

# STAT-615 Project I: Data Analysis

```
library(tidyverse)
```

## Introduction

```
#Load Data

"britain_species.dat" %>%
  read_delim(delim = " ", col_names = FALSE) %>%
  rename(island = X1,
         area = X2,
         elevation = X3,
         soil = X4,
         latitude = X5,
         distance = X6,
         species = X7) %>%
  mutate(area = as.double(area),
         elevation = as.integer(elevation),
         soil = as.integer(soil),
         latitude = as.double(latitude),
         distance = as.double(distance),
         species = as.integer(species)) -> species

species
```

```
## # A tibble: 42 x 7
##   island      area elevation  soil latitude distance species
##   <chr>      <dbl>    <int> <int>    <dbl>    <dbl>    <int>
## 1 Ailsa        0.8      340     1     55.3      14       75
## 2 Anglesey    712.       127     3     53.3       0.2     855
## 3 Arran      429.       874     4     55.6       5.2     577
## 4 Barra      18.4       384     2      57      77.4     409
## 5 Bressay     31.1       226     1     60.1     202.     177
## 6 Britain  229850.    1343    16     54.3       0    1666
## 7 Canna      12.7       210     1     57.1     40.6     300
## 8 Coll       74.1       103     3     56.6     14.5     443
## 9 Colonsay   44.8       143     1     56.1     31.1     482
## 10 Eigg       29        393     1     56.9     12.3     453
## # ... with 32 more rows
```

For our data analysis project, we were interested in exploring the diversity of species. The importance of this question subject is nontrivial, as preserving species diversity is incredibly important to prevent extinction of them. By finding out what factors are related or responsible for increased diversity, we can obtain general awareness of them, as well as understand their role in diversity and how to manage them properly to preserve variety of species.

Because this is a very broad subject, we decided to narrow in on a particular dataset. The University of Florida contains a data repository which contained a dataset which included information about bird species diversity in the islands, as well as mainland Britain.

The following variables are found in this dataset: - island: name of the island - area: measured in squared kilometers - elevation: highest peak, measured in meters - soil: number of different soil types - latitude - distance: from mainland britain - species: total number of bird species

```
head(species, 5)
```

```
## # A tibble: 5 x 7
##   island    area elevation  soil latitude distance species
##   <chr>    <dbl>    <int> <int>    <dbl>    <dbl>    <int>
## 1 Ailsa      0.8      340     1     55.3      14      75
## 2 Anglesey 712.      127     3     53.3      0.2     855
## 3 Arran    429.      874     4     55.6      5.2     577
## 4 Barra    18.4      384     2      57      77.4     409
## 5 Bressay  31.1      226     1     60.1     202.     177
```

```
ncol(species)
```

```
## [1] 7
```

```
nrow(species)
```

```
## [1] 42
```

The dataset contains 7 variables, and 42 observations. Although there is no missing data it should be noted that a limitation of this dataset and following analyses is that it is not include a particularly large number of observations, particularly with respect to the number of variables considered.

## Limitations

Since `latitude` is not particularly useful, it will not be used in any analyses. This will place the observations to variable ratio at 42:6. However this would only be the case for considered models in which all other 6 variables are included. Some limitations of the small dataset are that it does not produce strong statistical power, thus, any conclusions, null or alternative, should be interpreted with caution.

