

姚志勇 编著

# SAS 编程与 数据挖掘商业案例

- 从 PDV 角度详尽剖析 Base SAS 常用语句代码及应用
- 数据挖掘理论和商业应用紧密结合
- 原创朴素贝叶斯文本分类和 EM 迭代算法代码
- 三个典型的数据挖掘商业案例分析



信息科学与技术丛书

# SAS 编程与数据挖掘商业案例

姚志勇 编著

机械工业出版社

本书是作者多年来在企业实践工作中的经验总结，详细讲解了使用 SAS 进行商业数据挖掘的方法。其中包含了目前公开出版的诸多 SAS 教材没有的大量实战内容。

本书内容全面、新颖独创、综合性强，适合企业人员使用，也可作为数学、统计学、金融、电子商务、医药等专业的本科生、硕士生学习 SAS 编程和数据挖掘的参考资料。

读者可在 <http://www.cmpbook.com> 下载书中的 SAS 程序及相关测试数据集。

### 图书在版编目 (CIP) 数据

SAS 编程与数据挖掘商业案例 / 姚志勇编著 .—北京：机械工业出版社，  
2010.5  
(信息科学与技术丛书)  
ISBN 978-7-111-30535-4

I .①S… II .①姚… III .①商业统计 - 统计分析 - 应用软件，SAS  
IV .①F712.3

中国版本图书馆 CIP 数据核字 (2010) 第 078039 号

机械工业出版社 (北京市百万庄大街 22 号 邮政编码 100037)

策划编辑：车 忱

责任编辑：车 忱

责任印制：乔 宇

北京机工印刷厂印刷 (北京振兴源印务有限公司装订)

2010 年 5 月第 1 版·第 1 次印刷

184mm × 260mm · 22.25 印张 · 551 千字

0 001—4 000 册

标准书号： ISBN 978-7-111-30535-4

定价： 42.00 元

凡购本书，如有缺页、倒页、脱页，由本社发行部调换

电话服务

网络服务

社服务中心：(010) 88361066

门户网：<http://www.cmpbook.com>

销售一部：(010) 68326294

教材网：<http://www.cmpedu.com>

销售二部：(010) 88379649

封面无防伪标均为盗版

读者服务部：(010) 68993821

# 前言

当前国内的诸多数据挖掘书籍几乎都是基于理论说明，很少深入介绍数据挖掘实践，涉及 SAS 开发的更是少见。因此，从商业应用出发，基于实践而不是基于理论的数据挖掘书籍呼之欲出。本书作者从商业需求出发，以商业人士的眼光来看待企业数据挖掘，并给出大量的商业实践案例。把主流的数据挖掘技术用真实案例来实现是本书出版的初衷，同时为了满足初学者需求，作者也给出了数据挖掘必备的基础编程知识模块。

全书共分两部分。

第一部分是 SAS 编程：

第 1 章和第 2 章主要介绍 SAS 系统和编程基础，同时介绍 SAS 数据处理最核心的内容——数据指针和 PDV 流程。该核心内容贯穿第一部分，是已出版的其他 SAS 图书没有的。

第 3~9 章主要介绍 SAS 的数据处理技术，也是第一部分的主要内容，包括数据集处理、变量处理和观测处理等多种数据处理技术，同时也介绍了循环控制等稍难的内容，重要的是给出了诸多实际案例及商业应用。尽管第 3~9 章从表面上看和诸多已经出版的 SAS 图书没有什么大的不同，但是这些章节最大的亮点是作者对每一个示例和案例从数据指针和 PDV 流程的角度给予了最详细的程序解读，让读者真正读懂程序，而不是停留在程序的表面。

第 10 章是第一部分的难点。作者还是站在商业实践的角度逐一介绍宏最常用的部分，同时也给出了非常详细的程序解读。

第 11 章介绍 SQL 过程。有关内容在国内同类书中都出现过，但是作者独辟蹊径，融合了项目实践中诸多真正有用的语言，同时也给出了诸多开发建议和应注意的问题。

第 12 章介绍数据处理实践。该章共包括四个方面的内容，几乎都是目前国内没有出现过的，如 HASH 对象及商业应用、正则表达式等。随机抽样也是数据处理经常面临的问题，这里作者开发了在 SAS 系统中如何处理分层不等比例抽样的代码，这也是目前国内其他 SAS 图书没有介绍过的。

