

TP: is Batman somewhere?

Gabriele Degola

December 2021

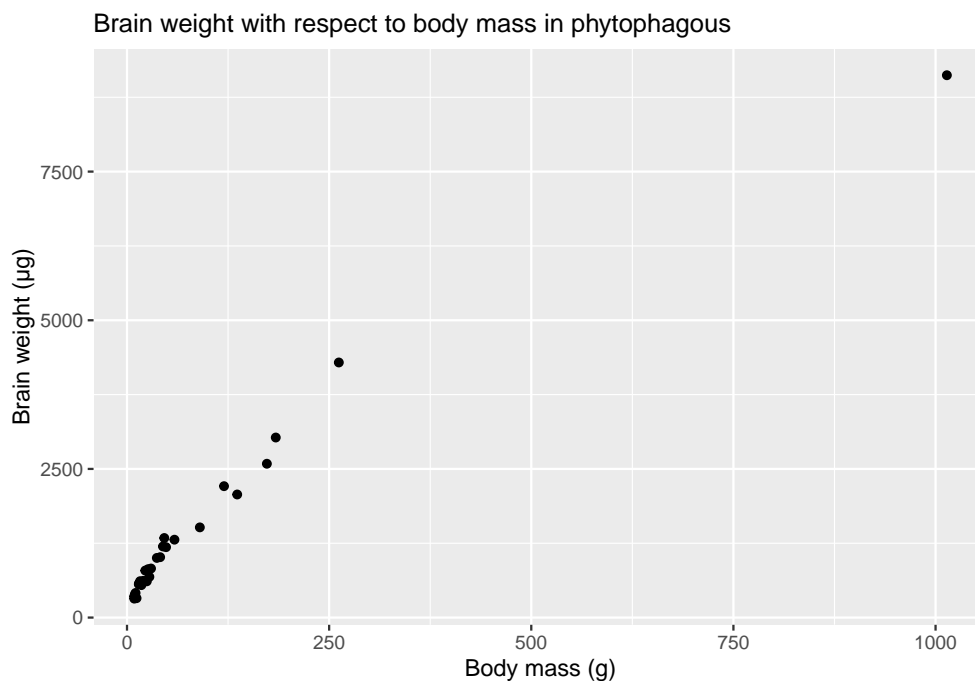
```
library(tidyverse)
library(corrplot)
```

```
myData <- read.table(file="../data/bats.csv", sep=";", skip=3, header=TRUE)
names(myData)
```

```
## [1] "Species" "Diet" "Clade" "BOW" "BRW" "AUD" "MOB"
## [8] "HIP"
```

Study of the relationship between brain weight and body mass

```
phyto <- myData[(myData$Diet == 1),]
ggplot(phyto, aes(x=BOW, y=BRW)) + geom_point() +
  xlab("Body mass (g)") + ylab("Brain weight (µg)") +
  ggtitle("Brain weight with respect to body mass in phytophagous") +
  theme(plot.title = element_text(size=12))
```



It looks like body mass (in grams) has a clear influence on brain weight (in micrograms). This relation can be analysed fitting a linear model of brain weight as a function of body mass. The fitted line is also shown in next plot.

```
reg1 <- lm(BRW ~ BOW, data=phyto)
ggplot(phyto, aes(x=BOW, y=BRW)) + geom_point() + geom_smooth(method="lm") +
  xlab("Body mass (g)") + ylab("Brain weight (µg)") +
  ggtitle("Brain weight with respect to body mass in phytophagous (with regression
  ↪ line)") +
  theme(plot.title = element_text(size=12))
```



With this expression, R estimates the β coefficients of the formula $Y = \beta_0 + \beta_1 X + \epsilon$, where X is the vector containing the values of the measured body masses and Y contains the measured brain weights. ϵ is random noise.

```
summary(reg1)
```

```
##
## Call:
## lm(formula = BRW ~ BOW, data = phyto)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -628.32 -233.94  -65.74  158.26 1308.59
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  623.4469    81.4762   7.652 3.14e-08 ***
## BOW           8.9999     0.3972  22.659 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 396.9 on 27 degrees of freedom
## Multiple R-squared:  0.95, Adjusted R-squared:  0.9482
## F-statistic: 513.4 on 1 and 27 DF, p-value: < 2.2e-16
```

