# STA 563 Project 2

Ryan McCollum

## I. Introduction

The purpose of this analysis will be to explore a data set containing 9 physical measurements and characteristics of 3893 farmed crabs from 2018-2019. The predictor variables provided in the data set include Sex, length, diameter, height, weight, shucked weight, viscera weight, and shell weight. The report will go through the model-building process to attempt to find the best model in terms of predicting the age of a crab based on the predictor variables in the data set.

Research hypothesis: Is the age of a farmed crab associated with the crab's physical measurements and/or the sex of the crab.

- It is worth noting that the Sex variable is a categorical variable with 3 categories (Male, Female, and Indeterminate) so it is separated into 2 indicator variables, using Male as the baseline category

The variables are defined as follows:

- $Y$ = Age (in months)
- $X1$ = Length (feet)
- $X2$ = Diameter (ft)
- $X3$ = Height (ft)
- $X4$ = Weight (ounces)
- $X5$ = Shucked.Weight (oz)
- $X6$ = Viscera.Weight (oz)
- $X7$ = Shell.Weight (oz)
- $X8$ = SexF (Indicates if a crab is female or not)
- $X9$ = SexI (Indicates if a crab's sex is indeterminate or not)

Some terms and definitions that may help with your understanding:

- The shucked weight of a crab is the crab's weight without their shell
- The viscera weight of a crab is the weight that wraps around the crab's abdominal organs.
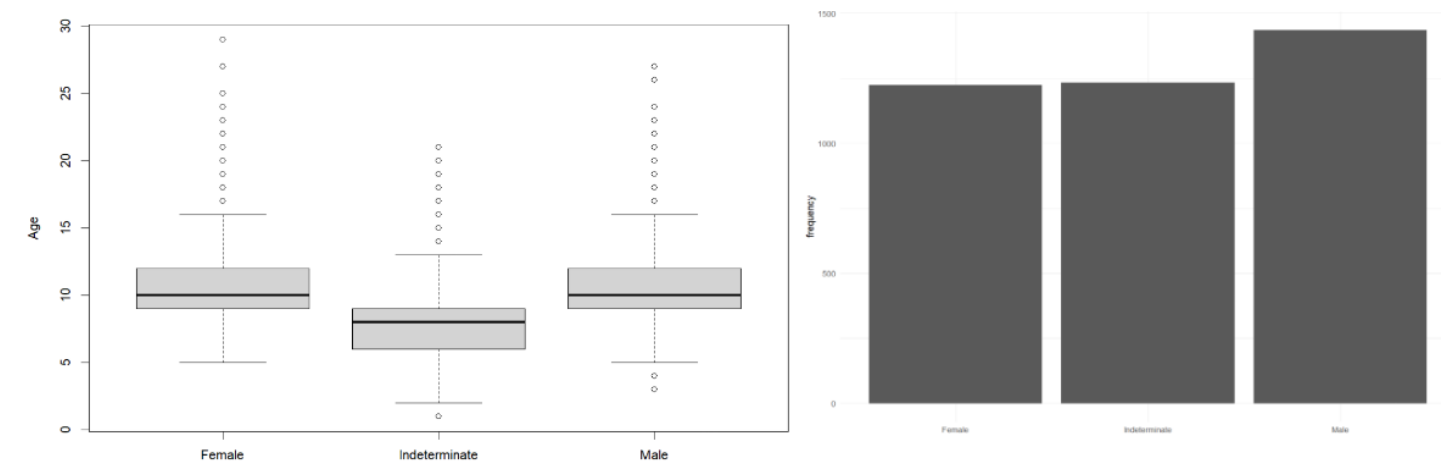
## II. Exploratory Data Analysis (EDA)

To start the analysis, view the first few observations of the raw data set (which can be viewed in the appendix). The first thing to notice is that almost all of the variables are numeric, except for sex. After reading the data description, it is known that there are 3 categories for the variable sex- female, male, and indicator. Since there are 3 categories it is necessary to make 2 indicator variables. Since the baseline is easy to compare, it would make sense for the baseline to be male or female so that we can compare the two basic sexes. For this report, male was made the baseline category. After adding the necessary indicator variables, we can look again at the format of the data.

