

**Black-Lip Abalone Linear Discrimination Study of Sex Types**

**MSDS 6372-402**

**Group Project II**

**July 13, 2017**

**James Park**

**Mooyoung Lee**



### Purpose/Question

Determine if a Blacklip Abalone is Male, Female, or an Infant (M, F, or I) based on various dimensional and weight parameters. The data set utilized in this study is from the Machine Learning Repository from the University of California, Irvine (UCI). We would like to investigate if it is possible to differentiate (M, F, or I), via a MANOVA and LDA. Specifically to differentiate M and F so that food preparers may quickly identify in order to determine how to prepare. In addition we will differentiate I from M and F so that fisheries may quickly identify and throw back the Infant Abalone back in the water to further mature.

### Background

There are many species of Abalone on coast lines across the world. Abalone are marine snails which are highly prized for their meat as well as their highly iridescent, nacre ('mother of pearl') shells used commonly in jewelry, buttons and inlay (figure 1). The flesh of abalones is widely considered to be a desirable food, and is consumed raw or cooked by a variety of cultures. Abalone mature to 9-10 centimeters in length in 3 to 6 years depending on the local conditions. Very few make it to maturity due to poaching and high demand.



Figure 1 (source: Wikipedia.org, wisegeek.com)

As in many fishing industries, demand is outstripping supply, and in order to maintain the species limits are placed by local authorities on the quantity and the size. In the 2016/2017 season in Australia an individual may obtain a 5 day (1 day per month over 5 months) fishing permit which allows up to 5 Abalone per day. The minimum size to keep is 13 centimeters in length, otherwise it must be thrown back.

The data investigated in this study is the Blacklip Abalone species found in the North Coast and Islands in the Bass Strait separating Tasmania from Australia (figure 2.)



Figure 2, (Source Google maps and pics)

The Sex of an Abalone is determined by orienting the shell so that the meat is facing up, then moving the meat to the side exposing the gonads. An off white color indicates a male and a bluish color indicates a female (figure 3). This process is very manual and could be sped up if Abalones could be first filtered by the variables collected prior to identification. Abalone vary in size from 20 millimeters (0.79 in) (*Haliotis pulcherrima*) to 200 millimeters (7.9 in) while *Haliotis rufescens* is the largest of the genus at 12 inches (30 cm).

Abalone has been an important staple in native cultures around the world, specifically in Africa and on the North American West coast. The meat was used as food, and the shell was used as currency for many tribes. More recently Japan has steadily increased demand for Abalones.



Figure 3 (source: fishconsult.org)

## The Data

The Abalone data set and variable descriptions may be found in table 1. In addition to the variables in table 1, an additional variable “Volume” (Height\*Diameter\*Length) was added in an attempt to simplify the data set. The data set contains one categorical variable (Sex) with three levels (Male, Female, and Infant), one integer variable (number of Rings) and the remaining are continuous relating to shell dimensions and weight. One item of measurement should be noted here regarding the measurement scale of the Length, Diameter, Height and Weights. There were 4,177 samples in the dataset. The source of the data indicates the measurements are represented in millimeters and grams although ALL

other sources of Blacklip Abalone biology indicate the measurements should be on the order of centimeters and kilograms. Given the original data collection was in 1995, and for the purposes of this project assignment, we will assume that the dimensions have been scaled to an obscure multiplier to protect the proprietary nature of the original study. This will still serve the purposes of this project.

Name	Data Type	Meas.	Description
---	-----	-----	-----
Sex	nominal		M, F, and I (infant)
Length	continuous	mm	Longest shell measurement
Diameter	continuous	mm	perpendicular to length
Height	continuous	mm	with meat in shell
Whole weight	continuous	grams	whole abalone
Shucked weight	continuous	grams	weight of meat
Viscera weight	continuous	grams	gut weight (after bleeding)
Shell weight	continuous	grams	after being dried
Rings	integer		+1.5 gives the age in years

Table 1 (source; uci.edu website)

Two potential outliers were identified in early diagnostic plots. The SAS GLM procedure revealed the data shown in Figure x. The outliers are easily visualized in the Residual plots, RStudent plots, as well as the qq plots. In looking more closely at the raw data, for these two points, it was an obvious data collection error in the "Height" variable. This is most likely a typo or recording error as similar data was either an order of magnitude different or a large amount off (>50% compared to all other data). These two point were removed leaving the total number of samples reduced to 4,175. Original data, corrections, and raw data excerpts may be seen in Figure 4.

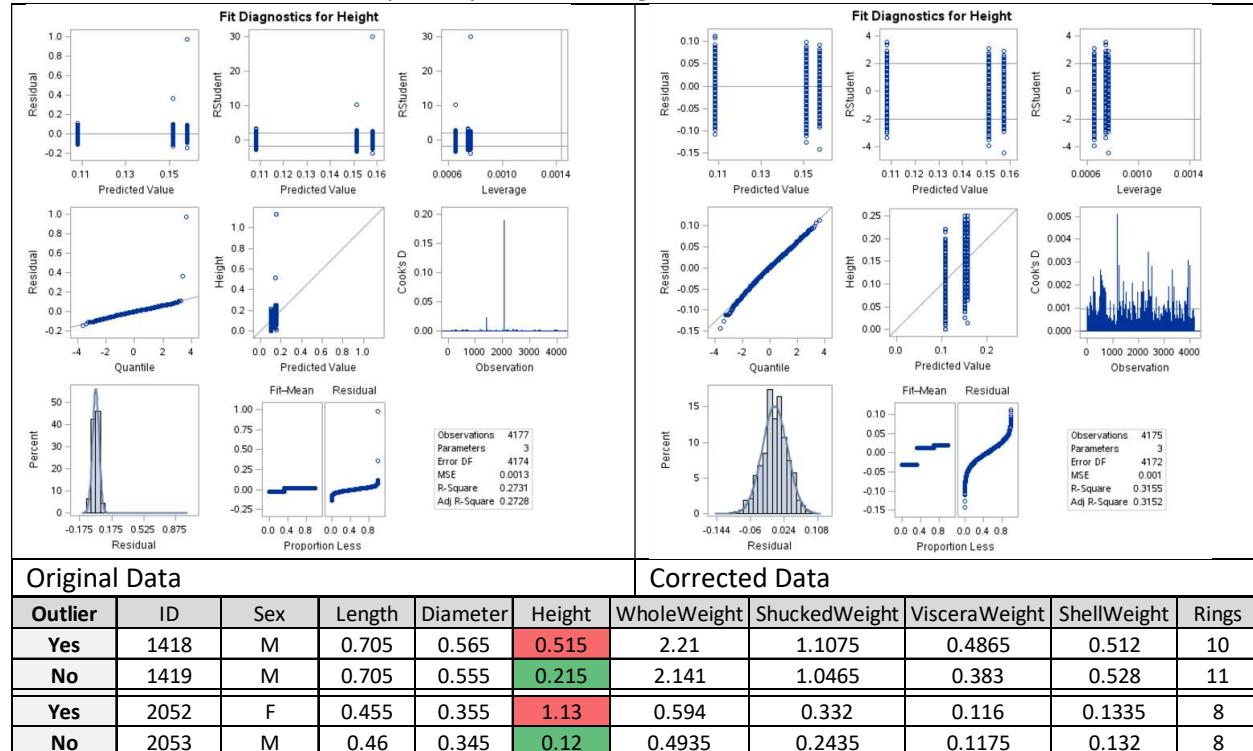
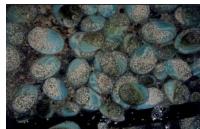


Figure 4



### Data Visualization (initial plots)

Prior to any deep data exploration the data was first visually explored on a high level to gain perspective and to provide a basic understanding of relative magnitudes and distributions. Exploration techniques included (See Appendix A for full results and code):

- Means Output: the calculation of means, minimums, maximums, range and standard deviations via a SAS means procedure (see Appendix A)
- Frequency and percent of Sex (M, F, I) via SAS frequency procedure. Frequencies are important in conducting and Linear Discriminant Analysis (LDA) so that the prior proportions may be taken into account. Frequency data is shown below in Figure 5. The data is relatively evenly distributed meaning about a third of the studied population is in each group.

Sex	Frequency	Percent	Cumulative Frequency	Cumulative Percent
F	1306	31.28	1306	31.28
I	1342	32.14	2648	63.43
M	1527	36.57	4175	100.00

Figure 5, frequency of Sex levels

- Boxplot of each variable grouped by Sex via SAS univariate procedure. As an example the “Length” variable by Sex (F, I, M) is shown below in Figure 6. The “Length” variable was typical of all other variables with respect to boxplot output, meaning the Infant was easily differentiated visually from the Male and Female samples and the Male and Female appeared fairly close. The remaining box-plots and associated code may be found in Appendix A.

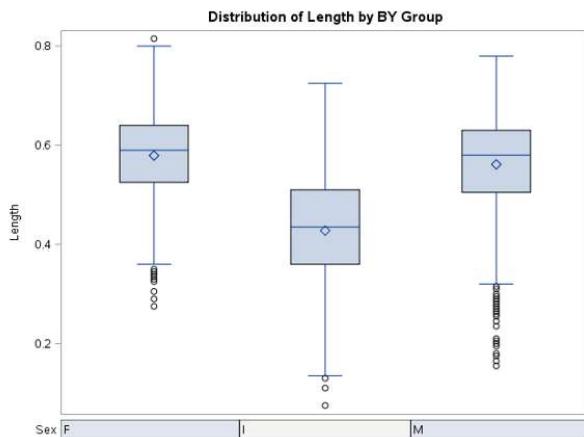


Figure 6, boxplot of example variable (Length) by Sex

- Scatter and Histogram plots of all continuous and integer variables involved via SAS sgscatter procedure were also viewed. As one can see in Figure 7, there is a fair amount of association of variables. Although this should not be an issue in conducting an LDA further residual examination will be explored in order meet LDA assumptions of normality of residuals.

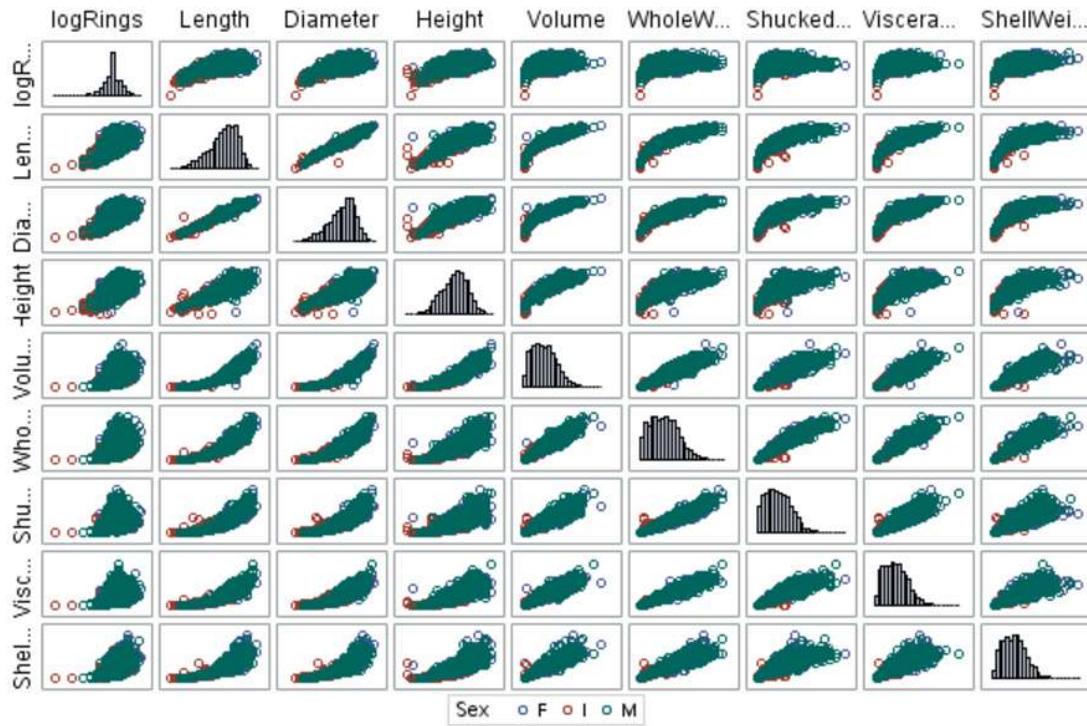


Figure 7, initial scatter plot

- Additionally scatter plots were more closely examined by Sex in order to visually examine if equal variance/covariance looked similar as part of the investigation if homogeneity is maintained. As an example in, in figure 8, the “Wholeweight” variable was compared to the “diameter” variable via scatter-plots and histograms by sex. The visual output would indicate that the variance/covariance may be similar. This was the case for all other variables investigated. While the full number of combinations is exhaustive several others and the associated SAS code used to generate are included in Appendix A.

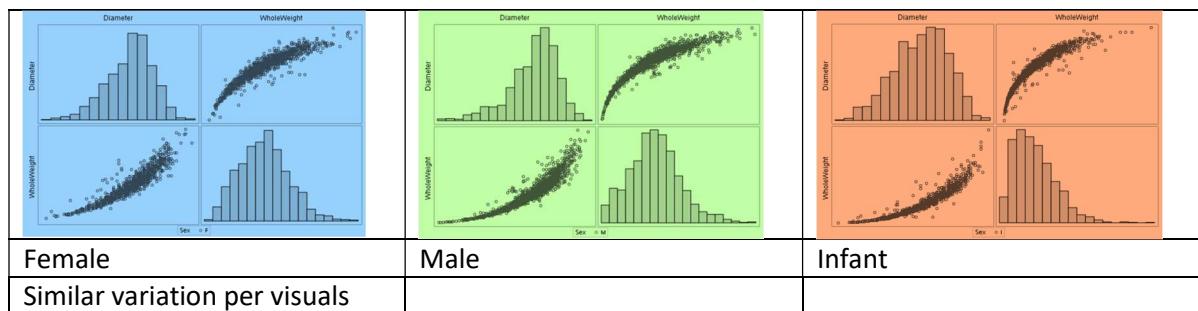


Figure 8, variable comparison by Sex

Next, 95% ellipse prediction plots were generated for variable pairs in order to assess potential differences in covariance matrices. Variables selected and the associated prediction ellipses may be viewed in figure 9. The ellipses indicate a difference in the covariance matrix between the Infant and other Sex categories (M and F). They also indicate similar shape and orientation between the M and F groups. Several variable pairs were randomly selected and they all exhibited similar trends. In the

interest of saving space and moving on in this study the other 95% prediction ellipses are not provided in this report.

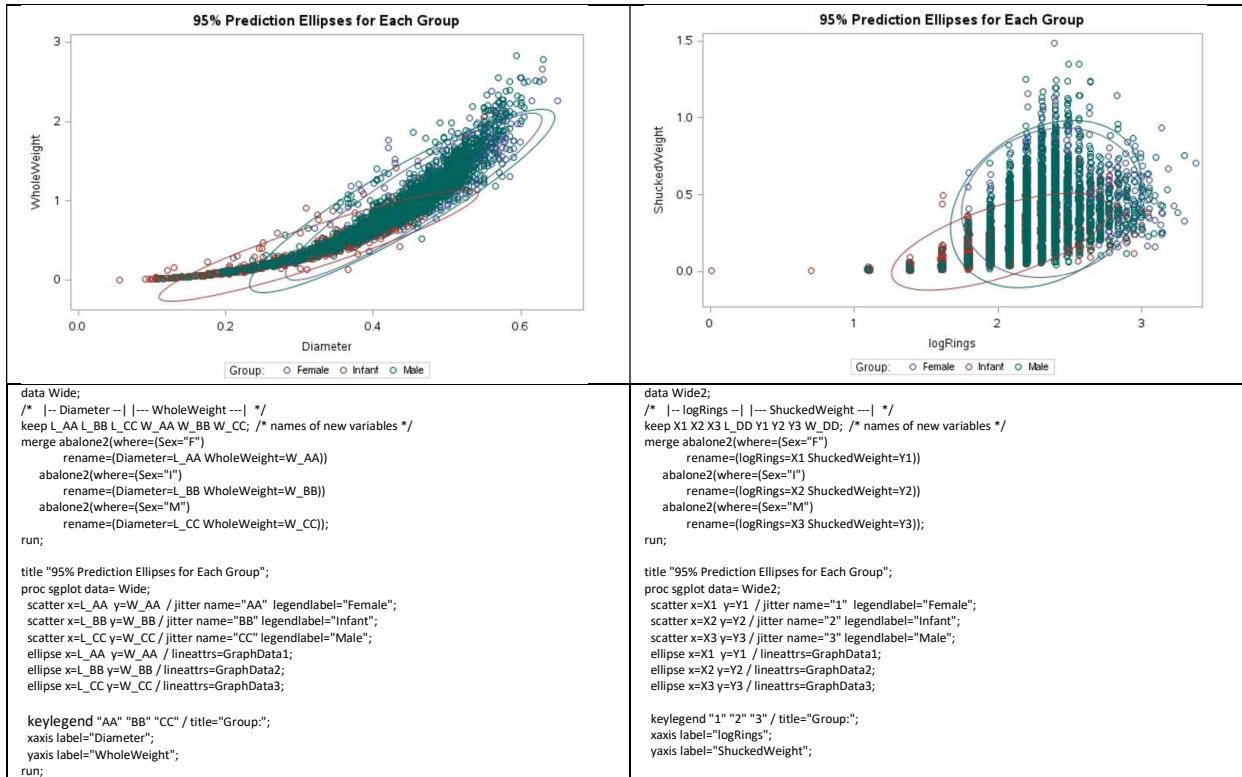


Figure 9, 95% prediction ellipses

### Assumptions (MANOVA/LDA)

Although visual indications were insightful in understanding the variables additional investigation was conducted to determine if the assumptions for a MANOVA/LDA were met. In principle the MANOVA will test whether or not the mean vectors of the three response levels are common and to what level. These investigations are outlined below and include rationale and some selected graphics output with associated code from SAS. The full SAS output and code reference may be found in Appendix B. The assumptions investigated include normality, variance/covariance, independence, and common means vector.