In this case, a few results can be observed:

- The intercept β_1 is estimated as 623.45.
- The p-value for the employed test statistics is extremely small, lower than 2.2×10^{-16} .
- The null hypothesis of this test is $H_0 : \beta_1 = \beta_2 = 0$. This hypothesis is therefore rejected.
- Body mass has a significative relation with brain weight. The intercept is also important.
- The coefficient of determination R^2 of the fitted model is 0.95, so a big portion of the global variation is explained by the model.

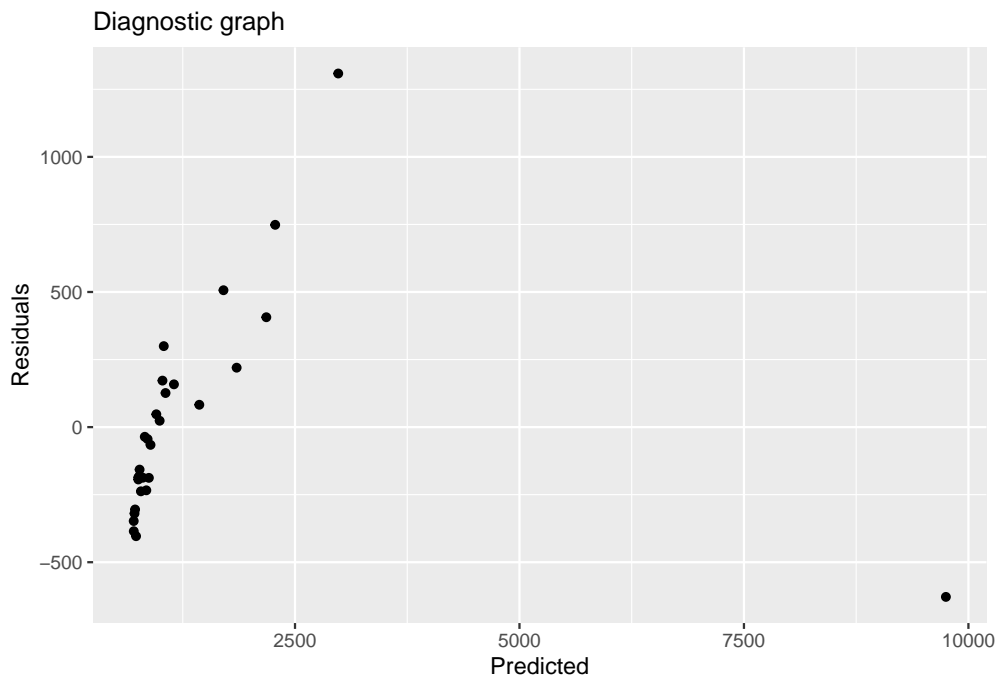
Next, analysis of variance can be performed.

```
anova(reg1)
```

```
## Analysis of Variance Table
##
## Response: BRW
##          Df  Sum Sq Mean Sq F value    Pr(>F)
## BOW         1 80888380 80888380  513.42 < 2.2e-16 ***
## Residuals  27  4253838   157550
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

This confirms that body weight is highly significative. In addition, the previous table contains information about the model residual, with the sum of residual squares being 4253838. The residuals have been computed during model fitting and can be plotted against the predicted values.

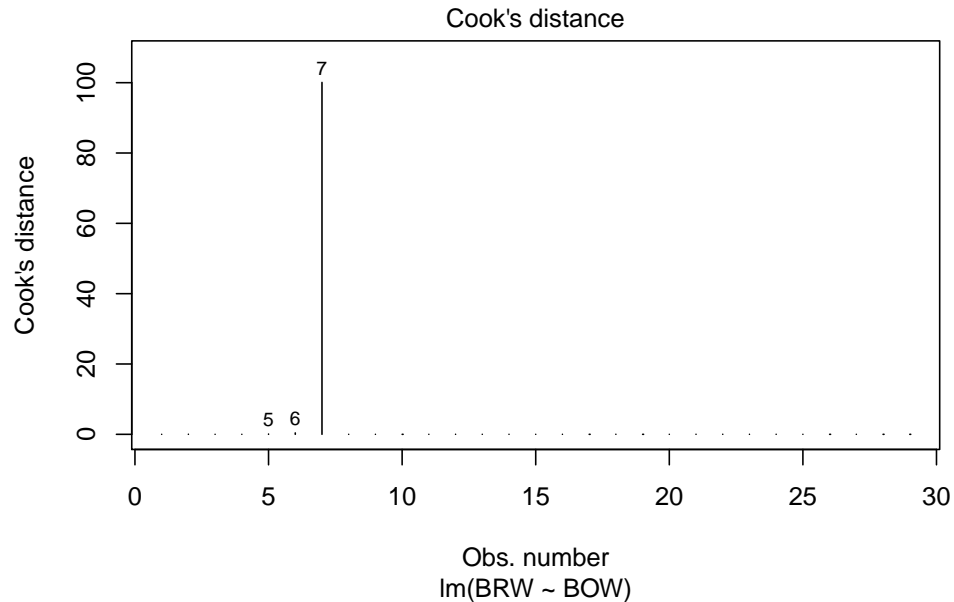
```
tmp <- data.frame(reg1$fitted.values, reg1$residuals)
ggplot(tmp, aes(x=reg1.fitted.values, y=reg1.residuals)) + geom_point() +
  xlab("Predicted") + ylab("Residuals") + ggtitle("Diagnostic graph") +
  theme(plot.title = element_text(size=12))
```



Clearly, something is not going well: most prediction values are below 3750 μg , except one with predicted brain weight around 10000 μg . As this last point is so far from the others, it influences the model results worsening the prediction for values in the middle. For now, we can consider it as an outlier.

Cook's distance can be used for outlier detection in the fitted model:

```
plot(reg1, 4)
```



Observation number 7 has distance much higher than the higher and probably corresponds to the previously identified point.

```
myData[7,]
```

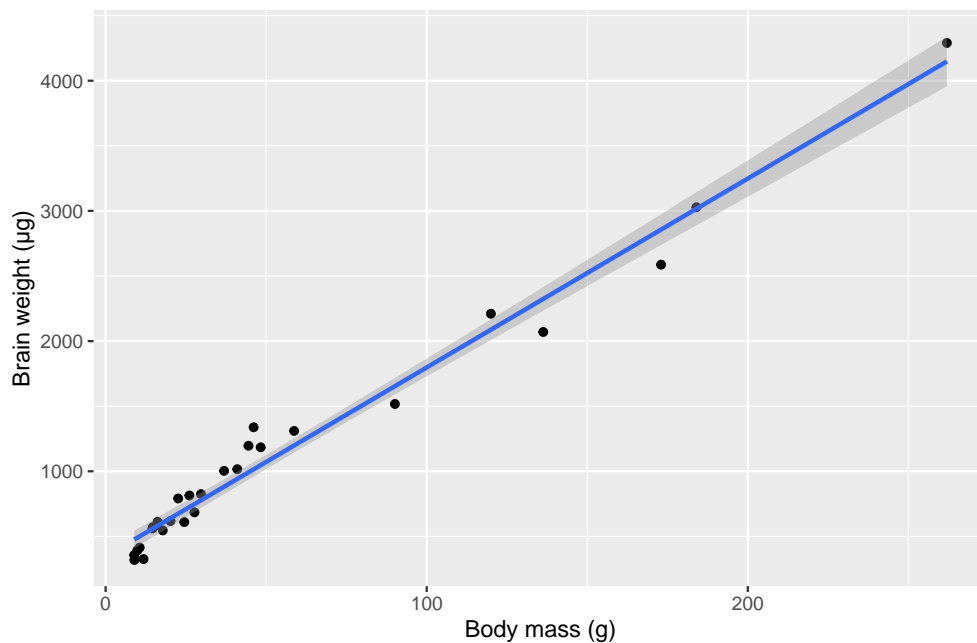
```
##           Species Diet Clade BOW BRW  AUD  MOB  HIP
## 7 Pteropus vampyrus    1     I 1014 9121 16.93 243.54 331.29
```

Indeed, Pteropus vampyrus has body mass 1014 g and brain weight 9121 μg .

We can therefore redo the previous analysis without taking it into account.

