

# P8106 Data Science II Midterm Project Code: Predicting COVID-19 Recovery Time and Identifying Important Risk Factors

Sarah Forrest - sef2183

4/5/2023

## Contents

Background	2
Data	2
Exploratory Analysis and Data Visualization:	3
Model Training	5
Results	17
Conclusion	19
Appendix	19

```
library(tidyverse)
library(caret)
library(mgcv)
library(earth)
```

## Background

To gain a better understanding of the factors that predict recovery time from COVID-19 illness, this study was designed to combine three existing cohort studies that have been tracking participants for several years. This study collects recovery information through questionnaires and medical records, and leverages existing data on personal characteristics prior to the pandemic. The ultimate goal is to develop a prediction model for recovery time and identify important risk factors for long recovery time.

## Data

This study uses data from the `recovery.RData` file. The dataset contains a variable for the time from COVID-19 infection to recovery (in days), which is the outcome of interest in this study. It also contains 14 predictor variables, including demographic characteristics, personal characteristics, vital measurements, and disease status. The predictors are a mix of continuous and categorical variables.

The `recovery.RData` file consists of data on 10000 participants. A random sample of 2000 participants was used for analysis. Additionally, the random sample of 2000 participants was further split into training and test datasets using the `createDataPartition()` function in R. The training dataset contained 70% of the sample and the test dataset contained the remaining 30%. All data partitions were conducted using a seed set to my UNI number (2183) for reproducibility.

```
# read in dataset
load("data/recovery.Rdata")

set.seed(2183) # for reproducibility - my UNI

# create a random sample of 2000 participants
dat <- dat[sample(1:10000, 2000),]

# remove the id variable from the dataset
dat <- dat %>%
  select(-id)

set.seed(2183)

# specify rows of training data (70% of the dataset)
trRows <- createDataPartition(dat$recovery_time, p = 0.7, list = FALSE)

# training data
dat_train <- dat[trRows, ]
## matrix of predictors
x <- model.matrix(recovery_time~.,dat)[trRows,-1]
## vector of response
y <- dat$recovery_time[trRows]

# test data
```

```

dat_test <- dat[-trRows, ]
## matrix of predictors
x2 <- model.matrix(recovery_time~.,dat)[-trRows,-1]
## vector of response
y2 <- dat$recovery_time[-trRows]

```

## Exploratory Analysis and Data Visualization:

Lattice plots were created using the `featurePlot()` function in the `caret` package to explore the multivariate dataset and identify patterns or relationships. A plot is created for each of the 14 predictors in the dataset in order to visualize each predictor's association with the outcome, COVID-19 recovery time.

```

# create dataset for exploratory analysis and data visualization
dat_train_viz <- dat_train %>%
  mutate(study = case_when( # turn study (character variable) into a numeric variable
    study == "A" ~ 1,
    study == "B" ~ 2,
    study == "C" ~ 3))

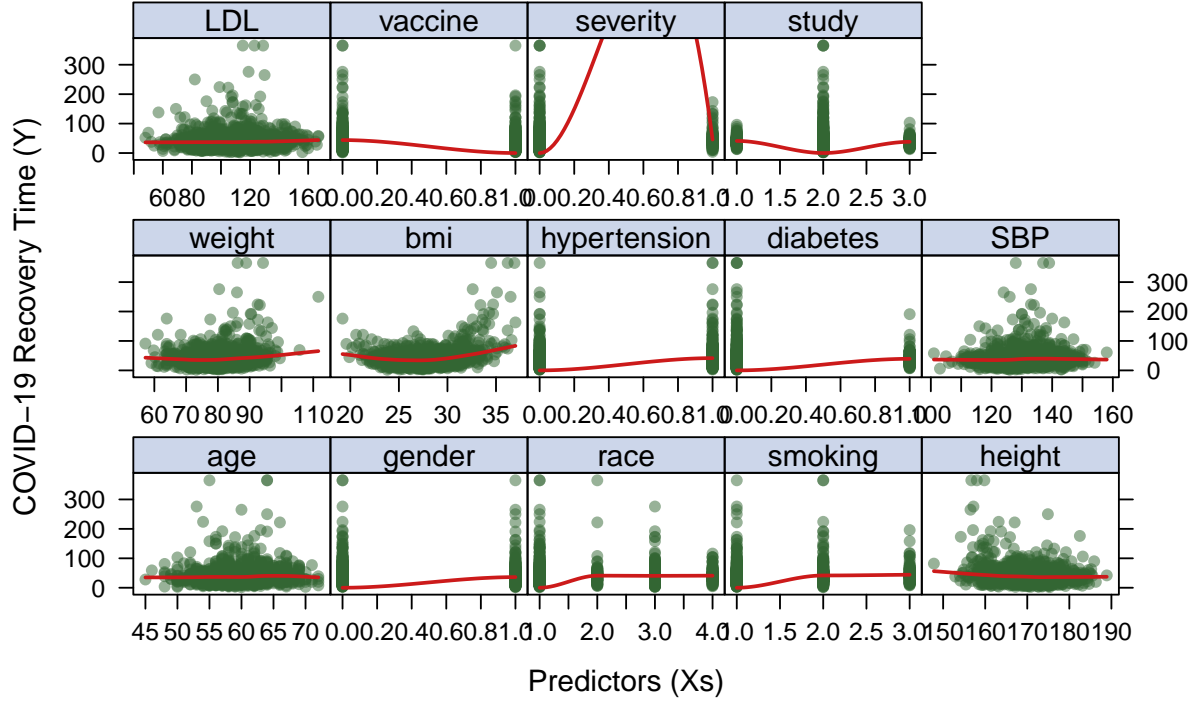
# Find the remaining non-numeric columns
non_numeric_cols <- sapply(dat_train_viz, function(x) !is.numeric(x))

# Convert non-numeric columns to numeric
dat_train_viz[, non_numeric_cols] <- lapply(dat_train_viz[, non_numeric_cols], as.numeric) # turn factors to numeric

# set various graphical parameters (color, line type, background, etc) to control the look of trellis d
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(x = dat_train_viz[,1:14],
            y = dat_train_viz[,15],
            plot = "scatter",
            span = .5,
            labels = c("Predictors (Xs)", "COVID-19 Recovery Time (Y)"),
            main = "Figure 1. Lattice Plot",
            type = c("p", "smooth"))

```

**Figure 1. Lattice Plot**

Using the lattice plot in Figure 1. the following patterns were observed:

