

Market Forecasting Assignment: Multiple Linear Regression

Teacher: Dr. Stefan Groesser, Dr. Christoph Imboden

Autors: Benedech Rodolpho, Buehler Pascal, Geeler Ken 04.12.2021

Short Theory Recap Multiple Linear Regression (MLR)

The model equation for MLR has the form $y = X\Theta + \epsilon$, with y being the target vector, X the explanatory Matrix, Θ the parameter vector and the residuals with the assumption that they are normally distributed $\epsilon \sim N(0, \sigma^2)$. The estimation of Θ is done with the Maximum likelihood $L(\theta) = P(X|\theta) = \frac{1}{\sqrt{(2\pi)\sigma}} * e^{-\frac{1}{2\sigma^2} ||y - \theta X||^2}$ by solving the equation $\Theta = (X^T X)^{-1} X^T y$.

Dataset

The data from 7 common fish species for fishmarket is provided from Kaggle. <https://www.kaggle.com/aungpyaeap/fish-market>

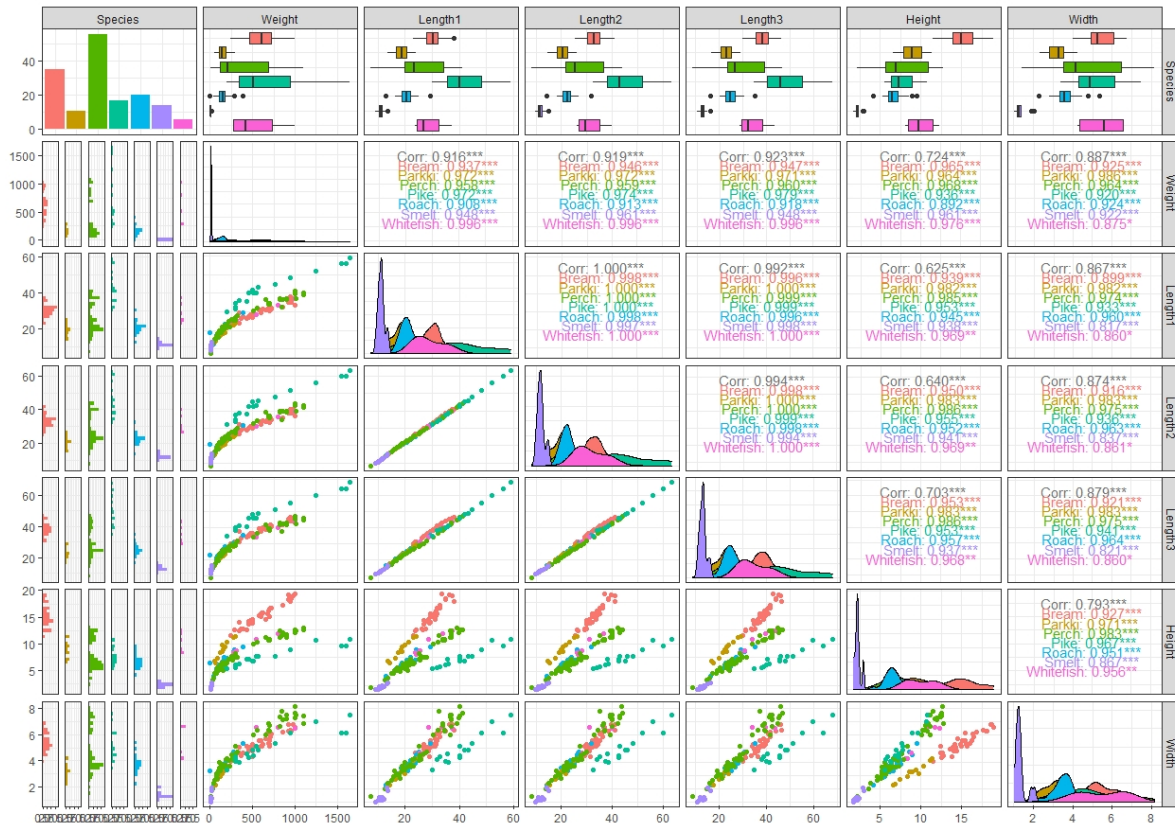


Figure 1: pairwise plot of the fish dataset

Data exploration

Information of the Dataset is gained with the following commands. *Note: the output is not shown to reduce space in the document full code and generation pf document can be found on https://github.com/buehlpa/TSM_MarkFor*

```
fish_dat=readr::read_csv("Fish.csv") # reading as tibble
GGally::ggpairs(fish_dat,aes(color = Species)) + theme_bw() # pairwise plot
summary(fish_dat);str(fish_dat) # statistics and structure
```

The pairwise plot Figure 1 provides a visual overview of the dataset. It consists of 159 observations measured for 6 numerical variables which measure properties of the fish and one categorical variable to describe the Species. The coloring in the pairwise plot allows a distinction of the Species, one can immediately observe that the species Perch has the most observations.

