

```
set.seed(42)
```

## Arbeitsblatt 5

### Aufgabe 1

Erinnerung. Wir nehmen die Gauss-Markov-Annahmen an (Fehlerterme haben Erwartungswert 0, gleiche Varianz und sind unkorreliert) und folglich gilt:

$$\text{Cov}(\epsilon_i, \epsilon_j) = 0 \quad \forall i \neq j$$

$$\text{E}(\epsilon_i) = \theta$$

$$\text{V}(\epsilon_i) = \sigma^2$$

$$Y_i = X_i b + \epsilon_i$$

$$\hat{Y}_i = (X\hat{b})_i$$

$$\hat{b} = (X^T X)^{-1} X^T Y$$

$$\hat{b}_i = ((X^T X)^{-1} X^T Y)_i$$

$$H = X(X^T X)^{-1} X^T$$

$$H = H \cdot H$$

$$H = H^T$$

$X_i$  ist hier eine Zeile der Designmatrix und  $H_i$  eine Zeile der Hatmatrix. Zuerst stellen wir fest, dass  $X_i b$  und  $\epsilon_i$  unabhängig voneinander sind und  $X_i b$  varianzlos ist. Dann erhalten wir

$$\text{V}(Y_i) = \text{V}(X_i b + \epsilon_i) = \text{V}(X_i b) + \text{V}(\epsilon_i) = \sigma^2.$$

Ferner können wir berechnen

$$\begin{aligned} \text{V}(\hat{Y}_i) &= \text{V}((X\hat{b})_i) = \text{V}((HY)_i) = H_i \text{V}(Y) H_i^T = H_i \text{V}(Xb + \epsilon) H_i^T \\ &= H_i \text{diag}(\sigma^2) H_i^T = h_{ii} \sigma^2. \end{aligned}$$

Ähnlich können wir berechnen

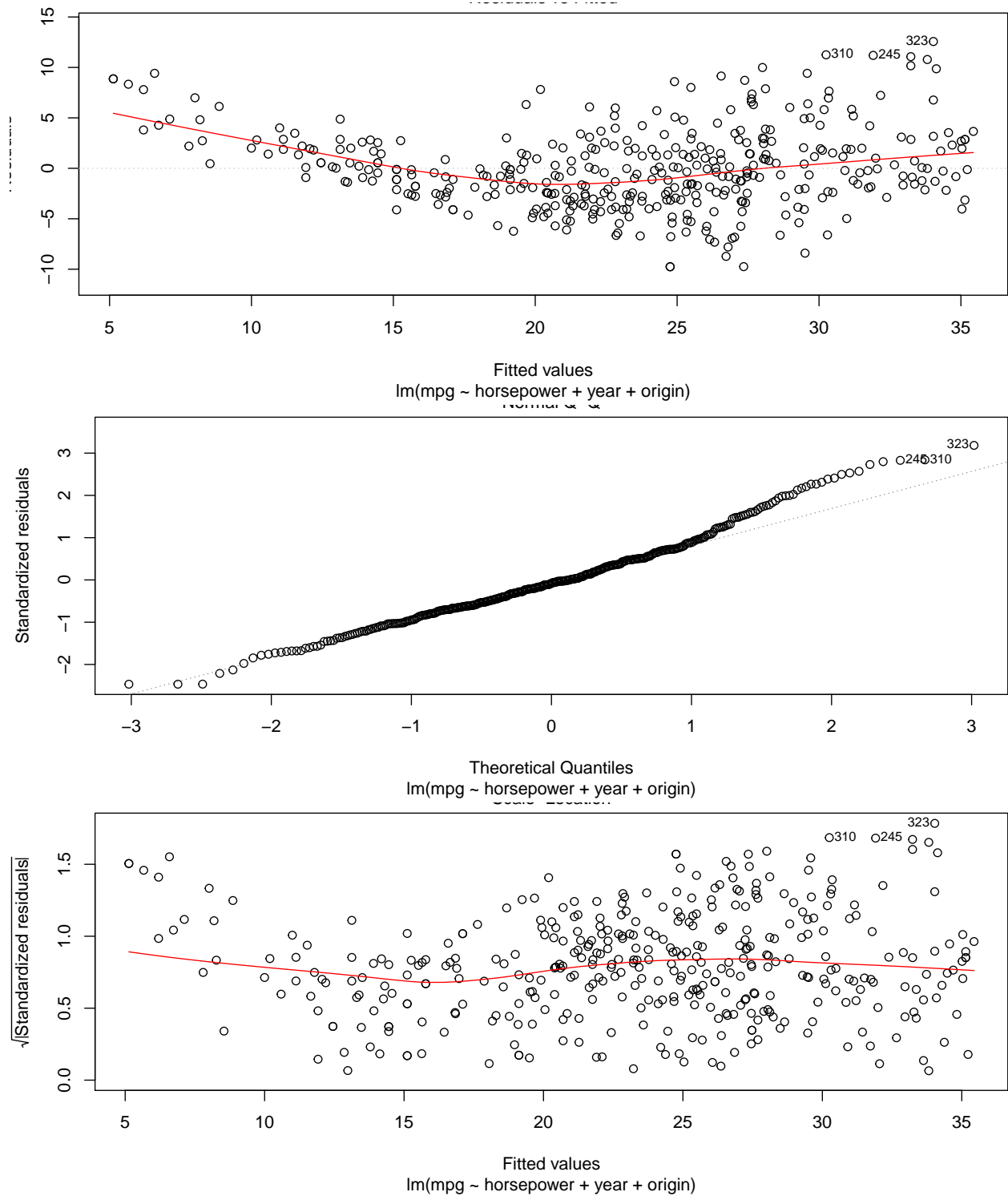
$$\begin{aligned} \text{Cov}(Y_i, \hat{Y}_i) &= \text{Cov}(Y_i, H_i Y) \\ &= H_i \text{Cov}(x_i b + \epsilon_i, (Xb + \epsilon)_i) H_i^T \\ &= h_{ii} \sigma^2. \end{aligned}$$

