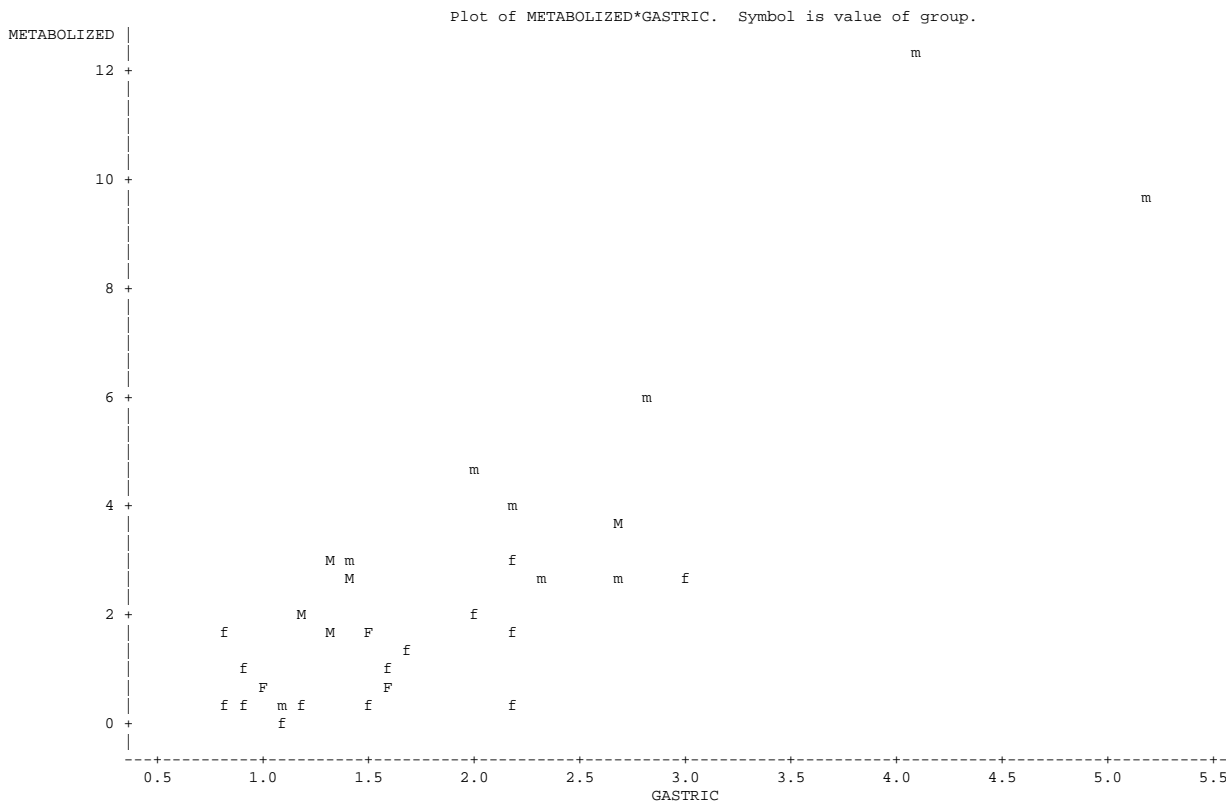**Chapter 11 : Model checking and refinement**

This chapter is primarily concerned with diagnostics for individual observations.  Previously, we have been concerned primarily with evaluating the model (ANOVA source table) and the individual variables, their significance (extra ss and various tests), multicolinearity (VIF and Condition number), relative measures of variation (standardized regression coefficient).   We have also looked at some tests and graphic tools for evaluating the assumptions.  So far we have mentioned outliers, evaluated with residual plots and box plots, but have not developed any formal tools or diagnostics.  That is one of the main themes of this chapter.

**An example:  Alcohol metabolism in men and women**

The study investigates first-pass metabolism of alcohol that is metabolized in the stomach before reaching the bloodstream.  To measure this, investigators administered equal amounts of alcohol orally and intravenously on a series of randomly selected days.  The differences in blood alcohol on these days provides a measure of how much alcohol is metabolized by the stomach (i.e first-pass metabolism).  In addition to this variable the amount of a key gastric enzyme "alcohol dehydrogenase" was measured.  This variable is called gastric AD activity.  Variables in the analysis included dummy variables for sex and for subjects categorized as "alcoholic".

Partial output for this problem is given below.  The complete program and handout are provided online.

