

Module 3: Recommended Exercises

alexaoh

21.01.2021

Problem 1 (Extension from Book Ex. 9)

This question involves the use of multiple linear regression on the `Auto` data set from `ISLR` package (you may use `?Auto` to see a description of the data). First we exclude from our analysis the variable `name` and look at the data summary and structure of the dataset.

```
library(ISLR)
Auto = subset(Auto, select = -name)
# Auto$origin = factor(Auto$origin)
summary(Auto)
```

```
#>      mpg      cylinders      displacement      horsepower      weight
#> Min.   : 9.00   Min.   :3.000   Min.    : 68.0   Min.    : 46.0   Min.    :1613
#> 1st Qu.:17.00   1st Qu.:4.000   1st Qu.:105.0   1st Qu.: 75.0   1st Qu.:2225
#> Median :22.75   Median :4.000   Median :151.0   Median : 93.5   Median :2804
#> Mean   :23.45   Mean   :5.472   Mean   :194.4   Mean   :104.5   Mean   :2978
#> 3rd Qu.:29.00   3rd Qu.:8.000   3rd Qu.:275.8   3rd Qu.:126.0   3rd Qu.:3615
#> Max.   :46.60   Max.   :8.000   Max.   :455.0   Max.   :230.0   Max.   :5140
#> acceleration      year      origin
#> Min.    : 8.00   Min.    :70.00   Min.    :1.000
#> 1st Qu.:13.78   1st Qu.:73.00   1st Qu.:1.000
#> Median :15.50   Median :76.00   Median :1.000
#> Mean    :15.54   Mean    :75.98   Mean    :1.577
#> 3rd Qu.:17.02   3rd Qu.:79.00   3rd Qu.:2.000
#> Max.    :24.80   Max.    :82.00   Max.    :3.000
```

```
str(Auto)
```

```
#> 'data.frame':   392 obs. of  8 variables:
#> $ mpg      : num  18 15 18 16 17 15 14 14 14 15 ...
#> $ cylinders : num   8  8  8  8  8  8  8  8  8  8 ...
#> $ displacement: num  307 350 318 304 302 429 454 440 455 390 ...
#> $ horsepower  : num  130 165 150 150 140 198 220 215 225 190 ...
#> $ weight      : num 3504 3693 3436 3433 3449 ...
#> $ acceleration: num  12 11.5 11 12 10.5 10 9 8.5 10 8.5 ...
#> $ year        : num  70 70 70 70 70 70 70 70 70 70 ...
#> $ origin      : num   1  1  1  1  1  1  1  1  1  1 ...
```

We obtain a summary and see that all variables are numerical (continuous). However, when we check the description of the data (again with `?Auto`) we immediately see that `origin` is actually encoding for either American (`origin=1`), European (`origin=2`) or Japanese (`origin=3`) origin of the car, thus the values 1, 2 and 3 do not have any actual numerical meaning. We therefore need to first change the data type of that variable to let R know that we are dealing with a qualitative (categorical) variable, instead of a continuous one (otherwise we will obtain wrong model fits). In R such variables are called *factor variables*, and before

