

Assignment3

Haojin Li (Declan)

2020/7/24

Contents

# 3.7 Exercises	1
Initiation	1
(a) Fit a multiple regression model to predict Sales using Price, Urban, and US.	1
(b) Provide an interpretation of each coefficient in the model. Be careful some of the variables in the model are qualitative!	2
(c) Write out the model in equation form, being careful to handle the qualitative variables properly.	2
(d) For which of the predictors can you reject the null hypothesis?	2
(e) On the basis of your response to the previous question, fit a smaller model that only uses the predictors for which there is evidence of association with the outcome.	2
(f) How well do the models in (a) and (e) fit the data ?	3
(g) Using the model from (e), obtain 95% confidence intervals for the coefficient(s).	3
(h) Is there evidence of outliers or high leverage observations in the model from (e) ?	3
# 3.6.1	4
# 3.6.2	4
# 3.6.3	8
# 3.6.4	10

3.7 Exercises

Initiation

```
library(ISLR)
data(Carseats)
```

(a) Fit a multiple regression model to predict Sales using Price, Urban, and US.

```
fit <- lm(Carseats$Sales ~ Carseats$Price+Carseats$Urban+Carseats$US)
summary(fit)
```

```
##
## Call:
## lm(formula = Carseats$Sales ~ Carseats$Price + Carseats$Urban +
##      Carseats$US)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.9206 -1.6220 -0.0564  1.5786  7.0581
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   13.043469    0.651012  20.036 < 2e-16 ***
## Carseats$Price -0.054459    0.005242 -10.389 < 2e-16 ***
## Carseats$UrbanYes -0.021916    0.271650  -0.081  0.936
## Carseats$USYes    1.200573    0.259042   4.635 4.86e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.472 on 396 degrees of freedom
## Multiple R-squared:  0.2393, Adjusted R-squared:  0.2335
## F-statistic: 41.52 on 3 and 396 DF, p-value: < 2.2e-16
```

(b) Provide an interpretation of each coefficient in the model. Be careful some of the variables in the model are qualitative!

When price increases by 1 and other variable stay unchanged, the sales will decrease by 0.054. When Urban increases by 1 and other variable stay unchanged, the sales will decrease by 0.021. The US sales on average 1.2 more than non-US.

(c) Write out the model in equation form, being careful to handle the qualitative variables properly.

$$\text{Sales} = 13.043469 + \text{Price} \times (-0.054459) + \text{Urban} \times (-0.021916) + \text{US} \times 1.200573$$

(d) For which of the predictors can you reject the null hypothesis?

I don't know what is null hypothesis

(e) On the basis of your response to the previous question, fit a smaller model that only uses the predictors for which there is evidence of association with the outcome.

```
fit1 <- lm(Sales ~ Price + US, data = Carseats)
summary(fit1)
```

```
##
## Call:
## lm(formula = Sales ~ Price + US, data = Carseats)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.9269 -1.6286 -0.0574  1.5766  7.0515
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 13.03079    0.63098  20.652 < 2e-16 ***
## Price       -0.05448    0.00523 -10.416 < 2e-16 ***
## USYes        1.19964    0.25846   4.641 4.71e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.469 on 397 degrees of freedom
## Multiple R-squared:  0.2393, Adjusted R-squared:  0.2354
## F-statistic: 62.43 on 2 and 397 DF,  p-value: < 2.2e-16
```

(f) How well do the models in (a) and (e) fit the data ?

We can measure the performance of models by Multiple R-squared values. Both values are around 0.23 which indicates models are not good since they are much smaller than 1

