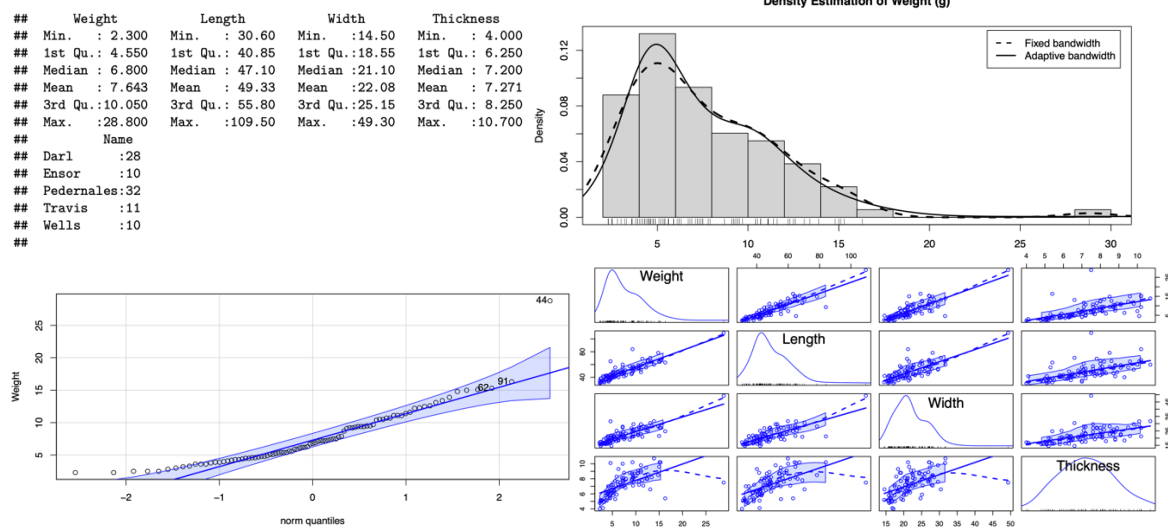


Excercise 1

(a) As a first step, we obtain summary statistics for the variables in Darts.csv; Secondly, we view distribution of the variable Weight;



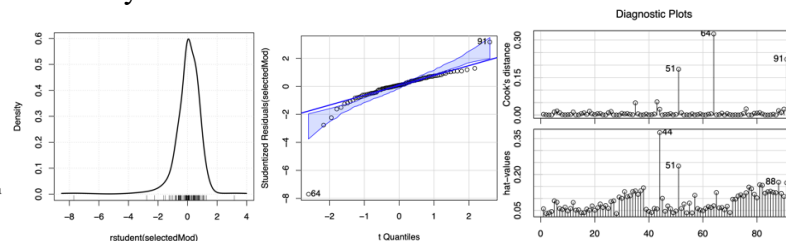
Both nonparametric density estimates and the histogram suggest a mode at around 5g, and all show that the distribution of weight is right-skewed. The fixed-bandwidth kernel estimate has more wiggle at the right where data are sparse, and the histogram is rough in this region, while the adaptive-kernel estimator is able to smooth out the density estimate in the low-density region. And because many points, especially at the left of the graph, are outside the confidence bounds, we have evidence that the distribution of weight is not like a sample from a normal population.

Then, we use scatterplots to provide summaries of the conditional distribution of a numeric response variable given a numeric predictor. The `scatterplotMatrix()` function produces scatterplots for all pairs of numeric variables.

From the above graphs, we apply the log transformation on the variable Weight, Length and Width.