- Patient age does not appear to have a clear association with COVID-19 recovery time. The mostly straight horizontal line shows that both low and high values for age have similar values for recovery time. The age variable appears to be normally distributed, with most data points in the middle between the minimum age (42) and maximum age (79).
- For the gender variable, male patients were assigned a value of 1 and female patients were assigned a value of 0. Therefore, the plot shows that males appear to have a longer recovery time than females in the dataset.
- For the race variable, White patients were assigned a value of 1, Asian patients were assigned 2, Black patients were assigned 3, and Hispanic patients were assigned 4. The plot shows that White patients appear to have a shorter recovery time compared to all other races.
- For the smoking variable, patients who never smoked were assigned a value of 0, former smokers were assigned 1, and current smokers were assigned 2. The plot shows that never smokers appear to have a shorter recovery time compared to former and current smokers. There is a positive association between smoking and recovery time.
- A slight negative association can be observed between patient height and recovery time. Patients with lower height values (shorter patients) appear to have a slightly longer recovery time than patients with higher height values (taller patients).
- A slight positive association can be observed between patient weight and recovery time. Patients with lower weight values appear to have a slightly shorter recovery time than patients with higher weight values.
- A slight U-shaped association can be observed between patient bmi and recovery time. Patients with bmi values near the minimum (19.2) and maximum (37.7) in the dataset have longer recovery times than patients with bmi values in the middle (25-30).
- For the hypertension variable, patients without hypertension were assigned a value of 0 and patients with hypertension were assigned a value of 1. The plot shows that patients with hypertension appear to have a longer recovery time compared to patients without hypertension. There is a positive association

between hypertension and recovery time.

- For the diabetes variable, patients without diabetes were assigned a value of 0 and patients with diabetes were assigned a value of 1. The plot shows that patients with diabetes appear to have a longer recovery time compared to patients without diabetes. There is a positive association between diabetes and recovery time.
- Patient systolic blood pressure (SBP) does not appear to have a clear association with recovery time. The mostly straight horizontal line shows that both low and high values for SBP have similar values for recovery time.
- Patient LDL cholesterol does not appear to have a clear association with recovery time. The mostly straight horizontal line shows that both low and high values for LDL have similar values for recovery time.
- For the vaccine variable, unvaccinated patients at the time of infection were assigned a value of 0 and vaccinated patients were assigned a value of 1. The plot shows that vaccinated patients appear to have a shorter recovery time compared to unvaccinated patients. There is a negative association between vaccination and recovery time.
- For the severity variable, patients without severe COVID-19 infection were assigned a value of 0 and patients with severe COVID-19 infection were assigned a value of 1. The plot shows that patients with severe infections appear to have a longer recovery time compared to patients without severe infections. There is a positive association between severity and recovery time.
- For the study variable, patients in study A were assigned a value of 1, patients in study B were assigned a value of 2, patients in study C were assigned a value of 3. The plot shows that patients in study B had a shorter recovery time compared to patients in studies A and C.

## Model Training

Several different models were fit using all predictors to predict time to recovery from COVID-19.

The `train()` function from the `caret` package was used to fit each model using either the training dataset or a matrix of all the predictors and a vector of the response variable, `recovery_time`, from the training dataset as inputs. Each model was also fit using cross-validation with 10 folds repeated 5-times. The `trainControl()` function was used to specify this cross-validation method, which was called on within the `train()` function using the `trControl` argument. Additionally, the `method` argument was used within the `train()` function to specify the type of model to fit. The resulting model object for each model contains the final model (`finalModel`) and information about the cross-validation performance.

The `predict()` function from the `caret` package was also used to generate predictions for the test dataset using each final model that was trained using the training dataset. The root mean squared error (RMSE) between the predicted and actual recovery times on the test dataset was calculated in order to evaluate each model's performance and support model comparison.

This approach results in a final model that has been trained using cross-validation, which can help to reduce overfitting and improve the generalization performance of the model on new, unseen data for future prediction purposes.

### 1. Linear model

A linear model that assumes a linear relationship between the predictor and response variables (linearity), and assumes that the errors are normally distributed (normality) and have constant variance (homoscedasticity), and that the observations are independent of each other (independence). The linear model is the most basic and assumes a linear relationship between the variables, while the other models allow for more flexible relationships. A `method = "lm"` argument was used within the `train()` function to specify a linear model fit.

```
set.seed(2183)
```

```

# 10-fold cross-validation repeated 5 times using the best rule
ctrl <- trainControl(method = "repeatedcv", number = 10, repeats = 5)

set.seed(2183)

# Fit a linear regression model using cross-validation on the training dataset
linear_model <- train(recovery_time ~ age + gender + race + smoking +
                      height + weight + bmi + hypertension +
                      diabetes + SBP + LDL +
                      vaccine + severity + study,
                      data = dat_train, # training dataset
                      method = "lm",
                      trControl = ctrl)

# view the model summary
summary(linear_model$finalModel)
##
## Call:
## lm(formula = .outcome ~ ., data = dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -64.119 -14.513  -1.103   9.605  252.860
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.109e+03  2.031e+02 -15.306 < 2e-16 ***
## age         -1.049e-01  1.839e-01  -0.570 0.568643
## gender      -7.061e+00  1.392e+00  -5.072 4.47e-07 ***
## race2        1.708e+00  2.928e+00   0.583 0.559903
## race3       -1.384e+00  1.770e+00  -0.782 0.434469
## race4       -1.644e+00  2.409e+00  -0.682 0.495145
## smoking1     4.899e+00  1.560e+00   3.140 0.001722 **
## smoking2     8.278e+00  2.372e+00   3.491 0.000497 ***
## height       1.804e+01  1.195e+00  15.103 < 2e-16 ***
## weight      -1.927e+01  1.266e+00 -15.227 < 2e-16 ***
## bmi          5.844e+01  3.617e+00  16.159 < 2e-16 ***
## hypertension  3.662e+00  2.330e+00   1.572 0.116236
## diabetes     -3.433e+00  1.977e+00  -1.737 0.082664 .
## SBP          1.001e-01  1.524e-01   0.657 0.511393
## LDL          2.597e-03  3.806e-02   0.068 0.945611
## vaccine      -8.170e+00  1.414e+00  -5.778 9.33e-09 ***
## severity      6.256e+00  2.396e+00   2.611 0.009120 **
## studyB        4.355e+00  1.777e+00   2.450 0.014398 *
## studyC       -1.339e+00  2.197e+00  -0.609 0.542396
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 25.88 on 1383 degrees of freedom
## Multiple R-squared:  0.2909, Adjusted R-squared:  0.2817
## F-statistic: 31.52 on 18 and 1383 DF,  p-value: < 2.2e-16

# view performance on the test set (RMSE)