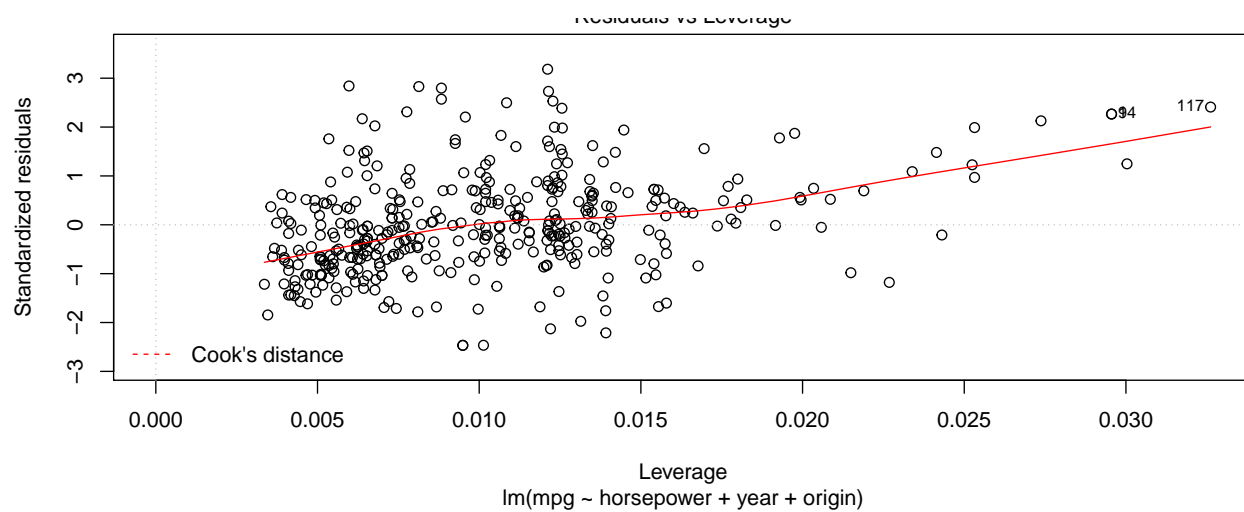
Und schließlich erhalten wir durch Einsetzen

