

Assignment3.2

Amrrial Khatib

6/3/2022

Question 2

```
## Import libraries
library(dplyr)

## Warning: package 'dplyr' was built under R version 4.1.3

##
## Attaching package: 'dplyr'

##
## The following objects are masked from 'package:stats':
##
##   filter, lag
##
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library(ggplot2)

## Warning: package 'ggplot2' was built under R version 4.1.3

library(lattice) # caret requirement
library(caret)

## Warning: package 'caret' was built under R version 4.1.3

library(tidyverse)

## Warning: package 'tidyverse' was built under R version 4.1.3

## -- Attaching packages ----- tidyverse 1.3.1 --

## v tibble 3.1.7   v purrr  0.3.4
## v dplyr  1.2.0   v rrrrr  1.4.0
## v readr  2.1.2   v forcats 0.5.1

## Warning: package 'libble' was built under R version 4.1.3

## Warning: package 'tidyr' was built under R version 4.1.3

## Warning: package 'forcats' was built under R version 4.1.3

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()   masks stats::lag()
## x purrr::lift()  masks caret::lift()

library(rpart)
library(rpart.plot)

## Warning: package 'rpart.plot' was built under R version 4.1.3

library(cmaplot)

## Warning: package 'cmaplot' was built under R version 4.1.3
```

Read Sacramento dataset

```
data(Sacramento)
sac_data = Sacramento
sac_data$limits = factor(ifelse(sac_data$city == "SACRAMENTO", "in", "out"))
sac_data = subset(sac_data, select = ~(city, zip, baths))
```

Data splitting

```
set.seed(42)
sac_trn_idx = sample(nrow(sac_data), size = 0.8 * nrow(sac_data))
sac_trn = sac_data[sac_trn_idx,]
sac_tst = sac_data[-sac_trn_idx,]

sac_est_idx = sample(nrow(sac_trn), size = 0.8 * nrow(sac_trn))
sac_est = sac_trn[sac_est_idx,]
sac_val = sac_trn[-sac_est_idx,]
```

RMSE function

```
rme = function(predicted, actual){sqrt(mean((actual - predicted)^2))}
```

Part A: Linear model

```
sac_reg_model_list = list()
model.1 = lm(formula = price ~ sqft + beds + limits, data = sac_est),
model.2 = lm(formula = price ~ ., data = sac_est),
model.3 = stepAIC(lm(formula = price ~ ., data = sac_est), trace = FALSE, direction = "forward"),
model.4 = lm(formula = price ~ (sqft ^ 2) + beds + limits + (latitude * longitude) ^ 2, data = sac_est),
model.5 = lm(formula = price ~ . - limits - latitude - longitude, data = sac_est)

sac_est_reg_predicted_list = lapply(sac_reg_model_list, predict, sac_est)
sac_val_reg_predicted_list = lapply(sac_reg_model_list, predict, sac_val)

sac_rmse_reg_est_list = sapply(sac_est_reg_predicted_list, rmse, sac_val$price)
sac_rmse_reg_val_list = sapply(sac_val_reg_predicted_list, rmse, sac_val$price)

report_reg_models = data.frame(data_type = c("estimation", "validation"),
                               model.1_rmse = c(sac_rmse_reg_est_list[1], sac_rmse_reg_val_list[1]),
                               model.2_rmse = c(sac_rmse_reg_est_list[2], sac_rmse_reg_val_list[2]),
                               model.3_rmse = c(sac_rmse_reg_est_list[3], sac_rmse_reg_val_list[3]),
                               model.4_rmse = c(sac_rmse_reg_est_list[4], sac_rmse_reg_val_list[4]),
                               model.5_rmse = c(sac_rmse_reg_est_list[5], sac_rmse_reg_val_list[5]))

report_reg_models
```

```
##      data_type model.1_rmse model.2_rmse model.3_rmse model.4_rmse model.5_rmse
## 1 estimation    77539.83    75872.25    75872.25    80431.33    75144.34
## 2 validation    91487.72    89761.84    89761.84    95228.69    82980.59
```

By fitting model 5 (price ~ . - limits - latitude - longitude), we didn't encounter with significant change in RMSE in comparison with the other models' RMSE value. Hence, we can conclude that the absence of **limits**, **latitude**, **longitude** do not causes a considerable shift in RMSE value of the model.

