

## STAT 380 – Assessing Model Accuracy Part 2 (Lecture 7)

EXAMPLE 1: (Sample data dictionary) L06\_Insurance\_m.csv contains information about a number of health insurance policies. In particular, the data set contains some attributes of the policy holder (such as age, sex, etc.) and the total charges billed by the health care provider. (While similar, this dataset is different than the Insurance.txt dataset used earlier in the course.) Here are the variables included:

- age: age of primary beneficiary
- sex: sex of primary beneficiary
- bmi: Body mass index, providing an understanding of body, weights that are relatively high or low relative to height, objective index of body weight ( $\text{kg} / \text{m}^2$ ) using the ratio of height to weight, ideally 18.5 to 24.9
- children: Number of children covered by the health insurance policy (i.e., the number of dependents)
- smoker: Status indicating whether the person is a smoker (options include 'yes' and 'no')
- region: the beneficiary's residential area in the US (options include northeast, southeast, southwest, northwest).
- charges: Individual medical costs as billed by health insurance

Read L06\_Insurance\_m.csv into R and name the result Ins.

- a. In Example 2 (Lecture 6), we saw that there is at least 1 NA in the dataset. Write a function called countNA that takes a single argument call dat, where dat represents a vector (column of a dataframe). The function should return the number of NA values in the column. (Code may be found in STAT380\_L7.Rmd.)

```
#This function will count the number of NA's in a vector data
countNA <- function(dat){
  numNA <- sum(is.na(dat))
  return(numNA)
}
```

b. Apply the function to each column of Ins. How much missing data is present?

Inefficient approach:

```
30 ▾ ## Example 1b - Apply function to one column at a time
31 ▾ ```{r}
32 #Apply countNA to each column of Ins
33 countNA(dat = Ins$age)
34 countNA(dat = Ins$sex)
35 countNA(dat = Ins$bmi)
36 countNA(dat = Ins$children)
37 countNA(dat = Ins$smoker)
38 countNA(dat = Ins$region)
39 countNA(dat = Ins$charges)
40 ▴ ```
```

Efficient approach:

```
42 ▾ ## Example 1b - Apply function to all columns at once
43 ▾ ```{r}
44 apply(X = Ins, MARGIN = 2, FUN = countNA)
45 ▴ ```
```

age	sex	bmi	children	smoker	region	charges
0	0	0	0	0	0	1

c. Remove any observations with missing data from the dataset. Name the new dataset Ins. (You should have 1337 observations after this step.)

```
47 ▾ ## Example 1c - Remove any observation (row) that includes an NA
48 ▾ ```{r}
49 Ins <-
50   Ins %>%
51   na.omit()
52 ▴ ```
```

d. To learn how to reproducibly generate random numbers in R, consider the following examples.

```
#Randomly generate 9 numbers from 1 to 20  
sample(x = 1:20, size = 9, replace = TRUE)
```

```
## [1]  2  7 19 13 13 16  9 16 12
```

```
59 #Randomly generate 9 numbers from 1 to 20 using a seed of 123  
60 set.seed(1234)  
61 sample(x = 1:20, size = 9, replace = TRUE) #16  5 12 15  9  5  6 16  4  
62 set.seed(NULL)  
63  
64 #Randomly generate 3 numbers from 1 to 20 3 times using a seed of 123  
65 set.seed(1234)  
66 sample(x = 1:20, size = 3, replace = TRUE) #16  5 12  
67 sample(x = 1:20, size = 3, replace = TRUE) #15  9  5  
68 sample(x = 1:20, size = 3, replace = TRUE) #6 16  4  
69 set.seed(NULL)  
70  
71 #What does set.seed(NULL) do?  
72 set.seed(1234)  
73 sample(x = 1:20, size = 3, replace = TRUE) #16  5 12  
74 sample(x = 1:20, size = 3, replace = TRUE) #15  9  5  
75 set.seed(NULL)  
76 sample(x = 1:20, size = 3, replace = TRUE)
```

e. Using a seed of 123, perform an 80/20 training/testing split.

IDEA: Randomly select 80% of the row numbers and store the result as train\_ind. Those row numbers determine the training set. The unselected row numbers, which correspond to the remaining 20%, determine the testing set.

```
84 set.seed(123)  
85 train_ind <- sample(x = 1:nrow(Ins), size = floor(0.8*nrow(Ins)))  
86 set.seed(NULL)  
87  
88 Train <- Ins[train_ind, ]  
89 Test <- Ins[-train_ind, ]
```

NOTE: There are other ways to do this in R; however, since there is randomness involved and we must agree on a correct answer for grading purposes, the expectation is that you use this method.

f. Build a regression model for predicting charges based on person's age using the training data. Write the estimated regression equation.

g. What proportion of the variation in charges is explained by the model containing age?

h. Using the model from f., predict the charges for each person in the testing data. HINT: Make sure the number of predictions matches the number of rows in Test.

i. Compute MSE and RMSE for the testing data.

NOTE: The holdout set (or validation set) approach works well on large datasets. On small datasets, the results are sensitive to the split (i.e., the result can vary a lot based on which observations were assigned to the training set and which are assigned to the testing set).

j. Explore some multiple linear regression models. For each model you try, complete a row in the table shown below. Which gives the best MSE (or best RMSE) based on the testing set?

Predictor Variables in the model	$R^2$	$R^2_{adj}$	Test Set MSE	Test Set RMSE