

HW1_Statistical Machine Learning

Moran Guo

2024-09-08

Question 1 (3.7-5)

We have $\hat{y}_i = x_i \hat{\beta} = x_i \left(\frac{\sum_{i'=1}^n x_{i'} y_{i'}}{\sum_{j=1}^n x_j^2} \right)$ so that $\hat{y}_i = \sum_{i'=1}^n \left(\frac{x_i x_{i'}}{\sum_{j=1}^n x_j^2} \right) y'_{i'}$, therefore

$$a'_i = \frac{x_i x'_i}{\sum_{j=1}^n x_j^2}.$$

Question 2 (3.7-6)

We have the least-square linear regression as $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$, we plug in $x = \bar{x}$ and $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ then we have

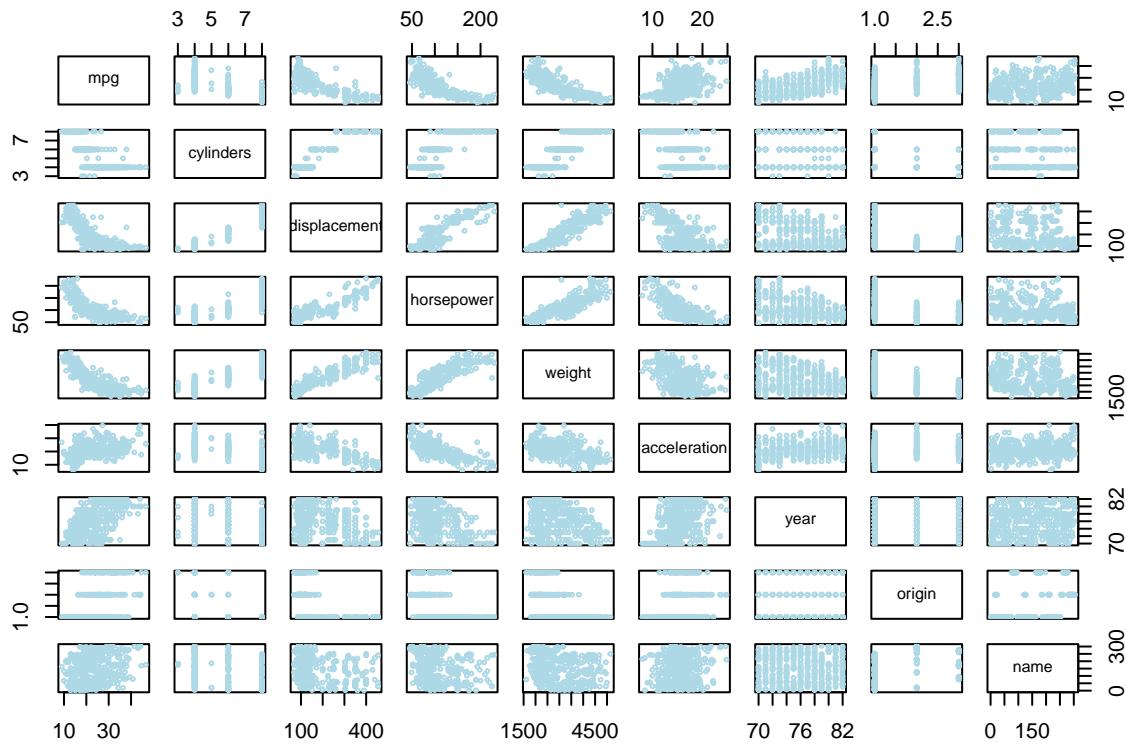
$$\hat{y} = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 \bar{x} = \bar{y}.$$

Therefore, the point (\bar{x}, \bar{y}) must be on the least-square line.

Question 3 (3.7-9)

(a)

```
# (a) Scatterplot Matrix
# install.packages("ISLR")
library(ISLR)
data(Auto)
plot(Auto, cex = 0.4, col = "lightblue")
```



(b)

```
# (b) correlation matrix
cor(Auto[, 1:8]) # Excluded the first variable `name`
```

```
##          mpg cylinders displacement horsepower      weight
## mpg      1.0000000 -0.7776175 -0.8051269 -0.7784268 -0.8322442
## cylinders -0.7776175  1.0000000  0.9508233  0.8429834  0.8975273
## displacement -0.8051269  0.9508233  1.0000000  0.8972570  0.9329944
## horsepower   -0.7784268  0.8429834  0.8972570  1.0000000  0.8645377
## weight       -0.8322442  0.8975273  0.9329944  0.8645377  1.0000000
## acceleration 0.4233285 -0.5046834 -0.5438005 -0.6891955 -0.4168392
## year         0.5805410 -0.3456474 -0.3698552 -0.4163615 -0.3091199
## origin        0.5652088 -0.5689316 -0.6145351 -0.4551715 -0.5850054
##           acceleration      year      origin
## mpg          0.4233285  0.5805410  0.5652088
## cylinders   -0.5046834 -0.3456474 -0.5689316
## displacement -0.5438005 -0.3698552 -0.6145351
## horsepower   -0.6891955 -0.4163615 -0.4551715
## weight        -0.4168392 -0.3091199 -0.5850054
## acceleration  1.0000000  0.2903161  0.2127458
## year          0.2903161  1.0000000  0.1815277
## origin        0.2127458  0.1815277  1.0000000
```

(c)

```
# (c) multiple linear regression
car_lm_fit <- lm(mpg ~ . - name, data = Auto)
summary(car_lm_fit)
```

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

```

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

