# Exercise 4

Tobias Raidl, 11717659

2023-10-30

## Contents

Setup dataset

```r
library(UsingR)
```

```
## Warning: Paket 'UsingR' wurde unter R Version 4.1.3 erstellt
```

```
## Lade nötiges Paket: MASS
```

```
## Lade nötiges Paket: HistData
```

```
## Lade nötiges Paket: Hmisc
```

```
## Warning: Paket 'Hmisc' wurde unter R Version 4.1.3 erstellt
```

```
##
## Attache Paket: 'Hmisc'
```

```
## Die folgenden Objekte sind maskiert von 'package:base':
##
##     format.pval, units
```

```r
library(dplyr)
```

```
## Warning: Paket 'dplyr' wurde unter R Version 4.1.3 erstellt
```

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

## Die folgenden Objekte sind maskiert von 'package:Hmisc':
##
##     src, summarize

## Das folgende Objekt ist maskiert 'package:MASS':
##
##     select

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

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

```r
df = select(fat, -c(case, body.fat.siri, ffweight, density))
X = select(df, -c(body.fat))
y = df$body.fat
```
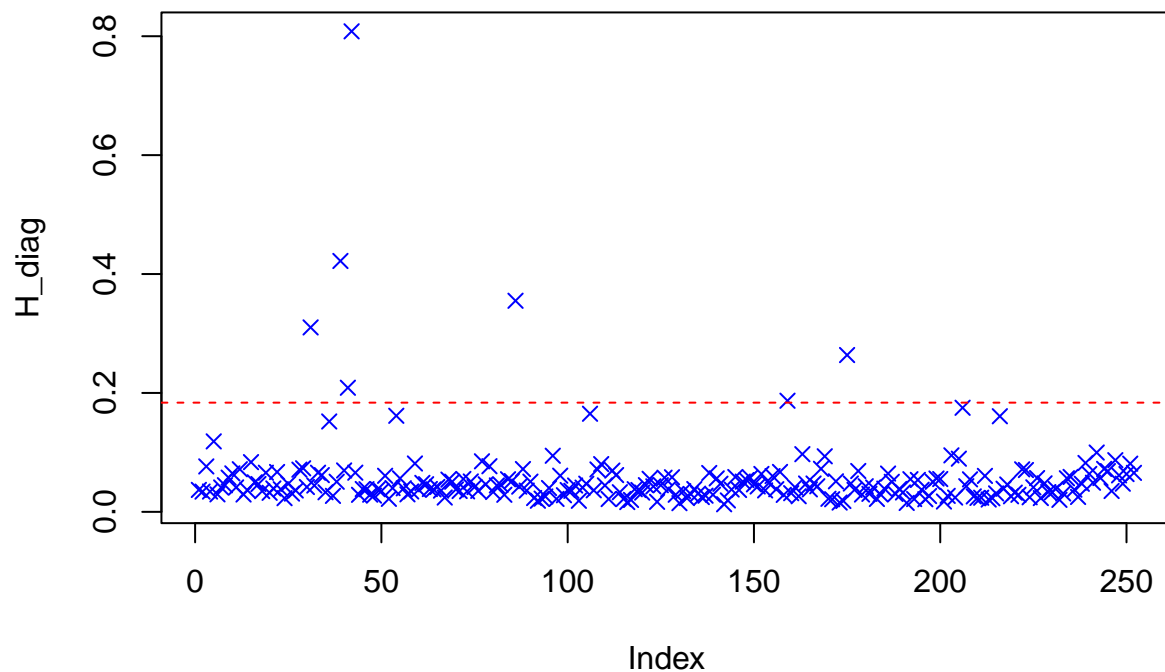
# 1

Investigate if there are leverage points by using ...

