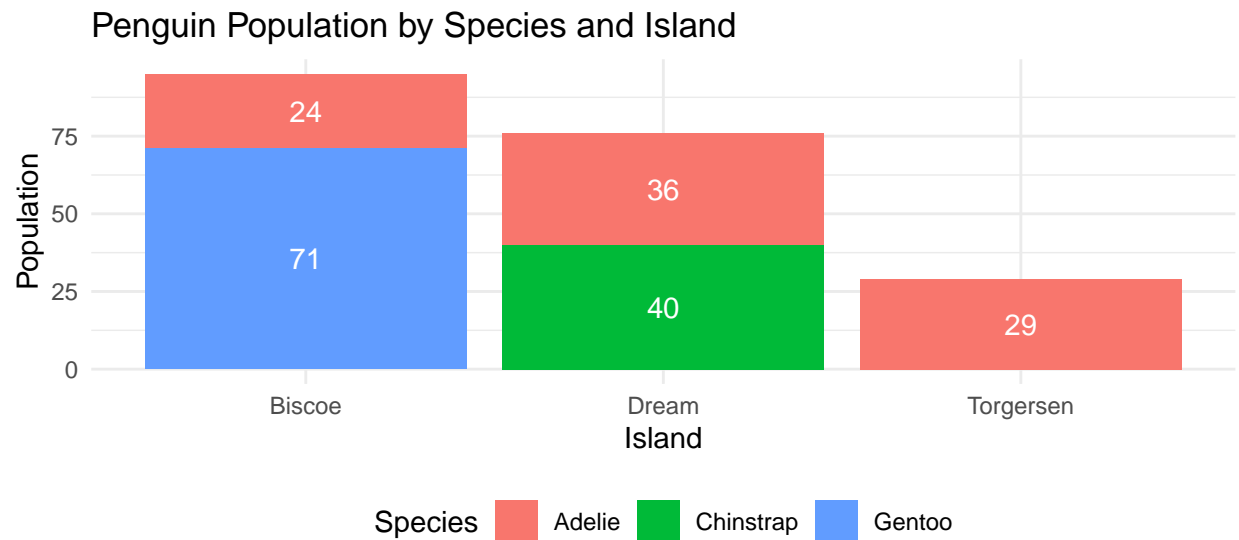# Palmer Penguins Analysis

Harish Sainath Sai Prasath (230164781)

## Introduction

In this report, we aim to analyse a subset of the famous *Palmer Penguins* data set. The subset data used for analysis has 200 entries and 8 different variables describing the penguins. These variables are *species, island, bill_length_mm, bill_depth_mm, flipper_length_mm, body_mass_g, sex, year* . Before diving deep into the data, let's try to get a better understanding about the penguins given in the data set by exploring all the variables, one at a time.
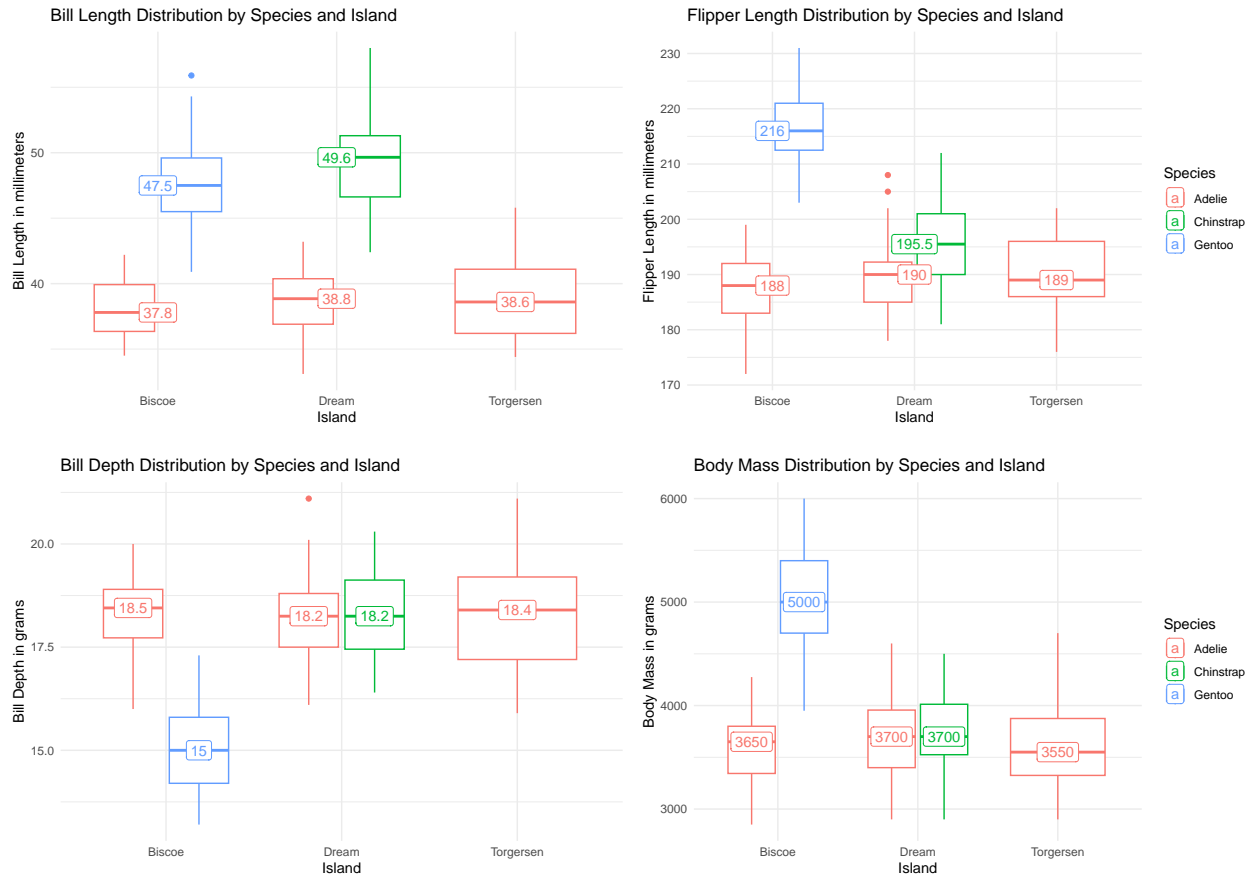
**Species and Islands:**

