

Week 12: Regression Analysis in R: Multiple Linear Regression

Jonathan Hernandez

November 12, 2018

Using R, build a multiple regression model for data that interests you. Include in this model at least

- For this discussion, I will look at the Kaggle Dataset “Medical Cost Personal Datasets” link can be found here to download: <https://www.kaggle.com/mirichoi0218/insurance>
- This dataset looks at medical insurance costs charges for various people based on several factors like number of children, region of residency, age etc.
- I will make a multiple linear regression model and make a best-fit line for computing medical costs.
- Start by loading the data and viewing attributes

```
insurance <- read.csv("insurance.csv")
dim(insurance)
```

```
## [1] 1338    7
```

```
summary(insurance)
```

```
##      age      sex      bmi      children      smoker
## Min.   :18.00  female:662  Min.   :15.96  Min.   :0.000  no :1064
## 1st Qu.:27.00  male  :676  1st Qu.:26.30  1st Qu.:0.000  yes: 274
## Median :39.00                      Median :30.40  Median :1.000
## Mean   :39.21                      Mean   :30.66  Mean   :1.095
## 3rd Qu.:51.00                      3rd Qu.:34.69  3rd Qu.:2.000
## Max.   :64.00                      Max.   :53.13  Max.   :5.000
##      region      charges
## northeast:324  Min.   : 1122
## northwest:325  1st Qu.: 4740
## southeast:364  Median : 9382
## southwest:325  Mean    :13270
##                3rd Qu.:16640
##                Max.    :63770
```

```
str(insurance)
```

```
## 'data.frame': 1338 obs. of 7 variables:
## $ age : int 19 18 28 33 32 31 46 37 37 60 ...
## $ sex : Factor w/ 2 levels "female","male": 1 2 2 2 2 1 1 1 2 1 ...
## $ bmi : num 27.9 33.8 33 22.7 28.9 ...
## $ children: int 0 1 3 0 0 0 1 3 2 0 ...
## $ smoker : Factor w/ 2 levels "no","yes": 2 1 1 1 1 1 1 1 1 1 ...
## $ region : Factor w/ 4 levels "northeast","northwest",...: 4 3 3 2 2 3 3 2 1 2 ...
## $ charges : num 16885 1726 4449 21984 3867 ...
```