```
phytobis <- phyto[which(phyto$BRW<8000),]
reg2 <- lm(BRW ~ BOW, data=phytobis)
ggplot(phytobis, aes(x=BOW, y=BRW)) + geom_point() + geom_smooth(method="lm") +
  xlab("Body mass (g)") + ylab("Brain weight ( $\mu\text{g}$ )") +
  ggtitle("Brain weight with respect to body mass in phytophagous (with regression  
↪ line)") +
  theme(plot.title = element_text(size=12))
```

Brain weight with respect to body mass in phytophagous (with regression line)



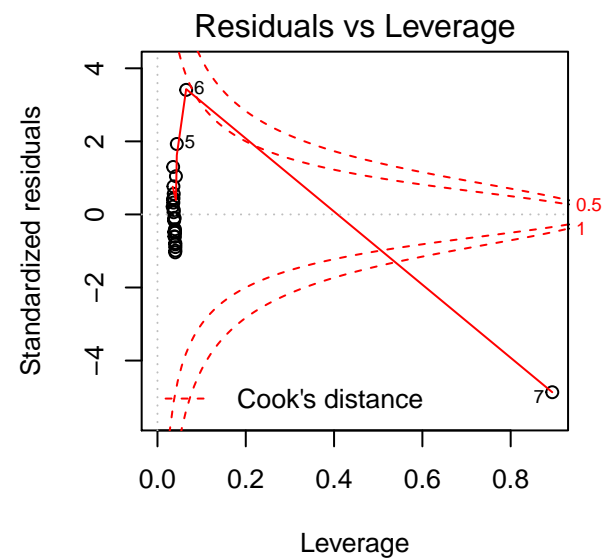
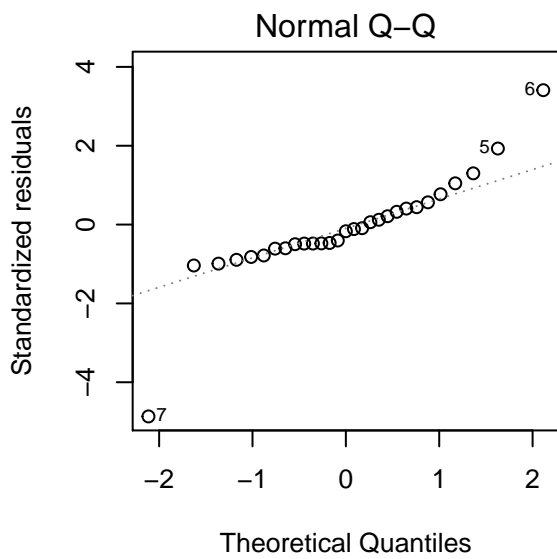
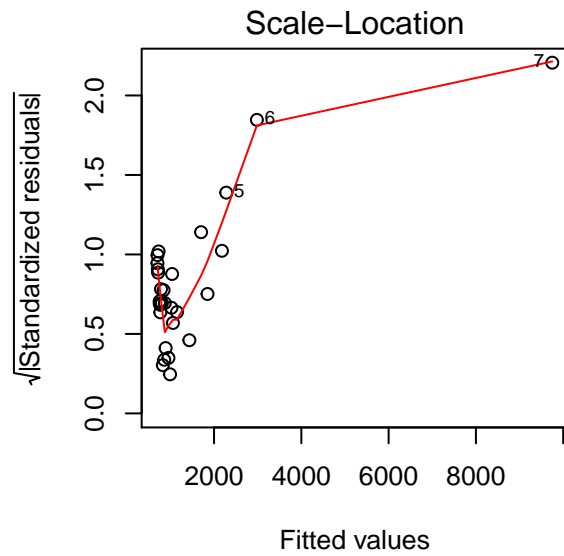
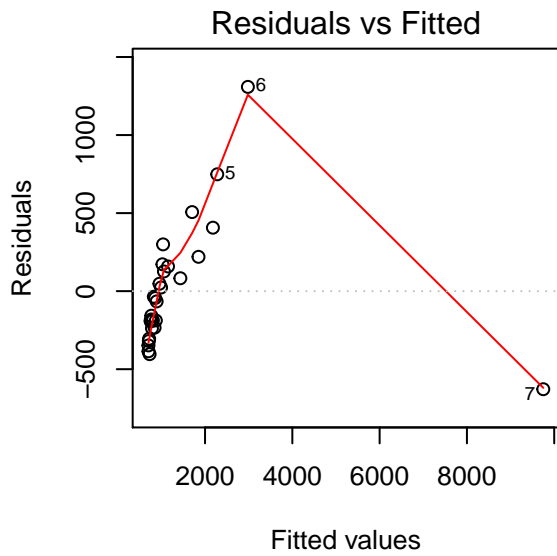
```
summary(reg2)
```

