

## Matching SAS® Data Sets with Hash Objects: If at First You Don't Succeed, Match, Match Again

Imelda C. Go, Questar Assessment, Inc.

### ABSTRACT

Two data sets can be matched using a number of techniques available in SAS. If the goal is to match as many records as possible between two data sets, then the matching process may have to be repeated several times between the two data sets by using different matching criteria. Although programmers new to hash objects will need to learn a few things before they can feel confident about using the technique, the resulting ease of use/maintenance of the matching code is well worth the effort. (This is the hash-object version of two prior papers on the same topic: the 2009 SESUG paper focused on PROC SQL and the 2004 SESUG paper focused on the DATA step.)

---

This paper does not provide details of hash object basics, but papers by Paul Dorfman and other authors on the topic are highly recommended. In general, this paper is for programmers who are aware of the pros and cons of the different matching techniques and the things you have to be cautious about when matching (e.g., do the data sets to be matched contain variables with exactly the same name that could be overwritten accidentally?). Although it is not difficult to get started with matching in SAS, it does take some experience to match correctly and without undesirable outcomes using the various techniques available in SAS.

In the example below, there is no multiplicative effect where the number of records increases as a result of matching. An increase in the number of records is fine if that does not conflict with the programmer's intentions. With the use of hash objects, the multiplicative problems/complications produced by one-to-many, many-to-one, and many-to-many matching are bypassed. When hash objects are used for matching, SAS will match the first record in the hash object that provides a match. It is your responsibility to make sure that this is the desired effect. If not, you will have to work on or prepare the data that produces the hash object so that you produce the desired matching results.

There are computer programming tasks that require that the number of records remain the same. For example, there might be a set of student test scores that need to be matched to demographic variables from student records. We do not want Jane Doe's test record that appeared once to appear 10 times after matching because Jane was enrolled in 10 different schools. (This unwanted duplication of records is often due to one-to-many, many-to-one, and many-to-many situations using the DATA step or PROC SQL to match records.) This is the scenario that will be presented in the programming example below.

### PROGRAMMING EXAMPLE

The programming example will attempt to match two data sets, which are:

- A data set named `testscores` that contains tests scores. This data set has `ctr` as its primary key or unique record identifier.
- A data set named `students` that contains demographic information for students. This data set has `ctrf` as its primary key or unique record identifier.

Step 1:

- Determine all the variables needed for the matching criteria.
- Clean up the data for the variables to match on. For example, make them as consistent as possible—remove spaces, capitalize all the letters, etc. This reduces unnecessary variation in the data values, which helps increase the number of matches (e.g., "John" does not match to "JOHN").
- All of these variables must be present in both data sets. The variable names need to be identical. For example, you might want to match on the first three letters of the first name and the first two letters of the last name. Define these substrings in the data set prior to matching.

Step 2:

- Identify the variables to be used per matching round. In general, the more variables used to match on, the more stringent the criterion is and the stringency of the criterion can decrease for each subsequent round. That is, the most stringent criterion can be used for the first round and subsequent rounds would be less stringent.

Therefore, let us suppose that you have prepared the two data sets you want to match and that you also decided on the rounds listed below as your matching strategy. The variables used for matching involve the school ID (`schoolid`); and the student's state ID (`stateid`), last name (`lname`), first name (`fname`), and date of birth (`dob`).

Each of the two data sets will contain all the variables listed in the table below. (Note for Variables Below: Variable names that begin with `m` mean they are the modified version of the original value [e.g., `mlname` is the modified version of variable `lname`--all values have consistent case, all unnecessary spaces removed, etc.]. Variable names that end in a number (`n`) mean the original variable is truncated to the first `n` characters. For example, `mlname3` contains the first 3 characters of the `mlname` variable's value. Variables names that begin with `sw` mean the first and last names were interchanged since that is a common data error with student names. At this point, it is not really that important what the variables below represent. What is important is that you decide which variables you will use when you match your data and what your matching strategy is.)

Round	Variables Used to Match On
1	<code>schoolid stateid mlname mfname dob</code>
2	<code>schoolid stateid mlname mfname</code>
3	<code>schoolid mlname mfname dob</code>
4	<code>schoolid stateid mlname3 mfname1 dob</code>
5	<code>schoolid stateid mlname3 mfname1</code>
6	<code>schoolid mlname3 mfname2 dob</code>
7	<code>schoolid stateid swmlname swmfname</code>
8	<code>schoolid swmlname swmfname dob</code>
9	<code>stateid mlname mfname mmname dob</code>
10	<code>stateid mlname mfname dob</code>
11	<code>stateid mlname dob</code>
12	<code>stateid mfname dob</code>
13	<code>stateid mlname3 mfname2 dob</code>
14	<code>stateid swmlname swmfname dob</code>
15	<code>stateid mlname mfname</code>
16	<code>stateid swmlname swmfname</code>
17	<code>schoolid mlname mfname mmname</code>
18	<code>schoolid mlname mfname</code>