```

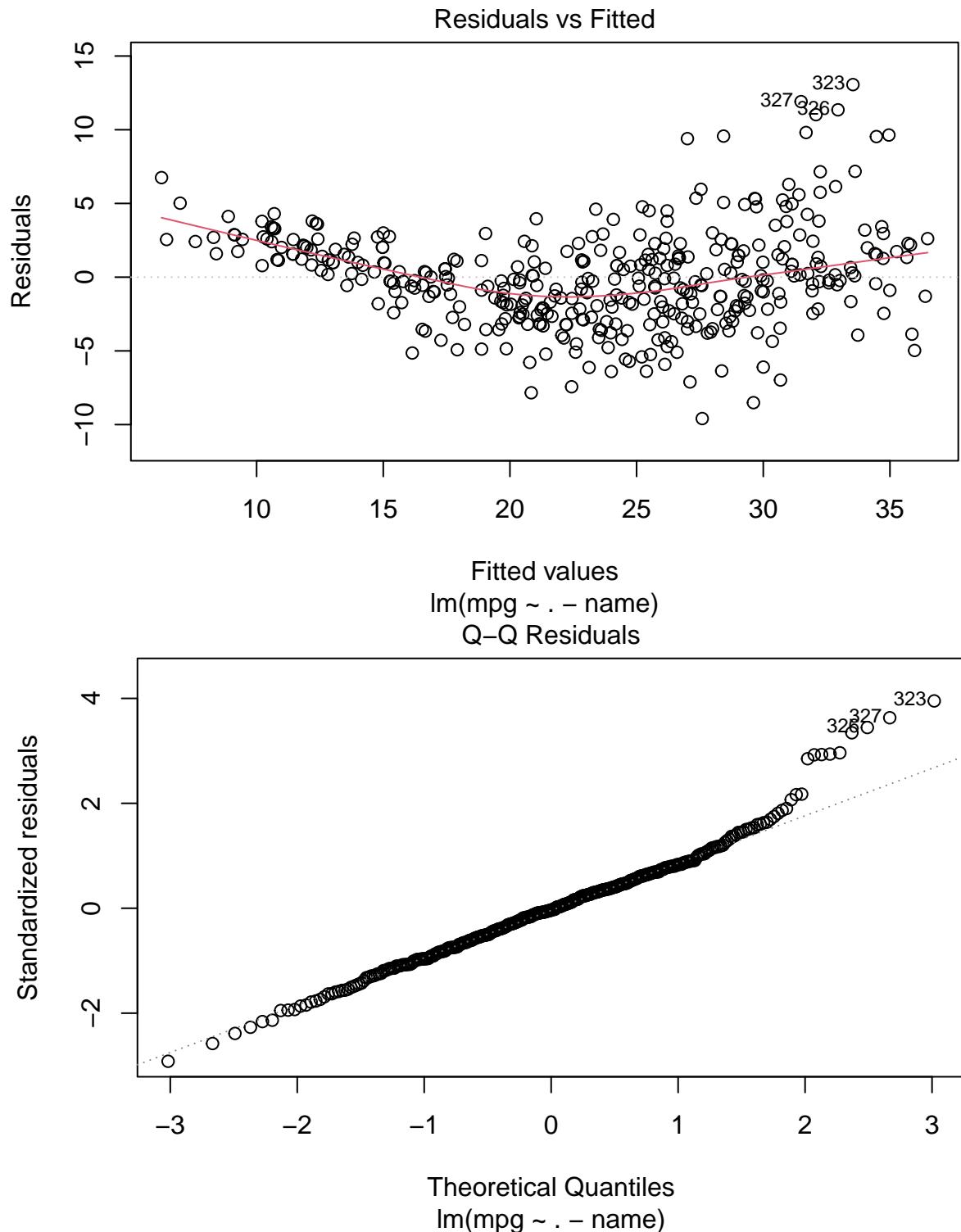
- (i) From the result, we see a significant p -value and a sufficiently large F statistics for the overall model fit, which suggests that there exist a linear relationship between mpg and other predictors.
- (ii) From the above linear model fit summary, we see variables `weight`, `year`, and `origin` are the most significant predictors.
- (iii) The coefficient of `year` is statistical significant and positive, which suggests that year is positively associated with the outcome variable mpg. This implies that as time goes by, the car manufacturers improved cars' mpg, possibly by innovating new technology.