$$\begin{aligned} \text{Cor}(Y_i, \hat{Y}_i) &= \frac{\text{Cov}(Y_i, \hat{Y}_i)}{\sqrt{\text{V}(Y_i) \text{V}(\hat{Y}_i)}} \\ &= \frac{h_{ii} \sigma^2}{\sqrt{\sigma^2 h_{ii} \sigma^2}} \\ &= \sqrt{h_{ii}}. \end{aligned}$$

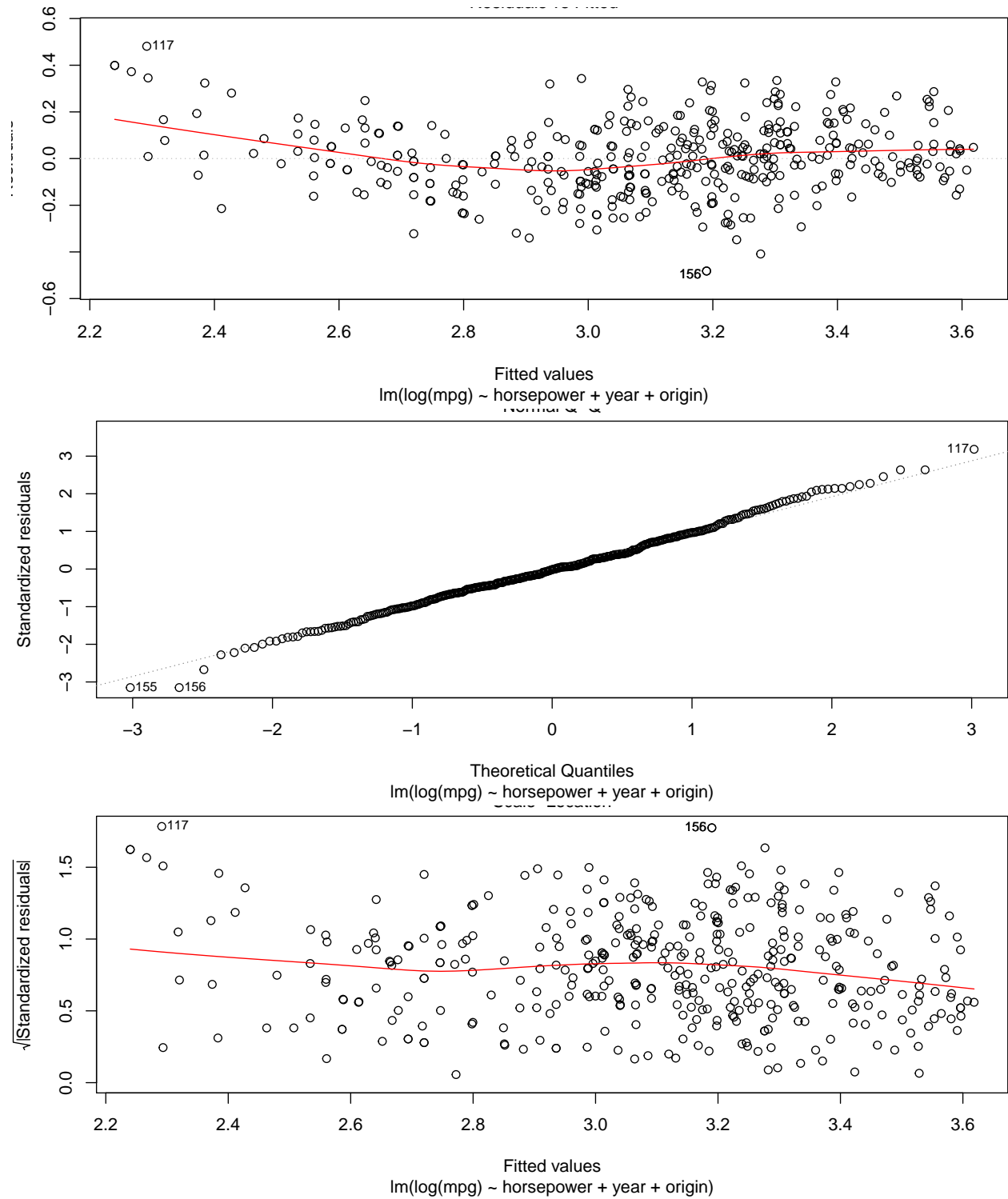
## Aufgabe 2

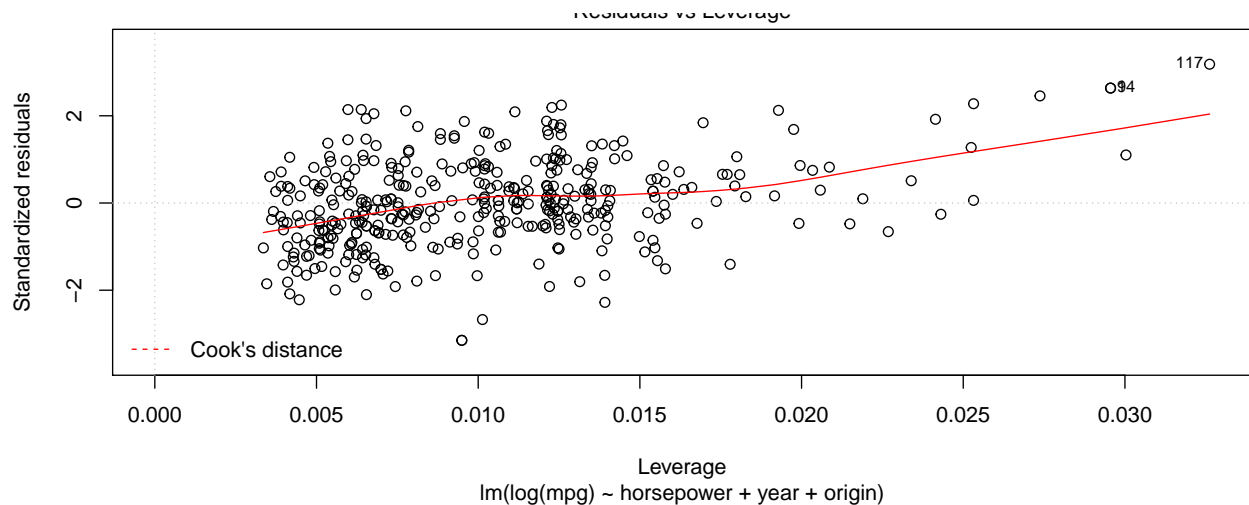
```
library("ISLR")
model_mpg <- lm(mpg~horsepower+year+origin, data = Auto)
plot(model_mpg)
```





```
model_log_mpg <- lm(log(mpg)~horsepower+year+origin, data = Auto)
plot(model_log_mpg)
```





a)

Wenn wir uns die Grafik “Residuals vs Fitted” ansehen, sehen wir, dass die Daten kein offensichtliches, eindeutiges Muster haben. Obwohl die Kurven leicht gewölbt sind, haben sie gleichmäßige Residuen um die horizontale Linie herum verteilt, ohne ein ausgeprägtes Muster. Dies ist ein guter Hinweis darauf, dass es sich nicht um eine nichtlineare Beziehung handelt.

Die Residuen sollten normalverteilt sein, und der Q-Q Plot ist ein gutes Instrument, um dies zu überprüfen. Für beide Modelle zeigt der Q-Q-Plot eine ziemlich gute Anpassung auf die Linie. Nach oben hin weichen die Punkte für das mpg-Modell deutlicher von der Linie ab, als im log(mpg)-Modell.

