

Machine Learning I Group Work

Adriana Ricklin, Yvonne Schaerli, Christina Sudermann, Carole Mattmann

11 June 2020

Contents

1 Packages	1
2 Import and data cleaning	2
3 Linear models	3
3.1 Basic analysis of continuous variables	3
3.2 Basic analysis of categorical variables	5
3.3 Basic analysis of categorical and continuous variables	8
3.4 Fitting a first Linear Model	14
3.5 Final Linear Model Development	17
4 Polynomial models	25
5 Generalised Additive Models	30
6 Generalised Linear Models	33
6.1 Generalised Linear Models for count data	33
6.2 Generalised Linear Models for binomial data	35
7 Decision Trees	37
7.1 Regression trees	37
7.2 Classification Trees	54
7.3 Pruning	57
7.4 Bagging	65
7.5 Random Forests	67
7.6 Boosting	68
8 Support Vector Machines	72
8.1 two classes linear	72
8.2 two classes polynomial	75
9 Neural Networks	77
9.1 Neural Networks using “neuralnet”	77
9.2 Neural Networks using “h2o”	86
9.3 Results and method comparison	89
10 Cross validation and discussion	89
10.1 Model overview	89
10.2 Cross validation	91
11 Final words	97

1 Packages

The following packages are used in the analysis:

```

library(gridExtra)
library(tree)      #for decision trees
library(dplyr)
library(readr)
library(randomForest) #for decision trees
library(e1071)     # for SVM
library(ISLR)
library(tibble)
library(MASS)
library(gbm)
library(ipred)
library(Metrics)
library(neuralnet) # for NN
library(h2o)       # for NN
library(bit64)     # for NN
library(tidyverse) # includes ggplot2
library(mgcv)      # for GAM

```

2 Import and data cleaning

The dataset “insurance” from <https://www.kaggle.com/mirichoi0218/insurance> is imported and transformed. The data contains information about 1338 people, their health care charges and further parameters such as age, number of children they have etc.

```

#getwd()
#setwd("~/GitHub/machinelearning/machinelearning/03_rmarkdown") #adrianas extrawurst
insurance <- read.csv("../01_data/insurance.csv", header=TRUE)
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 ...

# smoker = 1 / nonsmoker = 0
insurance$smoker <- as.character(insurance$smoker)
insurance$smoker[insurance$smoker == "yes"] <- "1"
insurance$smoker[insurance$smoker == "no"] <- "0"
insurance$smoker <- as.factor(insurance$smoker)
# female = 1 / male = 0
insurance$sex <- as.character(insurance$sex)
insurance$sex[insurance$sex == "female"] <- "1"
insurance$sex[insurance$sex == "male"] <- "0"
insurance$sex <- as.factor(insurance$sex)
# region / SE = 1 / SW = 0 / NE = 2 / NW = 3
insurance$region <- as.character(insurance$region)
insurance$region[insurance$region == "southwest"] <- "1"
insurance$region[insurance$region == "southeast"] <- "0"
insurance$region[insurance$region == "northeast"] <- "2"
insurance$region[insurance$region == "northwest"] <- "3"

```

```
insurance$region <- as.factor(insurance$region)
head(insurance)
```

```
##   age sex   bmi children smoker region   charges
## 1  19  1 27.900         0      1      1 16884.924
## 2  18  0 33.770         1      0      0 1725.552
## 3  28  0 33.000         3      0      0 4449.462
## 4  33  0 22.705         0      0      3 21984.471
## 5  32  0 28.880         0      0      3 3866.855
## 6  31  1 25.740         0      0      0 3756.622
```

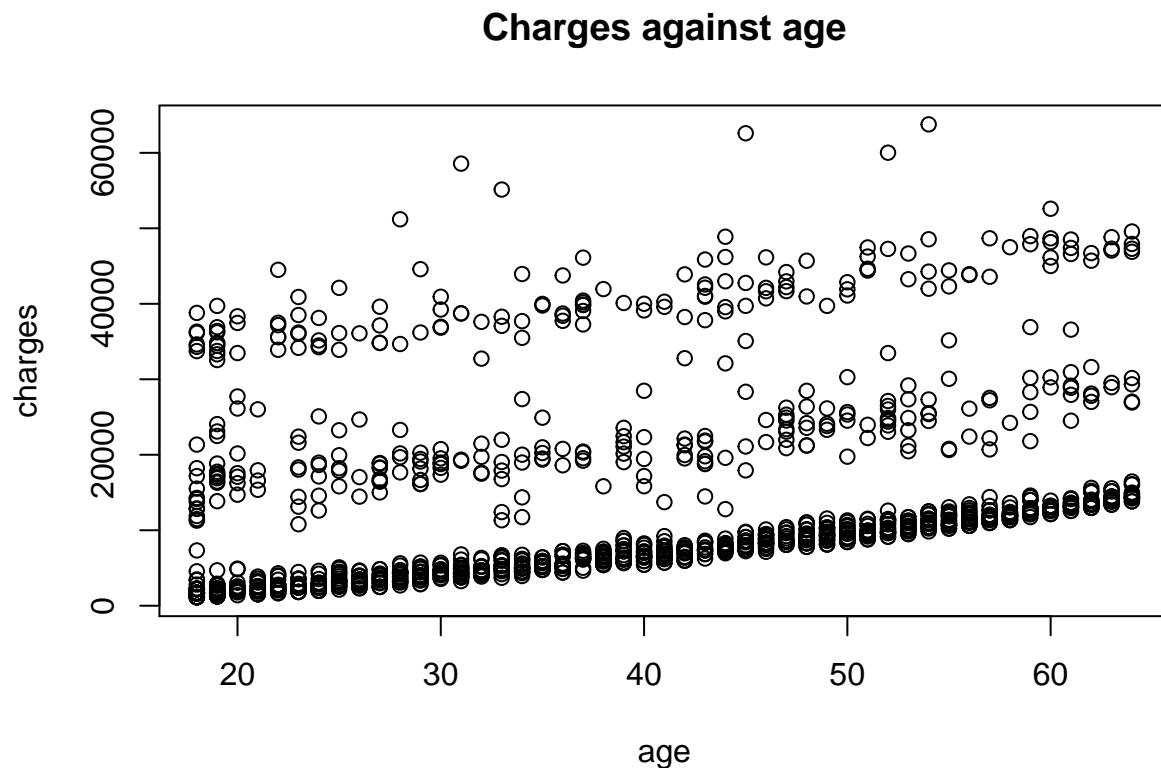
3 Linear models

The data set used for this project contains 1338 observations and seven variables containing the age, sex, bmi, children, smoker, region and charges. Those variables are covering continuous and categorical variables. The region as an categorical variable is covering four different levels: southeast (0), southwest (1), northeast (2) and northwest (3).

3.1 Basic analysis of continuous variables

We will begin with a graphical analysis. We will plot the response variable against the given predictors to gather first relationships, which can be used later in the modelling process.

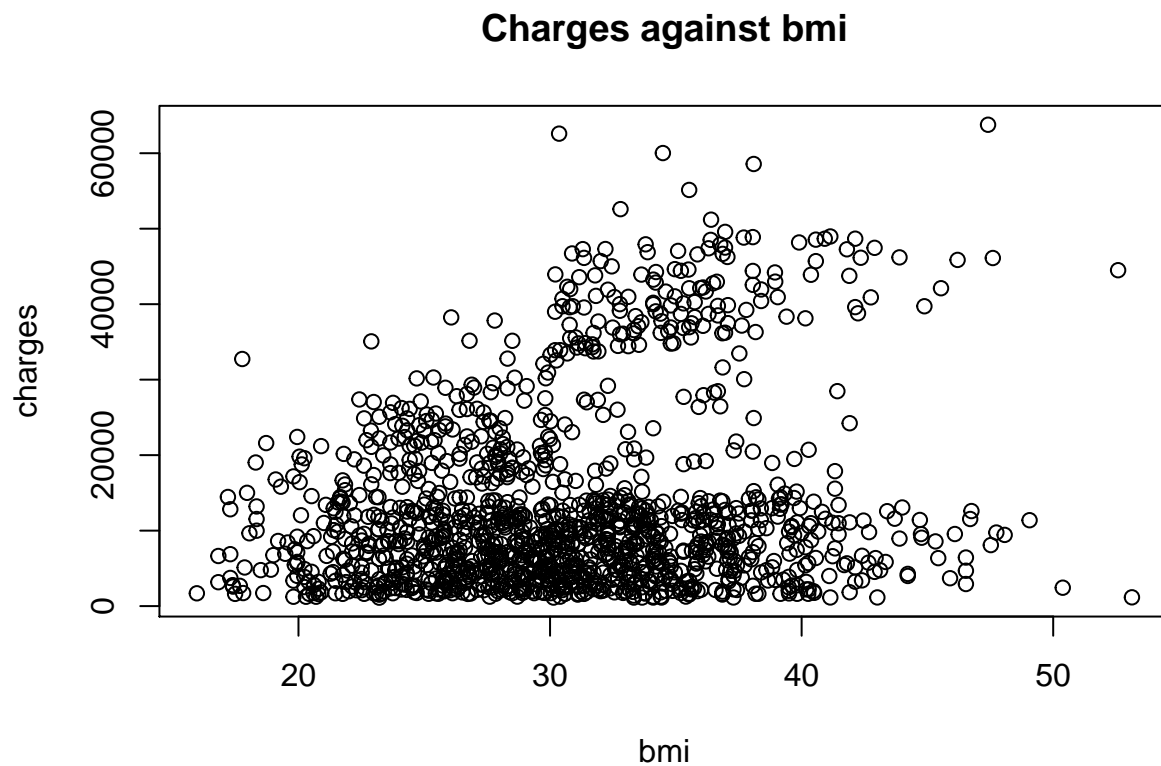
```
plot(charges ~ age, data = insurance, main = 'Charges against age')
```



There seems to be a positive relationship between age and charges as can be seen by the seemingly positive

relationship between ages and charges.

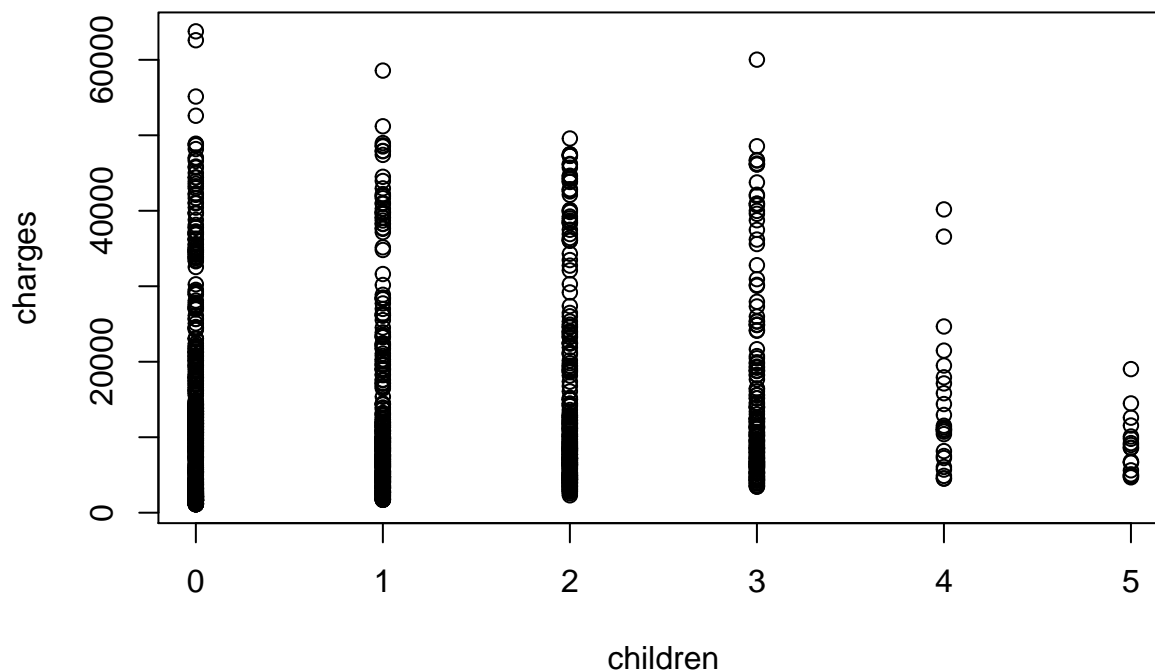
```
plot(charges ~ bmi, data =insurance, main = 'Charges against bmi')
```



The above plot is not showing a clear linear relationship between bmi and the corresponding charges.

```
plot(charges ~ children, data =insurance, main = 'Charges against children')
```

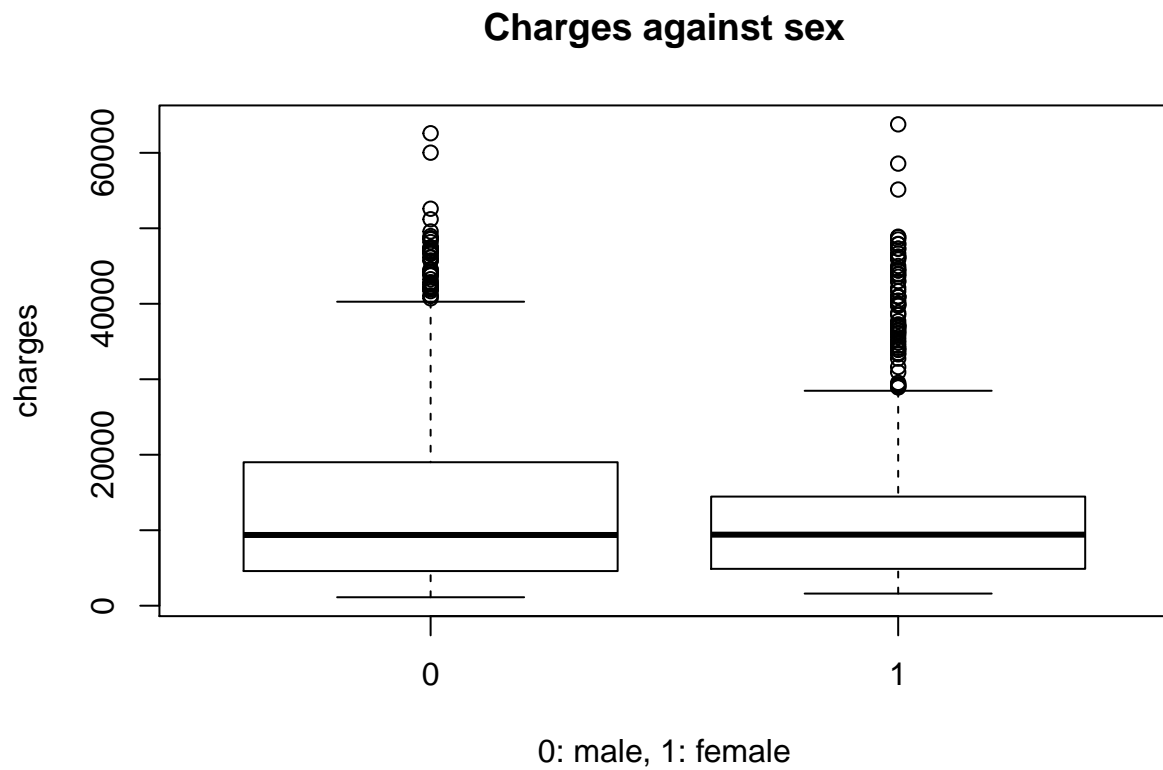
Charges against children



There might be a connection between the number of children and charges. As shown above, it might be possible to interpret that with a rising number of children the charges are decreasing. Still, this is not a clear relationship, more a wide interpretation of the plot. Not taken into consideration has been the number of participants per group, meaning that our interpretation of the plot leaves out that there are much more people with 0 children than there are with 4 or 5.

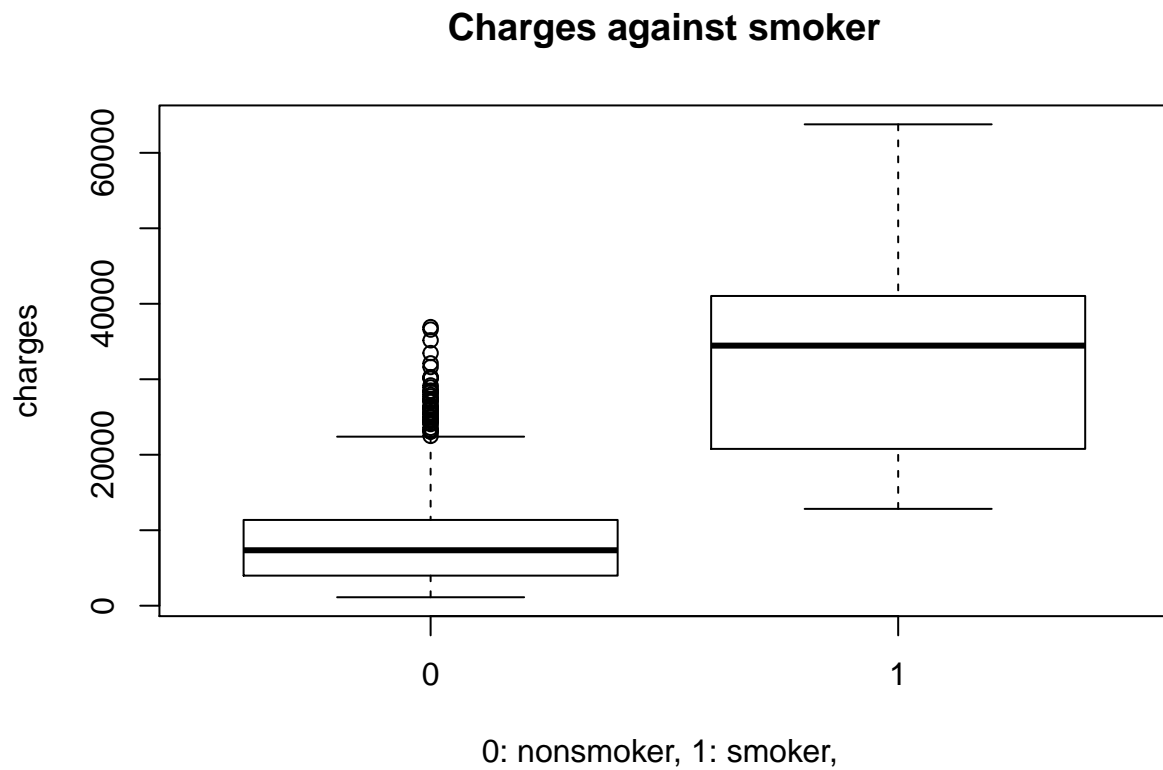
3.2 Basic analysis of categorical variables

```
plot(charges ~ sex, data = insurance, xlab = "0: male, 1: female", main = 'Charges against sex')
```



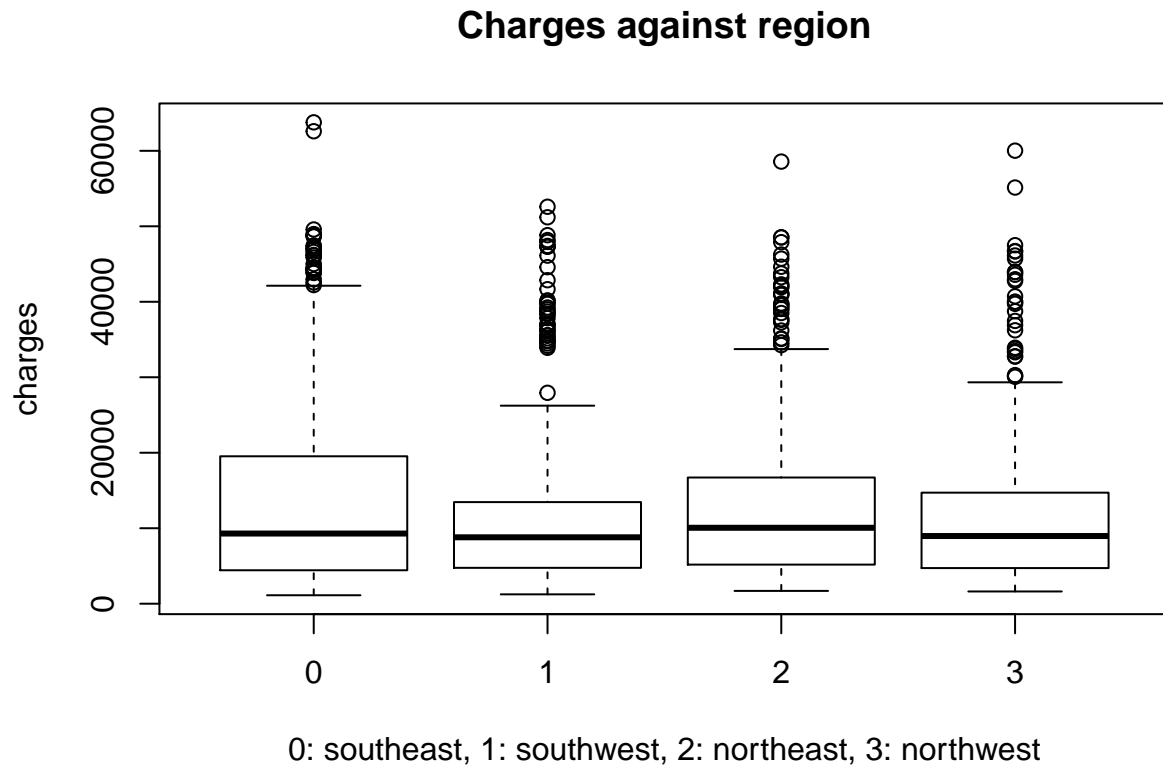
As shown above there are some differences in the charges, comparing the boxplot for the two given genders, whereas 0 stands for male, 1 stands for female. It seems that the range for 50% of the male observations is bigger than the one for the female group, meaning the distribution for males is broader than the one for females. Also the upper whisker is about 10'000 lower than the same for the male group.

```
plot(charges ~ smoker, data = insurance, xlab = "0: nonsmoker, 1: smoker, ", main = 'Charges against smoker')
```



This plot is showing a clear effect of smoking on the charges. As one can see, people who are not smoking (value = 0) have mean charges that are three times lower than the charges for smokers (value = 1). Even the outliers of the non-smokers are still within the range of the upper boxplot quartil of the smokers.

```
plot(charges ~ region, data =insurance, xlab = "0: southeast, 1: southwest, 2: northeast, 3: northwest")
```



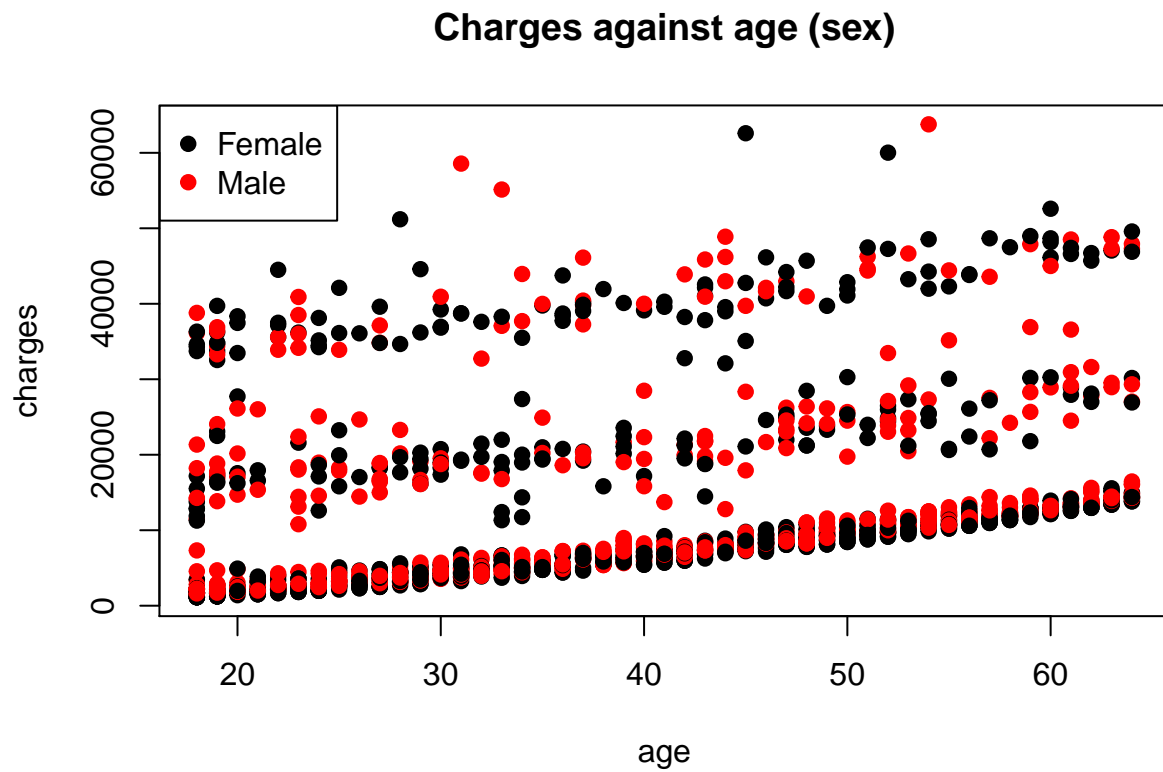
In this case, there seems not to be any clear relationship or trend visible between region and charges. Comparing the four regions with each other, it seems that there are slightly small differences between the width of the boxes, distribution of 50% of the observation as well the setting of the upper whisker. Those differences will be analysed later.

3.3 Basic analysis of categorical and continuous variables

In this part we want to show exemplary how the relationship between charges and age can be affected by adding a categorical variable. This can be used for further data explorations and can be adapted on all other variables. For this section, we will just perform this enlargement for the case of the age to illustrate some possible relationships, which also can be helpful for the following modelling process.

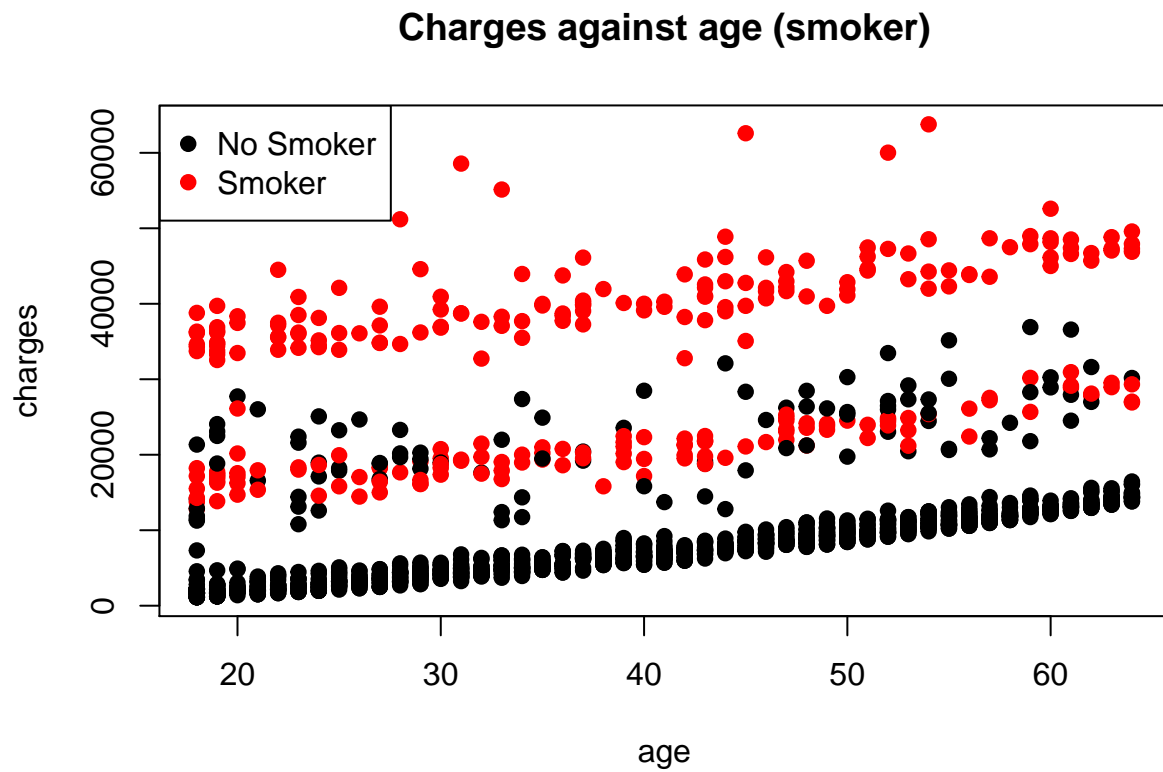
First, we add the variable “sex” to charges against age.

```
plot(charges ~ age, data = insurance,
     col = sex,
     pch = 19,
     main = "Charges against age (sex)")
legend("topleft",
     pch = 19,
     legend = c("Female", "Male"),
     col = c("black", "red"))
```

Also, we go a different approach and change the variable sex to variable smoker:

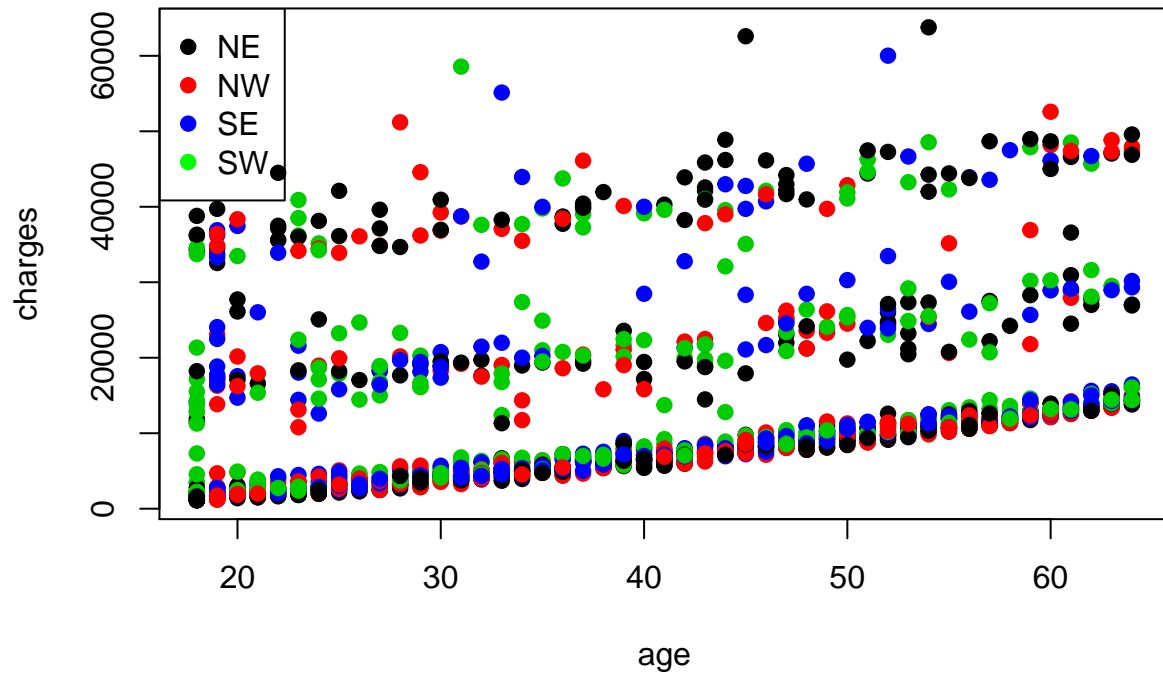
```
plot(charges ~ age, data = insurance,
     col = smoker,
     pch = 19,
     main = "Charges against age (smoker)")
legend("topleft",
     pch = 19,
     legend = c("No Smoker", "Smoker"),
     col = c("black", "red"))
```



Last but not least, we also try to fit “regions” to charges against age

```
plot(charges ~ age, data = insurance,
     col = region,
     pch = 19,
     main = "Charges against age (region)")
legend("topleft",
     pch = 19,
     legend = c("NE", "NW", "SE", "SW"),
     col = c("black", "red", "blue", "green"))
```

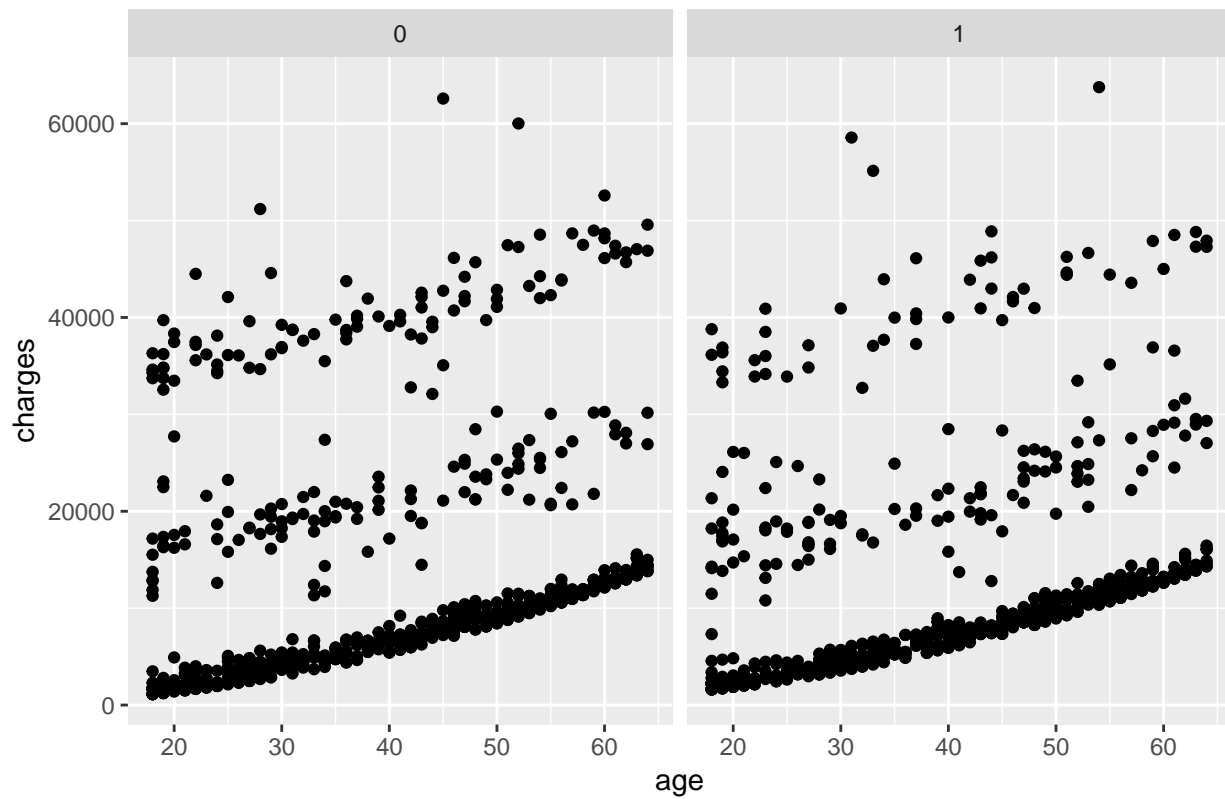
Charges against age (region)