**Normality:** Virgin, log-transformed, square-transformed, and mean-center-transformed data were all investigated in order to attempt to further improve the normality of the residual data. The SAS GLM procedure was used to generate individual ANOVAs in order to generate fit diagnostics plots for each variable. This was executed for virgin, log-transformed, mean-center transformed as well as mean-center transformed data. Primarily the residuals, histogram and qq plots were utilized to determine normality. Given the large sample size and the central limit theorem we will proceed with the raw data with the exception of the "Rings" variable which remains log transformed for the remainder of this study. The "Rings" variable output (pre and post log) is shown in figure 10, while several other investigations may be viewed in Appendix B. The improvement from "Rings" to "LogRings" is primarily

evident in the residual plots, qq plots and in the adjusted R-squared. For all other variables the qq-plots seemed to invert causing an equally “non-normal” condition. Mean-Center and Square transformations did not yield improvement to the normality of the residuals. In many cases the same inversion occurred on the qq plots or the fit simply did not improve. These efforts and associated code may also be viewed in Appendix B.

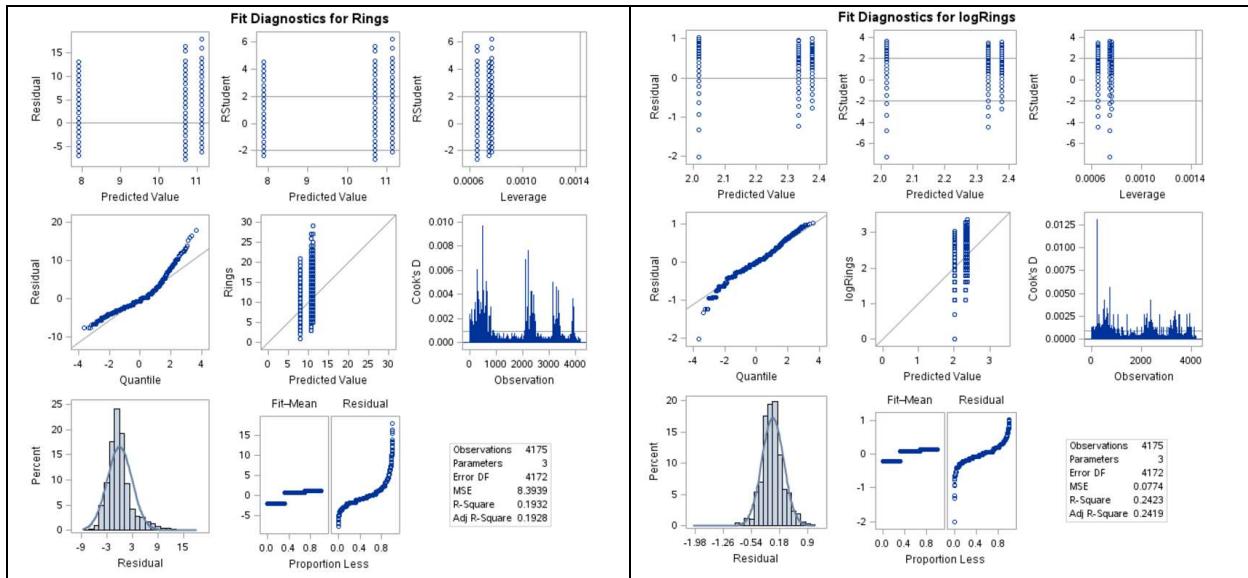


Figure 10, transformation example for Rings variable

Variance/Covariance: Once the state of variable transformation was determined the next assumption to be verified was the variance/covariance assumption. Prior 95% prediction ellipses indicated a difference in the Infant level to the two other levels (Male and Female). It also indicated similar variance/covariance between the Male and Female groups. The SAS code and output noted in figure 11 provide evidence that variance/covariance matrices may not be the same as indicated by the p-value on the Chi-Square test of less than 0.0001. Although the matrices cannot be confirmed to be the same (even though visual indications indicate they may be), this does not rule out the use of an LDA as SAS will use a quadratic, versus liner, method which may assume different matrices.

<pre>proc discrim data=abalone2 pool=test; class sex; var logRings Length Diameter Volume Height WholeWeight ShuckedWeight VisceraWeight ShellWeight; run;</pre>						
<b>The DISCRIM Procedure</b> <b>Test of Homogeneity of Within Covariance Matrices</b> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>Chi-Square</th> <th>DF</th> <th>Pr &gt; ChiSq</th> </tr> </thead> <tbody> <tr> <td>2465.893308</td> <td>20</td> <td>&lt;.0001</td> </tr> </tbody> </table>	Chi-Square	DF	Pr > ChiSq	2465.893308	20	<.0001
Chi-Square	DF	Pr > ChiSq				
2465.893308	20	<.0001				
<p>Since the Chi-Square value is significant at the 0.1 level, the within covariance matrices will be used in the discriminant function.  Reference: Morrison, D.F. (1976) Multivariate Statistical Methods p252.</p> <p>There is a difference among the covariance matrices. Also noted from SAS output is an indication that the within covariance matrices will be used in the discriminant function.</p>						

Figure 11, Variance/Covariance equality test

Independence: Since there were separate shells measured and none were repeated we will assume the independence assumption is made. All samples were gathered from the Bass Strait in an area north or Tasmania. Therefore the independence is limited to this region of data collection and may not be

necessarily extended outside of this region. For purposes of this study, all data is assumed to be randomly collected as there is no provided rationale to assume that they are not.

Although variance/covariance are not confirmed common we will proceed with MANOVA F-Tests, individual ANOVA's and LDA to further understand any differences between the Sex group mean vectors.

## MANOVA

A MANOVA was run to understand where the differences in the Sex group mean vectors. This knowledge of differences may be used to understand contributors in the LDA model. All of the test criteria noted in the SAS output in figure 12 indicate a rejected  $H_0$  (all p-values are  $<.0001$ ) that the mean vectors are the same. This means that there is difference in the Sex levels per the p-values.

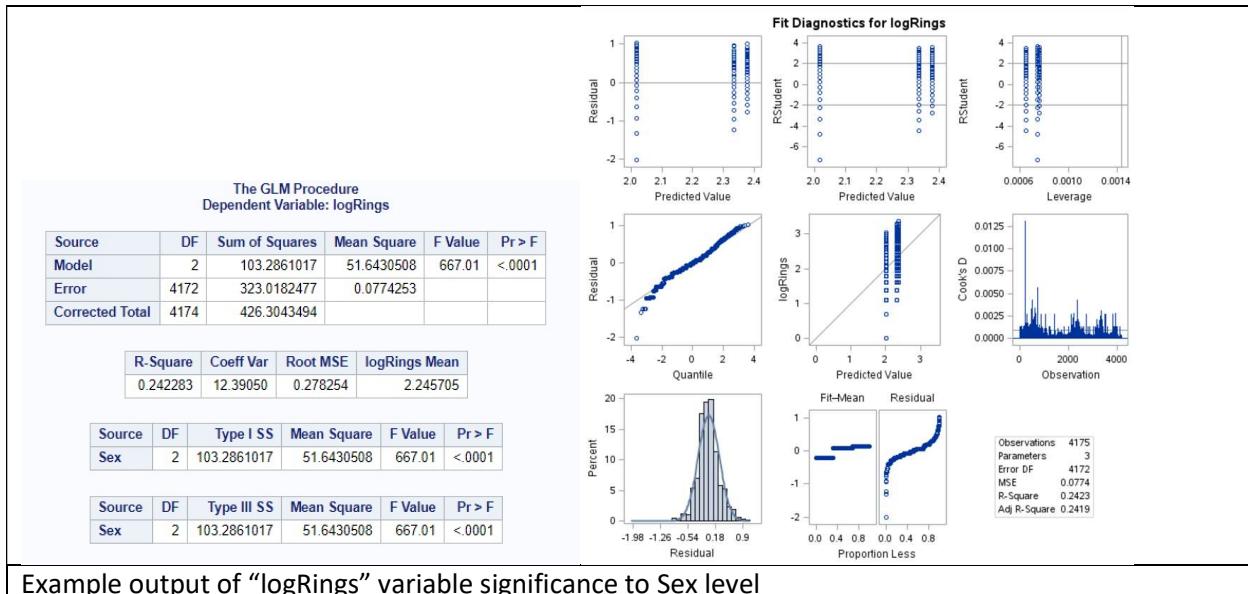
```
/* MANOVA after Log Transformation */
Proc GLM Data=abalone2 plots = all;
  Class Sex;
  Model logRings Length Diameter Height Volume WholeWeight
    ShuckedWeight VisceraWeight ShellWeight = Sex;
  Manova H=_All_ / PrintE PrintH Canonical;
Run;
```

MANOVA Test Criteria and F Approximations for the Hypothesis of No Overall Sex Effect					
$H$ = Type III SSCP Matrix for Sex $E$ = Error SSCP Matrix					
S=2 M=3 N=2081					
Statistic	Value	F Value	Num DF	Den DF	Pr > F
Wilks' Lambda	0.61800157	125.87	18	8328	<.0001
Pillai's Trace	0.38633119	110.79	18	8330	<.0001
Hotelling-Lawley Trace	0.61110793	141.34	18	6936.2	<.0001
Roy's Greatest Root	0.59941158	277.39	9	4165	<.0001

NOTE: F Statistic for Roy's Greatest Root is an upper bound.  
 NOTE: F Statistic for Wilks' Lambda is exact.

Figure 12, MANOVA comparing Sex level mean vectors. There is a difference.

Since there is a difference in mean vectors we proceed to several one-way ANOVA F-Tests to find which variables are potentially driving the differences. The SAS code and sample output may be seen in figure 13 (the SAS code and remaining output may be viewed in Appendix C). Per the output, all variables are significantly and individually influencing (p-value  $<0.0001$ ) the difference in response (Sex) levels.



Example output of “logRings” variable significance to Sex level

Figure 13, ANOVA F-Test example

Based on the MANOVA and individual ANOVA tests there must be a group mean difference so the individual mean test is performed using LSMEANS procedure. The SAS code and output may be viewed in figure 14. The Bonferroni adjustment is used in order to match Type-I error rate as 0.05 with multiple pair wise tests on a single set of data. In this test, every pair-wise test for every dependent variable shows that the differences are highly significant except one case. For the ShuckedWeight variable, Sex level 1 and 3 (which is Male and Female) shows no significant difference (p-value 0.1601). There is, however, a significant difference (p-value <0.0001) between the Infant and others (Male and Female).

<pre>Proc GLM Data=abalone2 plots = all plots=all; Class Sex; Model logRings Length Diameter Height Volume WholeWeight ShuckedWeight VisceraWeight ShellWeight = Sex; Manova H=_All_ / PrintE PrintH Canonical; Ismeans Sex / adjust=bon cl; Run;</pre>			
<b>Least Squares Means for effect Sex</b> $Pr >  t $ for $H_0: \text{LSMean}(i) = \text{LSMean}(j)$		<b>Least Squares Means for effect Sex</b> $Pr >  t $ for $H_0: \text{LSMean}(i) = \text{LSMean}(j)$	
<b>Dependent Variable: logRings</b>		<b>Dependent Variable: Length</b>	
i/j	1	2	3
1		<.0001	0.0002
2	<.0001		<.0001
3	0.0002	<.0001	
<b>Least Squares Means for effect Sex</b> $Pr >  t $ for $H_0: \text{LSMean}(i) = \text{LSMean}(j)$		<b>Least Squares Means for effect Sex</b> $Pr >  t $ for $H_0: \text{LSMean}(i) = \text{LSMean}(j)$	
<b>Dependent Variable: Diameter</b>		<b>Dependent Variable: Diameter</b>	
i/j	1	2	3
1		<.0001	<.0001
2	<.0001		<.0001
3	<.0001	<.0001	
<b>Least Squares Means for effect Sex</b> $Pr >  t $ for $H_0: \text{LSMean}(i) = \text{LSMean}(j)$		<b>Least Squares Means for effect Sex</b> $Pr >  t $ for $H_0: \text{LSMean}(i) = \text{LSMean}(j)$	
<b>Dependent Variable: Height</b>		<b>Dependent Variable: Volume</b>	
i/j	1	2	3
1		<.0001	<.0001
2	<.0001		<.0001
3	<.0001	<.0001	
<b>Least Squares Means for effect Sex</b> $Pr >  t $ for $H_0: \text{LSMean}(i) = \text{LSMean}(j)$		<b>Least Squares Means for effect Sex</b> $Pr >  t $ for $H_0: \text{LSMean}(i) = \text{LSMean}(j)$	
<b>Dependent Variable: WholeWeight</b>		<b>Dependent Variable: WholeWeight</b>	
i/j	1	2	3
1		<.0001	0.0007
2	<.0001		<.0001
3	0.0007	<.0001	

Least Squares Means for effect Sex Pr >  t  for H0: LSMean(i)=LSMean(j)				Least Squares Means for effect Sex Pr >  t  for H0: LSMean(i)=LSMean(j)				Least Squares Means for effect Sex Pr >  t  for H0: LSMean(i)=LSMean(j)			
Dependent Variable: ShuckedWeight				Dependent Variable: VisceraWeight				Dependent Variable: ShellWeight			
i/j	1	2	3	i/j	1	2	3	i/j	1	2	3
1		<.0001	0.1601	1		<.0001	<.0001	1		<.0001	<.0001
2	<.0001		<.0001	2	<.0001		<.0001	2	<.0001		<.0001
3	0.1601	<.0001		3	<.0001	<.0001		3	<.0001	<.0001	

Figure 14, SAS code and output for LSMEANS testing.

Based on the significance of the MANOVA, ANOVA, and LSMEANS tests and the indications in the visual inspection of variables we will use all the original variables in the initial LDA including Volume and log transformed "Rings".

Purely out of academic interest and to help visualize practical versus statistical significance between the Male and Female levels we also decided to perform two simple t-tests. The output may be viewed in figure 15. One t-test compares the "Volume" between the Male and Female levels. The second t-test compares the "Wholeweight" between the Male and Female levels. In both examples, the difference in means is apparent and statistically significant. However, there is significant overlap of the two histograms and boxplots indicating the difference may not be practical. For example, one could gather Wholeweight measurements that fall in the overlapping area and not be able to discern whether it came from the Male or Female. We will take this into consideration in looking at the LDA results.

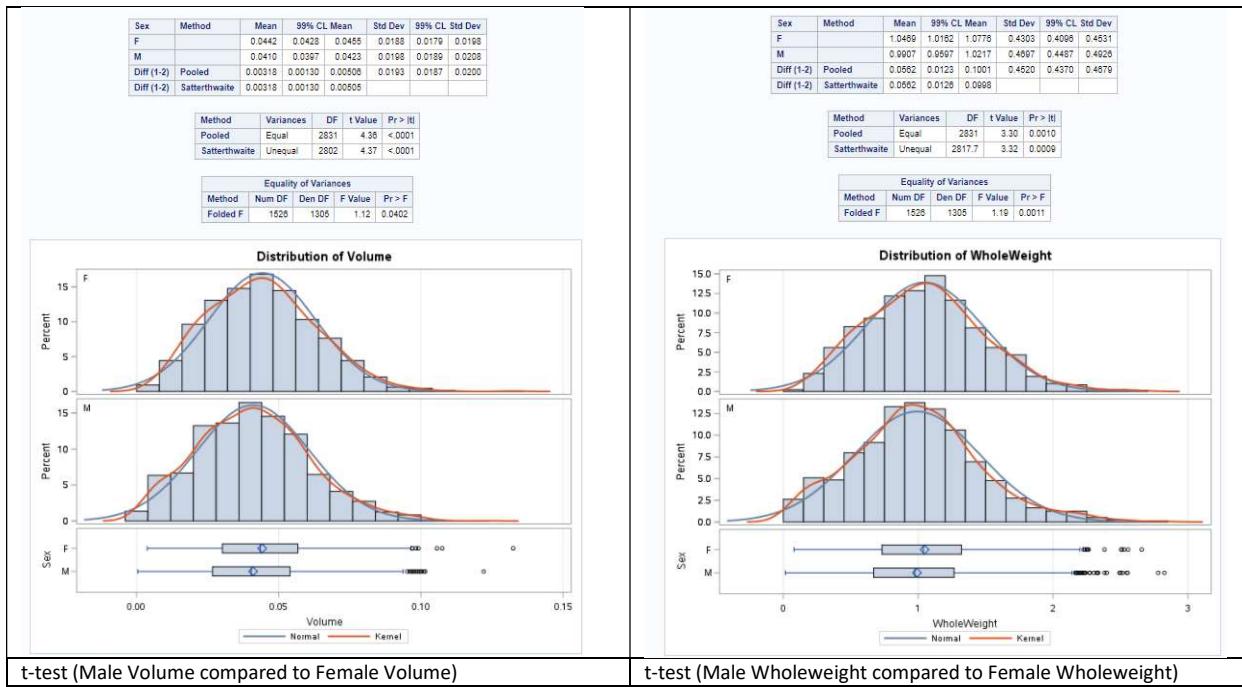


Figure 15, t-test for practical significance understanding

## Method for LDA

Based on the above visual and statistical examinations we will proceed with the Linear Discriminant Analysis (LDA) and cross-validation to determine if Male and Female Abalones may be distinguished and if the Infant Abalones may be distinguished by collecting measurements consistent with the data set described earlier in this paper. This Baseline data set will include the original data including "Volume" and a log transformed "Rings" variable. Once the LDA is run we will report the error rates and interpret the results with respect to the questions of interest. The LDA was run using the SAS DISCRIM procedure. The code and output may be seen in figure 16. The Baseline LDA indicates an overall error rate of 46.11%, meaning that it is correct, per the cross-validation, 53.89% of the time in the Bass Strait region of Ocean for the black-lip Abalone. As the t-test visuals and previous box-plots indicate there may be statistical significance in the mean vectors, however the large overlaps cause a practical concern where many of the variables collected are in the overlap region and may easily fall into the 46.11% error side of the discrimination.

