df30 = as.data.frame(pred.ss.30)[, "y"],
df40 = as.data.frame(pred.ss.40)[, "y"],
df50 = as.data.frame(pred.ss.50)[, "y"]
) |>
as.data.frame() |>
pivot_longer(
  !x,
  names_to = "df",
  values_to = "pred"
)

ggplot(aes(x = perc.alumni, y = Outstate), data = testing) +
  geom_point() +
  geom_line(aes(x = x, y = pred, col = df), linewidth = 0.75, data = combined) +
  theme_bw()

```



As visualized above, when fitting the predicted values for a given degree of freedom to the training dataset, as the degrees of freedom increases, the roughness of the fitted curve increases. With degrees of freedom = 10, the curve is at its smoothest, but with degrees of freedom = 50, the curve is most jagged.

```

ss <-
smooth.spline(
  training_x[, "perc.alumni"],
  training_y
)

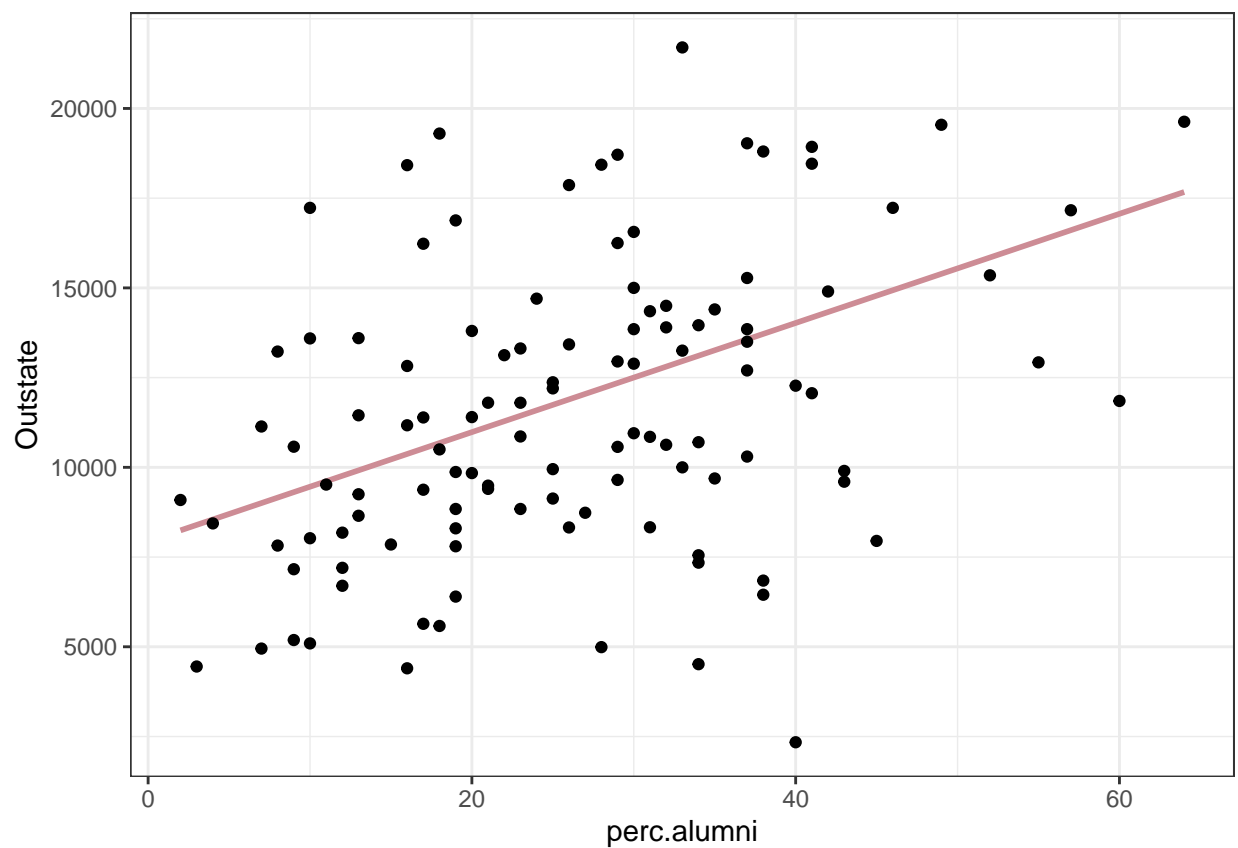
```

```

pred.ss <-
  predict(
    ss,
    x = testing_x[, "perc.alumni"]
  )

pred.ss |>
  as.data.frame() |>
  ggplot(aes(x = x, y = y)) +
  geom_line(linewidth = 1, color = "lightpink3") +
  geom_point(aes(x = testing_x[, "perc.alumni"], y = testing_y)) +
  labs(y = "Outstate", x = "perc.alumni") +
  theme_bw()

```



```
ss$df
```

```
## [1] 2.000232
```

When not passing a prespecified degree of freedom to `smooth.spline`, it automatically uses generalized cross validation to get the optimal degree of freedom, , which is used to get the optimal fit for smoothing spline. This is plotted above.

Fitting multivariate adaptive regression spline

To get the optimal number of terms and degree of interactions, ensure that the plotted RMSE has approximately a U shaped curve and the lowest RMSE is plotted.

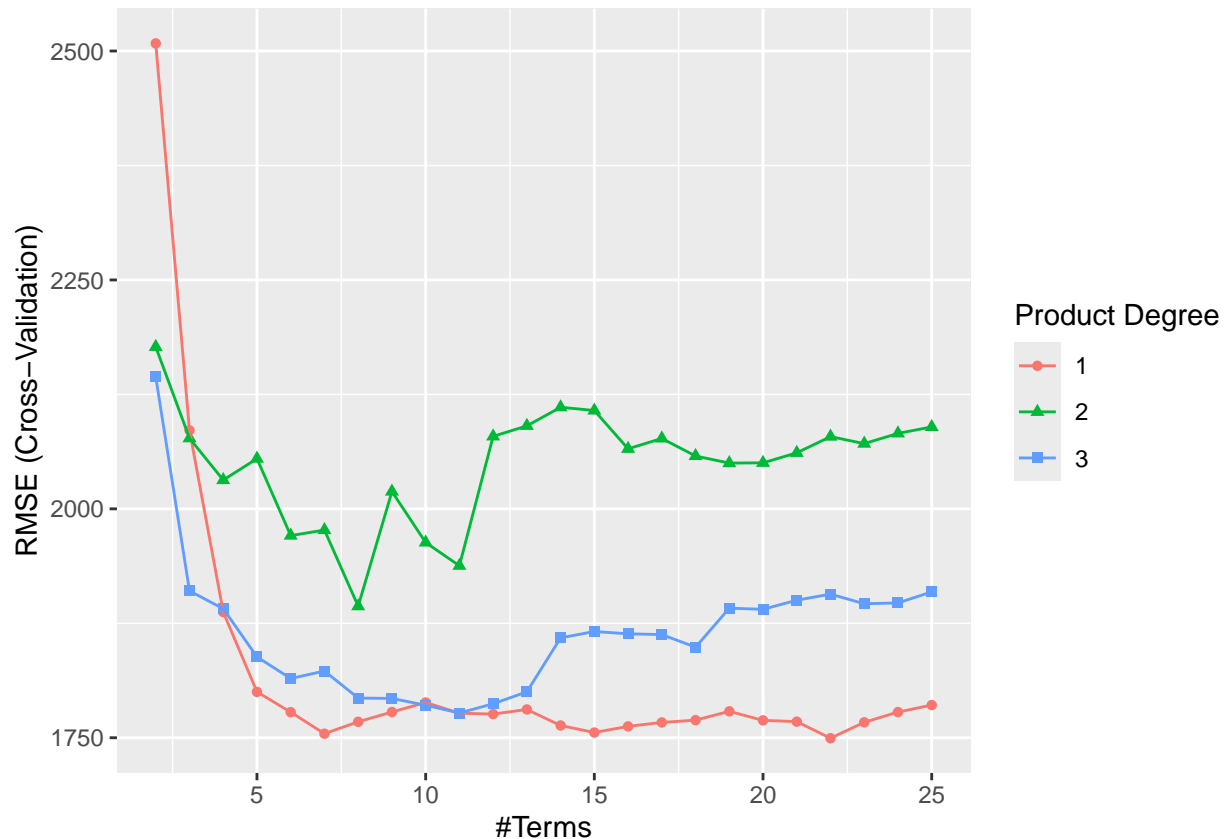
```
set.seed(81062)