Since it is possible to have overlapping observations, we will plot the same set-up using the `qplot`. Having separate facets will avoid overlapping and therefore might serve different results as already seen above.

```
qplot(y=charges, x=age,  
      main = "Charges against age per sex (0: male, 1: female)",  
      data = insurance,  
      facets = ~ sex  
    )
```

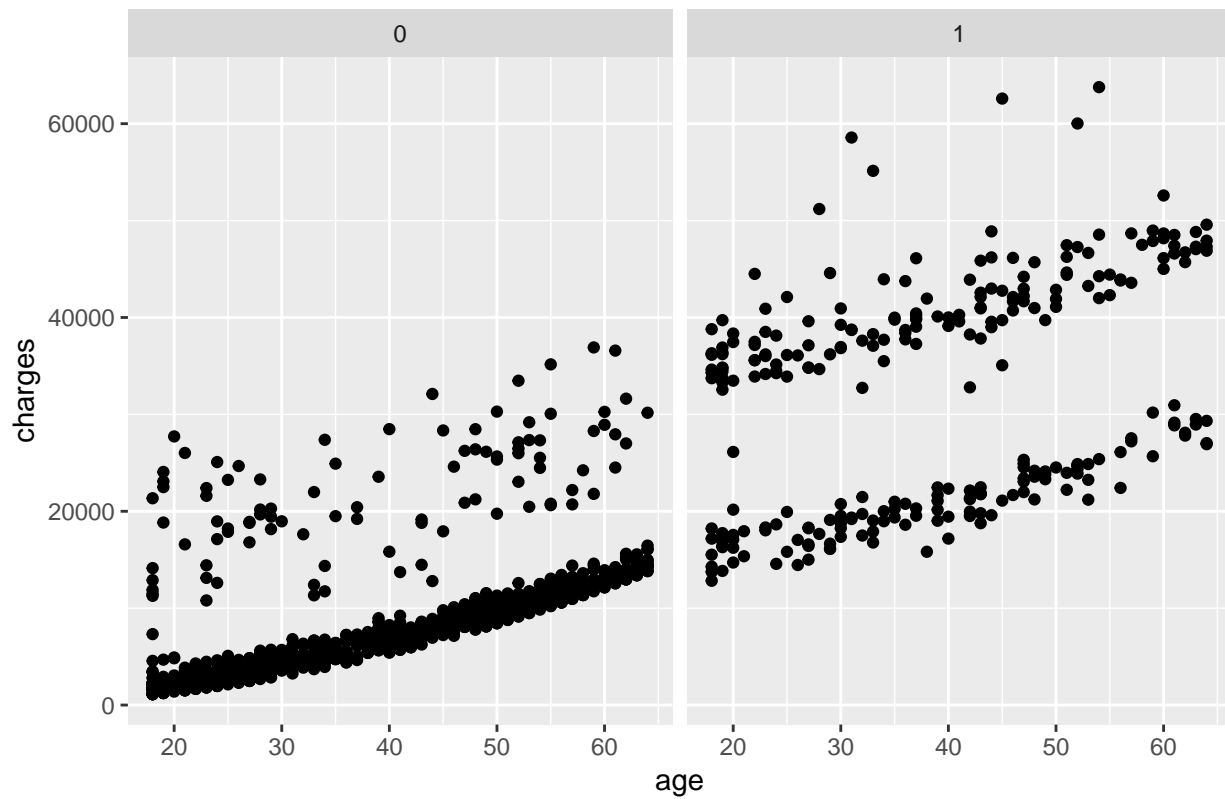
Charges against age per sex (0: male, 1: female)



```
lm.insurance <- lm(charges ~ age, data = insurance)
```

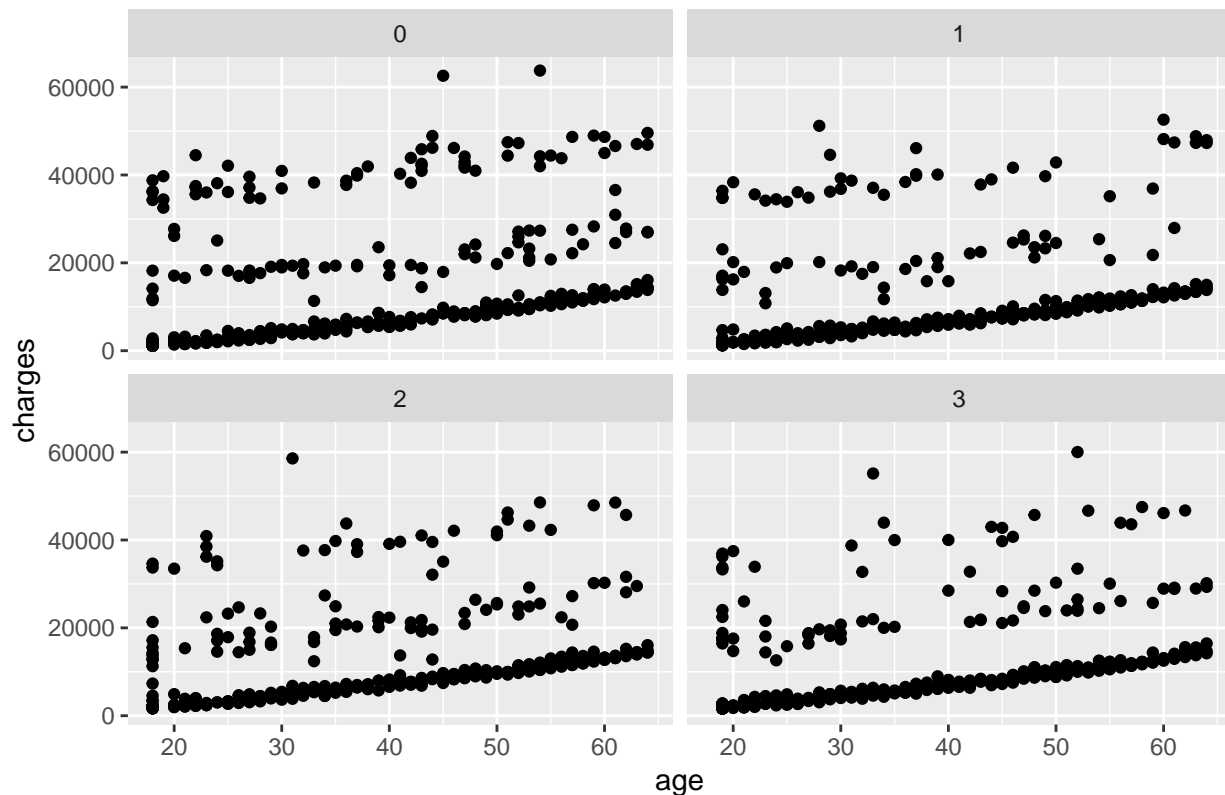
```
qplot(y=charges, x=age,  
      main = "Charges against age per smoker (0: non-smoker, 1: smoker)",  
      data = insurance,  
      facets = ~ smoker)
```

Charges against age per smoker (0: non-smoker, 1: smoker)



```
qplot(y=charges, x=age,  
      data = insurance,  
      main = "Charges against age per region (SE (0), SW (1), NE (2) and NW (3)",  
      facets = ~ region)
```

Charges against age per region (SE (0), SW (1), NE (2) and NW (3))



3.4 Fitting a first Linear Model

As usual, we will start by fitting a simple regression model to the insurance dataset. Therefore we will start with one variable and add additional complexity in each following subset. This step-by-step approach helps to explore the data slightly better and will simplify the final modelling.

3.4.1 Linear Model with one variable

For the first simple regression model we will use the age.

```
lm.insurance <- lm(charges ~ age, data = insurance)
summary(lm.insurance)
```

```
##
## Call:
## lm(formula = charges ~ age, data = insurance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8059   -6671   -5939    5440   47829
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3165.9      937.1     3.378 0.000751 ***
## age           257.7       22.5    11.453 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 11560 on 1336 degrees of freedom
## Multiple R-squared:  0.08941,    Adjusted R-squared:  0.08872
## F-statistic: 131.2 on 1 and 1336 DF,  p-value: < 2.2e-16
```

3.4.1.1 Coefficient and Interpretation

As shown in the R Output a person with the age of 0 will have a base charge of 3165.885. This intercept and its interpretation seems to be nonsensical, thus one needs to pay attention in what extend this model can be applied.

With each year, charges increase by 257.7226, which is the slope for the regression line.

3.4.1.2 P-values

With a value of 2e-16, age seems to have a very strong effect on charges. This means that the slope of charges is not flat, and the hypothesis of being zero has to be denied.

3.4.1.3 Including the gender as a second variable

In this subset, we will consider also the sex for modelling a simple regression model.

```
lm.insurance.2 <- lm(charges ~ age + sex, data = insurance)
summary(lm.insurance.2)

##
## Call:
## lm(formula = charges ~ age + sex, data = insurance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8821  -6947  -5511   5443  48203
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3882.46     980.49   3.960  7.9e-05 ***
## age          258.87      22.47  11.523 < 2e-16 ***
## sex1        -1538.83     631.08  -2.438  0.0149 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11540 on 1335 degrees of freedom
## Multiple R-squared:  0.09344,    Adjusted R-squared:  0.09209
## F-statistic: 68.8 on 2 and 1335 DF,  p-value: < 2.2e-16
```

As the output shows the additionally added variable of sex with a calculated p-value of 0.0149 seems to have an effect on charges and therefore should be considered within the modelling process.

3.4.1.4 Including interaction

Since a second variable was added, we want to explore if there might be significant interactions between age and sex, that might have to be considered for the modelling process.

```
lm.insurance.3 <-lm(charges ~ age * sex, data =insurance)
summary(lm.insurance.3)
```

```
##
## Call:
```

```

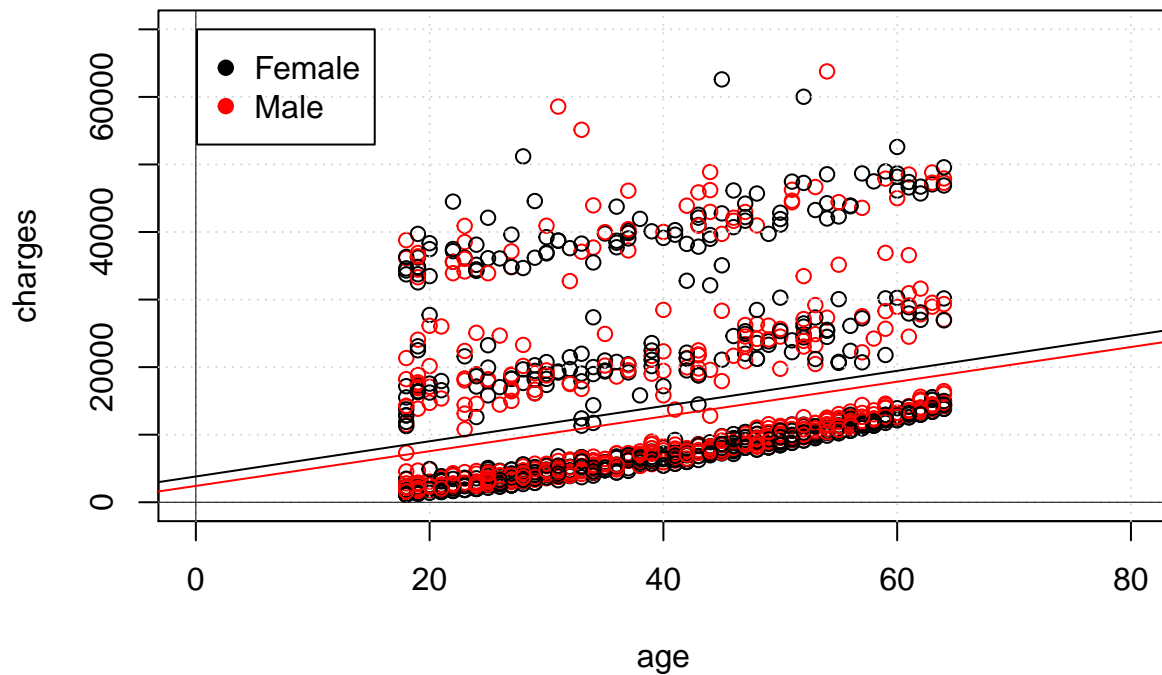
## lm(formula = charges ~ age * sex, data = insurance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8823  -6936  -5500   5456  48187
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3811.77    1308.29   2.914  0.00363 **
## age          260.68      31.62   8.244 3.96e-16 ***
## sex1        -1394.92    1872.30  -0.745  0.45638
## age:sex1      -3.67      44.95  -0.082  0.93494
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11540 on 1334 degrees of freedom
## Multiple R-squared:  0.09345,    Adjusted R-squared:  0.09141
## F-statistic: 45.84 on 3 and 1334 DF,  p-value: < 2.2e-16

plot(charges ~ age, data = insurance,
     main = "Model 'lm.insurance.3'",
     xlim = c(0, 80),
     ylim = c(0, 70000),
     col = sex)

##
grid()
##
abline(h = 0, lwd = 0.5)
abline(v = 0, lwd = 0.5)
##
abline(a = coef(lm.insurance.3)[1],
       b = coef(lm.insurance.3)["age"])
abline(a = coef(lm.insurance.3)[1] + coef(lm.insurance.3)["sex1"] ,
       b = coef(lm.insurance.3)["age"] + coef(lm.insurance.3)["age:sex1"],
       col = "red")
legend(x = 0.1, y = 70000,
       pch = 19,
       legend = c("Female", "Male"),
       col = c("black", "red"))

```


Model 'lm.insurance.3'



```
coef(lm.insurance.3)
```

```
## (Intercept)      age      sex1  age:sex1
## 3811.773852 260.681339 -1394.925331 -3.669849
```

```
summary(lm.insurance.3)$coefficients
```

```
##           Estimate Std. Error    t value    Pr(>|t|)
## (Intercept) 3811.773852 1308.29248  2.91354869 3.633010e-03
## age         260.681339   31.62249  8.24354130 3.958052e-16
## sex1        -1394.925331 1872.29663 -0.74503437 4.563822e-01
## age:sex1     -3.669849   44.95054 -0.08164194 9.349437e-01
```

The R-output shows that there isn't any significant interaction between age and sex that has to be considered in our model.

3.5 Final Linear Model Development

In this section we want to find an appropriate model, which accounts all relevant parameters and interactions. Afterwards we will then compare the fitted model with a base model and test its performance.

3.5.1 Linear Model with all variables

```
lm.insurance.all <- lm(charges ~ age + sex + bmi + children + smoker + region , data = insurance)
summary(lm.insurance.all)
```

```
##
## Call:
```

```
## lm(formula = charges ~ age + sex + bmi + children + smoker +
##     region, data = insurance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11304.9  -2848.1   -982.1   1393.9  29992.8
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -13104.87    1090.51  -12.017  < 2e-16 ***
## age          256.86      11.90    21.587  < 2e-16 ***
## sex1         131.31     332.95    0.394 0.693348
## bmi          339.19      28.60    11.860  < 2e-16 ***
## children     475.50     137.80    3.451 0.000577 ***
## smoker1     23848.53     413.15   57.723  < 2e-16 ***
## region1       74.97      470.64    0.159 0.873460
## region2     1035.02      478.69    2.162 0.030782 *
## region3       682.06      478.96    1.424 0.154669
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6062 on 1329 degrees of freedom
## Multiple R-squared:  0.7509, Adjusted R-squared:  0.7494
## F-statistic: 500.8 on 8 and 1329 DF,  p-value: < 2.2e-16
```

As the above R Output shows, not every variable seems to have a significant effect on the charges. Therefore we will separately test the variables of sex and region before dropping them from our model.

3.5.1.1 Testing sex before dropping

Before we remove the variable sex from our model we will first make a deeper analysis.

```
lm.sex <- lm(charges ~ sex, data=insurance)
summary(lm.sex)

##
## Call:
## lm(formula = charges ~ sex, data = insurance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12835  -8435  -3980   3476  51201
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  13956.8      465.2   30.003  <2e-16 ***
## sex1        -1387.2      661.3   -2.098   0.0361 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12090 on 1336 degrees of freedom
## Multiple R-squared:  0.003282, Adjusted R-squared:  0.002536
## F-statistic:  4.4 on 1 and 1336 DF,  p-value: 0.03613
coef(lm.sex)
```

```
## (Intercept)      sex1
## 13956.751    -1387.172
```

As seen in the output considering only sex is not a significant and standalone explaining variable for charges. Comparing the p-value of sex within the full model, including all possible variables, it is so high with a value of 0.693348 it can be dropped from our final model.

3.5.1.2 Testing region before dropping

```
lm.region.1 <- lm(charges ~ region, data = insurance)
aggregate(charges ~ region, FUN = mean, data = insurance)
```

```
##   region  charges
## 1      0 14735.41
## 2      1 12346.94
## 3      2 13406.38
## 4      3 12417.58
```

```
summary(lm.region.1)
```

```
##
## Call:
## lm(formula = charges ~ region, data = insurance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13614  -8463  -3793   3385  49035
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  14735.4      633.3   23.266  <2e-16 ***
## region1      -2388.5      922.2   -2.590   0.0097 **
## region2      -1329.0      922.9   -1.440   0.1501
## region3      -2317.8      922.2   -2.513   0.0121 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12080 on 1334 degrees of freedom
## Multiple R-squared:  0.006634,    Adjusted R-squared:  0.0044
## F-statistic:  2.97 on 3 and 1334 DF,  p-value: 0.03089
```

There is strong evidence that the mean charges for northeast is not equal to zero. But there isn't any evidence that all other regions differ from the reference region northeast. To have a better understanding we will make an anova test between the above shown model and a base model `lm.region.0` as shown below.

```
lm.region.0 <- lm(charges ~ 1 , data = insurance)
coef(lm.region.0)
```

```
## (Intercept)
## 13270.42
```

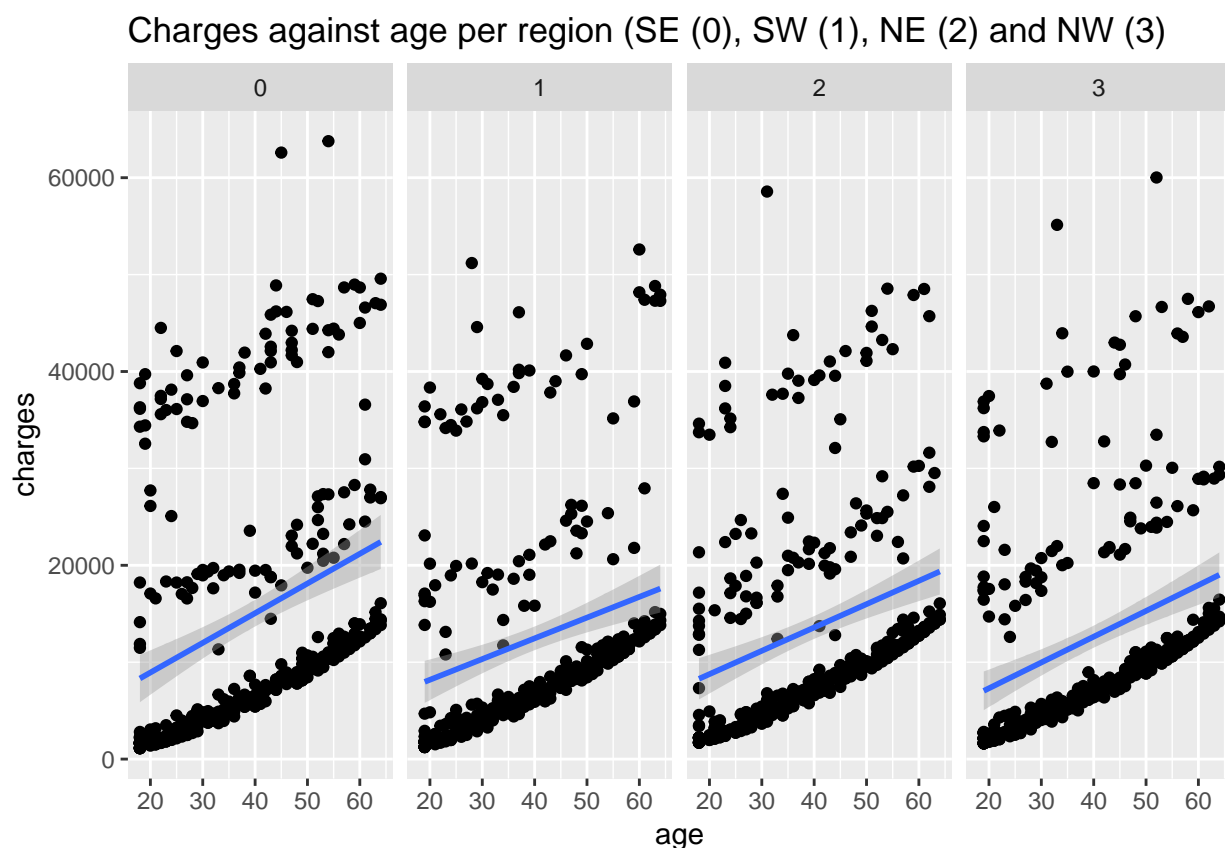
```
anova(lm.region.0, lm.region.1)
```

```
## Analysis of Variance Table
##
## Model 1: charges ~ 1
## Model 2: charges ~ region
##   Res.Df        RSS Df Sum of Sq    F Pr(>F)
```

```
## 1    1337 1.9607e+11
## 2    1334 1.9477e+11  3 1300759681 2.9696 0.03089 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The anova test shows low evidence that the model with more parameters (in this case only the region) fits better to the data. The F-value seems to be very small and also the p-value is not highly significant with 0.03089. Anyhow there is a drop in the RSS, by adding 3 additional parameters of the different regions. To have a better understanding, we can additionally perform posthoc contrasts to decide afterwards if we will drop region finally from our model. Before performing several posthoc test and repeating this exercise to all possible combinations we will use the ggplot to explore the data quickly upfront.

```
ggplot(data=insurance,
       mapping = aes(y = charges, x= age))+ geom_point()+geom_smooth(method = 'lm') + facet_grid(. ~region)
ggtitle("Charges against age per region (SE (0), SW (1), NE (2) and NW (3))")
```



Using the visualization, it seems that there seems to be no obvious difference and effect. Therefore, we decide to drop also the variable region from our linear model.

3.5.2 Dropping variables from the model

```
drop1(lm.insurance.all, test="F")

## Single term deletions
##
## Model:
## charges ~ age + sex + bmi + children + smoker + region
##      Df Sum of Sq      RSS      AIC    F value    Pr(>F)
```

```
## <none>          4.8840e+10 23316
## age            1 1.7124e+10 6.5964e+10 23717 465.9837 < 2.2e-16 ***
## sex            1 5.7164e+06 4.8845e+10 23315 0.1556 0.693348
## bmi            1 5.1692e+09 5.4009e+10 23449 140.6627 < 2.2e-16 ***
## children       1 4.3755e+08 4.9277e+10 23326 11.9063 0.000577 ***
## smoker         1 1.2245e+11 1.7129e+11 24993 3331.9680 < 2.2e-16 ***
## region         3 2.3343e+08 4.9073e+10 23317 2.1173 0.096221 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

lm.insurance.1 <- update(lm.insurance.all, .~. -region -sex)
formula(lm.insurance.1)

## charges ~ age + bmi + children + smoker

summary(lm.insurance.1)

##
## Call:
## lm(formula = charges ~ age + bmi + children + smoker, data = insurance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11897.9  -2920.8   -986.6   1392.2  29509.6
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -12102.77     941.98  -12.848 < 2e-16 ***
## age          257.85       11.90   21.675 < 2e-16 ***
## bmi          321.85       27.38   11.756 < 2e-16 ***
## children     473.50      137.79    3.436 0.000608 ***
## smoker1     23811.40     411.22   57.904 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6068 on 1333 degrees of freedom
## Multiple R-squared:  0.7497, Adjusted R-squared:  0.7489
## F-statistic: 998.1 on 4 and 1333 DF,  p-value: < 2.2e-16
```

Comparing the summary output of `lm.insurance.0` and `lm.insurance.1`, the latest model is only including variables with a p-value indicating a significant factor.

3.5.3 Considering Interactions

```
drop1(lm.insurance.1, test="F")

## Single term deletions
##
## Model:
## charges ~ age + bmi + children + smoker
##              Df Sum of Sq      RSS   AIC  F value    Pr(>F)
## <none>          4.9078e+10 23315
## age            1 1.7297e+10 6.6375e+10 23717 469.789 < 2.2e-16 ***
## bmi            1 5.0884e+09 5.4167e+10 23445 138.203 < 2.2e-16 ***
## children       1 4.3477e+08 4.9513e+10 23325 11.809 0.0006077 ***
## smoker         1 1.2345e+11 1.7253e+11 24995 3352.911 < 2.2e-16 ***
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

lm.insurance.2 <- update(lm.insurance.1, .~. +age:bmi + age:children + age:smoker + bmi:children + bmi:smoker)
summary(lm.insurance.2)

##
## Call:
## lm(formula = charges ~ age + bmi + children + smoker + age:bmi +
##     age:children + age:smoker + bmi:children + bmi:smoker + children:smoker,
##     data = insurance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13996.3  -1947.6  -1331.5   -406.4   29570.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -6.745e+02  2.106e+03  -0.320    0.749
## age             2.055e+02  4.963e+01   4.140 3.7e-05 ***
## bmi            -6.339e+01  6.756e+01  -0.938    0.348
## children        6.957e+02  6.341e+02   1.097    0.273
## smoker1       -1.983e+04  1.861e+03 -10.651 < 2e-16 ***
## age:bmi         1.912e+00  1.561e+00   1.225    0.221
## age:children    1.201e+00  8.527e+00   0.141    0.888
## age:smoker1     -9.141e-01  2.384e+01  -0.038    0.969
## bmi:children    -5.334e+00  1.863e+01  -0.286    0.775
## bmi:smoker1     1.437e+03  5.317e+01  27.029 < 2e-16 ***
## children:smoker1 -3.858e+02  2.841e+02  -1.358    0.175
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4874 on 1327 degrees of freedom
## Multiple R-squared:  0.8392, Adjusted R-squared:  0.838
## F-statistic: 692.8 on 10 and 1327 DF,  p-value: < 2.2e-16
```

As the above output shows, not all of the added interactions have to be considered in the model. Upfront we also tried to generate a linear model including all possible interactions. Since the output was not satisfying we skipped this analysis at this point. Therefore, the working assumption in this step was quite easy by editing only possible and simple interactions. Based on the results we will now drop the unnecessary interactions.

```
drop1(lm.insurance.2, test="F")
```

```
## Single term deletions
##
## Model:
## charges ~ age + bmi + children + smoker + age:bmi + age:children +
##     age:smoker + bmi:children + bmi:smoker + children:smoker
##              Df Sum of Sq      RSS   AIC  F value Pr(>F)
## <none>                 3.1519e+10 22735
## age:bmi                1 3.5624e+07 3.1555e+10 22734   1.4998 0.2209
## age:children            1 4.7151e+05 3.1520e+10 22733   0.0199 0.8880
## age:smoker              1 3.4909e+04 3.1519e+10 22733   0.0015 0.9694
## bmi:children            1 1.9467e+06 3.1521e+10 22733   0.0820 0.7747
## bmi:smoker              1 1.7353e+10 4.8872e+10 23319 730.5705 <2e-16 ***
## children:smoker         1 4.3796e+07 3.1563e+10 22734   1.8439 0.1747
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

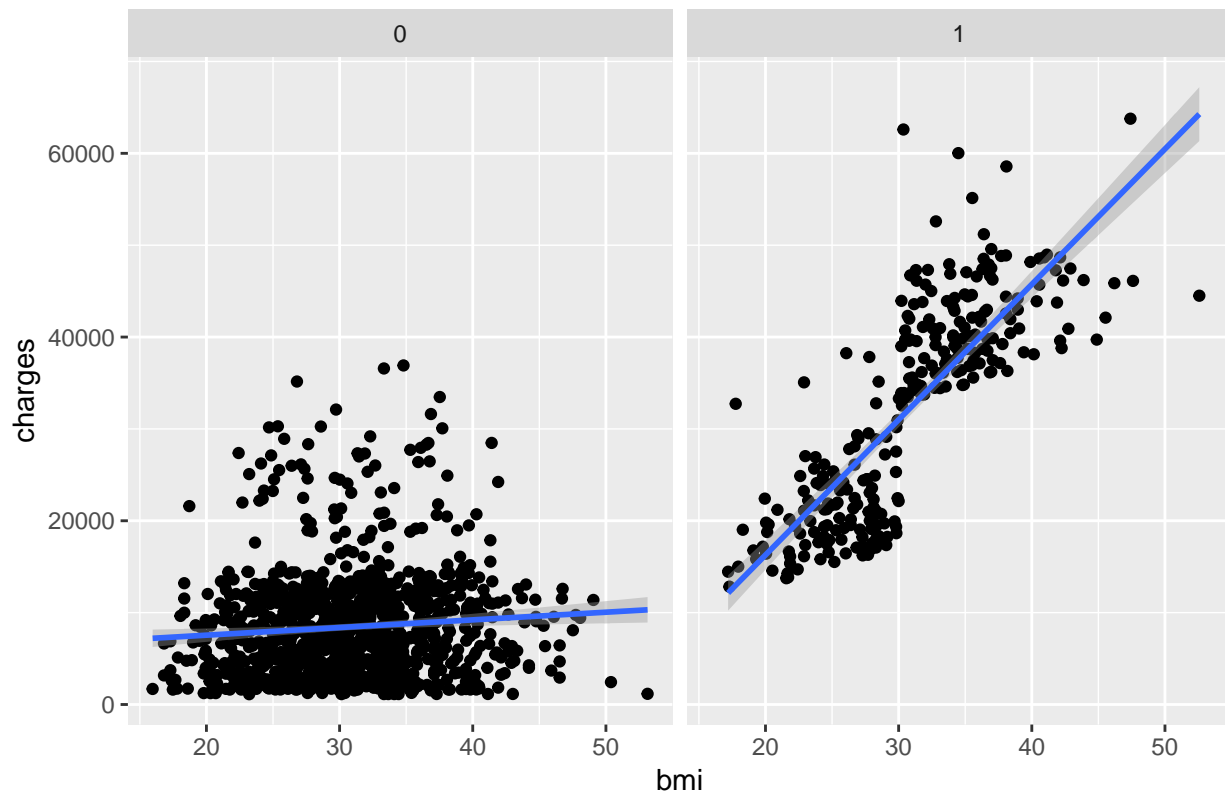
lm.insurance.3 <- update(lm.insurance.2, .~. -age:bmi - age:children - age:smoker - bmi:children -children
summary(lm.insurance.3)

##
## Call:
## lm(formula = charges ~ age + bmi + children + smoker + bmi:smoker,
##     data = insurance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -14598.6  -1924.4  -1321.4   -465.6   29892.4
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -2729.002     831.270  -3.283  0.00105 **
## age           264.948       9.553   27.735 < 2e-16 ***
## bmi           5.656       24.873    0.227  0.82014
## children     508.924     110.615    4.601 4.61e-06 ***
## smoker1     -20194.709    1654.505  -12.206 < 2e-16 ***
## bmi:smoker1   1433.788      52.823   27.143 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4871 on 1332 degrees of freedom
## Multiple R-squared:  0.8388, Adjusted R-squared:  0.8382
## F-statistic: 1387 on 5 and 1332 DF,  p-value: < 2.2e-16
```

After considering also the interaction between bmi and smoker, which seems to be significant with a small p-value of 2e-16. The variable bmi itself has now a p-value of 0.82014. Therefore let's have a look on the relationship between bmi and smoker to have a better understanding:

```
ggplot(data=insurance,
       mapping = aes(y = charges, x= bmi))+ geom_point()+ geom_smooth(method = 'lm') +facet_grid(. ~ smoker)
ggtitle("Charges against bmi per smoker (0: non-smoker, 1: smoker)")
```

Charges against bmi per smoker (0: non-smoker, 1: smoker)



As expected, there seems to be a clear relationship between charges and bmi, considering if a person is smoking or not. Since this relationship makes even sense from the domain perspective, we will definitely keep this in our model.

3.5.4 Measure of Fit

In this section, we were able to fit several models, having different levels of complexity. For the following analysis we have chosen our four models to measure the individual fit. Moreover we will compare them among each other to find the most appropriate model for the given insurance dataset.

