# An Introduction to the SAS® Survey Analysis PROCs

**Xiuhua Chen and Paul Gorrell**
**Social & Scientific Systems, Inc.**

## Abstract

This paper provides an introductory overview of the SAS survey analysis PROCS: SURVEYMEANS, SURVEYFREQ, SURVEYLOGISTIC and SURVEYREG.  For all procedures we will discuss the STRATA, CLUSTER and WEIGHT statements with respect to specifying the sample design and generating weighted estimates.  We will also discuss the importance of not deleting observations in order to analyze subpopulations, as well as analyzing data with missing values.  Other PROC-specific statements will be discussed, e.g. TABLES in SURVEYFREQ, DOMAIN in SURVEYMEANS and MODEL in SURVEYLOGISTIC.  In addition we will briefly discuss the Taylor linearization method for variance estimation used by SAS and how SAS (by default) handles situations where there is only one cluster in a stratum.  Concrete examples will be used (with SAS 9.1.3).  This paper will be of interest to SAS programmers who need to analyze the results of complex, i.e. not simple, random, surveys.

## Introduction

Before the introduction of the survey procedures, SAS could not be used for variance estimation with complex survey data because the existing statistical procedures, e.g. FREQ, MEANS, did not have the functionality for incorporating the design properties of complex surveys into the analysis, i.e. they assumed the design to be a simple random sample of an infinite population.  In our experience this assumption generally leads to an underestimation of variance.  This introductory paper will not answer all the questions you may have about using SAS for analyzing complex surveys, but we hope that reading this paper will put you in a better position to know what questions to ask when you first need to use a SAS survey procedure,

For a paper on variance estimation, this paper has a notable lack of statistical formulae.  This is because we have targeted SAS programmers rather than statisticians.  We have taken a 'functional' approach to discussing the survey procedures.  The SAS documentation contains excellent discussions of the statistical details which are built into the respective procedures.

## Variance Estimation with Complex Surveys

A complex survey design is any design that is not a simple random sample (where each person in the target population has an equal chance of being selected in the survey).  Complex designs involve dividing the population into groups and sampling the groups themselves (cluster sampling) and/or sampling from the groups (stratified sampling).

In SAS 9.1.3 the only available method for variance estimation is Taylor linearization.  Replication methods (e.g. jackknife and balanced repeated measures) are available in SAS 9.2 but will not be discussed in this paper (but see Mukhopadhyay *et al*, 2008).

One important point to emphasize when analyzing data from complex surveys is that you should resist the normally efficient practice of subsetting to just those records that will be used in the analysis.  This subverts the process of giving the procedures all the necessary survey design information with respect to strata and clusters.  That is, if you are interested in a particular analysis of city-dwellers within a nationally-representative survey and subset to just this subpopulation before your analysis, you may be deleting information about the full design

properties of the survey.  You are more or less likely to do this depending on the criteria for subsetting.  For example, subsetting to males from a typical national survey population is unlikely to exclude all records in a particular PSU and remove it from the analysis (and if it did you'd really want to know why!).  Below we discuss the DOMAIN option for analyzing subpopulations within a sample population.  Examples are shown for PROCs SURVEYMEANS and SURVEYLOGISTIC.

All examples used in this paper use as input the Medical Expenditure Panel Survey (MEPS), discussed in the next section.  No matter what data source is used as input, correct variance estimation for complex survey data requires knowing the design properties of the survey being analyzed.

### The Household Component of the Medical Expenditure Panel Survey (MEPS HC)

The MEPS HC is a nationally representative survey of the U.S. civilian noninstitutionalized population.  It collects medical expenditure data as well as information on demographic characteristics, access to health care, health insurance coverage, as well as income and employment data.  MEPS is cosponsored by the Agency for Healthcare Research and Quality (AHRQ) and the National Center for Health Statistics (NCHS).  For the comparisons reported here we used the MEPS 2005 Full Year Consolidated Data File (HC-097).  This is a public use file available for download from the MEPS web site (http://www.meps.ahrq.gov).  See also the MEPS Factsheet "Computing Standard Errors for MEPS Estimates", also available from the MEPS web site.

The MEPS is not a simple random sample, its design includes stratification, clustering, multiple stages of selection, as well as disproportionate sampling.  The MEPS public use files (such as HC-097) include variables for generating weighted national estimates and for use of the Taylor method for variance estimation.  These variables are: person-level weight (PERWT05F on HC-097); stratum (VARSTR on HC-097); and cluster/psu (VARPSU on HC-097).

### PROC SURVEYMEANS

In this section we will first review the basic syntax of SURVEYMEANS.  We will then generate some estimates and standard deviations and errors using the 2005 MEPS full year file (HC-097).  Note that we used SAS 9.1.3 for the comparisons reported here.

### SYNTAX

The following illustrates the relevant syntax and statements for PROC SURVEYMEANS:

```
(1)    PROC SURVEYMEANS DATA= HC097 SUM STD MEAN STDERR;
             STRATA VARSTR / LIST;
             CLUSTER VARPSU;
             WEIGHT PERWT05F;
             VAR TOTEXP05;
       RUN;
```

The options selected on the PROC SURVEYMEANS statement are SUM, STD (standard deviation of the sum), MEAN, and STDERR (standard error of the mean).  The STRATA statement (you can also use "STRATUM" when you have one variable and you're fussy about your Latin) is where you list the strata variables if you have a stratified sample.  The LIST option causes SAS to generate a *Stratum Information* table which lists the variables and sampling rates for each stratum, as well as the number of records and clusters for each stratum and analysis variable.  This can be very handy information.  For example, for the MEPS HC-097 file the Stratum

Information table shows that each stratum has at least 2 PSUs (clusters).  The head of this table shows the following information.

```
(2)    Number of Strata                                      203
       Number of Clusters                                    452
       Number of Observations                              33961
       Number of Observations Used                         32320
       Number of Obs with Nonpositive Weights               1641
       Sum of Weights                                  296185002
```

The 1,641 records with non-positive weights are excluded prior to analysis.  The N value that appears in the LST output will be 32,320 (number of observations used).  For TOTEXP05 the NMISS value will be zero because all records have a non-missing value for this analysis variable.

