

# Simple document

```
library(tidyverse)
library(caret)
library(visdat)
library(corrplot)
library(AppliedPredictiveModeling)
library(pROC)
library(rpart.plot)
library(vip)
library(ranger)
library(tidytext)
library(pdp)
library(lime)

ctrl <- trainControl(method = "cv",
                      summaryFunction = twoClassSummary,
                      classProbs = TRUE)

knitr::opts_chunk$set(
  fig.width = 6,
  out.width = "80%",
  fig.align = "center"
)
```

## Data pre-process

```
# Import data
dat_raw <- read.csv("airline.csv")

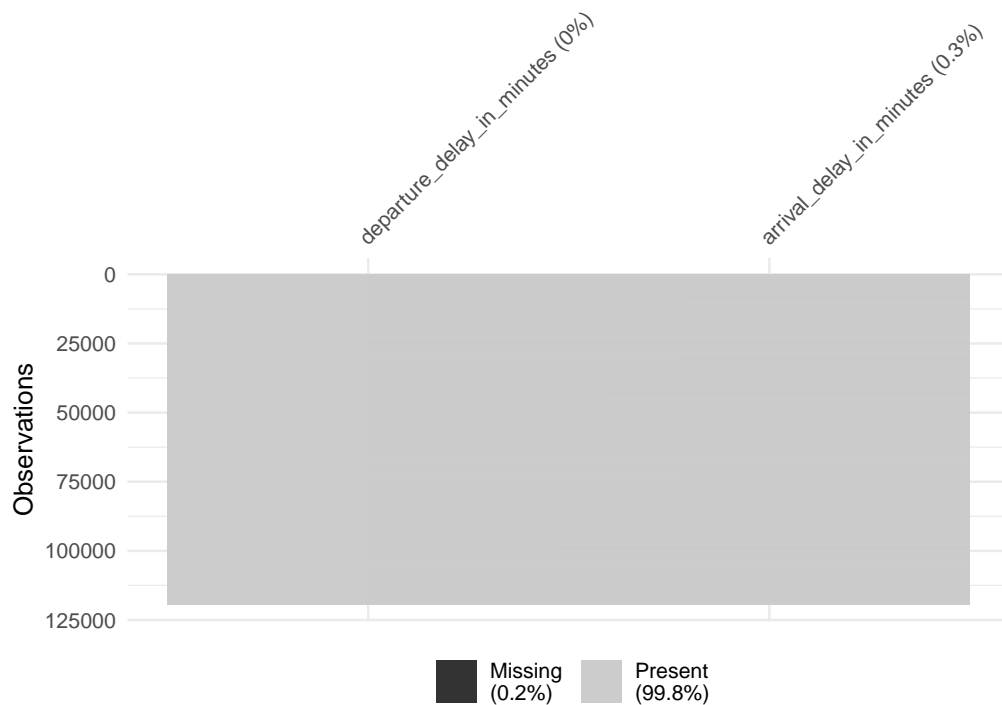
# data clean
dat <- dat_raw %>%
  janitor::clean_names() %>%
  select(-1) %>%
  mutate(satisfaction = recode(satisfaction,
                                "satisfied" = "yes",
                                "neutral or dissatisfied" = "no")) %>%
  filter_at(vars(7:20), all_vars(. > 0.5))

# Check missing value
sapply(dat, function(x) sum(is.na(x)))
```

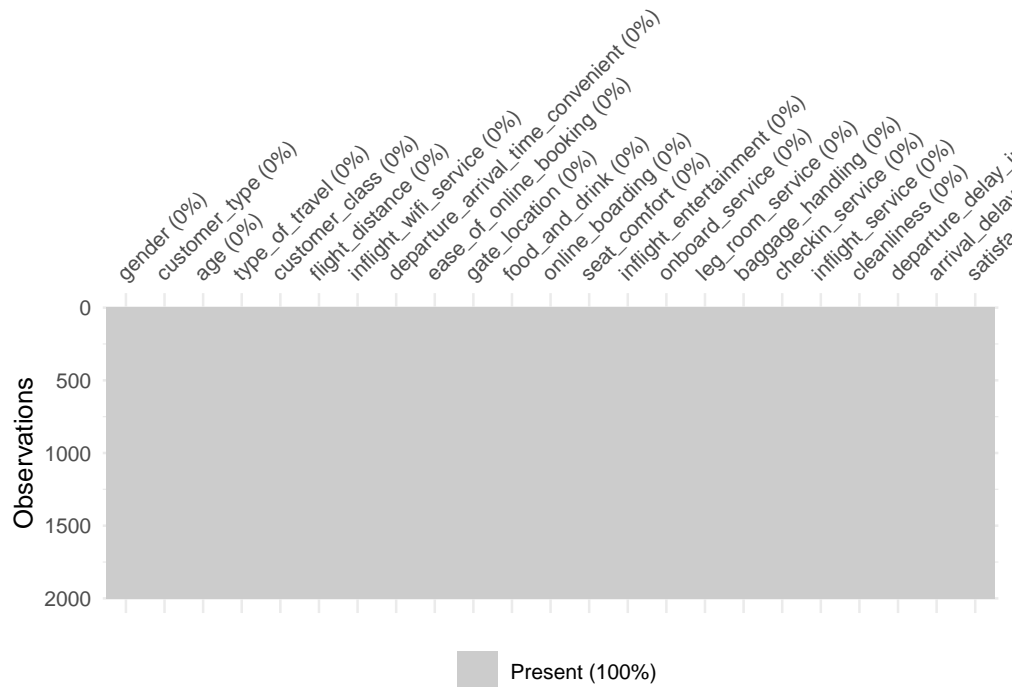
```
##                gender                customer_type
##                  0                          0
##                age                type_of_travel
##                  0                          0
##      customer_class                flight_distance
##                  0                          0
```

```
##          inflight_wifi_service departure_arrival_time_convenient
##                                0                                0
##          ease_of_online_booking          gate_location
##                                0                                0
##          food_and_drink          online_boarding
##                                0                                0
##          seat_comfort          inflight_entertainment
##                                0                                0
##          onboard_service          leg_room_service
##                                0                                0
##          baggage_handling          checkin_service
##                                0                                0
##          inflight_service          cleanliness
##                                0                                0
##          departure_delay_in_minutes          arrival_delay_in_minutes
##                                0                                363
##          satisfaction
##                                0
```

```
# deal with missing values
deal_mis <- dat[, 21:22]
bagImp = preProcess(deal_mis, method = "bagImpute")
dat = predict(bagImp, dat)
vis_miss(deal_mis)
```



```
# sample data
set.seed(1234)
dat <- dat[sample(1:nrow(dat), 2000, replace = FALSE), ]
vis_miss(dat) ## check
```