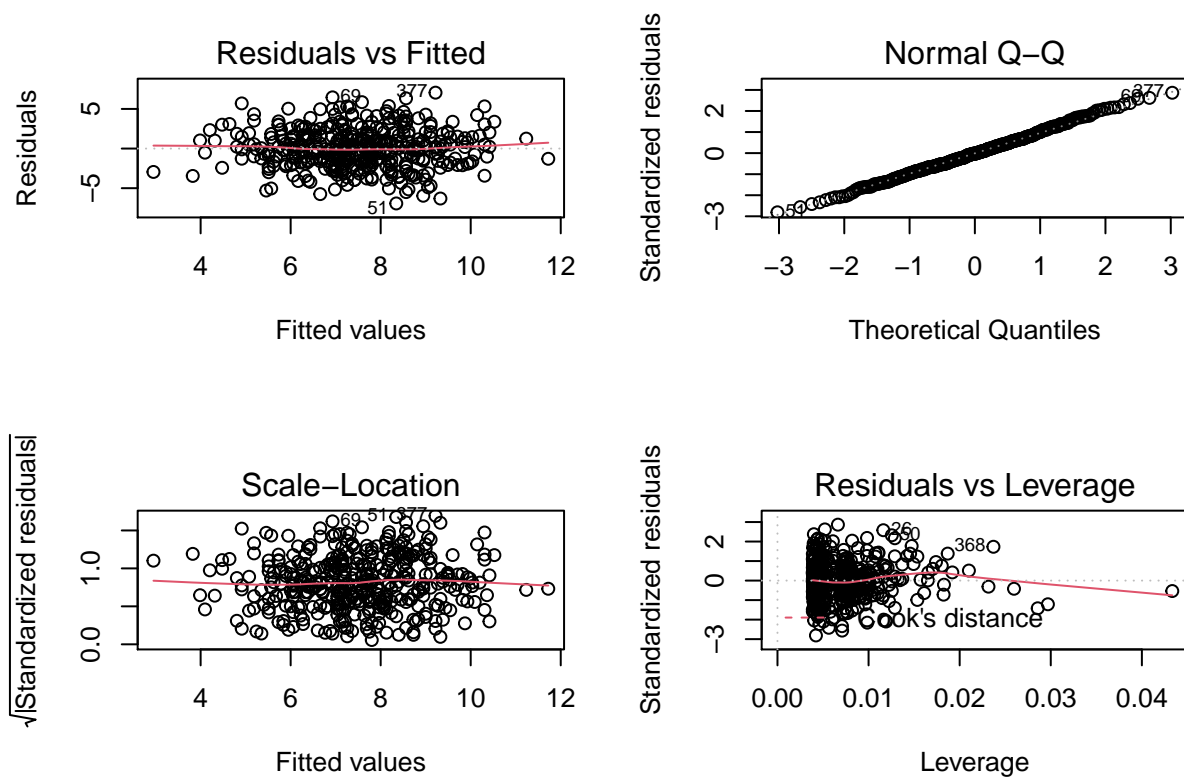
(g) Using the model from (e), obtain 95% confidence intervals for the coefficient(s).

```
confint(fit1)
```

```
##              2.5 %      97.5 %
## (Intercept) 11.79032020 14.27126531
## Price       -0.06475984 -0.04419543
## USYes        0.69151957  1.70776632
```

(h) Is there evidence of outliers or high leverage observations in the model from (e) ?

```
par(mfrow = c(2, 2))
plot(fit1)
```



3.6.1

```
library(ISLR)
library(MASS)
```

3.6.2

```
lm.fit <- lm(medv ~ lstat, data=Boston)
summary(lm.fit)
```

```
##
## Call:
## lm(formula = medv ~ lstat, data = Boston)
##
## 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
```

```
coef(lm.fit)
```

```
## (Intercept)      lstat
## 34.5538409  -0.9500494
```