**FIGURE 1:** View of Data Set

```
> head(CrabData_W_Dummys)
  Length Diameter Height    Weight Shucked.Weight Viscera.Weight Shell.Weight Age SexI SexF
1 1.4375   1.1750 0.4125 24.635715       12.332033       5.584852     6.747181   9    0    1
2 0.8875   0.6500 0.2125  5.400580        2.296310       1.374951     1.559222   6    0    0
3 1.0375   0.7750 0.2500  7.952035        3.231843       1.601747     2.764076   6    1    0
4 1.1750   0.8875 0.2500 13.480187        4.748541       2.282135     5.244657  10    0    1
5 0.8875   0.6625 0.2125  6.903103        3.458639       1.488349     1.700970   6    1    0
6 1.5500   1.1625 0.3500 28.661344       13.579410       6.761356     7.229122   8    0    1
```
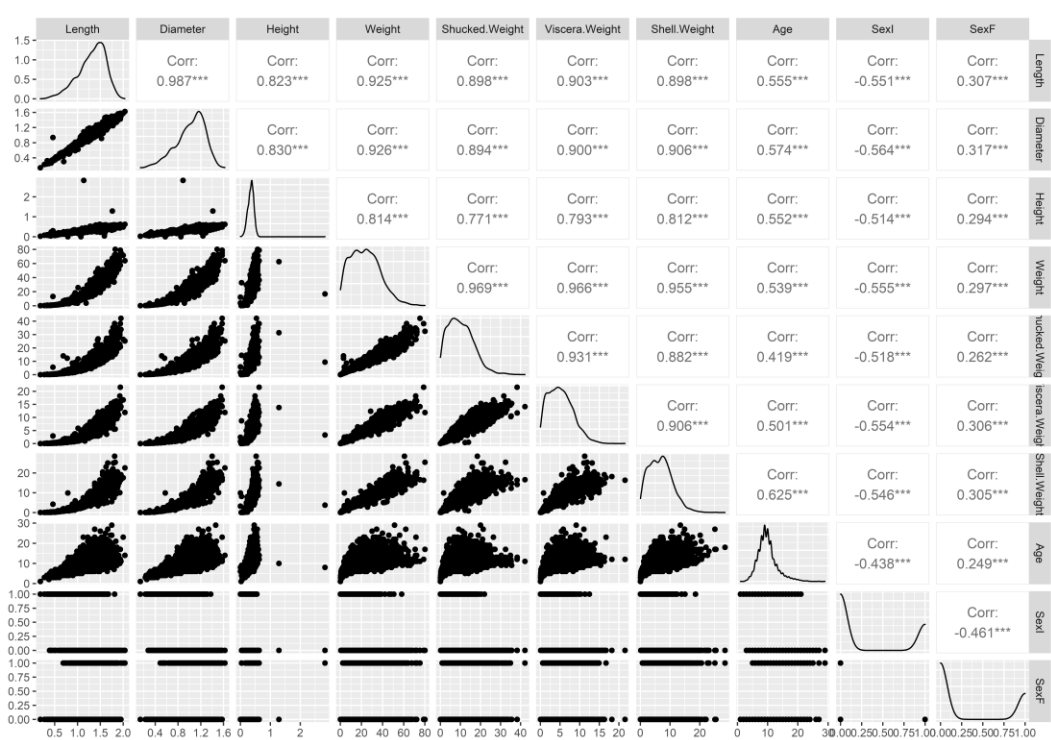
Another thing worth investigating is how common each of the sexes appears in the data set, and if the sex of the crab has a relationship with the average age of the crab. Looking at TABLE, the frequency chart shows that each of the sexes appears approximately the same number of times in the dataset, so it is appropriate to include all three in the evaluation.

**FIGURE 2:** Box plot of Sexes and Age and Frequency plot of Sexes



Hypothetically, this may suggest that the SexI indicator variable can be a significant predictor in the model. Also, looking at Figure 2, we can see that the average age of female and male crabs was about the same, but that the average age of the indeterminate crabs seems to be significantly lower than the other two groups. Lastly, looking at the correlation between variables in the data set may give insight into how the variables will behave through the model-building process. Figure 3 gives the scatterplot matrix of the data set.

