

ISLR Chapter3

Linear Regression: Labs

ACO

2025-02-13

Contents

Linear Regression	2
3.6.1 Libraries	2
3.6.2 Simple Linear Regression	2
3.6.3 Multiple Linear Regression	6
3.6.4 Interaction Terms	8
3.6.5 Non-linear Transformations of the Predictors	8
3.6.6 Qualitative Predictors	10
Applied	11
Q8	11
Q9	13
a	13
b	14
c	15
d	15
e	16
f	17

Linear Regression

3.6.1 Libraries

```
1 library(MASS)
2 library(ISLR2)
```

3.6.2 Simple Linear Regression

Predict `medv` using 12 predictors such as `rm` (average number of rooms per house), `age` (proportion of owner-occupied units built prior to 1940) and `lstat` (percent of households with low socioeconomic status).

```
1 str(Boston, 2)
```

```
1 ## 'data.frame':  506 obs. of  13 variables:
2 ## $ crim   : num  0.00632 0.02731 0.02729 0.03237 0.06905 ...
3 ## $ zn     : num  18 0 0 0 0 0 12.5 12.5 12.5 12.5 ...
4 ## $ indus  : num  2.31 7.07 7.07 2.18 2.18 2.18 7.87 7.87 7.87 7.87 ...
5 ## $ chas   : int   0 0 0 0 0 0 0 0 0 0 ...
6 ## $ nox    : num  0.538 0.469 0.469 0.458 0.458 0.458 0.524 0.524 0.524 0.524 ...
7 ## $ rm     : num  6.58 6.42 7.18 7 7.15 ...
8 ## $ age    : num  65.2 78.9 61.1 45.8 54.2 58.7 66.6 96.1 100 85.9 ...
9 ## $ dis    : num  4.09 4.97 4.97 6.06 6.06 ...
10 ## $ rad    : int   1 2 2 3 3 3 5 5 5 5 ...
11 ## $ tax    : num  296 242 242 222 222 222 311 311 311 311 ...
12 ## $ ptratio: num  15.3 17.8 17.8 18.7 18.7 18.7 15.2 15.2 15.2 15.2 ...
13 ## $ lstat  : num  4.98 9.14 4.03 2.94 5.33 ...
14 ## $ medv   : num  24 21.6 34.7 33.4 36.2 28.7 22.9 27.1 16.5 18.9 ...
```

```
1 lm_fit <- lm(medv ~ lstat, data = Boston)
2 summary(lm_fit)
```

```
1 ##
2 ## Call:
3 ## lm(formula = medv ~ lstat, data = Boston)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -15.168  -3.990  -1.318   2.034  24.500
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept) 34.55384    0.56263   61.41  <2e-16 ***
12 ## lstat       -0.95005    0.03873  -24.53  <2e-16 ***
13 ## ---
14 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15 ##
16 ## Residual standard error: 6.216 on 504 degrees of freedom
17 ## Multiple R-squared:  0.5441, Adjusted R-squared:  0.5432
18 ## F-statistic: 601.6 on 1 and 504 DF,  p-value: < 2.2e-16
```

confidence interval for coefficient estimates

```
1 confint(lm_fit)
```

```
1 ##                2.5 %    97.5 %  
2 ## (Intercept) 33.448457 35.6592247  
3 ## lstat      -1.026148 -0.8739505
```

confidence interval for predictions

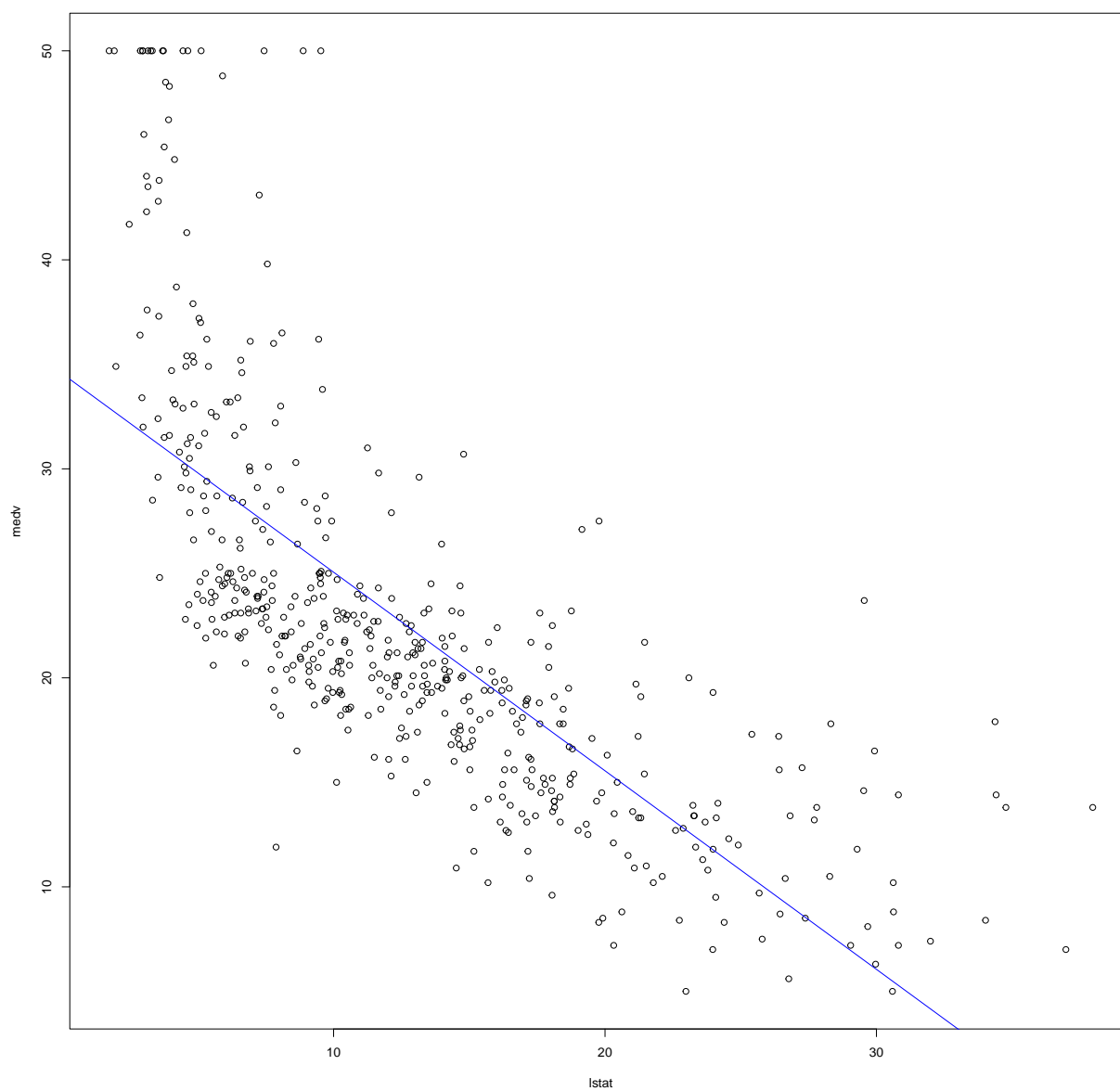
```
1 predict(  
2   lm_fit , data.frame(lstat = (c(5, 10, 15))), interval = "confidence")
```

```
1 ##      fit      lwr      upr  
2 ## 1 29.80359 29.00741 30.59978  
3 ## 2 25.05335 24.47413 25.63256  
4 ## 3 20.30310 19.73159 20.87461
```

