**Name:**      **Clayton Wells**
             Arkansas Foundation for Medical Care
             401 W Capitol Suite 410
             Little Rock AR 72201
             (501) 375-5700
             cwells@afmc.org

**Tips For Using PROC SQL to Extract Information From Medical Claims Data Sets**

## Introduction

This paper is aimed at beginning SAS users who are writing SAS code used to manipulate and analyze medical claims data sets.  Often, raw claims data sets contain one observation per claim.  The SQL procedure provides a quick and efficient way to manipulate data sets of this type and allow the programmer to extract lots of information with minimal coding.

## Problem

For simplicity, consider a data set that contains four variables – **id_number** (an identification number assigned to each unique person in the data set), **mem_month** (a date variable that represents for each person, the month the person was enrolled in a particular health care program), **provider** (health care provider where the person was enrolled for a given month), and **dob** (date of birth – this will obviously be the same within a given id_number).  This data set, **DATA_1**, is shown in Figure 1.  This data is fairly representative of a raw medical claims data set.  Each observation in the set represents a month of enrollment in a health care plan by a particular individual.  There are several questions that we can answer with this data:

1. How many unique members (ID numbers) are there in the data set?
2. How many months of membership does each unique member have?
3. How many unique providers are in the data set?
4. For each unique provider, how many unique members did the provider have enrolled for at least one month?
5. How many total "member months" did each unique provider have?

The SQL procedure is well suited to answer all of the questions above.  For comparison purposes, an alternative

(non-SQL) method will be demonstrated as well for questions 1 and 2.

## Method 1 – PROC TRANSPOSE

One way to answer questions 1 and 2 above if to transpose the data in Figure 1 on the **mem_month** variable by the **id_number** variable.  The result of this transformation is shown in Figure 1, **DATA_2**.  Since **id_number** 1001 and 1006 contain 6 months of membership, more than the other unique identification numbers in the data set, there are 6 new variables in the resulting data set, labeled COL1 through COL6.  The number of observations in **DATA_2** in Figure 1 is equal to the number of unique members and question 1 is answered.  So we must now answer question 2.  There are many ways to do this.  The only one discussed here is as follows: create 6 new variables corresponding to COL1 through COL6.  Call these new variables **m1** through **m6**.  Assign these variables a value of 1 if its corresponding COL variable is non-missing, 0 otherwise.  Then we can sum **m1** to **m12**.  The resulting sum is the number of membership months and question 2 is now answered.  The resulting data set containing the sum variable is called **DATA_3** in Figure 1.

The code for this method is as follows:

```
proc transpose data=lib.data_1 out=lib.data_2;
var mem_month; by id_number;
run;


data lib.data_3; set lib.data_2;
m1=0; m2=0; m3=0; m4=0; m5=0; m6=0;
if col1 ne . then m1=1;
if col2 ne . then m2=1;
if col3 ne . then m3=1;
if col4 ne . then m4=1;
if col5 ne . then m5=1;
if col6 ne . then m6=1;
sum=m1+m2+m3+m4+m5+m6;
run;
```

This method is not terribly inefficient (for this particular problem), but the SQL required to acquire the same information is definitely more compact.

## Method 2 – SQL

A technique using the count function in SQL requires less code.  The code shown below will create a new data set containing one observation for each unique member and will contain a variable indicating total number of membership months.  The resulting data can be seen in **DATA_4** in Figure 1.  The code needed to create the data set is as follows:

```
proc sql;
create table lib.data 4 as
select unique(id_number),
       count(mem_month) as member months
from lib.data_1 group by id_number;
quit;
```

To answer questions 3, 4, and 5, we can again use the SQL procedure to create a new data set containing one observation per unique provider that contains variable representing the total number or members enrolled at least one month at that particular provider and a variable to represent the total number of member months at the particular provider.  This data can be found in the **PROVIDER** data set in Figure 1.  The code to create this set is as follows:

```
proc sql;
create table lib.provider as
select provider,
       count(id number) as num months,
         count(unique(id number)) as num_recip
from lib.data_1 group by provider;
quit;
```

## Discussion

From inspection we are able to identify the answer to all of the questions listed above (since the raw data set only contains 26 observations).  For question 1, we can see that the raw data set contains 6 unique identification variables. Also note that the data sets **DATA_2** and **DATA_4** contain six observations, as they should.  For the next three questions, which pertain to the providers in the raw data set, we can note first that there are 3 unique providers, A, B, and C. Provider A had four unique members enrolled at one point – members 1001, 1002, 1003, and 1004.  Provider B only had three unique members enrolled at one point – members 1001, 1004, and 1006.  Hence the values of **num_recip** in the **PROVIDER** data set.  The values for the variable **num_months** in **PROVIDER** is obtained by counting the number of times a particular provider appears in the raw data set, since each observation represents a month of membership.  So provider A had 8 months of membership, provider B had 3, and so on.