```

```
test_pred <- predict(linear_model, newdata = dat_test) # test dataset
test_rmse <- sqrt(mean((test_pred - dat_test$recovery_time)^2))
test_rmse
## [1] 27.34561
```

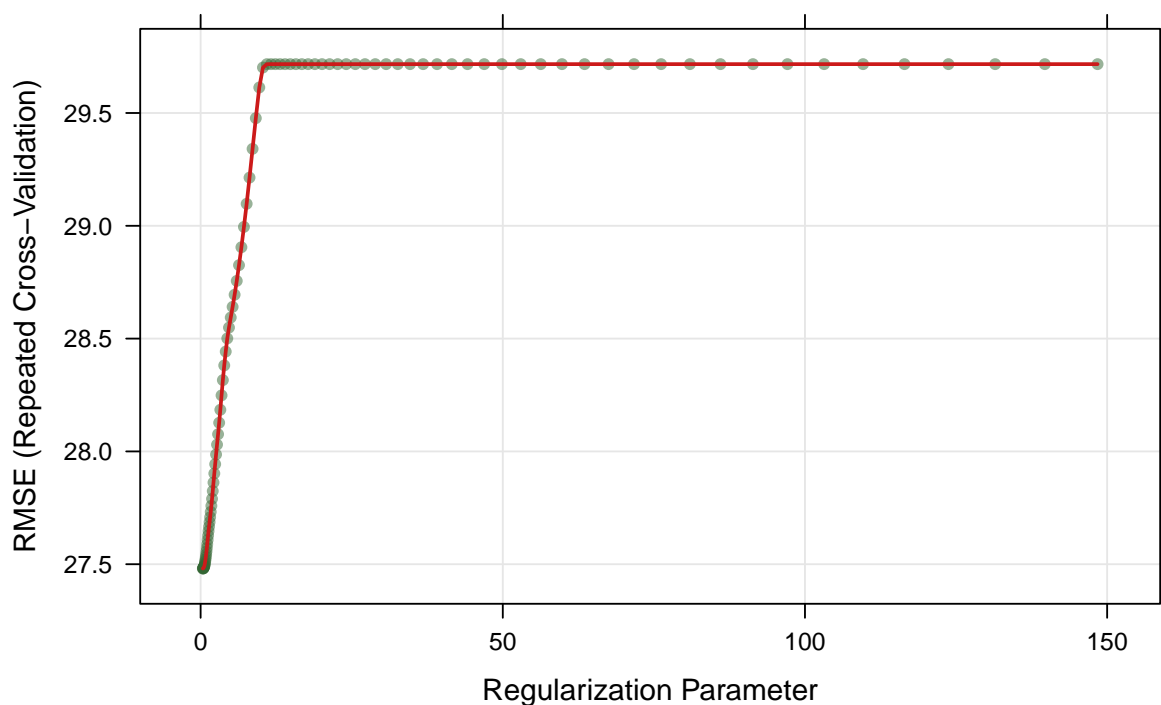
## 2. Lasso model

A lasso model is a linear regression model that adds a penalty term to the sum of absolute values of the coefficients to prevent overfitting, which can lead to sparse models by shrinking some coefficients to zero. It has the same assumptions as the linear model. A `method = "glmnet"` argument was used within the `train()` function to specify a lasso model fit. Additionally, a grid of tuning parameters for the model were set using the `tuneGrid` argument within the `train()` function. The grid contains 100 values of the `lambda` parameter, ranging from  $\exp(-1)$  to  $\exp(5)$  with equal intervals on the log scale. The `alpha` parameter is set to 1, indicating that we are only considering Lasso regularization. The `expand.grid()` function is used to generate all possible combinations of `alpha` and `lambda` in the grid for tuning.

```
set.seed(2183)

lasso_model <- train(x, y, # training dataset
  method = "glmnet",
  tuneGrid = expand.grid(alpha = 1,
    lambda = exp(seq(-1, 5, length = 100))),
  trControl = ctrl)

# view the model cross-validation (tuning parameter selection) plot and summary
print(plot(lasso_model))
```



```
summary(lasso_model$finalModel)
##           Length Class      Mode
## a0           94  -none-   numeric
## beta        1692 dgCMatrix S4
## df           94  -none-   numeric
## dim           2  -none-   numeric
## lambda       94  -none-   numeric
## dev.ratio     94  -none-   numeric
## nulldev        1  -none-   numeric
## npasses        1  -none-   numeric
## jerr           1  -none-   numeric
## offset         1  -none-   logical
## call           5  -none-   call
## nob           1  -none-   numeric
## lambdaOpt      1  -none-   numeric
## xNames         18  -none-   character
## problemType    1  -none-   character
## tuneValue       2 data.frame list
## obsLevels      1  -none-   logical
## param           0  -none-   list

# view performance on the test set (RMSE)
test_pred <- predict(lasso_model, newdata = x2) # test dataset
test_rmse <- sqrt(mean((test_pred - dat_test$recovery_time)^2))
test_rmse
## [1] 29.24864
```