- There are 3 different islands and these are *Biscoe, Dream, Torgersen.*
- We can observe from the below plot that, Biscoe is the most populated island followed by Dream and Torgersen.
- There are 3 species of penguins and they are *Adelie, Chinstrap, Gentoo* respectively.
- Population of each species in these island is given in the below plot.



**Bill Length, Bill Depth, Flipper Length and Body Mass:**

- To get an idea about the features of penguins, we will consider median as a the measure of central tendency because it is not affected if the data is skewed.
- The penguins have a median bill length of *43.5 mm.*
- From the below box plots we can observe that Chinstrap Penguins has the highest median Bill length. This is closely followed by Gentoo penguins.

- Apart from the median bill lengths, a few interesting points can also be observed. Adelie is the most common type of penguin that is present in all three islands. Chinstrap is present only in Dream island, and Gentoo is present only present in Biscoe island.
- All penguins have a median bill depth of *17.3 mm*. Median bill depth for each species can be observed from the below box plots.
- Gentoo Penguins have the highest median flipper length and Adelie has the overall lowest median flipper length.
- The penguins have a median body mass of *3950 g*.
- Gentoo penguins are the heaviest penguins in the dataset. Adelie penguins are the the least heavy of the group and chinstrap penguins are comparable to Adelie as per body mass is concerned.
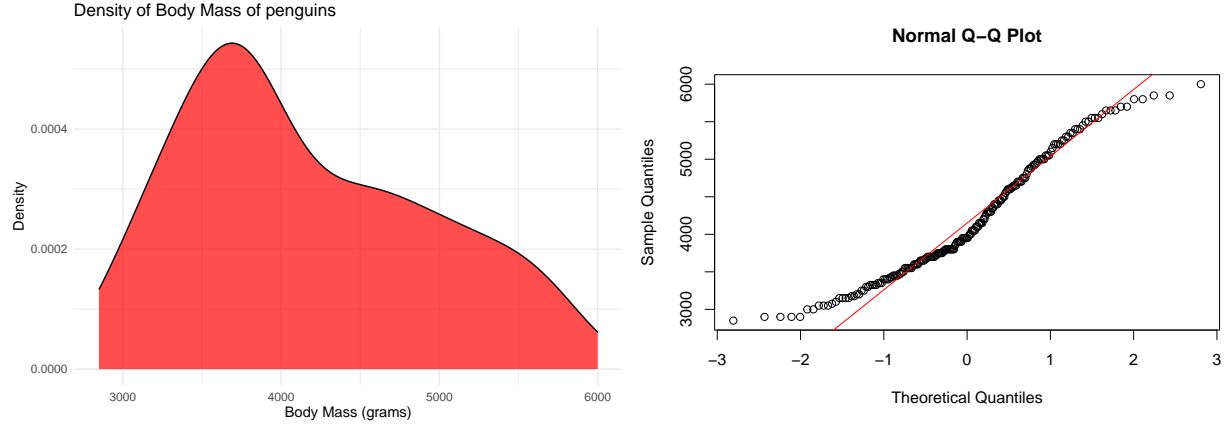


**Year and Sex :**

- The data is collected for 3 years and these are *2007, 2008, 2009* respectively. There are 98 female penguins and 102 male penguins in the dataset.