## a

... classical diagnostic based on the diagonal elements of the hat matrix

```r
Xm = data.matrix(X)
H = Xm %*% solve(t(Xm) %*% Xm) %*% t(Xm)
H_diag = diag(H)
plot(H_diag, pch = 4, col = "blue")
abline(h = quantile(H_diag, 0.975), lty = "dashed", col = "red")
```
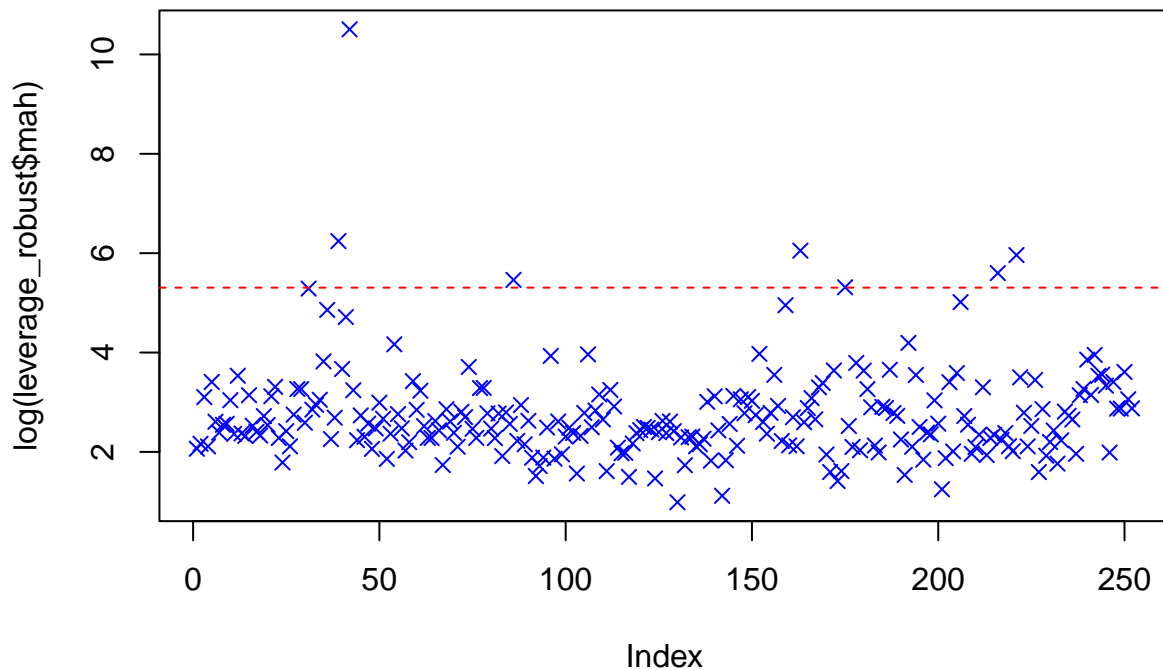
**b**

. . . robust diagnostics based on robust Mahalanobis distances. Use the MCD estimator (covMcd() from the package robustbase) for this purpose.

```r
library(robustbase)
```

```
## Warning: Paket 'robustbase' wurde unter R Version 4.1.3 erstellt
```

```r
leverage_robust = covMcd(Xm)
plot(log(leverage_robust$mah), pch = 4, col = "blue")
abline(h = quantile(log(leverage_robust$mah), 0.975), lty = "dashed", col = "red")
```

```r
leverage_points = which(leverage_robust$mah > quantile(leverage_robust$mah, 0.975))
print(paste("leverage points:", paste(leverage_points, collapse = ", ")))
```

```
## [1] "leverage points: 39, 42, 86, 163, 175, 216, 221"
```

What do you conclude? The classical non robust method results in 6 observations having a p-value greater than 0.975. The robust MCD gets 7 observations with a greater p-value. Therefore the difference seems to be negligible.
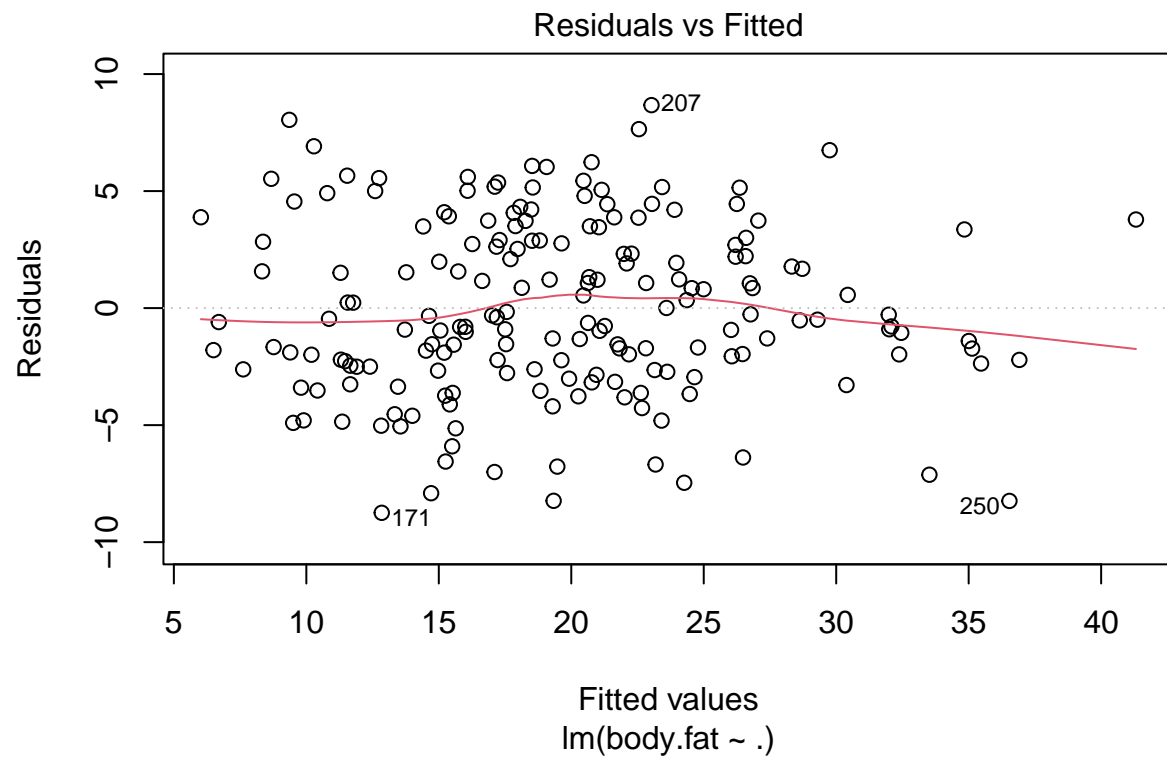
## 2

