# Exercise 8

## Tobias Raidl, 11717659

### 2023-12-18

```r
library(dplyr)
```

```
## Warning: Paket 'dplyr' wurde unter R Version 4.3.2 erstellt
```

```
##
## Attache Paket: 'dplyr'
```

```
## Die folgenden Objekte sind maskiert von 'package:stats':
##
##     filter, lag
```

```
## Die folgenden Objekte sind maskiert von 'package:base':
##
##     intersect, setdiff, setequal, union
```

```r
library(ISLR)
```

```
## Warning: Paket 'ISLR' wurde unter R Version 4.3.2 erstellt
```

```r
library(splines)
data(Auto)

df = select(Auto, -name)
df_trans = df
df_trans[, 1] = log(df[, 1])
df_trans[, 3] = log(df[, 3])
df_trans[, 4] = log(df[, 4])
df_trans[, 5] = log(df[, 5])



set.seed(11717659)
sample <- sample(c(TRUE, FALSE), nrow(df), replace = TRUE, prob = c(0.7, 0.3))
train <- df[sample, ]
test <- df[!sample, ]
```
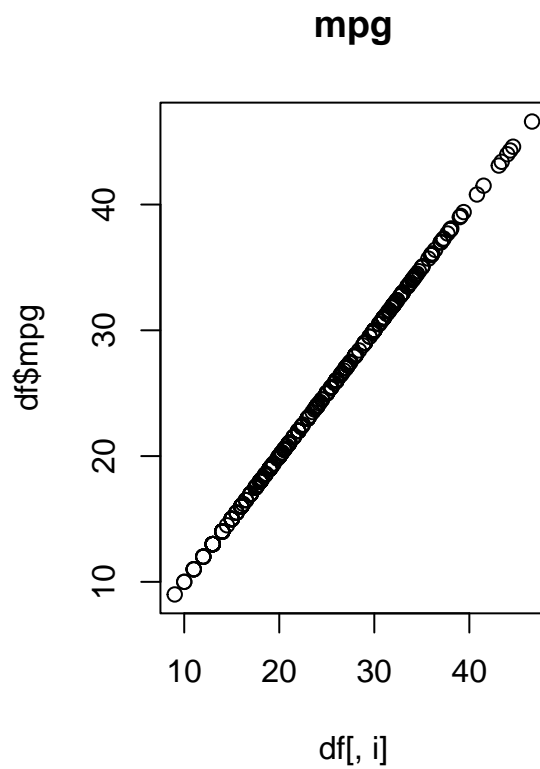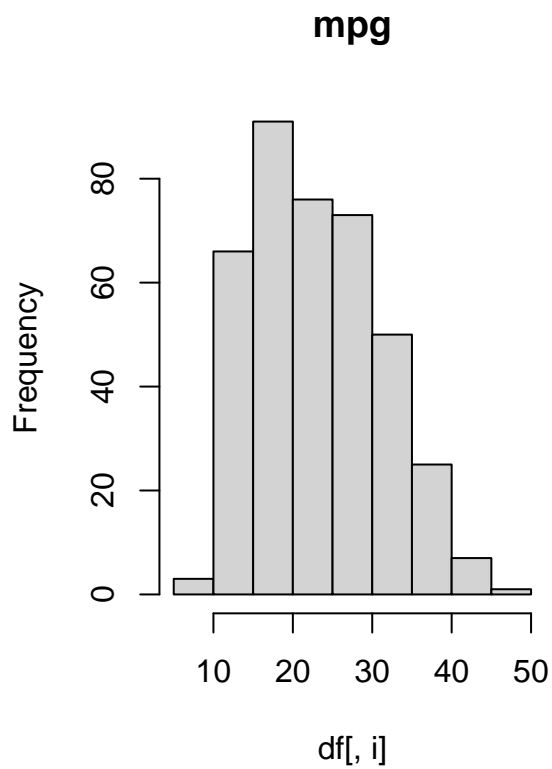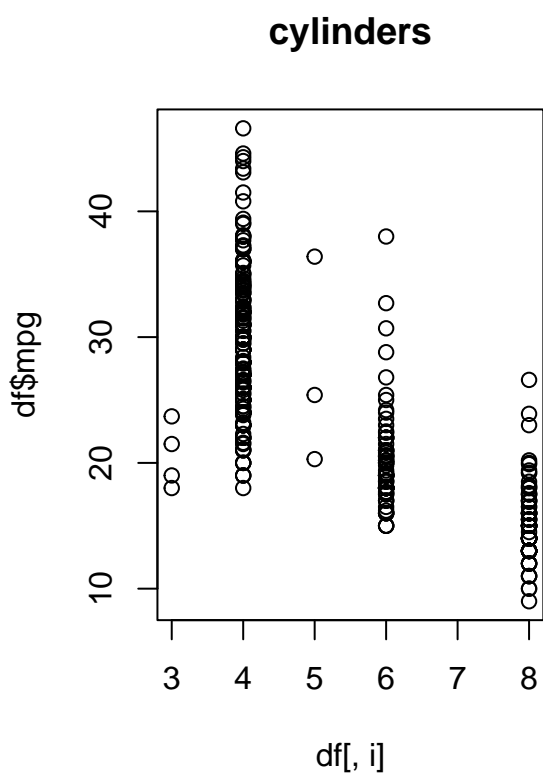
## Initial

```
library(ggplot2)
```
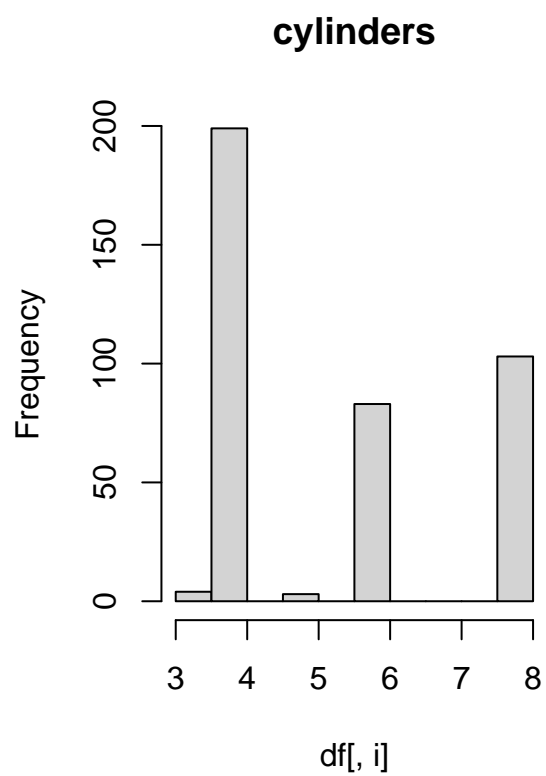
```
