

Anyone Can Learn PROC TABULATE, v2.0

Lauren Haworth
Ischemia Research & Education Foundation
San Francisco

➤ ABSTRACT

SAS® Software provides hundreds of ways you can analyze your data. You can use the DATA step to slice and dice your data, and there are dozens of procedures that will process your data and produce all kinds of statistics. But odds are that no matter how you organize and analyze your data, you'll end up producing a report in the form of a table.

This is why every SAS user needs to know how to use PROC TABULATE. While TABULATE doesn't do anything that you can't do with other PROCs, the payoff is in the output. TABULATE computes a variety of statistics, and it neatly packages the results in a single table.

Unfortunately, TABULATE has gotten a bad rap as being a difficult procedure to learn. This paper will prove that if you take things step by step, anyone can learn PROC TABULATE.

➤ INTRODUCTION

This paper will start out with the most basic one-dimensional table. We will then go on to two-dimensional tables, tables with totals, and finally, three-dimensional tables. By the end of this paper, you will be ready to build most basic TABULATE tables.

➤ ONE-DIMENSIONAL TABLES

To really understand TABULATE, you have to start very simply. The simplest possible table in TABULATE has to have three things: a PROC TABULATE statement, a TABLE statement, and a CLASS or VAR statement. In this example, we will use a VAR statement. Later examples will show the CLASS statement.

The PROC TABULATE statement looks like this:

```
PROC TABULATE DATA=TEMP;
```

The second part of the procedure is the TABLE statement. It describes which variables to use and how to arrange the variables. This first table will have only one variable, so you don't have to tell TABULATE where to put it. All you have to do is list it in the TABLE statement. When there is only one variable, you get a one-dimensional table.

```
PROC TABULATE DATA=TEMP;  
TABLE RENT;  
RUN;
```

If you run this code as is, you will get an error message because TABULATE can't figure out whether the variable RENT is intended as an analysis variable, which is used to compute statistics, or a classification variable, which is used to define row or column categories in the table.

In this case, we want to use rent as the analysis variable. We will be using it to compute a statistic. To tell TABULATE that RENT is an analysis variable, you use a VAR statement. The syntax of a VAR statement is simple: you just list the variables that will be used for analysis. So now the syntax for our PROC TABULATE is:

```
PROC TABULATE DATA=TEMP;  
VAR RENT;  
TABLE RENT;  
RUN;
```

The result is the table shown below. It has a single column, with the header RENT to identify the variable, and the header SUM to identify the statistic. There is just a single table cell, which contains the value for the sum of RENT for all of the observations in the dataset TEMP.

```
„ffffffffff†  
, RENT ,  
‡ffffffffff%  
, SUM ,  
‡ffffffffff%  
134582.00,  
$ffffffffff€
```

➤ ADDING A STATISTIC

The previous table shows what happens if you don't tell TABULATE which statistic to use. If the variable in your table is an analysis variable, meaning that it is listed in a VAR statement, then the statistic you will get by default is the sum. Sometimes the sum will be the statistic that you want. Most likely, sum isn't the statistic that you want.

To add a statistic to a PROC TABULATE table, you modify the TABLE statement. You list the statistic right after the variable name. To tell TABULATE that the statistic MEAN should be applied to the variable RENT, you use an asterisk to link the variable name to the statistic keyword. The asterisk is a TABULATE *operator*. Just as you use an asterisk as an operator when you want to multiply 2 by 3 (2*3), you use an asterisk when you want to apply a statistic to a variable.

TABLE statement, and following it with the variable name CITY, TABULATE knows that CITY will be used to categorize the mean values of RENT.

```
PROC TABULATE DATA=TEMP;
  CLASS CITY;
  VAR RENT;
  TABLE RENT*MEAN*CITY;
RUN;
```

The resulting table is shown below. Now the column headings have changed. The variable name RENT and the statistic name MEAN are still there, but under the statistic label there are now three columns. Each column is headed by the variable label "CITY" and the category name "Portland," "San Francisco," and "Indianapolis." The values shown in the table cells now represent subgroup means.

```
..+++++
,          RENT
,+++++
,          MEAN
,+++++
,          CITY
,          San
,  Portland  Franci sco  Indi anapoli s,
,+++++
,          859. 67,    1691. 79,    765. 98,
,+++++
,          $+++++
```

➤ **TWO-DIMENSIONAL TABLES**

You probably noticed that our example table is not very elegant in appearance. That's because it only takes advantage of one dimension. It has multiple columns, but only one row. It is much more efficient to build tables that have multiple rows and multiple columns. You can fit more information on a page, and the table looks better, too.