The CLUSTER statement is where (you guessed it!) you list the cluster variables.  For MEPS (as with all stratified samples), clusters are nested within strata.

This may be obvious but it's worth pointing out that SAS will not tell you the sample design of the data you're working with.  You need to know the design of the survey that produced the data and how those design properties are to be specified in the survey procedures.  The Federal agencies responsible for distributing public use files for data such as the MEPS do an excellent job of providing documentation of survey design and data file properties.

The WEIGHT statement is self-explanatory.  Records with nonpositive or missing weight values are excluded from the analysis.

The VAR statement lists the analysis variables.  Here TOTEXP05 is the 2005 MEPS total expenditure variable.  The HC-097ile is a person-level file with 33,961cords.  All records have a non-missing value for TOTEXP05  This is consistent with the situation described in (1a): all survey respondents are in the analytic population and all have a non-missing value for the analysis variable.

When we run (1) SAS generates the following output statistics (rounded values are given).  The full output listing is given in the Appendix.

```
(3)    SUM:      1,023,763,462,950
       STD:         33,630,954,198
       MEAN:                 3,457
       STDERR:                  90
```

If we run (2) but omit the STRATA and CLUSTER statements, SAS generates the same SUM and MEAN, but the standard deviation and standard error are those given in the parentheses in (4).

```
(4)    SUM:      1,023,763,462,950
       STD:         33,630,954,198 (23,916,914,371)
       MEAN:                 3,457
       STDERR:                  90 (79)
```

This comparison is consistent with our general experience that variance estimation where a simple random sample is assumed underestimates standard deviations and standard errors for complex samples.

**DOMAIN (Subgroup) Analysis**

In this section we show the use of the SURVEYMEAN's DOMAIN statement.  We will compute the mean and standard error again for the MEPS TOTEXP01 variable, but here we will also look at males and females as subgroups.

```
(5)     PROC SURVEYMEANS DATA= HC097 MEAN STDERR;
            STRATA VARSTR;
            CLUSTER VARPSU;
            WEIGHT PERWT05F;
            VAR TOTEXP05;
            DOMAIN SEX;
        RUN;
```

The DOMAIN statement shows the categorical variable SEX.  The output statistics are shown in (6).

(6)

```
                              Data Summary

                    Number of Strata                            203
                    Number of Clusters                          452
                    Number of Observations                     33961
                    Number of Observations Used                32320
                    Number of Obs with Nonpositive Weights      1641
                    Sum of Weights                         296185002


                                         Statistics

                                                               Std Error
     Variable               N         N Miss          Mean       of Mean
     ------------------------------------------------------------------------
     TOTEXP01           32320              0     3456.500008     90.433709
     ------------------------------------------------------------------------


                             Domain Analysis: SEX

                                                               Std Error
     SEX        Variable        N          N Miss          Mean      of Mean
     ------------------------------------------------------------------------
     FEMALE     TOTEXP05     15251              0     2954.576843    105.916815
     MALE       TOTEXP05     17069              0     3938.646448    135.136809
     ------------------------------------------------------------------------
```

**Missing Values for the Analysis Variable (one PSU in some strata)**

In this section we will first look at an example where records with missing values for the analysis variable create a situation where some strata only have one PSU (cluster) per stratum.  We create this example by modifying the MEPS example from the previous section so we can run an analysis just for expenses paid by Medicare (the MEPS variable TOTMCR05).  Further, for this analysis we are interested in computing the mean only for those persons with a workers' compensation expense, i..e. zero values are set to missing.

```
(7)     PROC SURVEYMEANS DATA= HC097 MEAN STDERR;
              STRATA VARSTR / LIST;
              CLUSTER VARPSU;
              WEIGHT PERWT05F;
              VAR TOTMCR05;
        RUN;
```

Running (7) generates the LOG note in (8) and the output statistics in (9):

```
(8)     Only one cluster in a stratum for variable TOTWCP05.  The variance in
        that stratum is estimated by zero.
```

```
(9)     MEAN:                   5331.97
        STDERR:                  255.53
```

The stratum information table would also show which of the strata have only one cluster (PSU). The NMISS value on the output statistics would be 28,168 and the N would be 4,152. NMISS indicates the number of records with a positive weight that have a missing value for the analysis variable TOTMCR05.

Note that if you run this type of analysis in a software package such as SUDAAN (using the default Taylor linearrization method settings), the generated standard error will be different. The reason for the difference in standard error values is the different assumptions made by SAS and SUDAAN concerning the treatment of missing values (see SAS FAQ # 1813 available online at SAS Technical Support). SAS survey analysis procedures assume that missing values for analysis variables are missing completely at random and delete them before running the analysis. In contrast, the SUDAAN default analysis assumes that these values are not missing completely at random and runs a domain analysis focused only on those records with non-missing values, i.e. the subset (or, domain) of respondents for the analysis variable. Note that this is a distinct situation from what we saw in the previous section. Here (as the stratum information table would show) all strata have at least 2 clusters.