Because of so many matching rounds, a macro will be written to accommodate the rounds. There will be two macro parameters: the matching round (`n`) and the matching variables (`matchvar`). The macro is described as follows:

SAS CODE	EXPLANATION
<code>%macro round (n,matchvar);</code>	There are 2 macro parameters. <code>N</code> is for the matching round number. <code>Matchvar</code> is for the variables used for matching.
<code>%let m=%eval(&amp;n-1);</code>	Macro variable <code>m</code> is just macro variable <code>n</code> minus one.
<pre>data _null_; length str \$150.; str="&amp;matchvar"; strlength=lengthn(str); cstrlength=lengthn(compress(str)); numvars=strlength-cstrlength+1; str="  tranwrd(str," ","','")  "; str=substr(str,1,strlength+numvars*2); lastvar=scan("&amp;matchvar",numvars);  call symput ('str',str); call symput ('lastvar',lastvar); run;</pre>	<p>This DATA _NULL_ step assigns to macro variable <code>str</code> the macro variable <code>matchvar</code> reformatted with punctuation for use in the <code>definekey</code> method of the hash object. For example, if <code>matchvar</code> is the value:</p> <p><code>a b c</code></p> <p>then the final result for the variable <code>str</code> is:</p> <p><code>"a","b","c"</code></p> <p>and the result for the variable <code>lastvar</code> is:</p> <p><code>c</code></p> <p>The <code>lastvar</code> variable is the last variable listed in the <code>str</code> variable value.</p>

SAS CODE	EXPLANATION
<pre>proc sort data=students (keep=&amp;matchvar ctrf) out=subset; by &amp;matchvar; run;</pre>	<p>PROC SORT will sort on the <code>matchvar</code> matching variables. Only the <code>matchvar</code> variables and the <code>ctrf</code> variable will be in the data set <code>subset</code>. Because the goal is to match the <code>ctrf</code> value to the matching <code>testscores</code> record, removing all other variables is a way to save on system resources since only the matching variables and <code>ctrf</code> are needed for the hash object.</p>
<pre>data keys dupkeys;   set subset; by &amp;matchvar;   if first.&amp;lastvar and last.&amp;lastvar then output keys; else output dupkeys; run;</pre>	<p>This DATA step divides data set <code>subset</code> into two. The first data set <code>keys</code> contains all records in <code>subset</code> that have no duplicates on the <code>matchvar</code> variables. None of those with duplicates appear in data set <code>keys</code>; duplicates will all appear in data set <code>dupkeys</code>.</p> <p>In this example, a conscious decision was made to exclude any record from <code>students</code> that has duplicate values on the matching variables.</p> <p>Preparing the data set for creating the hash object is critical to the success of the match. This is where the programmer has to know what he/she is doing and know how the hash object behaves during matching so that the outcome of the match is the desired correct result. After all, there are a number of ways to resolve duplicates. Depending on how you prepare the data, the duplicates may be included or excluded in the matching process.</p>
<pre>data matches; length ctrf \$30.; if _n_=1 then do;   declare hash round (dataset:'keys');   round.definekey(&amp;str);   round.definedata('ctrf');   round.definedone();   call missing (ctrf); end;</pre>	<p>These lines define the hash object named <code>round</code>. It will take the values of the hash object from the dataset called <code>keys</code>. For this hash object, the key consists of the <code>matchvar</code> variables expressed in the <code>str</code> macro variable. The data that the hash object provides is <code>ctrf</code>.</p> <p>You can add more variables to this declaration, but for now the priority is to find the <code>ctrf</code> value for the record in data set <code>students</code> that matches to a particular record in data set <code>testscores</code>.</p> <p>The LENGTH statement gives <code>ctrf</code> a length of 30. The CALL MISSING statement sets the variable <code>ctrf</code> to missing. Since it is a character variable in this case, it is set to blank.</p>
<pre>%if &amp;n&gt;1 %then   %do;     set matches;     if min (of rc1-rc&amp;m) ne 0 then rc&amp;n=round.find();   %end; %else %if &amp;n=1 %then   %do;     set testscores;     rc1=round.find();   %end;</pre>	<p>The rest of the macro specifies the data sets to be processed and the matching performed using the hash object.</p> <p>During the first round (<code>&amp;n=1</code>) of matching (round 1), the match begins with data set <code>testscores</code>.</p> <p>The return code for round 1 is variable <code>rc1</code>. Whenever the return code value is 0, there is a match.</p> <p>For subsequent rounds (<code>&amp;n&gt;1</code> for rounds 2, 3, etc.), the processing will start with data set <code>matches</code>. The condition <code>min(of rc1-rc&amp;m)</code> checks if any of the prior rounds produced a return code of 0. Note that the minimum for the 1<sup>st</sup> through the <i>m</i><sup>th</sup> (<i>n</i>-1<sup>th</sup>) return code is computed. A return code of 0 means there was a match. If there was a match found then the matching will not be attempted for that round.</p>
<pre>if rc&amp;n=0 then match=&amp;n;  run; %mend round;</pre>	<p>The last statement in the macro simply assigns the round to the <code>match</code> variable if the return code for corresponding return code variable (<code>rc&amp;n</code>) is 0.</p>

