# Exercise 1

## Nikolaus Czernin

```r
# install.packages("ISLR")
library("ISLR")
library("tidyverse")
```

```
## -- Attaching packages ------------------------------------- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.6      v purrr   0.3.4
## v tibble  3.1.7      v dplyr   1.0.9
## v tidyr   1.2.0      v stringr 1.4.0
## v readr   2.1.2      v forcats 0.5.1
```

```
## -- Conflicts ---------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```r
library("knitr")

data(College,package="ISLR")
saveRDS(College, file = "College.rds")
```

```r
?College
str(College)
```

```
## 'data.frame':    777 obs. of  18 variables:
##  $ Private    : Factor w/ 2 levels "No","Yes": 2 2 2 2 2 2 2 2 2 2 ...
##  $ Apps       : num  1660 2186 1428 417 193 ...
##  $ Accept     : num  1232 1924 1097 349 146 ...
##  $ Enroll     : num  721 512 336 137 55 158 103 489 227 172 ...
##  $ Top10perc  : num  23 16 22 60 16 38 17 37 30 21 ...
##  $ Top25perc  : num  52 29 50 89 44 62 45 68 63 44 ...
##  $ F.Undergrad: num  2885 2683 1036 510 249 ...
##  $ P.Undergrad: num  537 1227 99 63 869 ...
##  $ Outstate   : num  7440 12280 11250 12960 7560 ...
##  $ Room.Board : num  3300 6450 3750 5450 4120 ...
##  $ Books      : num  450 750 400 450 800 500 500 450 300 660 ...
##  $ Personal   : num  2200 1500 1165 875 1500 ...
##  $ PhD        : num  70 29 53 92 76 67 90 89 79 40 ...
##  $ Terminal   : num  78 30 66 97 72 73 93 100 84 41 ...
##  $ S.F.Ratio  : num  18.1 12.2 12.9 7.7 11.9 9.4 11.5 13.7 11.3 11.5 ...
##  $ perc.alumni: num  12 16 30 37 2 11 26 37 23 15 ...
##  $ Expend     : num  7041 10527 8735 19016 10922 ...
##  $ Grad.Rate  : num  60 56 54 59 15 55 63 73 80 52 ...
```

```r
summary(College)
```

```
##   Private        Apps           Accept          Enroll        Top10perc
##   No :212   Min.   :   81   Min.   :   72   Min.   :  35   Min.   : 1.00
##   Yes:565   1st Qu.:  776   1st Qu.:  604   1st Qu.: 242   1st Qu.:15.00
##             Median : 1558   Median : 1110   Median : 434   Median :23.00
##             Mean   : 3002   Mean   : 2019   Mean   : 780   Mean   :27.56
##             3rd Qu.: 3624   3rd Qu.: 2424   3rd Qu.: 902   3rd Qu.:35.00
##             Max.   :48094   Max.   :26330   Max.   :6392   Max.   :96.00
##    Top25perc      F.Undergrad     P.Undergrad        Outstate
##   Min.   :  9.0   Min.   :  139   Min.   :    1.0   Min.   : 2340
##   1st Qu.: 41.0   1st Qu.:  992   1st Qu.:   95.0   1st Qu.: 7320
##   Median : 54.0   Median : 1707   Median :  353.0   Median : 9990
##   Mean   : 55.8   Mean   : 3700   Mean   :  855.3   Mean   :10441
##   3rd Qu.: 69.0   3rd Qu.: 4005   3rd Qu.:  967.0   3rd Qu.:12925
##   Max.   :100.0   Max.   :31643   Max.   :21836.0   Max.   :21700
##    Room.Board      Books          Personal          PhD
##   Min.   :1780   Min.   :  96.0   Min.   : 250   Min.   :  8.00
##   1st Qu.:3597   1st Qu.: 470.0   1st Qu.: 850   1st Qu.: 62.00
##   Median :4200   Median : 500.0   Median :1200   Median : 75.00
##   Mean   :4358   Mean   : 549.4   Mean   :1341   Mean   : 72.66
##   3rd Qu.:5050   3rd Qu.: 600.0   3rd Qu.:1700   3rd Qu.: 85.00
##   Max.   :8124   Max.   :2340.0   Max.   :6800   Max.   :103.00
##    Terminal       S.F.Ratio       perc.alumni        Expend
##   Min.   : 24.0   Min.   : 2.50   Min.   : 0.00   Min.   : 3186
##   1st Qu.: 71.0   1st Qu.:11.50   1st Qu.:13.00   1st Qu.: 6751
##   Median : 82.0   Median :13.60   Median :21.00   Median : 8377
##   Mean   : 79.7   Mean   :14.09   Mean   :22.74   Mean   : 9660
##   3rd Qu.: 92.0   3rd Qu.:16.50   3rd Qu.:31.00   3rd Qu.:10830
##   Max.   :100.0   Max.   :39.80   Max.   :64.00   Max.   :56233
##    Grad.Rate
##   Min.   : 10.00
##   1st Qu.: 53.00
##   Median : 65.00
##   Mean   : 65.46
##   3rd Qu.: 78.00
##   Max.   :118.00
```