<pre>PROC DISCRIM data = abalone2 CROSSVALIDATE OUTSTAT=DIS_FUNC POOL=test wcov pcov; CLASS sex; VAR LogRings Length Diameter Height Volume WholeWeight ShuckedWeight VisceraWeight ShellWeight; PRIORS EQUAL; run; quit;</pre>																																																							
<b>The DISCRIM Procedure</b> Classification Summary for Calibration Data: WORK.ABALONE2 Cross-validation Summary using Quadratic Discriminant Function																																																							
<table border="1"> <thead> <tr> <th colspan="5">Number of Observations and Percent Classified into Sex</th> </tr> <tr> <th>From Sex</th> <th>F</th> <th>I</th> <th>M</th> <th>Total</th> </tr> </thead> <tbody> <tr> <td>F</td> <td>871 66.69</td> <td>328 25.11</td> <td>107 8.19</td> <td>1306 100.00</td> </tr> <tr> <td>I</td> <td>164 12.22</td> <td>1147 85.47</td> <td>31 2.31</td> <td>1342 100.00</td> </tr> <tr> <td>M</td> <td>916 59.99</td> <td>466 30.52</td> <td>145 9.60</td> <td>1527 100.00</td> </tr> <tr> <td>Total</td> <td>1951 46.73</td> <td>1941 48.49</td> <td>283 6.78</td> <td>4175 100.00</td> </tr> <tr> <td>Priors</td> <td>0.33333</td> <td>0.33333</td> <td>0.33333</td> <td></td> </tr> </tbody> </table> <table border="1"> <thead> <tr> <th colspan="5">Error Count Estimates for Sex</th> </tr> <tr> <th></th> <th>F</th> <th>I</th> <th>M</th> <th>Total</th> </tr> </thead> <tbody> <tr> <td>Rate</td> <td>0.3331</td> <td>0.1453</td> <td>0.9050</td> <td>0.4811</td> </tr> <tr> <td>Priors</td> <td>0.3333</td> <td>0.3333</td> <td>0.3333</td> <td></td> </tr> </tbody> </table>	Number of Observations and Percent Classified into Sex					From Sex	F	I	M	Total	F	871 66.69	328 25.11	107 8.19	1306 100.00	I	164 12.22	1147 85.47	31 2.31	1342 100.00	M	916 59.99	466 30.52	145 9.60	1527 100.00	Total	1951 46.73	1941 48.49	283 6.78	4175 100.00	Priors	0.33333	0.33333	0.33333		Error Count Estimates for Sex						F	I	M	Total	Rate	0.3331	0.1453	0.9050	0.4811	Priors	0.3333	0.3333	0.3333	
Number of Observations and Percent Classified into Sex																																																							
From Sex	F	I	M	Total																																																			
F	871 66.69	328 25.11	107 8.19	1306 100.00																																																			
I	164 12.22	1147 85.47	31 2.31	1342 100.00																																																			
M	916 59.99	466 30.52	145 9.60	1527 100.00																																																			
Total	1951 46.73	1941 48.49	283 6.78	4175 100.00																																																			
Priors	0.33333	0.33333	0.33333																																																				
Error Count Estimates for Sex																																																							
	F	I	M	Total																																																			
Rate	0.3331	0.1453	0.9050	0.4811																																																			
Priors	0.3333	0.3333	0.3333																																																				
46.11% Error Rate – Baseline variable set																																																							

Figure 16, LDA and percent error

Let's look a little closer at the questions of interest.

Can we discriminate Male from Female Black-Lip Abalones? The error of discriminating an Abalone in the Female group is 33.31% and in the Male group is 90.5%, meaning that the Female discrimination is correct 66.69% of the time and Male discrimination is correct only 9.5% of the time. This method of discriminating Male and Female Black-Lip Abalones in the Bass Strait region of ocean would not be very reliable or practical. In fact, the number of measurements, even if collected autonomously, would probably not rival the more manual method of checking the gonad color.

Can we discriminate Infant from Male and Female Black-Lip Abalones? The error of discriminating an Abalone in the Infant group is 14.53%, meaning that it is correct, per the cross-validation, 85.47% of the time in the Bass Strait region of ocean for the black-lip Abalone. This may be very useful to fisherman attempting to discriminate if they should throw the Abalone back or keep it. In fact, this is what

fisherman do (take a measurement and decide). Further investigation may likely reveal that the error occurs in Infants that are near the size and weight of the matured Male or Female Abalones. In this case many may simply be checked by eye and thrown back without question, thus speeding up the process and improving profits through eliminating non-value added and time consuming measurements. In the event the price of Abalones continues to rise automation may further improve this process. This level of discrimination is useful with respect to the original question posed.

### Auto-selection LDA

For academic purposes and in order to determine if a simplified LDA model could discriminate Male/Female as well or better with a reduced data set, the data set was reduced to only include Male and Female data entries and then the SAS STEPDISC procedure was utilized to auto-select variables for use in a revised LDA. The thought is to attempt to build either an improved model (less error) or a much simpler approach to future data collection utilizing fewer variables. For reference STEPDISC selection functions in the following manner (ref: support.sas.com website):

- selection methods include forward selection, backward elimination, and stepwise selection
- variables are chosen to enter or leave the model according to one of two criteria:
  - the significance level of an *F* test from an analysis of covariance, where the variables already chosen act as covariates and the variable under consideration is the dependent variable
  - the squared partial correlation for predicting the variable under consideration from the CLASS variable, controlling for the effects of the variables already selected for the model

The STEPDISC procedure resulted in a much smaller subset of variables (Diameter, Height, ShuckedWeight, and VisceraWeight) listed in the SAS output in figure 17.

<pre>Code reference for stepwise selection and results. *automatic variable selection using stepdisc for LDA variable selection; proc stepdisc data=abaloneMonly bscp tsscp; class Sex; var LogRings Length Diameter Height Volume WholeWeight ShuckedWeight VisceraWeight ShellWeight; run;</pre>	<p style="text-align: center;"><b>The STEPDISC Procedure</b> <b>Stepwise Selection: Step 5</b></p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th colspan="4" style="text-align: left;">Statistics for Removal, DF = 1, 2828</th> </tr> <tr> <th style="text-align: left;">Variable</th><th style="text-align: left;">Partial R-Square</th><th style="text-align: left;">F Value</th><th style="text-align: left;">Pr &gt; F</th></tr> </thead> <tbody> <tr> <td>Diameter</td><td>0.0077</td><td>21.82</td><td>&lt;.0001</td></tr> <tr> <td>Height</td><td>0.0008</td><td>2.29</td><td>0.1300</td></tr> <tr> <td>ShuckedWeight</td><td>0.0142</td><td>40.67</td><td>&lt;.0001</td></tr> <tr> <td>VisceraWeight</td><td>0.0017</td><td>4.94</td><td>0.0264</td></tr> </tbody> </table> <p style="text-align: center; margin-top: 10px;">No variables can be removed.</p>	Statistics for Removal, DF = 1, 2828				Variable	Partial R-Square	F Value	Pr > F	Diameter	0.0077	21.82	<.0001	Height	0.0008	2.29	0.1300	ShuckedWeight	0.0142	40.67	<.0001	VisceraWeight	0.0017	4.94	0.0264
Statistics for Removal, DF = 1, 2828																									
Variable	Partial R-Square	F Value	Pr > F																						
Diameter	0.0077	21.82	<.0001																						
Height	0.0008	2.29	0.1300																						
ShuckedWeight	0.0142	40.67	<.0001																						
VisceraWeight	0.0017	4.94	0.0264																						

Figure 17, Stepwise auto-selection of LDA variables. Final step shown.

The percent error rate comparing the original LDA to the auto-selected LDA may be viewed in figure 18. The error rate was degraded by 0.4% (from 46.11% to 46.51%) while the number of variables was reduced from ten down to the four auto-selected by SAS.

<pre>PROC DISCRIM data = abalone2 CROSSVALIDATE OUTSTAT=DIS_FUNC POOL=test wcov pcov; CLASS sex; VAR logRings Length Diameter Height Volume WholeWeight ShuckedWeight VisceraWeight ShellWeight; PRIORS EQUAL; run; quit;</pre>	<pre>*rerun LDA using auto-selected variables; PROC DISCRIM data = abaloneMonly CROSSVALIDATE OUTSTAT=DIS_FUNC POOL=test wcov pcov; CLASS sex; VAR Diameter Height ShuckedWeight VisceraWeight; PRIORS EQUAL;run;quit;</pre>																																																																																																												
<p style="text-align: center;"><b>The DISCRIM Procedure</b>  Classification Summary for Calibration Data: WORK.ABALONE2  Cross-validation Summary using Quadratic Discriminant Function</p> <table border="1" data-bbox="236 369 595 601"> <thead> <tr> <th colspan="5">Number of Observations and Percent Classified into Sex</th> </tr> <tr> <th>From Sex</th> <th>F</th> <th>I</th> <th>M</th> <th>Total</th> </tr> </thead> <tbody> <tr> <td>F</td> <td>871 85.69</td> <td>328 25.11</td> <td>107 8.19</td> <td>1306 100.00</td> </tr> <tr> <td>I</td> <td>164 12.22</td> <td>1147 85.47</td> <td>31 2.31</td> <td>1342 100.00</td> </tr> <tr> <td>M</td> <td>916 59.99</td> <td>466 30.52</td> <td>145 9.50</td> <td>1527 100.00</td> </tr> <tr> <td>Total</td> <td>1951 46.73</td> <td>1941 46.49</td> <td>283 6.78</td> <td>4175 100.00</td> </tr> <tr> <td>Priors</td> <td>0.33333</td> <td>0.33333</td> <td>0.33333</td> <td></td> </tr> </tbody> </table> <table border="1" data-bbox="277 633 546 728"> <thead> <tr> <th colspan="4">Error Count Estimates for Sex</th> </tr> <tr> <th></th> <th>F</th> <th>I</th> <th>M</th> <th>Total</th> </tr> </thead> <tbody> <tr> <td>Rate</td> <td>0.3331</td> <td>0.1453</td> <td>0.9050</td> <td>0.4611</td> </tr> <tr> <td>Priors</td> <td>0.3333</td> <td>0.3333</td> <td>0.3333</td> <td></td> </tr> </tbody> </table>	Number of Observations and Percent Classified into Sex					From Sex	F	I	M	Total	F	871 85.69	328 25.11	107 8.19	1306 100.00	I	164 12.22	1147 85.47	31 2.31	1342 100.00	M	916 59.99	466 30.52	145 9.50	1527 100.00	Total	1951 46.73	1941 46.49	283 6.78	4175 100.00	Priors	0.33333	0.33333	0.33333		Error Count Estimates for Sex					F	I	M	Total	Rate	0.3331	0.1453	0.9050	0.4611	Priors	0.3333	0.3333	0.3333		<p style="text-align: center;"><b>The DISCRIM Procedure</b>  Classification Summary for Calibration Data: WORK.ABALONE  Cross-validation Summary using Quadratic Discriminant Function</p> <table border="1" data-bbox="856 369 1199 601"> <thead> <tr> <th colspan="5">Number of Observations and Percent Classified into Sex</th> </tr> <tr> <th>From Sex</th> <th>F</th> <th>I</th> <th>M</th> <th>Total</th> </tr> </thead> <tbody> <tr> <td>F</td> <td>756 57.89</td> <td>380 27.57</td> <td>190 14.55</td> <td>1306 100.00</td> </tr> <tr> <td>I</td> <td>151 11.25</td> <td>1150 85.69</td> <td>41 3.06</td> <td>1342 100.00</td> </tr> <tr> <td>M</td> <td>782 51.21</td> <td>3189 31.89</td> <td>258 18.90</td> <td>1527 100.00</td> </tr> <tr> <td>Total</td> <td>1689 40.46</td> <td>1997 47.83</td> <td>489 11.71</td> <td>4175 100.00</td> </tr> <tr> <td>Priors</td> <td>0.33333</td> <td>0.33333</td> <td>0.33333</td> <td></td> </tr> </tbody> </table> <table border="1" data-bbox="889 623 1158 718"> <thead> <tr> <th colspan="4">Error Count Estimates for Sex</th> </tr> <tr> <th></th> <th>F</th> <th>I</th> <th>M</th> <th>Total</th> </tr> </thead> <tbody> <tr> <td>Rate</td> <td>0.4211</td> <td>0.1431</td> <td>0.8310</td> <td>0.4651</td> </tr> <tr> <td>Priors</td> <td>0.3333</td> <td>0.3333</td> <td>0.3333</td> <td></td> </tr> </tbody> </table>	Number of Observations and Percent Classified into Sex					From Sex	F	I	M	Total	F	756 57.89	380 27.57	190 14.55	1306 100.00	I	151 11.25	1150 85.69	41 3.06	1342 100.00	M	782 51.21	3189 31.89	258 18.90	1527 100.00	Total	1689 40.46	1997 47.83	489 11.71	4175 100.00	Priors	0.33333	0.33333	0.33333		Error Count Estimates for Sex					F	I	M	Total	Rate	0.4211	0.1431	0.8310	0.4651	Priors	0.3333	0.3333	0.3333	
Number of Observations and Percent Classified into Sex																																																																																																													
From Sex	F	I	M	Total																																																																																																									
F	871 85.69	328 25.11	107 8.19	1306 100.00																																																																																																									
I	164 12.22	1147 85.47	31 2.31	1342 100.00																																																																																																									
M	916 59.99	466 30.52	145 9.50	1527 100.00																																																																																																									
Total	1951 46.73	1941 46.49	283 6.78	4175 100.00																																																																																																									
Priors	0.33333	0.33333	0.33333																																																																																																										
Error Count Estimates for Sex																																																																																																													
	F	I	M	Total																																																																																																									
Rate	0.3331	0.1453	0.9050	0.4611																																																																																																									
Priors	0.3333	0.3333	0.3333																																																																																																										
Number of Observations and Percent Classified into Sex																																																																																																													
From Sex	F	I	M	Total																																																																																																									
F	756 57.89	380 27.57	190 14.55	1306 100.00																																																																																																									
I	151 11.25	1150 85.69	41 3.06	1342 100.00																																																																																																									
M	782 51.21	3189 31.89	258 18.90	1527 100.00																																																																																																									
Total	1689 40.46	1997 47.83	489 11.71	4175 100.00																																																																																																									
Priors	0.33333	0.33333	0.33333																																																																																																										
Error Count Estimates for Sex																																																																																																													
	F	I	M	Total																																																																																																									
Rate	0.4211	0.1431	0.8310	0.4651																																																																																																									
Priors	0.3333	0.3333	0.3333																																																																																																										
46.11% error rate using all variables	46.51% error rate using fewer variables via auto-selection procedure in SAS (STEPDISC)																																																																																																												

Figure 18, Comparison of error rate in original LDA to auto-selected LDA variables

## Conclusions

We initially set out to determine if Blacklip Abalones could be discriminated by Infant versus Male and Female and also if they could be discriminated between Male and Female based on several measurement and weight variables. Through a series of both visual and statistical methods it has been determined that the Infant may be discriminated 85.69% (14.31% error) of the time using a Cross-validated and auto-selected data set and SAS's DISCRIM LDA procedure. This means that the four auto-selected measurements would differentiate the Infant the majority of the time and could be useful, especially if the measurement and sorting was automated and occurred upon collection. The Male and Female discrimination was more problematic, meaning the error rates ranged between 33.31% and 90.5% for the base and auto-selected models attempted. Although the MANOVA and one-way ANOVA's indicated statistical significance, the difference was simply not practical as all of the one-way ANOVA's exhibited significant overlap areas. So, a measurement found in this area could potentially be either a Male or Female. Auto-selection techniques did not change the high combined error rates. Our recommendation is to simply continue to utilize current methods to determine the gender of the Blacklip Abalone.