ctrl1 <- trainControl(method = "cv", number = 10)

mars_grid <- expand.grid(degree = 1:3,
                        nprune = 2:25)

mars <-
  train(
    training_x, training_y,
    method = "earth",
    tuneGrid = mars_grid,
    trControl = ctrl1
  )

ggplot(mars)
```



```
mars$bestTune
```

```
##      nprune degree
```

```
## 21      22      1
```

```
coef(mars$finalModel)
```

```
##              (Intercept)
##              10059.9773425
##              h(Expend-15387)
##              -0.5850216
##              h(Grad.Rate-91)
##              -130.7595670
##              h(91-Grad.Rate)
##              -36.2601721
##              h(4340-Room.Board)
##              -1.2913130
##              h(2037-Accept)
##              -1.6626110
##              h(888-Enroll)
##              3.9677227
##              h(PhD-81)
##              66.1370472
##              h(F.Undergrad-1433)
##              -0.3114914
##              h(1433-F.Undergrad)
##              -1.3008168
## CollegeWentworth Institute of Technology
##              -6308.7429973
##              CollegeSpelman College
##              -5773.8421038
##              h(1320-Personal)
##              1.1441838
##              CollegeTrinity University
##              -6005.6546736
##              h(Apps-2302)
##              0.3217272
##              h(Expend-5524)
##              0.6574164
##              CollegeMorehouse College
##              -5323.7827976
##              h(8.8-S.F.Ratio)
##              -337.9561569
##              CollegeGreen Mountain College
##              3948.7247871
##              CollegeTuskegee University
##              -4994.0997311
##              CollegeArkansas College (Lyon College)
##              -4941.0946766
##              CollegeXavier University of Louisiana
##              -4052.3934151
```

Arbitrary predictors in the model of Room.Board and Grad.Rate were selected to include in the partial dependence plot.

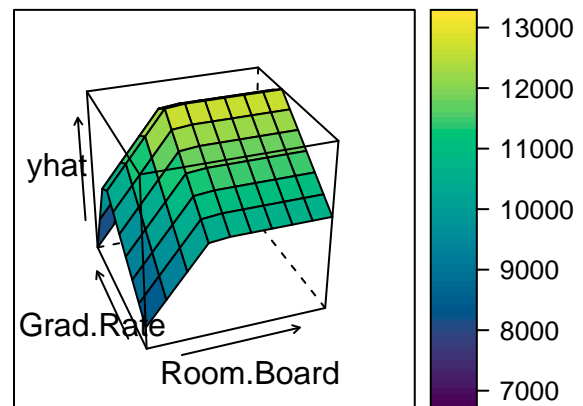
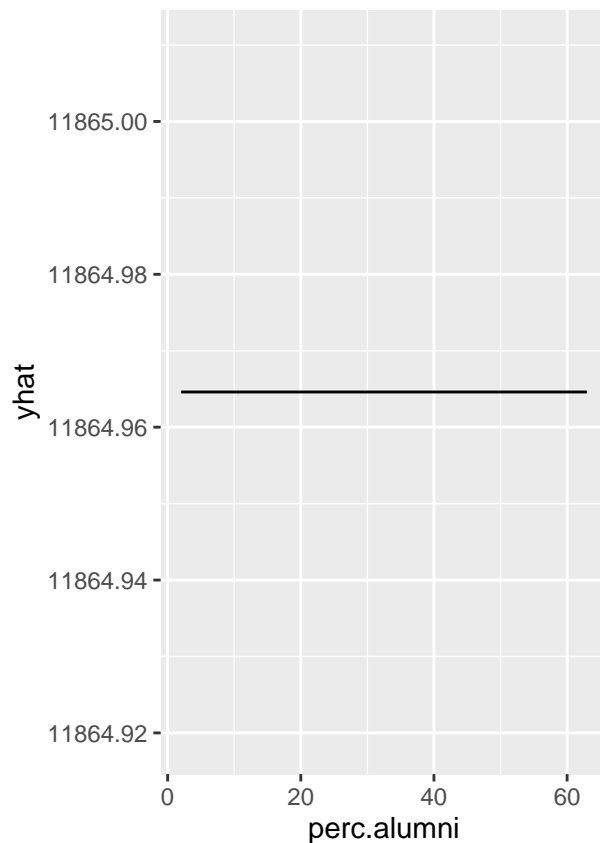
```

partial1 <-
  pdp::partial(
    mars, pred.var = c("perc.alumni"),
    grid.resolution = 10
  ) |>
  autoplot()

partial2 <-
  pdp::partial(
    mars, pred.var = c("Room.Board", "Grad.Rate"),
    grid.resolution = 10
  ) |>
  pdp::plotPartial(
    levelplot = FALSE,
    zlab = "yhat",
    drape = TRUE,
    screen = list(z = 20, x = -60)
  )

gridExtra::grid.arrange(partial1, partial2, ncol = 2)

```



Over the 10-fold cross validation, the mean CV RMSE across each fold is 1749.5082619.

```

mars_pred <- predict(mars, newdata = testing)