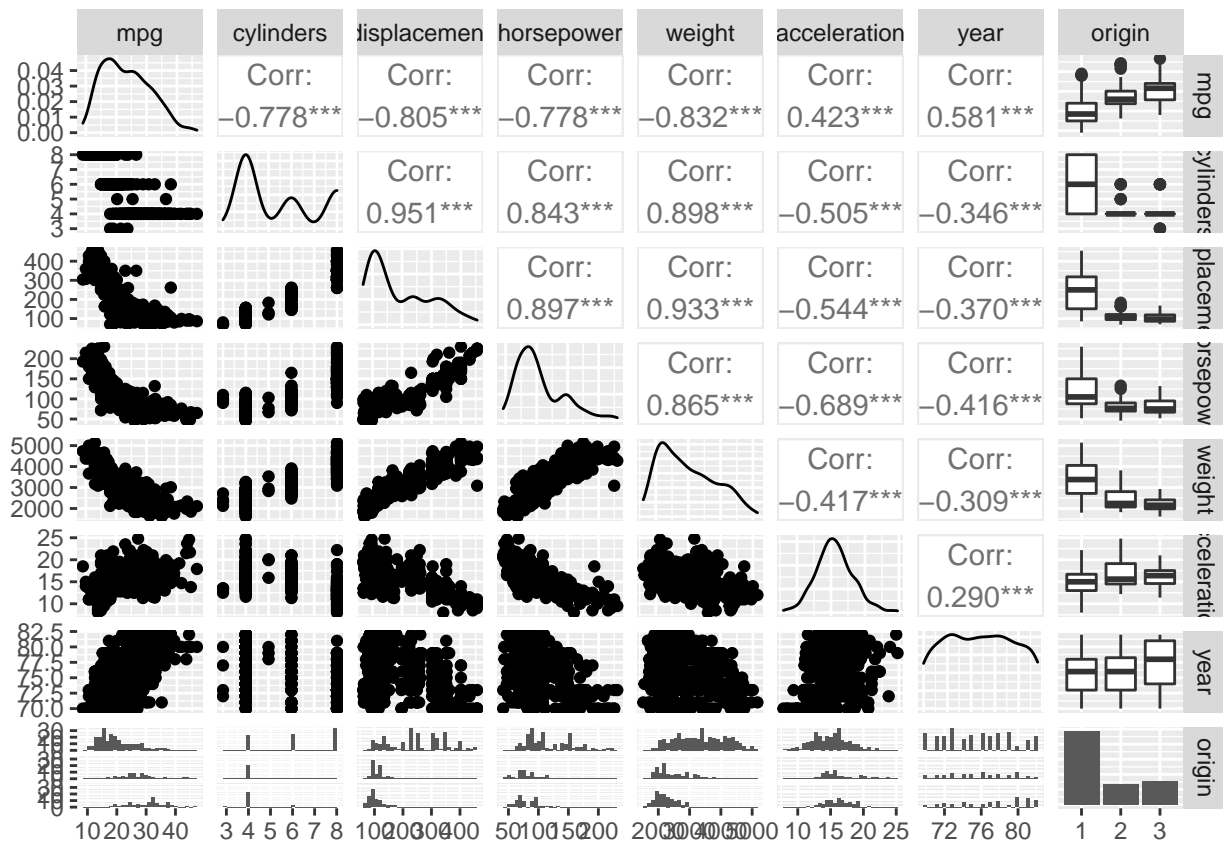
we continue to do any analyses we first need to convert `origin` into a factor variable (a synonymous for “qualitative predictor”):

```
Auto$origin = factor(Auto$origin)
```

a)

Use the function `ggpairs()` from `GGally` package to produce a scatterplot matrix which includes all of the variables in the data set.

```
library(GGally)
ggpairs(Auto)
```



b)

Compute the correlation matrix between the variables. You will need to remove the factor covariate `origin`, because this is no longer a continuous variable.

```
variables <- Auto[-c(8)] # Remove 'origin' from the data set.
Sigma <- cor(variables)
Sigma
```

```
#>           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
#>               acceleration      year
#> mpg                0.4233285  0.5805410
#> cylinders          -0.5046834 -0.3456474
#> displacement      -0.5438005 -0.3698552
#> horsepower        -0.6891955 -0.4163615
#> weight            -0.4168392 -0.3091199
#> acceleration      1.0000000  0.2903161
#> year              0.2903161  1.0000000
```

c)

Use the `lm()` function to perform a multiple linear regression with `mpg` (miles per gallon, a measure for fuel consumption) as the response and all other variables (except `name`) as the predictors. Use the `summary()` function to print the results. Comment on the output. In particular:

- i. Is there a relationship between the predictors and the response?
- ii. Is there evidence that the weight of a car influences `mpg`? Interpret the regression coefficient β_{weight} (what happens if a car weights 1000kg more, for example?).
- iii. What does the coefficient for the year variable suggest?

```
mreg <- lm(mpg ~ ., data = Auto)
summary(mreg)
```

```
#>
#> Call:
#> lm(formula = mpg ~ ., data = Auto)
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -9.0095 -2.0785 -0.0982  1.9856 13.3608
#>
#> Coefficients:
#>              Estimate Std. Error t value Pr(>|t|)
#> (Intercept) -1.795e+01  4.677e+00  -3.839 0.000145 ***
#> cylinders    -4.897e-01  3.212e-01  -1.524 0.128215
#> displacement  2.398e-02  7.653e-03   3.133 0.001863 **
#> horsepower   -1.818e-02  1.371e-02  -1.326 0.185488
#> weight       -6.710e-03  6.551e-04 -10.243 < 2e-16 ***
#> acceleration  7.910e-02  9.822e-02   0.805 0.421101
#> year          7.770e-01  5.178e-02  15.005 < 2e-16 ***
#> origin2       2.630e+00  5.664e-01   4.643 4.72e-06 ***
#> origin3       2.853e+00  5.527e-01   5.162 3.93e-07 ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 3.307 on 383 degrees of freedom
#> Multiple R-squared:  0.8242, Adjusted R-squared:  0.8205
#> F-statistic: 224.5 on 8 and 383 DF, p-value: < 2.2e-16
```

Comments on the output:

- i) Several of the p-values of the coefficients related to the predictors are significant. Moreover, the F-statistic is 224.5, with a p-value of less than 2.2e-6, which signals that there is a relationship between

the predictors and the response.