## Appendix A Data Visualization and Assumptions

### Variable Summary Stats

```
proc means data = abalone n mean max min range std fw=8;
var _numeric_;
output out = meansout mean = mean std = std;
title 'Summary Stat';
run;
```

Alphabetic List of Variables and Attributes					
#	Variable	Type	Len	Format	Informat
4	Diameter	Num	8	BEST12.	BEST32.
5	Height	Num	8	BEST12.	BEST32.
1	ID	Num	8	BEST12.	BEST32.
3	Length	Num	8	BEST12.	BEST32.
11	Rings	Num	8	BEST12.	BEST32.
2	Sex	Char	1	\$1.	\$1.
10	ShellWeight	Num	8	BEST12.	BEST32.
8	ShuckedWeight	Num	8	BEST12.	BEST32.
9	VisceraWeight	Num	8	BEST12.	BEST32.
6	Volume	Num	8	BEST12.	BEST32.
7	WholeWeight	Num	8	BEST12.	BEST32.

### The MEANS Procedure

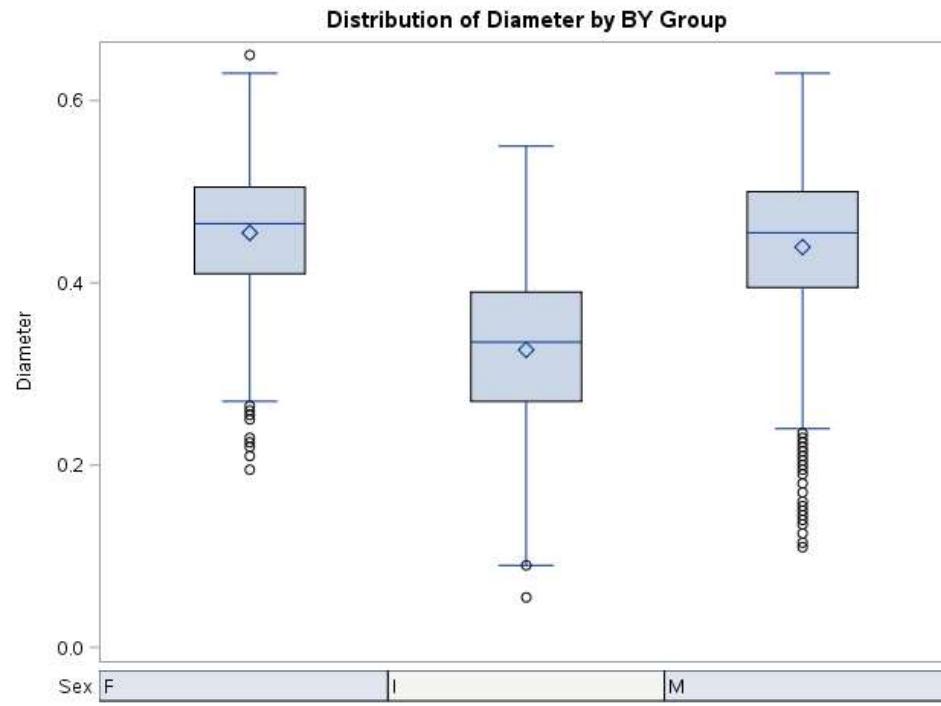
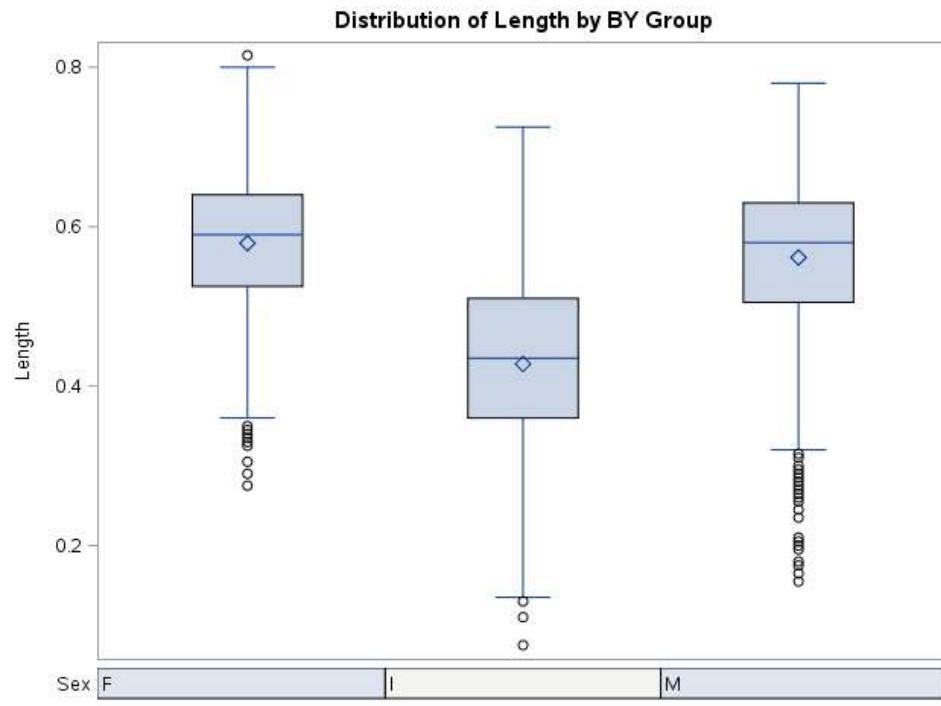
Variable	N	Mean	Maximum	Minimum	Range	Std Dev
ID	4175	2089.2	4177.0	1.0000	4176.0	1206.2
Length	4175	0.5240	0.8150	0.0750	0.7400	0.1201
Diameter	4175	0.4079	0.6500	0.0550	0.5950	0.0992
Height	4175	0.1392	0.2500	0	0.2500	0.0385
Volume	4175	0.0346	0.1324	0	0.1324	0.0209
WholeWeight	4175	0.8285	2.8255	0.00200	2.8235	0.4900
ShuckedWeight	4175	0.3592	1.4880	0.00100	1.4870	0.2217
VisceraWeight	4175	0.1805	0.7600	0.000500	0.7595	0.1095
ShellWeight	4175	0.2388	1.0050	0.00150	1.0035	0.1392
Rings	4175	9.9341	29.0000	1.0000	28.0000	3.2248

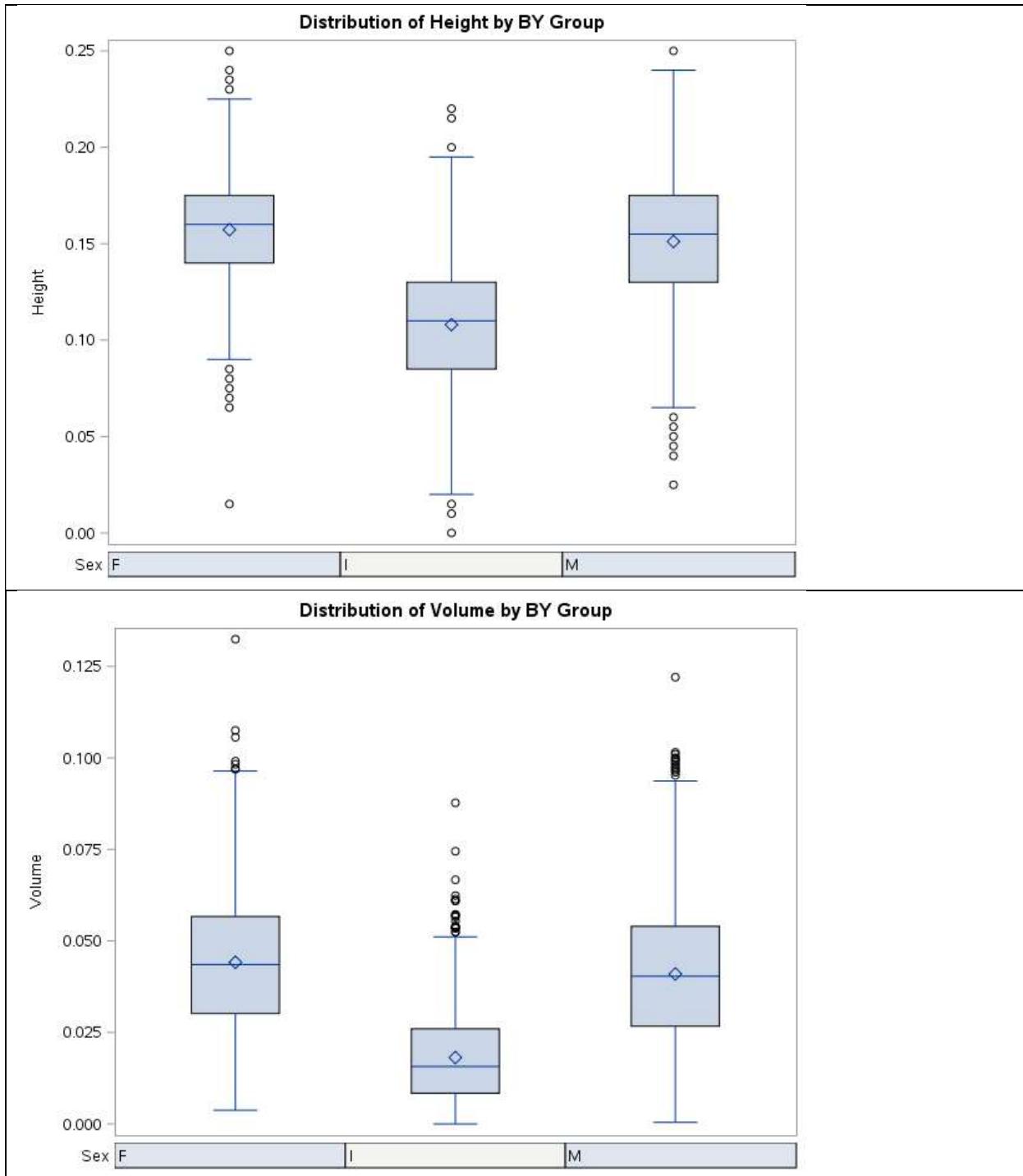
```
proc freq data = abalone;
run;
```

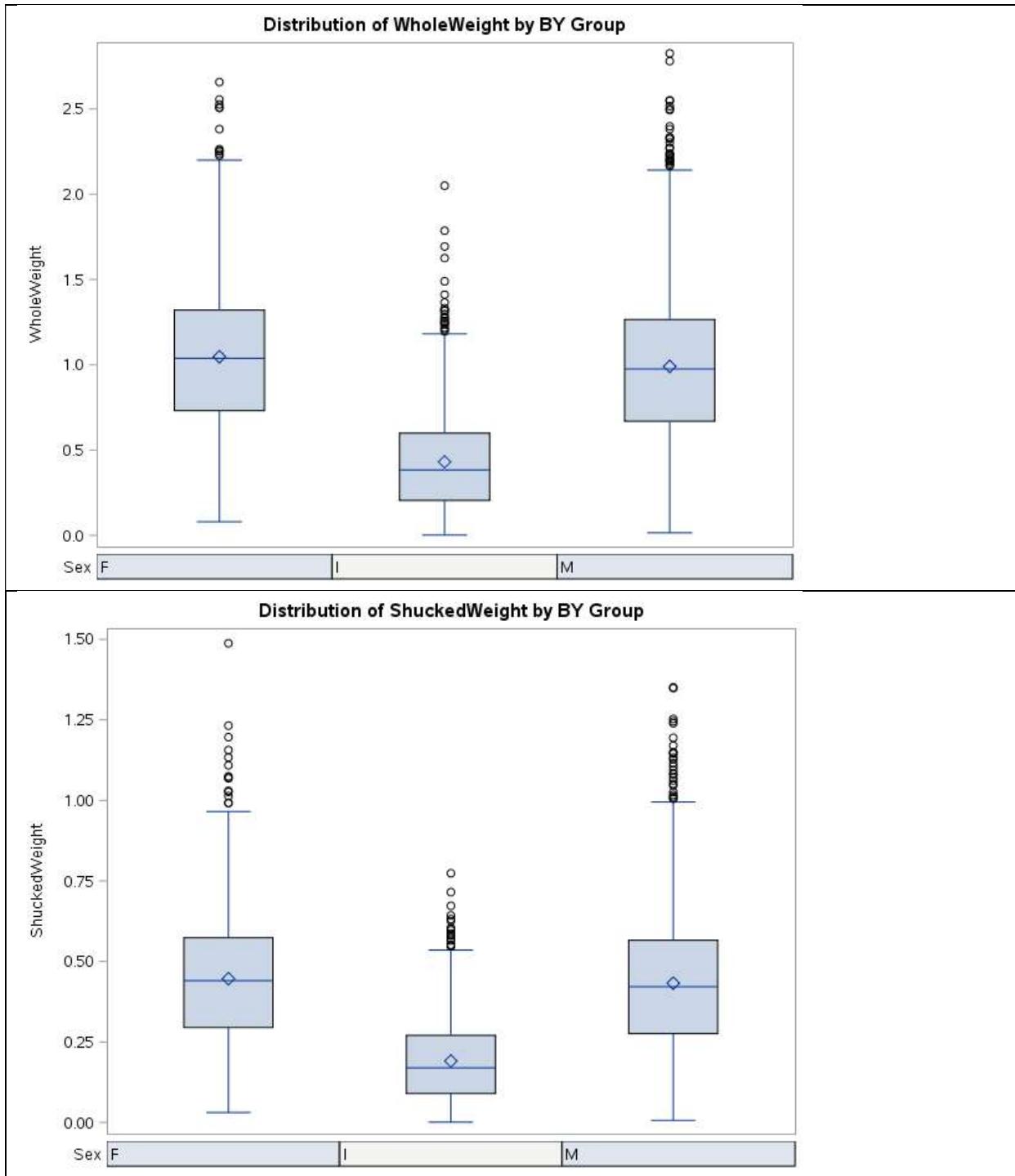
Sex	Frequency	Percent	Cumulative Frequency	Cumulative Percent
F	1306	31.28	1306	31.28
I	1342	32.14	2648	63.43
M	1527	36.57	4175	100.00

