

Handbook of  
**SAS® DATA Step  
Programming**



# Handbook of **SAS<sup>®</sup> DATA Step Programming**

Arthur Li



**CRC Press**

Taylor & Francis Group

Boca Raton London New York

---

CRC Press is an imprint of the  
Taylor & Francis Group, an **informa** business  
A CHAPMAN & HALL BOOK

CRC Press  
Taylor & Francis Group  
6000 Broken Sound Parkway NW, Suite 300  
Boca Raton, FL 33487-2742

© 2013 by Taylor & Francis Group, LLC  
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works  
Version Date: 20130226

International Standard Book Number-13: 978-1-4665-5239-5 (eBook - PDF)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access [www.copyright.com](http://www.copyright.com) (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

**Trademark Notice:** Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

**Visit the Taylor & Francis Web site at**  
**<http://www.taylorandfrancis.com>**

**and the CRC Press Web site at**  
**<http://www.crcpress.com>**

*To Dave*



---

# Contents

---

Preface.....	xiii
The Author .....	xix
Acknowledgments .....	xxi
<b>1. Introduction to SAS®.....</b>	<b>1</b>
1.1 SAS Program and Language .....	1
1.2 Reading Data into SAS .....	2
1.2.1 The SAS Data Set and SAS Library .....	2
1.2.2 Reading a SAS Data Set.....	3
1.2.3 Reading a Raw Data File with Fixed Fields .....	6
1.2.4 Reading Data Entered Directly into the Program.....	8
1.3 Creating and Modifying Variables .....	9
1.3.1 The Assignment Statement and SAS Expression .....	9
1.3.2 Creating Variables Conditionally .....	12
1.4 Base SAS Procedures .....	13
1.4.1 Common Statements in SAS Procedures: TITLE, BY, and WHERE Statements .....	13
1.4.2 The CONTENTS Procedure .....	14
1.4.3 The SORT Procedure .....	16
1.4.4 The PRINT Procedure .....	16
1.4.5 The MEANS Procedure.....	18
1.4.6 The FREQ Procedure.....	20
1.5 Subsetting Data by Selecting Variables .....	24
1.5.1 Selecting Variables with the KEEP= Data Set Option or KEEP Statement.....	25
1.5.2 Selecting Variables with the DROP= Data Set Option or DROP Statement .....	27
1.5.3 Where to Specify the DROP= and KEEP= Data Set Options and DROP/KEEP Statements .....	27
1.6 Changing the Appearance of Data .....	28
1.6.1 Labeling Variables .....	29
1.6.2 Formatting Variable Values Using SAS FORMATS .....	31
Exercises .....	33
<b>2. Creating Variables Conditionally .....</b>	<b>35</b>
2.1 The IF-THEN/ELSE Statement.....	35
2.1.1 Steps for Creating a Variable.....	35
2.1.2 Handling Missing Values When Creating Variables.....	37
2.1.3 TRUE and FALSE: Logical Expressions.....	39

2.1.4	The LENGTH Attribute .....	41
2.1.5	DO Group.....	43
2.2	Executing One of Several Statements.....	45
2.2.1	Multiple IF-THEN/ELSE Statements.....	45
2.2.2	Executing Statements Using the SELECT Group .....	48
2.3	Modifying the IF-THEN/ELSE Statement with the Assignment Statement.....	52
	Exercises .....	54
<b>3.</b>	<b>Understanding How the DATA Step Works .....</b>	<b>55</b>
3.1	DATA Step Processing Overview .....	55
3.1.1	DATA Step Compilation Phase.....	57
3.1.2	DATA Step Execution Phase.....	58
3.1.3	The Importance of the OUTPUT Statement.....	61
3.1.4	The Difference between Reading a Raw Data Set and a SAS Data Set.....	61
3.2	Retaining the Value of Newly Created Variables .....	62
3.2.1	The RETAIN Statement.....	62
3.2.2	The SUM Statement .....	64
3.3	Conditional Processing in the DATA Step .....	66
3.3.1	The Subsetting IF Statement .....	66
3.3.2	Detecting the End of a Data Set by Using the END= Option.....	68
3.3.3	Restructuring Data Sets from Wide Format to Long Format.....	68
3.4	Debugging Techniques .....	70
3.4.1	Using the PUT Statement to Observe the Contents of the PDV .....	70
3.4.2	Using the DATA Step Debugger.....	74
	Exercises .....	76
<b>4.</b>	<b>BY-Group Processing in the DATA Step .....</b>	<b>79</b>
4.1	Introduction to BY-Group Processing.....	79
4.1.1	The FIRST.VARIABLE and the LAST.VARIABLE .....	79
4.1.2	The Execution Phase of BY-Group Processing .....	81
4.2	Applications Utilizing BY-Group Processing .....	85
4.2.1	Calculating Mean Score within Each BY Group .....	87
4.2.2	Creating Data Sets with Duplicate or Non-Duplicate Observations.....	88
4.2.3	Obtaining the Most Recent Non-Missing Data within Each BY Group .....	89
4.2.4	Restructuring Data Sets from Long Format to Wide Format .....	91
	Exercises .....	92