### 3. Elastic net model

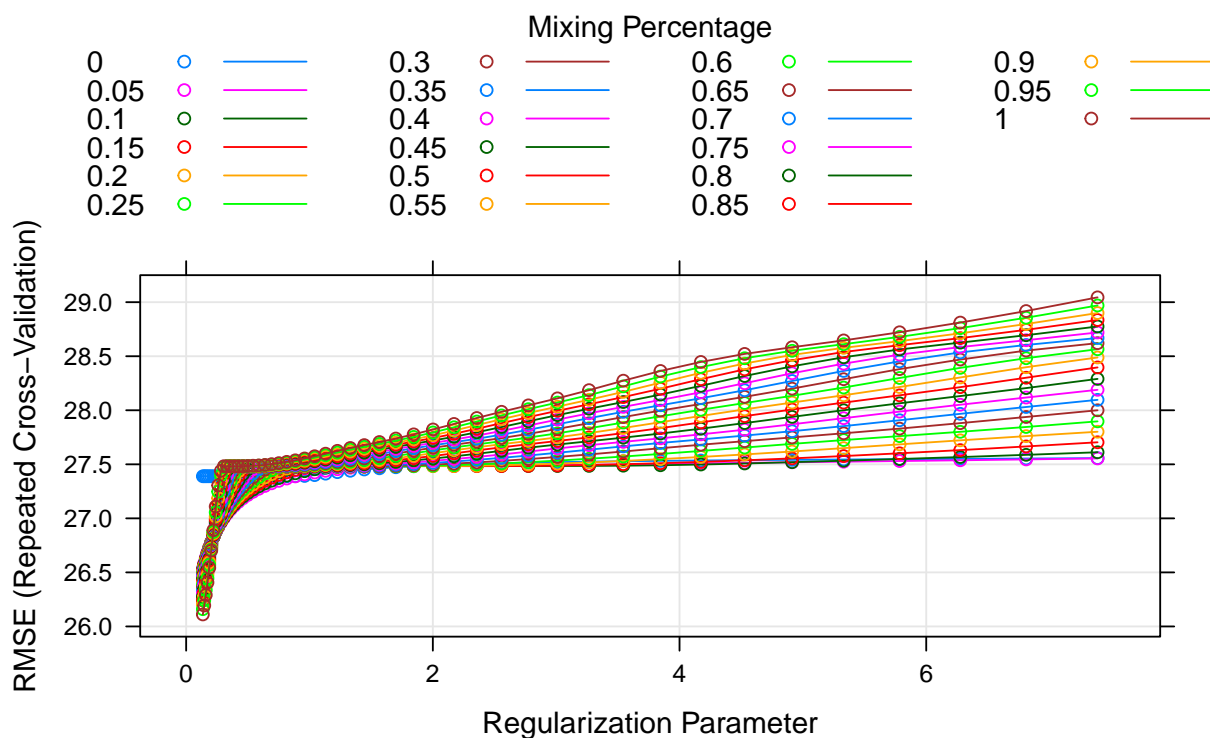
An elastic net model is linear regression model that combines the penalties of the lasso and ridge regression methods to prevent overfitting, which can result in better prediction accuracy than either method alone. It has the same assumptions as the linear model. A method = “glmnet” argument was used within the `train()` function to specify a elastic net model fit. Additionally, a grid of tuning parameters for the model were set using the `tuneGrid` argument within the `train()` function. The grid includes 21 values of the alpha parameter, ranging from 0 to 1 with equal intervals. The lambda parameter is set to a sequence of 50 values, ranging from  $\exp(2)$  to  $\exp(-2)$  with equal intervals on the log scale. The `expand.grid()` function generates all possible combinations of alpha and lambda in the grid for tuning.

```
set.seed(2183)

enet_model <- train(x, y, # training dataset
  method = "glmnet",
  tuneGrid = expand.grid(alpha = seq(0, 1, length = 21),
    lambda = exp(seq(2, -2, length = 50))),
  trControl = ctrl)

# view the model cross-validation (tuning parameter selection) plot and summary
print(plot(enet_model))
```





```
summary(enet_model$finalModel)
##           Length Class      Mode
## a0           94  -none-   numeric
## beta        1692 dgCMatrx  S4
## df           94  -none-   numeric
## dim           2  -none-   numeric
## lambda        94  -none-   numeric
## dev.ratio     94  -none-   numeric
## nulldev        1  -none-   numeric
## npasses        1  -none-   numeric
## jerr           1  -none-   numeric
## offset         1  -none-   logical
## call           5  -none-   call
## nobs           1  -none-   numeric
## lambdaOpt       1  -none-   numeric
## xNames         18 -none-   character
## problemType     1  -none-   character
## tuneValue       2 data.frame list
## obsLevels       1  -none-   logical
## param           0  -none-   list

# view performance on the test set (RMSE)
test_pred <- predict(enet_model, newdata = x2) # test dataset
test_rmse <- sqrt(mean((test_pred - dat_test$recovery_time)^2))
test_rmse
## [1] 27.77103
```

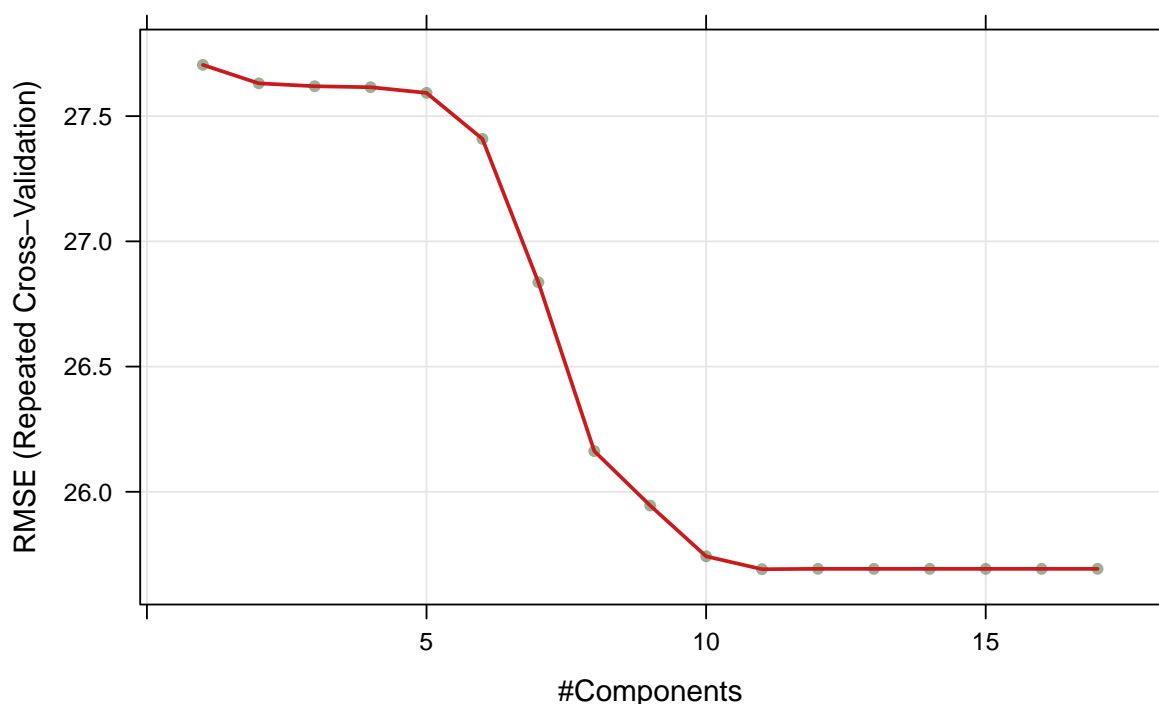
#### 4. Partial least squares (PLS) model