Regarding `Apps`, the mean is twice as high as the median, suggesting the data is left-skewed. By log-transforming the data, we reduce the effect of high numbers and hopefully make the model more robust.

```r
College_processed <- College %>%
  mutate(Apps = log(Apps)) %>%
  select(-Accept, -Enroll)

n <- nrow(College_processed)
idx <- sample(1:n, n%/%3*2)
train = College_processed[idx,]
test =  College_processed[-idx,]
```
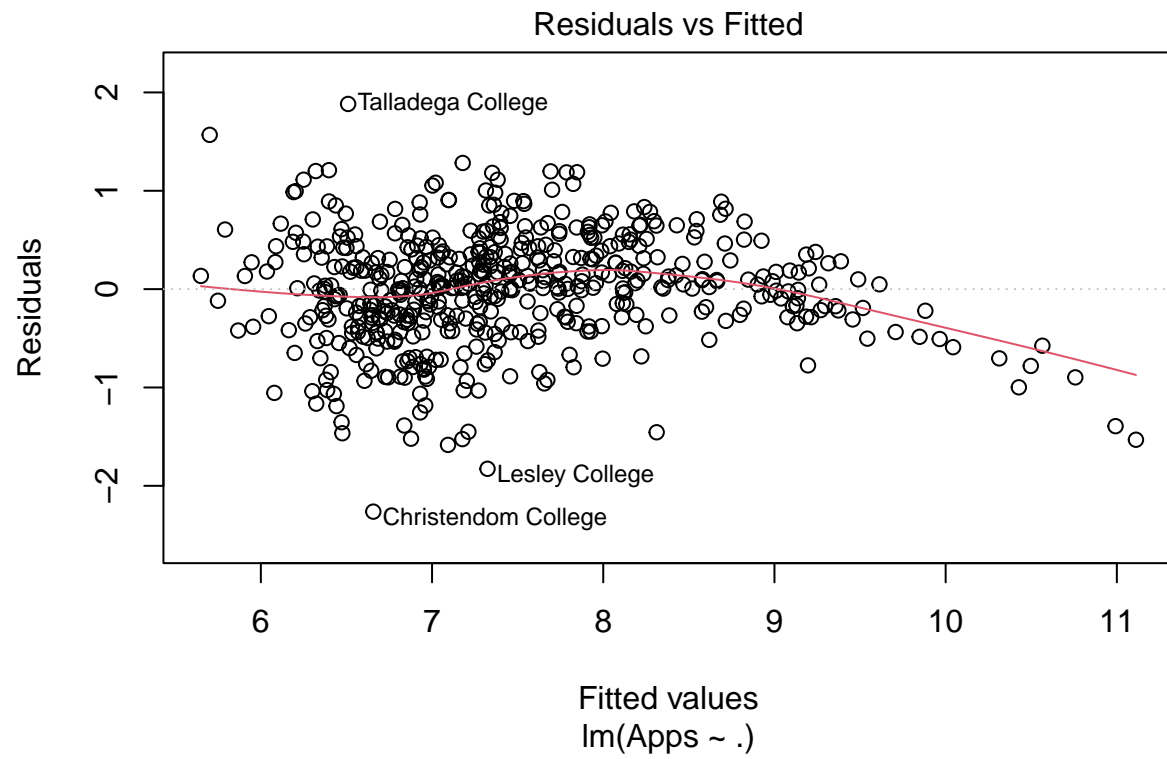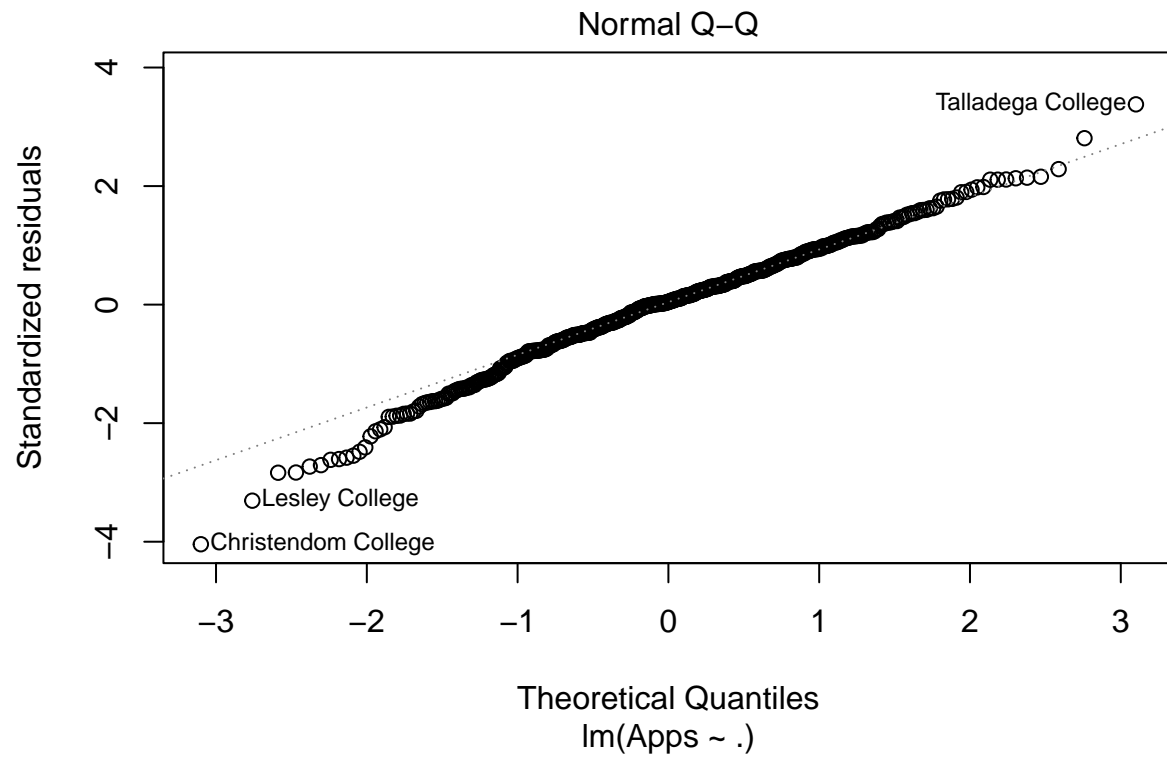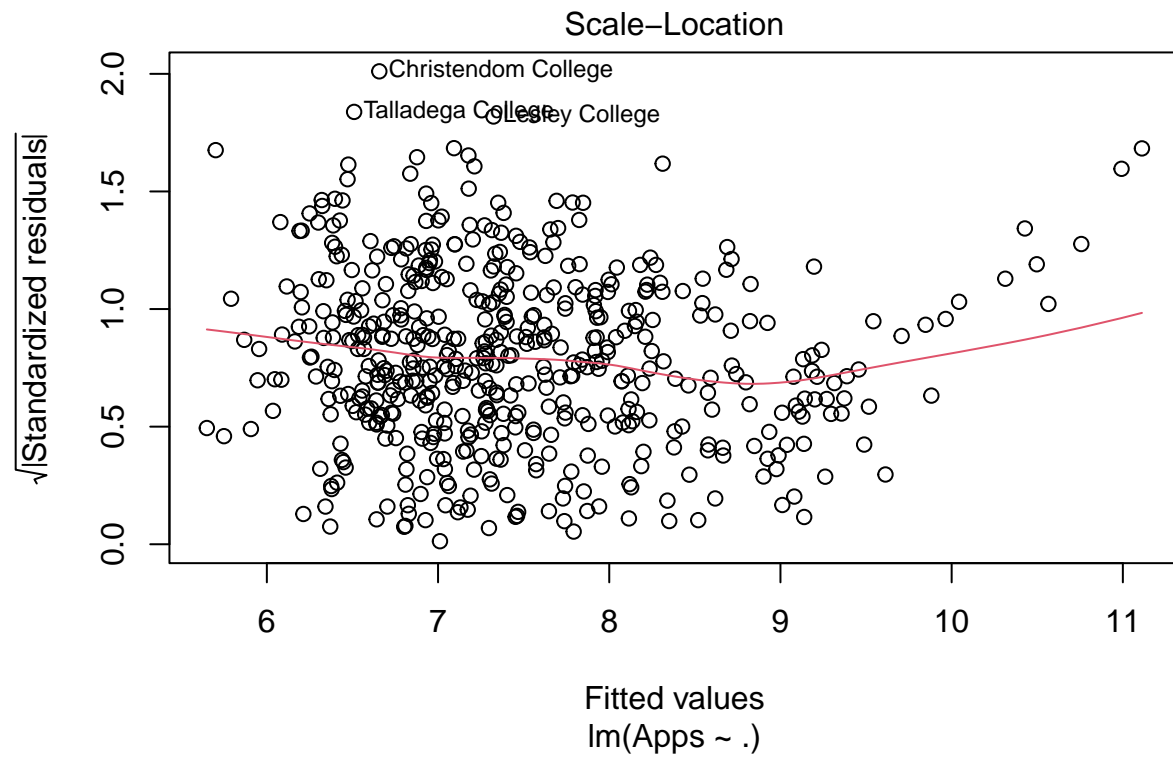
## 2. Full Model

```
full.lm <- lm(Apps ~ ., data = train)
full.lm %>% summary()
```
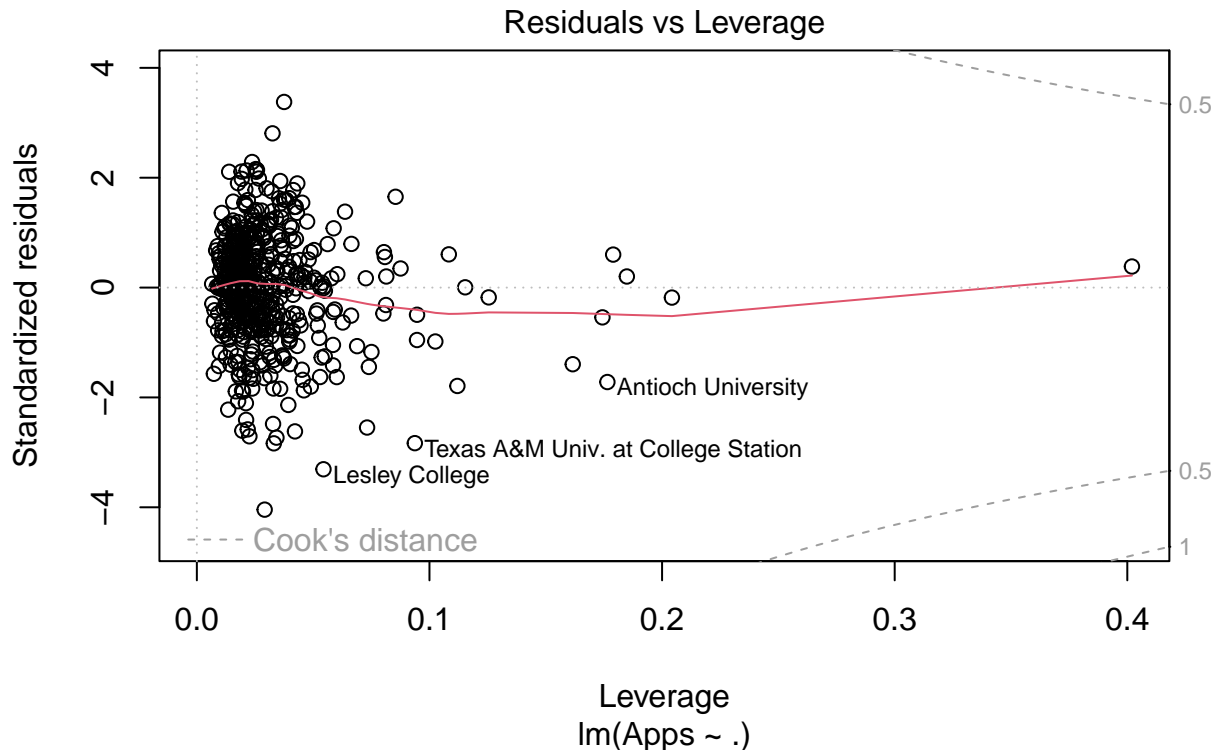
```
##
## Call:
## lm(formula = Apps ~ ., data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.26230 -0.31077  0.02488  0.35448  1.88278
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.633e+00  2.575e-01  17.990  < 2e-16 ***
