

ST309 week 6

2023-10-30

```
library(MASS)
names(Boston) # show the column names

## [1] "crim"    "zn"      "indus"   "chas"    "nox"     "rm"      "age"
## [8] "dis"     "rad"     "tax"     "ptratio" "black"   "lstat"   "medv"

# View(Boston)

boston = Boston

attach(boston)
lm1.Boston = lm(medv~lstat)
summary(lm1.Boston)

##
## Call:
## lm(formula = medv ~ lstat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.168  -3.990  -1.318   2.034  24.500
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 34.55384    0.56263   61.41  <2e-16 ***
## lstat       -0.95005    0.03873  -24.53  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.216 on 504 degrees of freedom
## Multiple R-squared:  0.5441, Adjusted R-squared:  0.5432
## F-statistic: 601.6 on 1 and 504 DF, p-value: < 2.2e-16

names(lm1.Boston)

## [1] "coefficients" "residuals"    "effects"      "rank"
## [5] "fitted.values" "assign"       "qr"           "df.residual"
## [9] "xlevels"      "call"        "terms"        "model"

lm1.Boston$rank # take a look at a specific item in the model

## [1] 2

confint(lm1.Boston)

##              2.5 %      97.5 %
## (Intercept) 33.448457 35.6592247
## lstat       -1.026148 -0.8739505
```

```
# calculate  $\hat{f}(5)$   $\hat{f}(10)$   $\hat{f}(15)$ 
predict(lm1.Boston, data.frame(lstat=c(5,10,15)), interval = 'confidence' )
```

```
##          fit      lwr      upr
## 1 29.80359 29.00741 30.59978
## 2 25.05335 24.47413 25.63256
## 3 20.30310 19.73159 20.87461
```

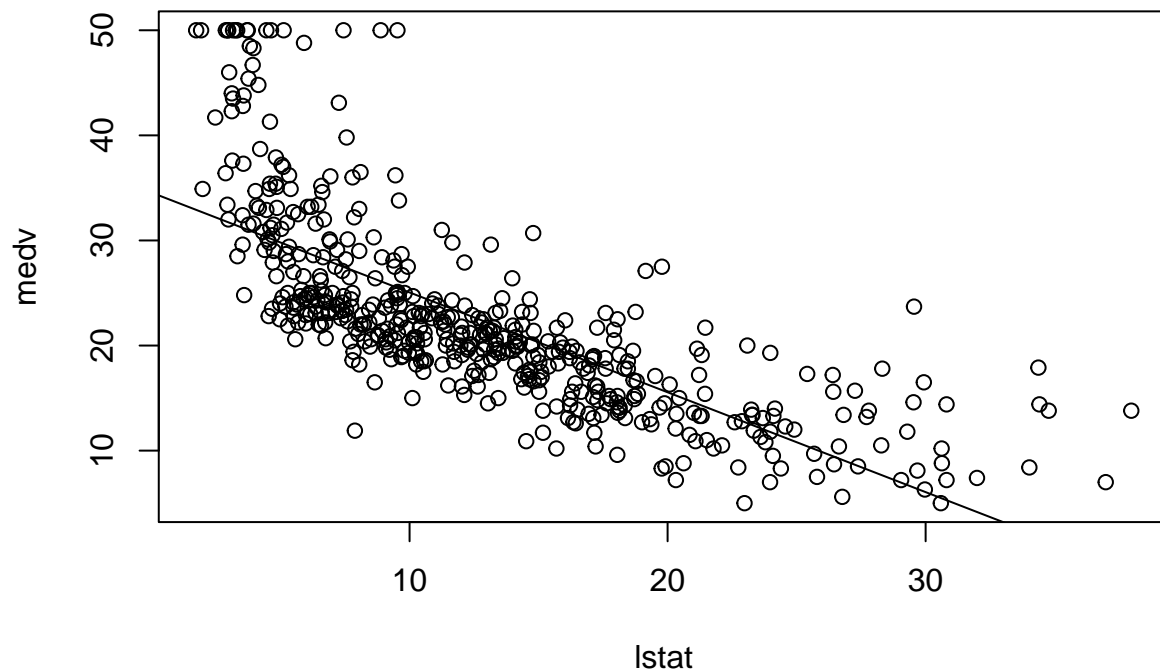
```
# the second argument of function 'predict' should be in dataframe format.
```

```
# difference between 'prediction interval' and 'confidence interval'
# confidence interval treat the true value of the mean as known.
# prediction interval has a larger range, as it first estimate the mean.
browseURL("https://online.stat.psu.edu/stat501/lesson/3/3.3")
```