The macro will be invoked as many times as the number of rounds. The following show the 18 macro calls.

```
%round(n=1,matchvar=schoolid stateid mlname mfname dob);
%round(n=2,matchvar=schoolid stateid mlname mfname);
%round(n=3,matchvar=schoolid mlname mfname dob);
%round(n=4,matchvar=schoolid stateid mlname3 mfname1 dob);
%round(n=5,matchvar=schoolid stateid mlname3 mfname1);
%round(n=6,matchvar=schoolid mlname3 mfname2 dob);
%round(n=7,matchvar=schoolid stateid swmlname swmfname);
%round(n=8,matchvar=schoolid swmlname swmfname dob);
%round(n=9,matchvar=stateid mlname mfname mmname dob);
%round(n=10,matchvar=stateid mlname mfname dob);
%round(n=11,matchvar=stateid mlname dob);
%round(n=12,matchvar=stateid mfname dob);
%round(n=13,matchvar=stateid mlname3 mfname2 dob);
%round(n=14,matchvar=stateid swmlname swmfname dob);
%round(n=15,matchvar=stateid mlname mfname);
%round(n=16,matchvar=stateid swmlname swmfname);
%round(n=17,matchvar=schoolid mlname mfname mmname);
%round(n=18,matchvar=schoolid mlname mfname);
run;
```

After the macro is invoked 18 times for the 18 rounds of matching, we can check to see how many successful matches were made.

```
proc freq data=matches;
tables
rc1*rc2*rc3*rc4*rc5*rc6*rc7*rc8*rc9*rc10*rc11*rc12*rc13*rc14*rc15*rc16*rc17*rc18*match/
list missing;
where ctrf>' ';
run;
```

rc1	rc2	rc3	rc4	rc5	rc6	rc7	rc8	rc9	rc10	rc11	rc12	rc13	rc14	rc15	rc16	rc17	rc18	match	Frequency	Percent	Cumulative Frequency	Cumulative Percent
0	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	1	3475	98.05	3475	98.05
160038	0	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	2	5	0.14	3480	98.19
160038	160038	0	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	3	17	0.48	3497	98.67
160038	160038	160038	0	.	.	.	.	.	.	.	.	.	.	.	.	.	.	4	8	0.23	3505	98.90
160038	160038	160038	160038	160038	160038	0	.	.	.	.	.	.	.	.	.	.	.	7	2	0.06	3507	98.96
160038	160038	160038	160038	160038	160038	160038	0	.	.	.	.	.	.	.	.	.	.	8	1	0.03	3508	98.98
160038	160038	160038	160038	160038	160038	160038	160038	0	.	.	.	.	.	.	.	.	.	9	34	0.96	3542	99.94
160038	160038	160038	160038	160038	160038	160038	160038	160038	0	.	.	.	.	.	.	.	.	10	1	0.03	3543	99.97
160038	160038	160038	160038	160038	160038	160038	160038	160038	160038	160038	160038	160038	160038	160038	160038	160038	0	18	1	0.03	3544	100.00

The above shows us that 3,475 matches were made during round 1, 5 matches for round 2, etc. We are essentially adding the variable `ctrf` to the data set `testscores`. To reduce the resources needed to find the matches, the most important thing is to find the value of `ctrf` that produces the match. After all the rounds are done, the variables of interest in data set `students` can be added to the `matches` data set since we have the variable `ctrf` that identifies which `students` record matches to the `testscores` record.

## CONCLUSION

SAS provides a number of ways to match data. An understanding of how SAS processes data and how things can go wrong during the matching process is essential to anticipating problems that need to be addressed when matching data. Being unaware of how problems can occur could potentially allow such problems to occur undetected, and thus endanger the integrity of the matched data.

## REFERENCES

Go, Imelda. 2004. *Matching SAS® Data Sets: If at First You Don't Succeed, Match, Match Again*. Proceedings CD of the Twelfth Annual Conference of the SouthEast SAS Users Group, Nashville, TN, USA, 12.

Go, Imelda. 2009. *Matching SAS® Data Sets with PROC SQL: If at First You Don't Succeed, Match, Match Again*. Proceedings CD of the 17th Annual Conference of the Southeast SAS Users Group, Birmingham, AL, USA, 17.

## CONTACT INFORMATION

Imelda C. Go

[igo@questarai.com](mailto:igo@questarai.com)

Working remotely from Columbia, SC

## TRADEMARK NOTICE

SAS is a registered trademark or trademark of the SAS Institute Inc. in the USA and other countries. ® indicates USA registration.