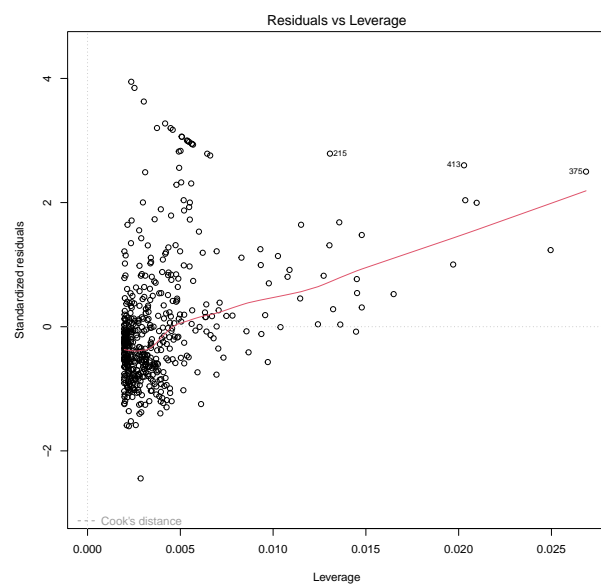
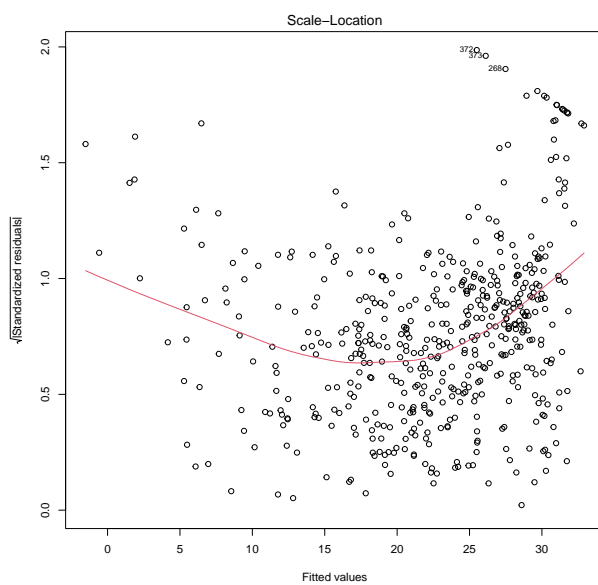
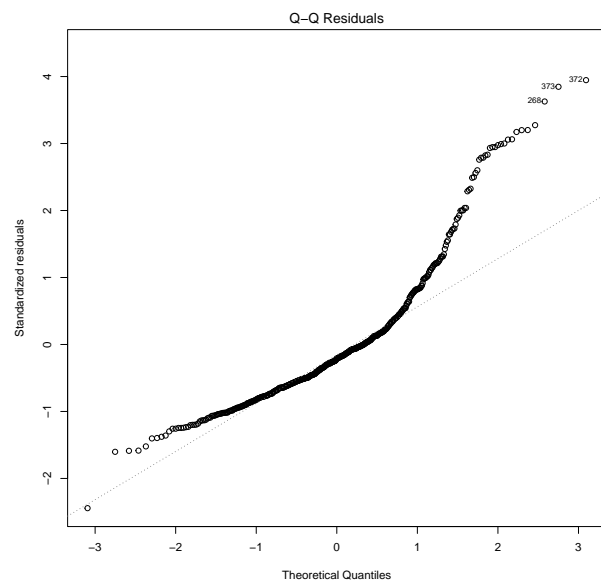
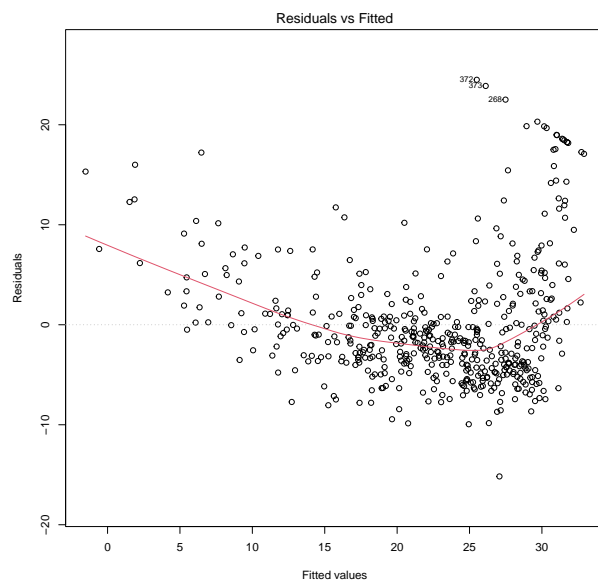
```
1 predict(  
2   lm_fit , data.frame(lstat = (c(5, 10, 15))), interval = "prediction")
```

```
1 ##      fit      lwr      upr  
2 ## 1 29.80359 17.565675 42.04151  
3 ## 2 25.05335 12.827626 37.27907  
4 ## 3 20.30310  8.077742 32.52846
```

```
1 plot(medv ~ lstat , data = Boston)  
2 abline(lm(medv ~ lstat , data = Boston), col = 'blue')
```