<b>5. Writing Loops in the DATA Step .....</b>	<b>95</b>
5.1 Implicit and Explicit Loops.....	95
5.1.1 Implicit Loops.....	95
5.1.2 Explicit Loops.....	96
5.1.3 Nested Loops.....	102
5.1.4 Combining Implicit and Explicit Loops .....	103
5.2 Utilizing Loops to Create Samples .....	103
5.2.1 Direct-Access Mode .....	104
5.2.2 Creating a Systematic Sample .....	105
5.2.3 Creating a Random Sample with Replacement.....	106
5.2.4 Creating a Random Sample without Replacement .....	108
5.3 Using Looping to Read a List of External Files .....	110
5.3.1 Using an Iterative DO Loop to Read an External File ....	110
5.3.2 Using an Iterative DO-Loop to Read Multiple External Files .....	111
Exercises .....	117
 <b>6. Array Processing.....</b>	 <b>121</b>
6.1 Introduction to Array Processing .....	121
6.1.1 Situations for Utilizing Array Processing .....	121
6.1.2 Defining and Referencing One-Dimensional Arrays.....	123
6.1.3 Compilation and Execution Phases for Array Processing .....	125
6.2 Functions and Operators Related to Array Processing .....	126
6.2.1 The DIM, HBOUND, and LBOUND Functions .....	126
6.2.2 Using the IN and OF Operator with an Array .....	129
6.3 Some Array Applications.....	130
6.3.1 Creating a Group of Variables by Using Arrays.....	130
6.3.2 Calculating Products of Multiple Variables .....	131
6.3.3 Restructuring Data Sets Using One-Dimensional Arrays .....	132
6.4 Applications That Use Multi-Dimensional Arrays .....	133
6.4.1 Calculating Average SBP for Pre- and Post-Treatment.....	133
6.4.2 Restructuring Data Sets by Using a Multi-Dimensional Array .....	135
Exercises .....	136
 <b>7. Combining Data Sets.....</b>	 <b>139</b>
7.1 Vertically Combining Data Sets.....	139
7.1.1 Concatenating Data Sets .....	139
7.1.2 Interleaving Data Sets .....	142
7.2 Horizontally Combining Data Sets .....	143
7.2.1 One-to-One Reading .....	143
7.2.2 One-to-One Merging.....	146