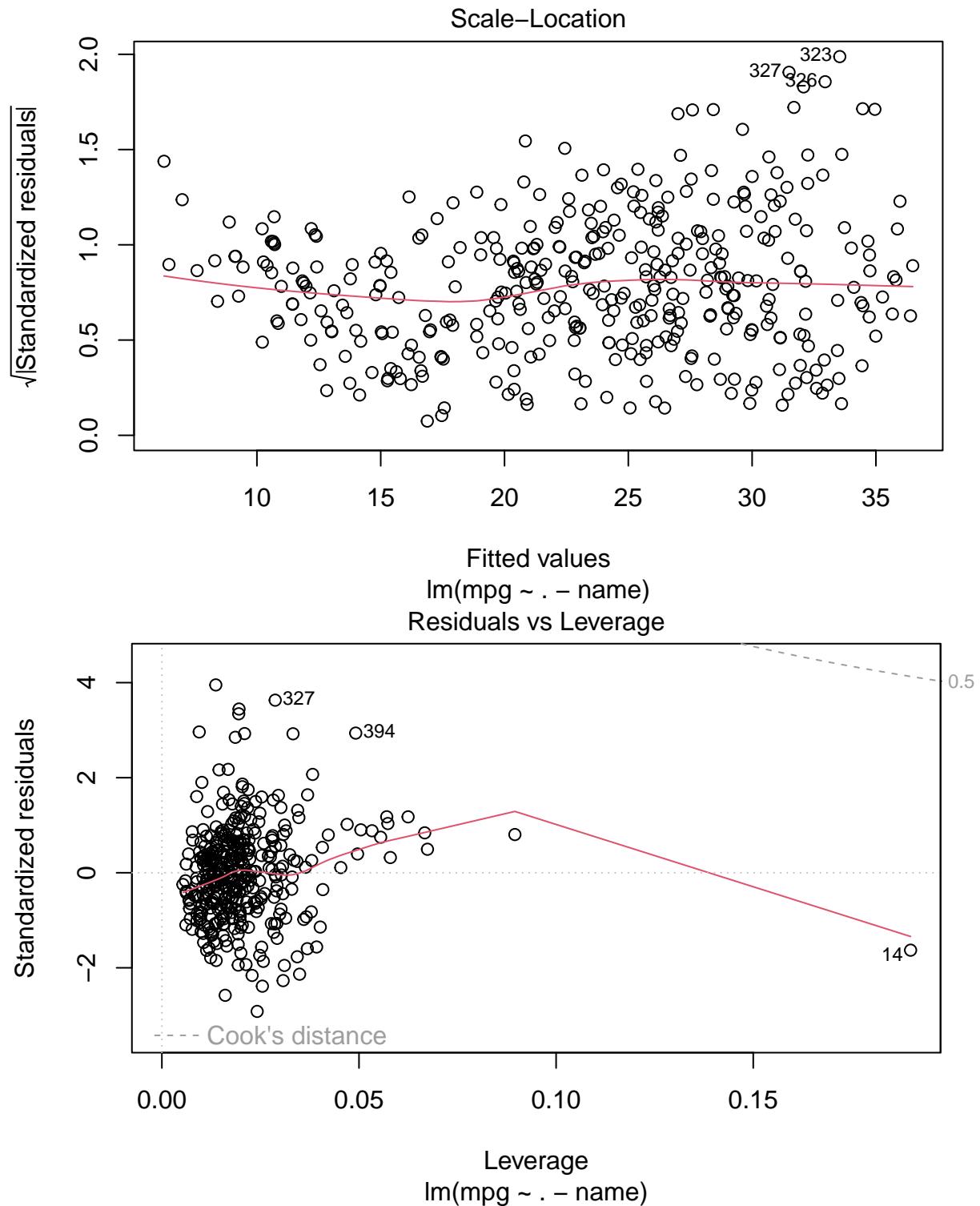
(d)

```

# (d) diagnostic plots
plot(car_lm_fit)

```





From the diagnostic plots (especially the Residuals vs. Leverage plot), we can see that observation 327, 394 seems to have a large residuals compared to the rest, which suggests that they are potentially outliers. Observation 14 has a large leverage value so that it is a point of large impact.

(e)

```

# (e) linear regression w/ interaction effects
summary(update(car_lm_fit, . ~ . + horsepower:acceleration))

##
## Call:
## lm(formula = mpg ~ cylinders + displacement + horsepower + weight +
##      acceleration + year + origin + horsepower:acceleration, data = Auto)
##
## Residuals:
##    Min     1Q Median     3Q    Max 
## -9.0329 -1.8177 -0.1183  1.7247 12.4870 
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)    
## (Intercept) -32.499820   4.923380 -6.601 1.36e-10 ***
## cylinders    0.083489   0.316913   0.263 0.792350    
## displacement -0.007649   0.008161  -0.937 0.349244    
## horsepower    0.127188   0.024746   5.140 4.40e-07 ***
## weight       -0.003976   0.000716  -5.552 5.27e-08 ***
## acceleration  0.983282   0.161513   6.088 2.78e-09 *** 
## year          0.755919   0.048179  15.690 < 2e-16 ***
## origin         1.035733   0.268962   3.851 0.000138 *** 
## horsepower:acceleration -0.012139   0.001772  -6.851 2.93e-11 ***
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.145 on 383 degrees of freedom
## Multiple R-squared:  0.841, Adjusted R-squared:  0.8376 
## F-statistic: 253.2 on 8 and 383 DF, p-value: < 2.2e-16

summary(update(car_lm_fit, . ~ . + horsepower:weight))

##
## Call:
## lm(formula = mpg ~ cylinders + displacement + horsepower + weight +
##      acceleration + year + origin + horsepower:weight, data = Auto)
##
## Residuals:
##    Min     1Q Median     3Q    Max 
## -8.589 -1.617 -0.184  1.541 12.001 
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)    
## (Intercept)  2.876e+00  4.511e+00   0.638 0.524147    
## cylinders   -2.955e-02  2.881e-01  -0.103 0.918363    
## displacement 5.950e-03  6.750e-03   0.881 0.378610    
## horsepower   -2.313e-01  2.363e-02  -9.791 < 2e-16 ***
## weight      -1.121e-02  7.285e-04 -15.393 < 2e-16 *** 
## acceleration -9.019e-02  8.855e-02  -1.019 0.309081    
## year         7.695e-01  4.494e-02  17.124 < 2e-16 *** 
## origin        8.344e-01  2.513e-01   3.320 0.000986 *** 
## horsepower:weight 5.529e-05  5.227e-06  10.577 < 2e-16 ***
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

## 
## Residual standard error: 2.931 on 383 degrees of freedom
## Multiple R-squared:  0.8618, Adjusted R-squared:  0.859
## F-statistic: 298.6 on 8 and 383 DF,  p-value: < 2.2e-16

```

Intuitively, we might think that the interactions between horsepower & acceleration **and** horsepower & car weight may have significant influence on car's mpg.

From the result of linear model fit with added interaction terms, we can see that as acceleration increases, the effect of horsepower on car's mpg becomes negative with statistical significance. Cars with higher acceleration and higher horsepower will tend to have even lower mpg than would be expected from each of these two variables individually.

For the interactions between horsepower and car weight, we see a significantly positive effect, whereas car weight & horsepower on their own have negative effects on car mpg. This suggests that as weight increases, the negative effect of horsepower on car's mpg is moderated. Or in other words, the drop of fuel efficiency of high horsepower car is less severe for high weight cars compared to lighter cars.