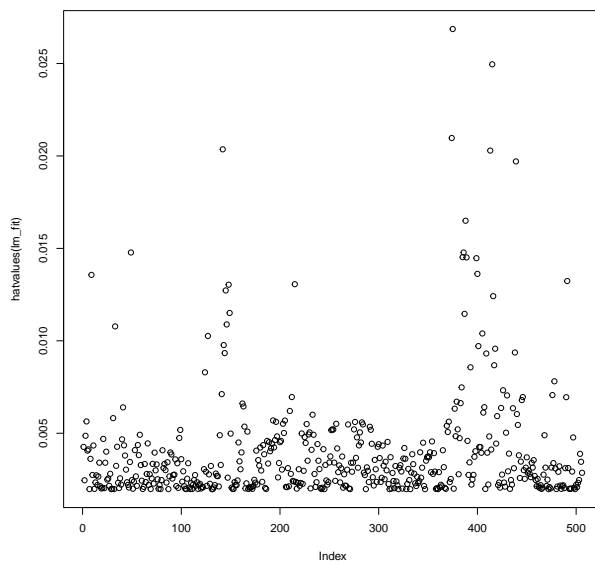
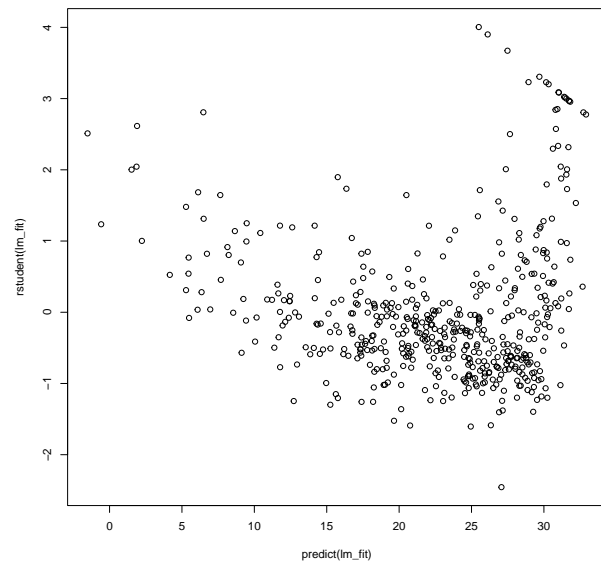
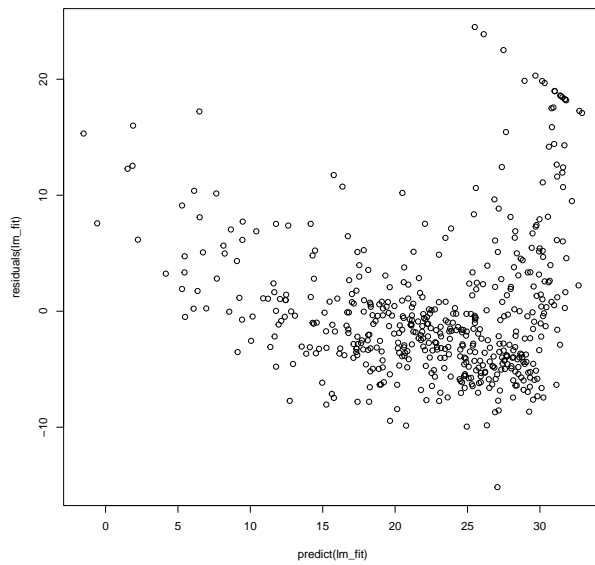


```
1 par(mfrow = c(2,2))
2 plot(lm(medv ~ lstat , data = Boston))
```



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

```
1 par(mfrow = c(2,2))
2 plot(predict(lm_fit), residuals(lm_fit))
3 plot(predict(lm_fit), rstudent(lm_fit))
4 plot(hatvalues(lm_fit))
5 par(mfrow = c(1,1))
```



```
1 which.max(hatvalues(lm_fit))
```

```
1 ## 375
```

```
2 ## 375
```

3.6.3 Multiple Linear Regression

```
1 lm_fit <- lm(medv ~ lstat + age , data = Boston)
2 summary(lm_fit)
```

```

1 ##
2 ## Call:
3 ## lm(formula = medv ~ lstat + age, data = Boston)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -15.981  -3.978  -1.283   1.968  23.158
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept) 33.22276    0.73085  45.458 < 2e-16 ***
12 ## lstat      -1.03207    0.04819 -21.416 < 2e-16 ***
13 ## age         0.03454    0.01223   2.826  0.00491 **
14 ## ---
15 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
16 ##
17 ## Residual standard error: 6.173 on 503 degrees of freedom
18 ## Multiple R-squared:  0.5513, Adjusted R-squared:  0.5495
19 ## F-statistic: 309 on 2 and 503 DF,  p-value: < 2.2e-16

```

```

1 lm_fit <- lm(medv ~ ., data = Boston)
2 summary(lm_fit)

```

```

1 ##
2 ## Call:
3 ## lm(formula = medv ~ ., data = Boston)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -15.1304  -2.7673  -0.5814   1.9414  26.2526
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept) 41.617270    4.936039   8.431 3.79e-16 ***
12 ## crim        -0.121389    0.033000  -3.678 0.000261 ***
13 ## zn          0.046963    0.013879   3.384 0.000772 ***
14 ## indus       0.013468    0.062145   0.217 0.828520
15 ## chas        2.839993    0.870007   3.264 0.001173 **
16 ## nox        -18.758022    3.851355  -4.870 1.50e-06 ***
17 ## rm          3.658119    0.420246   8.705 < 2e-16 ***
18 ## age         0.003611    0.013329   0.271 0.786595
19 ## dis        -1.490754    0.201623  -7.394 6.17e-13 ***
20 ## rad         0.289405    0.066908   4.325 1.84e-05 ***
21 ## tax        -0.012682    0.003801  -3.337 0.000912 ***
22 ## ptratio    -0.937533    0.132206  -7.091 4.63e-12 ***
23 ## lstat      -0.552019    0.050659 -10.897 < 2e-16 ***
24 ## ---
25 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
26 ##
27 ## Residual standard error: 4.798 on 493 degrees of freedom
28 ## Multiple R-squared:  0.7343, Adjusted R-squared:  0.7278
29 ## F-statistic: 113.5 on 12 and 493 DF,  p-value: < 2.2e-16

```

Hence R-Squared = `rsummary(lm_fit)$r.sq` and `RSE = summary(lm_fit)$sigma`

```
1 library(car)

1 vif(lm_fit)

1 ##      crim      zn      indus      chas      nox      rm      age      dis
2 ## 1.767486 2.298459 3.987181 1.071168 4.369093 1.912532 3.088232 3.954037
3 ##      rad      tax  ptratio      lstat
4 ## 7.445301 9.002158 1.797060 2.870777
```

3.6.4 Interaction Terms

```
1 # interaction, use : or *
2 lm_fit <- lm(medv ~ lstat:age , data = Boston)
3 summary(lm_fit)

1 ##
2 ## Call:
3 ## lm(formula = medv ~ lstat:age, data = Boston)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -13.347  -4.372  -1.534   1.914  27.193
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept) 30.1588631  0.4828240   62.46  <2e-16 ***
12 ## lstat:age    -0.0077146  0.0003799  -20.31  <2e-16 ***
13 ## ---
14 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15 ##
16 ## Residual standard error: 6.827 on 504 degrees of freedom
17 ## Multiple R-squared:  0.4501, Adjusted R-squared:  0.449
18 ## F-statistic: 412.4 on 1 and 504 DF,  p-value: < 2.2e-16
```