7.2.3	Match-Merging.....	147
7.2.4	Updating Data Sets.....	151
	Exercises.....	152
<b>8.</b>	<b>Data Input and Output.....</b>	<b>155</b>
8.1	Introduction to Reading and Writing Text Files.....	155
8.1.1	Steps for Reading Text Files.....	155
8.1.2	Steps for Writing Text Files.....	157
8.1.3	Data Informat.....	157
8.1.4	Data Format.....	158
8.1.5	SAS Date and Time Values.....	158
8.2	Reading Text Files.....	160
8.2.1	Column Input.....	161
8.2.2	Formatted Input.....	162
8.2.3	List Input.....	164
8.2.4	Modified List Input.....	168
8.2.5	Mixed Input.....	170
8.2.6	Creating Observations by Using the Line Pointer-Controls.....	171
8.2.7	Creating Observations by Using Line-Hold Specifiers.....	172
8.3	Creating Text Files.....	175
8.3.1	Column Output.....	175
8.3.2	Formatted Output.....	176
8.3.3	List Output.....	177
	Exercises.....	177
<b>9.</b>	<b>Data Step Functions.....</b>	<b>181</b>
9.1	Introduction to Functions and CALL Routines.....	181
9.1.1	Functions.....	181
9.1.2	CALL Routines.....	182
9.1.3	Categories of Functions and CALL Routines.....	184
9.2	Date and Time Functions.....	185
9.2.1	Creating Date and Time Values.....	185
9.2.2	Extracting Components from Date and Time Values.....	187
9.2.3	Date and Time Interval Functions.....	188
9.3	Character Functions.....	190
9.3.1	Functions for Changing Character Cases.....	190
9.3.2	Functions for Concatenating Character Strings.....	191
9.3.3	Functions for Searching, Exacting, and Replacing Character Strings.....	194
9.4	Functions for Converting Variable Types.....	198
9.4.1	The INPUT Function.....	198
9.4.2	The PUT Function.....	201
	Exercises.....	203

<b>10. Useful SAS® Procedures.....</b>	<b>205</b>
10.1 Using the SORT Procedure to Eliminate Duplicate Observations .....	205
10.1.1 Eliminating Observations with Duplicate BY Values .....	205
10.1.2 Eliminating Duplicate Observations .....	207
10.2 Using the COMPARE Procedure to Compare the Contents of Two Data Sets .....	208
10.2.1 Information Provided from PROC COMPARE .....	209
10.2.2 Comparing Observations with Common ID Values.....	212
10.3 Restructuring Data Sets Using the TRANSPOSE Procedure .....	215
10.3.1 Transposing an Entire Data Set .....	216
10.3.2 Introduction to Transposing BY Groups .....	219
10.3.3 Where the ID Statement Does Not Work for Transposing BY Groups .....	220
10.3.4 Where the ID Statement Is Essential for Transposing BY Groups .....	221
10.3.5 Handling Duplicated Observations Using the LET Option.....	222
10.3.6 Situations Requiring Two or More Transpositions .....	224
10.4 Creating the User-Defined Format Using the FORMAT Procedure .....	227
10.4.1 Creating User-Defined Formats.....	228
10.4.2 Retrieving User-Defined Formats .....	233
10.4.3 Creating Variables by Using User-Defined Formats.....	235
10.5 Using the OPTIONS Procedure to Modify SAS System Options .....	236
Exercises .....	239
 <b>References .....</b>	 <b>241</b>
<b>Index .....</b>	<b>243</b>



---

## *Preface*

---

A common statistical programmer's task generally begins with reading one or more raw data sets into a statistical software and performing data management and manipulation, such as checking and modifying values of variables to ensure they satisfy the analytic quality, transposing data into a desired shape for a certain type of statistical analysis, merging multiple data sets by common variables, etc. Once the data has been assured of possessing the desired quality, the transformed data is ready for statistical analysis by using procedures that are provided by the statistical software. Data manipulation is an essential step to obtaining a reliable, statistical, analytical result because an analytical result that is based on unreliable data is not trustworthy; this is often referred to as "garbage in, garbage out." Successfully creating reliable data solely depends upon writing an accurate computer program.

In SAS®, data manipulation and management are mostly performed in the DATA steps; conducting statistical analysis and creating reports is carried out by using SAS procedures (PROC steps). A computer program that is used to perform data manipulation and analysis consists of a series of DATA or PROC steps or both and is written within the SAS programming language. The DATA and PROC steps consist of a series of SAS statements that are created by following the SAS language syntax. A PROC step is often easy to write because it is purely syntax driven; therefore, simply knowing how to follow the syntax from the documentation will sometimes be sufficient to help you accomplish the task successfully. On the other hand, in order to write an accomplished program in the DATA step, a programmer must be able to understand programming logic and to know how to implement and even create his or her own programming algorithm.