第二部分是数据挖掘商业案例：

第 13 章主要介绍数据挖掘概念和流程。数据挖掘流程尤其是商业流程是本章的重点。该流程告诉读者一个真正的商业数据挖掘流程在商业环境中是如何实施的。

第 14 章重点介绍响应模型。响应模型是商业实践中最常用的预测模型，基于第 13 章的流程规范给出了一个具体的商业案例研究。

第 15 章是客户行为分析。该章有目前全球最流行的行为分析，包括“行为年龄”和“行为性别”（注意完全不同于具有自然属性特征的“真实年龄”和“真实性别”），作者运用 Naïve Bayesian 技术开发出一整套模型，并对该模型拥有完全自主知识产权。

第 16 章介绍文本挖掘。该章首先介绍了文本挖掘的流程，然后开发出基于 Naïve Bayesian 文本分类算法和 EM 迭代思想的大型代码，并成功应用于商业实践。

本书特色如下：

(1) 国内少有的由商界人士编写的 SAS 数据挖掘图书。书中有近一半的内容是目前国

## 出版说明

随着信息科学与技术的迅速发展，人类每时每刻都会面对层出不穷的新技术和新概念。毫无疑问，在节奏越来越快的工作和生活中，人们需要通过阅读和学习大量信息丰富、具备实践指导意义的图书来获取新知识和新技能，从而不断提高自身素质，紧跟信息化时代发展的步伐。

众所周知，在计算机硬件方面，高性价比的解决方案和新型技术的应用一直备受青睐；在软件技术方面，随着计算机软件的规模和复杂性与日俱增，软件技术不断地受到挑战，人们一直在为寻求更先进的软件技术而奋斗不止。目前，计算机在社会生活中日益普及，随着 Internet 延伸到人类世界的方方面面，掌握计算机网络技术和理论已成为大众的文化需求。由于信息科学与技术在电工、电子、通信、工业控制、智能建筑、工业产品设计与制造等专业领域中已经得到充分、广泛的应用。所以这些专业领域中的研究人员和工程技术人员越来越迫切需要汲取自身领域信息化所带来的新理念和新方法。

针对人们了解和掌握新知识、新技能的热切期待，以及由此促成的人们对语言简洁、内容充实、融合实践经验的图书迫切需要的现状，机械工业出版社适时推出了“信息科学与技术丛书”。这套丛书涉及计算机软件、硬件、网络和工程应用等内容，注重理论与实践的结合，内容实用、层次分明、语言流畅，是信息科学与技术领域专业人员不可或缺的参考书。

目前，信息科学与技术的发展可谓一日千里，机械工业出版社欢迎从事信息技术方面工作的科研人员、工程技术人员积极参与我们的工作，为推进我国的信息化建设作出贡献。

内没有出现过的。同时克服了诸多国内 SAS 书籍只翻译（如 SAS 帮助文档）而不予解释，或者即便解释也只是停留在程序表面的严重弊端，从程序的内在运行流程并基于作者多年的商业实践给予程序代码更详细到位的解析。

（2）国内首次用 SAS 编程开发 Naïve Bayesian 分类算法和 EM 迭代（原创代码）。

（3）国内首次把数据挖掘理论和商业应用相结合，给出三个非常典型的数据挖掘项目案例，并配有大型的算法代码。

（4）国内首次介绍非常实用的抽样技术。

（5）国内首次介绍 SAS 的一个重要组件——HASH 对象，并给出商业实践案例。

本书按照惯例，SAS 的关键字在代码中使用小写，而在正文中使用大写。

本书适合多层次多专业的人士阅读，如数学、统计学、经济学、保险和商业管理等专业的本科生、研究生及相关从业人员。希望这是一本让 SAS 用户和数据挖掘工作者都非常喜欢的书。