- ii) The p-value of the coefficient β_{weight} is $2\text{e}-16$ which could be evidence that the weight of a car influences **mpg**. The interpretation of the coefficient is that the **mpg** changes, on average, by $\beta_{\text{weight}} = 6.710\text{e}-3$ for every one-valued increase in the weight of the car, given that all the other predictors are fixed. This means that, e.g., if a car weighs 1000kg, the **mpg** is estimated to be reduced by 6.710\$.
- iii) The coefficient for the year variable suggests that **mpg** is increased by 0.770 for each increase in model year of the car.

d)

Look again at the regression output from question c). Now we want to test whether the **origin** variable is important. How does this work for a factor variable with more than only two levels?

We construct dummy variables such that the coefficients can be estimated in the regression. If we have k levels in the factor variable, we construct $k - 1$ dummy variables. In this way, a baseline level is made, and each coefficient tells something about the difference in the response with respect to each of the other levels in the factor variable. More details to come after learning more of the theory behind it.

Origin2 and Origin3 are two of the coefficients that came out of R. I am not quite sure what these mean, but I am guessing that they are the two differences in levels I was talking about above. And then the baseline is the level 1, which is given by the intercept? In this case we can see that both the estimations of European (2) and Japanese (3) are significant (from their p-values) and that they give positive values of the **mpg** when added to the intercept (the baseline, which is American (1)), since they represent the differences from the baseline. This seems to be correct! Have a look at the R-block below.

```
contrasts(Auto$origin)
```

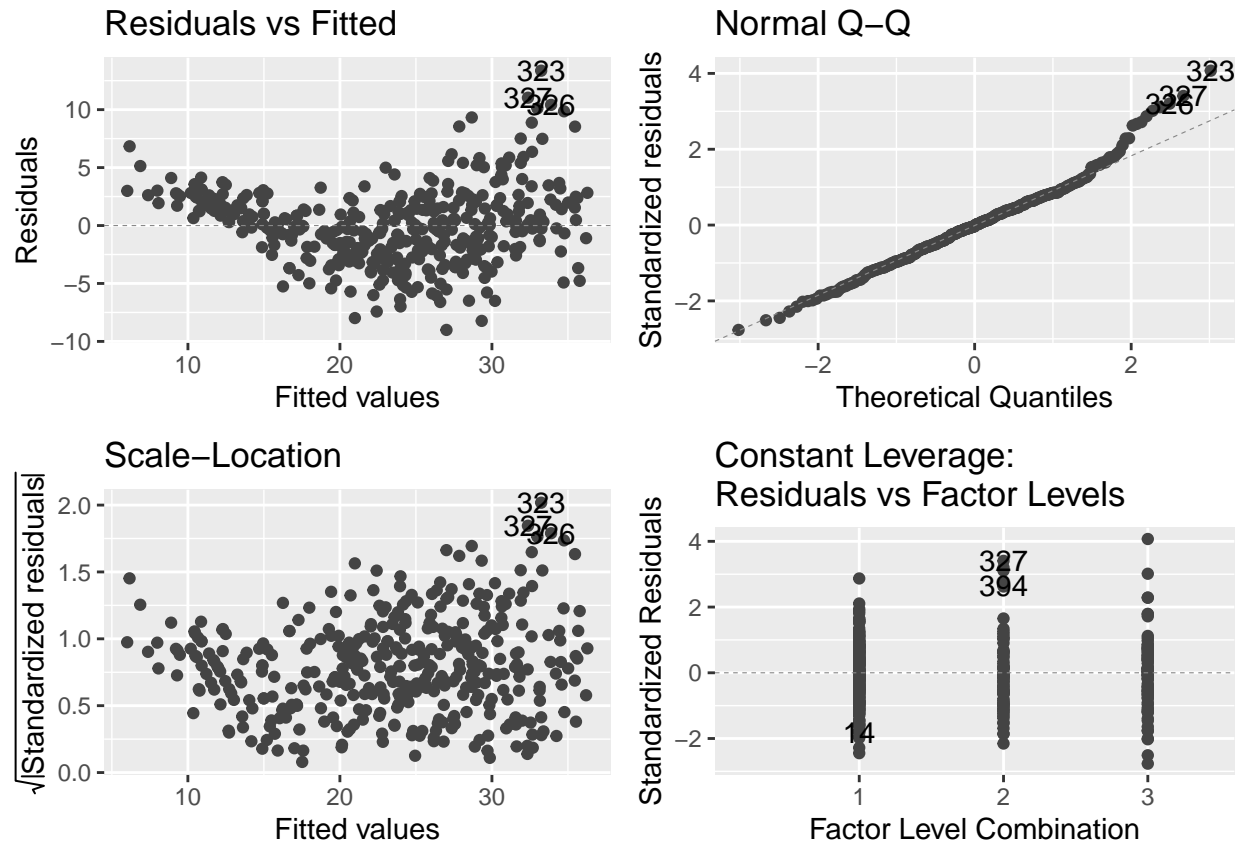
```
#>  2 3  
#> 1 0 0  
#> 2 1 0  
#> 3 0 1
```

