

Simple Parametric Tests
ANALYSING THE SPARROW DATA SET
-
BASIC STATISTICS FOR BIOLOGISTS

Erik Kusch
PhD Student
Aarhus University
Department of Bioscience
Section for Ecoinformatics & Biodiversity
Center for Biodiversity and Dynamics in a Changing World (BIOCHANGE)
Ny Munkegade 116, Building 1540
8000 Aarhus
Denmark
email: erik@i-solution.de

Summary:

Welcome to our sixth practical experience in R. Throughout the following notes, I will introduce you to a couple of simple parametric test. Whilst parametric tests are used extremely often in biological statistics, they can be somewhat challenging to fit to your data as you will see soon.

To do so, I will enlist the sparrow data set we handled in our first exercise. Additionally, today's seminar is showing plotting via base plot instead of `ggplot2` to highlight the usefulness of base plot and show you the base notation.

Contents

1	Preparing Our Procedure	2
1.1	Packages	2
1.2	Loading Data	2
2	t-Test (unpaired)	3
2.1	Testing For Normality And Homogeneity	3
2.2	Climate Warming/Extremes	4
2.2.1	Testing for Normality and Variance	4
2.2.2	Analyses	4
2.3	Sexual Dimorphism	7
2.3.1	Testing for Normality and Variance	7
3	t-Test (paired)	8
3.1	Preparing Data	8
3.2	Testing for Normality	8
3.3	Climate Warming/Extremes	9
4	One-Way ANOVA	11
4.1	Testing For Assumptions	11
4.2	Climate Warming/Extremes	12
4.2.1	Assumption Check	12
4.2.2	Analysis	13
4.3	Predation	15
4.3.1	Assumption Check	15
4.3.2	Analysis	16
5	Two-Way ANOVA	17
5.1	Testing For Assumptions	17
5.2	Sexual Dimorphism	18
5.2.1	Assumption Check	18
5.2.2	Analysis	19
6	ANCOVA	20
6.1	Climate Warming/Extremes	20
6.1.1	Assumption Check	20
6.1.2	Analysis	21
6.2	Sparrow Characteristics And Sites	23
6.2.1	Assumption Check	23
6.2.2	Analysis	24

1. Preparing Our Procedure

To ensure others can reproduce our analysis we run the following three lines of code at the beginning of our R coding file.

```
rm(list = ls()) # clearing environment
Dir.Base <- getwd() # soft-coding our working directory
Dir.Data <- paste(Dir.Base, "Data", sep = "/") # soft-coding our data directory
```

1.1 Packages

Using the following, user-defined function, we install/load all the necessary packages into our current R session.

```
# function to load packages and install them if they haven't been
# installed yet
install.load.package <- function(x) {
  if (!require(x, character.only = TRUE))
    install.packages(x)
  require(x, character.only = TRUE)
}
package_vec <- c("car") # needed for the Levene Test for Homogeneity
sapply(package_vec, install.load.package)
```

```
## car
## TRUE
```

1.2 Loading Data

During our first exercise (Data Mining and Data Handling - Fixing The Sparrow Data Set) we saved our clean data set as an RDS file. To load this, we use the `readRDS()` command that comes with base R.

```
Data_df_base <- readRDS(file = paste(Dir.Data, "/1 - Sparrow_Data_READY.rds",
  sep = ""))
Data_df <- Data_df_base # duplicate and save initial data on a new object
```

2. t-Test (unpaired)

Assumptions of the unpaired t-Test:

- Predictor variable is binary
- Response variable is metric and **normal distributed** within their groups
- Variable values are **independent** (not paired)

In addition, test whether variance of response variable values in groups are equal (`var.test()`) and adjust `t.test()` argument `var.equal` accordingly.

2.1 Testing For Normality And Homogeneity

We need to test the distribution of our response variables within each predictor variable group for their normality and variance. Since this involves two Shapiro tests and one variance test per variable for each response variable, we might want to write our own function to do so:

```
ShapiroTest <- function(Variables, Grouping) {  
  Output <- data.frame(x = Variables)  
  for (i in 1:length(Variables)) {  
  
    X <- Data_df[, Variables[i]]  
    Levels <- levels(Data_df[, Grouping])  
  
    Output[i, 2] <- shapiro.test(X[which(Data_df[, Grouping] == Levels[1])])$p.value  
    Output[i, 3] <- shapiro.test(X[which(Data_df[, Grouping] == Levels[2])])$p.value  
    Output[i, 4] <- var.test(x = X[which(Data_df[, Grouping] == Levels[1])],  
                           y = X[which(Data_df[, Grouping] == Levels[2])])$p.value  
  }  
  colnames(Output) <- c("Variable", "P.value1", "P.value2", "Var.Test")  
  return(Output)  
}
```

This function (`ShapiroTest()`) takes two arguments: (1) `Variables` - a vector of characters holding the names of the variables we want to have tested, and (2) `Grouping` - the binary variable by which to group our variables. The function returns a data frame holding the p-values of the Shapiro tests on each variable group values as well as the `var.test()` p-value.

2.2 Climate Warming/Extremes

Does sparrow morphology change depend on climate?

Using multiple different methods (i.e. Kruskal-Wallis and Mann-Whitney U Test), we have already identified climate (be it in its binary form or when recorded as a three-level variable) is a strong driving force of sparrow morphology. We expect the same results when using a t-Test.

Take note that we need to limit our analysis to our climate type testing sites again as follows (we include Manitoba this time as it is at the same latitude as the UK and Siberia and holds a semi-coastal climate type):

```
# prepare climate type testing data
Data_df <- Data_df_base
Index <- Data_df$Index
Rows <- which(Index == "SI" | Index == "UK" | Index == "RE" | Index ==
  "AU" | Index == "MA")
Data_df <- Data_df[Rows, ]
```

2.2.1 Testing for Normality and Variance

Before we can make use of our data with a t-Test, we need to do an **assumption check**. To this end, we first turn Climate records into a binary variable by turning records of a semi-coastal climate into a coastal one.

```
# Make climate binary
Data_df$Climate[which(Data_df$Climate == "Semi-Coastal")] <- "Coastal"
Data_df$Climate <- droplevels(Data_df$Climate)
```

Let's make sure our assumptions are met:

```
ShapiroTest(Variables = c("Weight", "Height", "Wing.Chord"), Grouping = "Climate")
```

```
##      Variable P.value1 P.value2 Var.Test
## 1      Weight    0.170    0.25    0.3262
## 2      Height    0.168    0.36    0.0106
## 3 Wing.Chord    0.054    0.17    0.0029
```

Luckily, all of our variables allow for the calculation of t-Test. Take note though that some need different specification of the `var.equal` argument than others.

2.2.2 Analyses

Sparrow Weight

Let's start with the weight of *Passer domesticus* individuals as grouped by the climate type present at the site weights have been recorded at:

```
t.test(Data_df$Weight ~ Data_df$Climate, var.equal = TRUE)
```

```
##
## Two Sample t-test
##
## data: Data_df$Weight by Data_df$Climate
## t = -15, df = 381, p-value <2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -2.4 -1.9
## sample estimates:
## mean in group Coastal mean in group Continental
##                      31                      33
```

According to our analysis, which has us **reject the null hypothesis**, we conclude that binary climate records are valuable information criteria for predicting sparrow weight with sparrows in coastal climates being lighter than sparrows in continental ones thus effectively varifying the results of our non-parametric approaches (Kruskal-Wallis, Mann-Whitney U).

Sparrow Height

Let's move on to the height of *Passer domesticus* individuals as grouped by the climate type present at the site weights have been recorded at:

```
t.test(Data_df$Height ~ Data_df$Climate, var.equal = FALSE)

##
##  Welch Two Sample t-test
##
## data:  Data_df$Height by Data_df$Climate
## t = -0.3, df = 366, p-value = 0.8
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -0.23  0.18
## sample estimates:
##      mean in group Coastal mean in group Continental
##                   14                      14
```

Confirming the results of our Mann-Whitney U Test, we **accept the null hypothesis**.

Sparrow Wing Chord

Lastly, we test the wing chords of *Passer domesticus* individuals as grouped by the climate type present at the site weights have been recorded at:

```
t.test(Data_df$Wing.Chord ~ Data_df$Climate, var.equal = FALSE)

##
##  Welch Two Sample t-test
##
## data:  Data_df$Wing.Chord by Data_df$Climate
## t = -0.1, df = 370, p-value = 0.9
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -0.040  0.035
## sample estimates:
##      mean in group Coastal mean in group Continental
##                   6.9                      6.9
```

Without confirming the results of our Mann-Whitney U Test, we **accept the null hypothesis**.