限于作者的水平，并且由于数据挖掘涉及领域非常广泛，所以本书只是把目前主流的数据挖掘付诸实践，今后如有机会，将不断更新和改进。

姚志勇

# 目 录

## 出版说明

## 前言

<b>第1章 SAS 系统简介</b>	1	2.6.2 变量列表	21
1.1 系统简介	1	2.6.3 自动变量	21
1.1.1 SAS 系统与商务智能系统	1	<b>第3章 数据获取与数据集操作</b>	23
1.1.2 SAS 系统与其他数据库的 数据交换	1	3.1 数据获取	23
1.1.3 SAS 语言与 SAS 系统	2	3.1.1 LIBNAME 方式	23
1.1.4 SAS 9 浏览窗口简介	3	3.1.2 PASSTHROUGH 方式	24
1.2 一个简单的编程实例	4	3.1.3 IMPORT 方式	25
1.2.1 编写一个 SAS 程序	4	3.1.4 INPUT 方式	25
1.2.2 提交一个 SAS 程序	5	3.2 SET 语句	26
1.2.3 保存和打开一个 SAS 程序	6	3.2.1 语法说明	26
1.3 DATA 步的数据指针和 PDV 流程	6	3.2.2 实例详解	26
1.3.1 数据指针和 PDV 流程	6	3.2.3 商业实践	33
1.3.2 DATA 步执行次数	7	3.3 BY 语句	36
<b>第2章 SAS 编程基础</b>	9	3.3.1 语法说明	36
2.1 SAS 逻辑库	10	3.3.2 实例详解	36
2.1.1 创建 SAS 逻辑库	10	3.4 MERGE 语句	38
2.1.2 删除 SAS 逻辑库	12	3.4.1 语法说明	38
2.1.3 永久逻辑库和临时逻辑库	13	3.4.2 实例详解	38
2.2 SAS 数据集	13	3.5 UPDATE 语句	41
2.2.1 SAS 数据集命名规则	13	3.5.1 语法说明	41
2.2.2 永久 SAS 数据集和临时 SAS 数据集	13	3.5.2 实例详解	41
2.2.3 SAS 数据集结构	15	3.6 MODIFY 语句	42
2.2.4 SAS 数据集形式	16	3.6.1 语法说明	42
2.3 SAS 索引	17	3.6.2 实例详解	44
2.3.1 创建索引	17	3.6.3 商业实践	47
2.3.2 删除索引	18	3.7 PUT 语句	49
2.4 SAS 目录	18	3.7.1 语法说明	49
2.5 数据字典	18	3.7.2 实例详解	51
2.6 SAS 变量	21	3.7.3 商业实践	54
2.6.1 变量属性	21	3.8 FILE 语句	55
		3.8.1 语法说明	56
		3.8.2 实例详解	57
		3.8.3 商业实践	59

3.9 INFILE 语句 .....	60	5.4.3 REPLACE、REMOVE 与 OUTPUT 应用 .....	95
3.9.1 语法说明 .....	60	5.5 DELETE 语句与 STOP 语句 .....	96
3.9.2 实例详解 .....	61	5.5.1 DELETE 语句 .....	96
3.9.3 商业实践 .....	62	5.5.2 STOP 语句 .....	97
<b>第 4 章 SAS 变量操作 .....</b>	<b>64</b>	<b>第 6 章 SAS 数据集管理 .....</b>	<b>98</b>
4.1 赋值语句和累加语句 .....	64	6.1 APPEND 过程 .....	98
4.1.1 赋值语句 .....	64	6.1.1 语法说明 .....	98
4.1.2 累加语句 .....	66	6.1.2 实例详解 .....	100
4.2 KEEP 语句和 DROP 语句 .....	67	6.2 SORT 过程 .....	101
4.2.1 KEEP 语句 .....	67	6.2.1 语法说明 .....	102
4.2.2 DROP 语句 .....	68	6.2.2 实例详解 .....	102
4.3 RETAIN 语句 .....	68	6.2.3 商业实践 .....	103
4.3.1 语法说明 .....	69	6.3 TRANSPOSE 过程 .....	104
4.3.2 实例详解 .....	69	6.3.1 语法说明 .....	104
4.3.3 商业实践 .....	70	6.3.2 实例详解 .....	105
4.4 ARRAY 语句 .....	75	6.4 CONTENTS 过程 .....	107
4.4.1 语法说明 .....	75	6.4.1 语法说明 .....	107
4.4.2 实例详解 .....	77	6.4.2 实例详解 .....	107
4.4.3 商业实践 .....	77	6.5 DATASETS 过程 .....	108
4.5 其他语句 .....	84	6.5.1 语法说明 .....	108
4.5.1 RENAME 语句 .....	84	6.5.2 实例详解 .....	110
4.5.2 LENGTH 语句 .....	85	<b>第 7 章 DATA 步循环与控制 .....</b>	<b>112</b>
4.5.3 LABEL 语句 .....	86	7.1 IF-THEN/ELSE 语句与 SELECT 语句 .....	112
<b>第 5 章 SAS 观测值操作 .....</b>	<b>87</b>	7.1.1 IF-THEN/ELSE 语句 .....	112
5.1 OUTPUT 语句 .....	87	7.1.2 SELECT 语句 .....	115
5.1.1 语法说明 .....	87	7.2 DO 语句 .....	118
5.1.2 实例详解 .....	88	7.2.1 DO 组语句 .....	118
5.2 子集 IF 语句 .....	89	7.2.2 DO 循环语句 .....	119
5.2.1 语法说明 .....	89	7.2.3 DO WHILE 语句 .....	121
5.2.2 实例详解 .....	90	7.2.4 DO UNTIL 语句 .....	121
5.2.3 子集 IF 与 OUTPUT 语句比较 .....	90	7.2.5 DO OVER 语句 .....	122
5.3 WHERE 语句 .....	92	7.2.6 商业实践 .....	123
5.3.1 语法说明 .....	92	7.3 各种控制语句 .....	127
5.3.2 实例详解 .....	92	7.3.1 GO TO 语句 .....	127
5.3.3 子集 IF 与 WHERE 语句比较 .....	94	7.3.2 CONTINUE 语句与 LEAVE 语句 .....	128
5.4 REPLACE 语句和 REMOVE 语句 .....	94		
5.4.1 REPLACE 语句 .....	94		
5.4.2 REMOVE 语句 .....	94		