Now split the observations randomly into training and test data, e.g. in a proportion 3:1. Then apply linear regression to the training set, with ## a the least-squares estimator (lm()) and the robust MM-estimator (lmrob() from library(robustbase)). Interpret the results of summary() and plot().

```r
set.seed(11717659)
sample <- sample(c(TRUE, FALSE), nrow(df), replace = TRUE, prob = c(3/4, 1/4))
train <- df[sample, ]
test <- df[!sample, ]

library(robustbase)

model = lm(body.fat ~ ., train)
```
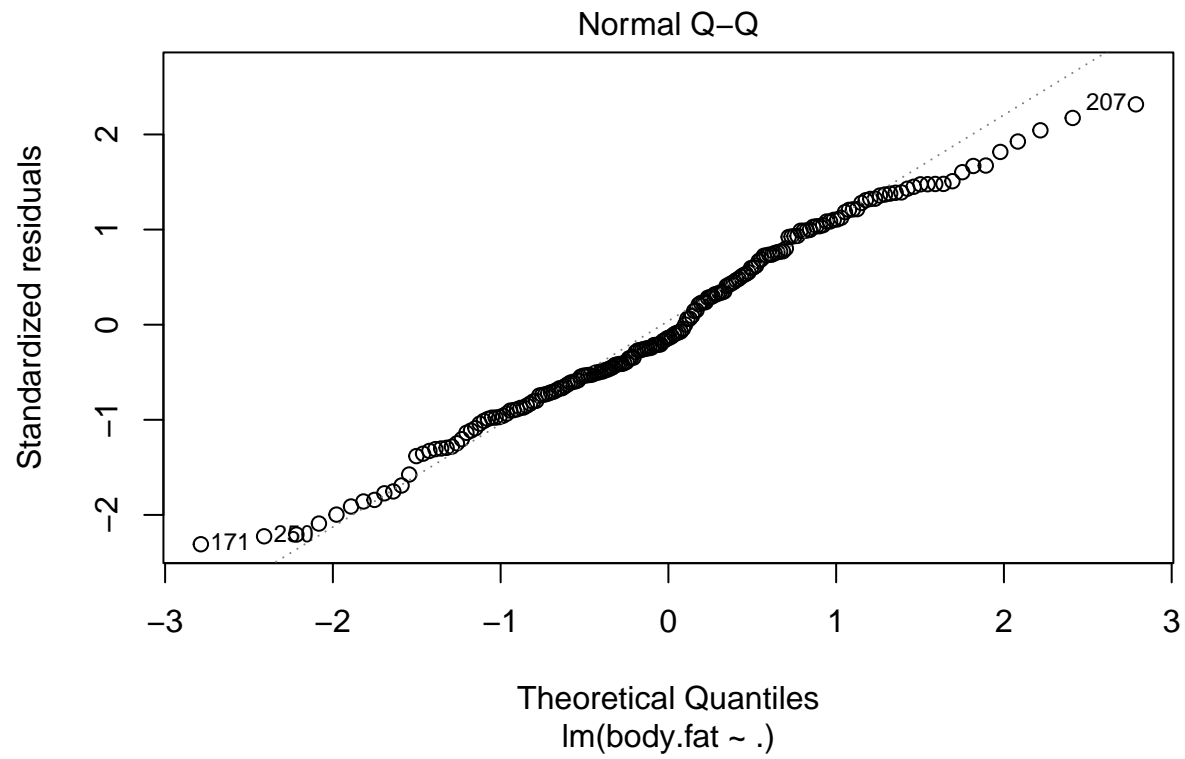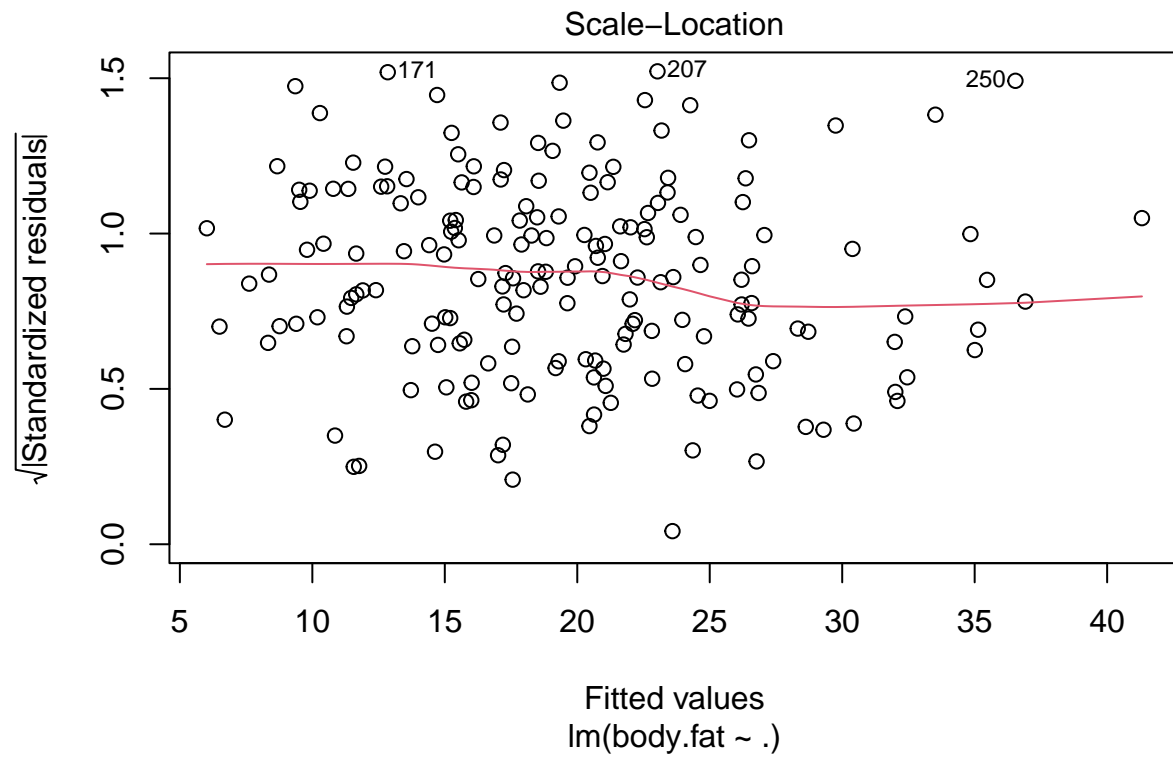
```
# summary(model)
plot(model)
```