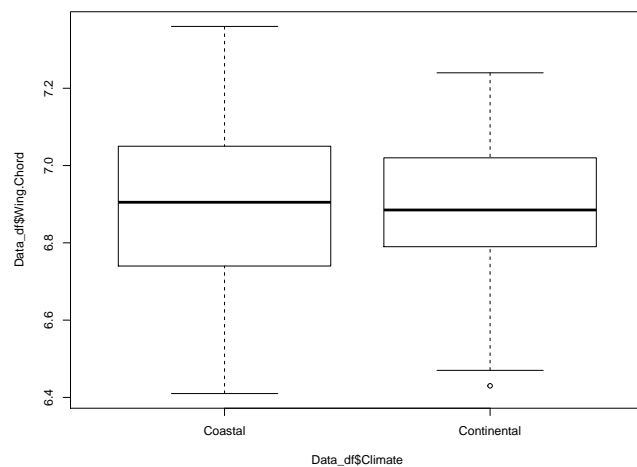
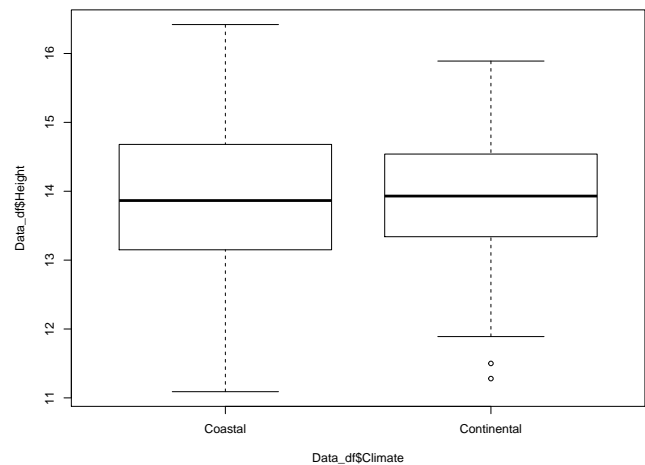
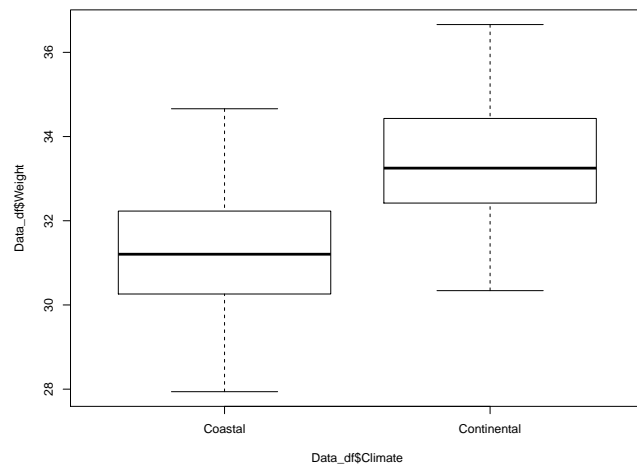
Conclusion

Here's what we've learned from the t-Test so far:

- Sparrow weight depends on (binary) climate types
- Sparrow height does not depend on (binary) climate types
- Sparrow wing chord does not depend on (binary) climate types

Let's end this by visualising all of the data:

```
par(mfrow = c(2, 2))
plot(Data_df$Weight ~ Data_df$Climate)
plot(Data_df$Height ~ Data_df$Climate)
plot(Data_df$Wing.Chord ~ Data_df$Climate)
```



2.3 Sexual Dimorphism

Does sparrow morphology change depend on Sex?

Using the Mann-Whitney U Test, we have already identified the sex of *Passer domesticus* is a good information criterion for understanding sparrow weight but not sparrow height or wing chord. Let's see if we can reproduce this using a t-Test approach.

We may wish to use the entirety of our data set again for this purpose:

```
Data_df <- Data_df_base
```

2.3.1 Testing for Normality and Variance

Again, before we can use our data in a t-Test for this purpose, we have to make sure that our assumptions are met. To this end, we can make use of our user defined `ShapiroTest()` function as follows:

```
ShapiroTest(Variables = c("Weight", "Height", "Wing.Chord"), Grouping = "Sex")
```

```
##      Variable P.value1 P.value2 Var.Test
## 1      Weight  2.9e-21  1.7e-21    0.75
## 2      Height  4.0e-17  6.6e-19    0.40
## 3 Wing.Chord  3.4e-25  1.6e-26    0.56
```

As it turns out, our data does not allow for any t-Test (this happens often in real studies). However, we can create sex-driven subgroups within each site and test whether these meet the requirements for our t-Test. This is out of the scope of this course though and so we will skip it. Spoiler alert: I have done this and the findings did not reveal anything we didn't uncover so far.

3. t-Test (paired)

Assumptions of the paired t-Test:

- Predictor variable is binary
- Response variable is metric
- *Difference of response variable pairs* is **normal distributed**
- Variable values are **dependent** (paired)

3.1 Preparing Data

For this purpose, we need an **additional data set with truly paired records** of sparrows and so we implement the same solution as we've used within our fourth seminar using the Wilcoxon Signed Rank Test. Within our study set-up, think of a **resettling experiment**, where you take *Passer domesticus* individuals from one site, transfer them to another and check back with them after some time has passed to see whether some of their characteristics have changed in their expression.

To this end, presume we have taken the entire *Passer domesticus* population found at our **Siberian** research station and moved them to the **United Kingdom**. Whilst this keeps the latitude stable, the sparrows *now experience a coastal climate instead of a continental one*. After some time (let's say: a year), we have come back and recorded all the characteristics for the same individuals again.

You will find the corresponding *new data* in **2b - Sparrow_ResettledSIUK_READY.rds**. Take note that this set only contains records for the transferred individuals in the **same order** as in the old data set.

```
Data_df_Resettled <- readRDS(file = paste(Dir.Data, "/2b - Sparrow_ResettledSIUK_READY.rds",  
    sep = ""))
```

Since earlier analysis such as the Wilcoxon Signed Rank test (fourth practical) and the Friedman Test (fifth practical) showed that height and wing chord records do not change when sparrows are resettled at all, we have excluded these here and **focus solely on sparrow weight**.

3.2 Testing for Normality

Before being able to run our paired t-Test, we must make sure that the *difference of response variable pairs* is **normal distributed**. We can do so using the `shapiro.test()` of base R as follows:

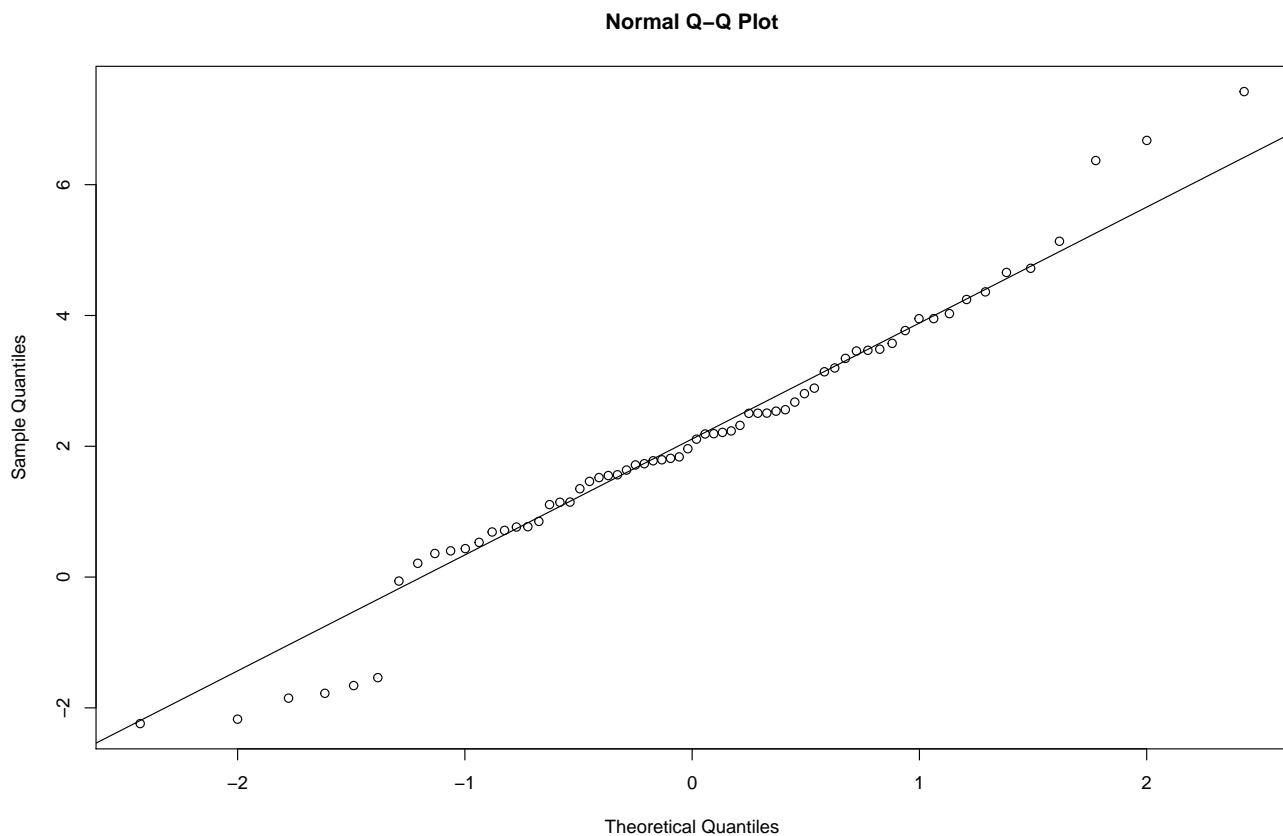
```
# selecting pre-resettling weights  
DataSI <- Data_df$Weight[which(Data_df$Index == "SI")]  
# calculating difference of before and after resettling weights  
WeightDiff <- DataSI - Data_df_Resettled$Weight  
# shapiro test  
shapiro.test(WeightDiff)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data:  WeightDiff  
## W = 1, p-value = 0.2
```

Thankfully, the **assumption of normality** is met.

Now let's visualise that using a qqplot:

```
qqnorm(WeightDiff)
qqline(WeightDiff)
```



3.3 Climate Warming/Extremes

Does sparrow morphology change depend on climate?

