

Model Fitting II

Carolina Alves de Lima Salge

3/16/2021

Slide Code

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.0 --
```

```
## v ggplot2 3.3.3    v purrr  0.3.4
## v tibble  3.0.6    v dplyr  1.0.4
## v tidyr   1.1.2    v stringr 1.4.0
## v readr   1.4.0    v forcats 0.5.1
```

```
## -- Conflicts ----- tidyverse_conflicts() --
```

```
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```
library(caret)
```

```
## Loading required package: lattice
```

```
##
```

```
## Attaching package: 'caret'
```

```
## The following object is masked from 'package:purrr':
```

```
##
```

```
## lift
```

```
insurance <- read_csv("insurance.csv")
```

```
##
```

```
## -- Column specification -----
```

```
## cols(
```

```
##   age = col_double(),
```

```
##   sex = col_character(),
```

```
##   bmi = col_double(),
```

```
##   children = col_double(),
```

```
##   smoker = col_character(),
```

```
##   region = col_character(),
```

```
##   charges = col_double()
```

```
## )
```

```
insurance
```

```
## # A tibble: 1,338 x 7
##   age sex    bmi children smoker region    charges
##   <dbl> <chr> <dbl>    <dbl> <chr>  <chr>    <dbl>
## 1    19 female  27.9        0 yes    southwest 16885.
## 2    18 male   33.8        1 no     southeast 1726.
## 3    28 male   33         3 no     southeast 4449.
## 4    33 male   22.7        0 no     northwest 21984.
## 5    32 male   28.9        0 no     northwest 3867.
## 6    31 female 25.7        0 no     southeast 3757.
## 7    46 female 33.4        1 no     southeast 8241.
## 8    37 female 27.7        3 no     northwest 7282.
## 9    37 male   29.8        2 no     northeast 6406.
## 10   60 female 25.8        0 no     northwest 28923.
## # ... with 1,328 more rows
```

```
# transform categories to numbers
```

```
insurance <- insurance %>%
```

```
  mutate(sexN = case_when(
    sex == "male" ~ 1,
    sex == "female" ~ 0
  )) %>%
```

```
  mutate(smokerN = case_when(
    smoker == "yes" ~ 1,
    smoker == "no" ~ 0
  )) %>%
```

```
  mutate(regionN = case_when(
    region == "southwest" ~ 1,
    region == "southeast" ~ 2,
    region == "northwest" ~ 3,
    region == "northeast" ~ 4
  ))
```

```
# only select numeric variables
```

```
df <- insurance %>%
```

```
  dplyr::select(charges, age, sexN, bmi, children, smokerN, regionN)
```

```
# drop missing values NAs
```

```
df1 <- drop_na(df)
```

```
# compute correlation between predictors
```

```
cor(df1[,2:7])
```

```
##           age           sexN           bmi    children    smokerN
## age      1.000000000 -0.020855872  0.109271882  0.04246900 -0.025018752
## sexN     -0.020855872  1.000000000  0.046371151  0.01716298  0.076184817
## bmi       0.109271882  0.046371151  1.000000000  0.01275890  0.003750426
## children  0.042468999  0.017162978  0.012758901  1.000000000  0.007673120
## smokerN  -0.025018752  0.076184817  0.003750426  0.00767312  1.000000000
## regionN  -0.002127313 -0.004588385 -0.157565849 -0.01656945  0.002180682
##           regionN
## age      -0.002127313
```

```
## sexN      -0.004588385
## bmi       -0.157565849
## children  -0.016569446
## smokerN   0.002180682
## regionN   1.000000000
```

```
# compute correlation between predictors and the target
cor(df1[,1:7])
```

```
##           charges           age           sexN           bmi           children
## charges  1.000000000  0.299008193  0.057292062  0.198340969  0.06799823
## age      0.299008193  1.000000000 -0.020855872  0.109271882  0.04246900
## sexN     0.057292062 -0.020855872  1.000000000  0.046371151  0.01716298
## bmi      0.198340969  0.109271882  0.046371151  1.000000000  0.01275890
## children 0.067998227  0.042468999  0.017162978  0.012758901  1.00000000
## smokerN  0.787251430 -0.025018752  0.076184817  0.003750426  0.00767312
## regionN  0.006208235 -0.002127313 -0.004588385 -0.157565849 -0.01656945
##           smokerN           regionN
## charges  0.787251430  0.006208235
## age      -0.025018752 -0.002127313
## sexN     0.076184817 -0.004588385
## bmi      0.003750426 -0.157565849
## children 0.007673120 -0.016569446
## smokerN  1.000000000  0.002180682
## regionN  0.002180682  1.000000000
```

```
# age, bmi, and smoking are highly correlated with health costs
```

```
set.seed(12L) # set a starting seed to be able to get reproducible results
```

```
# partition data
```

```
trainIndex <- createDataPartition(df1$charges, # target variable
                                   p = 0.8, # percentage that goes to training
                                   list = FALSE, # results will not be in a list
                                   times = 1) # number of partitions to create
```

```
charges_train <- df1[trainIndex, ] # data frame for training
```

```
## Warning: The 'i' argument of '['()' can't be a matrix as of tibble 3.0.0.
## Convert to a vector.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_warnings()' to see where this warning was generated.
```

```
charges_test <- df1[-trainIndex, ] # data frame for testing
```

```
# use training set to build model
```

```
model <- train(charges ~ age + bmi + smokerN,
               data = charges_train, # use training set
               method = "lm") # linear regression
```

```
# now predict outcomes in test set
```

```
p <- predict(model, charges_test)
```

```
# how did we do? calculate performance across resamples
# RMSE and R-squared
postResample(pred = p, obs = charges_test$charges)
```

```
##          RMSE      Rsquared      MAE
## 5808.0045894    0.7989742 4184.9721150
```

```
# on average, our prediction is off by $5,808.00
```

```
# how can we improve performance? Try a different method!
model2 <- train(charges ~ age + bmi + smokerN,
               data = charges_train, # use training set
               method = "ranger") # random forest
```

```
## note: only 2 unique complexity parameters in default grid. Truncating the grid to 2 .
```

```
# now predict outcomes in test set
p1 <- predict(model2, charges_test)
```

```
# how did we do? calculate performance across resamples
# RMSE and R-squared
postResample(pred = p1, obs = charges_test$charges)
```

```
##          RMSE      Rsquared      MAE
## 4205.4282707    0.8933816 2436.1432025
```

```
# on average, our prediction is off by $4,632.99
```

```
# first collect the resampling results of each model
resamps <- resamples(list(LM = model,
                        RF = model2))
resamps
```

```
##
## Call:
## resamples.default(x = list(LM = model, RF = model2))
##
## Models: LM, RF
## Number of resamples: 25
## Performance metrics: MAE, RMSE, Rsquared
## Time estimates for: everything, final model fit
```

```
# then use a simple t-test to evaluate the null hypothesis that there is no difference
summary(diff(resamps))
```

```
##
## Call:
## summary.diff.resamples(object = diff(resamps))
##
## p-value adjustment: bonferroni
```

```

## Upper diagonal: estimates of the difference
## Lower diagonal: p-value for H0: difference = 0
##
## MAE
##      LM          RF
## LM          1403
## RF < 2.2e-16
##
## RMSE
##      LM          RF
## LM          1156
## RF 2.139e-15
##
## Rsquared
##      LM          RF
## LM          -0.09466
## RF 1.17e-14

```