## PrivateYes  -7.296e-01  9.059e-02  -8.054 5.86e-15 ***
## Top10perc    1.390e-03  3.768e-03   0.369 0.712395
## Top25perc    5.105e-03  2.916e-03   1.751 0.080628 .
## F.Undergrad  1.073e-04  7.662e-06  14.010  < 2e-16 ***
## P.Undergrad  3.183e-06  1.919e-05   0.166 0.868364
## Outstate     5.164e-05  1.288e-05   4.008 7.06e-05 ***
## Room.Board   7.072e-05  3.272e-05   2.161 0.031144 *
## Books        2.777e-04  1.639e-04   1.694 0.090942 .
## Personal     1.071e-05  4.312e-05   0.248 0.803913
## PhD          2.685e-03  3.102e-03   0.866 0.387177
## Terminal     2.560e-03  3.300e-03   0.776 0.438317
## S.F.Ratio    3.259e-02  8.304e-03   3.924 9.91e-05 ***
## perc.alumni -9.705e-03  2.682e-03  -3.618 0.000327 ***
## Expend       2.665e-05  7.405e-06   3.598 0.000352 ***
## Grad.Rate    1.024e-02  2.004e-03   5.109 4.61e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5679 on 502 degrees of freedom
## Multiple R-squared:  0.7362, Adjusted R-squared:  0.7283
## F-statistic: 93.37 on 15 and 502 DF,  p-value: < 2.2e-16
```

```
full.lm %>% plot()
```

