# Appendix 8.C.  Empirical Regression Models and Goodness-of-Fit Diagnostics

Michael F. Coveney, SJRWMD

St. Johns River Water Management District
Palatka, Florida
2011

Appendix 8.C.

# Contents

# Introduction

The St. Johns River Water Management District's Water Supply Impact Study evaluated potential effects of water supply withdrawals from the St. Johns River (SJR) on biological and water resources. The Plankton Working Group was charged with the identification and quantification of possible environmental effects of water withdrawals on plankton communities and Total Maximum Daily Loads. Most of the potential effects that we investigated were causes or consequences of enhanced growth of phytoplankton. Consequently, phytoplankton blooms, primarily cyanobacteria in fresh water and dinoflagellates in brackish water, were a primary focus of our work. One important step in this work was to develop hydroecological regression models that link changes in hydrology to changes in plankton communities.

We aggregated the effects of algal blooms into four algal bloom metrics: marine algal blooms (dinoflagellate biovolume), increase in nitrogen (N) load due to $N_2$-fixation, magnitude of freshwater algal blooms (maximum chlorophyll-a concentration), and duration of freshwater algal blooms. We used water age (a measure of residence time) as the primary hydrologic variable and constructed empirical regression models to relate algal bloom metrics to water age variables. Water age was calculated in the Environmental Fluids Dynamic Code (EFDC) hydrodynamic model and was the average time that water resided in the model domain before reaching a specific site (model grid cell).

We constructed a total of eight regression models using both multiple linear and multiple logistic regression techniques (Table 1). These regression models covered the four algal bloom metrics in four of the five river segments that we assessed (Lake Poinsett in segment 8 was assessed separately; see Appendix 8.B.). Seven of the regression models were linear regressions with adjusted $R^2$ values that varied from 0.80 to 0.97. In each case, both the overall model and the individual independent variables were significant ($p < 0.05$). The remaining model, freshwater bloom duration in segments 3 and 4, was a logistic regression with two independent variables. In this case, we chose logistic regression because the best-fit linear regression model achieved only a modest $R^2$ value. Regression models were constructed using either stepwise (forward and backward selection) or $R^2$ selection of independent variables for linear regression and by highest likelihood score selection for logistic regression. SAS$^{©}$ was used in these analyses.

We assessed potential effects of water withdrawals by using the empirical regression models to predict algal bloom measures from water age variables. We used water age output from the EFDC hydrodynamic model for the baseline scenario and each withdrawal scenario and compared the predicted bloom measures. By focusing on the difference between predicted results for the baseline scenario and for withdrawal scenarios, we minimized effects of model bias since the bias would be present in both sets of results.

# Dependent and Independent Variables

Dependent variables for regression models were the quantities measured for each of the four algal bloom metrics (Table 1). These variables were tabulated annually and represented either the maximum (biovolume, chl-a, bloom duration) or the total (mass $N_2$-fixation) for each year. No variables were transformed except for maximum annual dinoflagellate biovolume, where logarithmic transformation appeared to be necessary (Chapter 8. Plankton).

Table 1. Regression type, source location(s) for data, and dependent variable names used in each regression model. See Chapter 8. Plankton for sampling sites corresponding to these locations.

| Bloom Measure | Segment | Type | Data Source Location | Dependent Variable |
|---|---|---|---|---|
| ln(Dinoflag Biovolume) | 2 | Linear | Mandarin Pt Doctors Lake | LNMaxDino |
| $N_2$-Fixation | 3,4 | Linear | Lake George | Est_N2_Fix |
| FW Max Chl-a | 2 | Linear | Doctors Lake | Max_Chl |
| FW Max Chl-a | 3,4 | Linear | Racy Pt Lake George | Max_Chl |
| FW Max Chl-a | 6 | Linear | Lake Monroe Lake Jesup Lake Harney | Max_Chla |
| FW Bloom Duration | 2 | Linear | Doctors Lake | Longest_d |
| FW Bloom Duration Event | 3,4 | Logistic | Racy Pt Palatka Lake George | Bloom_50d |
| FW Bloom Duration | 6 | Linear | Lake Monroe Lake Jesup | Longest_dur_Bloom |

Independent variables for regression models were minimum, mean, and maximum daily water age and the inverse of each water age calculated for seven periods, for a total of 42 possible variables (Table 2). Time periods were five quarters (A – E, starting with the last quarter of the previous calendar year) and two growth-season periods (Apr – Aug and Apr – Oct) (Table 2).

We included this wide range of water age variables as potential prediction variables in regression models because relationships between algal bloom patterns and hydrology were both positive and negative and showed complex seasonality (Chapter 8. Plankton).