The easiest way to build a two-dimensional table is to build it one dimension at a time. First we'll build the columns, and then we'll add the rows.

For this first table, we'll keep things simple. This is the table we built in a previous example. It has two columns: one showing the number of observations for RENT and another showing the mean of RENT.

```
PROC TABULATE DATA=TEMP;
  VAR RENT;
  TABLE RENT*( N MEAN );
RUN;
```

This table is shown below.

```
..+++++
,          RENT
,+++++
,          N          MEAN
,+++++
,          128. 00,    1051. 42,
,+++++
,          $+++++
```

To turn this table into a two-dimensional table, we will add another variable to the TABLE statement. In this

case, we want to add rows that show the N and MEAN of RENT for different sizes of apartments.

To add another dimension to the table, you use a comma as an operator. All you do is put a comma between the old part of the TABLE statement and the new part of the TABLE statement.

If a TABLE statement has no commas, then it is assumed that the variables and statistics are to be created as columns. If a TABLE statement has two parts, separated by a comma, then TABULATE builds a two-dimensional table using the first part of the TABLE statement as the rows and the second part of the TABLE statement as the columns.

So to get a table with rent as the columns and number of bedrooms as the rows, we just need to add a comma and the variable BEDROOMS. Since we want to add BEDROOMS as a row, we list it before the rest of the TABLE statement. If we wanted to add it as a column, we'd add it to the end of the TABLE statement.

```
PROC TABULATE DATA=TEMP;
  VAR RENT;
  CLASS BEDROOMS;
  TABLE BEDROOMS, RENT*N RENT*MEAN;
RUN;
```

This table is shown below.

```
..+++++
,          RENT          RENT
,          ++++++
,          N          MEAN
,+++++
,          BEDROOMS
,          1 Bedroom    64. 00,    792. 95,
,          2 Bedrooms   64. 00,    1309. 89,
,+++++
,          $+++++
```

By the way, there is a limit to which variables you can use in a two-dimensional table. You can't have a cross-tabulation of two analysis variables. A two-dimensional table must have at least one classification variable (i.e., you must have a CLASS statement). If you think about it, this makes sense. A table of mean rent by mean bedrooms would be meaningless, but a table of mean rent by categories of bedrooms makes perfect sense.

➤ **ADDING CLASSIFICATION VARIABLES ON BOTH DIMENSIONS**

The previous example showed how to reformat a table from one to two dimensions, but it did not show the true power of two-dimensional tables. With two dimensions, you can classify your statistics by two different variables at the same time.

To do this, you put one classification variable in the row dimension and one classification in the column dimension. The previous example had bedrooms displayed in rows, and rent as the column variable. In this new table, we will add city as an additional column variable. Instead

of just displaying the MEAN of the variable RENT for each number of bedrooms, we will display the statistic broken down city.

So in the following code, we leave BEDROOMS as the row variable, and we leave RENT in the column dimension. The only change is to add CITY to the column dimension using the asterisk operator. This tells TABULATE to break down each of the column elements into categories by city.

```
PROC TABULATE DATA=TEMP;
  VAR RENT;
  CLASS BEDROOMS CITY;
  TABLE BEDROOMS, RENT*CITY*MEAN;
RUN;
```

This table is shown below. Notice how the analysis variable RENT remains as the column heading, and MEAN remains as the statistic, but now there are additional column headings to show the three categories of CITY.

```
..+++++..+++++
      RENT
      +++++%
      CITY
      +++++%
      San
      Portland  Francisco  Indianapolis
      +++++%
      MEAN      MEAN      MEAN
+++++
BEDROOMS
+++++
1 Bedroom      620.63      1284.00      608.20
+++++
2 Bedrooms     1098.70      2099.59      923.75
+++++
PARKING
+++++
No              880.09      1479.00      636.50
+++++
Yes             845.63      1780.46      895.45
+++++
$+++++<+++++<+++++<+++++&
```

➤ **ADDING ANOTHER CLASSIFICATION VARIABLE**

The previous example showed how to add a classification variable to both the rows and columns of a two-dimensional table. But you are not limited to just one classification per dimension. This next example will show how to display additional subgroups of the data.