**FIGURE 3:** Scatterplot Matrix



From Figure 3 we can gather that there is a lot of correlation between the variables that analyze the physical measurements of the crab, which makes sense. However, it is interesting there is only a moderate correlation between the sex variables and the measurement variables, so I am unsure how the variables will interact in the model-building process. Although there is likely a lot of multicollinearity between the variables, I am unsure which measurement variable will act as the best predictor for the Age of the crab. Instead, the full model will be built using all of the measurement predictors, the sex indicator predictors, and then the interactions between the sex and measurement predictors. Since there is a high correlation between all the measurement predictors, it does not make sense to include their interaction effects.

There are some other noteworthy things that can be viewed from Figure 3. For one, there is only one crab measured that is taller than 2 feet. This crab seems to be an outlier in the dataset, and it is something to be made aware of. Another thing worth noting is that almost all of the scatterplots for the response variable Age tend to have a 'megaphone' shape, which might suggest non-constant variance. Although the goal of this project is not inference and model assumptions do not matter as much it is still worth noting. Lastly, the Length and Diameter scatterplots seem to have some curvature towards them, so in further analysis it may be worth looking into some quadratic terms, however, for this analysis only linear terms will be included to keep the model simple.
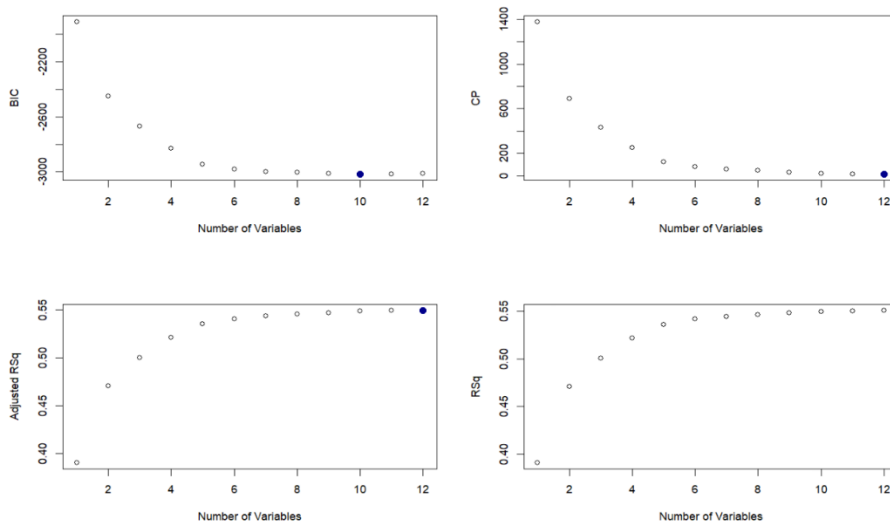
## III.    Analysis

Post-data exploration, we can start the predictive model-building process. After noting that there is probably multicollinearity between many of the measurement variables, the full model being built will not include interactions between the measurement variables, but interactions between the sex indicator variables and the measurement variables will be included, as mentioned earlier. The full model is as follows:

$$Y = \beta_0 + \sum_{i=1}^{9} \beta_i X_i + \sum_{j=1}^{7} \sum_{k=8}^{9} \beta_{j,k} X_j X_k + \varepsilon$$

Considering the fairly large data set, I feel it is appropriate to form multiple reduced models using automatic search procedures before comparing models through 10-fold cross-validation. To start, using the regsubsets function in R, models will be fit to the entire crab dataset, and for each unique number of predictors from 1 to 12, the best model in terms of $R^2$ will be kept. Then, using AIC and BIC, we will identify what the best number of predictors to use in the model. Maxing the possible chosen variables at 12 was selected to keep the model run time lower and prevent too complicated models being selected.

**FIGURE 4:** Suggested Number of Variables



As shown in Figure 4, BIC and CP (the same as AIC) do not agree on how many predictors to put in the model. Although adjusted $R^2$ and $R^2$ graphs are included, they are only there to help confirm the findings of BIC and AIC and their suggestions will not be used to create a model. One important thing to note is that even though BIC and CP suggest a different number of variables, their curves level out between the range 8-12. Since BIC and CP have different suggested numbers of predictors, it is possible to compare the model with the minimum value of each since they will be different models.