Die “Scale-Location”-Plots testen die lineare Regressionsannahme der gleichen Varianz (Homoscedastizität), d.h. dass die Residuen die gleiche Varianz entlang der Regressionslinie aufweisen.

In unseren beiden Modellen sind die Residuen relativ gut über und unter einer ziemlich horizontalen Linie verteilt.

Die “Residuals vs Leverage”-Plots können verwendet werden, um einflussreiche Datenpunkte im Datensatz zu finden. Ein einflussreicher Datenpunkt ist einer, der, wenn er entfernt wird, das Modell beeinflusst, so dass seine Einbeziehung oder sein Ausschluss in Betracht gezogen werden sollte.

Ein einflussreicher Datenpunkt kann ein Ausreißer sein oder auch nicht, und der Zweck dieser Grafik ist es, Fälle zu identifizieren, die einen hohen Einfluss auf das Modell haben. Ausreißer neigen dazu, Leverage und damit Einfluss auf das Modell auszuüben.

Ein einflussreicher Fall erscheint oben rechts oder unten links in der Grafik innerhalb einer roten Linie, die Cook’s Distance markiert.

In beiden Modellen gibt es keine offensichtlichen Ausreißer (laut Cook’s Distance).

Abschließend würden wir das log(mpg)-Modell bevorzugen, obwohl es kaum Unterschiede in den Residuen gibt. Alleine der Q-Q-Plot des log(mpg)-Modells deutet auf eine bessere Anpassung der Residuen auf die Normalverteilung hin als im mpg-Modell.

b)

Der Punkt 116 ist in beiden Modellen des Öfteren markiert. Schauen wir uns diesen näher an:

```
pt_116 <- Auto[116,]
summary <- summary(Auto)
c(pt_116)
#> $mpg
```

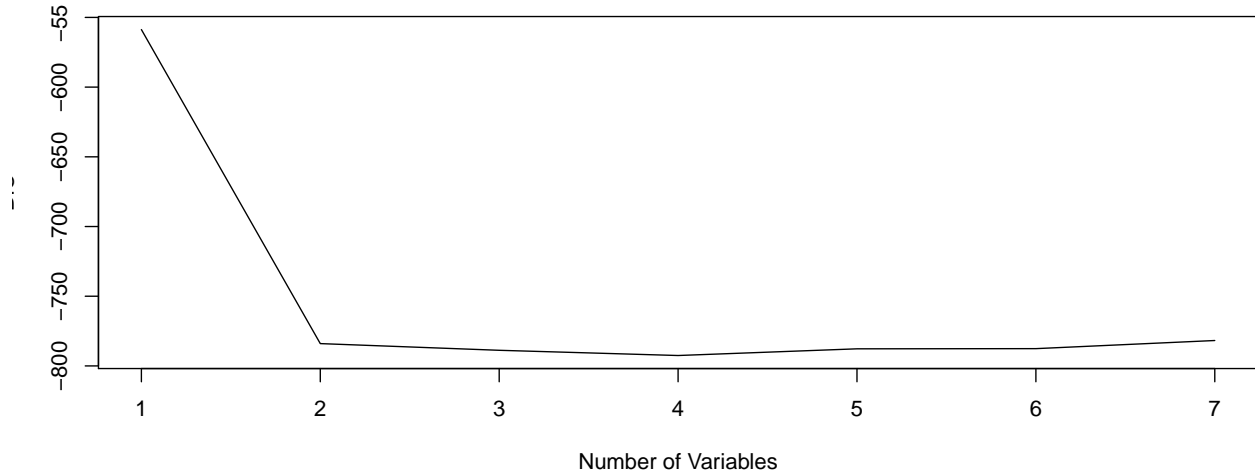
```
#> [1] 16
#>
#> $cylinders
#> [1] 8
#>
#> $displacement
#> [1] 400
#>
#> $horsepower
#> [1] 230
#>
#> $weight
#> [1] 4278
#>
#> $acceleration
#> [1] 9.5
#>
#> $year
#> [1] 73
#>
#> $origin
#> [1] 1
#>
#> $name
#> [1] pontiac grand prix
#> 304 Levels: amc ambassador brougham ... vw rabbit custom
summary
#>      mpg      cylinders  displacement  horsepower
#>  Min.   : 9.00   Min.    :3.000   Min.    : 68.0   Min.    : 46.0
#> 1st Qu.:17.00   1st Qu.:4.000   1st Qu.:105.0   1st Qu.: 75.0
#>  Median :22.75   Median :4.000   Median :151.0   Median : 93.5
#>  Mean    :23.45   Mean    :5.472   Mean    :194.4   Mean    :104.5
#> 3rd Qu.:29.00   3rd Qu.:8.000   3rd Qu.:275.8   3rd Qu.:126.0
#>  Max.    :46.60   Max.    :8.000   Max.    :455.0   Max.    :230.0
#>
#>      weight  acceleration      year      origin
#>  Min.   :1613   Min.    : 8.00   Min.    :70.00   Min.    :1.000
#> 1st Qu.:2225   1st Qu.:13.78   1st Qu.:73.00   1st Qu.:1.000
#>  Median :2804   Median :15.50   Median :76.00   Median :1.000
#>  Mean    :2978   Mean    :15.54   Mean    :75.98   Mean    :1.577
#> 3rd Qu.:3615   3rd Qu.:17.02   3rd Qu.:79.00   3rd Qu.:2.000
#>  Max.    :5140   Max.    :24.80   Max.    :82.00   Max.    :3.000
#>
#>      name
#> amc matador      : 5
#> ford pinto       : 5
#> toyota corolla   : 5
#> amc gremlin      : 4
#> amc hornet       : 4
#> chevrolet chevette: 4
#> (Other)          :365
```