mars_pred_error <- mean((mars_pred - testing_y)^2)

```

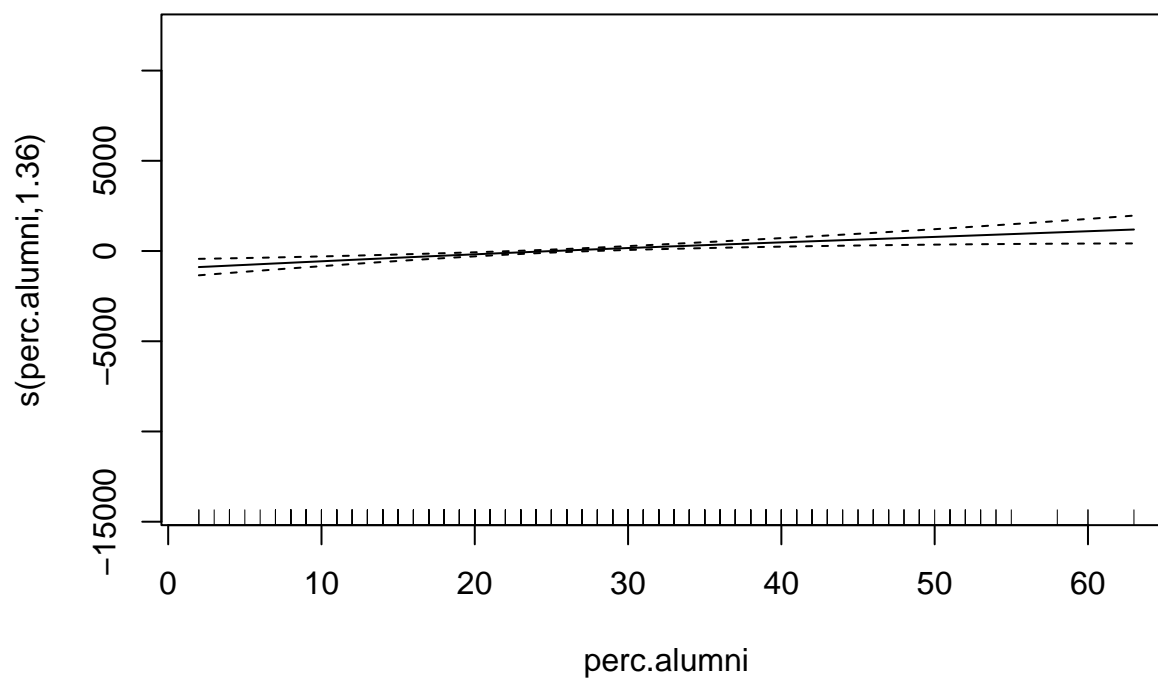
When fit to the testing dataset, the MARS model has a test MSE of NA.

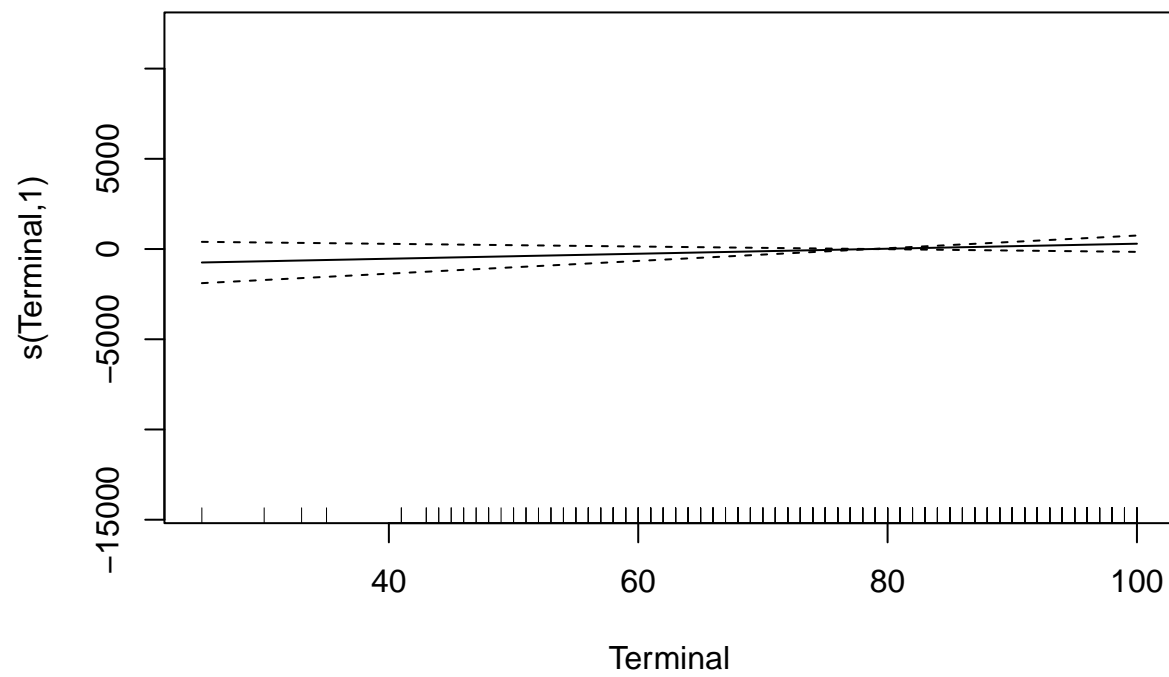
Fitting generalized additive model (GAM)

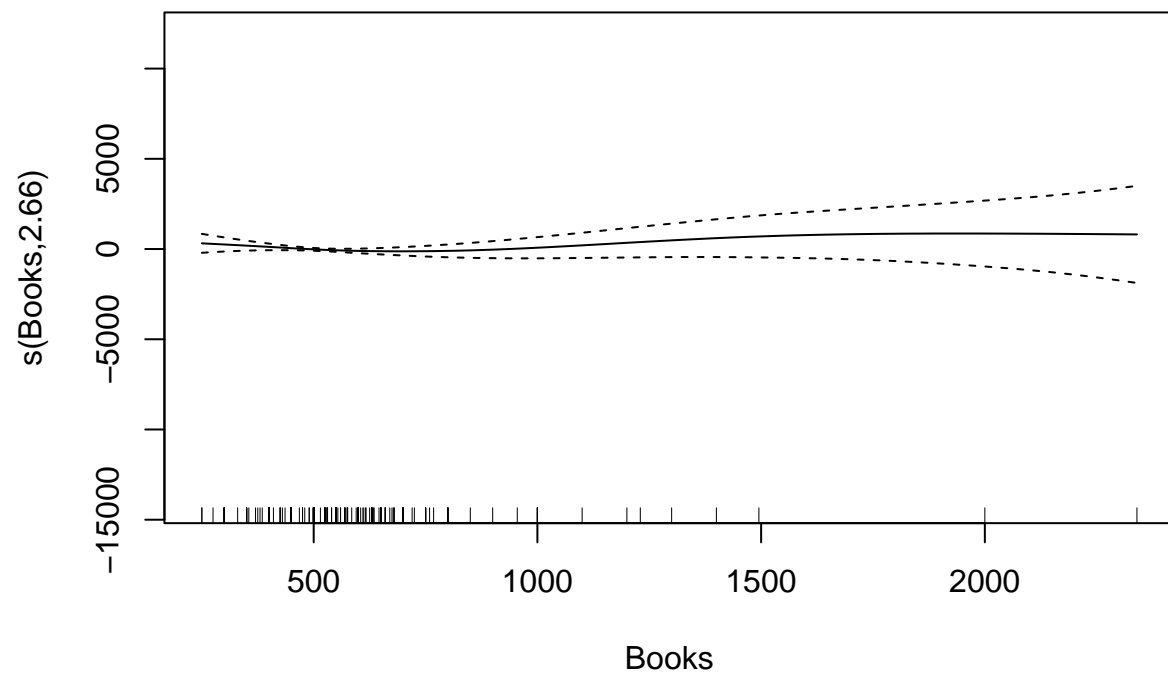
```
set.seed(81062)

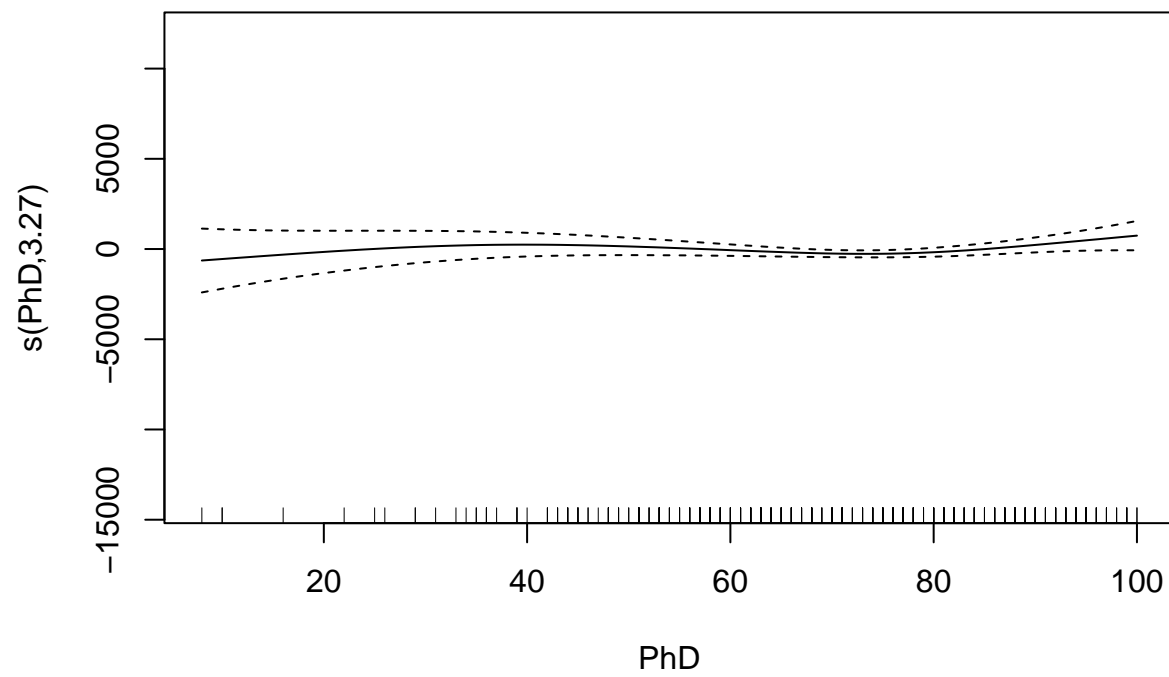
gam <- gam(Outstate ~ s(perc.alumni) + s(Terminal) + s(Books) + s(PhD) +
  s(Grad.Rate) + s(Top10perc) + s(Top25perc) + s(S.F.Ratio) +
  s(Personal) + s(P.Undergrad) + s(Room.Board) + s(Enroll) +
  s(Accept) + s(F.Undergrad) + s(Apps) + s(Expend),
  data = training)