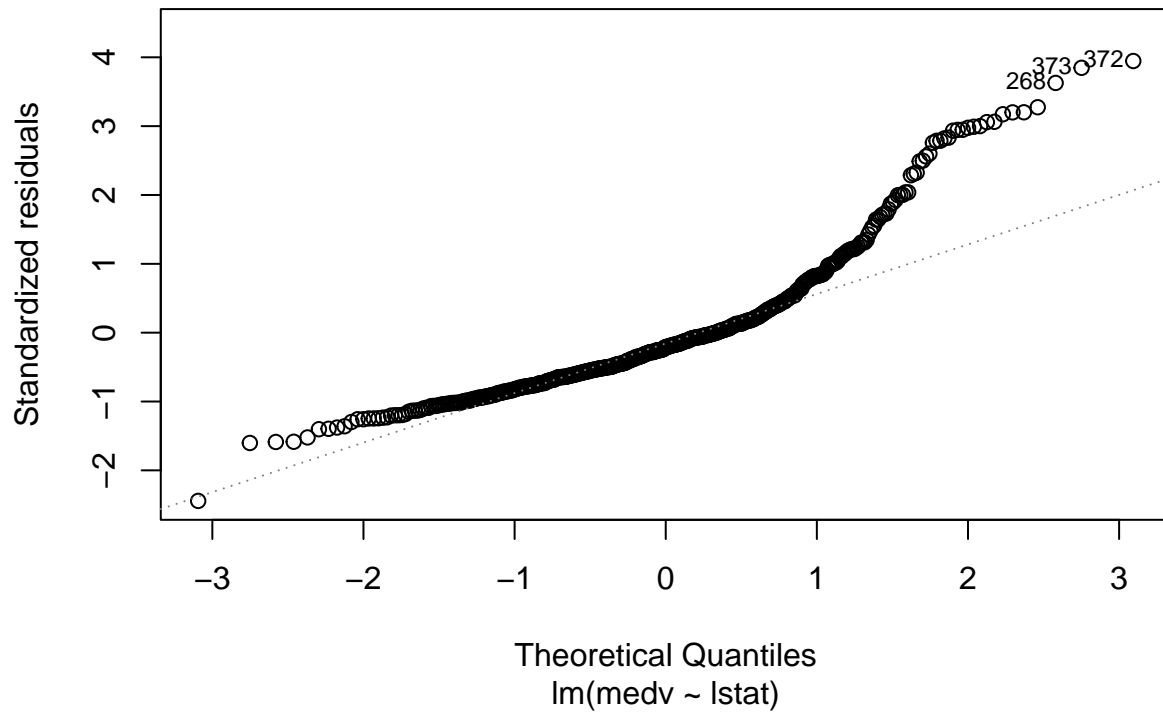
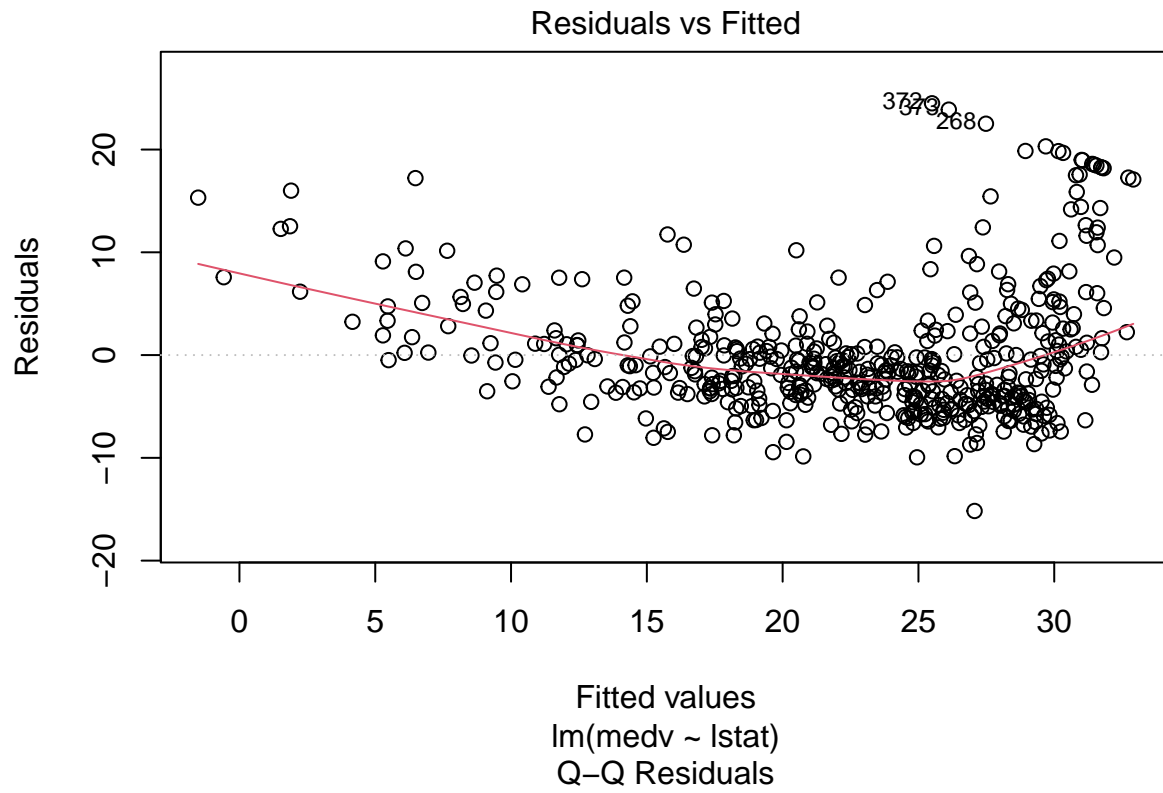
```
predict(lm1.Boston, data.frame(lstat=c(5,10,15)), interval="prediction")
```