```
formula(lm.insurance.all)
```

```
## charges ~ age + sex + bmi + children + smoker + region
```

```
summary(lm.insurance.all)$r.squared
```

```
## [1] 0.750913
```

```
summary(lm.insurance.all)$adj.r.squared
```

```
## [1] 0.7494136
```

```
formula(lm.insurance.1)
```

```
## charges ~ age + bmi + children + smoker
```

```
summary(lm.insurance.1)$r.squared
```

```
## [1] 0.7496945
```



```
summary(lm.insurance.1)$adj.r.squared

## [1] 0.7489434

formula(lm.insurance.2)

## charges ~ age + bmi + children + smoker + age:bmi + age:children +
##      age:smoker + bmi:children + bmi:smoker + children:smoker

summary(lm.insurance.2)$r.squared

## [1] 0.8392491

summary(lm.insurance.2)$adj.r.squared

## [1] 0.8380378

formula(lm.insurance.3)

## charges ~ age + bmi + children + smoker + bmi:smoker

summary(lm.insurance.3)$r.squared

## [1] 0.8388379

summary(lm.insurance.3)$adj.r.squared

## [1] 0.8382329
```

For the comparison we used R-Squared and the adjusted R-Squared to measure the performance of our models. Since the adjusted R-squared can provide a more precise view of that correlation by also taking into account how many independent variables are added to our particular models against we will base our conclusion on this parameter.

Therefore we are happy to state that the latest model number three is able to explain the charges with a percentage of 83.82329 % based on the independent variables.

Comparing the latest model with the first one there is an increase of 8.88% in the adjusted R-squared. Even if the latest model is performing the best compared to the others, it is always a trade off between the gain in the fit and the corresponding effort.

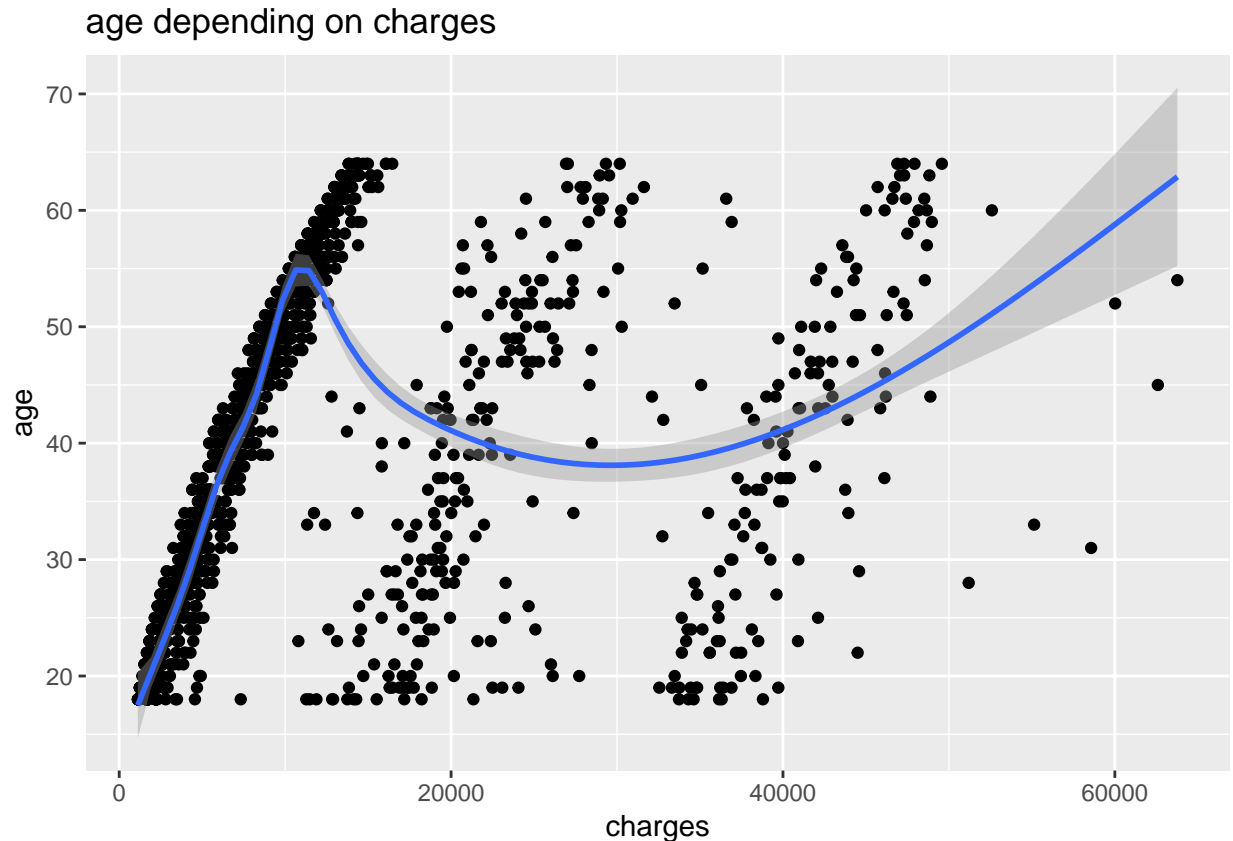
4 Polynomial models

As we can see from the following plot, the models are not always linear.

```
gg.age.charges <- ggplot(data = insurance,
                        mapping = aes (y = age,
                                      x = charges)) +

  geom_point()

gg.age.charges + geom_smooth() + ggtitle("age depending on charges")
```



Clearly visible is the increase in costs up to about 18'000 with increasing age. The dependence of amount and costs above 10,000 is more difficult to predict due to the non-linearity.

We add linearity to the model by adding a square term to charges.

```
# Model with a linear effect for charges
lm.insurance.4 <- lm(age ~ sex + bmi + children + charges, data = insurance)

#Model with a quadratic effect for charges
lm.insurance.5 <- update(lm.insurance.4, . ~ . + I(charges^2))
```

The F-test is used to test the quadrate term.

```
anova(lm.insurance.4, lm.insurance.5)

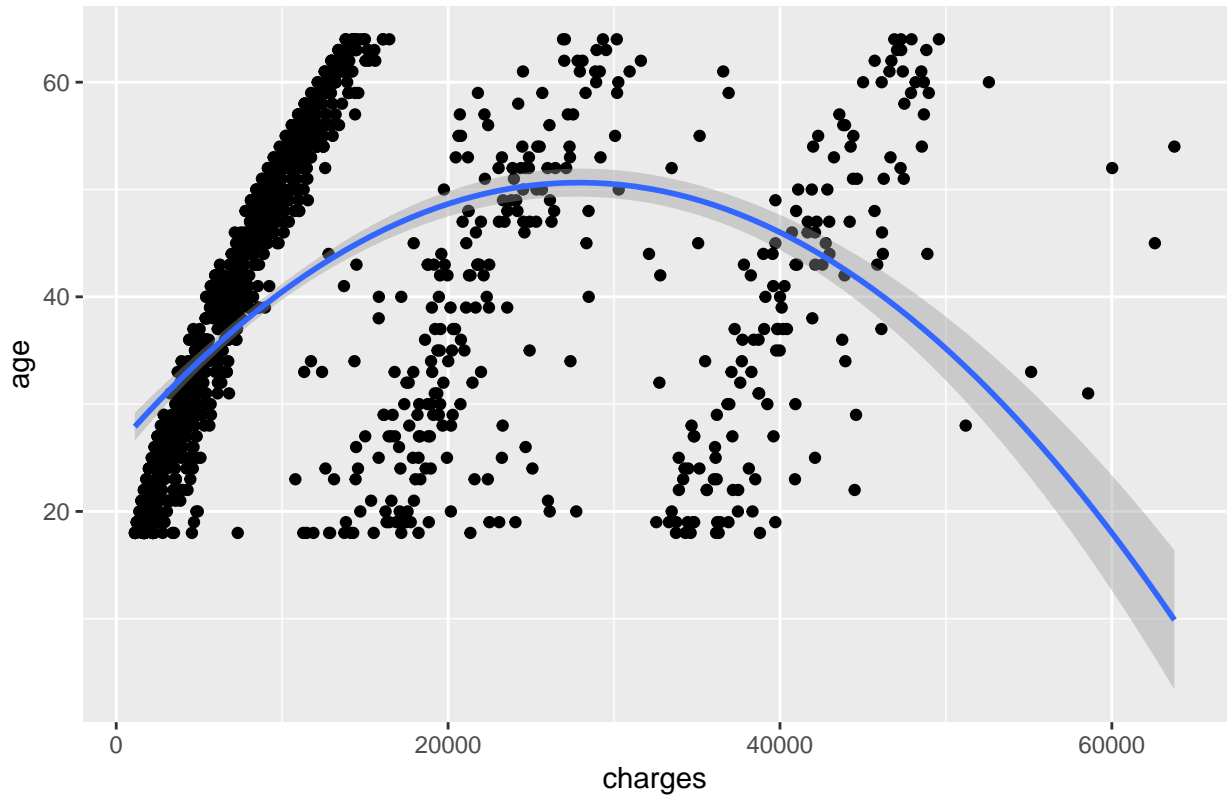
## Analysis of Variance Table
##
## Model 1: age ~ sex + bmi + children + charges
## Model 2: age ~ sex + bmi + children + charges + I(charges^2)
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1    1333 239086
## 2    1332 198860  1    40227 269.45 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The F-test shows a strong indication charges requires a quadratic term.

The graphical result is shown below.

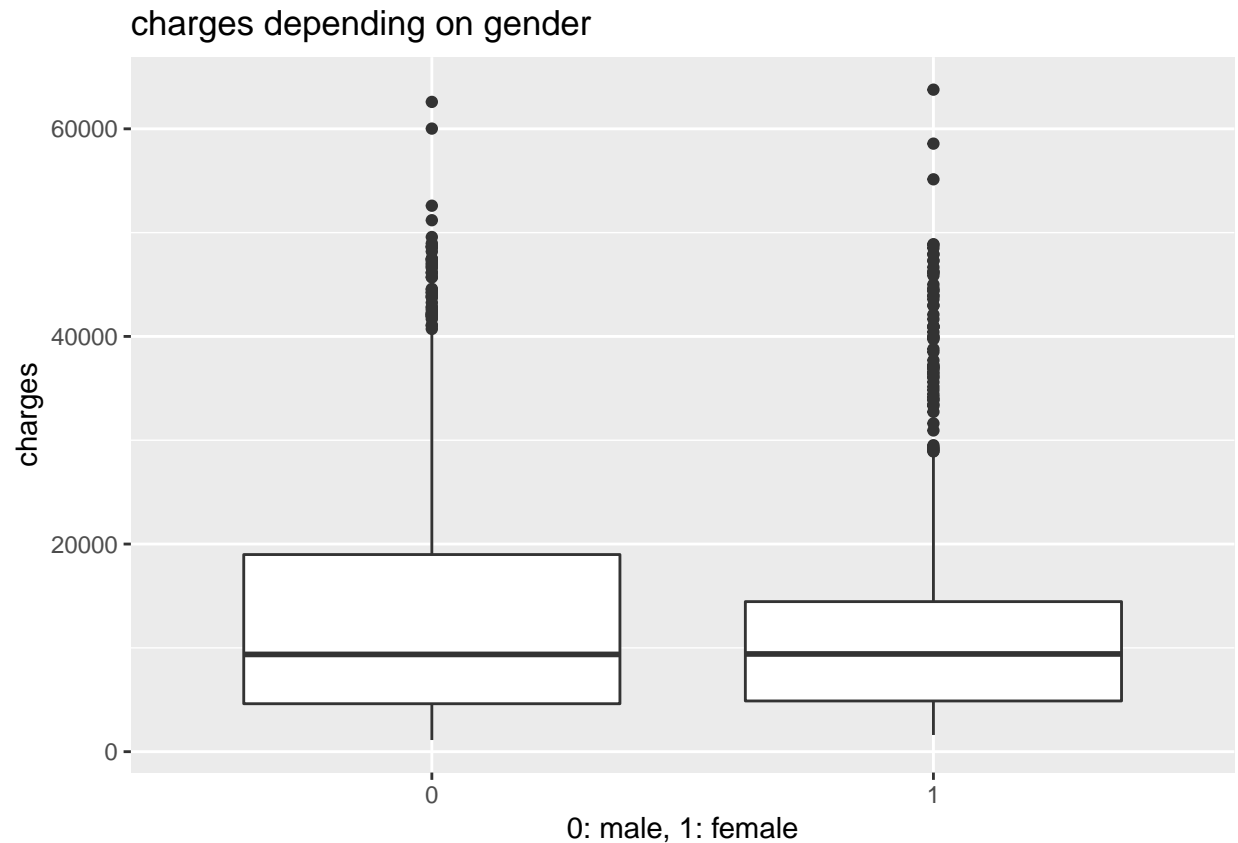
```
gg.age.charges +
  geom_smooth(method = "lm",
             formula = y ~ poly(x, degree = 2)) +
  ggtitle("age depending on charges with quadratic term")
```

age depending on charges with quadratic term



What is the result if you look at the men and women separately.

```
ggplot(data = insurance,
       mapping = aes(y = charges,
                     x = sex)) +
  geom_boxplot() + ggtitle("charges depending on gender") + xlab("0: male, 1: female")
```

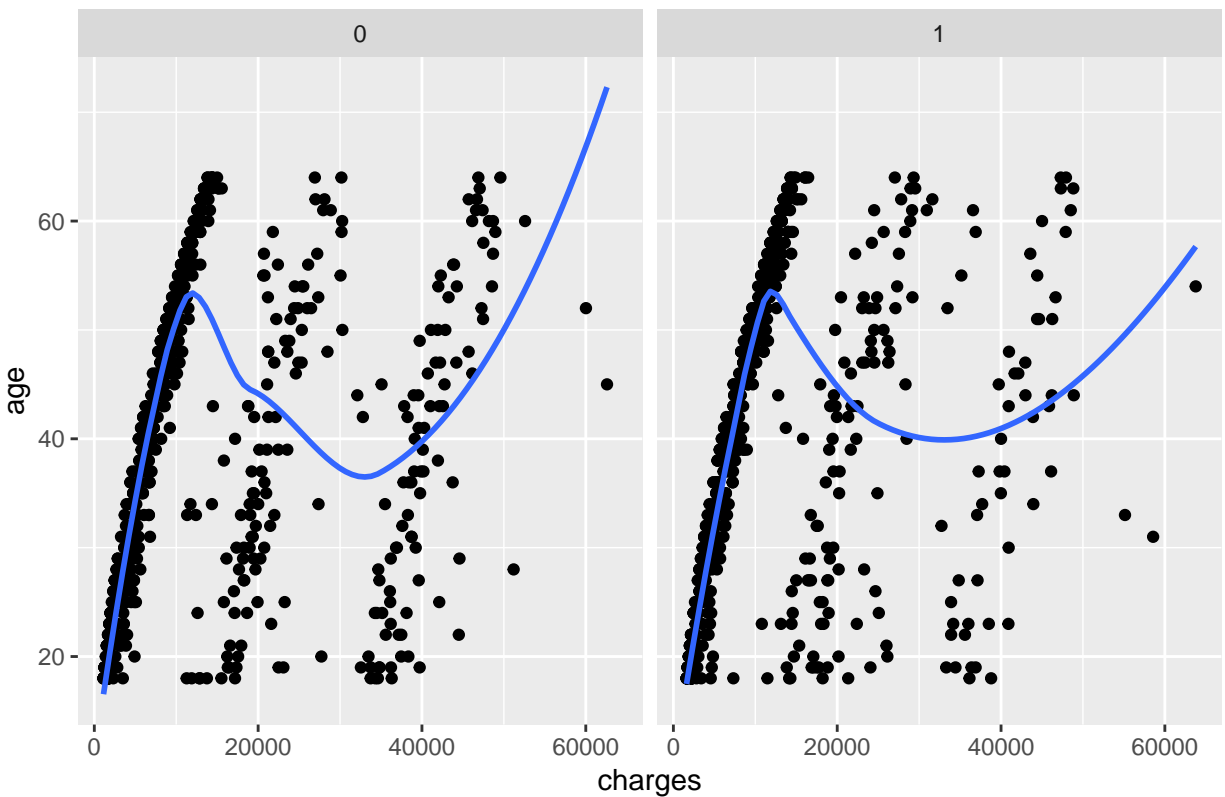


This graph shows that there is not really a difference in the median between women and men. The only difference is in the middle 50% range where the men have a larger range.

```
ggplot(data = insurance,
       mapping = aes(y = age,
                     x = charges)) +
  geom_point() +
  geom_smooth(se = FALSE) +
  facet_wrap(. ~ sex) + ggtitle("charges depending on gender (0: male, 1: female)")
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

charges depending on gender (0: male, 1: female)



The rough course of costs as a function of age is comparable for men and women. A closer look at the result reveals slight differences. However, this is not a reason to differentiate the models according to gender.

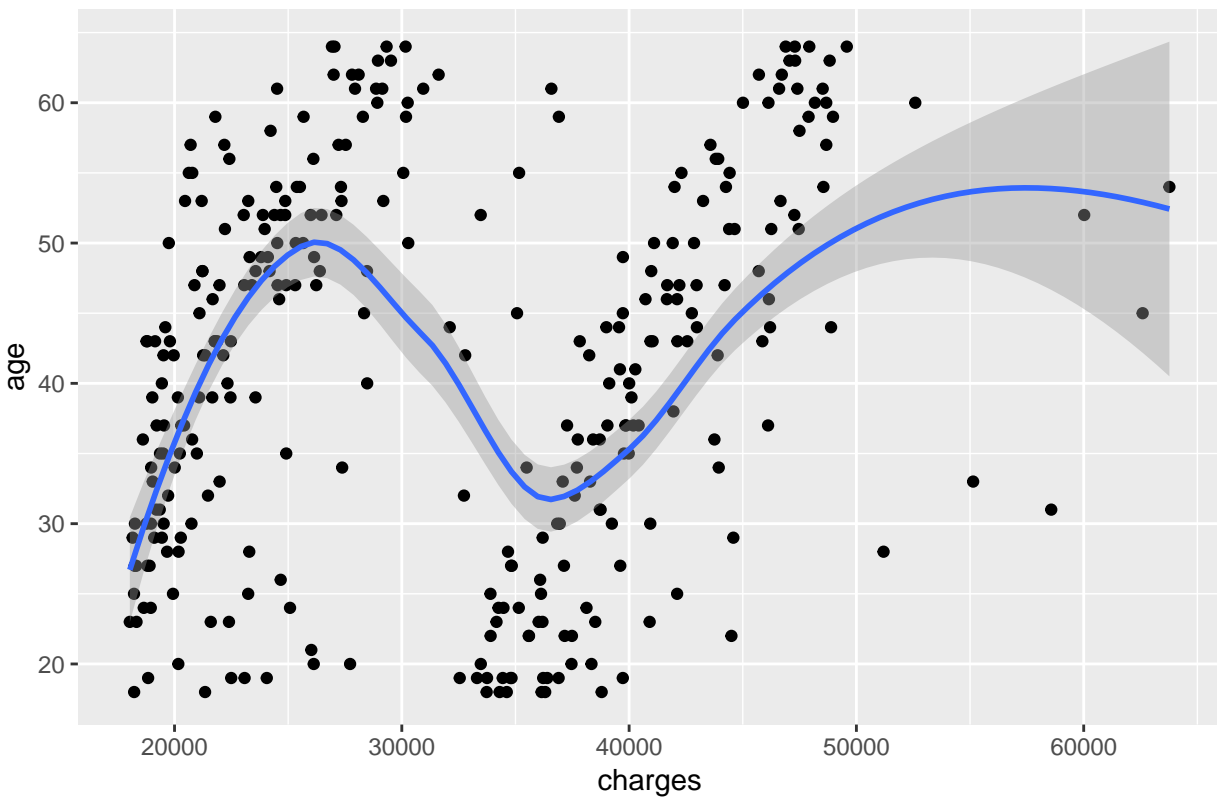
The first cost increase looks more like a linear model. What does the result look like if the costs are brought above 18'000? This will be examined in the following considerations.

```
insurance.part2 <- filter(insurance, charges > 18000 )
```

```
gg.age.charges.part2 <- ggplot(data = insurance.part2,
                               mapping = aes (y = age,
                                                x = charges)) +
  geom_point()
gg.age.charges.part2 + geom_smooth() + ggtitle("age depending on charges over 18'000")
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

age depending on charges over 18'000



```
lm.insurance.11 <- lm(age ~ sex + bmi + children + charges, data = insurance.part2)
lm.insurance.12 <- update(lm.insurance.11, . ~ . + I(charges^2))

anova(lm.insurance.11, lm.insurance.12)
```

```
## Analysis of Variance Table
##
## Model 1: age ~ sex + bmi + children + charges
## Model 2: age ~ sex + bmi + children + charges + I(charges^2)
##   Res.Df  RSS Df Sum of Sq    F Pr(>F)
## 1     308 53277
## 2     307 52872   1    405.12 2.3523 0.1261
```

The result shows that there is no strong evidence that charges needs a quadratic term. The model appears to be linear.

5 Generalised Additive Models

The previous data modelling is now performed with GAM.

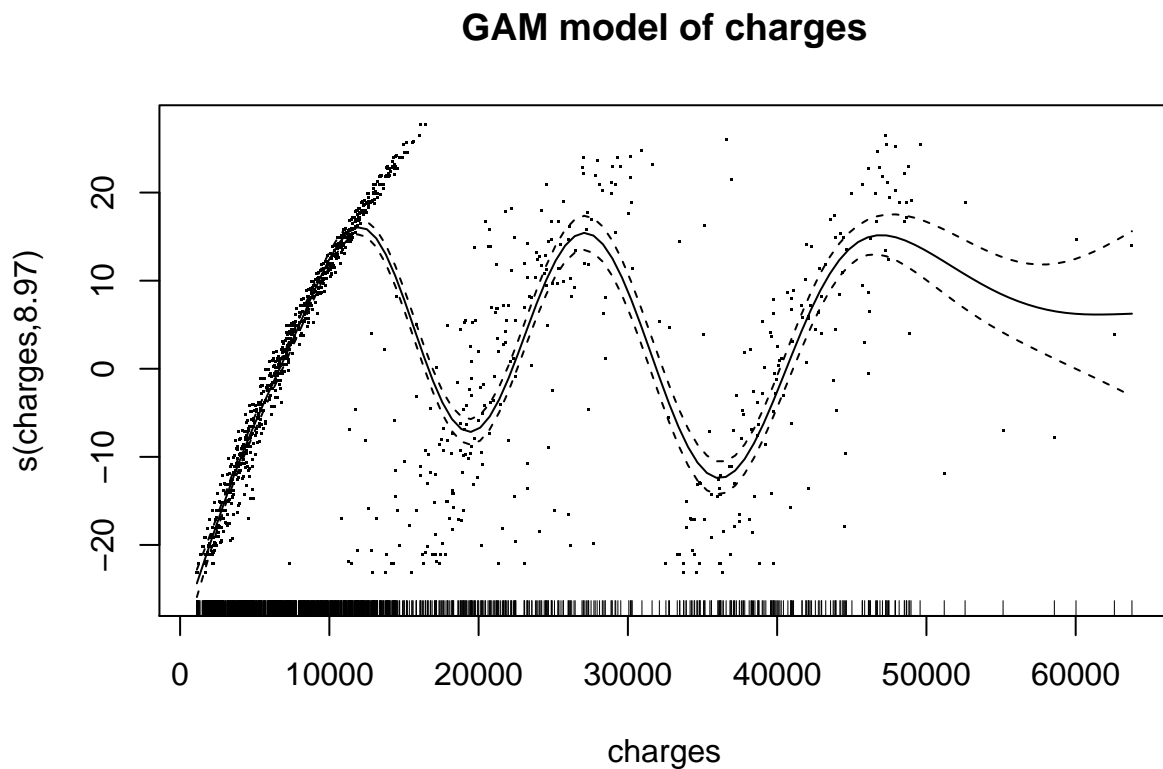
```
gam.insurance.1 <- gam(age ~ sex + s(charges) + children,
  data = insurance)
summary(gam.insurance.1)
```

```
##
## Family: gaussian
## Link function: identity
```

```
##
## Formula:
## age ~ sex + s(charges) + children
##
## Parametric coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  41.1364    0.3598 114.330 < 2e-16 ***
## sex1         -1.0849    0.4186  -2.591  0.00966 **
## children     -1.2719    0.1836  -6.926  6.72e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##             edf Ref.df    F p-value
## s(charges)  8.974     9 363.9 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.71   Deviance explained = 71.3%
## GCV = 57.673   Scale est. = 57.157    n = 1338
```

The $s(\text{charges})$ value of 8,974 is high and allows us to model the existing predictor with a smooth term. The corresponding visualization is shown below.

```
plot(gam.insurance.1, residuals = TRUE, select = 1, main = "GAM model of charges")
```



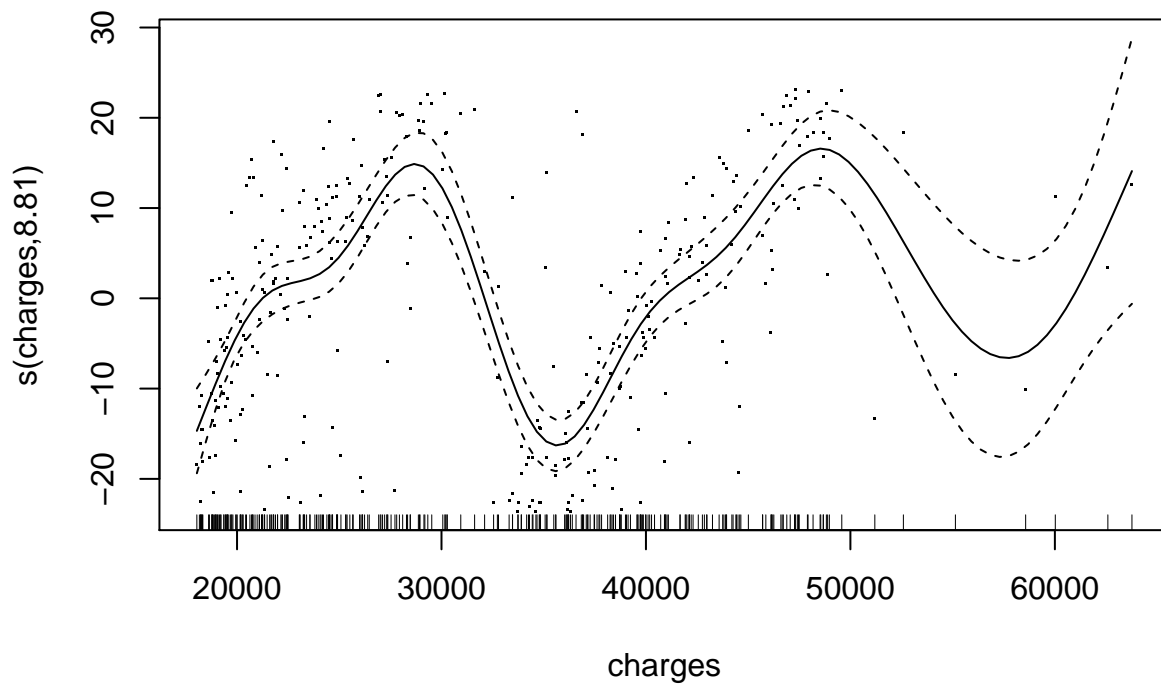
If you only consider the costs greater than 18'000, the result is identical with the previous consideration. A

linear model is also preferable in this case.

```
gam.insurance.2 <- gam(age ~ sex + s(charges) + children,
                        data = insurance.part2)
summary(gam.insurance.2)

##
## Family: gaussian
## Link function: identity
##
## Formula:
## age ~ sex + s(charges) + children
##
## Parametric coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  41.5861    0.9591  43.360  <2e-16 ***
## sex1         -0.1992    1.1457  -0.174    0.862
## children     -0.2804    0.4936  -0.568    0.570
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##              edf Ref.df      F p-value
## s(charges)  8.815   8.99 32.15  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.474   Deviance explained = 49.2%
## GCV = 101.7   Scale est. = 97.857      n = 313
plot(gam.insurance.2, residuals = TRUE, select = 1, main = "GAM model of charges over 18'000")
```


GAM model of charges over 18'000



6 Generalised Linear Models

6.1 Generalised Linear Models for count data

6.1.1 Original data

The number of children an insured person has is analysed. We have the following data on children per person. The number of children ranges from 0 to 5 with a median of 1.

```
summary(insurance$children)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  0.000   0.000   1.000   1.095   2.000   5.000
```

6.1.2 Poisson model

To model count data (number of children) the poisson model is used. An analysis performed beforehand showed that only the variables “charges” and “smoker” have a significant impact on the number of children.

```
glm.children <- glm(children ~ smoker + charges,
                    data=insurance,
                    family = "poisson")
summary(glm.children)
```

```
##
## Call:
## glm(formula = children ~ smoker + charges, family = "poisson",
##      data = insurance)
```

```
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8561  -1.4318  -0.1057   0.7768   2.9717
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -3.706e-02  4.213e-02  -0.880   0.3790
## smoker1     -3.239e-01  1.058e-01  -3.061   0.0022 **
## charges      1.419e-05  3.365e-06   4.217  2.48e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 2001.6  on 1337  degrees of freedom
## Residual deviance: 1984.1  on 1335  degrees of freedom
## AIC: 3879.4
##
## Number of Fisher Scoring iterations: 5
```

To get the coefficients, the log transformation needs to be reversed:

```
exp(coef(glm.children))
```

```
## (Intercept)      smoker1      charges
##   0.9636169    0.7233085    1.0000142
```

Smoker (factor): The model shows that for the factor smoker (yes/no), a smoker has on average 72% of the number of children a non-smoker has. The more common-sense interpretation might be the other way around, that people who have 1 or more children smoke less, but for the moment we have no proof of that.

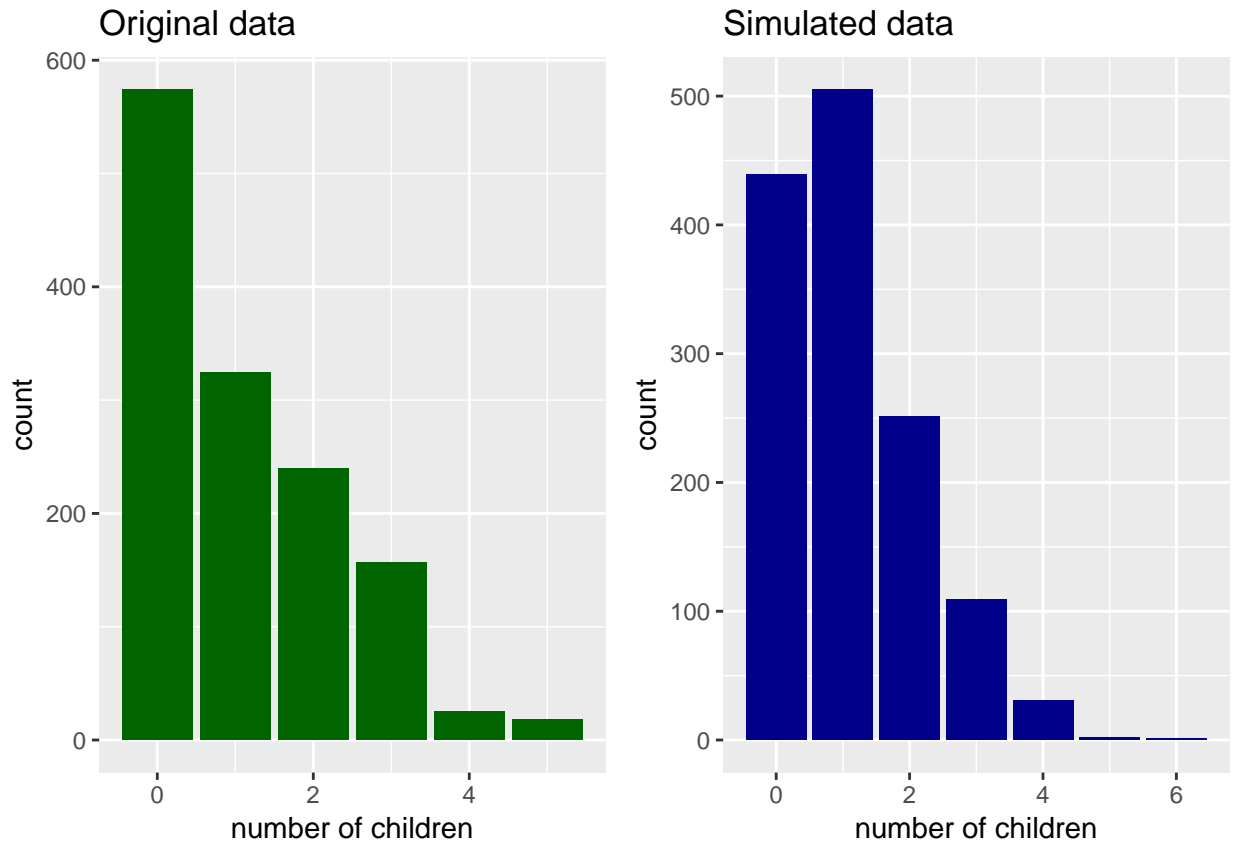
Charges: A person with higher charges will on average have more children. If charges are increased by 1000 dollars, the calculated number of children increases by 1.4%.

6.1.3 Simulation of data and comparison

With the calculated model, data is simulated:

```
##      sim_1
## Min.    :0.000
## 1st Qu.:0.000
## Median :1.000
## Mean    :1.102
## 3rd Qu.:2.000
## Max.    :6.000
```

The original and the simulated data are compared visually. The number of children from the simulated data (0-6) seem to be plausible. The distribution has a strong downwards trend starting at 1 like the original data. However the model does not seem to generate enough data with 0 children.



