

SAS

Getting Started  
with SAS Programming

# Using SAS Studio in the Cloud



Ron Cody



# Getting Started with SAS® Programming

# Using SAS® Studio in the Cloud



Ron Cody



**Getting Started  
with SAS® Programming**

# **Using SAS® Studio in the Cloud**



**Ron Cody**

The correct bibliographic citation for this manual is as follows: Cody, Ron. 2021. *Getting Started with SAS® Programming: Using SAS® Studio in the Cloud*. Cary, NC: SAS Institute Inc.

## **Getting Started with SAS® Programming: Using SAS® Studio in the Cloud**

Copyright © 2021, SAS Institute Inc., Cary, NC, USA

ISBN 978-1-953329-20-2 (Hardcover)

ISBN 978-1-953329-16-5 (Paperback)

ISBN 978-1-953329-17-2 (Web PDF)

ISBN 978-1-953329-18-9 (EPUB)

ISBN 978-1-953329-19-6 (Kindle)

All Rights Reserved. Produced in the United States of America.

**For a hard copy book:** No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, or otherwise, without the prior written permission of the publisher, SAS Institute Inc.

**For a web download or e-book:** Your use of this publication shall be governed by the terms established by the vendor at the time you acquire this publication.

The scanning, uploading, and distribution of this book via the Internet or any other means without the permission of the publisher is illegal and punishable by law. Please purchase only authorized electronic editions and do not participate in or encourage electronic piracy of copyrighted materials. Your support of others' rights is appreciated.

**U.S. Government License Rights; Restricted Rights:** The Software and its documentation is commercial computer software developed at private expense and is provided with RESTRICTED RIGHTS to the United States Government. Use, duplication, or disclosure of the Software by the United States Government is subject to the license terms of this Agreement pursuant to, as applicable, FAR 12.212, DFAR 227.7202-1(a), DFAR 227.7202-3(a), and DFAR 227.7202-4, and, to the extent required under U.S. federal law, the minimum restricted rights as set out in FAR 52.227-19 (DEC 2007). If FAR

52.227-19 is applicable, this provision serves as notice under clause (c) thereof and no other notice is required to be affixed to the Software or documentation. The Government's rights in Software and documentation shall be only those set forth in this Agreement.

SAS Institute Inc., SAS Campus Drive, Cary, NC 27513-2414

February 2021

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 trademarks of their respective companies.

SAS software may be provided with certain third-party software, including but not limited to open-source software, which is licensed under its applicable third-party software license agreement. For license information about third-party software distributed with SAS software, refer to <http://support.sas.com/thirdpartylicenses>.

# Contents

[About This Book](#)

[About The Author](#)

[Acknowledgements](#)

[\*\*Part I: Getting Acquainted with the SAS Studio Environment\*\*](#)

[\*\*Chapter 1: Introduction to SAS OnDemand for Academics\*\*](#)

[Introduction: An Overview of SAS OnDemand for Academics](#)

[Registering for ODA](#)

[Conclusion](#)

[\*\*Chapter 2: The SAS Studio Interface\*\*](#)

[Introduction](#)

[Exploring the Built-In Data Sets](#)

[Sorting Your Data](#)

[Switching between Column Names and Column Labels](#)

[Resizing Tables](#)

[Creating Filters](#)

[Conclusion](#)

[\*\*Chapter 3: Importing Your Own Data\*\*](#)

[Introduction](#)

[Uploading Data from Your Local Computer to SAS Studio](#)

[Listing the SAS Data Set](#)

[Importing an Excel Workbook with Invalid SAS Variable Names](#)

[Importing an Excel Workbook That Does Not Have Variable Names](#)

[Importing Data from a CSV File](#)

[Conclusion](#)

[\*\*Chapter 4: Creating Reports\*\*](#)

[Introduction](#)

[Using the List Data Task to Create a Simple Listing](#)

[Filtering Data](#)

[Sorting Data](#)

[Outputting HTML and PDF Files](#)

[Joining Tables \(Using the Query Window\)](#)

[Conclusion](#)

[\*\*Chapter 5: Summarizing Data Using SAS Studio\*\*](#)

[Introduction](#)

[Summarizing Numeric Variables](#)

[Adding a Classification Variable](#)

[Summarizing Character Variables](#)

## [Conclusion](#)

### [Chapter 6: Graphing Data](#)

[Introduction](#)

[Creating a Frequency Bar Chart](#)

[Creating a Bar Chart with a Response Variable](#)

[Adding a Group Variable](#)

[Creating a Pie Chart](#)

[Creating a Scatter Plot](#)

[Conclusion](#)

## [Part II: Learning How to Write Your Own SAS Programs](#)

### [Chapter 7: An Introduction to SAS Programming](#)

[SAS as a Programming Language](#)

[The SAS Studio Programming Windows](#)

[Your First SAS Program](#)

[DATA Statement](#)

[INILE Statement](#)

[INPUT Statement](#)

[Assignment Statement](#)

[How the DATA Step Works](#)

[How the INPUT Statement Works](#)

[Reading Delimited Data](#)

[How Procedures \(PROCs\) Work](#)

[How SAS Works: A Look Inside the “Black Box”](#)

[Conclusion](#)

### [Chapter 8: Reading Data from External Files](#)

[Introduction](#)

[Reading Data Values Separated by Delimiters](#)

[Reading Comma-Separated Values Files](#)

[Reading Data Separated by Other Delimiters](#)

[Reading Data in Fixed Columns](#)

[Reading Data in Fixed Columns Using Column Input](#)

[Reading Data in Fixed Columns Using Formatted Input](#)

[Conclusion](#)

[Problems](#)

### [Chapter 9: Reading and Writing SAS Data Sets](#)

[What's a SAS Data Set?](#)

[Temporary Versus Permanent SAS Data Sets](#)

[Creating a Library by Submitting a LIBNAME Statement](#)

[Using the Library Tab to Create a Permanent Library](#)

[Reading from a Permanent SAS Data Set](#)

[Conclusion](#)

[Problems](#)

### [Chapter 10: Creating Formats and Labels](#)

[What Is a SAS Format and Why Is It Useful?](#)

[Using SAS Built-in Formats](#)

[More Examples to Demonstrate How to Write Formats](#)  
[Describing the Difference between a FORMAT Statement in a Procedure and a FORMAT Statement in a DATA Step](#)  
[Making Your Formats Permanent](#)  
[Creating Variable Labels](#)  
[Conclusion](#)  
[Problems](#)

## [\*\*Chapter 11: Performing Conditional Processing\*\*](#)

[Introduction](#)  
[Grouping Age Using Conditional Processing](#)  
[Using Conditional Logic to Check for Data Errors](#)  
[Describing the IN Operator](#)  
[Using Boolean Logic \(AND, OR, and NOT Operators\)](#)  
[A Special Caution When Using Multiple OR Operators](#)  
[Conclusion](#)  
[Problems](#)

## [\*\*Chapter 12: Performing Iterative Processing: Looping\*\*](#)

[Introduction](#)  
[Demonstrating a DO Group](#)  
[Describing a DO Loop](#)  
[Using a DO Loop to Graph an Equation](#)  
[DO Loops with Character Values](#)  
[Leaving a Loop Based on Conditions \(DO WHILE and DO UNTIL Statements\)](#)  
[DO WHILE](#)  
[Combining an Iterative Loop with a WHILE Condition](#)  
[DO UNTIL](#)  
[Demonstrating That a DO UNTIL Loop Executes at Least Once](#)  
[Combining an Iterative Loop with an UNTIL Condition](#)  
[LEAVE and CONTINUE Statements](#)  
[Conclusion](#)  
[Problems](#)

## [\*\*Chapter 13: Working with SAS Dates\*\*](#)

[Introduction](#)  
[Reading Dates from Text Data](#)  
[Creating a SAS Date from Month, Day, and Year Values](#)  
[Describing a Date Constant](#)  
[Extracting the Day of the Week, Day of the Month, Month, and Year from a SAS Date](#)  
[Adding a Format to the Bar Chart](#)  
[Computing Age from Date of Birth: The YRDIF Function](#)  
[Conclusion](#)  
[Problems](#)

## [\*\*Chapter 14: Subsetting and Combining SAS Data Sets\*\*](#)

[Introduction](#)  
[Subsetting \(Filtering\) Data in a SAS Data Set](#)  
[Describing a WHERE= Data Set Option](#)  
[Describing a Subsetting IF Statement](#)

[A More Efficient Way to Subset Data When Reading Raw Data](#)  
[Creating Several Data Subsets in One DATA Step](#)  
[Combining SAS Data Sets \(Combining Rows\)](#)  
[Adding a Few Observations to a Large Data Set \(PROC APPEND\)](#)  
[Interleaving Data Sets](#)  
[Merging Two Data Sets \(Adding Columns\)](#)  
[Controlling Which Observations Are Included in a Merge \(IN= Data Set Option\)](#)  
[Performing a One-to-Many or Many-to-One Merge](#)  
[Merging Two Data Sets with Different BY Variable Names](#)  
[Merging Two Data Sets with One Character and One Numeric BY Variable](#)  
[Updating a Master File from a Transaction File \(UPDATE Statement\)](#)  
[Conclusion](#)  
[Problems](#)

## **[Chapter 15: Describing SAS Functions](#)**

[Introduction](#)  
[Describing Some Useful Numeric Functions](#)  
[Function Name: MISSING](#)  
[Function Name: N](#)  
[Function Name: NMISS](#)  
[Function Name: SUM](#)  
[Function Name: MEAN](#)  
[Function Name: MIN](#)  
[Function Name: MAX](#)  
[Function Name: SMALLEST](#)  
[Function Name: LARGEST](#)  
[Programming Example Using the N, NMISS, MAX, LARGEST, and MEAN Functions](#)  
[Function Name: INPUT](#)  
[CALL Routine: CALL SORTN](#)  
[Function Name: LAG](#)  
[Function Name: DIF](#)  
[Describing Some Useful Character Functions](#)  
[Function Names: LENGTHN and LENGTHC](#)  
[Function Names: TRIMN and STRIP](#)  
[Function Names: UPCASE, LOWCASE, and PROPCASE \(Functions That Change Case\)](#)  
[Function Name: PUT](#)  
[Function Name: SUBSTRN \(Newer Version of the SUBSTR Function\)](#)  
[Function Names: FIND and FINDC](#)  
[Function Names: CAT, CATS, and CATX](#)  
[Function Names: COUNT and COUNTC](#)  
[Function Name: COMPRESS](#)  
[Function Name: SCAN](#)  
[CALL Routine: CALL MISSING](#)  
[Function Names: NOTDIGIT, NOTALPHA, and NOTALNUM](#)  
[Function Names: ANYDIGIT, ANYALPHA, and ANYALNUM](#)  
[Function Name: TRANWRD](#)  
[Conclusion](#)

## [Problems](#)

### [\*\*Chapter 16: Working with Multiple Observations per Subject\*\*](#)

[Introduction](#)

[Useful Tools for Working with Longitudinal Data](#)

[Describing First and Last Variables](#)

[Computing Visit-to-Visit Differences](#)

[Computing Differences between the First and Last Visits](#)

[Counting the Number of Visits for Each Patient](#)

[Conclusion](#)

[Problems](#)

### [\*\*Chapter 17: Describing Arrays\*\*](#)

[Introduction](#)

[What Is an Array?](#)

[Describing a Character Array](#)

[Performing an Operation on Every Numeric Variable in a Data Set](#)

[Performing an Operation on Every Character Variable in a Data Set](#)

[Converting a Data Set with One Observation per Subject into a Data Set with Multiple Observations per Subject](#)

[Converting a Data Set with Multiple Observations per Subject into a Data Set with One Observation per Subject](#)

[Conclusion](#)

[Problems](#)

### [\*\*Chapter 18: Displaying Your Data\*\*](#)

[Introduction](#)

[Producing a Simple Report Using PROC PRINT](#)

[Using Labels Instead of Variable Names as Column Headings](#)

[Including a BY Variable in a Listing](#)

[Including the Number of Observations in a Listing](#)

[Conclusion](#)

[Problems](#)

### [\*\*Chapter 19: Summarizing Data with SAS Procedures\*\*](#)

[Introduction](#)

[Using PROC MEANS \(with the Default Options\)](#)

[Using PROC MEANS Options to Customize the Summary Report](#)

[Computing Statistics for Each Value of a BY Variable](#)

[Using a CLASS Statement Instead of a BY Statement](#)

[Including Multiple CLASS Variables with PROC MEANS](#)

[Statistics Broken Down Every Way](#)

[Using PROC MEANS to Create a Summary Data Set](#)

[Letting PROC MEANS Name the Variables in the Output Data Set](#)

[Creating a Summary Data Set with CLASS Variables](#)

[Using a Formatted CLASS Variable](#)

[Demonstrating PROC UNIVARIATE](#)

[Conclusion](#)

[Problems](#)

## **Chapter 20: Computing Frequencies**

[Introduction](#)

[Creating a Data Set to Demonstrate Features of PROC FREQ](#)

[Using PROC FREQ to Generate One-Way Frequency Tables](#)

[Creating Two-Way Frequency Tables](#)

[Creating Three-Way Frequency Tables](#)

[Using Formats to Create Groups for Numeric Variables](#)

[Conclusion](#)

[Problems](#)

# About This Book

## What Does This Book Cover?

This book has two goals: the first is to show readers how to use a free version of SAS called SAS OnDemand for Academics, including how to use point-and-click menus to view, summarize, and analyze data using built-in SAS Studio tasks. The second goal is to teach readers how to program using SAS.

The first part of the book shows readers how to register for SAS OnDemand for Academics. The remaining chapters in this section explore how to use the SAS Studio tasks to inspect, summarize, display, and, finally, how to create graphical representations of data.

The second part of the book is an introduction to SAS programming. Starting from basic concepts, this part of the book discusses conditional logic, looping, SAS functions, and some slightly more advanced topics such as how to analyze longitudinal data and transform SAS data sets.

Although this book covers basic and intermediate topics, more advanced topics such as SAS date interval functions and Perl regular expressions are not covered.

## Is This Book for You?

This book is appropriate for beginners as well as intermediate programmers. Even people with advanced SAS programming skills might find this book useful to learn how to use SAS Studio tasks in a cloud-based environment.

## What Are the Prerequisites for This Book?

There are essentially no prerequisites for people thinking about buying this book.

## What's New in This Edition?

Parts of this book are similar to an earlier book called *An Introduction to SAS® University Edition*. However, because that book used SAS University Edition while this book uses SAS OnDemand for Academics, it should rightfully be considered a new book and not a second edition of the older book.

## What Should You Know about the Examples?

This book includes tutorials for you to follow to gain hands-on experience with SAS. All the programs and data sets used in the text, as well as data used for the end-of-chapter problems, are included in a free download from the SAS author site. Every topic in the programming section is introduced by one or more examples.

## Software Used to Develop the Book's Content

Every program in the book was written and run using SAS OnDemand for Academics, the SAS cloud-based platform.

## Example Code and Data

You can access the example code and data for this book by linking to its author page at <https://support.sas.com/cody>.

## SAS OnDemand for Academics

 This book is compatible with SAS OnDemand for Academics. If you are using SAS OnDemand for Academics, then begin here: [https://www.sas.com/en\\_us/software/on-demand-for-academics.html](https://www.sas.com/en_us/software/on-demand-for-academics.html).

## Where Are the Exercise Solutions?

Solutions to the odd-numbered problems are included at the end of the book as well and in the free download from the author site. Self-learners or instructors can request the solutions to the even-numbered problems by contacting SAS Press.

## We Want to Hear from You

SAS Press books are written *by SAS Users for SAS Users*. We welcome your participation in their development and your feedback on SAS Press books that you are using. Please visit [sas.com/books](http://sas.com/books) to do the following:

- Sign up to review a book
- Recommend a topic
- Request information on how to become a SAS Press author
- Provide feedback on a book

Do you have questions about a SAS Press book that you are reading? Contact the author through [saspress@sas.com](mailto:saspress@sas.com) or [https://support.sas.com/author\\_feedback](https://support.sas.com/author_feedback).

SAS has many resources to help you find answers and expand your knowledge. If you need additional help, see our list of resources: [sas.com/books](http://sas.com/books).

# About The Author



Ron Cody, EdD, is a retired professor from the Rutgers Robert Wood Johnson Medical School who now works as a private consultant and national instructor for SAS. A SAS user since 1977, Ron's extensive knowledge and innovative style have made him a popular presenter at local, regional, and national SAS conferences. He has authored or co-authored numerous books, such as *Learning SAS by Example: A Programmer's Guide, Second Edition*; *A Gentle Introduction to Statistics Using SAS Studio*; *Cody's Data Cleaning Techniques Using SAS, Third Edition*; *SAS Functions by Example, Second Edition*; and several other books on SAS programming and statistical analysis. During his career at Rutgers Medical School, he authored or co-authored over 100 articles in scientific journals.

Learn more about this author by visiting his author page at <http://support.sas.com/cody>. There you can download free book excerpts, access example code and data, read the latest reviews, get updates, and more.

# Acknowledgements

Writing the acknowledgements is, for me, an emotional event. Not only did a talented group of people help to bring this book to fruition, I also realize I am almost at the end of the writing process.

As I have said in previous books, I thank my wife, Jan, for her inspiration and telling me to stop writing for a while and go for a walk! Once again, it was Jan who took the picture for the back cover.

I should start by thanking Suzanne Morgen, the developmental editor and coordinator of the entire project. There are a lot of moving parts in creating a book, and I thank Suzanne for keeping everything on track.

Next, a hearty thanks to the three reviewers: Paul Grant, Russell Lavery, and David Franklin. Paul has reviewed almost every book I have written and keeps coming back for more, even thanking me for the chance! This is the first time Russ Lavery has reviewed one of my books, and I see that his experience in reviewing several of Art Carpenter's books has made him very good at the task. I also needed an expert on SAS OnDemand for Academics. David Franklin was the go-to guy for this. So, thanks to the three of you.

Next are the copyeditors and folks who put everything together. I send in individual chapters, but the final book is so much more. Thanks to Catherine Connolly and Denise Jones for the very difficult task of reading every word, checking every figure caption, and putting everything together.

Last, but certainly not least, my sincere thanks to Robert Harris, artist extraordinaire! You only see one cover on this book, but Robert actually created four covers for me to choose from. The one you see here just grabs your eye to look at it. Covers are really important. Robert, it is a pleasure working with you!

Ron Cody

# Part I: Getting Acquainted with the SAS Studio Environment

[Chapter 1 Introduction to SAS OnDemand for Academics](#)

[Chapter 2 The SAS Studio Interface](#)

[Chapter 3 Importing Your Own Data](#)

[Chapter 4 Creating Reports](#)

[Chapter 5 Summarizing Data Using SAS Studio](#)

[Chapter 6 Graphing Data](#)

**Part I** shows you how to perform basic tasks such as producing a report, summarizing data, producing charts and graphs, and using the SAS Studio built-in tasks.

# Chapter 1: Introduction to SAS OnDemand for Academics

## Introduction: An Overview of SAS OnDemand for Academics

SAS is many things: A data analysis tool, a programming language, a statistical package, a tool for business intelligence, and more. Until recently, you could only get access to SAS by paying a license fee (this could be an individual license, or a license purchased by a company for as many users as necessary).

The really big news is that anyone can now obtain SAS for FREE! It's called SAS OnDemand for Academics—but don't be fooled by the name. **Anyone can use this free version of SAS, not just students enrolled in a class.**

Some of you might know about a SAS product called the SAS University Edition. This was another free version of SAS, but you had to download software to create a virtual computer on your real computer, then download the SAS software, and finally, set up a way to read and write files from your “real” computer to the “virtual computer.” This caused many people massive headaches (including this author). The great news about SAS OnDemand for Academics (hence forth called ODA – OnDemand for Academics) is that **you don't have to download anything!** You access SAS on a cloud platform. Also, reading data from your real computer is quite simple.

And now for the caveats: This product was developed so that people can use it to learn how to program and run tasks using SAS. It is not supposed to be used for commercial purposes. One final note: there is a 5-gigabyte limit for data files, but that is certainly not a problem for learning how to use SAS.

On many college campuses, students taking statistics courses, or any

course that needs a powerful analytic tool could access a computer language called R, for free. Since free is better than not free, these institutions sometimes choose to use R instead of SAS. That is fine, except that when these students graduate, they find that in the corporate world, SAS is by far the major package for powerful statistical analysis, data manipulation, and reporting. By offering a free version of SAS, users now have a choice between SAS or R and SAS Institute is hoping that the majority of users will choose SAS.

ODA uses SAS Studio as the interface. SAS Studio provides an environment that includes a point-and-click facility for performing many common tasks, such as producing reports, graphs, data summaries, and statistical tests. For those who either enjoy programming or have more complicated tasks, SAS Studio also enables you to write and run your own programs.

## Registering for ODA

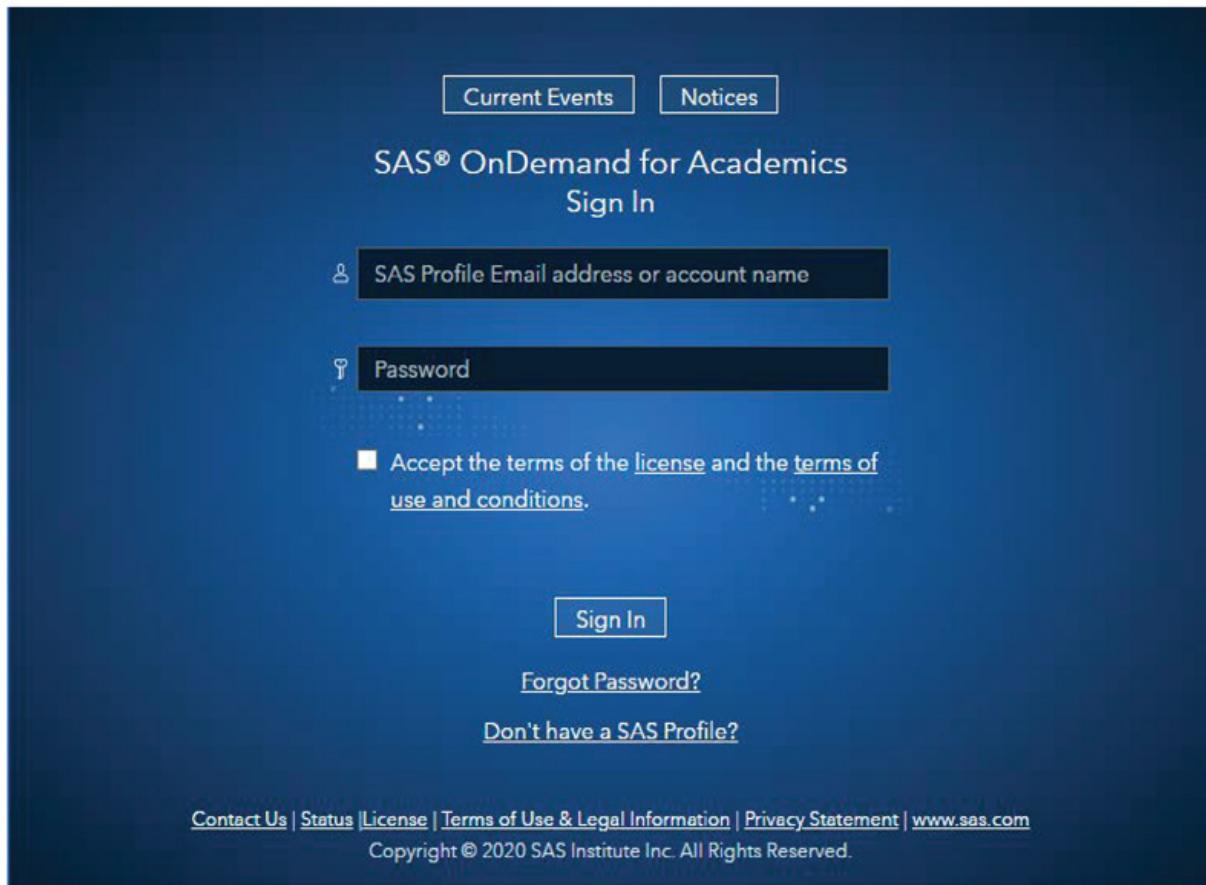
To gain access to ODA, you need to register with SAS Institute. Part of the registration process is to create a SAS profile. If you already have a SAS profile, skip that portion of the instructions.

To start, point your browser to:

<https://welcome.oda.sas.com>

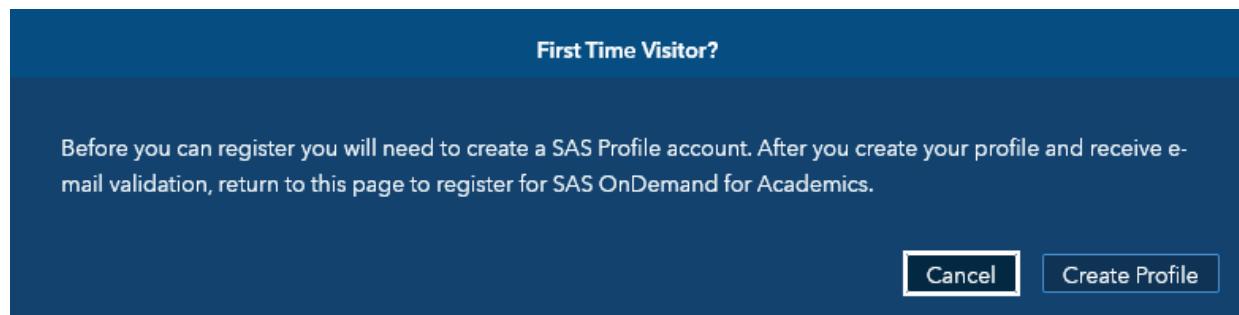
This brings up the screen shown here.

**Figure 1.1: Registration Screen for ODA**



If you do not have a SAS Profile, click “Don’t have a SAS Profile.” You will see the screen shown in Figure 1.2.

**Figure 1.2: First Time Visitor**



Click Create Profile.

**Figure 1.3: Enter Your Personal Information**



SAS Profile

Step 1 of 2: Tell us about yourself.

Preferred Language

First Name \*

Last Name \*

Email \*

Country/Region \*

Affiliation With SAS \*

Company/Organization \*

\*Required

By the way, if you are self-employed or “retired” as I am, just enter self-employed or retired, or whatever describes your situation.

Finally, click “Agree with terms” and click the box Create Profile.

#### **Figure 1.4: Agree with Terms and Continue**

Affiliation With SAS \*

Company/Organization \*

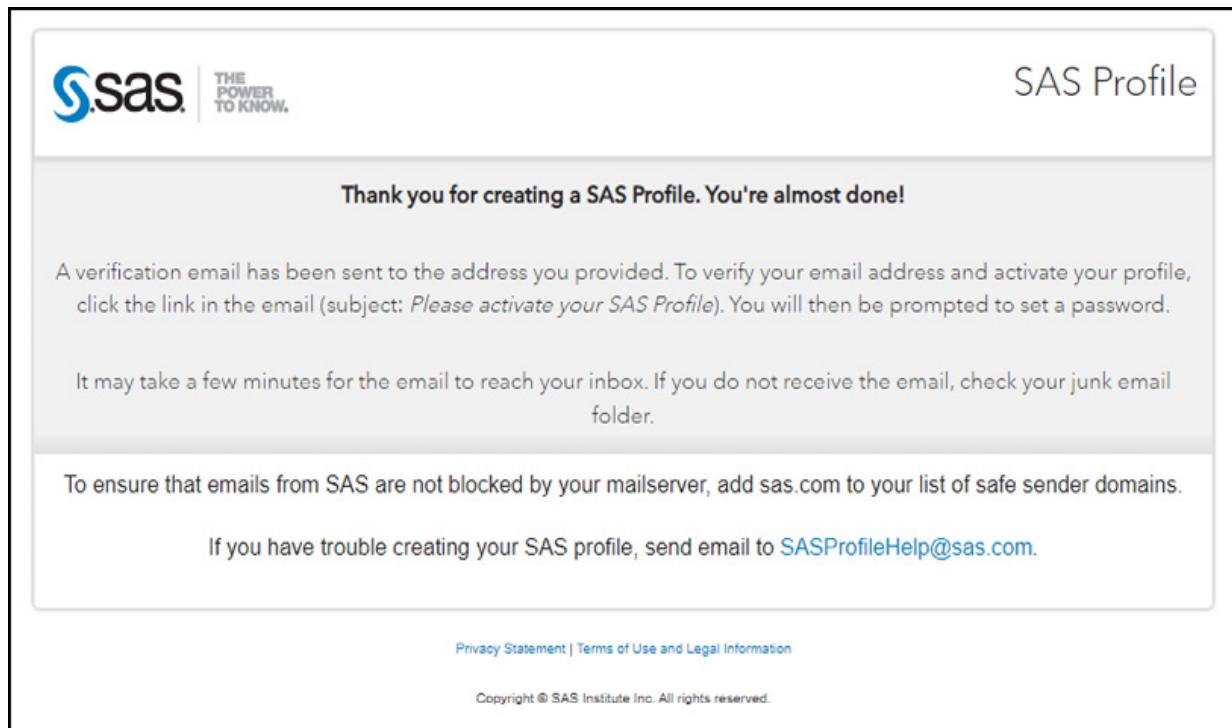
\*Required

Yes, I would like to receive occasional emails from SAS and its affiliates about SAS products, events, free white papers, special offers, etc. I understand that all personal information will be handled in accordance with the SAS Privacy Statement. [Learn more](#).

I agree to the [terms of use and conditions](#). \*

After clicking "Create profile," you will receive a verification email with instructions for setting your password and activating your profile.

**Figure 1.5: Verifying Your Email Address**



When your email arrives (this might take some time, perhaps a day), click "Activate your profile."

**Figure 1.6: Activate Your Profile**

replies-disabled@sas.com  
To:

---

**Please activate your SAS Profile**  
Today at 9:16 AM

Dear ODA User,

Thank you for creating a SAS Profile. Please verify your email address and activate your profile by clicking on the following link:

**[Activate your SAS Profile.](#)**

Once you've activated your profile, you can access free training, SAS communities, technical support and more. If you have any questions or have not requested a SAS Profile, please email [SASProfileHelp@sas.com](mailto:SASProfileHelp@sas.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. © SAS Institute Inc. All rights reserved.

Choose your password and confirm.

**Figure 1.7: Choose Your Password and Confirm**

The screenshot shows a web-based password reset interface. At the top left is the SAS logo with the tagline "THE POWER TO KNOW". To the right, the text "SAS Profile" is displayed. Below this, a header bar indicates "Step 2 of 2: Please set your new password." The main content area contains two input fields: "Password (show)" and "Confirm password", both represented by white input boxes with black outlines. A "Save changes" button is located below these fields. At the bottom of the page, there are links for "Privacy Statement | Terms of Use and Legal Information" and "Copyright © SAS Institute Inc. All rights reserved."

Note that your password must contain at least 8 characters and include at least 1 lowercase letter, 1 uppercase letter, 1 symbol, and 1 number.

**Figure 1.8: Continuing the Password Selection**

**SAS Profile**

Step 2 of 2: Please set your new password.

**Password (show)**

••••••••

Your password must contain at least:

8 characters 1 lowercase letter 1 uppercase letter 1 symbol 1 number

Do not use #, \, --, first name, last name, or part of your email address.

**Confirm password**

••••••••

[Privacy Statement](#) | [Terms of Use and Legal Information](#)

Copyright © SAS Institute Inc. All rights reserved.

**Figure 1.9: A Verification Email Will be Sent**

**SAS Profile**

Your profile is now active! You can now sign in using your new SAS Profile credentials.

[Privacy Statement](#) | [Terms of Use and Legal Information](#)

Copyright © SAS Institute Inc. All rights reserved.

You are now ready to sign in to OnDemand for Academics (ODA).

**Figure 1.10: Signing into OnDemand**

The screenshot shows the SAS OnDemand for Academics sign-in interface. At the top, there are two buttons: "Current Events" and "Notices". Below them is the title "SAS® OnDemand for Academics" followed by a "Sign In" button. The main form area contains two input fields: "SAS Profile Email address or account name" with a user icon and "Password" with a lock icon. Below these fields is a checkbox labeled "Accept the terms of the [license](#) and the [terms of use and conditions](#)". To the right of the checkbox is a note: "If you do not accept the terms, you will not be able to log in." At the bottom of the form are three links: "Sign In", "Forgot Password?", and "Don't have a SAS Profile?". At the very bottom of the page, there is a footer with links to "Contact Us", "Status", "License", "Terms of Use & Legal Information", "Privacy Statement", and the website "www.sas.com", followed by the copyright notice "Copyright © 2020 SAS Institute Inc. All Rights Reserved."

Be sure to accept terms and then click “Sign In.” You will be asked to select a region.

**Figure 1.11: Select a Region**

Need help selecting your home region?

In most cases, the site performance from your typical usage location will indicate the best choice. Below, you can perform a simple test to estimate that performance if you are currently at that location.

[Compare Performance of All Regions](#)

Alternatively, we can suggest a region based on your expected usage location if you select that location below. If you are currently at that location, you can select **(Current Location)** to detect it automatically.

You are done. Once you have received an email saying that your account is ready to use, you can select SAS Studio.

**Figure 1.12: Start SAS Studio**

SAS® OnDemand for Academics  
Dashboard

Planned Events Notices

Applications Enrollments Courses

**SAS® Studio**  
Write and run SAS code with a Web-based SAS development environment.  
Actions: [Clear my saved tabs](#).

Other Ways to Access SAS® OnDemand for Academics Resources

**SASPy access to SAS® hosted servers**  
Use Python applications to run SAS® from the SAS OnDemand for Academics hosted environments.  
[\(Configuration Steps Required\)](#)

Reference

[Support Site](#)  
[Step-by-Step Reference Guides](#)  
[Frequently Asked Questions](#)

Quotas [\(learn more\)](#)  
Home Directory (0.7MB/5120MB)  
0%

Click SAS Studio.

## Conclusion

Registering for OnDemand for Academics is really quick and easy. Just follow all the prompts and it should only take a few minutes of your time (you will need to wait for your profile to be accepted).

# Chapter 2: The SAS Studio Interface

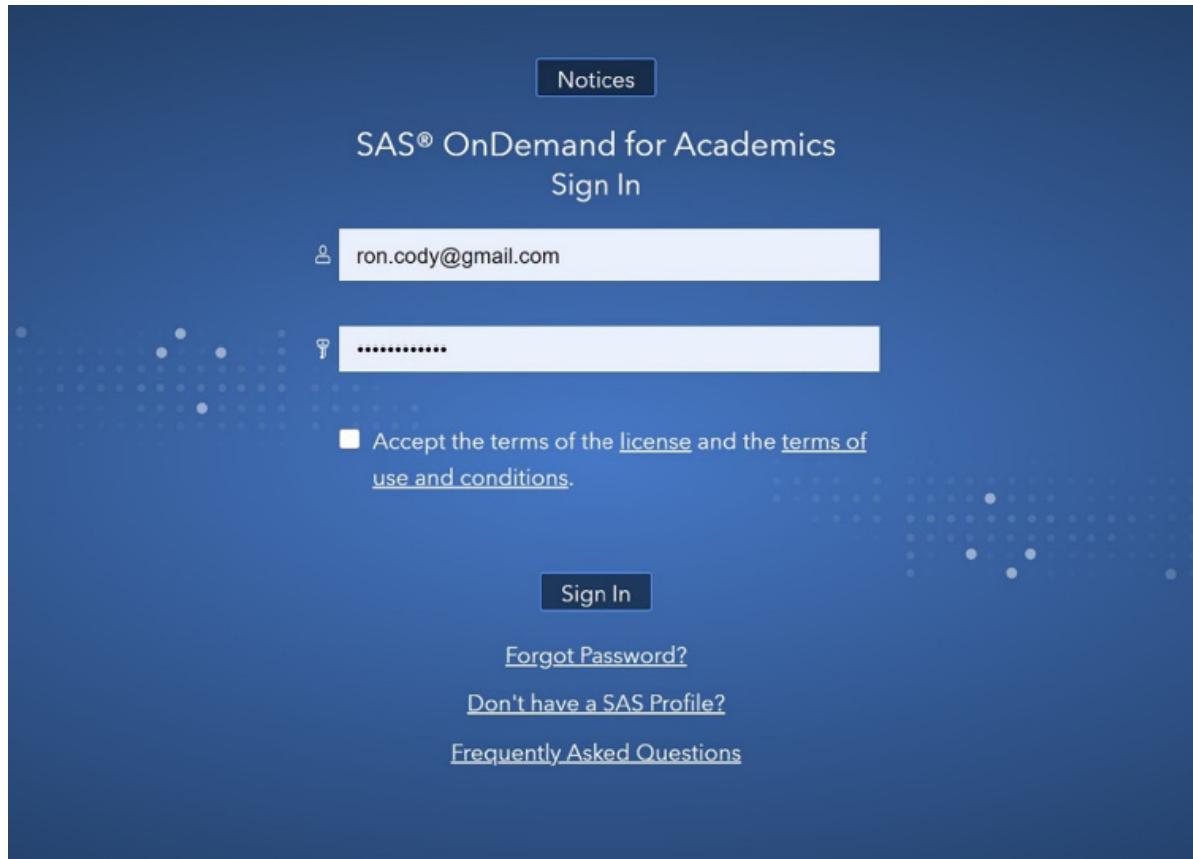
## Introduction

To begin using OnDemand for Academics, click the URL

<https://welcome.oda.sas.com>

When you first open OnDemand for Academics, you will see the following screen.

**Figure 2.1: Starting SAS OnDemand for Academics**



You need to check the box labeled “Accept the terms of the license and the terms of use and conditions.” Once that is completed, click

the box labeled Sign In.

The opening screen in SAS Studio looks like this.

**Figure 2.2: Opening Screen for SAS Studio**

The screenshot shows the SAS® OnDemand for Academics Dashboard. At the top, it says "SAS® OnDemand for Academics Dashboard". Below that are tabs for "Planned Events" and "Notices". A main section has three buttons: "Applications", "Enrollments", and "Courses", with "Applications" being the active tab. Under "Applications", there's a section for "SAS® Studio" which says "Write and run SAS code with a Web-based SAS development environment." and "Actions: Clear my saved tabs.". Below this, there's a section titled "Other Ways to Access SAS® OnDemand for Academics Resources" with links for "SASPy access to SAS® hosted servers" and "PY; Use Python applications to run SAS® from the SAS OnDemand for Academics hosted environments. (Configuration Steps Required)". To the right, there's a "Reference" sidebar with links for "Support Site", "Step-by-Step Reference Guides", "Frequently Asked Questions", and "Quotas (learn more)" which shows a progress bar at 0% and "Home Directory (0.3MB/5120MB)".

Click SAS Studio to begin your session.

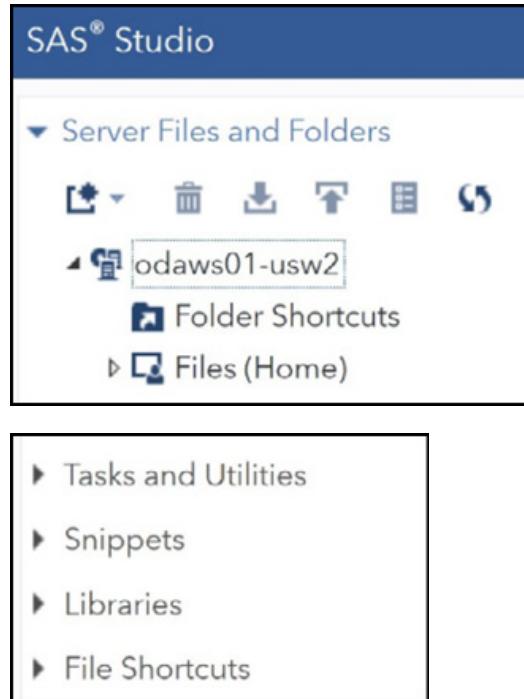
After a few seconds, you arrive at the home screen (yours might vary because this author has already created files and folders).

**Figure 2.3: Home Screen for SAS Studio**

The screenshot shows the SAS Studio work area. On the left is a navigation pane with sections for "Server Files and Folders" (containing "odaws02-usw2", "Folder Shortcuts", and "Files (lhome)"), "Tasks and Utilities", "Snippets", "Libraries", and "File Shortcuts". On the right is a work area titled "Program 1" with tabs for "CODE", "LOG", and "RESULTS". The "CODE" tab shows a single line of code: "1". The status bar at the bottom right shows "Line 1, Column 1", "UTF-8", "Messages", and "User: ronaldready".

On the left, you see the navigation pane—on the right, the work area. Here is a blow up of the navigation pane:

**Figure 2.4: Blow Up of the Navigation Pane**



## Exploring the Built-In Data Sets

SAS libraries are where SAS stores SAS data sets. SAS Studio ships with a number of data sets that you can play with. In the next chapter, you will see how to create your own SAS data sets from several different data sources and store them in a library of your own.

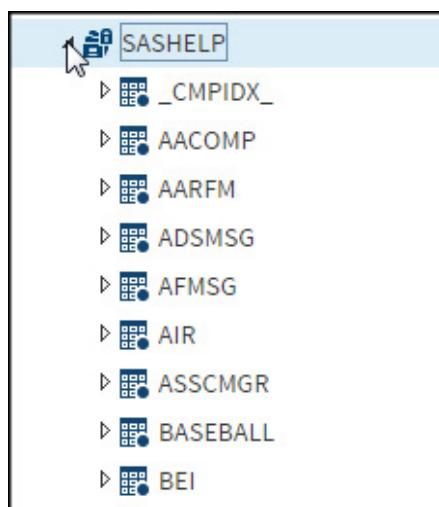
First click the **Libraries** tab. Next click the **My Libraries** tab, you will see something like this.

**Figure 2.5: Expanding the My Libraries Tab**



Some of the libraries that you see here were created by the author. Others, such as SASUSER and SASHelp, will always show up. The SASHelp library is where SAS has stored all of the demonstration data sets. Click the small triangle to the left of SASHelp to see a list of the SAS data sets stored there. It should look like Figure 2.6 below.

**Figure 2.6: Expanding the SASHelp Library**



You can click the small triangle to the left of any of these data sets to see a list of variables. As an alternative, double-click a data set of interest to display a list of variables and a partial listing of the data set. For this example, let's double-click the **BASEBALL** data set. Here's what happens.

**Figure 2.7: Opening the BASEBALL Data Set**

The screenshot shows the SAS Studio interface. On the left, the Libraries pane is open, displaying various datasets like BOOK, MYFMTS, SASHHELP, and BASEBALL. The BASEBALL dataset is selected. The main workspace shows the variable selection pane on the left and the data preview pane on the right.

**Variables Selection:**

- Selected variables: Name, Team, nAtBat, nHits, nHome, nRuns, nRBI, nBB, YrMajor, CrAtBat, CrHits, CrHome, CrRuns.
- Unselected variables: nRBI, nBB, CrAtBat, CrHits, CrHome, CrRuns.

**Data Preview:**

Name	Team	nAtBat	nHits
Allanson, Andy	Cleveland	293	66
Ashby, Alan	Houston	315	81
Davis, Alan	Seattle	479	130
Dawson, Andre	Montreal	496	141
Galaraga, Andres	Montreal	321	87
Griffin, Alfredo	Oakland	594	169
Newman, Al	Montreal	188	37
Salazar, Argenis	Kansas City	298	73
Thomas, Andre	Atlanta	323	81
Thornton, Andre	Cleveland	401	92
Trammell, Alan	Detroit	574	159
Trevino, Alex	Los Angeles	202	53
Van Slyke, Andy	St Louis	418	113
Wiggins, Alan	Baltimore	239	60
Almon, Bill	Pittsburgh	196	43
Beane, Billy	Minneapolis	183	39
Bell, Buddy	Cincinnati	568	158
Blancalana, Buddy	Kansas City	190	46
Bochte, Bruce	Oakland	407	104
Bochy, Bruce	San Diego	127	32
Bonds, Barry	Pittsburgh	413	92
Bonilla, Bobby	Chicago	426	109
Boone, Bob	California	442	98
Brenly, Bob	San Francisco	472	116

The middle pane shows a list of variables while the right pane shows a partial listing of the data. For a better view of the list of variables, the next figure shows an expanded view.

**Figure 2.8: Expanded View of the Variables Pane**

This image shows the expanded variable selection pane. All variables listed in Figure 2.7 are now checked, indicating they are selected for the analysis.

- Selected variables: Name, Team, nAtBat, nHits, nHome, nRuns, nRBI, nBB, YrMajor, CrAtBat, CrHits, CrHome, CrRuns.

When the data set is displayed, all the variables are selected. You can click any variable to deselect it or, alternatively, uncheck **Select All** and then, holding down the Ctrl key, select the variables that you want to display. If the variables that you want to select are in sequence, you can also click the first variable of interest, hold the Shift key down, and then click the last variable in the list to select all the variables from the first to the last.

As you select or deselect variables, the work area changes to reflect these selections. Let's look at the work area with the variable selection shown in Figure 2.8. To enlarge the work area, click the **Expand** icon (shown below) to expand this area.

**Figure 2.9: Expanding the Work Area**

The screenshot shows the SAS Studio interface. The top bar displays 'Program 1' and 'SASHHELP.BASEBALL'. Below the bar, there are buttons for 'View: Column names', 'Columns', 'Rows', 'Statistics', and 'Filter: (none)'. A status bar indicates 'Total rows: 322 Total columns: 24'. The 'Columns' pane has a 'Select all' checkbox and a 'Name' checkbox with a blue icon. The 'Work Area' pane shows a table with one row:

	Name	Team
1	Allanson, Andy	Cleveland

The work area is now enlarged and shown in Figure 2.10 below.

You can expand or collapse other panes by clicking on the **Expand** icon or, if already expanded, on the **Collapse** icon.

**Figure 2.10: Work Area with Selected Variables**

Total rows: 322 Total columns: 24

	Name	Team	nAtBat	nHits	nHome	nRuns
1	Aldrete, Mike	San Francisco	216	54	2	27
2	Allanson, Andy	Cleveland	293	66	1	30
3	Almon, Bill	Pittsburgh	196	43	7	29
4	Anderson, Dave	Los Angeles	216	53	1	31
5	Armas, Tony	Boston	425	112	11	40
6	Ashby, Alan	Houston	315	81	7	24
7	Backman, Wally	New York	387	124	1	67
8	Baines, Harold	Chicago	570	169	21	72
9	Baker, Dusty	Oakland	242	58	4	25
10	Balboni, Steve	Kansas City	512	117	29	54
11	Bando, Chris	Cleveland	254	68	2	28
12	Barfield, Jesse	Toronto	589	170	40	107
13	Barrett, Marty	Boston	625	179	4	94
14	Bass, Kevin	Houston	591	184	20	83
15	Baylor, Don	Boston	585	139	31	93
16	Beane, Billy	Minneapolis	183	39	3	20
17	Bell, Buddy	Cincinnati	568	158	20	89
18	Bell, George	Toronto	641	198	31	101
19	Bollard, Rafael	Pittsburgh	300	72	0	33
20	Beniquez, Juan	Baltimore	343	103	6	48
21	Bernazard, Tony	Cleveland	562	169	17	88
22	Biancalana, Buddy	Kansas City	190	46	2	24
23	Bilardello, Dann	Montreal	191	37	4	12
24	Bochte, Bruce	Oakland	407	104	6	57

You can use the scroll bars to scroll left or right, up or down. At the top of the work area, you see the number of rows (observations) and columns (variables) in the data set. Here is an enlarged view.

**Figure 2.11: Number of Rows and Columns**

Filter: (none)		
Total rows: 322 Total columns: 24		
	Name	Team
1	Allanson, Andy	Cleveland

You see that there are 322 rows and 24 columns. These numbers will change as you change your variable selections or create filters.

## Sorting Your Data

If you place the cursor on a column heading, an arrow appears. (See Figure 2.12.) [Figure 2.12: Placing the Cursor on a Column Heading](#)

**Figure 2.12: Placing the Cursor on a Column Heading**

nAtBat	nHits
293	66
315	81
479	130
496	141
321	87
594	169
185	37
298	73
323	81
401	92
574	159

By clicking anywhere in the column heading (the column heading nAtBat stands for the number of times at bat), the values are sorted from low to high (an ascending sort). See Figure 2.13 below.

**Figure 2.13: Ascending Sort**

nAtBat ▲	nHits
127	32
127	32
138	31
143	39
151	41
155	44
155	41
160	39
161	36

You can click the triangle to the right of the column name to request a descending sort (see Figure 2.14).

**Figure 2.14: Descending Sort**

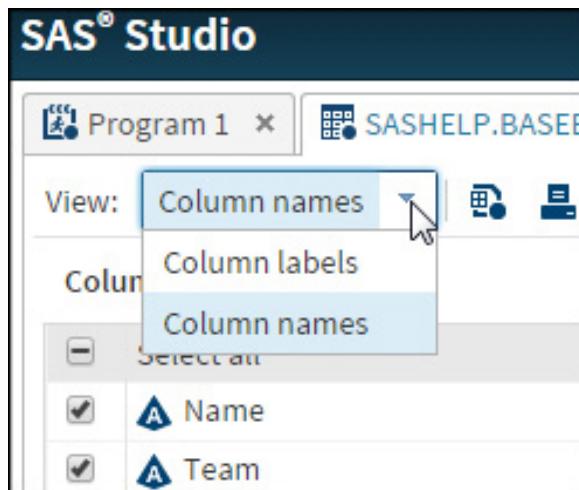
nAtBat ▾	nHits
687	213
680	223
677	238
663	200
642	211
641	198
637	174
633	210
631	170
629	168
627	178

## Switching between Column Names and Column Labels

SAS has been the premier data analysis language since computer variables were limited to only 8 characters (now SAS variable names can be 32 characters long). To fix the short variable name problem, long ago, SAS created labels that could be used to identify variables in addition to the variable names. You must use variable *names* in your programming statements, but you can ask SAS to display the longer, and informative *labels* (that you must create), when columns of data are presented in reports.

You can switch between variable names and variable labels by clicking on the **View** tab.

**Figure 2.15: Selecting Variable Names or Labels**



If you choose column labels, the column names are replaced by labels. (This only works, of course, if the person creating the data set created labels for the variables.) Figure 2.16 shows the effect of switching to column labels for the BASEBALL data set (that did contain labels):

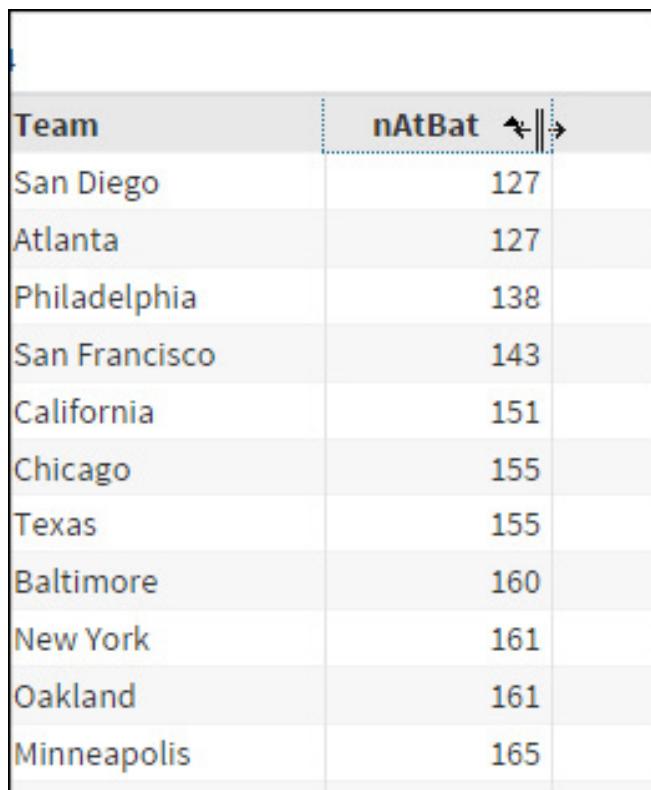
**Figure 2.16: Column Headings Changed to Labels**

Total rows: 322 Total columns: 24					
	Player's Name	Team at the End of 1986	Times at Bat in 1986	Hits in 1986	
1	Fernandez, Tony	Toronto	687	213	
2	Puckett, Kirby	Minneapolis	680	223	
3	Mattingly, Don	New York	677	238	
4	Carter, Joe	Cleveland	663	200	
5	Gwynn, Tony	San Diego	642	211	
6	Bell, George	Toronto	641	198	
7	Dykens, Dan	Chicago, IL	637	174	

## Resizing Tables

By placing the cursor on the dividing line between columns, the pointer changes to double vertical lines, enabling you to then drag the border of the column left or right. In Figure 2.17 below, the nAtBat column was resized (made smaller).

**Figure 2.17: Resizing a Column**



Team	nAtBat
San Diego	127
Atlanta	127
Philadelphia	138
San Francisco	143
California	151
Chicago	155
Texas	155
Baltimore	160
New York	161
Oakland	161
Minneapolis	165

## Creating Filters

Right-clicking a column brings up the following menu.

**Figure 2.18: Adding a Filter**

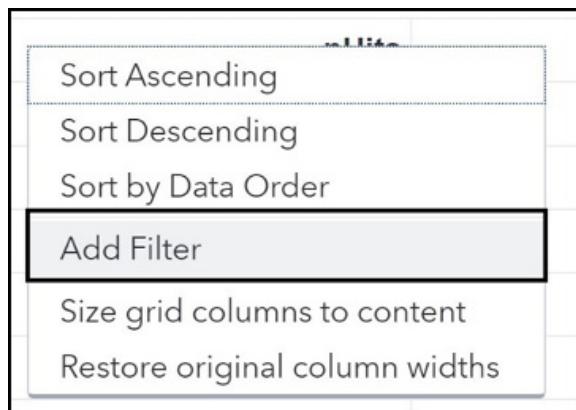
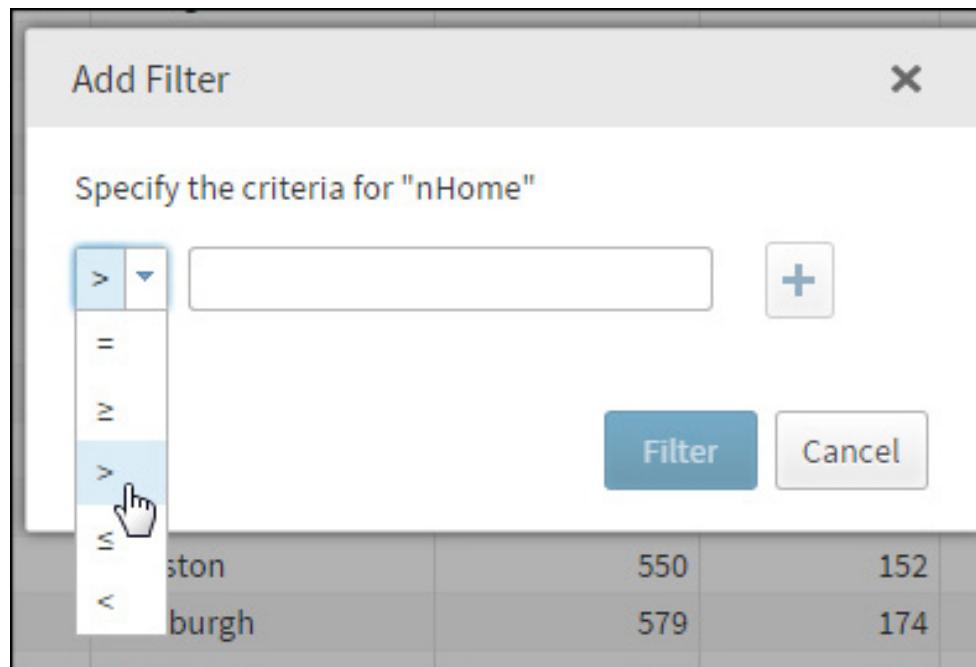


Figure 2.18 shows that you can perform ascending or descending sorts (the same as clicking a column heading). You can do things such as sort in data order, add a filter, or resize columns. *Filters* are very useful in exploring the data because they enable you to subset rows of the table by specifying criteria, such as displaying only rows

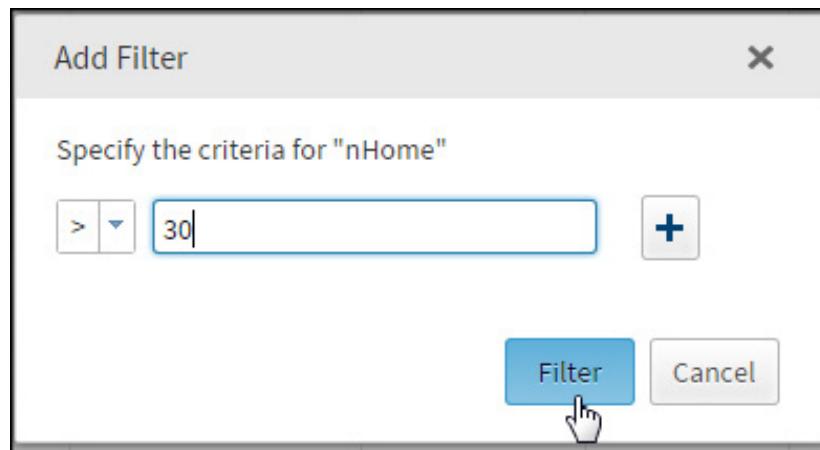
where the number of home runs was 30 or more. The following screen shows how to create a filter. Click **Add Filter** to display the screen shown next.

**Figure 2.19: Adding a Filter (Continued)**



In the drop-down menu, select a logical operator. In this example, you are choosing **greater than**.

**Figure 2.20: Adding a Filter (Continued)**



Next, enter a criterion for the variable that you have selected. If you choose a categorical variable like Team, then the filter dialog box gives you the list of values from which you can select. If you want to

add more conditions, click the + (plus) sign. If you are finished, click **Filter** to complete your request.

**Figure 2.21: Listing of Filtered Data**

Total rows: 322 Total columns: 24 Filtered rows: 11					
	Name	Team	nAtBat	nHits	nHome ▾
1	Barfield, Jesse	Toronto	589	170	40
2	Schmidt, Mike	Philadelphia	552	160	37
3	Kingman, Dave	Oakland	561	118	35
4	Gaetti, Gary	Minneapolis	596	171	34
5	Canseco, Jose	Oakland	600	144	33
6	Mattingly, Don	New York	677	238	31
7	Davis, Glenn	Houston	574	152	31
8	Baylor, Don	Boston	585	139	31
9	Bell, George	Toronto	641	198	31
10	Puckett, Kirby	Minneapolis	680	223	31
11	Parker, Dave	Cincinnati	637	174	31

In Figure 2.21, you see all rows in the BASEBALL data set where the player hit more than 30 home runs.

If you want to remove the filter, click the X next to the filter, as shown in Figure 2.22.

**Figure 2.22: Removing a Filter**

The screenshot shows the SAS Studio interface with the 'SASHelp.BASEBALL' data set open. In the top right, there is a filter bar with the text 'Filter: nHome > 30' and a circled 'X' button to its right. Below the filter bar, the status bar displays 'Total rows: 322 Total columns: 24 Filtered rows: 12'. On the left, a 'Columns' panel lists 'Name' and 'Team' with checkboxes checked. The main pane shows a filtered list of two rows:

	Name	Team
1	Baylor, Don	Boston
2	Kingman, Dave	Oakland

## Conclusion

In this chapter, you saw how to open one of the built-in SASHelp data sets and the various tasks that you can perform using simple

point-and-click operations supplied by SAS Studio. The next step is to create SAS data sets using your own data and then perform operations such as summarizing your data, creating plots and charts, and creating reports in HTML, PDF, or RTF format. You will learn how to accomplish all of these operations in the chapters to follow.

# Chapter 3: Importing Your Own Data

## Introduction

In the last chapter, you saw how to use some of the features of SAS Studio to manipulate data from built-in SAS data sets in SASHELP. In this chapter, you will see how easy it is to create SAS data sets by importing data from Excel workbooks, CSV files, and many other file formats such as Access and SPSS.

## Uploading Data from Your Local Computer to SAS Studio

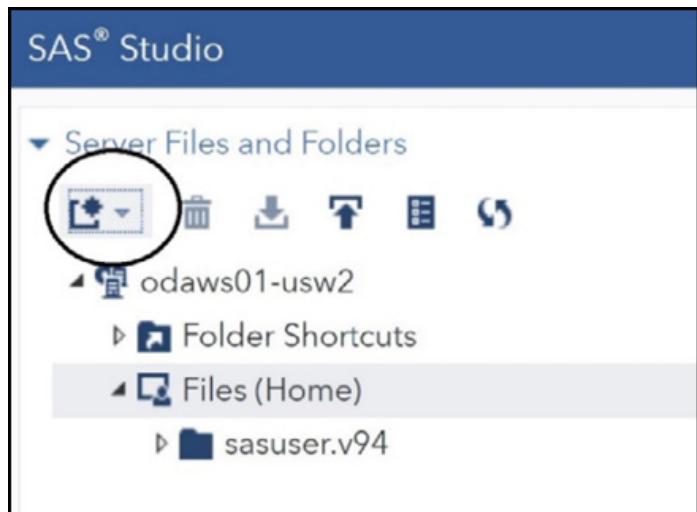
Before you start uploading data files and programs to the cloud, you need to create a folder to hold the files. This is similar to folders on a Windows platform. If you open up the **Server Files and Folders** tab in the navigation pane, you will see the following.

**Figure 3.1: The Server Files and Folders Tab in the Navigation Pane**



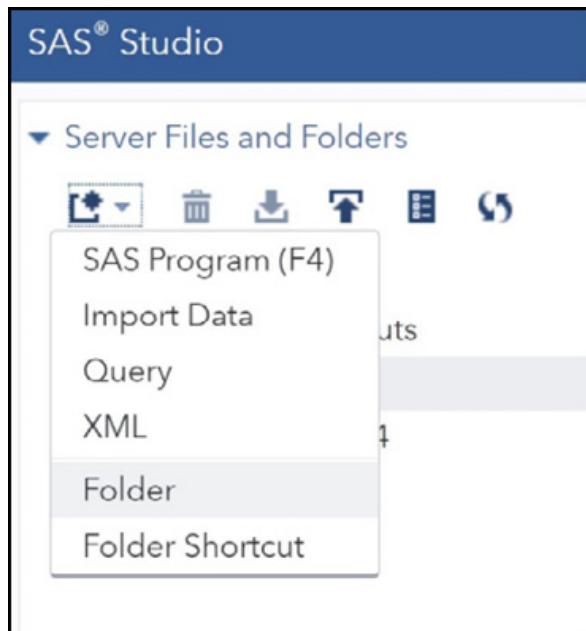
Next, click the **New** icon as shown circled in Figure 3.2.

**Figure 3.2: Click the New Icon**



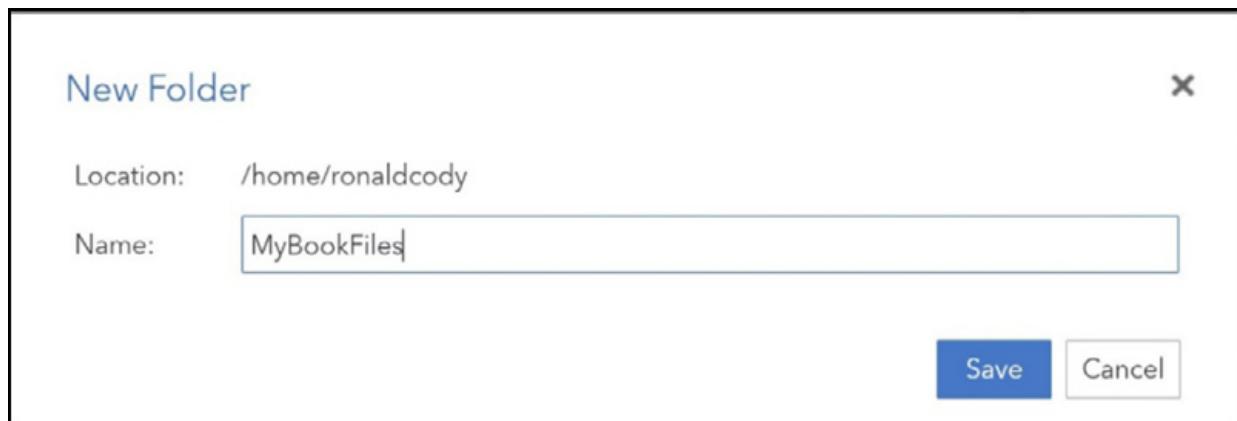
From the pull-down list, select **Folder** (Figure 3.3).

**Figure 3.3: Select Folder in the Pull-Down List**



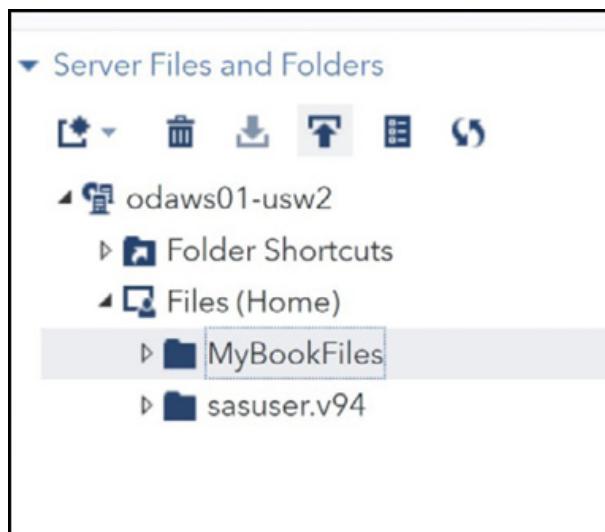
You can now name your new folder. In Figure 3.4, the new folder is called MyBookFiles.

**Figure 3.4: Naming Your Folder**



Click **Save** to finish creating your folder. Your folder now shows up in the list of folders as shown in Figure 3.5.

**Figure 3.5: Your New Folder is Now Included in the List of Files**



Notice that there is a small triangle to the left of the folder name. You can expand a list by clicking this triangle. When you expand the contents of a folder, the triangle points downward and you see the list of files. This action is a “toggle,” that is, it will either expand or collapse the lists.

For this example, we are going to upload an Excel workbook called *Grades.xlsx* to the *MyBookFiles* location. In Figure 3.6, you see a copy of the worksheet that contains student names, ID numbers, some Quiz and exam grades. The first row of the worksheet contains variable names (also known as *column names*). The remaining rows contain data on three students (yes, it was a very small class). The

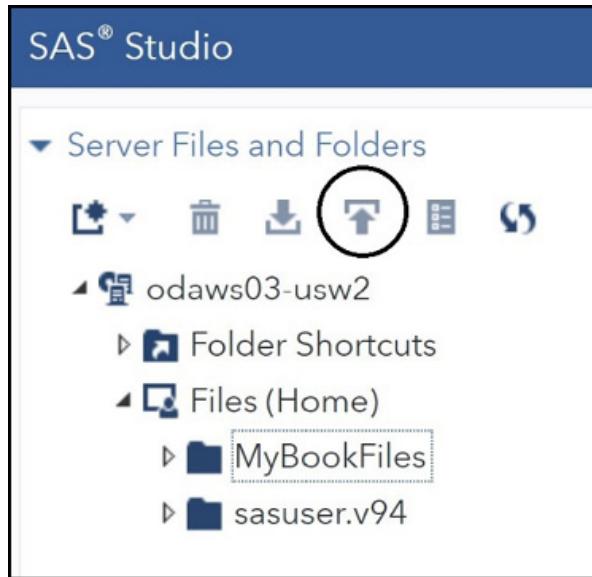
worksheet name was not changed, so it has the default name Sheet1.

**Figure 3.6: Contents of the Spreadsheet Grades.xlsx**

	A	B	C	D	E	F	G	H
1	Name	ID	Quiz1	Quiz2	Midterm	Quiz3	Quiz4	Final
2	Jones	12345	88	80	76	88	90	82
3	Hildebran	22222	95	92	91	94	90	96
4	O'Brien	33333	76	78	79	81	83	80
5								

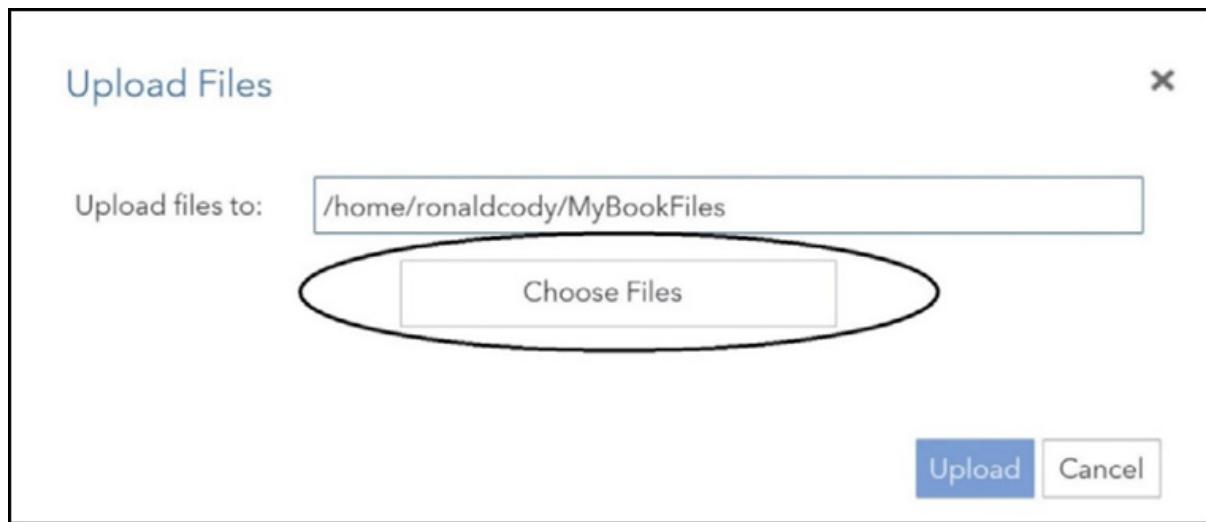
In order to permit SAS Studio to access data files stored on your hard drive, you first need to upload the files to SAS Studio. As shown in Figure 3.7, you click the **Upload** icon.

**Figure 3.7: Click the Upload Icon to Upload One or More Files**



This brings up a screen where you can choose which files you want to upload to SAS Studio.

**Figure 3.8: Choose Files Anywhere on Your Hard Drive**



Click the **Choose Files** box and then you can choose a single file or multiple files. There are two ways to choose multiple files, using methods familiar to Windows users. One way is to hold down the control key and click each file that you want to upload. The other is to click one file, hold down the Shift key, and then click another file. All files from the first to the last will be selected. In Figure 3.9 this latter method was used.