The data is plotted below.  The M and F are males and females respectively and the lower case letters are alcoholics.
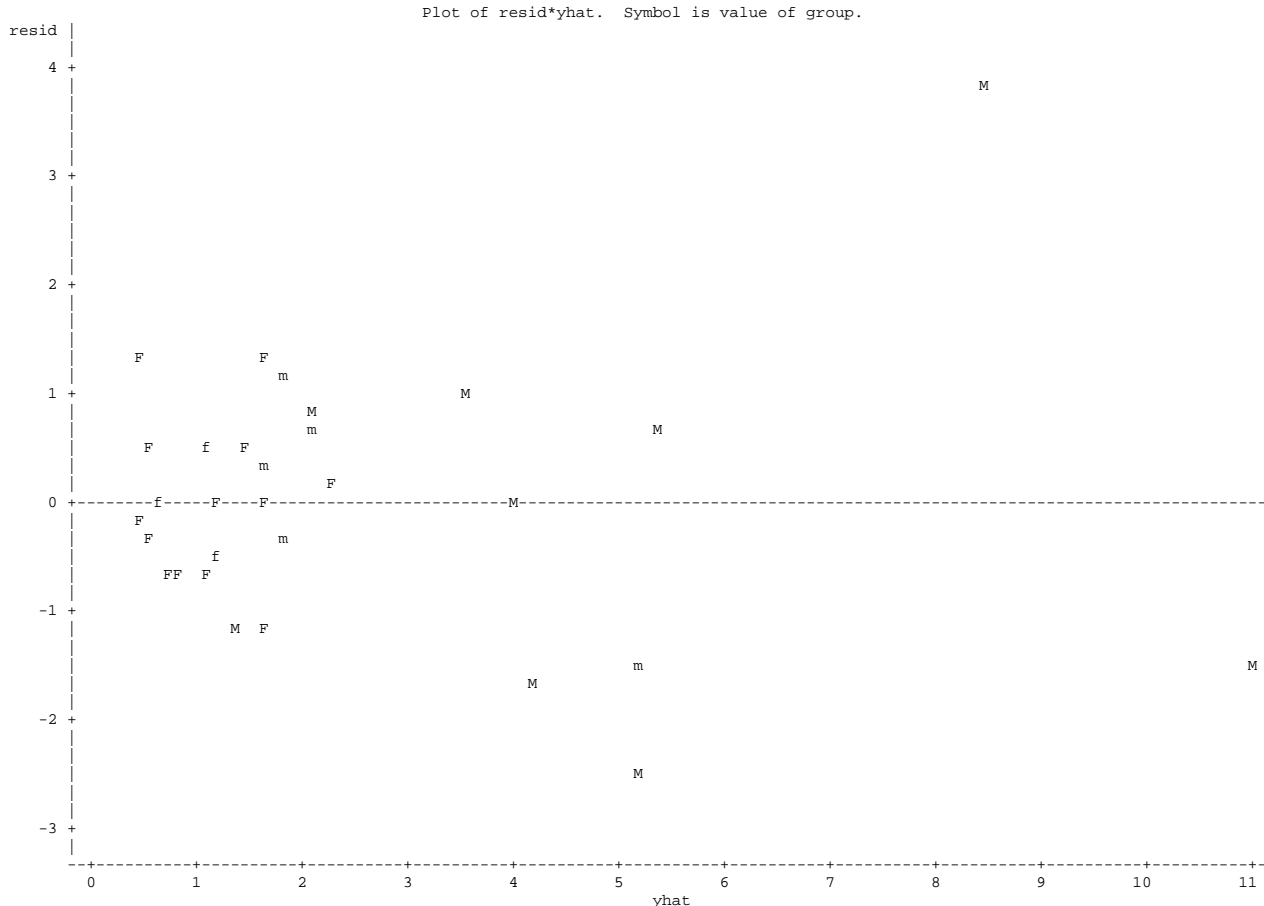
```
                                    Plot of METABOLIZED*GASTRIC.  Symbol is value of group.
   METABOLIZED |
               |                                                                        m
         12 +
               |
               |
               |
               |
         10 +
               |                                                                              m
               |
               |
               |
          8 +
               |
               |
               |
               |
          6 +                                              m
               |
               |
               |
               |                                   m
          4 +
               |                              m
               |                                   M
               |
               |          M m               f
               |            M                   m        m        f
          2 +              M              f
               |      f        M   F           f
               |                          f
               |     f              f
               |       F            F
               |    f f   m f       f              f
          0 +                f
               |
               ---+----------+----------+----------+----------+----------+----------+----------+----------+----------+----------+----------+--
                0.5        1.0        1.5        2.0        2.5        3.0        3.5        4.0        4.5        5.0        5.5
                                                          GASTRIC
NOTE: 1 obs hidden.
```

The analysis, done with PROC REG, is presented below.

```
Analysis of Variance
                                 Sum of            Mean
Source                   DF      Squares          Square     F Value    Pr > F
Model                     7    181.34065        25.90581      16.47     <.0001
Error                    24     37.75404         1.57309
Corrected Total          31    219.09469
```

```
                                   Parameter Estimates
                          Parameter        Standard
Variable          DF       Estimate           Error    t Value    Pr > |t|      95% Confidence Limits
Intercept          1       -1.65966         0.99965      -1.66      0.1099      -3.72283       0.40351
GASTRIC            1        2.51416         0.34337       7.32      <.0001       1.80548       3.22284
Female             1        1.46572         1.33255       1.10      0.2823      -1.28453       4.21597
Alcoholic          1        2.55210         1.94599       1.31      0.2021      -1.46421       6.56842
FEMxALC            1       -2.25171         4.39370      -0.51      0.6130     -11.31986       6.81644
ALCxGastric        1       -1.45874         1.05286      -1.39      0.1786      -3.63173       0.71425
FEMxGastric        1       -1.67344         0.62020      -2.70      0.0126      -2.95347      -0.39341
FEMxALCxGastric    1        1.19867         2.99783       0.40      0.6928      -4.98854       7.38588
```