We can simulate this approach in SURVEYMEANS by using the DOMAIN statement to explicitly distinguish missing and non-missing values for the analysis variable. First we need to create a categorical variable INDICATOR (see also the SURVEYMEANS documentation for analyzing data with missing values). Then we can run SURVEYMEANS with the code in (11).

```
(10)    DATA HC097I;
              SET HC097;
              IF TOTMCR05 GT 0
                     THEN INDICATOR = 'NOT_MISSING';
              ELSE INDICATOR = 'MISSING';
        RUN;
```

```
(11)    PROC SURVEYMEANS DATA= HC097I;
              STRATA VARSTR ;
              CLUSTER VARPSU;
              WEIGHT PERWT05F;
              VAR TOTMCR05;
              DOMAIN INDICATOR;
        RUN;
```

This will produce the output in (12).

```
                               Statistics


                                                           Std Error
Variable               N         N Miss           Mean      of Mean
-------------------------------------------------------------------
TOTMCR05            4152          28168     5331.970280    255.529494
-------------------------------------------------------------------



                      Domain Analysis: INDICATOR


                                                                Std Error
INDICATOR       Variable            N          N Miss         Mean      of Mean
-------------------------------------------------------------------------------
MISSING         TOTMCR05            0           28168            .            .
NOT_MISSING     TOTMCR05         4152               0  5331.970280   255.424552
-------------------------------------------------------------------------------

```

The standard error in the first statistics row is the one SURVEYMEANS computes for the mean under the assumption that missing values are missing completely at random.  The domain analysis section shows the standard error produced by the 'indicator' approach.  The value of 255.424552 is exactly what is produced by SUDAAN's PROC DESCRIPT.  Although intuitively the numeric difference between the two values is not great, it is important to understand that (i) the default assumptions which underlie the computation of standard errors in cases where there are cluster-level missing values, and (ii) that SURVEYMEANS, through the use of the DOMAIN statement, offers a way to incorporate into the analysis the non-default assumption that missing values are not missing completely at random.

The documentation for SURVEYMEANS lists the various keywords to be used for requesting additional statistical output, e.g. confidence limits for the mean or sum; coefficient of variation for the mean or sum.  For example, the RATIO statement requests ratio analysis for means or proportions of analysis variables.  These variables must also be listed on the VAR statement.


**PROC SURVEYFREQ**

If you are familiar with PROC FREQ then you have a clear jumpstart on using SURVEYFREQ.  In addition to the STRATA and CLUSTER statements, which have the same function we discussed for SURVEYMEANBS, the TABLES, WEIGHT and FORMAT  statements operate in a similar fashion to PROC FREQ.  Example (13) shows a typical two-way frequency.

```
(13)    PROC SURVEYFREQ DATA= HC097;
            TABLES HISPANX*INSCOV05;
            STRATA VARSTR;
            CLUSTER VARPSU;
            WEIGHT PERWT05F;
            FORMAT HISPANX HISPF. INSCOV05 INSF. ;
        RUN;
```

Here the variables on the TABLES statement are HISPANX (a demographic variable indicating self-reported race/ethnicity), and INSCOV05 (a variable indicating insurance status in 2005).  The output of (13) is given in (14).

(14)

```
                            The SURVEYFREQ Procedure

                                  Data Summary

                    Number of Strata                           203
                    Number of Clusters                         452
                    Number of Observations                   33961
                    Number of Observations Used              32320
                    Number of Obs with Nonpositive Weights    1641
                    Sum of Weights                        296185002


                            Table of HISPANX by INSCOV05

                                    Weighted    Std Dev of              Std Err of
        HISPANX        INSCOV05    Frequency    Frequency    Wgt Freq    Percent      Percent
        ------------------------------------------------------------------------------------
        HISPANIC    ANY PRIVATE        2887     18908300     1242584     6.3839       0.3674
                    PUBLIC ONLY        3512     12877938      733805     4.3479       0.2273
                      UNINSURED        2591     11790097      825434     3.9807       0.2465

                          Total        8990     43576335     2417073    14.7125       0.6899
        ------------------------------------------------------------------------------------
          BLACK     ANY PRIVATE       15423    186356680     4857520    62.9190       0.7401
                    PUBLIC ONLY        5401     41413223     1445538    13.9822       0.3970
                      UNINSURED        2506     24838765      932800     8.3862       0.2512

                          Total       23330    252608667     6097607    85.2875       0.6899
        ------------------------------------------------------------------------------------
          Total     ANY PRIVATE       18310    205264979     5264787    69.3030       0.5957
                    PUBLIC ONLY        8913     54291161     1755404    18.3302       0.4523
                      UNINSURED        5097     36628862     1371767    12.3669       0.3297

                          Total       32320    296185002     7072820   100.000
        ------------------------------------------------------------------------------------
```

One advantage of the default SURVEYFREQ output is that the table shows both unweighted and weighted frequencies.  This was not possible with PROC FREQ, where you had to run the PROC with and without the WEIGHT statement to generate both frequencies.

There are a number of TABLES statement options for selecting different output.  For example, if you want to generate confidence limits for weighted frequencies, you can use the CLWT option.  In (15) below we illustrate the ROW option.  The ouput is shown in (16).

(15)    **PROC SURVEYFREQ DATA= HC097;**
            **TABLES HISPANX*INSCOV05 / ROW;**
            **STRATA VARSTR;**
            **CLUSTER VARPSU;**
            **WEIGHT PERWT05F;**
            **FORMAT HISPANX HISPF. INSCOV05 INSF. ;**
        **RUN;**

7

(16)

```
                              The SURVEYFREQ Procedure

                                   Data Summary

                  Number of Strata                              203
                  Number of Clusters                            452
                  Number of Observations                      33961
                  Number of Observations Used                 32320
                  Number of Obs with Nonpositive Weights       1641
                  Sum of Weights                          296185002


                             Table of HISPANX by INSCOV05

                        Weighted    Std Dev of                 Std Err of      Row   Std Err of
HISPANX      INSCOV05    Frequency   Frequency    Wgt Freq   Percent    Percent  Percent  Row Percent
--------------------------------------------------------------------------------------------------
HISPANIC   ANY PRIVATE      2887     18908300     1242584    6.3839     0.3674   43.3912      1.3001
           PUBLIC ONLY      3512     12877938      733805    4.3479     0.2273   29.5526      1.1133
             UNINSURED      2591     11790097      825434    3.9807     0.2465   27.0562      0.8611

                Total       8990     43576335     2417073   14.7125     0.6899  100.000
--------------------------------------------------------------------------------------------------
  BLACK    ANY PRIVATE     15423    186356680     4857520   62.9190     0.7401   73.7729      0.5635
           PUBLIC ONLY      5401     41413223     1445538   13.9822     0.3970   16.3942      0.4533
             UNINSURED      2506     24838765      932800    8.3862     0.2512    9.8329      0.2952

                Total      23330    252608667     6097607   85.2875     0.6899  100.000
 ---------------------------------------------------------------------------------------------------
  Total    ANY PRIVATE     18310    205264979     5264787   69.3030     0.5957
           PUBLIC ONLY      8913     54291161     1755404   18.3302     0.4523
             UNINSURED      5097     36628862     1371767   12.3669     0.3297

                Total      32320    296185002     7072820  100.000
 ---------------------------------------------------------------------------------------------------
```

In (16), in addition to the frequencies shown in (14), the use of the ROW option causes SAS to output within-group percent estimates in the "Row Percent" column (the SEs for these estimates are shown in the "Std Err of Row Percent" column). For example, the first row of the table in (16), the "Percent" column shows that about 6.3% of the total population is Hispanic with Any Private insurance. The "Row Percent" column shows that about 43.3% of the Hispanic population has Any Private insurance.