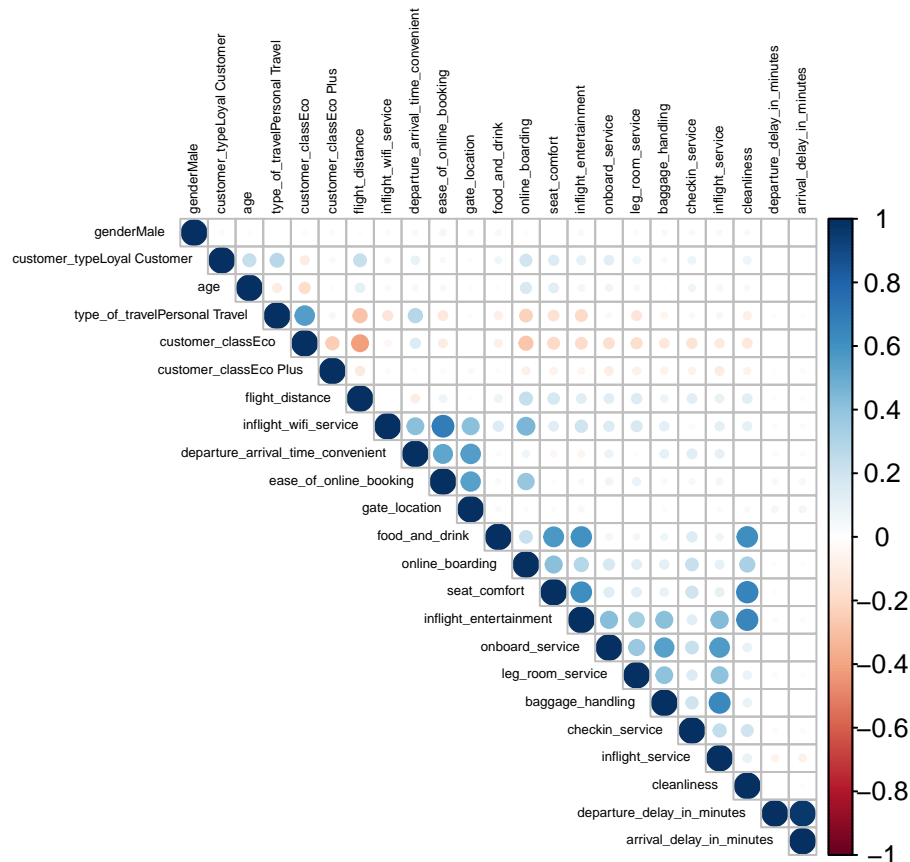
```
# --- Split data ---
set.seed(1234)
trRow <- createDataPartition(dat$satisfaction, p = 0.8, list = F)

# Train data
train <- dat[trRow, ]
x_train <- model.matrix(satisfaction ~., train)[,-1]
y_train <- train$satisfaction

# Test data
test <- dat[-trRow, ]
x_test <- model.matrix(satisfaction ~., test)[,-1]
y_test <- test$satisfaction
```

## EDA

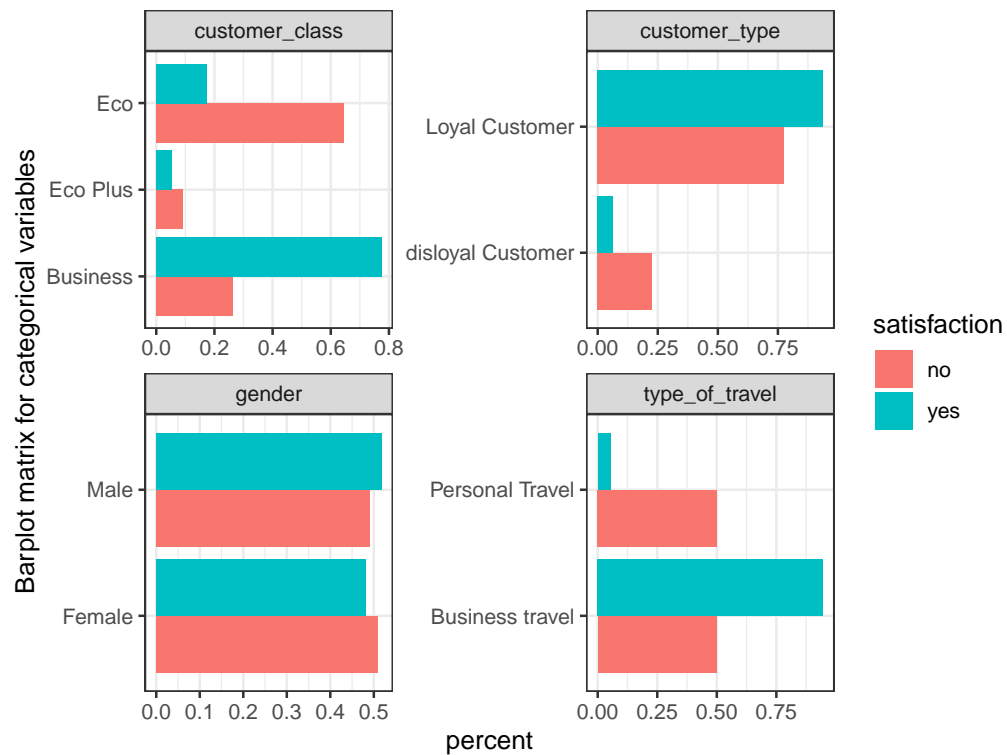
```
# Correlation plot
corrplot(cor(x_train),
  method = "circle",
  type = "upper",
  tl.col = "black",
  tl.cex = 0.4)
```



```
# Barplot matrix for categorical variables
```