The overall test of the model indicates a joint significance.  The only single variable in the model that is
significant is GASTRIC, indicating that the gastric activity does indeed correlate to first-pass
metabolism.  There are many variables in the model, and it is possible that removing one or more
of these would reveal other significant correlations.  Formal techniques for model selection will
be discussed in Chapter 12.  For the moment, accept that the only significant variables in this
model, after variables selection, are gastric and the gastric by sex interaction.  When an
interaction is significant both of its main effects are usually included in the model.  Therefore,
sex will be included and the final model is "MODEL METABOLIZED = GASTRIC FEMALE
FEMALExGASTRIC;".  This model is fitted and the residuals plotted below.  It appears there
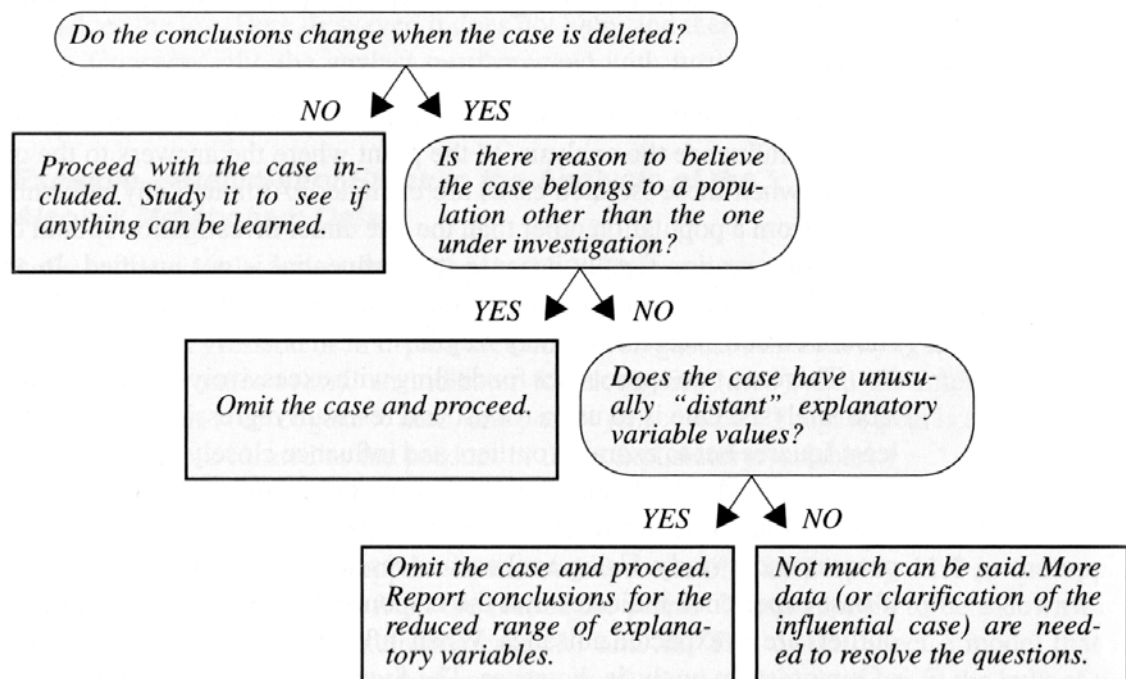may be problems with outliers.



Plot of resid*yhat.  Symbol is value of group.

NOTE: 2 obs hidden.

```
Parameter Estimates
                    Parameter      Standard
Variable     DF      Estimate        Error     t Value    Pr > |t|      95% Confidence Limits
Intercept    1       -1.18577       0.71168     -1.67      0.1068      -2.64359      0.27205
GASTRIC      1        2.34387       0.28015      8.37      <.0001       1.77001      2.91773
Female       1        0.98850       1.07239      0.92      0.3645      -1.20820      3.18519
FEMxGastric  1       -1.50692       0.55914     -2.70      0.0118      -2.65227     -0.36158
```

**Dealing with influential observations** – Outliers are a serious issue in least squares analysis. While the analysis is robust to some departures from the assumptions of X measured without error and normality, the analysis can be seriously and adversely influenced by the presence of outliers.

To make an evaluation on each observation we will need some additional diagnostics. Your textbook also suggests a strategy for dealing with outliers.

**Display 11.8**    A strategy for dealing with suspected influential cases



The book first points out the utility of some simple graphics like scatter plots and residual plots in developing an initial model. The book also suggests three general objectives paraphrased below.

1) *The model should estimate parameters that address research questions.* Obviously, if the model does not address research issues, the model is not really worth pursuing.

2) *Potentially confounding variables should be included.* Confounding, in the statistical sense, refers to closely related variables whose effects in a model often cannot be distinguished. The authors are suggesting that variables suspected of influencing the variable of interest should be included in the model.

3) *The model should take into consideration aspects seen in the initial graphics.* If the graphics display curvature or nonhomogeneous variance these issues should be addressed in the model development.

In order to assist with these decisions and developments we need some new diagnostics. The residual plot for the alcohol metabolism (above and from the book below) has two points that are potentially outliers. Outliers are points that often appear as if they do not fit in with the other observations in the data set because, graphically, they appear to be way out of line with the other points in the data set. The books version of the residual plot shows two observations (#31 and #32) that appear disjoint and separate from the remaining points.

**Display 11.7**  Residual plot from the regression of first-pass metabolism on gastric activity, sex indicator, alcoholism indicator, and all two- and three-factor interactions



There is often a temptation to simply remove an observation that does not fit in with the other observations. However, this action should only be done after careful consideration. The basic question is, "does this observation belong in the population we targeted or does it belong in some other population?" The strategy suggested by the book above will help in making the decision, but better tools are needed.

According to the strategy outlined in the book, if removal of the observation does not change the results of the analysis then it probably belongs in our target population and should not be removed. If the observation does influence our results, is there some reason to assume that it erroneously got included in our sample. Is it a coding error or a case that is influenced by some unusual factors we did not intend to include in our study?

To aid in answering these questions we will consider some additional diagnostics.

1) **Leverage values** – leverage values are produced in the process of doing the matrix algebra for multiple regression. In these calculations the predicted values of $Y_i$ are given by $X(X'X)^{-1}X'Y$ where the first portion "$X(X'X)^{-1}X'$" is known as the hat matrix and its main diagonal are the "hat values", denoted $h_{ii}$ and commonly referred to as leverage values.

Leverage values are between 1/n and 1 and sum to p/n, where p is the number of parameters in the model. Larger values are taken as an indicator of "unusual $X_i$ values". Note that a large value of $h_{ii}$ does not necessarily mean that the observation is a "bad" value, only that it is unusual.

2) **Studentized residuals** – It is not easy to look at residuals and determine which ones are excessively large. The values of residuals depend on the scale of the Y values, so they can be in the hundreds, thousands or thousandths. In order to create values on a recognizable scale, values are "standardized" to a Z scale or t scale ($Z = (Y_i - \mu) / \sigma$). The mean of the residuals is zero and the variances is estimated by the MSE (mean squared error), so the calculation is SemiStudent = $e_i$ / $\sqrt{MSE}$. A better calculation is given by $Student = e_i \Big/ \sqrt{MSE(1-h_{ii})}$. The use of the leverage value scales the variance so that values near the mean of X have smaller variances than those more distant from the center of the X values. This latter calculation is the one used in SAS.