(f)

```

# (f) different transformation of variables

# log transform of horsepower

car_lm_fit_log_hp <- lm (mpg ~ . + log(horsepower) - name, data = Auto)
summary(car_lm_fit_log_hp)

```

```

## 
## Call:
## lm(formula = mpg ~ . + log(horsepower) - name, data = Auto)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.5777 -1.6623 -0.1213  1.4913 12.0230
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 8.674e+01 1.106e+01  7.839 4.54e-14 ***
## cylinders  -5.530e-02 2.907e-01 -0.190 0.849230  
## displacement -4.607e-03 7.108e-03 -0.648 0.517291  
## horsepower   1.764e-01 2.269e-02  7.775 7.05e-14 ***
## weight      -3.366e-03 6.561e-04 -5.130 4.62e-07 ***
## acceleration -3.277e-01 9.670e-02 -3.388 0.000776 ***
## year        7.421e-01 4.534e-02 16.368 < 2e-16 ***
## origin      8.976e-01 2.528e-01  3.551 0.000432 ***
## log(horsepower) -2.685e+01 2.652e+00 -10.127 < 2e-16 ***
## ---      
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 2.959 on 383 degrees of freedom
## Multiple R-squared:  0.8592, Adjusted R-squared:  0.8562
## F-statistic: 292.1 on 8 and 383 DF,  p-value: < 2.2e-16

# quadratic transform of horsepower

```

```

car_lm_fit_sq_hp <- lm (mpg ~ . + I((horsepower)^2) - name, data = Auto)
summary(car_lm_fit_sq_hp)

```

```

## 
## Call:
## lm(formula = mpg ~ . + I((horsepower)^2) - name, data = Auto)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -8.5497 -1.7311 -0.2236  1.5877 11.9955 
## 
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 1.3236564  4.6247696  0.286 0.774872    
## cylinders   0.3489063  0.3048310  1.145 0.253094    
## displacement -0.0075649  0.0073733 -1.026 0.305550    
## horsepower   -0.3194633  0.0343447 -9.302 < 2e-16 ***  
## weight       -0.0032712  0.0006787 -4.820 2.07e-06 ***  
## acceleration -0.3305981  0.0991849 -3.333 0.000942 ***  
## year         0.7353414  0.0459918 15.989 < 2e-16 ***  
## origin        1.0144130  0.2545545  3.985 8.08e-05 ***  
## I((horsepower)^2) 0.0010060  0.0001065  9.449 < 2e-16 ***  
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 3.001 on 383 degrees of freedom
## Multiple R-squared:  0.8552, Adjusted R-squared:  0.8522 
## F-statistic: 282.8 on 8 and 383 DF,  p-value: < 2.2e-16

```

From the model fit results, we can see that both the natural log and the quadratic transform of variable `horsepower` showed significant effects. We can also observe from the F-statistics that the model improved with the addition of transformed variables. So we would expect that `horsepower` has a non-linear relationship with `mpg`, which is the reason that adding transformed variables improves the overall model.

Question 4 (3.7-15)

(a)

```

# (a) simple regression model fit

#install.packages("MASS")
library(MASS)
predictors_bos <- names(Boston)[names(Boston) != "crim"] # Exclude response variable `crim` 

#Empty lists to store statistics
coef_l <- c()
pval_l <- c()
r2_l <- c()

#For loop, looping through all variables
for (i in predictors_bos) {
  formula <- as.formula(paste("crim ~", i))
  model <- lm(formula, data=Boston)
  model_summary <- summary(model)

  #Extract estimates
  coef <- model_summary$coefficients[2, 1]
  pval <- model_summary$coefficients[2, 4]
  r2 <- model_summary$r.squared
}
```

```

#Store the results in lists
coef_1 <- c(coef_1, coef)
pval_1 <- c(pval_1, pval)
r2_1 <- c(r2_1, r2)
}

#Create dataframe
summary_tab <- data.frame(
  Predictor = predictors_bos,
  Coefficient = coef_1,
  P_value = pval_1,
  R_squared = r2_1
)

# Display the table
summary_tab

##      Predictor Coefficient      P_value   R_squared
## 1       zn -0.07393498 5.506472e-06 0.040187908
## 2     indus  0.50977633 1.450349e-21 0.165310070
## 3     chas -1.89277655 2.094345e-01 0.003123869
## 4     nox  31.24853120 3.751739e-23 0.177217182
## 5      rm -2.68405122 6.346703e-07 0.048069117
## 6     age  0.10778623 2.854869e-16 0.124421452
## 7     dis -1.55090168 8.519949e-19 0.144149375
## 8     rad  0.61791093 2.693844e-56 0.391256687
## 9     tax  0.02974225 2.357127e-47 0.339614243
## 10    ptratio  1.15198279 2.942922e-11 0.084068439
## 11    black -0.03627964 2.487274e-19 0.148274239
## 12    lstat  0.54880478 2.654277e-27 0.207590933
## 13    medv -0.36315992 1.173987e-19 0.150780469

```

From the integrated result summary table, we can observe that variables `rad` and `tax` have the most significant effect on the response variable, criminal rate per Capita.