The focus of this book is not about learning statistical procedures but rather learning how best to manage and manipulate data by using the DATA step. Beginning programmers often tend to focus on learning syntax without focusing on programming logic and algorithms, which often results in common problems when they create a SAS data set. For example, the data set that they created is not what they originally intended to create—that is, there are more or less observations than intended or the value of the newly created variable was not retained correctly. These types of mistakes are most commonly committed because programming novices don't understand fundamental and unique SAS programming concepts, such as understanding the compilation and execution phases of the DATA step, what happens in the program data vector (PDV) during the DATA step execution, etc. This book will provide insight to readers that simply learning syntax will not solve all the problems that they'll

encounter; instead, they need to understand SAS processing in order to be successful programmers.

Another common problem novice programmers face is a lack of programming strategies. Therefore, when SAS programming novices encounter a new programming task, often they don't know where to begin and what steps will be involved in solving the problem. Most of the examples in this book begin with discussing the strategies and steps for solving the problems, then providing a solution, and in the end, providing a more detailed explanation for the solution.

---

## An Overview of SAS Software

SAS was originally the acronym for Statistical Analysis System, which is an integrated software system that utilizes fourth-generation programming language to perform tasks like data management, report writing, statistical analysis, data warehousing, and application development. SAS has undergone various upgrades to its software system over the years. This book was written using Version 9.2.

The core component of the SAS system is Base SAS software, which consists of different modules such as DATA steps, SAS Base procedures, SAS macro facility, and Output Delivery System (ODS). Among the modules in Base SAS, this book covers the DATA step and some of the SAS Base procedures that relate to data management.

SAS provides multiple methods for starting and running your SAS program, which depends upon your operating system. The common method for most SAS users is to utilize the SAS Windowing Environment. Alternatively, you can also run your program by using an interactive or noninteractive line mode, as well as batch mode. Please refer to SAS documentation for these alternative methods.

SAS Windowing Environment (illustrated in [Figure 0.1](#)) consists of five windows where you can create and edit your SAS program and manage your SAS files, which includes Program Editor, Log, Output, Result, and Explorer windows. To create a SAS program in the SAS Windowing Environment, you can use the Program Editor window to write your SAS code. You can either submit your entire program at once or highlight only part of the program that you wish to run by clicking on the "Submit" button on the tool bar or by selecting "Run ◊ Submit" from the tool bar. Once the program is submitted, SAS will display messages such as the name of the newly created data set and error or warning messages in the Log window. If the program generates output, it will be displayed in the Output window. You can use the Result window to view and manage different output that is created from your program. To manage files that are stored in the SAS library, you can utilize the Explorer window.

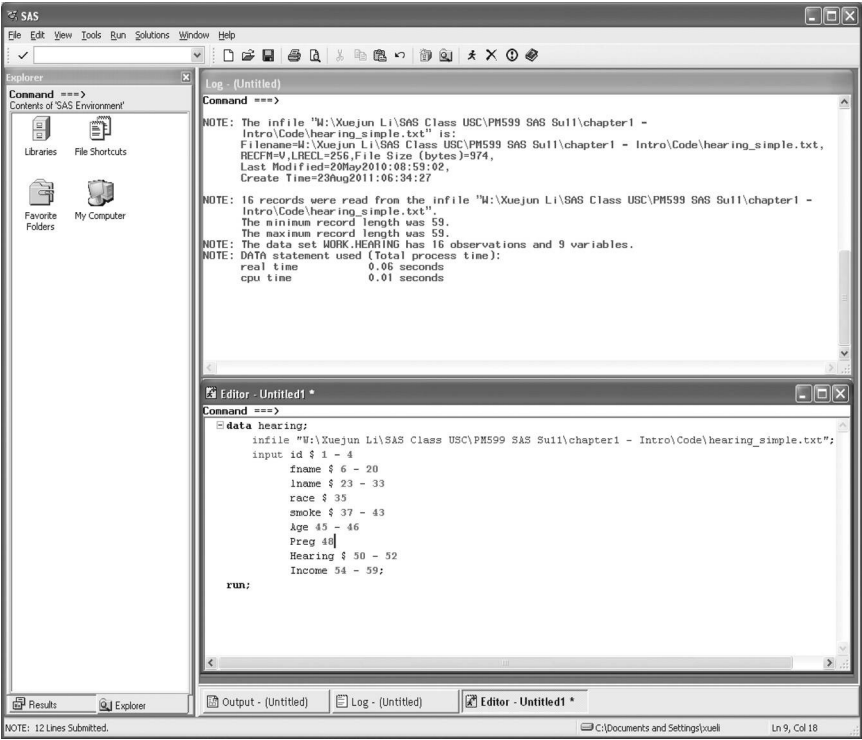


FIGURE 0.1  
SAS windowing environment.