The period for observed data used to build regression models was determined by the eleven-year simulation period for the hydrodynamic model: January 1995 through December 2005. Because model initiation ("spin-up") caused low water age values during the first three months, we did not use data from Jan – March 1995.

Table 2. Naming convention for water age (independent) variables. Variable names were combinations of two forms (value, inverse value), three statistics (mean, minimum, maximum), and seven periods, for a total of 42 variables. Examples of water age variables using this convention were MeanAgeD (mean daily water age for period D), Min_Age_Apr_Oct (minimum daily water age for the period April – October), and invMaxAgeA (the inverse of the maximum daily water age for period A). Variable names varied slightly between analysts; e.g., both "mean" and "avg" were used for mean, and position of underscore characters differed.

| Forms | Statistics | Periods |
|---|---|---|
| Water age value | Mean | April – October |
| Inverse of the water age value | Minimum | April – August |
| | Maximum | A (Oct-Dec of previous year) |
| | | B (Jan-Mar) |
| | | C (Apr-Jun) |
| | | D (Jul-Sep) |
| | | E (Oct-Dec) |

Appendix 8.C.

# SAS® Calculation of Regression Models

## A. Dinoflagellate Biovolume, Segment 2

```
                        The REG Procedure
                           Model: MODEL1
                  Dependent Variable: LNMaxDino LNMaxDino

                  Number of Observations Read          20
                  Number of Observations Used          20
```

### Analysis of Variance

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 5 | 36.50278 | 7.30056 | 17.36 | <.0001 |
| Error | 14 | 5.88607 | 0.42043 | | |
| Corrected Total | 19 | 42.38885 | | | |

| | | | | |
|---|---|---|---|---|
| Root MSE | 0.64841 | R-Square | 0.8611 | |
| Dependent Mean | 5.39282 | Adj R-Sq | 0.8115 | |
| Coeff Var | 12.02355 | | | |

### Parameter Estimates

| Variable | Label | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
|---|---|---|---|---|---|---|
| Intercept | Intercept | 1 | 7.31406 | 2.06093 | 3.55 | 0.0032 |
| Inv_Max_Age_E | Inv_Max_Age_E | 1 | 390.31362 | 124.20677 | 3.14 | 0.0072 |
| Inv_Min_Age_D | Inv_Min_Age_D | 1 | -540.09679 | 147.07539 | -3.67 | 0.0025 |
| Max_Apr_Oct | Max Apr_Oct | 1 | 0.15035 | 0.02471 | 6.08 | <.0001 |
| Max_Age_C | Max_Age_C | 1 | -0.12813 | 0.02268 | -5.65 | <.0001 |
| Min_Age_D | Min_Age_D | 1 | -0.04082 | 0.01245 | -3.28 | 0.0055 |

### Parameter Estimates

| Variable | Label | DF | Standardized Estimate |
|---|---|---|---|
| Intercept | Intercept | 1 | 0 |
| Inv_Max_Age_E | Inv_Max_Age_E | 1 | 0.80224 |
| Inv_Min_Age_D | Inv_Min_Age_D | 1 | -1.47481 |
| Max_Apr_Oct | Max Apr_Oct | 1 | 6.30995 |
| Max_Age_C | Max_Age_C | 1 | -5.05773 |
| Min_Age_D | Min_Age_D | 1 | -1.49375 |

Appendix 8.C.

## *B. N₂-Fixation, Segments 3 & 4*

```
                        The REG Procedure
                         Model: MODEL1
               Dependent Variable: Est_N2_Fix Est_N2_Fix

         Number of Observations Read                  22
         Number of Observations Used                  10
         Number of Observations with Missing Values   12


                      Analysis of Variance

                              Sum of        Mean
    Source             DF     Squares      Square   F Value   Pr > F

    Model               3     3689415     1229805     76.29   <.0001
    Error               6       96722       16120
    Corrected Total     9     3786137


              Root MSE          126.96616   R-Square    0.9745
              Dependent Mean    806.26123   Adj R-Sq    0.9617
              Coeff Var          15.74752


                      Parameter Estimates

                              Parameter    Standard
    Variable     Label      DF   Estimate      Error   t Value   Pr > |t|

    Intercept    Intercept   1  444.94101  213.20432      2.09   0.0819
    MeanAgeD     MeanAgeD    1    4.67005    1.15732      4.04   0.0068
    invMeanAgeB  invMeanAgeB 1    -259099      21012    -12.33   <.0001
    invMinAgeC   invMinAgeC  1     242239      16316     14.85   <.0001

                      Parameter Estimates

                                          Standardized
            Variable     Label      DF      Estimate

            Intercept    Intercept   1           0
            MeanAgeD     MeanAgeD    1     0.29719
            invMeanAgeB  invMeanAgeB 1    -2.16507
            invMinAgeC   invMinAgeC  1     2.63640
```

Appendix 8.C.