As one can see, the dummy variables 2 and 3 were made automatically by R, with 1 as a baseline, as predicted. This means that the dummy variable 'Origin2' takes on the value 1 if the the origin is 2 (European) and zero otherwise. Moreover, the dummy variable 'Origin3' takes on the value 1 if the origin is 3 (Japanese) and zero otherwise. The baseline, 1 (American), corresponds to when both dummy variables are zero-valued.

e)

Use the `autoplot()` function from the `ggfortify` package to produce diagnostic plots of the linear regression fit by setting `smooth.colour = NA`, as sometimes the smoothed line can be misleading. Comment on any problems you see with the fit. Do the residual plots suggest any unusually large outliers? Does the leverage plot identify any observations with unusually high leverage?

```
library(ggfortify)  
autoplot(mreg, smooth.colour = NA)
```



Comments on the plots:

From the residual plot in the upper left it looks like the residuals form a pattern closer to a quadratic, which suggests that there might be some problems with the fit. I cannot, however, identify any unusually large outliers in this plot.

Not sure how the rest of the plots, including the leverage plots should be interpreted (yet).

f)

For beginners, it can be difficult to decide whether a certain QQ plot looks “good” or “bad”, because we only look at it and do not test anything. A way to get a feeling for how “bad” a QQ plot may look, even when the normality assumption is perfectly ok, we can use simulations: We can simply draw from the normal distribution and plot the QQ plot. Use the following code to repeat this six times:

```
set.seed(2332)
n = 100
par(mfrow = c(2, 3))
for (i in 1:6) {
  sim = rnorm(n)
  qqnorm(sim, pch = 1, frame = FALSE)
  qqline(sim, col = "blue", lwd = 1)
}
```

The conclusion is that the Q-Q-plot looks good?

g)

Let us look at interactions. These can be included via the `*` or `:` symbols in the linear predictor of the regression function (see Section 3.6.4 in the course book).