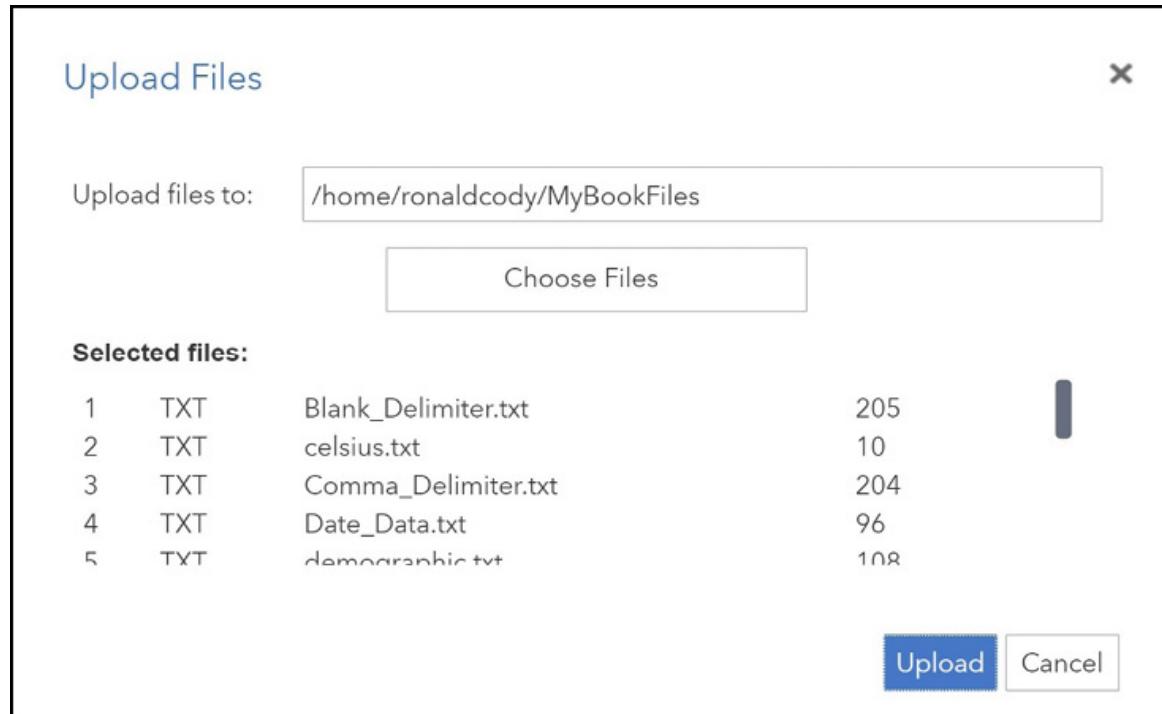
**Figure 3.9: Choose Your File(s)**

Name	Date modified	Type	Size
_Description of Files in this folder.txt	9/4/2015 1:03 PM	Text Document	1 KB
Blank_Delimiter.txt	12/30/2014 10:41 AM	Text Document	1 KB
celsius.txt	12/16/2014 4:05 PM	Text Document	1 KB
Comma_Delimiter.txt	2/27/2015 2:17 PM	Text Document	1 KB
Date_Data.txt	2/4/2015 9:22 AM	Text Document	1 KB
demographic.txt	12/26/2014 1:19 PM	Text Document	1 KB
<input checked="" type="checkbox"/> employee.csv	10/9/2007 9:05 AM	Microsoft Excel C...	1 KB
<input checked="" type="checkbox"/> Grades.csv	3/15/2015 12:57 PM	Microsoft Excel C...	1 KB
<input checked="" type="checkbox"/> Grades.xls	3/17/2015 9:18 AM	Microsoft Excel 97...	25 KB
<input checked="" type="checkbox"/> Grades.xlsx	3/17/2015 8:52 AM	Microsoft Excel W...	9 KB
<input checked="" type="checkbox"/> Grades2.xlsx	3/17/2015 8:38 AM	Microsoft Excel W...	9 KB
health.txt	1/22/2015 8:21 AM	Text Document	1 KB
HeightWeight.txt	1/11/2015 11:10 AM	Text Document	1 KB
<input checked="" type="checkbox"/> Programs.sas	9/4/2015 12:39 PM	SAS System Progr...	40 KB
Quick.txt	4/20/2015 11:57 AM	Text Document	1 KB
Tab_Delimiter.txt	1/6/2015 9:02 AM	Text Document	1 KB
<input checked="" type="checkbox"/> taxes.txt	1/21/2015 4:20 PM	Text Document	1 KB

For this example, we are uploading all the files that will be used in examples in this book. Notice the Excel workbook Grades.xlsx is included in this group of files. Once you have selected the files that you want to upload, click the box labeled **Upload**.

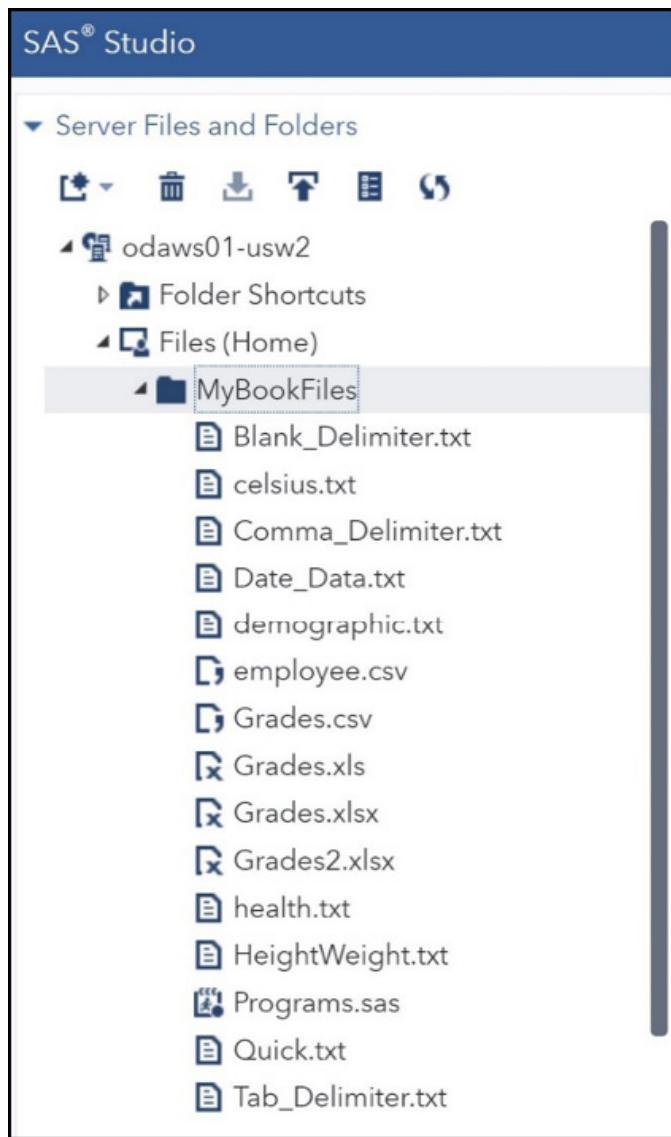
Although the maximum file space on the Studio cloud is 5 megabytes, you cannot upload more than 1 megabyte at a time.

**Figure 3.10: Getting Ready to Upload Selected Files**



In Figure 3.11, you see all the selected files in the MyBookFiles folder. Notice that the triangle to the left of MyBookFiles points downward because you clicked it to reveal the individual files.

**Figure 3.11: Files Uploaded to MyBookFiles**



You can now use SAS Studio to analyze any of these files or write a SAS program of your own using one or more of these files.

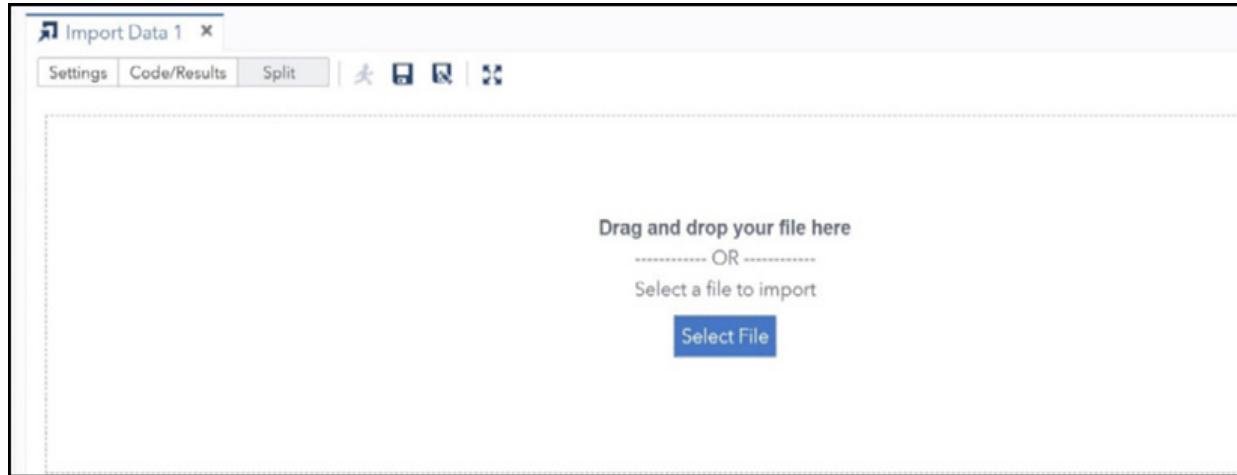
Before you can work with any data such as an Excel workbook or a CSV file, you need to convert it into a SAS data set. To accomplish this, you select the **Import** task that can be found under **Tasks and Utilities** task. (See Figure 3.12 below.)

**Figure 3.12: Selecting the Tasks and Utilities Task**



Double click **Import Data** to begin converting the Excel workbook Grades.xlsx into a SAS data set.

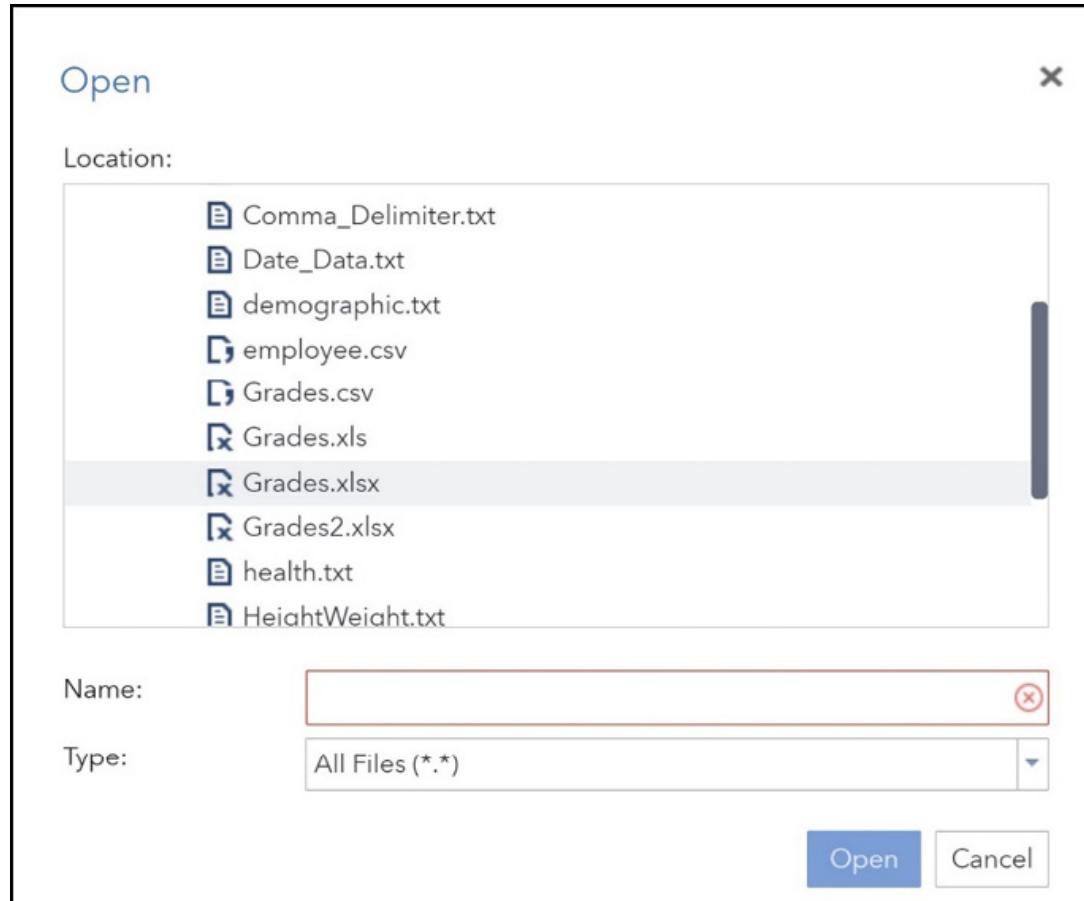
**Figure 3.13: Selecting a File to Import**



You have two ways to select which file you want to import. One is to click the **Select File** box on the right side of the screen—the other method is to click the **Server Files and Folders** tab in the navigation pane (on the left), find the file, and drag it to the drag and drop area.

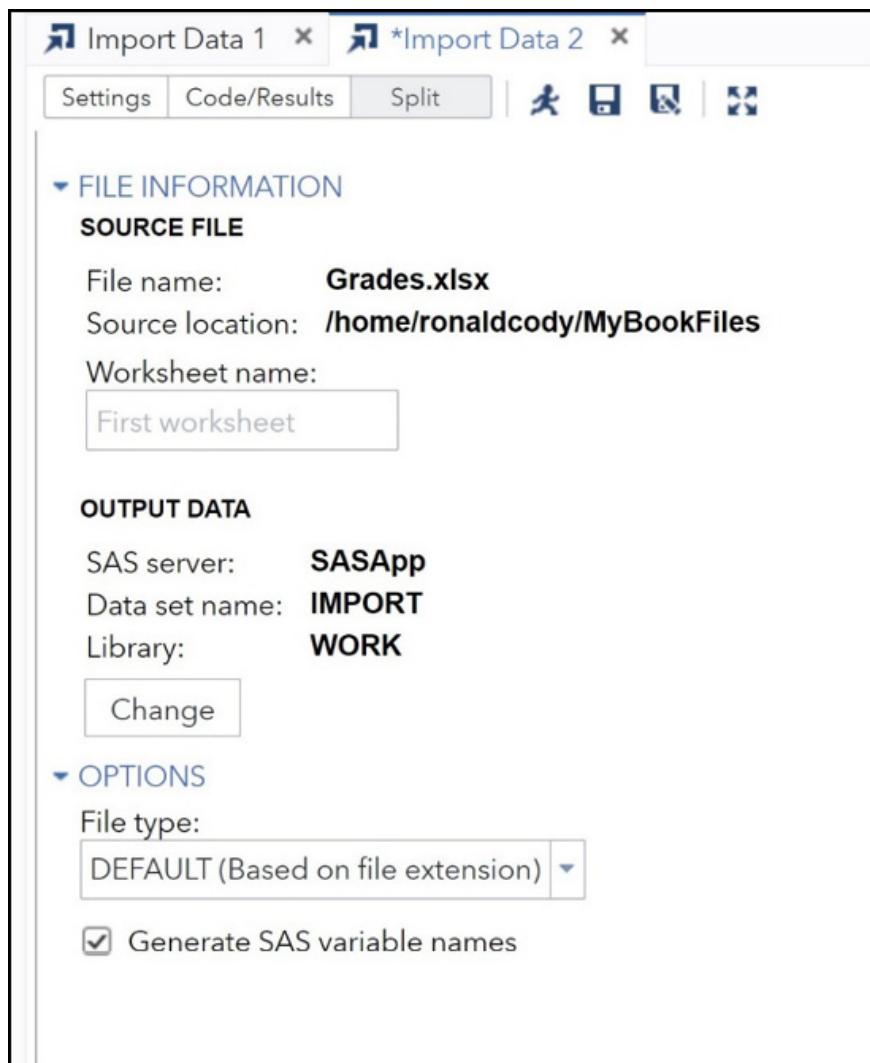
Clicking **Select File** brings up a window where you can select a file to import. Here it is.

**Figure 3.14: Selecting the Grades Workbook**



Click **Grades.xlsx** and then click **Open**. This brings up a split window showing options for the import on the top half and the SAS program that will be generated by these options in the bottom half.

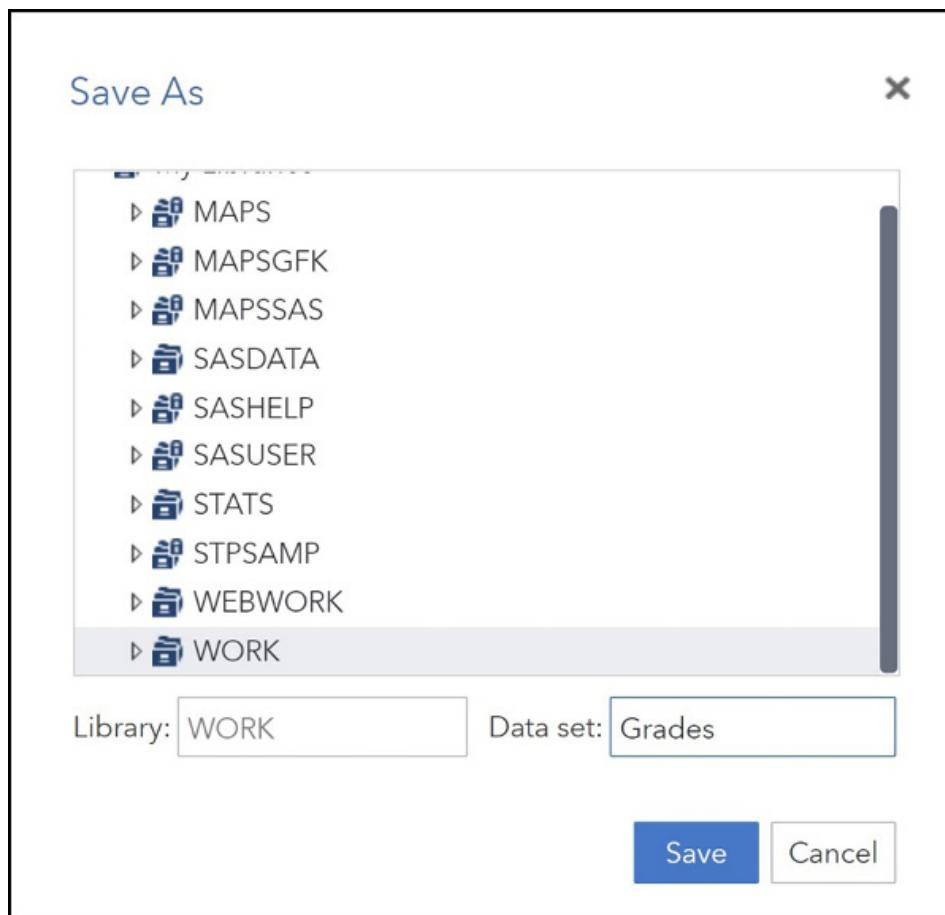
**Figure 3.15: File Information**



You might need to place your cursor at the bottom of this window to expand the window, enabling you see the selections for naming your output data set. There is an alternative. At the very top of the window are three tabs, labeled **Settings**, **Code/Results**, and **Split**. If you click the **Settings** tab, you do not have to expand the window manually.

Because you only have one worksheet, you do not have to enter a worksheet name. You probably want to change the name of the output (SAS) data set from the default name (IMPORT) to something more meaningful. You might also want to change the library where the SAS data set will be stored from the default WORK library. Clicking the **Change** button brings up a list of SAS libraries (below) and enables you to do both.

**Figure 3.16: Selecting a File Name and a SAS Library**



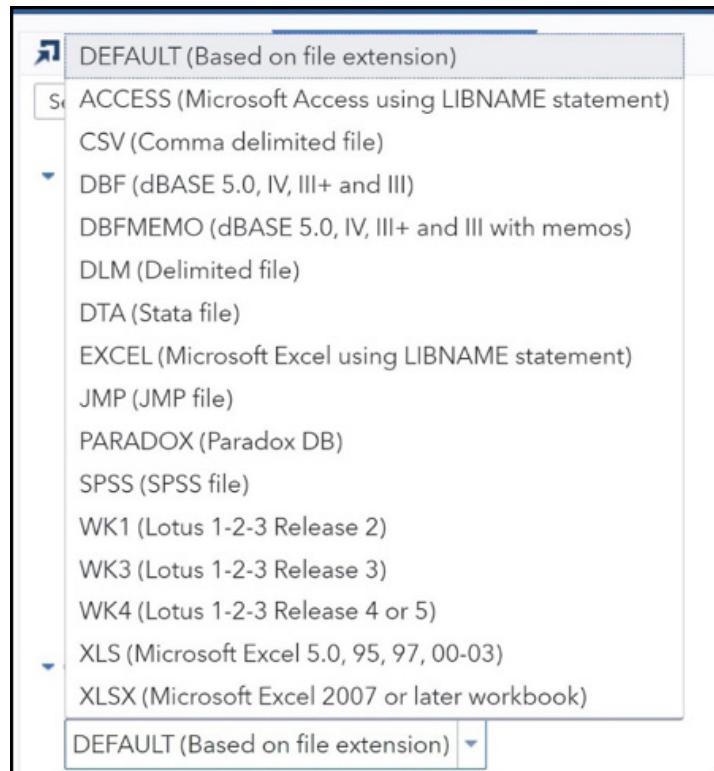
For this example, you are naming the file Grades and placing it in the default WORK library. This library is a temporary library whose contents are deleted when you exit your SAS session. Later in this chapter you will see how to create a permanent SAS library where your files will not be deleted when you close your SAS session.

Notice that you did not have to tell SAS Studio that you were importing an Excel workbook—it figured it out by the file extension (XLSX). SAS Studio will typically use the file extension to figure out how to import data. Because the first row of the spreadsheet contains variable names, leave the check on the **Generate SAS variable names** option. This tells the import utility to use the first row of the worksheet to generate variable names.

If you have a nonstandard file extension or if you prefer to manually select a file format, then you can use the drop-down list displayed in

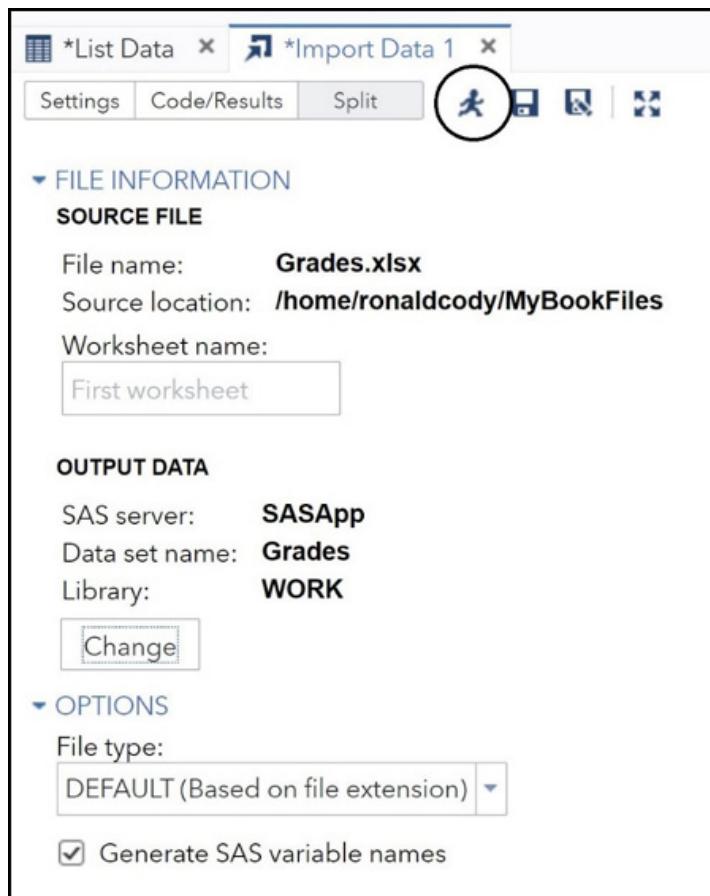
Figure 3.17 to instruct SAS how to convert your file.

### Figure 3.17: List of Supported File Types



When all is ready, click the **Run** icon (Figure 3.18).

### Figure 3.18: The Run Icon



SAS Studio will now list information about the Grades data set. Selections of the output are shown below.

**Figure 3.19: Edited Output from the Import Task**

The CONTENTS Procedure			
Data Set Name	WORK.GRADES	Observations	3
Member Type	DATA	Variables	8
Engine	V9	Indexes	0
Created	10/16/2020 15:48:52	Observation Length	72
Last Modified	10/16/2020 15:48:52	Deleted Observations	0
Protection		Compressed	NO
Data Set Type		Sorted	NO
Label			
Data Representation	SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64		
Encoding	utf-8 Unicode (UTF-8)		

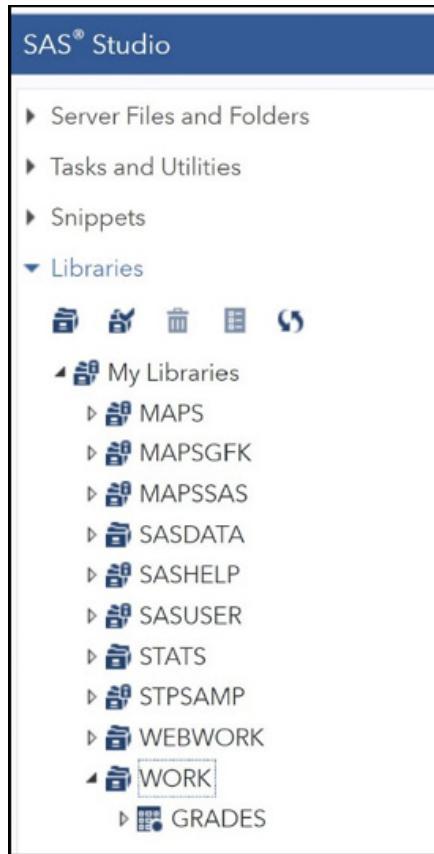
Alphabetic List of Variables and Attributes						
#	Variable	Type	Len	Format	Informat	Label
8	Final	Num	8	BEST.		Final
2	ID	Num	8	BEST.		ID
5	Midterm	Num	8	BEST.		Midterm
1	Name	Char	10	\$10.	\$10.	Name
3	Quiz1	Num	8	BEST.		Quiz1
4	Quiz2	Num	8	BEST.		Quiz2
6	Quiz3	Num	8	BEST.		Quiz3
7	Quiz4	Num	8	BEST.		Quiz4

At the bottom of this output, you see a list of the variable names, whether they are stored as numeric or character, along with some other information that we don't need at this time. Notice that the import utility correctly read Name as character and the other variables as numeric.

## Listing the SAS Data Set

A quick way to see a listing of the Grades data set is to select the **Libraries** tab in the navigation pane and select **My Libraries**.

**Figure 3.20: The Libraries Tab**



Expand the Work library and double-click **Grades**. It looks like Figure 3.21.

**Figure 3.21: Data Set Grades in the Work Library**

A screenshot of the SAS Studio interface showing the "WORK.GRADES" data set in a grid view. The top navigation bar includes "Program 1", "Import Data 1", and "WORK.GRADES". The grid view displays the following data:

Name	ID	Quiz1	Quiz2	Midterm	Quiz3	Quiz4
1 Jones	12345	88	80	76	88	90
2 Hildebrand	22222	95	92	91	94	90
3 O'Brien	33333	76	78	79	81	83

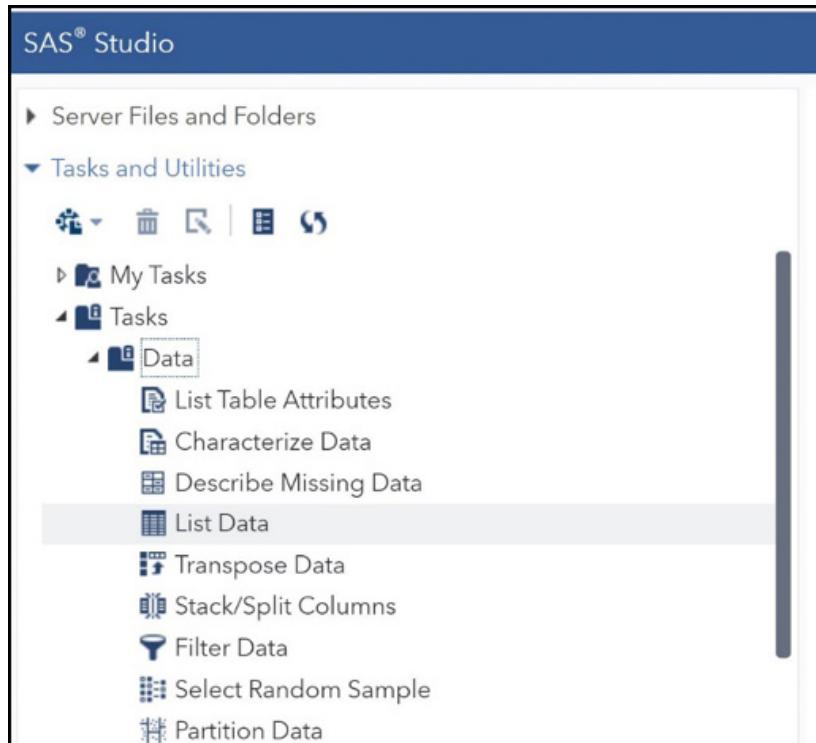
You can use your mouse to scroll to the right to see the rest of the table.

In this book, as well as in SAS Studio, you will see the terms SAS data set and table used interchangeably as well as these other equivalent terms: variables are also called columns and

observations are called rows.

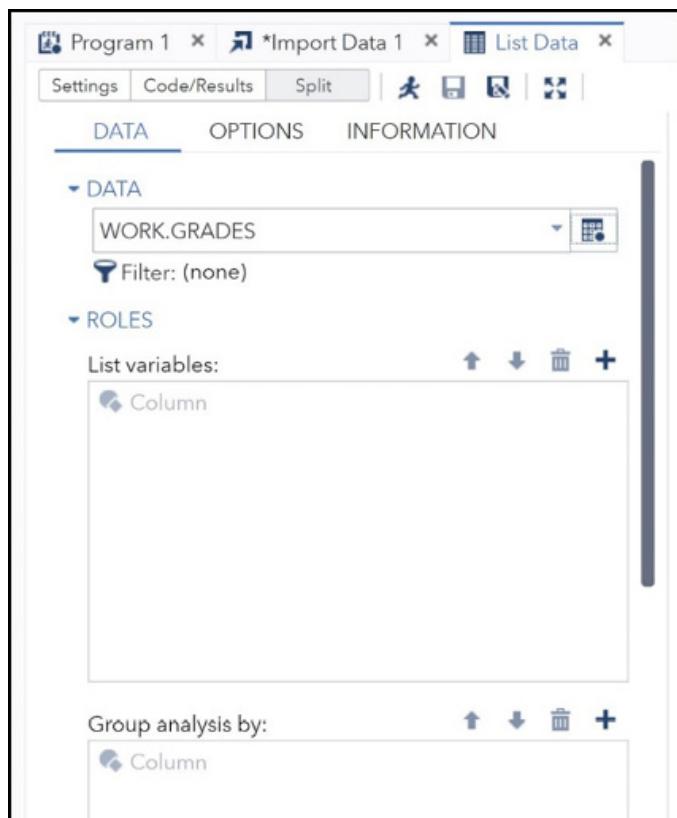
To create a nicer looking report, click the **Tasks and Utilities** tab of the navigation pane and select **Data** followed by **List Data**, as in Figure 3.22.

**Figure 3.22: The List Data Task**



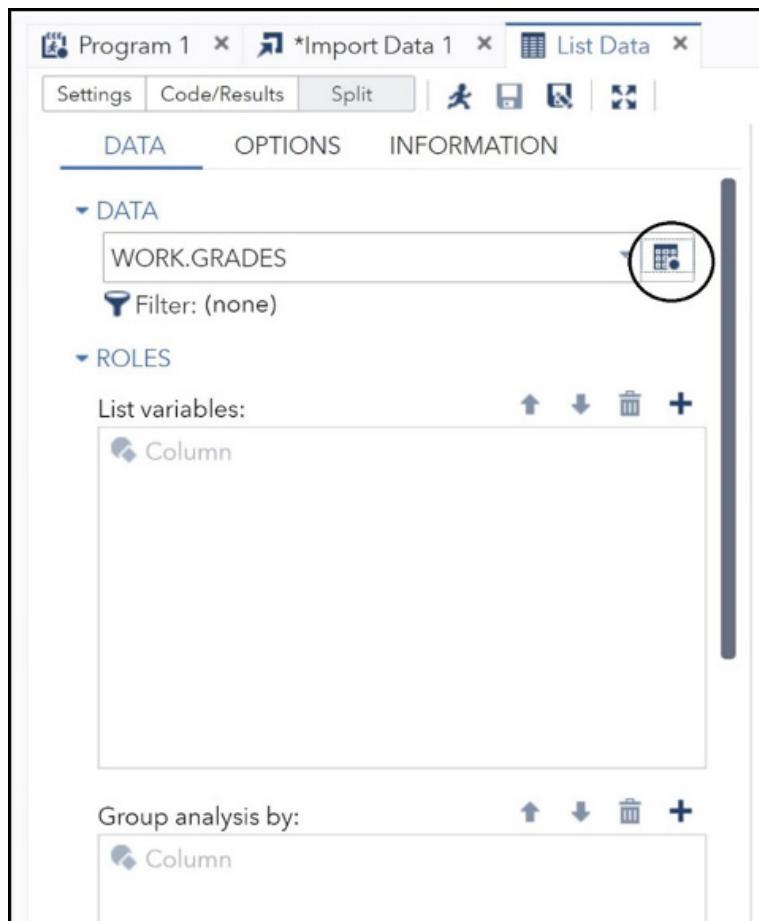
Double-click **List Data**. The screen that appears resembles Figure 3.23.

**Figure 3.23: List Data Task**

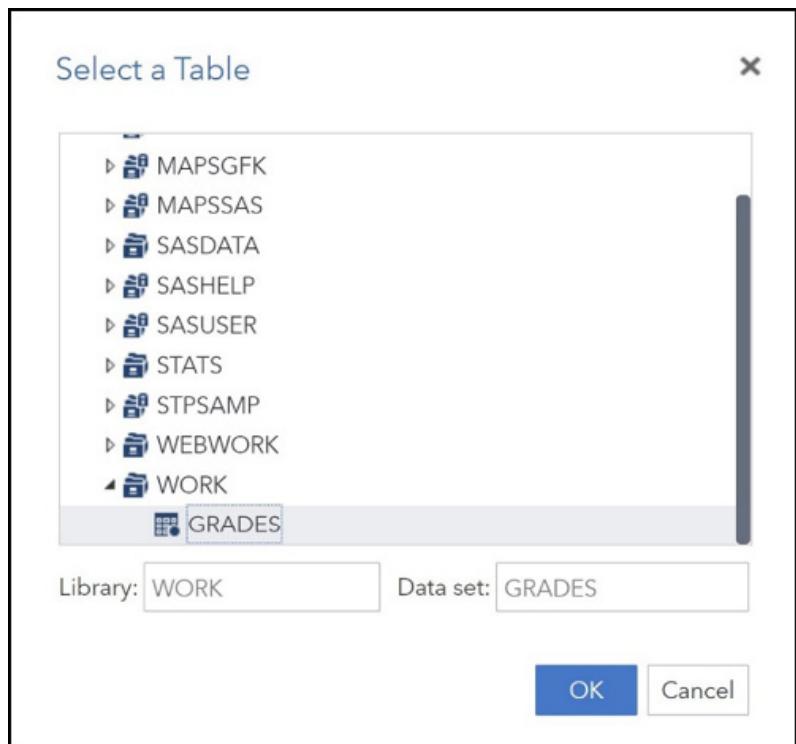


Click the icon at the far right side of the box labeled **DATA** (circled in the figure below).

**Figure 3.24: Select Table**

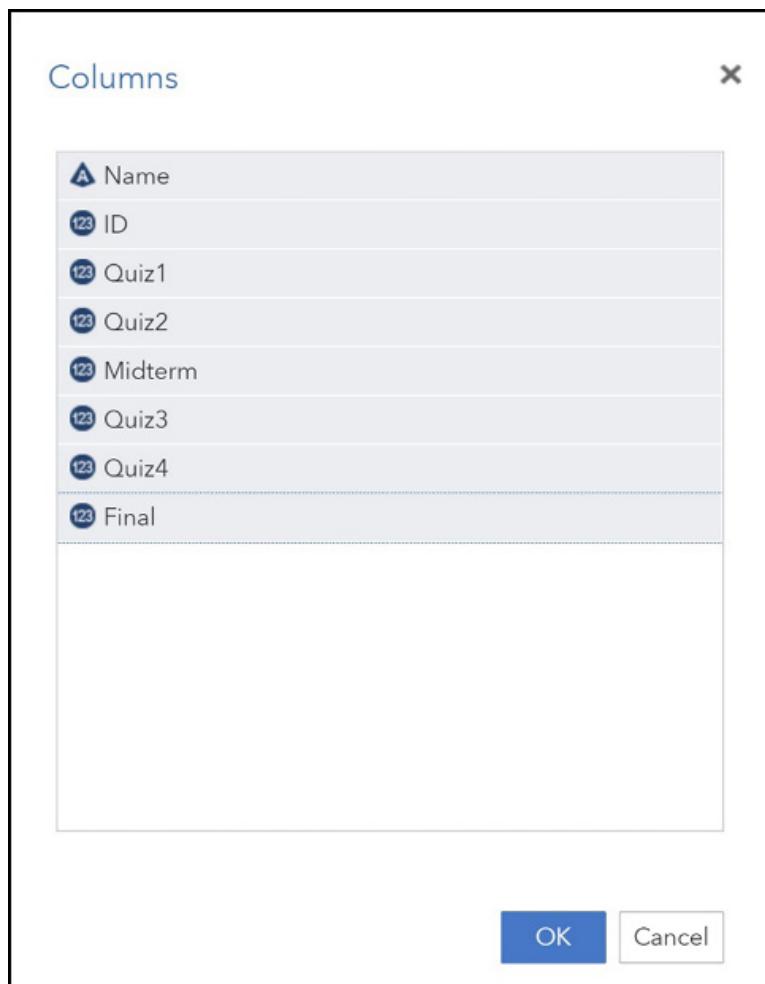


**Figure 3.25: Select Grades in the WORK Library**



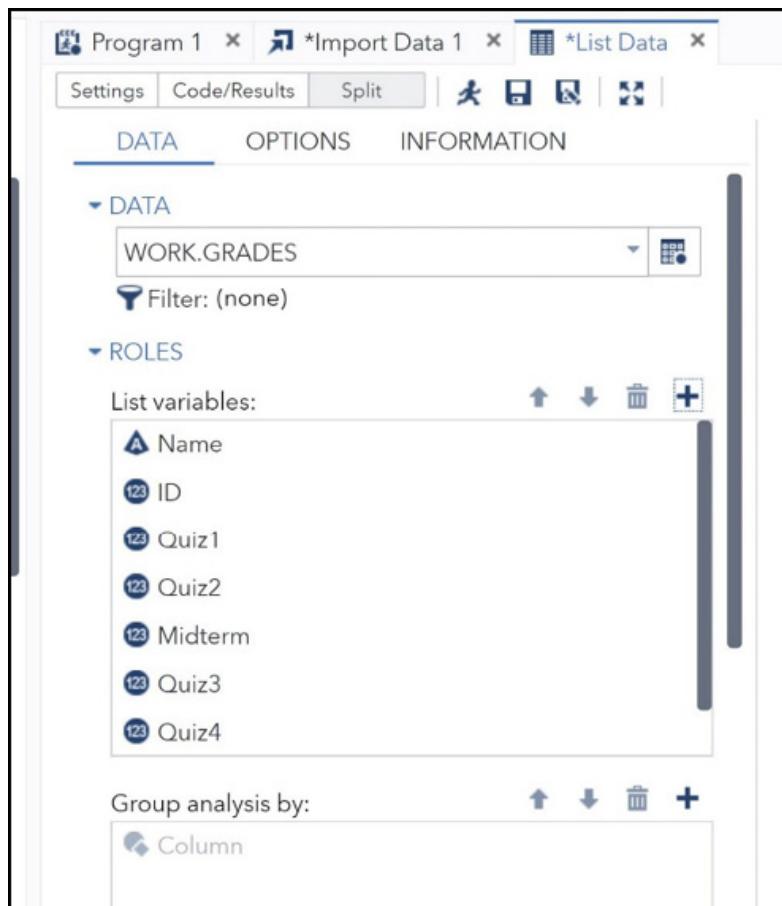
Click **OK**. Your next task is to select which variables you want in the listing. Under **ROLES**, you see four boxes. At this time, you are only interested in the first two: **List variables** and **Group analysis by**. Click the plus sign next to **List Variables**. A list of variables from the Grades data set will show up (Figure 3.26).

**Figure 3.26: Select Variables**



Select the variables that you want to see in your listing. If you want to see all of the variables, first click **Name**. Then hold down the Shift key and click **Final**. Click **OK**. Your screen now should look like Figure 3.27.

**Figure 3.27: Tables and Variables Selected**



Click the **Run** icon and the listing of the Grades data set appears as shown in Figure 3.28.

**Figure 3.28: Listing of Data Set Grades**

List Data for WORK.GRADES									
Obs	Name	ID	Quiz1	Quiz2	Midterm	Quiz3	Quiz4	Final	
1	Jones	12345	88	80	76	88	90	82	
2	Hildebrand	22222	95	92	91	94	90	96	
3	O'Brien	33333	76	78	79	81	83	80	

## Importing an Excel Workbook with Invalid SAS Variable Names

What if your Excel worksheet has column headings that are not valid SAS variable names?

Valid SAS data set names are up to 32 characters long. In many versions of SAS, the first character must be a letter or underscore—the remaining characters can be letters or digits. However, there is a system option called VALIDVARNAME=ANY that allows variable names that are invalid in many versions of SAS. This is the default in the current Studio version. It might be better to use the statement:

```
Options validvarname=V7;
```

so that your variable names are compatible with systems where this option is set to V7.

For example, take a look at the worksheet Grades2 shown in Figure 3.29.

**Figure 3.29: Listing of Excel Workbook Grades2**

	A	B	C	D	E	F	G	H
1	Stuent Name	ID	Quiz 1	Quiz 2	Mid Term	Quiz 3	Quiz 4	2015Final
2	Jones	12345	88	80	76	88	90	82
3	Hildebrand	22222	95	92	91	94	90	96
4	O'Brien	33333	76	78	79	81	83	80
5								

Most of the column headings in this spreadsheet are not valid SAS variable names when VALIDVARNAME is set to V7. Six of them contain a blank and the last column (2015Final) starts with a digit. What happens when you import this worksheet? Because you now know how to use the **Import Data** task, it is not necessary to describe the import task again. All you really need to see is the final list of variables in the data set. Here they are.

**Figure 3.30: Variable Names in the Grades2 SAS Data Set**

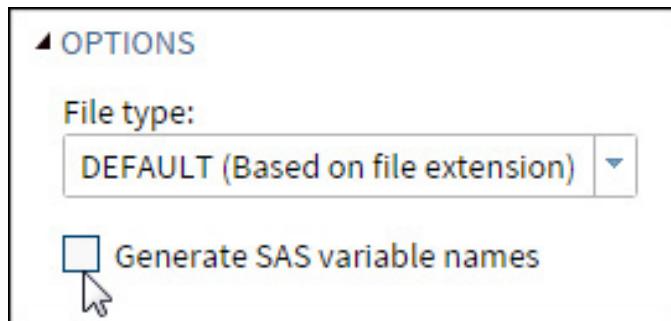
Alphabetic List of Variables and Attributes						
#	Variable	Type	Len	Format	Informat	Label
2	ID	Num	8	BEST.		ID
5	Mid_Term	Num	8	BEST.		Mid Term
3	Quiz_1	Num	8	BEST.		Quiz 1
4	Quiz_2	Num	8	BEST.		Quiz 2
6	Quiz_3	Num	8	BEST.		Quiz 3
7	Quiz_4	Num	8	BEST.		Quiz 4
1	Stuent_Name	Char	10	\$10.	\$10.	Stuent Name
8	_2015Final	Num	8	BEST.		2015Final

SAS replaced all the blanks with underscores and added an underscore as the first character in the 2015Final name to create valid SAS variable names. If you leave the ODA default setting of VALIDVARNAME=ANY, the underscores will not be added.

## Importing an Excel Workbook That Does Not Have Variable Names

What if the first row of your worksheet does not contain variable names (column headings)? You have two choices. First, you can edit the worksheet and insert a row with variable names (probably the best option). The other option is to uncheck **variable names** in the **OPTIONS** section in the **Import** window (see Figure 3.31), and let SAS create variable names for you.

**Figure 3.31: Uncheck the Variable Names Option**



Here is the result.

**Figure 3.32: Variable Names Generated by SAS**

Alphabetic List of Variables and Attributes						
#	Variable	Type	Len	Format	Informat	Label
1	A	Num	8	BEST.		A
2	B	Char	1	\$1.	\$1.	B
3	C	Num	8	MMDDYY10.		C
4	D	Num	8	BEST.		D
5	E	Num	8	BEST.		E
6	F	Char	1	\$1.	\$1.	F

SAS used the column identifiers (A through F) as variable names. You can leave these variable names as is or change them using DATA step programming. Another option is to use PROC DATASETS, a SAS procedure that enables you to alter various attributes of a SAS data set without having to create a new copy of the data set.

When you import a CSV file without variable names, you will see variable names generated by SAS such as, VAR1, VAR2, etc.

## Importing Data from a CSV File

CSV (*comma-separated values*) files are a popular format for external data files. As the name implies, CSV files use commas as data delimiters. Many websites enable you to download data as CSV files. As with Excel workbooks, your CSV file might or might not contain variable names at the beginning of the file. If the file does contain variable names, be sure the **Generate SAS variable names** options box is checked; if not, uncheck this option.

For example, take a look at the CSV file called Grades.csv in Figure 3.33 below.

**Figure 3.33: CSV File Grades.csv**

Grades.csv - Notepad						
File Edit Format View Help						
Name, ID, Quiz1, Quiz2, Midterm, Quiz3, Quiz4, Final						
Jones	12345	88	80	76	88	90
Hildebrand	22222	95	92	91	94	90
O'Brien	33333	76	78	79	81	83
						80

This CSV file contains the same data as the Excel workbook Grades.xlsx. Notice that variable names are included in the file. You can import this file and create a SAS data set, using the same steps that you used to import the Excel workbook. The import facility will automatically use the correct code to import this data file because of the CSV file extension. The resulting SAS data set is identical to the one shown in Figure 3.28.

## Conclusion

You have seen how easy it is to import data in a variety of formats and create SAS data sets. Even experienced programmers (at least the ones this author knows) would prefer to use the import data utility to convert external data to SAS data sets rather than writing their own code.

# Chapter 4: Creating Reports

## Introduction

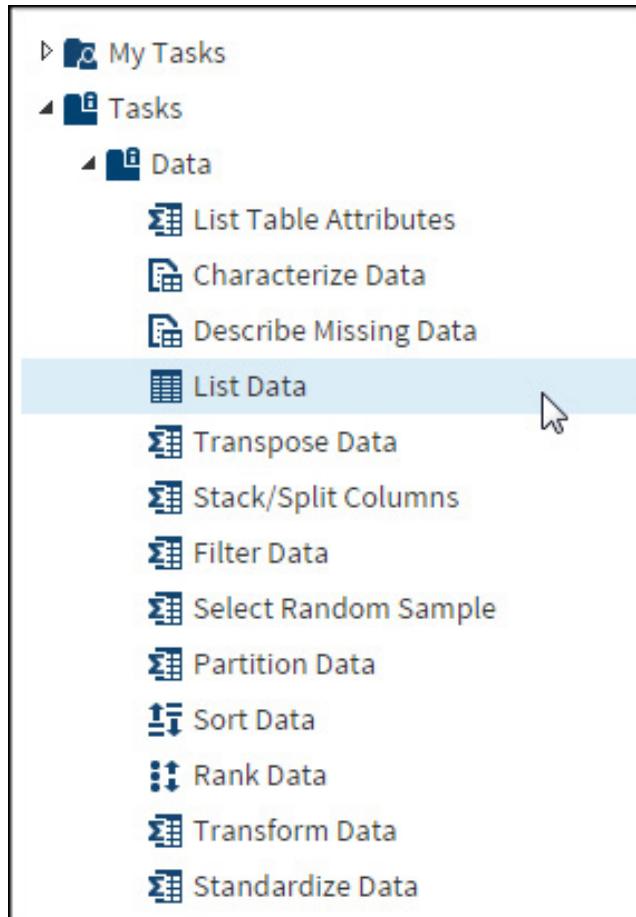
In the last chapter, you saw how to use the **Import Data** facility on the **Utilities** tab to import data. In this chapter, you will see how to use several of the most useful tasks as well as the Query tool on the **Utilities** tab.

It is worth repeating this note from the previous chapter: In this book, as well as in SAS Studio, you will see the terms SAS data set and table used interchangeably as well as these other equivalent terms: variables are also called columns and observations are called rows.

## Using the List Data Task to Create a Simple Listing

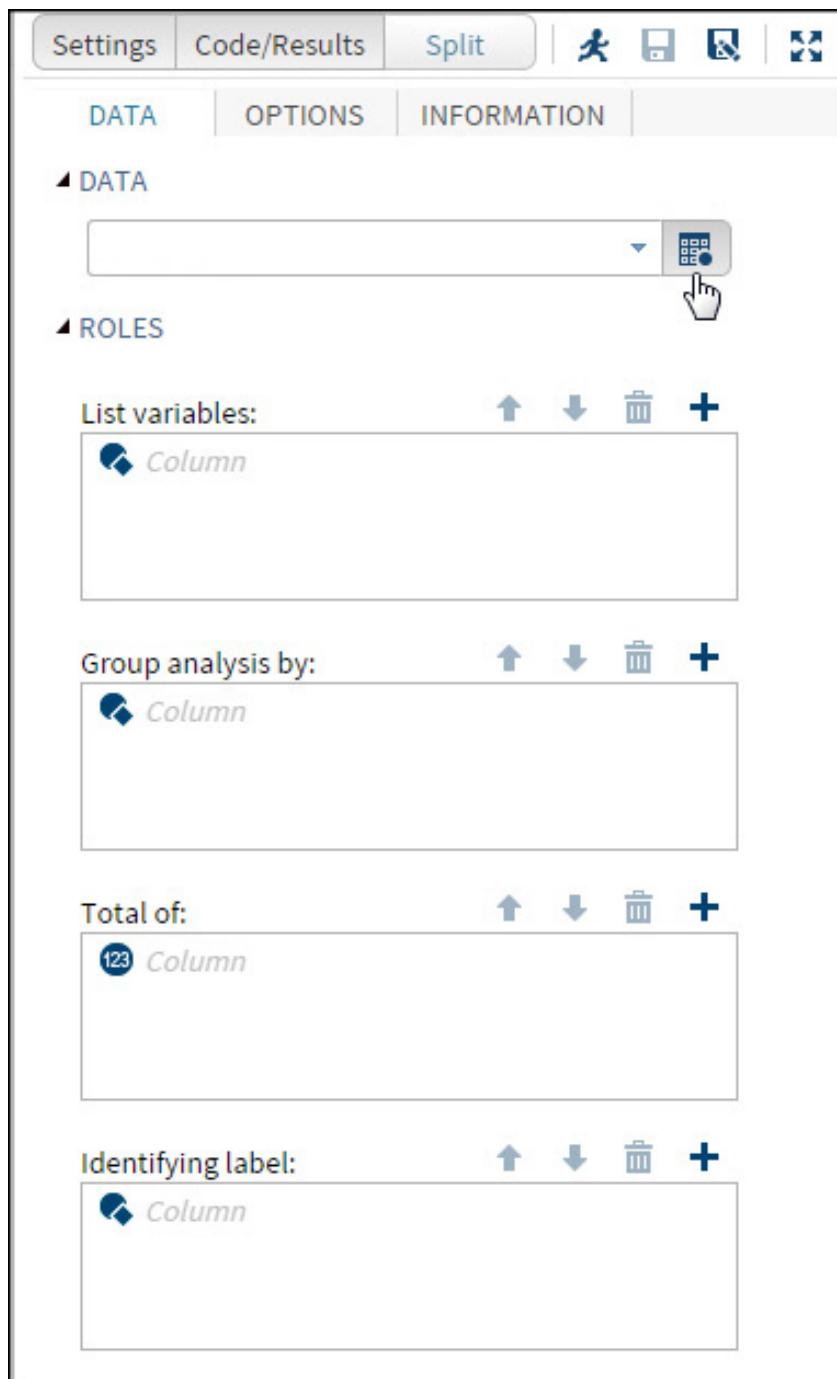
SAS Studio has dozens of built-in tasks. If you expand the **Tasks** tab and click the **Data** tab, you will see the following.

**Figure 4.1: Data Tasks**



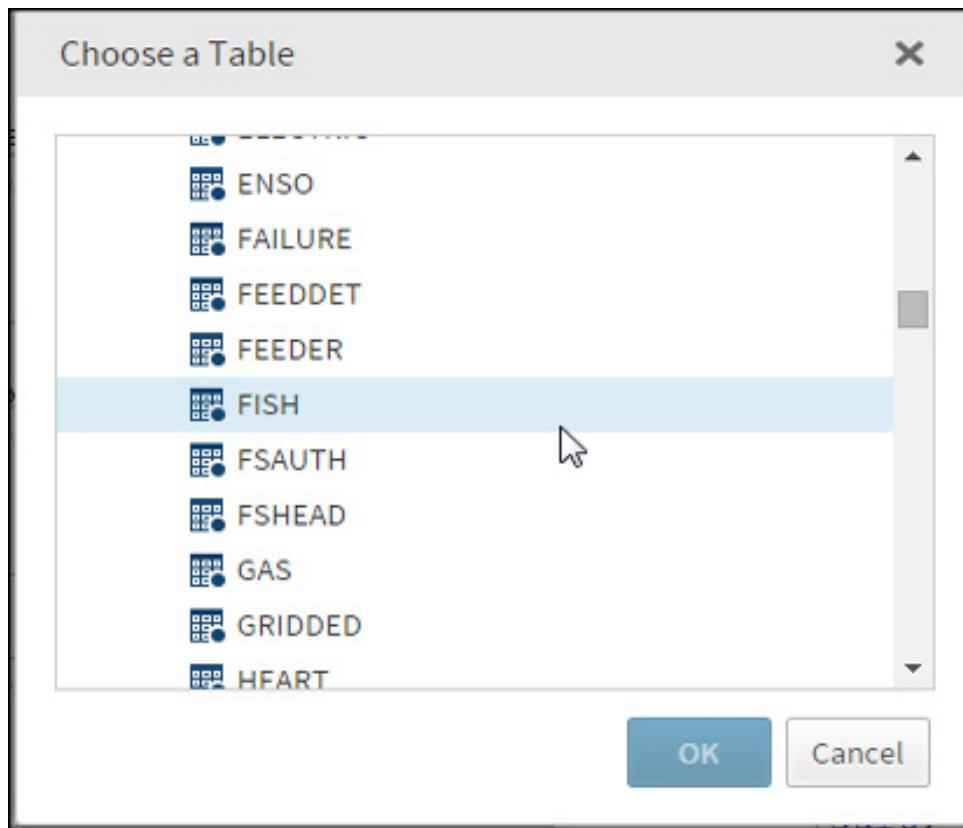
You can use the **List Data** task to create a listing of your data set. To demonstrate this, we are going to use the built-in SASHELP data set called Fish. This data set contains information about several species of fish, including weight, length, and width. To create a listing of this data set, expand the list of **Data** tasks and double-click **List Data**. This brings up the screen shown in Figure 4.2.

**Figure 4.2: The List Data Task Settings Screen**



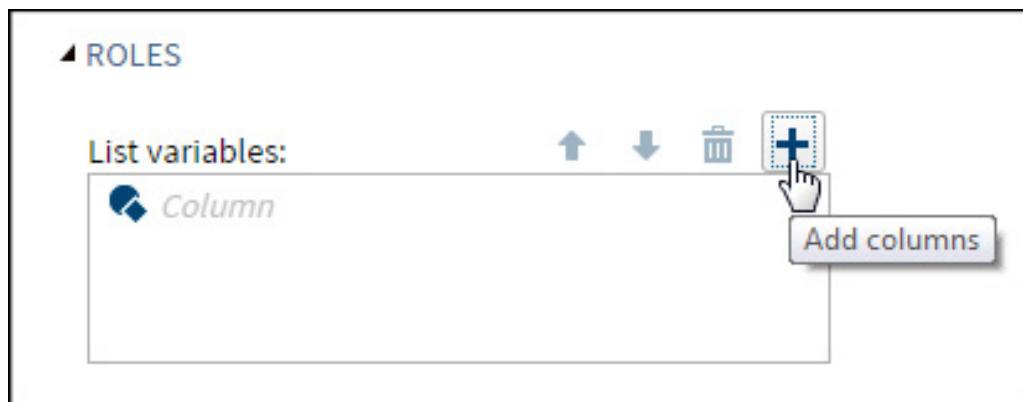
You can click the icon at the top right part of this screen (shown in the figure with a hand pointing to it) to select the library and data set you want to list. Because you want a listing of SASHELP.Fish, select this data set and click **OK**.

**Figure 4.3: Selecting the Fish Data Set in the SASHELP Library**



The next step is to click the plus sign to select which variables you want to include in your listing. (See Figure 4.4.)

**Figure 4.4: Adding Columns (Variables)**



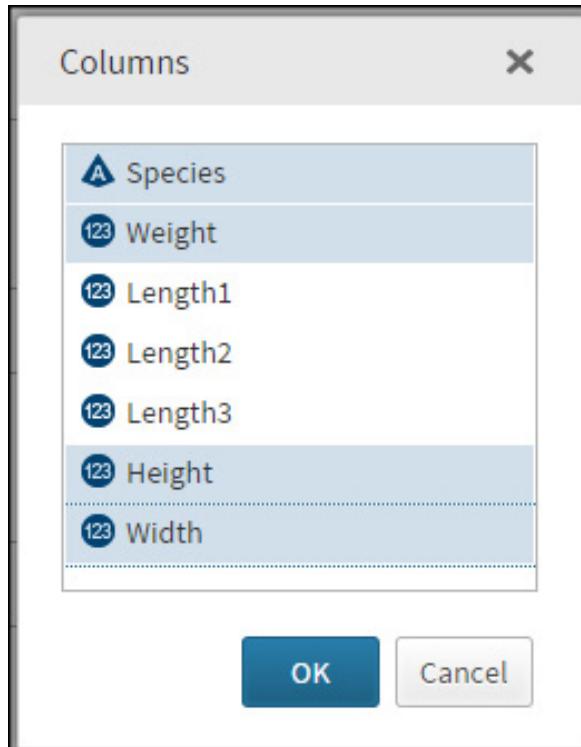
When you see the list of variables, you can select them in the usual way. (See the instructions below.)

To select variables from a list, use one of these two methods: 1) Hold the Ctrl key down and select the variables that you want; or 2) Click one variable, hold the Shift key down, and click another

variable—all the variables from the first to the last will be selected.

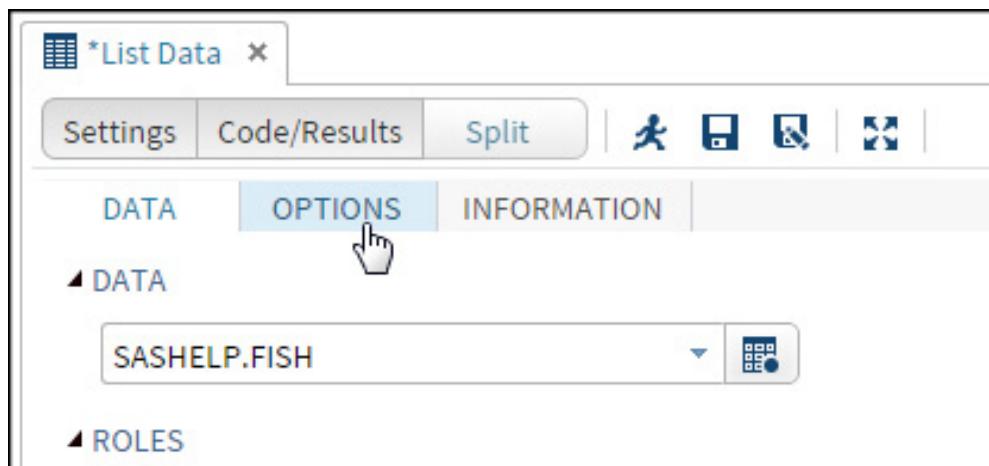
In this example, you are selecting **Species**, **Weight**, **Height**, and **Width**.

**Figure 4.5: Selecting Variables to List**



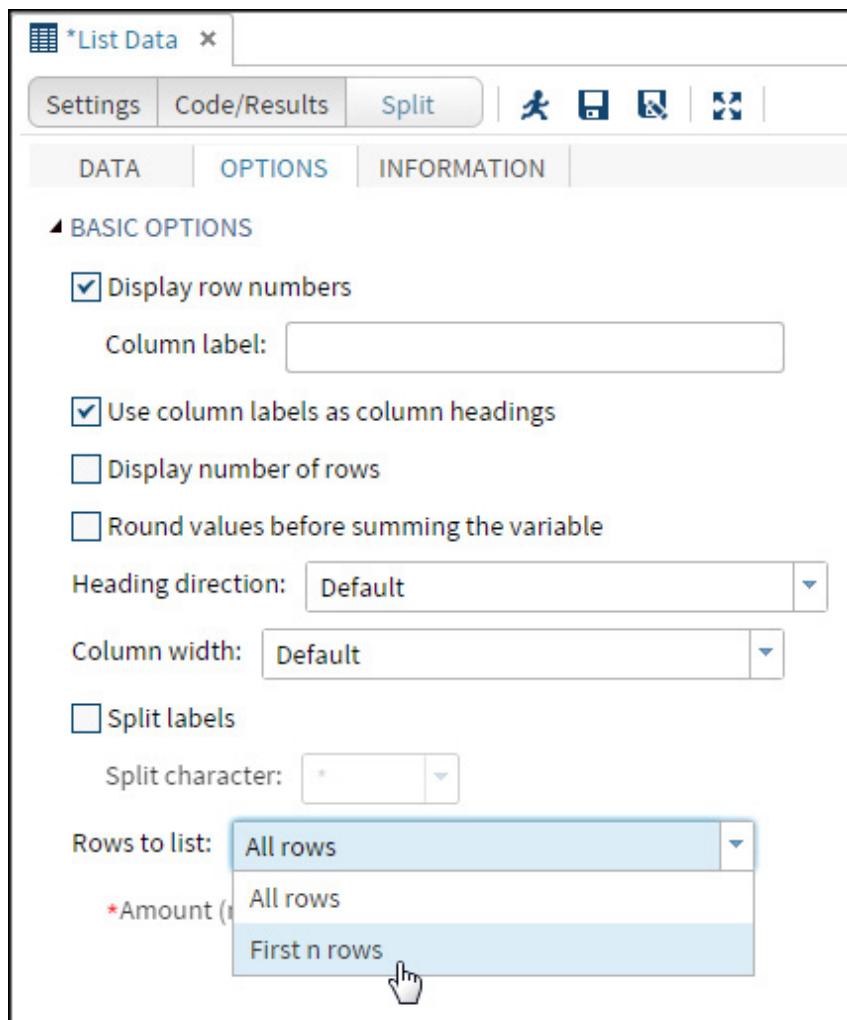
Click **OK** when you are finished. You can create the listing now or click the **OPTIONS** tab to customize the listing.

**Figure 4.6: Using the OPTIONS Tab to Customize the Listing**



Here is the list of options available for the **List Data** task.

**Figure 4.7: Options for the List Data Task**



You can check or uncheck the **Display row numbers** box to include the Obs column in the listing or to omit it. You have a choice of using column labels or column names in the listing. If you are a programmer, then you will probably want to see column names as headings. If you are creating the listing for a report, then you will probably want to see column labels.

The option displayed at the bottom of Figure 4.7 gives you the choice of listing all the rows of the table or the first  $n$  rows. In this example, you want to see the first seven rows of the Fish data set. This is shown in the next figure.

**Figure 4.8: Requesting the First Seven Rows to Be Displayed**

▲ BASIC OPTIONS

Display row numbers  
Column label:

Use column labels as column headings

Display number of rows

Round values before summing the variable

Heading direction:

Column width:

Split labels  
Split character:

Rows to list:

\*Amount (n):

Clicking the **Run** icon generates the following listing.

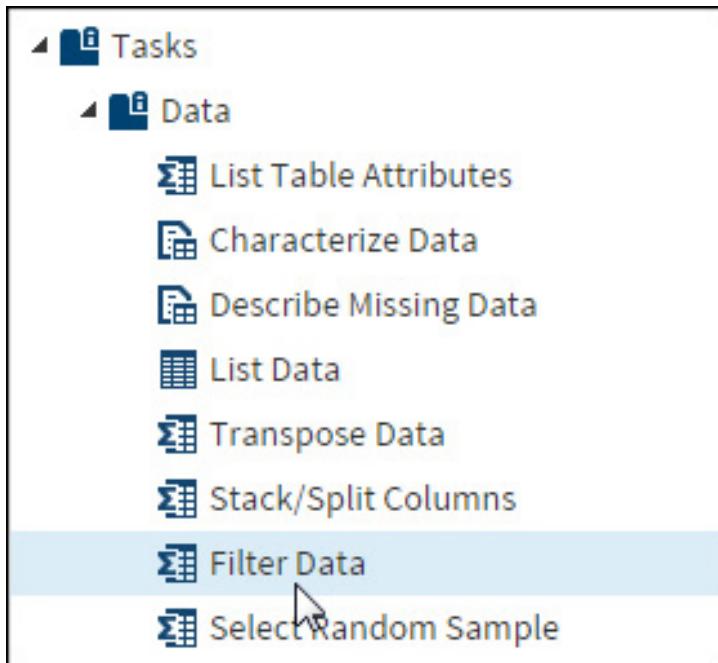
**Figure 4.9: Listing of the First Seven Rows of the Fish Data Set**

List Data for SASHELP.FISH				
Obs	Species	Weight	Height	Width
1	Bream	242	11.5200	4.0200
2	Bream	290	12.4800	4.3056
3	Bream	340	12.3778	4.6961
4	Bream	363	12.7300	4.4555
5	Bream	430	12.4440	5.1340
6	Bream	450	13.6024	4.9274
7	Bream	500	14.1795	5.2785

## Filtering Data

Another useful task is to filter the table—that is, you select rows that meet predefined criteria. To do this, double-click **Filter Data** in the task list.

**Figure 4.10: Selecting Filter Data from the Data List**



This brings up the following.

**Figure 4.11: Selections for Filtering Data**

The screenshot shows the SAS Studio interface with the following details:

- DATA:** SASHelp.FISH
- FILTER 1:**
  - \*Variable 1:** Weight
  - Comparison:** Less than (selected)
  - Value type:** Less than
  - \*Value:** (empty input field)
  - Logical:** (none)
- OUTPUT DATA:** (link)

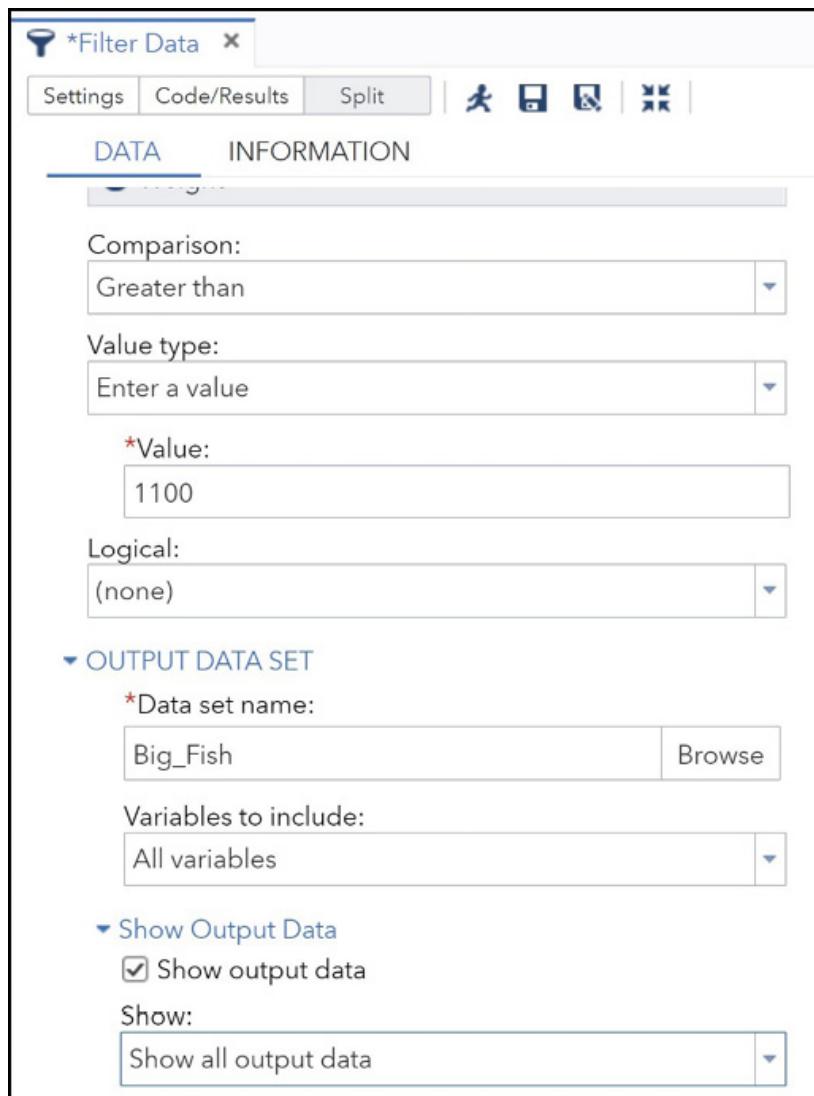
A dropdown menu is open under "Comparison" with the following options:

- Less than
- Less than or equal
- Equal
- Not equal
- Greater than (selected)
- Greater than or equal

A cursor icon is visible over the "Greater than or equal" option.

You select the data set as before. Next, you select a variable and a condition for your filter. In this example, you are selecting **Weight** as your variable and **Greater than** as your condition. You can now enter a value for the filter. In this example, you want to see rows in the table where the variable Weight is greater than 1,100.

**Figure 4.12: Selecting Rows Where the Weight Is Greater Than 1,100**



You can also expand the **OUTPUT DATA SET** option to override a default data set name. In most cases, you will want to supply your own data set name. Here you are naming the data set **Big\_Fish**. You can also check or uncheck the **Show Output Data** box. Selecting it (as in this example) generates a listing of the filtered data set.

Notice that there is no comma in the number 1100 in the figure above. A comma would generate an error.

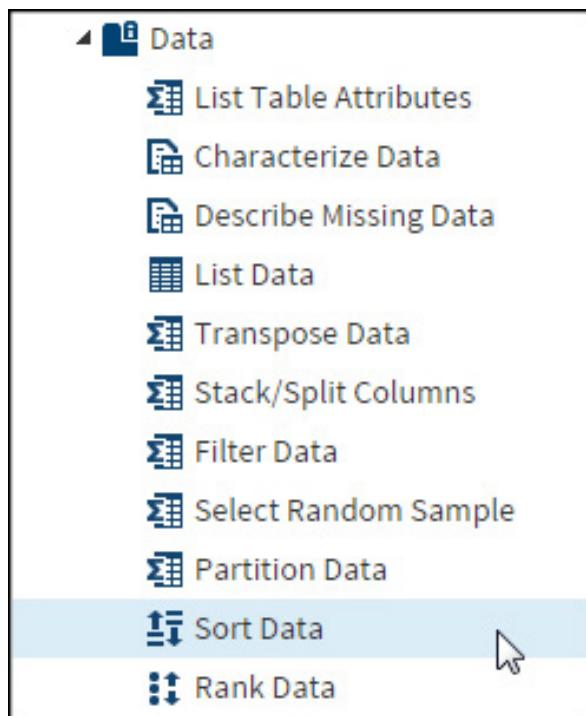
**Figure 4.13: Listing of the Filtered Data Set (Big\_Fish)**

List Data for WORK.BIG_FISH							
Obs	Species	Weight	Length1	Length2	Length3	Height	Width
1	Pike	1250	52	56.0	59.7	10.6863	6.9849
2	Pike	1600	56	60.0	64.0	9.6000	6.1440
3	Pike	1550	56	60.0	64.0	9.6000	6.1440
4	Pike	1650	59	63.4	68.0	10.8120	7.4800

## Sorting Data

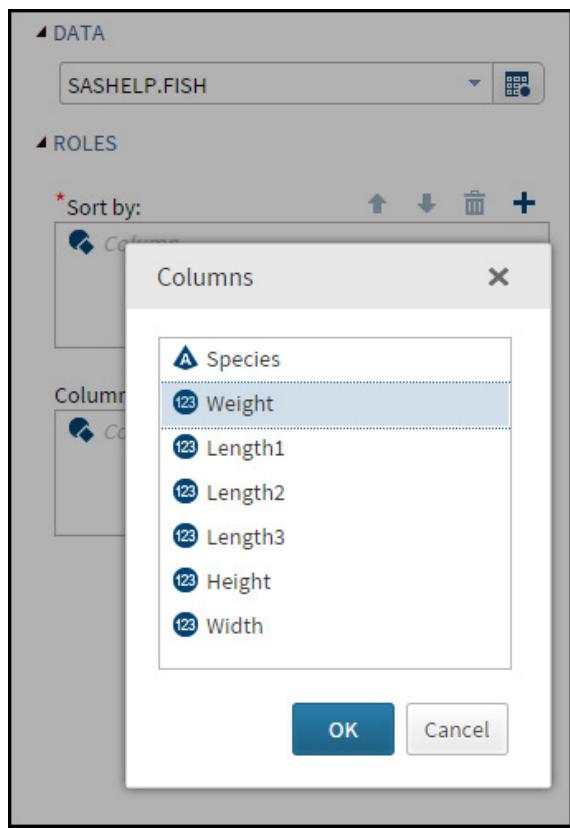
To sort a SAS data set by a variable, select **Sort Data** from the list.