## SAS Help and Documentation

One of the best resources for learning the SAS language and searching for SAS programming help is *SAS Help and Documentation*, which is part of your SAS software installation. To open *SAS Help and Documentation*, click on the “Help” icon on the tool bar, which will bring up the documents.

There are two windows within *SAS Help and Documentation*. The left window, which contains four tabs (Contents, Index, Search, and Favorites), can help you navigate the different topics. The contents of the selected topic are displayed on the right side of the window. The SAS documentations are grouped by categories under the Contents tab. The Contents tab can often be more useful compared to the other tabs if you know how each topic is categorized.

Most of the references in this book are based on *SAS Help and Documentation*. Every topic within this book is under the “SAS Products ♦ Base SAS” main branch. Within this main branch, if you are looking for documents for

non-statistical procedures, it will be under “Base SAS 9.2 Procedure Guide.” The topic under “SAS 9.2 Language Reference Concepts” provides detailed documentation about SAS system concepts, SAS DATA steps, and file concepts. Documents that are related to the language elements, such as SAS data set options, formats, functions, informats, statements, and system options, can be found under “SAS 9.2 Language Reference: Dictionary ◊ Dictionary of Language Elements.” This search path for SAS documentation through the HELP menu was introduced with SAS Version 9.2. If you are using another version of SAS, such as 9.3, the documentations might be organized differently. Thus, instead of listing the entire search path for a specific documentation, only the name of the last node, such as “Dictionary of Language Elements,” is provided as a reference in the book. If you cannot locate the referenced documentation from the HELP menu, an alternative and effective method is to perform a Google search on the name of the documentation (just make sure to include “SAS” in the search string).

---

## How Best to Navigate This Book

The contents of this book are grouped into ten chapters. These ten chapters can be grouped into three main sections.

The first section of the book includes Chapter 1 and Chapter 2, which serve as prerequisite reading for the main section of the book. Chapter 1 provides an introduction and overview of the SAS language, which includes reading data into the SAS system, some Base SAS procedures, creating variables, and subsetting data sets. Some of the topics in these chapters will also be expanded in detail in the later chapters. Chapter 2 discusses how to conditionally create variables based on existing variables, such as how best to use the IF-THEN/ELSE statement or the SELECT statement to create variables.

The second section of the book includes Chapters 3 through 6, which is also the core component of the entire book. These four chapters describe how essential it is to understand the PDV in order to write an accurate program in the DATA step. The topic covered in each chapter is built upon the concepts covered in a previous chapter. Chapter 3 provides an overview about DATA step processing, how to retain variables by using the RETAIN and SUM statements, conditional processing in the DATA step, and how to use the DATA step debugger. Chapter 4 covers BY-group processing within the DATA step and some applications relating to longitudinal data sets. Chapter 5 introduces topics on explicit loops in the DATA step and compares the differences between implicit and explicit loops. Chapter 6 covers array processing, which includes one- to multi-dimensional arrays. Many applications that utilize array processing are related to understanding the explicit loop.



The final section of the book includes Chapters 7 to 10. These four chapters are not built in sequence and can be read independently. Chapter 7 covers multiple methods for combining data sets vertically and horizontally. Although some examples in this chapter demonstrate an understanding of the PDV, readers should be able to grasp and understand most of the materials with or without having read the core section of the book. Chapter 8 covers data input and output. This chapter covers different input methods of reading and writing text files in the DATA step. Chapter 9 covers DATA step functions and call routines, which represent one of the many strengths of SAS software. This chapter introduces a few categories of SAS functions or call routines, such as date and time functions, character functions, and functions for converting variable types. Chapter 10 covers some useful SAS procedures that relate to data management, which include the SORT, COMPARE, TRANSPOSE, FORMAT, and OPTIONS procedures.