A PLS model is a model that seeks to find a low-dimensional representation of the predictor variables that explains the maximum variance in the response variable. It has the same assumptions as the linear model as well as a latent variable assumption, the assumption that the predictor variables are linearly related to the response variable via a set of underlying latent variables. A `method = "pls"` argument was used within the `train()` function to specify a PLS model fit. Additionally, a tuning parameter for the model were set using the `tuneGrid` argument within the `train()` function. The `tuneGrid` object includes a data frame with a single column `ncomp` that ranges from 1 to 17, representing the number of components used in the model. The `preProcess` argument is set to `"center"` and `"scale"`, which means that the training data will be centered and scaled prior to model fitting.

```
set.seed(2183)

pls_model <- train(x, y, # training dataset
  method = "pls",
  tuneGrid = data.frame(ncomp = 1:17),
  trControl = ctrl,
  preProcess = c("center", "scale"))

# view the model cross-validation (tuning parameter selection) plot and summary
print(plot(pls_model))
```



```
summary(pls_model$finalModel)
## Data:      X dimension: 1402 18
## Y dimension: 1402 1
## Fit method: oscorespls
## Number of components considered: 11
## TRAINING: % variance explained
```

```
##           1 comps  2 comps  3 comps  4 comps  5 comps  6 comps  7 comps
## X           9.853   17.56   27.14   32.78   38.53   43.47   47.92
## .outcome    15.537   16.50   16.65   16.80   17.26   19.02   23.30
##           8 comps  9 comps 10 comps 11 comps
## X           51.14   56.12   60.94   63.06
## .outcome    28.21   28.94   29.02   29.08

# view performance on the test set (RMSE)
test_pred <- predict(pls_model, newdata = x2) # test dataset
test_rmse <- sqrt(mean((test_pred - dat_test$recovery_time)^2))
test_rmse
## [1] 27.34478
```

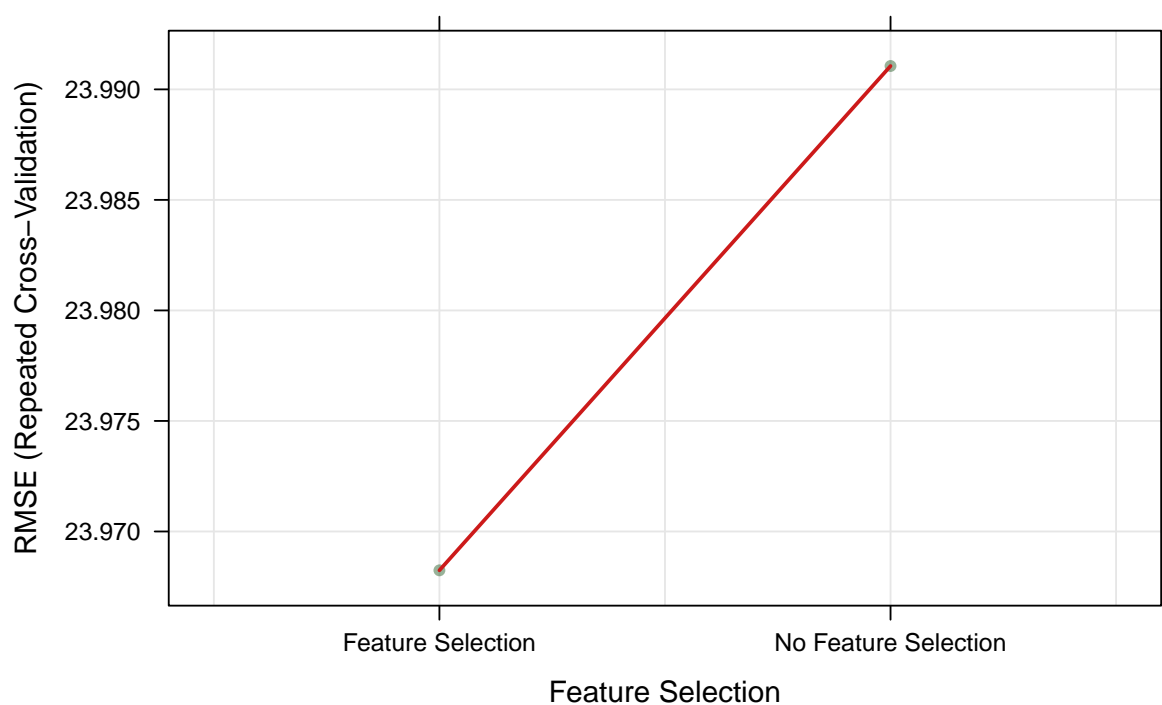
## 5. Generalized additive model (GAM)

A GAM model is a model that extends the linear model by allowing for nonlinear relationships between the predictor and response variables, using flexible functions called splines. It has the same assumptions as the linear model. A `method = "gam"` argument was used within the `train()` function to specify a GAM fit.

```
set.seed(2183)

# fit GAM model using all predictors
gam_model <- train(x, y, # training dataset
                  method = "gam",
                  trControl = ctrl,
                  control = gam.control(maxit = 200)) # Adjusted due to failure to converge at default s

# view the model cross-validation (tuning parameter selection) plot and summary
print(plot(gam_model))
```



```
summary(gam_model$finalModel)
##
## Family: gaussian
## Link function: identity
##
## Formula:
## .outcome ~ gender + race2 + race3 + race4 + smoking1 + smoking2 +
##   hypertension + diabetes + vaccine + severity + studyB + studyC +
##   s(age) + s(SBP) + s(LDL) + s(bmi) + s(height) + s(weight)
##
## Parametric coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  43.96343    1.93341   22.739  < 2e-16 ***
## gender       -7.51837    1.25233   -6.003  2.47e-09 ***
## race2         1.58941    2.62425    0.606  0.544841
## race3        -0.69381    1.59081   -0.436  0.662810
## race4        -1.95092    2.16143   -0.903  0.366895
## smoking1      4.64945    1.39753    3.327  0.000902 ***
## smoking2      8.73771    2.12834    4.105  4.28e-05 ***
## hypertension  4.43411    1.34730    3.291  0.001023 **
## diabetes     -0.65604    1.78312   -0.368  0.712993
## vaccine      -7.80343    1.26842   -6.152  1.00e-09 ***
## severity      6.61473    2.15179    3.074  0.002154 **
## studyB        4.44936    1.60182    2.778  0.005550 **
## studyC       -0.08194    1.97615   -0.041  0.966932
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Approximate significance of smooth terms:
##           edf Ref.df      F p-value
## s(age)      4.509e+00      9  1.554  0.00624 **
## s(SBP)      1.504e-06      9  0.000  0.41642
## s(LDL)      4.897e-07      9  0.000  0.69971
## s(bmi)      6.064e+00      9 50.348 < 2e-16 ***
## s(height)   8.241e+00      9  4.700 4.58e-07 ***
## s(weight)   4.979e+00      9  1.436  0.00783 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.427   Deviance explained = 44.2%
## GCV = 548.83   Scale est. = 534.42      n = 1402

# view performance on the test set (RMSE)
test_pred <- predict(gam_model, newdata = x2) # test dataset
test_rmse <- sqrt(mean((test_pred - dat_test$recovery_time)^2))
test_rmse
## [1] 25.45068
```