(b)

```

# (b) multiple linear regression model fit

model_fit_multi_bos <- lm(crim ~ . , data = Boston)
summary(model_fit_multi_bos)

##
## Call:
## lm(formula = crim ~ . , data = Boston)
##
## Residuals:
##      Min      1Q Median      3Q     Max 
## -9.924 -2.120 -0.353  1.019 75.051 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 17.033228   7.234903   2.354 0.018949 *  
## zn          0.044855   0.018734   2.394 0.017025 *  
## indus      -0.063855   0.083407  -0.766 0.444294

```

```

## chas      -0.749134  1.180147 -0.635 0.525867
## nox     -10.313535  5.275536 -1.955 0.051152 .
## rm       0.430131  0.612830  0.702 0.483089
## age      0.001452  0.017925  0.081 0.935488
## dis      -0.987176  0.281817 -3.503 0.000502 ***
## rad       0.588209  0.088049  6.680 6.46e-11 ***
## tax      -0.003780  0.005156 -0.733 0.463793
## ptratio   -0.271081  0.186450 -1.454 0.146611
## black    -0.007538  0.003673 -2.052 0.040702 *
## lstat     0.126211  0.075725  1.667 0.096208 .
## medv     -0.198887  0.060516 -3.287 0.001087 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.439 on 492 degrees of freedom
## Multiple R-squared:  0.454, Adjusted R-squared:  0.4396
## F-statistic: 31.47 on 13 and 492 DF, p-value: < 2.2e-16

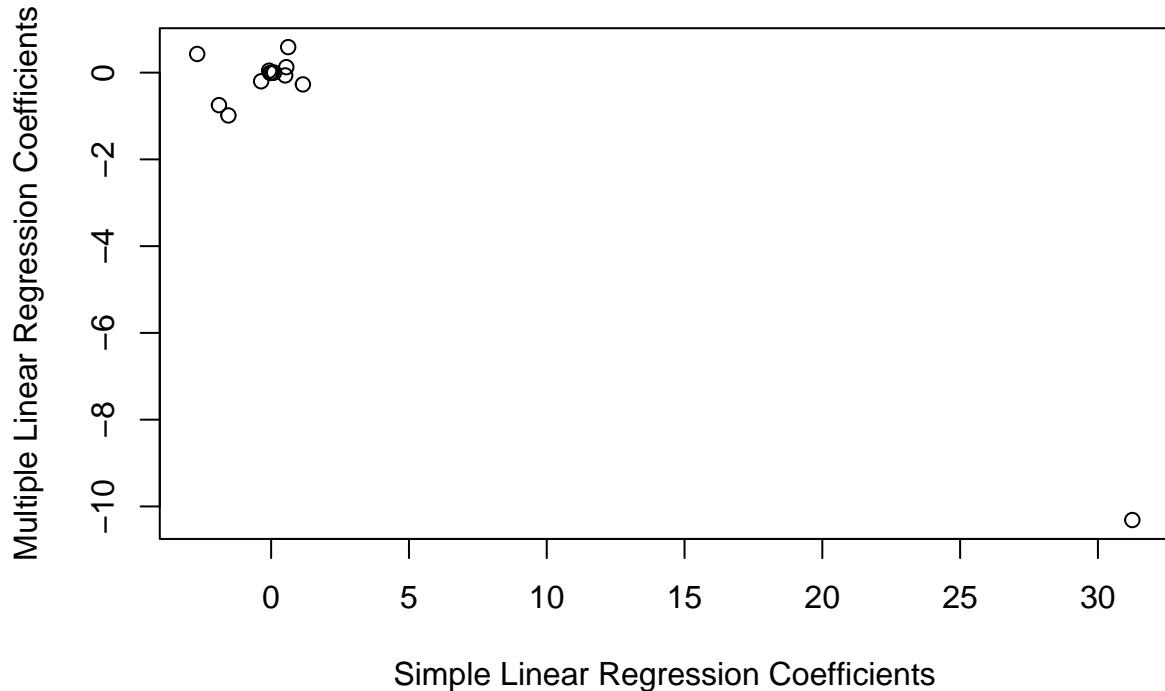
```

From the multiple linear regression result, it can be observed that variables `dis` and `rad` are the two most significant variables compared to the rest, with some of them have no significant effects on the response.

(c)

```
# (c) Compare coefficients of two models
```

```
plot(summary_tab$Coefficient, model_fit_multi_bos$coefficients[-1], xlab = "Simple Linear Regression Coefficients", ylab = "Multiple Linear Regression Coefficients")
```



From the result we can see that the coefficient of variable `nox` changed from around 31 in the simple linear regression model to around -10 in the multiple linear regression model. The coefficients of other variables also changed, but not as much as `nox`.

(d)

```

# (d) test for non-linear association

predictors_compare_bos <- names(Boston)[names(Boston) != "crim"]

predictor_names <- c()
f_stats <- c()
p_value <- c()

for (i in predictors_compare_bos) {
  #Fit a simple linear model
  linear_model <- lm(as.formula(paste("crim ~", i)), data = Boston)
  #Fit a model with quadratic & cubic terms
  nonlinear_model <- lm(as.formula(paste("crim ~", i, "+ I(", i, "^2)", "+ I(", i, "^3)")), data = Boston)
  #Use ANOVA to compare the two models
  anova_result <- anova(linear_model, nonlinear_model)
  #Extract statistics
  f_stats_temp <- anova_result$F[2]
  p_value_temp <- anova_result$`Pr(>F)`[2]

  predictor_names <- c(predictor_names, i)
  f_stats <- c(f_stats, f_stats_temp)
  p_value <- c(p_value, p_value_temp)
}

#Create a data frame to store the results
comparison_tab <- data.frame(
  Predictors = predictor_names,
  F_statistics = f_stats,
  P_value = p_value
)

comparison_tab[order(comparison_tab$P_value), ]

##      Predictors F_statistics      P_value
## 13      medv   116.6340058 2.504778e-42
## 7       dis    46.4603654 3.071837e-19
## 4       nox    42.7581707 7.122383e-18
## 2      indus   31.9869602 8.408754e-14
## 6       age    15.1400633 4.125056e-07
## 9       tax    11.6400227 1.144238e-05
## 10     ptratio   8.4155300 2.541647e-04
## 5        rm    5.3088168 5.229427e-03
## 1       zn     4.8118205 8.511995e-03
## 8       rad    3.6732699 2.607832e-02
## 12     lstat   3.3190437 3.698322e-02
## 11     black   0.4622222 6.301501e-01
## 3      chas      NA          NA

```

To compare the results between a simple linear regression model and a third order model suggested by the question, we fit the two models separately and use ANOVA test to see if the third order fit is sufficiently larger than the model we get from a simple linear regression.