## Residuals vs Fitted



Fitted values
lm(body.fat ~ .)

Normal Q–Q

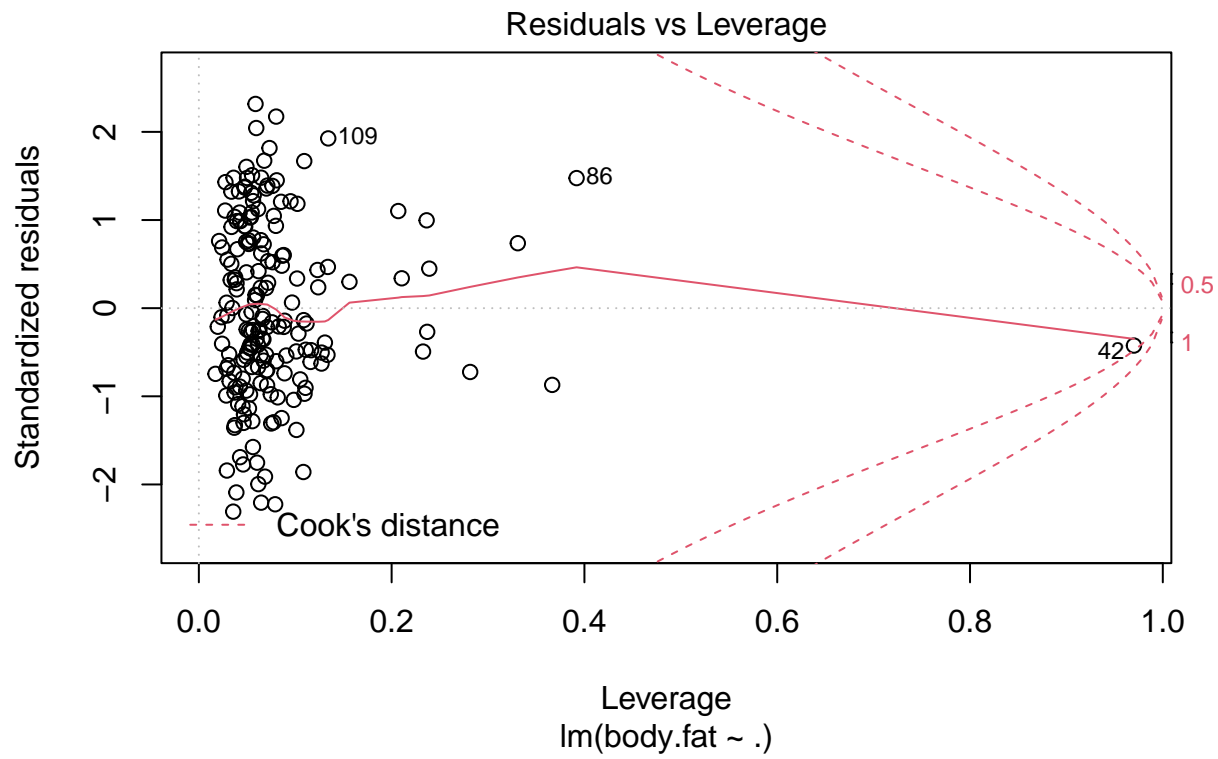Theoretical Quantiles
lm(body.fat ~ .)

