

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  3.1.0    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 == "northwest" ~ 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 -0.019119318  0.09570218 0.058405316 -0.02104324
## sexN     -0.019119318  1.000000000  0.06887488 0.011989296  0.07964590
## bmi       0.095702182  0.068874880  1.00000000 0.019267216  0.02191144
## children 0.058405316  0.011989296  0.01926722 1.000000000  0.01040992
## smokerN  -0.021043239  0.079645898  0.02191144 0.010409918  1.00000000
## regionN  -0.007367464 -0.004927365 -0.09209306 0.002040754  0.00000000
##           regionN
## age      -0.007367464
```

```
## sexN      -0.004927365
## bmi       -0.092093057
## children  0.002040754
## smokerN   0.000000000
## regionN   1.000000000
```

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

```
##          charges          age          sexN          bmi          children
## charges  1.000000000  0.298743694  0.062371225  0.19260362  0.068287662
## age      0.298743694  1.000000000 -0.019119318  0.09570218  0.058405316
## sexN     0.062371225 -0.019119318  1.000000000  0.06887488  0.011989296
## bmi      0.192603622  0.095702182  0.068874880  1.00000000  0.019267216
## children 0.068287662  0.058405316  0.011989296  0.01926722  1.000000000
## smokerN  0.802212827 -0.021043239  0.079645898  0.02191144  0.010409918
## regionN  0.002286109 -0.007367464 -0.004927365 -0.09209306  0.002040754
##          smokerN          regionN
## charges  0.80221283  0.002286109
## age      -0.02104324 -0.007367464
## sexN     0.07964590 -0.004927365
## bmi      0.02191144 -0.092093057
## children 0.01040992  0.002040754
## smokerN  1.00000000  0.000000000
## regionN  0.00000000  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
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
## 5779.1881082    0.7916937 4016.8698553
```

```
# on average, our prediction is off by $5,779.18
```

```
# 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  
## 4603.8843625    0.8683862 2625.9694526
```

```
# 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      1491  
## RF < 2.2e-16  
##
```

```
## RMSE
##      LM      RF
## LM      1220
## RF 4.433e-13
##
## Rsquared
##      LM      RF
## LM      -0.08187
## RF 5.742e-11
```