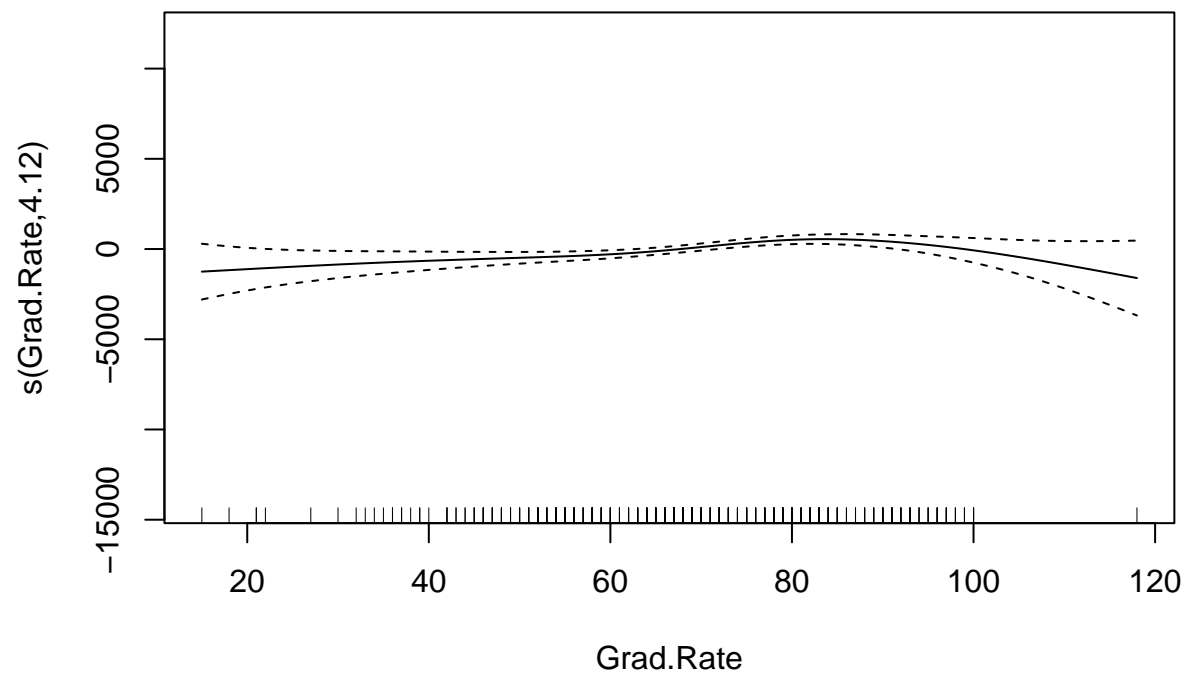
plot(gam)
```

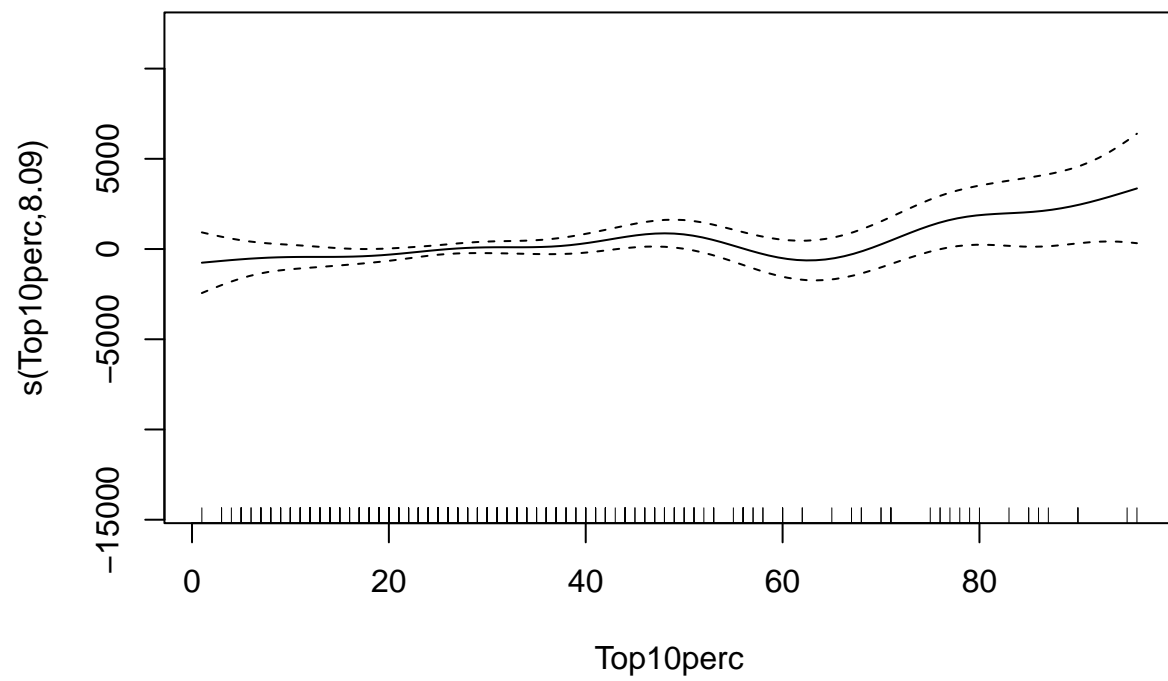


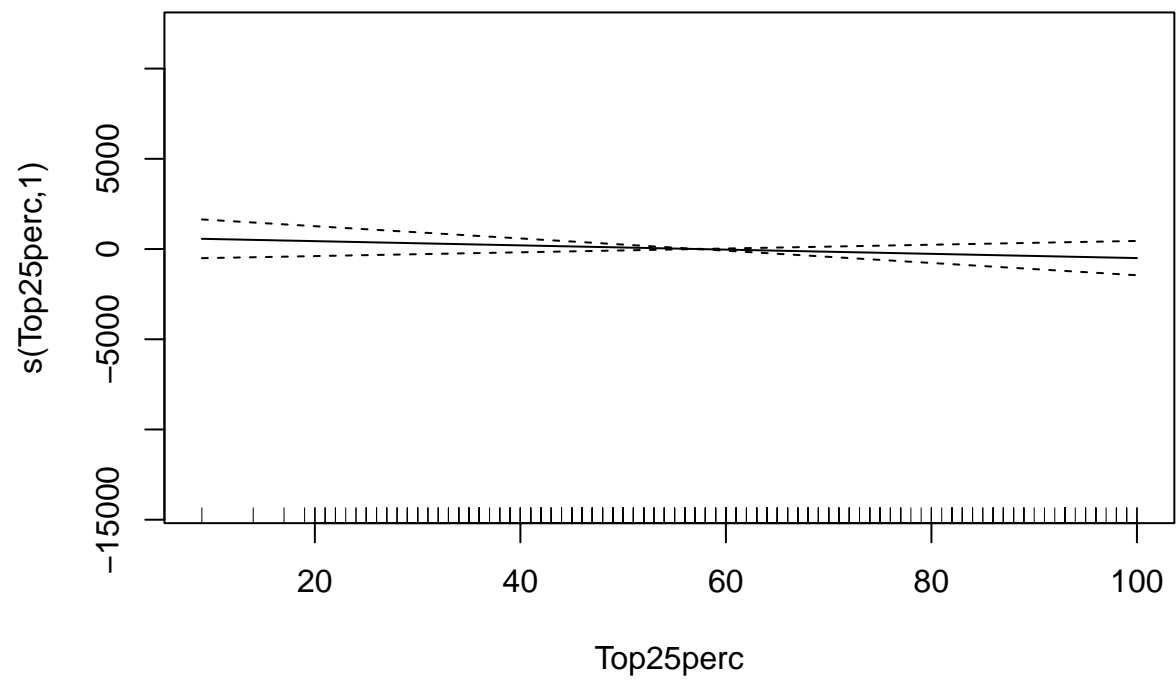