Residuals are assumed to follow a normal distribution, so Studentized residuals should follow a t-distribution or, for large samples, a z distribution. As a result the empirical rules apply; for large samples about 2/3 of the points should be between ±1, 95% between ±2 and 99% between ±2.6.

3) **Deleted studentized residuals** (not discussed in the text) – these values are similar to studentized residuals, but the residuals are calculated as deviations from a regression line that is calculated without the observation present. These residuals, denoted $e_{i(i)}$ are calculated the same way as the Studentized residuals and have similar properties and distributions.

4) **Cook's distance measure (D)** – this value is calculated as $Cook's\ D = \dfrac{\sum_{j=1}^{n}(\hat{Y}_j - \hat{Y}_{j(i)})^2}{pMSE}$. It is a measure of influence that is the change in the model calculated with an observation first included and then excluded. This statistic measures the overall change in the fit across all observations when each point is excluded. Cook's D values greater than 1 are considered large. As with the leverage values, large values are not necessarily "bad" values, but they do have considerable influence on the outcome of the calculations.

5) **DFFITS** – the measure is calculated as $DFFITS_i = \dfrac{\hat{Y}_i - \hat{Y}_{i(i)}}{MSE_{(i)}h_{ii}}$. It is also a measure of influence.

This statistic measures the change in the individual Yhat value when the observation is excluded. DIFFITS values greater than 1 are considered large.

SAS produces all of these values in PROC REG and PROC GLM.  An output statement was used it PROC REG to create these values for the
    Alcohol metabolism experiment.  The program and listing are produced below.

```
PROC REG DATA=Metabolism; Title2 'Fit of metabolized on indicator variables';
   MODEL metabolized = gastric female FEMxGastric / clb alpha=0.05;
   output out=next1 r=resid p=yhat lclm=lclm uclm=uclm lcl=lcli ucl=ucli
         student=student rstudent=rstudent cookd=cookd h=leverage dffits=dffits;
RUN;
proc print data=next1;
   var metabolized gastric female FEMxGastric yhat resid student rstudent
       cookd leverage dffits;
RUN;
```

```
Chapter 11 : Alcohol metabolism in men and women
Fit of metabolized on indicator variables with REG
```

| Obs | METABOLIZED | GASTRIC | Female | FEMx Gastric | yhat | resid | student | rstudent | cookd | leverage | dffits |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.6000 | 1.00000 | 1 | 1.00000 | 0.6397 | -0.03968 | -0.03472 | -0.03410 | 0.00004 | 0.10415 | -0.01163 |
| 2 | 0.6000 | 1.60000 | 1 | 1.60000 | 1.1418 | -0.54185 | -0.46192 | -0.45533 | 0.00316 | 0.05596 | -0.11086 |
| 3 | 1.5000 | 1.50000 | 1 | 1.50000 | 1.0582 | 0.44185 | 0.37667 | 0.37082 | 0.00210 | 0.05596 | 0.09028 |
| 4 | 0.4000 | 2.20000 | 1 | 2.20000 | 1.6440 | -1.24402 | -1.10056 | -1.10489 | 0.04264 | 0.12343 | -0.41460 |
| 5 | 0.1000 | 1.10000 | 1 | 1.10000 | 0.7234 | -0.62337 | -0.54070 | -0.53375 | 0.00706 | 0.08809 | -0.16589 |
| 6 | 0.2000 | 1.20000 | 1 | 1.20000 | 0.8071 | -0.60707 | -0.52288 | -0.51598 | 0.00556 | 0.07523 | -0.14717 |
| 7 | 0.3000 | 0.90000 | 1 | 0.90000 | 0.5560 | -0.25598 | -0.22646 | -0.22259 | 0.00181 | 0.12343 | -0.08352 |
| 8 | 0.3000 | 0.80000 | 1 | 0.80000 | 0.4723 | -0.17229 | -0.15441 | -0.15170 | 0.00102 | 0.14592 | -0.06270 |
| 9 | 0.4000 | 1.50000 | 1 | 1.50000 | 1.0582 | -0.65815 | -0.56106 | -0.55408 | 0.00466 | 0.05596 | -0.13490 |
| 10 | 1.0000 | 0.90000 | 1 | 0.90000 | 0.5560 | 0.44402 | 0.39281 | 0.38680 | 0.00543 | 0.12343 | 0.14514 |
| 11 | 1.1000 | 1.60000 | 1 | 1.60000 | 1.1418 | -0.04185 | -0.03567 | -0.03503 | 0.00002 | 0.05596 | -0.00853 |
| 12 | 1.2000 | 1.70000 | 1 | 1.70000 | 1.2255 | -0.02554 | -0.02181 | -0.02142 | 0.00001 | 0.05917 | -0.00537 |
| 13 | 1.3000 | 1.70000 | 1 | 1.70000 | 1.2255 | 0.07446 | 0.06358 | 0.06244 | 0.00006 | 0.05917 | 0.01566 |
| 14 | 1.6000 | 2.20000 | 1 | 2.20000 | 1.6440 | -0.04402 | -0.03894 | -0.03824 | 0.00005 | 0.12343 | -0.01435 |
| 15 | 1.8000 | 0.80000 | 1 | 0.80000 | 0.4723 | 1.32771 | 1.18997 | 1.19924 | 0.06048 | 0.14592 | 0.49569 |
| 16 | 2.0000 | 2.00000 | 1 | 2.00000 | 1.4766 | 0.52337 | 0.45396 | 0.44743 | 0.00498 | 0.08809 | 0.13906 |
| 17 | 2.5000 | 3.00000 | 1 | 3.00000 | 2.3136 | 0.18643 | 0.19825 | 0.19481 | 0.00637 | 0.39331 | 0.15685 |
| 18 | 2.9000 | 2.20000 | 1 | 2.20000 | 1.6440 | 1.25598 | 1.11115 | 1.11601 | 0.04346 | 0.12343 | 0.41877 |
| 19 | 1.5000 | 1.30000 | 0 | 0.00000 | 1.8613 | -0.36127 | -0.31925 | -0.31407 | 0.00352 | 0.12150 | -0.11680 |
| 20 | 1.9000 | 1.20000 | 0 | 0.00000 | 1.6269 | 0.27312 | 0.24287 | 0.23875 | 0.00225 | 0.13242 | 0.09327 |
| 21 | 2.7000 | 1.40000 | 0 | 0.00000 | 2.0957 | 0.60435 | 0.53110 | 0.52418 | 0.00886 | 0.11165 | 0.18583 |
| 22 | 3.0000 | 1.30000 | 0 | 0.00000 | 1.8613 | 1.13873 | 1.00631 | 1.00655 | 0.03501 | 0.12150 | 0.37432 |
| 23 | 3.7000 | 2.70000 | 0 | 0.00000 | 5.1427 | -1.44269 | -1.24695 | -1.25997 | 0.03456 | 0.08165 | -0.37569 |
| 24 | 0.3000 | 1.10000 | 0 | 0.00000 | 1.3925 | -1.09249 | -0.97829 | -0.97752 | 0.04039 | 0.14442 | -0.40161 |
| 25 | 2.5000 | 2.30000 | 0 | 0.00000 | 4.2051 | -1.70514 | -1.46571 | -1.49791 | 0.04136 | 0.07150 | -0.41566 |
| 26 | 2.7000 | 2.70000 | 0 | 0.00000 | 5.1427 | -2.44269 | -2.11128 | -2.26100 | 0.09908 | 0.08165 | -0.67418 |
| 27 | 3.0000 | 1.40000 | 0 | 0.00000 | 2.0957 | 0.90435 | 0.79474 | 0.78937 | 0.01985 | 0.11165 | 0.27984 |
| 28 | 4.0000 | 2.20000 | 0 | 0.00000 | 3.9708 | 0.02925 | 0.02514 | 0.02469 | 0.00001 | 0.07165 | 0.00686 |
| 29 | 4.5000 | 2.00000 | 0 | 0.00000 | 3.5020 | 0.99802 | 0.85960 | 0.85547 | 0.01502 | 0.07519 | 0.24393 |
| 30 | 6.1000 | 2.80000 | 0 | 0.00000 | 5.3771 | 0.72293 | 0.62663 | 0.61970 | 0.00934 | 0.08688 | 0.19115 |
| 31 | 9.5000 | 5.20000 | 0 | 0.00000 | 11.0024 | -1.50237 | -1.82581 | -1.91022 | 0.96070 | 0.53548 | -2.05094 |
| 32 | 12.3000 | 4.10000 | 0 | 0.00000 | 8.4241 | 3.87589 | 3.71413 | 5.12052 | 1.16725 | 0.25287 | 2.9790 |