Residuals vs Fitted

Talladega College

Lesley College

Christendom College

Residuals

Fitted values
lm(Apps ~ .)

4

Normal Q–Q

Standardized residuals

Theoretical Quantiles
lm(Apps ~ .)

Scale–Location

√|Standardized residuals|

Fitted values
lm(Apps ~ .)

## Residuals vs Leverage



The Residuals vs Fitted plot shows how poorly the model is performing. The red line not being straight is a sign, that the variance is not constant, which is also observableby the residual points being more spread out on the left hand side than the right end.

The QQ-Plot shows that the residuals are somewhat normally distributed.

## Manual computation of the coefficients

```r
X <- model.matrix(Apps ~ ., data = train)
# get the manual estimator
full.estimator <- solve(t(X) %*% X) %*% (t(X) %*% train$Apps)

# bind it to the coefficients of the lm function
summary(full.lm) %>% .$coefficients %>% .[,1] %>% cbind(full.estimator)
```

```
##                             .
## (Intercept)   4.632570e+00   4.632570e+00
## PrivateYes   -7.296089e-01  -7.296089e-01
## Top10perc     1.389893e-03   1.389893e-03
## Top25perc     5.104754e-03   5.104754e-03
## F.Undergrad   1.073479e-04   1.073479e-04
## P.Undergrad   3.182930e-06   3.182930e-06
## Outstate      5.163685e-05   5.163685e-05
## Room.Board    7.071469e-05   7.071469e-05
## Books         2.776634e-04   2.776634e-04
```

```
## Personal      1.071283e-05   1.071283e-05
## PhD           2.684861e-03   2.684861e-03
## Terminal      2.559598e-03   2.559598e-03
## S.F.Ratio     3.258802e-02   3.258802e-02
## perc.alumni  -9.704870e-03  -9.704870e-03
## Expend        2.664680e-05   2.664680e-05
## Grad.Rate     1.024037e-02   1.024037e-02
```
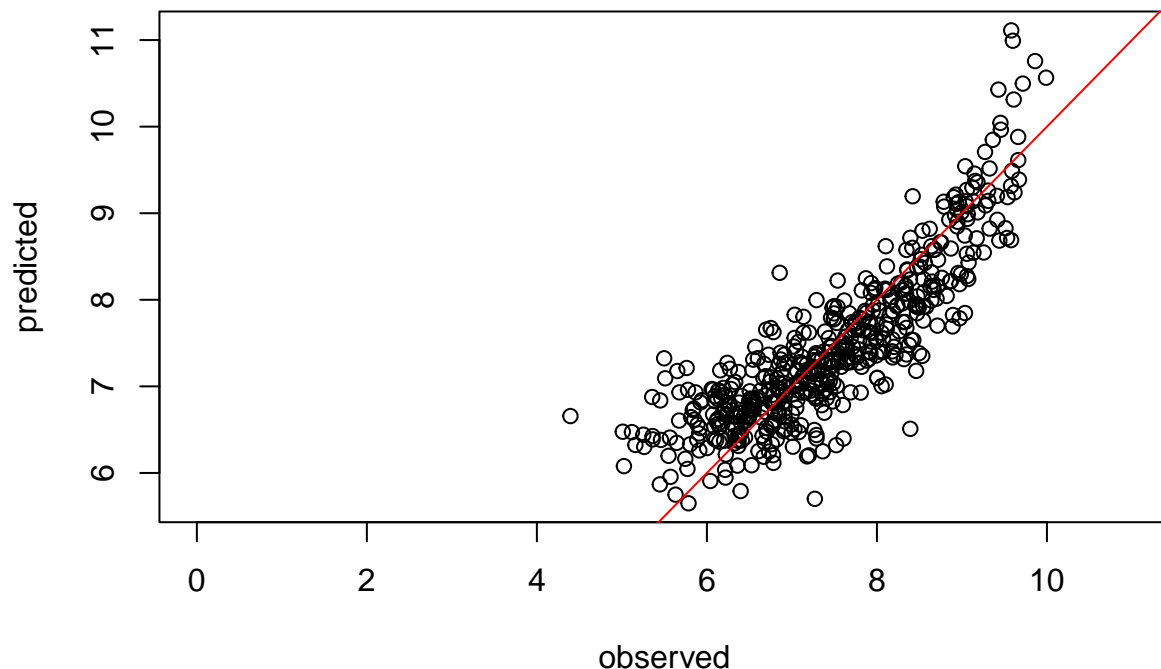
The coefficients of the lm() function and the manual estimation are equal.
`PrivateYes` is a variable with highly significant coefficient of ~-0.5, meaning that a value of "Yes" negatively influences the response.

**Predicting values**

```
plot(train$Apps, full.lm %>% predict(train) , xlim = c(0, 11),
     main="Observed vs predicted values (training data)", xlab="observed", ylab="predicted",
     # ylim=c(0, 11)
     )
abline(coef = c(0,1), col="red")
```
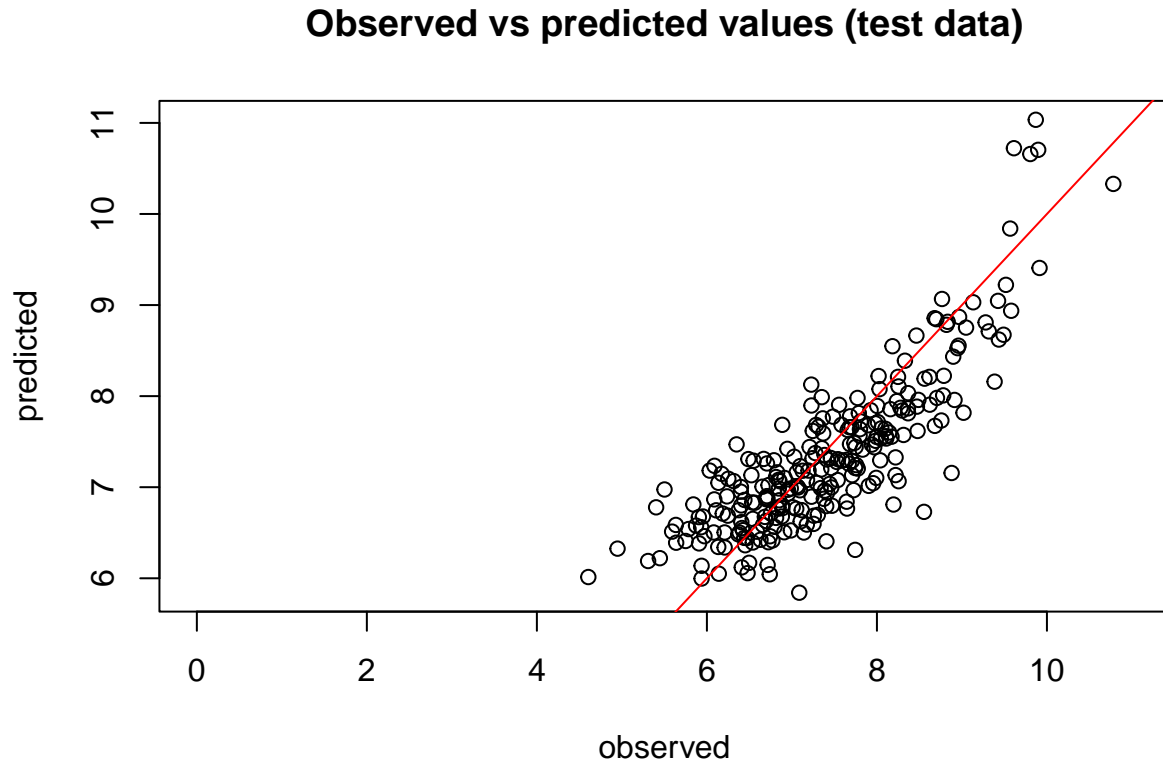
## Observed vs predicted values (training data)



```
plot(test$Apps, full.lm %>% predict(test) , xlim = c(0, 11),
     main="Observed vs predicted values (test data)", xlab="observed", ylab="predicted",
     # ylim=c(0, 11)
```

```
  )
abline(coef = c(0,1), col="red")
```