**Figure 4.14: The Sort Data Task**



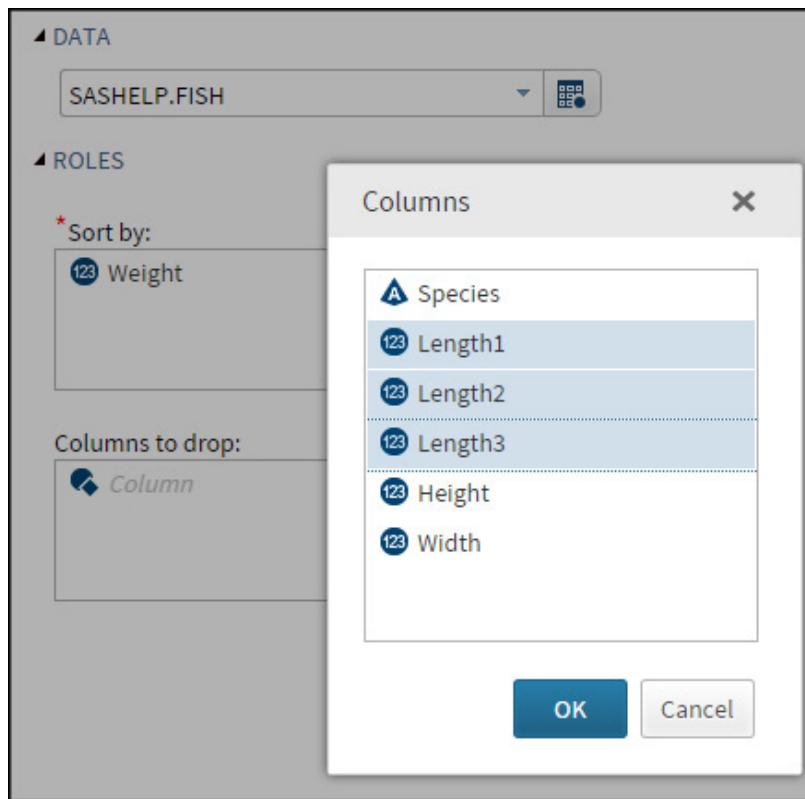
Just as in the previous tasks, you can now choose a data set and options.

**Figure 4.15: Selecting a Data Set and Variables for the Sort**



You are starting with the SASHHELP.Fish data set and requesting a sort based on the variable Weight. You can also choose columns to drop in this operation.

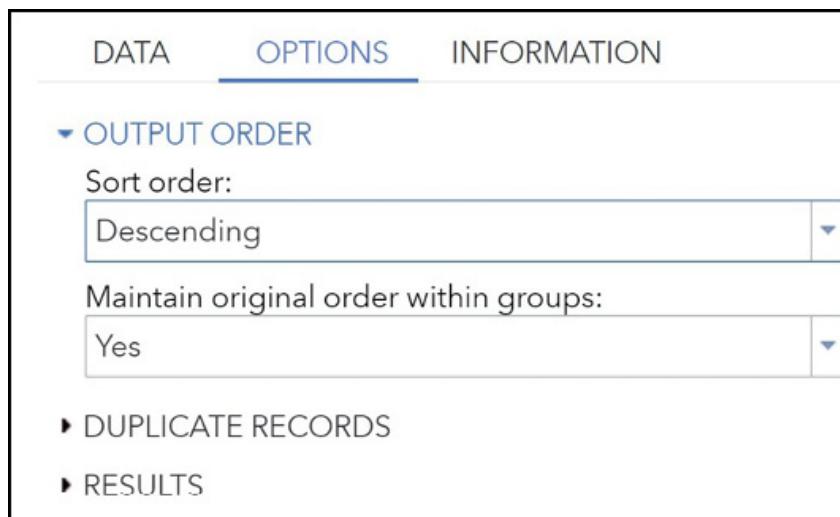
**Figure 4.16: Selecting Columns to Drop**



In this example, you are dropping the three Length variables.

Before you execute the sort, there are several options that you should consider. Click the **OPTIONS** tab to see these options. The default sort order is ascending (from smallest to largest). However, in this case, you want to see the heavier fish at the top of the list, so you choose **Descending** as the sort order.

**Figure 4.17: Selecting a Descending Sort**



Expand the **RESULTS** option to either sort in place or create a new data set with the sorted data.

**CAUTION:** Sorting in place replaces the original data set with the sorted data. If you drop columns, they will no longer be in the sorted data set.

In this example, you do not want to sort the original table—you want to create a new data set called `Sorted_Fish`:

**Figure 4.18: Naming the Output Data Set**



Click the **Run** icon to see the following screen.

**Figure 4.19: Result of Executing the Sort**

CODE | LOG | RESULTS | OUTPUT DATA

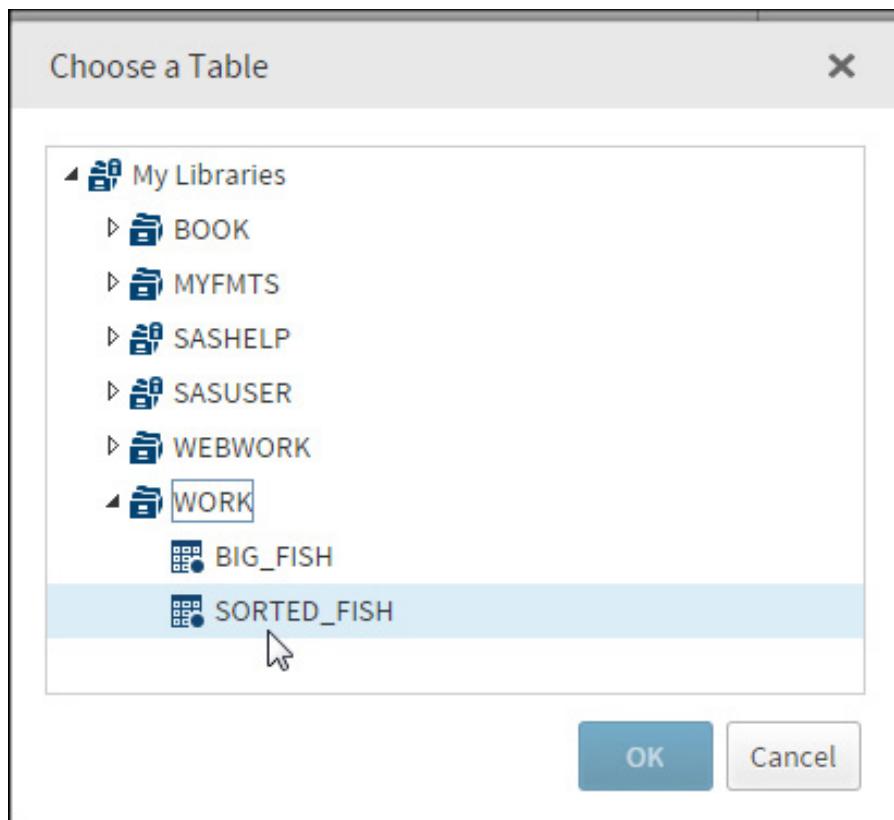
Table: WORK.SORTED\_FISH | View: Column names | Filter:

(none)

Columns		Total rows: 159 Total columns: 4	Rows 1-100
<input checked="" type="checkbox"/> Select all			
<input checked="" type="checkbox"/> Species		Pike	1650
<input checked="" type="checkbox"/> Weight		Pike	1600
<input checked="" type="checkbox"/> Height		Pike	1550
<input checked="" type="checkbox"/> Width		Pike	1250
		Perch	1100
		Perch	1100
		Perch	1015
		Bream	1000
		Whitefish	1000
		Perch	1000
		Perch	1000
		Perch	1000
		Bream	975
		Bream	955
		Bream	950
		Pike	950
		Bream	925
		Bream	920
		Perch	900

To see a nicer listing of this data set, go back to the **List Data** selection and proceed as you did in the first section of this chapter. Once you have opened the **List Data** task, choose a table to display and any options that you want. We will see that process below.

**Figure 4.20: Choose a Table to List**



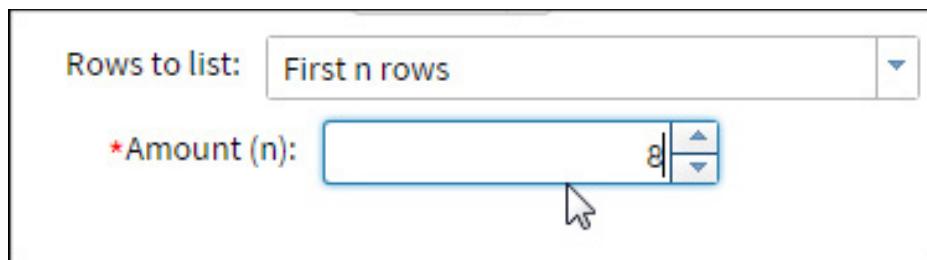
Next, choose which variables to include in the listing.

**Figure 4.21: Choosing Variables to List**



Finally, select any options that you want. In this example (see Figure 4.22), you want to see the first eight rows of the table.

**Figure 4.22: Option to List the First Eight Rows of the Table**



Here is the listing.

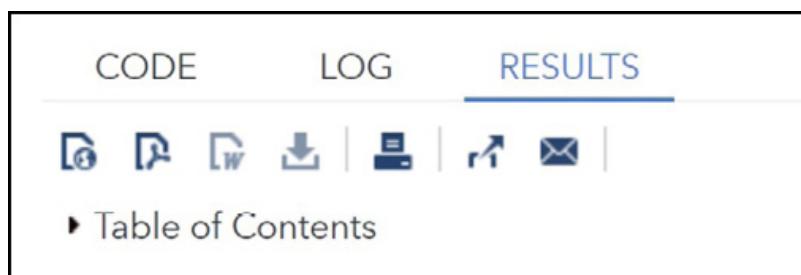
**Figure 4.23: First Eight Rows of the Sorted\_Fish Data Set**

List Data for WORK.SORTED_FISH				
Obs	Species	Height	Width	Weight
1	Pike	10.8120	7.4800	1650
2	Pike	9.6000	6.1440	1600
3	Pike	9.6000	6.1440	1550
4	Pike	10.6863	6.9849	1250
5	Perch	12.8002	6.8684	1100
6	Perch	12.5125	7.4165	1100
7	Perch	12.3808	7.4624	1015
8	Bream	18.9570	6.6030	1000

## Outputting HTML and PDF Files

You can download this listing as an HTML or PDF by clicking the appropriate icon at the top of the results window (as shown below).

**Figure 4.24: Icons to Download HTML or PDF Files**

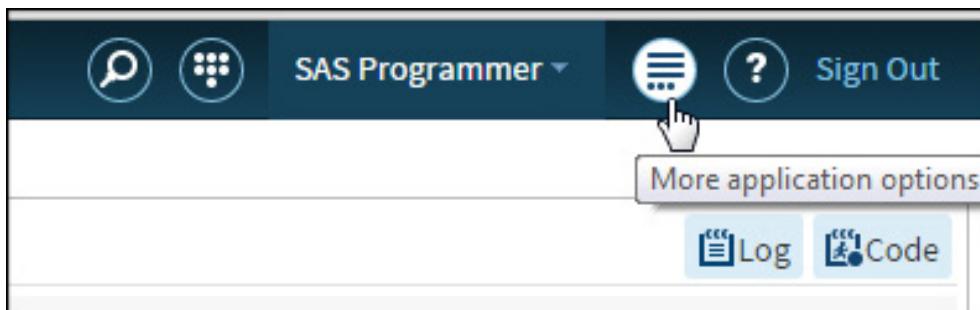


The left-most icon allows you to download the listing as an HTML file. The icon to the right of the HTML icon downloads the listing as a PDF file.

To create an RTF (rich text format) file, first click the SAS Studio

Options icon.

**Figure 4.25: Click More Application Options**



Click **Preferences**, then **Results**, and then check **RTF**. **HTML** and **PDF** are already chosen by default. (Note: These preferences remain in effect unless you change them later.)

In Figure 4.26, you are downloading a PDF file.

**Figure 4.26: Downloading a PDF File**

A screenshot of the SAS application interface showing the "RESULTS" tab selected. Below the tab are several icons: a magnifying glass, a double arrow, a downward arrow, a printer, and a refresh symbol. A tooltip "List D" is displayed above the printer icon, with the text "Download results as a PDF file" next to it. Below the icons is a table with the following data:

Species	Weight	Height	Width
Bream	242	11.5200	4.0200
Bream	290	12.4800	4.3056
Bream	340	12.3778	4.6961
Bream	363	12.7300	4.4555
Bream	430	12.4440	5.1340

Here is a listing of the PDF file.

**Figure 4.27: Listing of the PDF File**

### List Data for SASHELP.FISH

Species	Weight	Height	Width
Bream	242	11.5200	4.0200
Bream	290	12.4800	4.3056
Bream	340	12.3778	4.6961
Bream	363	12.7300	4.4555
Bream	430	12.4440	5.1340

## Joining Tables (Using the Query Window)

The last topic in this chapter describes how to use the **Query** utility to join two tables. Two data sets, ID\_Name and Grades, were created to explain how the joining process works. Figure 4.28 shows a listing of these two data sets.

**Figure 4.28: Listing of Data Sets ID\_Name and Grades**

List Data for WORK.ID_NAME		List Data for WORK.GRADES			
ID	Name	ID	Grade1	Grade2	Grade3
001	Ron	005	78	80	82
002	Jan	002	100	90	95
003	Peter	001	99	95	98
004	Paul	006	65	67	69
005	Mary	004	85	86	84

There are several features of these two data sets that are important for you to notice. Note the following: The data set ID\_Name is sorted in ID order—in the data set Grades, it is not. This does not cause a problem—the Query tool automatically sorts the data sets. Also, ID 003 is in data set ID\_Name and not in data set Grades; ID 006 is in Grades but not in ID\_Name. The goal is to join these two tables based on the ID column.

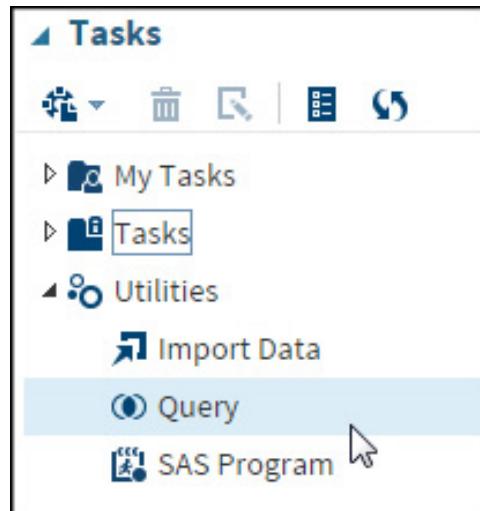
If you want to play along with this demonstration, you can run the program shown next to create these two tables.

### Program 4.1: Program to Create Data Sets ID\_Name and Grades

```
data ID_Name;
  informat ID $3. Name $12.;
  input ID Name;
datalines;
001 Ron
002 Jan
003 Peter
004 Paul
005 Mary
;
data Grades;
  informat ID $3. ;
  input ID Grade1-Grade3;
datalines;
005 78 80 82
002 100 90 95
001 99 95 98
006 65 67 69
004 85 86 84
;
```

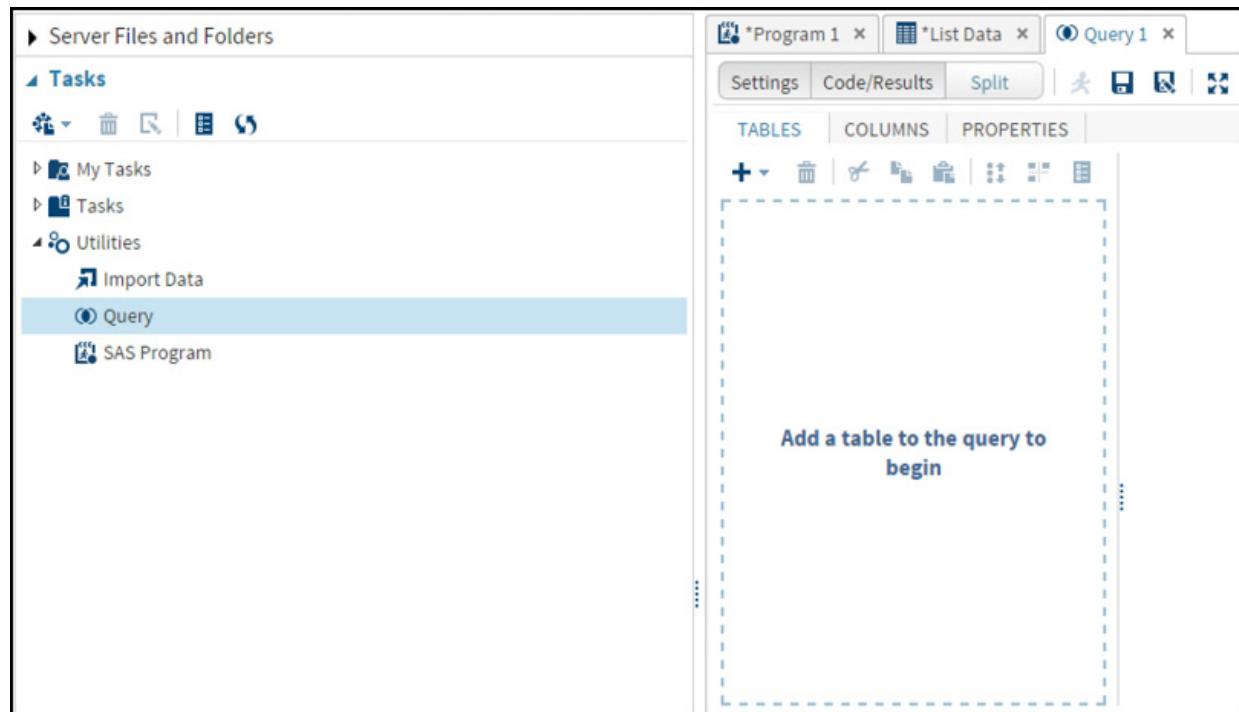
The first step in joining these two tables is to select **Query** from the **Utilities** tab.

**Figure 4.29: The Query Task in the Utilities Tab**



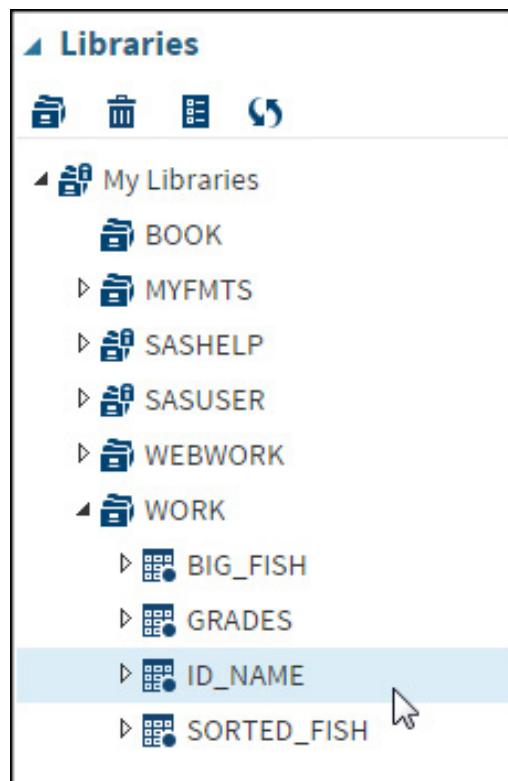
This brings up the following screen.

**Figure 4.30: The Query Window**



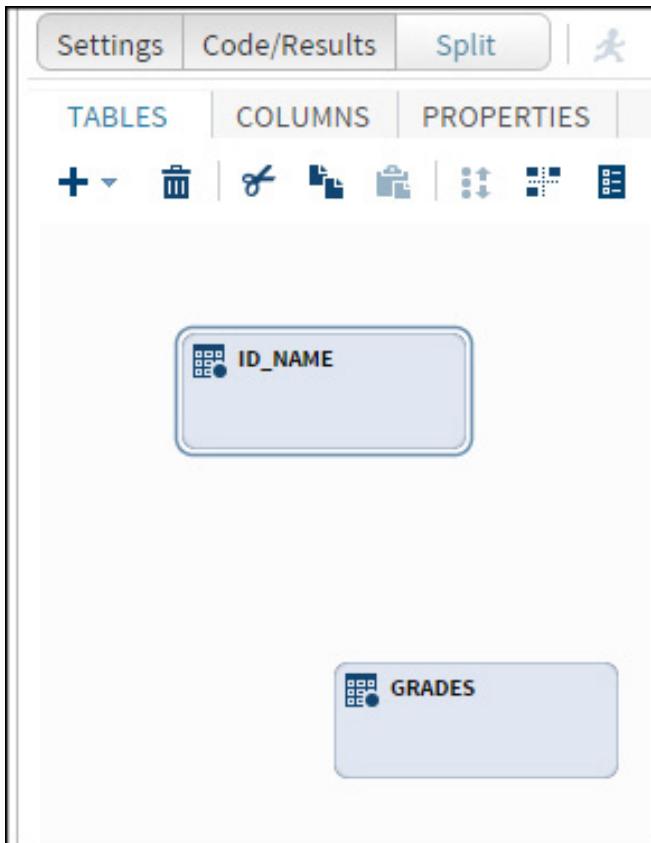
The next step is to open the **Libraries** tab and find the **Work** library.

**Figure 4.31: Locate the Two Tables in the Work Library**



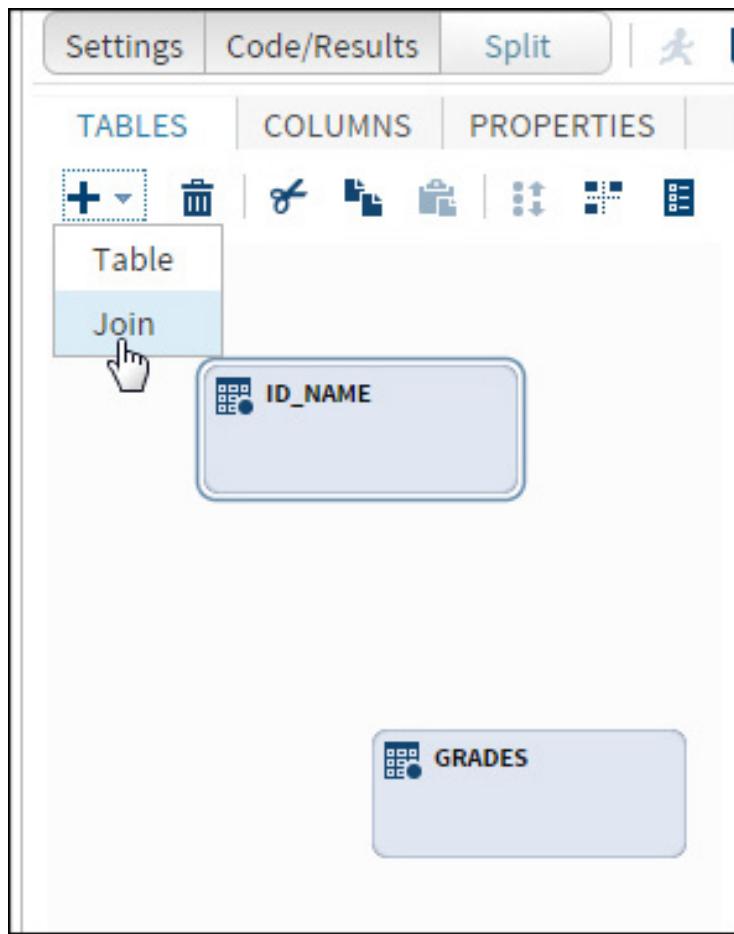
Click each file and drag it into the **Query** window. (If you drag the second file on top of the first file, SAS Studio automatically assumes that you want to perform a join operation.)

**Figure 4.32: Drag the Two Files into the Query Window**



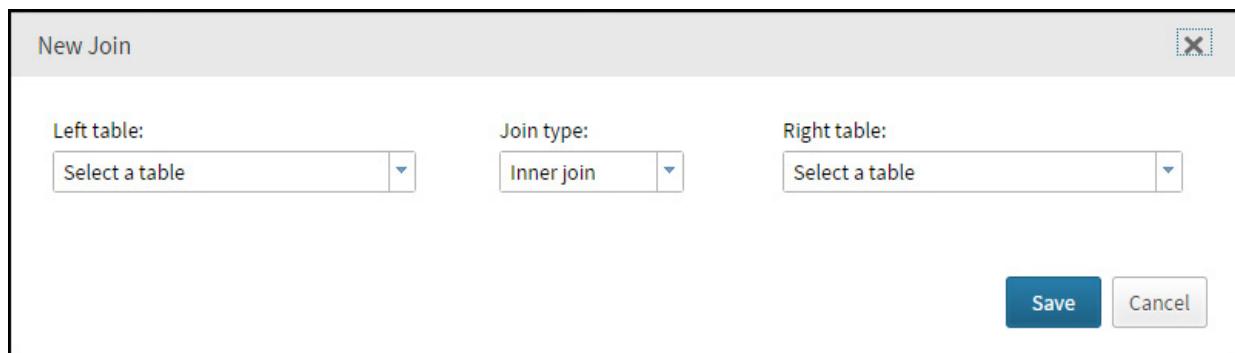
In the menu on the **TABLES** tab, select **Join** (unless you dragged the second file on top of the first, in which case the Query tool assumes that you want a join).

**Figure 4.33: Selecting Join in the Menu**



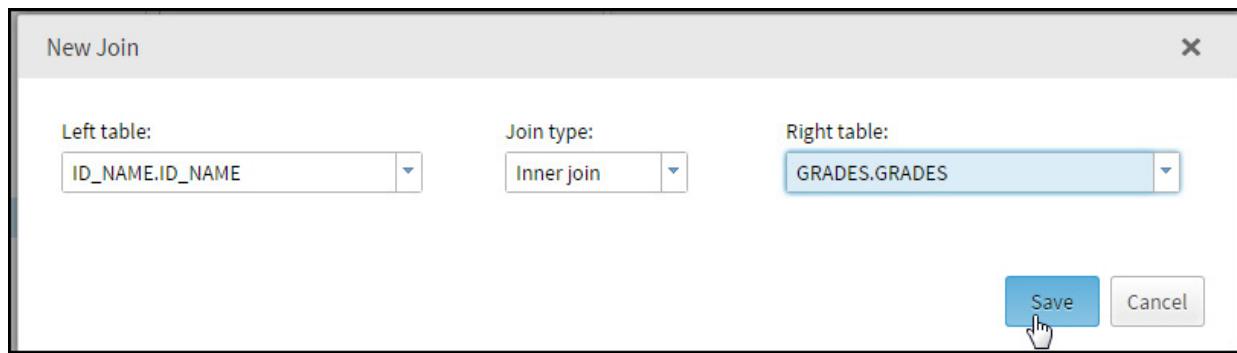
This brings up the following screen.

**Figure 4.34: Getting Ready to Join the Two Tables**



Select **ID\_Name** for the **Left table** and **Grades** for the **Right table**. For **Join type**, select **Inner join**.

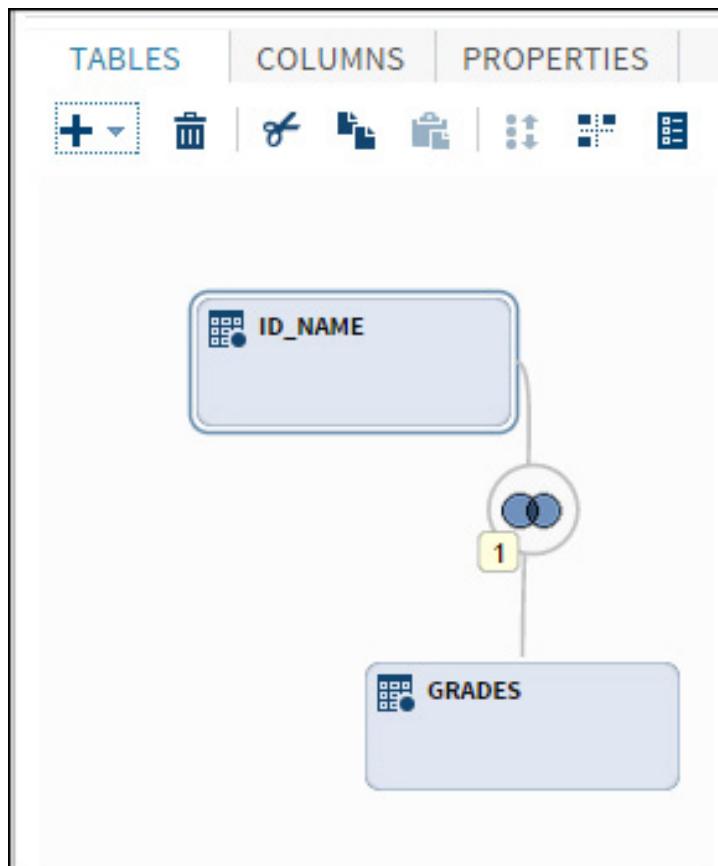
**Figure 4.35: Selecting the Two Tables and Inner Join**



**Click Save.** You now see the two tables with a Venn diagram that represents an inner join. If you are familiar with SQL, you already know the four types of joins. For those readers who are not, here is the explanation.

Because some IDs are only in one table, you have some decisions to make about how you want to handle the join. The most common join, selected in this example, is an *inner join*. This type of join includes only those rows where there is a matching ID in both files. An *outer join* includes all rows from both tables (with missing values in the rows from the table that does not contain an ID). Finally, the other two joins are a *left join* and a *right join*. In a left join, all IDs from the left table are included even if there isn't a matching ID in the right table. In a right join, all IDs from the right table are included even if there isn't a matching ID in the left table.

**Figure 4.36: Venn Diagram Showing an Inner Join of the Two Tables**



The next step is to name the columns that you want to use to join the two tables. In this example, because ID is in both tables, the Query tool automatically selects ID for the join variable. You are free to select any variable from each file to construct the join, even if the variable names are not the same in the two files.

**Figure 4.37: Selecting the Join Conditions**

The screenshot shows the 'Join conditions' section of the query tool. It includes fields for 'Left table: ID\_NAME', 'Join type: Inner join', and 'Right table: GRADES'. Below these, there is a 'Join conditions' section with a plus sign (+) button. Underneath is a table-like structure with two input fields: 'ID' on the left and '=' followed by another 'ID' field on the right. There are also delete and edit icons for these fields.

ID	=	ID
----	---	----

Next, select which columns you want in the joined table. To do this, click the **COLUMNS** tab. Select the columns in the usual way, and drag them to the **Add columns** area.

**Figure 4.38: Selecting Columns for the Final Table**

The figure shows the SAS Data step builder interface. The top navigation bar has tabs for TABLES, COLUMNS, and PROPERTIES. The COLUMNS tab is selected, indicated by a blue border. Below the tabs, there's a dropdown menu labeled 'View:' with options like 'Column names' and 'Table names'. On the left, a tree view lists tables: 'ID\_NAME' (selected) and 'GRADES'. Under 'ID\_NAME', columns 'ID' and 'Name' are listed. A cursor is hovering over 'Name'. On the right, a large text area is titled 'SELECT' and contains a dashed line indicating where to add columns. At the bottom right of this area, the text 'Add columns to include in the query (required)' is displayed.

The figure below shows the final list of columns in the joined table.

**Figure 4.39: Variables in the Joined Table**

Table	Source Column
ID_NAME	▲ ID
ID_NAME	▲ Name
GRADES	▲ ID
GRADES	123 Grade1
GRADES	123 Grade2
GRADES	123 Grade3

The last step is to click the **PROPERTIES** tab and indicate whether you want SAS to create a table or a report. If you choose a table, you can name the location (the Work library in this example) and the table name.

**Figure 4.40: Options in the PROPERTIES Tab**

TABLES | COLUMNS | PROPERTIES

**IDENTIFICATION**

Name: Query 1

Location:

**RESULTS**

Output type: Table

Output location: WORK

Output name: Combined

Clicking the **Run** icon finishes the join. A snapshot view of the resulting table is produced.

**Figure 4.41: View of the Resulting Table**

	ID	Name	ID2	Grade1	Grade2
1	005	Mary	005	78	80
2	002	Jan	002	100	90
3	001	Ron	001	99	95

You can use the **List Data** task to create a listing of the resulting table. It is shown in Figure 4.42.

**Figure 4.42: Listing of Combined Data Set**

List Data for WORK.COMBINED						
Obs	ID	Name	ID2	Grade1	Grade2	Grade3
1	005	Mary	005	78	80	82
2	002	Jan	002	100	90	95
3	001	Ron	001	99	95	98
4	004	Paul	004	85	86	84

Because this was an inner join, only those IDs that were in both

tables are listed in the final table.

## Conclusion

Only a few of the more popular data tasks were described in this chapter. Once you get the knack of running a task, you should feel confident in trying out some of the other data tasks in the list. The decision to use a task or write a SAS program is a personal choice. For those with programming experience, writing a program might be the way to go—for those folks who are new to SAS and just want to get things done, using the tasks is a great way to go. Or you can do both! Use a task or a utility to get the basic program written for you, and then take that and edit it to do more.

# Chapter 5: Summarizing Data Using SAS Studio

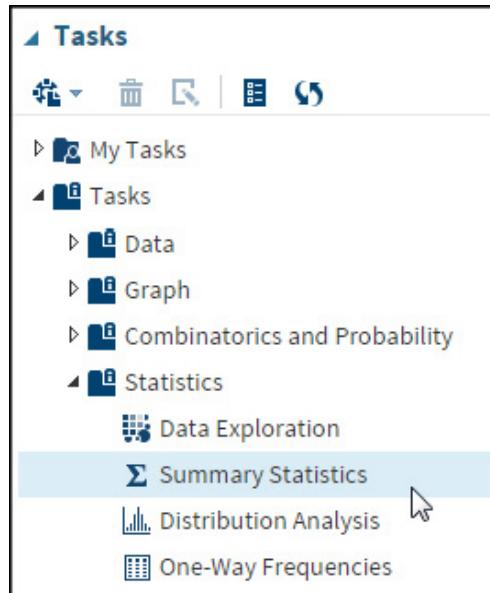
## Introduction

This chapter demonstrates how to summarize numeric and character data. For numeric variables, you will see how to compute statistics such as means and standard deviations, as well as histograms. For character variables, you will see how to generate frequency distributions and bar charts.

## Summarizing Numeric Variables

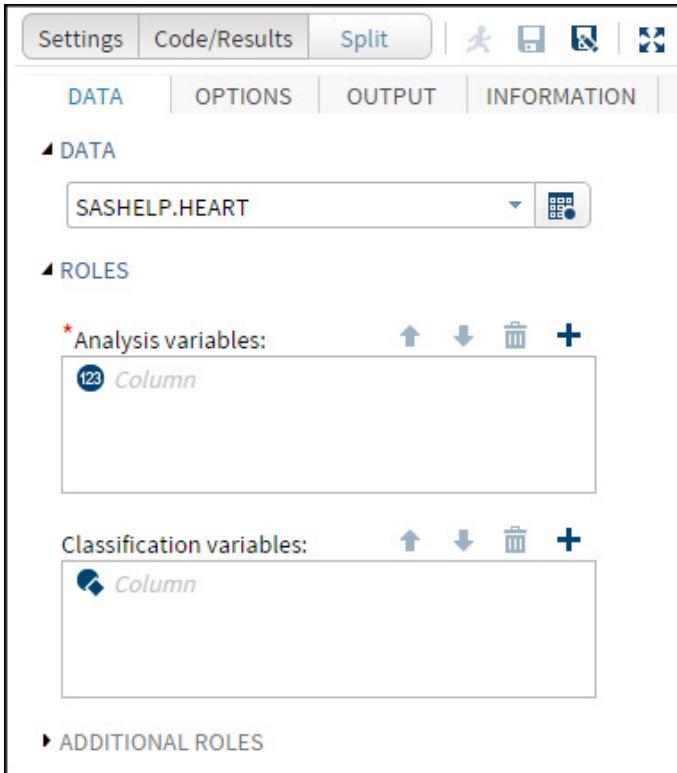
One of the most useful tasks for summarizing numeric data is found on the **Statistics** task list. Don't be alarmed that this task is listed under statistics—you don't need to be a statistician to understand how this works. Expand the **Statistics** task and select **Summary Statistics**, as shown in [Figure 5.1: Summary Statistics](#) below.

**Figure 5.1: Summary Statistics**



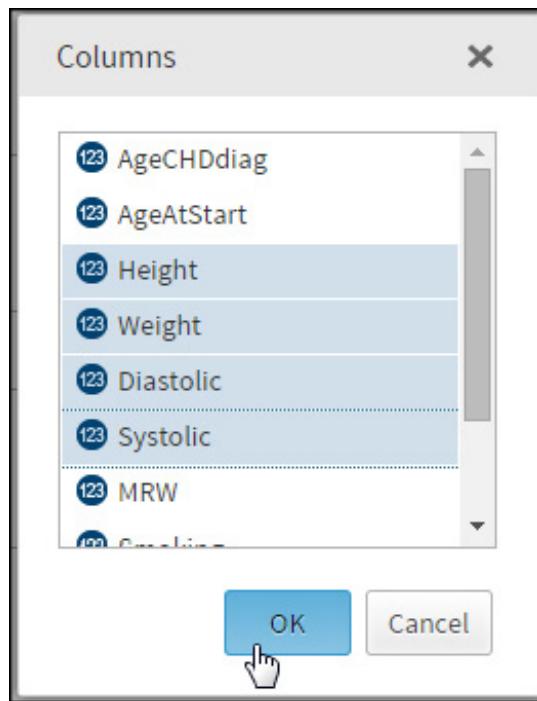
Double-click **Summary Statistics** to bring up the **DATA** and **OPTIONS** tabs (Figure 5.2).

**Figure 5.2: The Summary Statistics Task**



Let's choose the **Heart** data set in the SASHELP library to demonstrate how to use the **Summary Statistics** task. Looking at [Figure 5.2: The Summary Statistics Task](#), you see that SASHELP.Heart has already been chosen. You can click the **Data Table** icon to select a library and data set that you want to use. The next step is to select the variables that you want to summarize.

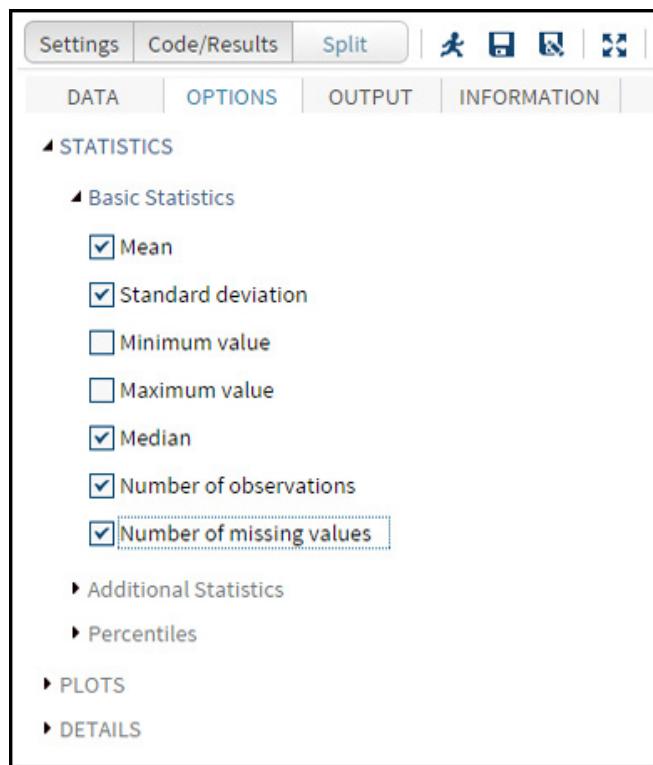
**Figure 5.3: Selecting Variables**



Click the plus sign in the **Analysis Variables** window and select the variables that you want to analyze (using the Ctrl or Shift keys, as described previously). Notice that the variable list contains only numeric variables. For this example, you have chosen the variables **Height**, **Weight**, **Diastolic** (diastolic blood pressure), and **Systolic** (systolic blood pressure). Click **OK** when you are done selecting variables.

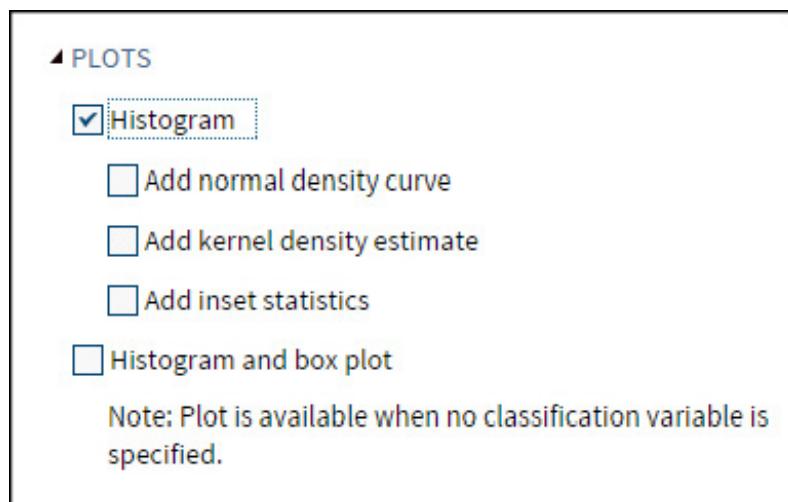
You can click the **Run** icon now or customize the report by clicking on the **OPTIONS** tab. Let's do that. It brings up the following screen.

**Figure 5.4: Selecting Options**



Select or deselect the options that you want. This author recommends that you select the two options—**Number of observations** and **Number of missing values**—as they are quite useful. At this point, you can run the task or continue on to request plots. Let's do that.

## Figure 5.5: Plots



You have the option to include a histogram or a histogram with a box

plot. In this example, you have chosen a histogram. If you are statistically minded, you can add a normal density curve and a kernel density estimate. The third option in this list places summary statistics in an inset window in the histogram.

It's time to run the task. The first part of the output shows the statistics that you requested in tabular form. It looks like this.

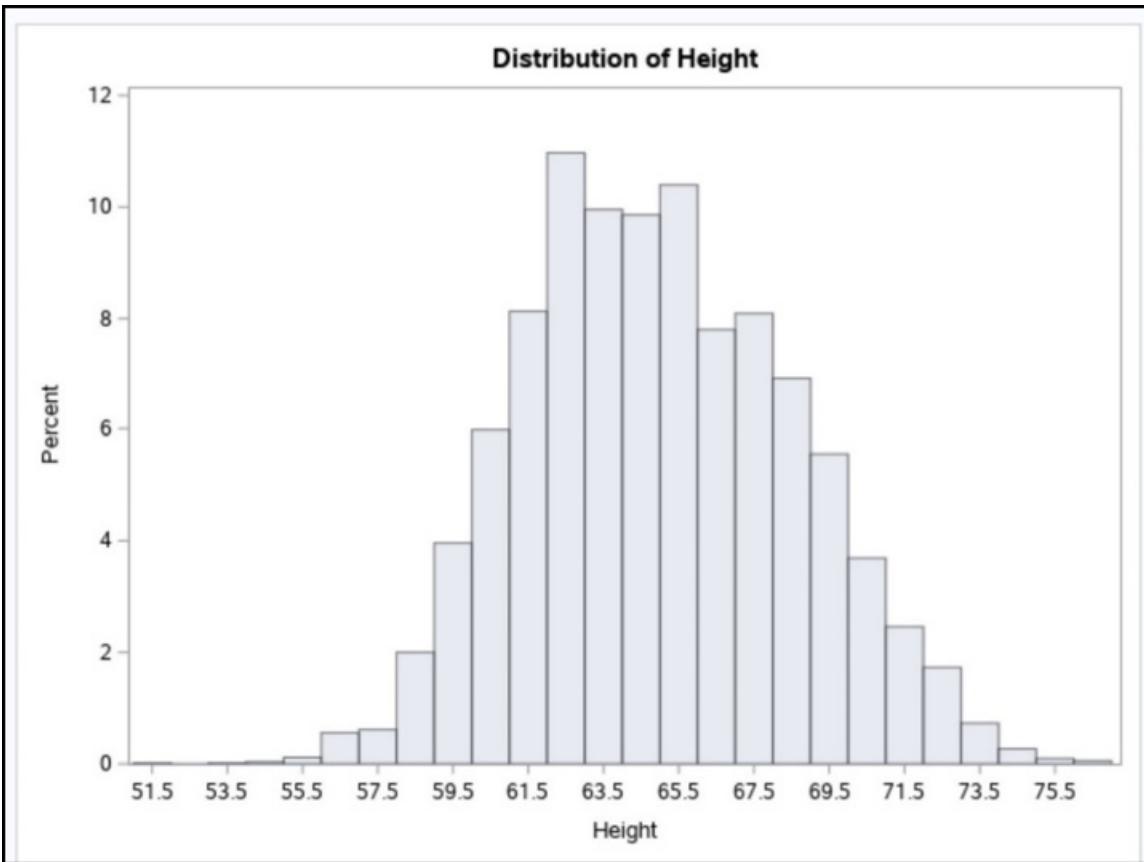
**Figure 5.6: Tabular Output**

Variable	Mean	Std Dev	Median	N	N Miss
Height	64.8131847	3.5827074	64.5000000	5203	6
Weight	153.0866808	28.9154261	150.0000000	5203	6
Diastolic	85.3586101	12.9730913	84.0000000	5209	0
Systolic	136.9095796	23.7395964	132.0000000	5209	0

Here you see the mean, standard deviation, and median for the selected variables. The last two columns, labeled N and N Miss, show the number of nonmissing observations and the number of missing values for all of your variables, respectively.

To economize on space in this book, only one histogram (Height) is displayed. It looks like this.

**Figure 5.7: Histogram**



You can use the scroll bars to move right or left, up or down. You can also click the **Expand** icon to see the entire histogram.

## Adding a Classification Variable

The statistics that you have seen so far are computed on the entire data set. To see statistics broken down by one or more classification variables, add those variables in the **Classification variables** box. To demonstrate this, let's see the statistics for Height, Weight, Diastolic, and Systolic broken down by the variable Sex. In the **DATA** tab, click the plus sign in the **Classification variables** box and select the variable **Sex** (Figure 5.8).

**Figure 5.8: Adding a Classification Variable**

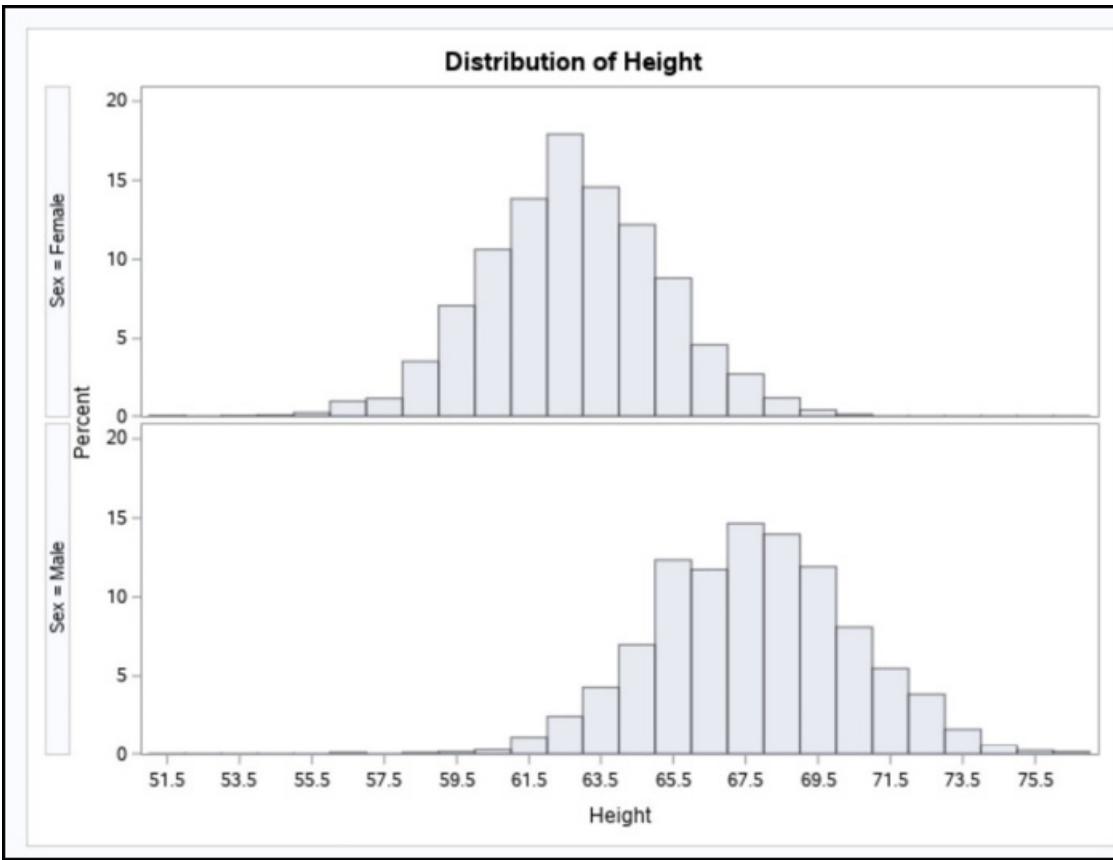
Now run the program to see the following table:

**Figure 5.9: Output Showing Classification Data**

Sex	N Obs	Variable	Mean	Std Dev	Minimum	N	N Miss
Female	2873	Height	62.5725863	2.4524112	51.5000000	2869	4
		Weight	141.3886372	26.2880439	67.0000000	2869	4
		Diastolic	84.6463627	13.3394548	50.0000000	2873	0
		Systolic	136.8861817	25.9835883	82.0000000	2873	0
Male	2336	Height	67.5673736	2.7321366	56.0000000	2334	2
		Weight	167.4661525	25.2907044	99.0000000	2334	2
		Diastolic	86.2345890	12.4548941	50.0000000	2336	0
		Systolic	136.9383562	20.6535522	90.0000000	2336	0

You see all of the statistics you originally requested for each value of Sex. If you requested plots, you will see separate histograms for each value of the classification variable. To save space, only the histogram for Height is displayed (Figure 5.10).

**Figure 5.10: Separate Histograms by Sex**



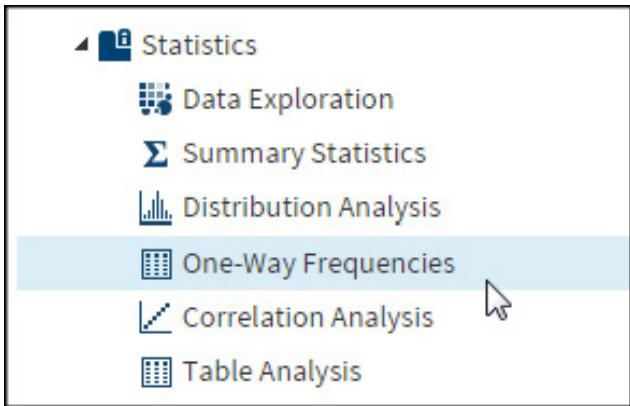
Seeing the two histograms juxtaposed like this is useful in determining if there are differences in the analysis variable for each level of the classification variable.

## Summarizing Character Variables

You can use the **One-Way Frequencies** task to compute counts and percentages for character or numeric variables. If you include any numeric variables in your selection, this task computes frequencies for every unique value of those variables. That is why this task is usually reserved for character variables or for numeric variables with very few unique values.

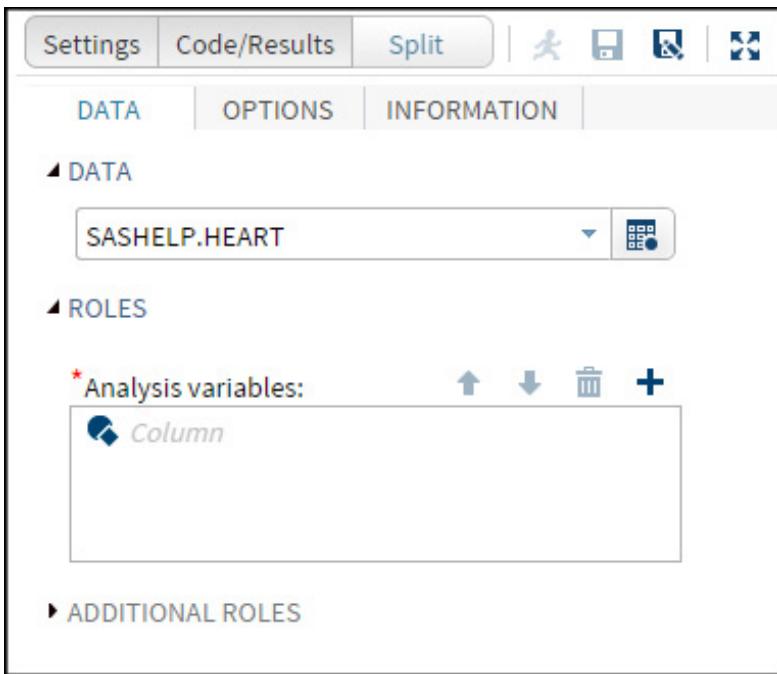
The first step is to double-click the **One-Way Frequencies** task.

**Figure 5.11: One-Way Frequencies**



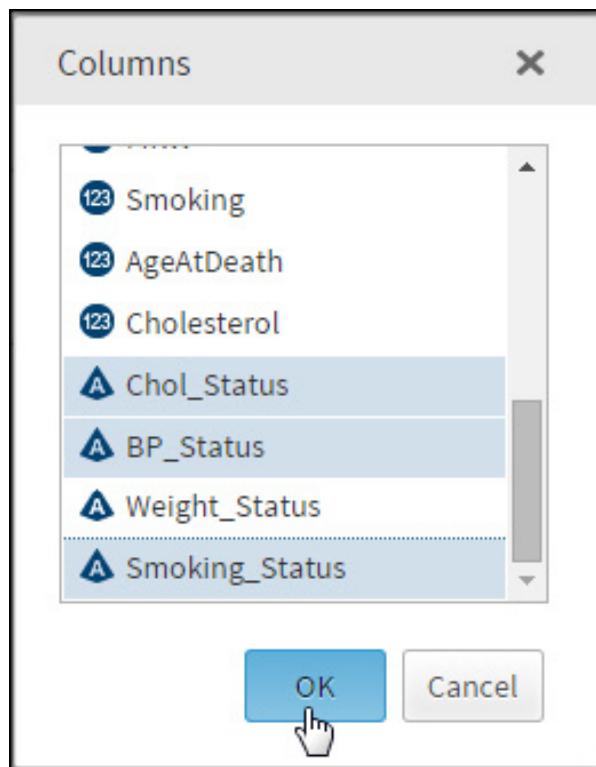
This brings up the screen shown below.

**Figure 5.12: One-Way Frequencies Task**



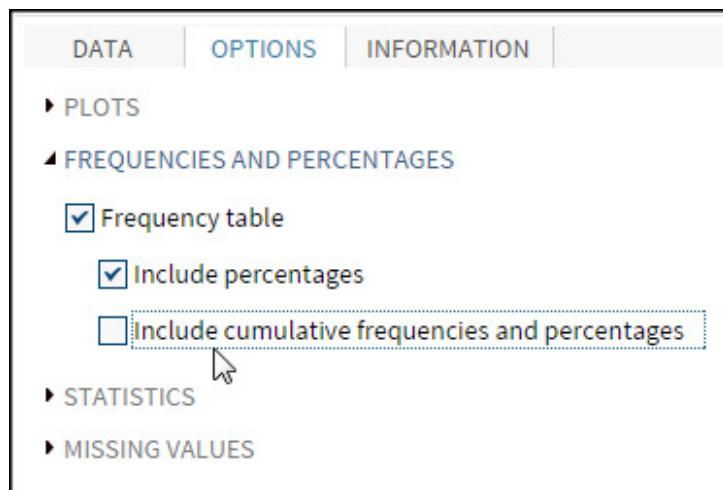
The SASHelp.Heart data set has already been selected. Click the plus sign attached to the **Analysis variables** box. Then select the variables for which you want to compute frequencies. For this demonstration, the variables **Sex**, **Chol\_Status**, **BP\_Status**, and **Smoking\_Status** were chosen (the variable **Sex** is farther up the list and does not appear in Figure 5.13).

**Figure 5.13: Selecting Variables**



Click **OK** to proceed. If you want to customize the frequency table, click the **OPTIONS** tab. This brings up the following:

**Figure 5.14: Frequency Options**



Here, you are deselecting the option to include cumulative frequencies and percentages. (The default is to include cumulative frequencies and percentages.) If you expand the **PLOTS** option, you will see that the default action is to produce plots. Select the option **Suppress plots** if you do not want bar charts. In this example, the

**Suppress plots** option is left unchecked.

**Figure 5.15: Suppress Plots**



Click the **Run** icon to complete the task. The output consists of frequency tables and bar charts. Figure 5.16 shows two of the four tables requested. Notice that cumulative frequencies and cumulative percentages are not included (because you unchecked the option to do this).

**Figure 5.16: Frequency Tables**

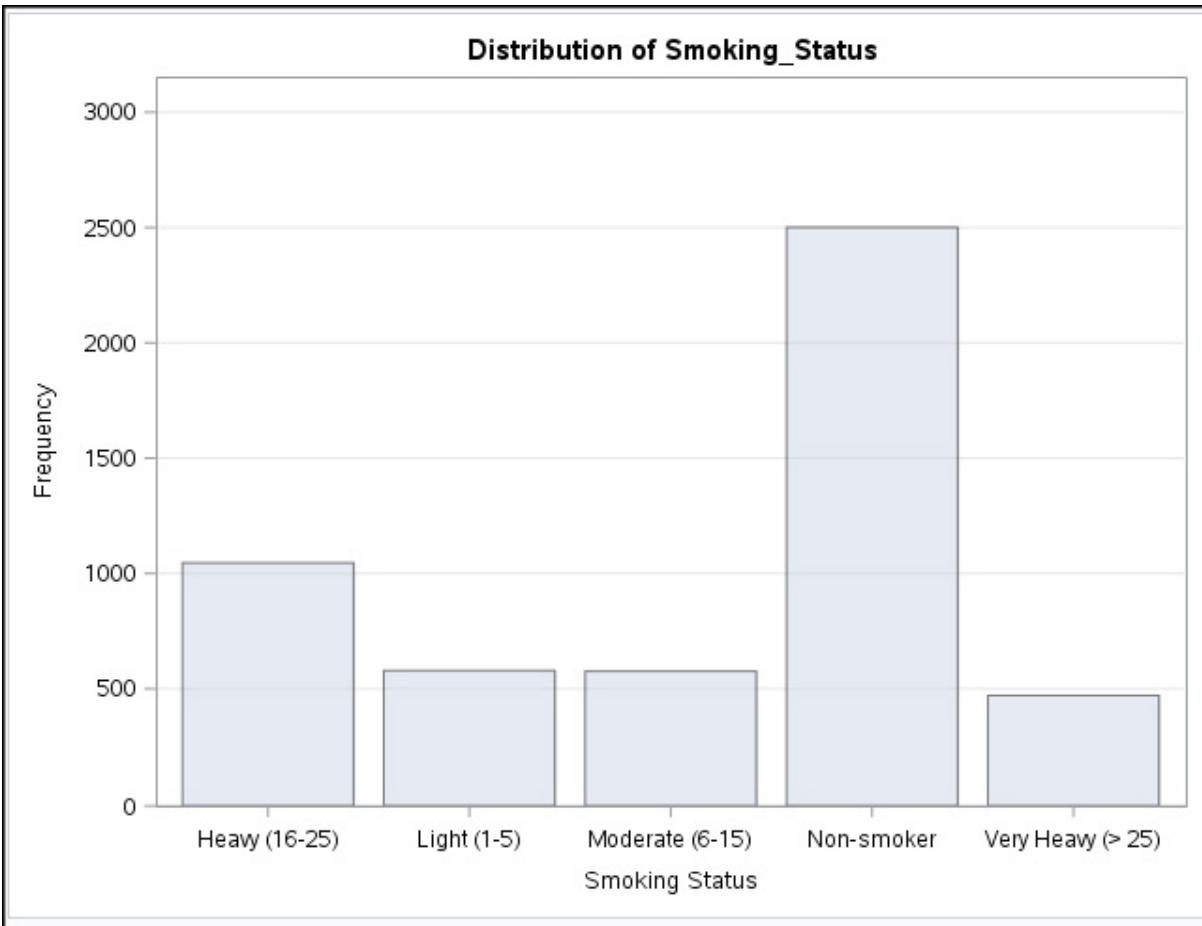
Sex	Frequency	Percent
Female	2873	55.15
Male	2336	44.85

Cholesterol Status		
Chol_Status	Frequency	Percent
Borderline	1861	36.80
Desirable	1405	27.78
High	1791	35.42
Frequency Missing = 152		

Only one bar chart (for Smoking\_Status) is displayed here (Figure 5.17).

**Figure 5.17: Bar Chart**



## Conclusion

The two statistics tasks, **Summary Statistics** and **One-Way Frequencies**, can be used to summarize numeric and character data, respectively. You can customize the tabular and graphical output by selecting options for both tasks.

# Chapter 6: Graphing Data

## Introduction

Creating charts and graphs is one of the more difficult programming tasks. Even veteran programmers need to pull out the manual or seek help online when attempting these tasks. Luckily, SAS Studio has a number of charting and graphing tasks that enable you to create beautiful, customized charts and graphs with ease.

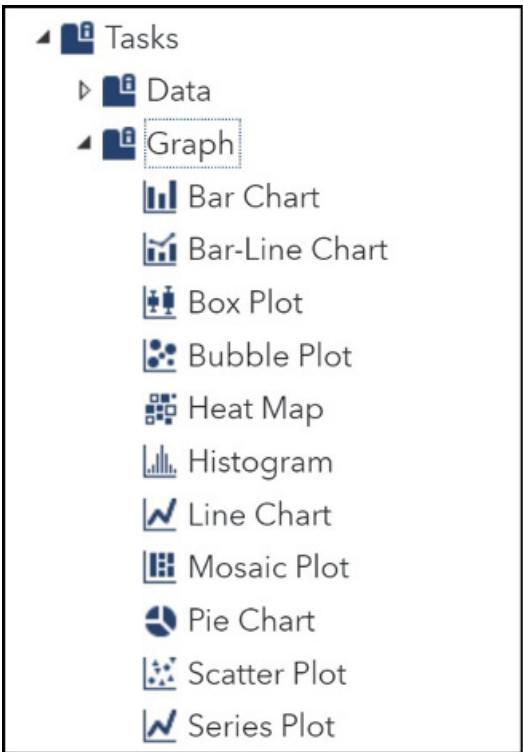
This chapter demonstrates a few of the more popular tasks, such as creating bar charts, pie charts, and scatter plots. Once you see how these tasks work, you will be able to use any of the other graph tasks offered by SAS Studio.

## Creating a Frequency Bar Chart

Let's start out with a simple bar chart where each bar represents a frequency on the Y axis. Because you are somewhat familiar with the SASHELP data set Heart (which is used in many of the previous chapters), let's start out using the variable Smoking\_Status to create a bar chart.

Open the **Tasks** tab in the navigation pane and expand the **Graph** tasks. It looks like this:

**Figure 6.1: Graph Tasks**



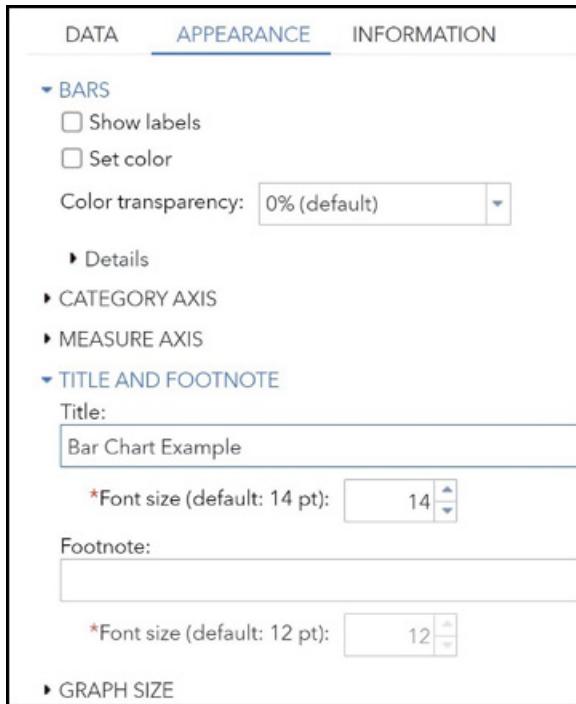
Double-click **Bar Chart** to get started. This opens the following screen.

**Figure 6.2: Requesting a Bar Chart**

This screenshot shows the configuration dialog for creating a bar chart. The top navigation bar has tabs for 'DATA', 'APPEARANCE', and 'INFORMATION'. The 'DATA' tab is active. In the 'DATA' section, 'SASHelp.Heart' is selected as the source. Under 'ROLES', 'Smoking\_Status' is assigned to the 'Category' role. The 'Measure' is set to 'Frequency count (default)'. The 'APPEARANCE' tab is visible at the bottom.

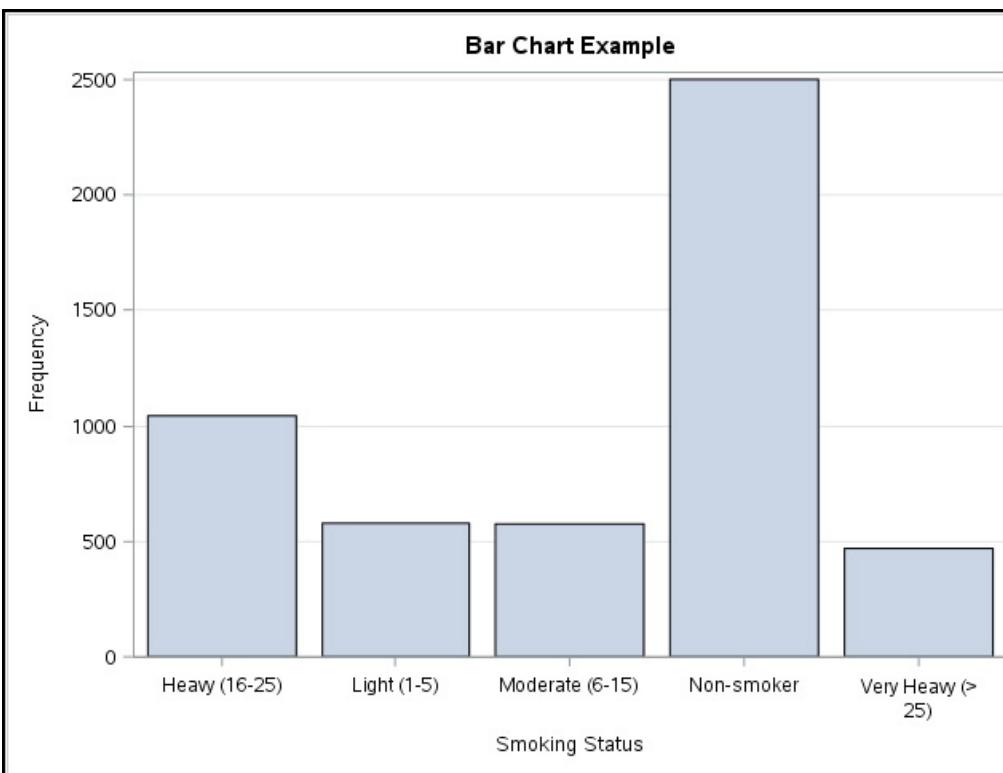
The Heart data set has already been entered and **Smoking\_Status** was chosen as the **Category variable**. The next step is to modify the appearance. Click the **APPEARANCE** tab to do this.

**Figure 6.3: Bar Chart Appearance**



You can enter titles and footnotes for this task and even choose the font size for each one. As you can see in Figure 6.3, you can also customize the bar details (change the bar color, for example), the bar labels, the two axes, the chart legend, and the graph size. You can also decide to accept all the defaults for these options and click the **Run** icon. You can always go back and modify these options later if you want. The figure below shows the resulting bar chart (accepting all the defaults).

**Figure 6.4: Final Bar Chart**



## Creating a Bar Chart with a Response Variable

The height of each bar in the previous chart represented frequency counts for each category of `Smoking_Status`. By choosing a response variable, you can create a bar chart where the height of each bar represents a statistic (the mean, for example) for a response variable at each level of smoking status.

Figure 6.5 shows the **DATA** tab for the **Bar Chart** task. Here you choose the **Heart** data set in the SASHELP library. First choose **Smoking\_Status** as your **Category** variable. Next, use the pull-down list under **Measure** to select **Variable** to represent the Y axis. Enter **Weight** as the variable that you want to plot. Finally, you use the menu of choices for statistics and choose **Mean** (the default is Sum). Note this task was significantly modified from previous versions of SAS Studio. Thanks to Paul Grant for pointing this out to me. Because of all these choices, the height of each bar will now represent the mean weight for each value of smoking status.