Interpretation of these statistics is often aided by a graphic presentation.  Graphics and the program statements that produced them are given below.  The reference lines are values suggested by the text and other sources.

```
Chapter 11 : Alcohol metabolism in men and women
Various plot with group variable

                          Plot of student*SUBJECT.  Symbol is value of group.
        4 +                                                                            M
          |
          |
          |
          |
        2 +--------------------------------------------------------------------------------
          |
          |
 student  |                                 F           F          m
          |                                                               M        M
          |                                      F                    m             M
          |           f                     F          F                               M
        0 +   f                         F  F  F  F                             M
          |                        F  F                     m
          |      f          F  F         F
          |                                                         M
          |         F                                 m          M
          |                                                        M                 M
       -2 +                                                   M
          |
          |
          |
          |
       -4 +
          ---+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--
             1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32
                                                      SUBJECT
```

```
Chapter 11 : Alcohol metabolism in men and women
Various plot with group variable

                          Plot of rstudent*SUBJECT.  Symbol is value of group.
 rstudent |
          |
        6 +
          |
          |                                                                               M
          |
          |
        4 +
          |
          |
          |
        2 +--------------------------------------------------------------------------------
          |
          |                                 F           F          m
          |                                                               M        M  M
          |           f                     F                   m  m
        0 +   f                         F         F  F  F  F          F                M
          |      f          F  F  F         F                          F
          |                                                         m
          |         F                                 m          M
          |                                                        M                 M
       -2 +                                                   M
          |
          |
          |
       -4 +
          |
          ---+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--
             1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32
                                                      SUBJECT
```

Chapter 11 : Alcohol metabolism in men and women
Various plot with group variable

```
               Plot of leverage*SUBJECT.  Symbol is value of group.
  leverage |
    0.6 +
        |
        |                                                                         M
    0.5 +-------------------------------------------------------------------------
        |
    0.4 +                                              F
        |
    0.3 +
        |                                                                       M
    0.2 +
        |                          F                          F        M
    0.1 +  f           F      F        F    F            F     F   F m m  m        M
        |       f  f          F           F          F F F          m    M M    M M M
    0.0 +
        |
        --+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--
          1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32
                                            SUBJECT
```

Chapter 11 : Alcohol metabolism in men and women
Various plot with group variable

```
               Plot of cookd*SUBJECT.   Symbol is value of group.
  cookd |
   1.25 +
        |                                                                           M
        |
   1.00 +---------------------------------------------------------------------------
        |                                                                           M
        |
   0.75 +
        |
        |
   0.50 +
        |
        |
   0.25 +
        |                                                             M
   0.00 +f   f   f     F       F  F  F  F  F  F  F  F  F  F    F    F     F    m  m  m    m m M M     M M M M
        |
        +---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
          1   2   3   4   5   6   7   8   9  10  11  12  13  14  15  16  17  18  19  20  21  22  23  24  25  26  27  28  29  30  31  32
                                            SUBJECT
```