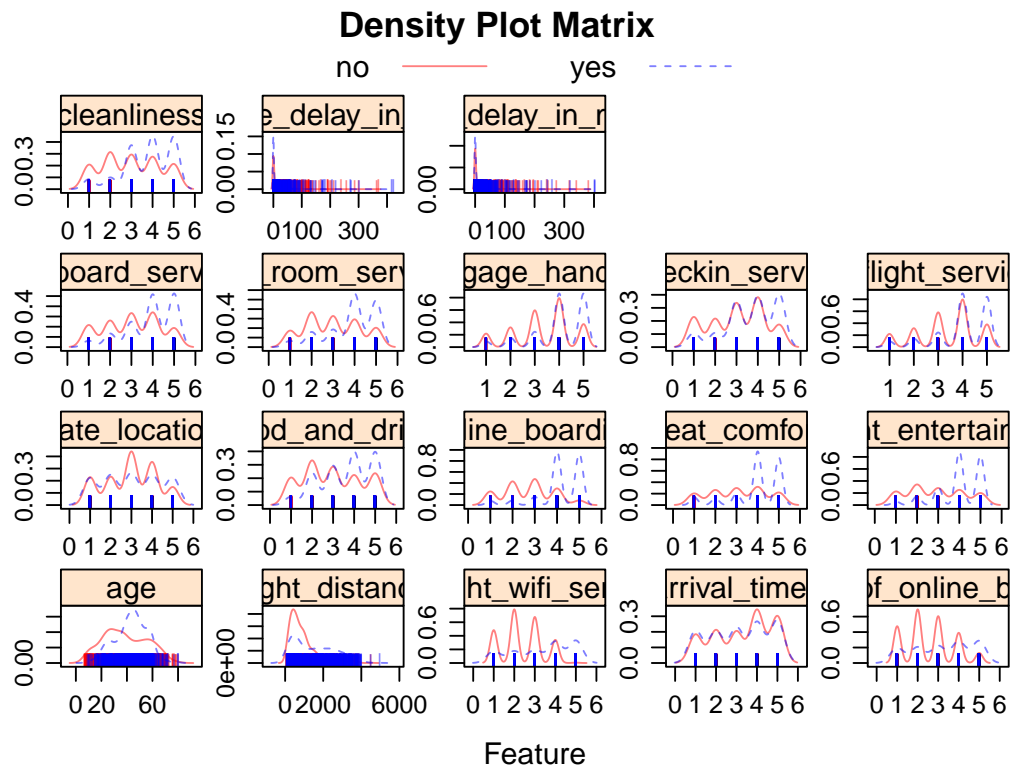
```
train %>%
  select(1:2, 4:5, 23) %>%
  pivot_longer(-5,
    names_to = "variable",
    values_to = "value") %>%
  group_by(variable, value, satisfaction) %>%
  summarize(num = n()) %>%
  ungroup() %>%
  group_by(variable, satisfaction) %>%
  mutate(percent = num / sum(num),
    indicator = case_when(value == "Eco" ~ 3,
      value == "Eco Plus" ~ 2,
      value == "Business" ~ 1,
      TRUE ~ 0)) %>%
  ggplot(aes(x = reorder_within(value, indicator, variable),
    y = percent, fill = satisfaction)) +
  geom_col(position = "dodge") +
  xlab("Barplot matrix for categorical variables") +
  coord_flip() +
  scale_x_reordered() +
  facet_wrap(~ variable, scales = "free") + theme_bw()
```

```
## `summarise()` has grouped output by 'variable', 'value'. You can override using the `.groups` argument
```



```
# Density plot matrix
theme1 <- transparentTheme(trans = .5)
trellis.par.set(theme1)

plt_feature <-
  featurePlot(x = x_train[, c(-1, -2, -4, -5, -6)],
             y = as.factor(y_train),
             plot = "density",
             scales = list(x = list(relation = "free"),
                           y = list(relation = "free")),
             pch = "|", auto.key = list(columns = 2))
update(plt_feature, main = "Density Plot Matrix")
```



## Model fitting

### Logistic regression

```
set.seed(1234)
model.glm <- train(x = x_train,
                   y = y_train,
                   method = "glm",
                   metric = "ROC",
                   trControl = ctrl)

# Test AUC and Misclassification error rate
pred_glm_auc <- predict(model.glm, newdata = x_test, type = "prob")[,2]
roc(y_test, pred_glm_auc)$auc[1]
```

```
## Setting levels: control = no, case = yes
## Setting direction: controls < cases
## [1] 0.9550372
```

```
# Confusion matrix
pred_glm <- predict(model.glm, newdata = x_test)
confusionMatrix(data = as.factor(pred_glm),
                reference = as.factor(y_test))
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction  no  yes
```

```
##      no  188  23
##      yes  28 161
##
##              Accuracy : 0.8725
##              95% CI : (0.8358, 0.9036)
##      No Information Rate : 0.54
##      P-Value [Acc > NIR] : <2e-16
##
##              Kappa : 0.7439
##
##      McNemar's Test P-Value : 0.5754
##
##      Sensitivity : 0.8704
##      Specificity : 0.8750
##      Pos Pred Value : 0.8910
##      Neg Pred Value : 0.8519
##      Prevalence : 0.5400
##      Detection Rate : 0.4700
##      Detection Prevalence : 0.5275
##      Balanced Accuracy : 0.8727
##
##      'Positive' Class : no
##
```

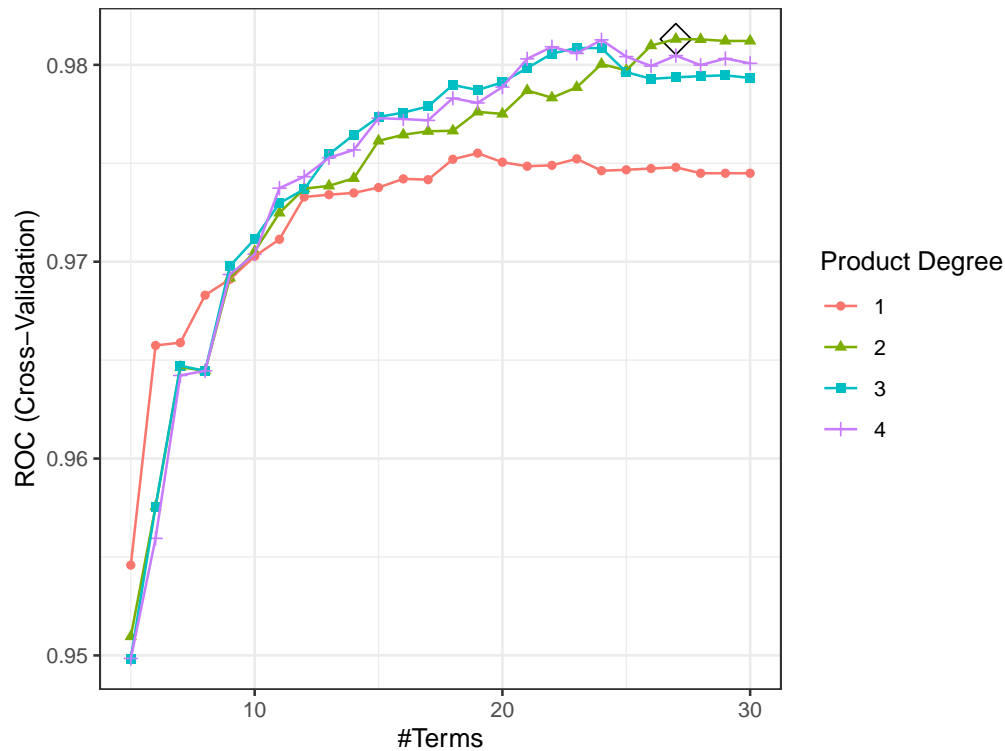
## MARS

```
set.seed(1234)
model.mars <- train(x = x_train,
                    y = y_train,
                    method = "earth",
                    tuneGrid = expand.grid(degree = 1:4,
                                           nprune = 5:30),
                    metric = "ROC",
                    trControl = ctrl)

model.mars$bestTune

##      nprune degree
## 49      27      2

ggplot(model.mars, highlight = T) +
  theme_bw()
```



```
## test auc and misclassification error rate
pred_mars_auc <- predict(model.mars, newdata = x_test, type = "prob")[,2]
roc(y_test, pred_mars_auc)$auc[1]
```

```
## [1] 0.9742728
```

```
pred_mars <- predict(model.mars, newdata = x_test)
pred.miserror_mars <- 1 - mean(pred_mars == y_test)
pred.miserror_mars
```

```
## [1] 0.0875
```

```
confusionMatrix(data = as.factor(pred_mars),
                 reference = as.factor(y_test))
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction no yes
```

```
##           no  201  20
```

```
##           yes   15 164
```

```
##
```

```
##           Accuracy : 0.9125
```

```
##           95% CI : (0.8804, 0.9383)
```

```
##           No Information Rate : 0.54
```

```
##           P-Value [Acc > NIR] : <2e-16
```

```
##
```

```
##           Kappa : 0.8235
```

```
##
```

```
##           McNemar's Test P-Value : 0.499
```

```
##
```

```
##           Sensitivity : 0.9306
```



```
##           Specificity : 0.8913
##           Pos Pred Value : 0.9095
##           Neg Pred Value : 0.9162
##           Prevalence : 0.5400
##           Detection Rate : 0.5025
##           Detection Prevalence : 0.5525
##           Balanced Accuracy : 0.9109
##
##           'Positive' Class : no
##
```