## Observed vs predicted values (test data)



observed

Visually, the variance of the predicted vs observed data points look similar in the plots of the training and the test data.

```
get_rmse <- function(y, yhat){
  N <- length(y)
  y_minus_yhat <- (y - yhat)^2
  avg_sum <- sum(y_minus_yhat)/N
  sqrt(avg_sum)
}

paste(
  "RMSE of training set:",
  get_rmse(train$Apps, full.lm %>% predict(train)) %>% round(4),
  " ---- RMSE of test set:",
  get_rmse(test$Apps, full.lm %>% predict(test)) %>% round(4)
)
```

```
## [1] "RMSE of training set: 0.5591  ---- RMSE of test set: 0.5784"
```

The RMSE of the training set being lower than that of the test set alsp checks out.

# 3. Slim model

Manually removing all insignificant variables from the full model, we are left with:
- An Intercept that is not zero
- `Private`
- `F.Undergrad`
- `Outstate`
- `Room.Board`
- `Books`
- `PhD`
- `S.F.Ratio`
- `perc.alumni`
- `Expend`
- `Grad.Rate`

```
slim.lm <- lm(Apps ~ Private + F.Undergrad + Outstate + Room.Board + Books + PhD + S.F.Ratio + perc.alu
slim.lm %>% summary()
```
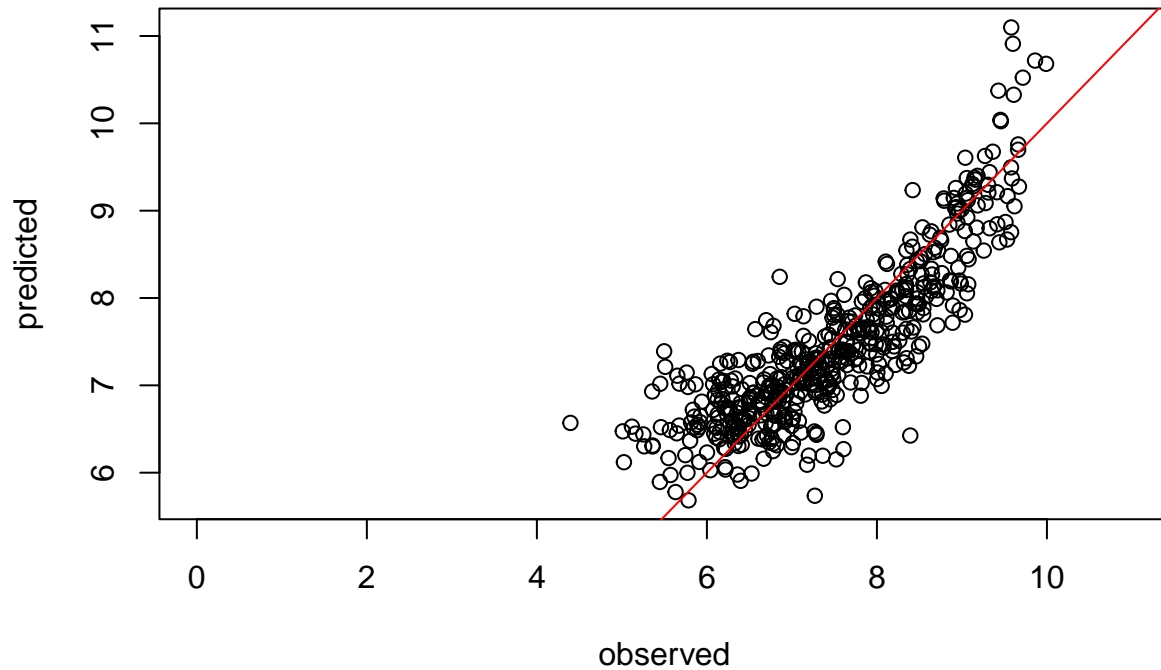
```
##
## Call:
## lm(formula = Apps ~ Private + F.Undergrad + Outstate + Room.Board +
##     Books + PhD + S.F.Ratio + perc.alumni + Expend + Grad.Rate,
##     data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.17549 -0.31134  0.01936  0.36590  1.96717
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.655e+00  2.294e-01  20.292  < 2e-16 ***
## PrivateYes  -7.297e-01  9.046e-02  -8.066 5.27e-15 ***
## F.Undergrad  1.133e-04  6.844e-06  16.550  < 2e-16 ***
## Outstate     5.757e-05  1.277e-05   4.510 8.05e-06 ***
## Room.Board   6.160e-05  3.201e-05   1.925  0.05482 .
## Books        3.806e-04  1.575e-04   2.416  0.01605 *
## PhD          6.396e-03  2.099e-03   3.047  0.00244 **
## S.F.Ratio    3.164e-02  8.319e-03   3.803  0.00016 ***
## perc.alumni -8.047e-03  2.639e-03  -3.049  0.00242 **
## Expend       3.065e-05  6.856e-06   4.471 9.62e-06 ***
## Grad.Rate    1.159e-02  1.911e-03   6.064 2.60e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5718 on 507 degrees of freedom
## Multiple R-squared:  0.7299, Adjusted R-squared:  0.7246
## F-statistic:   137 on 10 and 507 DF,  p-value: < 2.2e-16
```