From the results, we can see that variables `medv`, `dis`, `nox`, `indus`, `age`, `tax`, `ptratio`, etc., have significant

p-values which suggests a significant non-linear relationship.

Question 5 (4.7-1)

We work from Eq. (4.3) back to Eq. (4.2) to show that they are equivalent:

$$\frac{p(X)}{1-p(X)} = e^{\beta_0 + \beta_1 X} \Rightarrow p(X) = e^{\beta_0 + \beta_1 X} (1-p(X)) \Rightarrow p(X)(1+e^{\beta_0 + \beta_1 X}) = e^{\beta_0 + \beta_1 X} \Rightarrow p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1+e^{\beta_0 + \beta_1 X}}.$$

Question 6 (4.7-8)

Note that for the KNN with $N = 1$, it will have a zero error rate for the training data set. So we must have

$$\frac{0 + \text{test error rate}}{2} = 0.18 \Rightarrow \text{test error rate} = 0.36.$$

We have the information that logistic regression achieves a test error rate of 0.30, so we would prefer logistic regression because of small error rate.

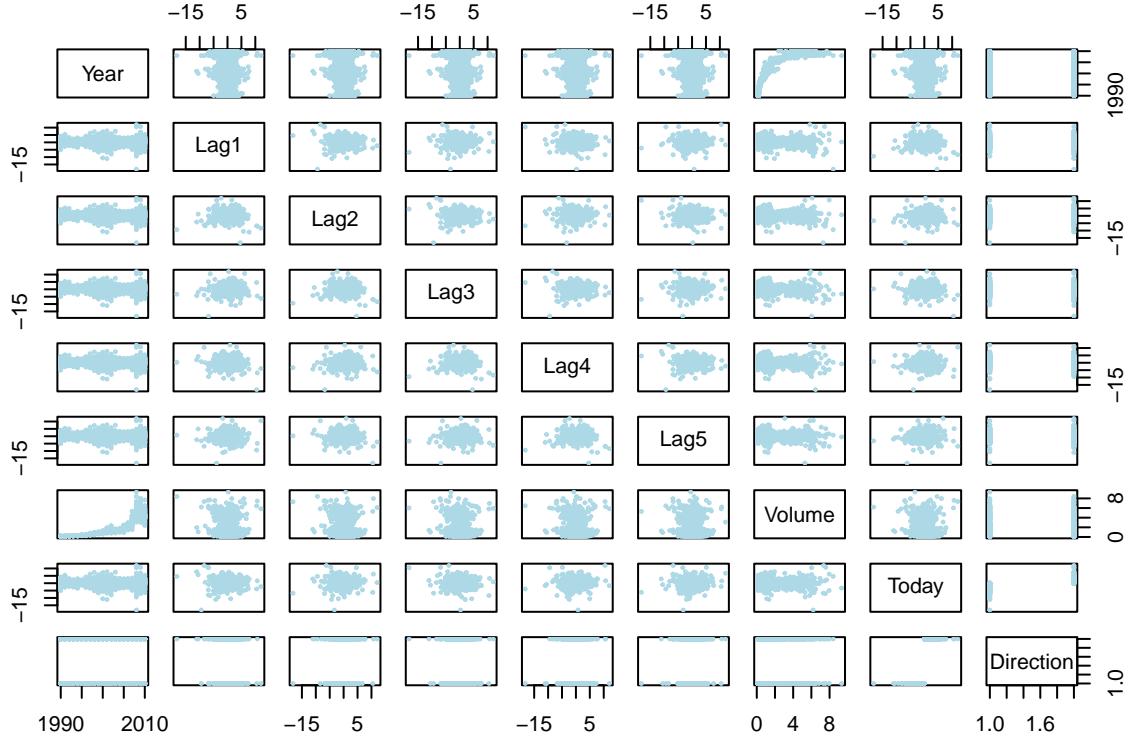
Question 7 (4.7-10)

(a)

```
# (a) Summary and plots
library(MASS)
summary(Weekly)

##      Year          Lag1          Lag2          Lag3
##  Min.   :1990   Min.   :-18.1950   Min.   :-18.1950
##  1st Qu.:1995   1st Qu.: -1.1540   1st Qu.: -1.1540
##  Median :2000   Median :  0.2410   Median :  0.2410
##  Mean   :2000   Mean   :  0.1506   Mean   :  0.1511
##  3rd Qu.:2005   3rd Qu.:  1.4050   3rd Qu.:  1.4090
##  Max.   :2010   Max.   : 12.0260   Max.   : 12.0260
##      Lag4          Lag5          Volume        Today
##  Min.   :-18.1950   Min.   :-18.1950   Min.   : 0.08747   Min.   :-18.1950
##  1st Qu.: -1.1580   1st Qu.: -1.1660   1st Qu.: 0.33202   1st Qu.: -1.1540
##  Median :  0.2380   Median :  0.2340   Median : 1.00268   Median :  0.2410
##  Mean   :  0.1458   Mean   :  0.1399   Mean   : 1.57462   Mean   :  0.1499
##  3rd Qu.:  1.4090   3rd Qu.:  1.4050   3rd Qu.: 2.05373   3rd Qu.:  1.4050
##  Max.   : 12.0260   Max.   : 12.0260   Max.   : 9.32821   Max.   : 12.0260
##      Direction
##  Down:484
##  Up :605
##
## 
## 
## 

plot(Weekly, cex = 0.3, col = "lightblue")
```



From the summary table and the plot, we can observe a relative balance distribution for all variables excluding `Year` and `Volume`, with mean values clustered around 0. The plot shows that there seems to be an exponential relationship between variable `Year` and `Volume`.