## LDA

```
set.seed(1234)
model.lda <- train(x = x_train,
                   y = y_train,
                   method = "lda",
                   metric = "ROC",
                   trControl = ctrl)

## test auc and misclassification error rate
pred_lda_auc <- predict(model.lda, newdata = x_test, type = "prob")[,2]
roc(y_test, pred_lda_auc)$auc[1]
```

```
## Setting levels: control = no, case = yes
```

```
## Setting direction: controls < cases
```

```
## [1] 0.953125
```

```
pred_lda <- predict(model.lda, newdata = x_test)
pred.miserror_lda <- 1 - mean(pred_lda == y_test)
pred.miserror_lda
```

```
## [1] 0.1275
```

```
confusionMatrix(data = as.factor(pred_lda),
                 reference = as.factor(y_test))
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction  no yes
```

```
##           no 187 22
```

```
##           yes 29 162
```

```
##
```

```
##           Accuracy : 0.8725
```

```
##           95% CI : (0.8358, 0.9036)
```

```
##           No Information Rate : 0.54
```

```
##           P-Value [Acc > NIR] : <2e-16
```

```
##
```

```
##           Kappa : 0.7441
```

```
##
```

```
##           McNemar's Test P-Value : 0.4008
```

```
##
```

```
##           Sensitivity : 0.8657
```

```
##          Specificity : 0.8804
##          Pos Pred Value : 0.8947
##          Neg Pred Value : 0.8482
##          Prevalence : 0.5400
##          Detection Rate : 0.4675
##          Detection Prevalence : 0.5225
##          Balanced Accuracy : 0.8731
##
##          'Positive' Class : no
##
```

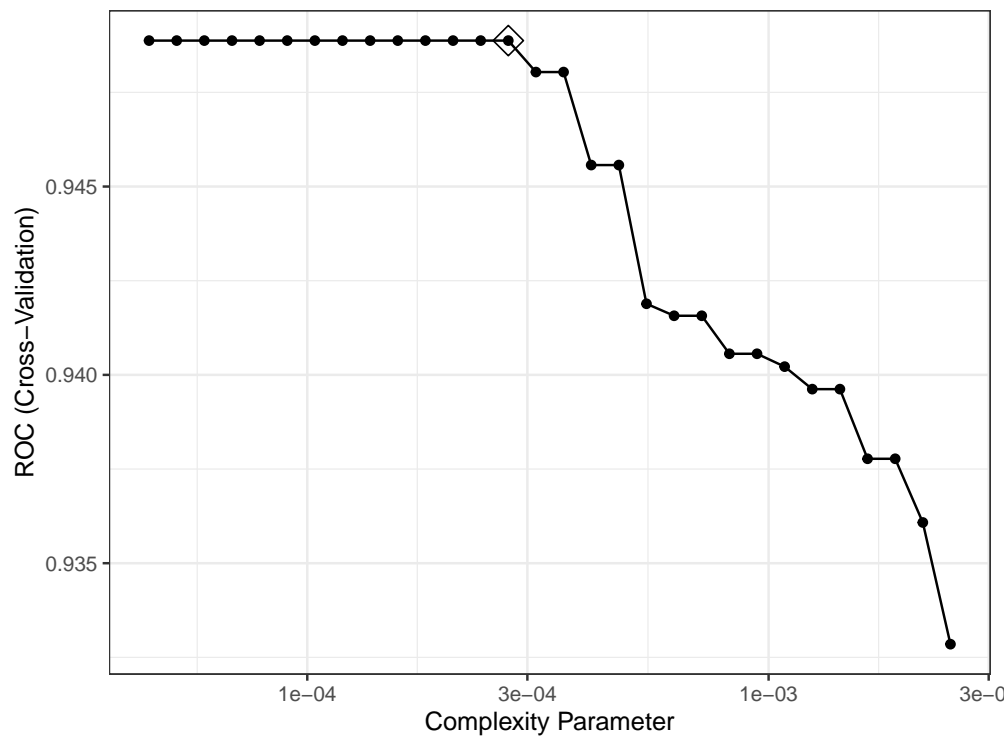
## Classification Tree

```
set.seed(1234)
model.tree <- train(x_train,
  y_train,
  method = "rpart",
  tuneGrid = data.frame(cp = exp(seq(-10, -6, length = 30))),
  trControl = ctrl,
  metric = "ROC")

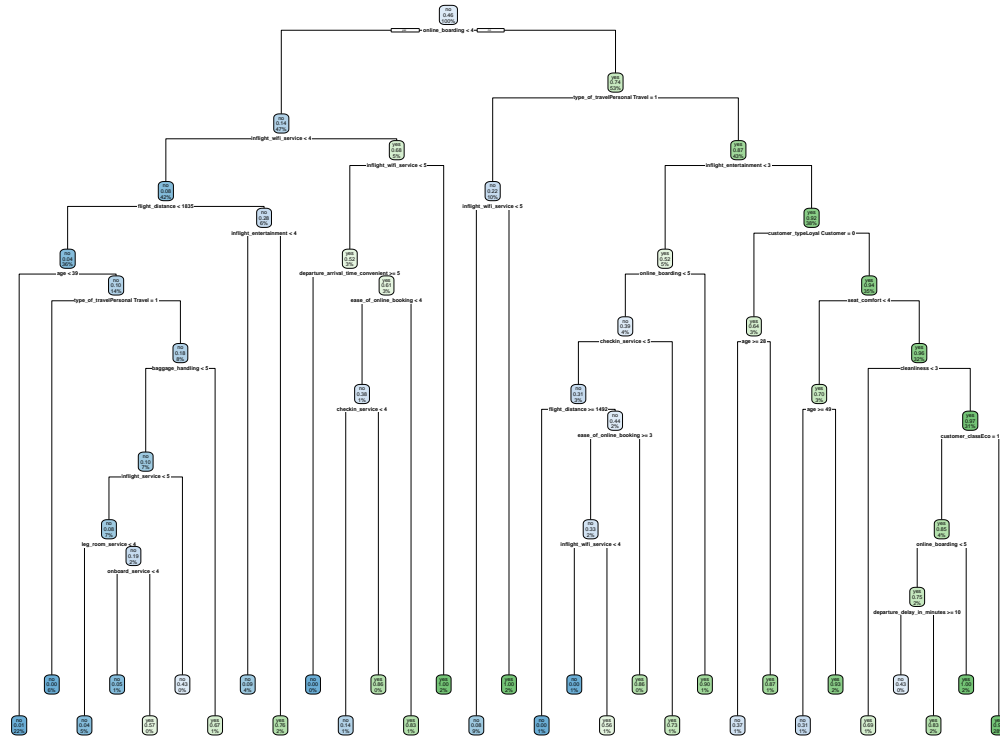
model.tree$bestTune
```

```
##          cp
## 14 0.0002727659

ggplot(model.tree, highlight = TRUE) +
  scale_x_continuous(trans = scales::log_trans(),
    breaks = scales::log_breaks()) +
  theme_bw()
```



```
rpart.plot(model.tree$finalModel)
```



```
## test auc and misclassification error rate
```

```
pred_tree_auc <- predict(model.tree, newdata = x_test, type = "prob")[,2]
roc(y_test, pred_tree_auc)$auc[1]
```

```
## [1] 0.9350594
```

```
pred_tree <- predict(model.tree, newdata = x_test)
pred.miserror_tree <- 1 - mean(pred_tree == y_test)
pred.miserror_tree
```

```
## [1] 0.1275
```

```
confusionMatrix(data = as.factor(pred_tree),
                 reference = as.factor(y_test))
```

```
## Confusion Matrix and Statistics
```