## 6. Multivariate adaptive regression spline (MARS) model

A MARS model is a model that uses piecewise linear or nonlinear functions to model the relationship between the predictor and response variables, which can capture complex nonlinear relationships. It has the same assumptions as the linear model. A method = “earth” argument was used within the `train()` function to specify a MARS model fit. Additionally, the `expand.grid()` function is used to generate a grid of tuning parameters. The `mars_grid` object includes two arguments: degree is set to 1, 2, and 3, representing the number of possible product hinge functions in a single term, and `nprune` is set to integers between 2 and 17, representing the upper bound on the number of terms in the model. The `tuneGrid` argument in the `train()` function uses the `mars_grid` object to specify the parameters for model tuning.

```
# create dummy variables for categorical variables
df_dummies <- data.frame(model.matrix(~ . - 1, data = dat[, c("gender", "race", "smoking", "hypertension",
  age = dat$age,
  height = dat$height,
  weight = dat$weight,
  bmi = dat$bmi,
  SBP = dat$SBP,
  LDL = dat$LDL,
  recovery_time = dat$recovery_time) # add continuous variables back to the data

# rename df_dummies dataset as dat
dat_mars <- df_dummies

# training data
dat_train_mars <- dat_mars[trRows, ]
## matrix of predictors
x_mars <- model.matrix(recovery_time~., dat_mars)[trRows, -1]
## vector of response
y_mars <- dat_mars$recovery_time[trRows]

# test data
dat_test_mars <- dat_mars[-trRows, ]
## matrix of predictors
```

```

x2_mars <- model.matrix(recovery_time~.,dat_mars)[-trRows,-1]
## vector of response
y2_mars <- dat_mars$recovery_time[-trRows]

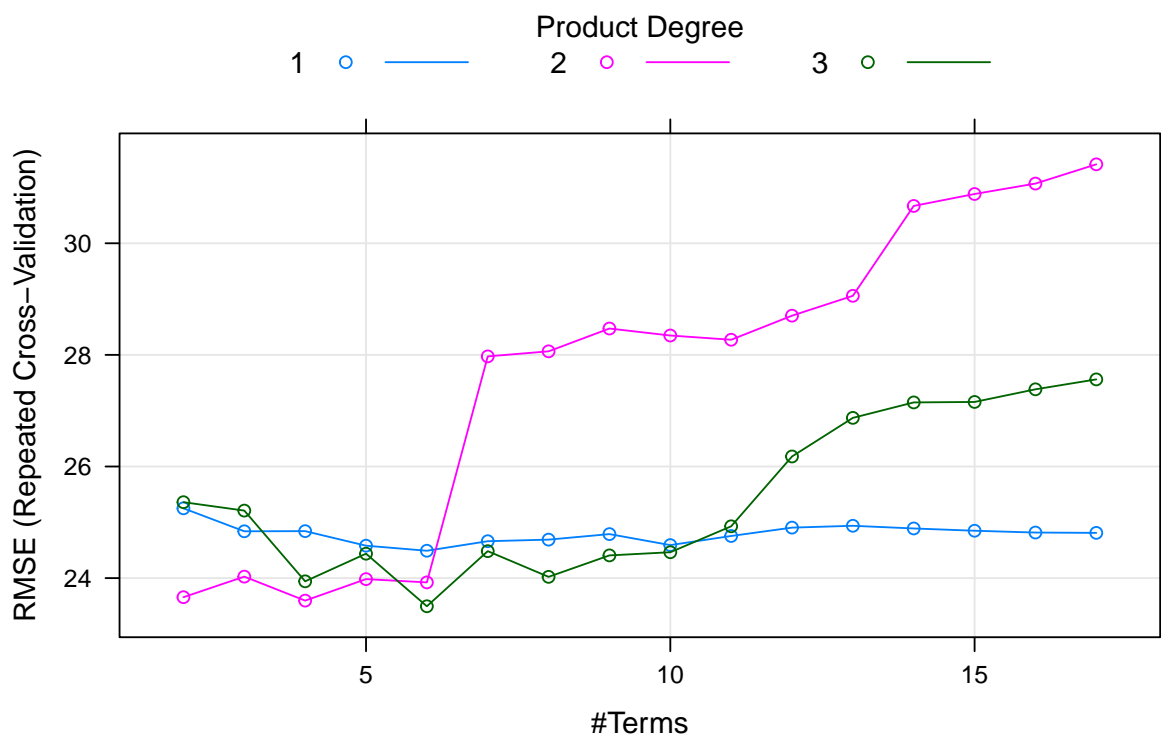
# create grid of all possible pairs that can take degree and nprune values
mars_grid <- expand.grid(degree = 1:3, # number of possible product hinge functions in 1 term
                        nprune = 2:17) # upper bound of number of terms in model

set.seed(2183)

mars_model <- train(x_mars, y_mars, # training dataset
                   method = "earth",
                   tuneGrid = mars_grid,
                   trControl = ctrl)

# view the model cross-validation (tuning parameter selection) plot and summary
print(plot(mars_model))

```



```

summary(mars_model$finalModel)
## Call: earth(x=matrix[1402,19], y=c(49,38,46,53,3...), keepxy=TRUE, degree=3,
##           nprune=6)
##
##
##               coefficients
## (Intercept)      -28.923070
## gender           -7.655526
## h(bmi-23.4)       9.511406
## h(30.7-bmi)       9.590069