# Distribution Fitting and Population Proportion Estimation

Let's take one Body mass of penguins and fit a distribution for proportion estimation. For getting initial impressions of how the data is distributed, let's plot a density plot. Also a lot of things in nature are observed to be falling under a normal distribution. So Let's use a q-q plot to check if body mass can be fitted using a Normal/Gaussian distribution.

Density of Body Mass of penguins / Normal Q–Q Plot

Two major things can be observed the above plots.

- q-q plot shows that fitting a normal distribution might not be the best choice for good parameter estimations.
- Distribution of body mass seems to right skewed.

Now we have to look for some other distribution to fit the data. The body mass is a positive value and the data is right skewed, *Log-Normal Distribution* might be a good choice to fit the data.

so let's fit Log-Normal Distribution to our data.

**Fitting the Data to a Log-Normal Distribution**

Probability density function for a Log-Normal distribution is given by

$$f(x) = \frac{1}{x\sigma\sqrt{2\pi}} e^{\left[-\frac{(ln(x)-\mu)^2}{2\sigma^2}\right]}$$

- $\mu$ is the mean of the normal distribution.
- $\sigma$ is the standard deviation of the normal distribution.
- $n$ is equal to 200 in our case.

Likelihood function:

$$L(\mu, \sigma | x_n) = \prod_{i=1}^{n}\left[\frac{1}{\sqrt{2\pi}\sigma x_i}e^{-\frac{1}{2\sigma^2}(lnx_i - \mu)^2}\right]$$

Rearranging the above equation to make the calculations convenient

$$L(\mu, \sigma | x_n) = \frac{1}{(2\pi)^{\frac{n}{2}}}\frac{1}{(\sigma^2)^{\frac{n}{2}}}\prod_{i=1}^{n}\left[(\frac{1}{x_i})e^{-\frac{1}{2\sigma^2}\Sigma(lnx_i - \mu)^2}\right]$$

taking ln on both sides

$$\ln L(\mu, \sigma^2 | x_n) = \ln 1 - \frac{n}{2}\ln \pi + \ln 1 - \frac{n}{2}\ln \sigma^2 + \ln(x_1.x_2....x_n) - \frac{1}{2\sigma^2}\Sigma(lnx_i - \mu)^2 \ln e$$

$$\ln L(\mu, \sigma^2 | x_n) = -\frac{n}{2}\ln \pi - \frac{n}{2}\ln \sigma^2 + \sum lnx_i - \frac{1}{2\sigma^2}\sum(lnx_i - \mu)^2 \qquad\qquad ...(1)$$

3

Where (1) is the log likelihood function

Now partially differentiating (1) we find equations for $\mu$ and $\sigma$.

- To estimate parameter $\mu$ (Mean) we differentiate the log-likelihood with respect to $\mu$

$$\frac{\partial}{\partial \mu} lnL(\mu, \sigma^2 | x_n) = \frac{1}{\sigma^2} \sum_{i=1}^{n} (ln(x_i) - \mu)$$

$$\frac{\partial}{\partial \mu} lnL(\mu, \sigma^2) = 0$$

Solving this equation we get,

$$\text{mean: } \hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} \ln(x_i)$$

-To estimate parameter for $\sigma^2$ (Variance) we differentiate the log-likelihood with respect to $\sigma^2$

$$\frac{\partial}{\partial \sigma^2} lnL(\mu, \sigma^2 | x_n) = -\frac{n}{2\sigma^2} - \frac{1}{2\sigma^4} \sum_{i=1}^{n} (ln(x_i) - \mu)^2$$
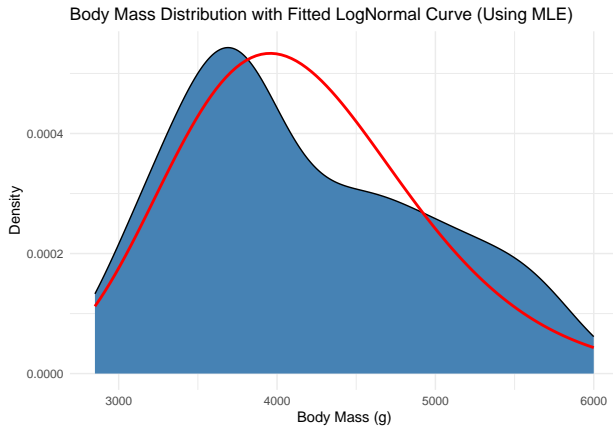
by solving this equation, we get