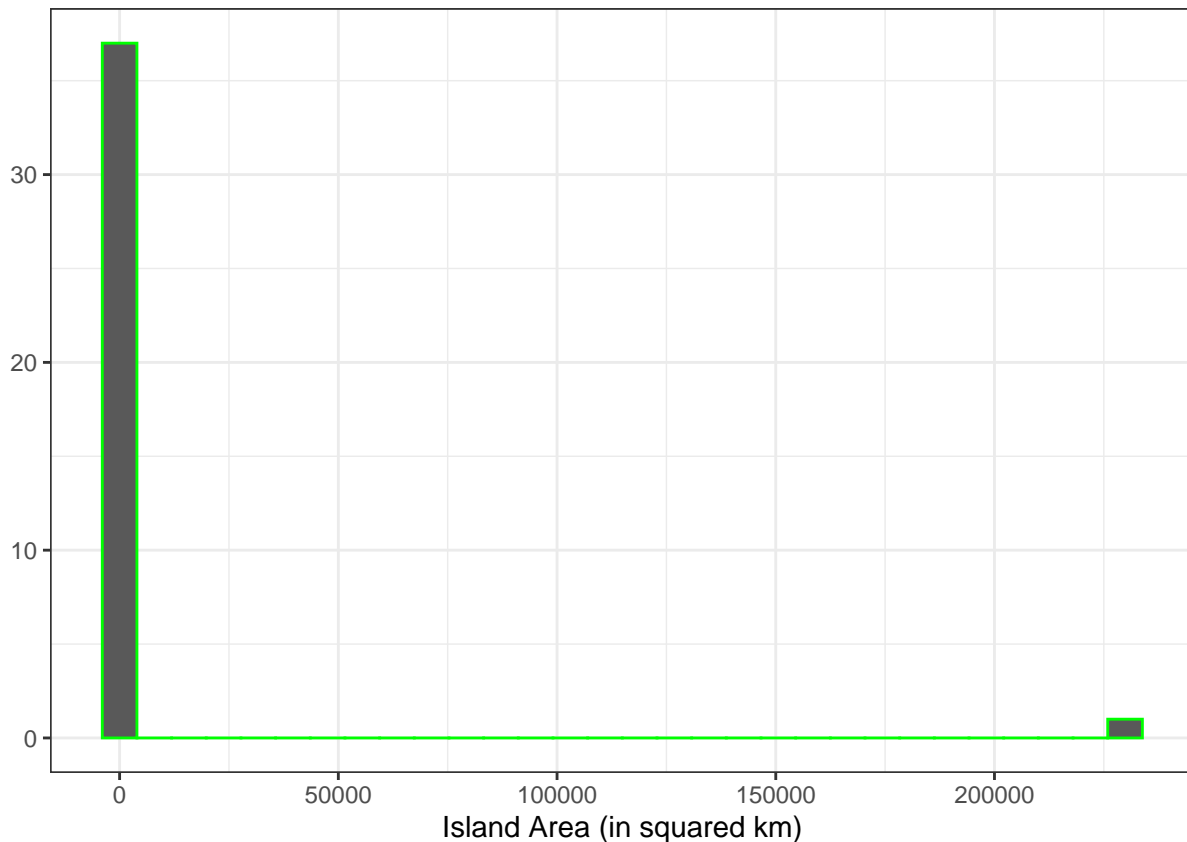
## Preliminary Data Analysis

```
summary(species)
```

```
##      island          area          elevation          soil
## Length:42      Min.   :    0.5      Min.   :    7.0      Min.   :    1.00
## Class :character 1st Qu.:   14.0      1st Qu.: 129.5      1st Qu.:    1.00
## Mode  :character Median :   52.8      Median : 232.0      Median :    2.00
##              Mean   : 6324.9      Mean   : 346.4      Mean   :   28.62
##              3rd Qu.:  417.2      3rd Qu.: 442.0      3rd Qu.:    3.75
##              Max.   :229849.8      Max.   :1343.0      Max.   :   620.00
##              NA's   :4
##      latitude      distance      species
## Min.   : 2.00      Min.   : 0.000      Min.   :    9.0
## 1st Qu.:55.62      1st Qu.:  9.825      1st Qu.: 159.5
## Median :56.80      Median : 33.950      Median : 346.0
## Mean   :51.97      Mean   : 63.119      Mean   : 368.2
## 3rd Qu.:59.08      3rd Qu.: 65.300      3rd Qu.: 450.8
```

```
## Max. :60.80 Max. :258.100 Max. :1666.0
##
```

```
ggplot(species, aes(x = area)) +
  geom_histogram(color = "green") +
  theme_bw() +
  ylab("") +
  xlab("Island Area (in squared km)")
```