Now let's go on to test whether sparrow weights change significantly per individual due to our relocation experiment (we expect this from future test in our practicals):

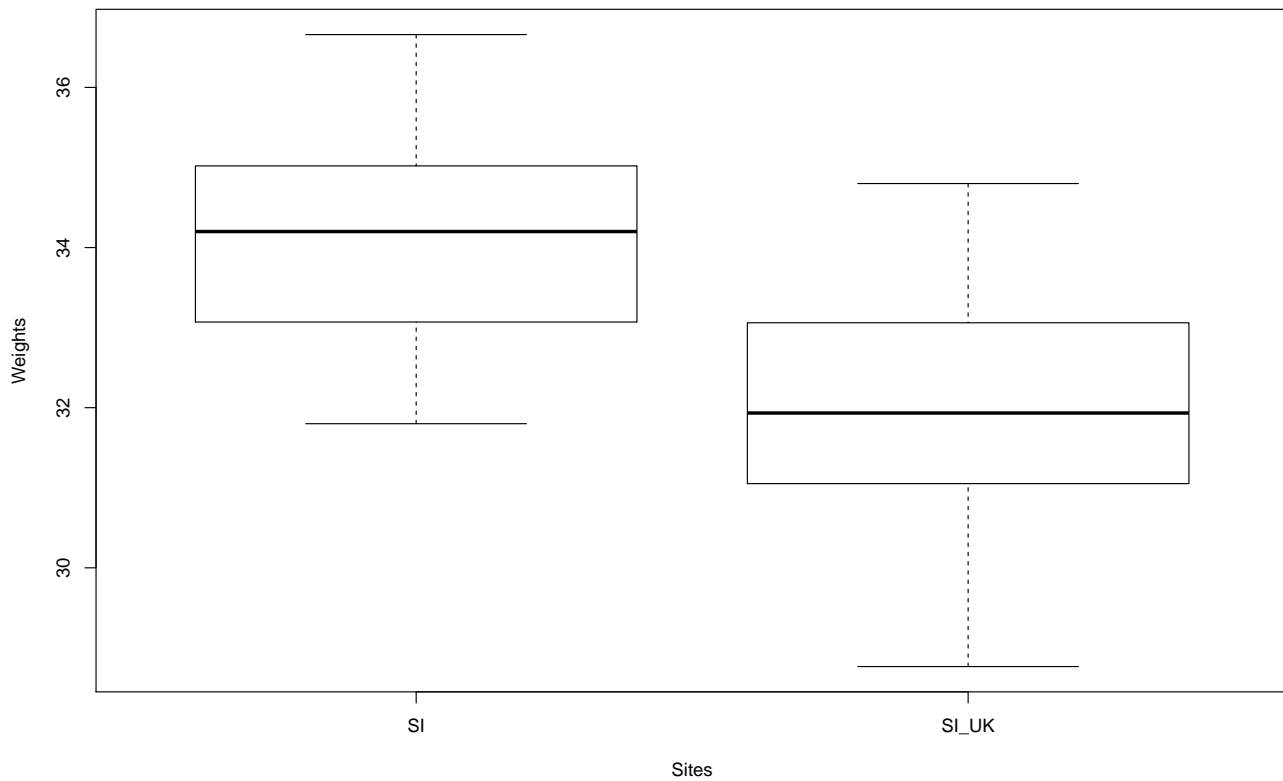
```
t.test(DataSI, Data_df_Resettled$Weight, paired = TRUE)
```

```
##
## Paired t-test
##
## data: DataSI and Data_df_Resettled$Weight
## t = 8, df = 65, p-value = 4e-12
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  1.6 2.6
## sample estimates:
## mean of the differences
##                2.1
```

We were right, individual sparrow weights change significantly after our relocation experiment and we **reject the null hypothesis**. This is in accordance with the results of the Wilcoxon Signed Rank Test as well as the Friedman Test.

Let's go on to visualise our data to make better sense of what is going on here:

```
# Select the sparrow weights
Weights <- c(DataSI, Data_df_Resettled$Weight)
# Select the sites
Sites <- factor(rep(c("SI", "SI_UK"), each = length(DataSI)))
# Plot
plot(Weights ~ Sites)
```



Quite obviously sparrows observed in Siberia are heavier than when they are resettled to the United Kingdom (this may be due to the more forgiving climate in the UK). Just like the test stated, the difference of the average weights is roughly 2g between the sparrows at the two sites.

4. One-Way ANOVA

Assumptions of the One-Way ANOVA:

- Predictor variable is categorical
- Response variable is metric
- *Response variable residuals* are **normal distributed**
- Variance of populations/samples are equal (**homogeneity**)
- Variable values are **independent** (not paired)

4.1 Testing For Assumptions

Firstly, we need to test the assumptions of our One-Way ANOVA. For this purpose, we write another user-defined function.

```
# User-defined function
ANOVACheck <- function(Variables, Grouping, data, plotting) {
  Output <- data.frame(x = NA)
  for (i in 1:length(Variables)) {
    # data
    Y <- as.numeric(data[, Variables[i]])
    X <- data[, Grouping]
    Levels <- levels(data[, Grouping])
    # Residuals?
    model <- lm(Y ~ X)
    Output[1, i] <- shapiro.test(residuals(model))$p.value
    # Homogeneity?
    Levene <- leveneTest(Y ~ X, center = median, data = data)
    Output[2, i] <- Levene[1, 3]
    # Plotting
    if (plotting == TRUE) {
      plot(model, 2) # Normality
      plot(model, 3) # Homogeneity
    }
  }
  colnames(Output) <- Variables
  rownames(Output) <- c("Residual Normality", "Homogeneity of Variances")
  return(Output)
}
```

This ANOVACheck() function takes four arguments: (1) **Variables** - a vector of characters holding the names of the variables we want to have tested, (2) **Grouping** - the categorical variable by which to group our variables, (3) **data** - the data frame which contains the **Variables** and the **Grouping** factor, and (4) **plotting** - a logical indicator of whether to produce plots visualising the test results or not.

The function returns a data frames containing the p-values indexing whether to accept or reject the notion of the normality of residuals per variable (**Residual Normality**), and the p-values indexing whether variances between groups are homogeneous or not (**Homogeneity of Variances**).

4.2 Climate Warming/Extremes

Does sparrow morphology change depend on climate?

Using the Kruskal-Wallis Test in our last exercise, we already identified climate to be an important factor in determining *Passer domesticus* morphology. Let's see if this holds true.

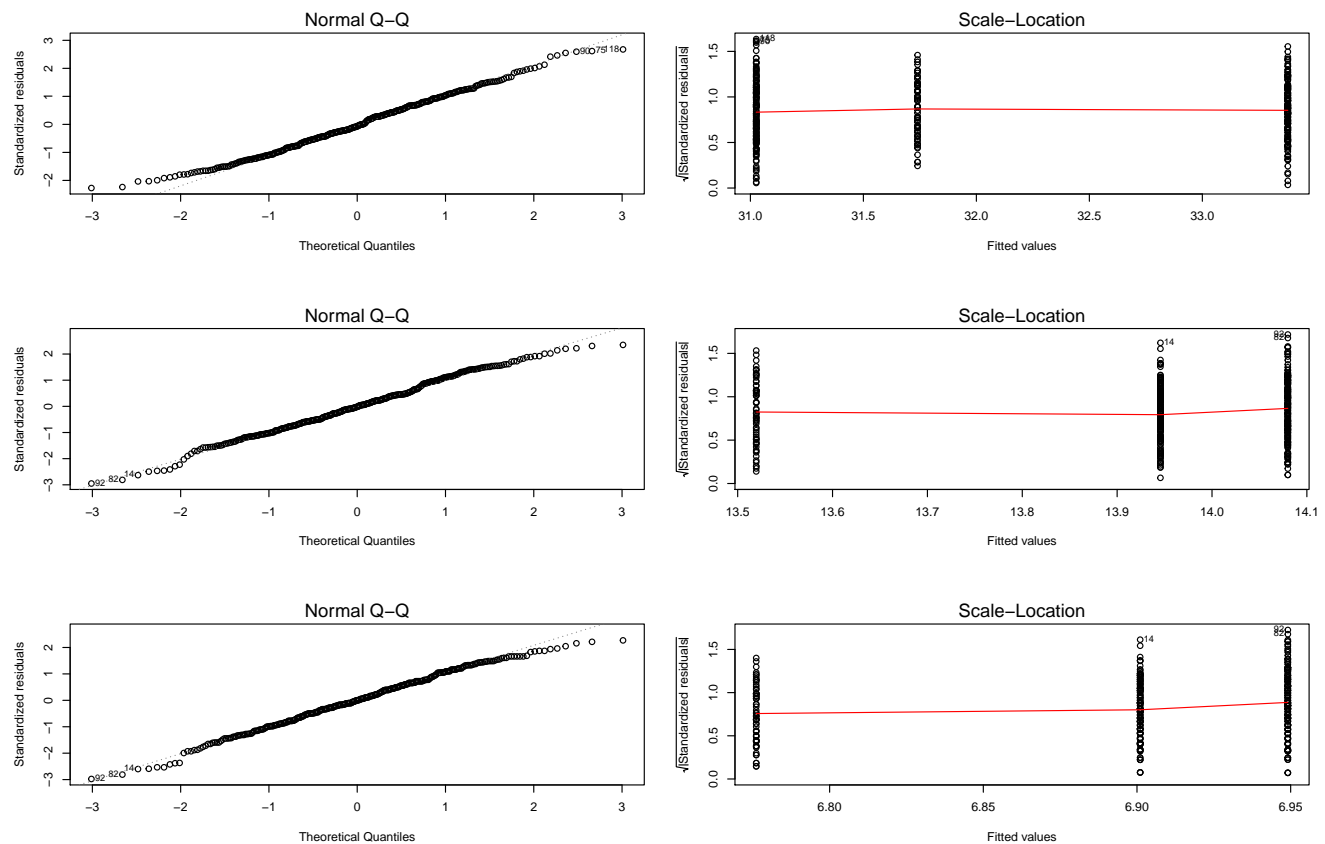
Take note that we need to limit our analysis to our climate type testing sites again as follows (we include Manitoba this time as it is at the same latitude as the UK and Siberia and holds a semi-coastal climate type):

```
# prepare climate type testing data
Data_df <- Data_df_base
Index <- Data_df$Index
Rows <- which(Index == "SI" | Index == "UK" | Index == "RE" | Index ==
  "AU" | Index == "MA")
Data_df <- Data_df[Rows, ]
```

4.2.1 Assumption Check

Let's use the ANOVACheck() function on our data:

```
par(mfrow = c(3, 2))
ANOVACheck(Variables = c("Weight", "Height", "Wing.Chord"), Grouping = "Climate",
  data = Data_df, plotting = TRUE)
```



```
##                               Weight Height Wing.Chord
## Residual Normality           0.045  0.099    0.0516
## Homogeneity of Variances      0.961  0.091    0.0064
```

Unfortunately, neither weight nor wing chord records fulfil our requirements.