6.2 Generalised Linear Models for binomial data

A model is fitted that predicts if a person is a smoker or not. Only the significant values age, bmi and charges are used.

```
glm.smoker.2 <- glm(smoker ~ age+bmi+charges,
                    data=insurance,
                    family = "binomial")
summary(glm.smoker.2)
```

```
##
## Call:
## glm(formula = smoker ~ age + bmi + charges, family = "binomial",
##      data = insurance)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -3.09442  -0.10998  -0.04475  -0.00970   1.53727
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  5.311e+00  1.029e+00   5.163 2.43e-07 ***
## age          -9.875e-02  1.300e-02  -7.597 3.02e-14 ***
## bmi          -3.481e-01  4.309e-02  -8.078 6.60e-16 ***
## charges       3.822e-04  2.917e-05  13.104 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

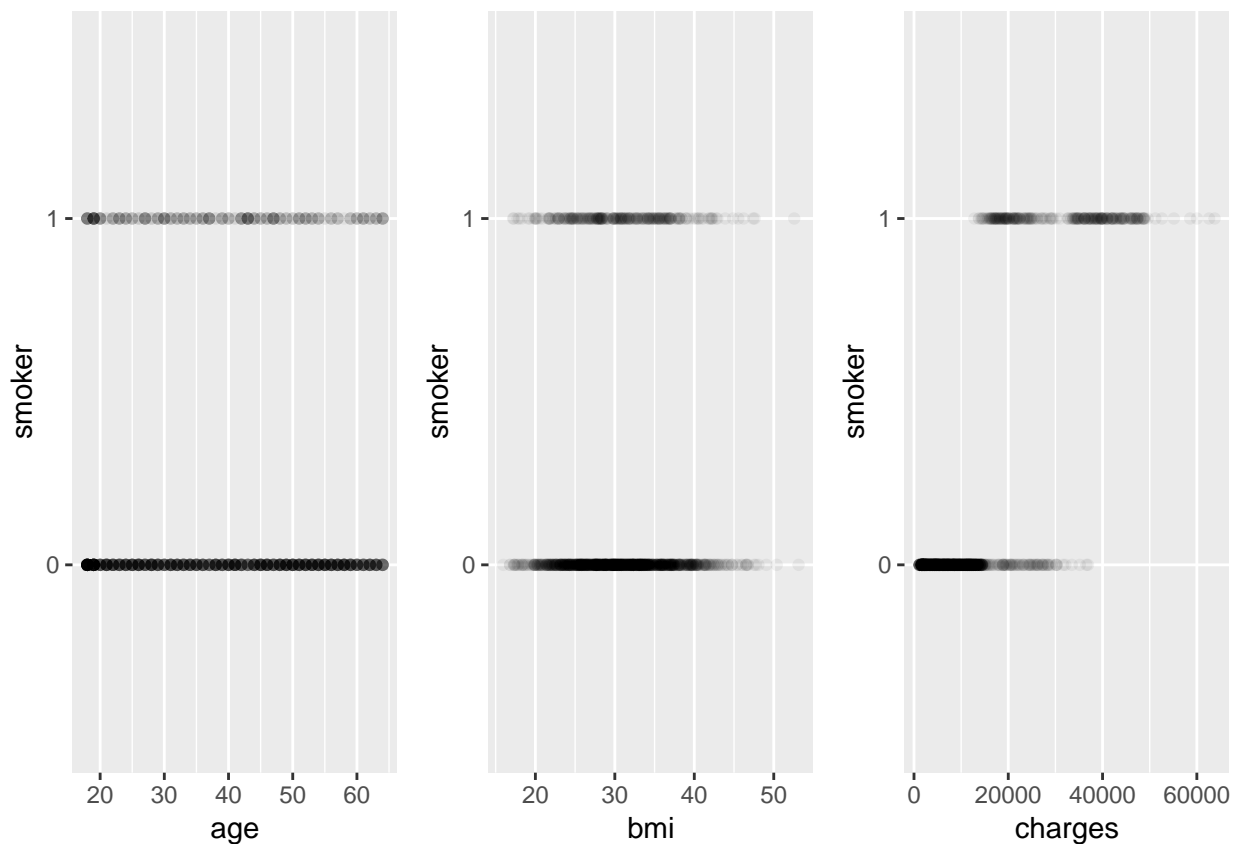
```
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 1356.63 on 1337 degrees of freedom
## Residual deviance: 311.98 on 1334 degrees of freedom
## AIC: 319.98
##
## Number of Fisher Scoring iterations: 8
exp(coef(glm.smoker.2))
```

```
## (Intercept)      age      bmi      charges
## 202.5686249    0.9059677    0.7060520    1.0003823
```

Age and bmi has a negative effect on smoker. This means the higher a persons bmi and age, the lower the probability that the person is a smoker. Charges has a positive effect. This means the higher a persons charges, the higher is the possibility that the person smokes.

6.2.1 Graphical analysis

This can also be explored graphically, at least for charges it is clearly visible that smokers have higher charges.



6.2.2 Estimating the model performance

The predicted values are transformed into binary (beforehand they indicated the probability) and compared with the actual data.

```
fitted.smoker.disc <- ifelse(fitted(glm.smoker.2) < 0.5,
                             yes = 0, no = 1)
head(fitted.smoker.disc)

## 1 2 3 4 5 6
## 1 0 0 1 0 0

d.obs.fit.smoker <- data.frame(obs = insurance$smoker,
                              fitted = fitted.smoker.disc)
head(d.obs.fit.smoker)

##   obs fitted
## 1    1      1
## 2    0      0
## 3    0      0
## 4    0      1
## 5    0      0
## 6    0      0
```

We observe the following fit:

```
##      fit
## obs    0    1
##   0 1028   36
##   1   23  251
```

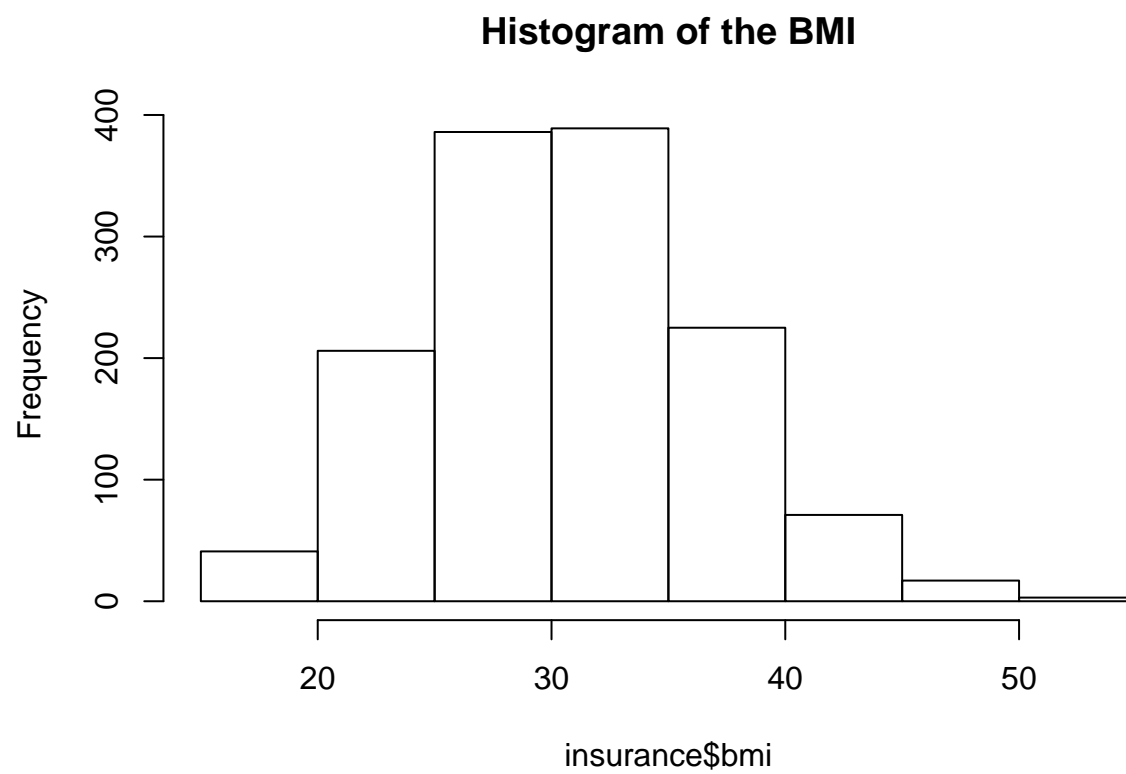
7 Decision Trees

7.1 Regression trees

7.1.1 Inspecting the Data

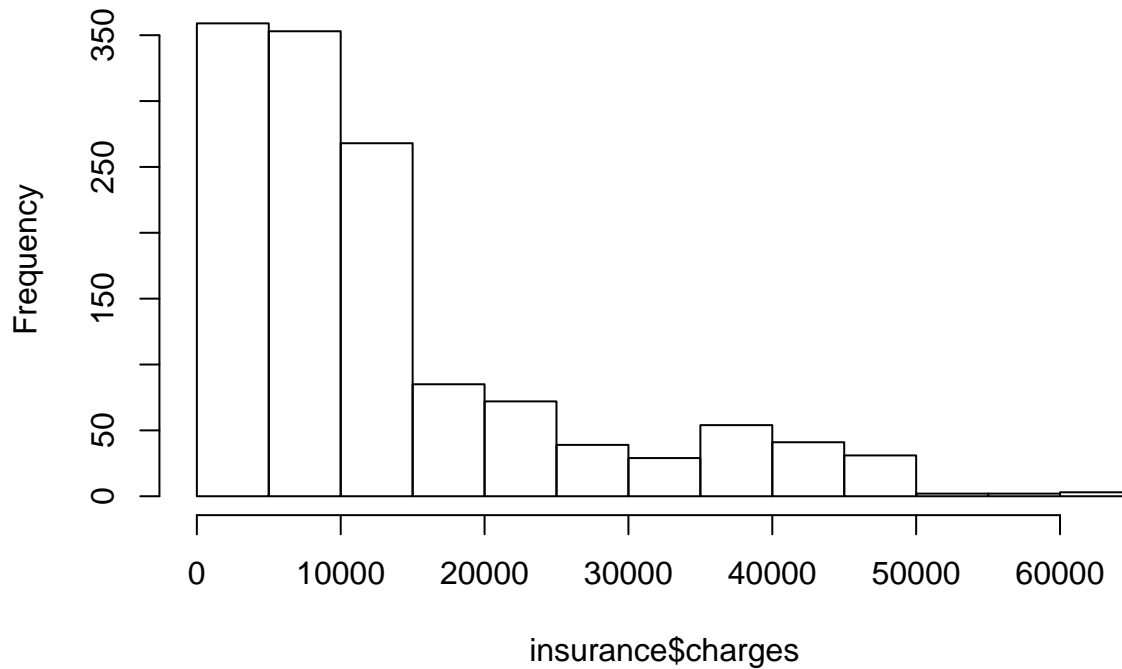
To start off with our regression trees, we throw a glance at the continuous variables:

```
hist(insurance$bmi, main = "Histogram of the BMI")
```



```
hist(insurance$charges, main = "Histogram of charges")
```

Histogram of charges



One can see that bmi seems to be normally distributed whereas charges follows the pattern of a $1/x$ -function.

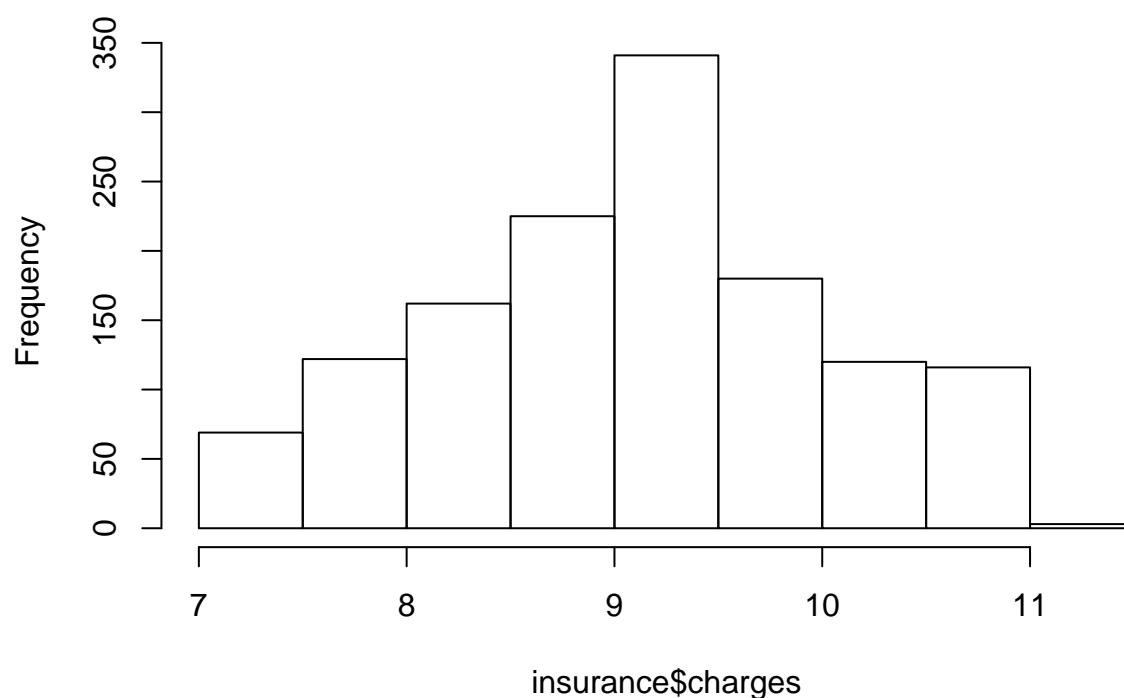
We thus normalize “charges” by logarithming it. (Attention: Please do only execute the code block (conversion to log) below once as executing it several times will apply the logarithmic function several times)

```
insurance$charges <- log(insurance$charges)
```

Please note that hence, all charges values for this chapter will strictly be indicated in log-value and not the true one.

```
hist(insurance$charges, main = "Histogram of the charges after applying logarithm")
```

Histogram of the charges after applying logarithm



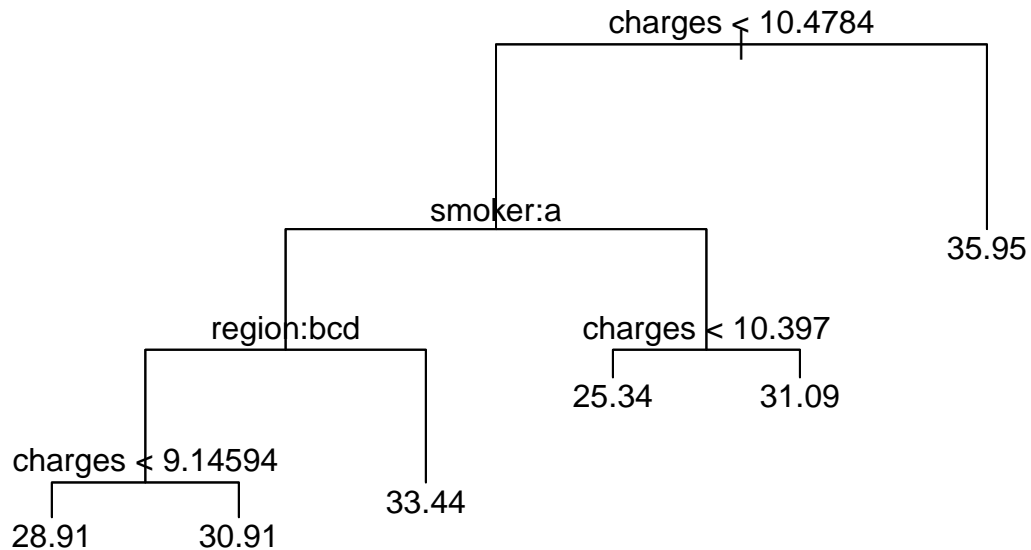
Applying a log-transformation solves this issue.

To begin with, we make the rookie mistake to use all the data, which we will then later change by applying a train-test approach. This is done to emphasize the importance of not using all the data to train one's model by showing the differences in methods.

7.1.2 Regression tree - BMI

```
set.seed(99999)
tree.regression.bmi <- tree(bmi~., data=insurance)
plot(tree.regression.bmi)
text(tree.regression.bmi)
mtext("Regression tree for BMI, using all data \n")
```


Regression tree for BMI, using all data



We here receive a tree with six terminal nodes. Thanks to the tree-function, the tree is already created by recursive binary splitting. As the variables (internal nodes) determining the bmi, “charges” have been applied as well as if one smokes and from which region one comes. This can also be seen in below summary:

```
summary(tree.regression.bmi)
```

```
##
## Regression tree:
## tree(formula = bmi ~ ., data = insurance)
## Variables actually used in tree construction:
## [1] "charges" "smoker" "region"
## Number of terminal nodes: 6
## Residual mean deviance: 29.23 = 38930 / 1332
## Distribution of residuals:
##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -13.09000 -3.62700 -0.07606  0.00000  3.26700  19.69000
```

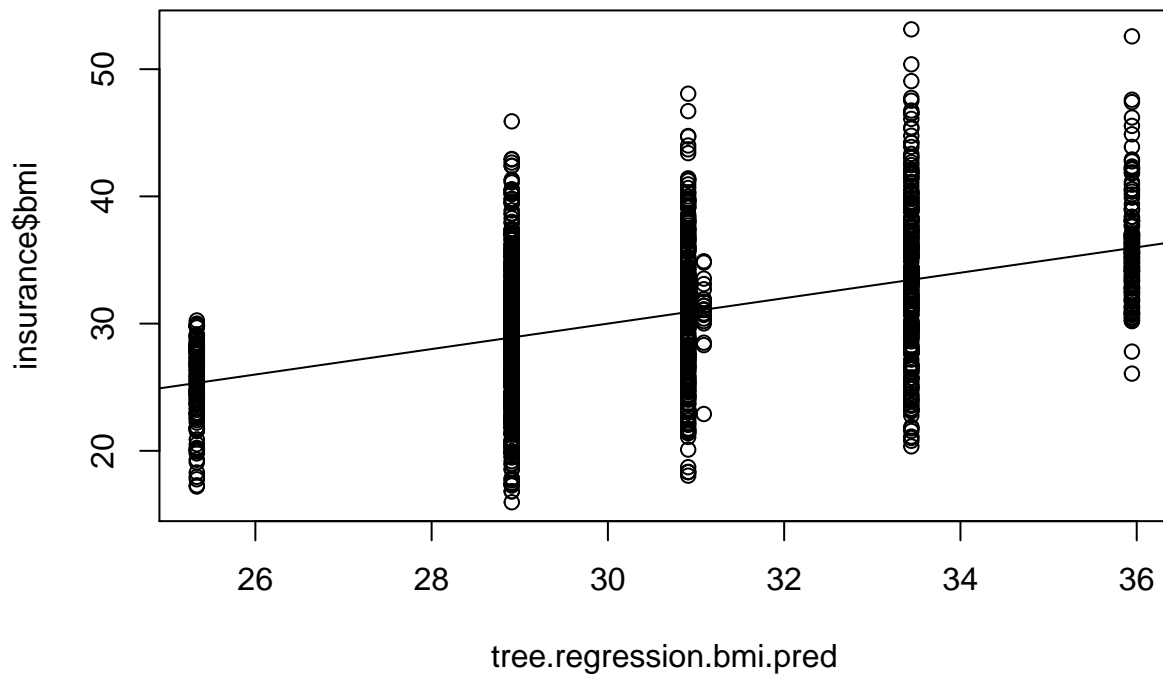
Our residual mean deviance, aka R-squared-error is at 29.23.

7.1.2.1 Prediction

We here establish a prediction for our bmi tree and also compare the predictions with the true value of the bmi:

```
tree.regression.bmi.pred <- predict(tree.regression.bmi, insurance, type="vector")
plot(tree.regression.bmi.pred, insurance$bmi, main = "Prediction of BMI against real BMI values")
abline(0,1)
```

Prediction of BMI against real BMI values



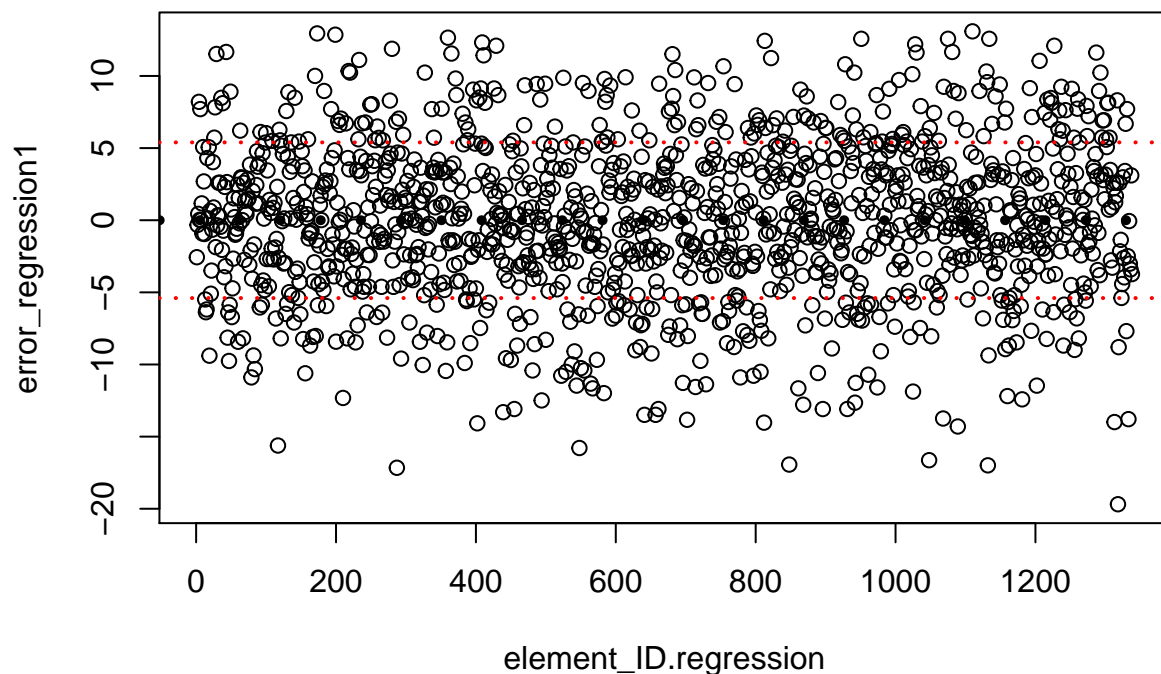
We here see how the prediction performs for the bmi values. As one can see, there seem to be rather a fair number of residuals.

7.1.2.2 Comparison of Prediction <-> True Values

Further, we examine the above said residuals:

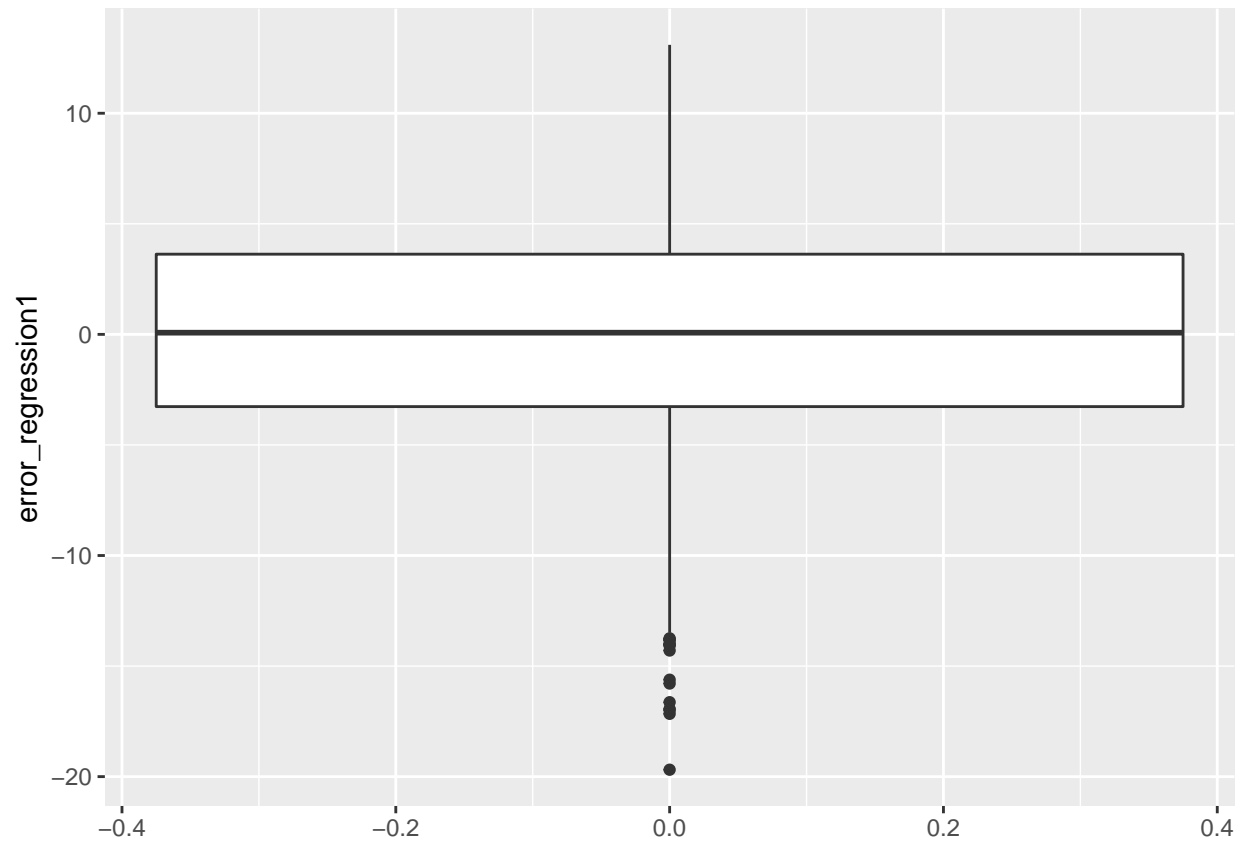
```
error_regression1 <- tree.regression.bmi.pred - insurance$bmi
element_ID_regression <- 1:length(error_regression1)
plot(element_ID_regression, error_regression1)
title(main="Analysis of the residuals")
abline(0, 0, lwd=5, lty="dotted")
abline(sd(error_regression1), 0, lwd=2, col="red", lty="dotted")
abline(-sd(error_regression1), 0, lwd=2, col="red", lty="dotted")
```

Analysis of the residuals



As can be seen, the residual analysis clearly supports the data to be normally distributed. So it is that most of the residuals are distributed fairly symmetrically around 0. A large portion of the residuals is furthermore within the margins of ± 1 standard deviation.

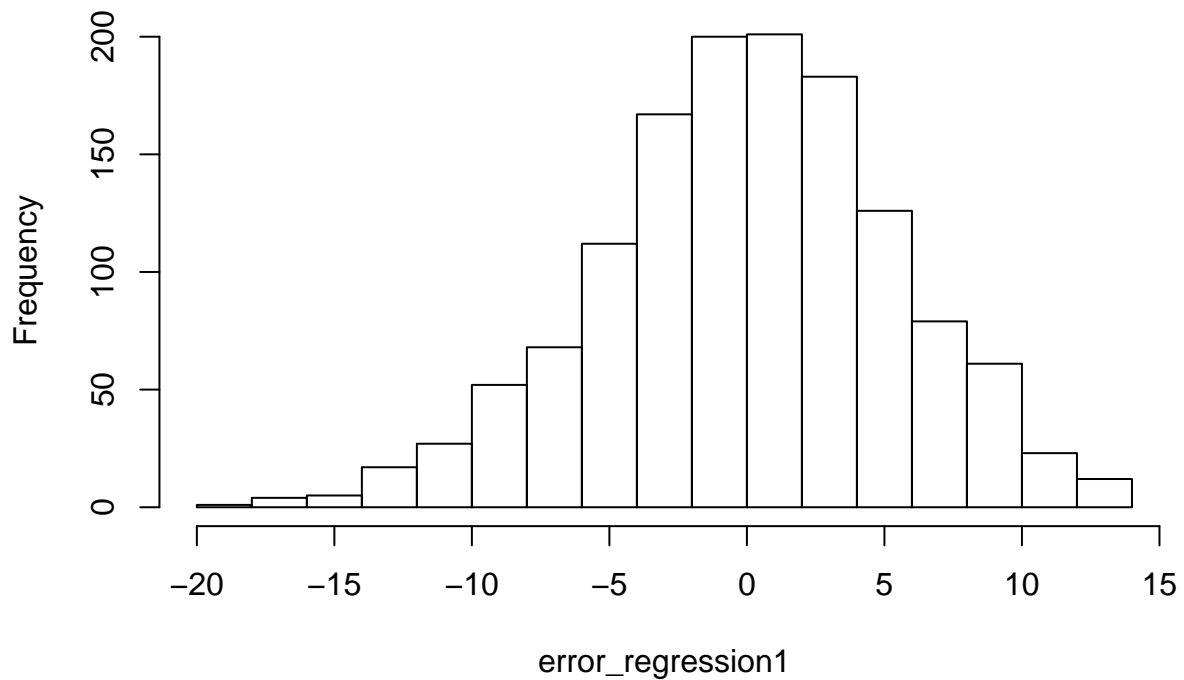
```
error_regression1_dataframe <- tibble(element_ID.regression, error_regression1)
ggplot(data=error_regression1_dataframe) + geom_boxplot(aes(y=error_regression1))
```



A boxplot further demonstrates the rather normally distributed data with mean around 0. It is apparent for there to be more outliers in the negative y-axis area than the positive y-axis area.

```
hist(error_regression1, main = "Histogram of error for the prediction of BMI")
```

Histogram of error for the prediction of BMI



```
RSS_bmi <- sum((insurance[3]-tree.regression.bmi.pred)^2)
MSE_bmi <- RSS_bmi/length(tree.regression.bmi.pred)
deviation_bmi <- sqrt(MSE_bmi)
cat(RSS_bmi)
```

```
## 38928.55
```

```
cat("\n")
```

```
cat(deviation_bmi)
```

```
## 5.39394
```

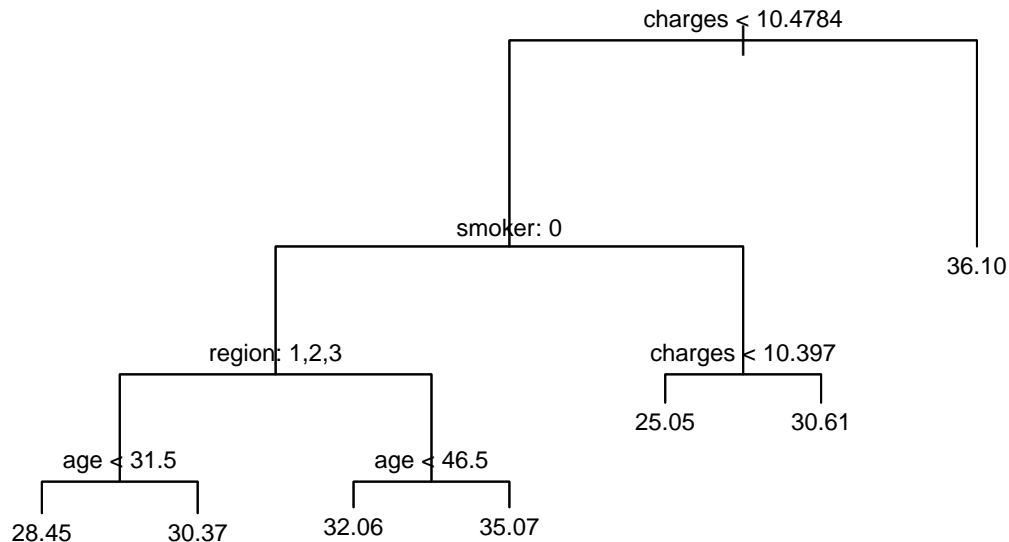
So the average estimation of a bmi is about 5.39 bmi points away from the true value.

We now switch to the train-test approach.

7.1.2.3 Train <-> Test Approach

```
ratio <- 0.7
total <- nrow(insurance)
train <- sample(1:total, as.integer(total * ratio))
tree.regression.bmi2 <- tree(bmi~., insurance, subset=train)
plot(tree.regression.bmi2)
text(tree.regression.bmi2, pretty=12, cex=0.75)
mtext("Regression tree for BMI, using train data \n")
```

Regression tree for BMI, using train data



By working with the test data, we can see at this point a change, meaning a possible overfit. Were there previously six end nodes, we now have seven end nodes, with variables used “charges”, “smoker”, “region” and “age”, as also can be seen here:

```
summary(tree.regression.bmi2)
```

```
##
## Regression tree:
## tree(formula = bmi ~ ., data = insurance, subset = train)
## Variables actually used in tree construction:
## [1] "charges" "smoker" "region" "age"
## Number of terminal nodes: 7
## Residual mean deviance: 27.84 = 25860 / 929
## Distribution of residuals:
##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -13.55000  -3.46800  -0.09678   0.00000   3.35700  16.48000
```

Our residual mean deviance has furthermore gone from 29.23 to 27.84.

7.1.2.4 Error Analysis

```
tree.regression.bmi2.pred <- predict(tree.regression.bmi2, insurance[-train,], type="vector")
(RSS.2.train <- mean(((insurance[train,][3]-predict(tree.regression.bmi2, insurance[train,], type="vector"))
```

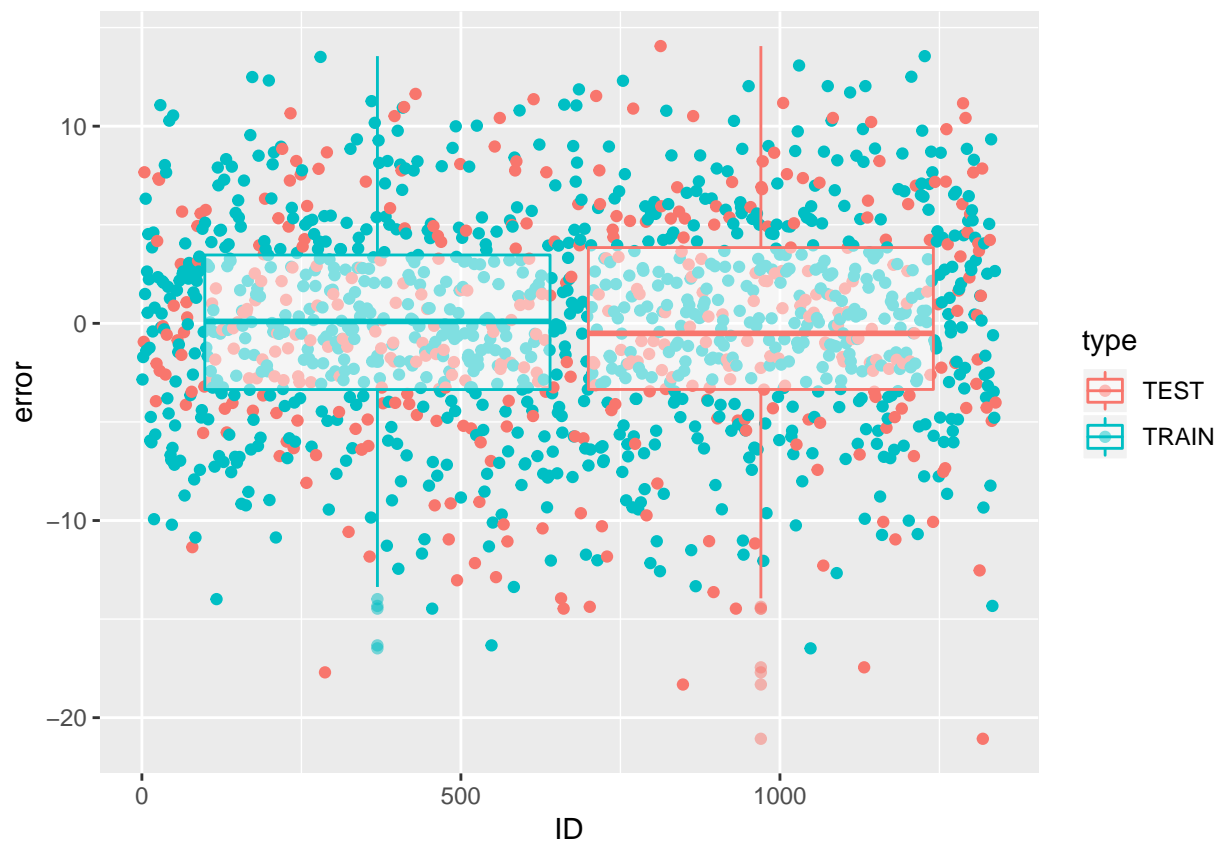
```
## [1] 27.63265
```

```
(RSS.2.test <- mean(((insurance[-train,][3]-tree.regression.bmi2.pred)^2)$bmi))
```

```
## [1] 32.85365
```

We receive an RSS for the train dataset of 27.21 whereas for the test set, the RSS is 34.55. We seem to have an overfit as our model performs considerably worse for the test data than the train data.