The suggested model with the lowest CP without its estimated model coefficients:

$$E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6 + \beta_7 X_7 + \beta_8 X_8 + \beta_9 X_9 + \beta_{3,9} X_3 X_9 + \beta_{4,8} X_4 X_8 + \beta_{5,9} X_5 X_9$$

The suggested model with the lowest BIC without its estimated model coefficients:

$$E(Y) = \beta_0 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6 + \beta_7 X_7 + \beta_8 X_8 + \beta_9 X_9 + \beta_{5,9} X_5 X_9 + \beta_{3,8} X_3 X_8$$

In the CP model, all of the individual variables are chosen, along with 3 interaction variables. In the BIC model, all of the individual variables except for Length are chosen, along with 2 interactions. Both models chose the interaction between Shucked.Weight and SexI, but the other 2 interactions chosen in the CP model differ from the second interaction term in the

BIC model. Since the models are relatively similar, another model or two to compare against may be helpful. Using forward and backward stepwise selection, it is possible to obtain more suggested models. Since there is already a 10 and 12 variable model, the best 8 predictor variable for both forward and backward selection will be chosen, if they differ. Since BIC, CP and Adjusted $R^2$ suggested models with between 8-12 variables, the simplest models from forward and backward selection that are in this range will be used to compare. After running both forward and backward stepwise selections the following 8 predictor models are selected.

The forward model with 8 variables selects the following model:
$$E(Y) = \beta_0 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6 + \beta_7 X_7 + \beta_8 X_8 + \beta_{8,3} X_3 : X_8$$

After viewing the following output from backward stepwise selection, we see that backward stepwise selection produced a different 8 variable model than forward, so this model can also be used to compare against the BIC and CP models.

The backward model with 8 variables selects the following model:

$$E(Y) = \beta_0 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_7 X_7 + \beta_8 X_8 + \beta_9 X_9 + \beta_{9,3} X_3 X_9 + \beta_{8,5} X_5 X_8$$

Now that 4 different models have been built through different model selection procedures, 10-fold cross-validation can be implemented to judge which model produces the best estimated test MSE. As mentioned earlier, the data set has almost 4000 observations, so each fold will have approximately 400 observations in it, which is large enough to avoid some of the possible complications we might encounter when trying to use 10-fold cross-validation on a smaller data set. After running the 10-fold cross-validation, test MSE can be found after a holdout set approach and leave one out cross-validation (LOOCV), however, these estimated test MSEs will only be used to support the findings from the 10-fold cross-validation, and won't be used to directly conclude the best model. The BIC, CP, Forward, and Backward models will be tested, along with the Full model to prove that the Full model likely overfits the data. Further, instead of outputting the test MSE, the RMSE will be analyzed, since it is on the same scale as the response variable and easier to interpret.

**TABLE 1:** The resulting RMSEs of the 10-fold cross-validation

| MODEL | BIC | CP | Forward | Backward | Full |
|---|---|---|---|---|---|
| RMSE ($\sqrt{\text{MSE}}$) | 2.171963 | 2.171047 | 2.193514 | 2.191208 | 2.198932 |

As evident by Table 1, the CP model produces the best RMSE, even though it is only slightly better than that of the BIC model. The forward and backward model are relatively similar, and as expected the full model performs the worst, likely due to overfitting. These results suggest that CP is the best model, however it is also worth considering that the BIC model is slightly simpler, so variables in the CP model that do not contribute as much to the model may be left out. The RMSE after LOOCV and holdout set approach is in Table 2.

**TABLE 2:** The resulting RMSEs of the holdout approach and LOOCV

| MODEL | BIC | CP | Forward | Backward | Full |
|---|---|---|---|---|---|
| HOLDOUT RMSE | 2.148398 | 2.150537 | 2.153326 | 2.157006 | 2.157903 |
| LOOCV RMSE | 2.17584 | 2.173386 | 2.185303 | 2.192786 | 2.189951 |

As mentioned previously, the values in TABLE are not going to be used to pick the best model, and were only found as a measure to support the findings of the 10-fold cross-validation. Since RMSE is in the same units as the response variable, the 10-fold cross-validation test error for the models ranges from 2.171 months to 2.199 months. Based on TABLE, the BIC and CP models performed the best out of the five models through both of the tests. This supports the findings from the 10-fold cross-validation, so it is safe to assume that one of the BIC and CP models is the best for modeling the age of a farmed crab.