4.2.2 Analysis

Let's run our analysis for height as grouped by the three-level climate variable:

```
model <- lm(Data_df$Height ~ Data_df$Climate)
anova(model)
```

```
## Analysis of Variance Table
##
## Response: Data_df$Height
##           Df Sum Sq Mean Sq F value    Pr(>F)
## Data_df$Climate      2      15      7.49      7.25 0.00081 ***
## Residuals          381      394      1.03
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

According to this, climate is a meaningful predictor of height of sparrows and we **reject the null hypothesis** thus confirming the results of our Kruskal-Wallis analysis.

Now, let's analyse the output a bit more in-depth:

```
summary(model)
```

```
##
## Call:
## lm(formula = Data_df$Height ~ Data_df$Climate)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.9899 -0.6981 -0.0147  0.6714  2.3704
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      14.0799     0.0796  176.80 < 2e-16 ***
## Data_df$ClimateContinental -0.1343     0.1143   -1.18  0.24060
## Data_df$ClimateSemi-Coastal -0.5604     0.1476   -3.80  0.00017 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1 on 381 degrees of freedom
## Multiple R-squared:  0.0367, Adjusted R-squared:  0.0316
## F-statistic: 7.25 on 2 and 381 DF, p-value: 0.000813
```

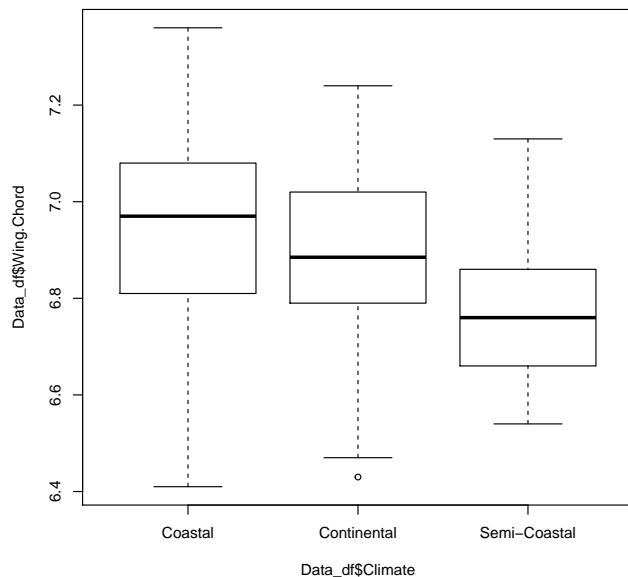
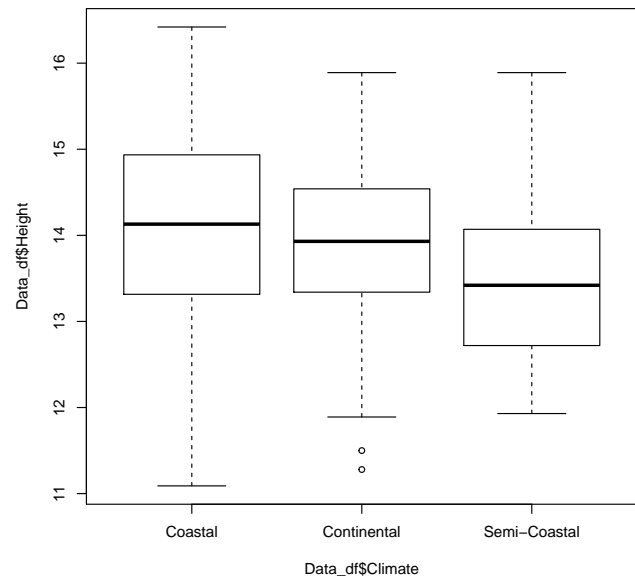
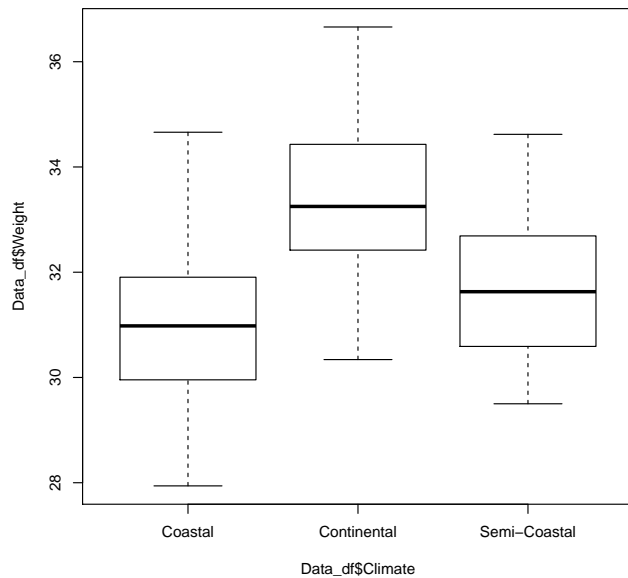
- The mean sparrow height in coastal climates is 14.08cm (this is our **Intercept**/*Baseline*)
- The mean sparrow height in continental climates is -0.13cm bigger than the **Intercept**
- The mean sparrow height in semi-coastal climates is -0.56cm bigger than the **Intercept**
- Only the estimates in coastal and semi-coastal climates are statistically significant

Personally, I would not place too much confidence in these results due to a couple of reasons:

- Our only semi-coastal site is on the northern hemisphere whereas two of our stations are located in the southern hemisphere
- Confounding factors such as population status might have an effect which we are not considering here

Let's end this by plotting all of our data:

```
par(mfrow = c(2, 2))
plot(Data_df$Weight ~ Data_df$Climate)
plot(Data_df$Height ~ Data_df$Climate)
plot(Data_df$Wing.Chord ~ Data_df$Climate)
```



As you can see, the variances are definitely not equal between our groups which explains why part of our assumption test failed.

4.3 Predation

Does nesting height depend on predator characteristics?

Again, using the Kruskal-Wallis Test in our last exercise, we already identified predator characteristics to be an important factor in determining *Passer domesticus* nesting height. Let's see if this holds true.

We may wish to use the entirety of our data set again for this purpose:

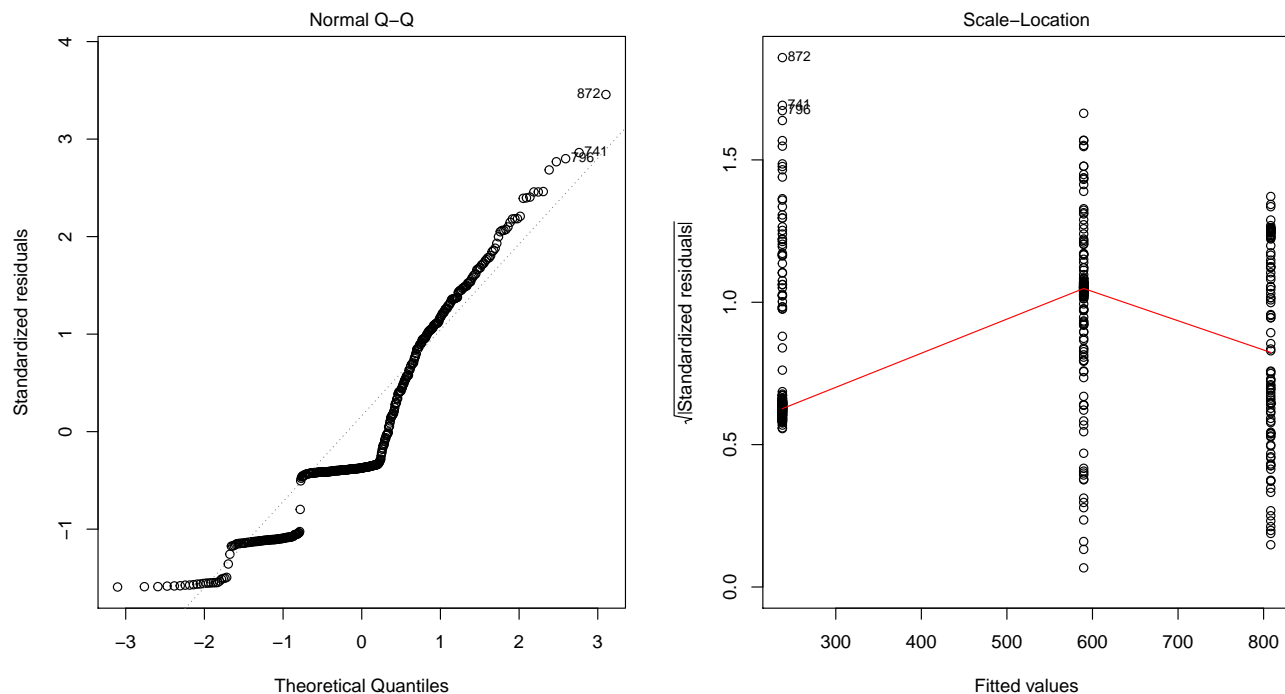
```
Data_df <- Data_df_base
```

4.3.1 Assumption Check

Let's use our `ANOVACheck()` function to test whether we can run our analysis. Before we can do so, however, we need to slightly adjust our predator type variable just like we did in our last exercise and as follows:

```
# changing levels in predator type
levels(Data_df$Predator.Type) <- c(levels(Data_df$Predator.Type), "None")
Data_df$Predator.Type[which(is.na(Data_df$Predator.Type))] <- "None"

# Assumption Check
par(mfrow = c(1, 2))
ANOVACheck(Variables = "Nesting.Height", Grouping = "Predator.Type", data = Data_df,
            plotting = TRUE)
```



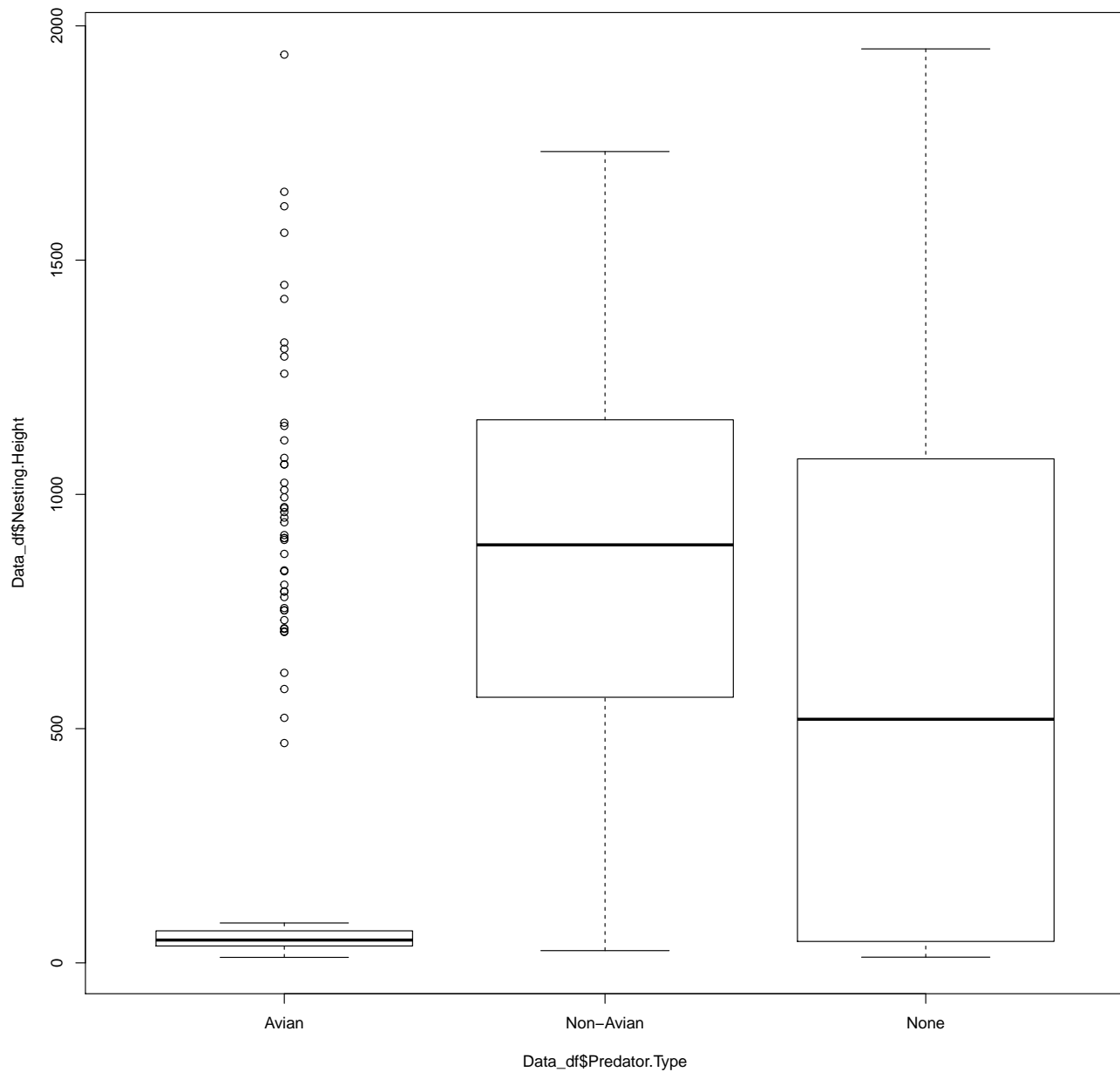
```
##                               Nesting.Height
## Residual Normality            1.2e-15
## Homogeneity of Variances      2.5e-20
```

Again, our data fails the assumption check. The residuals are definitely not normal distributed and the variance of nesting height records within our groups are not equal.

4.3.2 Analysis

Since none of our assumptions are met, we cannot run an ANOVA and therefore resort to data visualisation alone:

```
plot(Data_df$Nesting.Height ~ Data_df$Predator.Type)
```



Once more, we can see why our homogeneity of variances test failed.

5. Two-Way ANOVA

Assumptions of the Two-Way ANOVA:

- Predictor variables are categorical
- Response variable is metric
- *Response variable residuals* are **normal distributed**
- Variance of populations/samples are equal (**homogeneity**)
- Variable values are **independent** (not paired)

5.1 Testing For Assumptions

Yet again, we need to check if our assumptions are met first. Automating this procedure is definitely a good idea and only needs slight modification from our ANOVACheck() function.

```
# User-defined function
ANOVACheck_TWO <- function(Formulas, data, plotting) {
  Output <- data.frame(x = NA)
  for (i in 1:length(Formulas)) {
    # Check how many formulas there are
    if (length(Formulas) == 1) {
      Expression <- Formulas[[1]]
    } else {
      Expression <- Formulas[[i]]
    }
    # Residuals?
    model <- lm(formula = Expression, data = data)
    Output[1, i] <- shapiro.test(residuals(model))$p.value
    # Homogeneity?
    Levene <- leveneTest(Expression, center = median, data = data)
    Output[2, i] <- Levene[1, 3]
    # Plotting
    if (plotting == TRUE) {
      plot(model, 2) # Normality
      plot(model, 3) # Homogeneity
    }
  }
  colnames(Output) <- as.character(Formulas)
  rownames(Output) <- c("RN", "HoV")
  return(Output)
}
```

This ANOVACheck_TWO() function takes four arguments: (1) **Formulas** - a vector of formula specification for our ANOVA models we want to have tested, (2) **data** - the data frame which contains the variables and the grouping factor called upon in our **Formulas**, and (3) **plotting** - a logical indicator of whether to produce plots visualising the test results or not.

The function returns a data frames containing the p-values indexing whether to accept or reject the notion of the normality of residuals per variable (RN), and the p-values indexing whether variances between groups are homogeneous or not (HoV).

5.2 Sexual Dimorphism

Does sparrow morphology depend on population status and sex?

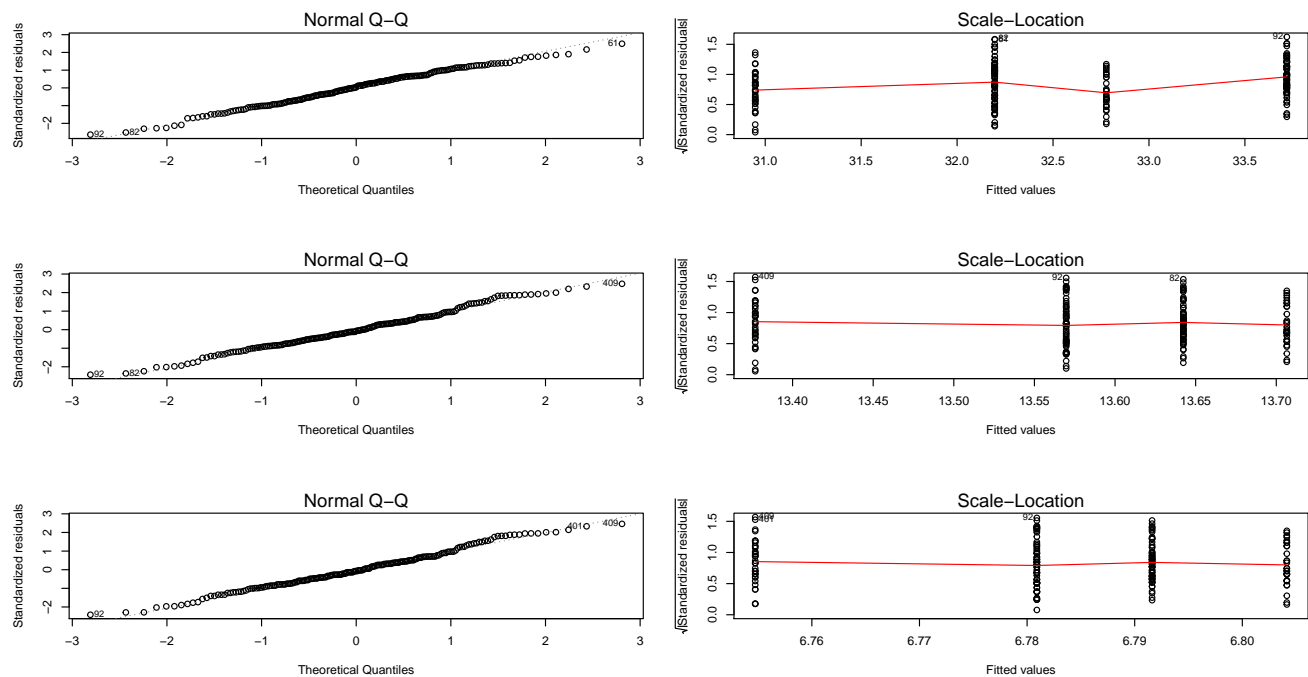
Given different factors affecting invasive species, we might expect different patterns of sexual dimorphism for invasive and native populations. Take note that we keep using the northern hemisphere subset of our climate testing sites as these present us with a nice set of invasive/native population records already whilst keeping confounding factors to a minimum.

5.2.1 Assumption Check

First, we need to check our assumptions:

```
# prepare climate type testing data
Data_df <- Data_df_base
Index <- Data_df$Index
Rows <- which(Index == "SI" | Index == "UK" | Index == "MA")
Data_df <- Data_df[Rows, ]

# analysis
par(mfrow = c(3, 2))
ANOVACheck_TWO(Formulas = c(Weight ~ Population.Status * Sex, Height ~
  Population.Status * Sex, Wing.Chord ~ Population.Status * Sex), data = Data_df,
  plotting = TRUE)
```



```
##      Weight ~ Population.Status * Sex Height ~ Population.Status * Sex
## RN              0.2880                      0.22
## HoV              0.0045                      0.91
##      Wing.Chord ~ Population.Status * Sex
## RN              0.19
## HoV              0.92
```

Again our assumptions are not met except for sparrow height and wing chord as a product of sex and population status.

5.2.2 Analysis

Let's run our analysis:

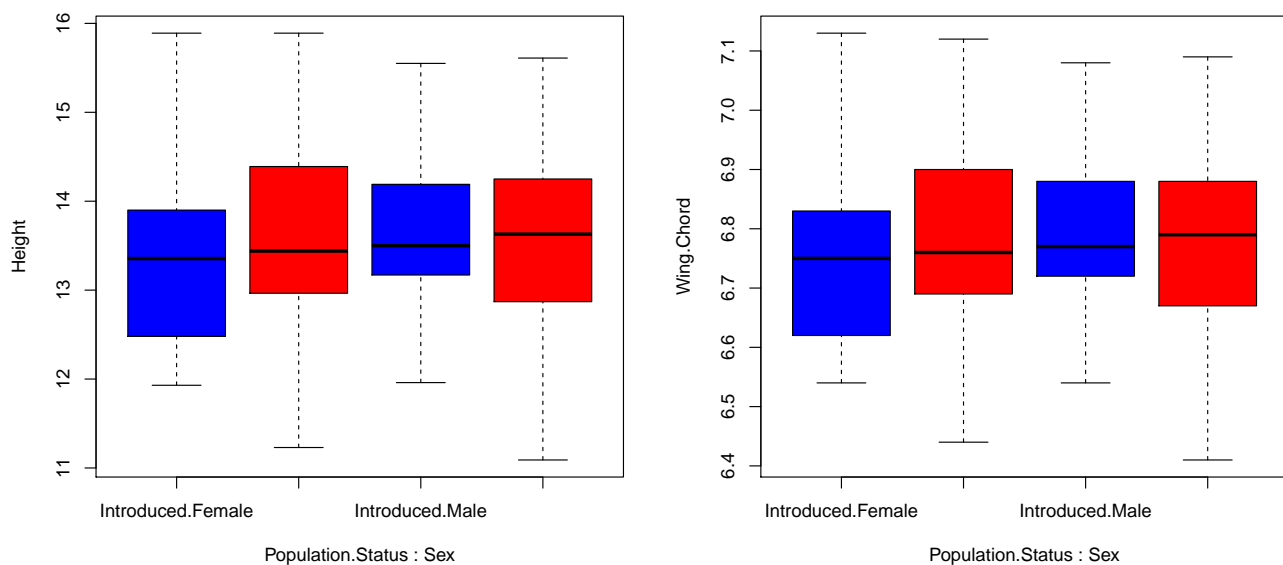
```
# height model
model <- lm(Height ~ Population.Status * Sex, data = Data_df)
anova(model)

## Analysis of Variance Table
##
## Response: Height
##
##              Df Sum Sq Mean Sq F value Pr(>F)
## Population.Status      1      0.3    0.338    0.32   0.57
## Sex                    1      0.2    0.179    0.17   0.68
## Population.Status:Sex    1      1.8    1.786    1.68   0.20
## Residuals             197    208.9    1.060

# wing chord model
model <- lm(Wing.Chord ~ Population.Status * Sex, data = Data_df)
anova(model)

## Analysis of Variance Table
##
## Response: Wing.Chord
##
##              Df Sum Sq Mean Sq F value Pr(>F)
## Population.Status      1    0.00  0.0047    0.20   0.66
## Sex                    1    0.00  0.0041    0.17   0.68
## Population.Status:Sex    1    0.04  0.0399    1.67   0.20
## Residuals             197    4.70  0.0239

# plotting
par(mfrow = c(1, 2))
boxplot(Height ~ Population.Status * Sex, data = Data_df, col = c("blue",
  "red"))
boxplot(Wing.Chord ~ Population.Status * Sex, data = Data_df, col = c("blue",
  "red"))
```



As it turns out, population status and sex are no viable predictors for sparrow height or wing chord and so we **accept the null hypothesis**.

6. ANCOVA

Assumptions of the ANCOVA:

- Predictor variables are categorical or continuous
- Response variable is metric
- *Response variable residuals* are **normal distributed**
- Variance of populations/samples are equal (**homogeneity**)
- Variable values are **independent** (not paired)
- Relationship between the response and covariate is **linear**.

6.1 Climate Warming/Extremes

Do sparrow characteristics depend on climate and latitude?

Latitude may have masked some effects of climate on sparrow morphology in our preceding analyses and vice-versa. At times, we have been able to account for this by including our site records, which can be seen as binned versions of latitude records. Let's test if the inclusion of raw latitude records are meaningful.

6.1.1 Assumption Check

Again, we need to do an assumption check. However, we need a new function for this, since we now need to test whether our response variable and the covariate are linear or not:

```
# overwriting prior changes in Data_df
Data_df <- Data_df_base
Data_df$Latitude <- abs(Data_df$Latitude)
# User-defined function
ANCOVACheck <- function(Variables, Grouping, Covariate, data, plotting) {
  Output <- data.frame(x = NA)
  for (i in 1:length(Variables)) {
    # data
    Y <- as.numeric(data[, Variables[i]])
    X <- data[, Grouping]
    Z <- data[, Covariate]
    # Residuals?
    model <- lm(Y ~ X * Z)
    Output[1, i] <- shapiro.test(residuals(model))$p.value
    # Homogeneity?
    Levene <- leveneTest(Y ~ X, center = median, data = data)
    Output[2, i] <- Levene[1, 3]
    # Plotting
    if (plotting == TRUE) {
      plot(model, 1) # Linearity
      plot(model, 2) # Normality
      plot(model, 3) # Homogeneity
    }
  }
  colnames(Output) <- Variables
  rownames(Output) <- c("RN", "HoV")
}
```

```

    return(Output)
}

```

This `ANCOVACheck()` function takes five arguments: (1) **Variables** - a vector of response variables used in our models, (2) **Grouping** - the categorical variable by which to group our variables, (3) **Covariate** - the covariate of our analysis, (4) **data** - the data frame which contains the variables, the grouping factor and our covariate, and (5) **plotting** - a logical indicator of whether to produce plots visualising the test results or not.

The function returns a data frames containing the p-values indexing whether to accept or reject the notion of the normality of residuals per variable (RN), and the p-values indexing whether variances between groups are homogeneous or not (HoV).

```

ANCOVACheck(Variables = c("Weight", "Height", "Wing.Chord", "Nesting.Height",
  "Egg.Weight", "Number.of.Eggs", "Home.Range"), Grouping = "Climate",
  Covariate = "Latitude", data = Data_df, plotting = FALSE)

```

```

##      Weight  Height Wing.Chord Nesting.Height Egg.Weight Number.of.Eggs Home.Range
## RN  6.2e-10 9.4e-05   1.8e-16    5.5e-18      7e-02      9.6e-14    2.9e-20
## HoV 5.2e-24 1.1e-24   3.2e-28    4.0e-01      1e-06      1.3e-13    1.2e-08

```

The assumptions aren't met. I have set the `plotting` argument to `FALSE` to suppress the plotting of model checking visualisation. This would be useful to judge linearity but not necessary here since the other two important assumptions (Homogeneity of variances and Normality of residuals) aren't met to begin with.

6.1.2 Analysis

Since none of our assumptions are met, we cannot run an ANOVA and therefore resort to data visualisation alone. We need a new function for this to do our plotting easily and automatically with some colours indicating our grouping factors whilst plotting response variables versus covariates.