## C. Freshwater Maximum Chl-a, Segment 2

```
                         The REG Procedure
                          Model: MODEL1
                  Dependent Variable: Max_Chl Max_Chl

        Number of Observations Read                    11
        Number of Observations Used                    10
        Number of Observations with Missing Values      1


                        Analysis of Variance

                              Sum of          Mean
        Source            DF   Squares        Square    F Value   Pr > F

        Model              4     10217    2554.17668      23.64   0.0019
        Error              5  540.17729     108.03546
        Corrected Total    9     10757


            Root MSE            10.39401   R-Square     0.9498
            Dependent Mean      82.34000   Adj R-Sq     0.9096
            Coeff Var           12.62328


                        Parameter Estimates

                               Parameter     Standard
        Variable       Label       DF    Estimate      Error   t Value  Pr > |t|

        Intercept      Intercept    1   100.21543   39.22541      2.55    0.0510
        Ave_Apr_Aug    Ave_Apr_Aug  1     3.62191    0.43635      8.30    0.0004
        Inv_Ave_Age_A  Inv_Ave_Age_A 1      57722      20280      2.85    0.0360
        Inv_Max_Age_A  Inv_Max_Age_A 1     -65902      20501     -3.21    0.0236
        Max_Apr_Oct    Max_Apr_Oct  1    -3.16151    0.40773     -7.75    0.0006

                        Parameter Estimates

                                            Standardized
            Variable       Label       DF      Estimate

            Intercept      Intercept    1             0
            Ave_Apr_Aug    Ave_Apr_Aug  1       5.89187
            Inv_Ave_Age_A  Inv_Ave_Age_A 1       2.44664
            Inv_Max_Age_A  Inv_Max_Age_A 1      -2.73423
            Max_Apr_Oct    Max_Apr_Oct  1      -5.57109
```

Appendix 8.C.

*D. Freshwater Maximum Chl-a, Segments 3 & 4*

```
                          The REG Procedure
                          Model: MODEL1
                    Dependent Variable: Max_Chl Max_Chl

           Number of Observations Read                    22
           Number of Observations Used                    20
           Number of Observations with Missing Values      2


                          Analysis of Variance

                               Sum of          Mean
   Source              DF      Squares        Square   F Value    Pr > F

   Model                7        18074    2582.06611     21.60    <.0001
   Error               12   1434.20270     119.51689
   Corrected Total     19        19509


              Root MSE              10.93238    R-Square     0.9265
              Dependent Mean        94.16500    Adj R-Sq     0.8836
              Coeff Var             11.60981


                          Parameter Estimates

                                           Parameter      Standard
Variable             Label              DF   Estimate        Error   t Value

Intercept            Intercept           1   24.63329     71.77171      0.34
MeanAgeD             MeanAgeD            1    1.00042      0.28955      3.46
MinAgeD              MinAgeD             1   -2.36030      0.29282     -8.06
MaxAgeE              MaxAgeE             1    0.86558      0.14935      5.80
invMean_Age_Apr_Oct  invMean_Age_Apr_Oct 1 6121.89276   2507.95714      2.44
invMeanAgeE          invMeanAgeE         1 3916.79838    965.69610      4.06
invMinAgeD           invMinAgeD          1 -7892.28591   1238.47933     -6.37
invMaxAgeA           invMaxAgeA          1 3676.65842   1016.17188      3.62

                          Parameter Estimates

                                                         Standardized
     Variable             Label              DF  Pr > |t|     Estimate

     Intercept            Intercept           1    0.7374            0
     MeanAgeD             MeanAgeD            1    0.0048      1.36713
     MinAgeD              MinAgeD             1    <.0001     -2.31831
     MaxAgeE              MaxAgeE             1    <.0001      1.29710
     invMean_Age_Apr_Oct  invMean_Age_Apr_Oct 1    0.0311      0.63618
     invMeanAgeE          invMeanAgeE         1    0.0016      0.75460
     invMinAgeD           invMinAgeD          1    <.0001     -1.76370
     invMaxAgeA           invMaxAgeA          1    0.0035      0.54011
```

Appendix 8.C.