The documentation for SURVEYFREQ shows additional options for requesting statistical tests (e.g. chi-square tests) and output (e.g. variances for weighted frequencies). Options for suppressing output (e.g. NOPERCENT, NOPRINT) are also listed.

## PROC SURVEYLOGISTIC

In this section we will first review the basic syntax of PROC SURVEYLOGISTIC. We will then discuss selected options and statements.

The STRATA and CLUSTER statements should be familiar by now, as are the WEIGHT and FORMAT statements.

In the model below, the dependent variable is DENTAL_VISIT a binary variable indicating a visit to a dentist's office in 2005 (1= YES, 2 = NO). This variable is based on the MEPS variable DVTOT05, a count of number

of dental visits in 2005.  The independent variables are: POVCAT05, a categorical variable with 5 values indicating poverty status (1= POOR – 5= HIGH INCOME); and INSCOV05, a categorical variable with three values indicating insurance-coverage status (1= ANY PRIVATE, 2= PUBLIC ONLY, 3= UNINSURED). INSCOV05 was used in the SURVEYFREQ examples above.

```
(17)    PROC SURVEYLOGISTIC DATA=MEPS_H97;
            STRATA VARSTR;
            CLUSTER VARPSU;
            WEIGHT PERWT05F;
            CLASS POVCAT05 INSCOV05 (REF='1 ANY PRIVATE')
                  / PARAM=REF ORDER=INTERNAL;
            MODEL DENTAL_VISIT (EVENT='1')= POVCAT05 INSCOV05/VADJUST=NONE;
            CONTRAST "Test poor, Near poor, and low income versus
                  Middle income" POVCAT05 0  0.5  0.5  -1  0 ;
            FORMAT  INSCOV05 INSCV05F. POVCAT05 POVCAT.;
        RUN;
```

For the CLASS statement, the PARAM=REF (reference) option overrides the default PARAM=EFFECT. When using the reference method, the highest variable value is, by default, used as reference.  This can be changed by the REF= option, as shown in (17).  Here default reference for INSCOV05 would be '3 UNINSURED'.  This is changed to '1 ANY PRIVATE' in the example.

Note also that, when the FORMAT statement is used, the SAS default is ORDER=FORMATTED.  Here this is changed to ORDER=INTERNAL, i.e. unformatted.

The MODEL option EVENT= specifies the event category for the binary variable DENTAL_VISIT.  The specification EVENT='1' specifies '1' as the event category value.  By default, SURVEYLOGISTIC models the lowest value.  The option VADJUST=NONE is used to suppress any variance adjustment associated with the Wald test.  Adjusting by degrees of freedom (VADJUST = DF) is the default.

The CONTRAST statement (which must appear after the MODEL statement) tests the hypothesis that the combined poor/near poor/low income group is different from the middle income group.  The ouput generated by this statement is shown in (18).

(18)

Contrast Test Results

| Contrast | DF | Wald Chi-Square | Pr > ChiSq |
|---|---|---|---|
| Test Negative/poor, Near poor, and low income versus middle income | 1 | 35.517 | <.0001 |

PROC SURVEYLOGISTIC outputs 11 tables by default. Here we will illustrate the ODS statements necessary for outputting the subset of these you may be interested in. One prerequisite to using the ODS statement is knowing the names of the output tables (e.g. ODDSRATIOS and PARAMETERESTIMATES in (19)).  Using the ODS TRACE ON/OFF statements before and after the program code for the relevant procedure will output these table names to the LOG.  Table names are also listed with the procedure documentation.

By adding the following ODS statement to the PROC SURVEYLOGISTIC above, two permanent SAS data sets, Dent_Care_Odds and Denta_Care_Beta, are output to the OUT directory. ODDSRATIOS and PARAMETERESTIMATES are table names defined by SAS.

```
(19)   ODS OUTPUT ODDSRATIOS   = OUT.DENT_CARE_ODDS
               PARAMETERESTIMATES = OUT.DENT_CARE_BETA;
```

The output below is from the output datasets DENT_VISIT_ODDS and DENT_VISIT_ODDS.

(20)    DENT_VISIT_BETA

| Variable | ClassVal0 | DF | Estimate | StdErr | WaldChiSq | ProbChiSq |
|----------|-----------|----|----------|--------|-----------|-----------|
| Intercept | | 1 | 0.367 | 0.035 | 111.120 | 5.57E-26 |
| | 1 | | | | | |
| POVCAT05 | POOR/NEGATIVE | 1 | -0.870 | 0.064 | 183.497 | 8.36E-42 |
| POVCAT05 | 2 NEAR POOR | 1 | -0.966 | 0.090 | 114.868 | 8.41E-27 |
| POVCAT05 | 3 LOW INCOME | 1 | -0.839 | 0.059 | 200.187 | 1.90E-45 |
| POVCAT05 | 4 MIDDLE INCOME | 1 | -0.520 | 0.046 | 128.661 | 8.04E-30 |
| INSCOV05 | 2 PUBLIC ONLY | 1 | -0.429 | 0.054 | 63.944 | 1.28E-15 |
| INSCOV05 | 3 UNINSURED | 1 | -1.356 | 0.067 | 411.707 | 1.56E-91 |