After pruning the variables that were not significant in the full model, all remaining variables' coefficients are significant in the pruned model.
Generally, this is not always the case, as highly correlated variables that are significant may not be significant anymore if you remove one.
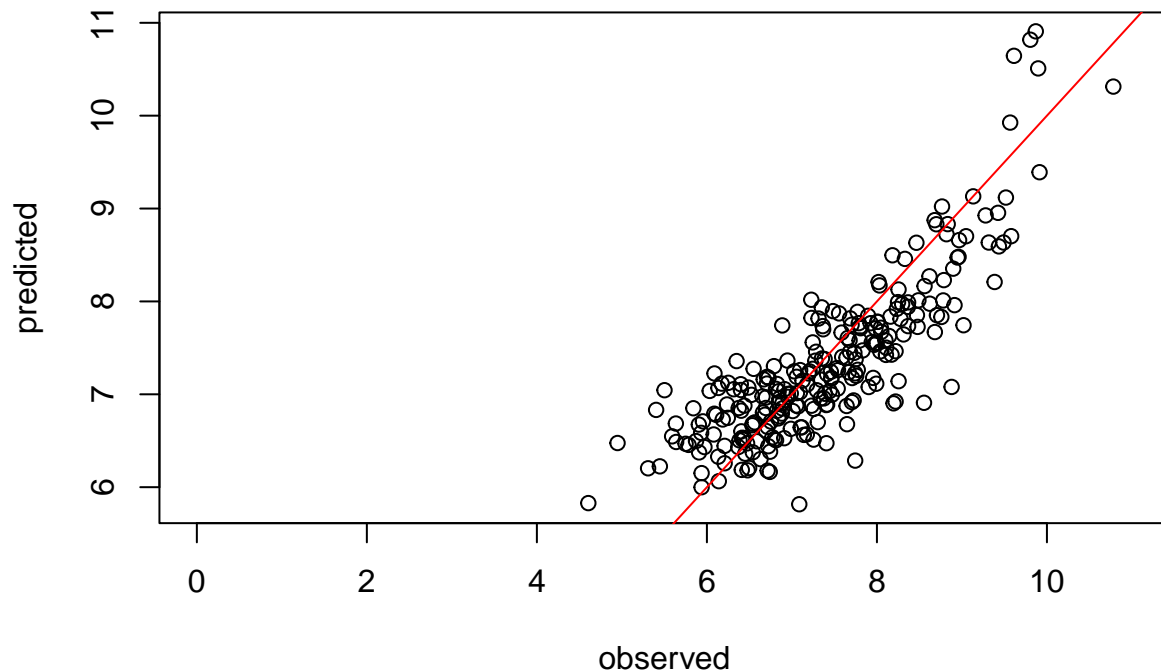
```r
plot(train$Apps, slim.lm %>% predict(train) , xlim = c(0, 11),
     main="Observed vs predicted values (training data)", xlab="observed", ylab="predicted",
     # ylim=c(0, 11)
     )
abline(coef = c(0,1), col="red")
```

## Observed vs predicted values (training data)



```r
plot(test$Apps, slim.lm %>% predict(test) , xlim = c(0, 11),
     main="Observed vs predicted values (test data)", xlab="observed", ylab="predicted",
     # ylim=c(0, 11)
     )
abline(coef = c(0,1), col="red")
```

## Observed vs predicted values (test data)



Visually, I don't see an obvious improvement of the slim model's performance to the full model.

```
paste(
  "RMSE of training set:",
  get_rmse(train$Apps, slim.lm %>% predict(train)) %>% round(4),
  " ---- RMSE of test set:",
  get_rmse(test$Apps, slim.lm %>% predict(test)) %>% round(4)
)
```

```
## [1] "RMSE of training set: 0.5657  ---- RMSE of test set: 0.5728"
```

On the test set, the RMSE has gotten marginally worse, which is to be expected when pruning predictors. On the other hand, also unsurprisingly, the RMSE on the test set has improved, though also by a bizmal amount.

```
anova(full.lm, slim.lm)
```

```
## Analysis of Variance Table
##
## Model 1: Apps ~ Private + Top10perc + Top25perc + F.Undergrad + P.Undergrad +
##     Outstate + Room.Board + Books + Personal + PhD + Terminal +
##     S.F.Ratio + perc.alumni + Expend + Grad.Rate
## Model 2: Apps ~ Private + F.Undergrad + Outstate + Room.Board + Books +
##     PhD + S.F.Ratio + perc.alumni + Expend + Grad.Rate
##   Res.Df    RSS Df Sum of Sq      F Pr(>F)
## 1    502 161.92
```

```
## 2     507 165.76 -5   -3.8427 2.3827 0.0375 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The p-value is >6%, so I would conservatively rule this slim model not to be significantly different from the full model.

# 4. Stepwise variable selection

```
step.fw.lm <- step(full.lm, direction = "forward")
```

```
## Start:  AIC=-570.36
## Apps ~ Private + Top10perc + Top25perc + F.Undergrad + P.Undergrad +
##     Outstate + Room.Board + Books + Personal + PhD + Terminal +
##     S.F.Ratio + perc.alumni + Expend + Grad.Rate
```

```
step.bw.lm <- step(full.lm, direction = "backward")
```

```
## Start:  AIC=-570.36
## Apps ~ Private + Top10perc + Top25perc + F.Undergrad + P.Undergrad +
##     Outstate + Room.Board + Books + Personal + PhD + Terminal +
##     S.F.Ratio + perc.alumni + Expend + Grad.Rate
##
##                 Df Sum of Sq    RSS     AIC
## - P.Undergrad  1     0.009 161.93 -572.33
## - Personal     1     0.020 161.94 -572.30
## - Top10perc    1     0.044 161.97 -572.22
## - Terminal     1     0.194 162.12 -571.74
## - PhD          1     0.242 162.16 -571.59
## <none>                     161.92 -570.36
## - Books        1     0.925 162.85 -569.41
## - Top25perc    1     0.988 162.91 -569.21
## - Room.Board   1     1.507 163.43 -567.56
## - Expend       1     4.177 166.10 -559.17
## - perc.alumni  1     4.222 166.15 -559.03
## - S.F.Ratio    1     4.968 166.89 -556.71
## - Outstate     1     5.181 167.10 -556.05
## - Grad.Rate    1     8.419 170.34 -546.10
## - Private      1    20.924 182.85 -509.41
## - F.Undergrad  1    63.311 225.23 -401.41
##
## Step:  AIC=-572.33
## Apps ~ Private + Top10perc + Top25perc + F.Undergrad + Outstate +
##     Room.Board + Books + Personal + PhD + Terminal + S.F.Ratio +
##     perc.alumni + Expend + Grad.Rate
##
##                 Df Sum of Sq    RSS     AIC
## - Personal     1     0.025 161.96 -574.25
## - Top10perc    1     0.041 161.97 -574.20
## - Terminal     1     0.195 162.13 -573.71
```