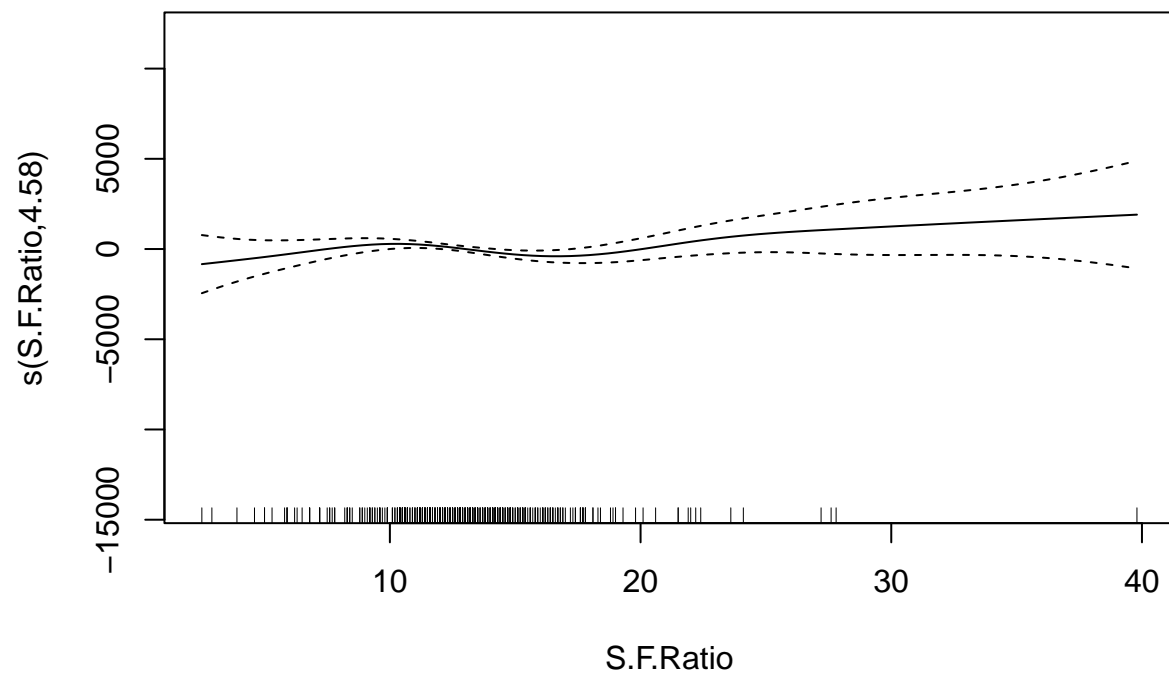


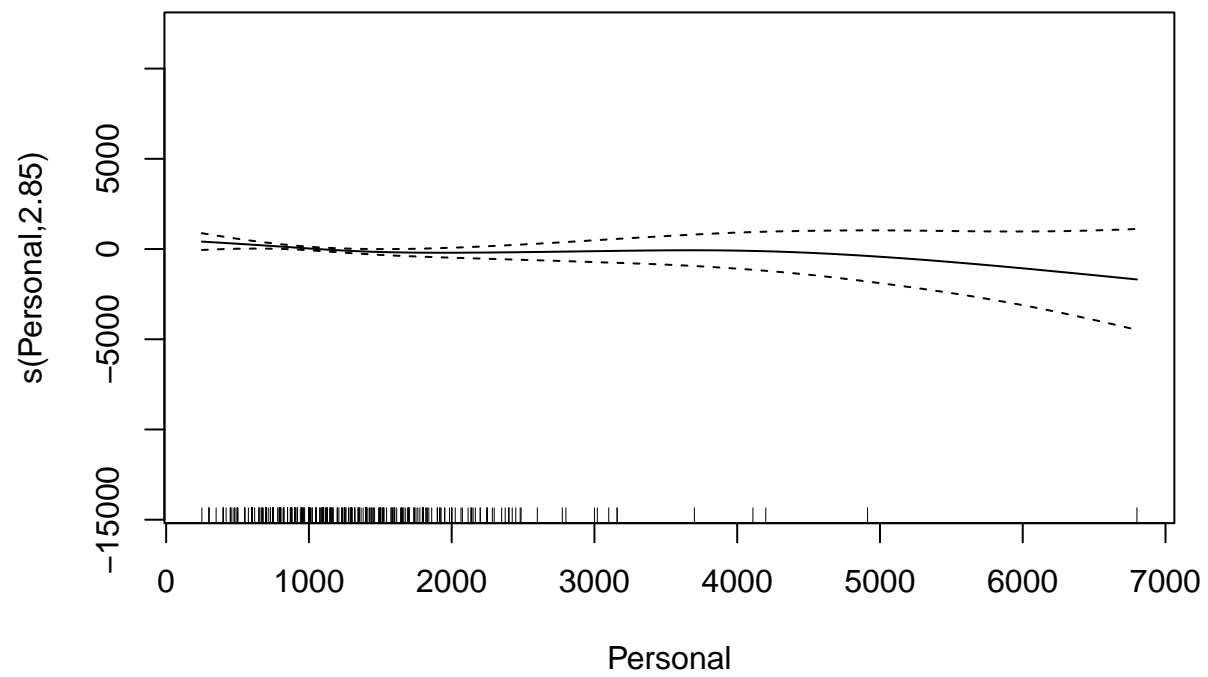


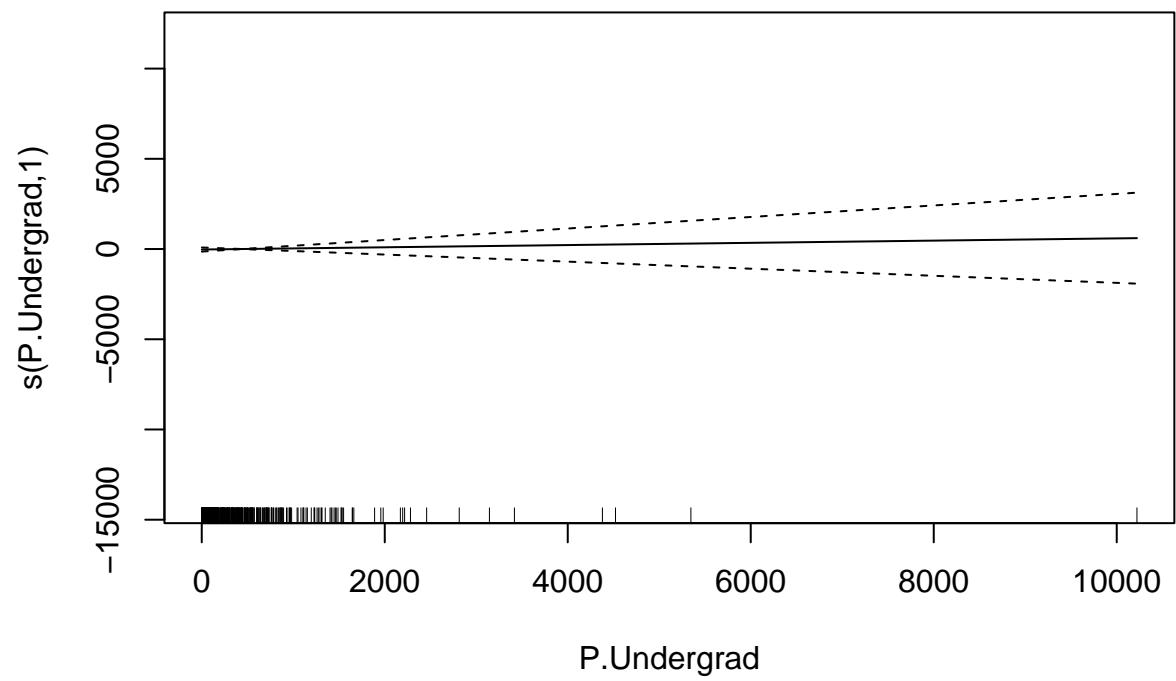