**Figure 6.5: Entering a Response Variable**

DATA APPEARANCE INFORMATION

▼ DATA  
SASHELP.HEART 

Filter: (none)

▼ CHART ORIENTATION  
 Vertical  
 Horizontal

▼ ROLES

\*Category: (1 item)    
A Smoking\_Status

Subcategory: (1 item)    
Column

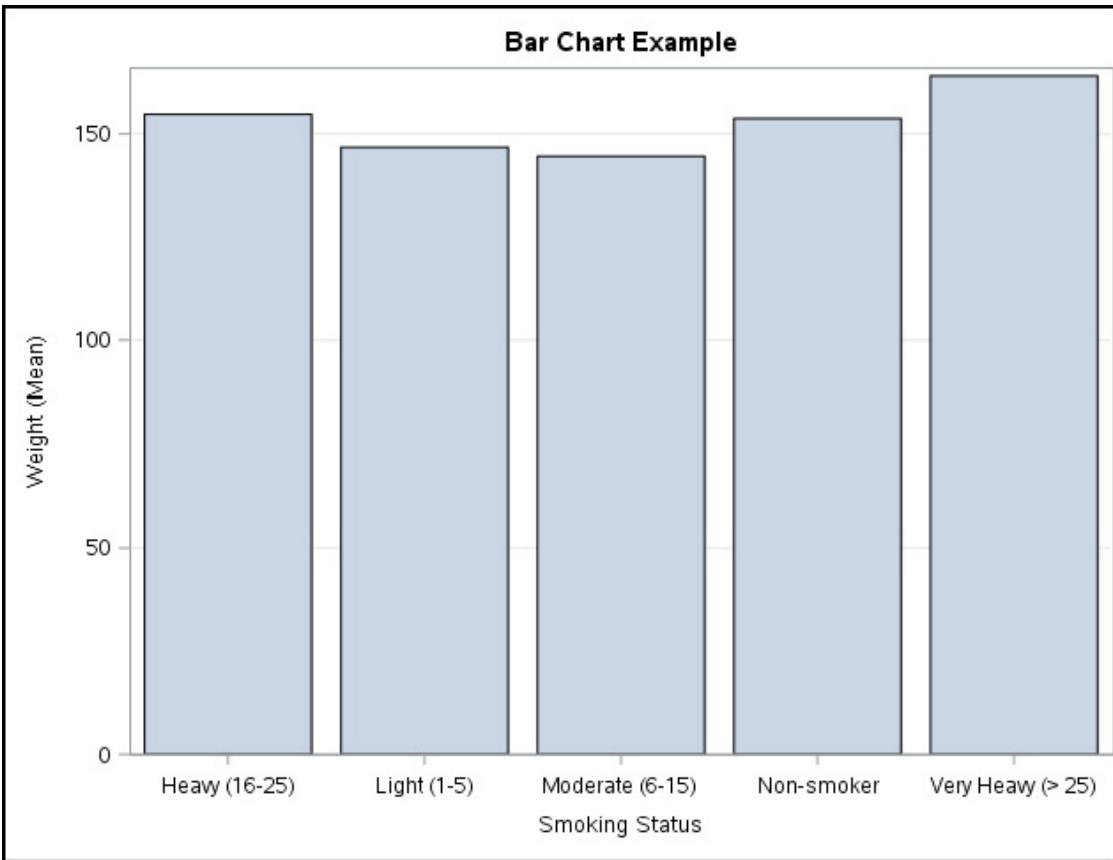
Measure: Variable  

\*Variable: (1 item)    
123 Weight

▼ Statistic:  
Mean 

Clicking the **Run** icon generates the chart shown in Figure 6.6.

**Figure 6.6: Bar Chart with a Response Variable**



The height of each bar represents the mean weight for each category of smoking status.

## Adding a Group Variable

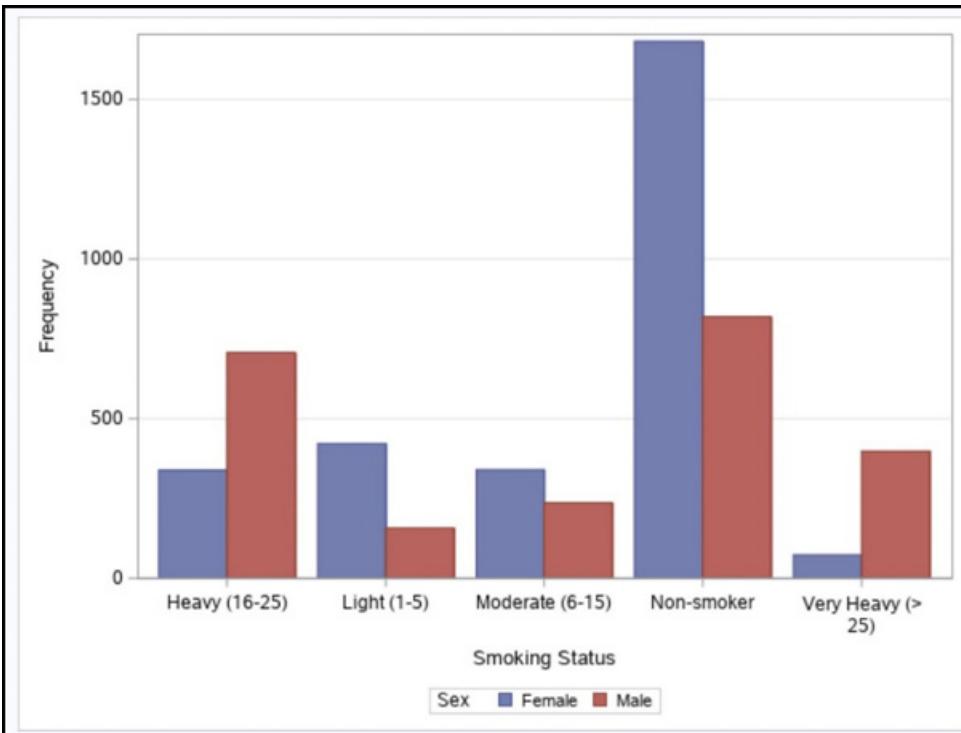
You can display more information in the bar chart by including a **Subcategory** variable in the **DATA** tab. Figure 6.7 shows that you want the variable Sex as a subgroup variable. At this point, you could choose to create a stacked chart (the male and female frequencies in a single bar with two different colors) or the default chart that shows side-by-side bars for men and women.

**Figure 6.7: Adding a Subgroup to a Bar Chart**

The screenshot shows the configuration interface for a bar chart. The top navigation bar has tabs for DATA, APPEARANCE, and INFORMATION, with DATA being the active tab. Under the DATA tab, there is a dropdown menu set to 'SASHELP.HEART' with a refresh icon. Below it is a 'Filter: (none)' button. The 'CHART ORIENTATION' section has two radio buttons: 'Vertical' (selected) and 'Horizontal'. The 'ROLES' section contains two categories: 'Category' (set to 'Smoking\_Status') and 'Subcategory' (set to 'Sex'). Both categories have '+' and '-' icons for managing items. Under 'Options:', there is a section for 'Display grouped bars:' with two radio buttons: 'Clustered side by side' (selected) and 'Stacked on one another'. There are also dropdown menus for 'Legend location' (set to 'Outside (default)') and 'Measure' (set to 'Frequency count (default)').

You can change the appearance of the plot by clicking the **Appearance** tab and make choices. For now, let's choose all the defaults and run the task. Figure 6.8 shows the resulting bar chart.

**Figure 6.8: Grouped Bar Chart**



It looks like there were many more female non-smokers in this sample, while in the heavy and very heavy categories, males predominated.

## Creating a Pie Chart

Pie charts represent a popular way to display frequencies. As with bar charts, the size of the slices can also represent a sum or mean of a response variable. Let's start out by generating a pie chart showing smoking status. Select **Pie Chart** from the task list.

**Figure 6.9: Requesting a Pie Chart**



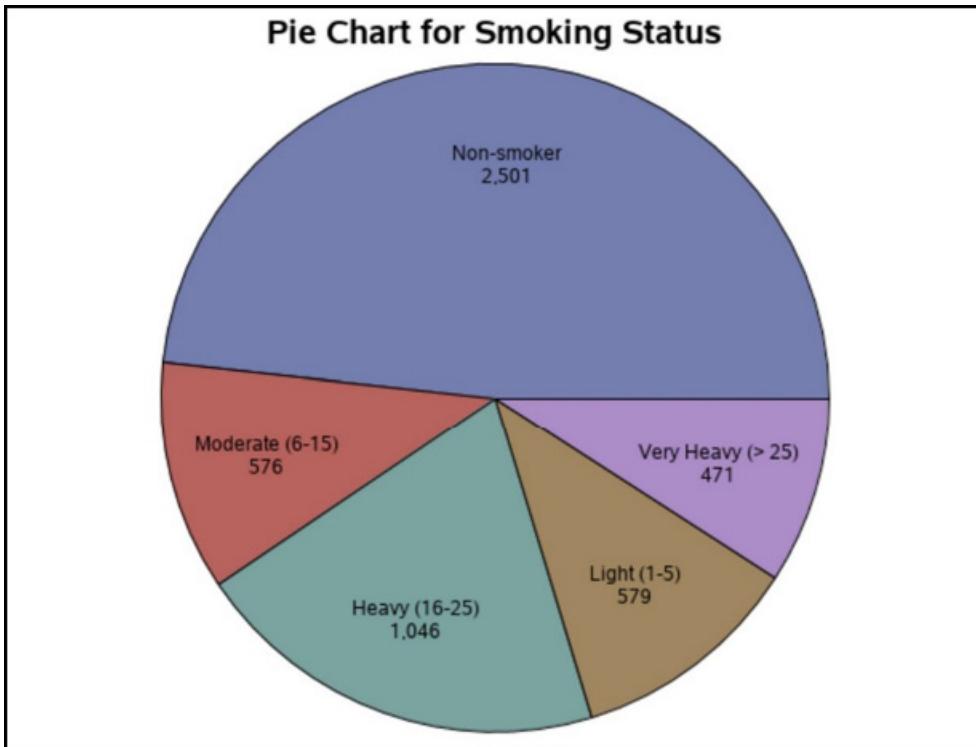
Double-click **Pie Chart** to get started. This brings up the following screen.

**Figure 6.10: Choosing a Variable for the Pie Chart**

The screenshot shows the configuration interface for a pie chart. The top navigation bar has tabs: DATA (underlined), APPEARANCE, and INFORMATION. The DATA section is expanded, showing the data source set to 'SASHELP.HEART'. Under the ROLES section, the 'Category' role is assigned to 'Smoking\_Status'. The 'Measure' setting is set to 'Frequency count (default)'. There are also sections for 'Subcategory' and 'Column' roles, each with a single item selected.

You can use the **Appearance** tab to modify the appearance of the pie chart. For this example, we will use this tab to add a title and leave all the other choices with default values. A chart with all the default options (except for the title) is shown in Figure 6.11.

**Figure 6.11: Final Pie Chart**



## Creating a Scatter Plot

A scatter plot shows the relationship between two variables by placing points on a set of X and Y axes. To get started, locate **Scatter Plot** in the **Graph** task list.

**Figure 6.12: Requesting a Scatter Plot**



Double-click **Scatter Plot** to bring up the following screen.

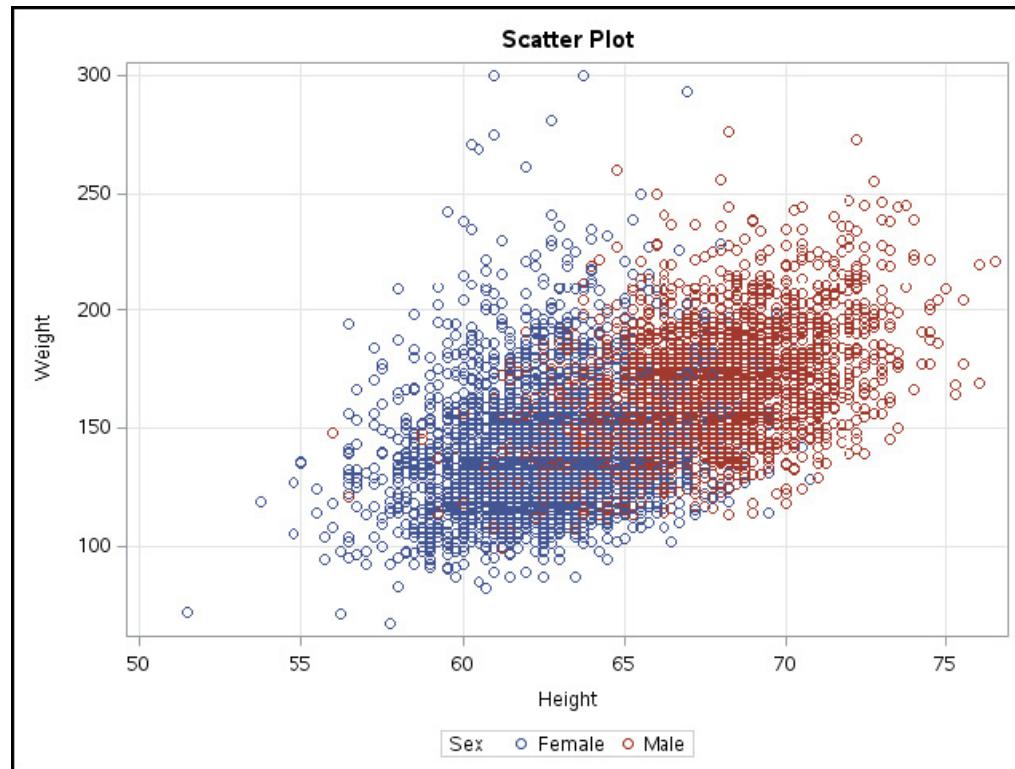
**Figure 6.13: Selecting the X, Y, and Grouping Variables**

A screenshot of a configuration dialog box for a scatter plot. The top navigation bar has tabs for DATA, APPEARANCE, and INFORMATION, with DATA being the active tab. The DATA section shows 'SASHelp.Heart' selected as the source. The ROLES section contains three items: 'X axis' with 'Height' assigned, 'Y axis' with 'Weight' assigned, and 'Group' with 'Sex' assigned. Below the roles, a 'Legend location' dropdown is set to 'Outside (default)'. There is also a section for 'ADDITIONAL ROLES' which is currently collapsed.

You enter the x- and y-variables as shown. Here you want to see **Height** on the X axis and **Weight** on the Y axis. You can also identify a group variable. For this example, you are choosing **Sex** as a group variable. Click the **Run** icon to generate the really impressive scatter

plot shown in Figure 6.14.

**Figure 6.14: Final Scatter Plot**



It appears that there is a positive relationship between Height and Weight. Also, the taller and heavier data is male-dominated.

## Conclusion

After seeing how to run the charts and graphs demonstrated in this chapter, you should have no trouble running any of the graph tasks included with SAS Studio. This author, for one, is immensely grateful that someone has already done all the behind-the-scenes work to make it so easy to create charts and graphs. Thank you, SAS!

# Part II: Learning How to Write Your Own SAS Programs

- [Chapter 7 An Introduction to SAS Programming](#)
- [Chapter 8 Reading Data from External Files](#)
- [Chapter 9 Reading and Writing SAS Data Sets](#)
- [Chapter 10 Creating Formats and Labels](#)
- [Chapter 11 Performing Conditional Processing](#)
- [Chapter 12 Performing Iterative Processing: Looping](#)
- [Chapter 13 Working with SAS Dates](#)
- [Chapter 14 Subsetting and Combining SAS Data Sets](#)
- [Chapter 15 Describing SAS Functions](#)
- [Chapter 16 Working with Multiple Observations per Subject](#)
- [Chapter 17 Describing Arrays](#)
- [Chapter 18 Displaying Your Data](#)
- [Chapter 19 Summarizing Your Data with SAS Procedures](#)
- [Chapter 20 Computing Frequencies](#)
- [Appendix Solutions to Odd-Numbered Problems](#)

**Part II** shows you how to write your own SAS programs and use SAS procedures to perform a variety of tasks. This section also explains how to read data from a variety of sources, including text files, Excel workbooks, and CSV files. In order to help you become familiar with the SAS Studio environment, the book also shows you how to access the dozens of interesting data sets that are included with the product.

# Chapter 7: An Introduction to SAS Programming

## SAS as a Programming Language

This section of the book is dedicated to teaching you how to write your own programs in SAS. Perhaps you have some programming experience with other languages, such as C+, Python, or Java. This is both an advantage and a possible disadvantage. The advantage is that you understand how to think logically and use conditional logic, such as IF-THEN-ELSE statements and DO loops. On the other hand, SAS is somewhat unique in the way it reads and processes data, so you need to “re-wire” your brain and start to think like a SAS programmer.

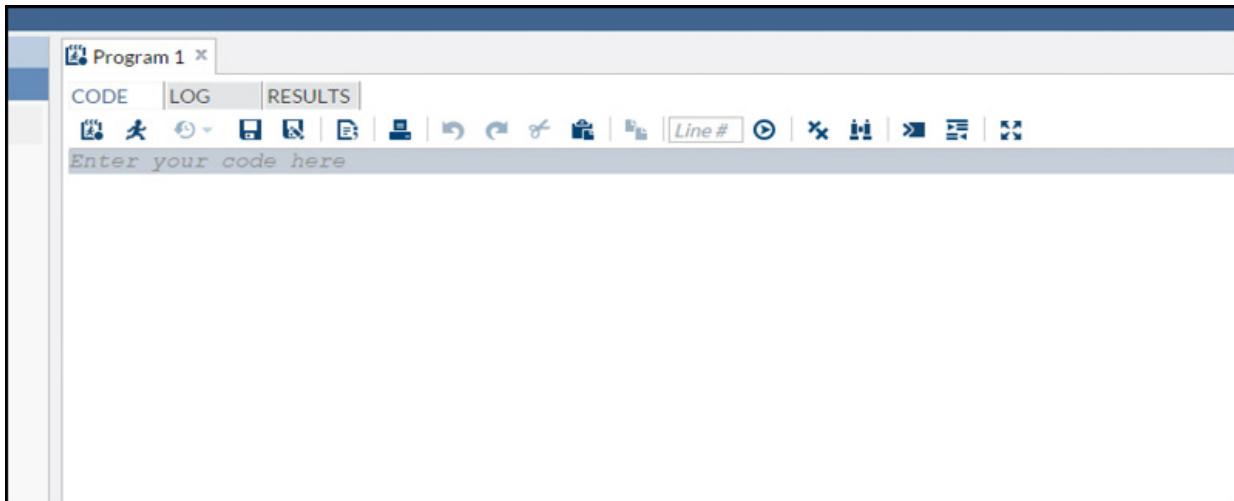
SAS programs consist of DATA steps, where you read, write, and manipulate data and PROC (short for procedure) steps, where you use built-in procedures to accomplish tasks such as writing reports, summarizing data, or creating graphical output. DATA steps begin with the keyword DATA and usually end with a RUN statement. PROC steps begin with the word PROC (did you guess that?) and end with either a RUN or QUIT statement (or both).

SAS statements all end with semicolons. This is a good place to mention that one of the most common programming mistakes, especially with beginning SAS programmers, is to forget a semicolon at the end of a statement. This sometimes leads to confusing error messages.

## The SAS Studio Programming Windows

When you open up SAS Studio, you see three tabs: **CODE**, **LOG**, and **RESULTS**. (See Figure 7.1 below.)

**Figure 7.1: The Three SAS Studio Windows**



The CODE window is where you enter your SAS program. When you run a SAS program, the LOG window displays your program, any syntax errors detected by SAS, information about data that was read or written out, and information about real time and CPU time used. The RESULTS window is where any SAS output appears. You can navigate among the three windows by clicking the appropriate tab.

## Your First SAS Program

Let's first write a simple program and then follow what happens when it runs. Suppose you want a program that converts temperatures from Celsius to Fahrenheit. It is a good idea to start your program by writing a comment statement. As part of the comment, you should, at a minimum, state the purpose of the program. In a more formal setting, you might also include information such as who wrote the program, the date on which it was written, and the location of input and/or output files. One way of writing comments in SAS programs is to start the comment with an asterisk and end it with a semicolon.

### DATA Statement

The next line is a DATA statement where you give a name to the data set you are going to create. Look what happens as you start to write the word DATA:

**Figure 7.2: Illustrating the Autocomplete Feature of SAS Studio**

The screenshot shows the SAS Studio interface with a code editor window titled "Program 1". The code window has tabs for "CODE", "LOG", and "RESULTS". Below the tabs is a toolbar with various icons. The main area contains the following code:

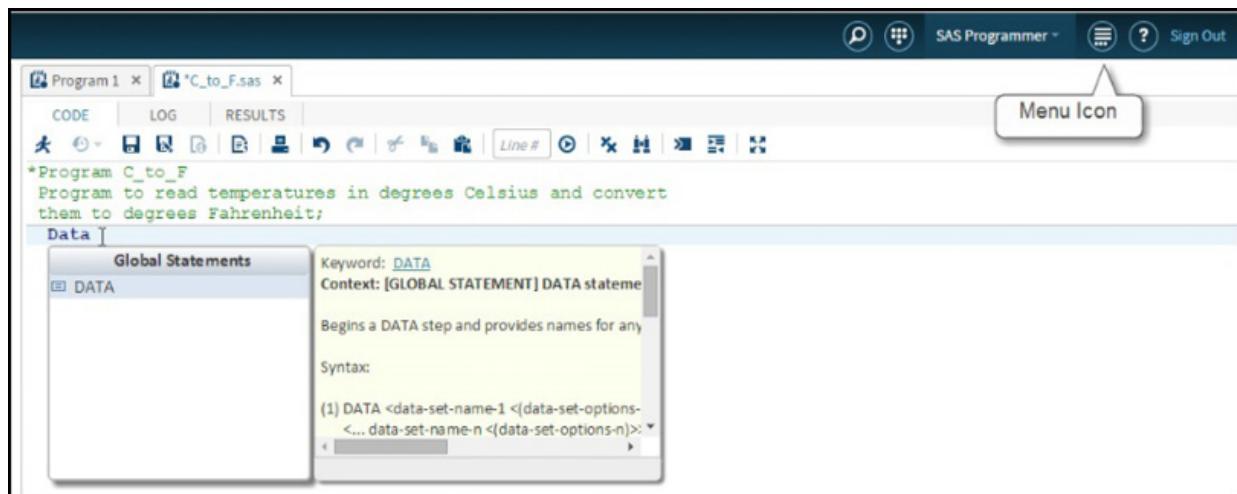
```
*Program C_to_F
Program to read temperatures in degrees Celsius and convert
them to degrees Fahrenheit;
```

As the user types "DATA", a context-sensitive help box appears. The box is titled "Global Statements" and "DATA". It provides the following information:

- Keyword:** [DATA](#)
- Context:** [GLOBAL STATEMENT] DATA statement
- Begins a DATA step and provides names for any output SAS data set
- Syntax:**
- (1) DATA <data-set-name-1 <(data-set-options-1)>>
<... data-set-name-n <(data-set-options-n)>>> </ <DEBUG> <NEST
- (2) DATA \_NULL\_ </ <DEBUG> <NESTING> <STACK = stack-size>> <N

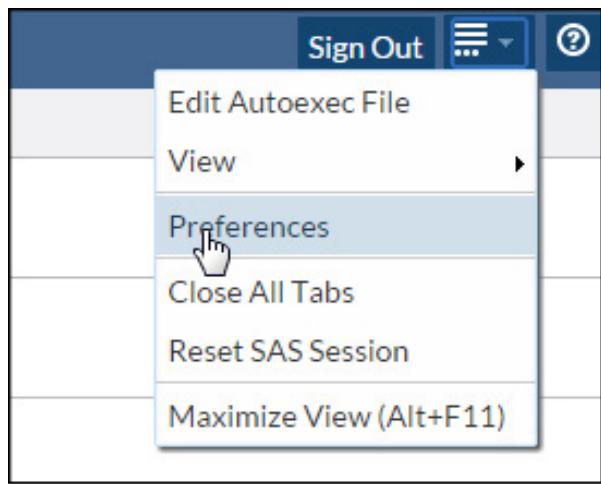
As you type certain keywords in the CODE window, context-sensitive boxes pop up to show you syntax and options that are available for the statement that you are writing. If you are a more advanced user who does not need this syntax help, you can turn it off by clicking the menu icon (on the top right of the screen) as follows:

**Figure 7.3: SAS Studio Options**



Choose the **Preferences** tab on this menu.

**Figure 7.4: Selecting SAS Studio Preferences**



Under the **Editor** tab, select (or deselect) **Enable autocomplete** and click **Save**.

**Figure 7.5: Select or Deselect Autocomplete**

A screenshot of the SAS Preferences dialog box. On the left, there's a sidebar with tabs: General, Start Up, Code and Log (which is selected and highlighted in light gray), Results, Tables, and Tasks. The main area is titled 'Editor options'. It contains several settings with checkboxes:

- Enable autocomplete (Ctrl+spacebar or Command+spacebar) - This checkbox is checked.
- Enable hint

Below these are input fields for 'Tab width' (set to 3 spaces) and 'Font size' (set to 16). There are also checkboxes for 'Substitute spaces for tabs', 'Enable color coding', and 'Show line numbers'. Further down are settings for 'Autosave': 'Enable autosave' (checked) and 'Autosave interval' (set to 30 seconds). The 'Log options' section includes checkboxes for 'Show generated code in the SAS log' (unchecked) and 'Stream log updates while a procedure is running' (checked). The 'With each submission' section has two radio buttons: 'Automatically clear log' (selected) and 'Append log'. At the bottom right is a dropdown menu set to 'Program and Task'.

The complete program is listed below.

### **Program 7.1: Program to Read Temperatures in Degrees Celsius and Convert Them to Fahrenheit**

```
*Program C_to_F
Program to read temperatures in degrees Celsius and convert
them to degrees Fahrenheit;

Data Convert;
  infile "~/MyBookFiles/celsius.txt";
  input Temp_C;
  Temp_F = 1.8*Temp_C + 32;
run;

title "Temperature Conversion Chart";

proc print data=Convert;
  var Temp_C Temp_F;
run;
```

This program starts with the keyword DATA. In this program you name the data set Convert. The rules for naming data sets and many other SAS names (such as variable names) are as follows:

In SAS, data set names and variable names must start with a letter or underscore. They can contain a maximum of 32 characters, and the remaining characters must be letters, digits, or underscores.

The following tables show examples of valid and invalid SAS names.

### Valid SAS Names

My\_Data

HeightWeight

---

X123

---

\_123

---

Price\_per\_pound

---

### Invalid SAS Names

My Data

Contains an invalid character  
(space)

---

---

123xyz	Starts with a digit
Temperature-Data	Contains an invalid character (-)
Group%	Contains an invalid character (%)

---

In SAS, variable names are not case-sensitive. However, the case that you use the first time you reference a variable is used in SAS output, regardless of how you write the variable name in other locations in the program.

## INFILE Statement

The INFILE statement tells the program where the raw data is located (use a different statement if you have cooked data). In this book, all the data from your hard drive was uploaded to the folder MyBookFiles (described in Chapter 3). In this author's account, the actual location for this file is /home/ronaldcody/MyBookFiles/Celsius.txt. A shortcut to refer to this file is to place a tilde (~) to represent the home directory, as shown in the INFILE statement above.

File names in Windows are not case sensitive; however SAS Studio is running in a Linux operating system on the SAS® Cloud platform. File names **are** case sensitive in Linux (also in UNIX).

For example, if you wrote Celsius.txt instead of celsius.txt the program would not be able to find the file. This author was just reminded of that fact when he ran a program and was momentarily

confused as to why the error message saying that the **physical file could not be found** was shown in the SAS LOG.

If you want to run the programs in this book, remember that you can download all of the programs and data sets from the following location:

[support.sas.com/cody](http://support.sas.com/cody)

Once you arrive at this location, locate this book and click the link "Example Code and Data." This action downloads a ZIP file containing all the programs and data files in the book to your hard drive. You will want to extract these files to your hard drive and upload them to the MyBookFiles folder in SAS Studio (or whatever name you choose). You will also want to perform a similar action on the files and programs needed for the end-of-chapter problems. The solutions at the end of this book use a folder called Problems for these files and programs. You can use this name or any other name of your choosing.

## INPUT Statement

The INPUT statement is an instruction to read data from the celsius.txt file. This text file contains one number per line and is listed below.

### File 7.1: celsius.txt

```
0  
100  
20
```

This INPUT statement uses one of three methods of reading data, called *list input*. When you use list input, you can read data values separated by blanks (the SAS default delimiter) or other delimiters such as commas. Information about how to process data with delimiters other than blanks is presented in the next chapter.

## Assignment Statement

Following the INPUT statement, you use an assignment statement to code the formula for the Celsius to Fahrenheit conversion. As with most programming languages, the calculation to the right of the equal

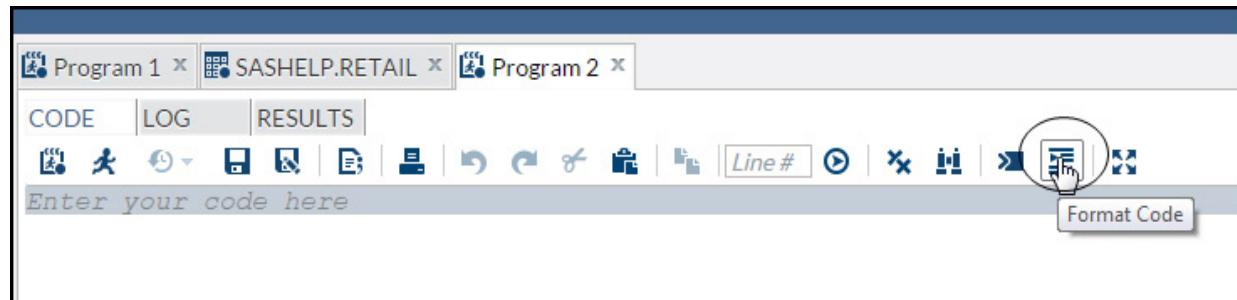
sign is assigned to the variable name to the left of the equal sign. In this example, the calculated value is assigned to the variable Temp\_F.

You use an asterisk to indicate multiplication, a forward slash for division, plus and minus signs for addition and subtraction, and two asterisks for exponentiation. Exponentiation is performed first, multiplication and division next, followed by addition or subtraction. You can always use parenthesis to determine the order of operations.

The DATA step ends with a RUN statement. In this program, the statements in the DATA step and PROC step are indented. This is not necessary, but it makes the program easier to read. It is also possible to place more than one SAS statement on a single line, as long as each statement ends with a semicolon. However, this practice is discouraged because it makes it difficult to read and understand the program. Finally, a SAS statement can use as many lines as necessary, such as the comment statement in this program. Just remember to end the statement with a semicolon.

SAS Studio has an auto formatting feature that you can use to automatically format your SAS programs. After you have written your program in the CODE window, click the **AUTOFORMAT** icon at the top of the **EDITOR** window (as shown in the figure below).

**Figure 7.6: Activating the Auto Formatting Feature of SAS Studio**



The result is a nicely formatted program. You can use CTRL+Z to undo the formatting in case you don't like the result.

## How the DATA Step Works

Unlike most programming languages, the SAS DATA step is actually

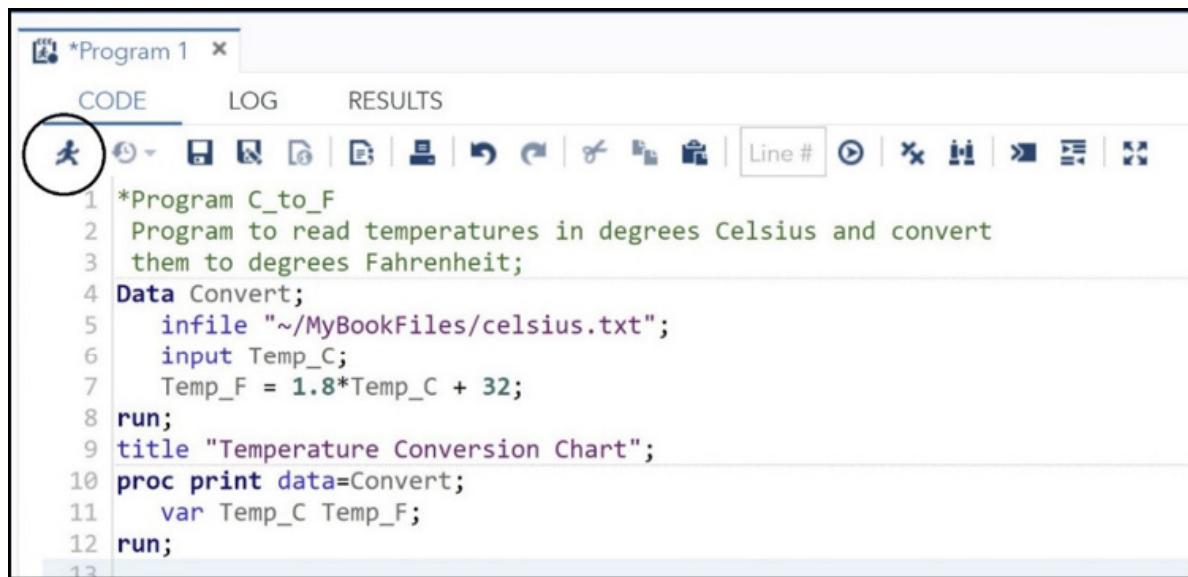
an automatic loop: The first time the INPUT statement executes, the program reads a value from the first line of data. It then computes the corresponding Fahrenheit temperature according to your formula and adds the newly created variable to the data set SAS is creating. Your next statement is a RUN statement that marks the end of the DATA step. Two things happen at this point: First, the program automatically outputs an observation (containing the variables Temp\_C and Temp\_F) to the output data set that you named Convert. Next, the program performs its implied loop by returning to the top of the DATA step to execute the INPUT statement again. On the second iteration of the DATA step, SAS reads data from the second line of data. Each time the DATA step iterates, the INPUT statement goes to a new line of data (unless you give it special instructions not to). In this example, the DATA step stops when it tries to read the fourth line of data from the file and encounters an end-of-file marker. At this point, your SAS data set Convert contains three observations.

You use PROC PRINT to list the contents of your SAS data set. In this example, you specify the name of the SAS data set using the procedure option DATA=. You can specify one or more title lines with a TITLE statement. You can place a TITLE statement before or after the PROC statement. When you specify a title, that title remains in effect until you replace it with another TITLE statement. In this program, you placed the TITLE statement between the DATA and PROC steps, a location referred to as *open code*.

You use a VAR statement to specify which variables you want to include in your report. The order you use to list the variables on the VAR statement is the order that PROC PRINT will use to print the results. If you leave out a VAR statement, PROC PRINT will print every variable in your data set in the order they are stored in the data set. You end the PROC step with a RUN statement.

You are now ready to run the program. You do this by clicking the RUN icon as shown in Figure 7.7.

### Figure 7.7: The RUN Icon

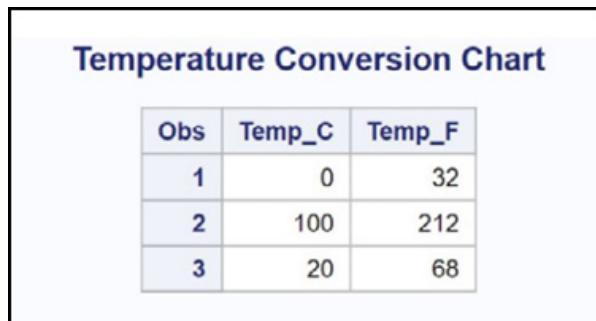


The screenshot shows the SAS Studio interface with the title bar "\*Program 1". Below it are three tabs: CODE (which is selected), LOG, and RESULTS. A toolbar with various icons follows. The main area contains the following SAS code:

```
1 *Program C_to_F
2 Program to read temperatures in degrees Celsius and convert
3 them to degrees Fahrenheit;
4 Data Convert;
5   infile "~/MyBookFiles/celsius.txt";
6   input Temp_C;
7   Temp_F = 1.8*Temp_C + 32;
8 run;
9 title "Temperature Conversion Chart";
10 proc print data=Convert;
11   var Temp_C Temp_F;
12 run;
13
```

SAS Studio now shows you the results.

**Figure 7.8: The RESULTS Window**

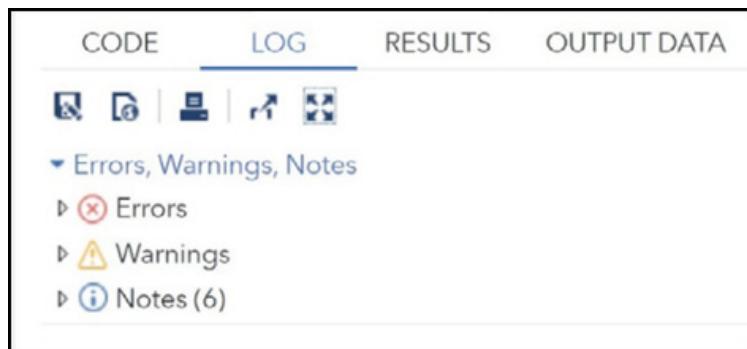


The screenshot shows the SAS Studio interface with the RESULTS tab selected. It displays a table titled "Temperature Conversion Chart" with the following data:

Obs	Temp_C	Temp_F
1	0	32
2	100	212
3	20	68

By default, SAS lists all rows (called *observations* in SAS terminology) and all columns (called *variables* by SAS). It also includes an Obs (observation number) column. A quick examination of the output shows that the program worked correctly. However, **you should always look at the SAS log**, even when you have output that seems to be correct. To examine the log, click the LOG tab. Below is a listing of the log.

**Figure 7.9: The LOG Window**



The first part of the log shows there were no errors or warnings. You can click any of these items to display any errors, warnings, or notes. Excerpts from the remaining log follow.

The first section shows your program, along with information about your input data file. Notice that SAS Studio added an OPTIONS statement to your program. This statement controls what information is displayed in the log and how the data appears in the RESULTS window.

```
1      OPTIONS NONOTES NOSTIMER NOSOURCE NOSYNTAXCHECK;
70
71      *Program C_to_F
72      Program to read temperatures in degrees Celsius and convert
73      them to degrees Fahrenheit;
74      Data Convert;
75          infile "~/MyBookFiles/celsius.txt";
76          input Temp_C;
77          Temp_F = 1.8*Temp_C + 32;
78      run;

NOTE: The infile "~/MyBookFiles/celsius.txt" is:
      Filename=/home/ronaldcody/MyBookFiles/celsius.txt,
      Owner Name=ronaldcody, Group Name=oda,
      Access Permission=-rw-r--r--,
      Last Modified=16Oct2020:15:42:16,
      File Size (bytes)=10
```

Next, you see that three records (lines) were read from the input file.

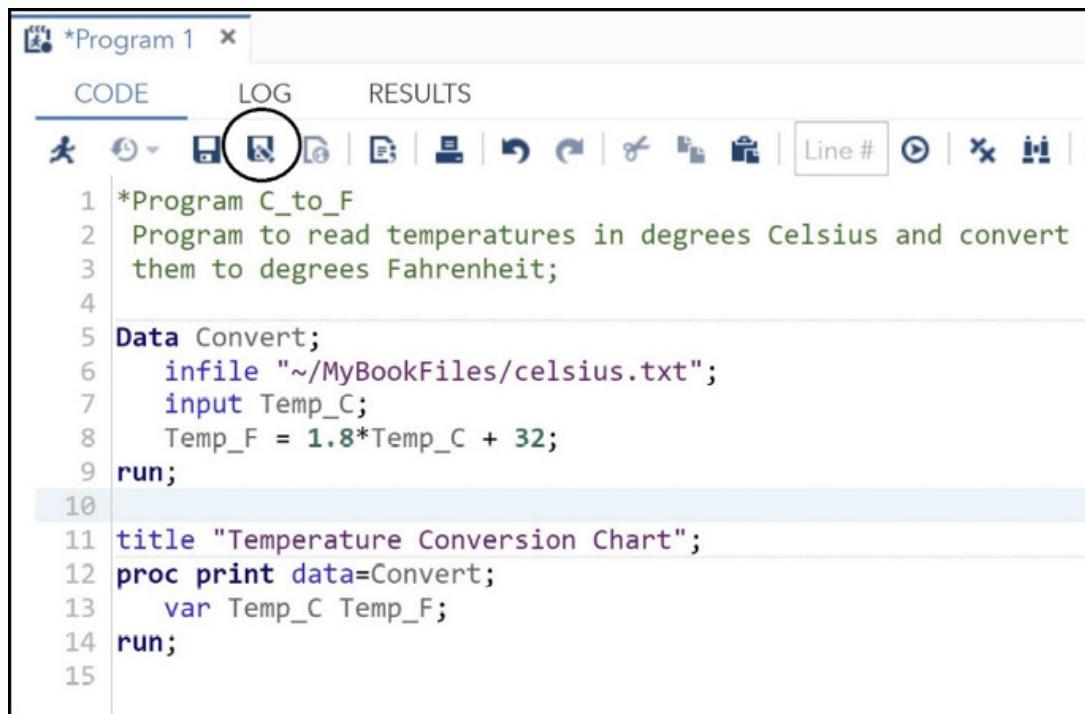
```
NOTE: 3 records were read from the infile "~/MyBookFiles/celsius.txt".  
      The minimum record length was 1.  
      The maximum record length was 3.  
NOTE: The data set WORK.CONVERT has 3 observations and 2 variables.  
NOTE: DATA statement used (Total process time):  
      real time          0.11 seconds  
      user cpu time     0.00 seconds  
      system cpu time   0.00 seconds
```

Finally, you see that PROC PRINT read three observations. You also see the real and CPU times.

```
79      title "Temperature Conversion Chart";  
80      proc print data=Convert;  
81          var Temp_C Temp_F;  
82      run;  
  
NOTE: There were 3 observations read from the data set WORK.CONVERT.  
NOTE: PROCEDURE PRINT used (Total process time):  
      real time          0.02 seconds  
      user cpu time     0.02 seconds  
      system cpu time   0.01 seconds
```

Before you close your SAS session, you should save your program so that you can work on the program later or, if it is finished, to save it. From the **CODE** tab, select the icon for **SAVE AS**:

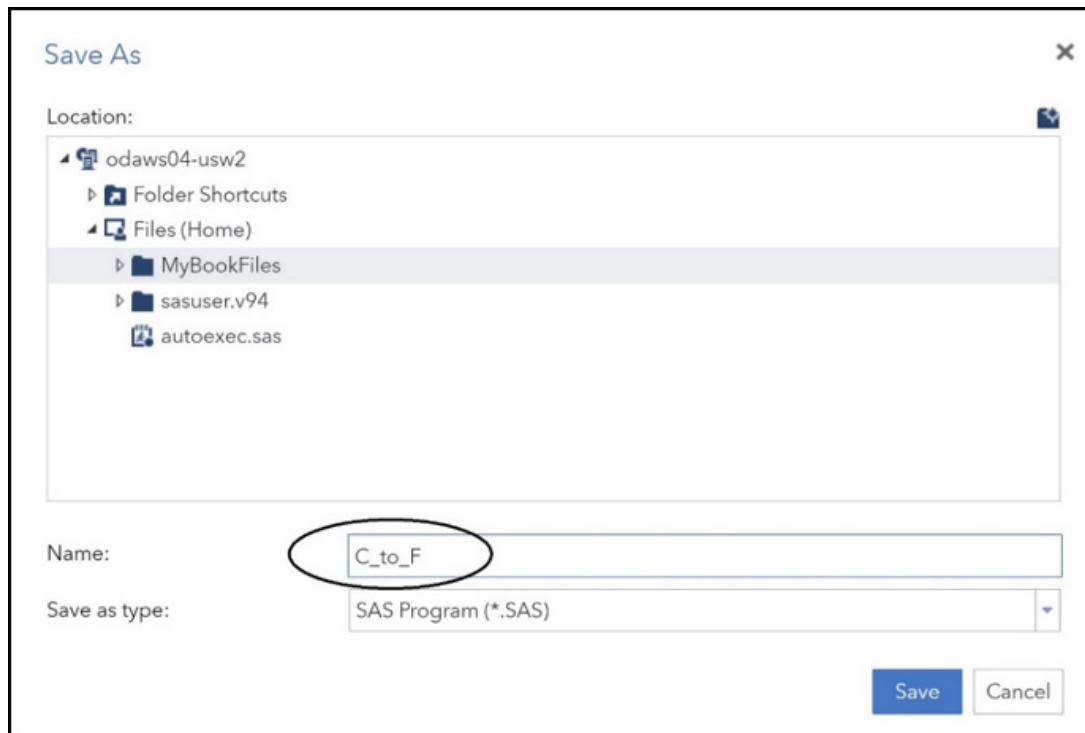
**Figure 7.10: The SAVE AS Icon**



```
*Program 1 x
CODE LOG RESULTS
1 *Program C_to_F
2 Program to read temperatures in degrees Celsius and convert
3 them to degrees Fahrenheit;
4
5 Data Convert;
6   infile "~/MyBookFiles/celsius.txt";
7   input Temp_C;
8   Temp_F = 1.8*Temp_C + 32;
9 run;
10
11 title "Temperature Conversion Chart";
12 proc print data=Convert;
13   var Temp_C Temp_F;
14 run;
15
```

This brings up the following screen.

**Figure 7.11: Saving Your Program in My Folders**

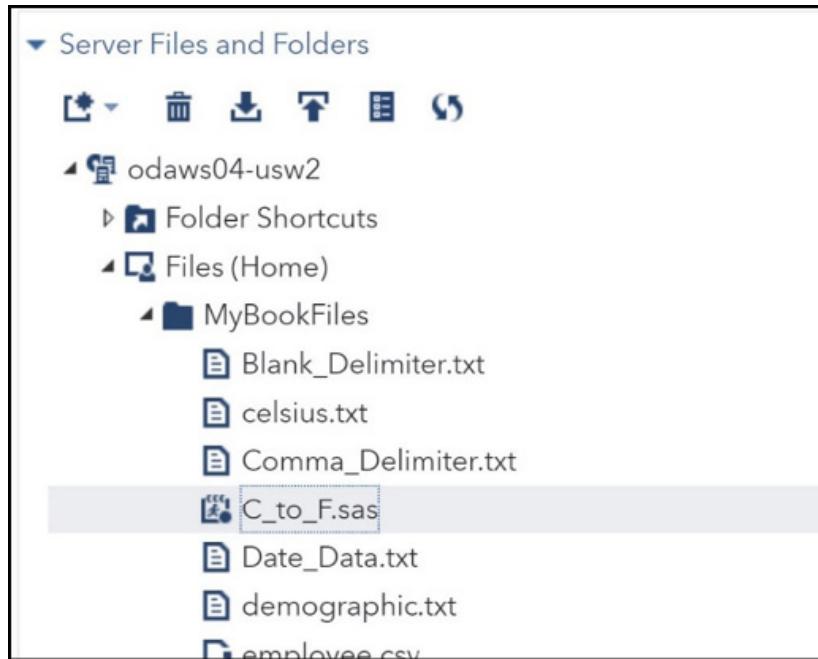


First, click the **MyBookFiles** folder. Next enter the program name (in this example, it is **C\_to\_F.sas**) and click **Save**. You can name your SAS

program anything you like (as long as it meets the naming conventions for files on your computer). The extension .SAS is automatically added to the file name.

If you look in the **Navigation** Pane, you will now see the program C\_to\_F.sas in the list (Figure 7.12).

**Figure 7.12: Contents of the MyBookFiles Folder**



## How the INPUT Statement Works

SAS programs can read just about any type of text data, whether it consists of numbers and letters, separated by delimiters (such as a CSV file), or whether the data values are arranged in fixed columns. This section just scratches the surface of the incredible versatility of this statement.

## Reading Delimited Data

Let's start out by reading a file where data values (either character or numeric) are separated by delimiters. The following file, called demographic.txt (stored in the MyBookFile folder), contains the following data values:

### Variables in the File demographic.txt

ID  
Gender (M or F)  
Age  
Height (in inches)  
Weight (in pounds)  
Party (political party affiliation (I=Independent, R=Republican, D=Democrat))

Here is the file.

### File 7.2: demographic.txt

```
012345 F 45 65 155 I
131313 M 28 70 220 R
987654 F 35 68 180 R
555555 M 64 72 165 D
172727 F 29 62 102 I
```

You want to create a SAS data set called Demo from this raw data file. Proceed as follows:

### Program 7.2: Reading Text Data from an External File

```
data Demo;
  infile "~/MyBookFiles/demographic.txt";
  input ID $ Gender $ Age Height Weight Party $;
run;
```

You simply list the variable names in the same order as the data values in the file. As you can probably guess, a dollar sign (\$) following a variable name indicates that you want to read and store this value as character data. Notice that ID is being stored as a character value so that any leading zeros will be maintained. This is a good time to mention that SAS has only two variable types: character and numeric. By default, all numeric values are stored in 8 bytes (64 bits). In most programs, you specify the length of character variables. In this first program, because no character variable lengths are specified, a default length of 8 bytes (characters) is used to store each of the character values (even though Gender and Party are only one character in length). You will see several methods of determining the storage length for character data later in this book.

Here is the SAS log that is produced when you run this program.

### Figure 7.13: The SAS Log from Program 7.2

```
71      data Demo;
72          infile "~/MyBookFiles/demographic.txt";
73              input ID $ Gender $ Age Height Weight Party $;
74      run;

NOTE: The infile "~/MyBookFiles/demographic.txt" is:
      Filename=/home/ronaldcody/MyBookFiles/demographic.txt,
      Owner Name=ronaldcody,Group Name=oda,
      Access Permission=-rw-r--r--,
      Last Modified=16Oct2020:15:42:16,
      File Size (bytes)=108

NOTE: 5 records were read from the infile "~/MyBookFiles/demographic.txt".
      The minimum record length was 20.
      The maximum record length was 20.
NOTE: The data set WORK.DEMO has 5 observations and 6 variables.
NOTE: DATA statement used (Total process time):
      real time            0.12 seconds
      user cpu time        0.00 seconds
      system cpu time     0.00 seconds
```

You see information about the input data file and note that 5 records were read (as expected) and the fact that a SAS data set, WORK.DEMO, was created. Unless you specify a storage location for your SAS data set, SAS places it in the WORK library. Data sets in the WORK library are temporary and disappear when you close your SAS session. You will see how to read and write permanent SAS data sets in Chapter 9. In the SAS log, you also see information about the real and CPU time used.

If you click the **OUTPUT DATA** tab, you will see a list of variables and your data like this:

**Figure 7.14: Contents in the OUTPUT DATA Tab**

Table: WORK.DEMO | View: Column names | Filter: (none)

Columns Total rows: 5 Total columns: 6

	ID	Gender	Age	Height	Weight	Party
1	012345	F	45	65	155	I
2	131313	M	28	70	220	R
3	987654	F	35	68	180	R
4	555555	M	64	72	165	D
5	172727	F	29	62	102	I

To adjust the column widths in this display, this author first right-clicked on one of the columns in the top row of the table and selected the option “Size grid columns to content”. (See Figure 7.15.)

**Figure 7.15: Resizing Columns**

Table: WORK.DEMO | View: Column names | Filter: (none)

Columns Total rows: 5 Total columns: 6

	ID	Gender	Age	Height	Weight	Party
1	012345	F	45	65	155	I
2	131313	M	28	70	220	R
3	987654	F	35	68	180	R
4	555555	M	64	72	165	D
5	172727	F	29	62	102	I

## How Procedures (PROCs) Work

As you saw in the first part of this book, you can now use the **TASKS** tab to summarize the data or produce plots. However, because this is

the programming section of the book, let's use PROC FREQ to compute frequencies and demonstrate how procedure options and statement options affect the results.

### Program 7.3: Computing Frequencies from the Demo Data Set

```
data Demo;
  infile("~/MyBookFiles/demographic.txt");
  input ID $ Gender $ Age Height Weight Party $;
run;

title "Computing Frequencies from the Demo Data Set";

proc freq data=Demo;
  tables Gender Party;
run;
```

You use PROC FREQ to compute frequencies on any of your variables. Use a TABLES statement to specify which variables you want to include in your results. These can include both character and numeric variables. If you include any numeric variables in the list, PROC FREQ will compute the frequency on every unique value—it does not group values into bins (there are other procedures that can produce histograms). If you do not use a TABLES statement, PROC FREQ will compute frequencies on all the variables in your data set.

The output from Program 7.3 is shown below.

**Figure 7.16: Output from Program 7.3**

## Computing Frequencies from the Demo Data Set

The FREQ Procedure

Gender	Frequency	Percent	Cumulative Frequency	Cumulative Percent
F	3	60.00	3	60.00
M	2	40.00	5	100.00

Party	Frequency	Percent	Cumulative Frequency	Cumulative Percent
D	1	20.00	1	20.00
I	2	40.00	3	60.00
R	2	40.00	5	100.00

In the frequency tables, you see frequency, percent, cumulative frequency, and cumulative percent for the two variables Gender and Party. The output you see results from the default settings for this procedure. In most cases, you will want to include procedure and statement options to control the output.

Most SAS procedures have what are called *procedure options*. These options affect how the procedure works, and they are placed between the procedure name and the semicolon. One popular procedure option used with PROC FREQ is ORDER=. There are several values that you can select for this option. For this example, if you use ORDER=FREQ, the frequencies are ordered from the most frequent to the least frequent.

You will most likely use *procedure statements* with most procedures as well. The TABLES statement in Program 7.3 is an example of a procedure statement. You can also add *statement options* to control how a statement works. The rule is that statement options are placed after a slash (/) following the statement. The TABLES option NOCUM (no cumulative statistics) is used in the next program to demonstrate a statement option. Let's also see how the procedure option ORDER=FREQ and the statement option NOCUM affect the output.

### Program 7.4: Adding Options to PROC FREQ

```
title "Adding Procedure Options and Statements";  
Proc freq order=freq;  
    tables Gender Party / nocum;  
run;
```

Your output now looks like this.

**Figure 7.17: Output from Program 7.4**

Adding Procedure Options and Statements		
Gender	Frequency	Percent
F	3	60.00
M	2	40.00

Party	Frequency	Percent
I	2	40.00
R	2	40.00
D	1	20.00

The tables are now displayed in decreasing frequency order (notice the change in the order of Party), and cumulative statistics are no longer included in the tables. In most cases, you will want to include the NOCUM option on your TABLES statement.

## How SAS Works: A Look Inside the “Black Box”

Although you can write basic SAS programs without understanding what goes on inside the “black box,” a more complete understanding of how SAS works will make you a better programmer. Furthermore, this knowledge is essential when you are writing more advanced programs.

SAS processes the DATA step in two stages: In the first stage, called the *compile stage*, several activities take place. The SAS complier reads each line of code from left to right, top to bottom. Each statement is broken up into tokens, with certain keywords such as DATA and RUN causing certain actions to take place. What is important to you as a programmer, is to know that this is the stage

where the data descriptor for each of your variables is written out. The first time SAS encounters a variable in a DATA step, it decides whether that variable is numeric or character. If it is numeric, SAS gives it a default length of 8 bytes. If it is character, it has rules that it uses to determine the storage length. SAS character values can be a maximum of 32,767 bytes in length. In this first stage, your source code is also compiled into machine language. Also, during the compile stage, SAS determines which variables will be written out to the new data set and which variables will be dropped (i.e., not written out). One way of determining whether a variable is kept or dropped is by explicitly writing a KEEP or DROP statement in the DATA step or by including a KEEP= or DROP= data set option (more on that later).

In the second stage, called the *execute stage*, the program performs its functions of reading data, performing logical actions, iterating loops, and so on. When you are reading raw data from a file or if you have variables defined in assignment statements (such as the variable Temp\_F in Program 7.1), SAS initializes each of these variables with a missing value at the top of the DATA step. During execution of the DATA step, these variables are usually given values, either from the raw data or from a computation. At the bottom of the DATA step (defined by the RUN statement), SAS performs an automatic output to one or more data sets. The DATA step continues its internal loop, reading data, performing calculations, and outputting observations to a data set. If you are reading data from a text file or from a SAS data set, the DATA step stops when you reach the end-of-file marker on any file.

## Conclusion

At this point, you understand how to write a simple program using the SAS OnDemand for Academics. You understand the various windows inside SAS Studio. You can read external data where blanks are used as delimiters and produce simple reports. The next chapter explores the INPUT statement, one of the most powerful statements in SAS.

# Chapter 8: Reading Data from External Files

## Introduction

This chapter describes three of the most common methods that you can use to read raw text data using SAS. It also demonstrates how to create SAS data sets from CSV files. If you typically receive data that is already in SAS data sets, you can skip this chapter (unless you are just curious).

## Reading Data Values Separated by Delimiters

One method of storing data in text files is to separate data values by a delimiter, usually blanks or commas. SAS refers to this as *list input*.

Let's start out with a file where blanks (spaces) are used as delimiters. This is a good starting place because a blank is the default delimiter in SAS. In this example, the data file that you want to read contains an ID, first name, last name, gender, age, and state abbreviation. These data lines are stored in a file called Blank\_Delimiter.txt and was uploaded to the MyBookFiles folder. Here is the listing.

### **File 8.1: Blank\_Delimiter.txt (Previously uploaded to the MyBookFiles folder)**

```
103-34-7654 Daniel Boone M 56 PA
676-10-1020 Fred Flintstone M
454-30-9999 Tracie Wortenberg F 34 NC
102-87-8374 Jason Kid M 23 NJ
888-21-1234 Patrice Marcella F . TX
788-39-1222 Margaret Mead F 77 PA
```

A program to read this data file is shown next:

### **Program 8.1: Reading Data with Delimiters (Blanks)**

```
/* This is another way to insert a comment */

data Blanks /* the data set name is Blanks */;
  infile "~/MyBookFiles/Blank_Delimiter.txt" missover;
  informat ID $11. First Last $15. Gender $1. State_Code $2. ;
  input ID First Last Gender Age State_Code;
run;

title "Listing of Data Set Blanks";
proc print data=Blanks;

run;
```

The first line of this program demonstrates another way to insert a comment in a SAS program. You begin the comment with a /\* and end the comment with a \*/. Unlike a comment statement that begins with an asterisk, this type of comment can be placed within a SAS statement (as demonstrated in the DATA statement of this program.).

The INFILE statement tells the SAS program where to find the input data. In this example, the data file was already uploaded to the folder MyBookFiles. As discussed briefly in the previous chapter, the actual location of the file is more complicated. It is:

/home/ronaldcody/MyBookFiles/Blank\_Delimiter.txt

Because this path is so complicated, you can use a notation common on many other operating systems and use a tilde (~) to specify your home directory.

Notice the keyword MISSOVER on this statement. This is not necessary if you have data values for every variable for every line of data. Line 2 of the data file (Fred Flintstone) is missing the last two values (Age and State\_Code). **Without the MISSOVER option, SAS would try to read these two values from the next line of data!** Obviously, you never want this to happen. The MISSOVER option comes into play when there are missing values at the end of the data line and you are using list input. This option tells the program to set each of these variables to a missing value. Notice line 5 (Patrice Marcella). To indicate that there is a missing value for Age, you use a period. Without the period, SAS would try to read the next value (TX) as the

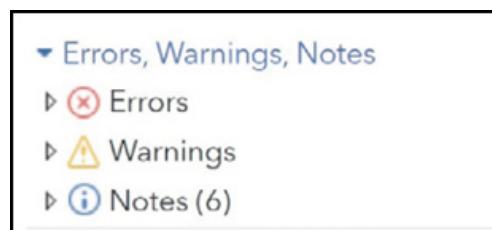
Age and really screw things up. This is handled differently in the next section describing CSV (comma-separated values) files. One last point: SAS interprets multiple blanks as a single delimiter.

You list each of the variable names in the INPUT statement in the order in which they appear in the raw data file. When you use list input (the type of input that reads delimited data), the default length for character variables is 8. To specify how to read your character variables, you use an INFORMAT statement. Following the keyword INFORMAT, you list one or more variables (usually character variables) and follow the variable or variables with a dollar sign (that indicates a character variable) and the number of characters (bytes) of storage that you want, followed by a period. In this example, \$11., \$15., \$1., and \$2. are called *character informats*. Informats are used in other styles of input as well.

Use of these informats also determines the storage length of each of the variables listed in the INFORMAT statement. In this example, you are specifying a length of 11 for ID, a length of 15 for both the variables First and Last, a length of 1 for Gender, and a length of 2 for State\_Code. You could also define the lengths of these character variables with a LENGTH statement. That statement would be similar to the INFORMAT statement—just replace the keyword INFORMAT with the keyword LENGTH and leave off the periods after the digits. An INFORMAT statement can also be used to tell SAS how to read other types of data such as dates and numbers with commas and dollar signs, therefore making it more flexible than a LENGTH statement for instructing SAS on how to read raw data.

To list the observations in the Blanks data set, you use PROC PRINT. Below are sections from the SAS log.

**Figure 8.1: SAS Log from Program 8.1**



```
71      /* This is another way to insert a comment */
72
73      data Blanks /* the data set name is Blanks */;
74          infile "~/MyBookFiles/Blank_Delimiter.txt" missover;
75          informat ID $11. First Last $15. Gender $1. State_Code $2.;
76          input ID First Last Gender Age State_Code;
77      run;
```

NOTE: The infile "~/MyBookFiles/Blank\_Delimiter.txt" is:  
Filename=/home/ronaldcody/MyBookFiles/Blank\_Delimiter.txt,  
Owner Name=ronaldcody,Group Name=oda,  
Access Permission=-rw-r--r--,br/>Last Modified=16Oct2020:15:42:16,  
File Size (bytes)=205

NOTE: 6 records were read from the infile "~/MyBookFiles/Blank\_Delimiter.txt"  
The minimum record length was 29.  
The maximum record length was 37.

NOTE: The data set WORK.BLANKS has 6 observations and 6 variables.

NOTE: DATA statement used (Total process time):  
real time 0.00 seconds  
user cpu time 0.00 seconds  
system cpu time 0.00 seconds

You see that there were no errors in the program, along with information about your input file. Finally, you see that your data set (Blanks) was created with 6 observations and 6 variables, along with the real and CPU time.

Even when you see the OUTPUT window after you submit a program, it is a good idea to check the SAS log for messages and warnings.

Here is the output.

**Figure 8.2: Output from Program 8.1**

### Listing of Data Set Blanks

Obs	ID	First	Last	Gender	State_Code	Age
1	103-34-7654	Daniel	Boone	M	PA	56
2	676-10-1020	Fred	Flintstone	M		.
3	454-30-9999	Tracie	Wortenberg	F	NC	34
4	102-87-8374	Jason	Kid	M	NJ	23
5	888-21-1234	Patrice	Marcella	F	TX	.
6	788-39-1222	Margaret	Mead	F	PA	77

Everything looks fine. Notice the missing values in observations 2 and 5. In observation 2, the missing values for State\_Code and Age result from the MISSOVER option in the INFILE statement. The missing value for Age in observation 5 results from the period in the input data.

## Reading Comma-Separated Values Files

CSV files (comma-separated values) are one of the most common file types for delimited data. One common use of CSV files is to output data from an Excel workbook. CSV files use commas to separate data values and, unlike the default behavior of SAS when you use blank delimiters, CSV files interpret two commas in a row to mean that there is a missing value for that data field.

The CSV file in this example contains the same data as the file Blank\_Delimiter.txt used in the last example. Here is a listing of the file:

### File 8.2: Comma\_Delimiter.txt

```
103-34-7654,Daniel,Boone,M,56,PA
676-10-1020,Fred,Flintstone,M
454-30-9999,Tracie,Wortenberg,F,34,NC
102-87-8374,Jason,Kid,M,23,NJ
888-21-1234,Patrice,Marcella,F,,TX
788-39-1222,Margaret,Mead,F,77,PA
```

Notice the two commas in a row for the fifth line of data (Patrice Marcella). Unlike the blank-delimited file where you needed a period to indicate a missing value for Age, the two commas in a row indicate

that this value is missing.

To read data from a CSV file, you add the DSD (Delimiter-Sensitive Data) option in the INFILE statement. This option does several things: First, it understands that the data values in the file are separated by commas. Next, it understands that two commas indicate a missing value. Also, if you have a data value in quotation marks (for example, a state name like ‘New York’ that consists of two words separated by a blank space), it ignores any delimiters inside the quotation marks and strips the quotation marks when it assigns the value to a SAS variable. The program to read the data from the CSV file Comma\_Delimiter.txt is shown next.

### Program 8.2: Reading CSV Files

```
data Commas;
  infile "~/MyBookFiles/Comma_Delimiter.txt" dsd missover;
  informat ID $11. First Last $15. Gender $1. State_Code $2.;
  input ID First Last Gender Age State_Code;
run;

title "Listing of Data Set Commas";
proc print data=Commas;
run;
```

If you see the following message in your SAS log:

**NOTE: SAS went to a new line when INPUT statement reached past the end of a line.**

you probably need to add the MISSOVER option to your INFILE statement.

Output from this program is identical to the output from Program 8.1.

### Reading Data Separated by Other Delimiters

By using the DLM= option in the INFILE statement, you can specify any delimiter you want, including non-printing characters such as tabs. Suppose you have a tab-delimited file that you want to read. Because a tab character is a non-printing character, you need to find

the hexadecimal code for a tab in ASCII (the coding system used in Windows, UNIX, and Linux systems). In ASCII, the Hex representation for a tab is 09. You can specify a Hex character anywhere in a SAS program by using a Hex constant in the form ‘nn’x, where *nn* is the Hex value that you want (placed in single or double quotation marks) and the ‘x’ is in lowercase or uppercase. Also, there are no spaces between the quoted Hex value and the ‘x’. The program below reads an ASCII file where a tab was used as the delimiter.

### Program 8.3: Reading Tab-Delimited Data

```
data Tabs;
  infile("~/MyBookFiles/Tab_Delimiter.txt" dlm='09'x missover;
  informat ID $11. First Last $15. Gender $1. State_Code $2.;
  input ID First Last Gender Age State_Code;
run;

title "Listing of Data Set Tabs";
proc print data=tabs;
run;
```

The output is identical to the output from the previous two programs. You can use the DLM= option in the INFILE statement to specify any delimiter you want. If you want two consecutive delimiters to indicate a missing value, include the DSD option as well as the DLM= option. For example, if you have data that uses a pipe symbol (vertical bar) as a delimiter and you want to interpret two bars together to indicate that there is a missing value, your INFILE statement would look like this:

```
infile "file-location" dlm='|' dsd;
```

## Reading Data in Fixed Columns

Another method of placing data in text files is to assign values to predefined columns.

This is much less common today than back in the “old days.” I’m talking about the time at which you entered your programs and data on punch cards. However, these types of file still exist and it’s a

good thing that SAS can read them with ease.

SAS provides you with two methods of reading column data. One is called *column input*. With this method, you follow a variable name with a starting and ending column. If the variable only occupies one column, you do not specify an ending column. If the variable will be defined as a character variable, you place a dollar sign (\$) between the variable name and the starting column number. This method is restricted to reading standard numeric data (numbers with or without decimal points) and character data.

The other method is called *formatted input*. With this method, you specify a starting column, the variable name, and what SAS calls an *informat*. SAS informats give SAS information about how to read data from one or more columns. Formatted input is much more flexible than column input because it can read data values such as dates and times. You will see examples of both of these methods in the sections that follow.

## Reading Data in Fixed Columns Using Column Input

For this example, you want to read data from a file called Health.txt. This file contains the following variables:

### Variables in the Health.txt Data File

Variable Name	Description	Starting Column	Ending Column
---------------	-------------	-----------------	---------------

Gender                    Gender (M or F) 4                    4

---

Age                        Age in years 5                        6

---

HR                        Heart rate 7                        8

---

SBP                      Systolic blood pressure 9                    11

---

DBP                      Diastolic blood pressure 12                14

---

Chol                     Total cholesterol 16                    18

---

The data file looks like this:

### File 8.3: health.txt

```
123456789012345678901234567890 (Ruler - this line is not in the  
file)  
001M2368120 90128  
002F5572180 90170  
003F1858118 72122  
004M8082      220  
005F3462128 80  
006F3878108 68220
```

Because the data values are in fixed columns, you can use column input to read it. Notice that blanks are used when there are missing values (although there are no blanks at the end of short lines). Below is a program that reads this data file using column input.

### Program 8.4: Reading Data in Fixed Columns Using Column Input

```
data Health;  
  infile '~/MyBookFiles/health.txt' pad;  
  input Subj    $ 1-3  
        Gender $ 4  
        Age     5-6  
        HR      7-8  
        SBP     9-11  
        DBP     12-14  
        Chol    15-17;  
  run;  
  
  title "Listing of Data Set Health";  
  proc print data=Health;  
    ID Subj;  
  run;
```

The INFILE statement tells the program where to find the health.txt data file. Following the file location, you see the keyword PAD. **This is very important.** Because some of the lines in the file are shorter than others (there is a carriage return after the last number in each line), SAS will not read the data correctly without it. This option pads each

line with blanks.

When programming was done on mainframe computers (often using punch cards as input), data lines were automatically padded with blanks (typically up to 80 columns). Because a lot of data is now entered on personal computers, data lines can be much longer than 80 columns and they are not typically padded with blanks.

Notice the file name “health.txt” is in lowercase. The reason is that the file that was uploaded to SAS Studio was also in lowercase. Whenever you see a “physical file not found” error message in your SAS log, check to see whether the file name in the INFILE statement matches the file name in your SAS Studio folder.

Here is a listing of data set Health:

**Figure 8.3: Output from Program 8.4**

Listing of Data Set Health						
Subj	Gender	Age	HR	SBP	DBP	Chol
001	M	23	68	120	90	128
002	F	55	72	180	90	170
003	F	18	58	118	72	122
004	M	80	82	.	.	220
005	F	34	62	128	80	.
006	F	38	78	108	68	220

Notice that the blanks in the raw data file result in missing values in the listing.

## Reading Data in Fixed Columns Using Formatted Input

Formatted input is the most common (and flexible) method for reading data in fixed columns. Let’s jump right to the program and then the explanation. Here it is.

### Program 8.5: Reading Data in Fixed Columns Using Formatted Input

```
data Health;
  infile '~/MyBookFiles/health.txt' pad;
```

```

input @1 Subj    $3.
      @4 Gender   $1.
      @5 Age      2.
      @7 HR       2.
      @9 SBP      3.
     @12 DBP      3.
     @15 Chol     3. ;
run;

title "Listing of Data Set Health";
proc print data=Health;
  ID Subj;
run;

```

The @ sign is called a *column pointer*. The number following the @ sign is the starting column for the value that you want to read.

Following the variable name is an informat. There are lots of informats for reading and interpreting things like dates, times, and values with dollar signs and commas. This program only uses two: The informat \$w. (w stands for width) reads w columns of character data; the informat w. reads w columns of numeric data.

The informat for numeric data is actually more general. You can specify a numeric informat that includes information about where to insert a decimal point in the value. For example, you might have numbers such as 10795 to represent 107 dollars and 95 cents. You could write an informat in the form of w.d, where w is the total number of columns to read and d indicates that there is a decimal place with d digits to the right. In the case of reading the dollars and cents value of 10795, you could use an informat of 5.2. The numeric value stored in the data set would be 107.95. When you use a w.d format, and the value that you are reading contains a decimal point, the d in the informat is ignored. The data set created by this program is identical to the data set produced by Program 8.4.

There are some shortcuts that you can use when employing formatted input. One of the most useful shortcuts uses variable lists and informat lists. If you have a group of variables that all share the same informat, you can list all the variables together (in a set of parentheses) and follow the list of variables by one or more informats (also in a set of parentheses). Here is Program 8.5, rewritten using

this feature.

### Program 8.6: Demonstrating Formatted Input

```
data Health;
  infile '~/MyBookFiles/health.txt' pad;
  input @1 Subj  $3.
        @4 Gender $1.
        @5 (Age HR) (2.)
        @9 (SBP DBP Chol) (3.);
run;

title "Listing of Data Set Health";
proc print data=Health;
  ID Subj;
run;
```

In this program, Age and HR both use the 2. informat; SBP, DBP, and Chol all use the 3. informat.

If you have as many informats as you have variables names the variables and informats will “pair up” on a one-to-one basis. If there are fewer informats than variable names, SAS will go back to the beginning of the informat list and reuse the informats in order. If there is only one informat, it will apply to every variable in the variable list. This is by far the most common way variable lists and informat lists are used.

In this example, this shortcut only saved a few lines of typing. However, imagine that you had 50 character values (called Ques1-Ques50), each occupying one column. You could write a very compact INPUT statement like this:

```
input (Ques1-Ques50) ($1.);
```

The variable list also demonstrates a convenient way to reference all the variables from Ques1 to Ques50. Anywhere in a SAS program where you need to name variables that have the same alphabetic root (Ques in this example), you can use a single hyphen to indicate that you are referencing all the variables from the first to the last.