```

PlotAncovas <- function(Variables, Grouping, Covariate, data){
  for(i in 1:length(Variables)){
    Y <- Data_df[,Variables[i]]
    X <- Data_df[,Covariate]
    G <- Data_df[, Grouping]
    plot(X, Y, col = G, xlab = Covariate, ylab = Variables[i])
    legend("top", # place legend at the top
      inset = -0.35, # move legend away from plot centre
      xpd = TRUE, # allow legend outside of plot area
      legend=levels(G), # what to include in legend
      bg = "white", col = unique(G), ncol=length(levels(G)), # colours
      pch = 1, # plotting symbols
      title = Variables[i] # title of legend
    )
  }
}

```

The `PlotAncovas()` returns a scatter plot and takes four arguments: (1) **Variables** - a vector of response variables, (2) **Grouping** - the name of the grouping factor according to which to colour the symbols in our plot, (3) **Covariate** - the covariate against which to plot individual variables, and (4) **data** - the data frame which holds our variables.

Let's use our function:

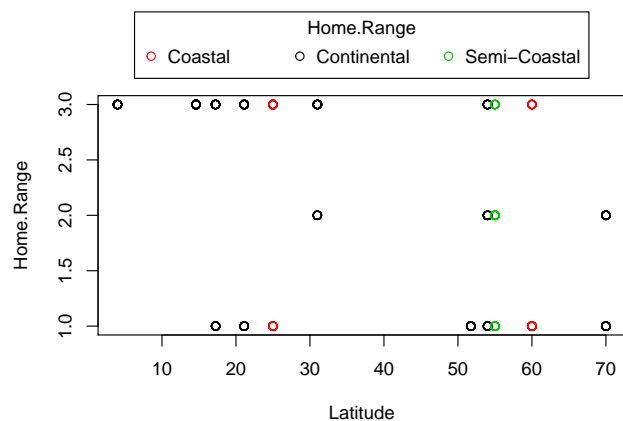
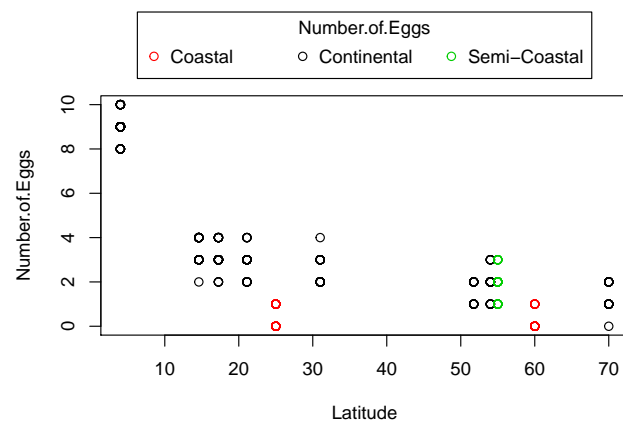
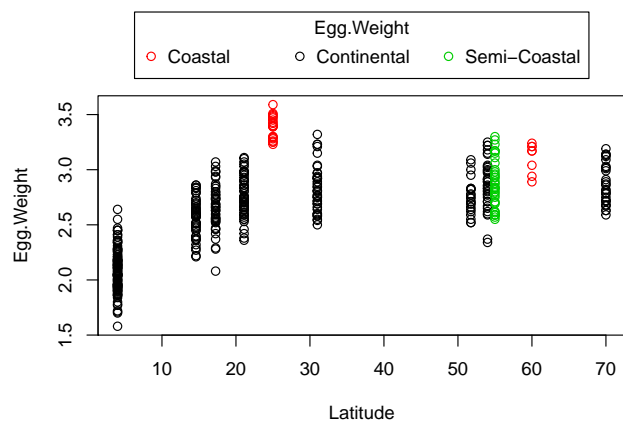
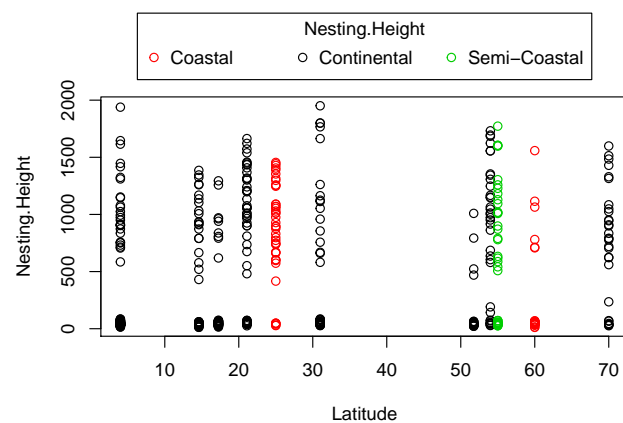
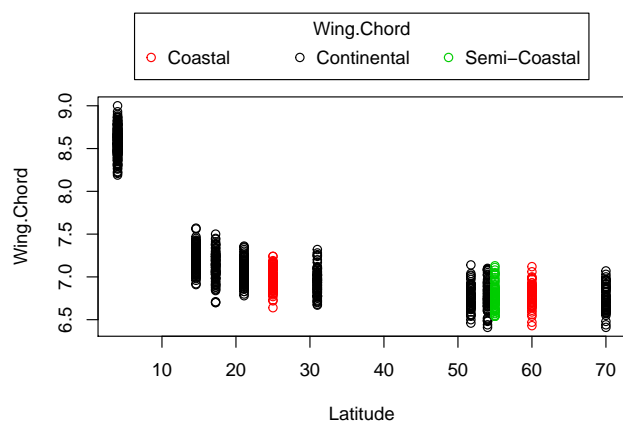
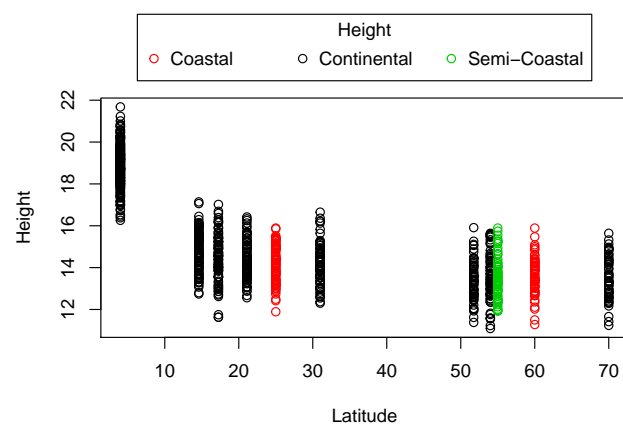
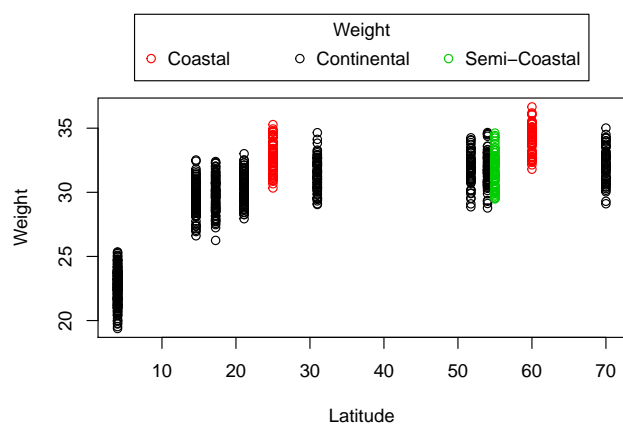
```

