

Paper 054-29

Colonoscopy for the SAS[®] Programmer

William C. Murphy
Howard M. Proskin & Associates, Inc., Rochester, NY

ABSTRACT

In order to maintain good program health, the colon should be used often and well in the SAS program. Although the colon is one of the most powerful and useful symbol operators available to the programmer, some of its features are rarely found in typical SAS code. In particular, its prefixing power, along with programming conventions, can produce smaller and more easily understood code. This, in turn, means fewer errors in the program and less typing for the programmer. Not only can you KEEP, DROP, or FORMAT variables quickly with the use of colons but you can readily SAVE or DELETE data sets. The colon prefix ability can be used in DATA steps, as well as in macro programs and PROCs. The colon can be employed in managing temporary variables, formatting similar variables, dealing with unknown numbers of variables, and cleaning up after macro execution.

INTRODUCTION

Colonoscopy is the examination of colons with a colonoscope. In lieu of a colonoscope, the SAS programmer can perform the examination of the colons (:) in his or her programs visually. If the colons are only sparingly used or not exploited to their fullest potential, the programs may not be as healthy as they could be. Colons can play a vital role in keeping code short and data sets clean.

BACKGROUND

As most SAS programmers know, the colon can appear many places in a program. If you label a line of code, the colon suffix indicates the label name. For arrays with a specified lower bound, the colon serves to separate the low bound from the high bound in the array statement. The colon is used as a FORMAT modifier in PUT statements and as an INFORMAT modifier in INPUT statements. And in PROC SQL, the colon prefix is used to create macro variables.

One of the most widely used functions of the colon is to compare strings with prefixes:

```
if variable =: 'A';  
if variable in: ('alpha', 'beta', 'c');
```

In this case, the colon transforms the meaning of '=' to 'begins with' and 'in' to 'begins with any string in'. It has a similar effect on other conditional operators.

The most useful feature of the colon is seldom found in a typical SAS program. The colon when used with variable and data set names can indicate any name starting with the string preceding the colon. This ability of the colon, along with some simple naming standards, enables the programmer to manage temporary variables better, format many variables quicker, determine unknown numbers of variables, clean up macro generated data sets, and shorten the code for a variety of PROCs.

TEMPORARY VARIABLES

When you process data within a DATA step, you often need to perform a variety of calculations. In this process, you may create variables that you do not wish to save in the output data set. For instance you might write

```
data two;  
  set one;  
  pi=3.14;  
  e=2.71;  
  c=2.99;  
  
  {... more statements ...}
```

where the variables pi, e, and c will be used for later calculations in the DATA step. However, in your new data set two, you don't really want to keep these variables. So you add a line to the DATA step:

```
drop pi e c;
```

Now for three temporary variables, this is no big hassle, but what if you write a program with dozens of temporary variables for holding intermediate calculations and constants. The task of keeping track of them and DROP them from the new data set becomes much more difficult. However, if we start all of our temporary variable names with the same distinct prefix, we would now write

```
data two;
  set one;
  tmppi=3.14;
  tmpe=2.71;
  tmpc=2.99;

  {... more statements ...}
```

where we have used the prefix tmp for temporary variables. Now to get rid of these variables in the output data set we write

```
drop tmp::;
```

where the expression tmp: is translated by the SAS system into a list of all variables starting with tmp. Thus we no longer have to keep track of our temporary variables; as long as we start the variable names with the same prefix, we can readily DROP them by the appropriate use of the colon.

This same strategy could also be used in the reverse way. You could prefix all the variables that you intended to KEEP and use the colon suffix to KEEP them in the output data set. The other advantage with this method is that you can go back and edit your data step, adding and removing variables, without having to change the DROP or KEEP statements.

FORMATTING SIMILAR VARIABLES

I often examine my SAS data sets in print-outs and in full-screen views. I have no problem in dealing with unformatted variables. In fact, when I'm doing a comparison with raw data, unformatted is better. That is, it's better in all cases but dates. If there is any chance at all that I understand them, dates must be formatted. I typically see variables like VisitDate, ScreenDate, StartDate, EndDate, etc. and of course I format them:

```
format VisitDate ScreenDate StartDate EndDate mmddyy8.;
```

Unfortunately, these names are used in the data sets provided to me. If I create the data sets myself, I usually prefer names like DateVisit, DateScreen, DateStart, DateEnd because then I can use the power of the colon to format them:

```
format Date: mmddyy8.;
```

As before, I can now readily add or remove date variables from my data set without having to worry about the content of the FORMAT statement. You can expand this idea to any group of variables that you format the same way. Of course, this technique would also apply to INFORMAT statements.