```
## - PhD           1     0.247 162.18 -573.54
## <none>                        161.93 -572.33
## - Books          1     0.927 162.86 -571.37
## - Top25perc      1     0.991 162.92 -571.17
## - Room.Board     1     1.535 163.47 -569.44
## - Expend         1     4.195 166.13 -561.08
## - perc.alumni    1     4.224 166.16 -560.99
## - S.F.Ratio      1     4.962 166.89 -558.70
## - Outstate       1     5.173 167.10 -558.04
## - Grad.Rate      1     8.503 170.44 -547.82
## - Private        1    21.011 182.94 -511.14
## - F.Undergrad    1    74.992 236.92 -377.20
##
## Step:  AIC=-574.25
## Apps ~ Private + Top10perc + Top25perc + F.Undergrad + Outstate +
##     Room.Board + Books + PhD + Terminal + S.F.Ratio + perc.alumni +
##     Expend + Grad.Rate
##
##              Df Sum of Sq    RSS     AIC
## - Top10perc   1     0.041 162.00 -576.12
## - Terminal    1     0.188 162.14 -575.65
## - PhD         1     0.255 162.21 -575.44
## <none>                    161.96 -574.25
## - Top25perc   1     0.986 162.94 -573.11
## - Books       1     1.067 163.02 -572.85
## - Room.Board  1     1.519 163.47 -571.42
## - Expend      1     4.235 166.19 -562.88
## - perc.alumni 1     4.314 166.27 -562.63
## - S.F.Ratio   1     4.937 166.89 -560.70
## - Outstate    1     5.148 167.10 -560.04
## - Grad.Rate   1     8.556 170.51 -549.58
## - Private     1    21.127 183.08 -512.74
## - F.Undergrad 1    78.015 239.97 -372.58
##
## Step:  AIC=-576.12
## Apps ~ Private + Top25perc + F.Undergrad + Outstate + Room.Board +
##     Books + PhD + Terminal + S.F.Ratio + perc.alumni + Expend +
##     Grad.Rate
##
##              Df Sum of Sq    RSS     AIC
## - Terminal    1     0.171 162.17 -577.57
## - PhD         1     0.281 162.28 -577.22
## <none>                    162.00 -576.12
## - Books       1     1.087 163.09 -574.66
## - Room.Board  1     1.479 163.48 -573.41
## - Top25perc   1     3.533 165.53 -566.94
## - perc.alumni 1     4.282 166.28 -564.60
## - S.F.Ratio   1     4.905 166.90 -562.67
## - Outstate    1     5.258 167.26 -561.57
## - Expend      1     5.284 167.28 -561.49
## - Grad.Rate   1     8.758 170.76 -550.85
## - Private     1    21.112 183.11 -514.66
## - F.Undergrad 1    78.254 240.25 -373.97
##
```
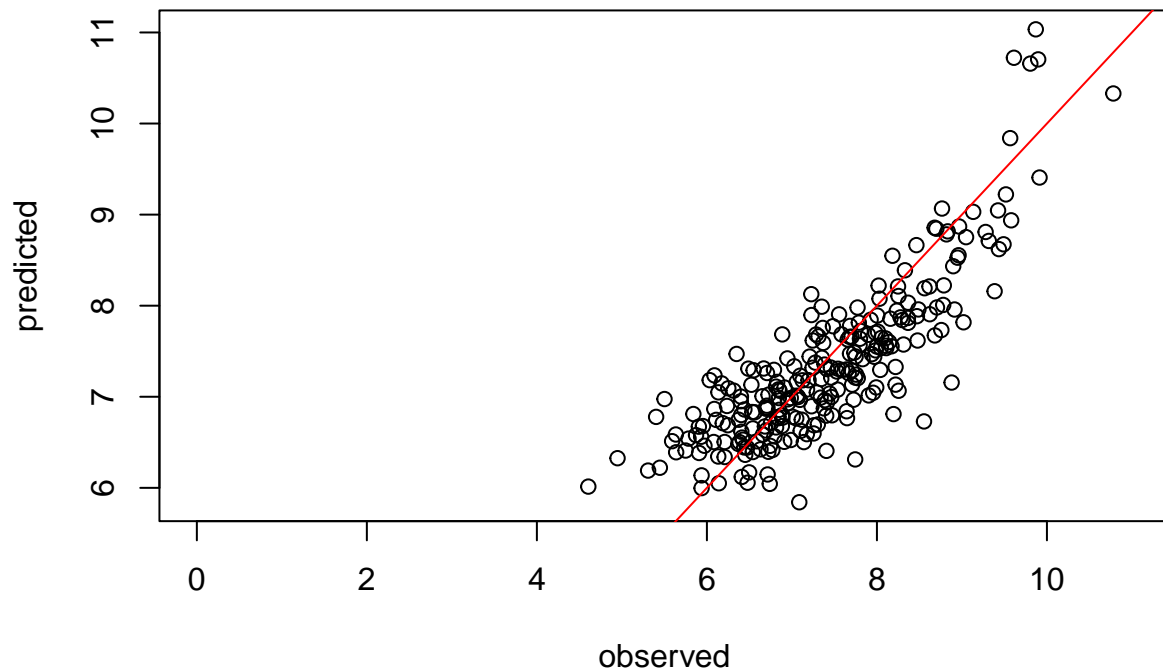
```
## Step:  AIC=-577.57
## Apps ~ Private + Top25perc + F.Undergrad + Outstate + Room.Board +
##     Books + PhD + S.F.Ratio + perc.alumni + Expend + Grad.Rate
##
##               Df Sum of Sq    RSS     AIC
## <none>                     162.17 -577.57
## - Books        1     1.159 163.33 -575.89
## - PhD          1     1.369 163.54 -575.22
## - Room.Board   1     1.659 163.83 -574.30
## - Top25perc    1     3.597 165.76 -568.21
## - perc.alumni  1     4.181 166.35 -566.39
## - S.F.Ratio    1     4.983 167.15 -563.90
## - Expend       1     5.329 167.50 -562.82
## - Outstate     1     5.446 167.61 -562.46
## - Grad.Rate    1     8.637 170.81 -552.70
## - Private      1    21.747 183.92 -514.39
## - F.Undergrad  1    78.712 240.88 -374.62
```