## ##

| ## | Reference |
|----|-----------|
|----|-----------|

```
## Prediction  no yes
```

```
##          no  189  24
```

```
##          yes  27 160
```

## ##

```
## Accuracy : 0.8725
```

```
##          95% CI : (0.8358, 0.9036)
```

```
##      No Information Rate : 0.54
```

```
##      P-Value [Acc > NIR] : <2e-16
```

##

```
##          Kappa : 0.7437
```

##

```
## McNemar's Test P-Value : 0.7794
```

```
##
##      Sensitivity : 0.8750
##      Specificity : 0.8696
##      Pos Pred Value : 0.8873
##      Neg Pred Value : 0.8556
##      Prevalence : 0.5400
##      Detection Rate : 0.4725
##      Detection Prevalence : 0.5325
##      Balanced Accuracy : 0.8723
##
##      'Positive' Class : no
##
```

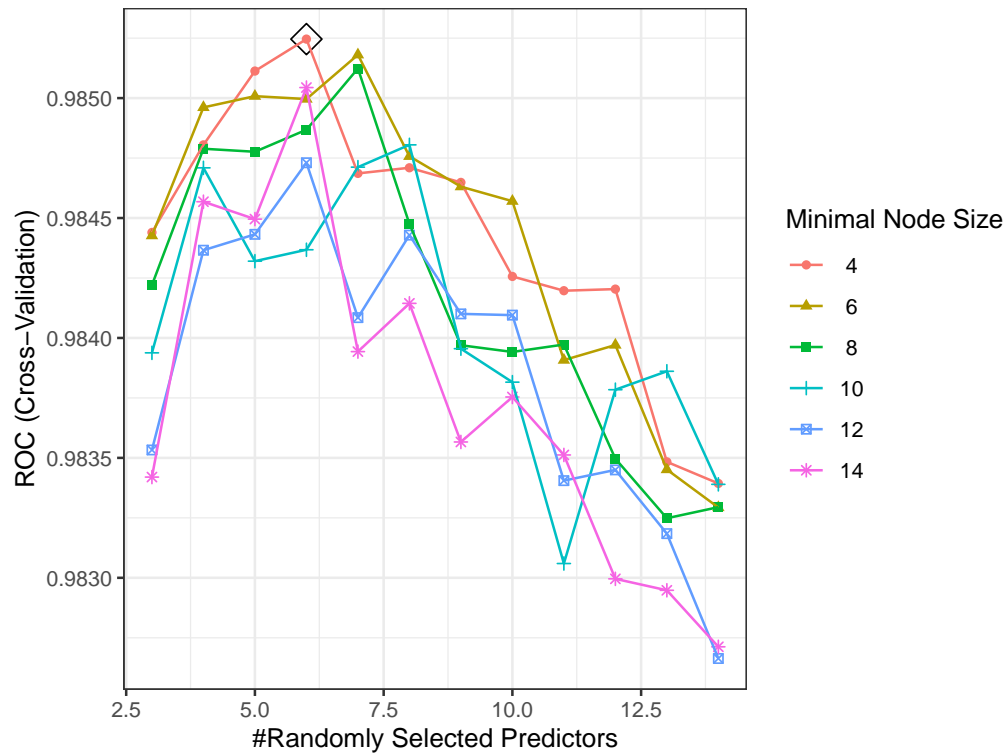
## Random forests

```
set.seed(1234)
model.rf = train(x_train,
                 y_train,
                 method = "ranger",
                 tuneGrid = expand.grid(mtry = 3:14,
                                       splitrule = "gini",
                                       min.node.size = seq(4, 14, by = 2)),
                 metric = "ROC",
                 trControl = ctrl)

model.rf$bestTune

##      mtry splitrule min.node.size
## 19      6      gini              4

ggplot(model.rf, highlight = T) +
  theme_bw()
```



```
## test auc and misclassification error rate
pred_rf_auc <- predict(model.rf, newdata = x_test, type = "prob")[,2]
roc(y_test, pred_rf_auc)$auc[1]
```

```
## Setting levels: control = no, case = yes
```

```
## Setting direction: controls < cases
```

```
## [1] 0.9764744
```

```
pred_rf <- predict(model.rf, newdata = x_test)
pred.miserror_rf <- 1 - mean(pred_rf == y_test)
pred.miserror_rf
```

```
## [1] 0.0775
```

```
confusionMatrix(data = as.factor(pred_rf),
                 reference = as.factor(y_test))
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction no yes
```

```
##           no 201 16
```

```
##           yes 15 168
```

```
##
```

```
##           Accuracy : 0.9225
```

```
##           95% CI : (0.8918, 0.9467)
```

```
##           No Information Rate : 0.54
```

```
##           P-Value [Acc > NIR] : <2e-16
```

```
##
```

```
##           Kappa : 0.8439
```