## Warning: Paket 'ggplot2' wurde unter R Version 4.3.2 erstellt
```
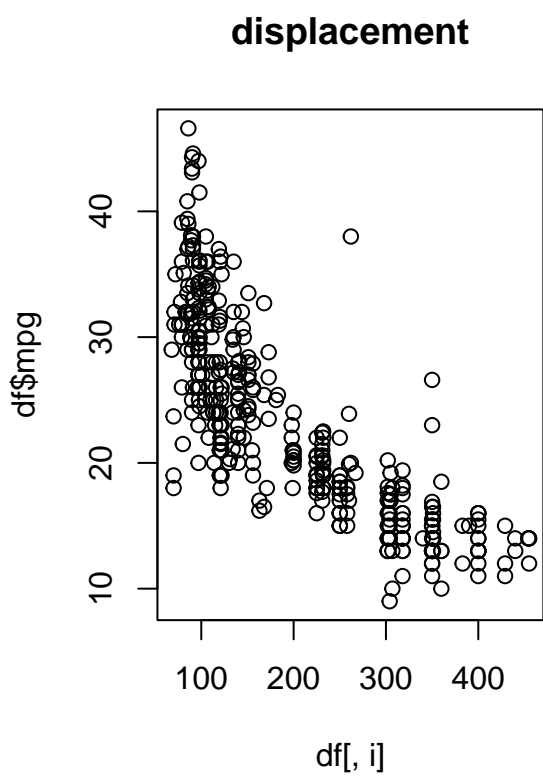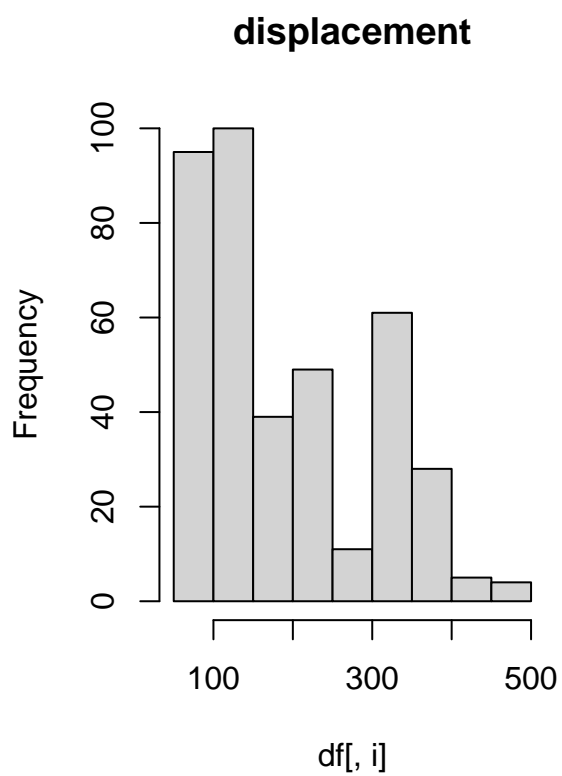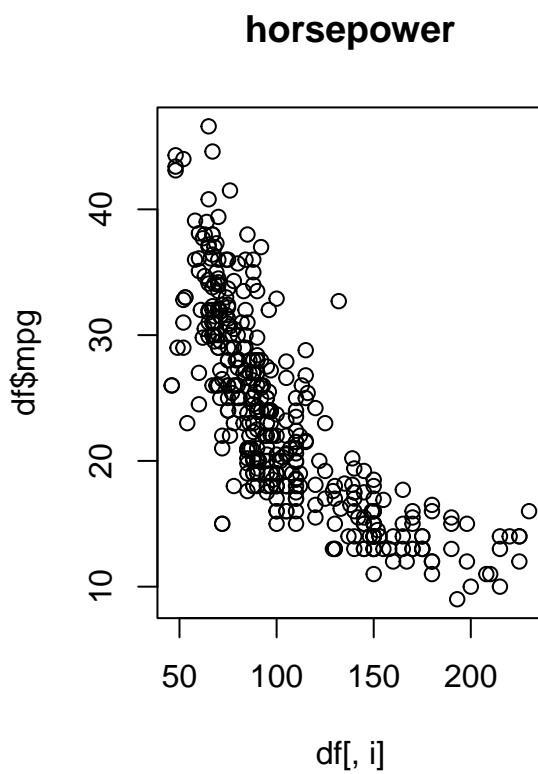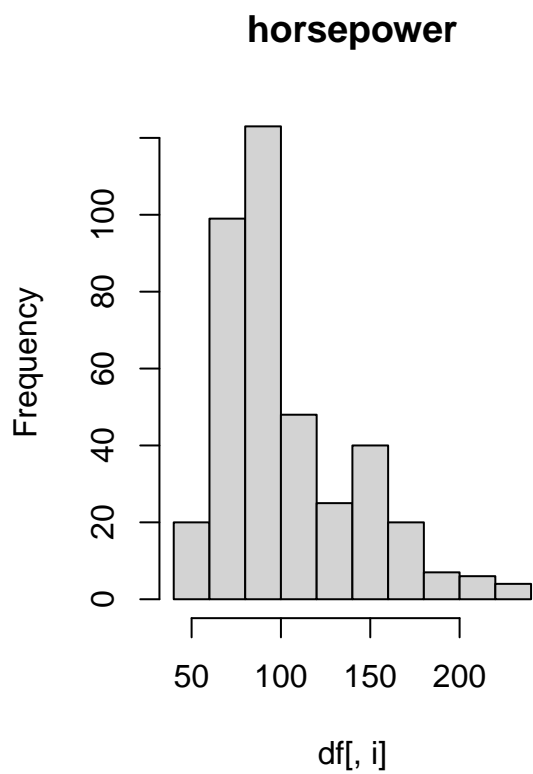
```
library(tidyr)
```

```
## Warning: Paket 'tidyr' wurde unter R Version 4.3.2 erstellt
```