Part B: KNN model

Feature scaling

```
## Estimation set
beds.s = scale(sac_est$beds)
beds.center = attr(beds.s, "scaled:center")
beds.scale = attr(beds.s, "scaled:scale")

latitude.s = scale(sac_est$latitude)
latitude.center = attr(latitude.s, "scaled:center")
latitude.scale = attr(latitude.s, "scaled:scale")

longitude.s = scale(sac_est$longitude)
longitude.center = attr(longitude.s, "scaled:center")
longitude.scale = attr(longitude.s, "scaled:scale")

sqft.s = scale(sac_est$sqft)
sqft.center = attr(sqft.s, "scaled:center")
sqft.scale = attr(sqft.s, "scaled:scale")

sac_est_scaled = data.frame(beds.s = beds.s, sqft.s = sqft.s,
                           type = sac_est$type, price = sac_est$price,
                           latitude.s = latitude.s,
                           longitude.s = longitude.s,
                           limits = sac_est$limits)

head(sac_est_scaled, 5)

##      beds.s    sqft.s    type price latitude.s longitude.s limits
## 1 -0.2841249 -0.71540995 Residential 181000 -0.3358432 -0.05052174   in
## 2  0.8226222 -0.43704621 Residential 209996  0.3695965 -0.52976228   in
## 3 -0.2841249 -0.83735978 Residential 120900 -2.0317650  0.36311540   out
## 4 -1.3899119 -0.28328337 Residential 195009  0.6980635  2.59948490   out
## 5 -1.3899119 -0.09903388 Residential 220702  0.2745569  0.16372781   out

## Validation set
beds.s_val = scale(sac_val$beds, center = beds.center, scale = beds.scale)
latitude.s_val = scale(sac_val$latitude, center = latitude.center, scale = latitude.scale)
longitude.s_val = scale(sac_val$longitude, center = longitude.center, scale = longitude.scale)
sqft.s_val = scale(sac_val$sqft, center = sqft.center, scale = sqft.scale)

sac_val_scaled = data.frame(beds.s = beds.s_val, sqft.s = sqft.s_val,
                           type = sac_val$type, price = sac_val$price,
                           latitude.s = latitude.s_val,
                           longitude.s = longitude.s_val,
                           limits = sac_val$limits)

head(sac_val_scaled, 5)

##      beds.s    sqft.s    type price latitude.s longitude.s limits
## 1 -0.2841249 -0.58033863 Residential 180900  1.8879680  4.8032884   out
## 2  0.8226222 -0.31111597 Residential 176805  1.8605236 -0.1036678   out
## 3  0.8226222  0.26283898 Residential 309000  0.6839732 -0.8022655   in
## 4 -0.2841249 -0.33806101 Residential 120900  0.7084474  0.8027248   out
## 5 -0.2841249 -0.92748958 Residential  90000 -0.7879952 -0.8396848   in
```

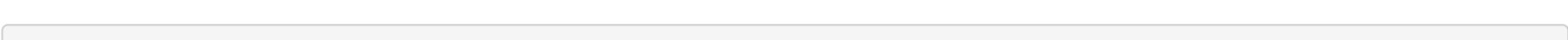
Train KNN models

```
sac_knn_model_list = list()
for (k in 1:180){
  sac_knn_model_list[[k]] = knnreg(formula = price ~ ., data = sac_est, k = k)
}
sac_knn_predicted_list = lapply(sac_knn_model_list, predict, sac_val)
sac_knn_rmse_list = sapply(sac_knn_predicted_list, rmse, sac_val$price)

sac_scaled_knn_model_list = list()
for (k in 1:180){
  sac_scaled_knn_model_list[[k]] = knnreg(formula = price ~ ., data = sac_est_scaled, k = k)
}
sac_scaled_knn_predicted_list = lapply(sac_scaled_knn_model_list, predict, sac_val_scaled)
sac_scaled_knn_rmse_list = sapply(sac_scaled_knn_predicted_list, rmse, sac_val$price)

report_knn_models = data.frame(k = as.numeric(1:180),
                               rmse = as.numeric(sac_knn_rmse_list),
                               rmse_scaled = as.numeric(sac_scaled_knn_rmse_list))

report_knn_models
```



trained with normalized data) is better than the one which was not trained with normalized data (for k in range 1 to 31) The model performance (which was not trained with normalized data) is better than the one which was trained with normalized data (for k in range 32 to 100)

Part C: Regression tree model

```
cp_list = c(1,0.1,0.01,0.001,0)
mingsplit_list = c(5,20)
sac_regtr_list = list()
cp_column = list()
mingsplit_column = list()
cp_mingsplit = list()
sac_rmse_regtr_list = list()
for (cp in cp_list){
  mingsplit_column = append(mingsplit_column, m)
  cp_column = append(cp_column, cp)
  cp_mingsplit = append(cp_mingsplit, paste(cp,"",m))
  regtr_model = rpart(formula = price ~ ., data = sac_est, cp = cp, mingsplit = m)
  predict_regtr = predict(regtr_model, sac_val)
  rmse_regtr = rmse(predict_regtr, sac_val$price)
  sac_rmse_regtr_list = append(sac_rmse_regtr_list, rmse_regtr)
}

report_regtr_models = data.frame(cp = as.numeric(cp_column),
                                mingsplit = as.numeric(mingsplit_column),
                                cp_mingsplit = as.character(cp_mingsplit),
                                rmse = as.numeric(sac_rmse_regtr_list))

head(report_regtr_models,5)

##      cp mingsplit    rmse
## 1 1.000      5      1 5 166067.02
## 2 0.100      5      0.1 5 118666.72
## 3 0.010      5      0.01 5 90820.59
## 4 0.001      5      0.001 5 10746.15
## 5 0.000      5      0 5 106650.04

## Plot line chart for regression tree line chart
ggplot(data = report_regtr_models, aes(x = cp_mingsplit, size = 1)) +
  geom_line(aes(y = rmse, color = "non_scaled"), size = 1) +
  geom_line(aes(y = rmse_scaled, color = "scaled"), size = 1) +
  labs(x = "k", y = "RMSE", color = "Scaling") +
  scale_color_manual(values = c("non_scaled" = "#E31A2C", "scaled" = "blue"))
```