Scale–Location

lm(body.fat ~ .)

```
## Warning in sqrt(crit * p * (1 - hh)/hh): NaNs wurden erzeugt

## Warning in sqrt(crit * p * (1 - hh)/hh): NaNs wurden erzeugt
```
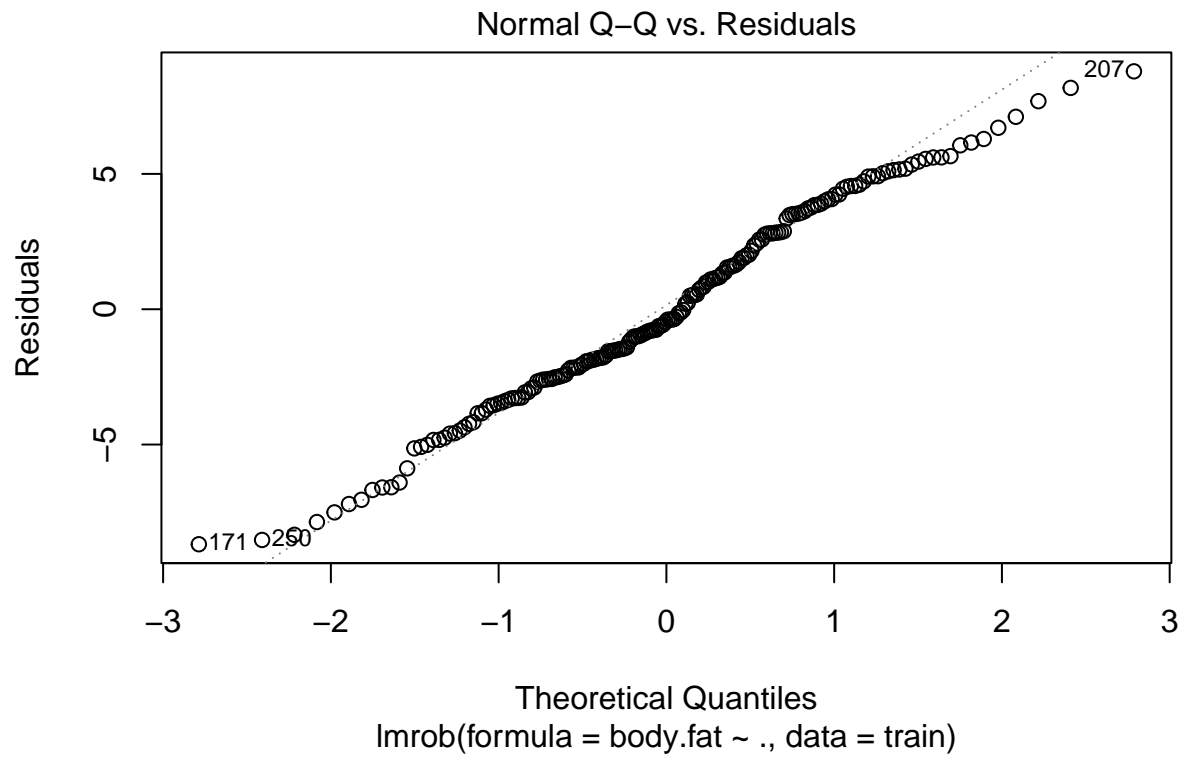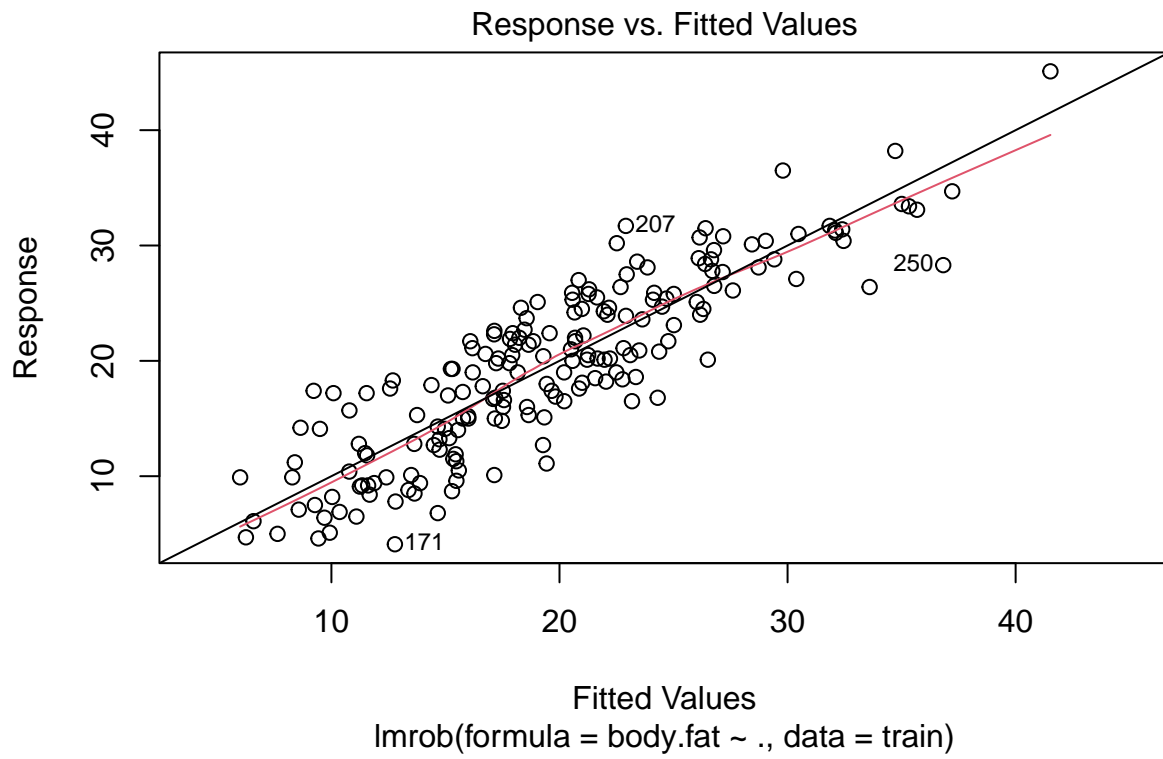
Residuals vs Leverage