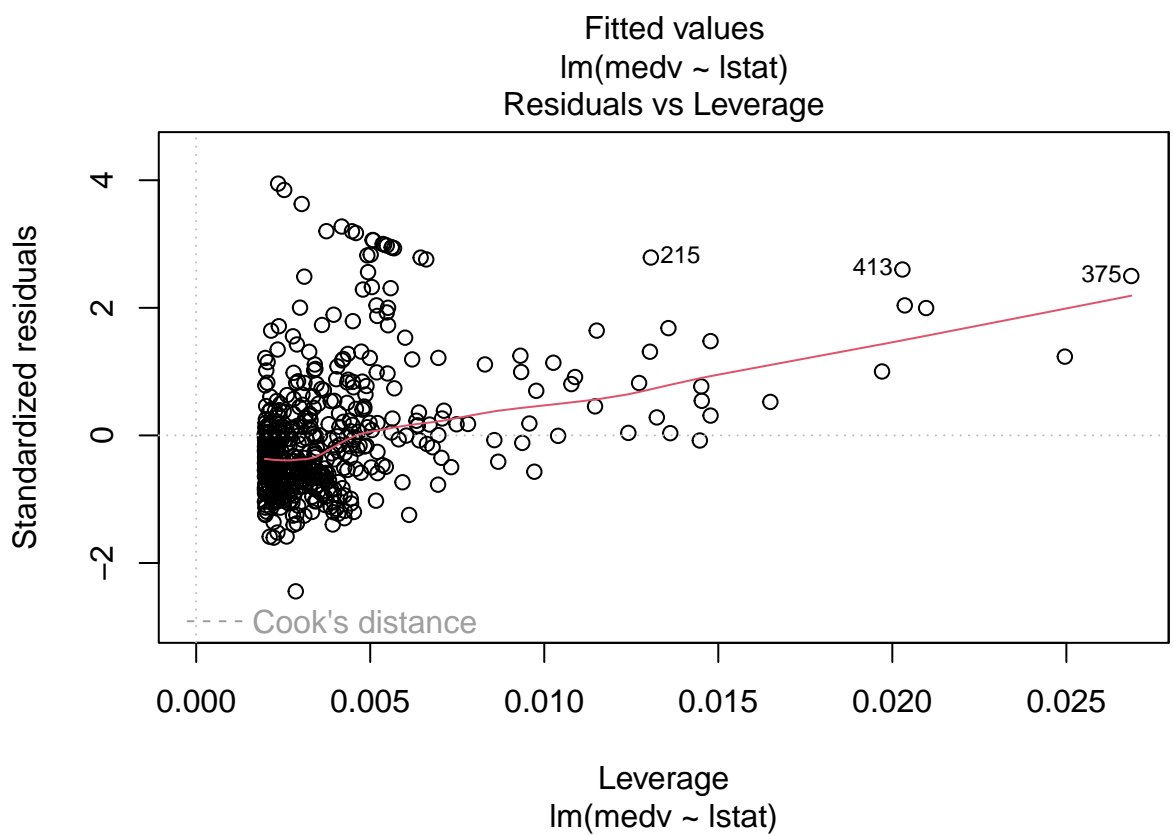
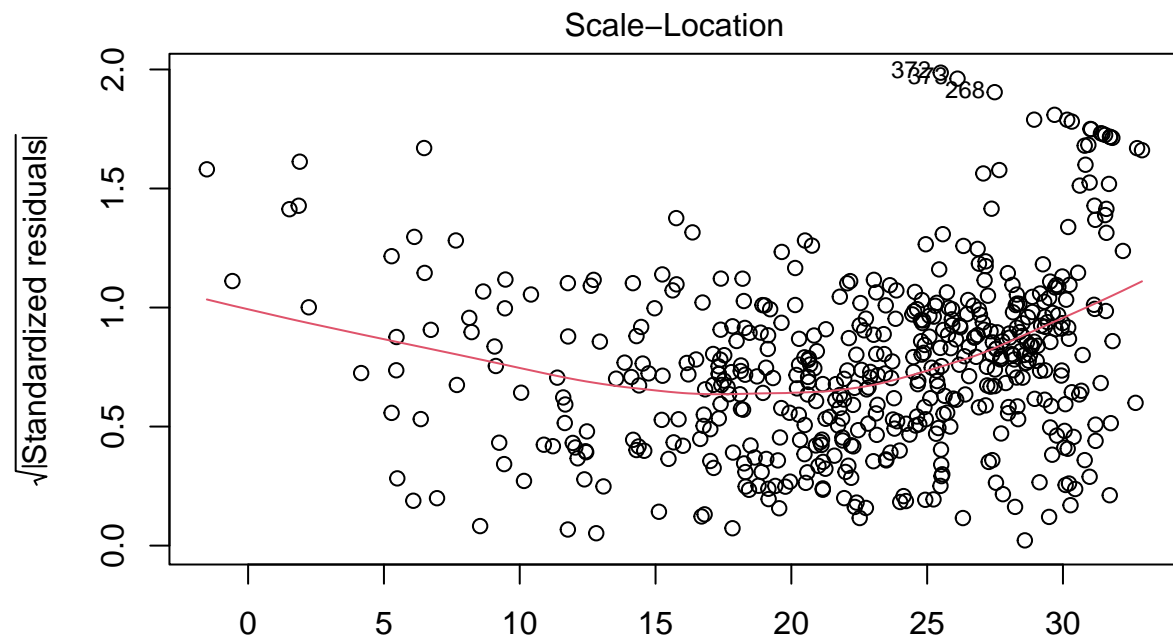
```
##          fit      lwr      upr
## 1 29.80359 17.565675 42.04151
## 2 25.05335 12.827626 37.27907
## 3 20.30310  8.077742 32.52846
```

```
plot(lstat, medv) # scatter plot
abline(lm1.Boston) # add a straight line
```



```
# par(mfrow=c(2,2))
plot(lm1.Boston) # model diagnostic checking
```