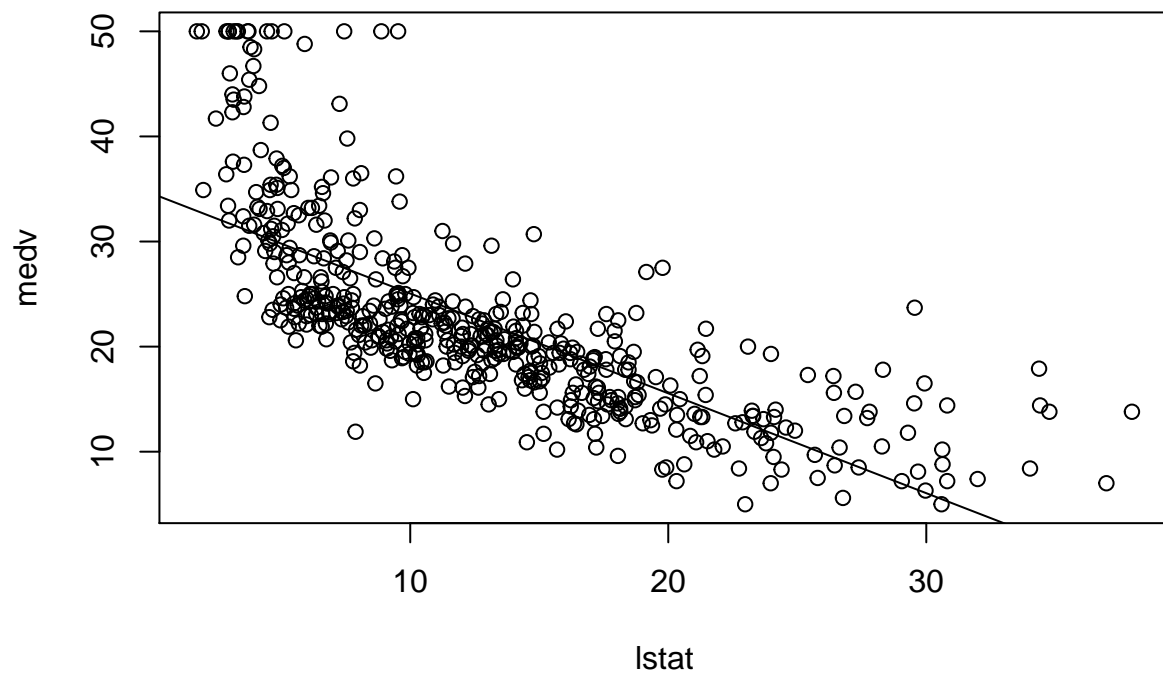
```
confint(lm.fit)
```

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

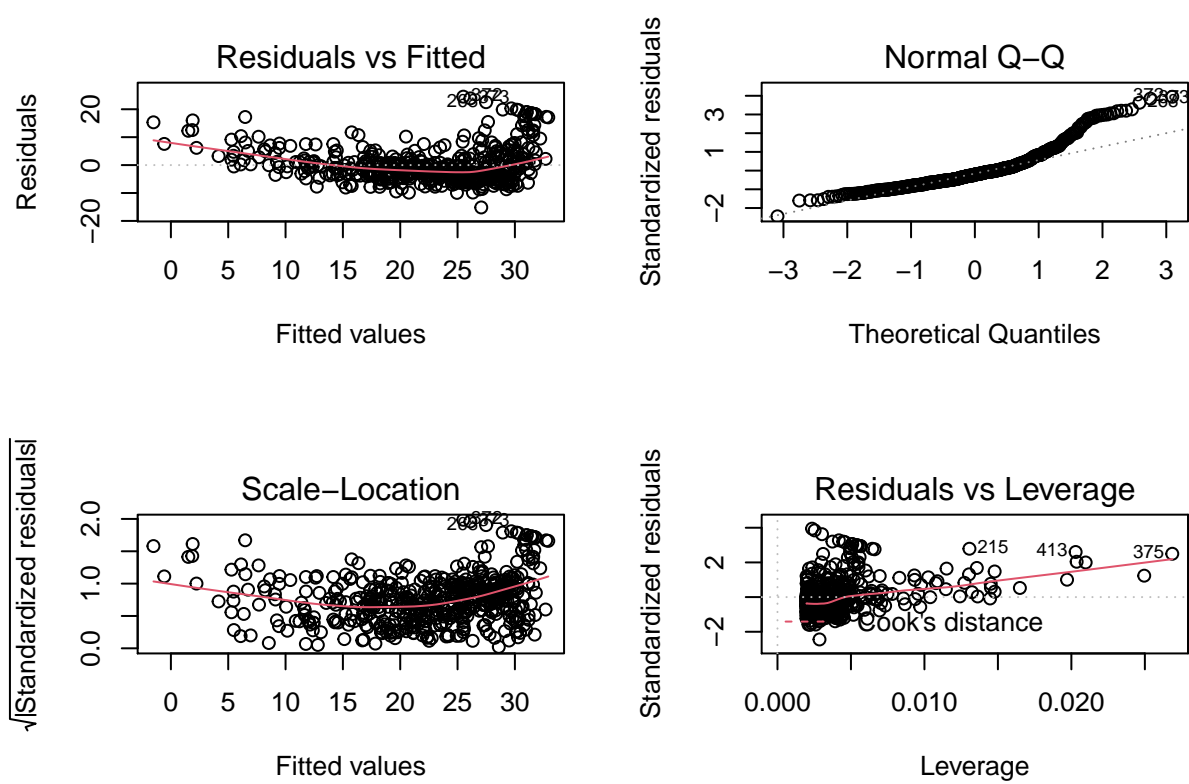
```
predict (lm.fit, data.frame(lstat=c(5,10 ,15)), interval ="confidence")
```