3.6.5 Non-linear Transformations of the Predictors

```
1 lm_fit <- lm(medv ~ lstat , data = Boston)
2 summary(lm_fit)

1 ##
2 ## Call:
3 ## lm(formula = medv ~ lstat, data = Boston)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -15.168  -3.990  -1.318   2.034  24.500
```



```

8 ##
9 ## Coefficients:
10 ##           Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept) 34.55384    0.56263   61.41  <2e-16 ***
12 ## lstat       -0.95005    0.03873  -24.53  <2e-16 ***
13 ## ---
14 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15 ##
16 ## Residual standard error: 6.216 on 504 degrees of freedom
17 ## Multiple R-squared:  0.5441, Adjusted R-squared:  0.5432
18 ## F-statistic: 601.6 on 1 and 504 DF,  p-value: < 2.2e-16

```

```

1 lm_fit2 <- lm(medv ~ lstat + I(lstat^2) , data = Boston)
2 summary(lm_fit)

```

```

1 ##
2 ## Call:
3 ## lm(formula = medv ~ lstat, data = Boston)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -15.168  -3.990  -1.318   2.034  24.500
8 ##
9 ## Coefficients:
10 ##           Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept) 34.55384    0.56263   61.41  <2e-16 ***
12 ## lstat       -0.95005    0.03873  -24.53  <2e-16 ***
13 ## ---
14 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15 ##
16 ## Residual standard error: 6.216 on 504 degrees of freedom
17 ## Multiple R-squared:  0.5441, Adjusted R-squared:  0.5432
18 ## F-statistic: 601.6 on 1 and 504 DF,  p-value: < 2.2e-16

```

```

1 anova(lm_fit, lm_fit2)

```

```

1 ## Analysis of Variance Table
2 ##
3 ## Model 1: medv ~ lstat
4 ## Model 2: medv ~ lstat + I(lstat^2)
5 ##   Res.Df  RSS Df Sum of Sq    F    Pr(>F)
6 ## 1     504 19472
7 ## 2     503 15347   1    4125.1 135.2 < 2.2e-16 ***
8 ## ---
9 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

1 lm_fit <- lm(medv ~ lstat + poly(lstat, 2, raw = TRUE) , data = Boston)
2 summary(lm_fit)

```

3.6.6 Qualitative Predictors

```
1 str(Carseats)

## 'data.frame': 400 obs. of 11 variables:
## $ Sales : num 9.5 11.22 10.06 7.4 4.15 ...
## $ CompPrice : num 138 111 113 117 141 124 115 136 132 132 ...
## $ Income : num 73 48 35 100 64 113 105 81 110 113 ...
## $ Advertising: num 11 16 10 4 3 13 0 15 0 0 ...
## $ Population : num 276 260 269 466 340 501 45 425 108 131 ...
## $ Price : num 120 83 80 97 128 72 108 120 124 124 ...
## $ ShelfLoc : Factor w/ 3 levels "Bad","Good","Medium": 1 2 3 3 1 1 3 2 3 3 ...
## $ Age : num 42 65 59 55 38 78 71 67 76 76 ...
## $ Education : num 17 10 12 14 13 16 15 10 10 17 ...
## $ Urban : Factor w/ 2 levels "No","Yes": 2 2 2 2 2 1 2 2 1 1 ...
## $ US : Factor w/ 2 levels "No","Yes": 2 2 2 2 1 2 1 2 1 2 ...
```

```
1 lm_fit <- lm(
2 Sales ~ . + Income:Advertising + Price:Age, data = Carseats)
3 summary(lm_fit)
```

```
1 ##
2 ## Call:
3 ## lm(formula = Sales ~ . + Income:Advertising + Price:Age, data = Carseats)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -2.9208 -0.7503  0.0177  0.6754  3.3413
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept)    6.5755654   1.0087470    6.519 2.22e-10 ***
12 ## CompPrice      0.0929371   0.0041183   22.567 < 2e-16 ***
13 ## Income         0.0108940   0.0026044    4.183 3.57e-05 ***
14 ## Advertising    0.0702462   0.0226091    3.107 0.002030 **
15 ## Population     0.0001592   0.0003679    0.433 0.665330
16 ## Price         -0.1008064   0.0074399  -13.549 < 2e-16 ***
17 ## ShelfLocGood   4.8486762   0.1528378   31.724 < 2e-16 ***
18 ## ShelfLocMedium 1.9532620   0.1257682   15.531 < 2e-16 ***
19 ## Age           -0.0579466   0.0159506   -3.633 0.000318 ***
20 ## Education     -0.0208525   0.0196131   -1.063 0.288361
21 ## UrbanYes       0.1401597   0.1124019    1.247 0.213171
22 ## USYes         -0.1575571   0.1489234   -1.058 0.290729
23 ## Income:Advertising 0.0007510  0.0002784    2.698 0.007290 **
24 ## Price:Age      0.0001068  0.0001333    0.801 0.423812
25 ## ---
26 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
27 ##
28 ## Residual standard error: 1.011 on 386 degrees of freedom
29 ## Multiple R-squared:  0.8761, Adjusted R-squared:  0.8719
30 ## F-statistic: 210 on 13 and 386 DF, p-value: < 2.2e-16
```

```
1 contrasts(Carseats$ShelveLoc)
```

```
1 ##           Good Medium
2 ## Bad           0      0
3 ## Good          1      0
4 ## Medium        0      1
```

Applied

Q8

```
1 q8_lm = lm(mpg~horsepower, data = Auto)
2 summary(q8_lm)
```

```
1 ##
2 ## Call:
3 ## lm(formula = mpg ~ horsepower, data = Auto)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -13.5710  -3.2592  -0.3435   2.7630  16.9240
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept) 39.935861   0.717499   55.66  <2e-16 ***
12 ## horsepower  -0.157845   0.006446  -24.49  <2e-16 ***
13 ## ---
14 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15 ##
16 ## Residual standard error: 4.906 on 390 degrees of freedom
17 ## Multiple R-squared:  0.6059, Adjusted R-squared:  0.6049
18 ## F-statistic: 599.7 on 1 and 390 DF, p-value: < 2.2e-16
```

- i. Yes
- ii. -0.7784268
- iii. R squared 60% iv

```
1 predict(q8_lm, newdata = data.frame(horsepower = 98))
```

```
1 ##           1
2 ## 24.46708
```