## E. Freshwater Maximum Chl-a, Segment 6

```
                          The REG Procedure
                            Model: MODEL1
                  Dependent Variable: Max_Chla Max_Chla

        Number of Observations Read                    33
        Number of Observations Used                    30
        Number of Observations with Missing Values      3


                        Analysis of Variance

                              Sum of          Mean
        Source          DF    Squares       Square    F Value   Pr > F

        Model            4     130729        32682      30.41   <.0001
        Error           25      26871   1074.85196
        Corrected Total 29     157600


             Root MSE            32.78494   R-Square    0.8295
             Dependent Mean      93.98230   Adj R-Sq    0.8022
             Coeff Var           34.88416


                        Parameter Estimates

                           Parameter     Standard
        Variable    Label         DF     Estimate        Error   t Value   Pr > |t|

        Intercept    Intercept     1      4.04120     10.67110      0.38    0.7081
        Min_Apr_Oct  Min_Apr_Oct   1     -4.81072      1.38030     -3.49    0.0018
        Max_Age_B    Max_Age_B     1      0.73032      0.17495      4.17    0.0003
        Min_Age_D    Min_Age_D     1      2.15419      0.65338      3.30    0.0029
        Min_Age_E    Min_Age_E     1      3.43669      1.07992      3.18    0.0039

                        Parameter Estimates

                                            Standardized
              Variable      Label      DF      Estimate

              Intercept     Intercept    1            0
              Min_Apr_Oct   Min_Apr_Oct  1     -1.79980
              Max_Age_B     Max_Age_B    1      0.50293
              Min_Age_D     Min_Age_D    1      0.90130
              Min_Age_E     Min_Age_E    1      1.33638
```

38

Appendix 8.C.

*F. Freshwater Bloom Duration, Segment 2*

```
                        The REG Procedure
                          Model: MODEL1
                 Dependent Variable: Longest_d Longest_d

          Number of Observations Read                    11
          Number of Observations Used                    10
          Number of Observations with Missing Values      1


                      Analysis of Variance

                             Sum of          Mean
     Source            DF    Squares        Square    F Value    Pr > F

     Model              4      23431    5857.75456      63.87    0.0002
     Error              5  458.58175      91.71635
     Corrected Total    9      23890


              Root MSE             9.57687    R-Square    0.9808
              Dependent Mean     103.80000    Adj R-Sq    0.9654
              Coeff Var            9.22627


                       Parameter Estimates

                                          Parameter    Standard
   Variable          Label           DF    Estimate       Error    t Value

   Intercept         Intercept        1    538.44157    44.83292      12.01
   Ave_Age_A         Ave_Age_A        1     -1.90085     1.06921      -1.78
   Ave_Age_D         Ave_Age_D        1     -2.34085     0.19796     -11.83
   Inv_Min_Apr_Oct   Inv_Min_Apr_Oct  1       -31891  3700.23837      -8.62
   Max_Age_A         Max_Age_A        1      3.13881     0.85479       3.67

                       Parameter Estimates

                                                      Standardized
      Variable          Label           DF   Pr > |t|     Estimate

      Intercept         Intercept        1    <.0001            0
      Ave_Age_A         Ave_Age_A        1    0.1356     -1.39338
      Ave_Age_D         Ave_Age_D        1    <.0001     -2.57233
      Inv_Min_Apr_Oct   Inv_Min_Apr_Oct  1    0.0003     -1.36967
      Max_Age_A         Max_Age_A        1    0.0144      2.70134
```

39

Appendix 8.C.

## G. Freshwater Bloom Duration, Segments 3 & 4

The LOGISTIC Procedure

Model Information

```
Data Set                      WORK.TEMP
Response Variable             Bloom_50d            Bloom_50d
Number of Response Levels     2
Model                         binary logit
Optimization Technique        Fisher's scoring
```

```
Number of Observations Read          39
Number of Observations Used          33
```

Response Profile

| Ordered Value | Bloom_50d | Total Frequency |
|---|---|---|
| 1 | 0 | 14 |
| 2 | 1 | 19 |

Probability modeled is Bloom_50d=1.

NOTE: 6 observations were deleted due to missing values for the response or explanatory variables.

Model Convergence Status

Convergence criterion (GCONV=1E-8) satisfied.

Model Fit Statistics

| Criterion | Intercept Only | Intercept and Covariates |
|---|---|---|
| AIC | 46.987 | 33.008 |
| SC | 48.484 | 37.497 |
| -2 Log L | 44.987 | 27.008 |

R-Square    0.4201    Max-rescaled R-Square    0.5645

Testing Global Null Hypothesis: BETA=0

| Test | Chi-Square | DF | Pr > ChiSq |
|---|---|---|---|
| Likelihood Ratio | 17.9794 | 2 | 0.0001 |
| Score | 13.3032 | 2 | 0.0013 |
| Wald | 7.2636 | 2 | 0.0265 |