```

```
## studyB * h(bmi-30.7)          17.390398
## smoking1 * studyB * h(bmi-30.7) 16.474027
##
## Selected 6 of 27 terms, and 4 of 19 predictors (nprune=6)
## Termination condition: Reached nk 39
## Importance: studyB, bmi, smoking1, gender, race1-unused, race2-unused, ...
## Number of terms at each degree of interaction: 1 3 1 1
## GCV 491.2801    RSS 675574    GRSq 0.47361    RSq 0.4829613

# view performance on the test set (RMSE)
test_pred <- predict(mars_model, newdata = x2_mars) # test dataset
test_rmse <- sqrt(mean((test_pred - dat_test_mars$recovery_time)^2))
test_rmse
## [1] 23.63058
```

## Model comparison

The models were assessed by comparing the RMSE for each model. Low RMSE indicates the best performing model (low prediction error), while high RMSE indicates the worst performing model (high prediction error).

```
set.seed(2183)

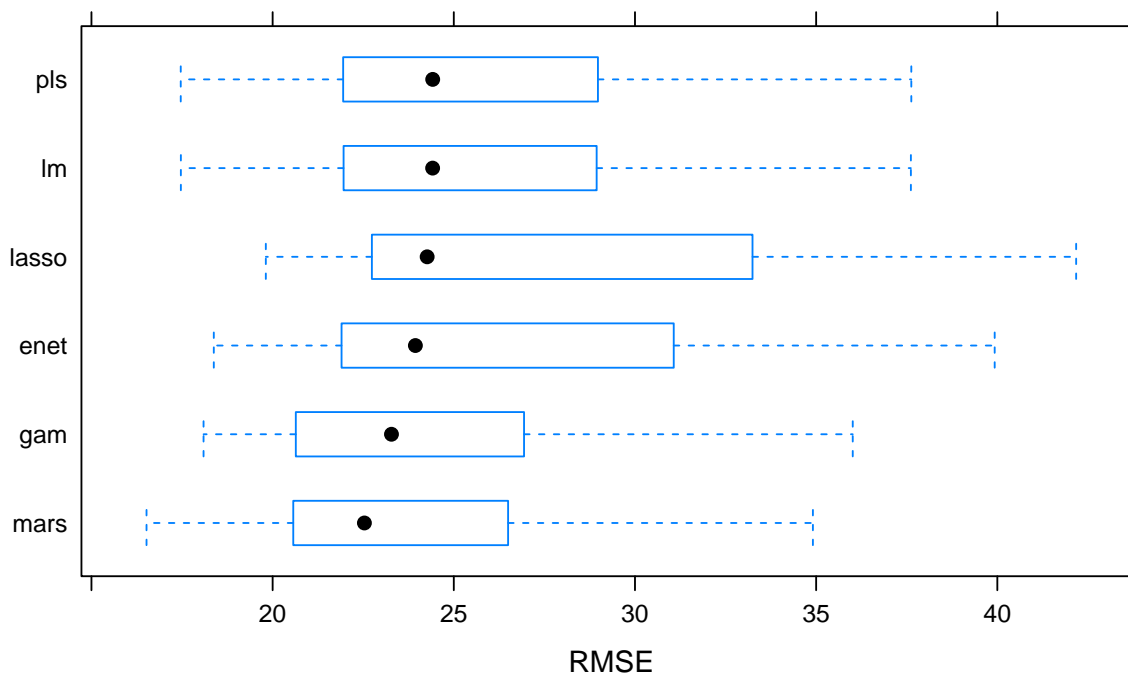
resamp <- resamples(list(
  lm = linear_model,
  lasso = lasso_model,
  enet = enet_model,
  pls = pls_model,
  gam = gam_model,
  mars = mars_model
))

summary(resamp)
##
## Call:
## summary.resamples(object = resamp)
##
## Models: lm, lasso, enet, pls, gam, mars
## Number of resamples: 50
##
## MAE
##           Min.   1st Qu.   Median     Mean   3rd Qu.     Max. NA's
## lm      13.06852 15.82773 16.65854 16.68009 17.52988 19.53885    0
## lasso   14.25843 15.92356 16.50435 16.89763 17.65061 20.10554    0
## enet    13.56431 15.46643 16.26234 16.31209 17.07245 19.30336    0
## pls     13.06709 15.77555 16.68261 16.67911 17.53197 19.49819    0
## gam     12.93158 15.19322 16.00336 15.87566 16.72480 19.49579    0
## mars    12.66352 14.32980 15.19281 15.40187 16.04710 18.87256    0
##
## RMSE
##           Min.   1st Qu.   Median     Mean   3rd Qu.     Max. NA's
## lm      17.46778 21.95953 24.41580 25.69236 28.94167 37.61437    0
## lasso   19.81271 22.77017 24.26439 27.48116 33.14234 42.17530    0
## enet    18.37685 21.91826 23.93991 26.11100 30.98068 39.92701    0
## pls     17.46405 21.95236 24.41786 25.69095 28.97255 37.62563    0
```

```
## gam 18.09395 20.67294 23.27995 23.96824 26.78507 36.01007 0
## mars 16.51977 20.58231 22.53261 23.49522 26.34327 34.91041 0
##
## Rsquared
##           Min.   1st Qu.   Median     Mean   3rd Qu.     Max. NA's
## lm      0.04722889 0.1947930 0.2788224 0.2727476 0.3461258 0.4918482 0
## lasso  0.01741077 0.1286390 0.1719990 0.1642228 0.2141429 0.2929154 0
## enet   0.04074479 0.1996787 0.2583386 0.2451645 0.3170416 0.4390000 0
## pls    0.04791543 0.1950316 0.2791061 0.2727537 0.3451321 0.4920658 0
## gam    0.13579308 0.2730556 0.3540842 0.3704064 0.4406017 0.6293796 0
## mars   0.03523834 0.2388844 0.3894285 0.3901751 0.5168117 0.7480162 0

bwplot(resamp,
       metric = "RMSE",
       main = "Figure 2. Model Comparison Plot Using RMSE")
```

**Figure 2. Model Comparison Plot Using RMSE**



According to Figure 2, the PLS model had the highest median training RMSE (i.e., worst performance), followed by the linear model, the lasso model, the elastic net model, the GAM model, and finally, The MARS model, which had the lowest median training RMSE (i.e., best performance).

The final model for predicting time to recovery from COVID-19 was selected by comparing the median training RMSE for all models created. The MARS model had the lowest value for mean and median RMSE out of all models, however, it only contained 4 of the predictors. A model with so few predictors would not be an ideal model for predicting recovery time and comprehensively identifying important risk factors for long recovery time, as so few risk factors were included. Therefore, the GAM model—the second best model in terms of mean and median RMSE—was selected as the final model for predicting time to recovery from COVID-19 in this study.