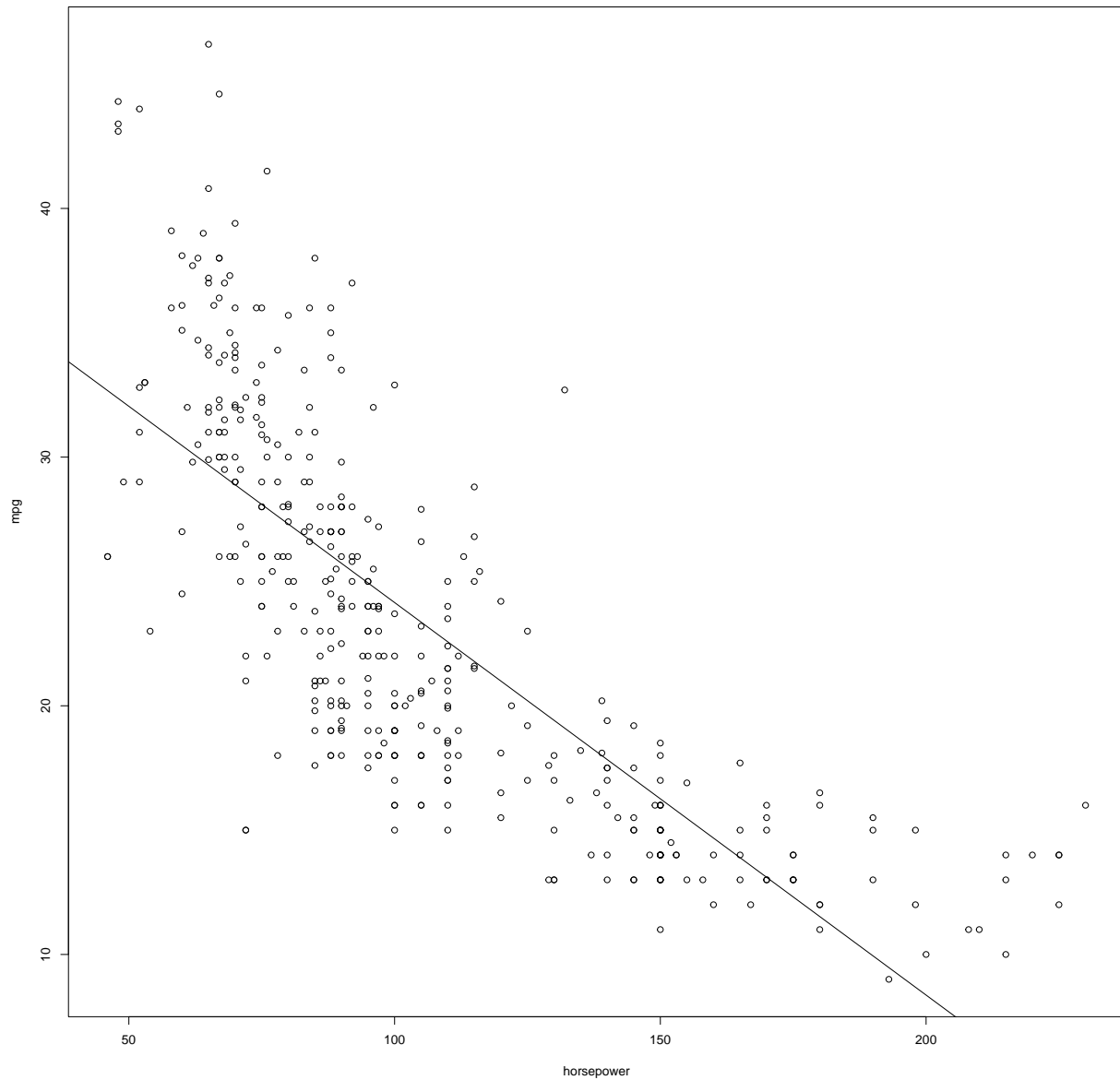
```
1 predict(q8_lm, newdata = data.frame(horsepower = 98), interval = "confidence")
```

```
1 ##           fit      lwr      upr
2 ## 1 24.46708 23.97308 24.96108
```

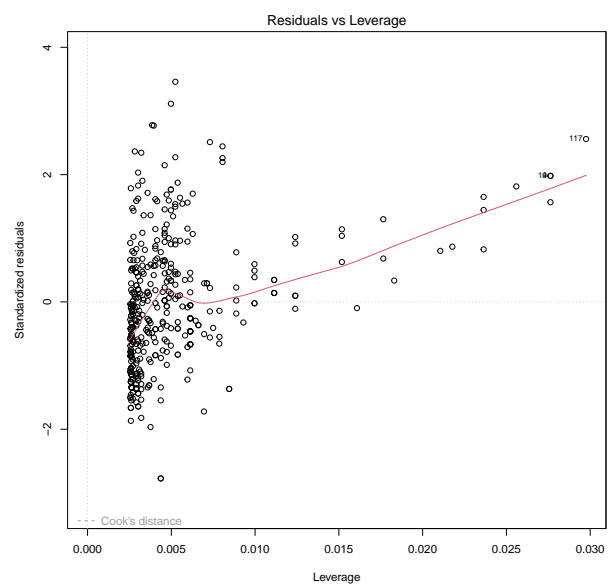
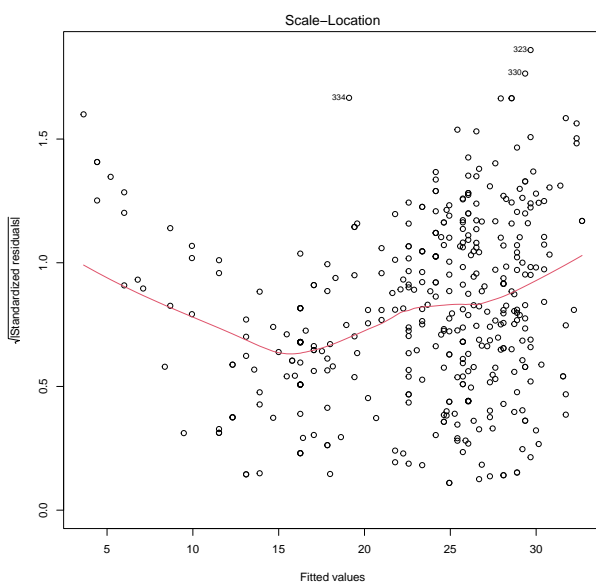
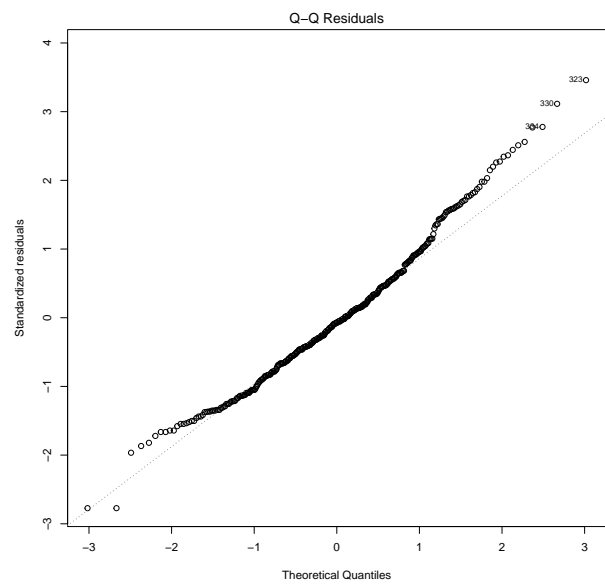
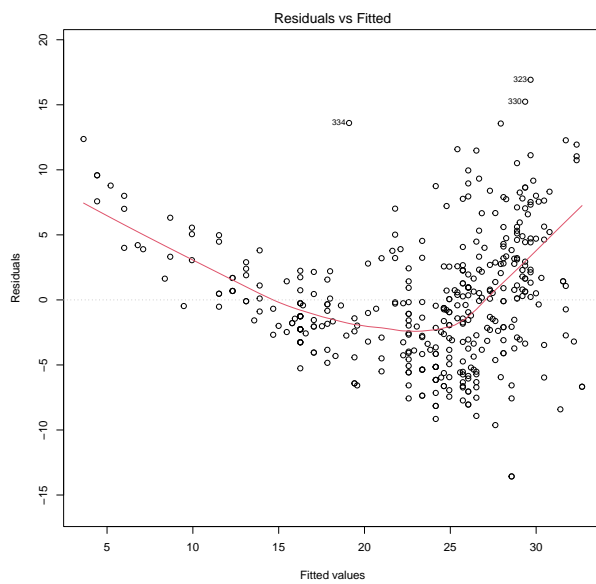
```
1 predict(q8_lm, newdata = data.frame(horsepower = 98), interval = "prediction")
```

```
1 ##          fit      lwr      upr  
2 ## 1 24.46708 14.8094 34.12476
```

```
1 q8_lm = lm(mpg~horsepower, data = Auto)  
2 plot(mpg~horsepower, data = Auto)  
3 abline(q8_lm)
```



```
1 par(mfrow = c(2,2))  
2 plot(q8_lm)
```

