

A GENERAL SAS MACRO FOR PERFORMING WEIGHTED LEAST SQUARES

Wanda H. Burton, Medical College of Virginia

The SAS procedure GLM provides an excellent means for performing least squares regression analysis when the usual model assumptions can be made. The model referred to is

$$Y = X\beta + \epsilon,$$

where ϵ is normally distributed with mean 0 and variance $\sigma^2 I$.

The case to be considered here is that in which the observations remain independent, but their variances are not all equal. The form of the variance-covariance matrix is $\sigma^2 V$ where V is a diagonal matrix with unequal diagonal elements,

$$\sigma^2 V = \begin{bmatrix} \sigma_1^2 & & & 0 \\ & \sigma_2^2 & & \\ & & \ddots & \\ 0 & & & \sigma_n^2 \end{bmatrix}$$

A unique nonsingular symmetric matrix P can be found such that

$$P^2 = V.$$

A transformation can then be made on our original model by premultiplying by P^{-1} , obtaining

$$P^{-1}Y = P^{-1}X\beta + P^{-1}\epsilon$$

or

$$Z = Q\beta + f$$

with obvious substitutions. This model satisfies the necessary assumptions for carrying out the usual least squares regression analysis; that is, $f \sim N(0, \sigma^2 I)$.

The MACRO to be presented provides a thorough analysis for a simple linear regression model,

$$E(y_1) = \beta_0 + \beta_1 x_1.$$

Let us use the following notation for the variance of Y :

$$\text{Var}(Y) = \sigma^2 V = \sigma^2 \begin{bmatrix} w_1 & & & 0 \\ & w_2 & & \\ & & \ddots & \\ 0 & & & w_n \end{bmatrix}$$

where the w 's are known weights.

For this situation,

$$P^{-1} = \begin{bmatrix} 1/\sqrt{w_1} & & & 0 \\ & 1/\sqrt{w_2} & & \\ & & \ddots & \\ 0 & & & 1/\sqrt{w_n} \end{bmatrix}$$

and a simple transformation of the variables is appropriate. These calculations are carried out in the MACRO WT REGR in statements 4-10, creating the variables Z, Q_0 and Q_1 .

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad P^{-1}Y = Z = \begin{bmatrix} y_1/\sqrt{w_1} \\ y_2/\sqrt{w_2} \\ \vdots \\ y_n/\sqrt{w_n} \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}$$

$$P^{-1}X = Q = (Q_0, Q_1) = \begin{bmatrix} 1/\sqrt{w_1} & X_1/\sqrt{w_1} \\ 1/\sqrt{w_2} & X_2/\sqrt{w_2} \\ \vdots & \vdots \\ 1/\sqrt{w_n} & X_n/\sqrt{w_n} \end{bmatrix}$$

These transformations are made after the regression variables, X and Y , and the weights, WT , are input separately and merged in statements 27-65. An identification variable, ID , is included in each of the input data sets and is necessary for the merger.

PROC GLM is used to compute the ordinary least squares regression of Z on Q_0 and Q_1 . Notice that the MODEL statement is written with NOINT option because the vector of ones normally present in a simple linear regression model has been transformed to a non-constant vector, Q_0 . The OUTPUT facilities are used for access to the predicted values and the residuals from the fitted line.

The MACRO WT_REGR includes a scatter of the original dependent variable (Y) against the independent variable (X). It provides plots which allow the user to examine the residuals from the analysis and thereby judge the effectiveness of applying the techniques of weighted regression to the data. A listing of the original variables with the scaled predicted values and residuals is included. This transformation back to the original scale of the data is carried out in statements 16-19.