## Conclusion

This chapter explored how to read external text data in almost any format. This is only the tip of the iceberg. For more information about the INPUT statement, please check out *Learning SAS by Example: A Programmer's Guide, Second Edition* (Cody, 2018), published by SAS Press.

# Problems

You can download all the files and programs you need for these problems from the author's website: [support.sas.com/cody](http://support.sas.com/cody). Programs and data for the problems are in a folder called Problems.

1. A quick survey was conducted, and the following data values were collected:

Variable	Description
----------	-------------

**Subj** Subject number (3 digits – stored as character)

Gender F=Female, M=Male

---

Height	Height in inches
Weight	Weight in pounds
Income_Group	Income group: L=Low, M=Medium, H=High

---

The data values were saved in a blank-delimited file called Quick.txt. Copy this file to your SAS Studio folder (for example, Problems) and write a SAS program to read this data file, create a temporary SAS data set, and produce a listing of the file. The file appears as follows:

#### File 8.4: Quick.txt

```
001 M 10/21/1950 68 150 H
002 F 9/11/1981 63 101 M
003 F 1/1/1983 62 120 L
004 M 5/17/2000 57 98 L
005 M 7/15/1970 79 220 H
006 F 6/1/1968 71 188 M
```

Use the mmddyy10. informat to read the DOB. Also, include the following statement in your program:

```
format DOB mmddyy10. ;
```

The reason for this is that, as you will see later in the chapter on dates, SAS stores dates as the number of days from January 1, 1960. The statement above is an instruction to print the DOB as a date and not the internal value (the number of days from 1/1/1960).

2. Use PROC FREQ and the data set in problem 1, to compute frequencies for the variables Gender and Income\_Group.
3. Add options to your program in Problem 2 to have the table list values in decreasing order of frequency and omit the cumulative statistics from the tables.
4. Modify the program in Problem 1 to compute a new variable called BMI (body mass index). BMI is defined as the weight in kilograms divided by the height (in meters) squared. Conversions are:  
1 kg. = 2.2 pounds  
1 meter = 39.3701 inches
5. The same data described in Problem 1 was saved as a CSV file called Quick.csv. Write a program to create a SAS data set from this CSV file. Be sure to include the FORMAT statement mentioned in that problem.
6. The same data described in Problem 1 in Chapter 8 was entered with the forward slash (/) as a delimiter. Because the dates also include slashes, the dates were placed in quotation marks. Here is a listing of the file:

#### **File 8.5: Quick\_Slash.txt**

```
001/M/"10/21/1950"/68/150/H
002/F//9/11/1981"/63/101/M
003/F/"1/1/1983"/62/120/L
004/M/"5/17/2000"/57/98/L
005/M/"7/15/1970"/79/220/H
006/F/"6/1/1968"/71/188/M
```

The variables in the file are ID (allow for 11 characters), First and Last name (each up to 15 characters), Gender (M or F), Age, and State code (2 characters). Write a program to read data from this file, create a temporary SAS data set (call it Slash), and produce a listing of the file.

Hints: 1) Some lines are missing values at the end of the line, and  
2) Two slashes in a row indicate there is a missing value (think  
DSD).

7. An Excel workbook called Grades.xlsx contains data on student grades. Use the IMPORT task on the Utilities menu to read the spreadsheet and create a SAS data set.
8. The data from the quick survey was entered into a file called Quick\_Cols.txt using fixed columns as follows:

Variable	Description	Columns

---

Subj	Subject number (3 digits)	1-3
------	---------------------------	-----

---

Gender	F=Female, M=Male	4
--------	------------------	---

---

DOB	Date of birth in <i>mm/dd/yyyy</i> form	5-14
-----	---	------

---

Height	Height in inches	15-16
<hr/>		

Weight	Weight in pounds	17-19
<hr/>		

Income_Group	Income group: L=Low, M=Medium, H=High	20
<hr/>		

Using column input, create a SAS data set from this file.

Important note: Because DOB is a date, you will have to read it as a character string.

9. Create a SAS data set from the file Quick\_Cols.txt using formatted input. Read the DOB with the mmddyy10. informat. Also, include the FORMAT statement mentioned in Problem 1.
10. What's wrong with this program?

```
1. data Names;  
2.   input Name $ Height Weight;  
3.   Height_CM = Height*2.54;  
4.   *Note 1 inch = 2.54 cm  
5.   datalines;  
    Zemlachenko 73 190  
    Holland 63 100  
    ;
```

# Chapter 9: Reading and Writing SAS Data Sets

## What's a SAS Data Set?

Besides reading raw data from text files, SAS can also read and write data from SAS data sets. Once you have created a SAS data set from raw data or you have been given a SAS data set, you are ready to run SAS procedures to analyze your data or to write a DATA step to create new variables or to further manipulate your data.

A SAS data set actually contains two parts: One is called the data descriptor, and the other is the data itself. The data descriptor contains your variable names, whether a variable is stored as character or numeric (the only two types allowed in a SAS data set), how many bytes of storage are used to store a variable, and other information about how to display the variable in reports or charts. A fancy word used to describe the data descriptor portion of a SAS data set is *metadata*—data about your data.

The format of a SAS data set is proprietary to SAS and only SAS can read and write SAS data sets directly (there are other programs that can read SAS data sets and convert them to other formats). If you try to open a SAS data set in Word or Notepad, it will look like gibberish.

If you want to examine the data descriptor (metadata) for a SAS data set, you have several options. Because this is the programming portion of the book, the method described here is to run a SAS procedure called PROC CONTENTS. Suppose you want to see the data descriptor for the data set called Demo described in Chapter 7.

The program shown next uses PROC CONTENTS to do this (the program to create the data set Demo is included so that you don't have to turn back to Chapter 7 to see it.)

### Program 9.1: Running PROC CONTENTS to Examine the Data

## Descriptor for Data Set Demo

```
data Demo;
  infile("~/MyBookFiles/demographic.txt");
  input ID $ Gender $ Age Height Weight Party $;
run;

*Program to display the data descriptor of data set Demo;

title "Data Descriptor for Data Set Demo";
proc contents data=Demo;
run;
```

You use a DATA= procedure option to specify which data set you want to examine. Here is the output from this procedure:

**Figure 9.1: Output from Program 9.1**

Data Descriptor for Data Set Demo			
Data Set Name	WORK.DEMO	Observations	5
Member Type	DATA	Variables	6
Engine	V9	Indexes	0
Created	10/25/2020 11:02:34	Observation Length	48
Last Modified	10/25/2020 11:02:34	Deleted Observations	0
Protection		Compressed	NO
Data Set Type		Sorted	NO
Label			
Data Representation	SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64		
Encoding	utf-8 Unicode (UTF-8)		

Engine/Host Dependent Information	
Data Set Page Size	131072
Number of Data Set Pages	1
First Data Page	1
Max Obs per Page	2722
Obs in First Data Page	5
Number of Data Set Repairs	0
Filename	/saswork/SAS_work63EA0001A890_odaws04-usw2.oda.sas.com/SAS_work8CDF0001A890_odaws04-usw2.oda.sas.com/demo.sas7bdat
Release Created	9.0401M6
Host Created	Linux
Inode Number	1610650715
Access Permission	rw-r--r--
Owner Name	ronaldcody
File Size	256KB
File Size (bytes)	262144

Alphabetic List of Variables and Attributes			
#	Variable	Type	Len
3	Age	Num	8
2	Gender	Char	8
4	Height	Num	8
1	ID	Char	8
6	Party	Char	8
5	Weight	Num	8

The first section of output shows global information about your data set: the number of observations, the number of variables, etc. You also see the date and time the data set was created or modified.

The next section of your output contains technical information about the size of the data set, the page size, the access permissions, and other details that you might (or might not) be interested in.

This last section of the output lists each of your variable names (in alphabetical order), the file type (Num for numeric, Char for character), and the number of bytes of storage used.

## Temporary Versus Permanent SAS Data Sets

All the data sets described in the previous chapter were temporary data sets. That is, once you end your SAS session, these data sets disappear. To create or read a permanent SAS data set, you need to specify a two-part name: the location where the data set is stored (referred to as a *library* in SAS terminology) and the actual data set name. You can specify a permanent SAS data set name by specifying the library name, a period, and then the data set name. Names of permanent SAS data sets are in the form:

*library-reference*.*data-set-name*

where *library-reference* (referred to as a *libref* in SAS terminology) is an alias for a folder and *data-set-name* is the name of the SAS data set. You might wonder about the data sets you created in the last chapter. Those data set names did not contain a period. It turns out that when you use a single name in a DATA statement (such as DATA

`Test ;`), SAS assumes that the data set is a temporary data set and stores it in the Work library (a temporary folder/directory located in your user area). The true name of the data set is `Work.Test`, where `work` is the built-in library reference to the Work library.

You can have SAS Studio create the library for you, using the **Libraries** tab or write a single line of SAS code. You will see both methods in this chapter. First, let's look at the programming method.

## Creating a Library by Submitting a LIBNAME Statement

There are two steps in creating or reading a permanent SAS data set. Step one is to create the libref by using a LIBNAME statement. For example, if you want to create a permanent SAS data set in the MyBookFile folder, you first create a library reference (the name Oscar is used in this example) like this:

```
libname Oscar "~/MyBookFiles";
```

Library references are a maximum of 8 characters and follow SAS naming conventions. If you want to create a permanent SAS data set called Demo in the Oscar library, your DATA statement looks like this:

```
data Oscar.Demo;
```

Data set Demo will remain even after you close your SAS session. If you look in the folder, MyBookFiles, you will see a file called `Demo.sas7bdat` (the extension refers to a SAS data set that is a binary data file compatible with SAS 7).

Here is a program to create the Demo data set in a permanent library called Oscar:

### Program 9.2: Using a LIBNAME Statement to Create a Permanent SAS Library

```
libname Oscar "~/MyBookFiles";

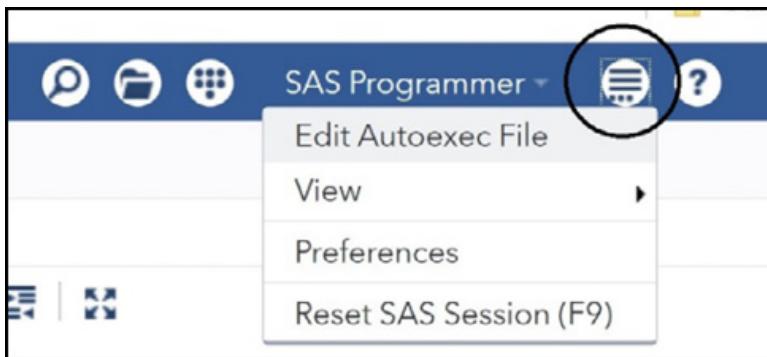
data Oscar.Demo;
  infile "~/MyBookFiles/demographic.txt";
  input ID $ Gender $ Age Height Weight Party $;
run;
```

Although Program 9.2 places the SAS data set Demo in the MyBookFiles folder, you will have to submit the LIBNAME statement each time you start a SAS Studio session. There are two approaches that make this process automatic. One is to add the LIBNAME statement in a SAS program called Autoexec.sas.

The reserved program name Autoexec.sas executes automatically every time you open a SAS session. Therefore, if you add your LIBNAME statement to this file it will re-create the Oscar library automatically.

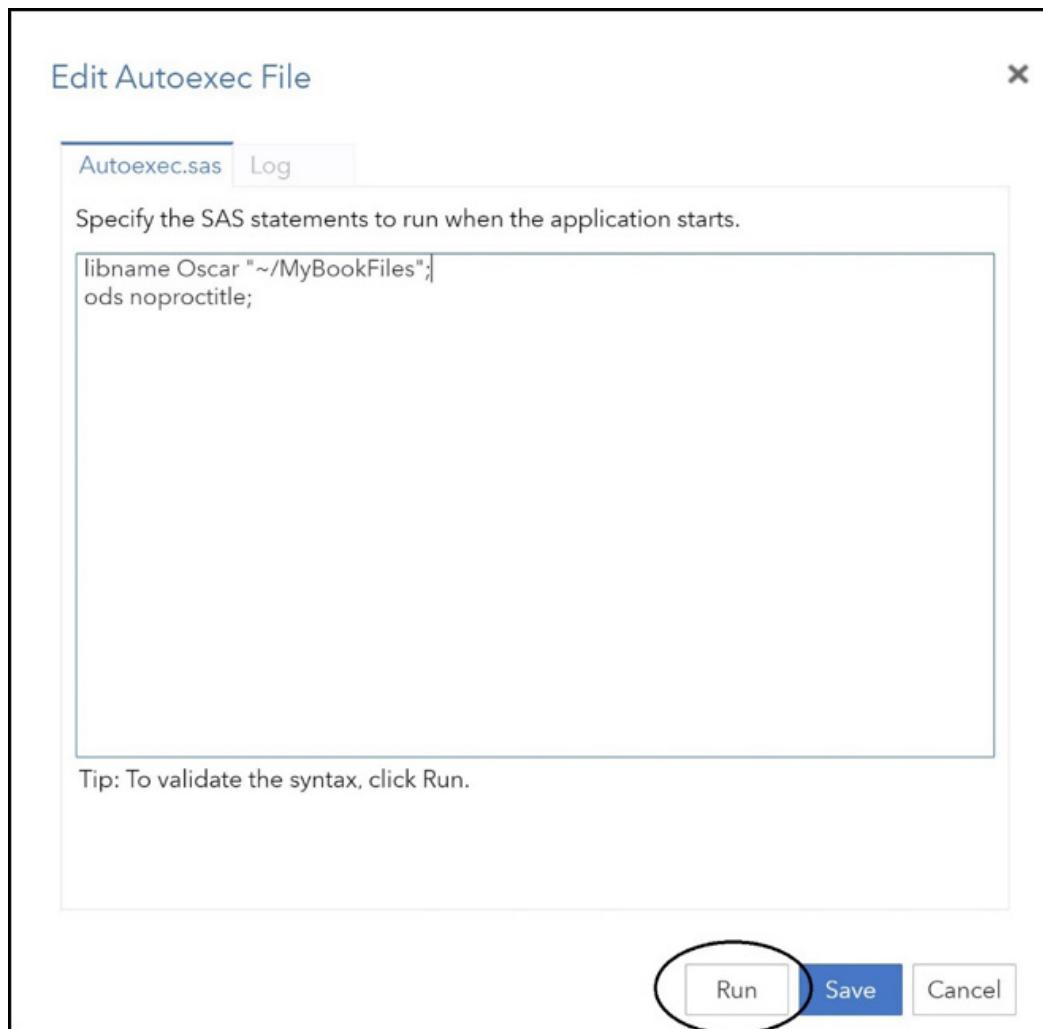
One way to edit the Autoexec.sas program is to have SAS Studio assist you. Click on the icon circled in Figure 9.2.

**Figure 9.2: Editing Autoexec.sas**



When you select **Edit Autoexec File**, you see the following.

**Figure 9.3: Editing Autoexec.sas**



The LIBNAME statement was entered. You will notice that this author previously added the statement:

```
ods noproctitle;
```

This statement removes the name of the procedure from all SAS output. It is a good idea to click the box labeled **RUN** before you save the file.

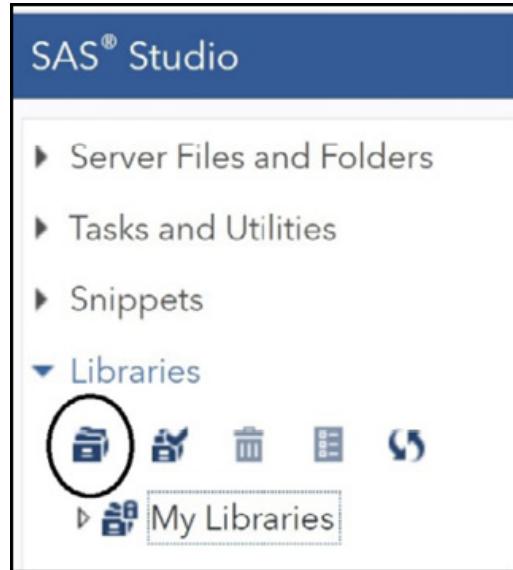
## Using the Library Tab to Create a Permanent Library

You can save yourself all this extra work if you simply let SAS Studio create a library for you. If you use this method, SAS Studio will ask you if you would like your library to be added the Autoexec.sas

automatically. To demonstrate how to use SAS Studio to create a library, let's create another library named Felix.

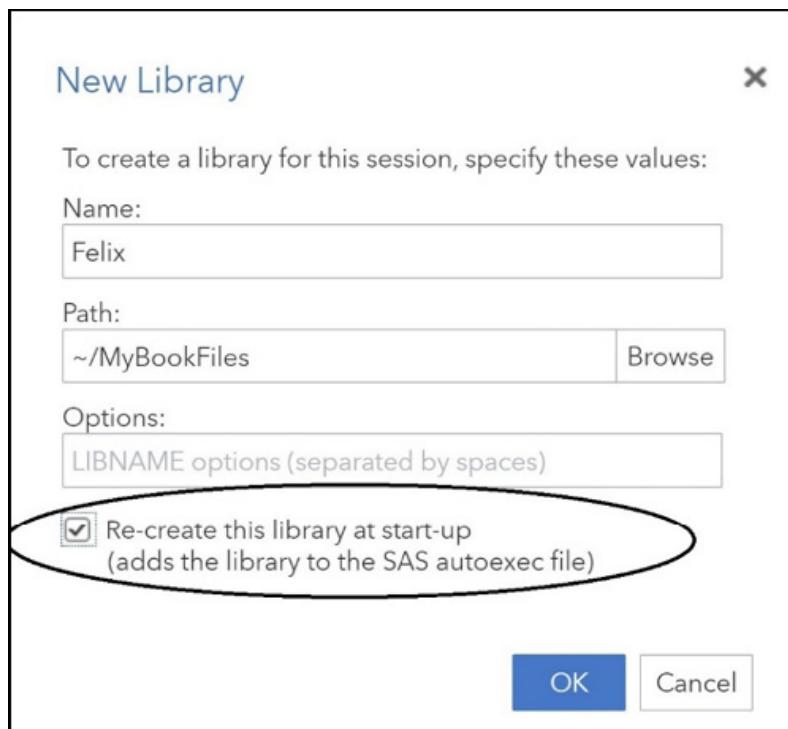
You start by selecting the **Library** tab in the Navigation Pane (Figure 9.4).

**Figure 9.4: Creating a Library Using the Library Tab**



Click the icon circled in this figure to create a new library. Here's what it looks like.

**Figure 9.5: New Library**



Be sure to check the box circled in Figure 9.5, and then click **OK**. When you do this, the library called Felix will be created every time you start a SAS Studio session.

## Reading from a Permanent SAS Data Set

You use an INPUT statement to read data from a raw data file, and you use a SET statement to read data from a SAS data set. The reason you do not need an INPUT statement is that the SAS data set contains all the information (called metadata) about the variables. In the example that follows, you want to use the SAS data set Demo as input to a DATA step that creates a new variable called Wt\_Kg, representing the person's weight in kilograms. In addition, let's assume that you are writing this program in a new SAS session. Here is the program:

### Program 9.3: Using a SAS Data Set as Input to a Program

```
data New_Demo;
  set Oscar.Demo;
  Wt_Kg = Weight / 2.2;
run;
title "Listing of Data Set New_Demo";
```

```
proc print data=New_Demo;  
run;
```

In this example, you plan to create a new, temporary data set called New\_Demo. Because the SAS data set Demo is stored in the Oscar library, you use the two-level name

Oscar.Demo

In the SET statement to read observations from the Oscar.Demo data set. You use an assignment statement (Wt\_Kg = Weight / 2.2;) to create the new variable Wt\_Kg.

Here is the output.

**Figure 9.6: Output from Program 9.2**

Listing of Data Set New_Demo							
Obs	ID	Gender	Age	Height	Weight	Party	Wt_Kg
1	012345	F	45	65	155	I	70.455
2	131313	M	28	70	220	R	100.000
3	987654	F	35	68	180	R	81.818
4	555555	M	64	72	165	D	75.000
5	172727	F	29	62	102	I	46.364

If you want the new data set New\_Demo to be permanent, replace the DATA statement with:

```
data Oscar.New_Demo;
```

## Conclusion

In this chapter, you saw how to create permanent SAS data sets as well as how to use existing SAS data sets to create either temporary or permanent SAS data sets. The previous version of this chapter, where SAS University Edition was described rather than OnDemand for Academics, was much longer and more complicated than what you see here. In this author's judgment, ODA is greatly improved over the University Edition.

## Problems

1. Use PROC CONTENTS to display the data descriptor for the data set Heart in the SASHELP library. Run it again with the VARNUM procedure option (remember, you place procedure options between the procedure name and the semicolon).
2. Create a new, temporary SAS data set called Heart\_Vars from the data set Heart in the SASHELP library. Include the variables BP\_Status, Chol\_Status, Systolic, Diastolic, and Status. Use a KEEP= data set option on the Heart data set to do this. Hint: First, click the **My Libraries** tab and then expand the list of the data sets in SASHELP.
3. **(For all the problems in this book where you need to create a permanent SAS data set, a folder called Problems was created in SAS Studio (the same way the MyBookFiles folder was created) and used in the solutions. You can use this name for your solutions or create a folder with whatever name you want to use).** Repeat Problem 2 except make the data set Heart\_Vars a permanent data set in your Problems folder (or whatever name you chose) folder.
4. Create a temporary SAS data set called Alive from the data set Heart in the SASHELP library. This data set should contain the variables BP\_Status, Chol\_Status, Systolic, and Diastolic. Use a WHERE= data set option to select only those observations where Status is equal to 'Alive'. Use the data set option (OBS=10) with PROC PRINT to list the Alive data set, like this:  

```
proc print data=Alive(obs=10);
```

Hint: You will need to include the variable Status in the KEEP= data set option because you need this variable to use in your WHERE= data set option. Use a DROP statement to drop Status so that it is not included in the Heart\_Vars data set.

5. Write the statements necessary to create a permanent SAS data set called Young\_Males in the Problems (or whatever name you chose) folder. Use as input the SASHELP data set Class and select only those observations where Gender is equal to 'M' and Age is 11 or 12.
6. What's wrong with this program?
  1. data New;
  2. set SASHELP.Fish(keep Species Weight);

```
3.      Wt_Kg = Weight/2.2;  
4.      *Note: 1 Kg = 2.2 Lbs *;  
5. run;
```

# Chapter 10: Creating Formats and Labels

## What Is a SAS Format and Why Is It Useful?

It is a common practice to store information in a database using codes rather than actual values. For example, you might have a questionnaire where the responses are strongly disagree, disagree, neutral, agree, and strongly agree. It would be unusual to store the actual values in your database. Rather, you would use coded values such as 1=strongly disagree, 2=disagree, and so on.

Even though you are storing codes in your database, you would like to see the actual labels printed in your output. SAS formats are the tool that allows this to happen.

To demonstrate how to create your own SAS formats, let's start with a SAS data set called Taxes, described in Table 10.1.

**Table 10.1: Formats in Taxes Data Set**

Variable	Description	Codes Used
SSN	Social Security Number	

---

SSN      Social Security Number

Gender

Gender

M=Male, F=Female

---

Question\_1

Do you pay taxes?

1=Yes, 0=No

---

Question\_2

Are you satisfied with  
this service?

1=strongly disagree,  
2=disagree, 3=neutral,  
4=agree, 5=strongly agree

---

Question\_3

How many phone calls  
did it take to resolve  
your problem?

A=0, B=1 or 2, C=3 to 5,  
D=More than 5

---

Question\_4

Was the person  
answering your call  
friendly?

Same as Question\_2

---

Question_5	How much did you pay?	The actual dollar amount
------------	-----------------------	--------------------------

---

The program to create the Taxes data set is shown next:

### Program 10.1: Creating the Taxes Data Set (and Demonstrating a DATALINES Statement)

```
data Taxes;
  informat SSN $11.
    Gender $1.
    Question_1 - Question_4 $1.;

  input SSN Gender Question_1 - Question_5;
datalines;
101-23-1928 M 1 3 C 4 23000
919-67-7800 F 9 2 D 2 17000
202-22-3848 M 0 5 A 5 57000
344-87-8737 M 1 1 B 2 34123
444-38-2837 F . 4 A 1 17233
763-01-0123 F 0 4 A 4 .
;

title "Listing of Data Set Taxes";
proc print data=Taxes;
  id SSN;
run;
```

This program uses a feature that you have not seen before: a DATALINES statement. When you want to write a short SAS program to test your logic or syntax, you can save the trouble of writing a text file (perhaps using Notepad or some similar editor) and then writing an INFILE statement telling the program where to find the data. Instead, you can actually include the lines of data in the program itself. You do this by writing a DATALINES statement and following this statement with your data. You end your data with a single semicolon or a RUN statement. SAS will read these lines of data as if they were in an external file. There is one additional new feature in this program—that is, the ID statement following PROC PRINT. You use an ID statement to do two things: First, you specify a variable that

you want to display in the first column of your listing. If your data set has more variables that can fit across one page of output, the ID variable is repeated on each new page. Second, when you specify an ID variable, the default Obs column is no longer printed (this is a good thing in this author's opinion).

Here is a listing of the Taxes data set.

**Figure 10.1: Output from Program 10.1**

Listing of Data Set Taxes							
SSN	Gender	Question_1	Question_2	Question_3	Question_4	Question_5	
101-23-1928	M	1	3	C	4	23000	
919-67-7800	F	9	2	D	2	17000	
202-22-3848	M	0	5	A	5	57000	
344-87-8737	M	1	1	B	2	34123	
444-38-2837	F		4	A	1	17233	
763-01-0123	F	0	4	A	4	.	

You would prefer to have values for Gender be listed as 'Male' and 'Female', values for Question\_1 to be listed as 'Yes' and 'No', values for Question\_2 and Question\_4 to display the agreement scale (called a *Likert scale* by psychometricians), and values for Question\_3 to show the number of calls. Finally, you would like to place the dollar amounts into four categories: 0 - \$10,000 = Low, \$10,001 - \$20,000 = Medium, \$20,001 - \$50,000 = High, and \$50,000+ = Very High.

The process of substituting these labels for the coded values is called *formatting* in SAS terminology. The first step is to create the formats. Once that is accomplished, you can associate one or more variables with these formats. You use PROC FORMAT to create your SAS formats, as shown in the following program.

### **Program 10.2: Creating Your Own Formats**

```
proc format;
  value $Gender 'M'='Male'
            'F'='Female';
  value $Yesno '0'='No'
             '1'='Yes'
             other='Did not answer';
  value $Likert '1'='Strongly Disagree'
```

```

      '2'='Disagree'
      '3'='No Opinion'
      '4'='Agree'
      '5'='Strongly Agree';
value $Calls 'A'='None'
      'B'='1 or 2'
      'C'='3 - 5'
      'D'='More than 5';

value Pay_group  low-10000 = 'Low'
      10001-20000 = 'Medium'
      20001-50000 = 'High'
      50001-high = 'Very High';

run;

```

You write VALUE statements to define your formats and follow the keyword VALUE with the name of the format that you want to create. Format names can be up to 32 characters long, and they follow the same naming conventions as other SAS names, except that they cannot end in a digit.

Formats that you plan to apply to character variables must begin with a dollar sign (\$) (leaving 31 characters for you to use). In this example, the first four formats will be associated with one or more character variables. The last format will be associated with the numeric variable Question\_5. Because the first four formats will be associated with character variables, their names all begin with a dollar sign. Following the format name, you list either a single value (such as 'M') or a range of values (such as low–10000). Notice that you place the character values in single or double quotation marks. There are several keywords that you can use in defining a value. As you can see in the \$Yesno format, the keyword OTHER will supply the label “Did not answer” for all values other than ‘0’ or ‘1’. You can also use the keywords LOW and HIGH to refer to the lowest and highest values, respectively. Note that for the Pay\_group format, you cannot include commas in the numerical ranges.

To demonstrate how formats work, let’s first run PROC FREQ to compute frequencies for each variable in the Taxes data set, without

formatting any of the variables. The PROC FREQ statements look like this:

### Program 10.3: Computing Frequencies on the Taxes Data Set (without Formats)

```
title "Frequencies for the Taxes Data Set";
proc freq data=Taxes;
    tables Gender Question_1 - Question_5 / nocum;
run;
```

The TABLES option NOCUM is an instruction to omit cumulative frequencies from the results, shown below.

**Figure 10.2: Output from Program 10.3**

Frequencies for the Taxes Data Set		
Gender	Frequency	Percent
F	3	50.00
M	3	50.00
Question_1	Frequency	Percent
0	2	40.00
1	2	40.00
9	1	20.00
Frequency Missing = 1		
Question_2	Frequency	Percent
1	1	16.67
2	1	16.67
3	1	16.67
4	2	33.33
5	1	16.67
Question_3	Frequency	Percent
A	3	50.00
B	1	16.67
C	1	16.67
D	1	16.67

Question_4	Frequency	Percent
1	1	16.67
2	2	33.33
4	2	33.33
5	1	16.67

Question_5	Frequency	Percent
17000	1	20.00
17233	1	20.00
23000	1	20.00
34123	1	20.00
57000	1	20.00

**Frequency Missing = 1**

PROC FREQ listed frequencies for each unique value of the variables. For numeric variables such as Question\_5 (How much did you pay, in dollars?), you see frequencies for each unique value—not something that is very useful. There are other SAS procedures (such as PROC SGLOT) that will automatically place numerical values into groups and produce histograms, etc.

It's time to see what the output from PROC FREQ looks like when you add a FORMAT statement, associating each of the variables in the Taxes data set with a format. Here is the code.

#### Program 10.4: Adding a FORMAT Statement to PROC FREQ

```
title "Frequencies for the Taxes Data Set";
proc freq data=Taxes;
  format Gender $Gender.
        Question_1 $Yesno.
        Question_2 Question_4 $Likert.
        Question_3 $Calls.
        Question_5 Pay_group.;
  tables Gender Question_1 - Question_5 / nocum;
run;
```

You use a FORMAT statement to associate your variables with the appropriate format. SAS programs recognize the difference between variables and formats because when you associate formats to a variable, you end each format name with a period. The two variables

Question\_2 and Question\_4 share the same format, so you list these two variables together and follow them with the format that you want to use (\$Likert). Output from this program is shown below.

**Figure 10.3: Output from Program 10.4**

Frequencies for the Taxes Data Set		
Gender	Frequency	Percent
Female	3	50.00
Male	3	50.00
Question_1	Frequency	Percent
No	2	50.00
Yes	2	50.00
Frequency Missing = 2		
Question_2	Frequency	Percent
Strongly Disagree	1	16.67
Disagree	1	16.67
No Opinion	1	16.67
Agree	2	33.33
Strongly Agree	1	16.67
Question_3	Frequency	Percent
None	3	50.00
1 or 2	1	16.67
3 - 5	1	16.67
More than 5	1	16.67

Question_4	Frequency	Percent
Strongly Disagree	1	16.67
Disagree	2	33.33
Agree	2	33.33
Strongly Agree	1	16.67

Question_5	Frequency	Percent
Medium	2	40.00
High	2	40.00
Very High	1	20.00
<b>Frequency Missing = 1</b>		

You now see the formatted values in the tables. You might wonder about the order of the values in the tables. Notice that Female comes before Male and No comes before Yes. Did you guess that PROC FREQ orders these values alphabetically? If so, you are correct. Later on, you will see how to change the order in the frequency tables.

If you are very observant, you will have noticed that the table for Question\_1 lists two missing values. However, when you look at the unformatted data, you see only one missing value and one value of 9. It turns out that the OTHER category was combined with the missing values (because it was not a valid value for this variable). If you want to separate the missing value from the OTHER category, you can add a format for a missing value like this:

```
value $Yesno '0'='No'
      '1'='Yes'
      ' '='Did not answer'
      other='Invalid value';
```

Because \$Yesno is a character format, you specify a missing value with a single space between the quotation marks. To specify a missing value for a numeric format, use a single period instead. When you use this format with PROC FREQ, the output frequencies for Question\_1 look like Figure 10.4.

**Figure 10.4: Using the Modified \$Yesno Format**

Question_1	Frequency	Percent
No	2	40.00
Yes	2	40.00
Invalid value	1	20.00
Frequency Missing = 1		

By using this modified format, you now see that there was one invalid value and one missing value.

## Using SAS Built-in Formats

SAS provides you with a large number of built-in formats that you can use, along with ones that you create yourself. One example of a numeric format is  $w.d$ , where  $w$  (stands for width) specifies the total number of spaces to use when writing a number and  $d$  specifies how many places to include after the decimal place. Table 10.2 shows some examples. For each of these examples,  $X=1234.567$ .

**Table 10.2: Built-in Numeric Format Examples**

Format	Display	Explanation
8.3	1234.567	The 8 is the total column width, including the decimal point.

---

8.3	1234.567	The 8 is the total column width, including the decimal point.
10.4	1234.5670	There is a leading space and a 0 is added to give four places to the right

of the decimal point.

---

8.2

1234.57

There is one leading blank and the decimal value is rounded.

---

4.

1235

When the width is shorter than the value, it gets truncated (and rounded).

---

10.1

1234.6

There are four leading blanks and the decimal is rounded.

---

Some other useful SAS formats are dollar*w.d* and commaw.*d*. The dollar format adds dollar signs, commas (if needed), and cents (if *d* is equal to 2). The value of *w* is the field width. This includes the dollar sign, any commas, a decimal point, and the digits to the right of the decimal point. If you leave out a value of *d* (the number after the period), the dollar value will be rounded. The comma format is similar to the dollar format, except that it does not include a dollar sign and the value of *d* can show as many decimal places as needed. Formats for SAS dates are particularly useful and will be discussed in the

chapter on SAS dates.

## More Examples to Demonstrate How to Write Formats

There is great flexibility in defining values or ranges when you create a SAS format. The following examples help illustrate this flexibility.

In Program 10.2, one of the formats placed numerical values into four categories as shown here:

```
value Pay_group  low-10000    = 'Low'  
           10001-20000 = 'Medium'  
           20001-50000 = 'High'  
           50001-high   = 'Very High';  
run;
```

This works fine as long as the values to be formatted are integers. However, suppose you tried to format a value of 10,000.50. This value falls between the Low and Medium ranges (and would not be formatted—it would print as 10000.50 in a listing). You might specify ranges like this:

```
value Pay_group  low-10000    = 'Low'  
           10000-20000 = 'Medium'  
           20000-50000 = 'High'  
           50000-high   = 'Very High';
```

Although this code works, it is confusing. What is the formatted value of 10,000? It turns out it would be formatted as ‘Low’. It is better to allow a value to only match a single range. You can exclude a value from the beginning or the end of a range by adding a less than (<) sign before or after the hyphen that specifies ranges. For example, to exclude 10,000 from the low range, you would use:

```
value Pay_group  low- <10000 = 'Low'
```

If you wanted to exclude 10,000 from the medium range, you would use:

```
value Pay_group  10000< -20000 = 'Medium'
```

You can also use hyphens and commas to specify character ranges. Here are some examples:

```
value $Grades 'A' , 'B' = 'Good'  
    'C' - 'E' = 'Passing'  
    'F'       = 'Fail'  
    other     = 'Error';
```

Values of A or B are formatted as ‘Good’; C, D, and E are formatted as ‘Passing’; F is formatted as ‘Fail’; and any characters not equal to any of these values is formatted as ‘Error’.

## Describing the Difference between a FORMAT Statement in a Procedure and a FORMAT Statement in a DATA Step

In the programs that you have seen thus far, the FORMAT statements have been used inside a SAS procedure. This creates an association between the variables and formats *only for that procedure*. If you place a FORMAT statement in a DATA step, the association between the variables and formats remains for the entire program. Then, if you write a PROC PRINT or PROC FREQ step, you will see formatted values for all the variables you listed in your FORMAT statement.

You will usually find it more convenient to associate your formats with variables, using a FORMAT statement in a DATA step, saving the trouble of having to rewrite (or copy) it for each procedure that you run. Keep in mind that even though you have assigned a format to a variable, many procedures such as PROC MEANS or other statistical procedures will still use the internal values of the variables when doing their calculations.

## Making Your Formats Permanent

All of the user-written formats previously described are temporary formats—that is, they only exist for the duration of your SAS session. If you have formats that you plan to use frequently, you can make them permanent. That way, every time you open up a SAS session, you can use any of your permanent formats without having to rerun PROC FORMAT. There are a few steps that you need to follow to make this happen.

First, you need to decide where you plan to store your formats. To keep this example simple, let's store your permanent formats in MyBookFiles. You run PROC FORMAT just as you did previously, except you add the procedure option LIBRARY=*libref* (where *libref* is the library reference you create using a LIBNAME statement). Here is an example.

### Program 10.5: Making a Permanent Format

```
libname Myfmnts "~/MyBookFiles";

proc format library=Myfmnts;
  value $Gender 1='Male'
        2='Female';
  value $Yesno  0='No'
        1='Yes';
run;
```

Figure 10.5: Log Window After Running Program 10.5

```
71      libname Myfmnts "~/MyBookFiles";
NOTE: Libref MYFMTS refers to the same physical library as FELIX.
NOTE: Libref MYFMTS was successfully assigned as follows:
      Engine:      V9
      Physical Name: /home/ronaldcody/MyBookFiles
72      proc format library=Myfmnts;
73        value $Gender 1='Male'
74          2='Female';
NOTE: Format $GENDER has been written to MYFMTS.FORMATS.
75        value $Yesno  0='No'
76          1='Yes';
NOTE: Format $YESNO has been written to MYFMTS.FORMATS.
77      run;
```

The two formats, \$Gender and \$Yesno, are now permanent formats and are stored in the Myfmnts library.

If you have variables associated with formats in a permanent SAS data set, keep in mind that the formats are now a permanent property of those variables, and the format definitions must be available in order to open and read the data set. That is, if you give someone a copy of one of your permanent SAS data sets, be sure to give them the format catalog (it will have the name formats.sas7bcat).

There is a way to open a data set that has formats associated with variables where you do not have the format library. The “trick” used by this author and many others is to include the following statement at the top of your program:

```
options nofmterr;
```

The option nofmterr stands for no format error.

A good way to keep track of your permanent formats is to include another PROC FORMAT option called FMTLIB. This option produces a listing of all the formats in the specified library, along with all of the defined values for these formats. To illustrate this, let's run [Program 10.5: Making a Permanent Format](#) again with the FMTLIB option added.

## Program 10.6: Adding the FMTLIB Option

```
libname Myfmts "~/MyBookFiles";  
  
proc format library=Myfmts fmtlib;  
    value $Gender 1='Male'  
                  2='Female';  
    value $Yesno  0='No'  
                  1='Yes';  
run;
```

Here is the output.

**Figure 10.6: Output from Program 10.6**

FORMAT NAME: \$GENDER LENGTH: 6 NUMBER OF VALUES: 2 MIN LENGTH: 1 MAX LENGTH: 40 DEFAULT LENGTH: 6 FUZZ: 0		
START	END	LABEL (VER. V7 V8 26OCT2020:10:00:46)
1	1	Male
2	2	Female

FORMAT NAME: \$YESNO LENGTH: 3 NUMBER OF VALUES: 2 MIN LENGTH: 1 MAX LENGTH: 40 DEFAULT LENGTH: 3 FUZZ: 0		
START	END	LABEL (VER. V7 V8 26OCT2020:10:00:46)
0	0	No
1	1	Yes

This output is an especially useful document for you to use or to share with others who want to use your formats.

Before you use these permanent formats in a new SAS session, you need to include two statements: One is a LIBNAME statement that defines your *libref*. The other statement is OPTIONS FMTSEARCH=*libref*. The option FMTSEARCH=*libref*, tells SAS to look in the location specified by your *libref* to find your permanent formats. So, beginning every SAS session, you would add these two statements:

```
libname Myfmts "~/MyBookFiles";
options fmtsearch=Myfmts;
```

If you plan to use permanent formats, you should put these two statements in the Autoexec file, as described in the last chapter.

## Creating Variable Labels

You can associate labels with your SAS variables. If you have variables such as Gender and Race, a label might not be necessary. However, for variable names such as Question\_1, Question\_2, etc., you might want to provide labels. You use a LABEL statement to associate labels with your variables. For example, the following program adds labels to the variables in the Taxes data set (described in the beginning of this chapter).

### Program 10.7: Creating Variable Labels

```
data Taxes;
  informat SSN $11.
    Gender $1.
    Question_1 - Question_4 $1.;

  input SSN Gender Question_1 - Question_5;
  label Question_1 = 'Do you pay taxes?'
    Question_2 = 'Are you satisfied with the service?'
    Question_3 = 'How many phone calls?'
    Question_4 = 'Was the person friendly?'
    Question_5 = 'How much did you pay?';
  datalines;
  101-23-1928 M 1 3 C 4 23000
```

```

919-67-7800 F 9 2 D 2 17000
202-22-3848 M 0 5 A 5 57000
344-87-8737 M 1 1 B 2 34123
444-38-2837 F . 4 A 1 17233
763-01-0123 F 0 4 A 4 .
;

title 'Frequencies for the Taxes Data Set';
proc freq data=Taxes;
    format Gender $Gender.
        Question_1 $Yesno.
        Question_2 Question_4 $Likert.
        Question_3 $Calls.
        Question_5 Pay_group. ;
    tables Gender Question_1 - Question_5 / nocum;
run;

```

Now that you have added variable labels to your program, let's see how it affects the output. Listed below is a partial output from this program.

**Figure 10.7: Output from Program 10.7**

Frequencies for the Taxes Data Set		
Gender	Frequency	Percent
Female	3	50.00
Male	3	50.00

Do you pay taxes?		
Question_1	Frequency	Percent
No	2	40.00
Yes	2	40.00
Invalid value	1	20.00
Frequency Missing = 1		

You now see the variable label listed at the top of each frequency table.

## Conclusion

Adding formats and labels to a SAS program will make the listings and tables much more readable. If you have formats that you use

frequently, be sure to create permanent formats so that you don't have to run PROC FORMAT every time you start a SAS session. If you have associated formats with variables in the DATA step, and plan to share your data set with others, be sure to include the format catalog (formats.sas7bcat) along with the SAS data set.

## Problems

In this chapter and several to follow, you will see programs with captions such as Program for Problem Sets 1, and so forth. These programs are all included in the download package in a folder called Problems. All these programs are stored in the file Problem Data Sets.sas. Other files in the Problems folder contain data that you can use in other problems.

1. Modify the following program to supply formats to the variables listed in the table below. (Use the SAS format MMDDYY8. to format Visit.)

Variable	Type	Formatted Values
Gender	Num	1=Male, 2=Female
Q1-Q4	Char	1='Strongly Disagree', 2='Disagree', 3='No Opinion', 4='Agree', 5='Strongly Agree'

---

Gender      Num      1=Male, 2=Female

Q1-Q4      Char      1='Strongly Disagree',  
2='Disagree', 3='No Opinion',  
4='Agree', 5='Strongly Agree'

---

Visit	Num	Month/Day/Year
Age	Num	0-20='Young', 21-40='Still young', 40-60='Middle', 61+='Older'

---

## Program for Problem Sets 1

```

data Questionnaire;
  informat Gender 1. Q1-Q4 $1. Visit date9.;
  input Gender Q1-Q4 Visit Age;
  format Visit date9.;
datalines;
1 3 4 1 2 29May2015 16
1 5 5 4 3 01Sep2015 25
2 2 2 1 3 04Jul2014 45
2 3 3 3 4 07Feb2015 65
;
title "Listing of Data Set Questionnaire";
proc print data=Questionnaire noobs;
run;

```

2. Using the program in Problem 1, create a format that places ages into the following categories:  
0-20='Group 1', 21-40='Group 2', 41-60='Group 3', 61-80='Group 5', 81+='Group 5'  
Use this format for the variable Age and the other formats described in Problem 1. Produce a listing showing the formatted values for all the variables.
3. You have a character variable called Grades. Values of Grades

are 'A', 'B', 'C', 'D', 'F', 'I', and missing. Write a format (call it \$Grades) that formats 'A' and 'B' as 'Good', 'C' as 'Average', 'D' as 'Poor', 'F' as 'Fail', 'I' as 'Incomplete', and missing values as 'Missing'. Also, include a format for any nonmissing value that is not one of the valid values (call them Invalid).

# Chapter 11: Performing Conditional Processing

## Introduction

All programming languages enable you to perform *conditional processing*—making decisions based on data values or other conditions. For example, you might want to create a new variable (Age\_Group) based on the values of age. Another common use of conditional logic is to check if data values are within a prescribed range.

## Grouping Age Using Conditional Processing

For the first example, you have data on gender, age, height, and weight. You want to create a new variable (Age\_Group) based on the variable Age. Here is a first attempt that runs but **has a logical flaw in regard to SAS missing values**.

### Program 11.1: First Attempt at Creating an Age Group Variable (Incorrect Program)

```
data People;
  input @1  ID      $3.
        @4  Gender  $1.
        @5  Age     3.
        @8  Height  2.
        @10 Weight  3.;

  if Age le 20 then Age_Group = 1;
  else if Age le 40 then Age_Group = 2;
  else if Age le 60 then Age_Group = 3;
  else if Age le 80 then Age_Group = 4;
  else if Age ge 80 then Age_Group = 5;

datalines;
001M 5465220
002F10161 98
003M 1770201
```

```

004M 2569166
005F   64187
006F 3567135
;
title "Listing of Data Set People";
proc print data=People;
  id ID;
run;

```

To indicate conditions such as less than, etc., you have a choice of two-letter abbreviations or symbols. The table below shows all the possible logical comparisons.

**Table 11.1: Logical Comparison Operators**

Logical Comparison	Mnemonic	Symbol
Equal to	EQ	=
Not equal to	NE	$\neq$ or $\sim$ or $\neg$
Less than	LT	<

Less than or equal to LE <=

Greater than GT >

Greater than or equal to GE >=

Equal to any value in a list IN

The new statements in this program are the IF and ELSE IF statements. They work like this: Following the IF or ELSE IF statement is a logical statement that is either true or false. If the statement is true, the following expression executes; if it is false, the following expression does not execute. Also, if the logical statement on an IF or ELSE IF statement is true, all the subsequent ELSE IF statements are skipped. For example, in the data set People, the first subject is 54 years old. The first ELSE IF statement that is true is

```
else if age le 60 then Age_Group = 3;
```

Because this statement is true, all the remaining ELSE IF statements

are skipped. This logic has the advantage of being more efficient than a series of IF statements; the program does not have to evaluate more IF statements than necessary.

Let's run the program and examine the output.

**Figure 11.1: Output from Program 11.1**

Listing of Data Set People					
ID	Gender	Age	Height	Weight	Age_Group
001	M	54	65	220	3
002	F	101	61	98	5
003	M	17	70	201	1
004	M	25	69	166	2
005	F	.	64	187	1
006	F	35	67	135	2

Most of the Age\_Group values are correct. However, there is a problem for ID 005. This person had a missing value for age but was placed in age group 1. Why?

In SAS, a numeric missing value is treated logically as the most negative number possible. Thus, a missing value is less than any real—positive or negative—number.

The first IF statement asks whether Age is less than or equal to 20. Person 005 has a missing value for Age and a missing value is less than 20, so this person is placed in age group 1. Here is one way to fix Program 11.1:

**Program 11.2: Corrected Version of Program 11.1**

```
data People;
  input @1 ID      $3.
        @4 Gender  $1.
        @5 Age     3.
        @8 Height  2.
        @10 Weight 3.;

  if missing(Age) then Age_Group = .;
  else if Age le 20 then Age_Group = 1;
  else if Age le 40 then Age_Group = 2;
  else if Age le 60 then Age_Group = 3;
  else if Age le 80 then Age_Group = 4;
```

```
else if Age ge 80 then Age_Group = 5;

datalines;
001M 5465220
002F10161 98
003M 1770201
004M 2569166
005F    64187
006F 3567135
;

title "Listing of Data Set People";
proc print data=People;
  id ID;
run;
```

The first IF statement tests if Age is a missing value. This is accomplished using the MISSING function. All SAS functions end with a set of parentheses. The values placed in the parentheses are called *arguments* to the function. The MISSING function returns a value of true if the argument is a missing value and false otherwise (The MISSING function works for both character and numeric arguments). When the program processes ID 005, the missing function returns a true value and the variable Age\_Group is set to a missing value (designated by a period). An alternative to testing for a missing value is the following line of code:

```
if Age = . then Age_Group = . ;
```

This author strongly recommends that you use the MISSING function to test for missing values.

Most SAS programmers would agree that failing to account for missing values in a DATA step is the most common logical error in SAS programming. Always be sure to consider the consequences of a missing value meeting your program logic.

Output from Program 11.2 results in a missing value for Age\_Group for ID 005.

## Using Conditional Logic to Check for Data Errors

You can use IF-THEN-ELSE logic to test if values of certain variables

are outside a predetermined range. For example, you might want to check if anyone in the People data set was heavier than 200 pounds or lighter than 100 pounds. The following program does just that.

### Program 11.3: Using Conditional Logic to Test for Out-of-Range Data Values

```
data _null_;
  set People;
  if Weight lt 100 and not missing(Weight) or Weight gt 200 then
    put "Weight for ID " ID "is " Weight;
run;
```

There are several new features in this program. First is the special data set name `_NULL_`. This is a reserved data set name that enables you to run a DATA step without actually creating a data set. The reason it is used in this program is that you are checking for out-of-range values, and you do not need a data set when you are finished checking—thus the use of DATA `_NULL_`. Because you are not creating a data set, the program is more efficient than one that does create a data set.

The SET statement brings in observations from the People data set. Notice that the condition in the IF statement checks for two things: First, is the value less than 100 and not missing (remember that a missing value is less than 100)? Second, is the Weight greater than 200? If either condition is true, the PUT statement executes. The PUT statement is an instruction to write out the text “Weight for ID” followed by the value of ID (note ID is not in quotation marks, so it represents a variable name) followed by the word “is” followed by the value of Weight. By default, the PUT statement writes its output to the SAS log. This is fine for programmers but not so fine for nonprogrammers. To tell SAS to write out the Weight values to the RESULTS window, add the line `file print;` before the PUT statement. Here is the output from Program 11.3.

**Figure 11.2: Output from Program 11.3**

```

71      data _null_;
72      set People;
73      if Weight lt 100 and not missing(Weight) or Weight gt 200 then
74          put "Weight for ID " ID " is " Weight;
75      run;

Weight for ID 001 is 220
Weight for ID 002 is 98
Weight for ID 003 is 201
NOTE: There were 6 observations read from the data set WORK.PEOPLE.

```

Three people had weights out-of-range.

## Describing the IN Operator

If you want to check if a value is any one of several values, you can use multiple OR operators or the IN operator. Suppose you want to check if values for Race are 'W', 'B', 'H', or 'O' (white, black, Hispanic, or other). Using the OR operator, you could write:

```

length Race_Value $ 7;
if Race = 'W' or Race = 'B' or Race = 'H' or Race = 'O' then
    Race_Value = 'Valid';
else Race_Value = 'Invalid';

length Race_Value $ 7;
if Race in ('W','B','H','O') then Race_Value = 'Valid';
else Race_Value = 'Invalid';

```

When you use the IN operator, you place the character values in quotation marks (single or double) and separate each value by a comma or space. The statement evaluates as true if Race matches any one of the listed values. It is useful to know that once a match is made, the IN operator stops looking for matches. When extreme efficiency is needed, programmers will place values most likely to be present in the data at the beginning of the list of values, saving the extra CPU time in searching the whole list.

You can use the IN operator with numeric data as well. For example, if you want to list all subjects in Age\_Group 3, 4, or 5, you could use the following statement:

```
if Age_Group in (3,4,5) then put "Older Folks";
```

You might wonder about the LENGTH statement in these code

segments. Why is it needed? Remember that the length of character variables is set at compile time (before any data values are read or any conditional logic is performed). Without the LENGTH statement, the first time the variable Race\_Value appears is where it is set equal to 'Valid'. Because 'Valid' is 5 characters long, SAS would set the storage length for Race\_Value to 5. The effect would be to truncate the other possible value for Race\_Group, 'Invalid', to 5 characters. By using a LENGTH statement before the assignment statement for Race\_Value, SAS assigns a storage length of 7 for this variable. A popular trick to avoid entering the LENGTH statement is to pad the first value of Race\_Value, 'Valid', with two extra blanks (for example, 'Valid '), so that SAS will assign the variable a length of 7 instead of 5. Using a LENGTH statement is considered more elegant and a better programming practice.

## Using Boolean Logic (AND, OR, and NOT Operators)

You can combine the three Boolean operators, AND, OR, and NOT, in logical expressions. Here is an example:

```
*Note: there are no missing values of LDL and HDL in the data;
if (LDL gt 100 or HDL lt 50) and Gender eq 'F' then
  Risk = 'High';
  else Risk = 'Low';
if (LDL gt 100 or HDL lt 40) and Gender = 'M' then
  Risk = 'High';
  else Risk = 'Low';
```

High values of LDL (low-density lipids) or low values of HDL (high-density lipids) are considered a risk for coronary artery disease. In addition, the Mayo Clinic uses different values for the HDL cutoff for men and women.

The order of precedence of the Boolean operators in decreasing order is:

1. NOT
2. AND
3. OR

You can use parentheses to force a different order of operation. For example:

```
if x and y or z;
```

is equivalent to

```
if (x and y) or z;
```

If you want to perform the OR operation before the AND operation, write the expression like this:

```
if x and (y or z);
```

Because the NOT operator has the highest precedence, the expression:

```
if x and not y or z;
```

is equivalent to

```
if x and (not y) or z;
```

Even though the Boolean operators have a built-in ordering, feel free to add parentheses in your logical expressions—it makes the logic easier to understand.

## A Special Caution When Using Multiple OR Operators

It is very easy to make a serious error when using multiple OR operators. Take a look at the following program:

### Program 11.4: A Common Error Using Multiple OR Operators

```
data Mystery;
  input x;
  if x = 3 or 4 then Match = 'Yes';
  else Match = 'No';
datalines;
3
4
9
.
-5
;
title "Listing of Data Set Mystery";
```

```
proc print data=Mystery noobs;  
run;
```

Before you look at the output, notice that the PROC PRINT option NOOBS was added. This option removes the Obs column from the output of PROC PRINT. If you use an ID statement with PROC PRINT, the NOOBS option is not necessary because the ID variable(s) replace the Obs column.

Here is the output.

**Figure 11.3: Output from Program 11.4**[Program 11.4: A Common Error Using Multiple OR Operators](#)

Listing of Data Set Mystery	
x	Match
3	Yes
4	Yes
9	Yes
.	Yes
-5	Yes

That is not what you expected, is it? First of all, you probably expected a syntax error because the statement:

```
if x = 3 or 4 then Match = 'Yes';
```

should have been written as:

```
if x = 3 or x = 4 then Match = 'Yes';
```

However, if you expressed the logic in this program verbally, you might very well say if x equals 3 or 4, then Match equals 'Yes'. Let's first look at the SAS log.

**Figure 11.4: SAS Log from Program 11.4**

▼ Errors, Warnings, Notes

► ✘ Errors

► ⚠ Warnings

► ⓘ Notes (4)

```
1      OPTIONS NONOTES NOSTIMER NOSOURCE NOSYNTAXCHECK;
70
71      data Mystery;
72          input x;
73          if x = 3 or 4 then Match = 'Yes';
74          else Match = 'No';
75      datalines;
```

NOTE: The data set WORK.MYSTERY has 5 observations and 2 variables.

There are no syntax errors, yet the output is obviously wrong. What is going on?

In SAS, any numerical value that is not equal to 0 or a missing value is considered true.

In this program, there are two conditions separated by an OR operator. One is ‘x = 3’; the other is ‘4’. ‘4’ is not equal to 0 or missing, so it is evaluated as true. An OR expression is true if either one (or both) of the expressions is true. Because ‘4’ is true, the logical expression is true regardless of the value of x.

Just to be sure this is clear, let’s rewrite Program 11.4 correctly like this.

### Program 11.5: A Corrected Version of Program 11.4

```
data Mystery;
    input x;
    if x = 3 or x = 4 then Match = 'Yes';
    else Match = 'No';
datalines;
3
4
9
.
-5
;
```

```
title "Listing of Data Set Mystery";  
  
proc print data=Mystery noobs;  
run;
```

Here is the output from the corrected program.

**Figure 11.5: Output from the Program 11.5**

### Listing of Data Set Mystery

x	Match
3	Yes
4	Yes
9	No
.	No
-5	No

Seeing how easy it is to make this error when using multiple OR operators, you should be convinced that it is better to write the logical test as:

```
if x in (3,4) then Match = 'Yes';
```

## Conclusion

Just about every program that you write will need to use conditional logic. Using IF-THEN-ELSE statements and the Boolean operators, NOT, AND, and OR, enables you to evaluate complex conditions. Finally, remember to consider missing values when you write your logical expressions.

## Problems

1. Using the SASHELP data set Fish, create a new, temporary SAS data set called Group\_Fish that contains the variables Species, Weight, and Height. Using IF-THEN-ELSE logic, create a new variable called Group\_Fish that places the fish into weight groups

as follows:

0 to 100=1, 101-200=2, 201-500=3, 501-1,000=4, and 1,001 and greater=5.

Use PROC PRINT to list the first 10 observations from data set Group\_Fish.

2. Run Program 11.2 in this chapter (listed below), adding statements to print one message if the variable Age is greater than 100 and another message if Age is missing. Include the Age value in the message for ages greater than 100.

```
data People;
  input @1  ID      $3.
        @4  Gender  $1.
        @5  Age     3.
        @8  Height  2.
        @10 Weight  3.;

if missing(Age) then Age_Group = .;
else if Age le 20 then Age_Group = 1;
else if Age le 40 then Age_Group = 2;
else if Age le 60 then Age_Group = 3;
else if Age le 80 then Age_Group = 4;
else if Age ge 80 then Age_Group = 5;

datalines;
001M 5465220
002F10161 98
003M 1770201
004M 2569166
005F   64187
006F 3567135
;
```

3. Create a new, temporary SAS data set called High\_BP containing subjects from the SASHELP.Heart data set who have systolic blood pressure greater than 250 or diastolic blood pressure greater than 180. The new data set should only have variables Diastolic, Systolic, and Status. Use PROC PRINT to list the contents of High\_BP.
4. If A is true, B is false, and C is false, what do each of these expressions evaluate as (true or false)?
  - a) A AND NOT B
  - b) NOT A OR NOT C
  - c) A AND NOT B AND C

d) A AND (NOT B OR NOT C)

5. What's wrong with this program?

```
1. data Weights;
2.   input Wt;
3.   if Wt lt 100 then Wt_Group = 1;
4.   if Wt lt 200 then Wt_Group = 2;
5.   if Wt lt 300 then Wt_Group = 3;
6. datalines;
50
150
250
;
```

6. Starting with the SASHELP data set Retail, write a program to create a new data set (Sales\_Status) with the following new variables:

If Sales is greater than or equal to 300, set Bonus equal to 'Yes' and Level to 'High'. Otherwise, if Sales is not a missing value, set Bonus to 'No' and level to 'Low'. Use PROC PRINT to list the observations in this data set.

# Chapter 12: Performing Iterative Processing: Looping

## Introduction

Many programs use a process called *looping*. You can execute a series of statements a fixed number of times or you can set a condition that, when satisfied, will cause the iteration to start or stop and the rest of the program to continue. SAS iterative statements all start with the keyword DO and are referred to as *DO loops*.

## Demonstrating a DO Group

Although this chapter deals mainly with DO loops, this is a good place to demonstrate a DO group. When you conduct a logical test, for example, if Gender = 'F', you might want to do several things. A DO group enables you to accomplish this.

In this example, you want to perform several actions when the value of LDL (low-density lipids, the “bad” cholesterol) is above 100.

```
if LDL gt 100 then do;
  Stoke_Risk = 'High';
  LDL_Group = 'Bad';
end;
```

Rather than repeating the IF statement twice, you create a DO group, starting with the keyword DO and ending with the keyword END. All the statements between DO and END execute when the value of LDL is greater than 100. It is standard programming practice to indent all the statements between the DO and END statements—it makes the program easier to read.

## Describing a DO Loop

Let's start off with an example, similar to the one in Chapter 7, that

converts temperatures in degrees Celsius to degrees Fahrenheit. Except this time, you want to produce a table of temperatures from 0 degrees Celsius to 100 degrees Celsius.

### Program 12.1: Creating a Table of Celsius and Fahrenheit Temperatures

```
data Convert_Temp;
  do Temp_C = 0 to 100;
    Temp_F = 1.8*Temp_C + 32;
    output;
  end;
run;

title "Listing of Data Set Convert_Temp";
proc print data=Convert_Temp noobs;
run;
```

The general form of a DO loop is:

```
do variable = lower-limit to upper-limit by increment;
   SAS statements;
end;
```

This is how it works. The SAS variable following the keyword DO is first given a value defined by the lower limit. Next, the statements between the DO and END execute. At the bottom of the loop (at the END statement), the value of the looping variable is incremented by the value of the increment. If an increment is not specified (that is, the BY variable is left out), the increment is assumed to be 1. If the value of the variable is greater than the upper limit, the loop is finished and the statement(s) following the END statement execute. If not, process control returns to the top of the loop and the statements in the DO loop execute again.

Let's follow the sequence of events in Program 12.1:

1. The variable Temp\_C is given a value of 0.
2. The variable Temp\_F is computed.
3. The OUTPUT statement writes an observation to the data set Convert\_Temp.
4. At the END statement, the variable Temp\_C is incremented by

- 1.
5. The program now loops back to the DO statement and checks to see whether the value of Temp\_C is greater than 100 (it is not).
6. Temp\_C is set equal to 1 and a value is computed for Temp\_F.
7. Another observation is written to the data set.
8. The loop continues until an observation with Temp\_C equal to 100 is written out.
9. The variable Temp\_C is incremented by 1 (now equal to 101).
10. Because 101 is greater than 100 (the upper limit), the DO loop terminates.
11. The program ends with the RUN statement.

It should be noted that because an OUTPUT statement was used in this DATA step, the automatic output at the bottom of the DATA step does not occur. This prevents the last observation in the data set from being written twice.

A portion of the output is shown below.

**Figure 12.1: Output from Program 12.1**

Listing of Data Set Convert_Temp	
Temp_C	Temp_F
0	32.0
1	33.8
2	35.6
3	37.4
4	39.2
5	41.0
6	42.8
7	44.6
8	46.4
9	48.2
10	50.0

If you want to change the interval to 10 degrees Celsius, change the DO statement to read:

```
do Temp_C = 0 to 100 by 10;
```

The output would then look like this:

**Figure 12.2: Output from Program 12.1 (with the Interval Set to 10)**

Listing of Data Set Convert_Temp	
Temp_C	Temp_F
0	32
10	50
20	68
30	86
40	104
50	122
60	140
70	158
80	176
90	194
100	212

## Using a DO Loop to Graph an Equation

Suppose you want to create a graph of an equation and can't find the graphing calculator that you used in high school. For example, let's look at the following equation.

You can easily compute multiple values of  $y$  for a family of  $x$  values. Let's generate  $y$  values for  $x$  values ranging from -5 to 5 in intervals of .01:

### Program 12.2: Graphing a Cubic Equation

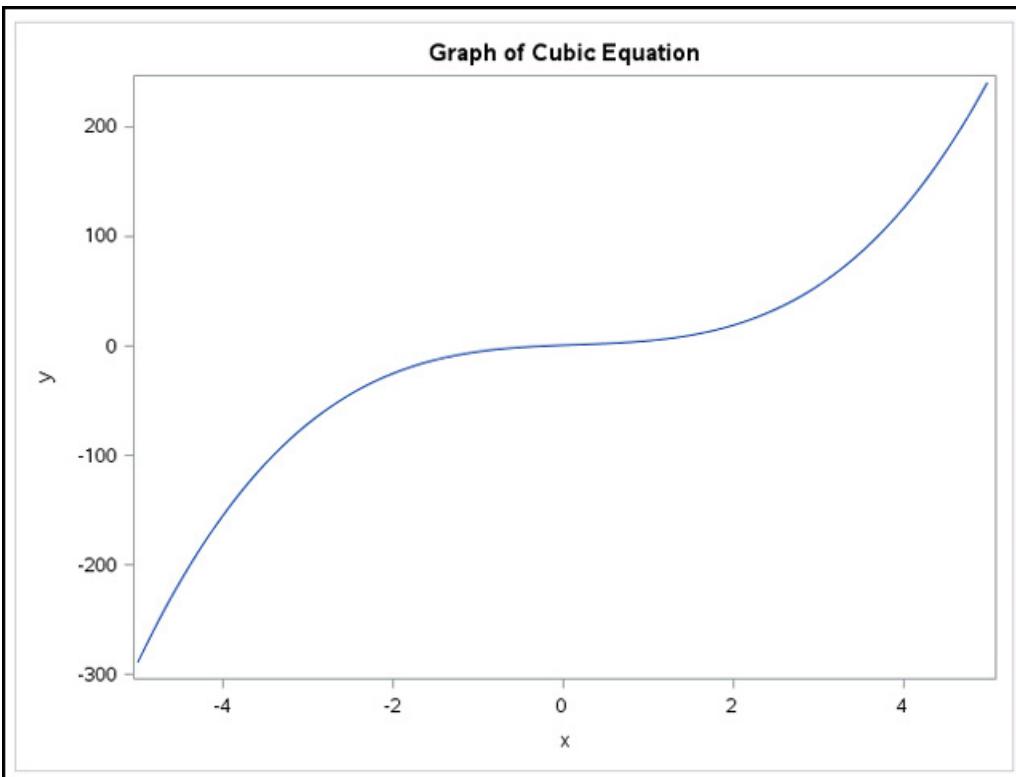
```
data Cubic;
  do x = -5 to 5 by .01;
    y = 2*x**3 - x**2 + 3*x;
    output;
  end;
run;

title "Graph of Cubic Equation";
proc sgplot data=Cubic;
  series x=x y=y;
```

```
run;
```

Data set Cubic contains the  $x$  and  $y$  values computed by the DO loop. At this point, you could use the built-in tasks to produce a series plot or to write the statements using PROC SGPlot to produce the plot. Here is the resulting plot.

**Figure 12.3: Output from Program 12.2**



And there it is!

## DO Loops with Character Values

Amazing as it might seem, you can write DO loops with character values. For example, suppose you have heart rates for subjects given one of three drugs: Placebo, Drug A, and Drug B. Each line of data contains three heart rates, one for each of the three drug groups. Here is what the data looks like.

### Data for Heart Rate Study

```
80 70 60  
82 77 63
```

```
76 74 70  
78 80 67
```

A program to read this data and assign values to a drug group variable is shown next.

### Program 12.3: Using a DO Loop with Character Values

```
data HR_Study;  
  do Drug_Group = 'Placebo', 'Drug A', 'Drug B';  
    input Heart_Rate @;  
    output;  
  end;  
datalines;  
80 70 60  
82 77 63  
76 74 70  
78 80 67  
;  
title "Listing of Data Set HR_Study";  
proc print data=HR_Study noobs;  
run;
```

Instead of a starting and ending value, you list each of the character values that you want to assign to Drug\_Group in quotation marks (single or double), separated by commas. The trailing @ sign in the INPUT statement prevents SAS from going to a new line each time the loop iterates, and the INPUT statement is executed. The trailing @ sign (“hold the line”) allows the program to read all three values from one line of data.

To help make this clear, here is a listing of data set HR\_Study.

**Figure 12.4: Listing of Data Set HR\_Study**

### Listing of Data Set HR\_Study

Drug_Group	Heart_Rate
Placebo	80
Drug A	70
Drug B	60
Placebo	82
Drug A	77
Drug B	63
Placebo	76
Drug A	74
Drug B	70
Placebo	78
Drug A	80
Drug B	67

## Leaving a Loop Based on Conditions (DO WHILE and DO UNTIL Statements)

There are two statements that perform a logical test to decide when to leave a DO loop—DO WHILE and DO UNTIL. Here's how they work.

### DO WHILE

The syntax for a DO WHILE loop is:

```
do while(expression);  
      SAS Statements  
end;
```

where *expression* is a SAS logical expression. This structure will loop from the DO statement to the END statement as long as the expression is true.

In a DO WHILE statement, the *expression* is evaluated at the top of the loop; therefore, if the expression remains false, the SAS statements in the DO WHILE loop will never execute.

Here are some examples.

### Program 12.4: Demonstrating a DO WHILE Loop

```
data Bank;
  Interest_Rate=.07;
  Amount = 1000;
  Goal = 2000;
  do while (Amount < Goal);
    Year + 1;
    Amount = Amount + Interest_Rate*Amount;
    output;
  end;
  format Amount Goal dollar9.2;
run;

title "Listing of Data Set Bank";
proc print data=Bank noobs;
run;
```

You start with \$1,000 in the bank. You want to see how many years it takes to reach (or exceed) your goal of \$2,000, with interest set at 7%. Because this bank compounds interest each year, you multiply the interest rate by the cumulative amount each year. The loop stops when Amount is greater than or equal to \$2,000. The statement Year + 1; is a SUM statement. It initializes Year at 0 and adds one each time the SUM statement executes. If this had been an assignment statement instead of a SUM statement, the variable Year would be set to a missing value at the top of the DATA step. Here is the listing.

**Figure 12.5: Output from Program 12.4**

### Listing of Data Set Bank

Interest_Rate	Amount	Goal	Year
0.07	\$1,070.00	\$2,000.00	1
0.07	\$1,144.90	\$2,000.00	2
0.07	\$1,225.04	\$2,000.00	3
0.07	\$1,310.80	\$2,000.00	4
0.07	\$1,402.55	\$2,000.00	5
0.07	\$1,500.73	\$2,000.00	6
0.07	\$1,605.78	\$2,000.00	7
0.07	\$1,718.19	\$2,000.00	8
0.07	\$1,838.46	\$2,000.00	9
0.07	\$1,967.15	\$2,000.00	10
0.07	\$2,104.85	\$2,000.00	11

When the value of Amount is greater than Goal, the DO WHILE loop ends. (It's been a long time since interest rates were at 7%, but this author didn't want the listing to be too long.)

### Combining an Iterative Loop with a WHILE Condition

**CAUTION:** If you are not careful, DO WHILE and (especially) DO UNTIL loops can cause infinite loops. This can occur if the condition for a DO WHILE loop is always true or the condition for a DO UNTIL loop is never achieved.

You can combine a traditional DO loop with a WHILE condition. You might want to do this only while testing your code, or you might have reason to use both conditions. Here is an example:

### Program 12.5: Combining an Iterative DO Loop with a WHILE Condition

```
data Bank;
  Interest_Rate=.07;
  Amount = 1000;
  Goal = 2000;
  do I = 1 to 20 while (Amount lt Goal);
    Year + 1;
    Amount = Amount + Interest_Rate*Amount;
    output;
```

```

end;
format Amount Goal dollar9.2;
drop I;
run;

title "Listing of Data Set Bank";
proc print data=Bank noobs;
run;

```

Even if the WHILE condition remains true, this loop will stop after 20 iterations.

## DO UNTIL

The companion to DO WHILE is DO UNTIL. There are some important differences. First of all, a DO UNTIL loop will exit the loop when the UNTIL condition becomes true. DO UNTIL loops are even more prone to infinite loops than DO WHILE loops. If the UNTIL expression is never true, the loop will continue indefinitely.

The UNTIL condition is evaluated at the bottom of the loop. Therefore, DO UNTIL loops will always execute at least one time.

The following program revisits the compound interest problem using a DO UNTIL. Here is the code.

### Program 12.6: Demonstrating a DO UNTIL Loop

```

data Bank;
Interest_Rate=.07;
Amount = 1000;
Goal = 2000;
do until (Amount gt Goal);
  Year + 1;
  Amount = Amount + Interest_Rate*Amount;
  output;
end;
format Amount Goal dollar9.2;
run;

title "Listing of Data Set Bank";

proc print data=Bank noobs;
run;

```

The output from this program is identical to the output from Program 12.4.