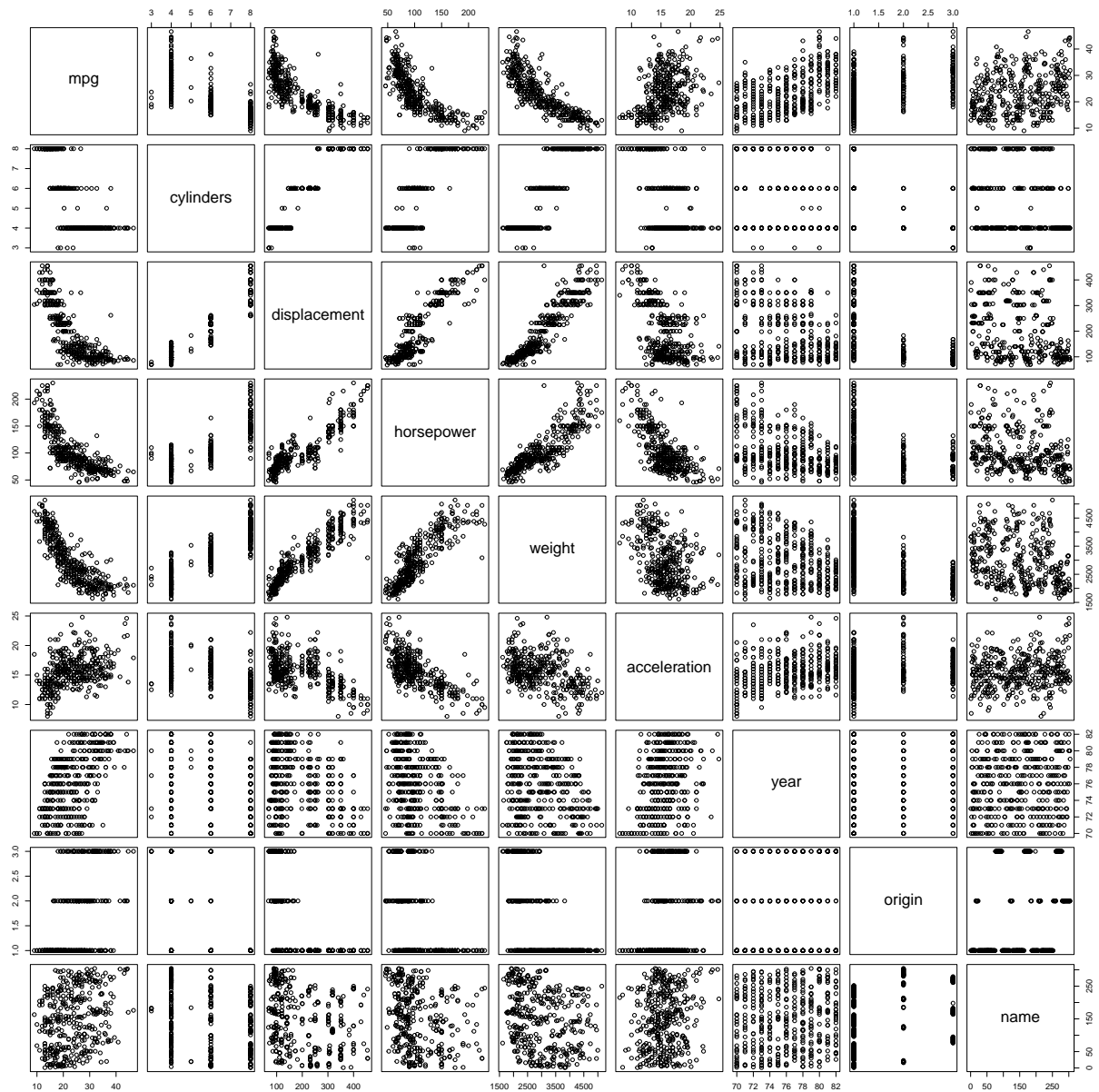


```
1 par(mfrow = c(1,1))
```

Q9

a

```
1 pairs(Auto)
```



b

```
1 cor(Auto[, sapply(Auto, FUN = \(x) is.numeric(x))])
```

```
1 ##           mpg  cylinders displacement horsepower    weight
2 ## mpg          1.000000  -0.7776175   -0.8051269  -0.7784268  -0.8322442
3 ## cylinders    -0.7776175   1.0000000    0.9508233   0.8429834   0.8975273
4 ## displacement -0.8051269   0.9508233    1.0000000   0.8972570   0.9329944
5 ## horsepower  -0.7784268   0.8429834    0.8972570   1.0000000   0.8645377
6 ## weight       -0.8322442   0.8975273    0.9329944   0.8645377   1.0000000
7 ## acceleration 0.4233285  -0.5046834   -0.5438005  -0.6891955  -0.4168392
8 ## year         0.5805410  -0.3456474   -0.3698552  -0.4163615  -0.3091199
```

```

9 ## origin      0.5652088 -0.5689316 -0.6145351 -0.4551715 -0.5850054
10 ##           acceleration      year      origin
11 ## mpg          0.4233285  0.5805410  0.5652088
12 ## cylinders    -0.5046834 -0.3456474 -0.5689316
13 ## displacement -0.5438005 -0.3698552 -0.6145351
14 ## horsepower   -0.6891955 -0.4163615 -0.4551715
15 ## weight       -0.4168392 -0.3091199 -0.5850054
16 ## acceleration  1.0000000  0.2903161  0.2127458
17 ## year          0.2903161  1.0000000  0.1815277
18 ## origin        0.2127458  0.1815277  1.0000000