In this case, we're going to add availability of off-street parking as an additional row classification. The variable is added to the CLASS statement and to the row dimension of the TABLE statement. It is added using a space as the operator, so we will get rows for number of bedrooms followed by rows for parking availability.

```
PROC TABULATE DATA=TEMP;
  VAR RENT;
  CLASS BEDROOMS CITY PARKING;
  TABLE BEDROOMS PARKING,
    RENT*CITY*MEAN;
RUN;
```

In the results shown below, you can see that we now have two two-dimensional "mini-tables" within a single table. First we have a table of rent by number of bedrooms and city, and then we have a table showing rent by parking availability and city.

```
..+++++..+++++
      RENT
      +++++%
      CITY
      +++++%
      San
      Portland  Francisco  Indianapolis
      +++++%
      MEAN      MEAN      MEAN
+++++
BEDROOMS
+++++
1 Bedroom      620.63      1284.00      608.20
+++++
2 Bedrooms     1098.70      2099.59      923.75
+++++
PARKING
+++++
No              880.09      1479.00      636.50
+++++
Yes             845.63      1780.46      895.45
+++++
$+++++<+++++<+++++<+++++&
```

This ability to stack multiple mini-tables within a single table can be a powerful tool for delivering large quantities of information in a user-friendly format.

➤ **NESTING THE CLASSIFICATION VARIABLES**

So far all we have done is added additional "tables" to the bottom of our first table. We used the space operator between each of the row variables to produce stacked tables.

By using a series of row variables, we can explore a variety of relationships between the variables. In previous table, we could see how rent varies by number of bedrooms and city, and we could see how rent varies by parking availability and city, but we could not see how number of bedrooms and parking availability interacted to affect rent for each city.

But we can learn even more from our data. The power of TABULATE comes from being able to look at combinations of categories within a single table. In the following example, we will build a table to look at rent by city for combinations of number of bedrooms and parking availability.

This code is the same as we used for the last example. The only change is that in the row definition, the asterisk operator is used to show that we want to nest the two row variables. In other words, we want to see the breakdown of rents by parking availability within each category of number of bedrooms.

```
PROC TABULATE DATA=TEMP;
  VAR RENT;
  CLASS BEDROOMS CITY PARKING;
  TABLE BEDROOMS*PARKING,
    RENT*CITY*MEAN;
RUN;
```

As you can see in the table below, this code produces nested categories within the row headings. The row headings are now split into two columns. The first column shows number of bedrooms and the second shows parking availability. It is now easier to interpret the interaction of the two variables.


```

PROC TABULATE DATA=TEMP
  FORMAT=DOLLAR12. ;
  CLASS BEDROOMS CITY;
  VAR RENT;
  TABLE BEDROOMS=' ',
    (CITY=' ' ALL=' Overall') *
    RENT=' '*MEAN=' '
  / BOX=' Average Rent' ;
RUN;

```

The revised output is shown below:

```

..-----..-----..-----..-----..-----..-----..-----..-----..
Average Rent          San
Rent      Portland  Francisco  Indianapolis  Overall
+-----+-----+-----+-----+-----+-----+-----+-----+
1 Bedroom          $621      $1,284      $608      $793
+-----+-----+-----+-----+-----+-----+-----+-----+
2 Bedrooms          $1,099    $2,100    $924      $1,310
$-----$-----$-----$-----$-----$-----$-----$-----$

```

➤ CREATING HTML OUTPUT

Once you know how to create TABULATE tables, it is a simple matter to turn them into HTML tables that can be posted on your web site.

Using Version 6, all you have to do is download and install the HTML conversion macros that are available on the SAS web site (www.sas.com). Then, you just add a macro call before and after your TABULATE code and SAS will generate the HTML output for you.

The Version 6 code is as follows:

```

%TAB2HTM (CAPTURE=ON, RUNMODE=B) ;
OPTIONS FORMCHAR=' 82838485868788898A8B8C' X;
PROC TABULATE DATA=TEMP
  FORMAT=DOLLAR12. ;
  CLASS BEDROOMS CITY;
  VAR RENT;
  TABLE BEDROOMS=' ',
    (CITY=' ' ALL=' Overall') *
    RENT=' '*MEAN=' '
  / BOX=' Average Rent' ;
RUN;
%TAB2HTM(CAPTURE=OFF, RUNMODE=B,
  OPENMODE=REPLACE, HTMLFILE=SAMPLE.HTML) ;

```

The output is shown below:

| Average Rent | Portland | San Francisco | Indianapolis | Overall |
|--------------|----------|---------------|--------------|---------|
| 1 Bedroom | \$621 | \$1,284 | \$608 | \$793 |
| 2 Bedrooms | \$1,099 | \$2,100 | \$924 | \$1,310 |