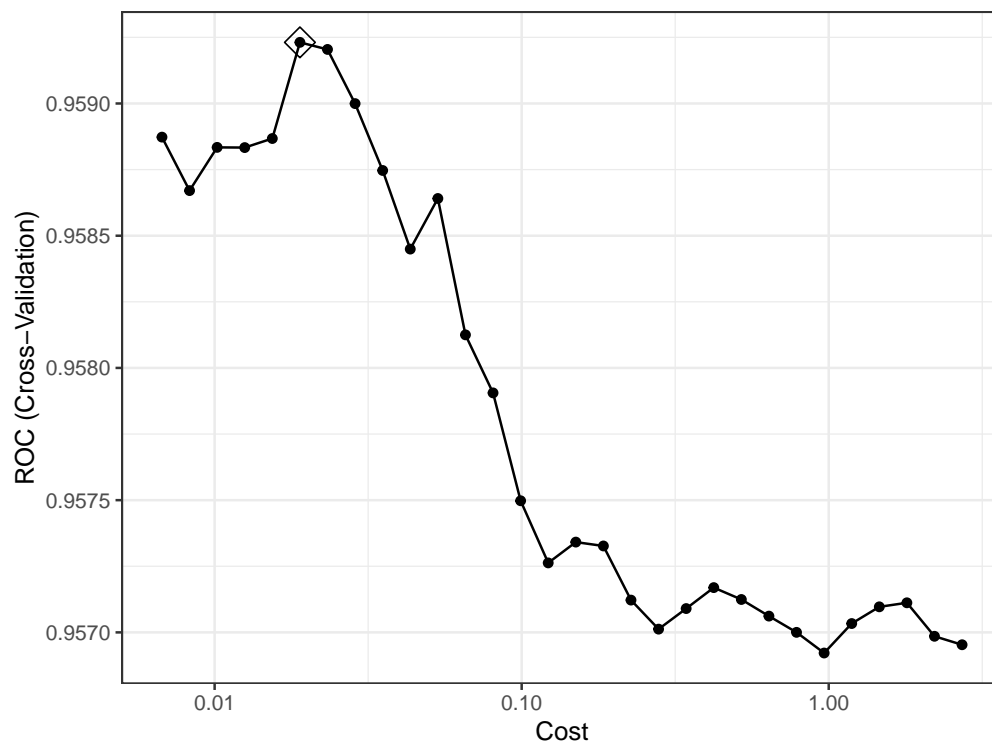
```
##
```

```
## McNemar's Test P-Value : 1
##
##      Sensitivity : 0.9306
##      Specificity : 0.9130
##      Pos Pred Value : 0.9263
##      Neg Pred Value : 0.9180
##      Prevalence : 0.5400
##      Detection Rate : 0.5025
##      Detection Prevalence : 0.5425
##      Balanced Accuracy : 0.9218
##
##      'Positive' Class : no
##
```

### Fit a support vector classifier (linear kernel)

```
set.seed(1234)
model.svm1 = train(x_train,
  y_train,
  method = "svmLinear",
  metric = "ROC",
  tuneGrid = data.frame(C = exp(seq(-5, 1, length = 30))),
  trControl = ctrl)

ggplot(model.svm1, highlight = TRUE) +
  scale_x_continuous(trans = scales::log_trans(),
    breaks = scales::log_breaks()) +
  theme_bw()
```



```
## test auc and misclassification error rate
pred_svml_auc <- predict(model_svml, newdata = x_test, type = "prob")[,2]
roc(y_test, pred_svml_auc)$auc[1]
```

```
## [1] 0.9557921
```

```
pred_svml <- predict(model_svml, newdata = x_test)
pred.miserror_svml <- 1 - mean(pred_svml == y_test)
pred.miserror_svml
```

```
## [1] 0.1275
```

```
confusionMatrix(data = as.factor(pred_svml),
                 reference = as.factor(y_test))
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction  no yes
```

```
##           no 189 24
```

```
##           yes 27 160
```

```
##
```

```
##           Accuracy : 0.8725
```

```
##           95% CI : (0.8358, 0.9036)
```

```
##           No Information Rate : 0.54
```

```
##           P-Value [Acc > NIR] : <2e-16
```

```
##
```

```
##           Kappa : 0.7437
```

```
##
```

```
##           Mcnemar's Test P-Value : 0.7794
```

```
##
```

```
##           Sensitivity : 0.8750
```

```
##           Specificity : 0.8696
```

```
##           Pos Pred Value : 0.8873
```

```
##           Neg Pred Value : 0.8556
```

```
##           Prevalence : 0.5400
```

```
##           Detection Rate : 0.4725
```

```
##           Detection Prevalence : 0.5325
```

```
##           Balanced Accuracy : 0.8723
```

```
##
```

```
##           'Positive' Class : no
```

```
##
```

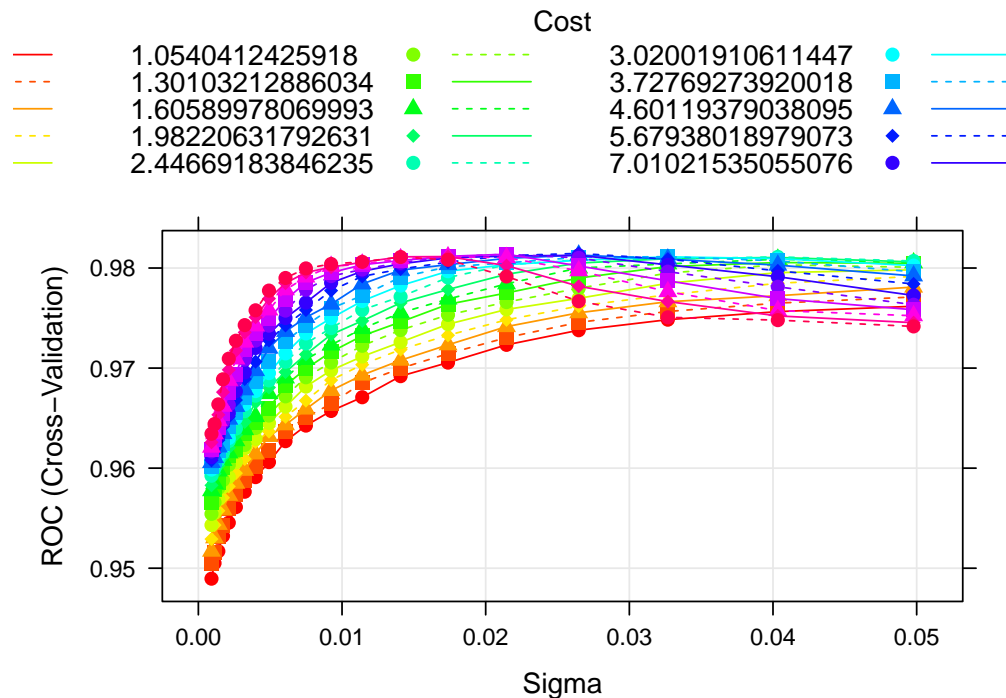
## Fit a support vector machine with a radial kernel

```
set.seed(1234)
model_svmr = train(x_train,
                   y_train,
                   method = "svmRadialSigma",
                   metric = "ROC",
                   tuneGrid = expand.grid(C = exp(seq(-1, 3, length = 20)),
                                          sigma = exp(seq(-7, -3, length = 20))),
                   trControl = ctrl)

model_svmr$bestTune
```

```
##          sigma      C
## 277 0.02647435 5.67938
```

```
myCol<- rainbow(20)
myPar <- list(superpose.symbol = list(col = myCol),
superpose.line = list(col = myCol))
plot(model.svmr, highlight = TRUE, par.settings = myPar)
```



```
## test auc and misclassification error rate
pred_svmr_auc <- predict(model.svmr, newdata = x_test, type = "prob")[,2]
roc(y_test, pred_svmr_auc)$auc[1]
```

```
## Setting levels: control = no, case = yes
```

```
## Setting direction: controls < cases
```

```
## [1] 0.9739583
```

```
pred_svmr <- predict(model.svmr, newdata = x_test)
pred.miserror_svmr <- 1 - mean(pred_svmr == y_test)
pred.miserror_svmr
```

```
## [1] 0.0875
```

```
confusionMatrix(data = as.factor(pred_svmr),
reference = as.factor(y_test))
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##          Reference
```

```
## Prediction no yes
```

```
##          no 198 17
```

```
##          yes 18 167
```

```
##
```

```
##          Accuracy : 0.9125
```