(b)

```
# (b) Logistic regression
weekly_lreg_fit <- glm(Direction ~ . - Today, data = Weekly, family = binomial)
summary(weekly_lreg_fit)
```

```
##
## Call:
## glm(formula = Direction ~ . - Today, family = binomial, data = Weekly)
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) 17.225822  37.890522  0.455   0.6494
## Year        -0.008500   0.018991 -0.448   0.6545
## Lag1        -0.040688   0.026447 -1.538   0.1239
## Lag2         0.059449   0.026970  2.204   0.0275 *
## Lag3        -0.015478   0.026703 -0.580   0.5622
## Lag4        -0.027316   0.026485 -1.031   0.3024
## Lag5        -0.014022   0.026409 -0.531   0.5955
## Volume      0.003256   0.068836  0.047   0.9623
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 1496.2  on 1088  degrees of freedom
## Residual deviance: 1486.2  on 1081  degrees of freedom
## AIC: 1502.2
```

```
##  
## Number of Fisher Scoring iterations: 4
```

From the result of logistic regression, only variable Lag2 is statistically significant.

(c)

```
# (c) Construct confusion matrix  
  
#Predict the outcomes based on predictors  
pred_probs <- predict(weekly_lreg_fit, type = "response")  
glm_pred <- rep("Down", length(pred_probs))  
glm_pred[pred_probs > 0.5] <- "Up"  
  
#Confusion matrix  
table(glm_pred, Weekly$Direction)  
  
##  
## glm_pred Down Up  
##      Down     56  47  
##      Up       428 558
```

From the table it seems that the model is predicting Up more frequently than predicting Down. So we could argue that this model tends to predict the market was up during this period of time. Another noticeable fact is that the overall correct classification rate is $\frac{56+558}{56+558+47+428} \approx 0.564$, which suggests a training error rate of 0.436, which is pretty high considering that the training error is often-time overly optimistic.

(d)

```
library(MASS)  
# (d) Logistic regression using Lag2  
  
#Specify test and train sets  
test_set <- Weekly[Weekly$Year > 2008, ]  
train_set <- Weekly[Weekly$Year <= 2008, ]  
  
lag2_lreg_fit <- glm(Direction ~ Lag2, family = binomial, data = train_set)  
  
lag_2_pred_probs <- predict(lag2_lreg_fit, newdata = test_set, type = "response")  
  
lag_2_preds <- rep("Down", length(lag_2_pred_probs))  
  
lag_2_preds[lag_2_pred_probs > 0.5] <- "Up"  
  
table(lag_2_preds, test_set$Direction)  
  
##  
## lag_2_preds Down Up  
##      Down     9  5  
##      Up       34 56
```

We can calculate the overall fraction of correct prediction is $\frac{9+56}{9+56+34+5} = 0.625$, improved from the previous model.

(e)

```
# (e) LDA model  
  
lag_2_lda_fit <- lda(Direction ~ Lag2, data = train_set)
```

```
lag_2_lda_pred <- predict(lag_2_lda_fit, newdata = test_set, type = "response")$class
table(lag_2_lda_pred, test_set$Direction)
```

```
##
## lag_2_lda_pred Down Up
##           Down    9  5
##           Up     34 56
```

Similarly, we can calculate the overall correct prediction rate: $\frac{9+56}{9+56+34+5} = 0.625$.

(f)

```
# (f) QDA model
```

```
lag_2_qda_fit <- qda(Direction ~ Lag2, data = train_set)
lag_2_qda_pred <- predict(lag_2_qda_fit, newdata = test_set, type = "response")$class
table(lag_2_qda_pred, test_set$Direction)
```

```
##
## lag_2_qda_pred Down Up
##           Down    0  0
##           Up     43 61
```

The overall fraction of correct prediction is: $\frac{61}{61+43} = 0.587$.

(g)

```
# (g) KNN model with K = 1
```

```
library(class)
train_set_KNN <- cbind(train_set$Lag2)
test_set_KNN <- cbind(test_set$Lag2)
train_response_KNN <- train_set$Direction

pred_probs_KNN <- knn(train_set_KNN, test_set_KNN, train_response_KNN) #Default value is k=1
table(pred_probs_KNN, test_set$Direction)

##
## pred_probs_KNN Down Up
##           Down   21 30
##           Up     22 31
```

The overall fraction of correct prediction for KNN with $K = 1$ is $\frac{21+31}{21+31+30+22} = 0.50$.

(h)

Among the models tested above, logistic regression and LDA performed the best based on the overall correct prediction rates.

(i)

We try to vary the value of K in the KNN model to see if the accuracy improves.

```
# (i) Modify KNN's K values
```

```
#Define a range of K values to test
k_values <- c(3, 5, 7, 9)
```

```
for (k in k_values) {
```

```

pred_probs_KNN <- knn(train_set_KNN, test_set_KNN, train_response_KNN, k = k)
confusion_matrix <- table(pred_probs_KNN, test_set$Direction)

cat("\nK =", k, "\n")
print(confusion_matrix)

accuracy <- sum(diag(confusion_matrix)) / sum(confusion_matrix)
cat("Accuracy for K =", k, ":", accuracy, "\n")
}

##  

## K = 3  

##  

## pred_probs_KNN Down Up  

##           Down   15 19  

##           Up     28 42  

## Accuracy for K = 3 : 0.5480769  

##  

## K = 5  

##  

## pred_probs_KNN Down Up  

##           Down   15 21  

##           Up     28 40  

## Accuracy for K = 5 : 0.5288462  

##  

## K = 7  

##  

## pred_probs_KNN Down Up  

##           Down   15 20  

##           Up     28 41  

## Accuracy for K = 7 : 0.5384615  

##  

## K = 9  

##  

## pred_probs_KNN Down Up  

##           Down   17 21  

##           Up     26 40  

## Accuracy for K = 9 : 0.5480769

#We use only the first three lags to see if we can improve the accuracy  

train_set_KNN <- cbind(train_set$Lag1, train_set$Lag2, train_set$Lag3)  

test_set_KNN <- cbind(test_set$Lag1, test_set$Lag2, test_set$Lag3)

train_response_KNN <- train_set$Direction

k_values <- c(1, 3, 5, 7, 9)

for (k in k_values) {
  pred_probs_KNN <- knn(train_set_KNN, test_set_KNN, train_response_KNN, k = k)
  confusion_matrix <- table(pred_probs_KNN, test_set$Direction)

  cat("\nK =", k, "\n")
  print(confusion_matrix)
}