Hier wird deutlich, dass der Punkt 116 ein Extremfall in unserem Datensatz ist.

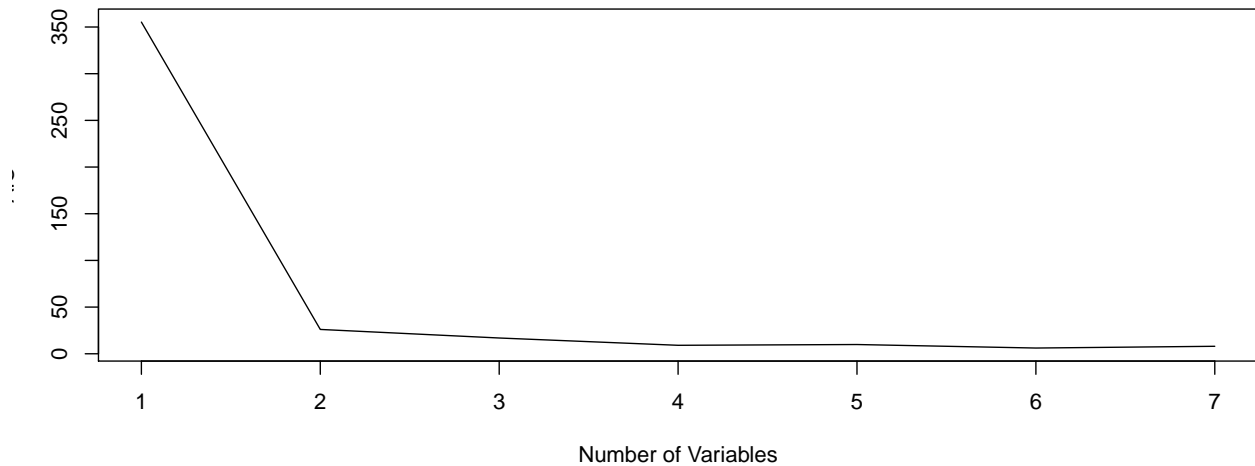
c)

```
library(leaps)
regfit.full = regsubsets(log(mpg)~cylinders+displacement+horsepower
                        +weight+acceleration+year+origin, data=Auto)
reg.summary = summary(regfit.full)
```