(21)    DENT_CARE_ODDS

| Effect | OddsRatioEst | LowerCL | UpperCL |
|--------|--------------|---------|---------|
| POVCAT05 1 POOR/NEGATIVE vs 5 HIGH INCOME | 0.419 | 0.369 | 0.475 |
| POVCAT05 2 NEAR POOR     vs 5 HIGH INCOME | 0.380 | 0.319 | 0.454 |
| POVCAT05 3 LOW INCOME    vs 5 HIGH INCOME | 0.432 | 0.385 | 0.486 |
| POVCAT05 4 MIDDLE INCOME vs 5 HIGH INCOME | 0.594 | 0.543 | 0.650 |
| INSCOV05 2 PUBLIC ONLY vs 1 ANY PRIVATE | 0.651 | 0.586 | 0.723 |

As discussed above with PROC SURVEYMEANS, the DOMAIN statement should be used for subpopulation, or domain analysis. The formation of these domains may be unrelated to the sample design. Therefore, the sample sizes for the domains are random variables. Domain incorporates this variability into the variance estimation.

The DOMAIN statement is different than a BY statement. The BY statement treats the sample sizes as fixed in each subpopulation, subsets the input data to each subpopulation group, and perform analysis within each BY group separately.  If the subpopulation do not include all strata and clusters (PSUs), the BY statement uses fewer degree of freedom for significance testing.

The code below would analyze how dental visits are associated with family income for each insurance category (DOMAIN INSCOV05).  The output for this is not shown.

10

```
(22)   PROC SURVEYLOGISTIC DATA=MEPS_H97;
           STRATA VARSTR;
           CLUSTER VARPSU;
           WEIGHT PERWT05F;
           DOMAIN INSCOV05;
           CLASS POVCAT05 INSCOV05 (REF='1') /PARAM=REF ORDER=INTERNAL;
           MODEL DENTAL_VISIT (EVENT= '1') = POVCAT05/VADJUST=NONE;
       RUN;
```

## PROC SURVEYREG

The SURVEYREG procedure performs regression analysis for sample survey data. The procedure fits linear models for survey data and computes regression coefficients and their variance-covariance matrix. The procedure also provides significance tests for the model effects and for any specified estimable linear functions of the model parameters. Using the regression model, the procedure can compute predicted values for the sample survey data.

The example below uses the MEPS data to illustrate how total dental expenditures are associated with family income and insurance coverage.  Variable DVTEXP05 is an expenditure variable indicating the amount spent per person in 2005.   The ODS statement requests two output datasets: D_EXP_TEST (from the EFFECTS table) and D_EXP_COEFFICIENTS (from the PARAMETERESTIMATES table).  This output is shown in (24) and (25), respectively.

```
(23)   PROC SURVEYREG DATA=MEPS_H97;
           STRATA VARSTR;
           CLUSTER VARPSU;
           WEIGHT PERWT05F;
           CLASS POVCAT05 INSCOV05;
           MODEL  DVTEXP05 = POVCAT05 INSCOV05/SOLUTION;
           CONTRAST "Test low income versus middle income"
                    POVCAT05 0 0 1 -1 0;
           ODS OUTPUT PARAMETERESTIMATES = out.d_exp_coefficient
                      EFFECTS = Out.D_Exp_test CONTRASTS = Econtrast;
       RUN;
```

As we have seen in previous examples, the poverty-status (POVCAT05) and insurance-coverage (INSCOV05) are used.  As we saw with the SURVEYLOGISTIC procedure, optional CLASS statement requests that they be used as classification variables.  The MODEL statement describes the requested linear model.  The SOLUTION option requests the regression coefficient estimates.  The CONTRAST statement requests hypothesis tests for linear combinations of the regression parameters (see (25) below for the output generated by this statement).

(24)    D_EXP_TEST

| Effect | NumDF | DenDF | FValue | ProbF |
|--------|-------|-------|--------|-------|
| Model | 6 | 249 | 97.0779 | 0 |
| Intercept | 1 | 249 | 834.929 | 0 |
| POVCAT05 | 4 | 249 | 13.2345 | 8.62E-10 |
| INSCOV05 | 2 | 249 | 100.651 | 0 |

11

(25)     D_EXP_COEFFICIENT

| Parameter | Estimate | StdErr | DenDF | tValue | Probt |
|---|---|---|---|---|---|
| Intercept | 166.508 | 14.965 | 249 | 11.126 | 0 |
| POVCAT05 1 POOR/NEGATIVE | -124.285 | 18.569 | 249 | -6.693 | 1.44E-10 |
| POVCAT05 2 NEAR POOR | -124.623 | 27.873 | 249 | -4.471 | 1.18E-05 |
| POVCAT05 3 LOW INCOME | -126.037 | 20.286 | 249 | -6.213 | 2.17E-09 |
| POVCAT05 4 MIDDLE INCOME | -67.282 | 15.272 | 249 | -4.406 | 1.57E-05 |
| POVCAT05 5 HIGH INCOME | 0.000 | 0.000 | 249 | . | . |
| INSCOV05 1 ANY PRIVATE | 177.040 | 13.079 | 249 | 13.536 | 0 |
| INSCOV05 2 PUBLIC ONLY | 83.982 | 12.444 | 249 | 6.749 | 1.04E-10 |
| INSCOV05 3 UNINSURED | 0.000 | 0.000 | 249 | . | . |

(26)     CONTRAST Statement Output

| ContrastLabel | NumDF | DenDF | FValue | ProbF |
|---|---|---|---|---|
| Test low income versus middle income | 1 | 249 | 11.8716 | 0.000669 |

## Summary

The SAS survey procedures are an important addition to the SAS statistical PROCS.  In addition to providing for variance estimation which incorporates the properties of complex surveys, they have the added benefit of allowing convenient output as SAS data sets.

As the paper by Mukhopadhyay *et al* shows, enhancements for SAS 9.2 will only add to the usefulness of these procedures for analysts working with data from complex surveys.