Under Versions 7-8, the process of outputting a TABULATE table to HTML is nearly identical. Instead of calling the macro before and after your code, you call on the Output Delivery System (ODS).

```

ODS HTML BODY=' SAMPLE.HTML' ;
PROC TABULATE DATA=TEMP
  FORMAT=DOLLAR12. ;
  CLASS BEDROOMS CITY;
  VAR RENT;
  TABLE BEDROOMS=' ',
    (CITY=' ' ALL=' Overall') *
    RENT=' '*MEAN=' '
  / BOX=' Average Rent' ;
RUN;
ODS HTML CLOSE;

```

The output is shown below:

| Average Rent | Portland | San Francisco | Indianapolis | Overall |
|--------------|----------|---------------|--------------|---------|
| 1 Bedroom | \$621 | \$1,284 | \$608 | \$793 |
| 2 Bedrooms | \$1,099 | \$2,100 | \$924 | \$1,310 |

You can see that the result is nearly identical to the Version 6 table. The only difference is that the Version 7 table uses more colors, and is more stylish.

➤ CHANGING THE STYLE

If you're using Version 7 or 8, you also have the option of changing the style. By adding a STYLE= option to your ODS statement, you can change from the default style to one of the following styles. The resulting output is shown below the code for each style

```

ODS HTML BODY=' SAMPLE.HTML'
STYLE=XXX;

```

The following examples show a few of the styles that ship with Versions 7-8.

STYLE=BARRETTSSBLUE

| Average Rent | Portland | San Francisco | Indianapolis | Overall |
|--------------|----------|---------------|--------------|---------|
| 1 Bedroom | \$621 | \$1,284 | \$608 | \$793 |
| 2 Bedrooms | \$1,099 | \$2,100 | \$924 | \$1,310 |

STYLE=BRI CK

| Average Rent | Portland | San Francisco | Indianapolis | Overall |
|--------------|----------|---------------|--------------|---------|
| 1 Bedroom | \$621 | \$1,284 | \$608 | \$793 |
| 2 Bedrooms | \$1,099 | \$2,100 | \$924 | \$1,310 |

STYLE=STATDOC

| Average Rent | Portland | San Francisco | Indianapolis | Overall |
|--------------|----------|---------------|--------------|---------|
| 1 Bedroom | \$621 | \$1,284 | \$608 | \$793 |
| 2 Bedrooms | \$1,099 | \$2,100 | \$924 | \$1,310 |

STYLE=BROWN

| Average Rent | Portland | San Francisco | Indianapolis | Overall |
|--------------|----------|---------------|--------------|---------|
| 1 Bedroom | \$621 | \$1,284 | \$608 | \$793 |
| 2 Bedrooms | \$1,099 | \$2,100 | \$924 | \$1,310 |

STYLE=THEME

| Average Rent | Portland | San Francisco | Indianapolis | Overall |
|--------------|----------|---------------|--------------|---------|
| 1 Bedroom | \$621 | \$1,284 | \$608 | \$793 |
| 2 Bedrooms | \$1,099 | \$2,100 | \$924 | \$1,310 |

STYLE=D3D

| Average Rent | Portland | San Francisco | Indianapolis | Overall |
|--------------|----------|---------------|--------------|---------|
| 1 Bedroom | \$621 | \$1,284 | \$608 | \$793 |
| 2 Bedrooms | \$1,099 | \$2,100 | \$924 | \$1,310 |

STYLE=MI NI MAL

| Average Rent | Portland | San Francisco | Indianapolis | Overall |
|--------------|----------|---------------|--------------|---------|
| 1 Bedroom | \$621 | \$1,284 | \$608 | \$793 |
| 2 Bedrooms | \$1,099 | \$2,100 | \$924 | \$1,310 |

➤ CONCLUSIONS

At this point, you should be comfortable with the basics of producing a table using PROC TABULATE. You should be able to produce a simple table with totals, be able to clean it up a bit, and be able to create HTML output.

This should be enough to get you going producing tables with your own data. And now that you're more comfortable with the procedure, you should be able to use the TABULATE manual and other books and papers to learn more advanced techniques.

➤ ACKNOWLEDGEMENTS

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

➤ CONTACTING THE AUTHOR

Please direct any questions or feedback to the author at:

Ischemia Research & Education Foundation
250 Executive Park Blvd, Suite 3400
San Francisco, CA 94134

E-mail: leh@iref.org