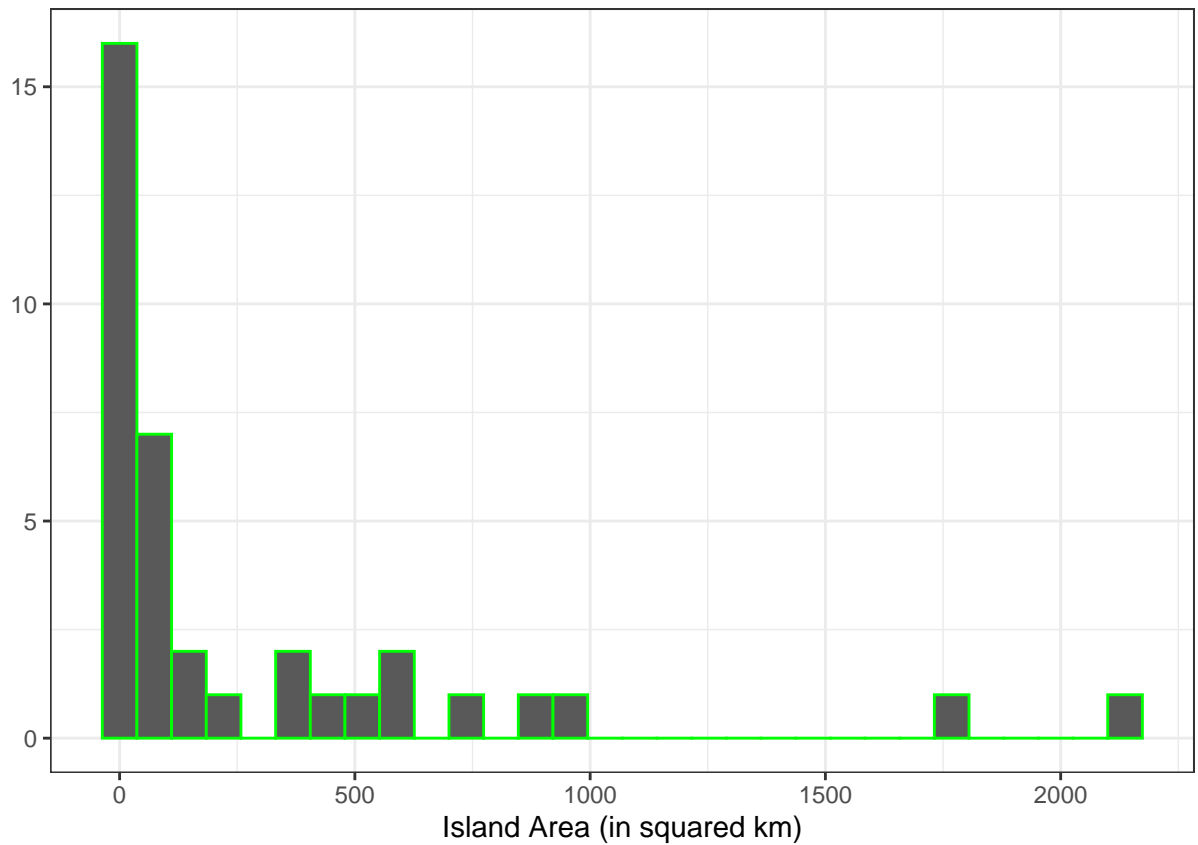
We note that for area there is a clear and extreme outlier.

```
species %>%
  filter(area > 200000)
```

```
## # A tibble: 1 x 7
##   island    area elevation  soil latitude distance species
##   <chr>    <dbl>    <int> <int>    <dbl>    <dbl>    <int>
## 1 Britain 229850.    1343   16     54.3      0    1666
```

Looking into the data we note that this outlier is clearly mainland Britain. To observe a more informative histogram, it is produced without this outlier.

```
species %>%
  filter(area < 200000) %>%
  ggplot(aes(x = area)) +
  geom_histogram(color = "green") +
  theme_bw() +
  ylab("") +
  xlab("Island Area (in squared km)")
```



After exclusion of other the outlier, we note that there are still a few more towards higher areas, and this variable in general appears to be strongly right skewed. Transformation of this variable might be necessary. We also note for further data visualization that the observation of mainland Britain will likely be an outlier as well.