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

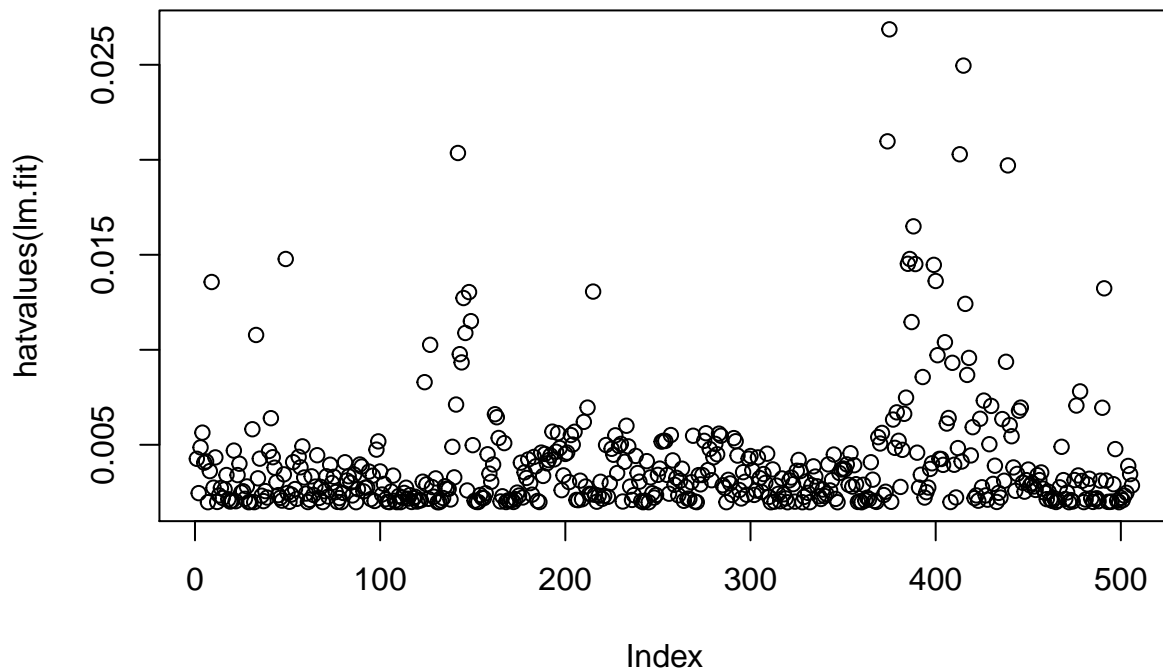
```
attach(Boston)
plot(lstat, medv)
abline(lm.fit)
```



```
par(mfrow=c(2,2))  
plot(lm.fit)
```



```
plot(hatvalues (lm.fit))
```



```
which.max(hatvalues (lm.fit))
```

```
## 375
## 375
```

3.6.3

```
lm.fit1 <- lm(medv ~ lstat+age, data=Boston )
summary(lm.fit1)
```

```
##
## Call:
## lm(formula = medv ~ lstat + age, data = Boston)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-15.981	-3.978	-1.283	1.968	23.158

```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	33.22276	0.73085	45.458	< 2e-16 ***
lstat	-1.03207	0.04819	-21.416	< 2e-16 ***
age	0.03454	0.01223	2.826	0.00491 **


```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.173 on 503 degrees of freedom
## Multiple R-squared:  0.5513, Adjusted R-squared:  0.5495
## F-statistic: 309 on 2 and 503 DF, p-value: < 2.2e-16
```

```
lm.fit2 <- lm(medv ~ ., data=Boston)
summary(lm.fit2)
```

```
##
## 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
```

```
lm.fit3 <- lm(medv ~ .-age, data=Boston)
summary(lm.fit3)
```

```
##
## 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
```

3.6.4

```
summary(lm(medv ~ lstat*age, data=Boston))
```

```
##
## Call:
## lm(formula = medv ~ lstat * age, data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.806  -4.045  -1.333   2.085  27.552
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 36.0885359  1.4698355  24.553 < 2e-16 ***
## lstat      -1.3921168  0.1674555  -8.313 8.78e-16 ***
## age       -0.0007209  0.0198792  -0.036  0.9711
## lstat:age   0.0041560  0.0018518   2.244  0.0252 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.149 on 502 degrees of freedom
## Multiple R-squared:  0.5557, Adjusted R-squared:  0.5531
## F-statistic: 209.3 on 3 and 502 DF,  p-value: < 2.2e-16
```