Chapter 11 : Alcohol metabolism in men and women
Various plot with group variable

```
               Plot of dffits*SUBJECT.  Symbol is value of group.
  dffits |
      3 +                                                                          M
        |
        |
      2 +
        |
        |
      1 +-------------------------------------------------------------------------
        |                                    F     F
        |                          F      F F    m            m M    M M
      0 +  f      f      F F       F F F F          m           M
        |     f        F F       F                 m        m M M
        |         F                                             M
     -1 +
        |
        |
     -2 +                                                          M
        |
        --+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--
          1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32
                                            SUBJECT
```

**Partial residual plots** – Recall that when each variable is entered in the model it alters the other variables in the model.  Scatter plots of the original variables ($Y_i$ on $X_i$) are of little use in examining potential relationships with $Y_i$ because of the effect of the $X_i$ variables on each other.  Partial residual plots provide a means of evaluating the relationship of the $Y_i$ variable and $X_i$ variable after adjusting for other variables.  Essentially the residuals of $Y_i$ and each $X_i$ are adjusted for all other $X_i$ variables and plotted on each other.  Although these are "residual" plots they can be used like scatter plots to examine for the strength of the relationship, curvature, homogeneity and outliers.  These are provided by PROC REG if the option "PARTIAL" is requested after the model statement.

```
options ps=60 ls=132;
PROC REG DATA=Metabolism; Title2 'Fit of metabolized on indicator variables with
      REG';
   MODEL metabolized = gastric female FEMxGastric / partial;
RUN;
```

```
Chapter 11 : Alcohol metabolism in men and women
Fit of metabolized on indicator variables with REG

The REG Procedure
Model: MODEL1
Partial Regression Residual Plot

            ---+----+----+----+----+----+----+----+----+----+----+----+----+----+----+----+----+---
       6 +                                                                                        +
         |                                                                                        |
         |                                                                                        |
         |                                                                                        |
         |                                                                                        |
         |                          1                                                             |
       4 +                                                                                        +
         |                                                                                        |
METABOLIZED |                                                                                     |
         |                                                                                        |
         |                                                                                        |
       2 +                                                                                        +
         |                                                                                        |
         |                                           2                                            |
         |                                          11           1                                |