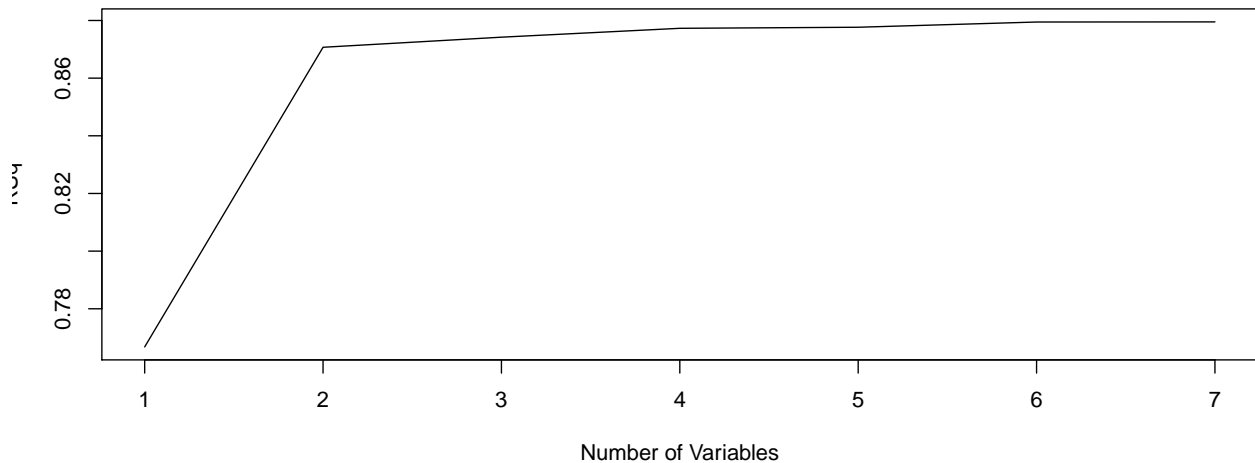
```
plot(reg.summary$bic, xlab="Number of Variables", ylab="BIC", type="l")
```



```
plot(reg.summary$cp, xlab="Number of Variables", ylab="AIC", type="l")
```



```
plot(reg.summary$rsq, xlab="Number of Variables", ylab="RSq", type="l")
```



BIC:

- BIC wird für das Modell mit nur 4 Kovariaten minimal.

AIC:

- Nach AIC ist die Entscheidungsfindung nicht ganz so eindeutig. Auch hier ist bei 4 Kovariaten der AIC sehr gering, jedoch steigt der AIC bei wachsender Anzahl an Kovariaten nicht mehr an, sinkt aber auch nicht deutlich weiter.

$R^2$  :

- Ein hohes  $R^2$  spricht für ein gutes Modell. Allerdings gibt es, anders als bei AIC und BIC, keinen Strafterm, welcher die Anzahl an Kovariaten "bestraft". Deshalb steigt das  $R^2$  mit Anzahl der Kovariaten (leicht) an.

Wir würden uns für ein Modell mit 4 Kovariaten entscheiden, da dort BIC und AIC sehr klein sind und  $R^2$  groß ist.

```
coef(regfit.full, 4)
#> (Intercept) horsepower weight year origin
#> 1.6584625886 -0.0010356091 -0.0002512191 0.0295010570 0.0346250328
```

d)

```
library(car)
model_all_log <- lm(log(mpg)~cylinders+displacement+horsepower+
                    weight+acceleration+year+origin, data=Auto)

model_4_log <- lm(log(mpg)~horsepower+weight+year+origin, data=Auto)
vif(model_all_log)
#> cylinders displacement horsepower weight acceleration
#> 10.737535 21.836792 9.943693 10.831260 2.625806
#> year origin
#> 1.244952 1.772386
cat("\n")
vif(model_4_log)
#> horsepower weight year origin
#> 4.462659 4.915735 1.227042 1.546807
```

Schaut man sich die VIF-Werte an, wird deutlich, dass im Modell mit allen Kovariaten definitiv Multikollinearität ein Problem darstellt.