## IV. Conclusion

After the analysis, the decision to go with either the BIC or CP model can be argued. Personally, although the CP model had a slightly better RMSE, and therefore was better at predicting values, I lean towards the BIC model. My main reasoning is that a

lot of the variables in the data set had very high correlations with each other. Therefore, it is probably not necessary for all of them to be in the model. When further research is done on more data through an explanatory approach, it is possible that some variables are removed from the model due to multicollinearity. Also, when finding the BIC and CP models, it is worth noting that BIC penalizes predictors that do not add very much to the quality of the model much more than CP does. Therefore, BIC models tend to be simpler and their variables contribute more to the overall prediction.

So to conclude, due to high correlation between variables, a slightly simple model will likely perform almost as effectively as a more complex model, without adding variables that contribute only slightly to the models effectiveness. The hypothesis that the age of a farmed crab is associated with the crab's physical measurements and/or the sex of the crab, held to be mostly true. The only measurement variable that was not included in the selected BIC model was Length, and the sex indicators were both chosen, as well as 2 interaction terms, one using SexF and one using SexI.

Since this was an exploratory approach to the data, it is not appropriate to test p-values of the newly selected model, however it would be interesting to see how the BIC model predicts a slightly larger indeterminate sex crab and to a smaller female crab, since in the EDA it was shown that the average age of the indeterminate crabs was slightly lower than that of the male and female crabs. The results are shown in Table 3.

**TABLE 3:** Prediction for Large Indeterminate and Small Female Crabs

|  | Diameter | Height | Weight | Shucked.Weight | Viscera.Weight | Shell.Weight | SexI | SexF | Pred Age |
|---|---|---|---|---|---|---|---|---|---|
| Indeter. | .95 | .4 | 18 | 8.5 | 4.5 | 5 | 1 | 0 | 9.04 |
| Female | 1 | .4 | 20 | 8.7 | 4.7 | 5.2 | 0 | 1 | 10.05 |

*Note:* the values of the Indeterminate crab were determined by roughly the approximate upper quartile of measurements for the Indeterminate crabs, as shown in the boxplots in the appendix. The same was done for the Female crab measurements, but in regard to their approximate lower quartile measurements.

With no other information, it is interesting that the Indeterminate sex crabs seem to be smaller than their male and female counterparts, as well as younger. If I had to guess, it may be harder to determine the sex of a crab before they reach a certain age or size range, which leads to the indeterminate crabs being predicted to being smaller than the other crabs, even when they have measurements that are larger than most of the other indeterminate crabs. Further analysis could be done using the selected model on new data to test this hypothesis more.

# V.  Appendices

**Appendix 1:**

This is a view of the actual raw data before indicator variables are made for the Sex variable.

```
> head(CrabData)
  Sex Length Diameter Height     Weight Shucked.Weight Viscera.Weight Shell.Weight Age
1   F 1.4375   1.1750 0.4125 24.635715      12.332033       5.584852     6.747181   9
2   M 0.8875   0.6500 0.2125  5.400580       2.296310       1.374951     1.559222   6
3   I 1.0375   0.7750 0.2500  7.952035       3.231843       1.601747     2.764076   6
4   F 1.1750   0.8875 0.2500 13.480187       4.748541       2.282135     5.244657  10
5   I 0.8875   0.6625 0.2125  6.903103       3.458639       1.488349     1.700970   6
6   F 1.5500   1.1625 0.3500 28.661344      13.579410       6.761356     7.229122   8
```