```
errors.2.in <- predict(tree.regression.bmi2, insurance[train,], type="vector")-insurance[train,]$bmi
element.2.in <- as.integer(names(errors.2.in))
errors.2.in_dataframe <- tibble(element.2.in,errors.2.in,"TRAIN")
colnames(errors.2.in_dataframe) <- c('ID','error','type')
errors.2 <- predict(tree.regression.bmi2, insurance[-train,], type="vector")-insurance[-train,]$bmi
element.2 <- 1:length(errors.2)
element.2 <- as.integer(names(errors.2))
errors.2.out_dataframe <- tibble(element.2,errors.2,"TEST")
colnames(errors.2.out_dataframe) <- c('ID','error','type')
errors.2_dataframe <- bind_rows(errors.2.in_dataframe,errors.2.out_dataframe)
errors.2_dataframe <- arrange(errors.2_dataframe, ID)
ggplot(data = errors.2_dataframe, mapping = aes(x = ID,y = error, color = type)) +
  geom_point() + geom_boxplot(alpha = 0.5)
```



Both dataframes show the about same boxplots in regard of the errors. Not surprisingly, the distribution of the test data is wider, considering that the model has been trained with

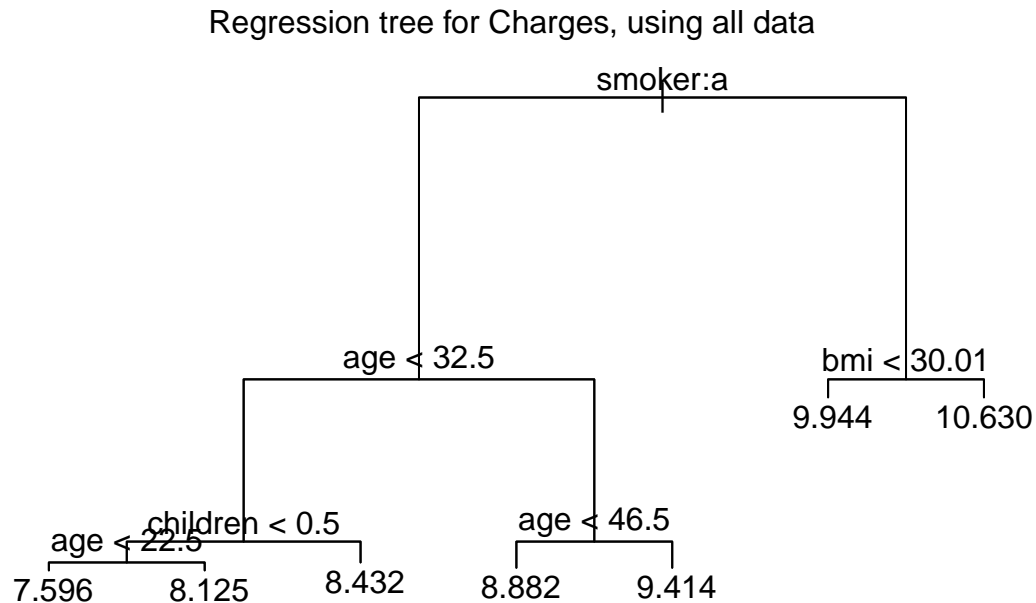
7.1.3 Regression tree - Charges

We now examine the behaviour of the variable “charges”. First off again with the entire data values.

```
summary(insurance$charges)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	7.023	8.464	9.147	9.099	9.720	11.063

```
tree.regression.charges <- tree(charges~., data=insurance)
plot(tree.regression.charges)
text(tree.regression.charges)
mtext("Regression tree for Charges, using all data \n")
```



If one smokes seems to be the crucial factor to start of the charges classification. We further receive a tree with seven end nodes and four influencing variables, as also the summary further demonstrates:

```
summary(tree.regression.charges)
```

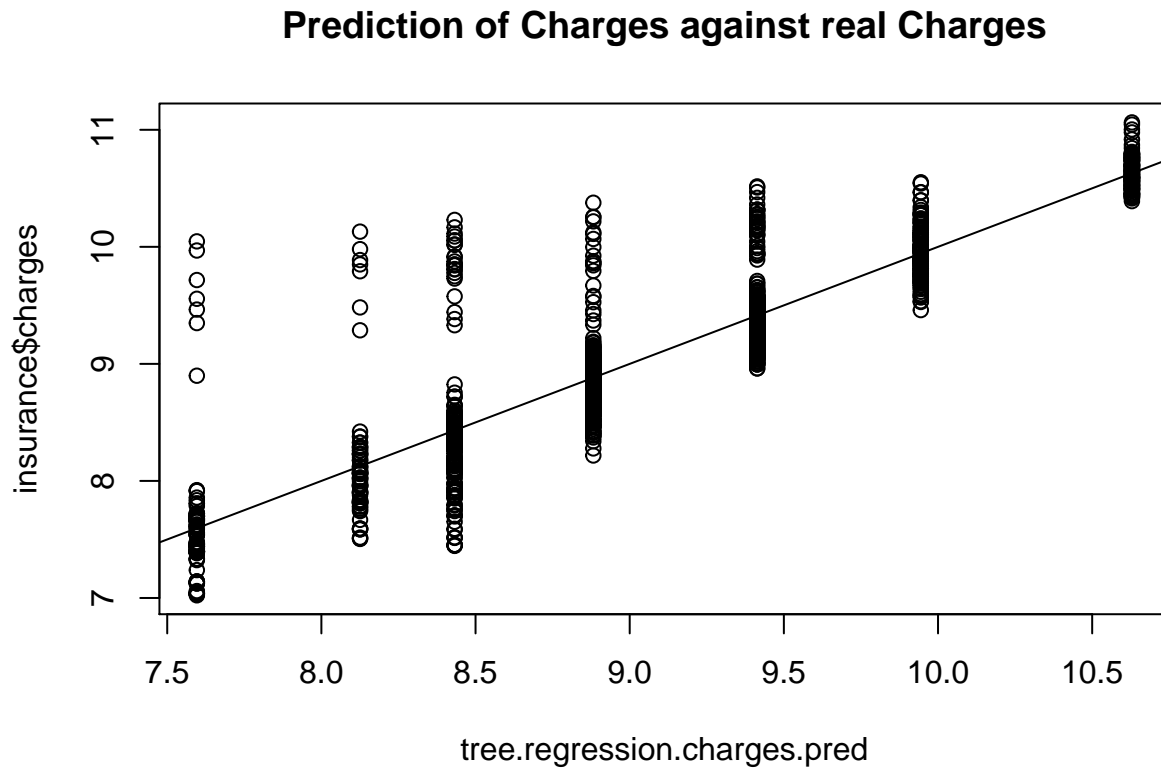
```
##
## Regression tree:
## tree(formula = charges ~ ., data = insurance)
## Variables actually used in tree construction:
## [1] "smoker" "age" "children" "bmi"
## Number of terminal nodes: 7
## Residual mean deviance: 0.1592 = 211.9 / 1331
## Distribution of residuals:
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -0.98730 -0.19690 -0.04999  0.00000  0.09942  2.45100
```

Our residual mean deviance (R-squared-error) is at 0.1592.

7.1.3.1 Prediction

We are now again predicting the data that has already been used for training the model.


```
tree.regression.charges.pred <- predict(tree.regression.charges, insurance, type="vector")
plot(tree.regression.charges.pred, insurance$charges, main="Prediction of Charges against real Charges",
abline(0 ,1)
```

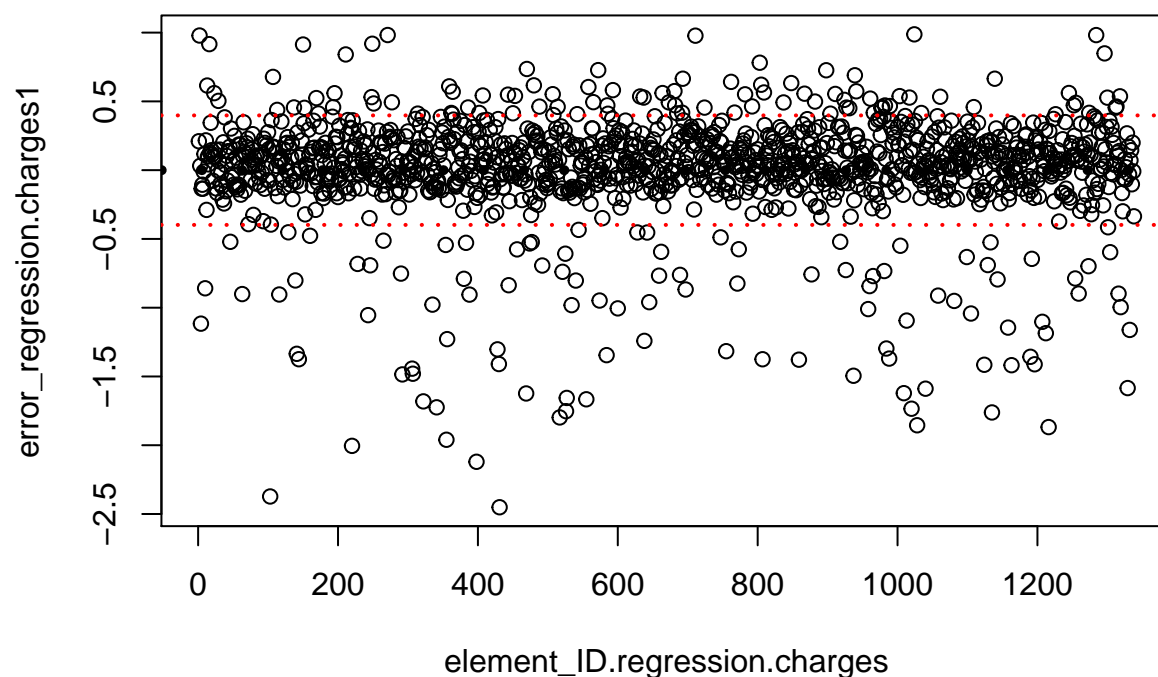


We here see how our model predicts the charges to by (x-axis) in comparison to what the charges actually are.

7.1.3.2 Comparison Prediction <-> True Values

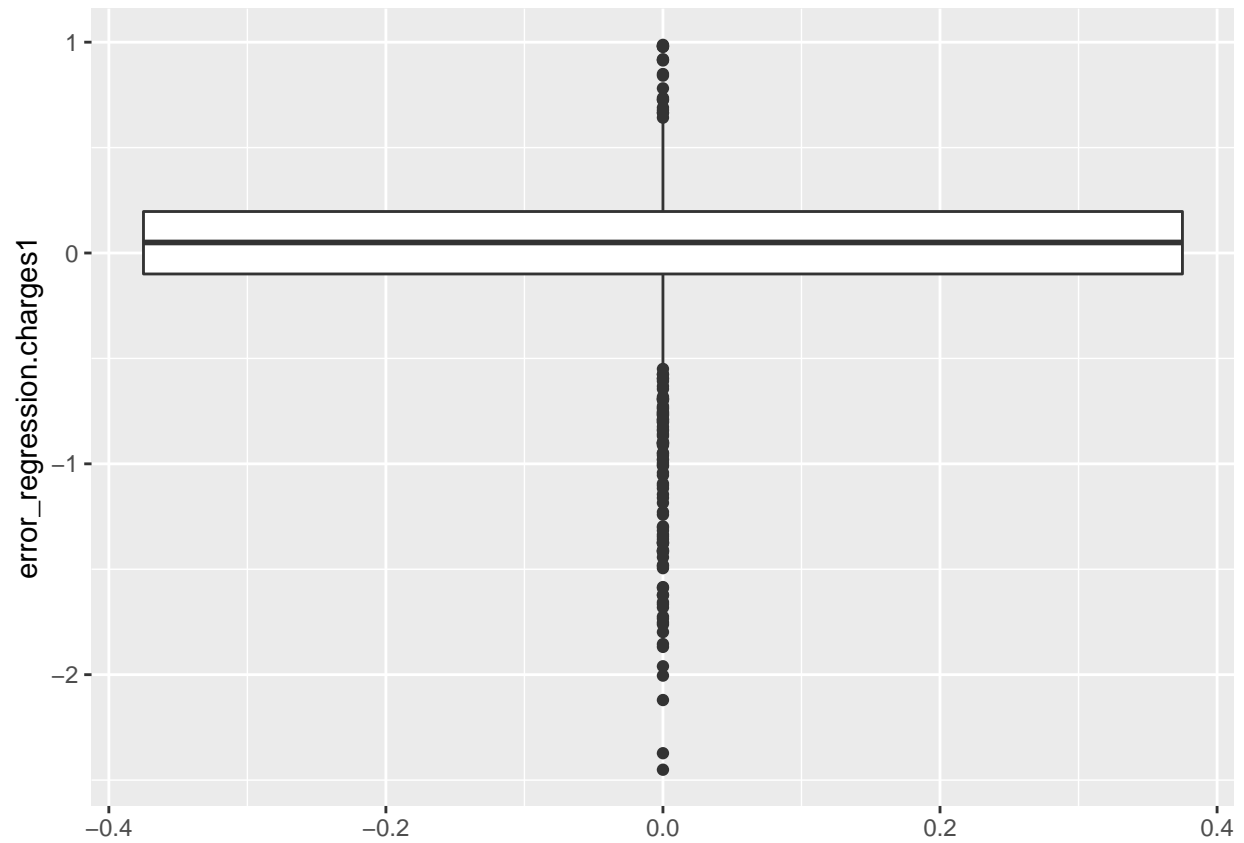
```
error_regression.charges1 <- tree.regression.charges.pred-insurance$charges
element_ID.regression.charges <- 1:length(error_regression.charges1)
plot(element_ID.regression.charges, error_regression.charges1)
title(main="Analysis of the residuals")
abline(mean(error_regression.charges1),0, lwd=5,lty="dotted")
abline(sd(error_regression.charges1) ,0, lwd=2, col="red", lty="dotted")
abline(-sd(error_regression.charges1) ,0, lwd=2, col="red", lty="dotted")
```

Analysis of the residuals



We can furthermore clearly see that the residuals are roughly normalized, with an approximate mean of 0. Outliers are however more distributed in the negative axis of y than in the positive.

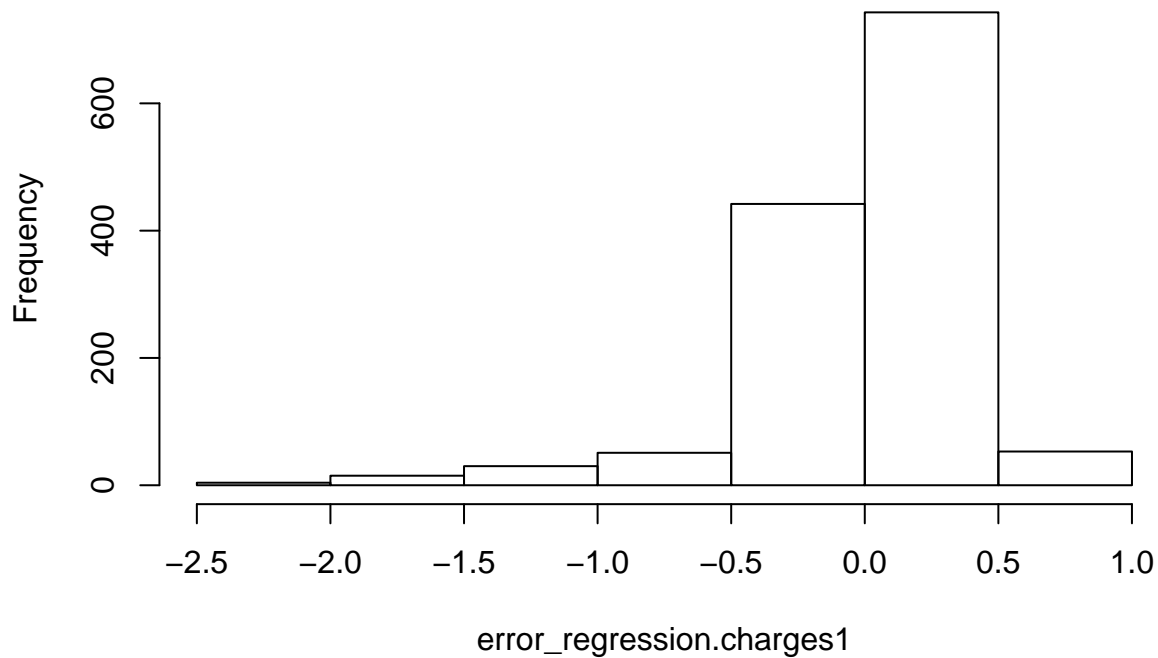
```
error_regression1_dataframe.charges <- tibble(element_ID.regression.charges, error_regression.charges1)
ggplot(data=error_regression1_dataframe.charges) + geom_boxplot(aes(y=error_regression.charges1))
```



The above boxplot further demonstrates that there is a slight bias in the distribution.

```
hist(error_regression.charges1, main = "Histogram of the Error for charges")
```

Histogram of the Error for charges



Here goes to further show that the data points are slightly skewed. However, mean of the residuals is still roughly at 1.

```
RSS_charges <- sum((insurance[7]-tree.regression.charges.pred)^2)
MSE_charges <- RSS_charges/length(tree.regression.charges.pred)
deviation_charges <- sqrt(MSE_charges)

cat(RSS_charges)
```

```
## 211.8598
```

```
cat("\n")
```

```
cat(deviation_charges)
```

```
## 0.3979204
```

RSS here is at 211.86, the deviation at 0.4.

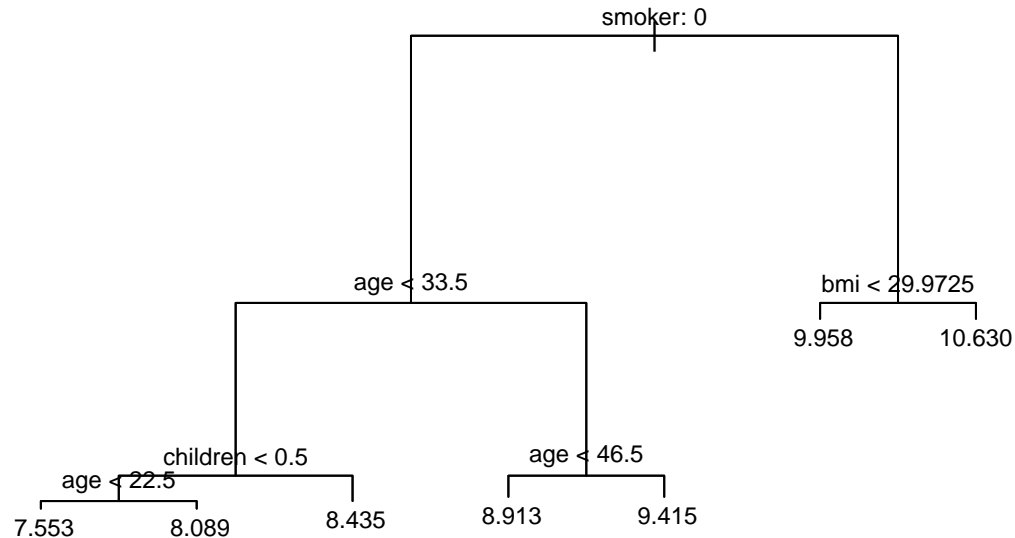
7.1.3.3 Train-Test-Approach

Again, we use the train-test approach to demonstrate the importance to not use all the data to train one's model.

```
ratio <- 0.7
total <- nrow(insurance)
train.charges <- sample(1:total, as.integer(total * ratio))
tree.regression.charges2 <- tree(charges~., insurance, subset=train)
plot(tree.regression.charges2)
```

```
text(tree.regression.charges2, pretty=12, cex=0.75)
mtext("Regression tree for Charges, using train data \n")
```

Regression tree for Charges, using train data



We now receive a tree with exactly the same number of nodes but where the criteria have different values, except for the first node, which remain the same (non-smoker → left branch, smoker → right branch).

```
summary(tree.regression.charges2)
```

```
##
## Regression tree:
## tree(formula = charges ~ ., data = insurance, subset = train)
## Variables actually used in tree construction:
## [1] "smoker" "age" "children" "bmi"
## Number of terminal nodes: 7
## Residual mean deviance: 0.1475 = 137 / 929
## Distribution of residuals:
## Min. 1st Qu. Median Mean 3rd Qu. Max.
## -0.98970 -0.18750 -0.04011 0.00000 0.10500 2.41600
```

Our residual mean deviance is now a bit lower at 0.1475 (previously: 0.1592)

```
tree.regression.charges2.pred <- predict(tree.regression.charges2, insurance[-train,], type="vector")
(RSS.2.train.charges <- mean(((insurance[train,][7]-predict(tree.regression.charges2, insurance[train,]
```

```
## [1] 0.146418
```

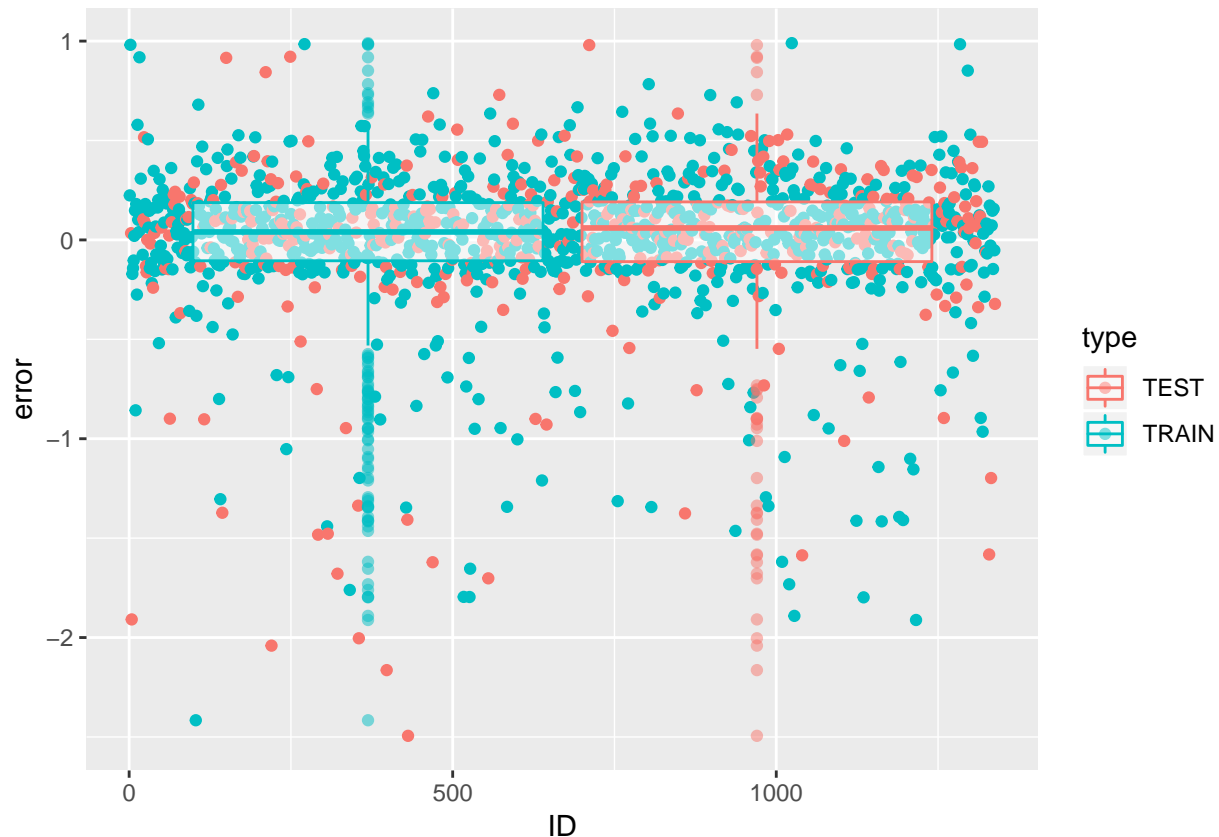
```
(RSS.2.test.charges <- mean(((insurance[-train,][7]-tree.regression.charges2.pred)^2)$charges))
```

```
## [1] 0.196484
```

However, our train data, despite still having the same number of terminal nodes, is biased and thus shows an overfit. Hence it is that RSS for the train set is at $e^{0.146418}$, whereas for the test set, we are at $e^{0.196484}$.

7.1.3.4 Errors

```
errors.2.in.charges <- predict(tree.regression.charges2, insurance[train,], type="vector")-insurance[tr
element.2.in.charges <- as.integer(names(errors.2.in.charges))
errors.2.in_dataframe.charges <- tibble(element.2.in.charges,errors.2.in.charges,"TRAIN")
colnames(errors.2.in_dataframe.charges) <- c('ID','error','type')
errors.2.charges <- predict(tree.regression.charges2, insurance[-train,], type="vector")-insurance[-tra
element.2.charges <- 1:length(errors.2.charges)
element.2.charges <- as.integer(names(errors.2.charges))
errors.2.out_dataframe.charges <- tibble(element.2.charges,errors.2.charges,"TEST")
colnames(errors.2.out_dataframe.charges) <- c('ID','error','type')
errors.2_dataframe.charges <- bind_rows(errors.2.in_dataframe.charges,errors.2.out_dataframe.charges)
errors.2_dataframe.charges <- arrange(errors.2_dataframe.charges, ID)
ggplot(data = errors.2_dataframe.charges, mapping = aes(x = ID,y = error, color = type)) +
  geom_point() + geom_boxplot(alpha = 0.5)
```



There are a few differences, especially in regard of the outliers between the test and the train data. Roughly, one can say however that the boxplot is about the same for train and test. Also, the outliers in the negative area are for both data sets more spread.

7.2 Classification Trees

We here examine a decision tree for a non continuous (hence categorical) variable.

```
tree.classification.region_all <- tree(region~., data=insurance)
summary(tree.classification.region_all)
```

```
##
## Classification tree:
## tree(formula = region ~ ., data = insurance)
## Variables actually used in tree construction:
## [1] "age" "bmi"
## Number of terminal nodes: 4
## Residual mean deviance: 2.604 = 3474 / 1334
## Misclassification error rate: 0.6547 = 876 / 1338
```

We here receive a classification tree with four end nodes. Furthermore our RMS is at 2.604 and the misclassification rate - which is basically the RSS for classification trees - is at 65.47%.

```
plot(tree.classification.region_all)
text(tree.classification.region_all, pretty=1, cex=0.75)
mtext("Classification tree for regions, using all data (0: SE, 1: SW, 2: NE, 3:NW 3)\n")
```

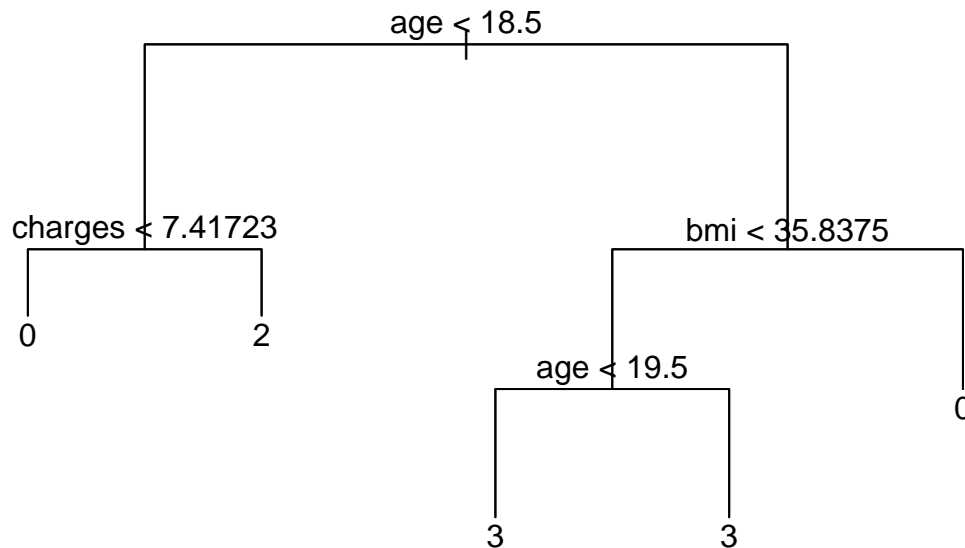
Classification tree for regions, using all data (0: SE, 1: SW, 2: NE, 3:NW 3)



It appears that the high misclassification rate is partly due to only 2 of the 4 factors being applied in the classification tree. Most important here is if one is under the age of 18.5 years (left branch) or above (right branch). We thus go into testing mode:

```
tree.classification.region <- tree(region~., data=insurance, subset=train)
plot(tree.classification.region)
text(tree.classification.region)
mtext("Classification tree for regions, using only train data (0: SE, 1: SW, 2: NE, 3:NW 3)\n")
```

Classification tree for regions, using only train data (0: SE, 1: SW, 2: NE, 3: NW)



Working with our training set, we see a clear change: we now have five terminal nodes as well as the factor of 2 -> “northeast” that is not part of the equation.

```
tree.classification.region.pred <- predict(tree.classification.region, insurance[train,], type="class")
```

7.2.1 Confusion Table for Train Data

We now examine the classification error on the train data by comparing the estimation with the true value to elaborate how exact our trained model works.

```

tree.classification.region.pred.ct <- table(tree.classification.region.pred, insurance[train,]$region)
tree.classification.region.pred.correct <- 0
tree.classification.region.pred.error <- 0
for (i1 in 1:4) {
  for (i2 in 1:4) {
    if (i1 == i2) {
      tree.classification.region.pred.correct <- tree.classification.region.pred.correct + tree.classification.region.pred.ct[i1,i2]
    }else{
      tree.classification.region.pred.error <- tree.classification.region.pred.error + tree.classification.region.pred.ct[i1,i2]
    }
  }
}
tree.classification.region.pred.rate <- tree.classification.region.pred.correct/sum(tree.classification.region.pred.ct)
tree.classification.region.pred.error.1 <- 1 - tree.classification.region.pred.rate
cat(tree.classification.region.pred.rate)

```

```
## 0.357906
```



```
cat("\n")
```

```
cat(tree.classification.region.pred.error.1)
```

```
## 0.642094
```

Our “new” tree manages to reduce our classification error by about 1% for the train data.

7.2.2 Confusion Table for Test Data

```
tree.classification.region.pred.test <- predict(tree.classification.region, insurance[-train,], type="c")
```

We now determine the classification error on our test data:

```
(tree.classification.region.pred.test.ct <- table(tree.classification.region.pred.test, insurance[-train,]))
```

```
##
## tree.classification.region.pred.test  0  1  2  3
##                                     0 41 14 16  9
##                                     1  0  0  0  0
##                                     2  3  0  8  0
##                                     3 73 81 79 78
```

```
tree.classification.region.pred.correct <- 0
tree.classification.region.pred.error <- 0
for (i1 in 1:3) {
  for (i2 in 1:3) {
    if (i1 == i2) {
      tree.classification.region.pred.correct <- tree.classification.region.pred.correct + tree.classification.region.pred.test.ct[i1, i2]
    } else {
      tree.classification.region.pred.error <- tree.classification.region.pred.error + tree.classification.region.pred.test.ct[i1, i2]
    }
  }
}
```

```
(tree.classification.region.pred.rate <- tree.classification.region.pred.correct/sum(tree.classification.region.pred.test.ct))
```

```
## [1] 0.1218905
```

```
(tree.classification.region.pred.error.2 <- 1 - tree.classification.region.pred.rate)
```

```
## [1] 0.8781095
```

Whereas the model performed a bit better with the train data in comparison with the entire data set, the model now clearly performs in an all time bad. The training model is thus a clear overfit.

It seems that region, considering the constant high error quote for entire data set, train set as well as test set, is not a variable that has a large connection with the other variables used.