Note in the problem presented here that the maximum number of member months per unique ID number was 6.  Often this number is 12 or more.  Method 1 requires at least one line

of additional code for each additional member month added to the data set.  The method utilizing PROC SQL, however, is independent of the number of member months.
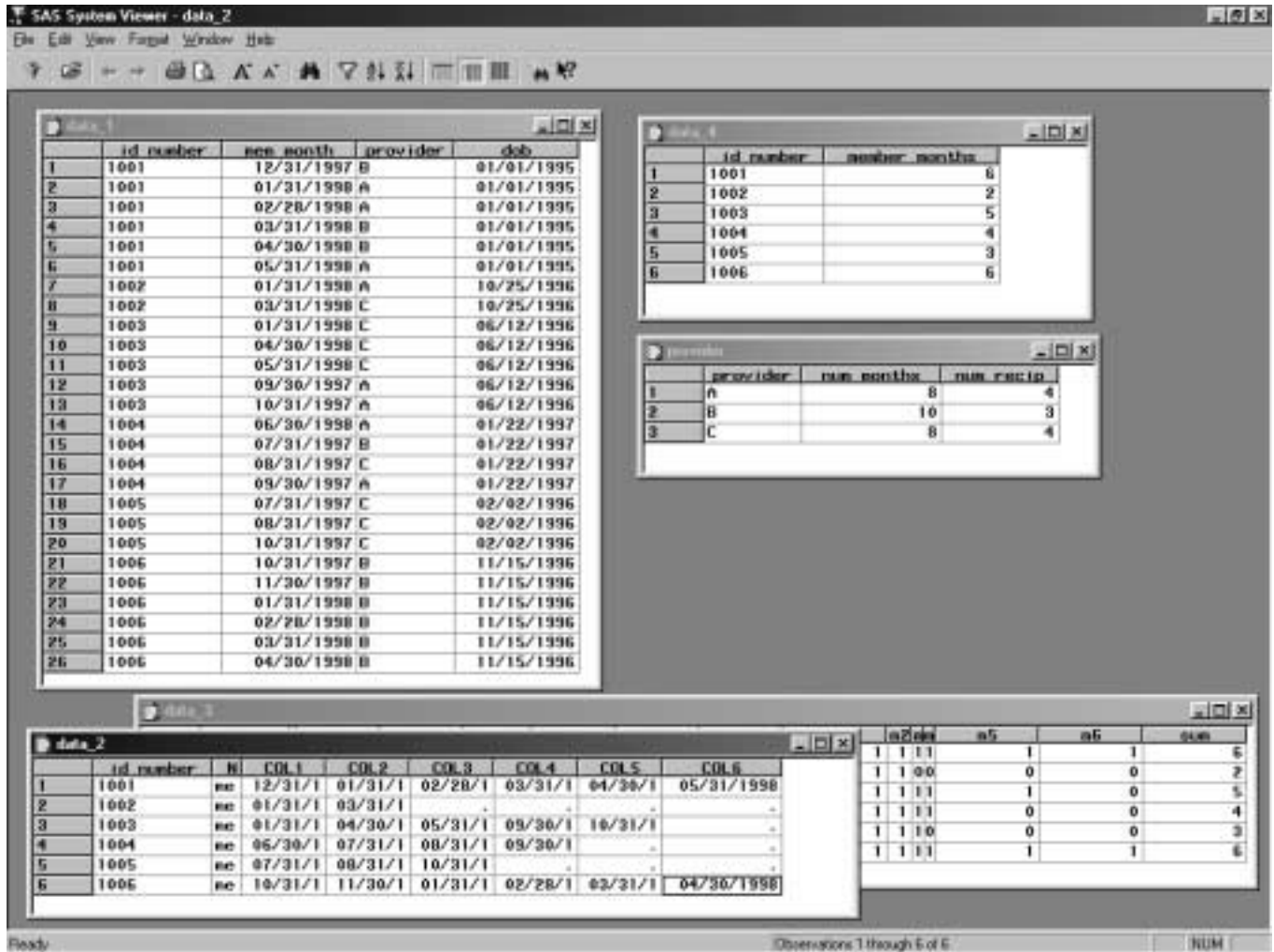


**Figure 1 – Data sets used by the SAS program**