## Appendix 2:

This is the summary output from the CP selected model with 12 predictors.

```
Call:
glm(formula = Age ~ Length + Diameter + Height + Weight + Shucked.Weight +
    Viscera.Weight + Shell.Weight + SexI + SexF + Height:SexF +
    Shucked.Weight:SexI + Height:SexI, data = CrabData_W_Dummys)

Deviance Residuals:
    Min      1Q   Median      3Q      Max
-8.4566  -1.3096  -0.2862   0.8375  14.4130

Coefficients:
                      Estimate Std. Error t value Pr(>|t|)
(Intercept)            5.56496    0.41817  13.308  < 2e-16 ***
Length                -1.83135    0.74619  -2.454  0.01416 *
Diameter               3.61900    0.91164   3.970 7.32e-05 ***
Height                 6.46088    0.94244   6.855 8.23e-12 ***
Weight                 0.33178    0.02586  12.830  < 2e-16 ***
Shucked.Weight        -0.69223    0.02961 -23.380  < 2e-16 ***
Viscera.Weight        -0.33863    0.04629  -7.316 3.10e-13 ***
Shell.Weight           0.31660    0.04051   7.815 7.03e-15 ***
SexI                  -2.56408    0.40619  -6.313 3.05e-10 ***
SexF                   1.99270    0.38045   5.238 1.71e-07 ***
Height:SexF           -5.14575    0.95784  -5.372 8.23e-08 ***
Shucked.Weight:SexI    0.09826    0.03341   2.941  0.00329 **
Height:SexI            3.85094    1.73378   2.221  0.02640 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 4.671514)

    Null deviance: 40378  on 3892  degrees of freedom
Residual deviance: 18125  on 3880  degrees of freedom
AIC: 17064

Number of Fisher Scoring iterations: 2
```

## Appendix 3:

This is the summary output from the BIC selected model with 10 predictors.

```
Call:
glm(formula = Age ~ Diameter + Height + Weight + Shucked.Weight +
    Viscera.Weight + Shell.Weight + SexI + SexF + Height:SexF +
    Shucked.Weight:SexI, data = CrabData_W_Dummys)

Deviance Residuals:
    Min      1Q   Median      3Q      Max
-8.4367  -1.2944  -0.2911   0.8343  14.5446

Coefficients:
                      Estimate Std. Error t value Pr(>|t|)
(Intercept)            4.85193    0.34819  13.935  < 2e-16 ***
Diameter               1.91460    0.46640   4.105 4.13e-05 ***
Height                 7.25723    0.85256   8.512  < 2e-16 ***
Weight                 0.33394    0.02587  12.907  < 2e-16 ***
Shucked.Weight        -0.70819    0.02913 -24.311  < 2e-16 ***
Viscera.Weight        -0.35447    0.04601  -7.703 1.67e-14 ***
Shell.Weight           0.30967    0.04034   7.677 2.04e-14 ***
SexI                  -1.73599    0.18659  -9.304  < 2e-16 ***
SexF                   2.20487    0.36231   6.086 1.27e-09 ***
Height:SexF           -5.71769    0.90574  -6.313 3.05e-10 ***
Shucked.Weight:SexI    0.14799    0.02174   6.806 1.16e-11 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 4.68118)

    Null deviance: 40378  on 3892  degrees of freedom
Residual deviance: 18172  on 3882  degrees of freedom
AIC: 17070

Number of Fisher Scoring iterations: 2
```

# Appendix 4:

This is the summary output of the forward model, highlighted are the variables selected by the 8 predictor model.

```
Selection Algorithm: forward
         Length Diameter Height Weight Shucked.Weight Viscera.Weight Shell.Weight SexI SexF Length:SexI Length:SexF Diameter:SexI Diameter:SexF Height:SexI
1  ( 1 ) " "    " "      " "    " "    " "            " "            "*"          " "  " "  " "         " "         " "           " "           " "
2  ( 1 ) " "    " "      " "    " "    "*"            " "            "*"          " "  " "  " "         " "         " "           " "           " "
3  ( 1 ) " "    "*"      " "    " "    "*"            " "            "*"          " "  " "  " "         " "         " "           " "           " "
4  ( 1 ) " "    "*"      " "    "*"    "*"            " "            "*"          " "  " "  " "         " "         " "           " "           " "
5  ( 1 ) " "    "*"      " "    "*"    "*"            " "            "*"          "*"  " "  " "         " "         " "           " "           " "
6  ( 1 ) " "    "*"      " "    "*"    "*"            " "            "*"          "*"  " "  " "         " "         " "           " "           "*"
7  ( 1 ) " "    "*"      " "    "*"    "*"            "*"            "*"          "*"  " "  " "         " "         " "           " "           "*"
8  ( 1 ) " "    "*"      "*"    "*"    "*"            "*"            "*"          "*"  " "  " "         " "         " "           " "           "*"
9  ( 1 ) " "    "*"      "*"    "*"    "*"            "*"            "*"          "*"  " "  " "         " "         " "           " "           "*"
10 ( 1 ) "*"    "*"      "*"    "*"    "*"            "*"            "*"          "*"  " "  " "         " "         " "           " "           "*"
11 ( 1 ) "*"    "*"      "*"    "*"    "*"            "*"            "*"          "*"  " "  " "         " "         " "           " "           "*"
12 ( 1 ) "*"    "*"      "*"    "*"    "*"            "*"            "*"          "*"  "*"  " "         " "         " "           " "           "*"

         Height:SexF Weight:SexI Weight:SexF Shucked.Weight:SexI Shucked.Weight:SexF Viscera.Weight:SexI Viscera.Weight:SexF Shell.Weight:SexI Shell.Weight:SexF
1  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
2  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
3  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
4  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
5  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
6  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
7  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
8  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
9  ( 1 ) " "         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
10 ( 1 ) " "         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
11 ( 1 ) " "         " "         " "         "*"                 "*"                 " "                 " "                 " "               " "
12 ( 1 ) " "         " "         " "         "*"                 "*"                 " "                 " "                 " "               " "
```