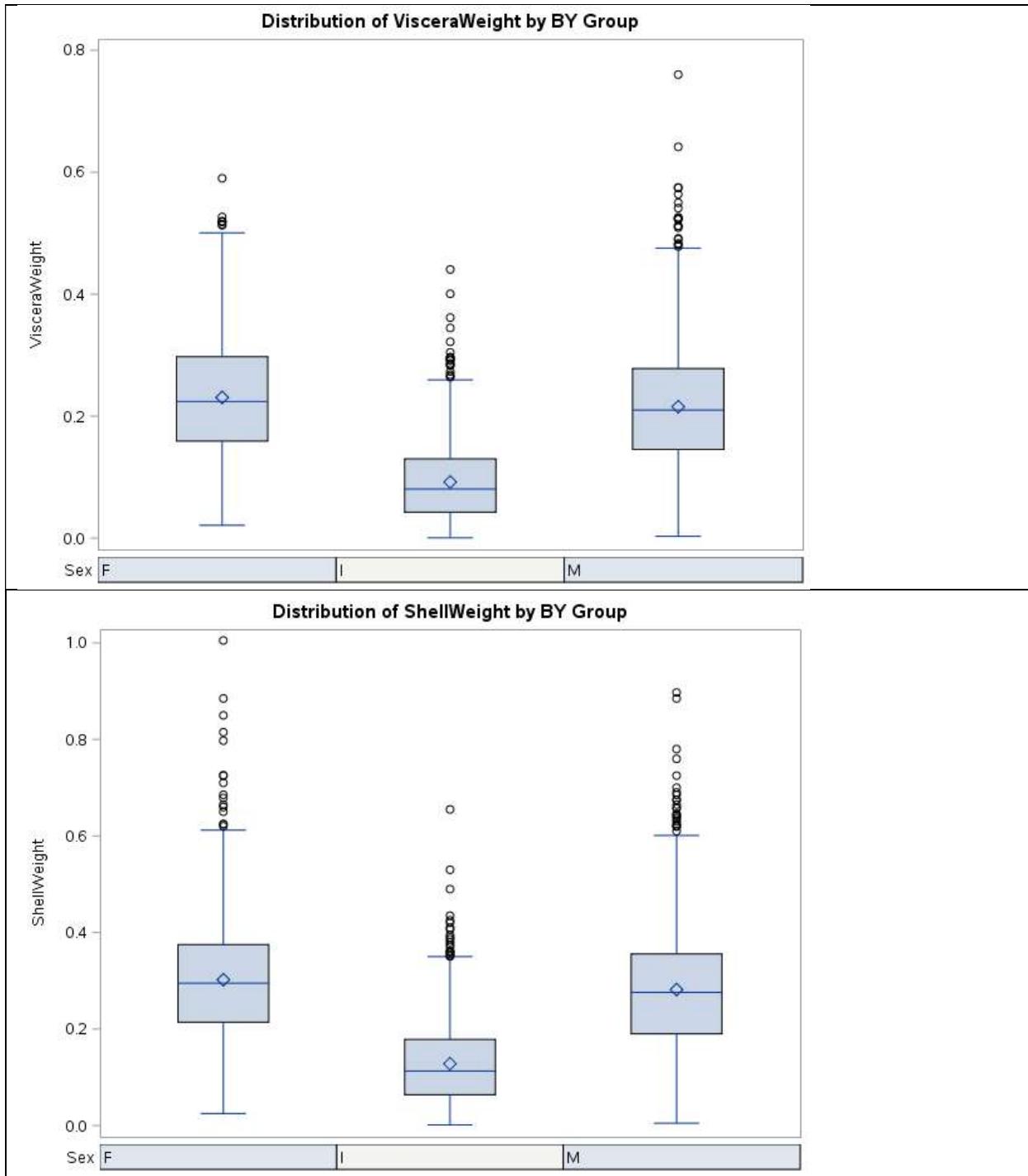
### Box Plot by Sex

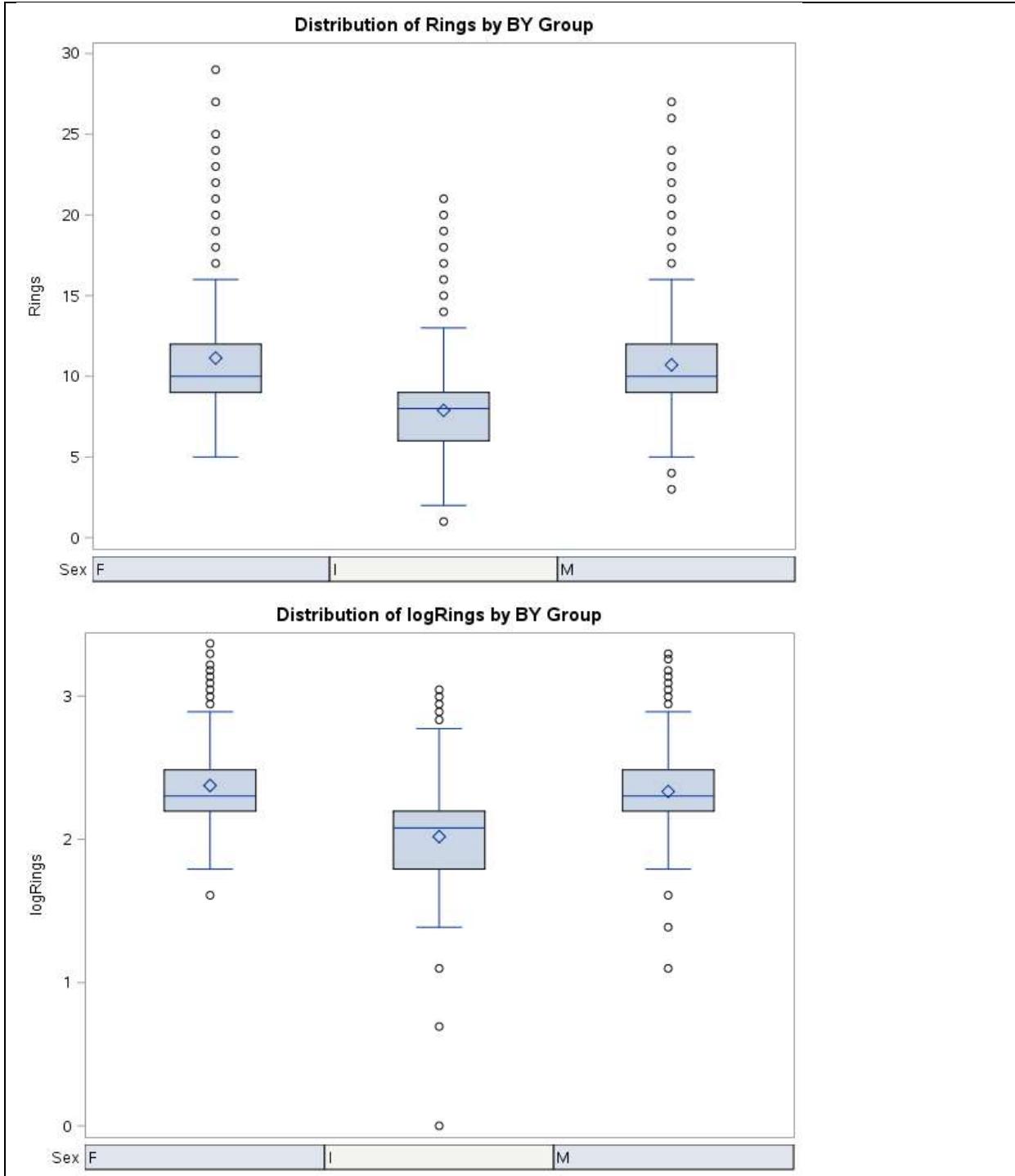
```
proc univariate data = abalone2 plots;  
by Sex;  
run;
```





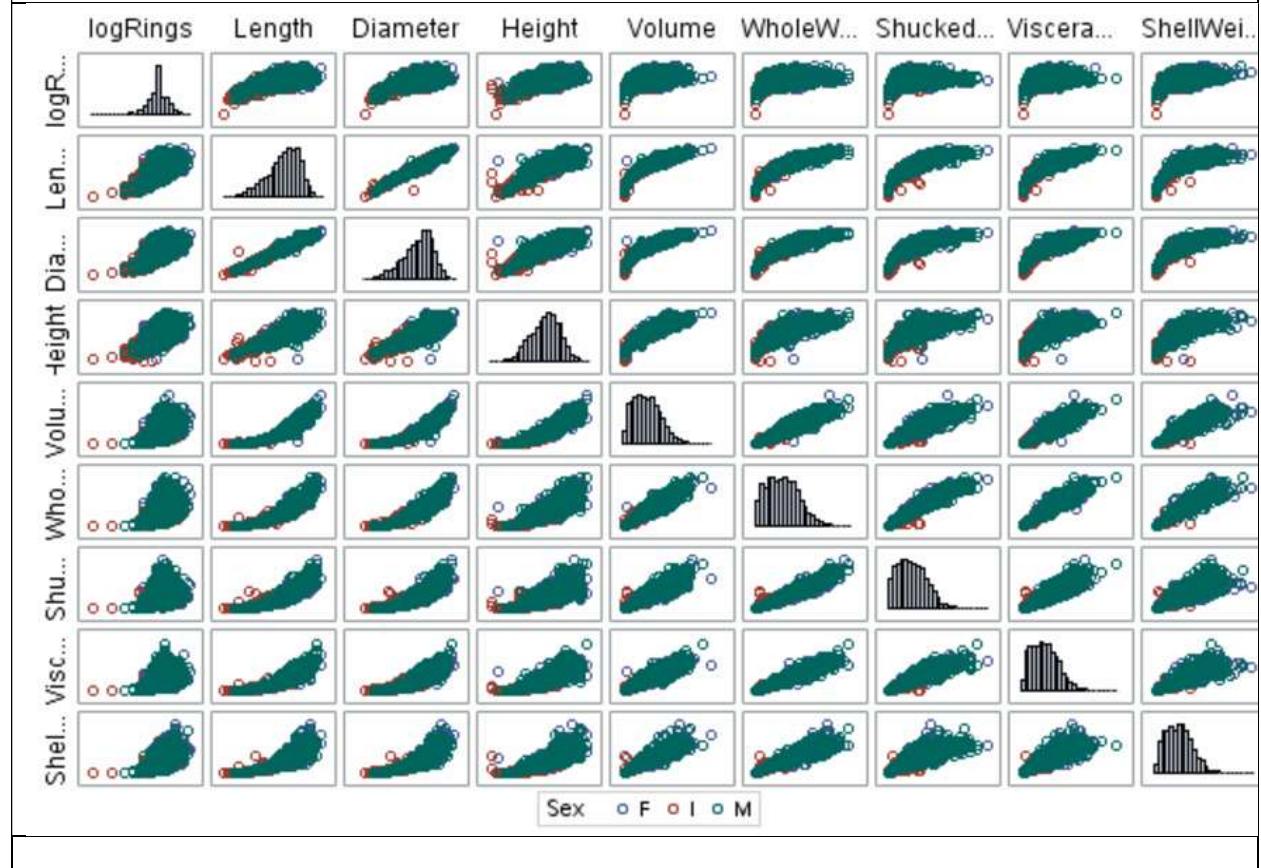






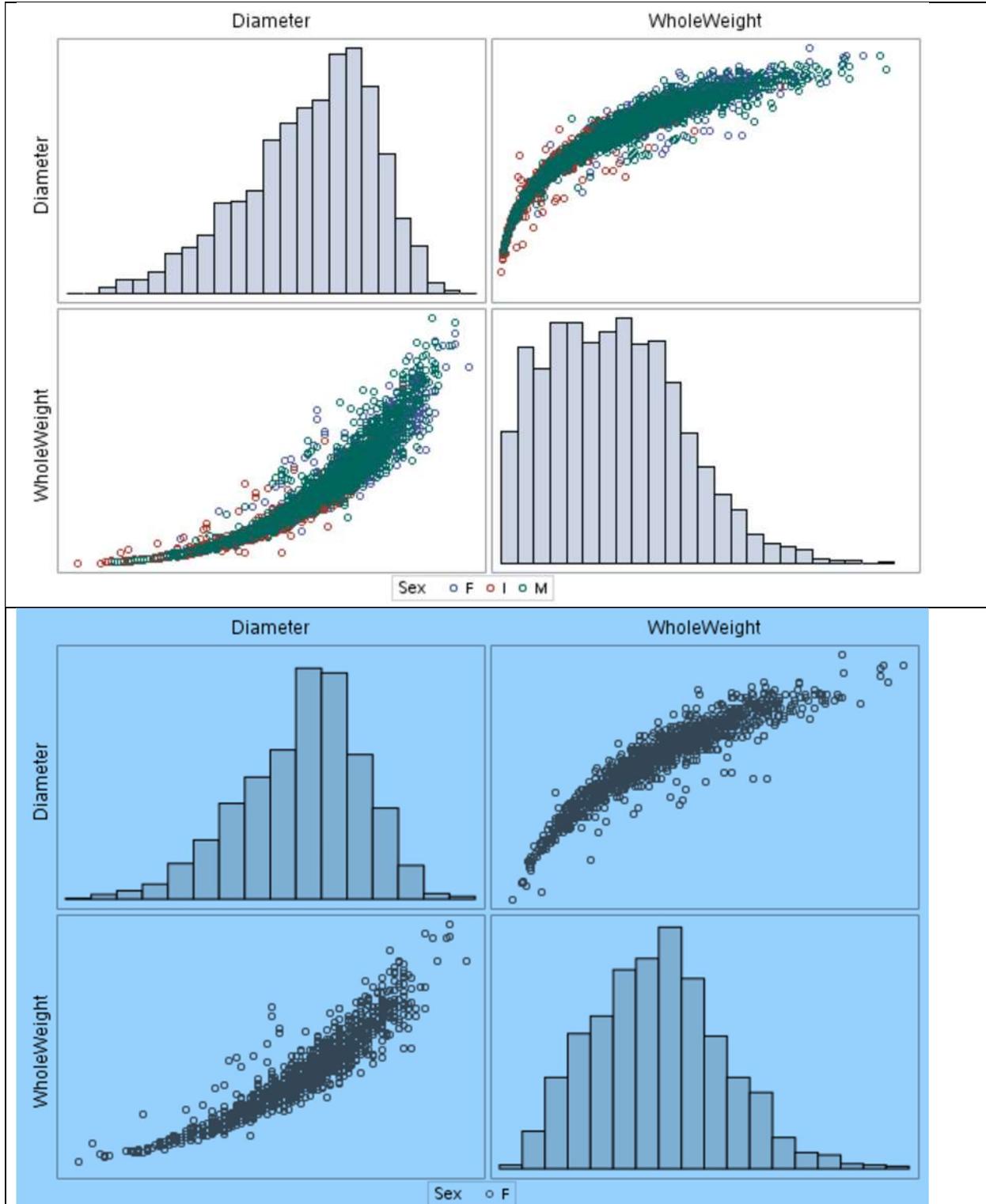
## Scatter plot by Sex: All Variables

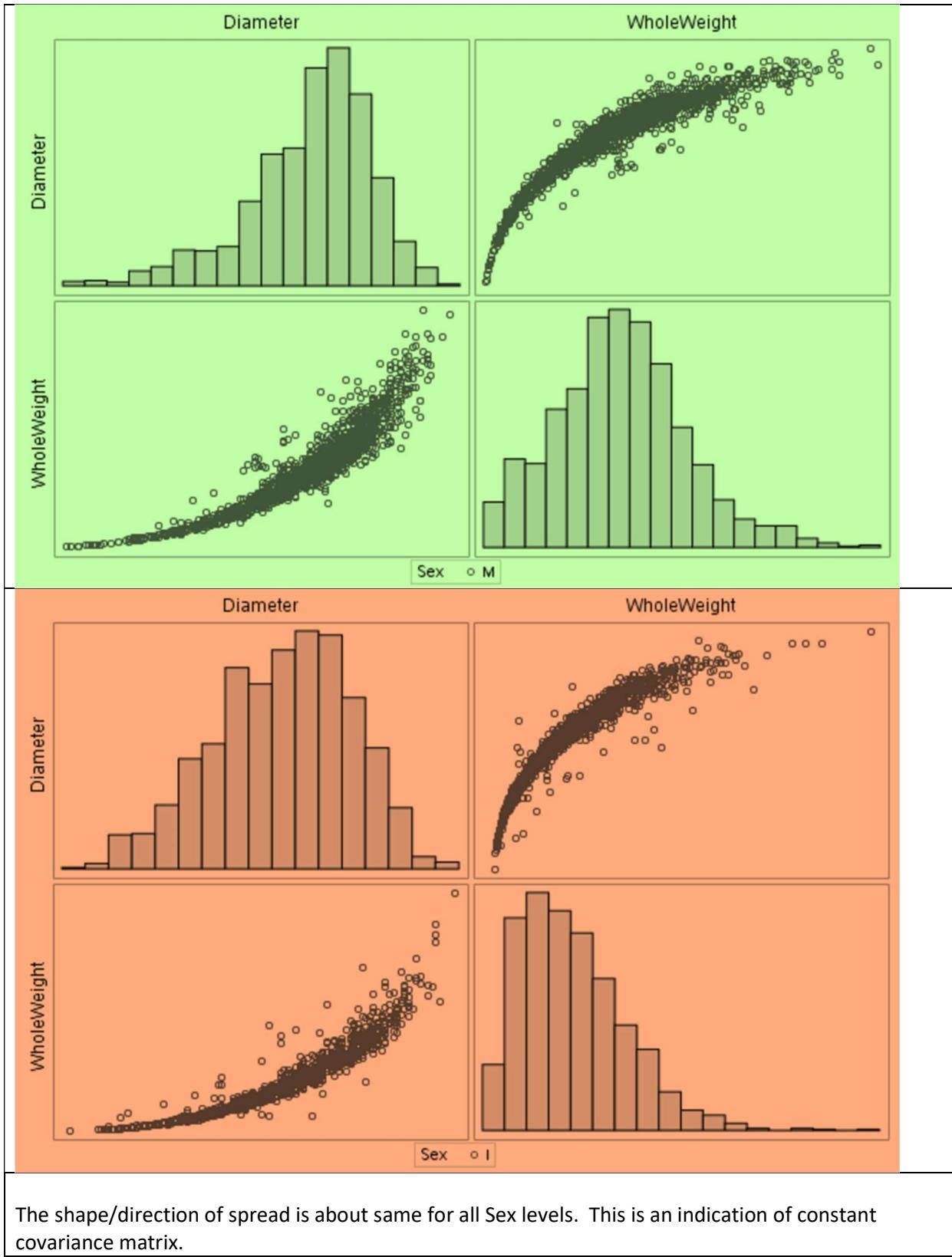
```
proc sgscatter data = abalone2;
matrix logRings Length Diameter Height Volume WholeWeight ShuckedWeight VisceraWeight
ShellWeight/ diagonal=(histogram) group = Sex;
run;quit;
```