## References

"Calculating Nationwide Inpatient Sample Variance," available at: http://www.hcup-us.ahrq.gov/db/nation/nis/reports/NISVariance2001_Revised%20031904.pdf.

"Computing Standard Errors for MEPS Estimates," available at: http://www.meps.ahrq.gov/FactSheets/FS_StandardErrors.HTM.

"Overview of Survey Data Analysis Using SAS Software".  Tony An.  Proceedings of NESUG 16 (2003).

"What assumptions are made about the treatment of missing values (nonresponse)? How are missing values handled," available at http://support.sas.com/faq/018/FAQ01813.html.

An, Anthony B. (2004).  "Performing logistic regression on survey data with the new SURVEYLOGISTIC procedure."  *Proceedings of the Twenty-Seven Annual SAS Users Group International Conference*, Paper 258.

Chen, X. (2006). "Survey logistic regression: some SAS®-SUDAAN comparisons." *Proceedings of NESUG 19,* Paper po12.

Chen, X. and Gorrell, P. (2004). "Variance estimation with complex surveys: some SAS®-SUDAAN comparisons." *Proceedings of NESUG 17,* Paper an02.

Binder, D.A (1983). "On the variances of asymptotically normal estimators from complex surveys," *International Statistical Review*, 51, 279-292.

Heeringa, S., and Liu, J. (1997), Complex sample design effects and inference for mental health survey data, *International Journal of Methods in Psychiatric Research, 7*, Whurr Publishers Ltd. – Pages 221 – 230

McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*, London: Chapman Hall.

Mukhopadhyay, P.K., An, A.B., Tobias, R.D., and Watts, D.L. (2008). "Try, Try Again: Replication-Based Variance Estimation Methods for Survey Data Analysis in SAS 9.2. *Proceedings of SAS Global Forum 2008*, Paper 367-2008.

Morel, G. (1989). "Logistic Regression under Complex Survey Designs," *Survey Methodology*, 15, 203–223.

Rust, K., (1985), Variance Estimation for Complex Estimators in Sample Surveys, *Journal of Official Statistics, 1 (4)*, Statistics Sweden Publishing Service. Pages 381 –397.

## Acknowledgements

## Contact Information

Xiuhua Chen
Paul Gorrell
Social & Scientific Systems, Inc.
8757 Georgia Avenue
Silver Spring, MD 20910
xchen@s-3.com
pgorrell@s-3.com

# Appendix

## PROC SURVEYMEANS Output for the Program Code in Example (1)

```
                              TheSURVEYMEANS Procedure

                                   Data Summary

                 Number of Strata                             203
                 Number of Clusters                           452
                 Number of Observations                     33961
                 Number of Observations Used                32320
                 Number of Obs with Nonpositive Weights      1641
                 Sum of Weights                          296185002


                               Stratum Information
```