par(mfrow = c(1, 2))
PlotAncovas(Variables = c("Weight", "Height", "Wing.Chord", "Nesting.Height",
  "Egg.Weight", "Number.of.Eggs", "Home.Range"), Grouping = "Climate",
  Covariate = "Latitude", data = Data_df)

```

I will not interpret these plots here in text and leave this to you.

Take note that this **could've been achieved much easier with ggplot2!**



6.2 Sparrow Characteristics And Sites

This was not part of what we set out to do according to the lecture slides but has been included as a logical conclusion to an earlier analysis.

Unfortunately, our previous attempt at an ANCOVA didn't work. So what other covariate do we have available for sparrow characteristics?

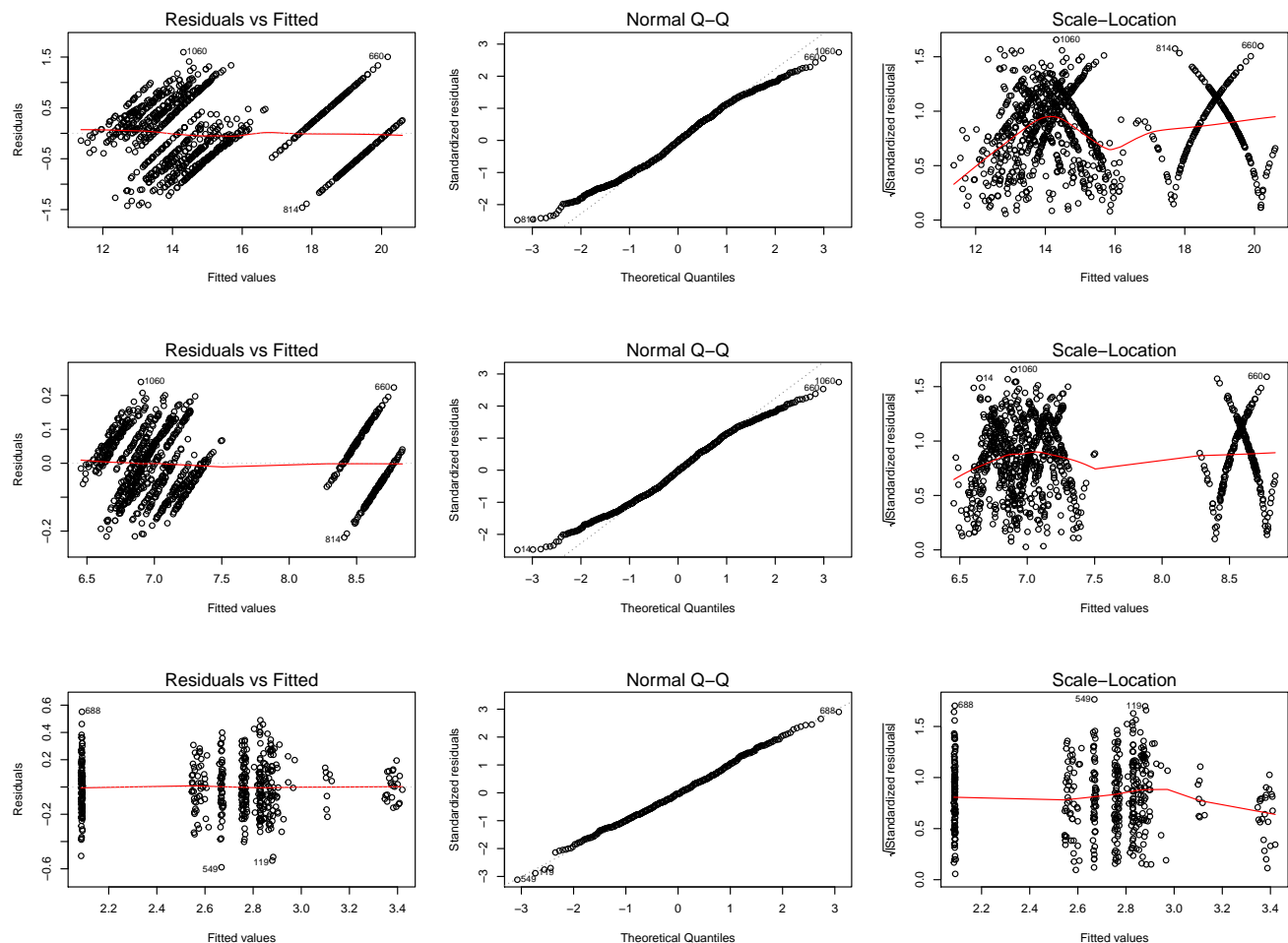
- *Latitude* doesn't make sense to include when grouping by site index as these two are synonymous - *Longitude* doesn't make sense to include when grouping by site index as these two are synonymous - *Weight* is well explained by other variables and we know the causal links - *Height* is not that well explained by other variables - *Wing.Chord* is not that well explained by other variables

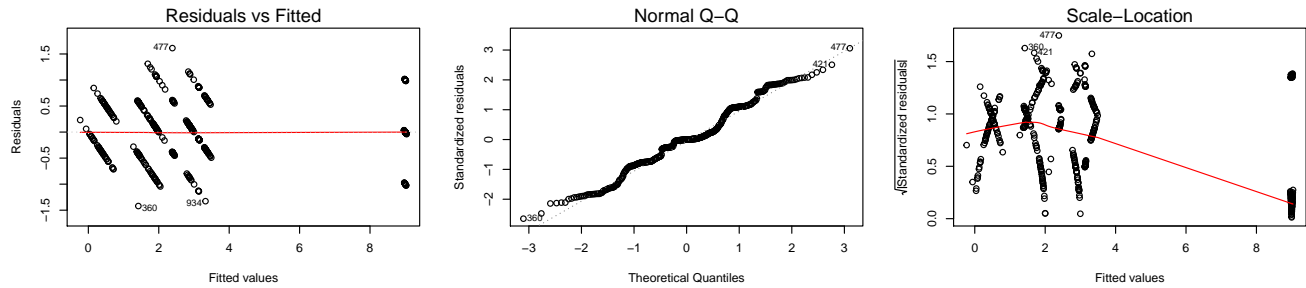
Of course, there are more within our data set but it has become apparent that **Weight** may make for an important covariate in our site-wise ANCOVA set-up. Using the Pearson correlation (third practical), we already identified a causal link between sparrow **Weight** and **Height** per site.

6.2.1 Assumption Check

Firstly, we test whether assumptions are met. For brevities sake, we only test four variables:

```
par(mfrow = c(1, 3))
ANCOVACheck(Variables = c("Height", "Wing.Chord", "Egg.Weight", "Number.of.Eggs"),
             Grouping = "Index", Covariate = "Weight", data = Data_df, plotting = TRUE)
```





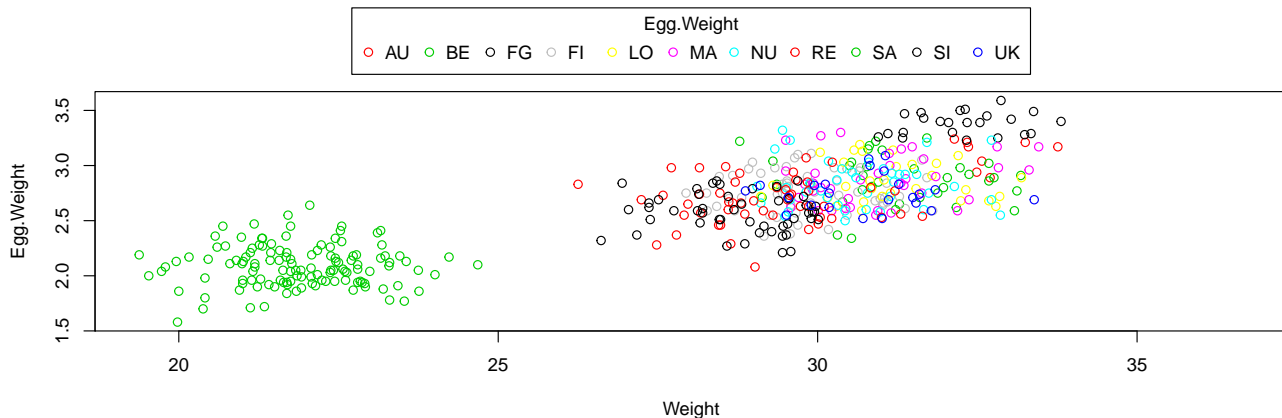
```
##      Height Wing.Chord Egg.Weight Number.of.Eggs
## RN   1.2e-07   1.2e-07   0.512      9.2e-07
## HoV   2.6e-01   2.5e-01   0.085      2.7e-02
```

As it turns out, we can run our ANCOVA on `Egg.Weight` when grouped by site `Index` and driven by `Weight`.

6.2.2 Analysis

First, let's visualise our data:

```
PlotAncovas(Variables = "Egg.Weight", Grouping = "Index", Covariate = "Weight",
  data = Data_df)
```



Quite obviously, Belize (BE) records are very different from the other stations, whose egg weight and weight records are grouped together. There seems to be some evidence for an overall linkage of sparrow weight and egg weight (a positive correlation).

Now we run the analysis:

```
LM_fit5 <- lm(Egg.Weight ~ Weight * Index, data = Data_df)
anova(LM_fit5)
```

```
## Analysis of Variance Table
##
## Response: Egg.Weight
##          Df Sum Sq Mean Sq F value Pr(>F)
## Weight      1   52.5    52.5 1442.56 <2e-16 ***
## Index     10    8.1     0.8   22.21 <2e-16 ***
## Weight:Index 10    0.1     0.0    0.35  0.97
## Residuals 455   16.6     0.0
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The above ANCOVA output tells us that there is no interaction effect between sites and sparrow weights when determining mean egg weight per nest of *Passer domesticus* and so we do another iteration of our model and remove

the postulated interaction:

```
LM_fit6 <- lm(Egg.Weight ~ Weight + Index, data = Data_df)
anova(LM_fit6)
```

```
## Analysis of Variance Table
##
## Response: Egg.Weight
##           Df Sum Sq Mean Sq F value Pr(>F)
## Weight      1   52.5    52.5   1462.9 <2e-16 ***
## Index      10    8.1     0.8    22.5 <2e-16 ***
## Residuals 465   16.7     0.0
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

By now, all of our model coefficients are significant and we can go on to interpret them:

```
summary(LM_fit6)

##
## Call:
## lm(formula = Egg.Weight ~ Weight + Index, data = Data_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.5889 -0.1315 -0.0062  0.1203  0.5514
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.34598    0.28083   11.91 < 2e-16 ***
## Weight        0.00108    0.00861    0.13  0.90020
## IndexBE      -0.70848    0.05486  -12.91 < 2e-16 ***
## IndexFG      -1.28117    0.09914  -12.92 < 2e-16 ***
## IndexFI      -0.62529    0.05744  -10.89 < 2e-16 ***
## IndexLO      -0.55014    0.05175  -10.63 < 2e-16 ***
## IndexMA      -0.51365    0.05135  -10.00 < 2e-16 ***
## IndexNU      -0.51702    0.05337   -9.69 < 2e-16 ***
## IndexRE      -0.61263    0.05142  -11.91 < 2e-16 ***
## IndexSA      -0.80605    0.05669  -14.22 < 2e-16 ***
## IndexSI      -0.27258    0.07787   -3.50  0.00051 ***
## IndexUK      -0.51167    0.05140   -9.95 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.19 on 465 degrees of freedom
## (590 observations deleted due to missingness)
## Multiple R-squared:  0.784, Adjusted R-squared:  0.779
## F-statistic: 153 on 11 and 465 DF, p-value: <2e-16
```