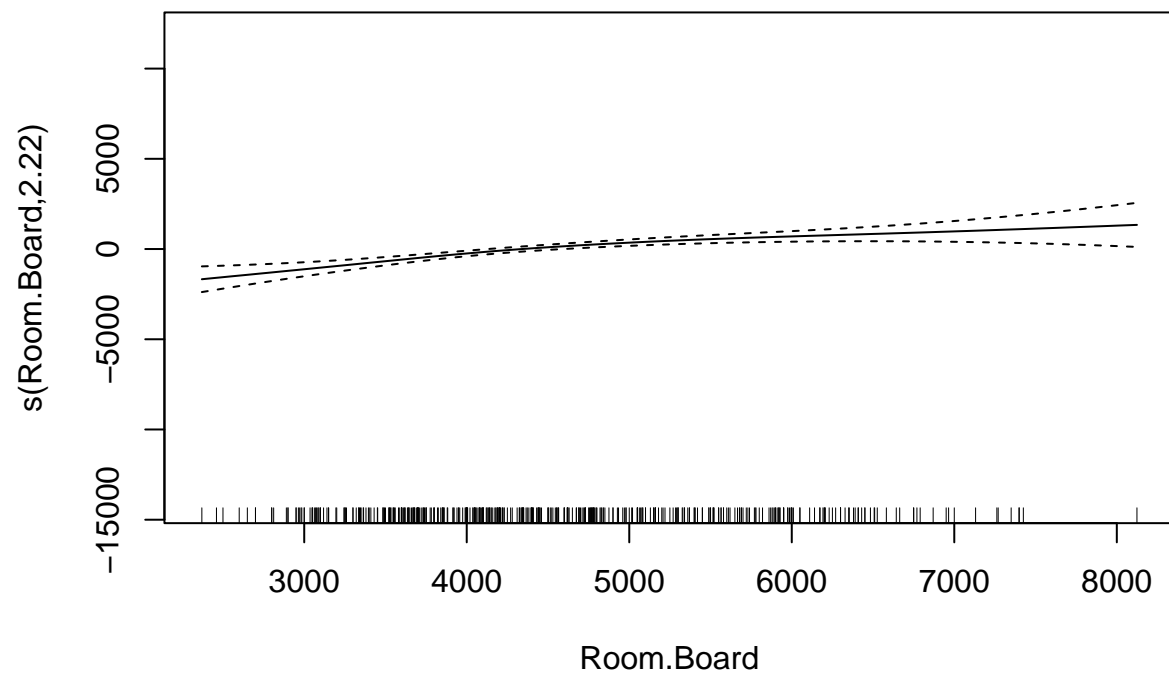


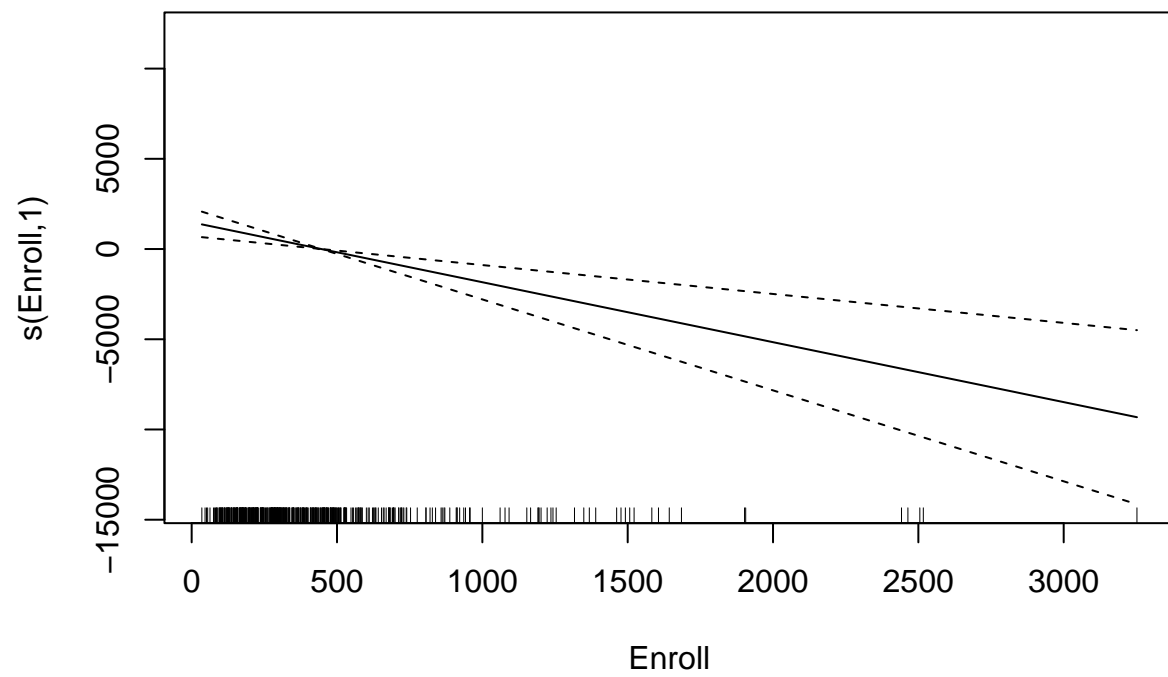


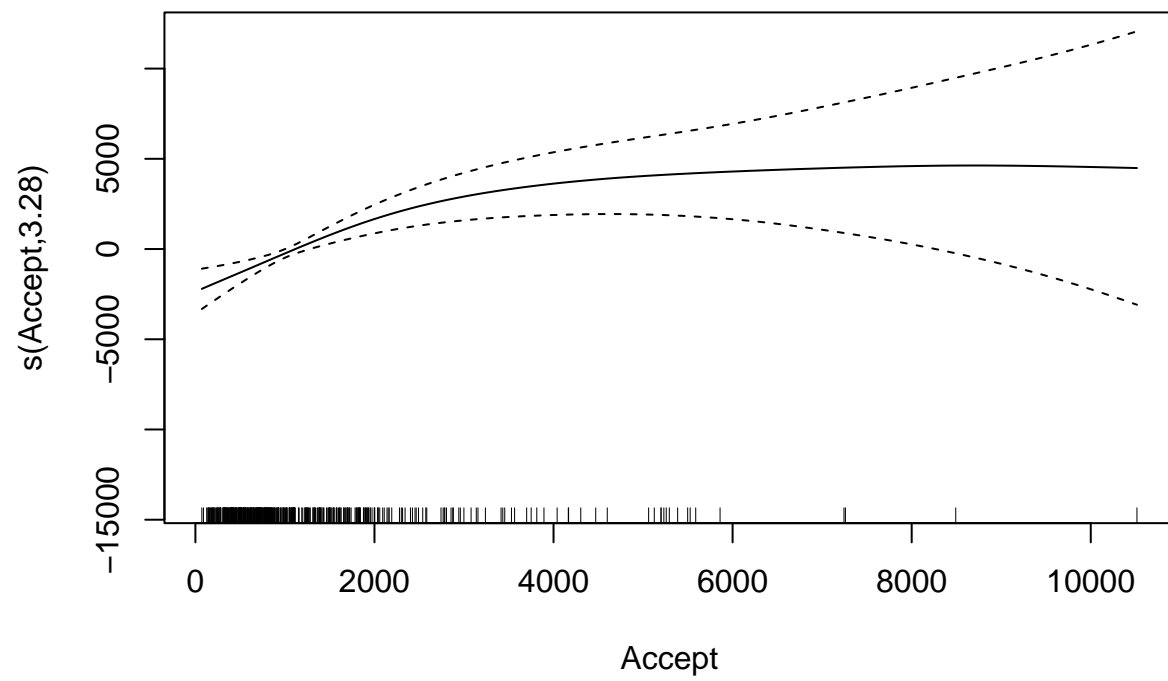