```r
data.frame(
  model = c("Full", "Slim", "Step.Forward", "Step.Backward"),
  rmse_train = c(
    get_rmse(train$Apps, predict(full.lm, train)),
    get_rmse(train$Apps, predict(slim.lm, train)),
    get_rmse(train$Apps, predict(step.fw.lm, train)),
    get_rmse(train$Apps, predict(step.bw.lm, train))
  ),
  rmse_test = c(
    get_rmse(test$Apps, predict(full.lm, test)),
    get_rmse(test$Apps, predict(slim.lm, test)),
    get_rmse(test$Apps, predict(step.fw.lm, test)),
    get_rmse(test$Apps, predict(step.bw.lm, test))
  )
) %>% kable()
```

| model | rmse_train | rmse_test |
|---|---|---|
| Full | 0.5590994 | 0.5783854 |
| Slim | 0.5656947 | 0.5727900 |
| Step.Forward | 0.5590994 | 0.5783854 |
| Step.Backward | 0.5595234 | 0.5775654 |

```r
plot(test$Apps, step.fw.lm %>% predict(test) , xlim = c(0, 11),
     main="Observed vs predicted values (forward stepwise model)", xlab="observed", ylab="predicted",
     # ylim=c(0, 11)
     )
abline(coef = c(0,1), col="red")
```
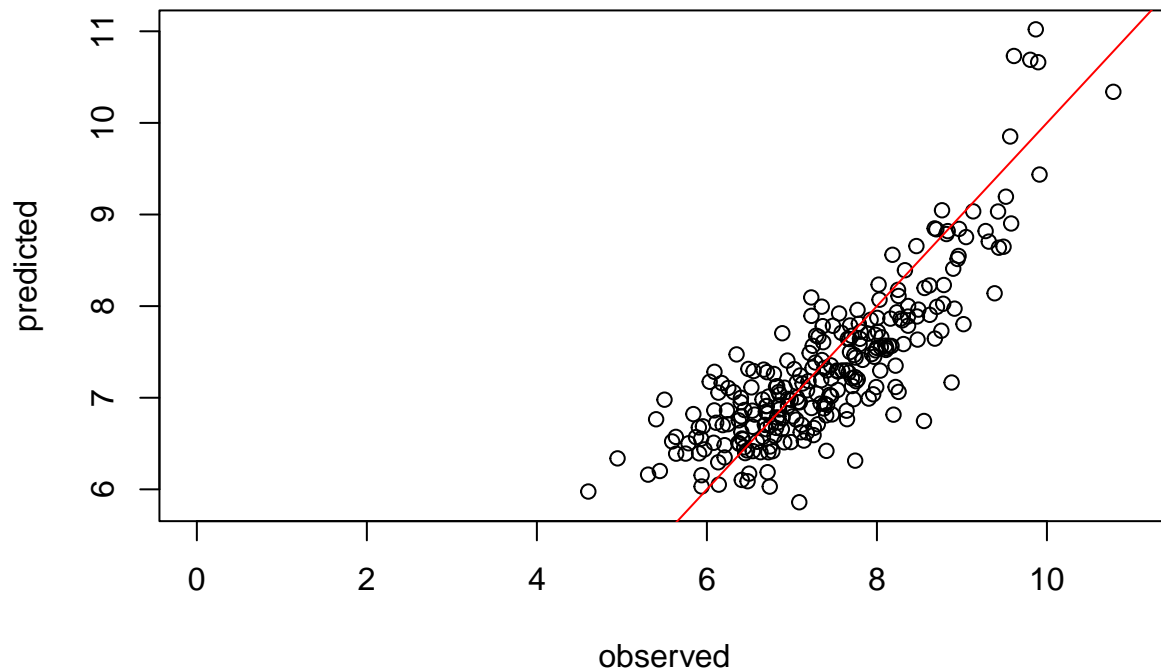
## Observed vs predicted values (forward stepwise model)



```r
plot(test$Apps, step.bw.lm %>% predict(test) , xlim = c(0, 11),
     main="Observed vs predicted values (backward stepwise model)", xlab="observed", ylab="predicted",
     # ylim=c(0, 11)
     )
abline(coef = c(0,1), col="red")
```

## Observed vs predicted values (backward stepwise model)



```
anova(full.lm, step.fw.lm)
```

```
## Analysis of Variance Table
##
## Model 1: Apps ~ Private + Top10perc + Top25perc + F.Undergrad + P.Undergrad +
##     Outstate + Room.Board + Books + Personal + PhD + Terminal +
##     S.F.Ratio + perc.alumni + Expend + Grad.Rate
## Model 2: Apps ~ Private + Top10perc + Top25perc + F.Undergrad + P.Undergrad +
##     Outstate + Room.Board + Books + Personal + PhD + Terminal +
##     S.F.Ratio + perc.alumni + Expend + Grad.Rate
##   Res.Df    RSS Df Sum of Sq F Pr(>F)
## 1    502 161.92
## 2    502 161.92  0         0
```

```
anova(full.lm, step.bw.lm)
```

```
## Analysis of Variance Table
##
## Model 1: Apps ~ Private + Top10perc + Top25perc + F.Undergrad + P.Undergrad +
##     Outstate + Room.Board + Books + Personal + PhD + Terminal +
##     S.F.Ratio + perc.alumni + Expend + Grad.Rate
## Model 2: Apps ~ Private + Top25perc + F.Undergrad + Outstate + Room.Board +
##     Books + PhD + S.F.Ratio + perc.alumni + Expend + Grad.Rate
##   Res.Df    RSS Df Sum of Sq      F Pr(>F)
```

```
## 1    502 161.92
## 2    506 162.17 -4  -0.24566 0.1904 0.9434
```

From looking at the resulting RMSE scores, the observed vs predicted plots and the results of the ANOVA tests, the stepwise models did not make a mentionable difference.