|              | VARIANCE ESTIMATION |       |          |                          |     |          |
|--------------|---------------------|-------|----------|--------------------------|-----|----------|
| Stratum Index | STRATUM – 2005      | N Obs | Variable | Label                    | N   | Clusters |
| 1            | 1                   | 127   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 122 | 3        |
| 2            | 2                   | 118   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 112 | 2        |
| 3            | 3                   | 110   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 104 | 2        |
| 4            | 4                   | 183   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 172 | 3        |
| 5            | 5                   | 434   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 414 | 2        |
| 6            | 6                   | 198   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 187 | 2        |
| 7            | 7                   | 123   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 123 | 2        |
| 8            | 8                   | 175   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 166 | 2        |
| 9            | 9                   | 106   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 104 | 2        |
| 10           | 10                  | 87    | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 84  | 3        |
| 11           | 11                  | 178   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 169 | 2        |
| 12           | 12                  | 212   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 196 | 2        |
| 13           | 13                  | 133   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 123 | 2        |
| 14           | 14                  | 332   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 309 | 2        |
| 15           | 15                  | 111   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 110 | 2        |
| 16           | 16                  | 88    | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 85  | 2        |
| 17           | 17                  | 210   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 196 | 3        |
| 18           | 18                  | 107   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 107 | 2        |
| 19           | 19                  | 163   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 160 | 2        |
| 20           | 20                  | 182   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 173 | 3        |
| 21           | 21                  | 81    | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 77  | 3        |
| 22           | 22                  | 100   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 94  | 3        |
| 23           | 23                  | 196   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 189 | 2        |
| 24           | 24                  | 71    | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 70  | 2        |
| 25           | 25                  | 222   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 213 | 2        |
| 26           | 26                  | 243   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 231 | 3        |
| 27           | 27                  | 103   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 100 | 3        |
| 28           | 28                  | 116   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 112 | 2        |
| 29           | 29                  | 105   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 99  | 2        |
| 30           | 30                  | 204   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 193 | 2        |
| 31           | 31                  | 141   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 137 | 2        |
| 32           | 32                  | 165   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 157 | 2        |
| 33           | 33                  | 236   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 228 | 2        |
| 34           | 34                  | 201   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 193 | 2        |
| 35           | 35                  | 161   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 151 | 2        |
| 36           | 36                  | 156   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 150 | 2        |
| 37           | 37                  | 87    | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 84  | 3        |
| 38           | 38                  | 83    | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 82  | 2        |
| 39           | 39                  | 112   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 98  | 3        |
| 40           | 40                  | 91    | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 84  | 2        |
| 41           | 41                  | 203   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 185 | 2        |
| 42           | 42                  | 129   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 118 | 3        |
| 43           | 43                  | 184   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 174 | 2        |
| 44           | 44                  | 162   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 158 | 2        |
| 45           | 45                  | 208   | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 198 | 3        |

```
46      46      606    TOTEXP05    TOTAL HEALTH CARE EXP 05    575    2
47      47       37    TOTEXP05    TOTAL HEALTH CARE EXP 05     34    2
48      48      109    TOTEXP05    TOTAL HEALTH CARE EXP 05    102    2
49      49      196    TOTEXP05    TOTAL HEALTH CARE EXP 05    183    2
50      50      106    TOTEXP05    TOTAL HEALTH CARE EXP 05    101    2
51      51      171    TOTEXP05    TOTAL HEALTH CARE EXP 05    165    2
52      52      249    TOTEXP05    TOTAL HEALTH CARE EXP 05    238    2
53      53       62    TOTEXP05    TOTAL HEALTH CARE EXP 05     58    2
54      54       64    TOTEXP05    TOTAL HEALTH CARE EXP 05     61    3
55      55      274    TOTEXP05    TOTAL HEALTH CARE EXP 05    258    2
56      56      115    TOTEXP05    TOTAL HEALTH CARE EXP 05    112    2
57      57      115    TOTEXP05    TOTAL HEALTH CARE EXP 05    106    2
58      58      116    TOTEXP05    TOTAL HEALTH CARE EXP 05    111    3
59      59      202    TOTEXP05    TOTAL HEALTH CARE EXP 05    186    2
60      60      103    TOTEXP05    TOTAL HEALTH CARE EXP 05    102    2
61      61      108    TOTEXP05    TOTAL HEALTH CARE EXP 05     97    2
62      62      345    TOTEXP05    TOTAL HEALTH CARE EXP 05    323    2
63      63      212    TOTEXP05    TOTAL HEALTH CARE EXP 05    208    2
64      64      159    TOTEXP05    TOTAL HEALTH CARE EXP 05    149    2
65      65      152    TOTEXP05    TOTAL HEALTH CARE EXP 05    149    2
66      66       96    TOTEXP05    TOTAL HEALTH CARE EXP 05     94    2
67      67      103    TOTEXP05    TOTAL HEALTH CARE EXP 05     97    2
68      68      150    TOTEXP05    TOTAL HEALTH CARE EXP 05    145    2
69      69      181    TOTEXP05    TOTAL HEALTH CARE EXP 05    175    3
70      70      117    TOTEXP05    TOTAL HEALTH CARE EXP 05    110    2
71      71      152    TOTEXP05    TOTAL HEALTH CARE EXP 05    144    2
72      72      147    TOTEXP05    TOTAL HEALTH CARE EXP 05    138    3
73      73      507    TOTEXP05    TOTAL HEALTH CARE EXP 05    478    2
74      74       84    TOTEXP05    TOTAL HEALTH CARE EXP 05     82    2
75      75      139    TOTEXP05    TOTAL HEALTH CARE EXP 05    128    2
76      76      154    TOTEXP05    TOTAL HEALTH CARE EXP 05    131    2
77      77      183    TOTEXP05    TOTAL HEALTH CARE EXP 05    178    2
78      78      157    TOTEXP05    TOTAL HEALTH CARE EXP 05    151    2
79      79       99    TOTEXP05    TOTAL HEALTH CARE EXP 05     92    2
80      80      446    TOTEXP05    TOTAL HEALTH CARE EXP 05    429    2
81      81      149    TOTEXP05    TOTAL HEALTH CARE EXP 05    147    2
82      82      197    TOTEXP05    TOTAL HEALTH CARE EXP 05    188    3
83      83      148    TOTEXP05    TOTAL HEALTH CARE EXP 05    140    2
84      84      130    TOTEXP05    TOTAL HEALTH CARE EXP 05    122    2
85      85      216    TOTEXP05    TOTAL HEALTH CARE EXP 05    205    3
86      86      126    TOTEXP05    TOTAL HEALTH CARE EXP 05    123    3
87      87      152    TOTEXP05    TOTAL HEALTH CARE EXP 05    148    2
88      88      119    TOTEXP05    TOTAL HEALTH CARE EXP 05    116    2
89      89       99    TOTEXP05    TOTAL HEALTH CARE EXP 05     98    2
90      90      139    TOTEXP05    TOTAL HEALTH CARE EXP 05    134    2
91      91      256    TOTEXP05    TOTAL HEALTH CARE EXP 05    248    3
92      92      139    TOTEXP05    TOTAL HEALTH CARE EXP 05    133    2
93      93      130    TOTEXP05    TOTAL HEALTH CARE EXP 05    129    2
94      94      136    TOTEXP05    TOTAL HEALTH CARE EXP 05    128    2
95      95      633    TOTEXP05    TOTAL HEALTH CARE EXP 05    602    2
96      96      137    TOTEXP05    TOTAL HEALTH CARE EXP 05    135    2
97      97      126    TOTEXP05    TOTAL HEALTH CARE EXP 05    125    2
98      98      134    TOTEXP05    TOTAL HEALTH CARE EXP 05    124    3
99      99       96    TOTEXP05    TOTAL HEALTH CARE EXP 05     92    2
100    100      123    TOTEXP05    TOTAL HEALTH CARE EXP 05    114    3
101    101      154    TOTEXP05    TOTAL HEALTH CARE EXP 05    147    2
102    102      175    TOTEXP05    TOTAL HEALTH CARE EXP 05    171    2
103    103       66    TOTEXP05    TOTAL HEALTH CARE EXP 05     59    2
104    104      178    TOTEXP05    TOTAL HEALTH CARE EXP 05    163    2
105    105       76    TOTEXP05    TOTAL HEALTH CARE EXP 05     71    2
106    106       84    TOTEXP05    TOTAL HEALTH CARE EXP 05     79    2
107    107      473    TOTEXP05    TOTAL HEALTH CARE EXP 05    452    2
108    108      144    TOTEXP05    TOTAL HEALTH CARE EXP 05    131    2
109    109       97    TOTEXP05    TOTAL HEALTH CARE EXP 05     95    3
110    110      132    TOTEXP05    TOTAL HEALTH CARE EXP 05    130    2
111    111      139    TOTEXP05    TOTAL HEALTH CARE EXP 05    132    2
112    112      100    TOTEXP05    TOTAL HEALTH CARE EXP 05    100    2
113    113      225    TOTEXP05    TOTAL HEALTH CARE EXP 05    209    2
114    114      227    TOTEXP05    TOTAL HEALTH CARE EXP 05    209    2
115    115      145    TOTEXP05    TOTAL HEALTH CARE EXP 05    137    2
```