```

c

```

1 q9_model = lm(mpg ~. -name, data = Auto)
2 summary(q9_model)

```

```

1 ##
2 ## Call:
3 ## lm(formula = mpg ~ . - name, data = Auto)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -9.5903 -2.1565 -0.1169  1.8690 13.0604
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept)  -17.218435   4.644294  -3.707  0.00024 ***
12 ## cylinders     -0.493376   0.323282  -1.526  0.12780
13 ## displacement  0.019896   0.007515   2.647  0.00844 **
14 ## horsepower    -0.016951   0.013787  -1.230  0.21963
15 ## weight       -0.006474   0.000652  -9.929 < 2e-16 ***
16 ## acceleration  0.080576   0.098845   0.815  0.41548
17 ## year          0.750773   0.050973  14.729 < 2e-16 ***
18 ## origin        1.426141   0.278136   5.127 4.67e-07 ***
19 ## ---
20 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
21 ##
22 ## Residual standard error: 3.328 on 384 degrees of freedom
23 ## Multiple R-squared:  0.8215, Adjusted R-squared:  0.8182
24 ## F-statistic: 252.4 on 7 and 384 DF,  p-value: < 2.2e-16

```

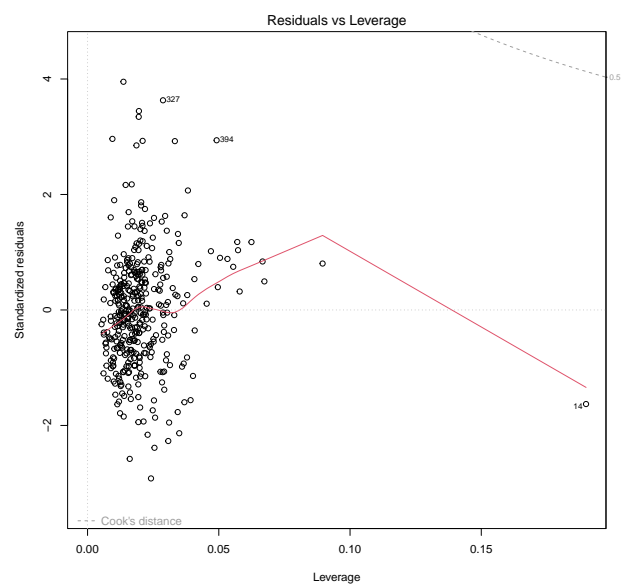
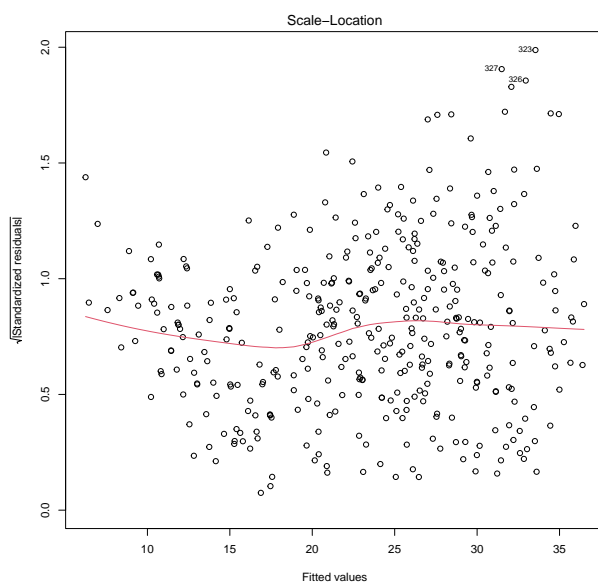
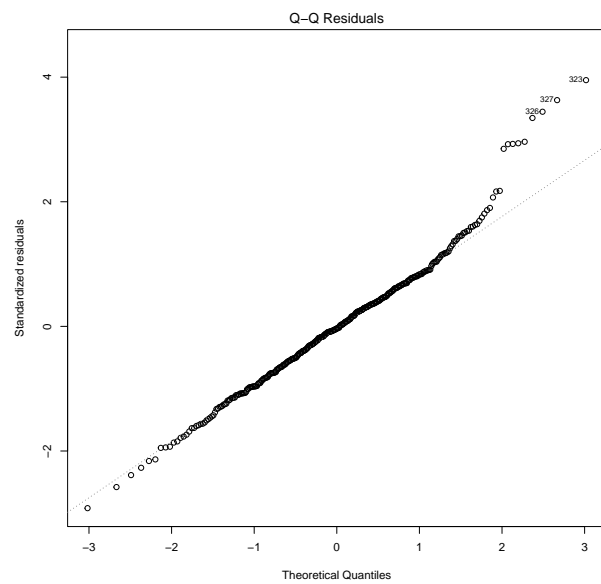
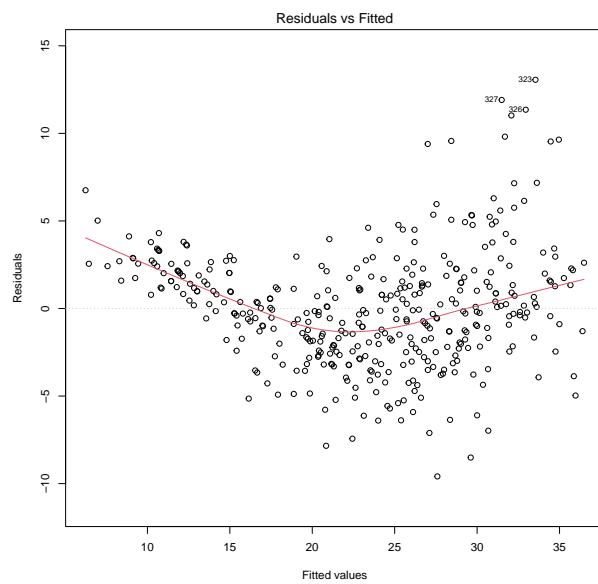
- i. Yes
- ii. displacement, weight, year, origin
- iii. Most recent cars have higher mpg

d

```

1 par(mfrow = c(2,2))
2 plot(q9_model)