```
##
## Call:
## lm(formula = BRW ~ BOW, data = phytobis)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -269.76  -93.33    8.73  112.93  322.55
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 346.5452    35.4920   9.764 3.48e-10 ***
## BOW         14.5099     0.4285  33.860 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 141.8 on 26 degrees of freedom
## Multiple R-squared:  0.9778, Adjusted R-squared:  0.977
## F-statistic: 1147 on 1 and 26 DF, p-value: < 2.2e-16
```

Clearly, now the regression line better fits intermediate point and is not affected by the removed extreme observation. The predicted coefficient for the body weight predictor is higher than before, as well as the coefficient of determination.

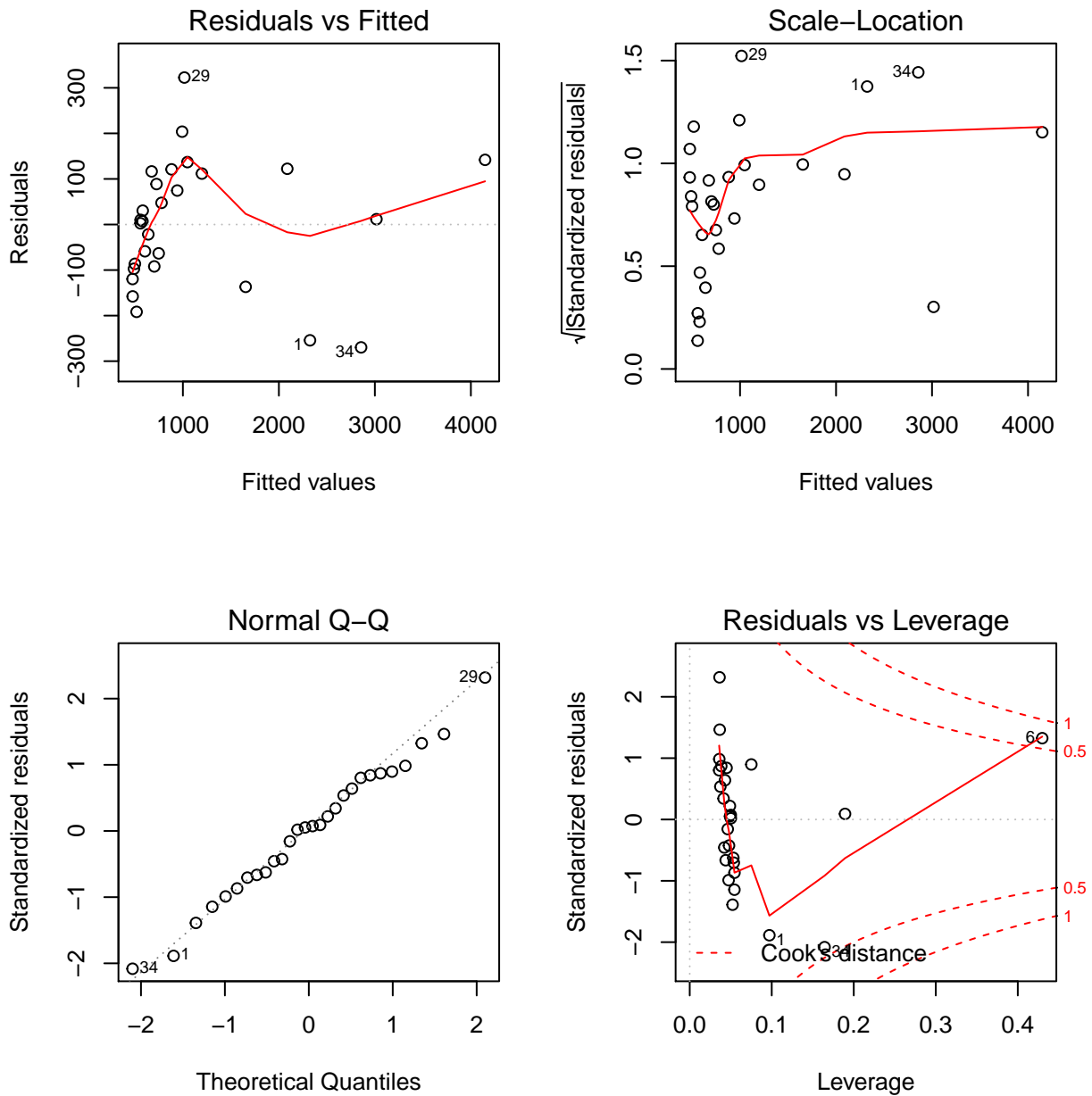
A few graphs can therefore plotted for model diagnosis, including the residuals against predicted shown before. For the first model, including the “outlier”:

```
par(mfcol=c(2,2))
plot(reg1)
```



If the “outlier” is removed:

```
par(mfcol=c(2,2))
plot(reg2)
```



Important differences can be noticed:

- The plots related to residuals vs fitted and scale-location are more “flat” for the second model, so it provides more accurate predictions than the first one.
- According to the Q-Q plot, it is more probable that the data used to fit the second model are drawn from a normal distribution, with respect to the first model.

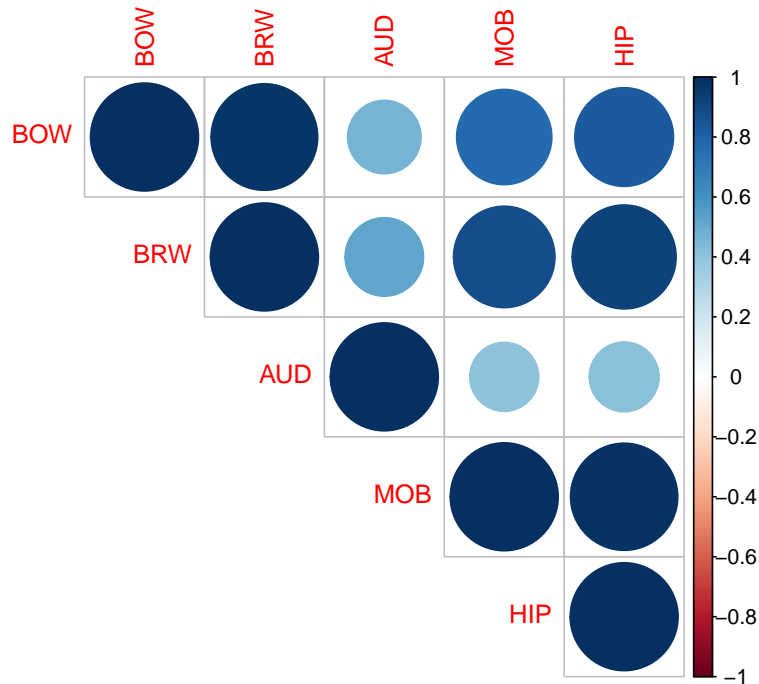
Therefore, the second model better satisfies common assumptions of linear regression.

Study of the contribution to the total weight of each part of the brain

Previous analysis can be expanded, using other variables (related to brain parts) to explain brain weight. The possible explanatory variables are AUD (auditory nuclei volume), MOB (main olfactory bulb volume)

and HIP (hippocampus volume).