Appendix 8.C.

```
                     Analysis of Maximum Likelihood Estimates

                                  Standard       Wald
Parameter            DF   Estimate    Error   Chi-Square   Pr > ChiSq

Intercept             1   -26.4944   9.9179      7.1362       0.0076
Mean_Age_Apr_Aug      1     0.1331   0.0494      7.2606       0.0070
invMeanAgeC           1     1155.3    448.4      6.6378       0.0100


                           Odds Ratio Estimates

                           Point        95% Wald
            Effect        Estimate   Confidence Limits

         Mean_Age_Apr_Aug   1.142     1.037      1.258
         invMeanAgeC     >999.999  >999.999   >999.999


     Association of Predicted Probabilities and Observed Responses

         Percent Concordant    91.0   Somers' D    0.823
         Percent Discordant     8.6   Gamma        0.826
         Percent Tied           0.4   Tau-a        0.415
         Pairs                  266   c            0.912


             Partition for the Hosmer and Lemeshow Test

                         Bloom_50d = 1         Bloom_50d = 0
     Group     Total   Observed   Expected   Observed   Expected

        1        3        1         0.14         2        2.86
        2        3        0         0.44         3        2.56
        3        3        0         0.79         3        2.21
        4        3        0         1.13         3        1.87
        5        3        1         1.31         2        1.69
        6        3        3         1.65         0        1.35
        7        3        3         2.18         0        0.82
        8        3        2         2.58         1        0.42
        9        3        3         2.81         0        0.19
       10        3        3         2.96         0        0.04
       11        3        3         2.99         0        0.01


             Hosmer and Lemeshow Goodness-of-Fit Test

            Chi-Square      DF      Pr > ChiSq

              13.6939        9         0.1336
```

Appendix 8.C.

## *H. Freshwater Bloom Duration, Segment 6*

```
                        The REG Procedure
                         Model: MODEL1
         Dependent Variable: Longest_Dur_Bloom Longest_Dur_Bloom

         Number of Observations Read                      22
         Number of Observations Used                      20
         Number of Observations with Missing Values        2


                       Analysis of Variance

                              Sum of           Mean
    Source              DF    Squares         Square    F Value    Pr > F

    Model                4      90152          22538      48.97    <.0001
    Error               15  6903.94924     460.26328
    Corrected Total     19      97056


             Root MSE            21.45375    R-Square     0.9289
             Dependent Mean     104.00000    Adj R-Sq     0.9099
             Coeff Var           20.62860


                       Parameter Estimates

                                     Parameter      Standard
    Variable          Label      DF   Estimate         Error    t Value

    Intercept         Intercept   1   -44.37433      35.92549      -1.24
    Inv_Ave_Apr_Oct   Inv_Ave_Apr_Oct  1  -7370.26368  1804.32954  -4.08
    Inv_Max_Apr_Oct   Inv_Max_Apr_Oct  1      15172    3532.35554   4.30
    Ave_Age_A         Ave_Age_A   1     0.51719       0.19136       2.70
    Max_Age_D         Max_Age_D   1     0.95232       0.18157       5.25

                       Parameter Estimates

                                                    Standardized
        Variable          Label        DF  Pr > |t|    Estimate

        Intercept         Intercept     1    0.2358           0
        Inv_Ave_Apr_Oct   Inv_Ave_Apr_Oct  1   0.0010   -1.01649
        Inv_Max_Apr_Oct   Inv_Max_Apr_Oct  1   0.0006    1.02284
        Ave_Age_A         Ave_Age_A     1    0.0164     0.26288
        Max_Age_D         Max_Age_D     1    <.0001     0.70370
```

# SAS® Goodness-of-Fit Diagnostics for Regression Models

## I. Introduction and Key

We show a set of diagnostic plots (SAS Institute Inc. 2010, page 6302) below for each of the seven linear regression models. Goodness-of-fit metrics for the logistic regression model (freshwater bloom duration, segments 3 & 4) are included in model calculation (above).

    1.   Key - Fit Diagnostic Plots



A. Residuals vs predicted values

B. Externally studentized residuals vs predicted values

C. Externally studentized residuals vs leverage

D. Normal quantile-quantile (Q-Q) plot of residuals

E. Observed values vs predicted values

F. Cook's D vs observation number

G. Histogram of residuals

H. Residual-Fit (RF): quantile plots of the centered fit and the residuals

Plots A, B, and H help diagnose the adequacy of models. Better models lack patterns in the plots of residuals or studentized residuals versus predicted values (A, B), and the spread of residuals is less than the spread of the centered fit in the RF plot (H). Better models also lack patterns in the spread about the 1:1 reference line in the plot of observed values versus predicted values (E). The Q-Q plot (D) and residual histogram plot (G) can indicate problems with lack of normality and with heteroscedasticity. Plot B provides a test for outliers (defined as observations with studentized residuals greater than 2). Plots C and F identify observations with high leverage values (indicated by index lines). Further information is found in the SAS/STAT 9.22 User's Guide (SAS Institute Inc. 2010, page 6301-2).