```



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

e

```
1 summary(lm(mpg ~ cylinders*displacement, data = Auto))
```

```
1 ##
2 ## Call:
3 ## lm(formula = mpg ~ cylinders * displacement, data = Auto)
4 ##
5 ## Residuals:
```



```

6  ##      Min      1Q   Median      3Q      Max
7  ## -16.0432 -2.4308 -0.2263   2.2048  20.9051
8  ##
9  ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept)      48.22040    2.34712  20.545 < 2e-16 ***
12 ## cylinders        -2.41838    0.53456  -4.524 8.08e-06 ***
13 ## displacement    -0.13436    0.01615  -8.321 1.50e-15 ***
14 ## cylinders:displacement 0.01182    0.00207   5.711 2.24e-08 ***
15 ## ---
16 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
17 ##
18 ## Residual standard error: 4.454 on 388 degrees of freedom
19 ## Multiple R-squared:  0.6769, Adjusted R-squared:  0.6744
20 ## F-statistic: 271 on 3 and 388 DF, p-value: < 2.2e-16

```

```

1  summary(lm(mpg ~ cylinders:displacement, data = Auto))

```

```

1  ##
2  ## Call:
3  ## lm(formula = mpg ~ cylinders:displacement, data = Auto)
4  ##
5  ## Residuals:
6  ##      Min      1Q   Median      3Q      Max
7  ## -11.705  -3.426  -0.450   2.704  17.715
8  ##
9  ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept)      30.9896203    0.3905111   79.36 <2e-16 ***
12 ## cylinders:displacement -0.0061177    0.0002462  -24.85 <2e-16 ***
13 ## ---
14 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15 ##
16 ## Residual standard error: 4.863 on 390 degrees of freedom
17 ## Multiple R-squared:  0.6128, Adjusted R-squared:  0.6119
18 ## F-statistic: 617.4 on 1 and 390 DF, p-value: < 2.2e-16

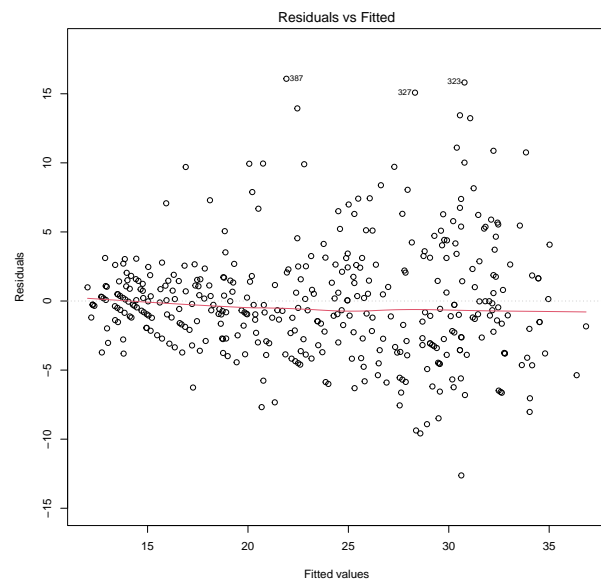
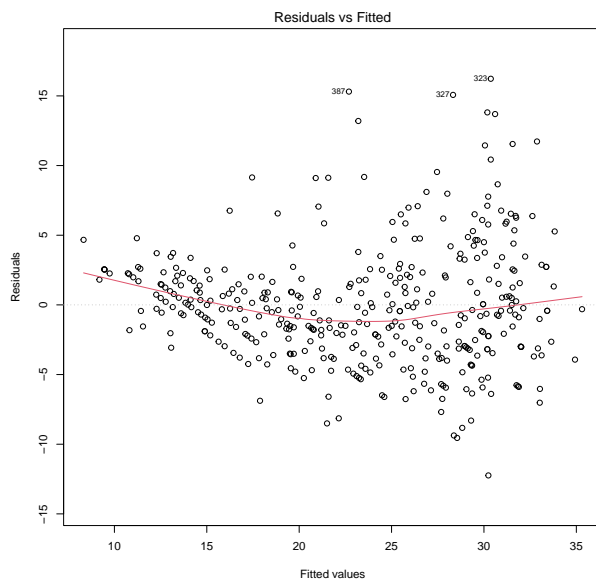
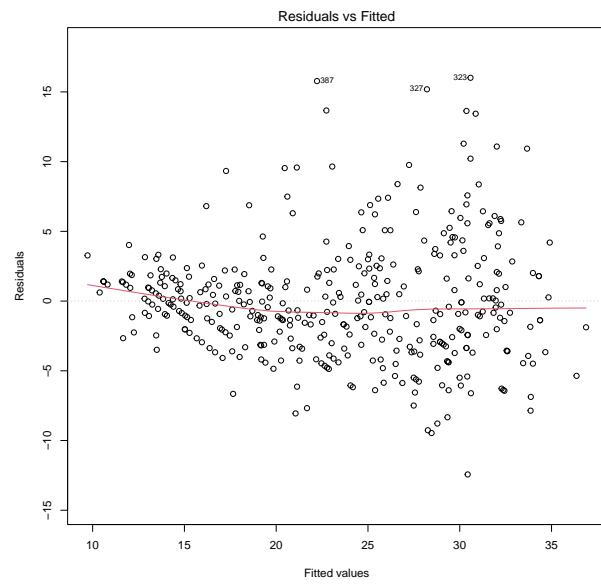
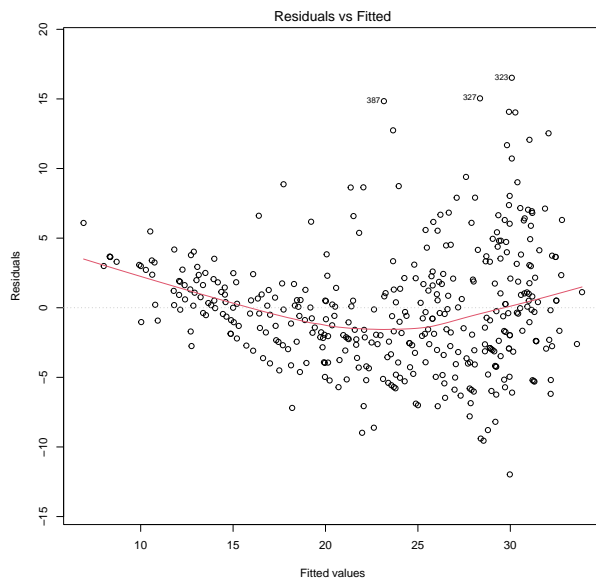
```

f

```

1  par(mfrow=c(2,2))
2  q9_lm_f = lm(mpg~weight, data= Auto)
3  plot(q9_lm_f, 1)
4
5  q9_lm_f = lm(mpg~log(weight), data= Auto)
6  plot(q9_lm_f, 1)
7
8  q9_lm_f = lm(mpg~sqrt(weight), data= Auto)
9  plot(q9_lm_f, 1)
10
11 q9_lm_f = lm(mpg~poly(weight, 2, raw = T), data= Auto)
12 plot(q9_lm_f, 1)

```



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