Im Modell, dass wir durch die Subset-Selektion gefunden haben, sind die VIF-Werte geringer. “Pi-Mal-Daumen” ist die Faustregel:

- 1 = unkorreliert
- zwischen 1 und 5 = moderat korreliert
- größer 5 = stark korreliert

Die VIF-Werte für horsepower und weight liegen beide bei circa 5. Nimmt man nun die Kovariate mit dem höchsten VIF-Wert aus dem Modell heraus, ergeben sich folgende VIF-Werte für das neue Modell:

```
model_3_log <- lm(log(mpg)~horsepower+year+origin, data=Auto)
vif(model_3_log)
#> horsepower      year      origin
#>  1.475700  1.209830  1.261445
```

Eindeutig wird der Einfluss von weight durch horsepower bereits gut “erklärt”. Die VIF-Werte liegen nun im annehmbaren Bereich. Das Modell mit horsepower, year und origin scheint das robustere zu sein.

### Aufgabe 3

a)

Likelihood-Ratio-Test:

Schätzstatistik:

$$T = 2 \cdot (LL_{M2} - LL_{M1}) \sim_{H0} \chi_q^2, \quad q = \text{Anzahl zusätzlicher Parameter im komplexeren Modell}$$

Hier:

$$T = 2 \cdot (LL_{M2} - LL_{M1}) \sim_{H0} \chi_1^2 \text{ mit } \alpha = 0.1572992$$

$\Rightarrow X_{p+1}$  wird in das Modell aufgenommen, wenn gilt:

$$T = 2 \cdot (LL_{M2} - LL_{M1}) > 2, \text{ da}$$

$$\chi_{1(1-0.1572992)}^2 = 2$$

AIC-Vergleich:

$$AIC = -2LL + 2k, \quad k = \text{Anzahl Modellparameter}$$

$$AIC_{M1} = -2LL_{M1} + 2 \cdot p$$

$$AIC_{M2} = -2LL_{M2} + 2 \cdot (p + 1)$$

$\Rightarrow X_{p+1}$  wird in das Modell aufgenommen, wenn gilt:

$$AIC_{M1} - AIC_{M2} > 0$$

$$= (-2LL_{M1} + 2p) - (-2LL_{M2} + 2(p + 1)) > 0$$

$$= -2LL_{M1} + 2p + 2LL_{M2} - 2p - 2 > 0$$

$$= -2(LL_{M1} - LL_{M2} - p + p + 1) > 0$$

$$= -2(LL_{M1} - LL_{M2} + 1) > 0$$

$$= 2(-LL_{M1} + LL_{M2} - 1) > 0$$

$$= 2(LL_{M2} - LL_{M1} - 1) > 0$$

$$= 2(LL_{M2} - LL_{M1}) - 2 > 0$$

$$= 2(LL_{M2} - LL_{M1}) > 2$$

b)

Modellvergleich:  $\log(\text{mpg}) \sim \text{acceleration}$  (model\_1) vs  $\log(\text{mpg}) \sim \text{acceleration} + \text{horsepower}$  (model\_2)

```
library(lmtest)
#> Error in library(lmtest): there is no package called 'lmtest'
model_1 = lm(log(mpg)~acceleration, data=Auto)
model_2 = lm(log(mpg)~acceleration+horsepower, data=Auto)
AIC(model_1)
#> [1] 184.1192
AIC(model_2)
#> [1] -223.4539
lmtest::lrtest(model_1, model_2)
#> Error in loadNamespace(name): there is no package called 'lmtest'
```

model\_2 ist sowohl im AIC-Vergleich als auch nach Likelihood-Ratio-Test besser als model\_1.