```

```

accuracy <- sum(diag(confusion_matrix)) / sum(confusion_matrix)
cat("Accuracy for K =", k, ":", accuracy, "\n")
}

## 
## K = 1
##
## pred_probs_KNN Down Up
##       Down   19 29
##       Up    24 32
## Accuracy for K = 1 : 0.4903846
##
## K = 3
##
## pred_probs_KNN Down Up
##       Down   19 27
##       Up    24 34
## Accuracy for K = 3 : 0.5096154
##
## K = 5
##
## pred_probs_KNN Down Up
##       Down   17 23
##       Up    26 38
## Accuracy for K = 5 : 0.5288462
##
## K = 7
##
## pred_probs_KNN Down Up
##       Down   17 18
##       Up    26 43
## Accuracy for K = 7 : 0.5769231
##
## K = 9
##
## pred_probs_KNN Down Up
##       Down   15 21
##       Up    28 40
## Accuracy for K = 9 : 0.5288462

```

To improve the model, I tried two strategies:

1. Increase the value of K in the KNN model. The result indicates that if we use all variables as predictors, $K = 3$ and $K = 9$ achieved the same accuracy around 0.548.
2. Use only the first three lags Lag1, Lag2, and Lag3 to predict the results. The KNN model was constructed based on modified data frame and varied K 's. From the result, we can see that $K = 7$ in this setting performed the best.

Question 8 (4.7-11)

(a)

```

# (a) Create variable

library(ISLR)
Auto$mpg01 <- ifelse(Auto$mpg > median(Auto$mpg), 1, 0)

```

(b)

```
# (b) Examine correlation

cor_mpg <- cor(Auto[, 1:8])[-1, 1]

##      cylinders displacement horsepower      weight acceleration      year
## -0.7776175   -0.8051269   -0.7784268   -0.8322442    0.4233285   0.5805410
##      origin
## 0.5652088
```

From the result, it seems that weight, displacement, cylinders, and horsepower have a strong negative correlation with cars' mpg; while car's year & origin have a moderate positive correlation. The former is kind of self-explanatory, and we could explain the latter by the advancement of technology in car manufacturing.

(c)

```
# (c) Split data

set.seed(817) #Set seed for reproducibility

n <- nrow(Auto)

indices <- sample(1:n)

train_size <- round(0.8 * n) #A 80/20 split

train_ind <- indices[1:train_size]
test_ind <- indices[(train_size + 1):n]

train_data <- Auto[train_ind, c(-1, -9)]
test_data <- Auto[test_ind, c(-1, -9)]
```

(d)

```
# (d) LDA model fit

library(MASS)

mpg01_lda_fit <- lda(mpg01 ~ . - acceleration, data = train_data)
lda_predictions <- predict(mpg01_lda_fit, newdata = test_data, type = "response")

predicted_classes <- lda_predictions$class
actual_classes <- test_data$mpg01

lda_test_error <- mean(predicted_classes != actual_classes)

print(paste("Test Error: ", round(lda_test_error, 4)))
```

```
## [1] "Test Error: 0.1282"
```

Therefore, the error rate on test set is 0.1282.

(e)

```
# (d) QDA model fit

library(MASS)
```

```

mpg01_qda_fit <- qda(mpg01 ~ . - acceleration, data = train_data)
qda_predictions <- predict(mpg01_qda_fit, newdata = test_data, type = "response")

predicted_classes <- qda_predictions$class
actual_classes <- test_data$mpg01

qda_test_error <- mean(predicted_classes != actual_classes)

print(paste("Test Error: ", round(qda_test_error, 4)))

```

[1] "Test Error: 0.1026"

The error rate of QDA on test set is 0.1026.

(f)

```

library(MASS)

mpg01_lreg_fit <- glm(mpg01 ~ . - acceleration, data = train_data, family = binomial)

lreg_pred <- ifelse(predict(mpg01_lreg_fit, test_data, type = "response") > 0.5, 1, 0)

actual_classes <- test_data$mpg01

lreg_test_error <- mean(lreg_pred != actual_classes)

print(paste("Test Error: ", round(lreg_test_error, 4)))

```

[1] "Test Error: 0.141"

The error rate of logistic regression on test set is 0.1410.

(g)

```

# (g) KNN model fit and output the best k

library(class)
set.seed(123)
result <- data.frame(k = 1:10, error = "") #Initialize dataframe

for (k in 1:10) {
  mpg_knn_pred <- knn(train = train_data[, c(-5, -6, -7, -8)],
                        test = test_data[, c(-5, -6, -7, -8)],
                        cl = train_data$mpg01, k = k)
  result[k, 2] <- mean(mpg_knn_pred != test_data$mpg01)
}

result

##      k          error
## 1    1 0.153846153846154
## 2    2 0.16666666666666667
## 3    3 0.115384615384615
## 4    4 0.115384615384615
## 5    5 0.128205128205128
## 6    6 0.115384615384615
## 7    7 0.128205128205128

```

```
## 8 8 0.115384615384615
## 9 9 0.128205128205128
## 10 10 0.128205128205128
```

Based on the result, $K = 3, 4, 6, 8$ yields the same lowest error rate. So a valid answer will be $K = 3$.