2. <u>Key - Residuals by Regressors (Independent Variables)</u>



Plots of residuals vs each of the individual independent variables (I) also help to diagnose the adequacy of models; better models do not show patterns in these scatter plots. Further information is found in SAS/STAT 9.22 User's Guide (SAS Institute Inc. 2010, page 6302).
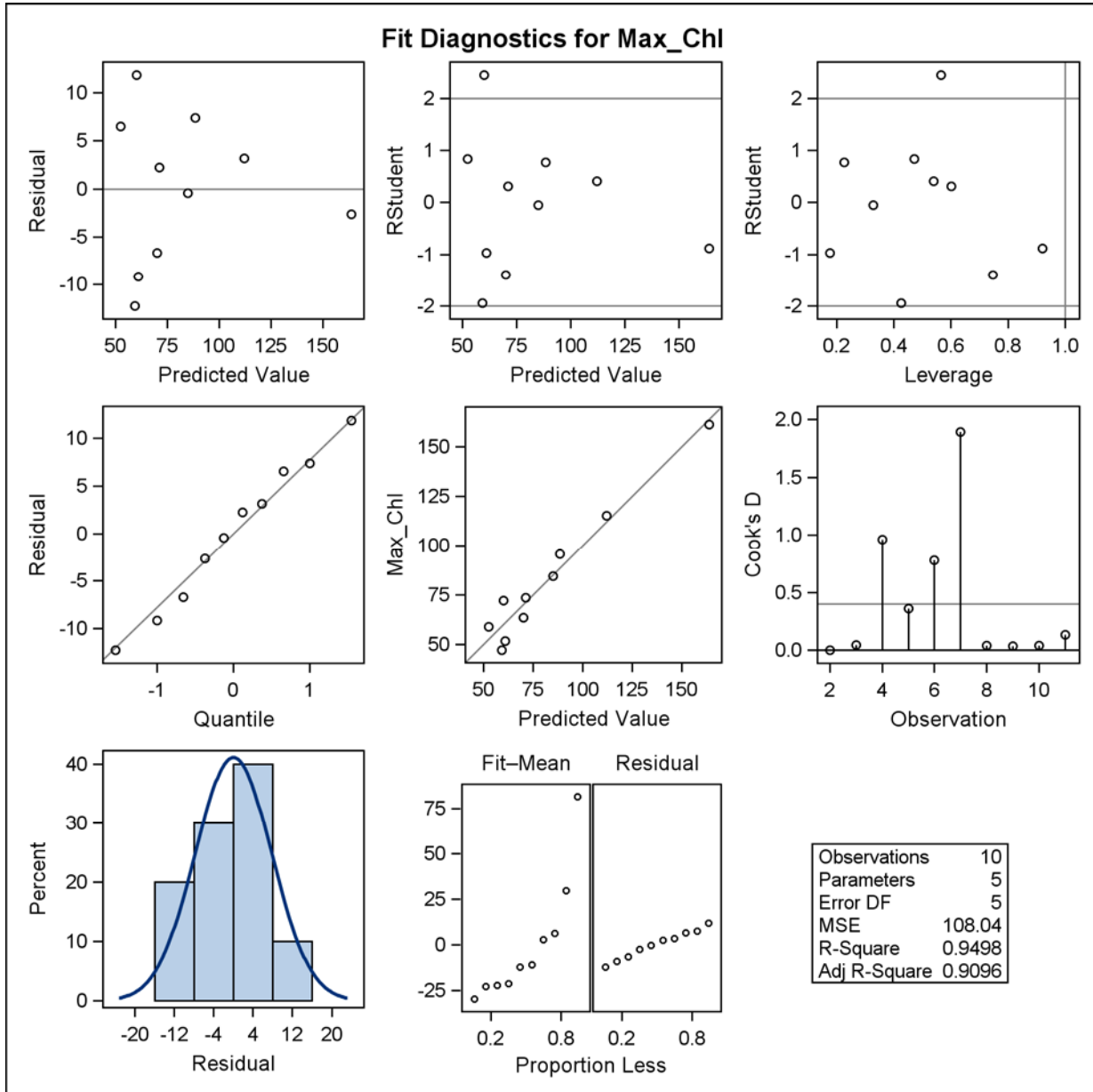
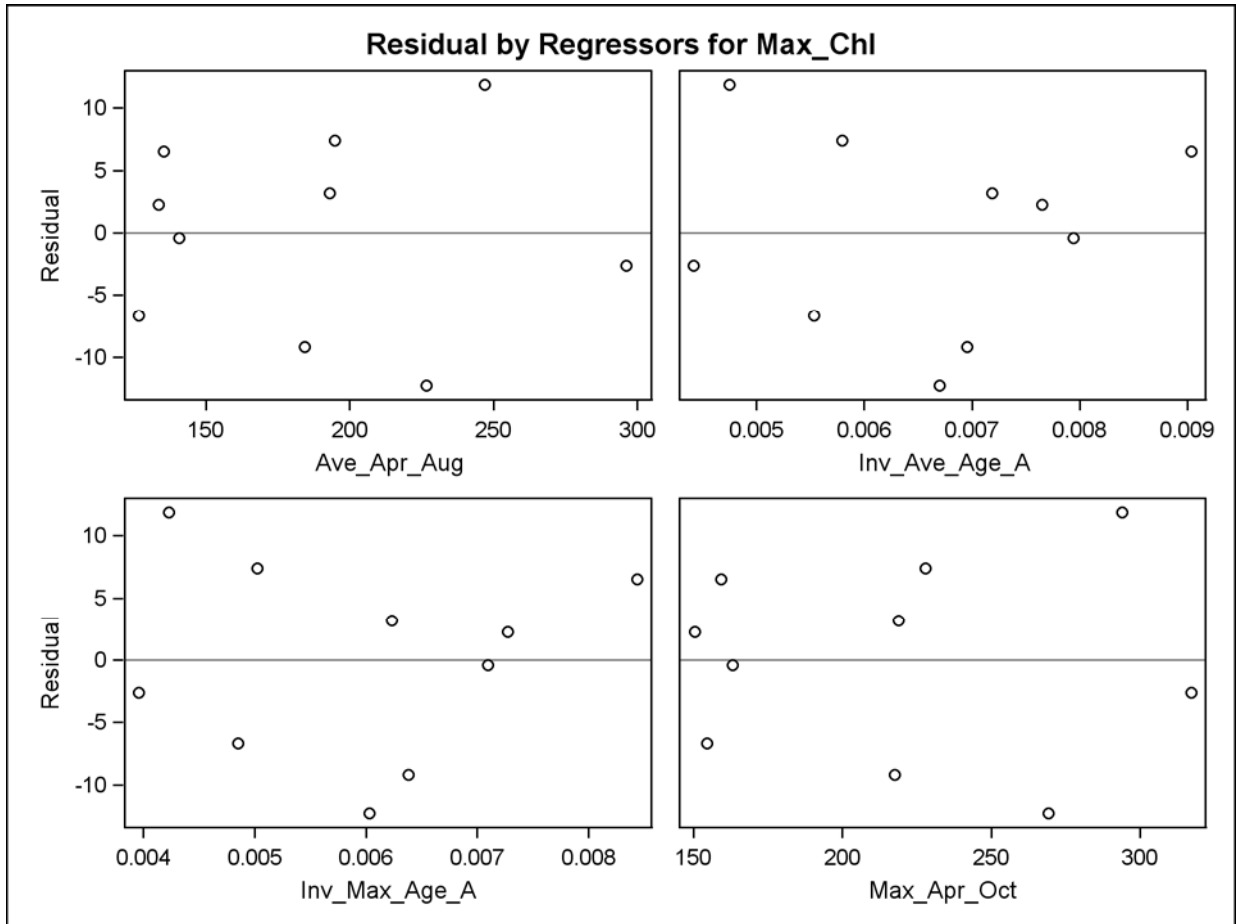## J. Dinoflagellate Biovolume, Segment 2