## Scatter plot by Sex: { Diameter WholeWeight }

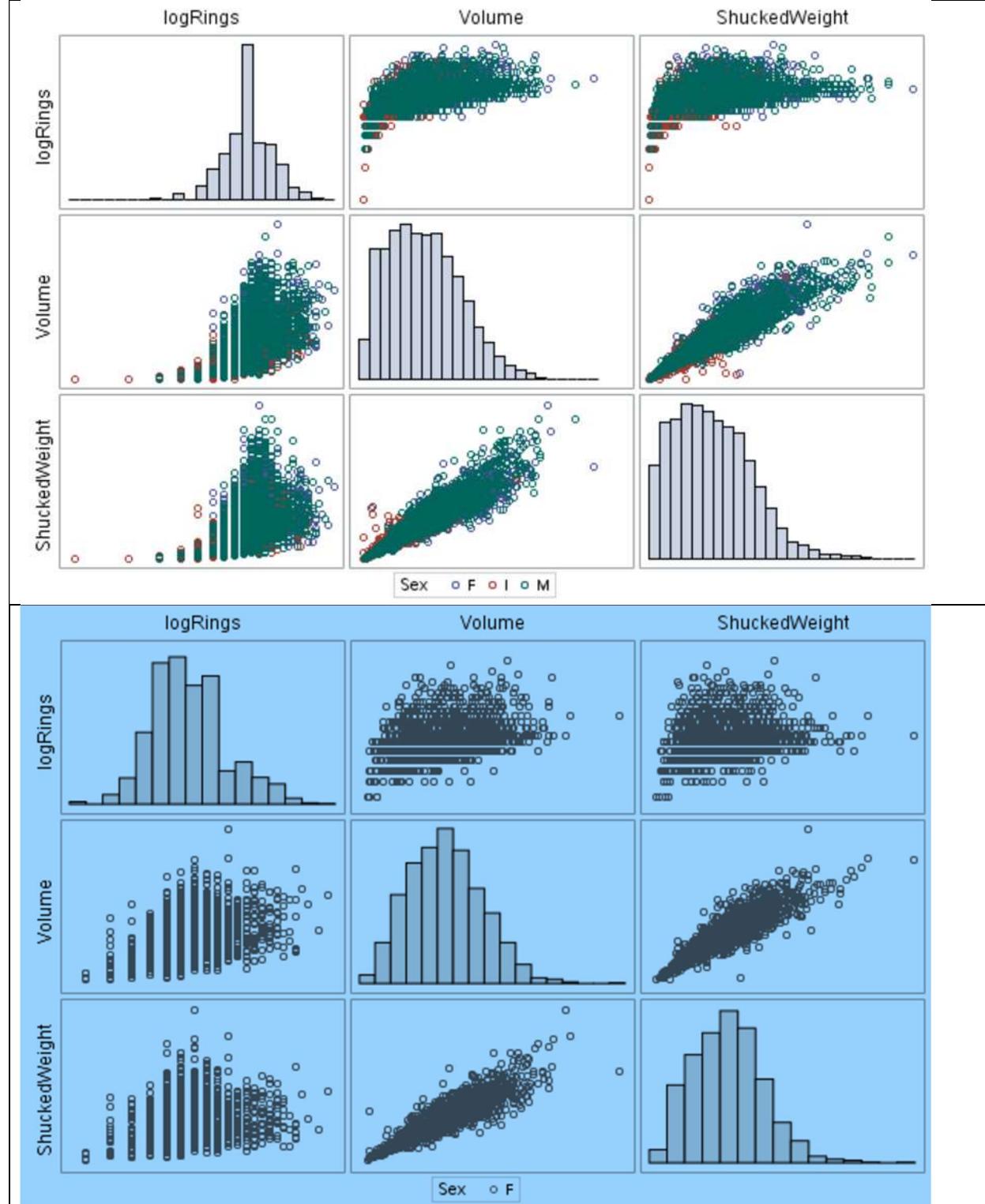
```
proc sgscatter data = abalone2;
matrix Diameter WholeWeight/ diagonal=(histogram) group = Sex;
run;quit;
```

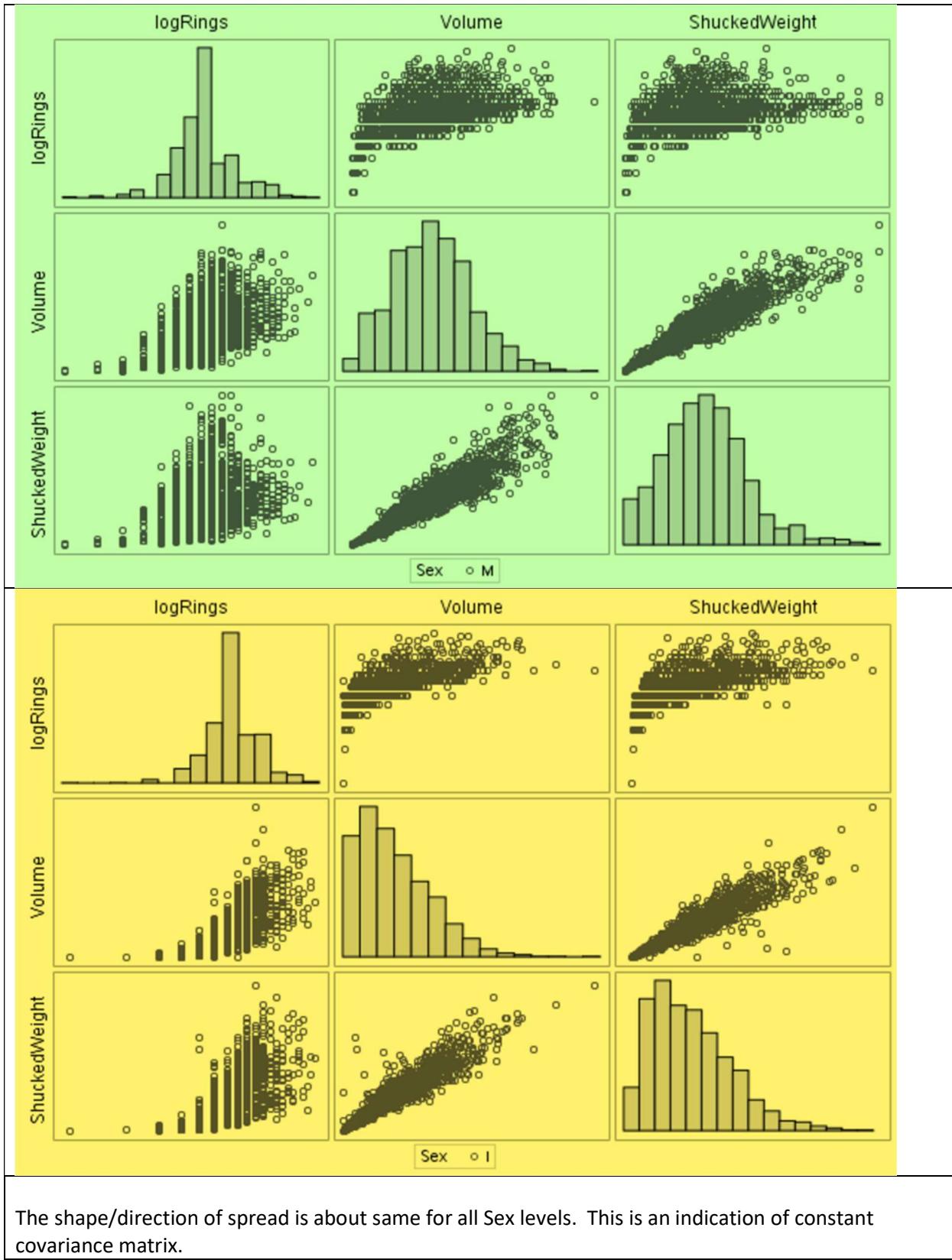




## Scatter Plot by Sex: { logRings Volume ShuckedWeight }

```
proc sgscatter data = abalone2;
matrix logRings Volume ShuckedWeight / diagonal=(histogram) group = Sex;
run;quit;
```







## Appendix B (Assumptions investigation output and SAS code)

MANOVA is used to see the model error distributions.

```
Proc GLM Data=abalone2 plots = all;
```

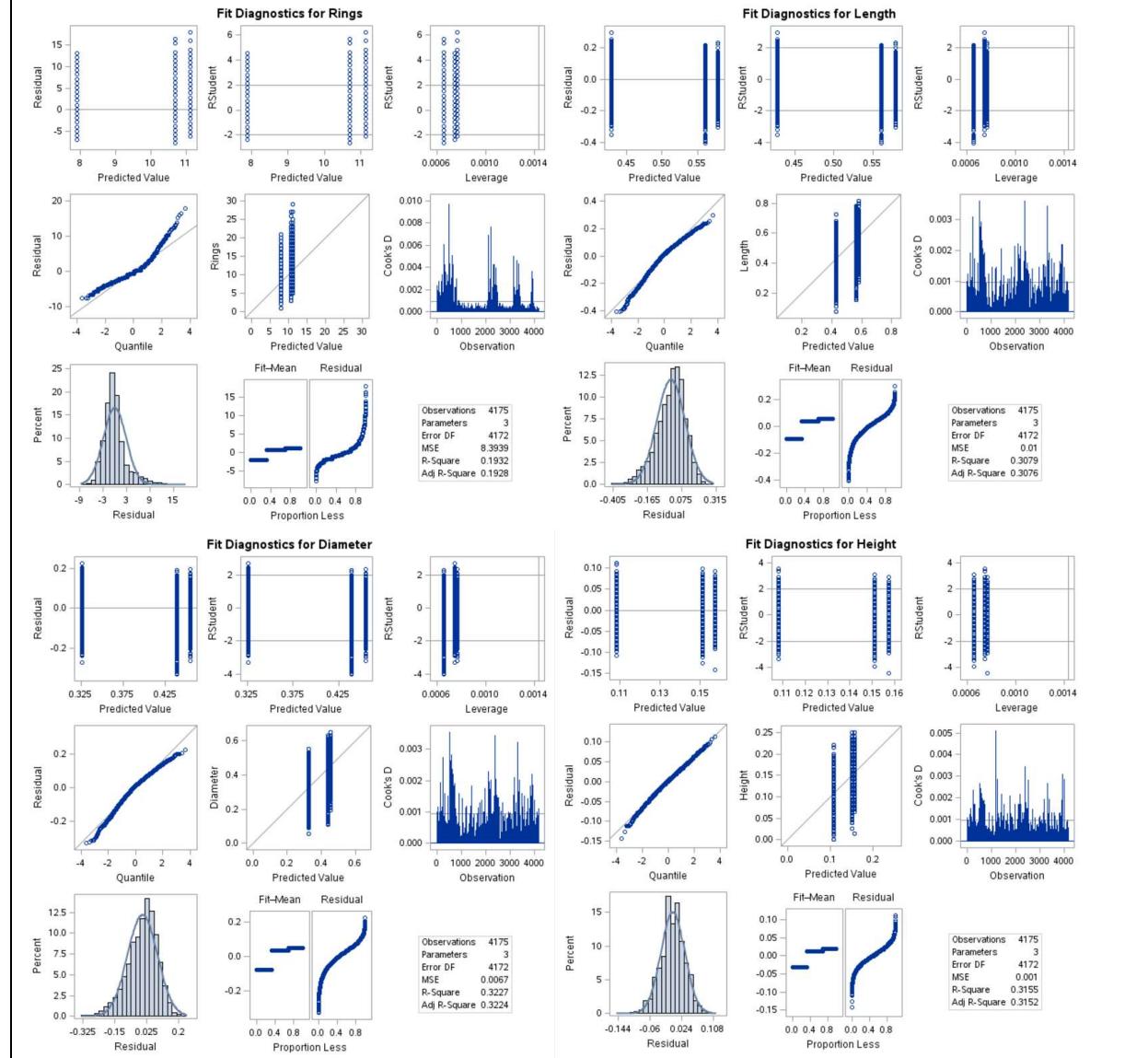
```
Class Sex;
```

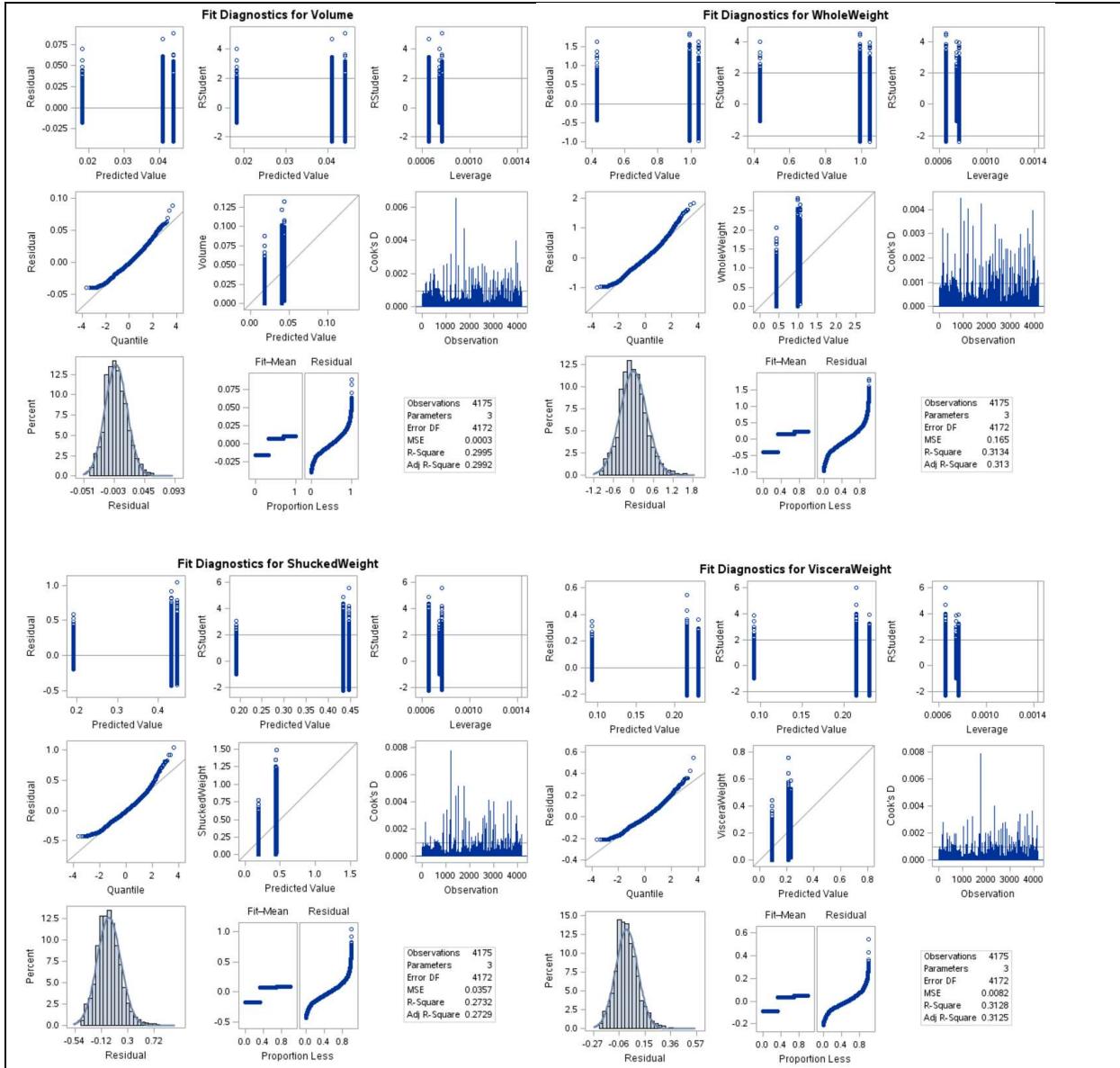
```
Model Rings Length Diameter Height Volume WholeWeight ShuckedWeight VisceraWeight
```

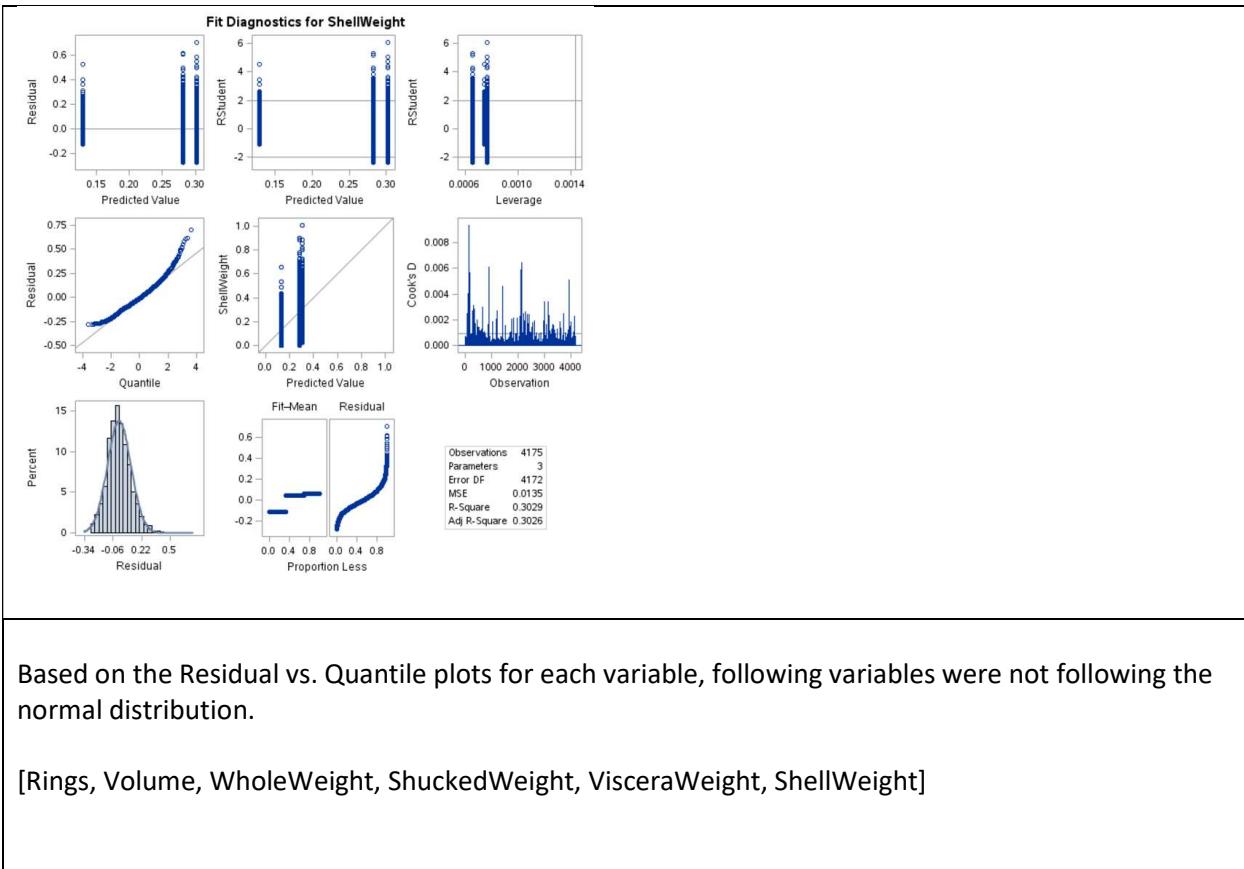
```
ShellWeight = Sex;
```

```
Manova H=_All_ / PrintE PrintH Canonical;
```

```
Run;
```