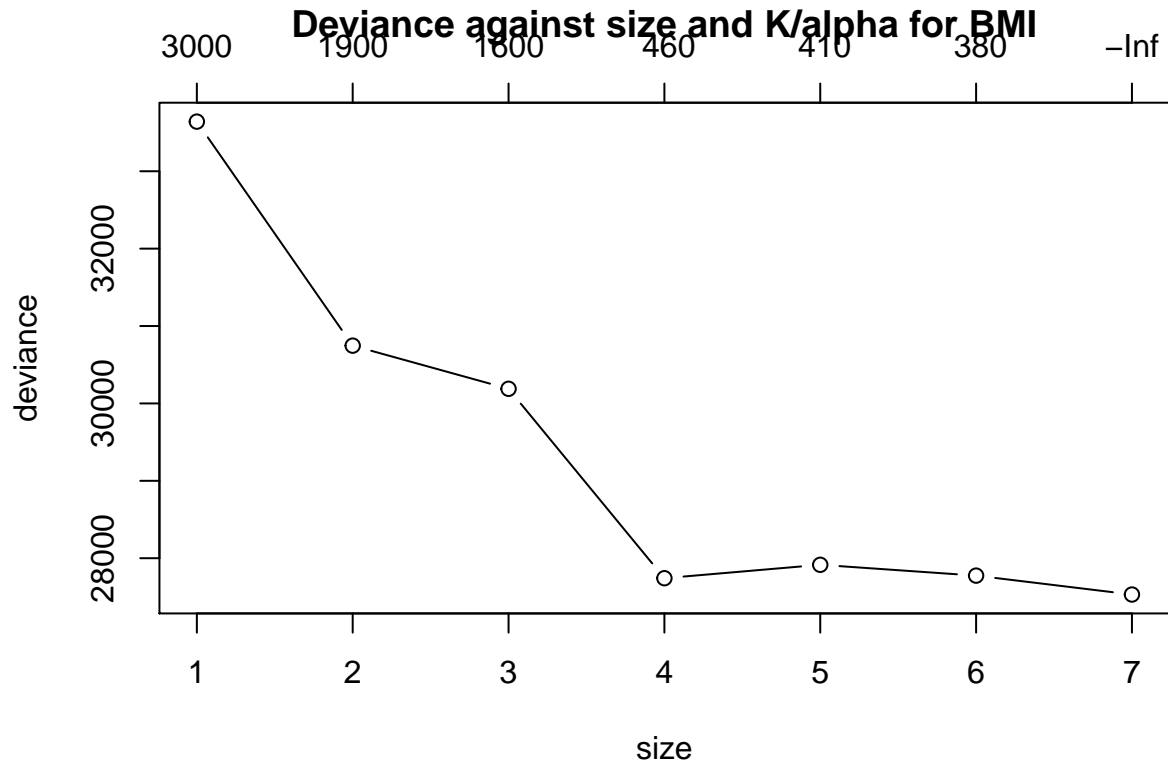
7.3 Pruning

We now continue to prune the trees as so far, we actually managed to create an overfit of our train data in comparison to our test data.

7.3.1 Pruning the continuous BMI variable

We recall: the bmi showed a clear overfit when comparing the RSS of train <-> test data: train = 27.21 test = 34.55

```
tree.regression.bmi.pruning = cv.tree(tree.regression.bmi2, FUN = prune.tree)
plot(tree.regression.bmi.pruning, type = "b")
title("Deviance against size and K/alpha for BMI")
```



With cvtree, the crossvalidation is done automatically for one. We here see that the tree with seven end nodes has the smallest cross-validation error of about 27'250. This is also further shown in below table. Four terminal end leaves is furthermore also very interesting as afterwards (up to 7), the changes are only small in comparison with the steep lines from 1 to 3.

```
tree.regression.bmi.pruning
```

```
## $size
## [1] 7 6 5 4 3 2 1
##
## $dev
## [1] 27530.68 27774.88 27914.67 27741.44 30188.86 30746.46 33641.26
##
## $k
## [1] -Inf 376.0091 414.7563 462.4535 1571.4867 1872.8382 3018.0179
##
## $method
## [1] "deviance"
##
## attr("class")
## [1] "prune" "tree.sequence"
```

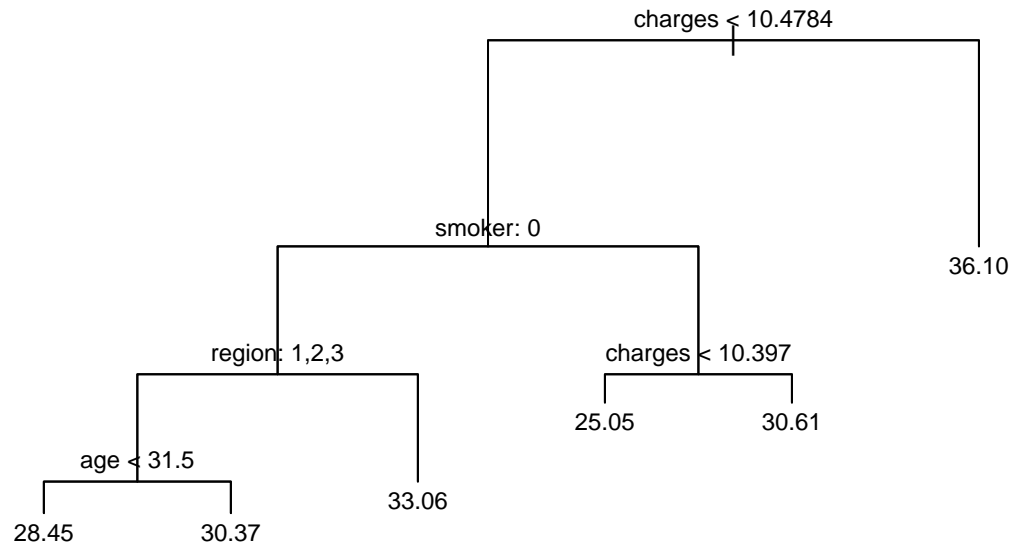
What is interesting to see is that between the number of end nodes 4 and 7, there is a slight growth as well

as decline in terms of variance, so the slope inbetween is both negative as well as positive.

```
tree.regression.bmi.pruned <- prune.tree(tree.regression.bmi2, best = 6)
```

```
plot(tree.regression.bmi.pruned)
text(tree.regression.bmi.pruned, pretty=1, cex=0.75)
mtext("Pruned regression tree for BMI, with best = 6 \n")
```

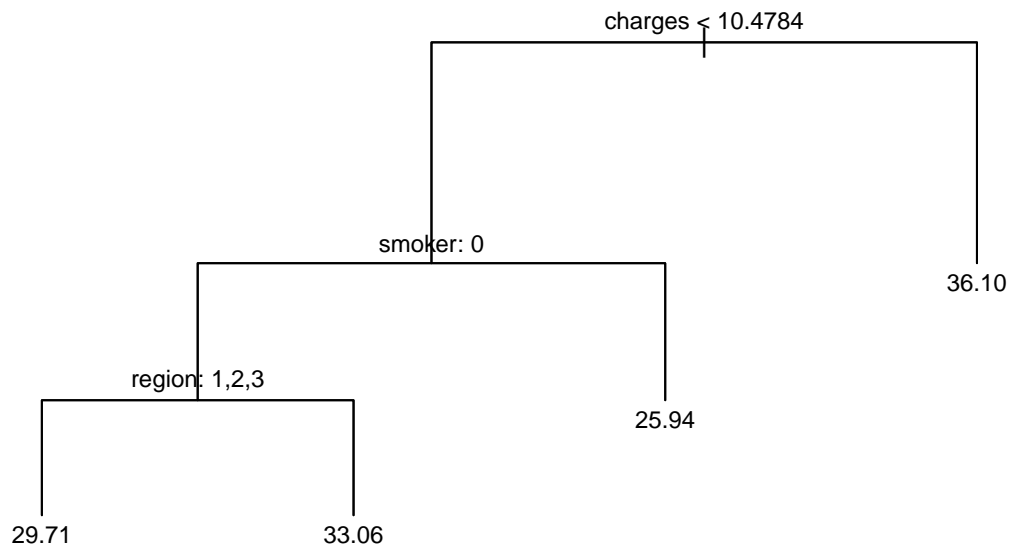
Pruned regression tree for BMI, with best = 6



```
tree.regression.bmi.pruned2 <- prune.tree(tree.regression.bmi2, best = 4)
```

```
plot(tree.regression.bmi.pruned2)
text(tree.regression.bmi.pruned2, pretty=1, cex=0.75)
mtext("Pruned regression tree for BMI, with best = 4 \n")
```

Pruned regression tree for BMI, with best = 4



```
summary(tree.regression.bmi.pruned)
```

```
##
## Regression tree:
## snip.tree(tree = tree.regression.bmi2, nodes = 9L)
## Variables actually used in tree construction:
## [1] "charges" "smoker" "region" "age"
## Number of terminal nodes: 6
## Residual mean deviance: 28.22 = 26240 / 930
## Distribution of residuals:
##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -13.55000 -3.57100 -0.09678  0.00000  3.33800  16.48000
```

```
summary(tree.regression.bmi.pruned2)
```

```
##
## Regression tree:
## snip.tree(tree = tree.regression.bmi2, nodes = c(9L, 5L, 8L))
## Variables actually used in tree construction:
## [1] "charges" "smoker" "region"
## Number of terminal nodes: 4
## Residual mean deviance: 29.1 = 27120 / 932
## Distribution of residuals:
##      Min. 1st Qu.  Median     Mean 3rd Qu.    Max.
## -13.750 -3.680 -0.110  0.000  3.635  16.990
```

In comparison, we can say that the tree with four end nodes more or less performs the same as the tree with six end nodes. Of course, for the tree with 6 leaves, the residual mean deviance is a bit lower than for four.

In terms of internal nodes however, the construction of the tree is roughly the same - with, of course, the exception of the two additional splits for which is also used an additional variable in comparison, namely “age”.

```
(RSS.cv.bmi.test <- sum((insurance[train,][3]-predict(tree.regression.bmi.pruned, insurance[train,], type="vector",
## [1] 28.03437
tree.regression.bmi.pruned.pred <- predict(tree.regression.bmi.pruned, insurance[-train,], type="vector",
(RSS.cv <- sum((insurance[-train,][3]-tree.regression.bmi.pruned.pred)^2)/length(tree.regression.bmi.pruned.pred))
## [1] 31.87446
```

Unsurprisingly, the error for the pruned tree (30.42) is still higher because with pruning, you basically have a fully grown tree and start cutting branches off that are not as necessary as others. Meaning pruning does not grant a perfect performance in an absolute sense.

7.3.2 Pruning regions

We now prune our classification tree, meaning the tree with the categorical variable. By setting FUN to prune.misclass, we communicate that this decision tree will be handling a continuous variable.

```
tree.classification.region.pruning <- cv.tree(tree.classification.region, FUN = prune.misclass)
summary(tree.classification.region.pruning)
```

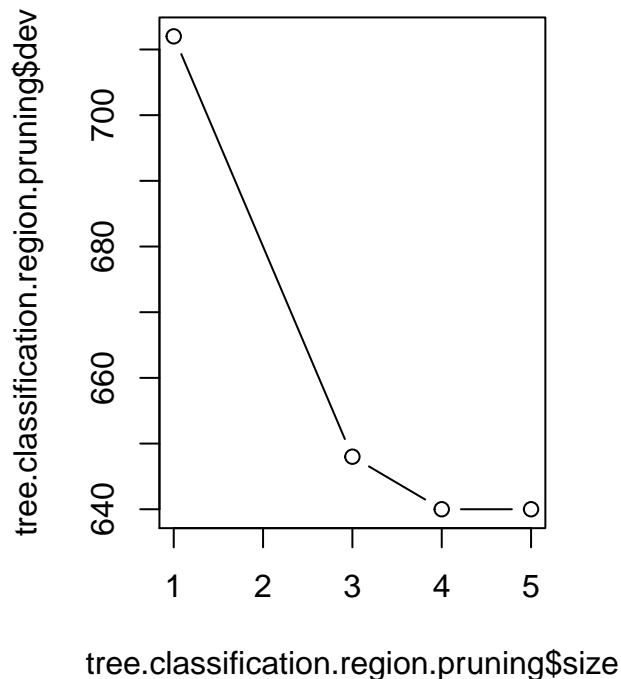
```
##           Length Class  Mode
## size      4      -none- numeric
## dev       4      -none- numeric
## k         4      -none- numeric
## method 1    -none- character
tree.classification.region.pruning

## $size
## [1] 5 4 3 1
##
## $dev
## [1] 640 640 648 712
##
## $k
## [1] -Inf    0   10   39
##
## $method
## [1] "misclass"
##
## attr("class")
## [1] "prune"          "tree.sequence"
```

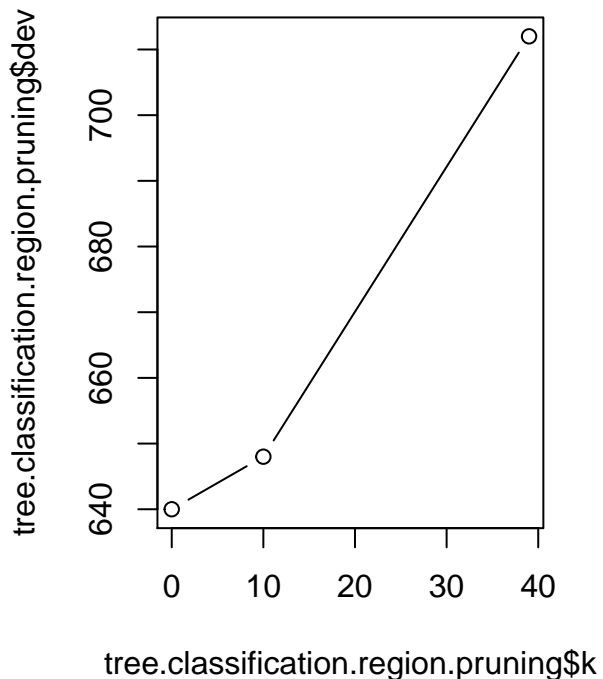
We have a length of four whereas the size of four is the optimal size.

```
par(mfrow=c(1,2))
plot(tree.classification.region.pruning$size, tree.classification.region.pruning$dev, type="b", main = "Size vs Dev")
plot(tree.classification.region.pruning$k, tree.classification.region.pruning$dev, type="b", main= "Eff vs k")
```

Effect of size on deviance



Effect of alpha on deviance



```
par(mfrow=c(1,1))
```

For around 3 and 4, the line is flat whereas before it is very steep and afterwards grows a little bit as well. For deviance, we receive good values for it inbetween 0 to 10.

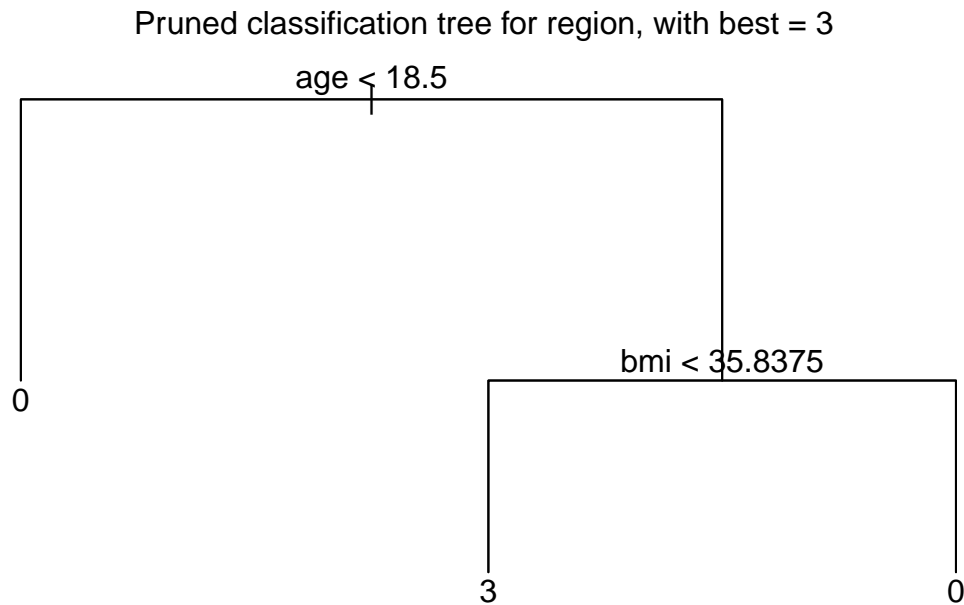
We go with three for the further processing of pruning this decision tree, as it combines low values in both illustrations above, thus is good in terms of a low deviance.

```
prune.tree.classification.region <- prune.misclass(tree.classification.region, best=3)
```

```
summary (prune.tree.classification.region)
```

```
##
## Classification tree:
## snip.tree(tree = tree.classification.region, nodes = c(6L, 2L
## ))
## Variables actually used in tree construction:
## [1] "age" "bmi"
## Number of terminal nodes: 3
## Residual mean deviance: 2.633 = 2456 / 933
## Misclassification error rate: 0.6528 = 611 / 936
```

```
plot(prune.tree.classification.region)
text(prune.tree.classification.region,pretty=0)
mtext("Pruned classification tree for region, with best = 3 \n")
```

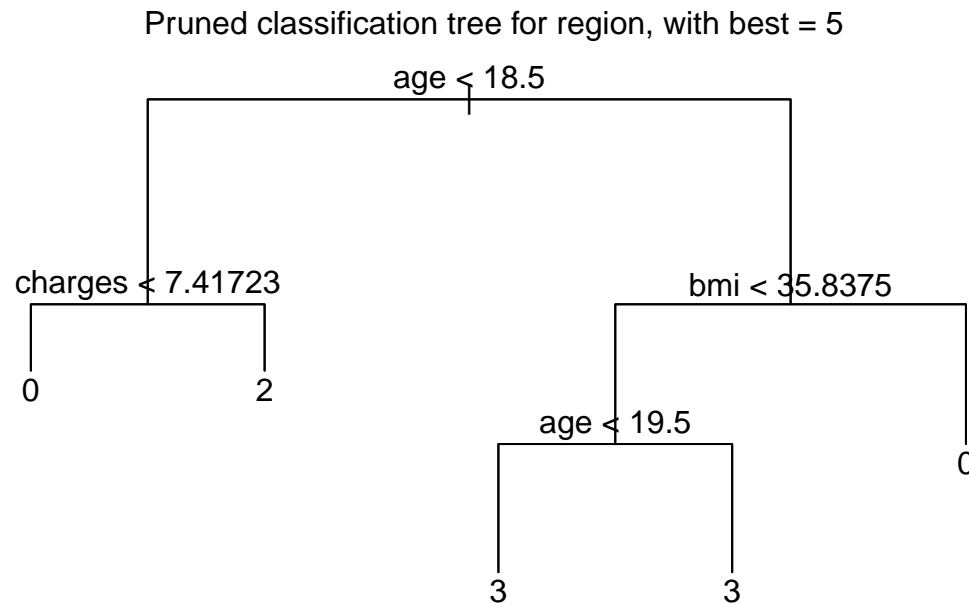


The residual mean deviance is at 2.633 whereas the misclassification error is at 65.28%.

```
prune.tree.classification.region2 <- prune.misclass(tree.classification.region, best=5)
summary(prune.tree.classification.region2)
```

```
##
## Classification tree:
## tree(formula = region ~ ., data = insurance, subset = train)
## Variables actually used in tree construction:
## [1] "age"      "charges"  "bmi"
## Number of terminal nodes:  5
## Residual mean deviance:  2.555 = 2379 / 931
## Misclassification error rate: 0.6421 = 601 / 936

plot(prune.tree.classification.region2)
text(prune.tree.classification.region2,pretty=0)
mtext("Pruned classification tree for region, with best = 5 \n")
```



The residual mean deviance is at 2.555 with a misclassification error rate of 0.6421.

We now make a prediction on our test data:

```

prune.tree.classification.region.pred <- predict(prune.tree.classification.region, insurance[-train,],
(prune.tree.classification.region.pred.ct <- table(prune.tree.classification.region.pred, insurance[-tr

##
## prune.tree.classification.region.pred  0  1  2  3
##               0 44 14 24  9
##               1  0  0  0  0
##               2  0  0  0  0
##               3 73 81 79 78
(prune.tree.classification.region.pred.correct <- sum(prune.tree.classification.region.pred==insurance[

## [1] 0.3034826
prune.tree.classification.region.pred.testError <- 1 - prune.tree.classification.region.pred.correct

As can be clearly seen, our rate to correctly classify the regions is still very low at 34.69%.
cat(tree.classification.region.pred.error)

## 0.8781095
cat("\n")

```



```
cat(prune.tree.classification.region.pred.testError)
```

```
## 0.6965174
```

However, despite the classification rate still being very low, in comparison, the error has drastically reduced in comparison with the unpruned tree.

7.4 Bagging

```
set.seed(5)
train <- sample(1:nrow(insurance), nrow(insurance)/2)
insurance.test=insurance[-train,]
```

We now are bagging our dataset. We furthermore say that we are happy with an estimation and not the real crossvalidation.

```
bag.insurance <- bagging(charges~., data=insurance, subset=train, nbagg=20, coob = TRUE)
print(bag.insurance)
```

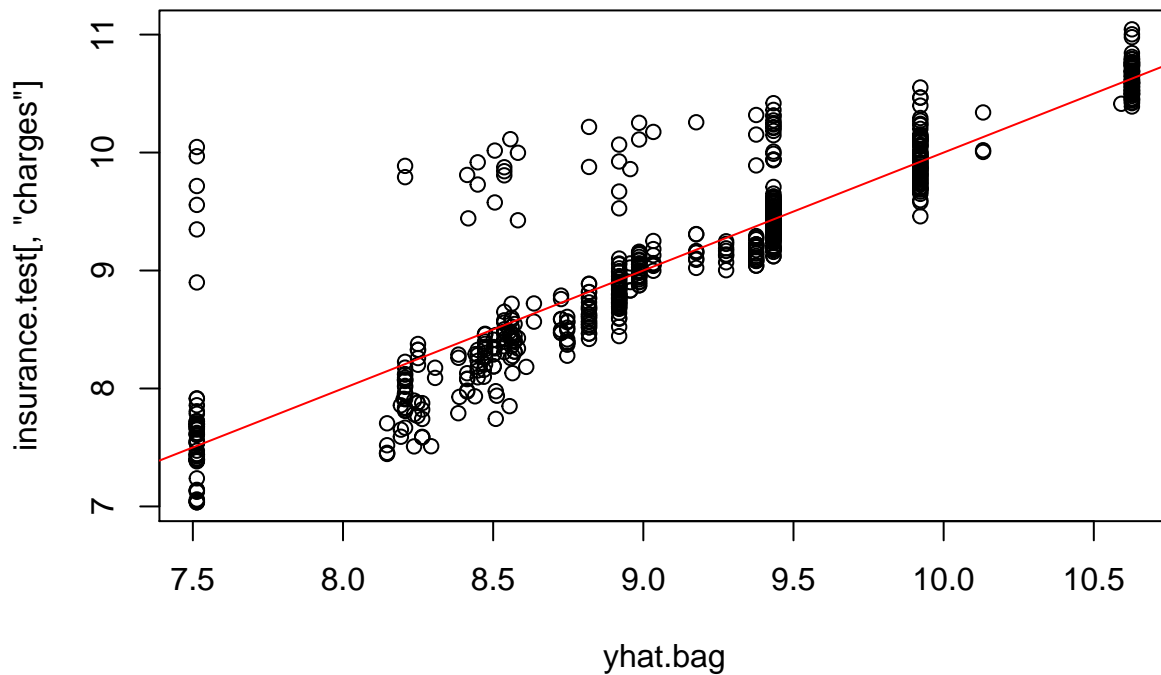
```
##
## Bagging regression trees with 20 bootstrap replications
##
## Call: bagging.data.frame(formula = charges ~ ., data = insurance, subset = train,
##      nbagg = 20, coob = TRUE)
##
## Out-of-bag estimate of root mean squared error: 0.391
```

We here receive a root mean squared error of 0.391 → mean squared error = 0.152881 (as the results one receives is rooted).

We now predict performance from our bagged model:

```
yhat.bag <- predict(bag.insurance, newdata=insurance[-train,])
plot(yhat.bag, insurance.test[, "charges"], main="Prediction against real data with test data")
abline(0,1, col="red")
```

Prediction against real data with test data



On the x-axis, we see the prediction, whereas on the y-axis, the real value (logarithmic).

We also compute the mean: in comparison with the previously calculated mean value, we can say they are rather related than not, however with some deviance.

```
mean((yhat.bag-insurance.test[, "charges"])^2)
```

```
## [1] 0.1575089
```

We also test this with the random forest way to do it.

```
## test with randomForest
```

```
bag.insurance.2 <- randomForest(charges~., data=insurance, subset=train, mtry=6, importance =TRUE)
print(bag.insurance.2)
```

```
##
```

```
## Call:
```

```
## randomForest(formula = charges ~ ., data = insurance, mtry = 6, importance = TRUE, subset = tr
```

```
## Type of random forest: regression
```

```
## Number of trees: 500
```

```
## No. of variables tried at each split: 6
```

```
##
```

```
## Mean of squared residuals: 0.1497648
```

```
## % Var explained: 82.09
```

```
importance(bag.insurance.2)
```

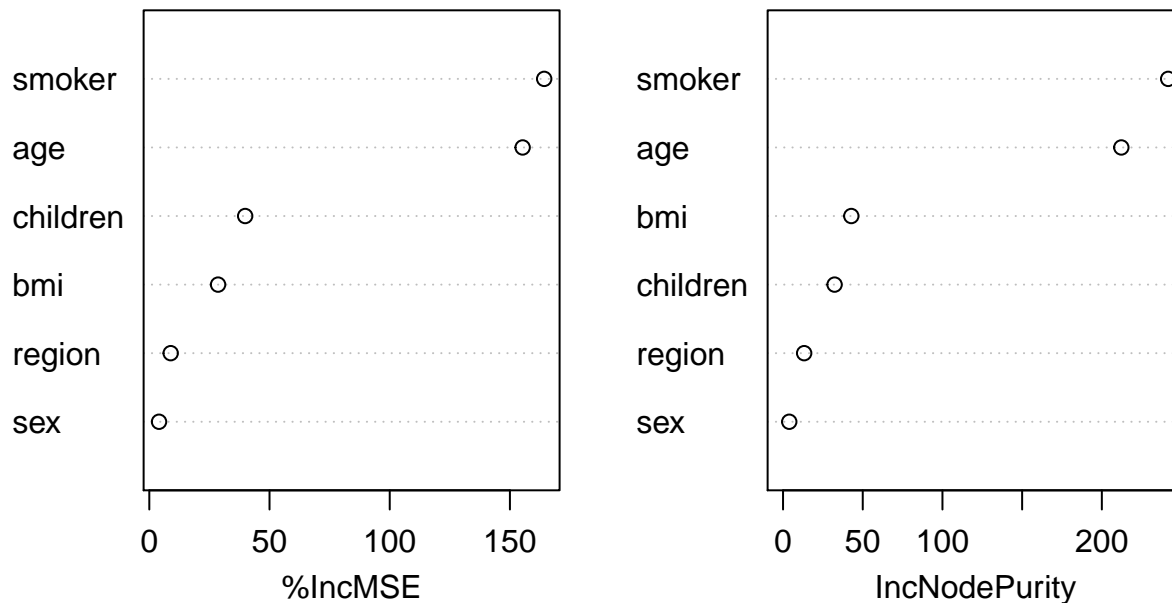
```
## %IncMSE IncNodePurity
```

```
## age 155.286361 212.223089
```

```
## sex      4.046307      3.899614
## bmi      28.599394     42.827102
## children 39.871574     32.346164
## smoker   164.251751    241.621328
## region    8.898203     13.246871
```

```
varImpPlot(bag.insurance.2)
```

bag.insurance.2



We here see that age and if one smokes are of biggest influence on charges. Age also seems fairly important.

7.5 Random Forests

```
set.seed(10)
rf.insurance <- randomForest(charges~.,data=insurance,subset=train, mtry=2,importance =TRUE)
yhat.rf <- predict(rf.insurance ,newdata=insurance[-train ,])
mean((yhat.rf-insurance.test[, "charges"])^2)
```

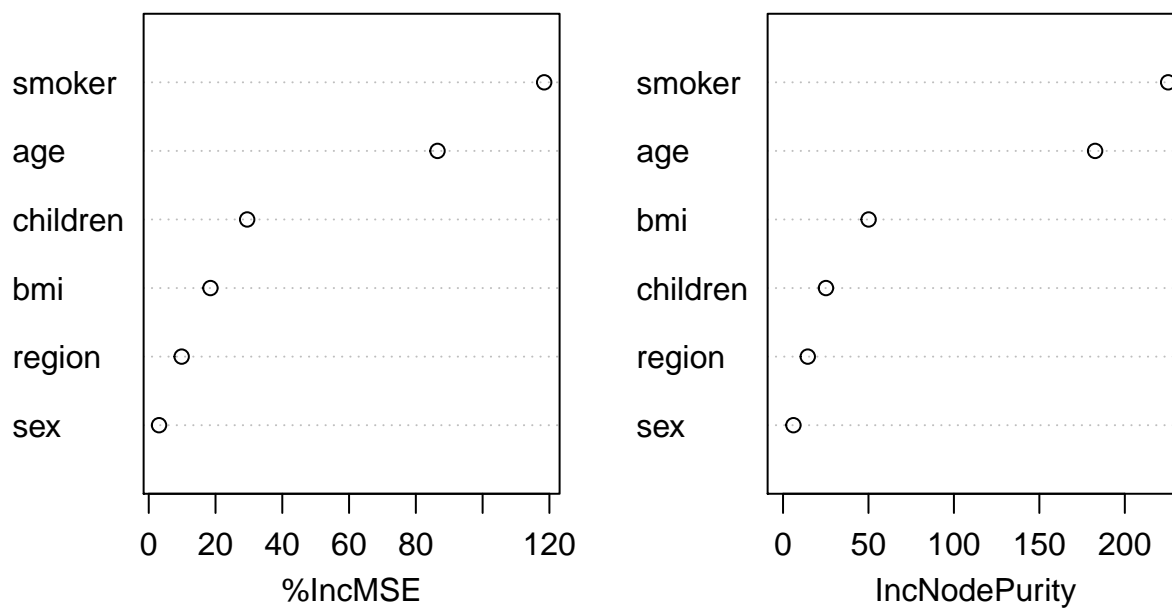
```
## [1] 0.158069
```

```
importance(rf.insurance)
```

```
##          %IncMSE IncNodePurity
## age      86.489219    182.623763
## sex       3.084226      6.073046
## bmi      18.476454     50.077296
## children 29.481898     25.129388
## smoker   118.440509    225.385251
## region    9.858480     14.461133
```

```
varImpPlot (rf.insurance)
```

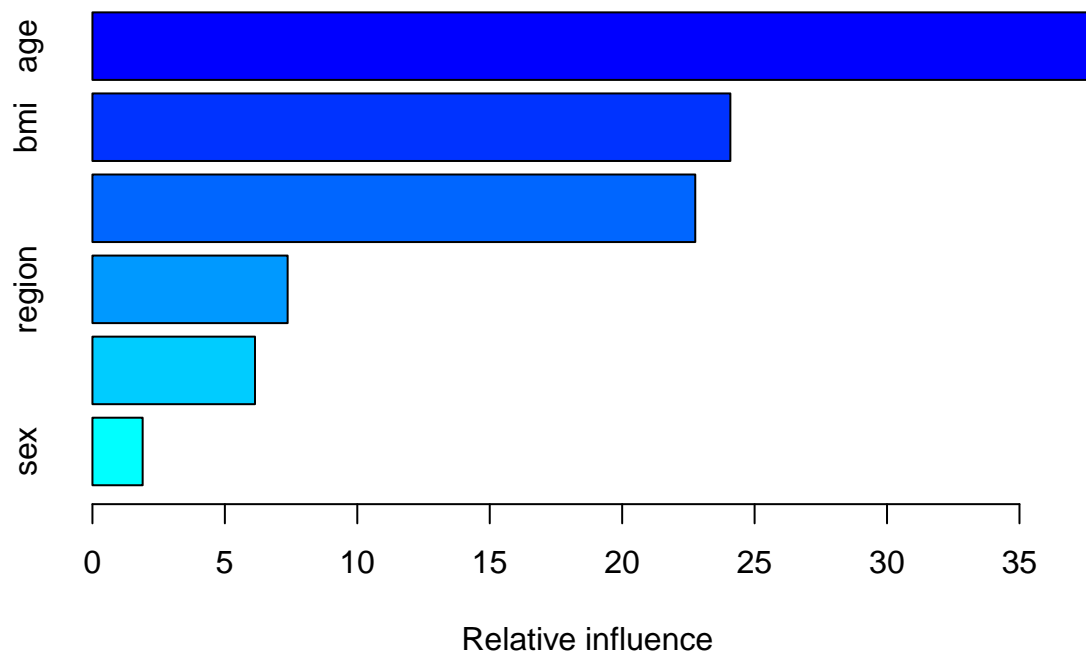
rf.insurance



7.6 Boosting

We now perform boosting by a generalized boosted regression modelling:

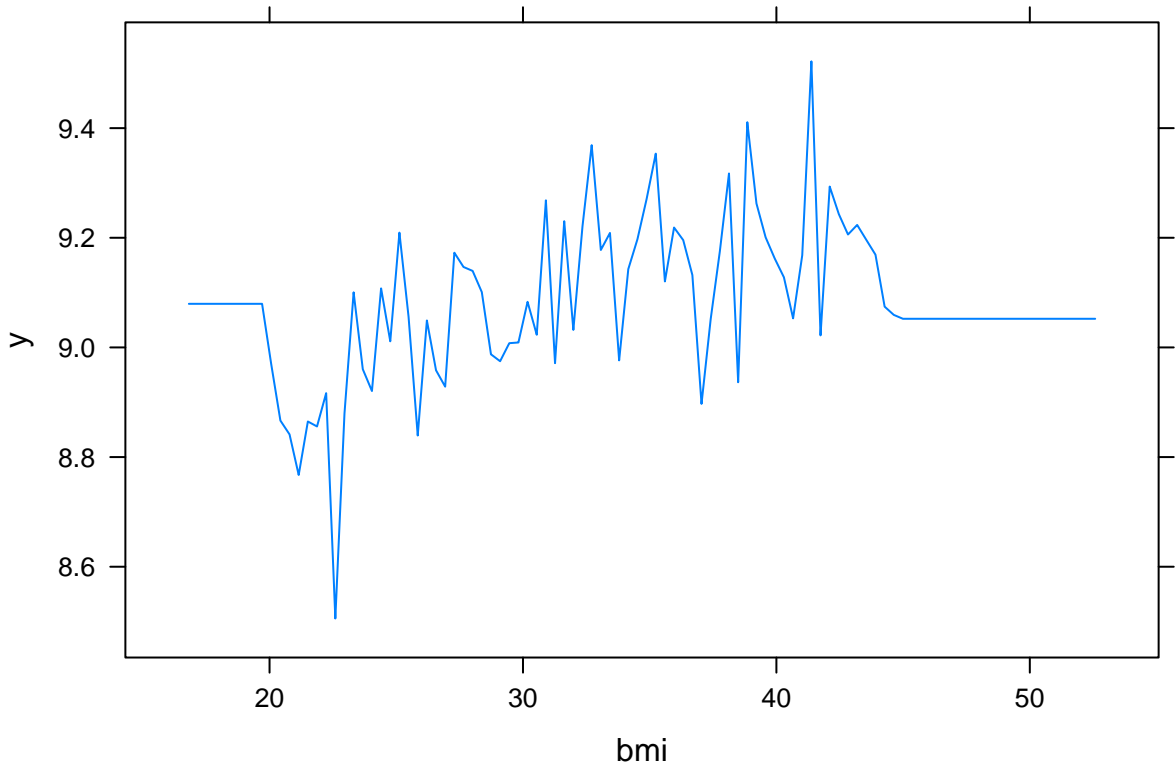
```
set.seed (1)
boost.insurance <- gbm(charges~.,data=insurance[train,],
                       distribution="gaussian",n.trees=5000, interaction.depth=4)
summary(boost.insurance)
```



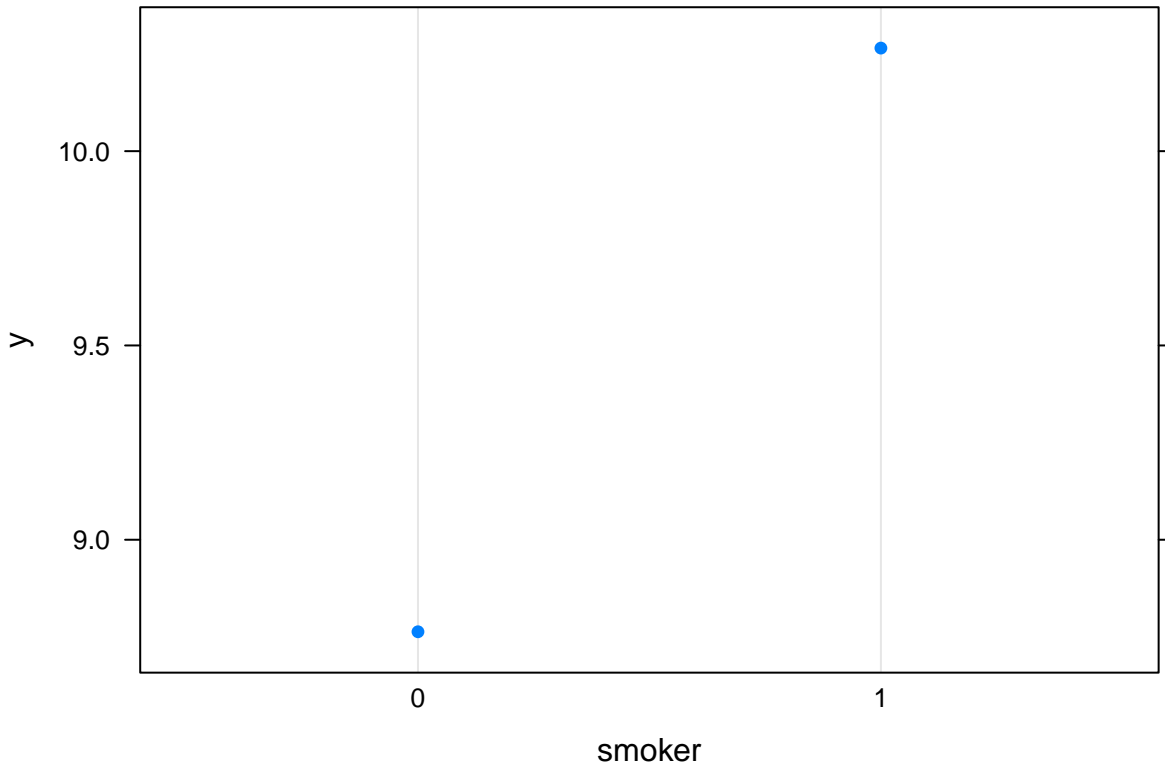
```
##           var    rel.inf
## age         age 37.751923
## bmi         bmi 24.084978
## smoker      smoker 22.759940
## region      region  7.370245
## children    children 6.137392
## sex         sex  1.895522
```

We now show partial dependence the four most important factors:

```
plot(boost.insurance ,i="bmi")
```



```
plot(boost.insurance ,i="smoker")
```



We now do a prediction with the boosted model.

```
yhat.boost <- predict(boost.insurance, newdata=insurance[-train,], n.trees=5000)
mean((yhat.boost - insurance.test[, "charges"])^2)
```

```
## [1] 0.2383602
```

The output is our MSE = 0.2383602.

We now are resetting shrinkage /lambda parameter which expresses basically that we have an overfit effect faster:

```
boost.insurance = gbm(charges~., data=insurance[train,], distribution="gaussian", n.trees=5000, interaction=1)
yhat.boost2=predict(boost.insurance, newdata=insurance[-train,], n.trees=5000)
mean((yhat.boost2 - insurance.test[, "charges"])^2)
```

```
## [1] 0.193239
```

We now receive a smaller error rate, so one could say this is a better model.

And we do the same thing again but on the other side:

```
boost.insurance2=gbm(charges~., data=insurance[train,], distribution="gaussian", n.trees=5000, interaction=1)
yhat.boost3=predict(boost.insurance2, newdata=insurance[-train,], n.trees=5000)
mean((yhat.boost3 - insurance.test[, "charges"])^2)
```

```
## [1] 0.2383602
```

We now have a higher test MSE which means this model is worse as overfitting has already started.

(Attention: Please do only execute the code block (applying the exponential function) below once as executing

it several times will apply the function several times)

```
insurance$charges <- exp(insurance$charges) #resetting the insurance variable
```

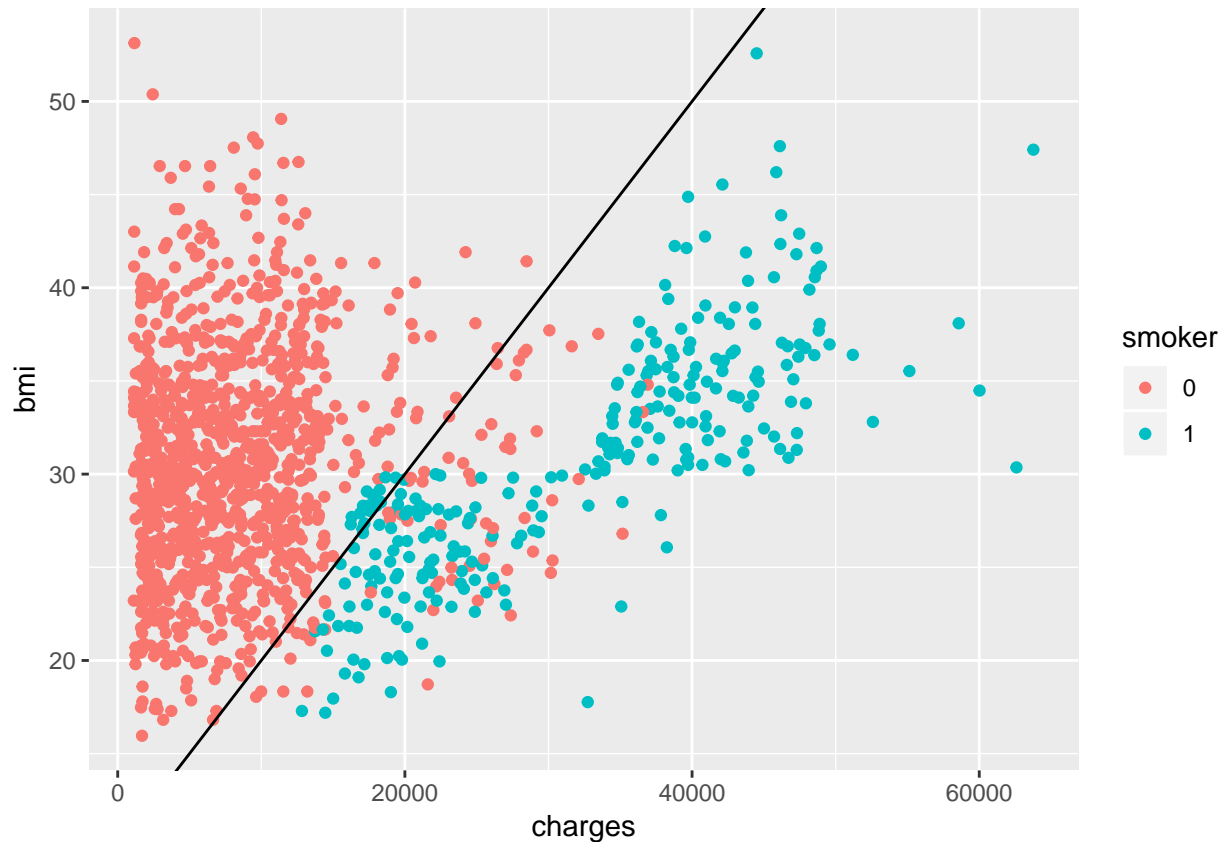
8 Support Vector Machines

In this chapter, Support Vector Machines are applied to the insurance dataset.

8.1 two classes linear

By visual exploration, we discover that smokers and non-smokers form two groups. We fit an SVM for those two groups. The lines in the graph are the first guess of how the data could be divided.

```
# guess of hyperplane
slope <- 0.001
intercept <- 10
ggplot(data = insurance) +
  geom_point(aes(x = charges, y = bmi, color=smoker)) +
  geom_abline(slope = slope, intercept = intercept)
```



A linear SVM is tuned using different cost ranges.

```
# prepare data

train.YES <- sample(x=c(TRUE,FALSE),
                    size=nrow(insurance),
                    replace = TRUE)
```



```

table(train.YES)

## train.YES
## FALSE  TRUE
##    688    650

insurance.train <- insurance[train.YES, ]
insurance.test  <- insurance[!train.YES, ]

set.seed(5)
#get optimal cost
cost_range <-
  c(1e-10,
    1e-7,
    1e-5,
    0.001,
    0.0025,
    0.005,
    0.0075,
    0.01,
    0.1,
    1,
    5,
    10,
    100,
    200)
tune.out.1 <- tune(
  svm,
  smoker ~ charges+bmi,
  data = insurance.train,
  kernel = "linear",
  ranges = list(cost = cost_range)
)

```

The following is the best model. A cost of 5 is used. There are 104 support vectors.

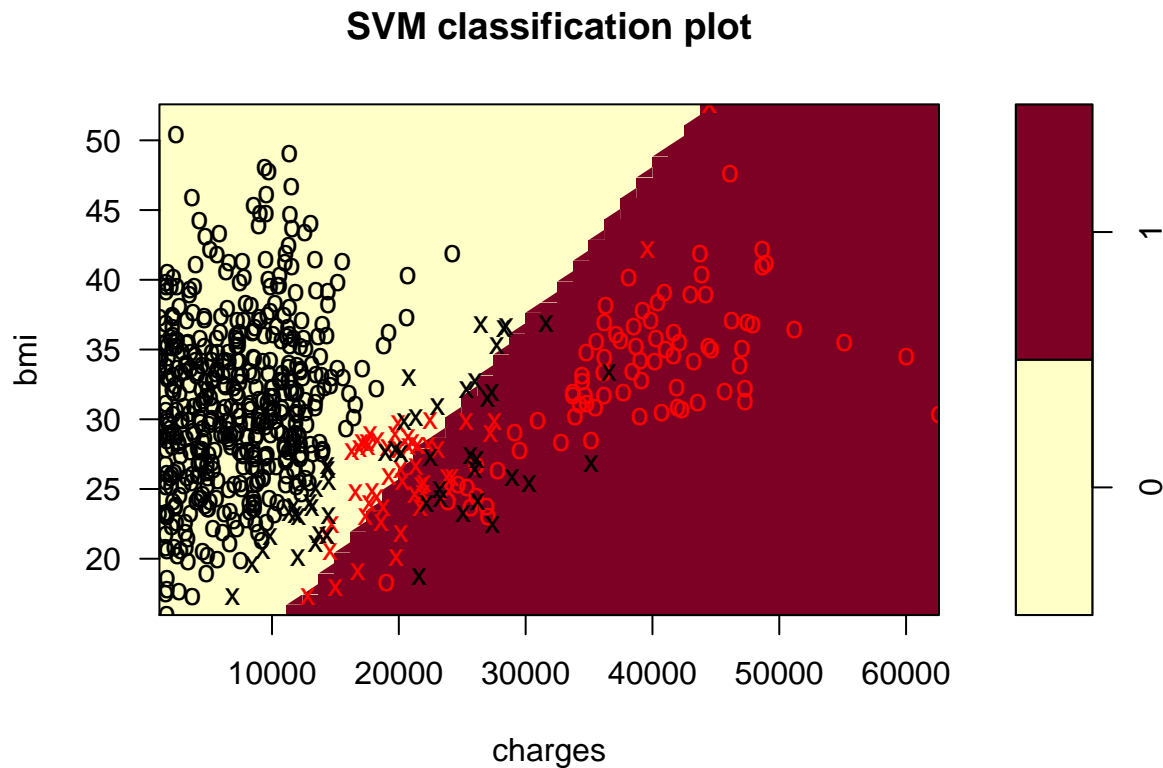
```

#best model
bestmod <- tune.out.1$best.model
summary(bestmod)

##
## Call:
## best.tune(method = svm, train.x = smoker ~ charges + bmi, data = insurance.train,
##   ranges = list(cost = cost_range), kernel = "linear")
##
##
## Parameters:
##   SVM-Type:  C-classification
##   SVM-Kernel: linear
##         cost: 10
##
## Number of Support Vectors: 100
##
## ( 50 50 )
##

```

```
##
## Number of Classes: 2
##
## Levels:
## 0 1
plot(bestmod, insurance.train, bmi~charges)
```



The confusion matrices and the classification error rate for train and test data are:

```
# confusion matrix train
confusion.lin.train <- table(predict = predict(bestmod, insurance.train),
                             truth = insurance.train$smoker)
confusion.lin.train

##          truth
## predict   0   1
##        0 501  20
##        1  20 109

#classification error rate train
(confusion.lin.train[1,2]+confusion.lin.train[2,1])/sum(confusion.lin.train[1:2,1:2])

## [1] 0.06153846

# confusion matrix linear test
confusion.lin.test <- table(predict = predict(bestmod, insurance.test),
                             truth = insurance.test$smoker)
confusion.lin.test
```

```
##      truth
## predict 0  1
##      0 523 26
##      1  20 119

#classification error rate test
(confusion.lin.test[1,2]+confusion.lin.test[2,1])/sum(confusion.lin.test[1:2,1:2])

## [1] 0.06686047
```

8.2 two classes polynomial

A polynomial model was also tested but the classification error rate was slightly higher:

```
set.seed(5)
#get optimal cost and degree for polynomial
cost_range <-
  c(1e-10,
    1e-7,
    1e-5,
    0.001,
    0.0025,
    0.005,
    0.0075,
    0.01,
    0.1,
    1,
    5,
    10,
    100)
degree_range <- 2:4
tune.out.poly <- tune(
  svm,
  smoker ~ charges+bmi,
  data = insurance.train,
  kernel = "polynomial",
  ranges = list(cost = cost_range, degree = degree_range))
#best model poly
tune.out.poly$best.parameters

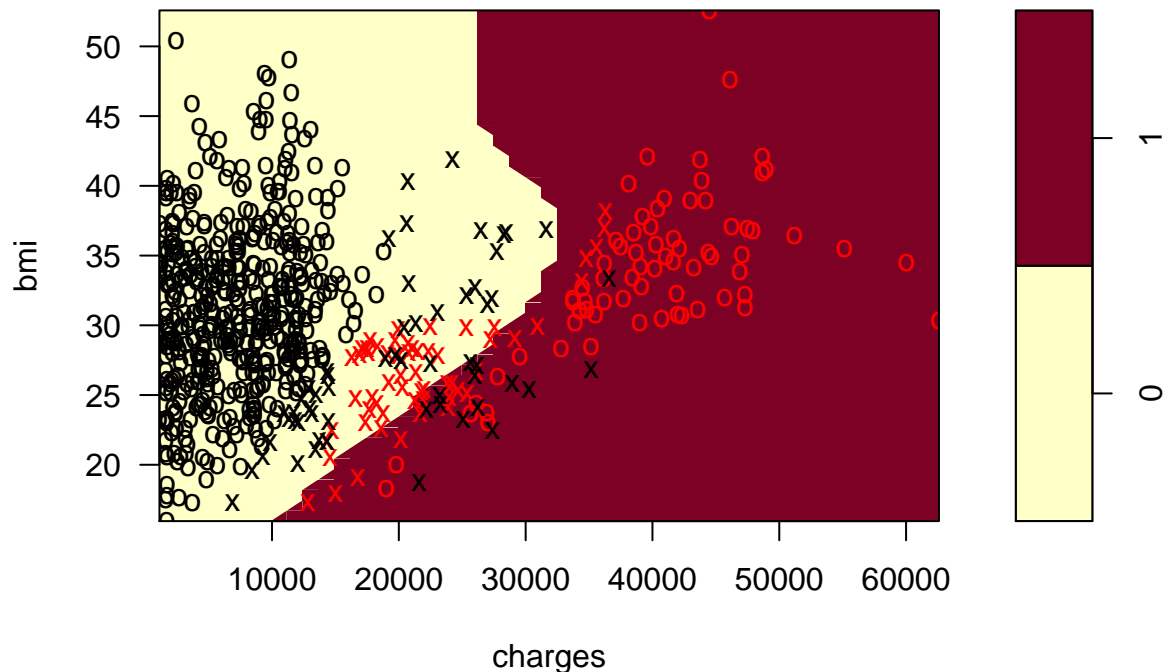
##      cost degree
## 24      5      3

bestmod.poly <- tune.out.poly$best.model
summary(bestmod.poly)

##
## Call:
## best.tune(method = svm, train.x = smoker ~ charges + bmi, data = insurance.train,
##      ranges = list(cost = cost_range, degree = degree_range),
##      kernel = "polynomial")
##
##
## Parameters:
##      SVM-Type:  C-classification
##      SVM-Kernel: polynomial
```

```
##      cost: 5
##     degree: 3
##    coef.0: 0
##
## Number of Support Vectors: 115
##
## ( 58 57 )
##
##
## Number of Classes: 2
##
## Levels:
## 0 1
plot(bestmod.poly, insurance.train, bmi~charges)
```

SVM classification plot



```
# confusion matrix
confusion.poly.train <- table(predict = predict(bestmod.poly, insurance.train),
                              truth = insurance.train$smoker)
confusion.poly.test <- table(predict = predict(bestmod.poly, insurance.test),
                              truth = insurance.test$smoker)
#classification error rate
(confusion.poly.train[1,2]+confusion.poly.train[2,1])/sum(confusion.poly.train[1:2,1:2])

## [1] 0.07230769
(confusion.poly.test[1,2]+confusion.poly.test[2,1])/sum(confusion.poly.test[1:2,1:2])
```

```
## [1] 0.07122093
```

9 Neural Networks

9.1 Neural Networks using “neuralnet”

In this section we will create and train a neural network which should be able to classify if a person is a smoker (output) based on the information of age, bmi and charges that are provided as input factors for the model. The assumption of the input variables is retrieved from our analysis already done in section 6. Generalized Linear Models.

9.1.1 Cleaning and Normalization

This part will be needed to ensure a cleansed and normalized database for a proper modelling of neural networks.

```
#setwd("~/GitHub/machinelearning/machinelearning/03_rmarkdown") #adrianas extrawurst
insurance <- read.csv("../01_data/insurance.csv", header=TRUE)
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 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 ...
```

```
# smoker = 1 / nonsmoker = 2
insurance$smoker <- as.character(insurance$smoker)
insurance$smoker[insurance$smoker == "yes"] <- "1"
insurance$smoker[insurance$smoker == "no"] <- "2"
insurance$smoker <- as.factor(insurance$smoker)
# female = 1 / male = 0
insurance$sex <- as.character(insurance$sex)
insurance$sex[insurance$sex == "female"] <- "1"
insurance$sex[insurance$sex == "male"] <- "0"
insurance$sex <- as.factor(insurance$sex)
# region / SE = 1 / SW = 0 / NE = 2 / NW = 3
insurance$region <- as.character(insurance$region)
insurance$region[insurance$region == "southwest"] <- "1"
insurance$region[insurance$region == "southeast"] <- "0"
insurance$region[insurance$region == "northeast"] <- "2"
insurance$region[insurance$region == "northwest"] <- "3"
insurance$region <- as.factor(insurance$region)
head(insurance)
```

```
## age sex bmi children smoker region charges
## 1 19 1 27.900 0 1 1 16884.924
## 2 18 0 33.770 1 2 0 1725.552
## 3 28 0 33.000 3 2 0 4449.462
## 4 33 0 22.705 0 2 3 21984.471
## 5 32 0 28.880 0 2 3 3866.855
## 6 31 1 25.740 0 2 0 3756.622
```

```

options(scipen = 99) # penalty for displaying scientific notation
options(digits = 4) # suggested number of digits to display

data <- dplyr::select(insurance, smoker, charges, age, bmi)

# Min-Max Normalization
data$charges <- (data$charges - min(data$charges))/(max(data$charges) - min(data$charges))
data$age <- (data$age - min(data$age))/(max(data$age) - min(data$age))
data$bmi <- (data$bmi - min(data$bmi))/(max(data$bmi) - min(data$bmi))

```

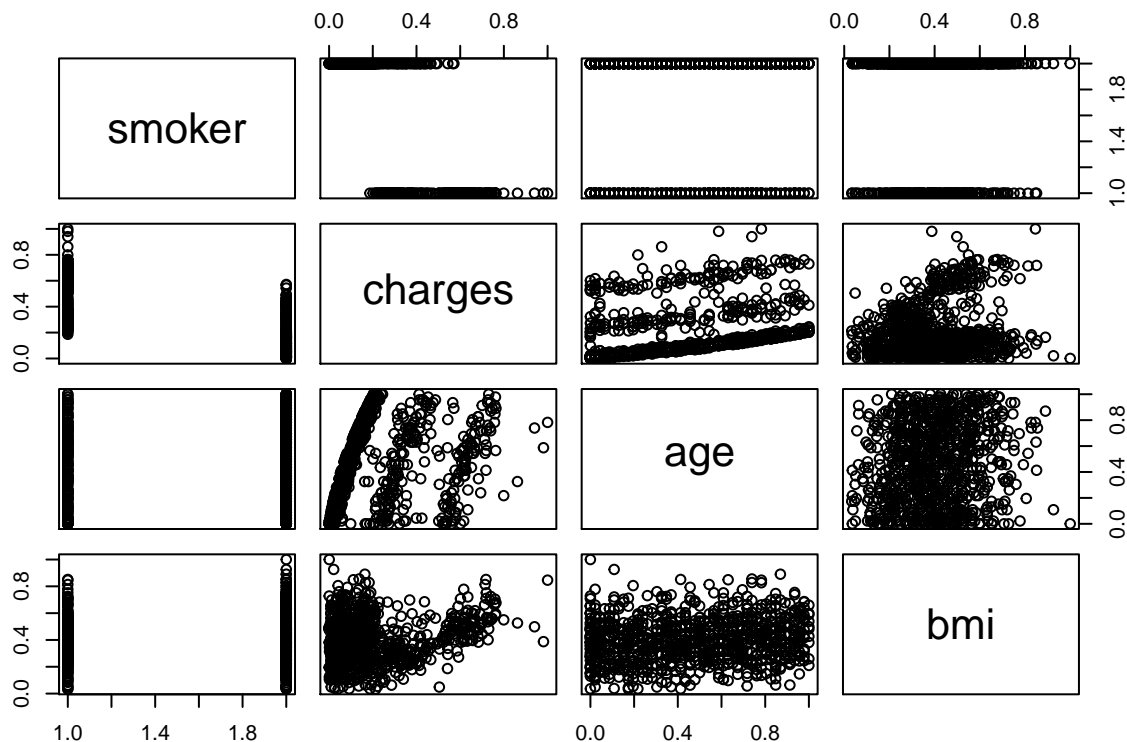
9.1.2 Partition into training and test subsets

As shown above we already cleansed and normalized data. Now we will create two subsets of the insurance data set. Therefore we will execute a partition, so we will have a training set covering 70% of the original dataset while the test set will cover the remaining 30%.

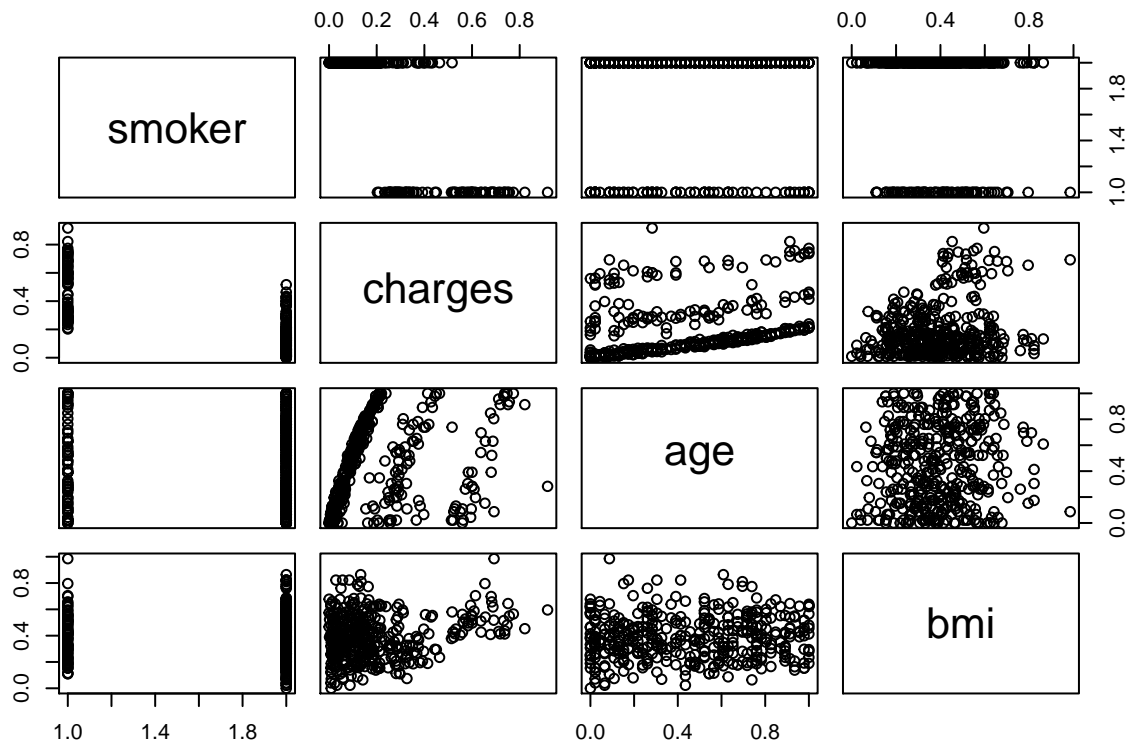
```

set.seed(222)
ind <- sample(2, nrow(data), replace = TRUE, prob = c(0.7, 0.3))
training <- data[ind==1,]
testing <- data[ind==2,]
par(mfrow = c(1, 1))
plot(training)

```



```
plot(testing)
```

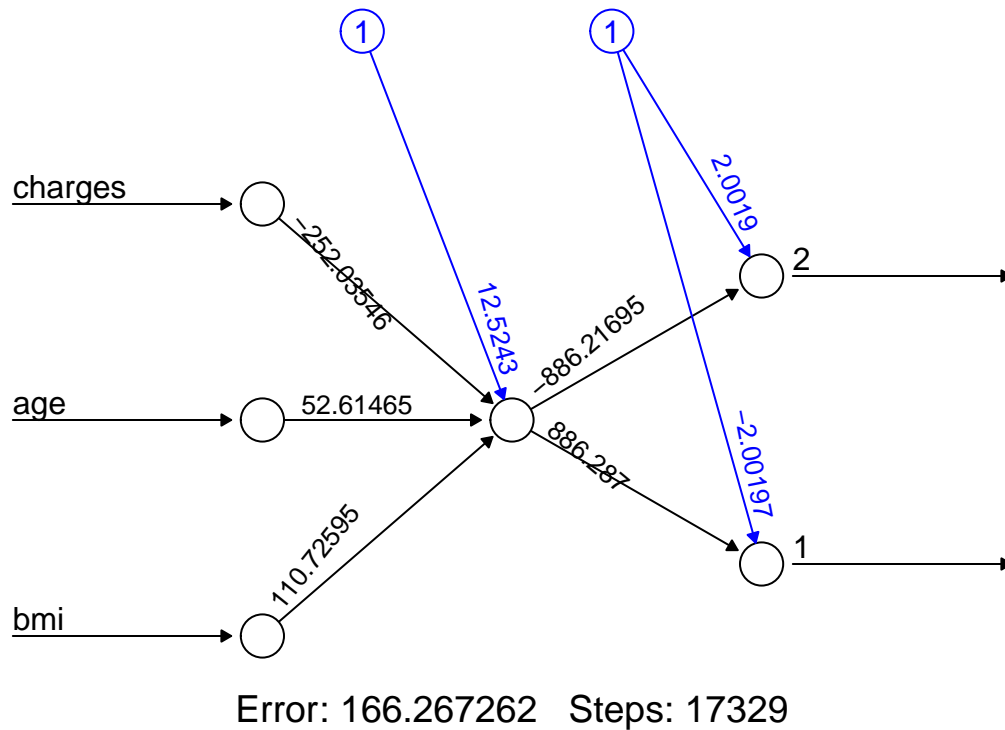


9.1.3 Creation and Training of the Neural Network

9.1.3.1 Not deep NN - one single hidden layer

9.1.3.1.1 Creation of NN_1

```
set.seed(444)
nn <- neuralnet(smoker~charges + age + bmi,
  data = training,
  hidden = 1,
  err.fct = "ce",
  linear.output = FALSE)
#summary(nn)
#print(nn)
plot(nn, rep='best')
```



```
##### Prediction on training subset
prediction_nn <- predict(nn,training)
#prediction_nn
class_nn <- apply(prediction_nn, 1, which.max)
```

9.1.3.1.2 Confusion Matrix for training

```
table(training$smoker, class_nn)
```

```
##      class_nn
##      1      2
##  1 197      0
##  2  27 711
```

```
#training$smoker[1]
#class_nn[1]
```

9.1.3.1.3 Performance for training

```
levels(training$smoker)[class_nn[1]]
```

```
## [1] "2"
```

```
corrects=sum(levels(training$smoker)[class_nn]==training$smoker)
errors=sum(levels(training$smoker)[class_nn]!=training$smoker)
(perf.nn=corrects/(corrects+errors))
```

```
## [1] 0.9711
```



```
print(paste('training accuracy: ',perf.nn*100,"%"))
```

```
## [1] "training accuracy: 97.1122994652406 %"
```

9.1.3.1.4 Prediction on test subset

```
prediction_nn <- predict(nn,testing)
#prediction_nn
class_nn <- apply(prediction_nn, 1, which.max)
```

9.1.3.1.5 Confusion Matrix for testing

```
table(testing$smoker, class_nn)
```

```
##      class_nn
##         1    2
##    1  75    2
##    2  13   313
```

```
#training$smoker[1]
#class_nn[1]
```

9.1.3.1.6 Performance for testing

```
levels(testing$smoker)[class_nn[1]]
```

```
## [1] "1"
```

```
corrects=sum(levels(testing$smoker)[class_nn]==testing$smoker)
errors=sum(levels(testing$smoker)[class_nn]!=testing$smoker)
(perf.nn=corrects/(corrects+errors))
```

```
## [1] 0.9628
```

```
print(paste('testing accuracy: ',perf.nn*100,"%"))
```

```
## [1] "testing accuracy: 96.2779156327543 %"
```