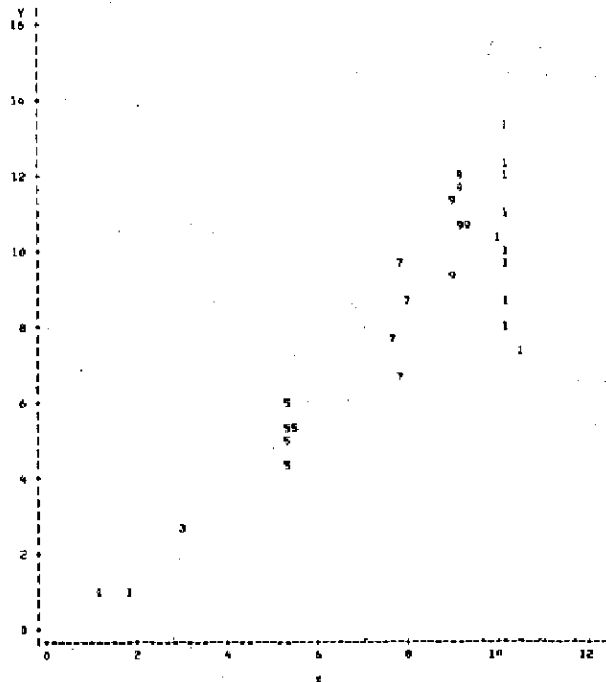
Acknowledgment

The author would like to thank Maria Sauer for the use of laboratory data from the Pulmonary Division of the Department of Medicine, Medical College of Virginia. This work was supported in part by grant 1 R01 01899-02 from the National Center for Health Services Research, HRA.

Reference

Draper, N. R. and Smith, H. Applied Regression Analysis, John Wiley & Sons, Inc., 1966, pp. 77-81.

STATISTICAL ANALYSIS 5
PLOT OF X*Y



1 STATISTICAL ANALYSIS 5
NOTE: THE JOB SORTS HAS BEEN RUN UNDER RELEASE 76.2 OF SAS AT VIRGINIA COMMONWEALTH

```

1 MACRO WT_REGR
2 PROC SORT DATA=DRIG BY ID;
3 PROC SORT DATA=WEIGHTS BY ID;
4 DATA TRANS;
5 MERGE ORIG WEIGHTS;
6 BY ID;
7 W=SORTINT(1);
8 Z=Y/W;
9 Q_0=1/W;
10 Q_1=Y/W;
11 PROC SCATTER;
12 PLOT X*Y=ID / MPOS=75;
13 PROC GLM;
14 MODEL Z = Q_0 Q_1 / NOINT S&I R;
15 OUTPUT OUT=PRED PREDICTED=Z_MAT RESIDUAL=Z_RESID;
16 DATA BACK;
17 SET PRED;
18 Y_MAT=Q_0*Z_MAT;
19 Y_RESID=Q_1*Z_RESID;
20 PROC PRINT;
21 VAR ID W X Y Y_MAT Y_RESID;
22 PROC SCATTER;
23 PLOT X*Y_RESID=Z' / MPOS=75;
24 PROC SCATTER;
25 PLOT X*Y_MAT=Y' / MPOS=75;
26 S;
27 DATA ORIG INPUT X 1-5 Y 6-10 ID 11-12;
28 CARDS;

```

NOTE: DATA SET WORK.ORIG HAS 35 OBSERVATIONS AND 3 VARIABLES.
NOTE: THE DATA STATEMENT USED 0.41 SECONDS AND 96K.

```

64 DATA WEIGHTS INPUT ID W;
65 CARDS;

```

NOTE: DATA SET WORK.WEIGHTS HAS 6 OBSERVATIONS AND 2 VARIABLES.
NOTE: THE DATA STATEMENT USED 0.16 SECONDS AND 96K.

```

72 S;
73 WT_REGR;

```

NOTE: DATA SET WORK.ORIG HAS 35 OBSERVATIONS AND 3 VARIABLES.
NOTE: THE PROCEDURE SORT USED 1.34 SECONDS AND 110K.