From the pairwise scatterplots a nonlinear relation between weight and the other numerical variables is observed, therefore the correlation *cannot* be correctly interpreted with these variablepairs. Further the other numerical variables are linearly correlated (the correlation statistic can be interpreted), interestingly the 3 length measurements have a very strong correlation to each other.

The density plots helps see how the variables themselves are distributed. For example the width of Smelt tends to have a bidmodal distribution, although this has to be taken with a grain of salt due to the sparse data, see on the inverted histograms on the left. Also the pairwise boxplots gives a good overview for comparison between the different variables. For example it can be seen that the variance of the species vary amongst the levels.

The distribution of the target variables is actually very important since MLR assumes multivariate normality of the data, for other distributions we might choose another model such as a Generalized-Linear-Model to improve the model validity, but for this analysis we stick to MLR.

Building a Model

The dataset provides many different options to do linear regression, nevertheless for MLR there is a problem with multicollinearity since a lot of the variables are strongly correlated, meaning that if multiple independent variables would be in the model it wouldn't be clear which one explains the effect the best for the target variable.

For a first approach we could just focus on the species Perch (in swiss-german it's called Egli) since we have the most data there. Lets find a Model for Predicting the weight of this particular species with the variables width height and length. Since the weight hasn't got a linear relation to the other variables, we'll transform it properly.

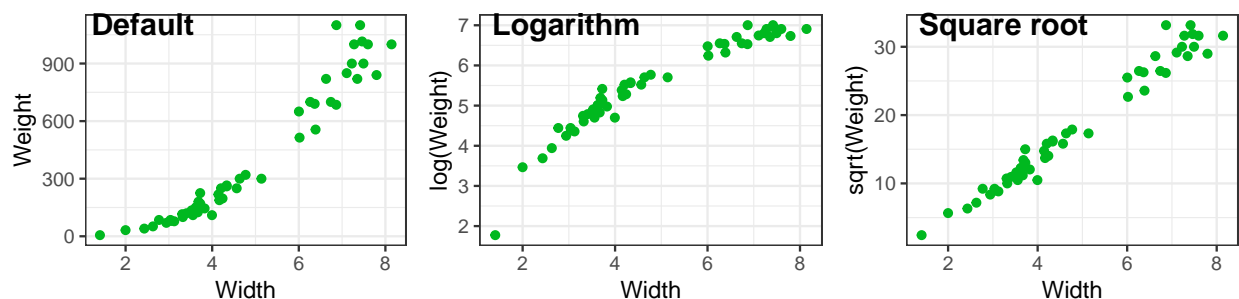


Figure 2: transformations of weight for Perch

If we transform the weight with the logarithm it looks a bit better but still isn't linear see Figure 2. The squareroot transformation seems to be the right one. Now the data seems much more linearly dependent. Therefore a first model can be set up

```
Perch_fit=lm(sqrt(Weight)~Width,data = Perch_dat) # Simple linear regression model
```

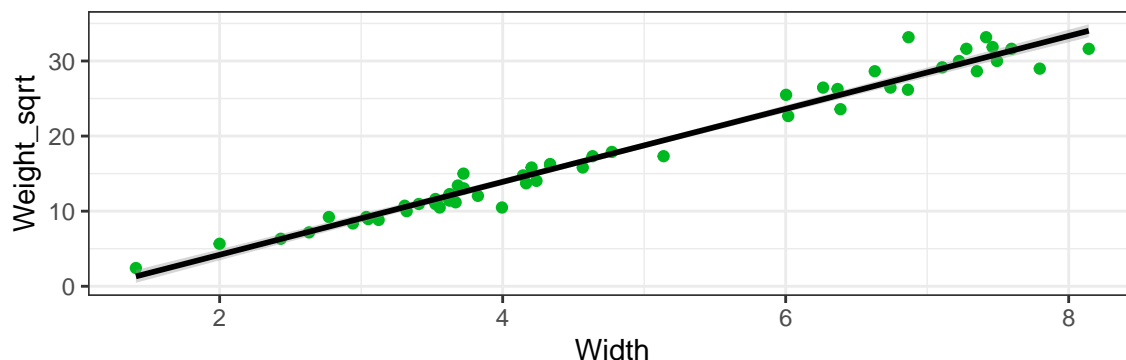


Figure 3: linear regression model for Perch weight vs width

From the summary output we get the following model $Weight_i = (-5.5245 + 4.8571 * Width_i)^2$. The explained Variance measurement R^2 is with 0.9719 very high. the AIC for this model is 206.76. The p-values for the both parameters are also significant. By inspecting the residuals in Figure 4 we can see that the model assumptions are met. The residuals seem to follow a normal distribution, have constant variance, constant mean and there are no points which are crucial to shift the model.