A novice programmer with minimum SAS background should read this book cover to cover in the order in which the book is presented. An intermediate SAS programmer who is interested in learning the concept of DATA step processing can begin reading Chapters 3 to 6 without reading the first two prerequisite chapters.

---

### A Note on SAS Output and Log Font

The output and logs that are generated from SAS procedures use the SAS Monospace font. However, the font used for SAS output and logs in this book is Lucida Sans Typewriter. The drawback of using the Lucida Sans Typewriter font is that grid lines for the table displays, such as the output from PROC FREQ, do not have the grid line “look.” Here’s an example of output from the PROC FREQ using the Lucida Sans Typewriter font:

*Output using* Lucida Sans Typewriter:

The FREQ Procedure				
Preg	Frequency	Percent	Cumulative Frequency	Cumulative Percent
0	19	63.33	19	63.33
1	11	36.67	30	100.00
Frequency Missing = 4				

The “f” symbols in the output above appear to be a straight line in the SAS Monospace font. In this situation, the “f” symbol is replaced with

a dash (“-”) to achieve an easier-to-comprehend display of the output, like the one below:

*Output using Lucida Sans Typewriter with modification:*

The FREQ Procedure				
Preg	Frequency	Percent	Cumulative Frequency	Cumulative Percent
0	19	63.33	19	63.33
1	11	36.67	30	100.00
Frequency Missing = 4				

---

### Additional Components of This Book

Here is a list of materials that readers can download from the publisher’s Web site (<http://www.crcpress.com/product/isbn/9781466552388>):

- All the data sets used in the book, as well as all the programs
- Exercise data sets and their solutions
- Slides for demonstrating the contents of the PDV for some examples in Chapters 3 through 6

---

## *The Author*

---

After receiving an MS in Biostatistics from the University of Southern California (USC), **Arthur Li** embarked on a career as a biostatistician at the City of Hope National Medical Center, in Los Angeles County, California. Li is also a part-time statistical programming instructor at USC. He has given numerous presentations and seminars on DATA step programming and statistical analysis using SAS software at SAS conferences throughout the United States.



---

## Acknowledgments

---

This book would not have been possible without the tremendous support of my mentor, Prof. Stanley Azen, who provided me the opportunity to teach my SAS class at the Department of Preventive Medicine at the University of Southern California (USC) for the past five years. The accumulated course materials bear an imprint upon this book. I would also like to express my sincere gratitude to Prof. Jim Gauderman, my first SAS teacher, who taught me much while I was his student and his teaching assistant. I also have to thank Prof. Roberta McKean-Cowdin, who enriched my hands-on SAS experience from the many research projects that were provided by her.

I would also like to thank Xiaolong Li, Han Tun, and Wuchen Zhao for reviewing the contents and testing all the programs in this book.

I would like to acknowledge the enthusiastic and delightful people from the SAS community who provided invaluable support throughout the years for my recognition in the SAS community: MaryAnne DePesquo, Perry Watts, Sunil Gupta, Ron Cody, Kirk Lafler, Peter Eberhardt, and my dear friend Nate Derby. Among them, Sunil Gupta, Ron Cody, and Peter Eberhardt provided many helpful technical suggestions for this book. I am especially deeply indebted to Perry Watts, who spent enormous chunks of her personal time to provide detailed critiques from cover to cover regarding organization, structure, references, and programs/exercises, resulting in a much more precise product.

I am also grateful for the help I received from Taylor & Francis: Rob Calver (Senior Acquisitions Editor), Rachel Holt (Senior Editorial Assistant), Kathryn Everett (Project Coordinator), and Amy Rodriguez (Production Editor).

I would lastly like to express my deepest thanks and gratitude to my beloved Zaccagnino family, especially my “crazy” but lovable in-laws, Dan and Sally Zaccagnino. The Zaccagnino family’s endless support makes my success possible. Finally, this book could not have been completed without my partner, David Zaccagnino, who not only edited the ESL (English as a second language) errors in this book but most importantly provided me with endless patience and encouragement.