The lower RMSE value a model have, the closer distance to the y = x line its points (predicted, actual) have.

We can deduce that the RMSE of above plotted model are approximately the same as each other.

```
final_rmse = paste("Linear model:", min(report_regtr_models$rmse),
                  "KNN:", min(report_knn_models$rmse),
                  "Regtr:", min(report_regtr_models$rmse))

final_rmse

## [1] "Linear model: 89761.8419277225 KNN: 94543.5808984794 KNN.S: 87550.4969798713 Regtr: 89127.3514891294"
```

Part D: Plot RMSE value of best model on test set

```
## Choose from linear models:
which_min(report_regtr_models$rmse, 2:5)

## [1] "0.001 20"

## model 2_rmse
rmse

## [1] 32

which_min(report_knn_models$rmse, 3)

## [1] 6

paste("non_scaled:", min(report_knn_models$rmse, 3),
      "scaled:", min(report_knn_models$rmse, 3))

## [1] "non_scaled: 94543.5808984794 scaled: 87550.4969798713"

final_knn_model = knnreg(formula = price ~ ., data = sac_est_scaled, k = 6)
final_knn_predict = predict(final_knn_model, sac_val_scaled)

## Choose from regression tree models:
report_regtr_models[which_min(report_regtr_models$rmse, 2)]

## [1] "0.001 20"

final_regtr_model = rpart(formula = price ~ ., data = sac_est, cp = 0.001, mingsplit = 20)
final_regtr_predict = predict(final_regtr_model, sac_val)

## Linear regression
plots = ggplot()
geom_point(aes(x = final_reg_predict, y = sac_val$price), color = "#3A3A3A")
geom_abline(intercept = 0, slope = 1)
labs(x = "Predicted", y = "Actual")

plot = ggplot()
geom_point(aes(x = final_knn_predict, y = sac_val$price), color = "#004080")
geom_abline(intercept = 0, slope = 1)
labs(x = "Predicted", y = "Actual")

plot3 = ggplot()
geom_point(aes(x = final_reg_predict, y = sac_val$price), color = "#004080")
geom_abline(intercept = 0, slope = 1)
labs(x = "Predicted", y = "Actual")

plot_grid(plot1, plot2, plot3, labels = c("Linear model", "KNN", "Regression tree"))
```

We can say:

The lower RMSE value a model have, the closer distance to the y = x line its points (predicted, actual) have.

We can deduce that the RMSE of above plotted model are approximately the same as each other.

final_rmse = paste("Linear model:", min(report_regtr_models\$rmse),

"KNN:", min(report_knn_models\$rmse),

"Regtr:", min(report_regtr_models\$rmse))

final_rmse

[1] "Linear model: 89761.8419277225 KNN: 94543.5808984794 KNN.S: 87550.4969798713 Regtr: 89127.3514891294"

Part E: RMSE of final model on test set

```
beds.s = scale(sac_tst$beds)
beds.center = attr(beds.s, "scaled:center")
beds.scale = attr(beds.s, "scaled:scale")

latitude.s = scale(sac_tst$latitude)
latitude.center = attr(latitude.s, "scaled:center")
latitude.scale = attr(latitude.s, "scaled:scale")

longitude.s = scale(sac_tst$longitude)
longitude.center = attr(longitude.s, "scaled:center")
longitude.scale = attr(longitude.s, "scaled:scale")

sqft.s = scale(sac_tst$sqft)
sqft.center = attr(sqft.s, "scaled:center")
sqft.scale = attr(sqft.s, "scaled:scale")

sac_trn_scaled = data.frame(beds.s = beds.s, sqft.s = sqft.s,
                           type = sac_trn$type, price = sac_trn$price,
                           latitude.s = latitude.s,
                           longitude.s = longitude.s,
                           limits = sac_trn$limits)

sac_val_scaled = data.frame(beds.s = beds.s, sqft.s = sqft.s,
                           type = sac_val$type, price = sac_val$price,
                           latitude.s = latitude.s,
                           longitude.s = longitude.s,
                           limits = sac_val$limits)

final_predict = predict(final_model, sac_tst_scaled)
final_rmse = rmse(final_predict, sac_tst$price)
final_rmse

## [1] 81899.39
```