7.3.3 RETURN 语句 .....	129
<b>第 8 章 常用全程语句</b> .....	131
8.1 COMMENT 语句 .....	131
8.2 X 语句 .....	131
8.3 FILENAME 语句 .....	132
8.4 %INCLUDE 语句 .....	134
8.5 TITLE 语句 .....	135
8.6 FOOTNOTE 语句 .....	136
<b>第 9 章 输出控制</b> .....	137
9.1 LOG 窗口输出控制 .....	137
9.2 OUTPUT 窗口输出控制 .....	138
9.3 常用 ODS 输出控制 .....	138
9.3.1 ODS LISTING .....	139
9.3.2 ODS RESULTS .....	141
9.3.3 ODS TRACE .....	142
9.3.4 ODS OUTPUT .....	144
9.3.5 ODS HTML .....	146
9.3.6 ODS CSVALL .....	148
9.3.7 ODS SELECT .....	149
9.3.8 ODS EXCLUDE .....	151
<b>第 10 章 SAS 宏变量</b> .....	152
10.1 宏运行的内在机制 .....	152
10.2 宏变量 .....	154
10.2.1 定义宏变量 .....	154
10.2.2 显示宏变量 .....	155
10.2.3 引用宏变量 .....	155
10.3 宏程序 .....	157
10.3.1 定义宏 .....	158
10.3.2 调用宏 .....	158
10.3.3 宏内宏 .....	158
10.3.4 宏存储 .....	158
10.4 宏参数 .....	159
10.4.1 创建参数 .....	159
10.4.2 参数赋值 .....	160
10.5 宏函数 .....	160
10.5.1 通配函数 .....	160
10.5.2 计算函数 .....	162
10.5.3 字符函数 .....	163
10.5.4 引用函数 .....	165
10.6 宏语句 .....	166
10.6.1 %IF-%THEN/%ELSE 语句 .....	167
10.6.2 %DO 组语句 .....	168
10.6.3 %DO 循环语句 .....	168
10.6.4 %DO%WHILE 循环语句 .....	169
10.6.5 %DO%UNTIL 循环语句 .....	170
10.7 宏应用 .....	171
10.7.1 创建宏变量的八种方法 .....	171
10.7.2 宏程序一般应用 .....	173
10.7.3 宏程序高级应用 .....	174
<b>第 11 章 SQL 过程</b> .....	178
11.1 单表操作 .....	178
11.2 多表操作 .....	180
11.2.1 多表关联 .....	180
11.2.2 子查询 .....	182
11.2.3 合并查询 .....	183
11.2.4 MERGE 与 SQL 比较 .....	184
11.3 创建、更新与删除表操作 .....	187
11.3.1 创建表 .....	187
11.3.2 行操作 .....	188
11.3.3 列操作 .....	190
11.3.4 删除表 .....	191
11.4 使用 SQL 注意的几个问题 .....	191
<b>第 12 章 数据处理实践</b> .....	192
12.1 随机抽样 .....	192
12.1.1 简单无重复随机抽样 .....	192
12.1.2 分层等比例随机抽样 .....	193
12.1.3 分层不等比例随机抽样 .....	194
12.1.4 随机抽样 MACRO .....	196
12.2 HASH 对象 .....	200
12.2.1 HASH 对象的引例 .....	201
12.2.2 HASH 对象的语法 .....	203
12.2.3 HITER 对象的引例 .....	204
12.2.4 HITER 对象的语法 .....	205
12.2.5 商业实践 .....	205
12.3 FORMAT 综述 .....	210
12.3.1 PROC 步创建 .....	210
12.3.2 DATA 步创建 .....	211
12.3.3 永久存储及调用 .....	212