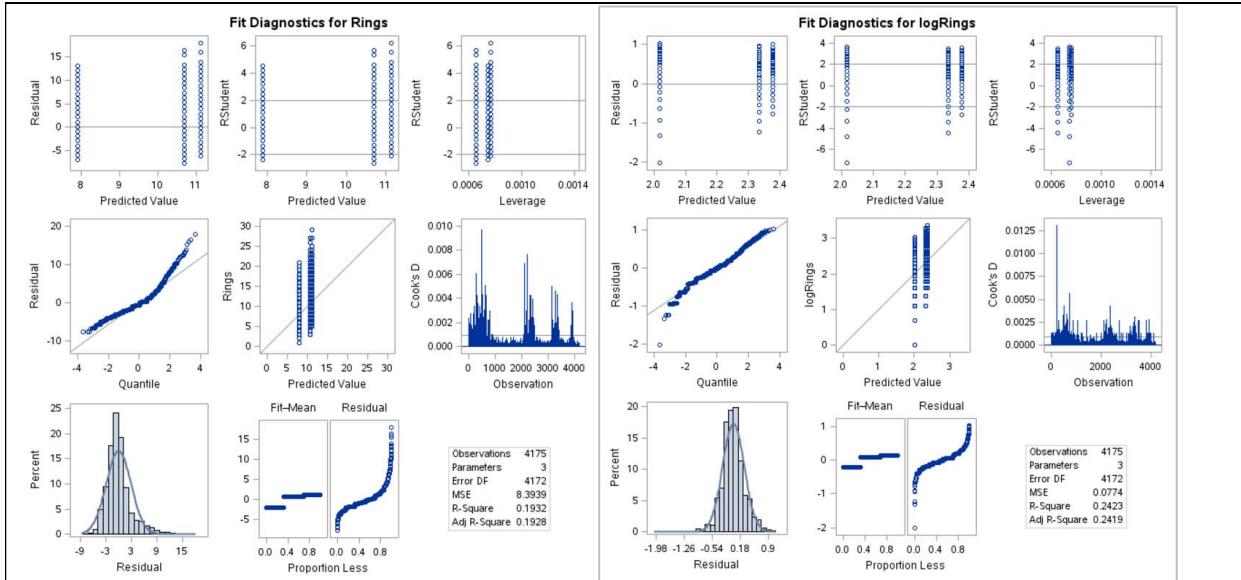




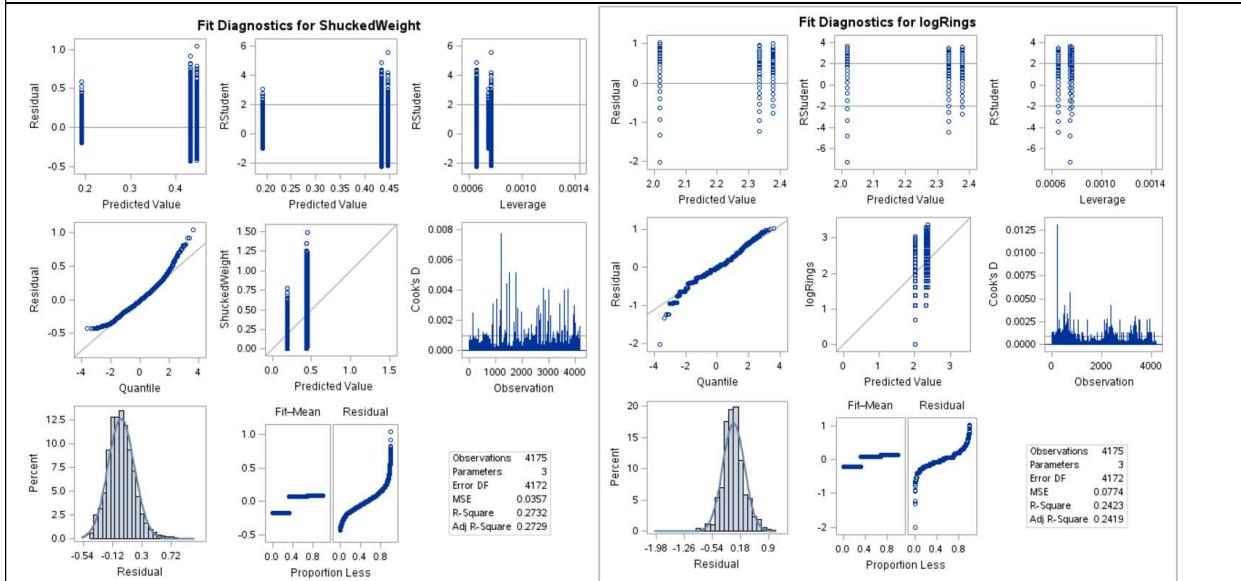


### Log Transformations

```
data abalone2;
set abalone;
logRings = log(Rings);
logVolume = log(Volume);
logWholeWeight = log(WholeWeight);
logShuckedWeight = log(ShuckedWeight);
logVisceraWeight = log(VisceraWeight);
logShellWeight = log(ShellWeight);
run;
```



'Rings' variable showed an improvement in error distribution after the log transformation.



Rest variables showed worse error distributions after the log transformation as shown above as an example. The scale of 'Residual vs. Quantile' plot is changed and the error is increased.

Other Transformations: Squared transformation and mean centering transformation were tried but there was no improvements in error distributions.

```
*squared transformation;
data abaloneSquare;
set abalone;
SRings = Rings*Rings;
SLength = Length*Length;
SDiameter = Diameter*Diameter;
SHeight = Height*Height;
```

```
SVolume = Volume*Volume;
SWholeWeight = WholeWeight*WholeWeight;
SShuckedWeight = ShuckedWeight*ShuckedWeight;
SVisceraWeight = VisceraWeight*VisceraWeight;
SShellWeight = ShellWeight*ShellWeight;
run;
proc print data = abaloneSqaure; run;

*Mean Centering Transformation;
data abaloneMC;
set abalone;
*MRings = Rings*Rings;
*MLength = Length*Length;
*MDiameter = Diameter*Diameter;
*MHeight = Height*Height;
*MVolume = Volume*Volume;
MWholeWeight = (WholeWeight-0.8285)**2;
MShuckedWeight = (ShuckedWeight-0.3591)**2;
MVisceraWeight = (VisceraWeight-0.1085)**2;
MShellWeight = (ShellWeight-0.2388)**2;
run;

*using Mean Centering;
proc means data = abalone n mean max min range std fw=8;
var _numeric_ ;
output out = meansout mean = mean std = std;
title 'Summary Stat';
run;

proc glm data=abaloneMC plots=all;
class sex;
model MWholeWeight      MShuckedWeight MVisceraWeight MShellWeight      = Sex;
output out=resides r=ral rfe rmg rca rna;
run;

* using squared transformatin;
proc glm data=abaloneSquare plots=all;
class sex;
model SRings SLength SDiameter SVolume SHeight SWholeWeight
SShuckedWeight SVisceraWeight SShellWeight      = Sex;
output out=resides r=ral rfe rmg rca rna;
run;
```

### Appendix C MANOVA and ANOVA output

```

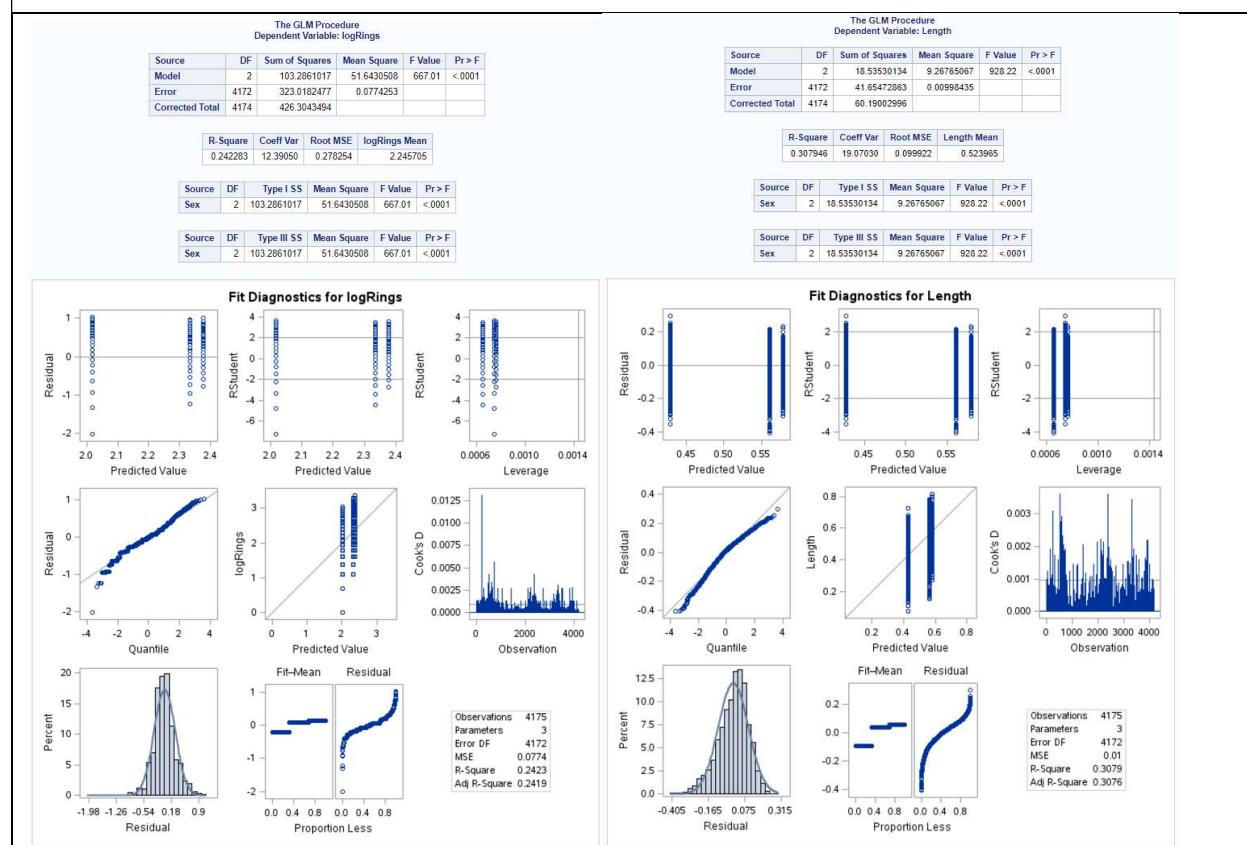
Proc GLM Data=abalone2 plots = all;
Class Sex;
Model logRings Length Diameter Height Volume WholeWeight ShuckedWeight VisceraWeight
ShellWeight = Sex;
Manova H=_All_ / PrintE PrintH Canonical;
Run;

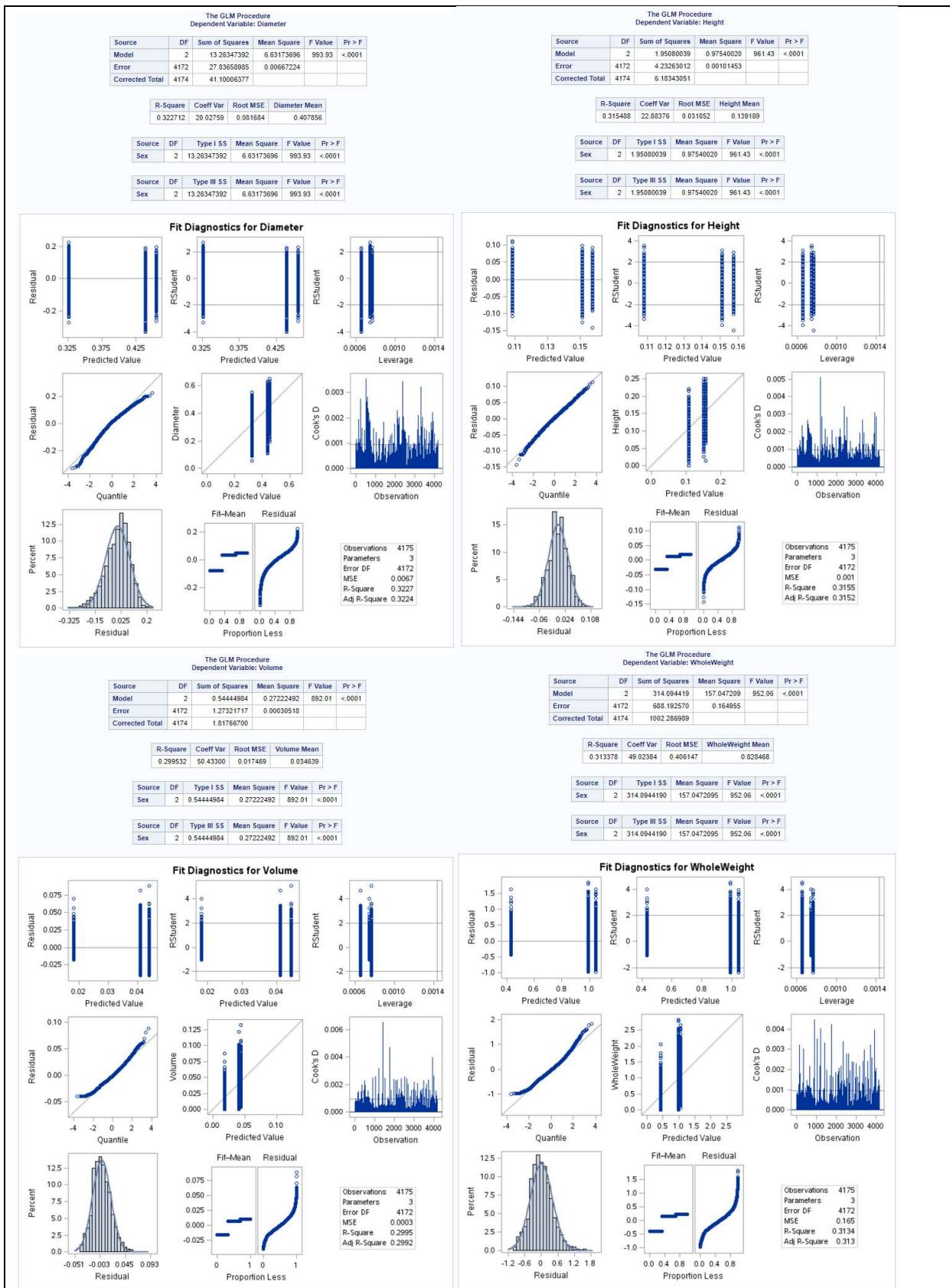
```