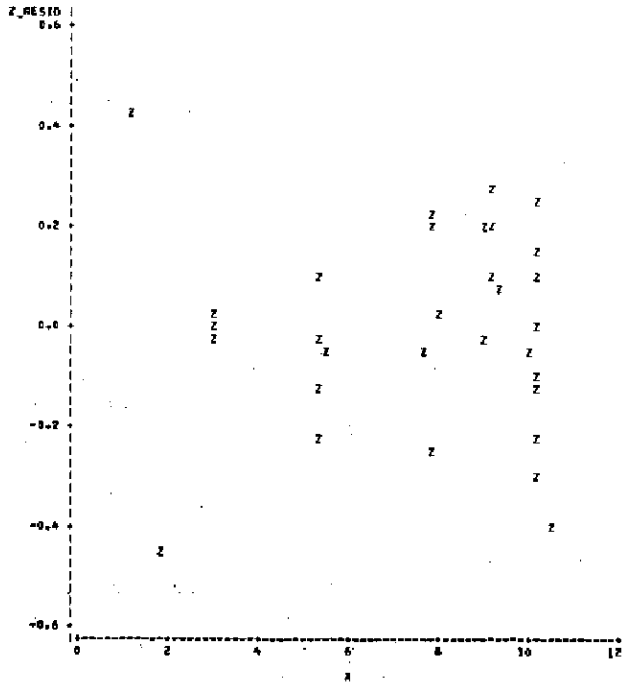
NOTE: DATA SET WORK.WEIGHTS HAS 6 OBSERVATIONS AND 2 VARIABLES.
NOTE: THE PROCEDURE SORT USED 1.36 SECONDS AND 110K.

NOTE: DATA SET WORK.TRANS HAS 35 OBSERVATIONS AND 8 VARIABLES.
NOTE: THE DATA STATEMENT USED 0.27 SECONDS AND 104K.

STATISTICAL ANALYSIS SYSTEM

OBS	ID	WT	X	Y	Y_MAT	Y_RESID
1	1	1	1.15	0.99	0.5738	0.4162
2	1	1	1.95	0.98	1.4354	-0.4554
3	3	9	3.05	2.60	2.4690	-0.8690
4	3	9	3.05	2.67	2.4690	-0.7990
5	3	9	3.85	2.66	2.4690	-0.8090
6	3	9	3.05	2.74	2.4690	-0.7290
7	3	9	3.00	2.80	2.4690	-0.3390
8	5	25	5.24	5.97	5.3872	0.5828
9	5	25	5.28	6.26	5.4331	0.8269
10	5	25	5.40	4.31	5.4561	-1.1461
11	5	25	5.40	4.89	5.4561	-0.5661
12	5	25	5.45	5.21	5.5135	-0.3035
13	7	49	7.78	7.66	8.0963	-0.4163
14	7	49	7.88	9.81	8.2132	1.5968
15	7	49	7.81	6.57	8.2946	-1.7246
16	7	49	7.85	9.71	8.2786	1.4314
17	7	49	7.87	9.82	8.2936	1.5264
18	7	49	7.91	9.81	8.3395	1.4705
19	7	49	7.96	8.50	8.3740	0.1260
20	9	81	8.03	9.47	9.4261	-0.1561
21	9	81	9.07	11.06	9.6781	1.3819
22	9	81	9.11	12.14	9.7186	2.4214
23	9	81	9.14	11.58	9.7525	1.8275
24	9	81	9.16	10.65	9.7755	0.8755
25	9	81	9.37	10.64	10.0167	0.6233
26	10	100	10.17	9.74	10.9257	-1.1857
27	10	100	10.18	12.39	10.9472	1.4428
28	10	100	10.22	11.03	10.9932	0.0368
29	10	100	10.22	8.00	10.9932	-2.9932
30	10	100	10.22	11.90	10.9932	0.9068
31	10	100	10.18	8.68	10.9472	-2.2672
32	10	100	10.50	7.25	11.3346	-4.0846
33	10	100	10.23	13.64	11.0047	2.6353
34	10	100	10.63	10.19	10.7769	-0.5869
35	10	100	10.22	9.94	11.0047	-1.0647

STATISTICAL ANALYSIS 5
 PLOT OF X*Z_RESID



STATISTICAL ANALYSIS 5
 PLOT OF X*Y_RESID