Delete 215-th 413-rd 375-th observations

```
Boston1 = Boston[-c(215,413,375), ]
dim(Boston1)
```

```
## [1] 503 14
```

```
lm11 = lm(Boston1$medv~Boston1$lstat )
summary(lm11)

##
## Call:
## lm(formula = Boston1$medv ~ Boston1$lstat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.268  -3.986  -1.302   1.968  24.470
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   34.98853    0.56008   62.47  <2e-16 ***
## Boston1$lstat -0.99246    0.03909  -25.39  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.099 on 501 degrees of freedom
## Multiple R-squared:  0.5627, Adjusted R-squared:  0.5618
## F-statistic: 644.6 on 1 and 501 DF,  p-value: < 2.2e-16
```

Remove covariates

```
lm2.Boston = lm(medv~., Boston)
summary(lm2.Boston)

##
## Call:
## lm(formula = medv ~ ., data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.595  -2.730  -0.518   1.777  26.199
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.646e+01  5.103e+00   7.144 3.28e-12 ***
## crim        -1.080e-01  3.286e-02  -3.287 0.001087 **
## zn           4.642e-02  1.373e-02   3.382 0.000778 ***
## indus        2.056e-02  6.150e-02   0.334 0.738288
## chas         2.687e+00  8.616e-01   3.118 0.001925 **
## nox         -1.777e+01  3.820e+00  -4.651 4.25e-06 ***
## rm           3.810e+00  4.179e-01   9.116 < 2e-16 ***
## age          6.922e-04  1.321e-02   0.052 0.958229
## dis         -1.476e+00  1.995e-01  -7.398 6.01e-13 ***
## rad          3.060e-01  6.635e-02   4.613 5.07e-06 ***
## tax         -1.233e-02  3.760e-03  -3.280 0.001112 **
## ptratio     -9.527e-01  1.308e-01  -7.283 1.31e-12 ***
## black        9.312e-03  2.686e-03   3.467 0.000573 ***
## lstat       -5.248e-01  5.072e-02 -10.347 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 4.745 on 492 degrees of freedom
## Multiple R-squared:  0.7406, Adjusted R-squared:  0.7338
## F-statistic: 108.1 on 13 and 492 DF,  p-value: < 2.2e-16

p-value of 'age' is 0.958229, remove it.

lm3.Boston = lm(medv ~ .-age, Boston)
summary(lm3.Boston)

##
## Call:
## lm(formula = medv ~ . - age, data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.6054  -2.7313  -0.5188   1.7601  26.2243
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  36.436927   5.080119   7.172 2.72e-12 ***
## crim        -0.108006   0.032832  -3.290 0.001075 **
## zn           0.046334   0.013613   3.404 0.000719 ***
## indus        0.020562   0.061433   0.335 0.737989
## chas         2.689026   0.859598   3.128 0.001863 **
## nox        -17.713540   3.679308  -4.814 1.97e-06 ***
## rm           3.814394   0.408480   9.338 < 2e-16 ***
## dis        -1.478612   0.190611  -7.757 5.03e-14 ***
## rad          0.305786   0.066089   4.627 4.75e-06 ***
## tax        -0.012329   0.003755  -3.283 0.001099 **
## ptratio    -0.952211   0.130294  -7.308 1.10e-12 ***
## black        0.009321   0.002678   3.481 0.000544 ***
## lstat      -0.523852   0.047625 -10.999 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

##
## Residual standard error: 4.74 on 493 degrees of freedom
## Multiple R-squared:  0.7406, Adjusted R-squared:  0.7343
## F-statistic: 117.3 on 12 and 493 DF,  p-value: < 2.2e-16
```

continue to remove 'indus'

```
lm4.Boston = update(lm3.Boston, ~ .-indus)
summary(lm4.Boston)

##
## Call:
## lm(formula = medv ~ crim + zn + chas + nox + rm + dis + rad +
##      tax + ptratio + black + lstat, data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.5984  -2.7386  -0.5046   1.7273  26.2373
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  36.436927   5.080119   7.172 2.72e-12 ***
## crim        -0.108006   0.032832  -3.290 0.001075 **
## zn           0.046334   0.013613   3.404 0.000719 ***
## chas         2.689026   0.859598   3.128 0.001863 **
## nox        -17.713540   3.679308  -4.814 1.97e-06 ***
## rm           3.814394   0.408480   9.338 < 2e-16 ***
## dis        -1.478612   0.190611  -7.757 5.03e-14 ***
## rad          0.305786   0.066089   4.627 4.75e-06 ***
## tax        -0.012329   0.003755  -3.283 0.001099 **
## ptratio    -0.952211   0.130294  -7.308 1.10e-12 ***
## black        0.009321   0.002678   3.481 0.000544 ***
## lstat      -0.523852   0.047625 -10.999 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```

## (Intercept) 36.341145 5.067492 7.171 2.73e-12 ***
## crim -0.108413 0.032779 -3.307 0.001010 **
## zn 0.045845 0.013523 3.390 0.000754 ***
## chas 2.718716 0.854240 3.183 0.001551 **
## nox -17.376023 3.535243 -4.915 1.21e-06 ***
## rm 3.801579 0.406316 9.356 < 2e-16 ***
## dis -1.492711 0.185731 -8.037 6.84e-15 ***
## rad 0.299608 0.063402 4.726 3.00e-06 ***
## tax -0.011778 0.003372 -3.493 0.000521 ***
## ptratio -0.946525 0.129066 -7.334 9.24e-13 ***
## black 0.009291 0.002674 3.475 0.000557 ***
## lstat -0.522553 0.047424 -11.019 < 2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.736 on 494 degrees of freedom
## Multiple R-squared: 0.7406, Adjusted R-squared: 0.7348
## F-statistic: 128.2 on 11 and 494 DF, p-value: < 2.2e-16

```

Multiple R-squared remains unchanged - 0.7406, but adjusted R-squared increase slightly.