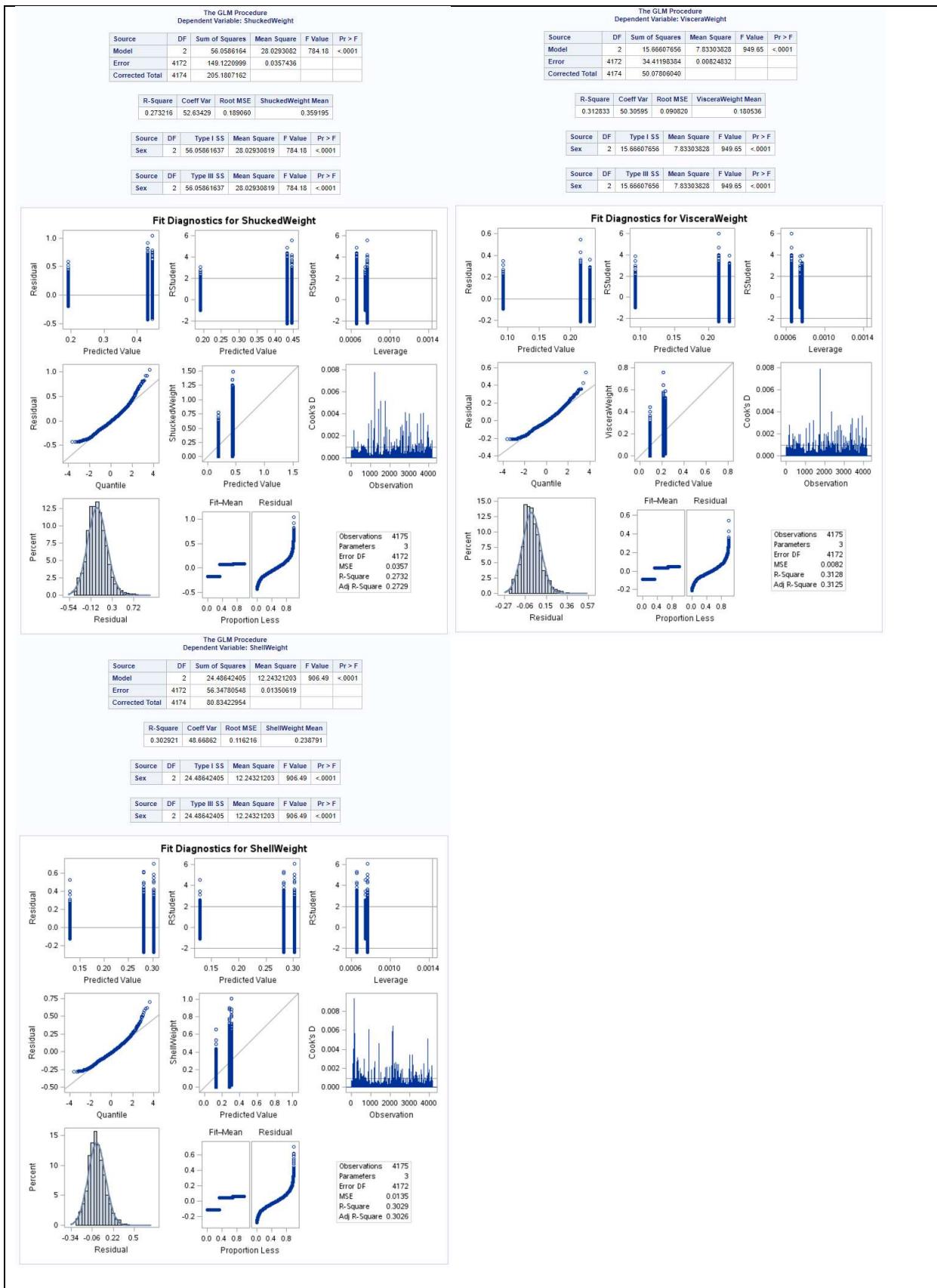
ANOVA F-test shows variables are significantly the sex levels for all variables. P-values from all F-tests are less than 0.0001.

Weight related variables shows slightly off-normal error distributions. At the same time the residual variance for each Sex level seem slightly different which is a violation of the ANOVA assumption.

... another assumption violation ... tried many transformations above but there was no effective method to fix the unequal variance and residual normality.









## Appending D – Full Code including all trials

```

/* MSDS6372 Project 2 Abalone Problem */
/* James Park 08JULY2017 */

/* Loading CSV */
FILENAME REFFILE '/home/jmpark1/Abalone Clean Data.csv';

PROC IMPORT DATAFILE=REFFILE
  DBMS=csv
  OUT=abalone;
  GETNAMES=yes;
RUN;
PROC CONTENTS DATA=abalone; RUN;
PROC PRINT DATA = abalone;run;

/* Data Exploring */

data abaloneM;
set abalone;
if Sex in ('M');
proc print data = abaloneM; run;

data abaloneF;
set abalone;
if Sex in ('F');
proc print data = abaloneF; run;

data abaloneI;
set abalone;
if Sex in ('I');
proc print data = abaloneI; run;

*log variables;
data abalone2;
set abalone;
logRings = log(Rings);
logLength = log(Length);
logDiameter = log(Diameter);
logHeight = log(Height);
logVolume = log(Volume);
logWholeWeight = log(WholeWeight);
logShuckedWeight = log(ShuckedWeight);
logVisceraWeight = log(VisceraWeight);
logShellWeight = log(ShellWeight);
Volume = Length*Diameter*Height;
run;
proc print data = abalone2; run;

*squared transformation;
data abaloneSquare;
set abalone;
SRings = Rings*Rings;
SLength = Length*Length;
SDiameter = Diameter*Diameter;
SHeight = Height*Height;
SVolume = Volume*Volume;
SWholeWeight = WholeWeight*WholeWeight;
SShuckedWeight = ShuckedWeight*ShuckedWeight;
SVisceraWeight = VisceraWeight*VisceraWeight;
SShellWeight = ShellWeight*ShellWeight;
run;
proc print data = abaloneSquare; run;

*Mean Centering Transformation;
data abaloneMC;
set abalone;
*MRings = Rings*Rings;
*MLength = Length*Length;
*MDiameter = Diameter*Diameter;
*MHeight = Height*Height;
*MVolume = Volume*Volume;
MWholeWeight = (WholeWeight-0.8285)**2;
MShuckedWeight = (ShuckedWeight-0.3591)**2;

```

```

MVisceraWeight = (VisceraWeight-0.1085)**2;
MShellWeight = (ShellWeight-0.2388)**2;
run;

*using Mean Centering;
proc glm data=abaloneMC plots=all;
class sex;
model MWholeWeight   MShuckedWeight      MVisceraWeight      MShellWeight      = Sex;
output out=resides r=rdf rfe rmg rca rna;
run;

* using squared transformatin;
proc glm data=abaloneSquare plots=all;
class sex;
model SRings SLength  SDiameter SVolume SHeight      SWholeWeight      SShuckedWeight      SVisceraWeight      SShellWeight
= Sex;
output out=resides r=rdf rfe rmg rca rna;
run;

proc means data = abalone n mean max min range std fw=8;
var _numeric_;
output out = meansout mean = mean std = std;
title 'Summary Stat';
run;

data abaloneFemaleMale;
set abalone2;
keep Sex WholeWeight Volume;
if Sex in ('M', 'F'); run;
proc print data = abaloneFemaleMale; run;

proc ttest data = abaloneFemaleMale ALPHA = .01;
CLASS sex;
var Volume WholeWeight;
run;

proc means data = abalone n mean max min range std fw=8;
var _numeric_;
by sex;
output out = meansout mean = mean std = std;
title 'Summary Stat by Sex';
run;

proc corr data = abalone;
run;

/* Scatter Plot */
proc sgscatter data=abalone ;
matrix Rings Length  Diameter Volume Height WholeWeight      ShuckedWeight      VisceraWeight      ShellWeight
/diagonal=(histogram);
run;

* Frequency distribution check;
proc freq data = abalone; run;

proc univariate data = abalone plots;
by sex;
Run;

proc univariate data = abalone alpha = 0.05 mu0=0 normal plot;
class sex;
var Rings Length  Diameter Volume Height WholeWeight      ShuckedWeight      VisceraWeight      ShellWeight;
histogram /normal
;inset mean std median min max;
run;quit;

*Initial Check of normality and covariance matrix;

proc glm data=abalone2 plots=all;
class sex;
model logRings Length  Diameter Volume Height WholeWeight      ShuckedWeight      VisceraWeight      ShellWeight = Sex;
output out=resides r=rdf rfe rmg rca rna;
run;

```

```

run;

proc glm data=abalone2 plots=all;
class sex;
model logRings logLength      logDiameter logHeight logVolume   logWholeWeight    logShuckedWeight   logVisceraWeight
      logShellWeight = Sex;
output out=resides r=r1 rfe rmg rca rna;
run;

* Check for equal variance and covariance matrix;
proc discrimin data=abalone2 pool=test;
class sex;
var logRings Length      Diameter Volume Height WholeWeight      ShuckedWeight      VisceraWeight      ShellWeight;
run;

* Rerun Check for equal variance and covariance matrix with fewer variables;
proc discrimin data=abalone pool=test;
class sex;
var Volume WholeWeight;
run;

/* Scatter Plot with log data*/
proc sgscatter data=abalone2 ;
matrix logRings logLength      logDiameter logHeight   logWholeWeight    logShuckedWeight   logVisceraWeight   logShellWeight
 / diagonal=(histogram);
run;

*PCA first attempt with all data;
ods rtf ; *places results in a word file;
ods graphics on;
proc princomp plots=all data=abalone cov out=pca;
  var Rings Length      Diameter Height      WholeWeight      ShuckedWeight      VisceraWeight      ShellWeight;
run;
proc print data=pca; run;
proc gplot data=pca;
  plot prin2*prin1=location;
run;
ods graphics off;
ods rtf close; *completes word file;

*PCA second attempt with only male data;
ods rtf ; *places results in a word file;
ods graphics on;
proc princomp plots=all data=abaloneM cov out=pcaM;
  var Rings Length      Diameter Height      WholeWeight      ShuckedWeight      VisceraWeight      ShellWeight;
run;
proc print data=pcaM; run;
proc gplot data=pcaM;
  plot prin2M*prin1M=locationM;
run;
ods graphics off;
ods rtf close; *completes word file;

*PCA third attempt with only female data;
ods rtf ; *places results in a word file;
ods graphics on;
proc princomp plots=all data=abaloneF cov out=pcaF;
  var Rings Length      Diameter Height      WholeWeight      ShuckedWeight      VisceraWeight      ShellWeight;
run;

data pcaF2;
set pcaF;
prin1F=prin1;
prin2F=prin2;
keep prin1F prin2F sex;
run;
proc print data=pcaF2; run;

data pcaM2;
set pcaM;
prin1M=prin1;
prin2M=prin2;
keep prin1M prin2M sex;
run;

```

```

proc print data=pcaM2; run;

proc gplot data=pcaF2;
  plot prin2F*prin1F=locationF;
  run;
ods graphics off;
ods rtf close; *completes word file;

*pca combined;
data pcaCombined;
set pcaF2 pcaM2;
run;
proc print pcaCombined; run;

data pcaCombined2;
set pcaCombined;
if prin1M='.' then prin1C = prin1F;
else if prin1F='.' then prin1C= prin1M;
if prin2M='.' then prin2C = prin2F;
else if prin2F='.' then prin2C= prin2M;
run;
proc print data=pcaCombined2; run;

*if we wanted to build a model to predict age;
proc glm data = abalone plots = all;
class sex;
model Rings = sex LengthDiameter Height WholeWeight ShuckedWeight VisceraWeight ShellWeight / solution;
run;quit;

*LDA STARTING CODE;
*First run a one-way ANOVA for each response;
Proc GLM Data=abalone plot=diagnostics;
Class Sex;
*Model height Length = Sex;
Model Rings Length Diameter Height WholeWeight ShuckedWeight VisceraWeight ShellWeight = Sex;
lsmeans sex / pdiff tdiff adjust=bon;
Output Out=Errs R=Exact Eant; *to look at residuals;
Means Sex / CLM T Alpha=0.0167; *code if we assume we are ok on assumptions. Adjusted Bonferonni;
Means Sex / CLDIFF T Alpha=0.0167; *code if we assume we are ok on assumptions. Adjusted Bonferonni;
*Contrast 'Sex contrast' sex -0.5 0 0.5; *order would be F, I, M;
*Estimate 'Sex estimate' sex -0.5 0 0.5; *order would be F, I, M;
Manova H=_All_ / PrintE PrintH Canonical; ***Must have this to actually function as a multiple response test;
Contrast 'Sex contrast F M' sex -0.5 0 0.5; *order would be F, I, M;
Contrast 'Sex contrast F I' sex -0.5 0.5 0; *order would be F, I, M;
ODS Output Cancoef=CanonicalCoefficients;
Run;
Quit;

*Second Run with log data one-way ANOVA for each response;
Proc GLM Data=abalone2 plot=diagnostics;
Class Sex;
*Model height Length = Sex;
Model logRings logLength logDiameter logHeight logWholeWeight logShuckedWeight logVisceraWeight logShellWeight = Sex;
lsmeans sex / pdiff tdiff adjust=bon;
Output Out=Errs R=Exact Eant; *to look at residuals;
Means Sex / CLM T Alpha=0.0167; *code if we assume we are ok on assumptions. Adjusted Bonferonni;
Means Sex / CLDIFF T Alpha=0.0167; *code if we assume we are ok on assumptions. Adjusted Bonferonni;
*Contrast 'Sex contrast' sex -0.5 0 0.5; *order would be F, I, M;
*Estimate 'Sex estimate' sex -0.5 0 0.5; *order would be F, I, M;
Manova H=_All_ / PrintE PrintH Canonical; ***Must have this to actually function as a multiple response test;
Contrast 'Sex contrast F M' sex -0.5 0 0.5; *order would be F, I, M;
Contrast 'Sex contrast F I' sex -0.5 0.5 0; *order would be F, I, M;
ODS Output Cancoef=CanonicalCoefficients;
Run;
Quit;

*using only principle components;
Proc GLM Data=pcaCombined2 plot=diagnostics;
Class Sex;
Model prin1C prin2C= Sex;
lsmeans sex / pdiff tdiff adjust=bon;
Output Out=Errs R=Exact Eant; *to look at residuals;
*Means Sex / CLM T Alpha=0.0125; *code if we assume we are ok on assumptions. Adjusted Bonferonni;
*Means Sex / CLDIFF T Alpha=0.0125; *code if we assume we are ok on assumptions. Adjusted Bonferonni;
*Contrast 'Sex contrast' sex -1 1; *order would be F, M;

```

```

*Estimate 'Sex estimate' sex -1 1; *order would be F, M;
Manova H=_All_ / PrintE PrintH Canonical; ***Must have this to actually function as a multiple response test;
Contrast 'Sex contrast F M' sex -1 1; *order would be F M;
*Contrast 'Sex contrast F I' sex -1 1; *order would be F M;
ODS Output Cancoef=CanonicalCoefficients;
Run;
Quit;

*using only Sex M and F data with wholeweight and volume;
Proc GLM Data=abaloneFemaleMale plot=diagnostics;
Class Sex;
Model WholeWeight Volume = Sex;
Ismeans sex / pdiff tdiff adjust=bon;
Output Out=Errs R=Exact Eant; *to look at residuals;
*Means Sex / CLM T Alpha=0.0125; *code if we assume we are ok on assumptions. Adjusted Bonferroni;
*Means Sex / CLDIFF T Alpha=0.0125; *code if we assume we are ok on assumptions. Adjusted Bonferroni;
*Contrast 'Sex contrast' sex -1 1; *order would be F, M;
*Estimate 'Sex estimate' sex -1 1; *order would be F, M;
Manova H=_All_ / PrintE PrintH Canonical; ***Must have this to actually function as a multiple response test;
Contrast 'Sex contrast F M' sex -1 1; *order would be F M;
*Contrast 'Sex contrast F I' sex -1 1; *order would be F M;
ODS Output Cancoef=CanonicalCoefficients;
Run;
Quit;

proc glm data=abalone; * order=data;
class sex;
model Rings Length Diameter Height WholeWeight ShuckedWeight VisceraWeight ShellWeight = sex;
Ismeans sex / pdiff tdiff adjust=bon;
means sex / lsd bon tukey scheffe;
*contrast 'lin' sex -1 1 0;*data order would be F, I, M;
estimate 'I M Comparison' Sex -1 0 1;
run;

PROC DISCRIM data = abalone2 CROSSVALIDATE OUTSTAT=DIS_FUNC POOL=test wcov pcov;
CLASS sex;
VAR logRings Length Diameter Height Volume WholeWeight ShuckedWeight VisceraWeight ShellWeight;
PRIORS EQUAL;
run;
quit;

data abaloneMFonly;
set abalone2;
if Sex in ('M' 'F'); run;
proc print data = abaloneMFonly; run;

PROC DISCRIM data = abaloneMFonly CROSSVALIDATE OUTSTAT=DIS_FUNC POOL=test wcov pcov;
CLASS sex;
VAR logRings Length Diameter Height Volume WholeWeight ShuckedWeight VisceraWeight ShellWeight;
PRIORS EQUAL;
run;
quit;

PROC DISCRIM data = abalone2 CROSSVALIDATE OUTSTAT=DIS_FUNC POOL=test;
CLASS sex;
VAR logRings logLength logDiameter logHeight logVolume logWholeWeight logShuckedWeight logVisceraWeight logShellWeight;
PRIORS EQUAL;
run;
quit;

PROC DISCRIM data = pcaCombined2 CROSSVALIDATE OUTSTAT=DIS_FUNC POOL=YES;
CLASS sex;
VAR Prin1C Prin2C;
PRIORS EQUAL;
run;
quit;
proc print data=dis_func; run;

*automatic variable selection using stepdisc for LDA variable selection;
proc stepdisc data=abaloneMFonly bsscp tsscp;
class Sex;
var LogRings Length Diameter Height Volume WholeWeight ShuckedWeight VisceraWeight ShellWeight;
run;

```

```
*rerun LDA using auto-selected variables;  
PROC DISCRIM data = abalone CROSSVALIDATE OUTSTAT=DIS_FUNC POOL=test wcov pcov;  
CLASS sex;  
VAR Diameter Height ShuckedWeight VisceraWeight;  
PRIORS EQUAL;  
run;  
quit;
```