## Demonstrating That a DO UNTIL Loop Executes at Least Once

The short program below demonstrates that the statements in a DO UNTIL loop **always** execute at least once. Here is the program.

### Program 12.7: Demonstrating That a DO UNTIL Loop Will Always Execute Once

```
data At_Least.Once;
  x = 5;
  do until (x = 5);
    put "This line is inside the loop";
  end;
run;
```

When you run this program, the message “This line is inside the loop” prints once in the SAS log.

## Combining an Iterative Loop with an UNTIL Condition

As mentioned previously, DO UNTIL loops are even more likely to cause infinite looping (a bad thing, especially if you are writing out to a file in the loop) than a DO WHILE loop. Therefore, you might want to “take out some insurance” by combining an iterative DO loop with an UNTIL condition. You can choose the number of iterations to be much larger than you need, but it is just there in case the UNTIL condition is never true. Here is an example.

### Program 12.8: Combining an Iterative DO Loop with an UNTIL Condition

```
data Bank;
  Interest_Rate=.07;
  Amount = 1000;
  Goal = 2000;
  do I = 1 to 100 until (Amount gt Goal);
    Year + 1;
    Amount = Amount + Interest_Rate*Amount;
    output;
```

```

end;
format Amount Goal dollar9.2;
drop I;
run;

title "Listing of Data Set Bank";

proc print data=Bank noobs;
run;

```

The output from this program is identical to the output from Program 12.6.

## LEAVE and CONTINUE Statements

The two statements, LEAVE and CONTINUE, cause two different actions to take place within a loop. LEAVE, as the name suggests, jumps to the first line following the loop; CONTINUE skips the remaining statements in the loop, but it does return control back to the top of the loop. You can add a conditional LEAVE statement inside an ordinary DO loop, and it can mimic a DO WHILE or DO UNTIL loop. Here is an example:

### Program 12.9: Demonstrating a LEAVE Statement

```

data Bank;
Interest_Rate=.07;
Amount = 1000;
Goal = 2000;
do I = 1 to 20;
    Year + 1;
    Amount = Amount + Interest_Rate*Amount;
    output;
    if Amount gt Goal then leave;
end;
format Amount Goal dollar9.2;
drop I;
run;

title "Listing of Data Set BANK";
proc print data=Bank noobs;
run;

```

Although this example produces the same output as Program 12.4, it

lacks the esthetics of DO WHILE or DO UNTIL loops.

This next program demonstrates how a CONTINUE statement can be used. In this example, you want to compute compound interest rates, but you only want to see the dollar amounts once the amount is greater than \$1,500. Here is the code:

### Program 12.10: Demonstrating a CONTINUE Statement

```
data Bank;
  Interest_Rate=.07;
  Amount = 1000;
  Goal = 2000;
  do until (Amount gt Goal);
    Year + 1;
    Amount = Amount + Interest_Rate*Amount;
    if Amount lt 1500 then continue;
    output;
  end;
  format Amount Goal dollar9.2;
run;

title "Listing of Data Set Bank";
proc print data=bank noobs;
run;
```

As long as the Amount is less than \$1,500, the OUTPUT statement is skipped, but the loop continues. Once the Amount is greater than or equal to \$1,500, the OUTPUT statement executes.

**Figure 12.6: Output from Program 12.10**

Listing of Data Set Bank			
Interest_Rate	Amount	Goal	Year
0.07	\$1,500.73	\$2,000.00	6
0.07	\$1,605.78	\$2,000.00	7
0.07	\$1,718.19	\$2,000.00	8
0.07	\$1,838.46	\$2,000.00	9
0.07	\$1,967.15	\$2,000.00	10
0.07	\$2,104.85	\$2,000.00	11

Compound interest was a convenient task to demonstrate the various types of DO loops in SAS. If you actually need to compute

compound interest in a program, you should be aware that there is a SAS function called COMPOUND that does this for you. This function can be found in the group of financial functions in SAS help.

## Conclusion

You have seen how to create a DO group and a DO loop as well as DO WHILE and DO UNTIL loops. It is important to remember that DO WHILE loops evaluate the logical expression at the top of the loop and DO UNTIL loops evaluate the logical expression at the bottom of the loop (thus, a DO UNTIL loop always executes at least once). Also, be careful to avoid infinite loops. They can possibly cause your computer to crash or a disk drive to fill to capacity (a really bad thing).

## Problems

1. Write a DATA step to create a conversion table for pounds and kilograms. The table should have one column showing pounds from 0 to 100 in units of 10. The second column in the table should show the kilogram equivalents. Note: 1 Kg. = 2.2 Lbs.
2. Write a program to create a table with column 1 containing the integers from 1 to 10, column 2 the square root of column 1, and column 3 the square of column 1. Hint: To take a square root, you can use an exponent of .5 or a function called SQRT. To use the function, you would write something like: Root\_x = sqrt(x);
3. You have data from three groups of subjects (Groups A, B, and C). The actual data looks like this:

Group	Score
A	85
B	92
C	78

A 10

---

B 11

---

C 12

---

A 20

---

B 21

---

C 22

---

However, the data was entered without the groups, like this:

**Data for Group Study**

10  
11  
12  
20  
21  
22

Write a program to read the six numbers shown here (use DATALINES) to create a table with the variables Group and Score.

4. You have a variable called Money initialized at 100. Write a DO WHILE loop that compounds this amount by 3 percent each year and computes the amount of money plus interest for each year. Stop when the total amount exceeds 200.

To help get you started, the beginning of the program should look like this:

### Program for Problem Sets 2

```
data Interest;  
    Money = 100;  
    do while (put something here);  
        Year + 1; *keep track of years;  
        *compute new amount;  
        output; *output an observation for each iteration  
               of the loop;  
    end;  
run;
```

5. Solve Problem 4 using DO UNTIL instead of DO WHILE.
6. What is the value of Y when you run this program?

### Program for Problem Sets 3

```
data Until;  
    X = 5;  
    Y = 10;  
    do until (X eq 5);  
        Y = 20;  
    end;  
run;
```

# Chapter 13: Working with SAS Dates

## Introduction

SAS can read and write dates in almost any format, such as 5/23/2015 or 23May2015. No matter how the date appears in the input file, **SAS converts all dates to the number of days from January 1, 1960**. Thus, January 1, 1960 is 0; January 2, 1960 is a 1; and so on. Dates before January 1, 1960 are converted to negative numbers. For example, December 31, 1959 is equal to -1. Almost all computer languages store dates as the number of days from a fixed date. SAS chose January 1, 1960—other languages use other dates. Although SAS dates are stored internally as numbers, SAS can display the date in any of the standard formats. For those readers interested in history, SAS does not compute dates before January 1, 1582. That is because Pope Gregory the VIII decreed that October 4, 1582 would be followed by October 15, 1582 (the beginning of the Gregorian calendar).

Remember INFORMATS help read data into SAS and change how the data is stored in the SAS data set. Formats do not change how the data is stored and are used to change how the data is displayed.

## Reading Dates from Text Data

SAS uses informats to convert dates expressed in any one of the possible date formats into its internal representation of dates (the number of days from January 1, 1960). Here is an example.

You have a text file Date\_Data.txt containing three dates, as follows:

**File 13.1: Date\_Data.txt**

```
123456789012345678901234567890 (Ruler, not part of the file)
05/23/2015 23May2015 23/05/2015
10/21/1950 21Oct1950 21/10/1950
5/7/2013 7Jul2013 7/5/2013
```

The first date, starting in column 1, is in the form *mm/dd/yyyy* (month, day, year). The second date, starting in column 12, is in the form *ddMonyyyy* (day of the month, month abbreviation, year). The third date, starting in column 22, is in the form *dd/mm/yyyy* (day, month, year). The following program reads each of these dates:

### Program 13.1: Reading Dates in a Variety of Date Formats

```
data Read_Dates;
  infile "~/MyBookFiles/Date_Data.txt" pad;
  input @1 Date1 mmddyy10.
        @12 Date2 date9.
        @22 Date3 ddmmmyy10.;

run;

title "Listing of Data Set Dates";
proc print data=Read_Dates noobs;
run;
```

You use the appropriate informat to read each of the dates. The two informats mmddyy10. and ddmmmyy10. are pretty obvious. By the way, the 10 at the end specifies that the program is reading 10 columns of data. If you had a date such as 5/6/2013 anywhere in the specified 10 columns, SAS would still read it correctly. The informat DATE9. reads dates in this form: day of the month (one or two digits), month abbreviation (not case-sensitive), and four-digit year.

Here is the listing.

**Figure 13.1: Output from Program 13.1**

### **Listing of Data Set Dates**

Date1	Date2	Date3
20231	20231	20231
-3359	-3359	-3359
19485	19546	19485

This looks pretty strange. What you are seeing are the internal values for each of the dates—the number of days from January 1, 1960. To see these values as dates, you need to associate a SAS date format with each variable. Here is a list of some of the more popular date formats:

**Table 13.1: List of Popular Formats for Displaying Dates**

How Date is Displayed	Date Format
10/21/1950	mmddyy10.
10-21-1950	mmddyyd10.
10 21 1950	mmddyyb10.

---

10/21/1950                    mmddyy10.

---

10-21-1950                    mmddyyd10.

---

10 21 1950                    mmddyyb10.

---

10:21:1950 mmddyy10.

---

21/10/1950 ddmmyy10.

---

21Oct1950 DATE9.

---

Saturday, October 21, 1950 WEEKDATE.

---

The formats that end in 10d, 10b, and 10c (rows 3-5) replace the default slash with hyphens, blanks, or colons, respectively. You can add d's, b's, or c's to the *ddmmyy* formats as well.

You can use any date format you want to display any of these date variables. The program below is identical to Program 13.1 [Program 13.1: Reading Dates in a Variety of Date Formats](#), except that each of the dates is now formatted.

## Program 13.2: Adding Formats to Display the Date Values

```

data Read_Dates;
  infile("~/MyBookFiles/Date_Data.txt") pad;
  input @1 Date1 mmddyy10.
        @12 Date2 date9.
        @22 Date3 ddmmyy10. ;
  format Date1 mmddyy10. Date2 Date3 date9. ;
run;

title "Listing of Data Set Dates";
proc print data=Read_Dates noobs;
run;

```

The listing below shows the effect of these formats.

**Figure 13.2: Output from Program 13.2**

Listing of Data Set Dates		
Date1	Date2	Date3
05/23/2015	23MAY2015	23MAY2015
10/21/1950	21OCT1950	21OCT1950
05/07/2013	07JUL2013	07MAY2013

The date values are now displayed properly. Because the month-day-year format is popular in the United States and the day-month-year format is popular in most of the world outside the United States, many companies that do global business prefer the DATE9. format (used for Date2 and Date3 in the listing).

## Creating a SAS Date from Month, Day, and Year Values

If you have month, day, and year values as separate variables, you can use the MDY function to create a SAS date. This function takes three numeric arguments and, as the function name suggests, the three arguments provide values for month, day, and year.

Here is an example: You have a SAS data set with variables Month, Day, and Year. To compute a SAS date, use the following statement:

Date = MDY(Month, Day, Year);

The variable Date will be a SAS date (i.e., the number of days from January 1, 1960). You will most likely want to add a FORMAT statement to associate one of the date formats with the variable Date.

## Describing a Date Constant

If you need to refer to a specific date in a DATA step, you could always enter that date as the number of days from January 1, 1960. However, computing that value is inconvenient and a program that referred a date by its internal value would also be difficult to read. The answer: The date constant (also known as a date literal).

Suppose you want to test a date to see whether it is earlier than January 1, 2020, or later than December 31, 2021. Here's how to do it:

### Program 13.3: Demonstrating a Date Constant

```
data _null_;
  title "Checking for Out of Range Dates";
  input @1 Date mmddyy10.;
  file print;
  if Date lt '01Jan2020'd and not missing(Date) or
    Date gt '31Dec2021'd then put "Date " Date "is out of
range";
  format Date mmddyy10.;
  datalines;
  10/13/2020
  5/1/2012
  1/1/2015
  6/5/2020
  1/1/2000
;
```

As demonstrated in this program, a date constant consists of a date in the form *ddMonyyyy*, where *dd* is the day of the month, *Mon* is a three-character month abbreviation, and *yyyy* is the year. You place this value in single or double quotation marks and follow it (no spaces) with an uppercase or lowercase d. At compile time, SAS converts the date constants to SAS dates.

Before we get to the output, let's discuss a few features of the program. First of all, it uses the reserved data set name `_NULL_`. Data `_NULL_` was discussed earlier in this book, but in case you didn't read the chapters of this book in order, here is the explanation again.

Giving a data set this special name has the effect of not creating any data set at all. That is, after the DATA step runs, there is no data set left behind. Using a `_NULL_` statement saves the computer all the overhead of creating the data set and writing observations to it. Because of this, you can't use a PROC PRINT (or any other procedure) to list the contents of the data set, **because there is no data set**. The solution is to use a PUT statement. Following the keyword PUT, you can enter text (in quotation marks) and variable names. SAS will print out the quoted text and the value of the variables listed in the PUT statement. Finally, the statement FILE PRINT is an instruction to print the results to the RESULTS window. If you leave off this statement, SAS, by default, sends the results of the PUT statement to the SAS log.

As a reminder, a missing value (logically the most negative value that you can have) is going to be less than any date. By using the MISSING function in this program, you are giving instructions not to print the out-of-range message if the date is a missing value.

Here is the output from Program 13.3.

**Figure 13.3: Output from Program 13.3**

**Checking for Out of Range Dates**

Date 05/01/2012 is out of range
Date 01/01/2015 is out of range
Date 01/01/2000 is out of range

The program worked as advertised.

## Extracting the Day of the Week, Day of the

## Month, Month, and Year from a SAS Date

You can use the four functions, WEEKDAY, DAY, MONTH and YEAR, to compute the day of the week (a number from 1 to 7, with 1=Sunday), day of the month (a number from 1 to 31), month of the year (a number from 1 to 12), and year from a SAS date. Let's look at an example.

You have a SAS data set with a variable called Date. You would like to generate bar charts showing frequencies for day of the week, day of the month, and year.

To save room, this program only produces a bar chart for the day of the week. You can substitute the day of the month or year to obtain bar charts for these variables. The DATA step that creates the data set containing the Date variable is included so that you can try running the program yourself.

### Program 13.4: Extracting the Day of the Week, Day of the Month, and Year from a SAS Date

```
data Extract;
  informat Date mmddyy10.;
  input Date @@; 
  Day_of_Week = weekday(Date); 
  Day_of_Month = day(Date);
  Year = year(Date);
  format Date mmddyyd10.;
datalines;
1/5/2000 2/8/2000 4/23/2000 4/12/2000 8/21/2000 8/21/2000
8/22/2000
12/12/2000 12/15/2000 12/18/2000
2/22/2001 2/1/2001 4/18/2001 4/18/2001 4/18/2001 9/17/2001
12/25/2001
12/22/2001 3/3/2001 3/6/2001 3/7/2001
;
title "Listing of the First Eight Observations from Extract";
proc print data=Extract (obs=8); 
run;

title "Frequencies for Day of the Week";
proc sgplot data=Extract; 
  vbar Day_of_Week;
```

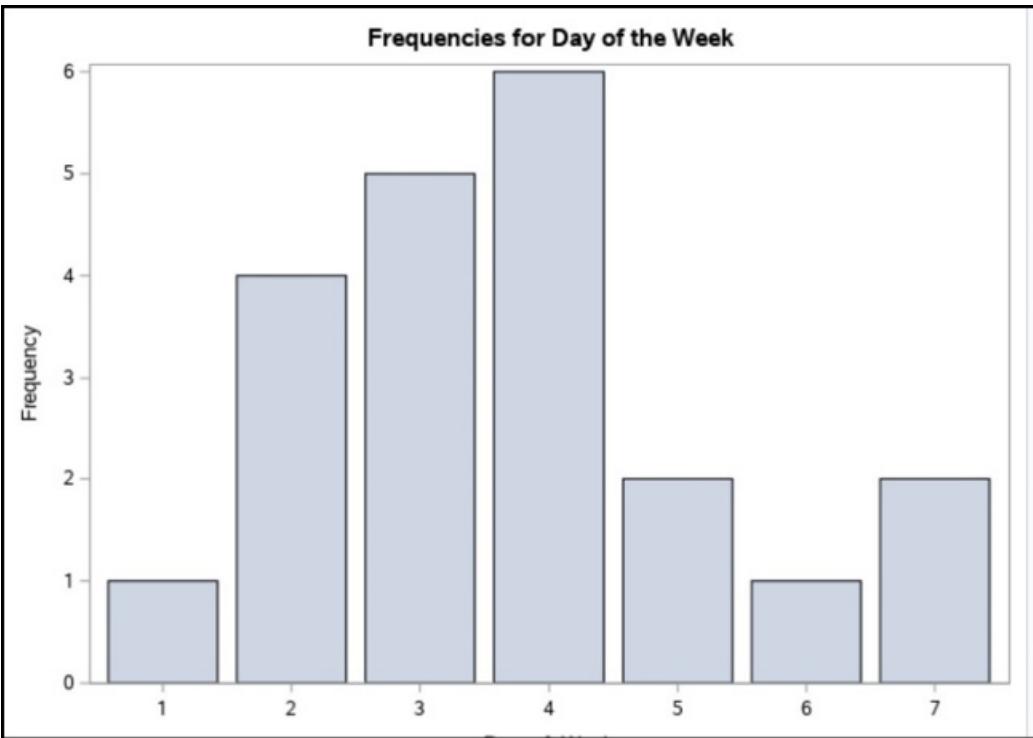
```
run;
```

- The two @ signs at the end of the INPUT statement prevent SAS from going to a new line when the DATA step iterates. It enables you to place data for more than one observation on a single line of data.
- The three functions, WEEKDAY, DAY, and YEAR, are used to extract the day of the week, the day of the month, and the year from the variable Date.
- The data set option OBS= is an instruction to stop processing when you reach observation 8. It is a convenient way to list the first  $n$  observations in a data set.
- You use PROC SGPlot to generate a vertical bar chart (VBAR) for the variable Day\_of\_Week.

Here is the output from Program 13.4.

**Figure 13.4: Output from Program 13.4**

Listing of the First Eight Observations from Extract				
Obs	Date	Day_of_Week	Day_of_Month	Year
1	01-05-2000	4	5	2000
2	02-08-2000	3	8	2000
3	04-23-2000	1	23	2000
4	04-12-2000	4	12	2000
5	08-21-2000	2	21	2000
6	08-21-2000	2	21	2000
7	08-22-2000	3	22	2000
8	12-12-2000	3	12	2000



The first eight observations are printed. Notice the date format (it uses hyphens because of the mmddyyd10. format) and the values for the other three variables.

The bar chart shows frequencies for the days of the week.

## Adding a Format to the Bar Chart

To make the bar chart more readable, you can write a format that will substitute the day abbreviations (Mon, Tue, etc.) for the numbers 1–7. The program that follows does just that:

### Program 13.5: Creating a Variable Representing Day of the Week Abbreviations

```

proc format;
  value DOW 1='Sun' 2='Mon' 3='Wed' 4='Thu'
        5='Fri' 6='Sat' 7='Sun';
run;

data Extract;
  informat Date mmddyy10.;
  input Date @@;
  Day_of_Week = weekday(Date);

```

```

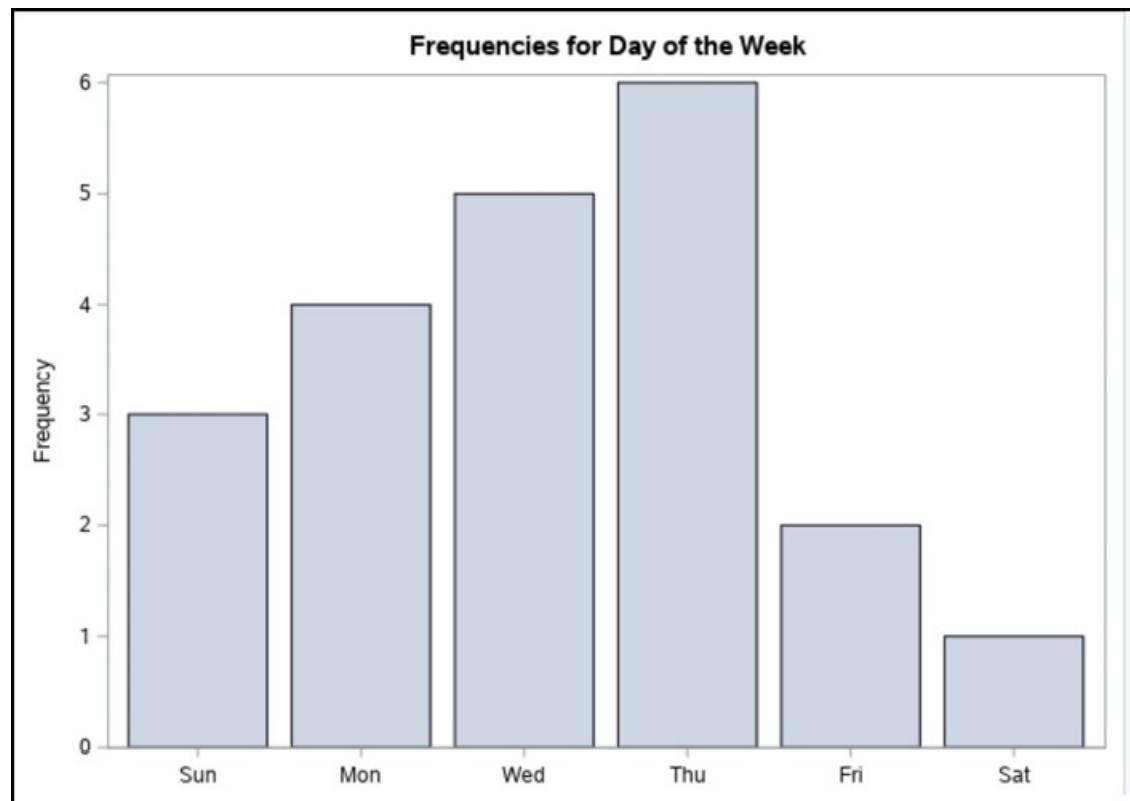
Day_of_Month = day(Date);
Year = year(Date);
format Date mmddyyd10. Day_of_Week DOW. ;
datalines;
1/5/2000 2/8/2000 4/23/2000 4/12/2000 8/21/2000 8/21/2000
8/22/2000
12/12/2000 12/15/2000 12/18/2000
2/22/2001 2/1/2001 4/18/2001 4/18/2001 4/18/2001 9/17/2001
12/25/2001
12/22/2001 3/3/2001 3/6/2001 3/7/2001
;
title "Listing of the First Eight Observations from EXTRACT";
proc print data=Extract (obs=8);
run;

title "Frequencies for Day of the Week";
proc sgplot data=Extract;
    vbar Day_of_Week;
run;

```

The bar chart now looks like this.

**Figure 13.5: Output from Program 13.5**



This is a substantial improvement over the previous bar chart.

# Computing Age from Date of Birth: The YRDIF Function

Computing a person's age, given the date of birth, is a problem faced by just about every programmer. Luckily, SAS has a function called YRDIF (year difference) that computes the difference between two SAS dates, in years. For example, suppose you have a variable called DOB (date of birth) and you want to compute a person's age as of January 1, 2020. The calculation is:

```
Age = yrdif(DOB, '01Jan2020'd);
```

The two arguments to the YRDIF function are the first and last dates from which you want to compute the interval. In this example, the first date is a SAS date and the second date is a date constant.

There are many applications where you need to compute age as of the last birthday. For example, you cannot vote even if your 18<sup>th</sup> birthday is the day after Election Day. You can use the INT (integer) function to extract the integer part of a number. The expression to compute a person's age on January 1, 2020, as of his or her last birthday is:

```
Age_Last = int(yrdif(DOB, '01Jan2020'd));
```

Notice that it is OK to have a SAS function as an argument to another SAS function. Just be careful with parentheses when you do this. You might want to draw the line at two functions in a single SAS statement. More than this can be a bit confusing to read.

If you need to round the result of the YRDIF function, use the ROUND function in place of the INT function.

## Conclusion

SAS stores dates as the number of days from January 1, 1960. SAS has informats to read and interpret dates in almost any form. Once you have a SAS date, you can use the WEEKDAY, DAY, MONTH, and YEAR functions to extract any of these values from the date. To learn more about SAS date functions, please refer to one (or both) of the following references:

Cody, Ron. 2010. *SAS Functions by Example, Second Edition*. Cary, NC: SAS Institute Inc.

Morgan, Derek. 2014. *The Essential Guide to SAS Dates and Times, Second Edition*. Cary, NC: SAS Institute Inc.

## Problems

1. You have a raw data file with the following data:

```
12345678901234567890 Ruler - not part of the data  
10/21/2015 12Jun2015  
12/25/2015 9Apr2014
```

The date starting in column 1 is in the form month/day/year. The date starting in column 12 is in the form 2-digit day of the month, 3-character month abbreviation, and 4-digit year.

Write an INPUT statement to read these two lines of data (you can use DATALINES). Call the first date Date1 and the second date Date2. Format both dates with the mmddyy10. format.

2. Run the program below to create a data set called Date\_Test:

### Program for Problem Sets 4

```
data Date_Test;  
    input Month Day Year;  
datalines;  
10 21 1988  
3 4 2015  
1 1 1960  
;
```

Modify this program so that you have a variable called Date that is a SAS date. Format this date using the DATE9. format.

3. The SASHELP data set Retail contains the variables Month, Day, and Year. Create a new, temporary data set called Dates that has these three variables plus one other called SAS\_Date that is a true SAS date. Format this variable using the mmddyy10. format and list the first five observations from this data set. The Retail data set contains other variables besides Month, Day, and Year. You do not want any of these other variables in your Dates data set.

4. Using the data from Problem 1, compute three new variables,

Month, Day (day of the week), and Year based on the date starting in column 1. Compute frequencies for these variables.

5. Run the following program to create a data set called Study:

### **Program for Problem Sets 5**

```
data Study;
  call streaminit(13579);
  do Subj = 1 to 10;
    Date = '01Jan2015'd + int(rand('uniform')*300);
    output;
  end;
  format Dates date9. ;
run;
```

Write a DATA step that will print out all dates in the Study data set that are before January 1, 2015 or after July 4, 2015. You can either use a DATA \_NULL\_ DATA step with a PUT statement or create a data set of out-of-range dates and use PROC PRINT to print it.

6. Run the program in Problem 5 except change the line that computes dates to:

```
Dates = '01Jan1950'd + int(rand('uniform')*15000);
```

Assuming these dates represent the date of birth, compute the age of each subject as of January 1, 2015. Print out the subject numbers (variable Subj) and ages for each person. Add a FORMAT statement to your PROC PRINT to assign the format 4.1 to the variable Age (which will list each age, rounded to a tenth of a year).

# Chapter 14: Subsetting and Combining SAS Data Sets

## Introduction

This chapter discusses ways to subset (filter) a SAS data set and to combine data from several SAS data sets. For those readers who understand SQL (structured query language), you should know that SAS supports SQL in PROC SQL. Because SQL is covered elsewhere in many texts, this chapter does not include a discussion of PROC SQL. The methods described in this chapter, SET, MERGE, and UPDATE, are DATA step statements that provide an alternative to PROC SQL.

## Subsetting (Filtering) Data in a SAS Data Set

As you saw in the first section of this book, you can subset or filter a SAS data set using interactive point-and-click operations provided by SAS Studio. This section shows you how to subset a SAS data set using DATA step programming.

To demonstrate how to create a subset starting with observations from an existing SAS data set, let's use one of the built-in data sets that comes with SAS Studio. There is a library called SASHELP that ships with SAS Studio and contains many data sets to be used for training.

To see a list of data sets in this library, perform the following steps: First, click the **Libraries** tab in the SAS Studio window.

**Figure 14.1: Select the Libraries Tab**



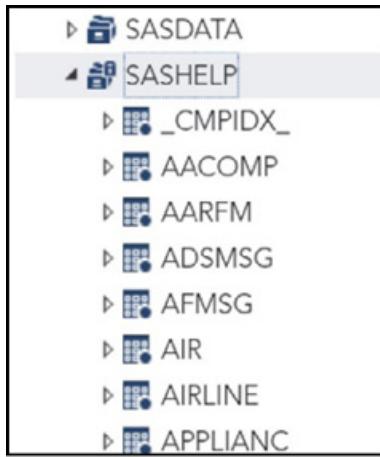
After you click the **Libraries** tab, select **My Libraries** and you will then see a list of libraries. (Note: your list might differ slightly from the list in Figure 14.2.)

**Figure 14.2: Opening the Libraries Tab**



If you click the triangle to the left of SASHelp, you will see a list of all the built-in data sets. A partial list looks like this:

**Figure 14.3: Partial list of SASHelp Data Sets**



We are going to select the Retail data set. Scroll down the list of data sets and select (double click) **RETAIL**.

**Figure 14.4: Opening the Retail Data Set**

View: Column names							Filter: (none)
Columns	Total rows: 58 Total columns: 5						
<input checked="" type="checkbox"/> Select all		SALES	DATE	YEAR	MONTH	DAY	
<input checked="" type="checkbox"/>	1	\$220	80Q1	1980	1	1	
<input checked="" type="checkbox"/>	2	\$257	80Q2	1980	4	1	
<input checked="" type="checkbox"/>	3	\$258	80Q3	1980	7	1	
<input checked="" type="checkbox"/>	4	\$295	80Q4	1980	10	1	
<input checked="" type="checkbox"/>	5	\$247	81Q1	1981	1	1	
<input checked="" type="checkbox"/>	6	\$292	81Q2	1981	4	1	
<input checked="" type="checkbox"/>	7	\$286	81Q3	1981	7	1	

You see the variables in this data set along with a listing of the data set.

Suppose your goal is to concentrate on data from the year 1980. You proceed as follows:

### Program 14.1: Using a WHERE Statement to Subset a SAS Data Set

```
data Year1980;
  set SASHelp.Retail;
  where Year = 1980;
run;
```

```
title "Listing of Data Set Year1980";
proc print data=Year1980 noobs;
run;
```

You use a SET statement to read the observations in data set SASHELP.Retail. Notice that you do not need to submit a LIBNAME statement. The SASHELP library is automatically made available to you every time you open SAS Studio. To subset the observations in data set Retail, you use a WHERE statement. Following the keyword WHERE, you specify the subsetting condition. In this example, you specify that you only want to see data for the year 1980. Using PROC PRINT, you obtain the following listing:

**Figure 14.5: Output from Program 14.1**

Listing of Data Set Year1980				
SALES	DATE	YEAR	MONTH	DAY
\$220	80Q1	1980	1	1
\$257	80Q2	1980	4	1
\$258	80Q3	1980	7	1
\$295	80Q4	1980	10	1

Only four observations met your selection criteria.

Let's look at a more complicated subset. Suppose you want data for the years 1980, 1982, and 1985, and you only want to look at observations where the Sales were greater than or equal to \$250. Here is the program.

### **Program 14.2: Demonstrating a More Complicated Query**

```
data Complicated;
  set SASHELP.retail;
  where Year in (1980, 1983, 1985) and Sales ge 250;
run;

title "Listing of Data Set Complicated";
proc print data=Complicated noobs;
run;
```

This WHERE statement is more complicated. You use an IN operator

to select any observations where Year is equal to 1980, 1983, or 1985. Because you also want to restrict observations where Sales are greater than \$250, you use the AND Boolean operator to add this condition. Here is the listing.

**Figure 14.6: Output from Program 14.2**

Listing of Data Set Complicated				
SALES	DATE	YEAR	MONTH	DAY
\$257	80Q2	1980	4	1
\$258	80Q3	1980	7	1
\$295	80Q4	1980	10	1
\$299	83Q1	1983	1	1
\$351	83Q2	1983	4	1
\$359	83Q3	1983	7	1
\$384	83Q4	1983	10	1
\$337	85Q1	1985	1	1
\$399	85Q2	1985	4	1
\$412	85Q3	1985	7	1
\$448	85Q4	1985	10	1

## Describing a WHERE= Data Set Option

An alternative to a WHERE statement is a WHERE= data set option. This data set option is one of many possible data set options available to you. You place all of the data set options that you want to use in a set of parentheses following the data set name. For example, you can rewrite Program 14.1 using a WHERE= data set option like this:

### Program 14.3: Rewriting Program 14.1 Using a WHERE= Data Set Option

```
data Year1980;
  set SASHELP.retail (where=(Year = 1980));
run;
```

You need to be careful with parentheses. The outermost set of parentheses encloses any data set options that you select—the

inner set of parentheses encloses the WHERE condition. You use an equal sign following the keyword WHERE when you are subsetting using a data set option. You do not use an equal sign when you are writing a WHERE statement. The resulting data set (Year1980) is identical to the one created in Program 14.1.

By the way, you can use a WHERE statement or a WHERE= data set option in any SAS procedure. If you want a listing of the 1980s data from the SASHELP.Retail data set, you can use the following statements.

### Program 14.4: Using a WHERE= Data Set Option in a SAS Procedure

```
proc print data=SASHELP.Retail (where=(Year = 1980));
run;
```

## Describing a Subsetting IF Statement

Suppose you have some raw data consisting of ID, gender, age, height, and weight. Earlier in this book, you developed a simple DATA step to read this collection of data. Suppose you want to create a SAS data set from the raw data, but you only want to see the female subjects. You have to use an IF statement because WHERE statements can only be used to subset SAS data sets. A demonstration of a subsetting IF is shown below.

### Program 14.5: Demonstrating the Subsetting IF Statement

```
data Females;
  input @1  ID      $3.
        @4  Gender  $1.
        @5  Age     3.
        @8  Height  2.
        @10 Weight  3. ;
  if Gender = 'F';
datalines;
001M 5465220
002F10161 98
003M 1770201
004M 2569166
005F    64187
```

```

006F 3567135
;

title "Listing of Data Set Females";
proc print data=Females;
  id ID;
run;

```

Notice that the IF statement does not have a THEN clause. This special IF statement, known as a subsetting IF, works like this: If the statement is true, the program continues; if the statement is false, the program returns to the top of the DATA step to read another line of data. Because an implicit output is performed at the bottom of the DATA step, only females will be in the resulting data set. Here is the listing.

**Figure 14.7: Output from Program 14.5**

Listing of Data Set Females				
ID	Gender	Age	Height	Weight
002	F	101	61	98
005	F	.	64	187
006	F	35	67	135

## A More Efficient Way to Subset Data When Reading Raw Data

Being a compulsive programmer, this author can't leave this section without showing you a more efficient way to write Program 14.5. Why read all of the data for males when you are not going to keep it? It is better to read just the single byte of data for Gender and only if it is equal to 'F' read the remaining data. Here is the program—the explanation follows:

**Program 14.6: Demonstrating a Trailing @**

```

data Females;
  input @4  Gender $1. @;
  if Gender = 'F' then

```

```
      input @5    Age     3.
          @8    Height   2.
          @10   Weight   3.;

      else delete;
datalines;
001M 5465220
002F10161 98
003M 1770201
004M 2569166
005F    64187
006F 3567135
;
```

You might ask, “What is the @ sign doing at the end of the first INPUT statement?” Here’s the explanation: If you have more than one INPUT statement in a DATA step, SAS will go to the next line of data each time it encounters another INPUT statement. The @ sign at the end of the line is called a *trailing @* and it tells SAS to “hold the line” and do not go to the next line of data when you encounter the next INPUT statement. **Without the trailing @, the program would read the value of Gender from one line and input the other values from the next line.**

The trailing @ on the first INPUT statement enables you to test the value of Gender and if it is equal to ‘F’ to read the values of Age, Height, and Weight from the **same line of data**. If Gender is not equal to ‘F’, the program executes the DELETE statement. This causes a return to the top of the DATA step without having to read the other values on the input record. **Keep the trailing @ in mind whenever you need to read a portion of your data to decide how to read other values from the same line.**

## Creating Several Data Subsets in One DATA Step

If you need to create several subsets from a single SAS data set, you can reduce processing time by creating multiple SAS data sets in one DATA step. Doing this reduces processing time because you read through the data just once instead of multiple times. This is

especially important when you are dealing with large files. For example, suppose you want to create three data sets from the SASHELP.Retail data set, each for a separate year. Here's how to do it.

### Program 14.7: Creating Several SAS Data Sets in One DATA Step

```
data Year1980 Year1981 Year1982;
  set SASHELP.Retail;
  if Year = 1980 then output Year1980;
  else if Year = 1981 then output Year1981;
  else if Year = 1982 then output Year1982;
run;
```

You name each of the data sets you want to create in the DATA statement. Next, you test the value of Year and use an OUTPUT statement, outputting observations to the appropriate data set. It is important to name the data set in the OUTPUT statement. If you use an OUTPUT statement without naming a data set, the program will output an observation to each of the data sets named in the DATA statement. One additional and important feature of this program is that when you include an OUTPUT statement in a DATA step, SAS does not perform an automatic output at the bottom of the DATA step.

When you run Program 14.7, the three data sets, Year1980, Year1981, and Year1982, are created.

## Combining SAS Data Sets (Combining Rows)

One way to combine two or more data sets is to “stack them up” one on top of the other (also referred to as *concatenation*). For example, you might have data sets for each of the four quarters of the year and want to put them together to create a data set with all the data for the year. The simple example that follows shows how to do this using a SAS DATA step.

For this example, you have two data sets (cleverly called One and Two). They are shown below.

**Figure 14.8: Data Sets One and Two**

Data Set One					Data Set Two				
ID	DOB	Gender	Height	Weight	ID	DOB	Gender	Height	Weight
001	10/21/60	M	68	160	004	05/13/78	M	75	190
002	11/11/81	F	62	120	005	08/23/88	F	59	99
003	01/05/83	M	72	220					

You use a SET statement to combine these two data sets, one after the other, like this:

### **Program 14.8: Using a SET Statement to Combine Two SAS Data Sets**

```
data Both;
  set One Two;
run;
```

You list each of the data sets you want to combine in the SET statement. The program first reads all the observations from data set One. When it reaches the end of the file, it switches to data set Two and continues to read observations from data set Two until it reaches the end of that file. The result is all the observations from the two data sets.

Here is a listing of data set Both.

**Figure 14.9: Listing of Data Set Both**

Listing of Data Set Both				
ID	DOB	Gender	Height	Weight
001	10/21/60	M	68	160
002	11/11/81	F	62	120
003	01/05/83	M	72	220
004	05/13/78	M	75	190
005	08/23/88	F	59	99

In case you want to run this program as is or with modifications, the program to create data sets One and Two is listed here.

## Program 14.9: Program to Create Data Sets One and Two

```
data One;
  informat ID $3. DOB mmddyy10. Gender $1.;
  input ID DOB Gender Height Weight;
  format DOB mmddyy10.;
datalines;
001 10/21/1950 M 68 160
002 11/11/1981 F 62 120
003 1/5/1983 M 72 220
;

data Two;
  informat ID $3. DOB mmddyy10. Gender $1.;
  input ID DOB Gender Height Weight;
  format DOB mmddyy10.;
datalines;
004 5/13/1978 M 70 190
005 8/23/1988 F 59 98
;
```

## Adding a Few Observations to a Large Data Set (PROC APPEND)

If your goal is to add observations (typically from a relatively small data set) to a large data set, you can, of course, use a SET statement and list the names of the two data sets, like this:

### Program 14.10: Using a SET Statement to Solve the Problem

```
Data Both;
  Set Big Small;
run;
```

When you run this program, SAS will first read all the observations from data set Big and, when it reaches the end of the file, it will then read all the observations from data set Small.

If this is something that you need to do often, you might be concerned about efficiency, especially if data set Big is really big (millions or tens of millions of observations). A much more efficient method is to use PROC APPEND. A program that accomplishes the

same goal as Program 14.10 (almost) is demonstrated in Program 14.11.

### Program 14.11: Using PROC APPEND to Add Observations from One Data Set to Another

```
proc append base=Big Data=Small;  
run;
```

You name the first data set in the BASE= procedure option and the data set to be added in the DATA= procedure option. This program **does not need to read any observations from data set Big**—it strips the end-of-file marker from the end of data set Big and then appends the new observations. If data set Big is large, this is a huge improvement in efficiency over using a SET statement. As with most things that look really quick and easy, there is a catch: In this program, you are **replacing** data set Big with the contents of Big and Small. If you have errors or bad data in data set Small, it is not easy to undo the process. It would be a good idea to have a backup copy of Big somewhere.

This method should only be used when you are confident that: One, the data that you are adding does not contain errors; and two, data set Small has all the same variables and attributes (especially character variable lengths) as data set Big. The reason for this is that the data descriptor (containing all the information such as variable types and lengths) is taken from data set Big.

If, for example, you had a character variable in data set Small that had a longer length than a variable of the same name in data set Big, that variable would be truncated. As a matter of fact, PROC APPEND would not even run in this situation unless you added an option called FORCE. The opinion of this author is that if you need to use the FORCE option, you had better know exactly what you are doing. One final point: If you have variables that have different names in the two data sets but are otherwise equivalent, you can use a RENAME= data set option to rename variables in one data set to match the variable names in the other data set.

## Interleaving Data Sets

If you have two or more data sets that are so large that, when put together, it would be difficult or impossible to sort the result, you can sort each of the data sets first and then interleave them. That is, you can add observations to the resulting data set from each of the constituent data sets so that the result is in sorted order. All that is necessary to accomplish this is to follow the SET statement with a BY statement, listing the variables that define the sorted order. Here is an example:

### Program 14.12: Demonstrating How to Interleave Two or More Data Sets

```
*Note: data sets One, Two, and Three are sorted by ID;  
data Combined;  
  set One Two Three;  
  by ID;  
run;
```

The resulting data set (Combined) will be in ID order.

## Merging Two Data Sets (Adding Columns)

By *merging*, we usually mean adding extra columns (variables) to a SAS data set based on one or more matching variables (such as ID or Name) from another data set. In SQL terms, you would be doing a *join* (inner, left, right, or, full). Imagine you are collecting clinical data on some patients. In one data set, you have ID, gender, and date of birth—in another data set, you have measurements such as weight, heart rate, and blood pressure. When it is time to analyze your data, you want to combine data from these two data sets. The program below creates two data sets, Patients and Visits.

### Program 14.13: Creating the Patients and Visits Data Sets

```
data Patients;  
  informat ID $4. Gender $1. DOB mmddyy10.;  
  input ID Gender DOB;  
  format DOB date9.;  
datalines;
```

```

0001 M 10/10/1980
0023 F 1/2/1977
1243 M 6/17/2000
0002 M 8/23/1981
4535 F 2/25/1967
;

data Visits;
  informat ID $4. Visit_Date mmddyy10.;
  input ID Visit_Date Weight HR SBP DBP;
  format Visit_Date date9.;
datalines;
0023 2/10/2015 122 76 122 78
4535 10/21/2014 155 78 138 88
0001 11/11/2014 210 68 118 78
;

```

A listing of these two data sets is shown next:

**Figure 14.10: Listing of Data Sets Patients and Visits**

**Listing of Data Set Patients**

ID	Gender	DOB
0001	M	10OCT1980
0023	F	02JAN1977
1243	M	17JUN2000
0002	M	23AUG1981
4535	F	25FEB1967

**Listing of Data Set Visits**

ID	Visit_Date	Weight	HR	SBP	DBP
0023	10FEB2015	122	76	122	78
4535	21OCT2014	155	78	138	88
0001	11NOV2014	210	68	118	78

In order to combine (merge) these two data sets, you must first sort each of them by the variable (or variables) that you plan to use to match the two data sets. Here is the program, followed by an explanation:

**Program 14.14: Merging Two SAS Data Sets**

```

proc sort data=Patients;
```

```

    by ID;
run;

proc sort data=Visits;
    by ID;
run;

data Merged;
    merge Patients Visits;
    by ID;
run;

title "Listing of Data Set MERGED";
proc print data=Merged;
    id ID;
run;

```

You use PROC SORT to sort both data sets by ID. Next, you use a MERGE statement to combine the patient and visit data and a BY statement to instruct the program to use ID as the matching variable. Here is the resulting data set.

**Figure 14.11: Output from Program 14.4**

Listing of Data Set Merged								
ID	Gender	DOB	Visit_Date	Weight	HR	SBP	DBP	
0001	M	10OCT1980	11NOV2014	210	68	118	78	
0002	M	23AUG1981	.	.	.	.	.	
0023	F	02JAN1977	10FEB2015	122	76	122	78	
1243	M	17JUN2000	.	.	.	.	.	
4535	F	25FEB1967	21OCT2014	155	78	138	88	

What you see here is all the observations from the Patients data set matched with all the observations in the Visits data set (known as a *full join* in SQL), even if there is no corresponding visit for a particular ID.

This is probably not what you want. In this example, you most likely want to see a listing of only those patients who were in the Visits data set. Luckily, SAS has a way to control which observations will be included in the merged data set. Enter the IN= data set option.

## Controlling Which Observations Are Included in a Merge (IN= Data Set Option)

You have previously seen a WHERE= data set option. You can now add to your collection of data set options by learning about the IN= data set option. Here's how it works:

When you are merging two data sets, you might find a matching ID in both data sets or you might only find an ID in one but not the other data set. To determine whether a data set is making a contribution on a particular merge, you add an IN= *variable\_name* data set option following each of the two data sets in the MERGE statement. Let's examine the value of these variables (called *contributor variables* by some SAS programmers) in the following program.

### Program 14.15: Demonstrating the IN= Data Set Option

```
*Note: both data set previously sorted by ID;
data Merged;
  merge Patients(in=In_Patients) Visits(In=In_Visits);
  by ID;
  put ID= In_Patients= In_Visits=;
run;
```

The variable name following the IN= data set option can be any valid SAS name. However, this author likes to use names that help you remember which variable goes with which data set. Hence, the names In\_Patients and In\_Visits.

**These two variables are temporary variables. That is, they are not included in the output data set.** They are available only during the DATA step to help you determine how to perform the merge. A PUT statement was added to Program 14.15 so that you can examine these two variables. Here is the relevant section of the SAS log:

**Figure 14.12: Partial Listing of the SAS Log from Program 14.15**

```

83      *Note: both data set previously sorted by ID;
84      data Merged;
85          merge Patients(in=In_Patients) Visits(In=In_Visits);
86          by ID;
87          put ID= In_Patients= In_Visits=;
88      run;

ID=0001 In_Patients=1 In_Visits=1
ID=0002 In_Patients=1 In_Visits=0
ID=0023 In_Patients=1 In_Visits=1
ID=1243 In_Patients=1 In_Visits=0
ID=4535 In_Patients=1 In_Visits=1

```

In each observation, the IN= variable is a 1 when that data set is contributing to the merge and 0 otherwise. Patient 0001 is in both data sets, so the value of the two IN= variables is 1. Patient 0002 is in the Patients data set but not the Visits data set, so In\_Patients is a 1 and In\_Visits is a 0 for this observation.

You can use the IN= variables to control which observations you want to include in the output data set. For example, to include only observations where there is a contribution from both data sets (referred to as an *inner join* in SQL), you would use:

```
if In_Patients=1 and In_Visits=1;
```

This can also be written as:

```
if In_Patients and In_Visits;
```

You do not need to include the equals 1 because SAS understands that these variables are either 1 or 0 (interpreted as true or false).

In this example, every patient in the Visits data set was matched to a record in the Patients data set. In the real world, this would not always be the case. The next program creates a data set containing patients in the Visits data set who are matched with the data in the Patients data set.

To make the program more general, statements were added to print an error message for any patient in the Visits data set who is missing from the Patients data set. It is good programming to expect the unexpected! Here is the program.

## Program 14.16: Using the IN= Data Set Option to Control the Merge Operation

```
*Note: Both data sets already sorted by ID;  
  
title "Listing of Patients with Visits Who are Not in the  
Patients Data Set";  
data Only_Visit_Patients;  
    file print;  
    merge Patients(in=In_Patients) Visits(in=In_visits);  
    by ID;  
    if In_Visits then output;  
    if In_Visits and not In_Patients then  
        put "Patient " ID "not found in the Patients data set.";  
run;  
  
title "Listing of Data Set Only_Visit_Patients";  
proc print data=Only_Visit_Patients;  
    id ID;  
run;
```

You test the value of the variable In\_Visits, and only output observations where In\_Visits is true. If there is an ID in the Visits data set that does not have an observation in the Patients data set, a message is written. By default, PUT statements write data to the SAS log. By including the statement FILE PRINT, the output from this PUT statement will be sent to the RESULTS window. It turns out that there were no patients who were in the Visits data set and not in the Patients data set, so no messages were printed by the PUT statement. Here is the output.

**Figure 14.13: Output from Program 14.16**

Listing of Data Set Only_Visit_Patients								
ID	Gender	DOB	Visit_Date	Weight	HR	SBP	DBP	
0001	M	10OCT1980	11NOV2014	210	68	118	78	
0023	F	02JAN1977	10FEB2015	122	76	122	78	
4535	F	25FEB1967	21OCT2014	155	78	138	88	

## Performing a One-to-Many or Many-to-One

## Merge

If you want to merge two data sets where one data set has more than one observation for your choice of BY variable(s) and the other data set has exactly one observation for the same BY variable(s), you can still perform a merge. For this example, you want to merge patient data (ID, Gender, and DOB) with visit data (ID, date of visit, HR, SBP, and DBP). To demonstrate this merge, we will use the Patients data set (see Program 14.13) and a new data set, Many\_Visits, which contains several visits for each patient. The program to create the Many\_Visits data set is shown in Program 14.17.

### Program 14.17: Creating Another Data Set to Demonstrate a One-to-Many Merge

```
data Many_Visits;
  informat ID $4. Visit_Date mmddyy10.;
  input ID Visit_Date HR SBP DBP;
  format Visit_Date date9.;

  datalines;
  0023 2/10/2015 122 76 122 78
  0023 3/10/2015 120 74 120 76
  4535 10/21/2014 155 78 138 88
  0001 11/11/2014 210 68 118 78
  0001 12/20/2014 210 68 120 82
  0001 1/5/2015 212 70 210 80
  ;
```

For reference, both data sets (Patients and Many\_Visits) are displayed below.

**Figure 14.14: Listing of Data Sets Patients and Many\_Visits**

Listing of Data Set Patients			Listing of Data Set Many_Visits				
ID	Gender	DOB	ID	Visit_Date	HR	SBP	DBP
0001	M	10OCT1980	0023	10FEB2015	122	76	122
0002	M	23AUG1981	0023	10MAR2015	120	74	120
0023	F	02JAN1977	4535	21OCT2014	155	78	138
1243	M	17JUN2000	0001	11NOV2014	210	68	118
4535	F	25FEB1967	0001	20DEC2014	210	68	120
			0001	05JAN2015	212	70	210

You can now sort and merge the two data sets like this.

### Program 14.18: Performing a One-to-Many Merge

```

proc sort data=Patients;
  by ID;
run;

proc sort data=Many_Visits;
  by ID;
run;

data One_to_Many;
  merge Patients Many_Visits;
  by ID;
run;

title "Listing of data set One_To_Many";
proc print data=One_to_Many;
  id ID;
run;

```

When you run this program, the variables Gender and DOB from the Patients data set are combined with matching observations from the Many\_Visits data set (based on ID) to create the data set One\_To\_Many, as shown in the listing below.

### Figure 14.15: Output from Program 14.18

Listing of data set One_To_Many						
ID	Gender	DOB	Visit_Date	HR	SBP	DBP
0001	M	10OCT1980	11NOV2014	210	68	118
0001	M	10OCT1980	20DEC2014	210	68	120
0001	M	10OCT1980	05JAN2015	212	70	210
0002	M	23AUG1981	.	.	.	.
0023	F	02JAN1977	10FEB2015	122	76	122
0023	F	02JAN1977	10MAR2015	120	74	120
1243	M	17JUN2000	.	.	.	.
4535	F	25FEB1967	21OCT2014	155	78	138

In this example, you could reverse the order of the two data sets (performing a many-to-one merge) and the result would be identical to the one obtained here.

**CAUTION:** Do not attempt to merge two data sets where there are multiple observations for each BY variable in both data sets and the number of multiples is not the same in both files.

## Merging Two Data Sets with Different BY Variable Names

You might find yourself needing to merge two data sets, but the BY variable has a different variable name in each of the two files (this is corollary number 17 of Murphy's Law). The solution to this problem is actually very simple. Let's redo Program 14.14, but make a slight change to the data set Visits. The new data set (Visits\_2) has all the same values of data set Visits except that the ID variable is named Pt.

Here is a listing of data set Visits\_2:

**Figure 14.16: Listing of Data Set Visits\_2**

**Listing of Data Set Visits\_2**

Pt	Visit_Date	Weight	HR	SBP	DBP
0023	10FEB2015	122	76	122	78
4535	21OCT2014	155	78	138	88
0001	11NOV2014	210	68	118	78

Before we continue, here is the program that was used to create the data set Visits\_2 (in case you want to try this yourself):

### **Program 14.19: Program to Create Data Set Visits\_2**

```
data Visits_2;
  informat Pt $4. Visit_Date mmddyy10.;
  input Pt Visit_Date Weight HR SBP DBP;
  format Visit_Date date9.;
datalines;
0023 2/10/2015 122 76 122 78
4535 10/21/2014 155 78 138 88
0001 11/11/2014 210 68 118 78
;
```

In order to merge the two files Patients and Visits\_2, you have to rename one of the variables (either ID or Pt) so that both BY variables have the same name. You do this with a RENAME= data set option. In this example, you are going to rename the variable Pt in the Visits\_2 data set to ID.

### **Program 14.20: Using a RENAME= Data Set Option to Rename the Variable Pt to ID**

```
proc sort data=Patients;
  by ID;
run;

proc sort data=Visits_2;
  by Pt;
run;

data Merged;
  merge Patients
    Visits_2(rename=(Pt = ID));
by ID;
run;
```

```
title "Listing of Data Set Merged";
proc print data=Merged;
  id ID;
run;
```

This is another opportunity to get confused with the parentheses. Remember, the outer-most set holds all of the data set options, and the inner-most set is a list of old variable names and new variable names. The form of the RENAME= data set option is:

*Data-Set-Name*( rename=(*Old-name1*=*New\_Name1* *Old\_Name2*=*New\_Name2* . . . ));

You can rename as many variables as needed using the RENAME= data set option. However, the two variables on each side of the equal sign must be the same type (either character or numeric).

Here is the output from Program 14.20.

**Figure 4.17: Output from Program 14.20**

Listing of Data Set Merged							
ID	Gender	DOB	Visit_Date	Weight	HR	SBP	DBP
0001	M	10OCT1980	11NOV2014	210	68	118	78
0002	M	23AUG1981	.	.	.	.	.
0023	F	02JAN1977	10FEB2015	122	76	122	78
1243	M	17JUN2000	.	.	.	.	.
4535	F	25FEB1967	21OCT2014	155	78	138	88

If you have a situation where the two variables that you would like to use in a merge are not the same type, you have to do a bit more work. The solution to this problem is discussed in the next section.

## Merging Two Data Sets with One Character and One Numeric BY Variable

Because you cannot use a RENAME= data set option with two variables of different types, you need to create a new variable with

the same name and type as the matching variable in the other data set. In this example, you have a variable called SS (Social Security number) in both files; however, one is stored as a character string and the other as a numeric value. Here is a listing of the two data sets you want to merge.

**Figure 14.18: Listings of Data Set One\_Char and Two\_Num**

Listing of Data Set One_Char			Listing of Data Set Two_Num		
SS	Gender	Age	Visit_Date	Fee_Paid	SS
123-45-6789	M	45	10/14/2015	Yes	123456789
088-54-1950	F	23	02/10/2015	No	88541950
321-43-7766	M	68	03/23/2015	Yes	321437766

The program to create these two data sets is listed next, in case you want to try it out for yourself.

**Program 14.21: Program to Create Data Sets One\_Char and Two\_Num**

```
data One_Char;
    informat SS $11. Gender $1. ;
    input SS Gender Age;
datalines;
123-45-6789 M 45
088-54-1950 F 23
321-43-7766 M 68
;
data Two_Num;
    informat Visit_Date mmddyy10. Fee_Paid $3. ;
    input SS Visit_Date Fee_Paid;
    format Visit_Date mmddyy10. ;
datalines;
123456789 10/14/2015 Yes
088541950 2/10/2015 No
321437766 3/23/2015 Yes
;
```

Before you conduct a merge, it is a good idea to run PROC CONTENTS (or examine the variables using SAS Studio) to ensure that the BY variables are the same type (and length). With this in mind, here is a section of PROC CONTENTS output for each of the

two data sets.

**Figure 14.19: Output from PROC CONTENTS for Data Set One\_Char**

Alphabetic List of Variables and Attributes				
#	Variable	Type	Len	Informat
3	Age	Num	8	
2	Gender	Char	1	\$1.
1	SS	Char	11	\$11.

**Figure 14.20: Output from PROC CONTENTS for Data Set Two\_Num**

Alphabetic List of Variables and Attributes					
#	Variable	Type	Len	Format	Informat
2	Fee_Paid	Char	3		\$3.
3	SS	Num	8		
1	Visit_Date	Num	8	MMDDYY10.	MMDDYY10.

You can see that the variable SS is stored as character in data set One\_Char and numeric in data set Two\_Num.

It's time to merge these two data sets. You have a choice—convert the character variable to a numeric variable or vice versa. Let's choose the latter—changing the numeric variable to character. Here is the program.

### **Program 14.22: Merging Two Data Sets with One Character and One Numeric BY Variable**

```
data Two_Char;
  set Two_Num(rename=(SS = SS_Num));
  SS = put(SS_Num,SSN11.);
  drop SS_Num;
run;

proc sort data=One_Char;
  by SS;
run;
```

```
proc sort data=Two_Char;
  by SS;
run;

Data One_and_Two;
  merge One_Char Two_Char;
  by SS;
run;

title "Listing of Data Set One_and_Two";
proc print data=One_and_Two noobs;
run;
```

There is quite a lot going on here. To start, you need to create a new data set (called Two\_Char in this example) with a variable called SS that is stored as a character value. You use a SET statement to read the observations from data set Two\_Num. Because you need the final character variable to be called SS, you use the RENAME= data set option to rename the original, numeric variable SS to SS\_Num.

The next step is to use the PUT function to perform a numeric-to-character conversion. Here's how it works: There is a SAS format called SSN11. that writes out numeric representations of Social Security numbers as 11-byte character strings, complete with leading zeros and hyphens. The PUT function takes two arguments. The first argument is the variable that you want to convert (SS\_Num in this example), and the second argument is a format that you want to use to format this value. The result is a character string in the same form as in data set One\_Char. Because you no longer need (or want) the variable SS\_Num in your data set, you use a DROP statement. This statement says that when SAS writes out observations to a data set, it should not include any variables listed in the DROP statement.

Output from Program 14.22 is listed next:

**Figure 14.21: Output from Program 14.22**

### **Listing of Data Set One\_and\_Two**

SS	Gender	Age	Visit_Date	Fee_Paid
088-54-1950	F	23	02/10/2015	No
123-45-6789	M	45	10/14/2015	Yes
321-43-7766	M	68	03/23/2015	Yes

That wasn't easy, but it is very likely that some time in your programming career, you will face this problem. Now you know the solution.

## **Updating a Master File from a Transaction File (UPDATE Statement)**

The final (yeah!) section of this chapter describes how to update values in a master file, using transaction data in a transaction file. In this example, you have a file of item numbers, descriptions, and prices, and you want to modify the file to change the price of a few items. Here is a program to create the master Price file and a listing of the file.

### **Program 14.23: Program to Create the Master File**

```
data Price;
    input @1 Item_Number $5.
        @7 Description $10.
        @18 Price;
datalines;
12345 Hammer      11.98
22222 Saw          25.89
44010 Nails 10p   17.95
44008 Nails 8p    15.56
;

title "Listing of Data Set Price";
proc print data=Price;
    id Item_Number;
run;
```

Here is the listing:

**Figure 14.22: Data Set Price**

Listing of Data Set Price		
Item_Number	Description	Price
12345	Hammer	11.98
22222	Saw	25.89
44010	Nails 10p	17.95
44008	Nails 8p	15.56

You want to change the price of the hammer (item number 12345) to \$12.98 and the price of 8 penny nails (item 44008) to \$16.50. You first create a data set with these two items (you only need the item number and the new price). This will be your transaction data set.

**Program 14.24: Creating the Transaction Data Set**

```
data Transact;
  informat Item_Number $5.;
  input Item_Number Price;
datalines;
12345 12.98
44008 16.50
;

title "Listing of Data Set Transact";
proc print data=Transact;
  id Item_Number;
run;
```

This is the listing:

**Figure 14.23: Output from Program 14.24**

Listing of Data Set Transact	
Item_Number	Price
12345	12.98
44008	16.50

Here is the program to update the prices in your master file:

**Program 14.25: Updating Your Master File Using a Transaction**

## Data Set

```
proc sort data=Price;
  by Item_Number;
run;

proc sort data=Transact;
  by Item_Number;
run;

data Price_100ct2020;
  update Price Transact;
  by Item_Number;
run;

title "Listing of Data Set Price_100ct2020";
proc print data= Price_100ct2020;
  id Item_Number;
run;
```

You first sort both data sets by Item\_Number. Next, you decide to give the new Price data set a different name—Price10Oct2020. Giving the new Price data set a new name is a good idea as it helps prevent confusion. Finally, you use an UPDATE statement to update the prices in the original Price data set. UPDATE is similar to MERGE. However, in a MERGE, if you have a variable in the second data set with the same name as a variable in the first data set, the value in the second data set will replace the value in the first data set, even if that value is a missing value. When you use the UPDATE statement, a missing value in the second data set **does not** replace the value in the first data set—just what you want to happen. Here is the listing:

**Figure 14.24: Output from Program 14.25**

### **Listing of Data Set Price\_10Oct2020**

Item_Number	Description	Price
12345	Hammer	12.98
22222	Saw	25.89
44008	Nails 8p	16.50
44010	Nails 10p	17.95

You see that the two prices are updated. Keep the UPDATE statement in mind whenever you need to replace values in one data set using values from another data set. It is one of the often-forgotten SAS statements.

## **Conclusion**

This was clearly one of the more difficult chapters. However, many, if not most, programming problems require you to manipulate data from multiple data sets. If you know and love SQL, you can use that knowledge to combine your SAS data sets. The choice of SET, MERGE, and UPDATE versus SQL is not always easy. We often use what we know best. Because PROC SQL came much later in the evolution of SAS, many of the “older” SAS programmers choose to use DATA step processing.

## **Problems**

1. Starting with the SASHELP data set Fish, create a data set called Small\_Perch that contains only perch that weigh less than 50 (whatever the weigh units are). Do this using a WHERE statement.
2. Repeat Problem 1 using a WHERE= data set option.
3. Run the program shown below. If you don’t want to enter it, you will find this program included in a program called “Programs for Data Sets.sas”, included in your download of data for this book. Use a subsetting IF statement to include only those subjects where the sum of Q1–Q3 is greater than using equal to 6.

### **Program for Problem Sets 1**

```

data Questionnaire;
  informat Gender 1. Q1-Q4 $1. Visit date9. ;
  input Gender Q1-Q4 Visit Age;
  format Visit date9. ;
datalines;
1 3 4 1 2 29May2015 16
1 5 5 4 3 01Sep2015 25
2 2 2 1 3 04Jul2014 45
2 3 3 3 4 07Feb2015 65
;

```

4. Starting with the SASHELP data set Cars, create two temporary data sets. The first one (Cheap) should include all the observations from Cars where the MSRP (manufacturer's suggested retail price) is less than or equal to \$11,000. The other (Expensive) should include all the observations from Cars where the MSRP is greater than or equal to \$100,000. Use a KEEP= data set option to include only the variables Model, Type, Origin, and MSRP from the Cars data set. Create these two data sets in one DATA step. Use PROC PRINT to list the observations in Cheap and Expensive. Even though there are no missing values for the variable MSRP, write your program so that any observation with a missing value for MSRP will not be written to data set Cheap.
5. Run Program for Problem Sets 6 to create two data sets (FirstQtr and SecondQtr). Then create a new data set (FirstHalf) that contains all the observations from FirstQtr and SecondQtr.

### **Program for Problem Sets 6**

```

data FirstQtr;
  input Name $ Quantity Cost;
datalines;
Fred 100 3000
Jane 90 4000
April 120 5000
;
data SecondQtr;
  input Name $ Quantity Cost;
datalines;
Ron 200 9000
Jan 210 9500
Steve 177 5400
;
```

6. Repeat Problem 5, except use PROC APPEND to combine the observations from data sets FirstQtr and SecondQtr. Because you want the resulting data set to be called FirstHalf, you will first need to make a copy of FirstQtr that is called First\_Half.
7. Run Program for Problem Sets 7. Then create a new data set (Both) that contains ID, X, Y, Z, and Name. Include only those IDs that are in both data sets.

### **Program for Problem Sets 7**

```

data First;
  input ID $ X Y Z;
datalines;
001 1 2 3
004 3 4 5
002 5 7 8
006 8 9 6
;
data Second;
  input ID $ Name $;
datalines;
002 Jim
003 Fred
001 Susan
004 Jane
;

```

8. Repeat Problem 7, except include all observations from each data set, even if there is no corresponding ID in one of the files.
9. Run Program for Problem Sets 8. Write a program to create a new data set called New\_Prices, where the price of item X200 is \$410 and the price of item A123 is \$121. Caution: Item\_Number is a character variable and the data set is not sorted by Item\_Number.

### **Program for Problem Sets 8**

```

data Prices;
  input Item_Number $ Price;
datalines;
A123 $123
B76 4.56
X200 400
D88 39.75
;

```

# Chapter 15: Describing SAS Functions

## Introduction

You have already seen a number of SAS functions in earlier chapters of this book. This chapter explores some of the most useful SAS functions that work with numeric and character data.

To review, SAS functions either return some system value or perform some calculation and return a value. Here are some useful facts about SAS functions:

- All SAS function names are followed by zero or more arguments, placed in parentheses following the function name.
- Although there are some exceptions, arguments to SAS functions can be variable names (the most common type of argument), constants (numbers for numeric functions, character values in quotation marks for character arguments), expressions (such as arithmetic expressions), or even other functions.
- Functions can only return a single value and, in most cases, the arguments to SAS functions do not change after you execute the function.

You will also see two CALL routines in this chapter. CALL routines are something like functions except for two important differences: One, some or all of the arguments can change value after the call; and two, you do not use a CALL routine in an assignment statement.

The following examples, using the WEEKDAY function (which returns the day of the week given a SAS date), demonstrate the variety of arguments that are valid with most SAS functions:

- Day = weekday(Date) where Date is the name of a SAS variable
- Day = weekday(20013) a numeric constant

- Day = weekday('01Jan2015'd) a date constant
- Day = weekday(today()) where Today is a SAS function that returns the current date
- Day = weekday(Date + 1) an arithmetic expression giving the day of the week one day after Date

## Describing Some Useful Numeric Functions

Each function in this chapter will be described and you will see either a full program or a statement showing how the function works.

### Function Name: MISSING

What it does: The MISSING function returns a value of true (1) if the argument is a missing value and false (0) otherwise.

Arguments: A character or numeric value

Examples:

### Program 15.1: Demonstrating the MISSING Function

```
data Old_Miss;
  input ID $ Age;
  if missing(Age)      then Age_group = .;
    else if Age le 50 then Age_group = 1;
    else                Age_group = 2;
datalines;
001 15
002 .
003 78
004 26
;

title "Listing of Data Set Old_Miss";
proc print data=Old_Miss noobs;
run;
```

Explanation: You use the MISSING function to test if Age is a missing value. If so, you assign a missing value to the variable Age\_Group. Here is the listing from Program 15.1:

**Figure 15.1: Output from Program 15.1**

### **Listing of Data Set Old\_Miss**

ID	Age	Age_group
001	15	1
002	.	.
003	78	2
004	26	1

Notice that Age\_Group has a missing value for ID 002.

### **Function Name: N**

What it does: Returns the number of nonmissing values in the list of arguments

Arguments: One or more numeric values. If any of the arguments are in the form *Variable1-VariableN*, the list must be preceded by the keyword OF.

Examples:

```
X1=1;      X2=. ;      X3=3;      Age=27;      Wt=. ;      Ht=68;
```

- `n(of X1-X3) = 2`
- `n(Age, Wt, Ht) = 2`

Explanation:

- In the list of variables X1–X3, there are two nonmissing values.
- Among the three variables Age, Wt, and Ht, there are two nonmissing values.

### **Function Name: NMISS**

What it does: Returns the number of missing values in the list of arguments

Arguments: One or more numeric values. If any of the arguments are in the form *Variable1-VariableN*, the list must be preceded by the keyword OF.

Examples:

```
X1=1;      X2=. ;      X3=3;      Age=27;      Wt=. ;      Ht=68;
```

- `nmiss(of X1-X3) = 1`
- `nmiss(Age, Wt, Ht) = 1`

Explanation:

- In the list of variables X1–X3, there is one missing value.
- Among the three variables Age, Ht, and Wt, there is one missing value.

## Function Name: SUM

What it does: Returns the sum of the arguments. In performing this operation, missing values are ignored. If all the arguments are missing, the result is a missing value.

Arguments: One or more numeric values. If any of the arguments are in the form *Variable1-VariableN*, the list must be preceded by the keyword OF.

Examples:

```
X1=1;      X2=. ;      X3=3;      Age=27;      Wt=. ;      Ht=68;
```

- `sum(of X1-X3) = 4`
- `sum(Age, Wt, Ht) = 95`

Explanation:

- Because X2 is a missing value, the result is the sum of 1 and 3.
- The sum of Age (27) and Ht (68) is 95. The missing value is ignored.

Because the SUM function returns a missing value when all of its arguments are missing values, a popular trick to return a value of 0 in this situation is to include a 0 as one of the arguments. For example,

If Cost1–Cost5 are all missing, the expression

```
sum(0,of Cost1-Cost5)
```

will return a value of 0.

## Function Name: MEAN

What it does: Returns the mean (average) of the arguments. In performing this operation, missing values are ignored. If all the arguments are missing, the result is a missing value.

Arguments: One or more numeric values. If any of the arguments are in the form *Variable1-VariableN*, the list must be preceded by the keyword OF.

Examples:

```
X1=1;      X2=. ;      X3=3;      Age=27;      Wt=. ;      Ht=68;
```

- `mean(of X1-X3) = 2`
- `mean(Age, Wt, Ht) = 47.5`

Explanation:

- X2 is a missing value and is ignored. The mean is computed as  $(1 + 3) / 2$ .
- The mean of Age (27) and Ht (68) is 47.5. The missing value is ignored.

## Function Name: MIN

What it does: Returns the smallest nonmissing value of its arguments. If all the arguments are missing, the result is a missing value.

Arguments: One or more numeric values. If any of the arguments are in the form *Variable1-VariableN*, the list must be preceded by the keyword OF.

Examples:

```
X1=1;      X2=. ;      X3=3;      Age=27;      Wt=. ;      Ht=68;
```

- `min(of X1-X3) = 1`
- `min(Age, Wt, Ht) = 27`

Explanation:

- 1 is the smallest nonmissing value of the arguments.
- 27 is the smallest nonmissing value.

## Function Name: MAX

What it does: Returns the largest nonmissing value of its arguments. If all the arguments are missing, the result is a missing value.

Arguments: One or more numeric values. If any of the arguments are in the form *Variable1-VariableN*, the list must be preceded by the keyword OF.

Examples:

```
X1=1;      X2=. ;      X3=3;      Age=27;      Wt=. ;      Ht=68;
```

- `max(of X1-X3) = 3`
- `max(Age, Wt, Ht) = 68`

Explanation:

- 3 is the largest value of the arguments.
- The largest value is 68.

## Function Name: SMALLEST

What it does: Returns the *n*th smallest nonmissing value of its arguments. If there is no *n*th nonmissing value, the function returns a missing value.