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (\ln(x_i) - \hat{\mu})^2$$

$$\text{Standard deviation: } \hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (\ln(x_i) - \hat{\mu})^2}$$

```
## MLE for mu (mean of log body mass): 8.317716
```

```
## MLE for sigma (standard deviation of log body mass): 0.1857341
```



Body Mass Distribution with Fitted LogNormal Curve (Using MLE)

**Key Observations:**

- While Log-Normal represents skewness, it still does not match the actual data. But this is far better than Normal distribution where skewness is zero.
- A more sophisticated model like gamma distribution would better represent the data that is presented.

# Determination of Sex

Determination of sex among penguins is often a difficult process where a lot of distress is caused to the penguins. Here, we look to use the available data to determine the sex of the penguins, thereby avoiding any invasive procedures.

Now, we will be looking at the following variables and compare the means of specific variable between two sexes. The t-test that we will be performing will give us an idea if there is a significant difference between the variable for each sex.

- Bill Length
- Bill Depth
- Flipper Length
- Body Mass

Let's perform t-test

```
options(scipen = 999)
options(width = 100)
t.test(bill_length_mm ~ sex, data = my.penguins,var.equal=TRUE)
```

**Note: We are assuming variances of two groups are equal**

```
##
##   Two Sample t-test
##
## data:  bill_length_mm by sex
## t = -4.802, df = 198, p-value = 0.000003093
## alternative hypothesis: true difference in means between group female and group male is not equal to
## 95 percent confidence interval:
##  -5.116182 -2.137399
## sample estimates:
## mean in group female    mean in group male
##             41.94184              45.56863
```

```
t.test(bill_depth_mm ~ sex, data = my.penguins,var.equal=TRUE)
```

```
##
##   Two Sample t-test
##
## data:  bill_depth_mm by sex
## t = -5.7102, df = 198, p-value = 0.00000004084
## alternative hypothesis: true difference in means between group female and group male is not equal to
## 95 percent confidence interval:
##  -1.8922373 -0.9207679
## sample estimates:
## mean in group female    mean in group male
##             16.40918              17.81569
```

```r
t.test(flipper_length_mm ~ sex, data = my.penguins,var.equal=TRUE)
```

```
## 
##  Two Sample t-test
## 
## data:  flipper_length_mm by sex
## t = -4.0131, df = 198, p-value = 0.00008494
## alternative hypothesis: true difference in means between group female and group male is not equal to
## 95 percent confidence interval:
##  -11.76401  -4.01190
## sample estimates:
## mean in group female    mean in group male
##              196.3571              204.2451
```

```r
t.test(body_mass_g ~ sex, data = my.penguins,var.equal=TRUE)
```

```
## 
##  Two Sample t-test
## 
## data:  body_mass_g by sex
## t = -6.1672, df = 198, p-value = 0.000000003838
## alternative hypothesis: true difference in means between group female and group male is not equal to
## 95 percent confidence interval:
##  -832.5006 -429.0940
## sample estimates:
## mean in group female    mean in group male
##              3845.918              4476.716
```

## Difference in Physical Characterstics of Penguins between Islands

```r
# Perform a one-way ANOVA for each physical characteristic
# Bill length by island
anova_bill_length <- aov(bill_length_mm ~ island, data = my.penguins)
summary(anova_bill_length)
```

```
##              Df Sum Sq Mean Sq F value      Pr(>F)    
## island        2    979   489.4   18.11 0.0000000603 ***
## Residuals   197   5324    27.0                      
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```r
# Bill depth by island
anova_bill_depth <- aov(bill_depth_mm ~ island, data = my.penguins)
summary(anova_bill_depth)
```

```
##              Df Sum Sq Mean Sq F value                Pr(>F)    
## island        2  294.1  147.07   71.52 <0.0000000000000002 ***
## Residuals   197  405.1    2.06                              
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# Flipper length by island
anova_flipper_length <- aov(flipper_length_mm ~ island, data = my.penguins)
summary(anova_flipper_length)
```

```
##              Df Sum Sq Mean Sq F value              Pr(>F)
## island        2  15243    7622   57.53 <0.0000000000000002 ***
## Residuals   197  26098     132
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# Body mass by island
anova_body_mass <- aov(body_mass_g ~ island, data = my.penguins)
summary(anova_body_mass)
```

```
##              Df   Sum Sq  Mean Sq F value              Pr(>F)
## island        2 48985030 24492515   64.82 <0.0000000000000002 ***
## Residuals   197 74432216   377829
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```