```
phytoNum <- phyto[, c(4:8)]
mat.cor <- cor(phytoNum)
corrplot(mat.cor, type="upper")
```



We see high positive correlation in the couples, for example, (BOW, BRW) and (MOB, HIP), while the AUD variable is not particularly correlated with any other. Let's statistically analyse the correlation between some variables, using the Pearson test.

```
cor.test(phyto$BRW, phyto$HIP)
```

```
##
## Pearson's product-moment correlation
##
## data: phyto$BRW and phyto$HIP
## t = 12.91, df = 27, p-value = 4.574e-13
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## 0.8502663 0.9658107
## sample estimates:
## cor
## 0.9276811
```

```
cor.test(phyto$BRW, phyto$MOB)
```

```
##
## Pearson's product-moment correlation
##
## data: phyto$BRW and phyto$MOB
## t = 9.7964, df = 27, p-value = 2.203e-10
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
```



```
## 0.7644185 0.9442114
## sample estimates:
##      cor
## 0.8834215
```

```
cor.test(phyto$BRW, phyto$AUD)
```

```
##
## Pearson's product-moment correlation
##
## data: phyto$BRW and phyto$AUD
## t = 3.2338, df = 27, p-value = 0.003215
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## 0.2007495 0.7497021
## sample estimates:
##      cor
## 0.5283792
```

Resulting p-values are very low if BRW is tested against HIP and MOB, while it is bigger (0.003) against AUD. What is shown in the previous figure is confirmed: correlation is important with HIP and MOB, less with AUD.

```
regm <- lm(BRW ~ AUD + MOB + HIP, data=phytobis)
summary(regm)
```

```
##
## Call:
## lm(formula = BRW ~ AUD + MOB + HIP, data = phytobis)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -268.55  -68.84    9.88   61.66  375.34
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -312.692     76.628  -4.081  0.00043 ***
## AUD           47.989      6.067   7.910  3.85e-08 ***
## MOB          -2.444      3.257  -0.750  0.46034
## HIP           15.981      2.960   5.399  1.52e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 158.5 on 24 degrees of freedom
## Multiple R-squared:  0.9744, Adjusted R-squared:  0.9712
## F-statistic: 304.5 on 3 and 24 DF,  p-value: < 2.2e-16
```

```
anova(regm)
```

```
## Analysis of Variance Table
##
## Response: BRW
##      Df  Sum Sq Mean Sq F value  Pr(>F)
## AUD    1 6817133 6817133 271.210 1.397e-14 ***
## MOB    1 15409397 15409397 613.040 < 2.2e-16 ***
## HIP    1  732653  732653  29.148 1.519e-05 ***
```

```
## Residuals 24    603265    25136
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The estimated model predicts the β coefficient of the regression model $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$, where Y is the BRW variable and X_1 , X_2 and X_3 are AUD, MOB and HIP respectively. The objective is therefore to predict the brain mass of a bat from the volumes of its auditory nuclei, main olfactory bulb and hippocampus. Theoretically, they may all have an impact on the brain mass, even though the main olfactory bulb can be present in the bat's nose and not in its brain.

The coefficient associated to the three variables are, respectively, 47.989, -2.444 and 15.981. The coefficient related to MOB is not significant, showing a high p-value, and could therefore be removed from the analysis.

Additional analysis can be conducted on this model. Here, the “batman” is taken into account.