DEALING WITH UNKNOWN NUMBERS OF VARIABLES

Often we receive data with multiple observations per subject. For our analysis and reporting, we usually have to reduce these multiple observations to one per subject. Most of the time, this means averaging the variables or taking some simple statistic (minimum, maximum, median, etc.) and this can be readily accomplished with PROC SUMMARY. However, our client often asks for the values nearest a reference, or the second nearest, or some other algorithm. In these cases, the first thing I do is transpose the data:

```
proc transpose data=one out=two;
  by Subject;
  var score;
run;
```

I would now have one observation per Subject with the value of score occupying a different variable for each of the original multiple observations. By default, PROC TRANSPOSE names these variables Col1, Col2, Col3... up to the maximum

number of multiple observations for any subject in data set one. We could now process the various Col# variables using an ARRAY statement and applying the appropriate algorithm in a DO loop. But there is a problem. In general we don't know how many Col# variables there are. Some subjects have only one observation, while others may have dozens. We could just use a very large guess and write

```
array col{999} col1-col999;
do i=1 to 999;
    {... more statements ...};
end;
```

This approach will work but it can lead to unnecessary computation and memory usage. And what happens if your ridiculously large guess for the number of multiple observations (999) is too small? This problem however can quickly be resolved by writing

```
array col{*} col:;
do i=1 to dim(col);
    {... more statements ...};
end;
```

where we use the colon to list all variables starting with col in the array. Using the array argument '*' causes the SAS system to count the elements in the array and the function DIM calculates the total number of Col# variables. So the colon has enabled us to perform the necessary computations without every knowing the number of observations per subject.

CLEANING UP AFTER MACRO EXECUTION

I rarely write a large SAS program without employing at least one macro program. Such macro programs often generate several temporary data sets for use by the macro. Once the macro is finished executing, I no longer need or want these data sets. So I might write my typical macro as

```
%macro doit;
    data _1;
        set inputdata1;
        {... more statements...}
    run;
    data _2;
        set inputdata2;
        {... more statements...}
    run;

    {... more macro statements...}
    {...create data sets _3, _4, _5 ...}
%mend;
```

You will note that I start all of the names of my macro-generated data sets with an underscore. This is so I have no conflict with data sets produced by the calling program and can easily distinguish the data sets created by the macro. At the conclusion of the macro, good etiquette calls for you to clean up your temporary data sets, so you write

```
proc datasets library=work;
    delete _1 _2 _3 _4 _5;
quit;
```

But if we really have a large macro with many temporary data sets, we could have a big chore keeping track of them and listing them in the PROC DATASETS. However, here again we can exploit the power of the colon:

```
proc data sets library=work;
    delete _: ;
quit;
```

In our previous examples, the use of the colon was restricted to lists of variables, but the colon works just as well on list of data sets. So we can now readily clean up our macro temporary data sets with little effort. In general, if you want to manage groups of data sets, choose appropriate prefixes for their names and you will be able to use PROC DATASETS and the colon to manage them.

AND SO FORTH

The use of the colon for prefixing can be use in many other cases. For example, if you are constantly programming with the same BY-variables, perhaps you should consider shortening the list by using the colon with the appropriate prefix. You would first name all the common BY-variables something like byName, byLocation, byDate, *etc.*; then you could write the BY statements as `'by by: ;'`. Or if you are constantly using the same CLASS-variables, you could prefix the variables and write the CLASS statement `'class class: ;'`. If you are using PROC MEANS or a PROC UNIVARIATE on the same group of analysis variables, you could try `'var anal: ;'`. You could also do a PROC FREQ with the statement `'table tab: ;'`. In short, the prefixing power of the colon can be used in many, if not most, SAS PROCs. (There is unfortunately one notable exception: the colon shortening ability is not available in PROC SQL!).

CONCLUSION

The colon is a very powerful tool for maintaining clear, concise code. It eases the problem of editing old code and enables quicker construction of new programs. The colon facilitates the clearing of temporary variables from data sets and interim data sets from the libraries. In short, the colon can be used across the SAS system to simplify programs and increase productivity.

CONTACT

Your criticisms and witticisms are invited. Contact the author at

William C. Murphy
Howard M. Proskin and Associates, Inc.
300 Red Creek Drive, Suite 220
Rochester, NY 14623
Phone (585) 359-2420
Fax (585) 359-0465
eMail wmurphy@hmproskin.com or wcmurphy@usa.net
web www.hmproskin.com

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are registered trademarks or trademarks of their respective companies.