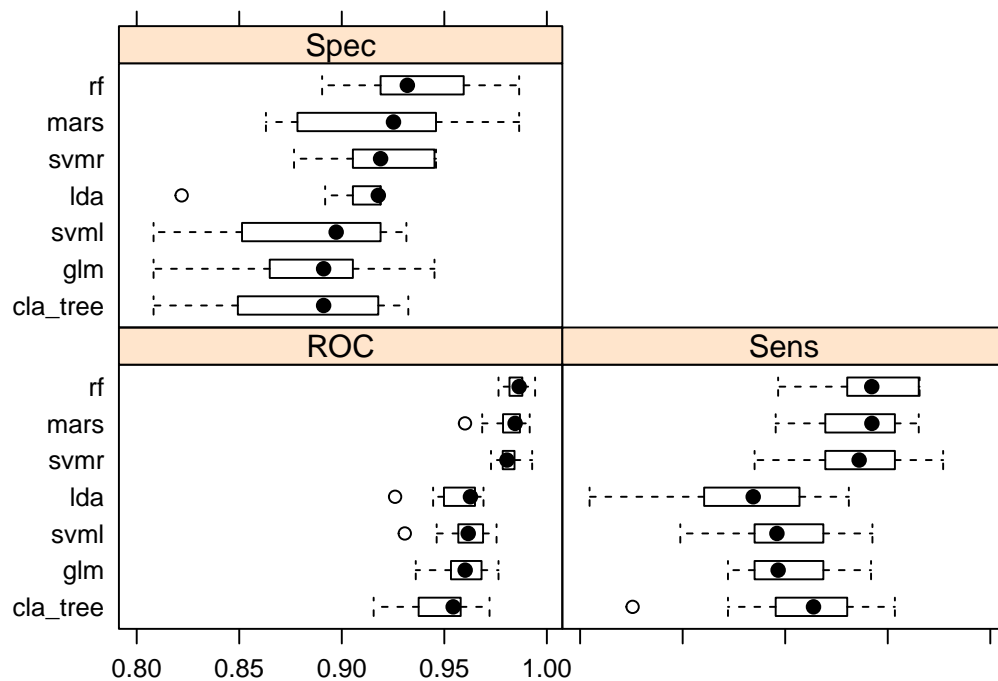


```
##          95% CI : (0.8804, 0.9383)
##    No Information Rate : 0.54
##    P-Value [Acc > NIR] : <2e-16
##
##          Kappa : 0.8239
##
##    Mcnemar's Test P-Value : 1
##
##          Sensitivity : 0.9167
##          Specificity : 0.9076
##          Pos Pred Value : 0.9209
##          Neg Pred Value : 0.9027
##          Prevalence : 0.5400
##          Detection Rate : 0.4950
##    Detection Prevalence : 0.5375
##          Balanced Accuracy : 0.9121
##
##          'Positive' Class : no
##
```

## Resample

```
resamp <- resamples(list(glm = model.glm, mars = model.mars,
                        lda = model.lda, cla_tree = model.tree,
                        rf = model.rf, svm1 = model.svm1,
                        svmr = model.svmr))
```