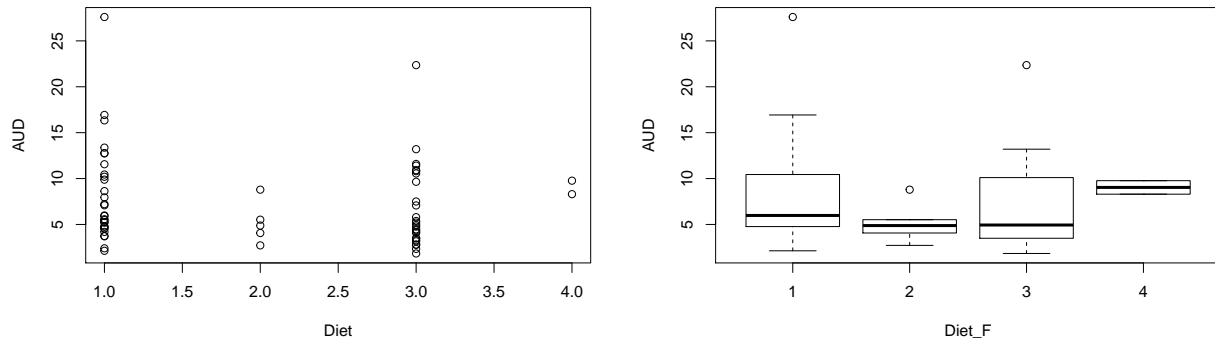
```
reg0 <- lm(BRW ~ 1, data=phyto)
step(reg0, scope=BRW ~ AUD + MOB + HIP, direction="forward")
```

```
## Start:  AIC=433.88
## BRW ~ 1
##
##           Df Sum of Sq      RSS      AIC
## + HIP      1  73272731 11869487 378.74
## + MOB      1  66447848 18694370 391.92
## + AUD      1  23770396 61371823 426.39
## <none>                    85142218 433.88
##
## Step:  AIC=378.74
## BRW ~ HIP
##
##           Df Sum of Sq      RSS      AIC
## + MOB      1   2846939  9022548 372.79
## + AUD      1   2013783  9855704 375.35
## <none>                    11869487 378.74
##
## Step:  AIC=372.79
## BRW ~ HIP + MOB
##
##           Df Sum of Sq      RSS      AIC
## + AUD      1   1910121  7112426 367.89
## <none>                    9022548 372.79
##
## Step:  AIC=367.89
## BRW ~ HIP + MOB + AUD
##
## Call:
## lm(formula = BRW ~ HIP + MOB + AUD, data = phyto)
##
## Coefficients:
## (Intercept)           HIP           MOB           AUD
##   -1003.95         44.35        -29.24         52.82
```

According to the `help`, `step` performs a stepwise search for a best model according to the AIC criterion. In this case, intercept only is considered at the first step and all explanatory variables are included at the end. In sequence, the models `BRW ~ 1`, `BRW ~ HIP`, `BRW ~ HIP + MOB` and `BRW ~ HIP + MOB + AUD` are evaluated, with the latter having the lowest information loss.

Link between volume of the auditory part and diet

```
myData$Diet_F <- as.factor(myData$Diet)
with(myData, plot(AUD~Diet))
with(myData, plot(AUD~Diet_F))
```



The two graphs basically show the same information, but in different ways. The first displays the single data points and it is clear that few examples are available for diets 2 and 4. The second displays the boxplots instead and shows data distributions for each different diet. Boxplots for diets 2 and 4 are condensed around the mean, while the ones with more data points are more spread.

```
lm <- lm(AUD~Diet_F, data=myData)
anova(lm)
```

```
## Analysis of Variance Table
##
## Response: AUD
##          Df Sum Sq Mean Sq F value Pr(>F)
## Diet_F    3  66.07  22.023  0.9293 0.4323
## Residuals 59 1398.26  23.699
```

The p-value of the ANOVA test is quite high, 0.432, so it is not possible to conclude that there is a significant difference between the different diets. The bats' diets don't have therefore an impact on the volume of their auditory part of the brain.