## Results

### Model formula

```
# view the model summary for the GAM model using all predictors
summary(gam_model$finalModel)
##
## Family: gaussian
## Link function: identity
##
## Formula:
## .outcome ~ gender + race2 + race3 + race4 + smoking1 + smoking2 +
##      hypertension + diabetes + vaccine + severity + studyB + studyC +
##      s(age) + s(SBP) + s(LDL) + s(bmi) + s(height) + s(weight)
##
## Parametric coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  43.96343    1.93341  22.739  < 2e-16 ***
## gender       -7.51837    1.25233  -6.003  2.47e-09 ***
## race2         1.58941    2.62425   0.606  0.544841
## race3        -0.69381    1.59081  -0.436  0.662810
## race4        -1.95092    2.16143  -0.903  0.366895
## smoking1      4.64945    1.39753   3.327  0.000902 ***
## smoking2      8.73771    2.12834   4.105  4.28e-05 ***
## hypertension  4.43411    1.34730   3.291  0.001023 **
## diabetes     -0.65604    1.78312  -0.368  0.712993
## vaccine       -7.80343    1.26842  -6.152  1.00e-09 ***
## severity      6.61473    2.15179   3.074  0.002154 **
## studyB        4.44936    1.60182   2.778  0.005550 **
## studyC       -0.08194    1.97615  -0.041  0.966932
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##              edf Ref.df      F  p-value
## s(age)       4.509e+00     9  1.554  0.00624 **
## s(SBP)       1.504e-06     9  0.000  0.41642
## s(LDL)       4.897e-07     9  0.000  0.69971
## s(bmi)       6.064e+00     9 50.348  < 2e-16 ***
## s(height)    8.241e+00     9  4.700  4.58e-07 ***
## s(weight)    4.979e+00     9  1.436  0.00783 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.427   Deviance explained = 44.2%
## GCV = 548.83   Scale est. = 534.42     n = 1402
```

The final GAM model formula uses the `recovery_time` variable as the outcome (Y), and the following terms as predictors. Note that “White” (`race = 0`) was set as the reference category for the `race` variable, “Never Smoker” (`smoking = 0`) was set as the reference category for the `smoking` variable, and “Study A” (`study = A`) was set as the reference category for the `study` variable. An `s()` around the variable name indicates that a smoothing function was applied to the variable. An asterisk (\*) next to the term indicates that the term was statistically significant at the 5% level of significance.

- `gender*`
- Race: Asian (`race2`)
- Race: Black (`race3`)
- Race: Hispanic (`race4`)
- Smoking: former smoker (`smoking1`)\*
- Smoking: current smoker (`smoking2`)\*
- `hypertension*`
- `diabetes`
- `vaccine*`
- `severity*`
- Study: B (`studyB`)\*
- Study: C (`studyC`)
- `s(age)*`
- `s(SBP)`
- `s(LDL)`
- `s(bmi)*`
- `s(height)*`
- `s(weight)*`

### Model interpretation

- On average, being male is associated with a statistically significantly shorter predicted COVID-19 recovery time compared to being female. More specifically, holding all else constant, being male is associated with a 7.51837-day shorter predicted recovery time than being female.
- On average, being a former smoker is associated with a statistically significantly longer predicted COVID-19 recovery time compared to being a never smoker. More specifically, holding all else constant, former smoking is associated with a 4.64945-day longer predicted recovery time than never smoking.
- On average, being a current smoker is associated with a statistically significantly longer predicted COVID-19 recovery time compared to being a never smoker. More specifically, holding all else constant, current smoking is associated with a 8.73771-day longer predicted recovery time than never smoking.
- On average, having hypertension is associated with a statistically significantly longer predicted COVID-19 recovery time compared to not having hypertension. More specifically, holding all else constant, a hypertension diagnosis is associated with a 4.43411-day longer predicted recovery time than no hypertension diagnosis.
- On average, being vaccinated is associated with a statistically significantly shorter predicted COVID-19 recovery time compared to not being vaccinated. More specifically, holding all else constant, vaccination is associated with a 7.80343-day shorter predicted recovery time than no vaccination.
- On average, severe COVID-19 infection is associated with a statistically significantly longer predicted COVID-19 recovery time compared to non severe infection. More specifically, holding all else constant, infection severity is associated with a 6.61473-day longer predicted recovery time than no infection severity.
- On average, being in study B is associated with a statistically significantly longer predicted COVID-19 recovery time compared to study A. More specifically, holding all else constant, study B is associated with a 4.44936-day longer predicted recovery time than study A.

The GAM model also includes several significant smooth terms to capture the nonlinear relationships between the predictors and COVID-19 recovery time. There is a statistically significant relationship between age, BMI, weight, and height, and predicted COVID-19 recovery time.

### Model performance

```
# print the GAM model training RMSE
gam_model$results$RMSE
## [1] 23.99106 23.96824
```

```
set.seed(2183)

# calculate and print the GAM model test RMSE
test_pred <- predict(gam_model, newdata = x2)
test_rmse <- sqrt(mean((test_pred - dat_test$recovery_time)^2))
test_rmse
## [1] 25.45068
```

All of these factors together explain 44.2% of the deviance in COVID-19 recovery time. Additionally, the model's training error (RMSE using the training dataset) of about 24.0 (rounded) indicates that, on average, the model's predictions on the training data deviate from the actual values by 24 units. Meanwhile, the test error (RMSE using the test dataset) of 25.5 (rounded) suggests that the model's performance on unseen data is slightly worse than its performance on the training data.

## Conclusion

The final GAM model found that several factors were statistically significant in predicting time to recovery from COVID-19. On average, having a history of former or current smoking, having hypertension, and experiencing severe COVID-19 infection were all significantly associated with a longer predicted recovery time. Additionally, being in study B was associated with a longer recovery time compared to being in study A. On the other hand, being male and being vaccinated was associated with a shorter recovery time. Finally, age, BMI, height, and weight are also significantly associated with predicted COVID-19 recovery time. The model did not find significant associations with race or diabetes. These insights can be useful in predicting recovery time and informing clinical decisions in the management of COVID-19 patients.

## Appendix

Predictor plots for the final GAM model:

```
set.seed(2183)

# formula based on final GAM model
gam_final_model <- gam(recovery_time ~ gender + race + smoking +
                      hypertension + diabetes + vaccine + severity + study +
                      s(age) + s(SBP) + s(LDL) + s(bmi) + s(height) +
                      s(weight),
                      data = dat_train) # training dataset

plot(gam_final_model)
```