Arguments: If the first argument is a 1, the function returns the smallest (nonmissing) value in the list of numeric values; if it is a 2, the function returns the second smallest value; and so on. The second to last arguments are numeric values. If any of the arguments are in the form *Variable1-VariableN*, the list must be preceded by the keyword OF.

Examples:

```
X1= 1;      X2=. ;      X3=3;      X4=2;
```

- `smallest(1,of X1-X4) = 1`
- `smallest(2,of X1-X4) = 2`
- `smallest(3,of X1-X4) = 3`
- `smallest(4,of X1-X4) = . (missing value)`

Explanation:

- When the first argument is equal to 1, the result is identical to the MIN function.
- When it is equal to 2, it returns the second (nonmissing) smallest number.
- The third nonmissing smallest value is 3.
- When you are asking for the fourth smallest number, the function returns a missing value because there is no fourth smallest number.

## Function Name: LARGEST

What it does: Returns the  $n$ th largest nonmissing value of its arguments. If there is no  $n$ th nonmissing value, the function returns a missing value.

Arguments: If the first argument is a 1, the function returns the largest (nonmissing) value in the list of numeric values; if it is a 2, the function returns the second largest value; and so on. The second through last arguments are numeric values. If any of the arguments are in the form *Variable1-VariableN*, the list must be preceded by the keyword OF.

Examples:

X1= 1;	X2=. ;	X3=3;	X4=2;
--------	--------	-------	-------

- |  |
|--|
| □ <code>largest(1,of X1-X4) = 3</code>                 |
| □ <code>largest(2,of X1-X4) = 2</code>                 |
| □ <code>largest(3,of X1-X4) = 1</code>                 |
| □ <code>largest(4,of X1-X4) = . (missing value)</code> |

Explanation:

- When the first argument is equal to 1, the result is identical to the MAX function.
- When it is equal to 2, it returns the second (nonmissing) largest number.
- The third largest number is 1.
- When you are asking for the fourth largest number, the function returns a missing value because there is no fourth largest number.

## Programming Example Using the N, NMISS, MAX, LARGEST, and MEAN Functions

The program that follows demonstrates how several of the functions just described can work in concert to produce a very useful result. Here is the problem:

You are given the results of a psychological study. There are 10 variables, Q1 to Q10. You are asked to compute several scores as follows:

1. Score1: the mean of the first five questions (Q1–Q5). Compute this value only if there are three or more nonmissing values.
2. Score2: the mean of Q6–Q10. Compute this mean if there are two or fewer missing values.
3. Score3: the highest score of Q1–Q10.
4. Score4: the sum of the three highest scores of Q1–Q10.

### Program 15.2: Program Demonstrating Several Functions (N, NMISS, MAX, LARGEST, and MEAN)

```
data Score;
  input ID $ Q1-Q10;
  if n(of Q1-Q5) ge 3 then Score1 = mean(of Q1-Q5);
  if nmiss(of Q6-Q10) le 2 then Score2 = mean(of Q6-Q10);
  Score3 = max(of Q1-Q10);
  Score4 = sum(largest(1,of Q1-Q10),
                largest(2,of Q1-Q10),
                largest(3,of Q1-Q10));
datalines;
001 9 7 8 6 7 6 . . 9 2
002 . . . 9 8 7 8 9 9
003 6 7 6 7 6 . . . 9 9
;
title "Listing of Data Set Score";
proc print data=Score noobs;
run;
```

Explanation: The combination of N and MEAN or NMISS and MEAN is very useful for solving problems of this nature. Also, determining the second or third largest value in a list of values is very difficult without using the LARGEST function. Here is the output.

**Figure 15.2: Output from Program 15.2**

Listing of Data Set Score														
ID	Q1	Q2	Q3	Q4	Q5	Q6	Q7	Q8	Q9	Q10	Score1	Score2	Score3	Score4
001	9	7	8	6	7	6	.	.	9	2	7.4	5.66667	9	26
002	.	.	.	.	9	8	7	8	9	9	.	8.20000	9	27
003	6	7	6	7	6	.	.	.	9	9	6.4	.	9	25

There are missing values for Score1 and Score2 because the criteria for computing a score were not met.

## Function Name: INPUT

What it does: Takes the first argument (character value) and “reads” it as if it were being read from text data in a file using an INFORMAT that you supply as the second argument of the function. The most popular use of this function is to perform **character-to-numeric conversion**.

Arguments: The first argument is a character value (typically a character variable). The second argument is the informat that you want to use to associate with the first argument.

Examples:

```
C_Num = '123';      C_Date = '10/21/1950';      C_Money =  
'$12,345.54';
```

```
Num = input(C_num,12.);  Num is the number 123  
Date = input(C_Date,mmddyy10.); Date is a SAS date  
Money = input(C_Money,Dollar12.); Money is a numeric variable  
Explanation: In the first example, notice that the informat (12.) is  
much larger than you need (you only needed 3.). However, unlike  
reading data from a text file, the INPUT function will not read past the  
end of a character value, so there is no harm in choosing a large  
number for the numeric informat.
```

The second example shows how to convert a date, stored as a character string, into a true SAS date.

The last example uses the Dollar12. informat that strips dollar signs

and commas from a value. Note that the Comma. informat also works for this example.

## CALL Routine: CALL SORTN

What it does: After the CALL statement executes, the values of all of the calling arguments are in ascending order. That is, this CALL routine **can sort within an observation**.

Arguments: One or more numeric variables. If any of the arguments are in the form *Variable1-VariableN*, the list must be preceded by the keyword OF.

Examples:

```
X1=7;          X2=3;          X3=9;          X4=. ;          X5=2;
```

```
call sortn(of X1-X5);
```

The values of X1 to X5 are now:

```
X1=.          X2=2          X3=3          X4=7          X5=9
```

```
call sortn(of X5 - X1);
```

The values of X1 to X5 are now:

```
X1=9          X2=7          X3=3          X4=2          x5=.
```

**If you enter the arguments in reverse order, you can perform a descending sort of the values.**

As a practical example, you have 10 test scores (Score1 to Score10) for each student (Stud\_ID). You want to compute the mean of the eight highest scores.

### Program 15.3: Using CALL SORTN to Compute the Mean of the Eight Highest Scores

```
data Test;
  input Stud_ID $ Score1-Score10;
  call sortn(of Score1-Score10);
  Mean_Top_8 = mean(of Score3-Score10);
  datalines;
  001 90 90 80 78 100 95 90 92 88 82
  002 50 55 60 65 70 75 80 85 90 95
```

```

;
title "Listing of Data Set TEST";
proc print data=Test;
  id Stud_ID;
  var Score1-Score10 Mean_Top_8;
run;

```

Here is the listing from Program 15.3.

**Figure 15.3: Output from Program 15.3**

Listing of Data Set TEST											
Stud_ID	Score1	Score2	Score3	Score4	Score5	Score6	Score7	Score8	Score9	Score10	Mean_Top_8
001	78	80	82	88	90	90	90	92	95	100	90.875
002	50	55	60	65	70	75	80	85	90	95	77.500

Explanation: Notice that the Score variables are now in ascending order. If you prefer, you can reverse the order of the arguments in the CALL SORTN routine and then compute the mean as `mean(of Score1-Score8)`.

## Function Name: LAG

What it does: If you execute the LAG function **for every iteration of the DATA step**, it returns the value of its argument from the previous observation.

The true definition of the LAG function is that it returns the value of its argument the last time the function **executed**.

Arguments: A character or numeric variable

Examples:

Because SAS processes one observation at a time, you need special tools to be able to retrieve a value from an earlier observation. The LAG function is one of those tools. For this first example, you have daily stock prices, and you want to compare the current day's price with that of the previous day. Here is the program.

## Program 15.4: Demonstrating the LAG Function

```

data Stocks;
  informat Date mmddyy10.;

```

```

input Date Price;
Up_Down = Price - lag(Price);
format Date mmddyy10.;
datalines;
1/1/2015 100
1/2/2015 98
1/3/2015 96
1/4/2015 101
1/5/2015 101
1/6/2015 104
;

title "Listing of Data Set Stocks";
proc print data=Stocks noobs;
run;

```

Explanation: Because this program is executing the LAG function for every iteration of the DATA step, the variable Up\_Down will be the current day's price minus the price from the previous day. Here is the listing.

**Figure 15.4: Output from Program 15.4**

Listing of Data Set Stocks		
Date	Price	Up_Down
01/01/2015	100	.
01/02/2015	98	-2
01/03/2015	96	-2
01/04/2015	101	5
01/05/2015	101	0
01/06/2015	104	3

The value of Up\_Down is missing in the first observation because that was the first time the LAG function executed and there was no previous value.

There is actually a whole family of LAG functions (LAG, LAG2, LAG3, etc.). LAG2 returns the value of its argument from two previous executions of the function, LAG3 three times, and so on. For example, you can use the Stocks data set to compute a three-day moving average as follows.

## Program 15.5: Using the Family of LAGn Functions to Compute a Moving Average

```
data Moving;
  set Stocks;
  Yesterday = lag(Price);
  Two_Days_Ago = lag2(Price);
  Moving = mean(Price, Yesterday, Two_Days_Ago);
run;

title "Listing of Data Set Moving";
proc print data=Moving noobs;
run;
```

Here is the output.

**Figure 15.5: Output from Program 15.5**

Listing of Data Set Moving					
Date	Price	Up_Down	Yesterday	Two_Days_Ago	Moving
01/01/2015	100	.	.	.	100.000
01/02/2015	98	-2	100	.	99.000
01/03/2015	96	-2	98	100	98.000
01/04/2015	101	5	96	98	98.333
01/05/2015	101	0	101	96	99.333
01/06/2015	104	3	101	101	102.000

Explanation: The variable Moving represents a three-day moving average. You can decide not to output this value until you reach day three—your choice.

## Function Name: DIF

What it does: The value of DIF(X) is equal to X – LAG(X). Because one of the most common uses of the LAG function is to compute inter-observation differences, SAS created the DIF function to save you a few keystrokes when you write your programs.

Arguments: A numeric variable

Examples:

You can substitute the line:

```
Up_Down = dif(Price);
```

for the line that computes the difference of the current day price and the price from the day before in Program 15.4.

## Describing Some Useful Character Functions

The functions discussed in this section all deal with character values. Because there are so many useful character functions, you will see them grouped into logical categories, such as functions that extract substrings, functions that combine strings, and functions that take strings apart.

### Function Names: LENGTHN and LENGTHC

What it does: LENGTHN returns the length of a character value, not counting trailing blanks (blanks to the right of the string). If the argument is a missing value, the function returns a 0. LENGTHC returns the storage length of a character variable.

Arguments: A character value or a character variable

Examples:

### Program 15.6: Demonstrating the Two Functions LENGTHN and LENGTHC

```
data How_Long;
  length String $ 5 Miss $ 4;
  String = 'Abe';
  Miss = ' ';
  Length_String = lengthn(String);
  Store_String = lengthc(String);
  Display = ':' || String || ':';
  Length_Miss = lengthn(Miss);
  Store_Miss = lengthc(Miss);
run;

title "Listing of Data Set HOW_LONG";
proc print data=How_Long noobs;
run;
```

Explanation: This program demonstrates how the two functions LENGTHN and LENGTHC work. Looking at the output (below) should

help clarify the difference between these two functions. The variable Display uses the concatenation operator (||) to place a colon on each side of String. This enables you to see any leading or trailing blanks in the value. Here is the output.

**Figure 15.6: Output from Program 15.6**

Listing of Data Set How_Long						
String	Miss	Length_String	Store_String	Display	Length_Miss	Store_Miss
Abe		3	5	:Abe :	0	4

The variable String is assigned a length of 5, using a LENGTH statement. The LENGTHN function returns a 3, the length of String with the (2) trailing blanks not counted. The storage length, returned by the LENGTHC function, correctly shows that String has a length of 5. The variable Display is a colon, followed by 'Abe', followed by two blanks, followed by a colon.

Character missing values are represented by blanks. The number of blanks is equal to the storage length for that variable. However, when you are either assigning a missing value to a variable (as shown in the assignment statement for Miss) or testing if a character value is missing, you use a single blank in single or double quotation marks (you can also test for, or assign a missing value by using two quotation marks together, with no space—but a single blank is the standard). The LENGTHN function returns a length of 0 for the variable Miss while LENGTHC returns the storage length (4).

For nonmissing character values, an older function called LENGTH is identical to the newer LENGTHN function. However, if the argument is a missing value, the older LENGTH function returns a 1 instead of a 0. The LENGTHN function was introduced with SAS version 9, the same version where strings of zero length were allowed. By the way, the 'N' in the function name LENGTHN stands for *null string*.

## Function Names: TRIMN and STRIP

What it does: TRIMN removes trailing blanks and STRIP removes both leading and trailing blanks. Note that the TRIMN function replaces the older TRIM function. The difference is that the TRIMN function returns a null string (a string of zero length) if its argument is a missing value, and the TRIM function returns a single blank when its argument is a missing value.

Both of these functions are particularly useful when you have character values that might contain leading or trailing blanks. Other functions that search strings for values such as digits or letters (see the ANY and NOT functions later in this chapter) search every position of a character value, including leading or trailing blanks. You usually want to remove leading and trailing blanks before using these searching functions.

Arguments: A character value

Examples:

```
Trail='ABC      ' ;           Lead= '      ABC' ;           Both='      ABC      ' ;
```

Expression	Value
'::'    Trail    '::'	:ABC:
'::'    trimn(Trail)    '	:ABC:

```
'::' || Lead || '::' : ABC:
```

---

```
'::' || strip(Lead) || '::' :ABC:
```

---

```
'::' || trimn(Both) || '::' : ABC:
```

---

```
'::' || strip(Both) || '::' :ABC:
```

---

Explanation: The first example uses the concatenation operator to place a colon on each side of the variable Trail. Therefore, there are blanks between the 'C' and the colon.

In the second example, you use the TRIMN function to remove the trailing blanks before concatenating the final colon. Therefore, there are no blanks between the 'C' and the colon.

The third example concatenates a colon on each side of the variable Lead. Therefore, there are blanks between the first colon and the 'A'.

In the fourth example, the STRIP function removes leading and trailing blanks from the argument. Therefore, there are no blanks between the colons and 'ABC'.

In the fifth example, the TRIMN function removes the trailing blanks from Both so that there are blanks between the first colon and the ‘A’ and no blanks between the ‘C’ and the final colon.

The last example uses the STRIP function to remove leading and trailing blanks from Both, resulting in no blanks between the first colon and the ‘A’ or the ‘C’ and the final colon.

It is unusual for a SAS character value to contain leading blanks because the \$W. informat left-justifies character values. However, if you are not sure if a variable contains leading blanks, use the STRIP function instead of the TRIMN function, just to be sure.

Before we leave these two functions, let’s take a look at the following statements:

```
string = 'ABC      ';
string = trimn(string);
```

What is the value of String? The answer is ‘ABC’. The reason is that although the TRIMN function removed the trailing blanks, when you then assign the trimmed value to a string of length 6, the trailing blanks return.

## Function Names: UPCASE, LOWCASE, and PROPCASE (Functions That Change Case)

What it does: The three functions in this group all change the case of their argument. This is especially useful when you have character data that is entered in different cases (some upper, some lower, some mixed).

Arguments: LOWCASE and UPCASE: A character value.

PROPCASE: A character value and, optionally, a second argument where you list your choice of delimiters.

You can probably guess the purpose of the first two functions, UPCASE and LOWCASE. UPCASE converts all letters in its argument to uppercase, and LOWCASE converts letters to lowercase. PROPCASE (stands for *proper case*) capitalizes the first letter of each “word” and converts the remaining letters to lowercase. The reason that “word” is in quotation marks is that besides the

default blank delimiter, you can specify characters (in addition or instead of blanks) that you want to act as delimiters. The examples that follow will make this clear.

Examples:

```
Name1='rOn Cody';    Name2="D'amore"    String='AbC123xYZ';
```

```
upcase(Name1) = RON CODY
upcase(Name2) = D'AMORE
upcase(String) = ABC123XYZ
lowcase(Name1) = ron cody
lowcase(Name2) = d'amore
lowcase(String) = abc123xyz
propcase(Name1) = Ron Cody
propcase(Name2) = D'amore
propcase(Name2, " ") = D'Amore
```

Explanation: The last two PROPCASE examples need some discussion. You probably want the letter following a single quotation mark in a name to be displayed in uppercase. Because the default delimiter is a blank, using PROPCASE with a single argument results in the value D'amore. By including the second argument containing a blank and a single quotation mark, the result is D'Amore. One final comment: Because you want to include a single quotation mark in the list of delimiters, you need to use double quotation marks to enclose the second argument in this example.

There are some names that do not print properly after you use the PROPCASE function. For example, the name McDonald will become Mcdonald (lowercase 'd').

## Function Name: PUT

What it does: Typically performs numeric-to-character conversion.

Arguments: First argument is a numeric or character value. The second argument is a format (either a built-in SAS format or one that you wrote). The PUT function takes the first argument, formats it using the second argument, and assigns the result to a character value.

Examples:

## Program 15.7: Examples Using the PUT Function

```
proc format;
  value Agegrp low-50='Young'
    51-high='Older'
    . ='Missing'
    other ='Error';
run;

data Put_Eg;
  informat Date mmddyy10.;
  input SS_Num Date Age;
  SS = put(SS_Num,ssn11.);
  Day = put(Date,downname3.);
  Age_Group = put(Age,agegrp.);
  format Date date9.;
datalines;
123456789 10/21/1950 42
890001233 11/12/2015 86
987654321 1/1/2015 15
;
title "Listing of Data Set Put_Eg";
proc print data=Put_Eg noobs;
run;
```

Explanation: This example demonstrates several uses of the PUT function. You first want to create a character variable (SS) from the numeric variable SS\_Num. The built-in SAS format Ssn11. formats a numeric value in Social Security form (i.e., adds leading zeros and hyphens). Next, by using the Downname3. format (Downname formats SAS dates to the names for the days of the week), you create a character variable containing the first three letters for each of the days. Finally, you create a format that groups Age into categories, and then you use the PUT function to create the variable Age\_Group. This variable has values of 'Young' and 'Older' (plus 'Error' and 'Missing' if there are missing values and errors for Age). You would normally drop the variable SS\_Num, keeping only SS, but it was not dropped so that you can see the original values of SS\_Num in the listing.

**Figure 15.7: Output from Program 15.7**

### Listing of Data Set Put\_Eg

Date	SS_Num	Age	SS	Day	Age_Group
21OCT1950	123456789	42	123-45-6789	Sat	Young
12NOV2015	890001233	86	890-00-1233	Thu	Older
01JAN2015	987654321	15	987-65-4321	Thu	Young

## Function Name: SUBSTRN (Newer Version of the SUBSTR Function)

What it does: Extracts a substring from a string.

Arguments: The first argument is a character value (one for which you want to extract a substring). The second argument is the starting position in the character value where you want to begin the substring. The third argument is the length of the substring that you want to extract. This last argument is optional. If you leave it out, the function returns characters from the starting position to the last non-blank character in the string.

Examples:

```
ID='123NJ456';    Amount='$12,456';
```

```
substrn(ID,4,2) = NJ
substrn(ID,4) = NJ456
substrn(Amount,2) = 12,456
substr(ID,4,2) = NJ
```

Explanation: In the first example, you want to extract the state codes from the variable ID. The starting position is 4 and you want to extract two characters. In the second example, you leave off the third argument (the length) and the SUBSTRN function returns all the remaining characters in the string.

The variable Amount starts with a dollar sign, and you want to extract the digits (and commas) without the dollar sign.

Finally, the older SUBSTR function, used in this example, returns the same substring as the SUBSTRN function. There are some more advanced features that are available to you with the SUBSTRN

function compared to the SUBSTR function, but these features are not used in most programs.

Very important point: If you do not define the length of the result of the SUBSTRN function, SAS gives it a default length equal to the length of the first argument (you cannot extract a substring longer than the original string). Every time you create a new variable using the SUBSTRN function, go back to the top of the DATA step and add a LENGTH statement, specifying the length of the substring variable.

## Function Names: FIND and FINDC

What it does: FIND searches a string for a specific string of characters and returns the position in the string where this substring starts. If the substring that you are searching for is not found, FIND returns a 0.

FINDC searches a string for **any one of the characters that you specify**. It returns the position of the first character that it finds. As with the FIND function, if the search fails, the function returns a 0.

Arguments: The first argument in both of these functions is the character value that you are searching. For the FIND function, the second argument is the substring that you are looking for. For the FINDC function, the second argument is a list of individual characters that you are searching for. Both functions allow a third, optional argument (called a *modifier*) that allows these two functions to ignore case when doing their search.

There is also an optional third or fourth argument where you can specify a starting position for the search. If you use a modifier and a starting position, it doesn't matter what order you place these two arguments. If you only use a modifier or a starting position, place it as the third argument. SAS figures out which is a modifier and which is a starting position because modifiers are always character values and starting positions are always numeric values.

Examples:

```
String1='Good bad good';    String2='NNyYnYNN';
```

- `find(String1,'good')` = 10 (position of 'good' - lowercase)
- `find(String1,'good','i')` = 1 (the 'i' modifier says to ignore case)
- `find(String1,'ugly','i')` = 0 (did not find string 'ugly')
- `findc(String2,'Y')` = 4 (position of first uppercase 'Y')
- `findc(String2,'Y','i')` = 3 (the 'i' modifier says to ignore case)
- `findc(String2,'X')` = 0 (no 'X' in string)
- `findc(String1,'abcd')` = 4 (position of 'd')

Explanation:

- FIND returns a 10 because the search is case-sensitive and the string 'good' starts in column 10.
- FIND uses the 'i' (ignore case) modifier, so it finds the string 'Good' in the first position.
- There is no string 'ugly', so the function returns a 0.
- FINDC is searching for an uppercase 'Y' and finds it in position 4.
- The 'i' modifier is used and the lowercase 'y' in position 3 is chosen.
- There is no 'x' in the string, so the function returns a 0.
- You are looking for an 'a', 'b', 'c', or 'd'. The first character it finds in the string is the 'd' in position 4.

Note that these two functions replace the older INDEX and INDEXC functions (with added capability). There is also a third function, FINDW, that searches for words (strings bounded by word boundaries or other delimiters). Because this is not used as much as the other two functions, it was not included. You can find information about FINDW in the documentation.

## Function Names: CAT, CATS, and CATX

What it does: You have already seen the concatenation operator (either `||` or `!!`) earlier in this chapter. Why do you need functions to concatenate strings? There are several advantages that these functions provide, as you will see in the examples that follow. CAT works in a similar manner to the concatenation operator. It takes a list of arguments and puts them together. CATS (pronounced Cat – S by

most) is similar to CAT except that it strips leading and trailing blanks from each of the strings before putting them together (hence, the 'S' in the name). Finally, CATX does everything CATS does with one additional feature—it uses the first argument to the function as a separator between the strings. Note that the first argument to CATX can be more than one character.

Very important point: If you assign the result of a concatenation function to a variable and you do not define a length for this variable, SAS will assign a length of 200. This is different from the rule for the concatenation operator. When you use the operator, the length of the result is the sum of the lengths of each of the strings that you are concatenating.

Arguments: CAT and CATS: One or more **character or numeric** values. If you have variables in the form *Variable1-VariableN*, precede the list with the keyword OF. CATX: The first argument specifies one or more characters to use as separators between the strings. The second through last arguments are identical to the other CAT functions.

Examples:

```
Trail='ABC      ';      Lead= '      ABC';      Both='      ABC      ';
C1='A';      C2='B';      C3='C';      C4='D';      C5='E';
```

```
cat('::', Trail, Lead, '::') = :ABC      ABC:
cats('::', Trail, Lead, '::') = :ABCABC:
catx('-', '::', Trail, Lead, '::') = ABC-ABC:
catx('.', 908, 782, 1323) = 908.782.1323
cats(of C1-C5) = ABCDE
```

Explanation: In the first example, you are concatenating a colon, the two variables Trail and Lead, and another colon. Notice that Trail has three trailing blanks and Lead has three leading blanks. The result has six blanks between the letters.

The second example uses the CATS function. Therefore, there are no spaces between the letters.

In the third example, the first argument is a hyphen, so the resulting string is similar to the CATS example except you now have a hyphen

between the letters.

The next CATX example demonstrates that the arguments to all the CAT functions can be numeric. Remember that the resulting string is a character value.

The last example takes the five variables (C1–C5) and creates a single string.

In practice, the two functions CATS and CATX are used the most. Notice that in the fourth example, the arguments are numeric values. CAT, CATS, and CATX can all take either character or numeric arguments. When you use numeric arguments, an automatic conversion from numeric to character is performed and there are **no conversion messages in the log**.

You will see later how the CATS function and the COUNTC function make for a very powerful duo.

## Function Names: COUNT and COUNTC

What it does: COUNT counts the number of substrings in a character value. COUNTC counts the number of individual characters in a character value.

Arguments: The first argument in both functions is a character value. For the COUNT function, the second argument is the substring that you are searching for. For the COUNTC function, the second argument is a list of individual characters that you want to count. As with the CAT functions, if you have a list of character variables in the form *Variable1-VariableN*, precede the list with the keyword OF. Both functions can take a third, optional argument, the most popular being 'i' that means ignore case. (See the examples.)

Examples:

```
String='Good bad good Good';      String2='AABBxxxxccc';
```

```
count(String,'good') = 1 (the search is case-sensitive)
count(String,'good','i') = 3 ('i' is the ignore case modifier)
countc(String2,'ABC') = 4 (2 A's and 2 B's - the c's are
lowercase)
countc(String2,'ABC','i') = 8 (the 'i' modifier is specified)
```

Explanation: The first example is searching for the string ‘good’ and it only finds it once.

Because the ‘i’ modifier is used in the second example, the function finds three occurrences of ‘good’.

In the first COUNTC example, there are a total of four uppercase As and Bs.

In the last COUNTC example, the ‘i’ modifier allows the four lowercase c’s to be counted as well.

In the example shown next, you have five variables (Q1–Q5), and you want to count the number of uppercase or lowercase Ys in these five variables. The classical approach to solving this problem is to create an array of the five ‘Q’ variables, set a counter to 0, and use a DO loop to test each of the five ‘Q’ variables for an uppercase or lowercase ‘Y’. Each time you find one, you increment the counter.

You can greatly simplify this problem by using the CATS function to place each of the five ‘Q’ values in a single character string. In the first observation, the result of the CATS function is the string ‘YyYnn’. You use this as the argument of the COUNTC function to count the number of uppercase or lowercase Y’s.

### **Program 15.8: Demonstrating the Combination of CATS and COUNTC**

```
data Survey;
  input (Q1-Q5)($1.);
  Number_Y = countc(cats(of Q1-Q5), 'Y', 'i');
datalines;
YyYnn
NNnnn
NYNyy
;
title "Listing of Data Set Survey";
proc print data=Survey noobs;
run;
```

Next time you find yourself creating an array and DO loops, consider if the problem can be solved using the method described in this example.

The output from Program 15.8 is shown below.

**Figure 15.8: Output from Program 15.8**

Listing of Data Set Survey					
Q1	Q2	Q3	Q4	Q5	Number_Y
Y	y	Y	n	n	3
N	N	n	n	n	0
N	Y	N	y	y	3

Imagine, a one-line solution to this problem!

### Function Name: COMPRESS

What it does: Traditionally, it removes characters from a string. With the 'k' modifier, you can use this function to extract characters (for example, all digits) from a string. **This is one of the most powerful character functions in the SAS arsenal—don't skip this section!**

Arguments: If you provide only one argument (a character value), this function removes blanks from the string. An optional second argument is a string of characters that you want to remove from the first argument. If you provide an optional third argument, you can specify character classes to remove from the first argument (such as all letters or all digits).

One of the most useful applications of this function is to use a 'k' (keep) modifier as the third argument. When you do this, the characters you list as the second argument and/or the characters that you specify with modifiers in the third argument are kept and all others are removed. The examples that follow should make this much clearer.

A list of some of the more useful modifiers is shown here:

Modifier	Description
	Removes all blank characters
k	Keeps only the characters listed in the second argument
~	Keeps all characters except those listed in the second argument
-	Keeps all characters except those listed in the second argument, including the '-' character itself
^	Keeps all characters except those listed in the second argument, including the '^' character itself



'a'

All upper- and lowercase letters

---

'd'

All digits

---

'p'

All punctuation (such as periods, commas, etc.)

---

's'

All whitespace characters (spaces, tabs, linefeeds, carriage returns)

---

'i'

Ignore case

---

'k'	Keep the specified characters; remove all others (very useful)
-----	--

---

### Examples:

```
String1='abc def 123';   String2='(908)782-1234';   String3='120  
Lbs.';
```

- **compress(String1)** = abcdef123 (one argument, default action remove blanks)
- **compress(String1,'0123456789')** = abc def (remove digits 0-9)
- **compress(String1,, 'd')** = abc def ('d' modifier means remove all digits)
- **compress(String2,, 'kd')** = 9087821234 (keep the digits)
- **compress(String3,, 'kd')** = 120 (this is a character string, not a number)
- **input(compress(String3,, 'kd'))** = 120 (a numeric value)
- **compress(String1, ' ', 'kadp')** = abc def 123

### Explanation:

- Because there is only one argument, the COMPRESS function removes all blanks.
- The second argument (all the digits 0 to 9) are removed. Back in SAS®8, the COMPRESS function did not have a third (modifier) argument, and this was the way you would need to remove all digits. You might see examples like this if you are looking at programs written before SAS®9.
- It is important to notice that there are two commas in a row following the first argument. With only one comma, the function would think that 'd' is the second argument and you are trying to remove all 'd's from the string. Using two commas tells the function that 'd' is the third argument (modifier) and you want to remove all digits.
- One of my favorites—**keep the digits and throw everything else away.**
- Again, keep the digits. Here it is used to remove units from a value.
- This example shows how to combine the INPUT function and the

COMPRESS function to extract the numeric value when there are units (or other non-digit values) included. The result of the COMPRESS function is a character value, and the INPUT function performs the character-to-numeric conversion.

- If you have character data that includes non-printing ASCII or EBCDIC characters, you can use an expression similar to this. Here you are keeping spaces, uppercase and lowercase letters, digits and punctuation, and throwing everything else away.

To further illustrate the power of the COMPRESS function, the next program shows how to deal with a quantity that contains different units (pounds and kilograms) and wind up with numeric values, all in the same units:

### Program 15.9: Using the COMPRESS Function to Read Data That Includes Units

```
data Weight;
  input Wt $ @@;
  Wt_Kg = input(compress(Wt,, 'kd'), 12.);
  if findc(Wt, 'L', 'i') then Wt_Kg = Wt_Kg / 2.2;
datalines;
120lbs. 90Kg 80Kgs. 200Lb
;
title "Listing of Data Set Weight";
proc print data=Weight noobs;
run;
```

You have weights that include units. Notice that the case of the units is not consistent, and the units might or might not include periods. The double trailing @ in the INPUT statement enables you to read multiple observations on one line of data (it prevents the program from going to a new line each time that the DATA step iterates). You use the COMPRESS function with the 'kd' modifier to extract the digits from each weight and the INPUT function to perform the character-to-numeric conversion. Although you call this Wt\_Kg, it might actually be in pounds. Next, you check to see whether there is an upper- or lowercase 'L' in the units of the Wt variable with the FINDC function. If you find an 'L', you know the value is in pounds, so you divide Wt\_Kg by 2.2 to convert the value to kilograms.

Here is the listing.

**Figure 15.9: Output from Program 15.9**

Listing of Data Set Weight	
Wt	Wt_Kg
120lbs.	54.5455
90Kg	90.0000
80Kgs.	80.0000
200Lb	90.9091

## Function Name: SCAN

What it does: Takes a string apart (parses a string). The most common use is to extract the first and last name from a variable that contains the entire name (and possibly middle initial). The length of the result will be the same as the length of the first argument if you do not use a LENGTH statement to specify the length of the result.

Arguments: The first argument is the string that you want to take apart. The second argument specifies which “word” you want. The reason “word” is in quotation marks is because if you do not list one or more delimiters in the third argument, the default action of the SCAN function is to use blanks plus other characters such as commas, periods, hyphens, etc., as delimiters. In addition, the list of default delimiters is slightly different between ASCII and EBCDIC character sets. **It is a good idea to specify exactly what delimiters you want as the optional, third argument.** If the second argument is negative, the scan proceeds from right to left.

Examples:

```
String1='Alfred E. Newman';   String2='Cody, Ronald';
String3='12-34-56';
```

- `scan(String1,1,' ')` = Alfred (Specifying a blank delimiter)
- `scan(String1,2,' ')` = E. (The second word)
- `scan(String1,3)` = Newman (Using the default delimiters that include blanks)
- `scan(String1,-1,' ')` = Newman (a negative second argument means

```
scan right to left)
□ scan(String1,4) = ' ' (missing value - there is no fourth word)
□ scan(String2,1,', ') = Cody (Specifying blanks and commas as
delimiters)
□ scan(String3,2,'-') = 34 (Specifying a hyphen as the delimiter)
```

Explanation:

- By specifying a blank delimiter, the first word is 'Alfred'.
- The second word is 'E'.
- Because one of the default delimiters is a blank, the third word is 'Newman'. It's probably better to specify a blank as the delimiter if that is the only delimiter you want to use.
- Using a -1 as the second argument is extremely useful in situations where you sometimes have first name and last name and at other times your name variable contains a middle name or initial. Using the -1 causes a right-to-left scan so that you can easily extract the last name in either of these situations.
- Here you are looking for the fourth word in a string that only contains three words, so the function returns a missing value. This provides you with a useful tool for extracting words from strings when you don't know how many words there are. You can extract each word in a DO loop and you know you are finished when the SCAN function returns a missing value.
- Because blanks and commas are both used as delimiters in String2, you can specify both in the third argument. Any combination of commas and blanks will be treated as a single delimiter.
- The SCAN function has many other uses besides extracting words from a string. By specifying the correct delimiter, you can parse other types of data.

When you use the SCAN function, you will probably want to use a LENGTH statement to declare a length for the result so that the result does not have the default length equal to the length of the first argument.

## CALL Routine: CALL MISSING

What it does: Sets all the arguments (character and/or numeric) to a missing value after the CALL statement executes. This routine is especially useful when you are writing a program where you need to initialize a large number of character or numeric variables to a missing value.

Arguments: As many character and/or numeric variables as you want. If you have a variable list in the form *Variable1*-*VariableN*, precede the list by the keyword OF.

Examples:

```
X1=1;  X2=2;  X3=3;  C1='ABC';  C2='D';  C3='Fred';  C4='***';
```

```
call missing(of X1-X3, of C1-C4);
```

After the call, the variables X1–X3 are all numeric missing values and the variables C1–C4 are all character missing values.

## Function Names: NOTDIGIT, NOTALPHA, and NOTALNUM

There are many more NOT functions (NOTSPACE, NOTPUNCT, etc.), but these three are the most useful.

What it does: Determines the first position in a string that is not a digit (NOTDIGIT), an upper- or lowercase letter (NOTALPHA), a letter or digit (NOTALNUM). The functions all return the position of the first character in a string that is not in one of the specified categories. If all the characters in the string fit the category, the functions return a 0. One of the most useful applications of these functions is for data cleaning. **If you have rules concerning allowable characters in a character variable (all digits, for example), you can use the appropriate NOT function to test if there are any characters that violate your condition.**

Arguments: The first argument is a character value that you want to test. There is an optional second argument that is the starting position to begin the search. If the starting position is a negative number, go to the absolute value of the starting position and search from right to left.

Examples:

```
String1='1234NJ';  String2='abc123';  String3='123456';
```

```
notdigit(String1) = 5 (the position of the 'N')
notalpha(String2) = 4 (the position of the '1')
notalnum(String2) = 0 (all characters are alphanumeric)
notdigit(String3) = 0 (all characters are digits)
notalpha(String2, -3) = 0 (start searching at position 3 and search
right to left)
```

Explanation: In the first example, the 'N' is the first non-digit in the string. In the second example, the '1' is the first non-letter in the string. In the third example, you get a 0 because NOTALNUM returns the position of the first character that is not a letter or digit. Because String2 contains only digits and letters, the function returns a 0. String3 contains all digits. Therefore, NOTDIGIT returns a 0. In the last example, a starting position of -3 is an instruction to go to position 3 in String2 and search from right to left. There are only digits from the third position to the beginning of the string, so NOTALPHA returns a 0.

## Function Names: ANYDIGIT, ANYALPHA, and ANYALNUM

What it does: Searches a string for the first occurrence of a digit (ANYDIGIT), an upper- or lowercase letter (ANYALPHA), any digit or letter (ANYALNUM).

Arguments: The first argument is a character value that you want to test. There is an optional second argument that is the starting position to begin the search. If the starting position is a negative number, go to the absolute value of the starting position and search from right to left.

Examples:

```
String1='1234NJ';  String2='abc123';  String3='123456';
```

```
anydigit(String2) = 4 (position of the '1')
anyalpha(String1) = 5 (position of the 'N')
anyalpha(String3) = 0 (no letters in the string)
anyalnum(String2) = 1 (position of the 'a')
```

Explanation: In the first example, the first digit in String2 is the '1' in

position 4. In the second example, the ‘N’ is the first letter and it is in position 5. In the third example, String3 does not contain any letters and the function returns a 0. In the last example, the first character ‘a’ is a letter or digit, so the function returns a 1.

## Function Name: TRANWRD

What it does: Performs a find-and-replace operation. One popular use of this function is to help with address standardization, converting words like “Street” to “St.” and “Road” to “Rd.”

Arguments: The three arguments to the TRANWRD function are: 1) the string that you want to modify, 2) the ‘find’ string, and 3) the ‘replace’ string. Because the ‘replace’ string might be longer than the ‘find’ string, the default length for the result is 200. This fits with the general SAS rule that if the length returned by a function is longer than the length of the string argument, the default length will be 200.

Examples:

```
String1='123 First Street';    String2='Mr. Frank Jones Jr.';
```

```
tranwrd(String1,'Street','St.') = 123 First St.  
tranwrd(String1,'Road','Rd.') = 123 First Street  
tranwrd(String2,'Jr.',' ') = Mr. Frank Jones  
tranwrd(String2,'Mr.',' ') =   Frank Jones (2 leading blanks)
```

Explanation: In the first example, “Street” is replaced by the abbreviation “St.” In the second example, because there is no “Road” in String1, nothing is changed. In the third example, you are substituting a blank for “Jr. “, removing “Jr.” from the string. In the last example, the resulting string has two leading blanks, one from converting “Mr.” to a blank and the other is the original blank that was between “Mr.” and “Frank.”

## Conclusion

Having a knowledge of the SAS functions described in this chapter can save you immense time and effort. Some of the more obscure or complex SAS functions, including Perl regular expressions, were not covered here. Please refer to SAS for help online or one of the SAS

Press books for information about these functions.

## Problems

1. You have a SAS data set called Questionnaire2. It contains variables Subj (subject) and Q1–Q20 (numeric variables). Some of the values of Q1–Q20 might be missing values. Compute the following scores:

Score1 is the mean of Q1–Q10. Only compute Score1 if there are seven or more nonmissing values of Q1–Q10.

Score2 is the median of Q11–Q20. Compute Score2 only if there are five or fewer missing values in Q11–A20.

Score3 is the largest value in variables Q1–Q10.

Score4 is the sum of the two largest values in variables Q1–Q10.

To test your program, run the program below to create the Questionnaire2 data set:

### Program for Problem Sets 9

```
data Questionnaire2;
    input Subj $ Q1-Q20;
    datalines;
001 1 2 3 4 5 1 2 3 4 5 1 2 3 4 5 1 2 3 4 5
002 . . . 3 2 3 1 2 3 4 3 2 3 4 3 5 4 4 4
003 1 2 1 2 1 2 12 3 2 3 . . . . 4 5 5 4
004 1 4 3 4 5 . 4 5 4 3 . 1 1 1 1 1 1 1 1
;
```

2. Modify the program from problem 1 to compute two new variables as follows: Mean\_3\_Large is the mean of the three largest values from Q1 to Q10. Mean\_3\_Low is the mean of the three lowest nonmissing values in variables Q1–Q10. To ensure a good bulletproof solution to this problem, do not compute a mean if there are fewer than three nonmissing values for variables Q1–Q10. That is not the case with the current data set, but it could happen. (Remember that the MEAN function ignores missing values.)
3. Run Program for Problem Sets 10 to create the data set Char\_Data. Create a new SAS data set called Num\_Data that has the three variables Date, Weight, and Height, which are

numeric. Format Date with the DATE9. format. You will need to rename the variables derived from Char\_Date so that those names can be used in the new data set. Hint: The RENAME= option is:

```
data-set-name(rename=(old1=new1 old2=new2, . . .));
```

### **Program for Problem Sets 10**

```
data Char_Data;
  length Date $10 Weight Height $ 3;
  input Date Weight Height;
datalines;
10/21/1966 220 72
5/6/2000 110 63
;
```

4. Working with the data set Questionnaire2 (Problem 1), use CALL SORTN to compute the mean of the 12 highest scores in variables Q1–Q20. It is OK if some of the scores have missing values. You can perform a sort either in ascending order by listing the variables as Q1–Q20 or in descending order by listing the variables as Q20–Q1. Don't forget the keyword OF in the CALL routine.
5. Run Program for Problem Sets 11 to create data set Oscar (which will be used for several of the remaining problems in this section).

### **Program for Problem Sets 11**

```
data Oscar;
  length String $ 10 Name $ 20 Comment $ 25 Address $ 30
    Q1-Q5 $ 1;
  infile datalines dsd dlm=" ";
  *Note: the DSD option is needed to strip the quotes from
  the variables that contain blanks;
  input String Name Comment Address Q1-Q5;
  datalines;
AbC "jane E. MarPle" "Good Bad Bad Good" "25 River Road" y
n N Y Y
12345 "Ron Cody" "Good Bad Ugly" "123 First Street" N n n n
N
98x "Linda Y. d'amore" "No Comment" "1600 Penn Avenue" Y Y
y y y
. "First Middle Last" . "21B Baker St." . . . Y N
; 
```

Using the two length functions, compute the length of String, not counting trailing blanks and the storage length of String. Call these two variables L1 and L2.

6. Modify Program for Problem Sets 11 so that String is in uppercase and Name is in proper case. Use appropriate delimiters so that the name d'amore is spelled D'Amore.
7. Modify Program for Problem Sets 11 to create a new variable called Two\_Three that contains the second and third characters in String. Be sure the length of Two\_Three is 2.
8. Modify Program for Problem Sets 11 to include a new variable Yes\_No. This variable should have a value of "Y" if the variable Comment contains the string "Good" (ignore case) and "N" if Comment does not contain the string "Good". Yes\_No should be a missing value if Comment is a missing value. In addition, create another variable called Count\_Y that represents the total number of Ys in variables Q1–Q5. Hint: Use the CATS function to concatenate Q1–Q5.
9. First, run Program for Problem Sets 12. Modify this program to create a new numeric variable called Height that is the height in inches, computed from the variable Ht. Note: 1 inch = 2.54 cm.

### **Program for Problem Sets 12**

```
Data How_Tall;  
    input Ht $ @@;  
    *Note: the @@ at the end of the INPUT statement allow you  
    to place several observations on one line of data;  
    datalines;  
    65inches 200cm 70In. 220Cm. 72INCHES  
    ;
```

10. Using Program for Problem Sets 11, create a new variable called Last\_Name (length 10) that is the last name computed from the variable Name. Hint: Remember that a minus value representing which "word" you want causes the function to scan from right to left. Also, convert Name to proper case before computing Last\_Name.
11. Using Program for Problem Sets 11, convert the variable Address so that the words "Street" and "Road" are replaced by their abbreviations "St." and "Rd."

# Chapter 16: Working with Multiple Observations per Subject

## Introduction

Because SAS processes one observation at a time, you need special programming tools when you want to analyze data where there are multiple observations per subject or any other grouping variable. For example, suppose you have a data set in which each observation represents a patient visit to a clinic. You record the patient ID, the date of the visit, and some health-related values. Some common programming tasks would include computing differences in values from visit to visit or differences between the first visit and the last visit for each patient. This chapter describes methods for solving problems of this type.

## Useful Tools for Working with Longitudinal Data

Three useful programming tools for working with longitudinal data are **first. and last. variables**, the **LAG and DIF functions**, and **retained variables**. You will use all of these tools in the programs that follow.

Data set Clinic, shown below, is a collection of data on patients visiting a medical clinic. Some patients have only one visit—others have multiple visits. Data for each visit is stored in a separate observation. The following program creates the Clinic data set as well as a listing.

### Program 16.1: Creating the Clinic Data Set

```
data Clinic;
  informat Date mmddyy10. PtNum $3. ;
  input PtNum Date Height Weight Heart_Rate SBP DBP;
  format Date date9. ;
  datalines;
  001 10/21/2015 68 190 68 120 80
```

```

001 11/25/2015 68 195 72 122 84
002 9/1/2015 72 220 76 140 94
003 5/6/2015 63 101 78 118 66
003 7/8/2015 63 106 76 122 70
003 9/1/2015 63 105 77 116 68
;
title "Listing of Data Set Clinic";
proc print data=Clinic;
  id PtNum;
run;

```

Here is the listing.

**Figure 16.1: Output from Program 16.1**

Listing of Data Set Clinic						
PtNum	Date	Height	Weight	Heart_Rate	SBP	DBP
001	21OCT2015	68	190	68	120	80
001	25NOV2015	68	195	72	122	84
002	01SEP2015	72	220	76	140	94
003	06MAY2015	63	101	78	118	66
003	08JUL2015	63	106	76	122	70
003	01SEP2015	63	105	77	116	68

Patient 001 had two visits, patient 002 had one visit, and patient 003 had three visits.

## Describing First. and Last. Variables

With this type of data, one of the first things you want to do is determine when you are processing the first observation and when you are processing the last observation for each patient. For this example, you want to first sort the data set by PtNum and Date. (Note: Sorting by date is not necessary here because visits are in date order for each patient. It is just a precaution in case a visit is out of order.) Next, you create a new data set (let's call it Clinic\_New), as follows.

**Program 16.2: Creating First. and Last. Variables**

```

proc sort data=Clinic;
  by PtNum Date;
run;

```

```
data Clinic_New;
  set Clinic;
  by PtNum;
  file print;
  put PtNum= Date= First.PtNum= Last.PtNum=;
run;
```

The key is to follow the SET statement with a BY statement. Doing this creates two temporary variables, First.PtNum and Last.PtNum. These variables are automatically dropped from the output data set. They only exist as the DATA step is processing. That is the reason that you use a PUT statement to examine these variables. Using PROC PRINT will not work because these two variables are temporary and are not in data set Clinic\_New.

The statement FILE PRINT is an instruction to send the results of the PUT statement to the RESULTS window, not to the SAS log, which is the default location.

Here is the output from Program 16.2.

**Figure 16.2: Output from Program 16.2**

```
PtNum=001 Date=21OCT2015 FIRST.PtNum=1 LAST.PtNum=0
PtNum=001 Date=25NOV2015 FIRST.PtNum=0 LAST.PtNum=1
PtNum=002 Date=01SEP2015 FIRST.PtNum=1 LAST.PtNum=1
PtNum=003 Date=06MAY2015 FIRST.PtNum=1 LAST.PtNum=0
PtNum=003 Date=08JUL2015 FIRST.PtNum=0 LAST.PtNum=0
PtNum=003 Date=01SEP2015 FIRST.PtNum=0 LAST.PtNum=1
```

The variable First.PtNum is equal to 1 for the first visit for each patient and 0 otherwise—the variable Last.PtNum is equal to 1 for the last visit for each patient and 0 otherwise. Let's see how you can use the First. and Last. variables in a program.

## Computing Visit-to-Visit Differences

The first task is to compute visit-by-visit changes in the heart rate and the systolic and diastolic blood pressure for each patient. Here is the program.

## Program 16.3: Computing Visit-to-Visit Differences in Selected Variables

```
proc sort data=Clinic;
  by PtNum Date;
run;

data Diff;
  set Clinic;
  by PtNum;
  if First.PtNum and Last.PtNum then delete;
  Diff_HR = dif(Heart_Rate);
  Diff_SBP = dif(SBP);
  Diff_DBP = dif(DBP);
  if not First.PtNum then output;
run;

title "Listing of Data Set Diff";
proc print data=Diff;
  id PtNum;
run;
```

It makes no sense to compute differences for patients with only one visit. Patients with one visit have both variables, First.PtNum and Last.PtNum, equal to 1 (it is both the first and last visit), so you delete these observations. You use the DIF function to compute differences in heart rate, systolic blood pressure, and diastolic blood pressure for each patient (you can review the DIF function in Chapter 15). For the first visit for each patient (not including the first patient), you are actually computing the difference in each variable based on the **last value from the previous patient and the current value**. This is OK. You are only outputting observations for the second through last visit for each patient. You might be tempted to compute the difference values conditionally (that is, only for the second through last visits). **This does not work.** Remember that the DIF function returns the value of its argument the last time the function executed. If you do not execute the DIF function for the first visit for each patient, when you finally execute it, the function will return the value from the last time the function executed, giving you the last value from the previous patient. If you like, on the first visit for each patient, after you execute the DIF function, you can set the difference variables to a missing value.

Here is the output from Program 16.3.

**Figure 16.3: Output from Program 16.3**

Listing of Data Set Diff										
PtNum	Date	Height	Weight	Heart_Rate	SBP	DBP	Diff_HR	Diff_SBP	Diff_DBP	
001	25NOV2015	68	195	72	122	84	4	2	4	
003	08JUL2015	63	106	76	122	70	-2	4	4	
003	01SEP2015	63	105	77	116	68	1	-6	-2	

Each observation shows differences between the current visit and the previous visit.

## Computing Differences between the First and Last Visits

The next task is to compute the difference in the three variables (Heart\_Rate, SBP, and DBP) from the first visit to the last visit. The primary tool to solve this problem is to use retained variables. As you (hopefully) remember, variables defined in assignment statements are automatically set to a missing value at the top of the DATA step. By using a RETAIN statement, you can instruct the program not to set these variables to a missing value—you are retaining them. This is one way to have the program “remember” a value from some previous observation. If you set a retained variable equal to a value from the first visit, then that retained value will be available when you are processing the last visit, enabling you to compute a difference score. The next program uses this strategy.

### Program 16.4: Computing Differences Between the First Visit and the Last Visit

```
proc sort data=Clinic;
  by PtNum Date;
run;

data First_Last;
  set Clinic;
  retain First_Heart_Rate First_SBP First_DBP; □
  by PtNum;
  if First.PtNum and Last.PtNum then delete; □
```

```

if First.PtNum then do; 
    First_Heart_Rate = Heart_Rate;
    First_SBP = SBP;
    First_DBP = DBP;
End;

if Last.PtNum then do; 
    Diff_HR = Heart_Rate - First_Heart_Rate;
    Diff_SBP = First_SBP - SBP;
    Diff_DBP = First_DBP - DBP;
    output;
end;

run;

title "Listing of Data Set First_Last";
proc print data=First_Last;
    id PtNum;
run;

```

- You use a RETAIN statement for the three variables that will contain the value of heart rate, SBP, and DBP from the first visit.
- Delete patients with only one visit.
- When you are processing the first visit for each patient, set each of the retained variables equal to the value of heart rate, SBP, and DBP, respectively.
- When you reach the last visit for each patient, compute the difference of the current value minus the value from the first visit. Also, output an observation to data set First\_Last.

Here is the listing.

**Figure 16.4: Output from Program 16.4**

Listing of Data Set First_Last													
PtNum	Date	Height	Weight	Heart_Rate	SBP	DBP	First_Heart_Rate	First_SBP	First_DBP	Diff_HR	Diff_SBP	Diff_DBP	
001	25NOV2015	68	195	72	122	84	68	120	80	4	-2	-4	
003	01SEP2015	63	105	77	116	68	78	118	66	-1	2	-2	

Each observation represents the last visit for each patient (who had more than one visit) and the differences between the first visit and the last visit.

## Counting the Number of Visits for Each Patient

Another common task is to count the number of observations (visits, in this example) for each subject or other grouping variable. One of the most straightforward ways to accomplish this is to initialize a counter to 0 when you are processing the first visit for each patient, increment the counter for each visit, and output the patient number and visit count when you are processing

the last visit for each patient. The next program uses this technique to create a data set of patient numbers and visit counts:

### Program 16.5: Counting the Number of Visits for Each Patient

```
proc sort data=Clinic;
  by PtNum;
run;

data Counts;
  set Clinic;
  by PtNum;
  if First.PtNum then N_Visits=0;
  N_Visits + 1;
  if Last.PtNum then output;
run;

title "Listing of Data Set Counts";
proc print data=Counts;
  id PtNum;
run;
```

Because you follow the SET statement with a BY statement, you have the temporary variables First.PtNum and Last.PtNum at your disposal. When you are processing the first visit for each patient, you set the variable N\_Visits equal to 0. Then, for every visit (observation), you use a SUM statement to add one to the counter.

Because it is so important, let's take a minute to review the properties of the SUM statement. If you had used an assignment statement like this:

```
N_Visits = N_Visits + 1;
```

The variable N\_Visits would be set to a missing value for each

iteration of the DATA step. A SUM statement differs from an assignment statement because it is in the form:

*Variable* + *Expression*;

There is no equal sign in this expression. There are three very important properties of a SUM statement:

- First, the *variable* is automatically retained.
- Second, the *variable* is initialized to 0.
- Third, if the *expression* results in a missing value, it is ignored.

The bottom line is that the SUM statement as used here is simply a counter.

When you are processing the last visit for each patient (Last.PtNum is true), you output an observation. This observation contains the heart rate, the systolic and diastolic blood pressure, and the number of visits for each patient. Here is the listing.

**Figure 16.5: Output from Program 16.5**

Listing of Data Set Counts								
PtNum	Date	Height	Weight	Heart_Rate	SBP	DBP	N_Visits	
001	25NOV2015	68	195	72	122	84	2	
002	01SEP2015	72	220	76	140	94	1	
003	01SEP2015	63	105	77	116	68	3	

The variable N\_Visits represents the number of visits for each patient.

## Conclusion

You need special programming tools to analyze data for which you have multiple observations for each value of a BY variable. This chapter described some of these tools.

## Problems

1. First run Program for Problem Sets 13 to create the Clinic data set. Next, write a DATA step that computes the visit-to-visit differences in Heart\_Rate and Weight. Be sure to omit any

patient with only one visit and only output an observation for the second through last visits for each patient. Careful, the data set is not sorted by Subj or Date.

### **Program for Problem Sets 13**

```
data Clinic;
  informat Date mmddyy10. Subj $3. ;
  input Subj Date Heart_Rate Weight;
  format Date date9. ;
datalines;
001 10/1/2015 68 150
003 6/25/2015 75 185
001 12/4/2015 66 148
001 11/5/2015 72 152
002 1/1/2014 75 120
003 4/25/2015 80 200
003 5/25/2015 78 190
003 8/20/2015 70 179
;
```

2. Using the Clinic data set from Problem 1, create a new data set (New\_Clinic) that has one observation per subject. This one observation should include the number of visits for each subject (N\_Visits) and the change in heart rate and weight from the first visit to the last visit. Omit any subject with only one visit.
3. Try to answer this question without running Program for Problem Sets 14 first. What is the value of Last\_x in each of the five observations?

### **Program for Problem Sets 14**

```
data Tricky;
  input x;
  if x gt 5 then Last_x = lag(x);
datalines;
6
7
2
10
11
;
```

4. What's wrong with this program (assume that you have just run Program for Problem Sets 13)?

```
1  data New;
2    set Clinic;
```

```
3      if first.Subj then First_Wt = Weight;
4      if last.Subj then Diff-Wt = Weight - First_Wt;
5 run;
```

# Chapter 17: Describing Arrays

## Introduction

Cody's law of SAS programming states that if you are writing a program and it is becoming very tedious, you should stop what you are doing and ask yourself, "Is there a non-tedious way to write this program?" The answer is often "yes." Perhaps there is a function that will save you time and effort, perhaps a SAS macro. SAS arrays are often the tool to save you huge amounts of time and effort writing your program. This chapter demonstrates the enormous power of SAS arrays.

## What Is an Array?

A SAS *array* (with the exception of a temporary array) is a collection of SAS variables. Using the array name and a subscript, an array element can represent any one of the variables included in the array. The best way to understand how an array works is to look at a program that doesn't use arrays and then see how array processing can make the program much less tedious to write. So, with that in mind, let's look at a program that is just crying out for an array.

You have been given data on a group of subjects that includes age, height, weight, heart rate, systolic blood pressure, and diastolic blood pressure. Missing values were coded as 999 for each of these variables. Because SAS doesn't treat values of 999 as missing values, you need to

convert every value of 999 to a SAS missing value. Here is a program that doesn't use arrays to solve this problem.

### **Program 17.1: Program to Convert 999 to a Missing Value (Without Using Arrays)**

```
data Health_Survey;
```

```

input ID $ Age Height Weight Heart_Rate SBP DBP;
if Age = 999 then Age = .;
if Height = 999 then Height = .;
if Weight = 999 then Weight = .;
if Heart_Rate = 999 then Heart_Rate = .;
if SBP = 999 then SBP = .;
if DBP = 999 then DBP = .;
datalines;
001 23 68 190 68 120 999
002 56 72 220 76 140 88
003 37 999 999 80 132 78
004 82 60 110 80 999 999
;
title "Listing of Data Set Health_Survey";
proc print data=Health_Survey noobs;
run;

```

Here is the output.

**Figure 17.1: Output from Program 17.1**

Listing of Data Set Health_Survey						
ID	Age	Height	Weight	Heart_Rate	SBP	DBP
001	23	68	190	68	120	.
002	56	72	220	76	140	88
003	37	.	.	80	132	78
004	82	60	110	80	.	.

The program works fine, but it is tedious. Imagine if there were 30 or 40 variables that needed processing. The program below uses arrays to save you from writing so many lines of code:

### Program 17.2: Rewriting Program 17.1 Using Arrays

```

data Health_Survey;
input ID $ Age Height Weight Heart_Rate SBP DBP;
array Miss[6] Age Height Weight Heart_Rate SBP DBP;
do i = 1 to 6;
  if Miss[i] = 999 then Miss[i] = .;
end;
drop i;
datalines;
001 23 68 190 68 120 999
002 56 72 220 76 140 88

```

```
003 37 999 999 80 132 78  
004 82 60 110 80 999 999  
;  
  
title "Listing of Data Set Health_Survey";  
proc print data=Health_Survey noobs;  
run;
```

You create an array using an ARRAY statement. Array names follow the same rules as SAS variable names. You follow the array name with a set of brackets. You are free to use either ( ), { }, or [ ] when defining an array. Because SAS functions all use regular parentheses, it is better to use either the curly brackets { } or the square brackets [ ] when defining an array. Most SAS documentation uses the curly variety—this author likes the square brackets. Take your pick.

You place the number of variables included in the array inside the brackets. In later examples, you will see that you can also use an asterisk instead of the number, and SAS will do the counting for you. Following the array name and the brackets, you list all the variables that you want to include in the array. These variables must be either all numeric or all character.

Once you have defined your array, you can use the array name followed by a number in the brackets (referred to as a *subscript*) anywhere in the DATA step. The array reference Miss[1] is the same as writing the variable name Age, the array reference Miss[2] is the same as writing the variable name Height, and so on. **This enables you to use a DO loop to process all of the variables in the array.**

The first iteration of the DO loop is the following statement:

```
if Miss[1] = 999 then Miss[1] = .;
```

which becomes

```
if Age = 999 then Age = .;
```

Once the DO loop has finished, each of the variables in the array have been processed. You don't need (or want) the DO loop counter (i) in the output data set, so you use a DROP statement to remove it

from the data set. The listing of data set Health\_Survey is identical to the output from Program 17.1.

## Describing a Character Array

You can create an array of character variables in much the same way that you did for numeric variables. If the variables in your array have already been defined as character variables elsewhere in the DATA step, you can use the same syntax that you used for numeric arrays. If you are creating the array at a point in the DATA step where the variables have not yet been defined, you can either precede the ARRAY statement with a LENGTH statement or, better yet, define the variables as character as well as defining a storage length right in the ARRAY statement. For example, an ARRAY statement to define 10 character variables, all of length 2, is:

```
array Chars[10] $ 2 Char1-Char10;
```

You follow the array name with a dollar sign and a length, followed by your list of variables.

As an example of a program that uses a character array, the following converts all of the character variables in the data set to uppercase:

### Program 17.3: Converting Character Variables to Uppercase

```
data Uppity;
  informat Name $15. Q1-Q5 $1. ;
  input Name Q1-Q5;
  array Up[6] Name Q1-Q5;
  do i = 1 to 6;
    Up[i] = upcase(Up[i]);
  end;
  drop i;
  datalines;
fred a B c D e
Sue a b c d D
;
title "Listing of Data Set Uppity";
proc print data=Uppity noobs;
run;
```

Each of the variables, Name and Q1–Q5, are now in uppercase.  
(See listing below.)

**Figure 17.2: Output from Program 17.3**

Listing of Data Set Uppity					
Name	Q1	Q2	Q3	Q4	Q5
FRED	A	B	C	D	E
SUE	A	B	C	D	D

## Performing an Operation on Every Numeric Variable in a Data Set

Suppose you want to repeat the process demonstrated in Program 17.1, replacing all values of 999 with a SAS missing value, except you want to do this for every numeric variable in a SAS data set. A very useful strategy is to use the keyword `_NUMERIC_` for your variable list when you define your array. When used in a DATA step, `_NUMERIC_` refers to all the numeric variables **defined up to that point** in the DATA step. You will see how this works in the program that follows. But first, you should know that you can use the keyword `_NUMERIC_` in VAR, KEEP, DROP, and all statements that enable you to specify a list of variables, even in procedures.

Suppose you are given a SAS data set called Big that contains many numeric and character variables. Suppose, further, that the value 999 was used to represent missing values. The following program creates a new data set (Big\_New) where all values of 999 are replaced by a SAS missing value.

### Program 17.4: Converting 999 to a Missing Value for All Numeric Variables in a Data Set

```
data Big_New;
  set Big;
  array All_Nums[*] _numeric_;
  do i = 1 to dim(All_Nums);
    if All_Nums[i] = 999 then All_Nums[i] = .;
  end;
```

```
drop i;  
run;
```

There are some important features in this program. First, the ARRAY statement uses the keyword \_NUMERIC\_. Because the ARRAY statement follows the SET statement, the array All\_Nums is an array of all the numeric variables in data set Big. **If you placed the ARRAY statement before the SET statement, the array would not contain any variables.** To save you the time and effort of counting the number of numeric variables in data set Big, you place an asterisk in the brackets. But what value do you use in your DO loop? The DIM function (stands for *dimension*) takes an array name as its argument and returns the number of variables in the array. You can use the method in this example any time you need to perform an operation on every numeric variable in a SAS data set. If you need to perform an operation on every character variable in a data set, use the keyword \_CHARACTER\_ when you define your array.

## Performing an Operation on Every Character Variable in a Data Set

This section describes a program to convert every character variable in a data set to uppercase:

### Program 17.5: Converting Every Character Variable in a Data Set to Uppercase

```
data Big_New;  
Set Big;  
array All_Chars[*] _character_;  
do i = 1 to dim(All_Chars);  
    All_Chars[i] = upcase(All_Chars[i]);  
end;  
drop i;  
run;
```

When this program executes, all the character variables will be in uppercase. This program and the previous program that operated on all numeric variables will save you immense time and effort

whenever you need to operate on every numeric or character variable in a data set.

## Converting a Data Set with One Observation per Subject into a Data Set with Multiple Observations per Subject

A fairly common programming problem is to transform or restructure a SAS data set. For example, suppose you are given the following SAS data set.

### Program 17.6: Program to Create One Observation per Subject Data Set

```
data Wide;
    input Subj $ Wt1-Wt5;
datalines;
001 120 122 124 123 128
002 200 190 188 180 173
003 115 114 113 110 90
;
title "Listing of Data Set Wide";
proc print data=Wide noobs;
run;
```

Here is the listing.

**Figure 17.3: Output from Program 17.6**

Listing of Data Set Wide					
Subj	Wt1	Wt2	Wt3	Wt4	Wt5
001	120	122	124	123	128
002	200	190	188	180	173
003	115	114	113	110	90

You want to create a data set with five observations per subject like this.

**Figure 17.4: Data Set You Want to Create**

**Listing of data set Thin**

Subj	Time	Weight
001	1	120
001	2	122
001	3	124
001	4	123
001	5	128
002	1	200
002	2	190
002	3	188
002	4	180
002	5	173
003	1	115
003	2	114
003	3	113
003	4	110
003	5	90

There are several ways to go about this. One is to use PROC TRANSPOSE. A straightforward DATA step approach using arrays will also do the trick. Although you can program this task without using arrays, the solution shown here is the standard way of solving this problem. First the program, then the explanation:

### **Program 17.7: Converting a Data Set with One Observation per Subject into a Data Set with Multiple Observations per Subject**

```
data Thin;
  set Wide;
  array Wt[5];
  do Time = 1 to 5;
    Weight = Wt[Time];
    output;
  end;
  drop Wt1-Wt5;
run;

title "Listing of data set Thin";
proc print data=Thin noobs;
run;
```

First, a word (or two) about the ARRAY statement: Notice that there are no variables listed after the array name. When you omit the list of

variables, SAS uses the array name as the base and creates variables using this base and the numbers from 1 to the number of variables in the array. The ARRAY statement in this program is equivalent to:

```
array Wt[5] Wt1-Wt5;
```

Because you want a variable called Time in the new data set, you use that as the variable for your DO loop counter. Inside the loop, you set a new variable that you call Weight equal to each of the Wt variables and output an observation. **It is important that the OUTPUT statement is inside the DO loop.** For each observation coming in from the Wide data set, you are creating five observations in the Thin data set. You no longer need the original Wt variables, so you drop them. The listing of this data was shown in Figure 17.4.

## Converting a Data Set with Multiple Observations per Subject into a Data Set with One Observation per Subject

This section describes a restructuring in the reverse direction—going from a data set with many observations per subject to a data set with one observation per subject. You are probably more likely to go from one to many, as described in the previous section, but the reverse process is sometimes needed. This program is a bit trickier. If you need to do this, be sure to use the program listed here as the basis for your own program. Of course, there is always PROC TRANSPOSE to consider.

You can use the data set Thin, created in the previous section, as the starting point to create one that is identical to the original Wide data set. Once more, first the program, then the explanation:

### Program 17.8: Converting a Data Set with Multiple Observations per Subject into a Data Set with One Observation per Subject

```
proc sort data=Thin;
  by Subj;
run;
```

```
data Wide;
  set Thin;
  by Subj;
  array Wt[5];
  retain Wt1-Wt5;
  if first.Subj then call missing(of Wt1-Wt5);
  Wt[Time] = Weight;
  if last.Subj then output;
  keep Subj Wt1-Wt5;
run;
```

You first sort data set Thin because you plan to follow the SET statement in the DATA step with a BY statement. In the previous chapter, you saw that this process creates two new variables: First.Subj and Last.Subj. You will need both of these variables to make this program work. When you are reading the first weight for each subject, you want to set the variables Wt1–Wt5 to a missing value. You need to do this because of the RETAIN statement. Why do you need to retain these five variables? Without the RETAIN statement, each of the variables Wt1–Wt5 would be set to missing at the top of the DATA step. The RETAIN statement is an instruction not to set these variables to missing. Initializing the variables Wt1–Wt5 is not necessary in this example. However, suppose data set Thin looked like this:

**Figure 17.5: Data Set Thin with a Missing Observation for Subject 002**

**Data Set Thin with a Missing Observation for Subject 002**

Subj	Time	Weight
001	1	120
001	2	122
001	3	124
001	4	123
001	5	128
002	1	200
002	2	190
002	3	188
002	5	173
003	1	115
003	2	114
003	3	113
003	4	110
003	5	90

There is a missing observation for subject 002 at time 4. If you had this situation and did not initialize Wt1–Wt5 to missing for each subject, the missing observation would wind up with the weight at time 4 for the previous subject. Initializing Wt1–Wt5 is just a good precaution.

Back to the program. You assign the value of Weight for each of the array variables. When all the variables Wt1–Wt5 have been given a value, you output a single observation. The resulting data set is identical to the original data set Wide shown previously.

## Conclusion

Arrays can save you time and effort in writing programs. Remember to use the keywords `_NUMERIC_` and `_CHARACTER_` to define arrays when you want to perform an operation on all numeric variables or all character variables in a data set. Finally, you can use arrays to transform or restructure data sets.

## Problems

1. Rewrite Program for Problem Sets 15 using arrays.

## **Program for Problem Sets 15**

```
data Prob1;
  length Char1-Char5 $ 8;
  input x1-x5 Char1-Char5;
  x1 = round(x1);
  x2 = round(x2);
  x3 = round(x3);
  x4 = round(x4);
  x5 = round(x5);
  Char1 = upcase(Char1);
  Char2 = upcase(Char2);
  Char3 = upcase(Char3);
  Char4 = upcase(Char4);
  Char5 = upcase(Char5);
datalines;
1.2 3 4.4 4.9 5 a b c d e
1.01 1.5 1.6 1.7 1.8 frank john mary jane susan
;
title "Listing of Data Set Prob1";
proc print data=Prob1 noobs;
run;
```

2. Starting with the SASHELP data set Fish, create a new, temporary data set (Inches) where each of the variables Height, Width, and Length1–Length3 are converted by dividing their values by 2.54. Use an array to do this.
3. Modify Program for Problem Sets 16 so that all values of 999 are converted to a SAS numeric missing value and all character values of 'NA' (uppercase or lowercase) are converted to a SAS character missing value. Even though there are only four numeric and three character variables, define the arrays with the keywords \_NUMERIC\_ and \_CHARACTER\_. Remember that the DIM function returns the number of variables in an array.

## **Program for Problem Sets 16**

```
data Missing;
  input w x y z c1 $ c2 $ c3 $;
datalines;
999 1 999 3 Fred NA Jane
8 999 10 20 Michelle Mike John
11 9 8 7 NA na Peter
;
```

4. Use the SASHELP data set Heart and create a new, temporary

SAS data set UpHeart that contains all the variables from SASHELP.Heart, but convert all of the character variables to uppercase. Use the keyword \_CHARACTER\_ to define your array and the DIM function to define the upper limit of your DO loop. Because the SASHELP.Heart data set is so large, use the data set option OBS=10 to read only the first 10 observations from this data set. Your SET statement should look like this:

```
set SASHELP.Heart(obs=10);
```

# Chapter 18: Displaying Your Data

## Introduction

The chapter describes how to use PROC PRINT to list all or part of a SAS data set. PROC PRINT can create decent-looking reports that include column labels and summary data. However, if you need a fancier report or want more control over its appearance, you can use PROC REPORT to produce your report. There is a tradeoff: PROC REPORT can produce more sophisticated, custom reports, but it is more complicated to use. As a matter of fact, there are entire books dedicated to PROC REPORT. For many applications, especially with the added tricks described in this chapter, you might find that the reports created by PROC PRINT will do just fine.

## Producing a Simple Report Using PROC PRINT

This first example uses PROC PRINT to list the contents of one of the SASHELP data sets called Shoes. This first program does not use any options or statements—it is “bare bones.”

### Program 18.1: Demonstrating PROC PRINT without any Options or Statements

```
title "Listing of Data Set SHOES";
title2 "In the SASHELP Library";
title3 "-----";
proc print data=SASHELP.Shoes;
run;
```

Three TITLE statements were added. You can include up to 10 TITLE statements in a program. Remember, all the TITLE statements remain in effect until you change them. If you write a new TITLE statement, say TITLE2 following the PROC PRINT statement in Program 18.1, the original second title line will be replaced by the new TITLE2 text, and all title lines higher than two will be deleted.

That is, anytime you submit a new TITLE $n$  statement, all title lines greater than  $n$  are deleted.