Fit another model for `mpg`, including only `displacement`, `weight`, `year` and `origin` as predictors, plus an interaction between `year` and `origin` (interactions can be included as `year*origin`; this adds the main effects and the interaction at once). Is there evidence that the interactions term is relevant? Give an interpretation of the result.

```
interactions.fit <- lm(mpg ~ displacement + weight + year * origin, data = Auto)
summary(interactions.fit)
```

```
#>
#> Call:
#> lm(formula = mpg ~ displacement + weight + year * origin, data = Auto)
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -8.7710 -2.0204 -0.0207  1.7045 13.0017
#>
#> Coefficients:
#>              Estimate Std. Error t value Pr(>|t|)
#> (Intercept) -5.117e+00  5.259e+00  -0.973  0.331220
#> displacement  4.803e-03  5.032e-03   0.955  0.340420
#> weight       -6.685e-03  5.543e-04 -12.060 < 2e-16 ***
#> year          6.152e-01  6.614e-02   9.302 < 2e-16 ***
#> origin2      -3.735e+01  1.026e+01  -3.642  0.000307 ***
#> origin3      -2.532e+01  9.441e+00  -2.682  0.007631 **
#> year:origin2  5.187e-01  1.342e-01   3.865  0.000130 ***
#> year:origin3  3.564e-01  1.213e-01   2.937  0.003514 **
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 3.257 on 384 degrees of freedom
#> Multiple R-squared:  0.829, Adjusted R-squared:  0.8259
#> F-statistic: 265.9 on 7 and 384 DF, p-value: < 2.2e-16
```

```
fit.without.interactions <- lm(mpg ~ displacement + weight + year + origin,
                                data = Auto)
summary(fit.without.interactions)
```

```
#>
#> Call:
#> lm(formula = mpg ~ displacement + weight + year + origin, data = Auto)
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -9.2633 -2.1572  0.0205  1.8226 13.5061
#>
#> Coefficients:
#>              Estimate Std. Error t value Pr(>|t|)
#> (Intercept) -1.959e+01  4.064e+00  -4.822  2.05e-06 ***
#> displacement  9.321e-03  5.001e-03   1.864   0.0631 .
#> weight       -6.820e-03  5.637e-04 -12.099 < 2e-16 ***
#> year          7.980e-01  5.081e-02  15.705 < 2e-16 ***
```

```
#> origin2      2.383e+00  5.606e-01  4.251 2.67e-05 ***
#> origin3      2.438e+00  5.309e-01  4.592 5.94e-06 ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 3.327 on 386 degrees of freedom
#> Multiple R-squared:  0.8206, Adjusted R-squared:  0.8183
#> F-statistic: 353.2 on 5 and 386 DF,  p-value: < 2.2e-16
```

The R-squared is slightly bigger than the value obtained when only fitting the main effects, while the F-statistic is slightly smaller. However, both all the coefficients in the first fit above are significant, including the interaction terms, which might signal that the interaction terms are relevant? However, based on the small differences in R-squared and F-statistic between the two fits, I would say that there is no clear evidence that the interaction terms are relevant.

h)

Try a few different transformations of the variables, such as $\log(X)$, \sqrt{X} , X^2 . See Section 3.6.5 in the course book for how to do this. Perhaps you manage to improve the residual plots that you got in e)? Comment on your findings.

```
# log-transformation.

# sqrt-transformation.

# x^2-transformation.
```

Problem 2

a)

A core finding for the least-squares estimator $\hat{\beta}$ of linear regression models is

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y},$$

with $\hat{\beta} \sim N_p(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$.

- Show that $\hat{\beta}$ has this distribution with the given mean and covariance matrix.
- What do you need to assume to get to this result?
- What does this imply for the distribution of the j th element of $\hat{\beta}$?
- In particular, how can we calculate the variance of $\hat{\beta}_j$?

b)

What is the interpretation of a 95% confidence interval? Hint: repeat experiment (on Y), on average how many CIs cover the true β_j ? The following code shows an interpretation of a 95% confidence interval. Study and fill in the code where is needed

- Model: $Y = 1 + 3X + \varepsilon$, with $\varepsilon \sim N(0, 1)$.

```
beta0 = ...
beta1 = ...
true_beta = c(beta0, beta1) # vector of model coefficients
true_sd = 1 # choosing true sd
X = runif(100, 0, 1) # simulate the predictor variable X
Xmat = model.matrix(~X, data = data.frame(X)) # create design matrix
```

```

ci_int = ci_x = 0 # Counts how many times the true value is within the confidence interval
nsim = 1000
for (i in 1:nsim) {
  y = rnorm(n = 100, mean = Xmat %*% true_beta, sd = rep(true_sd, 100))
  mod = lm(y ~ x, data = data.frame(y = y, x = X))
  ci = confint(mod)
  ci_int[i] = ifelse(..., 1, 0) # if true value of beta0 is within the CI then 1 else 0
  ci_x[i] = ifelse(..., 1, 0) # if true value of beta_1 is within the CI then 1 else 0
}
c(mean(ci_int), mean(ci_x))

```

c)

What is the interpretation of a 95% prediction interval? Hint: repeat experiment (on Y) for a given \mathbf{x}_0 . Write R code that shows the interpretation of a 95% PI. Hint: In order to produce the PIs use the data point $x_0 = 0.4$. Furthermore you may use a similar code structure as in b).

d)

Construct a 95% CI for $\mathbf{x}_0^T \beta$. Explain what is the connections between a CI for β_j , a CI for $\mathbf{x}_0^T \beta$ and a PI for Y at \mathbf{x}_0 .

e)

Explain the difference between *error* and *residual*. What are the properties of the raw residuals? Why don't we want to use the raw residuals for model check? What is our solution to this?