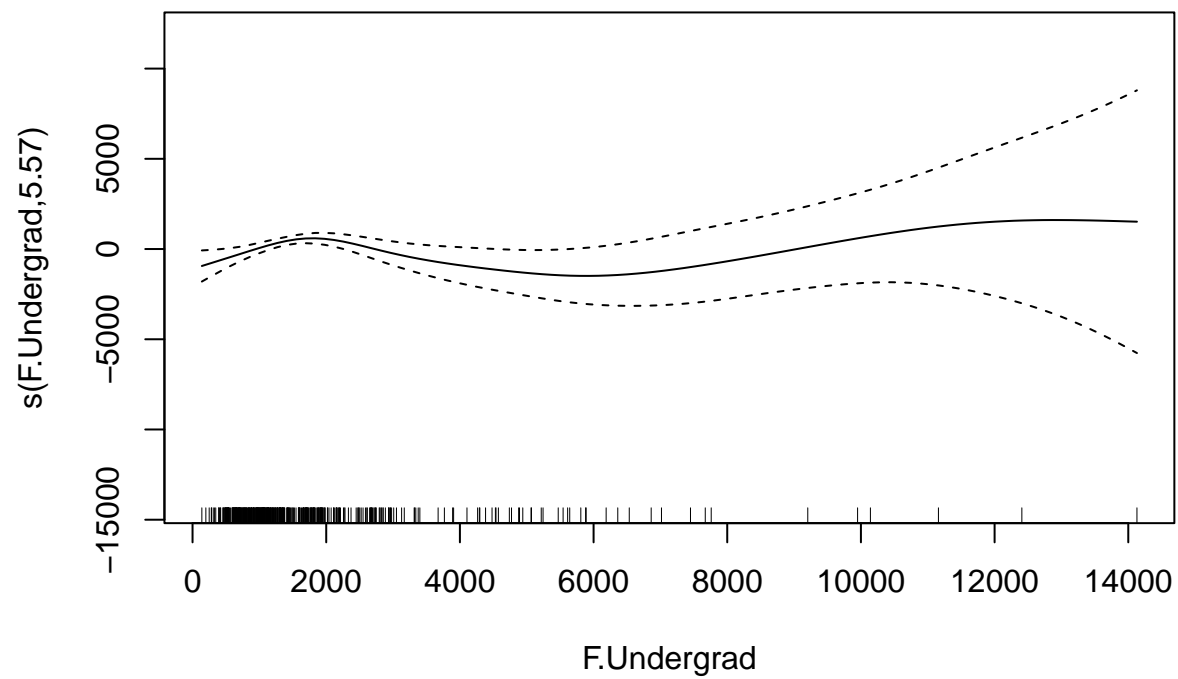


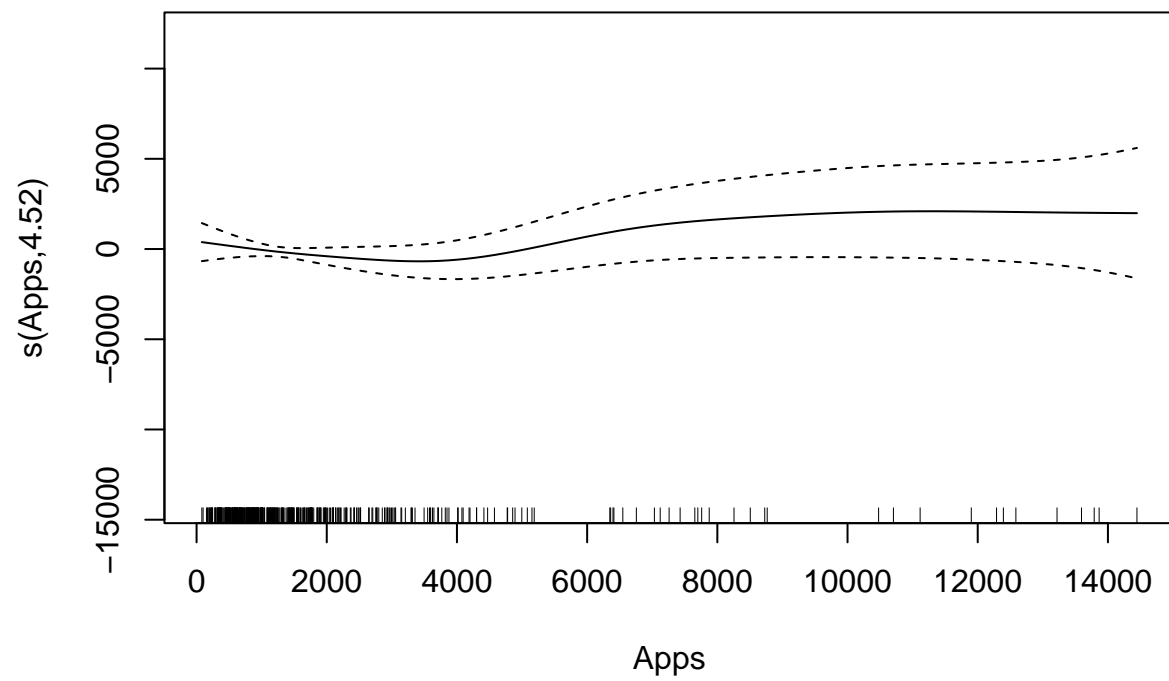


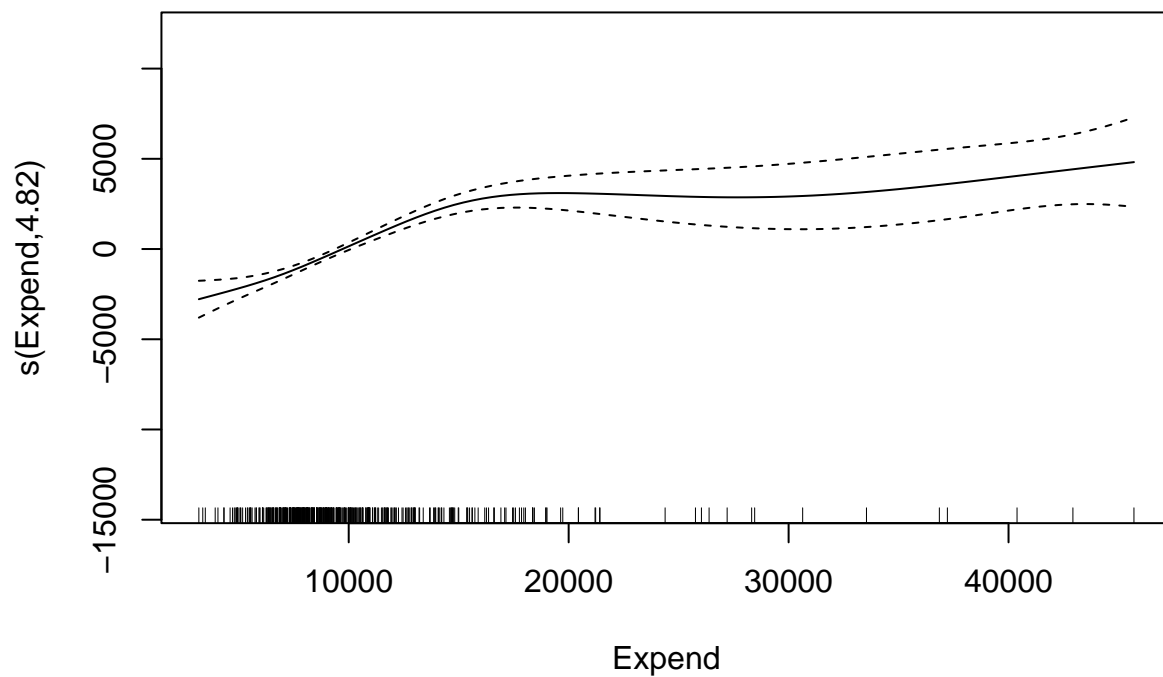