lm(body.fat ~ .)

```r
model_rob = lmrob(body.fat ~ ., train)
# summary(model_rob)
plot(model_rob)
```
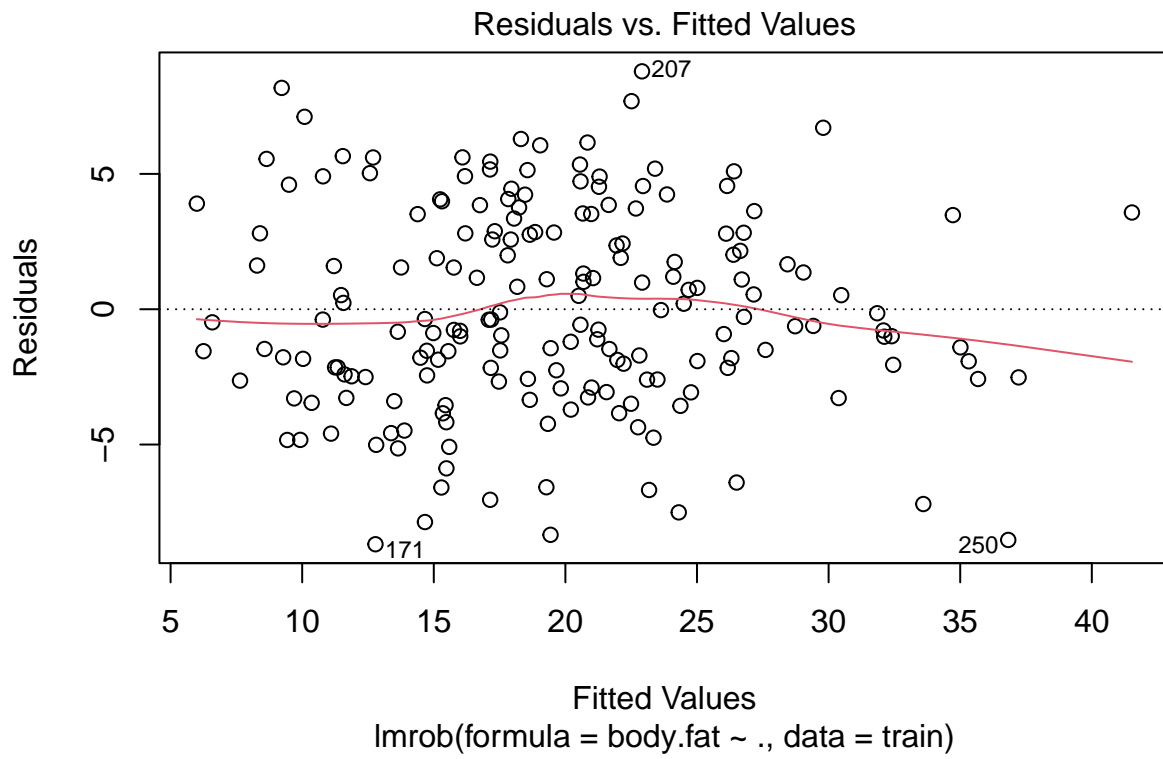
```
## recomputing robust Mahalanobis distances
```

```
## saving the robust distances 'MD' as part of 'model_rob'
```

Standardized residuals vs. Robust Distances

Robust Standardized residuals

Robust Distances
lmrob(formula = body.fat ~ ., data = train)

Normal Q–Q vs. Residuals

Residuals

Theoretical Quantiles
lmrob(formula = body.fat ~ ., data = train)

Response vs. Fitted Values

Response

207
250
171

Fitted Values
lmrob(formula = body.fat ~ ., data = train)

Residuals vs. Fitted Values

Fitted Values
lmrob(formula = body.fat ~ ., data = train)

## Sqrt of abs(Residuals) vs. Fitted Values



Fitted Values
lmrob(formula = body.fat ~ ., data = train)

Is robustness recommendable? I dont think so because the non robust estimator, detects nearly as many outliers as the non robust estimator.