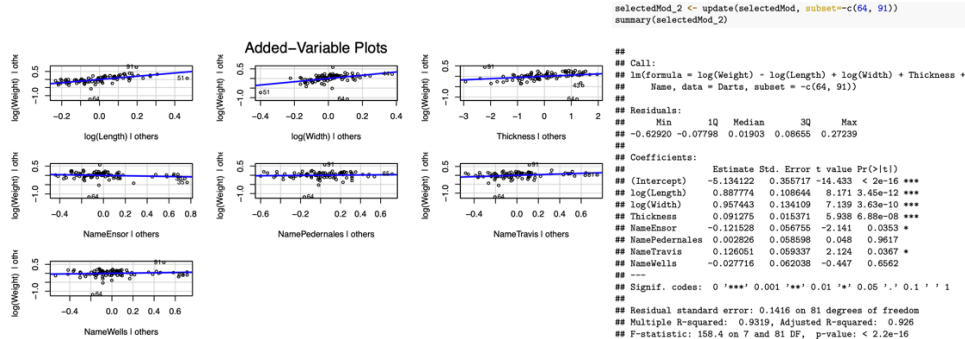
We use `thelm()` function to fit a linear regression model to the data. We use stepwise regression to select best model for the variable Weight. We pass the full model to step function. It iteratively searches the full scope of variables in backwards directions by default, if scope is not given. It performs multiple iterations by dropping one X variable at a time. In each iteration, multiple models are built by dropping each of the X variables at a time. The AIC of the models is also computed and the model that yields the lowest AIC is retained for the next iteration.

```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.17289 -0.08479  0.01706  0.11497  0.56577
##
## Coefficients:
## (Intercept)      Estimate Std. Error t value Pr(>|t|)
## log(Length)    -5.04406    0.50942   -9.902 1.05e-15 ***
## log(Width)      1.00401    0.15590    6.440 7.34e-09 ***
## log(Thickness)  0.84842    0.19394    4.375 3.51e-05 ***
## Thickness      0.05734    0.02149    2.668 0.00919 **
## NameDarl       -0.09834    0.02044   -1.193 0.23634
## NamePedernales  0.04303    0.08504    0.506 0.61421
## NameTravis      0.18990    0.08586    2.186 0.03325 *
## NameWells       0.07776    0.08787    0.887 0.37768
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.206 on 83 degrees of freedom
## Multiple R-squared:  0.8579, Adjusted R-squared:  0.8469
## F-statistic: 71.59 on 7 and 83 DF, p-value: < 2.2e-16
```



The `rstudent()` function returns studentized residuals, and the `densityPlot()` function fits an adaptive kernel density estimator to the distribution of the studentized residuals. A `qqPlot()` can be used as a check for nonnormal errors, comparing the studentized residuals to a t-distribution. This next function tests for outliers in the regression. This graph displays influence measures in

index plots; Added-variable plots for the regression, looking for influential cases; Removing the 64th and 91th rows:



(b) The 90% prediction interval for a new observation of Weight for a Dart of type Pedernales with Length = 50, Width = 20 and Thickness = 6 is (1.460714, 2.054469).

```
new_obs <- tibble(
  Length = 50,
  Width = 20,
  Thickness = 6,
  Name = c('Pedernales')
)
predict(selectedMod_2, newdata = new_obs, interval = "prediction", level = .9)
```

```
##          fit          lwr          upr
## 1 1.757592 1.509327 2.005857
```

Exercise 2

(a) As a first step, we obtain summary statistics for the dataset wheat; Because the data follow the binomial distribution, the objective is to model the success probability p as a function of the covariates, i.e., to predict the species of the wheat seed based on the measurements of area, perimeter, compactness and asymmetry. We choose logistic regression from a series of generalised linear models.

```
model <- glm(species ~ ., data = wheat, family = "binomial")
model >% summary()
```

```
## Call:
## glm(formula = species ~ ., family = "binomial", data = wheat)
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.42987 -0.07678  0.00023  0.07125  2.37811
## Coefficients:
## (Intercept)      -79.86844    605.46579   -0.132 0.89505
## area           -0.04656    18.66997   -0.002 0.99904
## perimeter       5.97347    39.89812   0.150 0.86099
## compactness    -15.29775    343.34608  -0.045 0.96446
## asymmetry       1.10838    0.42658   2.598 0.00937 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## (Dispersion parameter for binomial family taken to be 1)
## Null deviance: 194.081 on 139 degrees of freedom
## Residual deviance: 32.469 on 135 degrees of freedom
## AIC: 42.469
## Number of Fisher Scoring iterations: 8
```

```
new_seed <- tibble(
  area = 13,
  perimeter=10,
  compactness=0.75,
  asymmetry=2
)
predict(model, newdata = new_seed, type = 'response')
```

```
##          1
## 9.472809e-14
```

(b) The probability that a seed with area = 13, perimeter=10, compactness=0.75, asymmetry=2 is of species Rosa is 9.472809e-14. Perimeter=10 and compactness=0.75 are less than the minimum of these two variables in the data, respectively, which are not used in modelling the logistic regression. This could be harm the confidence of the prediction.