| 116 | 116 | 198 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 184 | 3 |
| 117 | 117 | 182 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 170 | 2 |
| 118 | 118 | 195 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 185 | 2 |
| 119 | 119 | 123 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 115 | 2 |
| 120 | 120 | 113 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 105 | 2 |
| 121 | 121 | 150 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 140 | 2 |
| 122 | 122 | 146 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 142 | 3 |
| 123 | 123 | 143 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 139 | 2 |
| 124 | 124 | 111 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 106 | 2 |
| 125 | 125 | 121 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 114 | 2 |
| 126 | 126 | 174 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 154 | 2 |
| 127 | 127 | 135 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 128 | 3 |
| 128 | 128 | 182 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 176 | 2 |
| 129 | 129 | 120 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 110 | 2 |
| 130 | 130 | 90 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 87 | 2 |
| 131 | 131 | 147 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 140 | 2 |
| 132 | 132 | 68 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 67 | 2 |
| 133 | 133 | 140 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 123 | 3 |
| 134 | 134 | 100 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 99 | 3 |
| 135 | 135 | 380 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 363 | 2 |
| 136 | 136 | 198 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 193 | 2 |
| 137 | 137 | 168 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 165 | 2 |
| 138 | 138 | 72 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 63 | 3 |
| 139 | 139 | 116 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 112 | 2 |
| 140 | 140 | 422 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 408 | 3 |
| 141 | 141 | 182 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 172 | 3 |
| 142 | 142 | 75 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 73 | 2 |
| 143 | 143 | 209 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 193 | 2 |
| 144 | 144 | 97 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 93 | 2 |
| 145 | 145 | 218 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 200 | 2 |
| 146 | 146 | 177 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 170 | 2 |
| 147 | 147 | 199 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 193 | 2 |
| 148 | 148 | 97 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 91 | 2 |
| 149 | 149 | 142 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 130 | 2 |
| 150 | 150 | 176 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 164 | 2 |
| 151 | 151 | 514 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 490 | 2 |
| 152 | 152 | 133 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 127 | 2 |
| 153 | 153 | 138 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 137 | 2 |
| 154 | 154 | 127 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 121 | 3 |
| 155 | 155 | 224 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 217 | 3 |
| 156 | 156 | 119 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 115 | 2 |
| 157 | 157 | 192 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 179 | 2 |
| 158 | 158 | 100 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 91 | 3 |
| 159 | 159 | 74 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 74 | 2 |
| 160 | 160 | 198 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 190 | 3 |
| 161 | 161 | 95 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 94 | 2 |
| 162 | 162 | 211 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 207 | 2 |
| 163 | 163 | 116 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 113 | 3 |
| 164 | 164 | 130 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 128 | 2 |
| 165 | 165 | 131 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 122 | 2 |
| 166 | 166 | 58 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 53 | 2 |
| 167 | 167 | 248 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 229 | 2 |
| 168 | 168 | 146 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 145 | 2 |
| 169 | 169 | 178 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 166 | 2 |
| 170 | 170 | 129 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 121 | 3 |
| 171 | 171 | 102 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 93 | 2 |
| 172 | 172 | 325 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 301 | 2 |
| 173 | 173 | 88 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 83 | 3 |
| 174 | 174 | 178 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 176 | 2 |
| 175 | 175 | 268 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 262 | 2 |
| 176 | 176 | 143 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 134 | 2 |
| 177 | 177 | 725 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 677 | 2 |
| 178 | 178 | 172 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 161 | 2 |
| 179 | 179 | 281 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 267 | 2 |
| 180 | 180 | 280 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 257 | 3 |
| 181 | 181 | 72 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 70 | 2 |
| 182 | 182 | 219 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 214 | 2 |
| 183 | 183 | 180 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 170 | 2 |
| 184 | 184 | 187 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 182 | 2 |
| 185 | 185 | 160 | TOTEXP05 | TOTAL HEALTH CARE EXP 05 | 157 | 2 |

```
    186          186          126    TOTEXP05     TOTAL HEALTH CARE EXP 05          124         3
    187          187          199    TOTEXP05     TOTAL HEALTH CARE EXP 05          197         3
    188          188          134    TOTEXP05     TOTAL HEALTH CARE EXP 05          133         2
    189          189          208    TOTEXP05     TOTAL HEALTH CARE EXP 05          194         2
    190          190          164    TOTEXP05     TOTAL HEALTH CARE EXP 05          161         3
    191          191          113    TOTEXP05     TOTAL HEALTH CARE EXP 05          110         2
    192          192          202    TOTEXP05     TOTAL HEALTH CARE EXP 05          195         2
    193          193          115    TOTEXP05     TOTAL HEALTH CARE EXP 05          103         2
    194          194           99    TOTEXP05     TOTAL HEALTH CARE EXP 05           90         2
    195          195          110    TOTEXP05     TOTAL HEALTH CARE EXP 05          107         3
    196          196          185    TOTEXP05     TOTAL HEALTH CARE EXP 05          179         3
    197          197          224    TOTEXP05     TOTAL HEALTH CARE EXP 05          216         2
    198          198           60    TOTEXP05     TOTAL HEALTH CARE EXP 05           58         2
    199          199          175    TOTEXP05     TOTAL HEALTH CARE EXP 05          172         2
    200          200          145    TOTEXP05     TOTAL HEALTH CARE EXP 05          140         2
    201          201           80    TOTEXP05     TOTAL HEALTH CARE EXP 05           80         2
    202          202          115    TOTEXP05     TOTAL HEALTH CARE EXP 05          112         3
    203          203          147    TOTEXP05     TOTAL HEALTH CARE EXP 05          142         2
----------------------------------------------------------------------------------------------


                                      Statistics


                                                                Std Error
             Variable    Label                          Mean     of Mean
             ------------------------------------------------------------------
             TOTEXP05    TOTAL HEALTH CARE EXP 05    3456.500008   90.433709
             ------------------------------------------------------------------
```