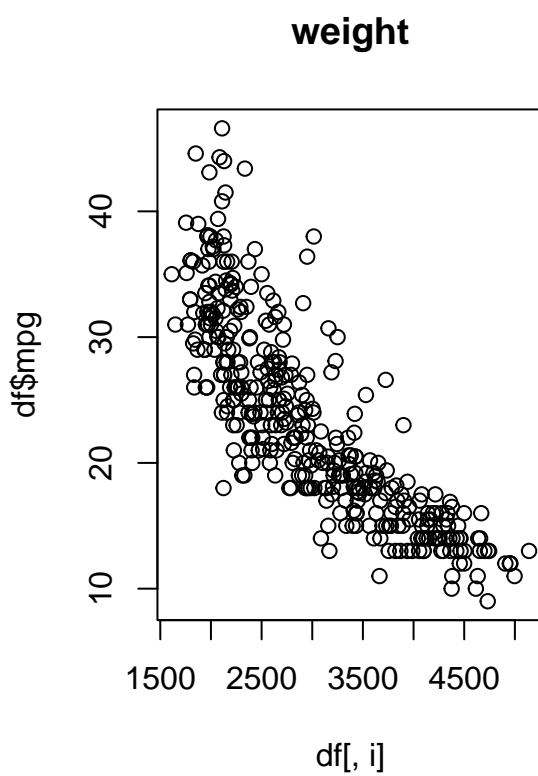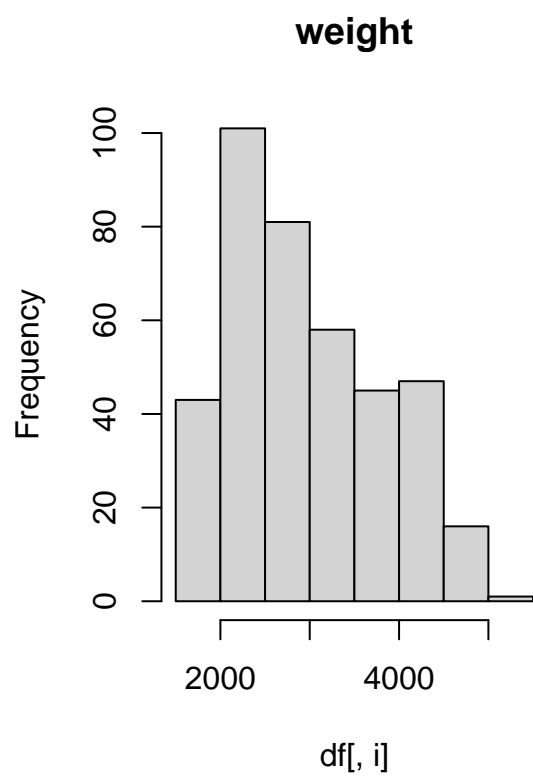
```
for (i in 1:ncol(df)) {
    par(mfrow = c(1, 2))
    hist(df[, i], main = colnames(df)[i])
    plot(df[, i], df$mpg, main = colnames(df)[i])
}
```
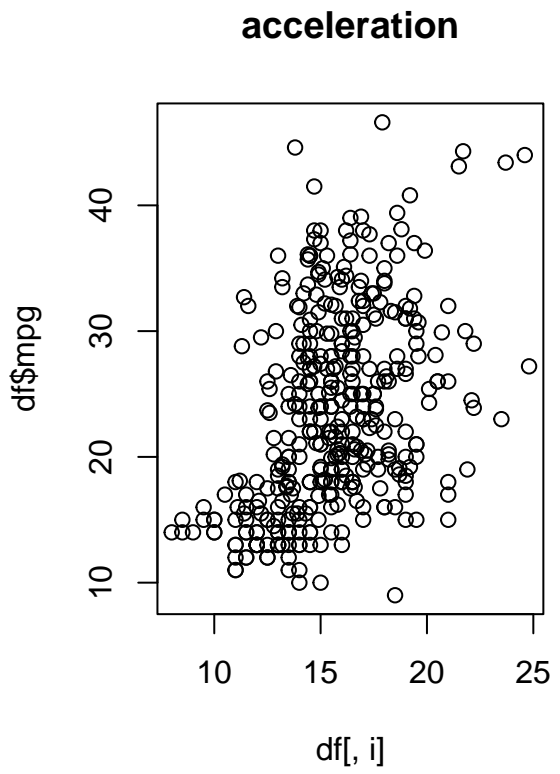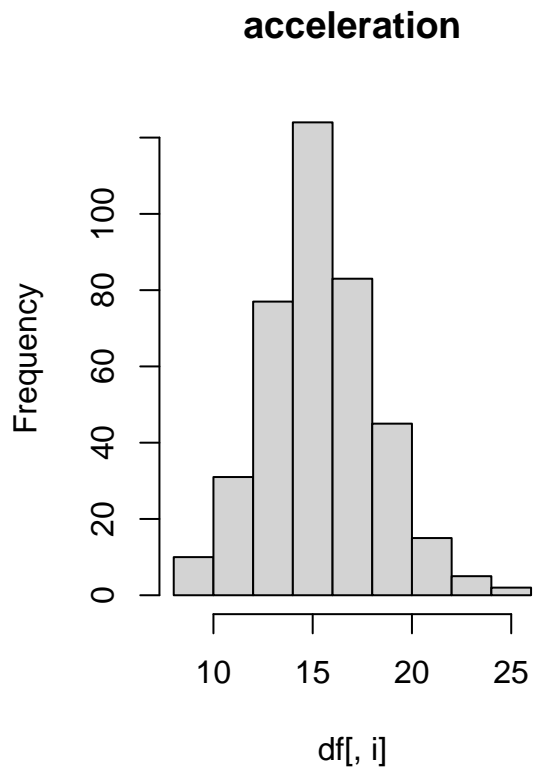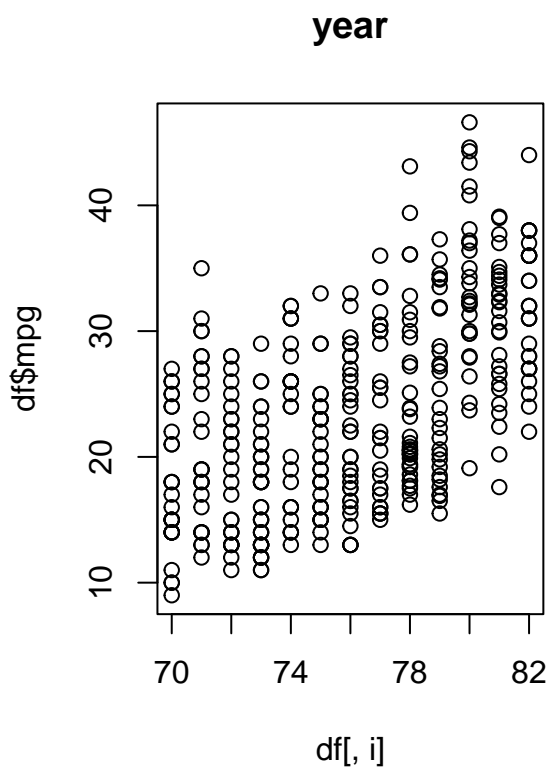
**cylinders**