         |                                          3                    1 1                       |
       0 +                                          5                    1                         +
         |                                          2         1             1                      |
         |            1                             4                                              |
         |                                                            1                            |
         |                                          1                                              |
         |                                            1                         1                  |
      -2 +                                                  1                                      +
         |                                                                                         |
         |                                          1                                              |
         |                                                                                         |
         |                                                                                         |
      -4 +                                                                                        +
            ---+----+----+----+----+----+----+----+----+----+----+----+----+----+----+----+----+---
             -0.9 -0.8 -0.7 -0.6 -0.5 -0.4 -0.3 -0.2 -0.1  0.0  0.1  0.2  0.3  0.4  0.5  0.6  0.7
                                              Intercept
```

```
Parameter Estimates
                      Parameter       Standard
Variable       DF     Estimate         Error      t Value    Pr > |t|      95% Confidence Limits
Intercept      1      -1.18577        0.71168      -1.67      0.1068      -2.64359       0.27205
GASTRIC        1       2.34387        0.28015       8.37      <.0001       1.77001       2.91773
Female         1       0.98850        1.07239       0.92      0.3645      -1.20820       3.18519
FEMxGastric    1      -1.50692        0.55914      -2.70      0.0118      -2.65227      -0.36158
```

```
       -------+------+------+------+------+------+------+------+------+------+-------
    8 +                                              1
      |
      |
      |
      |
    6 +                                                                           +
      |
      |                                                      1
      |
    4 +                                                                           +
      |
METABOLIZED
      |
      |
    2 +                                    1                                      +
      |
      |                        2
      |                    1   3
    0 +                        17                                                 +
      |                        2       1
      |                        3
      |            11          1
      |             1          1   1
   -2 +                                                                           +
      |        1
      |         1
      |
      |
   -4 +     1                                                                     +
       -------+------+------+------+------+------+------+------+------+------+-------
            -1.5    -1.0    -0.5    0.0     0.5     1.0    1.5     2.0    2.5     3.0
                                        GASTRIC
```

```
METABOLIZED  ---+---------+---------+---------+---------+---------+---------+---------+---
    4 +                                         1                                   +
      |
      |
      |
      |
    2 +                                                          1                  +
      |
      |                      1