9.1.3.2 Deep NN - two hidden layers

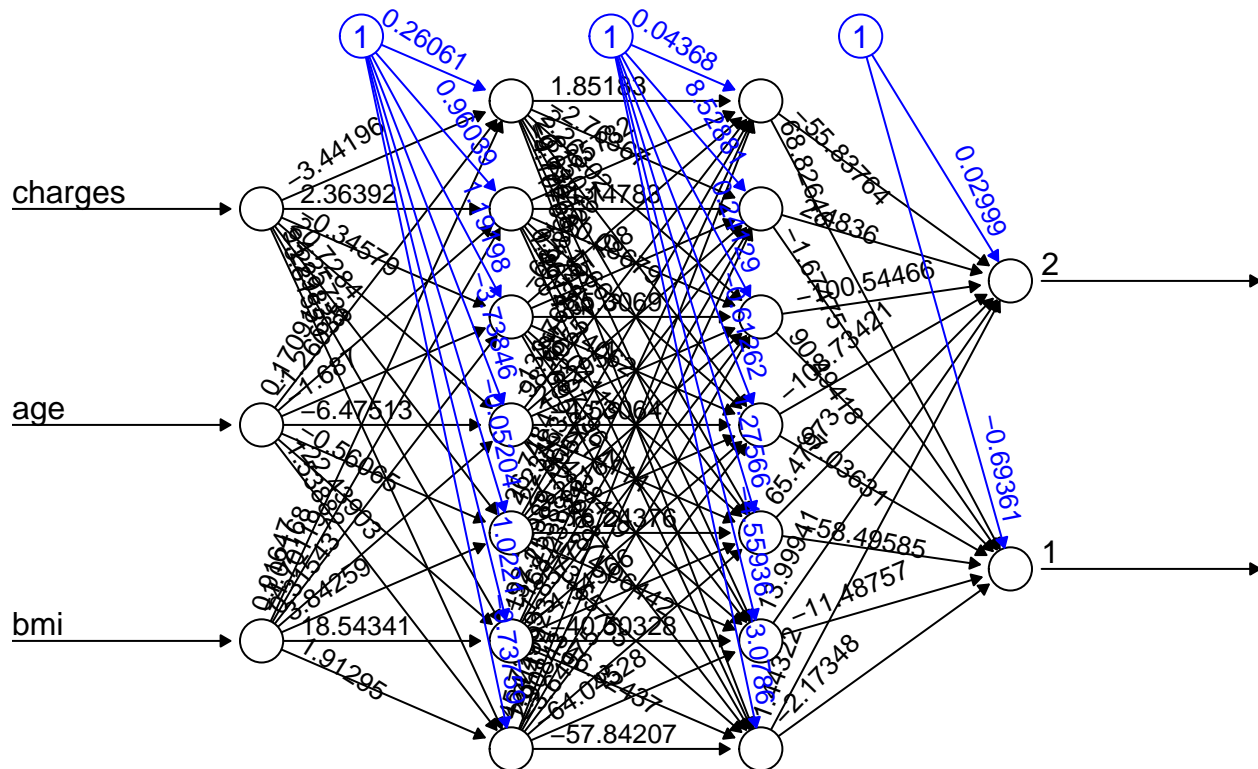
9.1.3.2.1 Creation of NN_2

```
set.seed(444)
nn2 <- neuralnet(smoker~charges + age + bmi,
  data = training,
  hidden = c(7,7),
  threshold = 0.0025,
  lifesign = 'full',
  lifesign.step = 500,
  linear.output = FALSE)
```

```
## hidden: 7, 7      thresh: 0.0025      rep: 1/1      steps:      500 min thresh: 0.11825064758526
##                                                         1000 min thresh: 0.112617356849393
##                                                         1500 min thresh: 0.112617356849393
##                                                         2000 min thresh: 0.107500337348156
##                                                         2500 min thresh: 0.0500175621736521
##                                                         3000 min thresh: 0.0219775031983741
```

```
## 3500 min thresh: 0.0129098584782895
## 4000 min thresh: 0.00367818486599488
## 4441 error: 19.00295 time: 7.28 secs
```

```
#summary(nn)
#print(nn)
plot(nn2, rep='best')
```



Error: 19.002947 Steps: 4441

9.1.3.2.2 Prediction on training subset

```
prediction_nn2 <- predict(nn2,training)
#prediction_nn
class_nn <- apply(prediction_nn2, 1, which.max)
```

9.1.3.2.3 Confusion Matrix for training

```
table(training$smoker, class_nn)
```

```
##   class_nn
##      1   2
## 1 197   0
## 2  19 719
```

```
#training$smoker[1]
#class_nn[1]
```

9.1.3.2.4 Performance for training

```

levels(training$smoker)[class_nn[1]]

## [1] "2"

corrects=sum(levels(training$smoker)[class_nn]==training$smoker)
errors=sum(levels(training$smoker)[class_nn]!=training$smoker)
(perf.nn=corrects/(corrects+errors))

## [1] 0.9797

print(paste('training accuracy: ',perf.nn*100,"%"))

## [1] "training accuracy: 97.9679144385027 %"

```

9.1.3.2.5 Prediction on test subset

```

prediction_nn2 <- predict(nn2,testing)
#prediction_nn
class_nn <- apply(prediction_nn2, 1, which.max)

```

9.1.3.2.6 Confusion Matrix for testing

```

table(testing$smoker, class_nn)

```

```

##      class_nn
##      1      2
##  1  74      3
##  2  12    314

#training$smoker[1]
#class_nn[1]

```

9.1.3.2.7 Performance for testing

```

levels(testing$smoker)[class_nn[1]]

## [1] "1"

corrects=sum(levels(testing$smoker)[class_nn]==testing$smoker)
errors=sum(levels(testing$smoker)[class_nn]!=testing$smoker)
(perf.nn=corrects/(corrects+errors))

## [1] 0.9628

print(paste('testing accuracy: ',perf.nn*100,"%"))

## [1] "testing accuracy: 96.2779156327543 %"

```

9.1.3.3 Deep NN - three hidden layers

9.1.3.3.1 Creation of NN_3

```

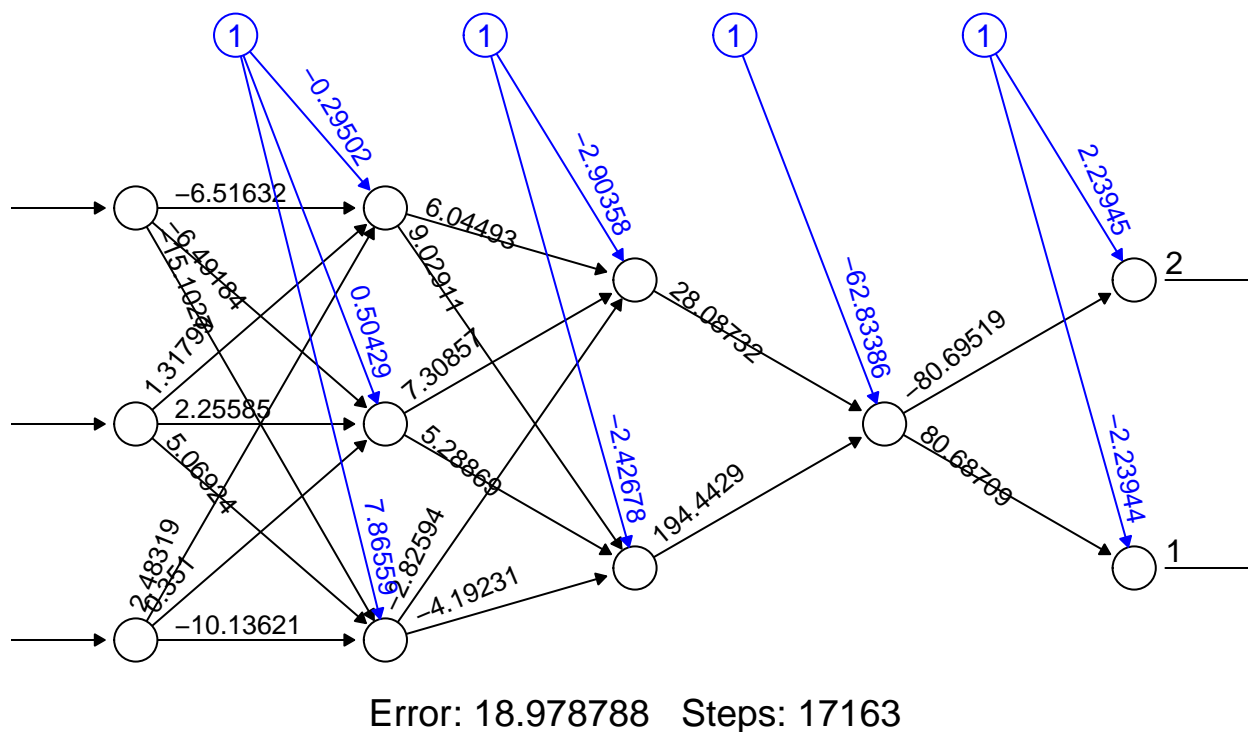
set.seed(444)
nn3 <- neuralnet(smoker~charges + age + bmi,
  data = training,
  hidden = c(3,2,1),
  threshold = 0.0025,
  lifesign = 'full',

```

```
lifesign.step = 500,
linear.output = FALSE)
```

```
## hidden: 3, 2, 1    thresh: 0.0025    rep: 1/1    steps:      500 min thresh: 0.0938559883189038
##                                                           1000 min thresh: 0.0938559883189038
##                                                           1500 min thresh: 0.0938559883189038
##                                                           2000 min thresh: 0.0938559883189038
##                                                           2500 min thresh: 0.0779042619497504
##                                                           3000 min thresh: 0.0477729129630322
##                                                           3500 min thresh: 0.0477729129630322
##                                                           4000 min thresh: 0.0477729129630322
##                                                           4500 min thresh: 0.0477729129630322
##                                                           5000 min thresh: 0.0477729129630322
##                                                           5500 min thresh: 0.0477729129630322
##                                                           6000 min thresh: 0.0477729129630322
##                                                           6500 min thresh: 0.0477729129630322
##                                                           7000 min thresh: 0.0477729129630322
##                                                           7500 min thresh: 0.0401919383687264
##                                                           8000 min thresh: 0.0319441973901314
##                                                           8500 min thresh: 0.0271607588021149
##                                                           9000 min thresh: 0.0231670821049527
##                                                           9500 min thresh: 0.0231670821049527
##                                                           10000 min thresh: 0.0160246751435115
##                                                           10500 min thresh: 0.0143015470518307
##                                                           11000 min thresh: 0.00978614370100106
##                                                           11500 min thresh: 0.00978614370100106
##                                                           12000 min thresh: 0.00813136952436977
##                                                           12500 min thresh: 0.00670291691221209
##                                                           13000 min thresh: 0.00588460075826878
##                                                           13500 min thresh: 0.00568340157694122
##                                                           14000 min thresh: 0.00431292899953402
##                                                           14500 min thresh: 0.00405437903080155
##                                                           15000 min thresh: 0.00396128911212355
##                                                           15500 min thresh: 0.00308676494427534
##                                                           16000 min thresh: 0.00289723264713905
##                                                           16500 min thresh: 0.00274623242108184
##                                                           17000 min thresh: 0.00274623242108184
##                                                           17163 error: 18.97879 time: 12.78 secs
```

```
#summary(nn)
#print(nn)
plot(nn3, rep='best')
```



9.1.3.3.2 Prediction on training subset

```
prediction_nn3 <- predict(nn3,training)
#prediction_nn
class_nn <- apply(prediction_nn3, 1, which.max)
```

9.1.3.3.3 Confusion Matrix for training

```
table(training$smoker, class_nn)
```

```
##      class_nn
##      1      2
##  1 197    0
##  2  21 717
```

```
#training$smoker[1]
#class_nn[1]
```

9.1.3.3.4 Performance for training

```
levels(training$smoker)[class_nn[1]]
```

```
## [1] "2"
```

```
corrects=sum(levels(training$smoker)[class_nn]==training$smoker)
errors=sum(levels(training$smoker)[class_nn]!=training$smoker)
(perf.nn=corrects/(corrects+errors))
```

```
## [1] 0.9775
print(paste('training accuracy: ',perf.nn*100,"%"))
```

```
## [1] "training accuracy: 97.7540106951872 %"
```

9.1.3.3.5 Prediction on test subset

```
prediction_nn3 <- predict(nn3,testing)
#prediction_nn
class_nn <- apply(prediction_nn3, 1, which.max)
```

9.1.3.3.6 Confusion Matrix for testing

```
table(testing$smoker, class_nn)
```

```
##      class_nn
##      1      2
## 1  76      1
## 2   9    317
```

```
#training$smoker[1]
#class_nn[1]
```

9.1.3.3.7 Performance for testing

```
levels(testing$smoker)[class_nn[1]]
```

```
## [1] "1"
```

```
corrects=sum(levels(testing$smoker)[class_nn]==testing$smoker)
errors=sum(levels(testing$smoker)[class_nn]!=testing$smoker)
(perf.nn=corrects/(corrects+errors))
```

```
## [1] 0.9752
```

```
print(paste('testing accuracy: ',perf.nn*100,"%"))
```

```
## [1] "testing accuracy: 97.5186104218362 %"
```

9.2 Neural Networks using “h2o”

9.2.1 Creation NN

```
set.seed(99)
h2o.init()
```

```
## Connection successful!
##
```

```
## R is connected to the H2O cluster:
```

```
##      H2O cluster uptime:      3 hours 5 minutes
##      H2O cluster timezone:    Europe/Berlin
##      H2O data parsing timezone: UTC
##      H2O cluster version:     3.30.0.1
##      H2O cluster version age:  2 months and 6 days
##      H2O cluster name:        H2O_started_from_R_CaroleM_pal554
##      H2O cluster total nodes: 1
##      H2O cluster total memory: 1.94 GB
```

```
##      H2O cluster total cores:      4
##      H2O cluster allowed cores:    4
##      H2O cluster healthy:          TRUE
##      H2O Connection ip:             localhost
##      H2O Connection port:           54321
##      H2O Connection proxy:          NA
##      H2O Internal Security:         FALSE
##      H2O API Extensions:            Amazon S3, Algos, AutoML, Core V3, TargetEncoder, Core V4
##      R Version:                     R version 3.6.1 (2019-07-05)
```

```
train_hf <- as.h2o(training)
```

```
##
|
|                                     | 0%
|
|=====| 100%
```

```
model_dl <- h2o.deeplearning(x = 2:4, y = 1, training_frame = train_hf, nfolds = 10, seed=545, verbose = TRUE)
```

```
##
|
|                                     | 0%
|
|=====| 19%
## Scoring History for Model DeepLearning_model_R_1591727398193_201 at 2020-06-09 23:35:33
## [1] "Model Build is 19.122022% done..."
##
|
|=====| 39%
## Scoring History for Model DeepLearning_model_R_1591727398193_201 at 2020-06-09 23:35:34
## [1] "Model Build is 39.338842% done..."
##
|
|=====| 51%
## Scoring History for Model DeepLearning_model_R_1591727398193_201 at 2020-06-09 23:35:39
## [1] "Model Build is 50.7438% done..."
##
|
|=====| 65%
## Scoring History for Model DeepLearning_model_R_1591727398193_201 at 2020-06-09 23:35:42
## [1] "Model Build is 65.376765% done..."
##
|
|=====| 77%
## Scoring History for Model DeepLearning_model_R_1591727398193_201 at 2020-06-09 23:35:43
## [1] "Model Build is 76.84978% done..."
##
|
|=====| 81%
## Scoring History for Model DeepLearning_model_R_1591727398193_201 at 2020-06-09 23:35:48
## [1] "Model Build is 80.95576% done..."
##
|
|=====| 100%
```

```
## Scoring History for Model DeepLearning_model_R_1591727398193_201 at 2020-06-09 23:35:50
## [1] "Model Build is 99.80457% done..."
##
## Scoring History for Model DeepLearning_model_R_1591727398193_201 at 2020-06-09 23:35:51
## [1] "Model Build is 100% done..."
```

9.2.2 Prediction on training

```
predicted_h2o_deep_neural_network_probabilites_train <- h2o.predict(model_dl, train_hf)
```

```
##
|
|
|
|=====| 100%
```

```
table( as.data.frame(predicted_h2o_deep_neural_network_probabilites_train$predict))
```

```
##
##      1      2
## 235 700
```

```
predicted_h2o_deep_neural_network_class_train <- as.data.frame(predicted_h2o_deep_neural_network_probabilites_train$predict)
```

9.2.3 Performance on training

```
table(training$smoker, as.integer(predicted_h2o_deep_neural_network_class_train[,1]))
```

```
##
##      1      2
##      1 197      0
##      2  38 700
```

```
predicted_h2o_deep_neural_network_performance_train <- mean(as.integer(predicted_h2o_deep_neural_network_class_train[,1]) == training$smoker)
print(paste('training accuracy: ', predicted_h2o_deep_neural_network_performance_train, "%"))
```

```
## [1] "training accuracy: 95.9358288770053 %"
```

9.2.4 Prediction on testing

```
test_hf <- as.h2o(testing)
```

```
##
|
|
|
|=====| 100%
```

```
predicted_h2o_deep_neural_network_probabilites_test <- h2o.predict(model_dl, test_hf)
```

```
##
|
|
|
|=====| 100%
```



```
table( as.data.frame(predicted_h2o_deep_neural_network_probabilites_test$predict))
```

```
##
##      1      2
##    90 313
```

```
predicted_h2o_deep_neural_network_class_test <- as.data.frame(predicted_h2o_deep_neural_network_probabilites_test$predict)
```

9.2.5 Performance on testing

```
table(testing$smoker, as.integer(predicted_h2o_deep_neural_network_class_test[,1]))
```

```
##
##      1      2
##    1 76      1
##    2 14 312
```

```
predicted_h2o_deep_neural_network_performance_test <- mean(as.integer(predicted_h2o_deep_neural_network_class_test[,1]))
print(paste('testing accuracy: ', predicted_h2o_deep_neural_network_performance_test, "%"))
```

```
## [1] "testing accuracy: 96.2779156327543 %"
```

9.3 Results and method comparison

Neural Network	Training Performance	Testing Performance
NN	97.11 %	96.28 %
NN2	97.97 %	96.28 %
NN3	97.75 %	97.52 %
h2o	96.68 %	96.52 %

In the above table you can find the latest results for all created models which are nearly the same. Based on the sample results, there is no huge difference between using the neuralnet package or the h2o package. Moreover the parameters for the hidden layer were chosen by a trial and error approach. We tried to create different levels for the hidden parameter, to explore if the performance will increase with the given complexity. Since the results are not reflecting this assumption and all performance results are in the same range, all models can be used for an adequate classification, whether a person is a smoker or not.

10 Cross validation and discussion

10.1 Model overview

Throughout the analysis, many different model are applied. Here is a quick overview of the models used that performed the best:

10.1.1 Linear models (standard lm)

```
lm.insurance.3$call #best lm
```

```
## lm(formula = charges ~ age + bmi + children + smoker + bmi:smoker,
##      data = insurance)
```

```
lm.insurance.all$call # lm with all variables
```

```
## lm(formula = charges ~ age + sex + bmi + children + smoker +
```

```
##      region, data = insurance)
```

10.1.2 Polynomial and GAM

```
lm.insurance.5$call
```

```
## lm(formula = age ~ sex + bmi + children + charges + I(charges^2),  
##      data = insurance)
```

```
gam.insurance.1$call
```

```
## gam(formula = age ~ sex + s(charges) + children, data = insurance)
```

10.1.3 GLM

```
glm.children$call
```

```
## glm(formula = children ~ smoker + charges, family = "poisson",  
##      data = insurance)
```

```
glm.smoker.2$call
```

```
## glm(formula = smoker ~ age + bmi + charges, family = "binomial",  
##      data = insurance)
```

10.1.4 Decision Trees

```
tree.regression.bmi$call #regression tree
```

```
## tree(formula = bmi ~ ., data = insurance)
```

```
tree.regression.charges$call #regression tree
```

```
## tree(formula = charges ~ ., data = insurance)
```

```
tree.regression.bmi.pruned$call #regression tree pruned
```

```
## snip.tree(tree = tree.regression.bmi2, nodes = 9L)
```

```
bag.insurance$call #regression tree bagged
```

```
## bagging.data.frame(formula = charges ~ ., data = insurance, subset = train,  
##      nbagg = 20, coob = TRUE)
```

```
bag.insurance.2$call #regression tree bagged random forest
```

```
## randomForest(formula = charges ~ ., data = insurance, mtry = 6,  
##      importance = TRUE, subset = train)
```

```
rf.insurance$call #random forest
```

```
## randomForest(formula = charges ~ ., data = insurance, mtry = 2,  
##      importance = TRUE, subset = train)
```

10.1.5 Support Vector Machines

```
bestmod$call #svm
```

```
## best.tune(method = svm, train.x = smoker ~ charges + bmi, data = insurance.train,  
##      ranges = list(cost = cost_range), kernel = "linear")
```

```
bestmod.poly$call #svm polynomial
```

```
## best.tune(method = svm, train.x = smoker ~ charges + bmi, data = insurance.train,  
##         ranges = list(cost = cost_range, degree = degree_range),  
##         kernel = "polynomial")
```

10.1.6 Neural Networks

```
nn3$call
```

```
## neuralnet(formula = smoker ~ charges + age + bmi, data = training,  
##         hidden = c(3, 2, 1), threshold = 0.0025, lifesign = "full",  
##         lifesign.step = 500, linear.output = FALSE)
```

10.2 Cross validation

As the models have different dependent variables and are different types (regression/classification) it is difficult to compare all numerically. To do a cross validation, the regression models with charges as dependent variables are compared and the same for classification models with dependent variable smoker.

10.2.1 Regression (dependent variable “charges”)

The in sample performance, using R Squared as measure is the following for the 3 selected models.

```
summary(lm.insurance.3)$r.squared
```

```
## [1] 0.8388
```

```
summary(lm.insurance.all)$r.squared
```

```
## [1] 0.7509
```

```
summary(tree.regression.charges)
```

```
##  
## Regression tree:  
## tree(formula = charges ~ ., data = insurance)  
## Variables actually used in tree construction:  
## [1] "smoker" "age" "children" "bmi"  
## Number of terminal nodes: 7  
## Residual mean deviance: 0.159 = 212 / 1330  
## Distribution of residuals:  
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
## -0.9870 -0.1970 -0.0500  0.0000  0.0994  2.4500
```

As decision trees do not work with an R-squared value per se, we take the residual mean deviance for the decision trees (here at 0.1592). The residual mean deviance (0.1592) can be considered as R-squared error. Interpreting error as 1 - R-squared, we receive roughly similar values for the three variables.

The out of sample performance is computed by using 50:50 training and test data and repeating the process 100 times.

```
set.seed(5)
```

```
RMSE.1 <- c()  
RMSE.2 <- c()  
RMSE.3 <- c()
```

```

RMSE.4 <- c()

for(i in 1:100){

  train.YES <- sample(x = c(TRUE, FALSE),
                     size = nrow(insurance),
                     replace = TRUE)

  insurance.train <- insurance[train.YES, ]

  insurance.test <- insurance[!train.YES, ]

  # fit model with train data

  lm.insurance.3.train <- lm(formula = formula(lm.insurance.3),
                             data = insurance.train)

  lm.insurance.all.train <- lm(formula = formula(lm.insurance.all),
                               data = insurance.train)

  tree.regression.charges.train <- lm(formula = formula(tree.regression.charges),
                                      data = insurance.train)

  rf.insurance.train <- lm(formula = formula(rf.insurance),
                           data = insurance.train)

  # make prediction on test data

  lm.insurance.3.predict <- predict(lm.insurance.3.train,
                                    newdata = insurance.test)

  lm.insurance.all.predict <- predict(lm.insurance.all.train,
                                      newdata = insurance.test)

  tree.regression.charges.predict <- predict(tree.regression.charges.train,
                                             newdata = insurance.test)

  rf.insurance.predict <- predict(rf.insurance.train,
                                  newdata = insurance.test)

  #RMSE

  RMSE.1[i] <- rmse(lm.insurance.3.predict, insurance.test$charges)
  RMSE.2[i] <- rmse(lm.insurance.all.predict, insurance.test$charges)
  RMSE.3[i] <- sqrt((sum(tree.regression.charges.predict - insurance.test[7])^2)/length(tree.regression.charges.predict))

}

```

The out of sample performance, using RMSE as measure is the following:

```
#lm.3
mean(RMSE.1)
```

```
## [1] 4890
```

```
#lm.all
mean(RMSE.2)
```

```
## [1] 6102
```

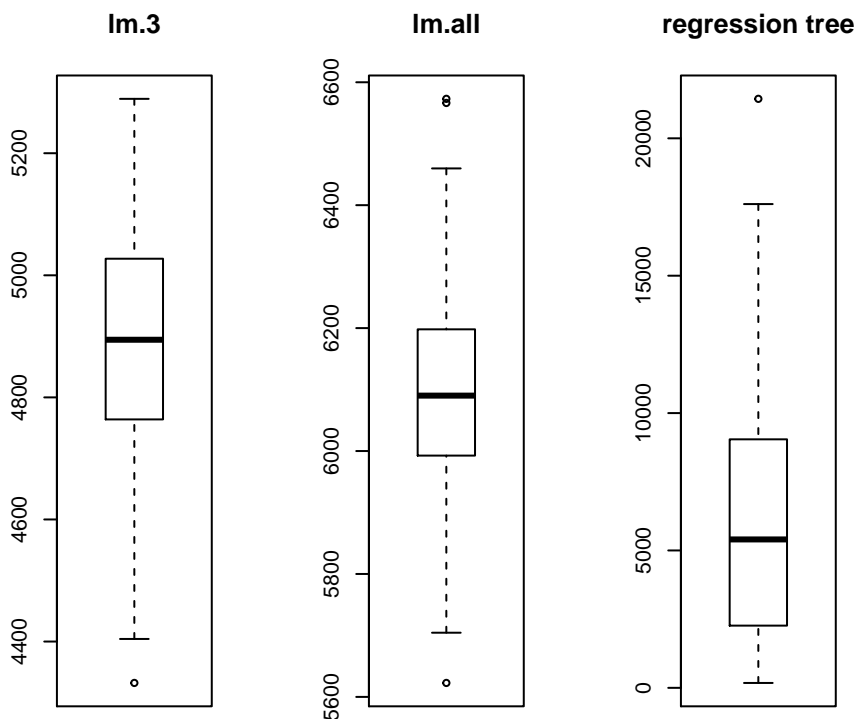
```
#tree regression
mean(RMSE.3)
```

```
## [1] 6129
```

As one can see, the best linear model also performs in cross-validation as the better in comparison with the linear model with all variables as well as the decision tree. This is not very astonishing, as with more variables, the more variance once gets.

What is interesting to see, is that the lm modell with all the variables and the regression tree perform very similarly, albeit not with the best values.

```
par(mfrow=c(1,4))
boxplot(RMSE.1, main="lm.3")
boxplot(RMSE.2, main="lm.all")
boxplot(RMSE.3, main="regression tree")
```



When having a closer look at the boxplots however, one can see the difference in the methods. The linear models perform very regularly and very symmetrically, whereas decision trees have a much larger distribution.

Here, one can also directly see the importance of cross-validation: as one practically plays with the chances when performing a train <-> test approach, one cannot take such results as true. It can be that one ends up on either end of the scale / boxplot.

10.2.2 Classification (Dependent Variable “smoker”)

Throughout the analysis, 3 models have been used to make a classification for smoker/non-smoker. Those models will be cross validated:

```
glm.smoker.2$call #glm
```

```
## glm(formula = smoker ~ age + bmi + charges, family = "binomial",  
##      data = insurance)
```

```
bestmod$call #svm
```

```
## best.tune(method = svm, train.x = smoker ~ charges + bmi, data = insurance.train,  
##           ranges = list(cost = cost_range), kernel = "linear")
```

```
nn3$call #neural network
```

```
## neuralnet(formula = smoker ~ charges + age + bmi, data = training,  
##           hidden = c(3, 2, 1), threshold = 0.0025, lifesign = "full",  
##           lifesign.step = 500, linear.output = FALSE)
```

```
set.seed(5)
```

```
MC.1 <- c()
```

```
MC.2 <- c()
```

```
MC.3 <- c()
```

```
for(i in 1:10){
```

```
  train.YES <- sample(x = c(TRUE, FALSE),  
                      size = nrow(insurance),  
                      replace = TRUE)
```

```
  insurance.train <- insurance[train.YES, ]
```

```
  insurance.test <- insurance[!train.YES, ]
```

```
  insurance.train.normalized <- dplyr::select(insurance.train, smoker, charges, age, bmi)
```

```
  insurance.test.normalized <- dplyr::select(insurance.test, smoker, charges, age, bmi)
```

```
  # Min-Max Normalization for neural networks
```

```
  insurance.train.normalized$charges <- (insurance.train.normalized$charges - min(insurance.train.normalized$charges)) / (max(insurance.train.normalized$charges) - min(insurance.train.normalized$charges))
```

```
  insurance.train.normalized$age <- (insurance.train.normalized$age - min(insurance.train.normalized$age)) / (max(insurance.train.normalized$age) - min(insurance.train.normalized$age))
```

```
  insurance.train.normalized$bmi <- (insurance.train.normalized$bmi - min(insurance.train.normalized$bmi)) / (max(insurance.train.normalized$bmi) - min(insurance.train.normalized$bmi))
```

```
  insurance.test.normalized$charges <- (insurance.test.normalized$charges - min(insurance.test.normalized$charges)) / (max(insurance.test.normalized$charges) - min(insurance.test.normalized$charges))
```

```
  insurance.test.normalized$age <- (insurance.test.normalized$age - min(insurance.test.normalized$age)) / (max(insurance.test.normalized$age) - min(insurance.test.normalized$age))
```

```
  insurance.test.normalized$bmi <- (insurance.test.normalized$bmi - min(insurance.test.normalized$bmi)) / (max(insurance.test.normalized$bmi) - min(insurance.test.normalized$bmi))
```

```

# fit model with train data

mod.1 <- glm(formula = smoker ~ age + bmi + charges, family = "binomial",
             data = insurance.train)

mod.2 <- svm(smoker~charges+bmi, data = insurance.train, cost =5)

nn3 <- neuralnet(smoker~charges + age + bmi,
                data = insurance.train.normalized,
                hidden = c(3,2,1),
                threshold = 0.0025,
                lifesign = 'full',
                lifesign.step = 500,
                linear.output = FALSE)

# make prediction on test data

#mod1
mod.1.pred <- predict(mod.1, newdata = insurance.test, type="response")
mod.1.pred.binomial <- ifelse(mod.1.pred > 0.5, yes = 1, no = 0)
df.mod.1.pred <- data.frame(obs = insurance.test$smoker, fitted = mod.1.pred.binomial)
mc.table.1 <- table(obs = df.mod.1.pred$obs, fit = df.mod.1.pred$fitted)
MC.1[i] <- (mc.table.1[2]+mc.table.1[3])/(mc.table.1[2]+mc.table.1[3]+mc.table.1[1]+mc.table.1[4])

#mod2
mod.2.pred <- predict(mod.2, newdata = insurance.test)
df.mod.2.pred <- data.frame(obs = insurance.test$smoker, fitted = mod.2.pred)
mc.table.2 <- table(obs = df.mod.2.pred$obs, fit = df.mod.2.pred$fitted)
MC.2[i] <- (mc.table.2[2]+mc.table.2[3])/(mc.table.2[2]+mc.table.2[3]+mc.table.2[1]+mc.table.2[4])

#mod3
prediction_nn3 <- predict(nn3, insurance.test.normalized)
class_nn <- apply(prediction_nn3, 1, which.max)
df.mod.3.pred <- data.frame(obs = insurance.test$smoker, fitted = class_nn)
mc.table.3 <- table(obs = df.mod.3.pred$obs, fit = df.mod.3.pred$fitted)
MC.3[i] <- (mc.table.3[2]+mc.table.3[3])/(mc.table.3[2]+mc.table.3[3]+mc.table.3[1]+mc.table.3[4])
}

```

The following are the misclassification rates of our models:

```
#misclassification rates of the models
```

```
#glm
```

```
mean(MC.1)
```

```
## [1] 0.0462371
```

```
#svm
```

```
mean(MC.2)
```

```
## [1] 0.042029
```

```
#neuralnet
```

```
mean(MC.3)
```

```
## [1] 0.05506842
```

We can observe that all models have a misclassification rate of around 4-5.5%. The model that performs best is the svm. The model that performs worst, but still reasonably well is the neural network.

10.2.3 Cross-Validation for Classification Trees

We recall the classification errors in the above code:

```
(tree.classification.region.pred.error.1)
```

```
## [1] 0.642094
```

```
(tree.classification.region.pred.error.2)
```

```
## [1] 0.8781095
```

We now recheck if misclassification improves when cross-validating by a 10fold crossvalidations our model:

```
n <- nrow(insurance)
K <- 10
size <- n/10
set.seed(5)
alea <- runif(n)
rang <- rank(alea)
bloc <- (rang-1)/size+1
bloc <- as.factor(bloc)
print(summary(bloc)) #Anzahl Datensätze und Daten pro Dign
```

```
##      1      2      3      4      5      6      7      8      9     10     11
## 133 133 133 133 133 133 133 133 133 133 8
```

```
library(rpart)
pred <- predict(tree.classification.region_all, newdata = insurance, type = "class")
mc <- table(insurance$region, pred) #missclassification
print(mc)
```

```
##      pred
##      0      1      2      3
## 0 211      0      0 153
## 1 112      0      0 213
## 2 105      0      0 219
## 3  74      0      0 251
```

```
all.err <- numeric(0)
for (k in 1:K) {
  baum <- rpart(region ~., data = insurance[bloc!=k,], method="class")
  pred <- predict(baum, newdata=insurance[bloc==k,], type = "class")
  mc <- table(insurance$region[bloc==k], pred)
  err <- 1.0 -(mc[1,1]+mc[2,2]+mc[3,3]+mc[4,4])/sum(mc)
  #print(err)
  all.err <- rbind(all.err,err)
}

#print(all.err)

err.cv <- mean(all.err)
print(err.cv)
```



```
## [1] 0.6699248
```

We see here that with 66.99%, our misclassification error has not much improved with crossvalidation. In fact, it even has gotten a bit worse in comparison with `tree.classification.region.pred.error.1`.

11 Final words

In the analysed data set, we got a lot of different insights on various sides. We learned that not every variable might be of importance for certain forecasts - or at least that not enough data is available at the time in order to do so. We think here of the variable “region” which performed not very good in terms of prediction for most of the models.

Then there is the other side of performance: models that can predict at a very good rate if somebody is a smoker or not, with an astonishingly low error margin of 5%.

All in all, it was very interesting to see the theory performing in practice and to get to know what kind of insights one receives by these machine learning technique hands on.

And “most” importantly, what we learned in this project and which is clearly documented: smoking is not only bad for the health, it is also bad for ones wallet. Never smoke, kids!