**cylinders**

**displacement**

**displacement**

**horsepower**



**horsepower**



5

## weight



## weight

**acceleration**



**acceleration**

## year

## year

## origin



**Transformed**

```r
library(ggplot2)
library(tidyr)
for (i in 1:ncol(df_trans)) {
    par(mfrow = c(1, 2))
    hist(df_trans[, i], main = colnames(df_trans)[i])
    plot(df_trans[, i], df_trans$mpg, main = colnames(df_trans)[i])
}
```

## mpg



## mpg

## cylinders



## cylinders

**displacement**



**displacement**

## horsepower



## horsepower

**weight**

**weight**

## acceleration



## acceleration

## RMSE function

```
get_rmse = function(y_pred, y_gt) {
    return(sqrt(mean((y_pred - y_gt)^2)))
}
```

### 1

Cylinders, year and origins should be treated as factors. Using log transformation all variables could be represented linearly. Not using log transform displacement, horsepower, weight, cylinders and year should be represented by splines.

### 2 & 3

```
degfree = 2
lin_model = lm(mpg ~ ., data = train)
spl_model = lm(mpg ~ (ns(displacement, df = degfree) + ns(horsepower, df = degfree) +
    ns(weight, df = degfree) + ns(cylinders, df = degfree) + ns(year, df = degfree)),
    data = train)
summary(spl_model)
```
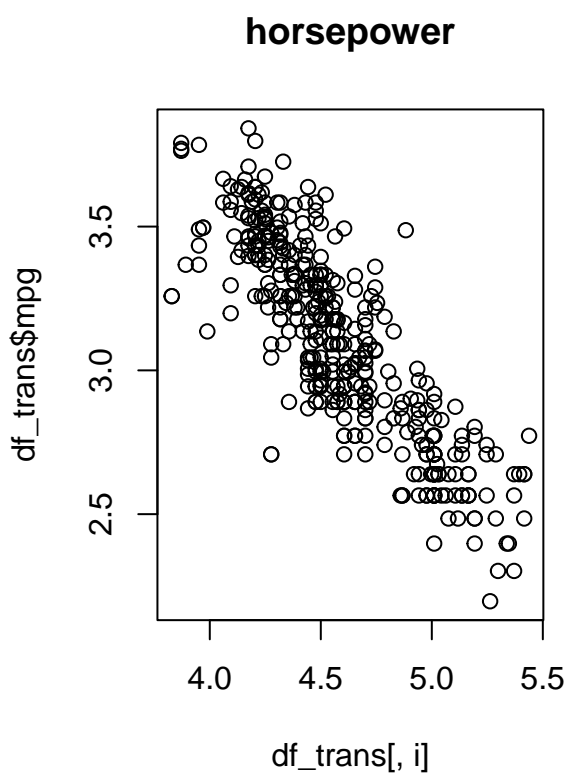
```
##
## Call:
## lm(formula = mpg ~ (ns(displacement, df = degfree) + ns(horsepower,
##     df = degfree) + ns(weight, df = degfree) + ns(cylinders,
##     df = degfree) + ns(year, df = degfree)), data = train)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -7.9958 -1.5826 -0.0154  1.3109 11.8227
##
## Coefficients:
##                              Estimate Std. Error t value Pr(>|t|)
## (Intercept)                   32.8342     1.1943  27.492  < 2e-16 ***
## ns(displacement, df = degfree)1  -7.4566     3.9936  -1.867  0.06300 .
## ns(displacement, df = degfree)2   1.7562     2.6124   0.672  0.50202
## ns(horsepower, df = degfree)1   -13.4381     2.7779  -4.838 2.25e-06 ***
## ns(horsepower, df = degfree)2    -6.6600     2.0873  -3.191  0.00159 **
## ns(weight, df = degfree)1       -17.5720     3.5293  -4.979 1.16e-06 ***
## ns(weight, df = degfree)2        -8.7383     1.8540  -4.713 3.96e-06 ***
## ns(cylinders, df = degfree)1      3.4320     3.1731   1.082  0.28043
## ns(cylinders, df = degfree)2      1.2350     1.5485   0.798  0.42587
## ns(year, df = degfree)1           6.9496     1.1528   6.029 5.60e-09 ***
## ns(year, df = degfree)2           9.0863     0.5933  15.315  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.906 on 262 degrees of freedom
## Multiple R-squared:  0.8707, Adjusted R-squared:  0.8658
## F-statistic: 176.4 on 10 and 262 DF,  p-value: < 2.2e-16
```

## 5 & 6

```
y_pred_lin = predict(lin_model, test)
y_pred_spl = predict(spl_model, test)
get_rmse(y_pred_lin, test$mpg)
```

```
## [1] 3.440262
```

```
get_rmse(y_pred_spl, test$mpg)
```

```
## [1] 2.532837
```

```
preds = data.frame(lin = y_pred_lin, spl = y_pred_spl, gt = test$mpg)
ggplot() + geom_point(data = preds, aes(gt, lin), col = "red") + geom_point(data = preds,
    aes(gt, spl), col = "steelblue") + geom_abline(intercept = 0, slope = 1)
```