```
gam_2 <- train(
  training_x, training_y,
  method = "gam",
  trControl = ctrl1
)
```

```
gam_2$bestTune
```

```
## select method
## 1 FALSE GCV.Cp
```

```
gam_2$finalModel
```

```
##
## Family: gaussian
## Link function: identity
##
## Formula:
## .outcome ~ s(perc.alumni) + s(Terminal) + s(Books) + s(PhD) +
## s(Grad.Rate) + s(Top10perc) + s(Top25perc) + s(S.F.Ratio) +
## s(Personal) + s(P.Undergrad) + s(Room.Board) + s(Enroll) +
## s(Accept) + s(F.Undergrad) + s(Apps) + s(Expend)
##
## Estimated degrees of freedom:
## 1.36 1.00 2.66 3.27 4.12 8.09 1.00
```

```
## 4.58 2.85 1.00 2.22 1.00 3.28 5.57
## 4.52 4.82 total = 52.34
##
## GCV score: 2545855
```

```
gam_pred <- predict(gam, newdata = testing)

gam_pred_error <- mean((gam_pred - testing_y)^2)
```

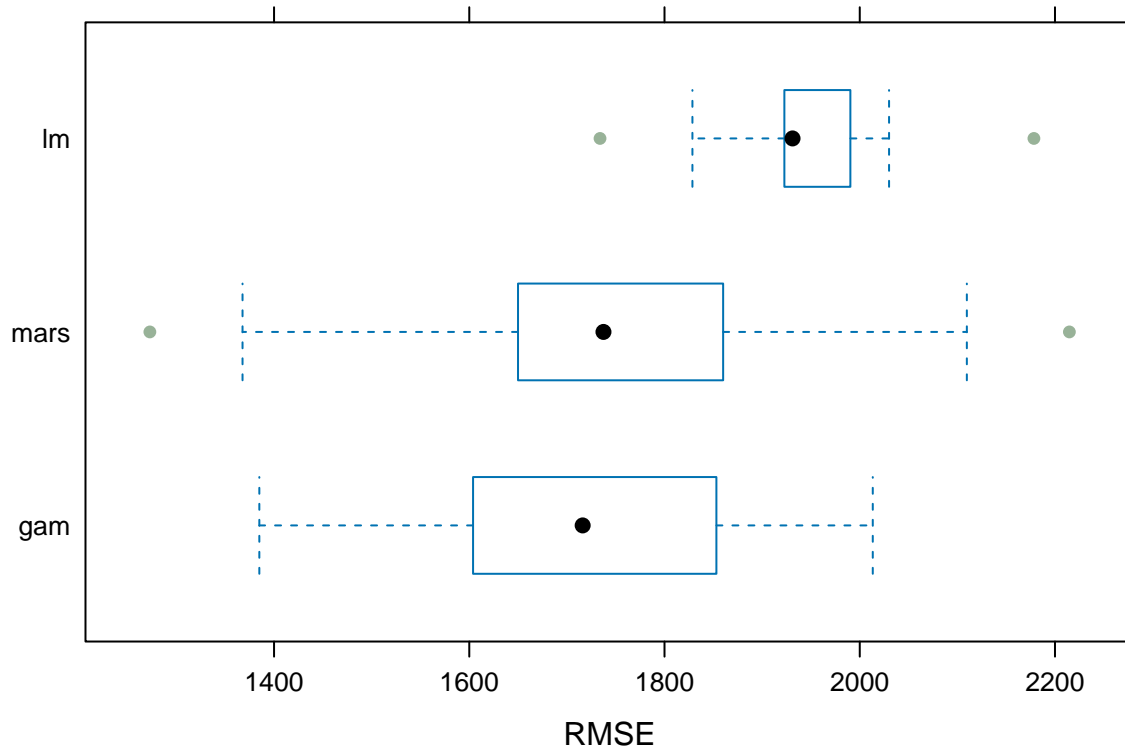
When fit to the testing dataset, the GAM has a test error MSE of 4.2132774×10^6 .

```
lm <-
  train(
    Outstate ~ perc.alumni + Terminal + Books + PhD +
    Grad.Rate + Top10perc + Top25perc + S.F.Ratio +
    Personal + P.Undergrad + Room.Board + Enroll +
    Accept + F.Undergrad + Apps + Expend,
    data = training,
    method = "lm",
    trControl = ctrl1
  )

resamp <-
  resamples(
    list(
      gam = gam_2,
      mars = mars,
      lm = lm
    )
  )

mars_rmse <- mean(resamp$values$`mars~RMSE`)
gam_rmse <- mean(resamp$values$`gam~RMSE`)
lm_rmse <- mean(resamp$values$`lm~RMSE`)

bwplot(resamp, metric = "RMSE")
```



```
lm_pred <- predict(lm, newdata = testing)
lm_pred_error <- mean((lm_pred - testing_y)^2)
```

For this particular dataset, I would favor MARS over a linear model to predict out of state tuition. When comparing the CV mean RMSEs, the MARS model has a mean RMSE of 1749.5082619, while the linear model has a mean RMSE of 1942.6879676. While only slightly higher, this is likely from the non-linear relations out of state tuition has with some predictors like `perc.alumni`, `P.Undergrad`, and `F.Undergrad`. Additionally, MARS is better able to account for interaction.

More broadly, it cannot be generalized whether MARS or linear models are preferred. It depends on a given dataset. If the dataset can be simply predicted with a linear relationship, without much interaction, then a linear model may be preferable, especially since it is easily interpretable and provides 95% CI. However, if there is considerable interaction and non-linear relations, then MARS may be preferred.