# Appendix 5:

This is the summary output of the backward model, highlighted are the variables selected by the 8 predictor model.

```
Selection Algorithm: backward
         Length Diameter Height Weight Shucked.Weight Viscera.Weight Shell.Weight SexI SexF Length:SexI Length:SexF Diameter:SexI Diameter:SexF Height:SexI
1  ( 1 ) " "    " "      " "    " "    " "            " "            "*"          " "  " "  " "         " "         " "           " "           " "
2  ( 1 ) " "    " "      " "    " "    "*"            " "            "*"          " "  " "  " "         " "         " "           " "           " "
3  ( 1 ) " "    " "      " "    " "    "*"            " "            "*"          "*"  " "  " "         " "         " "           " "           " "
4  ( 1 ) " "    " "      " "    " "    "*"            " "            "*"          "*"  " "  " "         " "         " "           " "           " "
5  ( 1 ) " "    " "      " "    "*"    "*"            " "            "*"          "*"  " "  " "         " "         " "           " "           " "
6  ( 1 ) " "    " "      "*"    "*"    "*"            " "            "*"          "*"  " "  " "         " "         " "           " "           " "
7  ( 1 ) " "    " "      "*"    "*"    "*"            " "            "*"          "*"  " "  " "         " "         " "           " "           " "
8  ( 1 ) " "    " "      "*"    "*"    "*"            " "            "*"          "*"  "*"  " "         " "         " "           " "           " "
9  ( 1 ) " "    " "      "*"    "*"    "*"            "*"            "*"          "*"  "*"  " "         " "         " "           " "           " "
10 ( 1 ) " "    "*"      "*"    "*"    "*"            "*"            "*"          "*"  "*"  " "         "*"         " "           " "           " "
11 ( 1 ) " "    "*"      "*"    "*"    "*"            "*"            "*"          "*"  "*"  " "         "*"         " "           " "           "*"
12 ( 1 ) "*"    "*"      "*"    "*"    "*"            "*"            "*"          "*"  "*"  " "         "*"         " "           " "           "*"

         Height:SexF Weight:SexI Weight:SexF Shucked.Weight:SexI Shucked.Weight:SexF Viscera.Weight:SexI Viscera.Weight:SexF Shell.Weight:SexI Shell.Weight:SexF
1  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
2  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
3  ( 1 ) " "         " "         " "         " "                 " "                 " "                 " "                 " "               " "
4  ( 1 ) " "         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
5  ( 1 ) " "         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
6  ( 1 ) " "         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
7  ( 1 ) "*"         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
8  ( 1 ) "*"         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
9  ( 1 ) "*"         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
10 ( 1 ) "*"         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
11 ( 1 ) "*"         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
12 ( 1 ) "*"         " "         " "         "*"                 " "                 " "                 " "                 " "               " "
```

**Appendix 6:**

Following are box plots comparing the differences in sexes across the 6 measurement variables selected by the BIC model