By default, PROC PRINT will list all variables and all observations. Also, the order of the variables will be the order that they are stored in the SAS data set. Here is the listing.

**Figure 18.1: Output from Program 18.1 (Partial Listing)**

Listing of Data Set Shoes In the SASHELP Library							
Obs	Region	Product	Subsidiary	Stores	Sales	Inventory	Returns
1	Africa	Boot	Addis Ababa	12	\$29,761	\$191,821	\$769
2	Africa	Men's Casual	Addis Ababa	4	\$67,242	\$118,036	\$2,284
3	Africa	Men's Dress	Addis Ababa	7	\$76,793	\$136,273	\$2,433
4	Africa	Sandal	Addis Ababa	10	\$62,819	\$204,284	\$1,861
5	Africa	Slipper	Addis Ababa	14	\$68,641	\$279,795	\$1,771
6	Africa	Sport Shoe	Addis Ababa	4	\$1,690	\$16,634	\$79
7	Africa	Women's Casual	Addis Ababa	2	\$51,541	\$98,641	\$940
8	Africa	Women's Dress	Addis Ababa	12	\$108,942	\$311,017	\$3,233
9	Africa	Boot	Algiers	21	\$21,297	\$73,737	\$710

The column labeled Obs is created by the procedure and it is an observation number. If you modify your data set (by sorting, adding, or deleting observations, for example), the observation number associated with a particular line of data might change. Most programmers want to replace the Obs column with a variable that identifies the observation, such as an ID or name.

Therefore, the next step is to replace the Obs column with a variable of your choice. You do this by including an ID statement. You can also select which variables to include in the listing by adding a VAR statement. Using a VAR statement also specifies the order of the columns in the report. Here is an example.

### **Program 18.2: Adding ID and VAR Statements to the Procedure**

```
title "Listing of Data Set Shoes";
title2 "In the SASHELP Library";
title3 "-----";
proc print data=SASHELP.Shoes;
```

```

id Region;
var Product Stores Sales Inventory;
run;

```

An ID statement and a VAR statement were added.

The result is Program 18.2 listed below.

Here is the top portion of the revised listing.

**Figure 18.2: Output from Program 18.2 (Partial Listing)**

Listing of Data Set Shoes In the SASHELP Library				
Region	Product	Stores	Sales	Inventory
Africa	Boot	12	\$29,761	\$191,821
Africa	Men's Casual	4	\$67,242	\$118,036
Africa	Men's Dress	7	\$76,793	\$136,273
Africa	Sandal	10	\$62,819	\$204,284
Africa	Slipper	14	\$68,641	\$279,795
Africa	Sport Shoe	4	\$1,690	\$16,634
Africa	Women's Casual	2	\$51,541	\$98,641
Africa	Women's Dress	12	\$108,942	\$311,017
Africa	Boot	21	\$21,207	\$73,737

Region has replaced the Obs columns and only the variables listed in the VAR statement are included in the listing.

You can use a data set option, OBS=, to list the first  $n$  observations in the data set. This is a very useful technique when you only want to see a few observations (from a possibly very big data set) to check that your program is working OK. To see the first eight observations from data set Shoes, you can add the OBS= option to the program, like this.

### **Program 18.3: Using the OBS= Data Set Option to List the First Eight Observations**

```

title "Listing of Data Set Shoes (First 8 Obs)";
title2 "In the SASHELP Library";
title3 "-----";

```

```

proc print data=SASHELP.Shoes(obs=8);
  id Region;
  var Product Stores Sales Inventory;
run;

```

The data set option is placed in parentheses following the data set name. The listing will now show the first eight observations (shown below).

**Figure 18.3: Output from Program 18.3**

Listing of Data Set Shoes (First 8 Obs) In the SASHELP Library				
Region	Product	Stores	Sales	Inventory
Africa	Boot	12	\$29,761	\$191,821
Africa	Men's Casual	4	\$67,242	\$118,036
Africa	Men's Dress	7	\$76,793	\$136,273
Africa	Sandal	10	\$62,819	\$204,284
Africa	Slipper	14	\$68,641	\$279,795
Africa	Sport Shoe	4	\$1,690	\$16,634
Africa	Women's Casual	2	\$51,541	\$98,641
Africa	Women's Dress	12	\$108,942	\$311,017

This listing stops at observation 8. If you want even more control over which observations to list, you can include the two data set options OBS= and FIRSTOBS=. You use FIRSTOBS= to specify the first observation to list and OBS= to specify the last observation to list. For example, to list observations 20 through 25, you would use:

```
proc print data=SASHELP.Shoes(firstobs=20 obs=25);
```

You can use these data set options anywhere you specify a SAS data set name in a program or procedure.

## Using Labels Instead of Variable Names as Column Headings

By default, PROC PRINT uses variable names, not variable labels, as column headings. If you would like variable labels to head your columns, use the LABEL procedure option. The program that follows

creates a small data set (Dimensions) where several variables are given labels. You can then use PROC PRINT with and without the LABEL option to see its effect.

### Program 18.4: PROC PRINT with and without a LABEL Option

```
data Dimensions;
  input Subj $ Ht Wt Waist;
  label Subj = "Subject"
        Ht   = "Height in Inches"
        Wt   = "Subject's Weight";
datalines;
001 68 180 35
002 75 220 40
003 60 101 28
;
title "PROC PRINT Without a LABEL Option";
proc print data=Dimensions;
  id Subj;
  Var Ht Wt Waist;
run;

title "PROC PRINT with LABEL Option";
proc print data=Dimensions label;
  id Subj;
  Var Ht Wt Waist;
run;
```

Before we look at the listings, take a look at the label for the variable Wt. Because you want to include a single quotation mark in the label, you need to enclose the label in double quotation marks. In most places where you need to enclose a string in quotation marks, either single or double quotation marks will work fine. This label is an exception. For more advanced programmers, you should also know that if you want SAS to resolve a macro variable, you need double quotation marks.

Here is the output.

**Figure 18.4: Output from Program 18.4**

### PROC PRINT Without a LABEL Option

Subj	Ht	Wt	Waist
001	68	180	35
002	75	220	40
003	60	101	28

### PROC PRINT with LABEL Option

Subject	Height in Inches	Subject's Weight	Waist
001	68	180	35
002	75	220	40
003	60	101	28

Each listing has a use. The first one (without the LABEL option) is most useful for a programmer, and the second one (with the LABEL option) is most useful for presenting the data to a non-programmer. Note that the variable Waist was not given a label, and its name is used in the second listing.

## Including a BY Variable in a Listing

You can include a BY statement with PROC PRINT to break the output down by that variable. For example, the program to produce the Health data set is reproduced here. To generate a listing broken down by Gender, you must first sort the data set by Gender and then include a BY statement with PROC PRINT. Here is the program:

### Program 18.5: Listing the Health Data Set Broken Down by Gender

```
data Health;
  infile("~/MyBookFiles/health.txt") pad;
  input Subj $ 1-3
        Gender $ 4
        Age   5-6
        HR    7-8
        SBP   9-11
        DBP   12-14
        Chol  15-17;
run;
```

```

proc sort data=Health;
  by Gender;
run;

title "Listing of Data Set Health - by Gender";
proc print data=Health;
  ID Subj;
  by Gender;
run;

```

The output looks like this.

**Figure 18.5: Output from Program 18.5**

Listing of Data Set Health - by Gender					
Gender=F					
Subj	Age	HR	SBP	DBP	Chol
002	55	72	180	90	170
003	18	58	118	72	122
005	34	62	128	80	.
006	38	78	108	68	220

Gender=M					
Subj	Age	HR	SBP	DBP	Chol
001	23	68	120	90	128
004	80	82	.	.	220

You now have separate listings for females and males.

Notice that some of the file names in the folder MyBookFiles are in lowercase and others are in Proper case. Once again, this author forgot that SAS Studio is running in a LINUX environment where file names are case sensitive and used the syntax “~/MyBookFiles/Health.txt” and received the error message “Physical file does not exist” in the log. When you see this error message, the first thing to do is check if you matched the correct file name.

## Including the Number of Observations in a

## Listing

You can include a count of the number of observations in the entire data set or for a BY group by adding the option N='label' to PROC PRINT. To add a count to a listing of the Health data set, you could use the following program.

### Program 18.6: Adding a Count of the Number of Observations to the Listing

```
title "Listing of Data Set Health";
proc print data=health n = "Total Observations =" ;
  ID Subj;
run;
```

By adding the N= option to PROC PRINT, you now see the total number of observations in the listing, as shown next.

**Figure 18.6: Output from Program 18.6**

Listing of Data Set Health						
Subj	Gender	Age	HR	SBP	DBP	Chol
001	M	23	68	120	90	128
002	F	55	72	180	90	170
003	F	18	58	118	72	122
004	M	80	82	.	.	220
005	F	34	62	128	80	.
006	F	38	78	108	68	220
Total Observations =6						

If you use the N= option along with a BY statement, you can see the observation count for each BY group (not shown).

## Conclusion

There are several PROC PRINT options that are not described in this chapter. You can refer to the SAS Studio **Help** menu to see a complete list. However, if you need even more control over the format of your report, you might need to investigate PROC REPORT. I highly recommend the following books:

Carpenter, Art. 2007. *Carpenter's Complete Guide to the PROC REPORT Procedure*. Cary, NC: SAS Institute Inc.

Eslinger, Jane. 2016. *The SAS Programmer's Handbook*. Cary, NC: SAS Institute Inc.

Fine, Lisa. 2013. *PROC REPORT by Example: Techniques for Building Professional Reports Using SAS*. Cary, NC: SAS Institute Inc.

## Problems

1. Use PROC PRINT to list the first 10 observations in the SASHELP data set Fish. Replace the Obs column with the variable Species, and do not include the three variables Length1–Length3. Add three TITLE statements to produce the report heading as follows:

Listing of the First 10 Observations in Data Set Fish

Prepared by: (your name here)

---

2. Prepare a report similar to the one in Problem 1 except break down the report by Species and include all the observations in the data set. Omit the three variables Length1–Length3 and use the PROC PRINT option NOOBS to omit the Obs column.
3. Repeat Problem 2 except use Species as a BY variable and as an ID variable. How does the listing differ from the listing in Problem 2?
4. Create a report from the SASHELP data set Retail. The report should include a title (of your choice) and show sales broken down by year. Include the number of observations for each year, and use variable labels instead of variable names as column headings. The beginning of your report should look like this:

**Listing of Retail Data by Year**

YEAR	Retail sales in millions of \$
1980	\$220
	\$257
	\$258
	\$295
Number=4	

YEAR	Retail sales in millions of \$
1981	\$247
	\$292
	\$286
	\$323
Number=4	

# Chapter 19: Summarizing Data with SAS Procedures

## Introduction

This chapter describes how to summarize numeric data (means, standard deviations, etc.), and how to create SAS data sets containing this summary information. The main tools are PROC MEANS, PROC SUMMARY, and PROC UNIVARIATE.

## Using PROC MEANS (with the Default Options)

PROC MEANS is one of the most useful procedures for summarizing data. As you will see in the programs that follow, you can run this procedure without specifying any options and obtain useful information, including the number of nonmissing observations, the mean, the standard deviation, and the minimum and the maximum values for all the numeric variables in the data set. Later on, in this chapter, you will see how to add procedure options and statements to customize the summary report.

To demonstrate the various ways that you can use PROC MEANS to summarize data, run Program 19.1 to create a data set called Blood\_Pressure. The program uses a random number generator to create a data set containing data from a fictitious drug trial.

The details of how this program works are not described in detail here, but some general information about this program might be interesting. First, remember that the library Oscar was created in a statement that you added to the Autoexec.sas file. The RAND function is used to generate random numbers that are normally distributed. The way the program is designed, if you run it yourself, you will obtain exactly the same data set used in the examples that follow. The program also generates a subject number, drug group,

and gender, as well as the systolic and diastolic blood pressure (adjusted so that the drug group pressures are lower than the placebo group pressures). The program also outputs missing values for some of the variables, to make the data set more realistic. Here is the program.

### Program 19.1: Creating the Blood\_Pressure Data Set

```
Libname Oscar "~/MyBookFiles";
/* Not necessary if already in your Autoexec.sas file */

data Oscar.Blood_Pressure;
  call streaminit(37373);
  do Drug = 'Placebo', 'Drug A', 'Drug B';
    do i = 1 to 20;
      Subj + 1;
      if mod(Subj,2) then Gender = 'M';
      else Gender = 'F';
      SBP = rand('normal',130,10) +
        7*(Drug eq 'Placebo') - 6*(Drug eq 'Drug B');
      SBP = round(SBP,2);
      DBP = rand('normal',80,5) +
        3*(Drug eq 'Placebo') - 2*(Drug eq 'Drug B');
      DBP = round(DBP,2);
      Age = int(rand('normal',50,10) + .1*SBP);
      if Subj in (5,15,25,55) then call missing(SBP, DBP);
      if Subj in (4,18) then call missing(Gender);
      output;
    end;
  end;
  drop i;
run;

title "Listing of Data Set Drug_Trial (first 10 observations)";
proc print data=Oscar.Blood_Pressure(obs=10);
  id Subj;
run;
```

The first 10 observations of this data set are listed next.

### Figure 19.1: Output from Program 19.1

### **Listing of Data Set Drug\_Trial (first 10 observations)**

Subj	Drug	Gender	SBP	DBP	Age
1	Placebo	M	138	86	50
2	Placebo	F	134	90	40
3	Placebo	M	136	84	61
4	Placebo		132	80	65
5	Placebo	M	.	.	71
6	Placebo	F	146	82	60
7	Placebo	M	138	88	70
8	Placebo	F	140	80	74
9	Placebo	M	134	80	59
10	Placebo	F	152	90	62

You are now ready to run PROC MEANS (without specifying any options) on this data set.

### **Program 19.2: Running PROC MEANS on the Blood\_Pressure Data Set (Using All the Default Options)**

```
title "Running PROC MEANS with all the Defaults";
proc means data=Oscar.Blood_Pressure;
run;
```

After running this program, the RESULTS window looks like this:

**Figure 19.2: Output from Program 19.2**

Running PROC MEANS with all the Defaults					
Variable	N	Mean	Std Dev	Minimum	Maximum
Subj	60	30.5000000	17.4642492	1.0000000	60.0000000
SBP	56	129.6785714	11.0389511	100.0000000	152.0000000
DBP	56	81.3571429	5.0755334	72.0000000	92.0000000
Age	60	64.0166667	10.3063115	40.0000000	88.0000000

Because you did not include a VAR statement, PROC MEANS summarized every numeric variable in the Blood\_Pressure data set (including the subject number). Let's see how to produce a more meaningful report.

## Using PROC MEANS Options to Customize the Summary Report

Most of the time, you will want to override the default output and select the statistics that you want in your report. The following table shows some of the most popular statistics that you can request using PROC MEANS:

**Table 19.1: Table of PROC MEANS Options**

Option	Description
N	The number of nonmissing observations used to compute the statistics
NMISS	The number of missing observations
MEAN	The mean

**STD** The standard deviation

---

**CV** The coefficient of variation

---

**CLM** The 95% confidence interval for the mean

---

**STDERR** The standard error

---

**MIN** The minimum value

---

**MAX** The maximum value

---

MEDIAN

The median

---

MAXDEC=*n*

The maximum number of decimal places in all the table values

---

You will probably want to include the MAXDEC= option every time you run PROC MEANS, whether you override the default options or not. This option controls the number of digits to the right of the decimal place for most of the statistics in the output (values for N and NMISS are always integers).

In most cases, you will also want to use the N and NMISS options. These two options add the number of nonmissing and missing values in the output table.

What other options you select will vary depending on what information you need for your particular project. The next program uses MAXDEC, N, and NMISS along with options to print the mean, the standard deviation, the coefficient of variation, the standard error, and the 95% confidence interval for the mean. These statistics are commonly used in reports and journal articles. Here is the program, followed by the output.

### Program 19.3: Adding Options to Control PROC MEANS Output

```
title "Running PROC MEANS with Selected Options";
proc means data=Oscar.Blood_Pressure n nmiss mean std cv stderr
  clm      maxdec=3;
  var SBP DBP;
run;
```

A VAR statement was also included in this program to request statistics only for the two variables SBP and DBP (systolic and diastolic blood pressure).

Most of the time, you will want to include a VAR statement when you use PROC MEANS because there might be several numeric variables in your data set for which it makes no sense to compute means and other statistics (such as the subject variable in this data set).

Here is the output.

**Figure 19.3: Output from Program 19.3**

Running PROC MEANS with Selected Options								
Variable	N	N Miss	Mean	Std Dev	Coeff of Variation	Std Error	Lower 95% CL for Mean	Upper 95% CL for Mean
SBP	56	4	129.679	11.039	8.513	1.475	126.722	132.635
DBP	56	4	81.357	5.076	6.239	0.678	79.998	82.716

Once you specify any PROC MEANS options (with the exception of MAXDEC=), the output will contain statistics only for the options that you specify. Notice that the minimum and maximum values, part of the default set of statistics, are not present in this output.

This summary table displays statistics for SBP and DBP for all subjects in the trial. The next step is to see a breakdown of these statistics by Gender and Drug.

## Computing Statistics for Each Value of a BY Variable

It is often useful to compute statistics for each value of some other variable. For example, one of the variables in the Blood\_Pressure data set is Gender. You might want to see selected statistics for males and females. One way to do this is to sort the data set by Gender and then include a BY statement with PROC MEANS.

The next program does just this.

### Program 19.4: Adding a BY Statement with PROC MEANS

```
proc sort data=Oscar.Blood_Pressure out=Blood_Pressure;
  by Gender;
run;
```

```

title "Statistics for Blood Pressure Study Broken Down by
Gender";
proc means data=Blood_Pressure n nmiss mean std maxdec=2;
  by Gender;
  var SBP DBP;
run;

```

You add an OUT= PROC SORT option so that the data in the original (permanent) data set is not affected. You can use the same data set name (Blood\_Pressure) for the output data set in the Work library as you used in the permanent library.

Most of the time, when you run PROC SORT, you should specify an OUT= procedure option, especially if you are subsetting either observations or variables. This prevents you from damaging the original data set.

Because you added the BY statement to PROC MEANS, you now have your statistics broken down by Gender. Here is a portion of the output.

**Figure 19.4: Output from Program 19.4**

Statistics for Blood Pressure Study Broken Down by Gender				
Gender=' '				
Variable	N	N Miss	Mean	Std Dev
SBP	2	0	133.00	1.41
DBP	2	0	84.00	5.66
Gender=F				
Variable	N	N Miss	Mean	Std Dev
SBP	28	0	130.86	11.73
DBP	28	0	81.29	5.56
Gender=M				
Variable	N	N Miss	Mean	Std Dev
SBP	26	4	128.15	10.70
DBP	26	4	81.23	4.63

You now see the N, NMISS, mean, and standard deviation for males

and females (as well as the two observations where Gender was missing).

If you want to omit the two observations where Gender is missing, add a WHERE= data set option as part of your PROC SORT, like this:

```
proc sort data=Oscar.Blood_Pressure(where=(Gender is not missing))
          out=Blood_Pressure;
```

The output data set, Blood\_Pressure, is now sorted by Gender and does not include any observations with missing values for Gender. This also shows why it is usually advantageous to include an OUT= option with PROC SORT. Without it, you would have removed two observations from your permanent Blood\_Pressure data set.

## Using a CLASS Statement Instead of a BY Statement

You can use a CLASS statement to generate the same information that you produced with a BY statement. One major difference between using a BY statement versus a CLASS statement is that **you do not have to sort the data when you use a CLASS statement**.

If you have very large data sets and several CLASS variables, there is a possibility that the program will run out of memory, and you will have to use PROC SORT and a BY statement instead of a CLASS statement. Using cloud based programs, this situation is highly unlikely.

The output that you obtain when you use a CLASS statement instead of a BY statement has a slightly different, more compact format, but the numbers are the same. The next program shows how to use a CLASS statement:

### Program 19.5: Using a CLASS Statement to See Statistics Broken Down by Region

```
title "Using a CLASS Statement with PROC MEANS";
proc means data=Oscar.Blood_Pressure n nmiss mean std maxdec=2;
  class Gender;
  var SBP DBP;
run;
```

Notice that this program is using the original permanent data set and no sorting is necessary. Here is the result:

**Figure 19.5: Output from Program 19.5**

Using a CLASS Statement with PROC MEANS						
Gender	N Obs	Variable	N	N Miss	Mean	Std Dev
F	28	SBP	28	0	130.86	11.73
		DBP	28	0	81.29	5.56
M	30	SBP	26	4	128.15	10.70
		DBP	26	4	81.23	4.63

Using a CLASS statement produces a slightly more compact and easy-to-read report. Also, it does not include the observations where Gender is missing.

## Including Multiple CLASS Variables with PROC MEANS

Because this was a drug study, you will want to see statistics on SBP and DBP broken down by Gender and Drug. This is easily accomplished by listing these two variables in the CLASS statement. Here is the program and the output.

**Program 19.6: Using Two CLASS Variables with PROC MEANS**

```
title "Using a CLASS Statement with Two CLASS Variables";
proc means data=Oscar.Blood_Pressure n nmiss mean std maxdec=2;
  class Gender Drug;
  var SBP DBP;
run;
```

**Figure 19.6: Output from Program 19.6**

### Using a CLASS Statement with Two CLASS Variables

Gender	Drug	N Obs	Variable	N	N Miss	Mean	Std Dev
F	Drug A	10	SBP	10	0	130.60	12.86
			DBP	10	0	80.20	6.43
	Drug B	10	SBP	10	0	123.00	5.10
M			DBP	10	0	79.80	3.94
	Placebo	8	SBP	8	0	141.00	8.88
			DBP	8	0	84.50	5.42
M	Drug A	10	SBP	9	1	131.33	9.22
			DBP	9	1	80.44	4.22
	Drug B	10	SBP	9	1	119.33	8.12
			DBP	9	1	78.89	4.14
	Placebo	10	SBP	8	2	134.50	8.80
			DBP	8	2	84.75	3.85

You now see blood pressures broken down by Gender and Drug.

## Statistics Broken Down Every Way

You can add the PROC MEANS option PRINTALLTYPES to output statistics broken down by every combination of the CLASS variables. To show how this works, here is Program 19.6 rewritten with this option included.

### Program 19.7: Adding the PRINTALLTYPES Option to PROC MEANS

```
proc means data=Oscar.Blood_Pressure n nmiss mean std maxdec=2
  printalltypes;
  class Gender Drug;
  var SBP DBP;
run;
```

The output now looks like this.

**Figure 19.7: Output from Program 19.7**

N Obs	Variable	N	N Miss	Mean	Std Dev
58	SBP	54	4	129.56	11.22
	DBP	54	4	81.26	5.08

Drug	N Obs	Variable	N	N Miss	Mean	Std Dev
Drug A	20	SBP	19	1	130.95	10.98
		DBP	19	1	80.32	5.34
Drug B	20	SBP	19	1	121.26	6.77
		DBP	19	1	79.37	3.95
Placebo	18	SBP	16	2	137.75	9.18
		DBP	16	2	84.63	4.54

Gender	N Obs	Variable	N	N Miss	Mean	Std Dev
F	28	SBP	28	0	130.86	11.73
		DBP	28	0	81.29	5.56
M	30	SBP	26	4	128.15	10.70
		DBP	26	4	81.23	4.63

Gender	Drug	N Obs	Variable	N	N Miss	Mean	Std Dev
F	Drug A	10	SBP	10	0	130.60	12.86
			DBP	10	0	80.20	6.43
	Drug B	10	SBP	10	0	123.00	5.10
M	Placebo	8	SBP	8	0	141.00	8.88
			DBP	8	0	84.50	5.42
	Drug A	10	SBP	9	1	131.33	9.22
	Drug B	10	SBP	9	1	119.33	8.12
			DBP	9	1	78.89	4.14
	Placebo	10	SBP	8	2	134.50	8.80
			DBP	8	2	84.75	3.85

You now see statistics for every combination of the CLASS variables. This is a very useful option and you should always consider using it when you have one or more CLASS variables.

## Using PROC MEANS to Create a Summary Data Set

Besides producing printed output, you can use PROC MEANS to create a data set containing the same data that was in the printed reports. Let's start out by computing the mean SBP and DBP for the entire data set. To do this, you add an OUTPUT statement. On this

statement, you name the output data set and specify what statistics you want in that data set. For this first example, you will name the data set Grand\_Mean and request that the mean and the number of nonmissing observations be included in the data set. Here is the program.

### Program 19.8: Using PROC MEANS to Create a Data Set Containing the Grand Mean

```
proc means data=Oscar.Blood_Pressure noprint;  
  var SBP DBP;  
  output out=Grand_Mean mean=Grand_SBP Grand_DBP  
        n=Nonmiss_SBP Nonmiss_DBP;  
run;
```

Because you want to create a summary data set and do not want to print a report, you use the PROC MEANS option NOPRINT. This option suppresses any printed output from the procedure.

PROC SUMMARY produces data sets identical to those produced by PROC MEANS. The only difference between the two procedures is that PROC SUMMARY automatically includes a NOPRINT option. So, use PROC MEANS with the NOPRINT option or PROC SUMMARY—your choice.

You use an OUTPUT statement where you specify the name of the output data set (with the keyword OUT=), and you use keywords to specify which statistics you want in the output data set. These keywords are the same ones that you used as PROC MEANS options (see the keywords for PROC MEANS earlier in this chapter) followed by an equal sign, followed by the variable names that you want to use. These variable names are in the same order as the variable names on the VAR statement. The variable Grand\_SBP is the mean SBP for all observations in the data set; the variable Grand\_DBP is the mean DBP for all observations in the data set.

Even if you are an experienced SAS programmer, it is a good idea to use PROC PRINT to list the contents of the summary data set. Program 19.9 does this.

### Program 19.9: Listing of Data Set Grand\_Mean

```
title "Listing of Data Set Grand_Mean";
proc print data=Grand_Mean noobs;
run;
```

Here is the listing.

**Figure 19.8: Output from Program 19.9**

Listing of Data Set Grand_Mean					
_TYPE_	_FREQ_	Grand_SBP	Grand_DBP	Nonmiss_SBP	Nonmiss_DBP
0	60	129.679	81.3571	56	56

This data set contains the four variables that you specified in the OUTPUT statement (two for SBP and two for DBP). SAS has added two additional variables to this data set: \_TYPE\_ and \_FREQ\_. The variable \_TYPE\_ is useful when you have one or more CLASS variables. The variable \_FREQ\_ is the total number of observations available for use in the data set. There were 60 observations in data set Blood\_Pressure, but because there were four missing values for SBP and DBP, both variables Nonmiss\_SBP and Nonmiss\_DBP are equal to 56.

## Letting PROC MEANS Name the Variables in the Output Data Set

A nice option that you can use with an OUTPUT statement is called AUTONAME. When you include this option, you do not have to name any of the variables in the output data set—PROC MEANS will name them for you. It uses the names on the VAR statement and adds an underscore, followed by the statistic that you requested. For example, if you want the mean for SBP and DBP, the output data set will use the variable names SBP\_Mean and DBP\_Mean. Here is Program 19.8 rewritten to use the AUTONAME option.

### Program 19.10: Using AUTONAME to Name the Variables in the Output Data Set

```
proc means data=Oscar.Blood_Pressure noprint;
```

```

var SBP DBP;
output out=Grand_Mean mean= n= / autoname;
run;

title "Listing of Data Set Grand_Mean";
title2 "Using the AUTONAME Output Option";
proc print data=Grand_Mean noobs;
run;

```

You use the keywords for the desired statistics, followed by an equal sign and no variable names. Because AUTONAME is a statement option, it follows a slash in the OUTPUT statement. Here is a listing of data set Grand\_Mean where the AUTONAME option was used.

**Figure 19.9: Output from Program 19.10**

**Listing of Data Set Grand\_Mean  
Using the AUTONAME Output Option**

_TYPE_	_FREQ_	SBP_Mean	DBP_Mean	SBP_N	DBP_N
0	60	129.679	81.3571	56	56

This author recommends using AUTONAME for two reasons: One, it saves entering and creating variable names for all your statistics; and two, it creates consistent and predictable variable names for all of your statistics.

## Creating a Summary Data Set with CLASS Variables

You might want to create a summary data set that contains statistics broken down by one or more variables. One way to do this is to add a CLASS statement to PROC MEANS, along with an OUTPUT statement. This produces an especially useful data set. To demonstrate this process, let's output several statistics for the variables SBP and DBP and include both Gender and Drug as CLASS variables. Here is the program.

**Program 19.11: Creating a Summary Data Set with Two CLASS Variables**

```

proc means data=Oscar.Blood_Pressure noprint;
  class Gender Drug;
  var SBP DBP;
  output out=Summary mean= n= std= / autoname;
run;

title "Listing of Data Set Summary";
proc print data=Summary noobs;
run;

```

You are requesting values for the mean, the number of nonmissing values, and the standard deviation for the variables SBP and DBP, broken down by Gender and Drug. Here is a listing of the data set.

**Figure 19.10: Output from Program 19.11**

Listing of Data Set Summary									
Gender	Drug	_TYPE_	_FREQ_	SBP_Mean	DBP_Mean	SBP_N	DBP_N	SBP_StdDev	DBP_StdDev
		0	58	129.556	81.2593	54	54	11.2244	5.08488
	Drug A	1	20	130.947	80.3158	19	19	10.9822	5.34429
	Drug B	1	20	121.263	79.3684	19	19	6.7728	3.94702
	Placebo	1	18	137.750	84.6250	16	16	9.1761	4.54423
F		2	28	130.857	81.2857	28	28	11.7307	5.55683
M		2	30	128.154	81.2308	26	26	10.7021	4.63299
F	Drug A	3	10	130.600	80.2000	10	10	12.8599	6.42564
F	Drug B	3	10	123.000	79.8000	10	10	5.0990	3.93841
F	Placebo	3	8	141.000	84.5000	8	8	8.8802	5.42481
M	Drug A	3	10	131.333	80.4444	9	9	9.2195	4.21637
M	Drug B	3	10	119.333	78.8889	9	9	8.1240	4.13656
M	Placebo	3	10	134.500	84.7500	8	8	8.7994	3.84522

This might not be what you expected. The first observation (\_TYPE\_ = 0) is the mean (along with the other requested statistics) for all the observations. The next three observations are statistics broken down by Drug; the next two observations are statistics broken down by Gender; and, finally, the last six observations show the statistics for every combination of Drug and Gender. Notice that the suffix added for the two variables holding the standard deviations is StdDev, not Std (the option that you used). The reason is that Std is an abbreviation for the actual term StdDev.

You probably want only the last six observations in this data set. The easiest way to do this is to add the PROC MEANS option NWAY. This option provides your requested statistics broken down by all of your CLASS variables. Here is Program 19.11 rewritten to include the NWAY option.

### Program 19.12: Adding the NWAY Option to PROC MEANS

```
proc means data=Oscar.Blood_Pressure noprint nway;  
  class Gender Drug;  
  var SBP DBP;  
  output out=Summary mean= n= std= / autoname;  
run;  
  
title "Listing of Data Set Summary";  
title2 "NWAY Option Added";  
proc print data=Summary noobs;  
run;
```

Here is the listing of the Summary data set.

**Figure 19.11: Output from Program 19.12**

Listing of Data Set Summary NWAY Option Added									
Gender	Drug	_TYPE_	_FREQ_	SBP_Mean	DBP_Mean	SBP_N	DBP_N	SBP_StdDev	DBP_StdDev
F	Drug A	3	10	130.600	80.2000	10	10	12.8599	6.42564
F	Drug B	3	10	123.000	79.8000	10	10	5.0990	3.93841
F	Placebo	3	8	141.000	84.5000	8	8	8.8802	5.42481
M	Drug A	3	10	131.333	80.4444	9	9	9.2195	4.21637
M	Drug B	3	10	119.333	78.8889	9	9	8.1240	4.13656
M	Placebo	3	10	134.500	84.7500	8	8	8.7994	3.84522

You now have the data set you want.

If you sorted the Blood\_Pressure data set by Gender and Drug and used a BY statement instead of a CLASS statement, you would not need the NWAY option. The data set would be identical to the one above.

It is very important to remember the NWAY option when all you want are the statistics broken down by all of the CLASS variables. This author uses the mental “trick” of pairing NOPRINT and NWAY

together in his mind.

## Using a Formatted CLASS Variable

The Blood\_Pressure data set also contains the age of each subject in the study. You might want to see the mean systolic and diastolic blood pressures for two or more age groups. Conveniently, if you use a continuous variable such as Age in a CLASS statement and you also include a FORMAT statement, **PROC MEANS will use the formatted values of the CLASS variables to compute statistics**. Here is an example:

You want to see the mean and standard deviation of SBP and DBP for three age groups:

- Low–50
- 51–70
- 71–High

All you need to do is create a format and include CLASS and FORMAT statements when you run PROC MEANS. The next program demonstrates this.

### Program 19.13: Using a Formatted CLASS Variable

```
proc format;
  value AgeGroup low-50  = '50 and Lower'
        51-70   = '51 to 70'
        71-high = '71 +';
run;

title "Using a Formatted CLASS Variable";
proc means data=Oscar.Blood_Pressure n nmiss mean std maxdec=2;
  class Age;
  format Age AgeGroup.;
  var SBP DBP;
run;
```

Here is the output.

### Figure 19.12: Output from Program 19.13

### Using a Formatted CLASS Variable

Age	N Obs	Variable	N	N Miss	Mean	Std Dev
50 and Lower	6	SBP	6	0	130.67	7.66
		DBP	6	0	80.33	6.25
51 to 70	40	SBP	38	2	130.16	11.01
		DBP	38	2	81.79	5.24
71 +	14	SBP	12	2	127.67	13.01
		DBP	12	2	80.50	4.10

Using formatted CLASS variables enables you to see all of your statistics broken down by any grouping of continuous variables. It saves time and effort.

## Demonstrating PROC UNIVARIATE

One of the most popular procedures for summarizing data, especially for statistical purposes, is PROC UNIVARIATE. This procedure has a lot in common with PROC MEANS. However, you can use statements to produce histograms and probability plots that are not available with PROC MEANS.

To demonstrate this procedure, the next program uses PROC UNIVARIATE to analyze SBP and produce a histogram and Q-Q plot.

### Program 19.14: Demonstrating PROC UNIVARIATE

```
title "Demonstrating PROC UNIVARIATE";
proc univariate data=Oscar.Blood_Pressure;
  id Subj;
  var SBP;
  histogram;
  qqplot / normal (mu=est sigma=est);
run;
```

If you have a variable such as Subj or ID, be sure to include an ID statement including that variable. It will be useful in some of the output. You use a VAR statement just as you did with PROC MEANS. Two additional statements, HISTOGRAM and QQPLOT, were added. HISTOGRAM, as the name implies, generates a histogram for all the variables on the VAR statement. QQPLOT requests a quantile-

quantile plot, that is used by statisticians to help determine deviations from normality. When you request a Q-Q plot, you can add the option NORMAL. This option draws a straight line representing what a normal distribution would look like on the plot. Following this option, you can specify a mean (mu) and a standard deviation (sigma) for your theoretical normal plot. In this example, you want to use the data to estimate these two values. You use the term EST (stands for *estimated*) to request this.

Here is the result.

**Figure 19.13: Output from Program 19.14**

Demonstrating PROC UNIVARIATE			
Variable: SBP			
Moments			
N	56	Sum Weights	56
Mean	129.678571	Sum Observations	7262
Std Deviation	11.0389511	Variance	121.858442
Skewness	0.02568868	Kurtosis	-0.2313177
Uncorrected SS	948428	Corrected SS	6702.21429
Coeff Variation	8.51254836	Std Error Mean	1.47514189

This first section shows you the mean, the standard deviation, and several other statistics. For those who are interested, the skewness and kurtosis are values that help determine whether the data values are normally distributed. For both of these statistics, values close to 0 indicate a distribution close to normal.

Basic Statistical Measures			
Location		Variability	
Mean	129.6786	Std Deviation	11.03895
Median	128.0000	Variance	121.85844
Mode	120.0000	Range	52.00000
		Interquartile Range	16.00000

Here you see other measures of central tendency, including the mean and median, along with measures of spread or dispersion.

Tests for Location: Mu0=0				
Test	Statistic		p Value	
Student's t	t	87.90922	Pr >  t	<.0001
Sign	M	28	Pr >=  M	<.0001
Signed Rank	S	798	Pr >=  S	<.0001

Here you see statistical tests to test the null hypothesis that the mean is 0.

Quantiles (Definition 5)	
Level	Quantile
100% Max	152
99%	152
95%	150
90%	146
75% Q3	138
50% Median	128
25% Q1	122
10%	118
5%	114
1%	100
0% Min	100

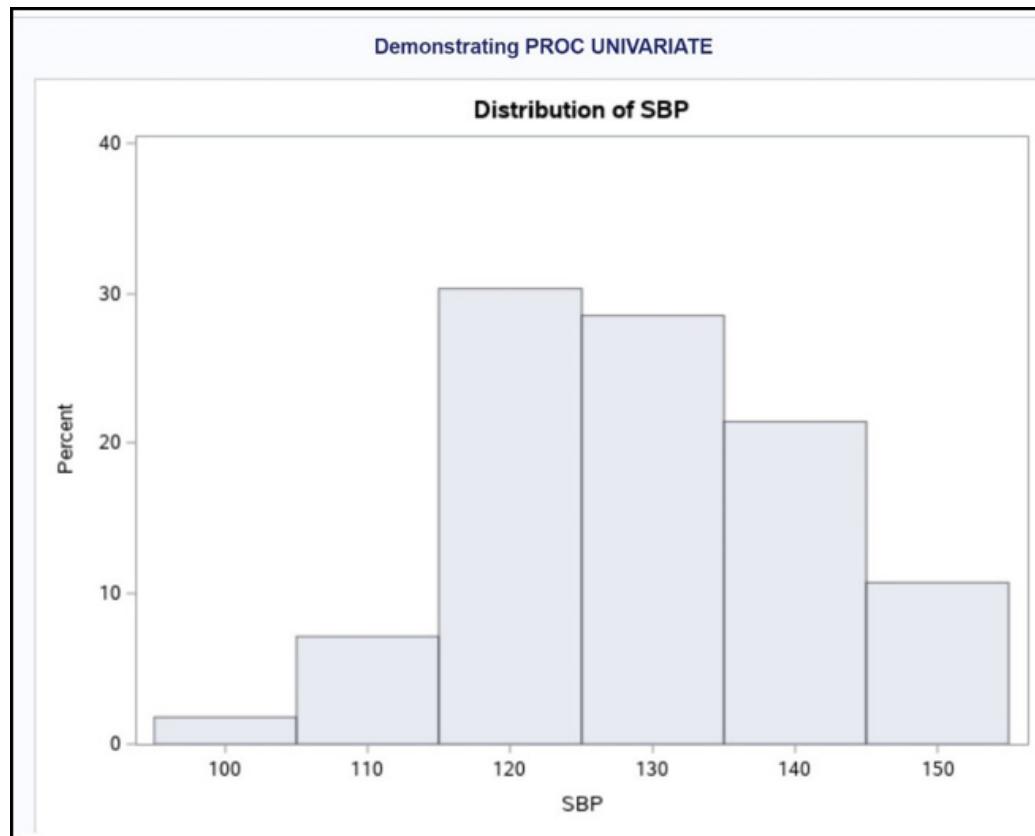
This section of output shows values of your variable at several different quantiles. In this section, you see that the largest value for SBP was 152 and the lowest was 100. The 25<sup>th</sup> percentile and the 75<sup>th</sup> percentile (122 and 138, respectively) are also popular measures.

Extreme Observations					
Lowest			Highest		
Value	Subj	Obs	Value	Subj	Obs
100	53	53	146	20	20
112	32	32	148	17	17
114	49	49	150	14	14
114	48	48	150	28	28
114	34	34	152	10	10

The section shows the five lowest and five highest observations in the data set. The section is especially useful to check if you have some extreme values, possibly data errors. You can have PROC UNIVARIATE print out more than five extreme observations by using a procedure option called NEXTROBS=N, where  $N$  is the number of high and low values that you want listed in the table. NEXTROBS stands for *Number of EXTReme OBServations*.

Missing Values			
Missing Value	Count	Percent Of	
		All Obs	Missing Obs
.	4	6.67	100.00

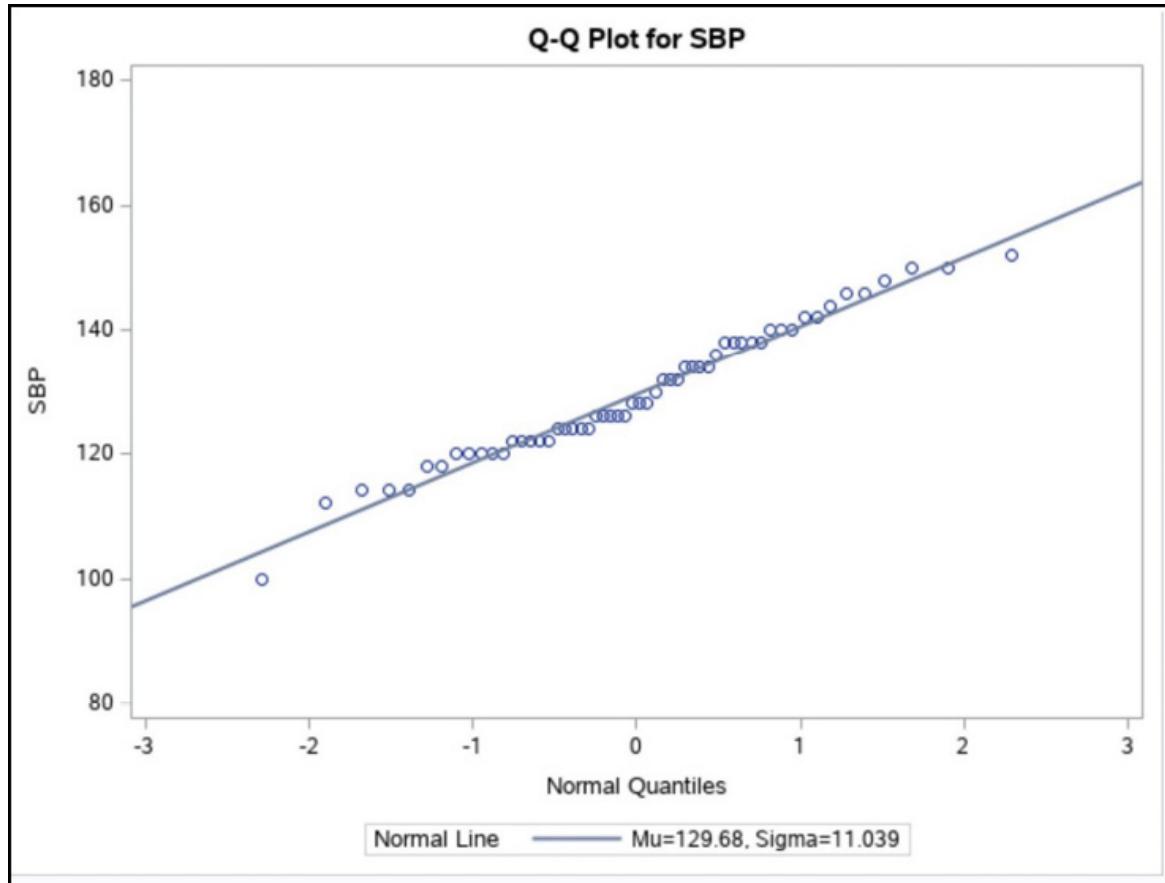
Here you see the number of missing values as a count and also as a percentage of all observations (6.67% in this data set).



This histogram was produced by the HISTOGRAM statement. If you want to change the bin sizes, you can use a MIDPOINTS= option in the HISTOGRAM statement. For example, to see bars from 100 to

150 but with a bin width of 5 instead of 10, the HISTOGRAM statement would look like this:

```
histogram / midpoints=100 to 150 by 5;
```



Finally, this is the Q-Q plot. You can see that the data points fall along a straight line, indicating that the values of SBP are close to normally distributed. If you prefer, you can substitute PROBPLOT for QQPLOT to obtain a similar plot called a *probability plot*.

## Conclusion

This chapter covered a lot of ground. Besides showing you how to generate summary statistics, you saw how to use PROC MEANS to create a data set containing summary information. This process blurs the distinction between DATA steps and PROC steps. You now know how to create a data set by running a procedure.

There are some topics that were not covered in detail. For example,

there are ways to use and interpret the `_TYPE_` variable placed in the output data set by PROC MEANS. At some point, you might also want to combine the summary data generated by PROC MEANS with the original raw data. For more information about these topics, I recommend the following book:

Cody, Ron. 2011. *SAS Statistics by Example*. Cary, NC: SAS Institute Inc.

## Problems

1. Starting with the SASHELP data set Heart, compute the following statistics for Height and Weight:
  1. Number of nonmissing values
  2. Number of missing values
  3. Mean
  4. Standard deviation
  5. Minimum
  6. MaximumUse the appropriate option to print statistics 3 through 6 with two places to the right of the decimal point.
2. Repeat Problem 1 except use a BY statement to compute the statistics broken down by Status.
3. Repeat Problem 2 except use a CLASS statement instead of a BY statement.
4. Repeat Problem 3 except use the two CLASS variables Sex and Status. In addition, use the procedure option to show the statistics broken down by every one of the CLASS variables.
5. Repeat Problem 1 except produce an output data set (Summary) instead of printed output. Use the AUTONAME option in the OUTPUT statement to name all of the variables in the Summary data set.
6. Repeat Problem 5 except produce all of the statistics broken down by Status. Use a CLASS statement and only include observations for each level of Status (i.e., do not include the statistics for all the subjects together).

7. Use PROC UNIVARIATE to analyze the variables Height and Weight from the SASHELP data set Heart. Include statements to produce a histogram for each of these variables.
8. Use PROC MEANS to compute the mean and standard deviation of systolic blood pressure (variable name Systolic) for men only (use a WHERE= data set option to do this) for two groups of men: one group comprising men with weights less than or equal to 150 and the other group for men weighing 151 pounds or more. Do this by using a CLASS statement and writing a format to place men in the two weight groups.

# Chapter 20: Computing Frequencies

## Introduction

You can use PROC FREQ to create one-way, two-way (row and column), and three-way (page, row, and column) tables. You can output counts and percentages as well as statistics such as chi-square and Fisher's Exact test. This chapter shows you examples of these tasks.

## Creating a Data Set to Demonstrate Features of PROC FREQ

The first step is to have some test data to demonstrate features of PROC FREQ. To this end, you can run Program 20.1 to generate a data set called Risk that includes the following variables:

Variable	Description
Subj	Subject number

---

Subj                              Subject number

Gender 1=Male, 2=Female

## Chol Cholesterol level

BP\_Status Blood pressure status (High or Low)

## Chol\_Status Cholesterol status (High or Low)

Heart Attack 1-Yes or 2-No

As in the previous chapter, this program randomly generates all of

the data. Here is the program.

### Program 20.1: Program to Generate Test Data Set Risk

```
proc format;
  value yesno 1 = '1-Yes'
            0 = '2-No';
run;
data Risk;
  call streaminit(12345678);
  length Age_Group $ 7;
  do Subj = 1 to 250;
    do Gender = 'F','M';
      Age = round(rand('uniform')*30 + 50);
      if Age lt 60 then Age_Group = '1:< 60';
      else if Age le 70 then Age_Group = '2:60-70';
      else Age_Group = '3:71+';
      if rand('uniform') lt .3 then BP_Status = 'High';
      else BP_Status = 'Low';
      Chol = rand('normal',200,30) +
             rand('uniform')*8*(Gender='M');
      Chol = round(Chol);
      if Chol gt 240 then Chol_Status = 'High';
      else Chol_Status = 'Low';
      Score = .1*Chol + age + 8*(Gender eq 'M') +
               10*(BP_Status = 'High');
      Heart_Attack = (Score gt 100)*(rand('uniform') lt .6);
      output;
    end;
  end;
  keep Subj Gender Age Chol Chol_Status BP_Status Heart_Attack;
  format Heart_Attack yesno.;
run;

title "Listing of Data Set Risk (first 10 observations)";
proc print data=Risk(obs=10);
  id Subj;
run;
```

Figure 20.1 is a listing of the first 10 observations from data set Risk.

### Figure 20.1: Output from Program 20.1

### Listing of Data Set Risk (first 10 observations)

Subj	Gender	Age	BP_Status	Chol	Chol_Status	Heart_Attack
1	F	68	Low	237	Low	2-No
1	M	72	High	210	Low	1-Yes
2	F	50	Low	157	Low	2-No
2	M	71	Low	179	Low	2-No
3	F	64	Low	184	Low	2-No
3	M	71	Low	202	Low	2-No
4	F	76	Low	253	High	2-No
4	M	57	Low	208	Low	2-No
5	F	75	High	207	Low	1-Yes
5	M	57	High	185	Low	2-No

## Using PROC FREQ to Generate One-Way Frequency Tables

To create one-way frequency tables, specify the list of variables in a TABLES statement as follows.

### Program 20.2: Using PROC FREQ to Generate One-Way Frequency Tables

```
title "One-way Frequency Tables";
proc freq data=Risk;
  tables Gender Heart_Attack;
run;
```

In this example, you are requesting one-way frequencies for Gender and Heart\_Attack. PROC FREQ can also compute frequencies for numeric variables. However, beware, if you have a variable such as Age with many different values, PROC FREQ will compute frequencies for every unique value. Here is the output.

**Figure 20.2: Output from Program 20.2**

### One-way Frequency Tables

Gender	Frequency	Percent	Cumulative Frequency	Cumulative Percent
F	250	50.00	250	50.00
M	250	50.00	500	100.00

Heart_Attack	Frequency	Percent	Cumulative Frequency	Cumulative Percent
2-No	415	83.00	415	83.00
1-Yes	85	17.00	500	100.00

You see frequency, percent, cumulative frequency, and cumulative percent in the two tables. The order of the values is based, by default, on the internal values of the variables. For numeric variables, lower numbers come before higher numbers; for character variables, the values are sorted alphabetically. Because Heart\_Attack was coded as 0 and 1, the first entry in the table is 2-No, the formatted value for the internal value of 0.

You might wonder why the formats ‘2-No’ and ‘1-Yes’ were created. Why not just ‘Yes’ and ‘No’? The answer is that you can request the order in any tables to be based on the formatted value of a variable by using the PROC FREQ option ORDER=formatted. Alphabetically, ‘1-Yes’ comes before ‘2-No’. That will force the ‘Yes’ frequencies to be listed before the ‘No’ frequencies. For the 2 x 2 tables coming up next, it is preferable (for statistical purposes) to have the ‘Yes’ values come before the ‘No’ values.

It is rarely useful to see cumulative frequencies or percentages. You can eliminate these values from the tables with a TABLES option called NOCUM. The program that follows uses both the ORDER= procedure option and the NOCUM statement option.

### Program 20.3: Changing the Table Order and Removing the Cumulative Statistics

```

title "One-way Frequency Tables";
proc freq data=Risk order=formatted;
  tables Gender Heart_Attack / nocum;

```

```
run;
```

Figure 20.3 shows how these two options affect the tables.

### Figure 20.3: Output from Program 20.3

One-way Frequency Tables		
Gender	Frequency	Percent
F	250	50.00
M	250	50.00

Heart_Attack	Frequency	Percent
1-Yes	85	17.00
2-No	415	83.00

The order of the frequencies for the variable Heart\_Attack is now controlled by the formatted values for Gender and Heart\_Attack, and the cumulative frequencies are no longer included in the tables.

Before we leave this section, you should know one other ordering option: ORDER=freq. This option arranges the frequencies from the most frequent to the least frequent. At times, this can be extremely useful.

## Creating Two-Way Frequency Tables

You can create a row by column table by placing one or more variables in the TABLES statement, followed by an asterisk, followed by another list of variables. The variables before the asterisk form the rows of the tables, and the variables after the asterisk form the columns of the tables. If you select more than one variable for the rows or columns list, that list must be placed in parentheses. The syntax is:

```
tables (list of row variables) * (list of column variables)
/ options;
```

PROC FREQ will create a table for every combination of variables in

the rows list with every variable in the columns list. If you write the following TABLES statement:

```
tables (A B) * (C D E);
```

PROC FREQ will create tables A by C, A by D, A by E, B by C, B by D, and B by E.

To demonstrate a two-way frequency table, let's create a table using the two variables BP\_Status and Heart\_Attack. Traditionally, epidemiologists like the outcome variable (Heart\_Attack, in this example) to form the columns of the table. Here is the program.

#### Program 20.4: Creating a Two-Way Frequency Table

```
title "Two-way Frequency Table of BP_Status by Heart_Attack";  
proc freq data=Risk order=formatted;  
  tables BP_Status * Heart_Attack;  
run;
```

This table request produces a table with BP\_Status forming the rows of the table and Heart\_Attack forming the columns of the table. Here is the listing.

Figure 20.4: Output from Program 20.4

Two-way Frequency Table of BP_Status by Heart_Attack				
Frequency Percent Row Pct Col Pct	Table of BP_Status by Heart_Attack			
	BP_Status	Heart_Attack		
		1-Yes	2-No	Total
	High	44 8.80 28.57 51.76	110 22.00 71.43 26.51	154 30.80
	Low	41 8.20 11.85 48.24	305 61.00 88.15 73.49	346 69.20
	Total	85 17.00	415 83.00	500 100.00

The key to this table is displayed in the upper left-hand part of the output. The top number in each box is a frequency count. For

example, there were 44 subjects with high blood pressure who had a heart attack. The second number in the table is a percent: These 44 subjects represent 8.8% of all the subjects in the table (500). The third number in the table is a row percentage. Of the 154 subjects who had high blood pressure, 28.57% of them had a heart attack. Finally, the bottom number in the box is a column percentage. Of the 85 subjects who had a heart attack, 51.76% of them had high blood pressure.

The primary numbers of interest to a medical researcher looking at this table would be the percent of people in the high blood pressure group who had a heart attack (28.57%) compared to the percent of people in the low blood pressure group who had a heart attack (11.85%).

One popular statistical test to decide whether the difference in these two proportions is statistically significant is called *chi-square*. You can request a chi-square test by adding the CHISQ TABLES option to your program. The modified program is shown next.

### Program 20.5: Adding a Request for a Chi-Square Test

```
title "Two-way Frequency Table of BP_Status by Heart_Attack";
proc freq data=Risk order=formatted;
  tables BP_Status * Heart_Attack / chisq;
run;
```

The output contains the same table as above, with the following added information.

**Figure 20.5: Partial Output from Program 20.5**

Statistics for Table of BP\_Status by Heart\_Attack

Statistic	DF	Value	Prob
Chi-Square	1	21.1184	<.0001
Likelihood Ratio Chi-Square	1	19.7866	<.0001
Continuity Adj. Chi-Square	1	19.9500	<.0001
Mantel-Haenszel Chi-Square	1	21.0762	<.0001
Phi Coefficient		0.2055	
Contingency Coefficient		0.2013	
Cramer's V		0.2055	

The CHISQ option produces other statistics that are not shown here. Of primary interest in this table is the first row where you see a chi-square of 21.1184 with a *p*-value of <.0001. This very small *p*-value tells you that if blood pressure status had no relationship to having a heart attack, the probability of getting such a large difference in the proportion between the two blood pressure groups by chance alone is less than .0001. This is considered highly significant by statisticians.

## Creating Three-Way Frequency Tables

You can create three-way frequency tables by specifying the page, row, and column variables, separated by an asterisk. You should be cautious when you do this as it might generate a large volume of output if your page variable(s) have large numbers of values. To keep the output small, the example for a three-way table that follows, uses Gender as the page variable, BP\_Status as the row variable, and Heart\_Attack as the column variable. Here is the program.

### Program 20.6: Creating a Three-Way Table

```
title "Three-way Table of Gender by BP_Status by Heart_Attack";
proc freq data=Risk order=formatted;
  tables Gender * BP_Status * Heart_Attack;
run;
```

Here is the output.

**Figure 20.6: Output from Program 20.6**

Three-way Table of Gender by BP_Status by Heart_Attack			
	Table 1 of BP_Status by Heart_Attack		
	Controlling for Gender=F		
BP_Status	Heart_Attack		Total
High	18 7.20 23.38 72.00	59 23.60 76.62 26.22	77 30.80
Low	7 2.80 4.05 28.00	166 66.40 95.95 73.78	173 69.20
Total	25 10.00	225 90.00	250 100.00

Table 2 of BP_Status by Heart_Attack			
	Controlling for Gender=M		
	Heart_Attack		
BP_Status	1-Yes	2-No	Total
High	26 10.40 33.77 43.33	51 20.40 66.23 26.84	77 30.80
Low	34 13.60 19.65 56.67	139 55.60 80.35 73.16	173 69.20
Total	60 24.00	190 76.00	250 100.00

Inspection of these tables suggests that high blood pressure is related to the incidence of heart attack in females and males. You could verify this by adding the CHISQ option in the TABLES statement.

## Using Formats to Create Groups for Numeric Variables

You can use formats to group values to be displayed in the frequency tables. One of the unique features of PROC FREQ is that it automatically uses formatted values for any variable that is associated with a format. For example, suppose you want to compute frequencies for the variable Age, and you want age groups of 20 years. The program that follows demonstrates how you can use a format to accomplish this task.

### Program 20.7: Using Formats to Group a Numeric Variable

```
proc format;
  value Agegroup low-19 = '<20'
                20-39  = '20 to 39'
                40-59  = '40 to 59'
                60-79  = '60 to 79'
                80-high= '80+'
                .      = 'Missing';
run;

title "Using a Format to Group a Numeric Variable";
proc freq data=Risk;
  tables Age / nocum;
  format Age Agegroup. ;
run;
```

Because you included a FORMAT statement in the procedure, frequencies are computed for the age groups, not the actual age values. Here is the output.

**Figure 20.7: Output from Program 20.7**

#### Using a Format to Group a Numeric Variable

Age	Frequency	Percent
40 to 59	157	31.40
60 to 79	330	66.00
80+	13	2.60

Frequencies are displayed for age groups. Notice the order of the frequencies. Because the option ORDER= was omitted in the program, PROC FREQ used its default ordering called INTERNAL. So, regardless of the values of the formatted ranges, the order in the

table is by increasing age.

## Conclusion

PROC FREQ is one of the main procedures for computing frequencies. It can construct one-way, two-way, and three-way tables. In addition, by adding options, you can request a variety of statistical tests to be performed on the resulting tables.

## Problems

1. Using the SASHELP data set Heart, compute one-way frequencies and percentages for the variables Status, BP\_Status, and Smoking\_Status. Use options to omit cumulative frequencies from the resulting tables.
2. Compute frequencies and percentages for Smoking\_Status from the SASHELP data set Heart. Use a procedure option so that the table is in frequency order (with the most frequent category listed first).
3. Compute frequencies and percentages for Status from the SASHELP data set Heart. Create formats and use the appropriate PROC FREQ option so that the category “Dead” comes before “Alive” in the resulting table.
4. Using the SASHELP data set Heart, construct a two-way table for the variables Sex (rows of the table) and Status (columns of the table). Arrange the table so that “Dead” comes before “Alive” and “Male” comes before “Female”. If you are statistically inclined, add an option to compute the chi-square for this table.
5. Using the SASHELP data set Heart, construct a three-way table with Sex as the page dimension, Weight\_Status as the row dimension, and Status as the column dimension. There are three categories for Weight\_Status. Write the necessary code to eliminate all observations where Weight\_Status is equal to “Underweight”.
6. Using the SASHELP data set Heart, construct a two-way table of Systolic (rows) by Diastolic (columns) blood pressures. Both

of these variables are numeric. Use formats to divide Systolic and Diastolic into two groups (below 200 versus 200 and above for Systolic; below 120 versus 120 and above for Diastolic). If you are interested, include a request for the chi-square.

# Ready to take your SAS® and JMP® skills up a notch?



Be among the first to know about new books,  
special events, and exclusive discounts.

**[support.sas.com/newbooks](https://support.sas.com/newbooks)**

Share your expertise. Write a book with SAS.

**[support.sas.com/publish](https://support.sas.com/publish)**

Continue your skills development with free online learning.

**[www.sas.com/free-training](https://www.sas.com/free-training)**



**sas.com/books**  
*for additional books and resources.*

**sas**  
THE POWER TO KNOW®

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 trademarks of their respective companies. © 2020 SAS Institute Inc. All rights reserved. M2043821 US1120

# Contents

1. [Contents](#)
2. [About This Book](#)
3. [About The Author](#)
4. [Acknowledgements](#)
5. [Part I: Getting Acquainted with the SAS Studio Environment](#)
6. [Chapter 1: Introduction to SAS OnDemand for Academics](#)
  1. [Introduction: An Overview of SAS OnDemand for Academics](#)
  2. [Registering for ODA](#)
  3. [Conclusion](#)
7. [Chapter 2: The SAS Studio Interface](#)
  1. [Introduction](#)
  2. [Exploring the Built-In Data Sets](#)
  3. [Sorting Your Data](#)
  4. [Switching between Column Names and Column Labels](#)
  5. [Resizing Tables](#)
  6. [Creating Filters](#)
  7. [Conclusion](#)
8. [Chapter 3: Importing Your Own Data](#)
  1. [Introduction](#)
  2. [Uploading Data from Your Local Computer to SAS Studio](#)
  3. [Listing the SAS Data Set](#)
  4. [Importing an Excel Workbook with Invalid SAS Variable Names](#)
  5. [Importing an Excel Workbook That Does Not Have Variable Names](#)
  6. [Importing Data from a CSV File](#)
  7. [Conclusion](#)

9. [Chapter 4: Creating Reports](#)

1. [Introduction](#)
2. [Using the List Data Task to Create a Simple Listing](#)
3. [Filtering Data](#)
4. [Sorting Data](#)
5. [Outputting HTML and PDF Files](#)
6. [Joining Tables \(Using the Query Window\)](#)
7. [Conclusion](#)

10. [Chapter 5: Summarizing Data Using SAS Studio](#)

1. [Introduction](#)
2. [Summarizing Numeric Variables](#)
3. [Adding a Classification Variable](#)
4. [Summarizing Character Variables](#)
5. [Conclusion](#)

11. [Chapter 6: Graphing Data](#)

1. [Introduction](#)
2. [Creating a Frequency Bar Chart](#)
3. [Creating a Bar Chart with a Response Variable](#)
4. [Adding a Group Variable](#)
5. [Creating a Pie Chart](#)
6. [Creating a Scatter Plot](#)
7. [Conclusion](#)

12. [Part II: Learning How to Write Your Own SAS Programs](#)

13. [Chapter 7: An Introduction to SAS Programming](#)

1. [SAS as a Programming Language](#)
2. [The SAS Studio Programming Windows](#)
3. [Your First SAS Program](#)
  1. [DATA Statement](#)
  2. [INILE Statement](#)
  3. [INPUT Statement](#)
  4. [Assignment Statement](#)
4. [How the DATA Step Works](#)

5. [How the INPUT Statement Works](#)
6. [Reading Delimited Data](#)
7. [How Procedures \(PROCs\) Work](#)
8. [How SAS Works: A Look Inside the “Black Box”](#)
9. [Conclusion](#)
14. [Chapter 8: Reading Data from External Files](#)
  1. [Introduction](#)
  2. [Reading Data Values Separated by Delimiters](#)
    1. [Reading Comma-Separated Values Files](#)
    2. [Reading Data Separated by Other Delimiters](#)
  3. [Reading Data in Fixed Columns](#)
    1. [Reading Data in Fixed Columns Using Column Input](#)
    2. [Reading Data in Fixed Columns Using Formatted Input](#)
  4. [Conclusion](#)
  5. [Problems](#)
15. [Chapter 9: Reading and Writing SAS Data Sets](#)
  1. [What’s a SAS Data Set?](#)
  2. [Temporary Versus Permanent SAS Data Sets](#)
  3. [Creating a Library by Submitting a LIBNAME Statement](#)
  4. [Using the Library Tab to Create a Permanent Library](#)
  5. [Reading from a Permanent SAS Data Set](#)
  6. [Conclusion](#)
  7. [Problems](#)
16. [Chapter 10: Creating Formats and Labels](#)
  1. [What Is a SAS Format and Why Is It Useful?](#)
  2. [Using SAS Built-in Formats](#)
  3. [More Examples to Demonstrate How to Write Formats](#)
  4. [Describing the Difference between a FORMAT Statement in a Procedure and a FORMAT Statement in a DATA Step](#)
  5. [Making Your Formats Permanent](#)
  6. [Creating Variable Labels](#)

- 7. [Conclusion](#)
- 8. [Problems](#)
- 17. [Chapter 11: Performing Conditional Processing](#)
  - 1. [Introduction](#)
  - 2. [Grouping Age Using Conditional Processing](#)
  - 3. [Using Conditional Logic to Check for Data Errors](#)
  - 4. [Describing the IN Operator](#)
  - 5. [Using Boolean Logic \(AND, OR, and NOT Operators\)](#)
  - 6. [A Special Caution When Using Multiple OR Operators](#)
  - 7. [Conclusion](#)
  - 8. [Problems](#)
- 18. [Chapter 12: Performing Iterative Processing: Looping](#)
  - 1. [Introduction](#)
  - 2. [Demonstrating a DO Group](#)
  - 3. [Describing a DO Loop](#)
  - 4. [Using a DO Loop to Graph an Equation](#)
  - 5. [DO Loops with Character Values](#)
  - 6. [Leaving a Loop Based on Conditions \(DO WHILE and DO UNTIL Statements\)](#)
    - 1. [DO WHILE](#)
    - 2. [Combining an Iterative Loop with a WHILE Condition](#)
    - 3. [DO UNTIL](#)
    - 4. [Demonstrating That a DO UNTIL Loop Executes at Least Once](#)
    - 5. [Combining an Iterative Loop with an UNTIL Condition](#)
    - 7. [LEAVE and CONTINUE Statements](#)
    - 8. [Conclusion](#)
    - 9. [Problems](#)
- 19. [Chapter 13: Working with SAS Dates](#)
  - 1. [Introduction](#)

2. [Reading Dates from Text Data](#)
3. [Creating a SAS Date from Month, Day, and Year Values](#)
4. [Describing a Date Constant](#)
5. [Extracting the Day of the Week, Day of the Month, Month, and Year from a SAS Date](#)
6. [Adding a Format to the Bar Chart](#)
7. [Computing Age from Date of Birth: The YRDIF Function](#)
8. [Conclusion](#)
9. [Problems](#)
20. [Chapter 14: Subsetting and Combining SAS Data Sets](#)
  1. [Introduction](#)
  2. [Subsetting \(Filtering\) Data in a SAS Data Set](#)
  3. [Describing a WHERE= Data Set Option](#)
  4. [Describing a Subsetting IF Statement](#)
  5. [A More Efficient Way to Subset Data When Reading Raw Data](#)
  6. [Creating Several Data Subsets in One DATA Step](#)
  7. [Combining SAS Data Sets \(Combining Rows\)](#)
  8. [Adding a Few Observations to a Large Data Set \(PROC APPEND\)](#)
  9. [Interleaving Data Sets](#)
  10. [Merging Two Data Sets \(Adding Columns\)](#)
  11. [Controlling Which Observations Are Included in a Merge \(IN= Data Set Option\)](#)
  12. [Performing a One-to-Many or Many-to-One Merge](#)
  13. [Merging Two Data Sets with Different BY Variable Names](#)
  14. [Merging Two Data Sets with One Character and One Numeric BY Variable](#)
  15. [Updating a Master File from a Transaction File \(UPDATE Statement\)](#)
  16. [Conclusion](#)

## 17. Problems

### 21. Chapter 15: Describing SAS Functions

1. Introduction

2. Describing Some Useful Numeric Functions

1. Function Name: MISSING

2. Function Name: N

3. Function Name: NMISS

4. Function Name: SUM

5. Function Name: MEAN

6. Function Name: MIN

7. Function Name: MAX

8. Function Name: SMALLEST

9. Function Name: LARGEST

10. Programming Example Using the N, NMISS, MAX, LARGEST, and MEAN Functions

11. Function Name: INPUT

12. CALL Routine: CALL SORTN

13. Function Name: LAG

14. Function Name: DIF

3. Describing Some Useful Character Functions

1. Function Names: LENGTHN and LENGTHC

2. Function Names: TRIMN and STRIP

3. Function Names: UPCASE, LOWCASE, and PROPCASE (Functions That Change Case)

4. Function Name: PUT

5. Function Name: SUBSTRN (Newer Version of the SUBSTR Function)

6. Function Names: FIND and FINDC

7. Function Names: CAT, CATS, and CATX

8. Function Names: COUNT and COUNTC

9. Function Name: COMPRESS

10. Function Name: SCAN

11. [CALL Routine: CALL MISSING](#)
  12. [Function Names: NOTDIGIT, NOTALPHA, and NOTALNUM](#)
  13. [Function Names: ANYDIGIT, ANYALPHA, and ANYALNUM](#)
  14. [Function Name: TRANWRD](#)
  4. [Conclusion](#)
  5. [Problems](#)
22. [Chapter 16: Working with Multiple Observations per Subject](#)
1. [Introduction](#)
  2. [Useful Tools for Working with Longitudinal Data](#)
  3. [Describing First. and Last. Variables](#)
  4. [Computing Visit-to-Visit Differences](#)
  5. [Computing Differences between the First and Last Visits](#)
  6. [Counting the Number of Visits for Each Patient](#)
  7. [Conclusion](#)
  8. [Problems](#)
23. [Chapter 17: Describing Arrays](#)
1. [Introduction](#)
  2. [What Is an Array?](#)
  3. [Describing a Character Array](#)
  4. [Performing an Operation on Every Numeric Variable in a Data Set](#)
  5. [Performing an Operation on Every Character Variable in a Data Set](#)
  6. [Converting a Data Set with One Observation per Subject into a Data Set with Multiple Observations per Subject](#)
  7. [Converting a Data Set with Multiple Observations per Subject into a Data Set with One Observation per Subject](#)
  8. [Conclusion](#)
  9. [Problems](#)
24. [Chapter 18: Displaying Your Data](#)

1. [Introduction](#)
  2. [Producing a Simple Report Using PROC PRINT](#)
  3. [Using Labels Instead of Variable Names as Column Headings](#)
  4. [Including a BY Variable in a Listing](#)
  5. [Including the Number of Observations in a Listing](#)
  6. [Conclusion](#)
  7. [Problems](#)
25. [Chapter 19: Summarizing Data with SAS Procedures](#)
1. [Introduction](#)
  2. [Using PROC MEANS \(with the Default Options\)](#)
  3. [Using PROC MEANS Options to Customize the Summary Report](#)
  4. [Computing Statistics for Each Value of a BY Variable](#)
  5. [Using a CLASS Statement Instead of a BY Statement](#)
  6. [Including Multiple CLASS Variables with PROC MEANS](#)
  7. [Statistics Broken Down Every Way](#)
  8. [Using PROC MEANS to Create a Summary Data Set](#)
  9. [Letting PROC MEANS Name the Variables in the Output Data Set](#)
  10. [Creating a Summary Data Set with CLASS Variables](#)
  11. [Using a Formatted CLASS Variable](#)
  12. [Demonstrating PROC UNIVARIATE](#)
  13. [Conclusion](#)
  14. [Problems](#)
26. [Chapter 20: Computing Frequencies](#)
1. [Introduction](#)
  2. [Creating a Data Set to Demonstrate Features of PROC FREQ](#)
  3. [Using PROC FREQ to Generate One-Way Frequency Tables](#)
  4. [Creating Two-Way Frequency Tables](#)

5. [Creating Three-Way Frequency Tables](#)
6. [Using Formats to Create Groups for Numeric Variables](#)
7. [Conclusion](#)
8. [Problems](#)

## Landmarks

1. [Cover](#)
2. [Table of Contents](#)