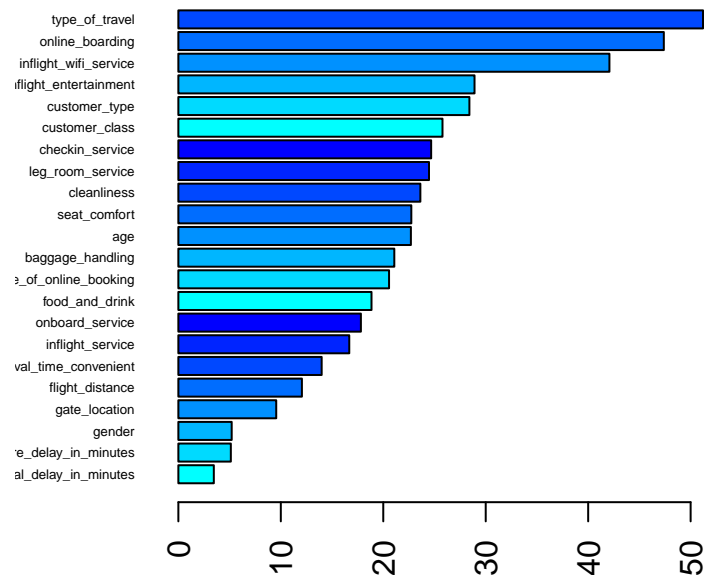
```
bwplot(resamp)
```



## Select the rf model and interpret

```
## importance variable
set.seed(1234)
rf2.final.per <- ranger(factor(satisfaction) ~ .,
  data = train,
  mtry = model.rf$bestTune[[1]],
  min.node.size = model.rf$bestTune[[3]],
  splitrule = "gini",
  importance = "permutation",
  scale.permutation.importance = TRUE)

barplot(sort(ranger::importance(rf2.final.per), decreasing = FALSE),
  las = 2,
  horiz = TRUE,
  cex.names = 0.4,
  col = colorRampPalette(colors = c("cyan", "blue"))(8))
```



```
pdp1.rf <- model.rf %>%
  partial(pred.var = c("inflight_wifi_service")) %>%
  autoplot(train = train, rug = TRUE)

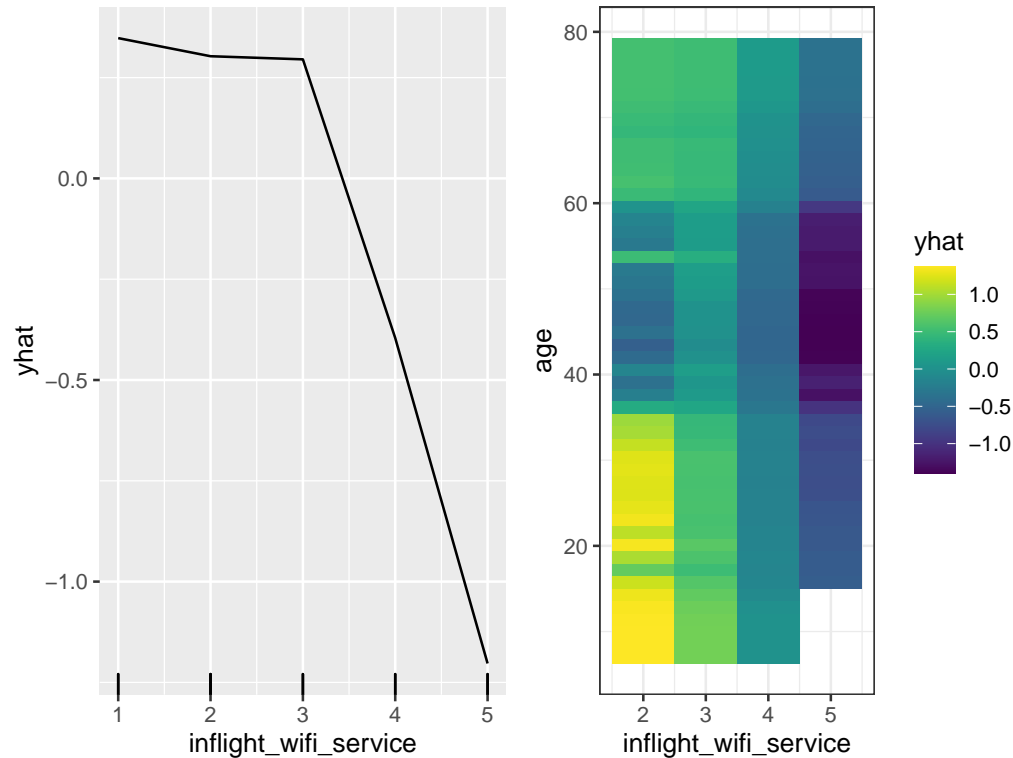
pdp2.rf <- model.rf %>%
  partial(pred.var = c("inflight_wifi_service", "age"), chull = TRUE) %>%
  autoplot(train = train, rug = TRUE)

grid.arrange(pdp1.rf, pdp2.rf, nrow = 1)
```

```
## Warning: Use of `object[[1L]]` is discouraged. Use `.data[[1L]]` instead.
## Warning: Use of `object[["yhat"]]` is discouraged. Use `.data[["yhat"]]`
## instead.
## Warning: Use of `x.rug[[1L]]` is discouraged. Use `.data[[1L]]` instead.
## Warning: Use of `object[[1L]]` is discouraged. Use `.data[[1L]]` instead.
## Warning: Use of `object[[2L]]` is discouraged. Use `.data[[2L]]` instead.
```

```
## Warning: Use of `object[["yhat"]}` is discouraged. Use `.data[["yhat"]}`  
## instead.
```

```
## Warning: Use of `object[["yhat"]}` is discouraged. Use `.data[["yhat"]}`  
## instead.
```



```
ice.rf <- model.rf %>%  
  partial(pred.var = "inflight_wifi_service",  
    grid.resolution = 100,  
    ice = TRUE) %>%  
  autoplot(train = dat, alpha = .1,  
    center = TRUE)
```

```
## Warning: `fun.y` is deprecated. Use `fun` instead.
```

```
ice.rf
```

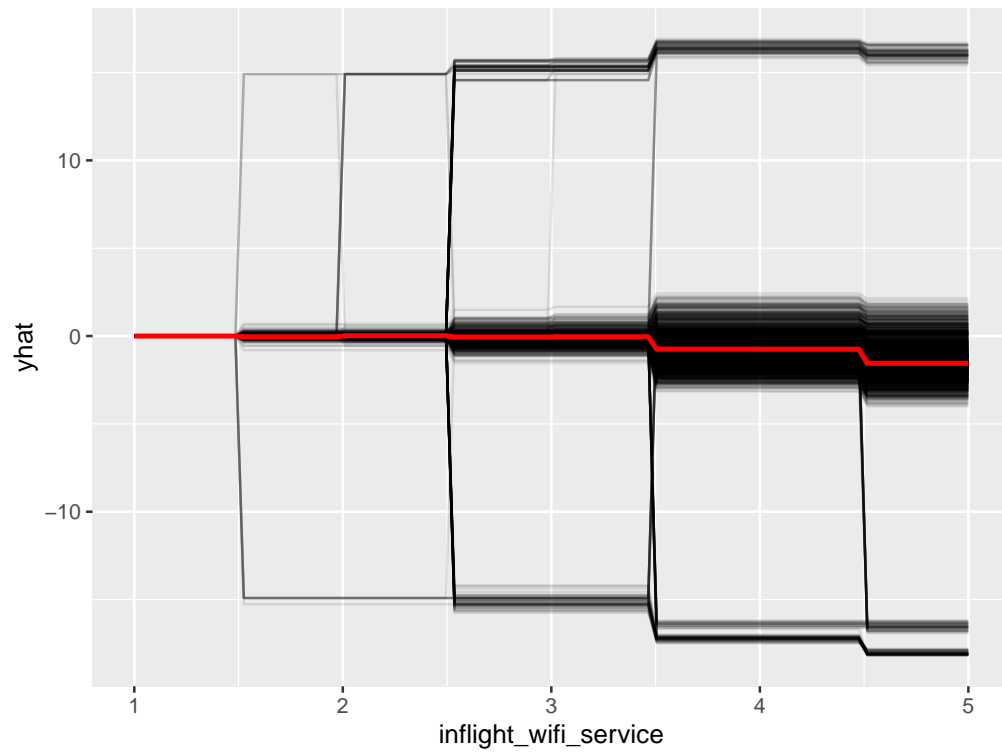
```
## Warning: Use of `object[["yhat.id"]}` is discouraged. Use `.data[["yhat.id"]}`  
## instead.
```

```
## Warning: Use of `object[[1L]]` is discouraged. Use `.data[[1L]]` instead.
```

```
## Warning: Use of `object[["yhat"]}` is discouraged. Use `.data[["yhat"]}`  
## instead.
```

```
## Warning: Use of `object[[1L]]` is discouraged. Use `.data[[1L]]` instead.
```

```
## Warning: Use of `object[["yhat"]}` is discouraged. Use `.data[["yhat"]}`  
## instead.
```



```
##lime
#explain.rf <- lime(data.frame(x_train), model.rf)

#new_obs = x_train[1:10,]
#explana.obs = explain(data.frame(new_obs),
#                        explain.rf,
#                        n_features = 10)

#plot_features(explana.obs)
```