## b

Compute the Cook distances from the least-squares solution.
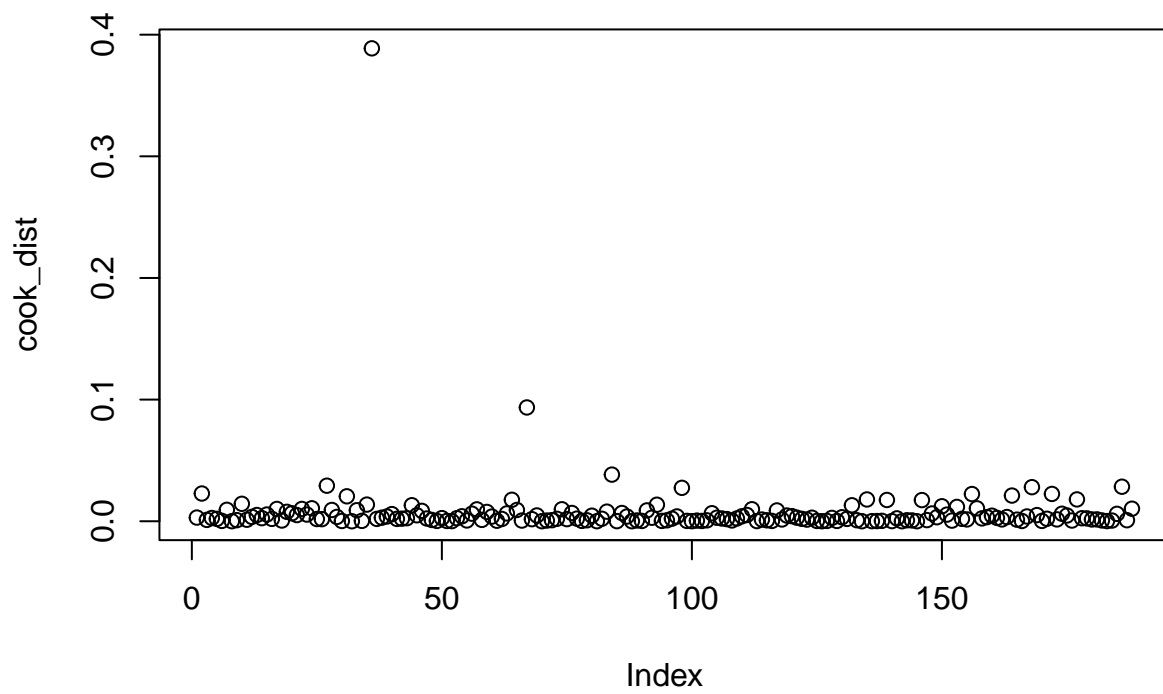
```
cook_dist = cooks.distance(model)
summary(cook_dist)
```

```
##      Min.  1st Qu.   Median     Mean  3rd Qu.      Max.
## 0.0000000 0.0007219 0.0023485 0.0074450 0.0063865 0.3887606
```
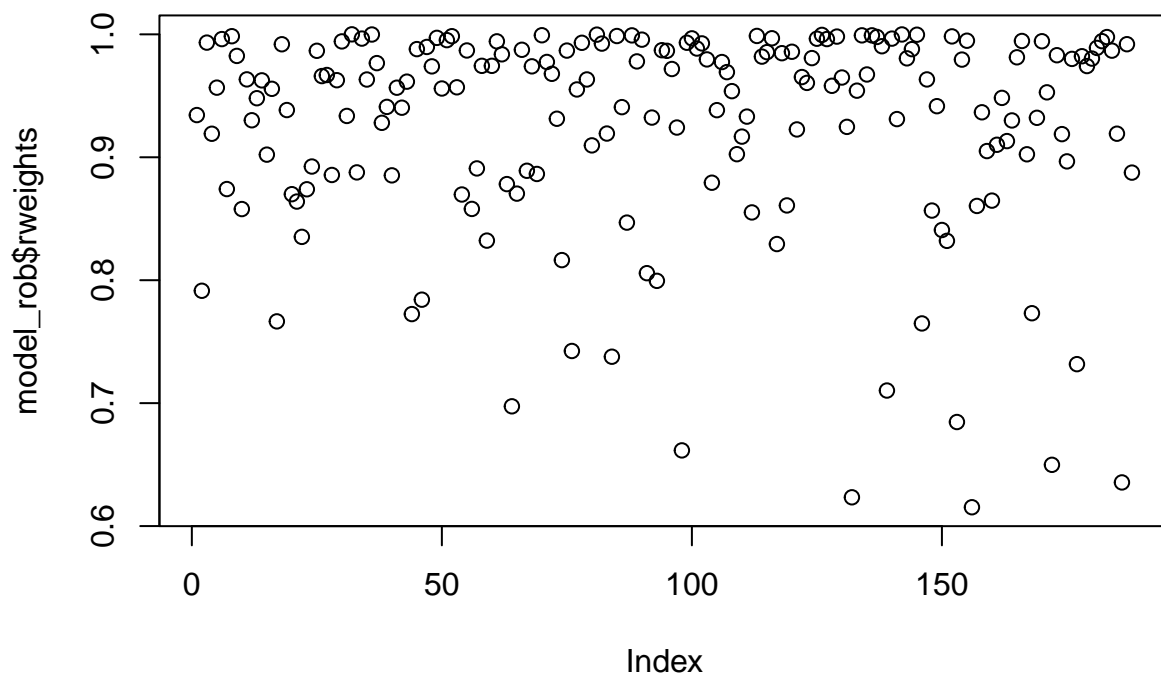
```
summary(model_rob$weights)
```

```
## Length  Class   Mode
##      0   NULL   NULL
```

```
plot(cook_dist)
```
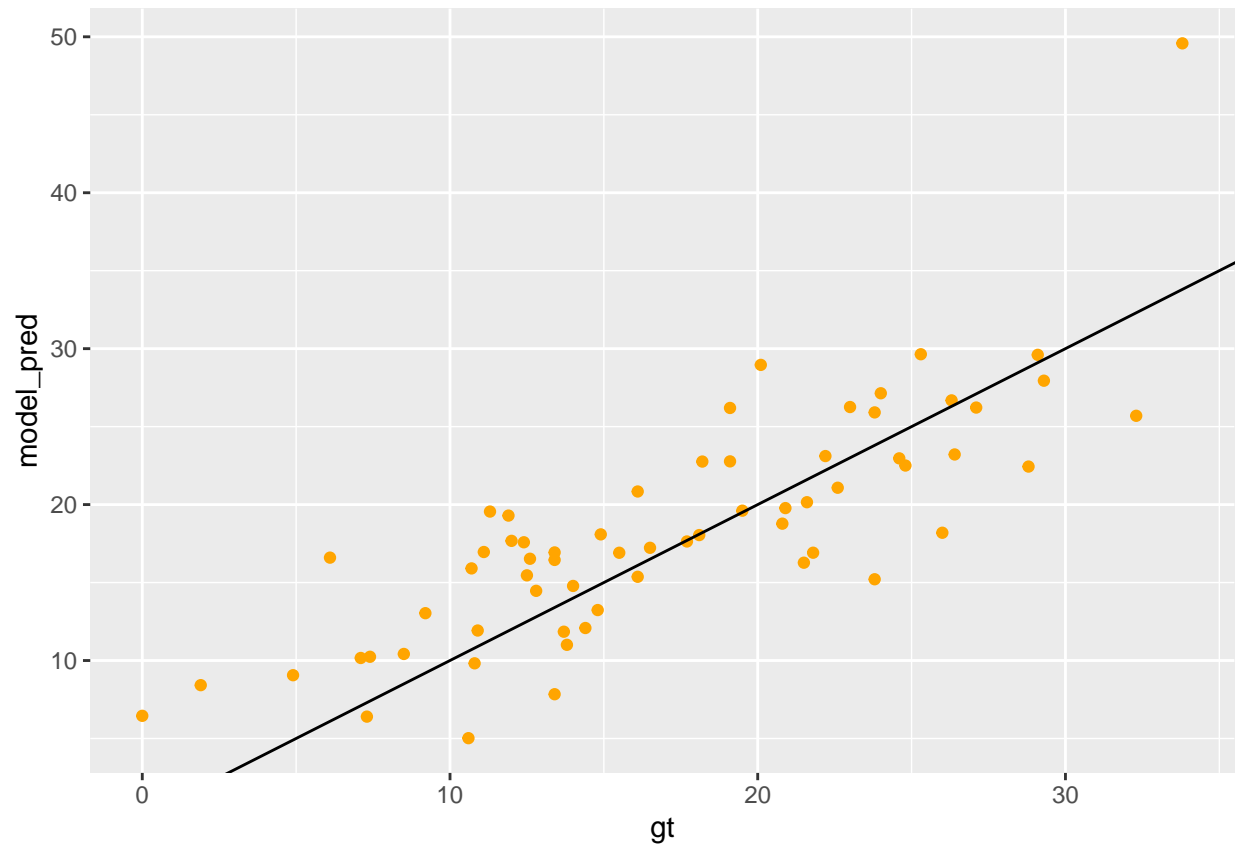
```r
plot(model_rob$rweights)
```

**c**

Use the models to predict the response of the test set. Compare the classical and robust predictions graphically and numerically using an appropriate measure of prediction accuracy.

```r
library(ggplot2)
gt = test$body.fat
model_pred = predict(model, select(test, -body.fat))
model_rob_pred = predict(model_rob, select(test, -body.fat))

res = data.frame(gt, model_pred, model_rob_pred)

ggplot(res) + geom_point(aes(x = gt, y = model_pred), color = "orange") +
↪   geom_abline(intercept = 0,
    slope = 1)
```

```
ggplot(res) + geom_point(aes(x = gt, y = model_rob_pred), color = "orange") +
↪  geom_abline(intercept = 0,
   slope = 1)
```