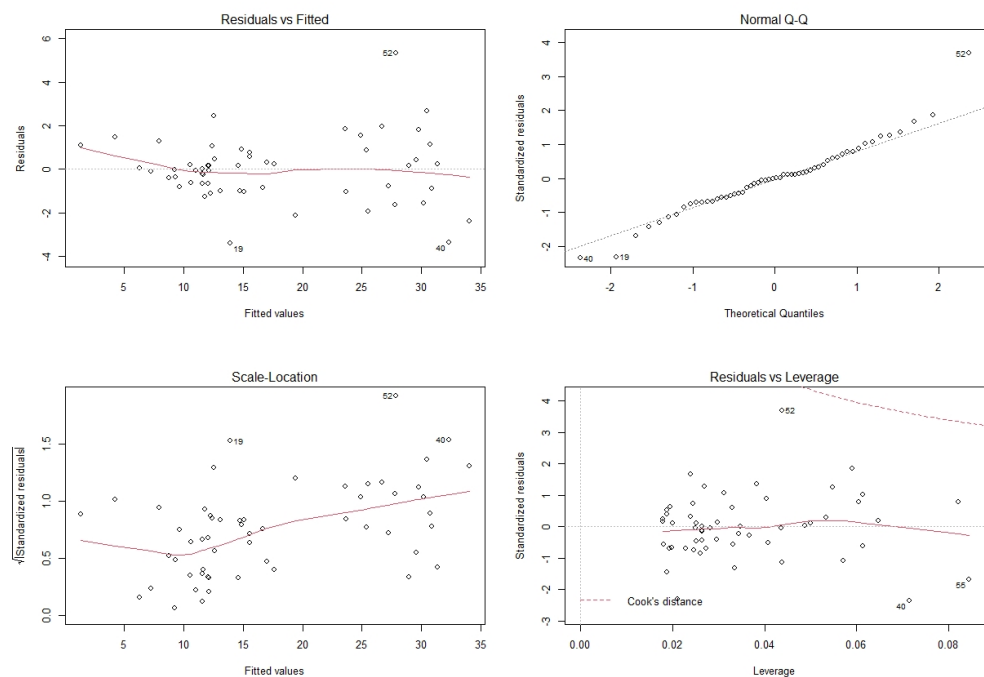


Figure 4: Residuals for the Perchfit

The way we have chosen is starting with a sparse model. Lets make now a full model and test if a larger model is more beneficial. The multicollinearity is checked with Variance Inflation Factor and a F test checks the importance of the variables

```
Perch_fit_full=lm(sqrt(Weight)~.,data = Perch_dat[,-1]) # Multiple linear regression model
summary(Perch_fit_full)
```

```
##
## Call:
## lm(formula = sqrt(Weight) ~ ., data = Perch_dat[, -1])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.238e-15 -7.163e-17  2.680e-17  1.494e-16  6.605e-16
##
## Coefficients:
##              Estimate Std. Error    t value Pr(>|t|)
## (Intercept)  3.357e-15  5.021e-16  6.685e+00  2.04e-08 ***
## Length1     -2.028e-16  2.649e-16 -7.660e-01  0.44762
## Length2      2.778e-16  4.045e-16  6.870e-01  0.49545
## Length3     -2.377e-16  2.747e-16 -8.650e-01  0.39115
## Height      -4.814e-16  1.538e-16 -3.131e+00  0.00294 **
## Width       -5.623e-16  1.866e-16 -3.014e+00  0.00408 **
## Weight_sqrt  1.000e+00  5.986e-17  1.671e+16  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.967e-16 on 49 degrees of freedom
## Multiple R-squared:  1, Adjusted R-squared:  1
## F-statistic: 4.454e+33 on 6 and 49 DF, p-value: < 2.2e-16
```

```
car::vif(Perch_fit_full)
```

```
##      Length1      Length2      Length3      Height      Width Weight_sqrt
## 1798.09443 4655.27365 2396.05622 68.45427 38.31884 95.74095
```

```
drop1(Perch_fit_full)
```

```
## Single term deletions
##
## Model:
## sqrt(Weight) ~ Length1 + Length2 + Length3 + Height + Width +
##      Weight_sqrt
##      Df Sum of Sq    RSS    AIC
## <none>                 0.000 -3965.4
## Length1      1      0.000  0.000 -4061.0
## Length2      1      0.000  0.000 -3998.4
## Length3      1      0.000  0.000 -4026.3
## Height       1      0.000  0.000 -4017.9
## Width        1      0.000  0.000 -3943.2
## Weight_sqrt  1     43.916 43.916    -1.6
```

Vif tells us that all the variables have very high kollinearity (as expected) whereas the vif of the lengths is massively hiigher than of the other variables.