      |           1          1          1                     1
      |             1            1             1
      |             1
    0 +           1            1         2  1            1  1    1  1               +
      | 1                      1                1  1                                +
      |             1
      |                                                          1
      |        1            1          1
      |                      1
   -2 +                                                                            +
      |                         1
      |
      |
      |
   -4 +                                                                            +
      |
       ---+---------+---------+---------+---------+---------+---------+---------+---
        -0.4      -0.3      -0.2      -0.1      0.0       0.1       0.2      0.3     0.4
                                        Female
```

```
Parameter Estimates    Parameter        Standard
Variable      DF         Estimate          Error      t Value    Pr > |t|      95% Confidence Limits
GASTRIC       1          2.34387         0.28015        8.37      <.0001       1.77001       2.91773
Female        1          0.98850         1.07239        0.92      0.3645      -1.20820       3.18519
```
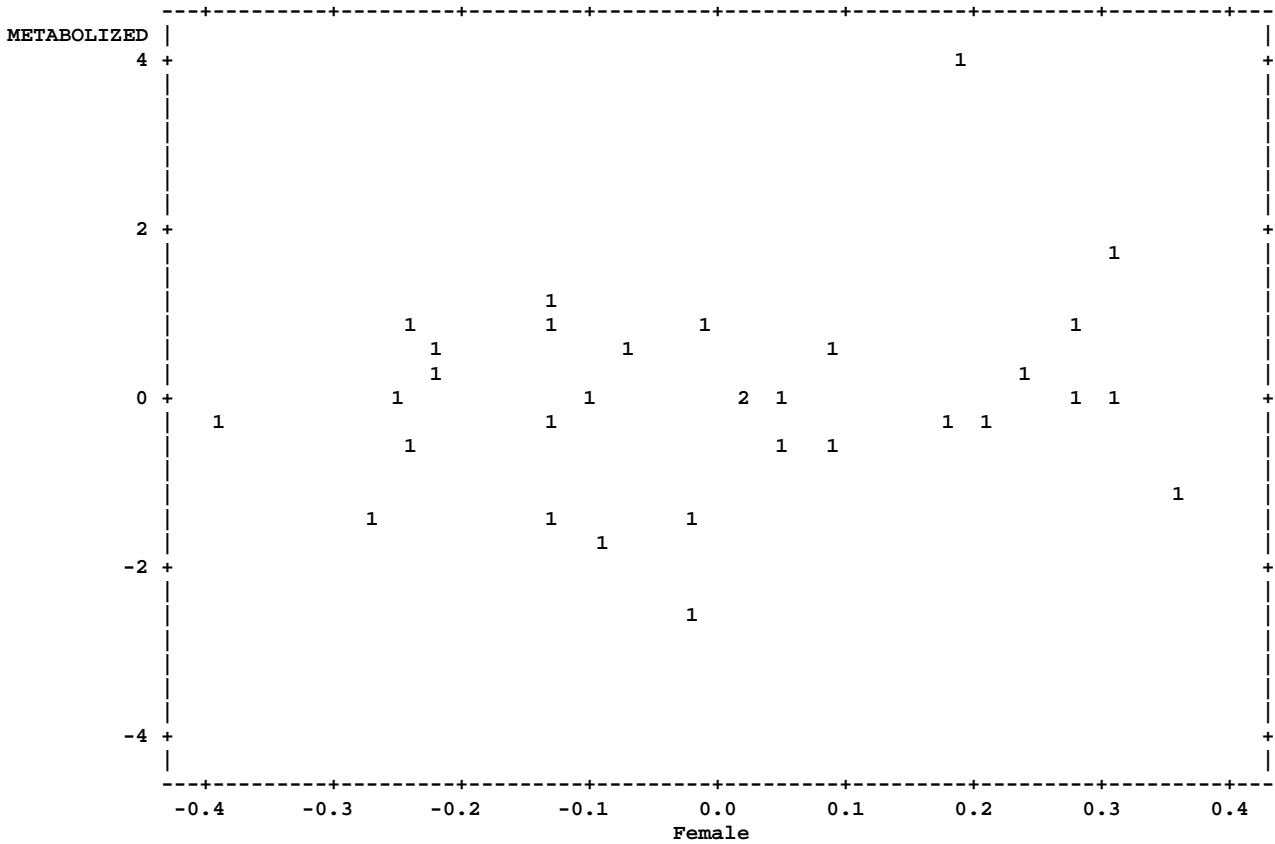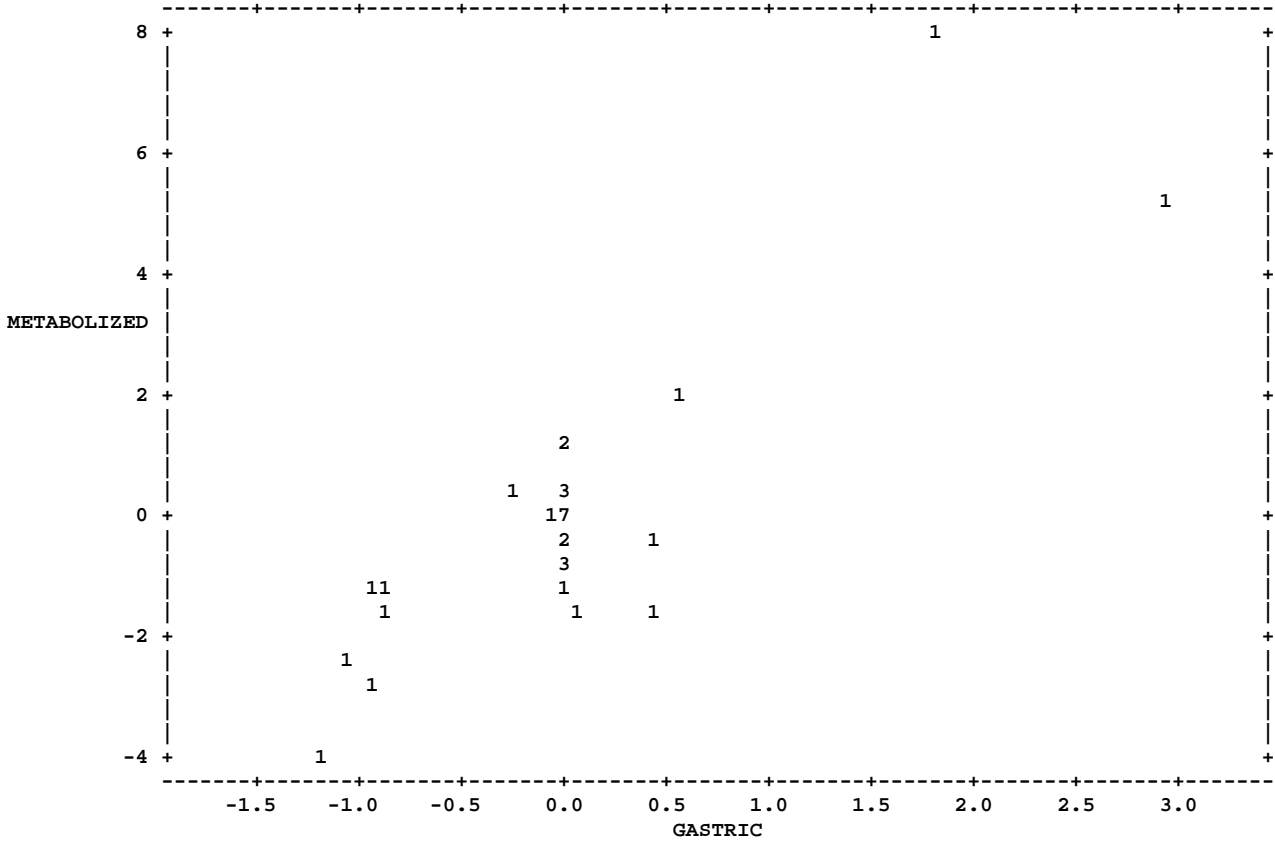
```
Chapter 11 : Alcohol metabolism in men and women
Fit of metabolized on indicator variables with REG

The REG Procedure
Model: MODEL1
Partial Regression Residual Plot
```
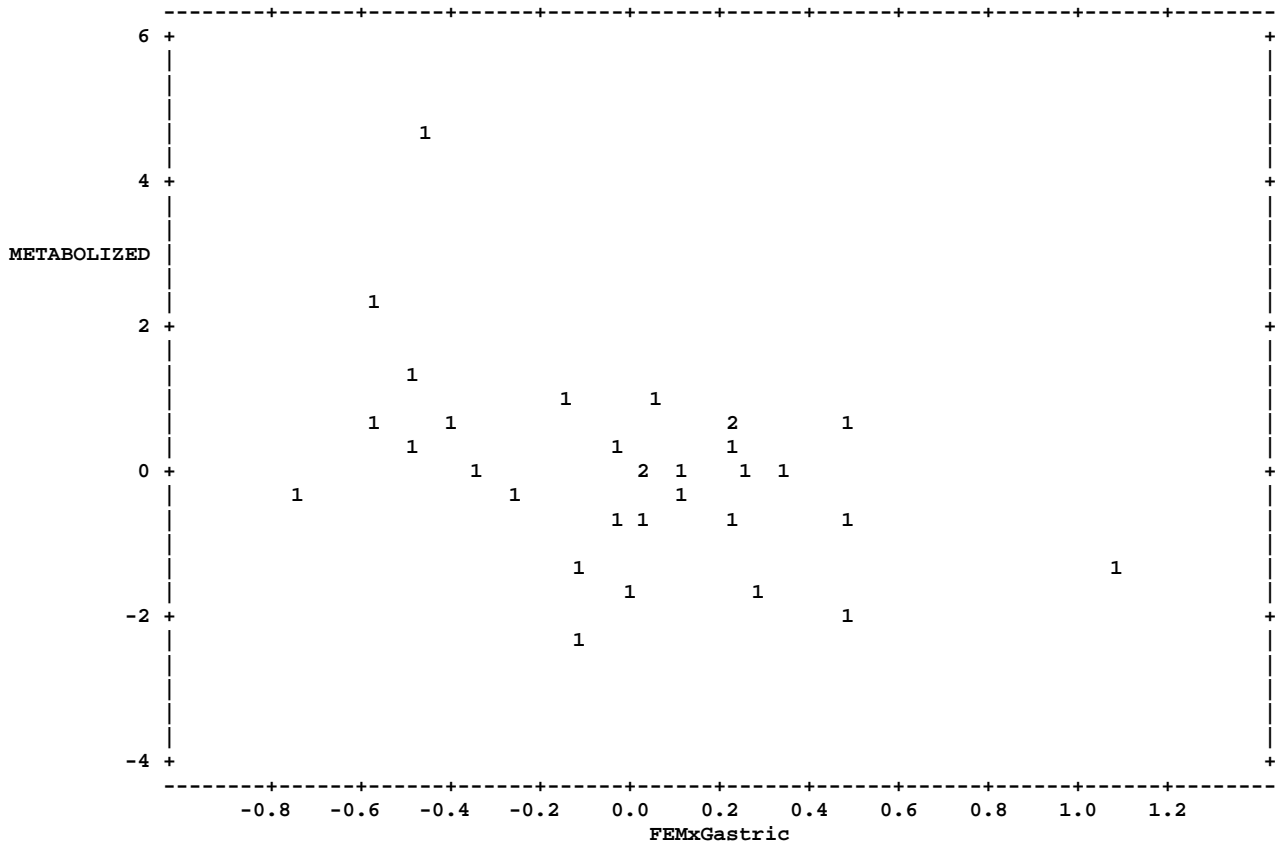
```
            --------+------+------+------+------+------+------+------+------+------+------+--------
         6 +                                                                                      +
           |
           |
           |
           |                          1
           |
         4 +                                                                                      +
           |
           |
METABOLIZED|
           |                      1
         2 +                                                                                      +
           |
           |                     1
           |                              1       1
           |                 1     1                   2         1
           |                   1            1          1
         0 +               1          1        2  1      1  1                                     +
           |           1              1           1
           |                              1 1       1         1
           |
           |                        1                                       1
           |                          1            1
        -2 +                                          1                                           +
           |                        1
           |
           |
           |
        -4 +                                                                                      +
            --------+------+------+------+------+------+------+------+------+------+------+--------
                 -0.8   -0.6   -0.4   -0.2    0.0    0.2    0.4    0.6    0.8    1.0    1.2
                                              FEMxGastric
```

```
Parameter Estimates
                      Parameter      Standard
Variable      DF      Estimate         Error    t Value    Pr > |t|      95% Confidence Limits
Intercept      1      -1.18577       0.71168      -1.67      0.1068      -2.64359       0.27205
GASTRIC        1       2.34387       0.28015       8.37      <.0001       1.77001       2.91773
Female         1       0.98850       1.07239       0.92      0.3645      -1.20820       3.18519
FEMxGastric    1      -1.50692       0.55914      -2.70      0.0118      -2.65227      -0.36158
```