Fit Diagnostics for LNMaxDino

Appendix 8.C.



Residual by Regressors for LNMaxDino

46

Appendix 8.C.

*K. N$_2$-Fixation, Segments 3 & 4*

Appendix 8.C.



Residual by Regressors for Est_N2_Fix

Appendix 8.C.

*L. Freshwater Maximum Chl-a, Segment 2*



Fit Diagnostics for Max_Chl

Appendix 8.C.



Residual by Regressors for Max_Chl

Appendix 8.C.

*M. Freshwater Maximum Chl-a, Segments 3 & 4*



Fit Diagnostics for Max_Chl

Appendix 8.C.



**Residual by Regressors for Max_Chl**

Appendix 8.C.

*N. Freshwater Maximum Chl-a, Segment 6*



Fit Diagnostics for Max_Chla

Appendix 8.C.



Residual by Regressors for Max_Chla

Appendix 8.C.

*O. Freshwater Bloom Duration, Segment 2*



Fit Diagnostics for Longest_d
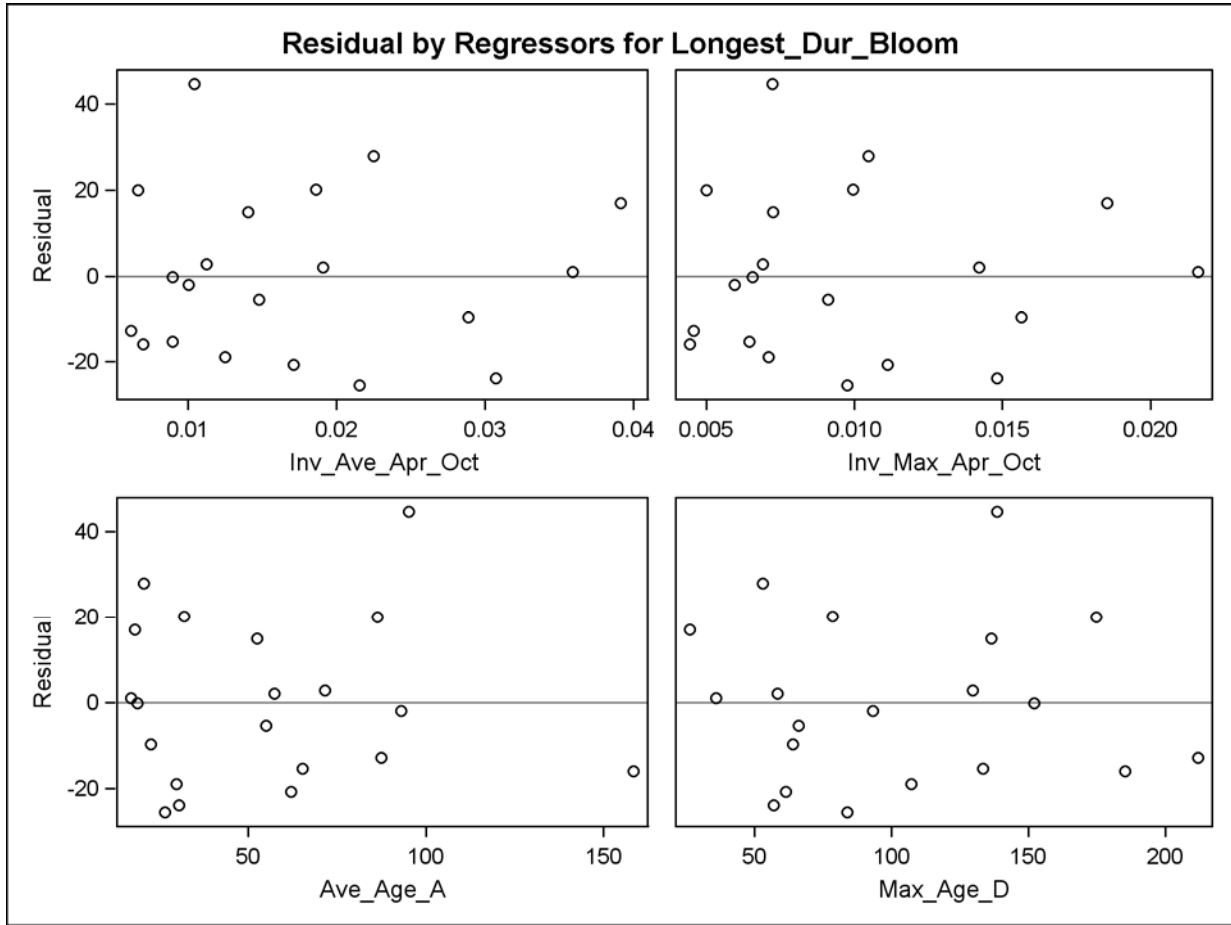
Appendix 8.C.



Residual by Regressors for Longest_d

*P. Freshwater Bloom Duration, Segments 3 & 4*

This regression model used logistic rather than linear regression, and goodness-of-fit metrics are included in model calculation (above) rather than summarized in graphics. The figure below shows observations (coded 1 = occurrence of bloom >50 d, or 0 = no occurrence) versus the probability of a bloom >50 d predicted by the binomial logistic regression. The overall regression was significant (likelihood ratio test of global null hypothesis, $p < 0.001$), and both regression coefficients were significant (Wald maximum likelihood chi-square, $p < 0.05$). The Hosmer-Lemeshow goodness-of-fit test confirmed that the observed frequency of events did not differ significantly from that predicted by the regression (chi-square, $p > 0.05$) (see model calculation results above).

Appendix 8.C.

*Q. Freshwater Bloom Duration, Segment 6*



Fit Diagnostics for Longest_Dur_Bloom

Appendix 8.C.



**Residual by Regressors for Longest_Dur_Bloom**

# References

SAS Institute Inc. 2010. SAS/STAT® 9.22 User's Guide. Cary, NC: SAS Institute Inc.