12.4 正则表达式 .....	213	14.3 模型开发 .....	271
12.4.1 语法说明 .....	214	14.3.1 全模型法选择所有候选 模型 .....	271
12.4.2 常用函数 .....	215	14.3.2 逐步回归法筛选候选模型 .....	272
12.4.3 实例详解 .....	216	14.3.3 创建两个重要数据集 .....	273
12.5 宏在 SAS 与 Excel 转换 中的应用 .....	220	14.3.4 创建 LIFT 图 .....	274
12.5.1 SAS 数据集转换成 Excel .....	220	14.3.5 创建评分卡文件 .....	278
12.5.2 Excel 转换成 SAS 数据集 .....	221	14.4 模型验证 .....	279
<b>第 13 章 数据挖掘概念、任务和     流程 .....</b>	<b>223</b>	14.4.1 评分卡文件导入 .....	280
13.1 数据挖掘概念 .....	223	14.4.2 LIFT 图比较 .....	281
13.2 数据挖掘任务 .....	224	14.4.3 模型确认 .....	283
13.3 数据挖掘流程 .....	225	14.5 模型实施与监控 .....	283
13.3.1 定义商业目标 .....	225	14.5.1 模型实施 .....	283
13.3.2 编制需求文档 .....	228	14.5.2 模型监控 .....	284
13.3.3 选择数据源 .....	231	14.6 小结 .....	285
13.3.4 建模流程图 .....	232	<b>第 15 章 行为建模：客户行为     属性分析 .....</b>	<b>286</b>
13.4 LOGISTIC 建模及结果 详解 .....	233	15.1 前期准备 .....	286
13.4.1 数学模型 .....	233	15.1.1 商业需求 .....	286
13.4.2 参数估计 .....	234	15.1.2 定义目标 .....	286
13.4.3 模型评价指标 .....	235	15.1.3 选择建模方法 .....	288
13.4.4 回归系数 .....	237	15.2 数据获取与处理 .....	288
13.4.5 变量筛选方法 .....	238	15.3 模型开发 .....	294
13.4.6 应用举例及输出结果详解 .....	239	15.4 模型验证 .....	296
13.4.7 多值 LOGISTIC 模型 .....	242	15.5 模型打分 .....	296
<b>第 14 章 响应模型：定位新客户 .....</b>	<b>244</b>	15.6 模型预测 .....	298
14.1 前期准备 .....	244	15.7 模型实施 .....	301
14.1.1 商业需求 .....	245	15.8 小结 .....	302
14.1.2 定义目标 .....	245	<b>第 16 章 文本挖掘：Web 文本     分析 .....</b>	<b>303</b>
14.1.3 选择变量 .....	245	16.1 文本挖掘概念与流程 .....	303
14.2 数据获取与数据处理 .....	246	16.1.1 文本挖掘概念 .....	303
14.2.1 创建建模数据集 .....	248	16.1.2 文本挖掘流程 .....	303
14.2.2 变量首次筛选 .....	249	16.2 商业案例 .....	308
14.2.3 数据探索 .....	252	16.2.1 商业需求 .....	308
14.2.4 数据清洗 .....	254	16.2.2 建模框架设计 .....	308
14.2.5 变量二次筛选 .....	259	16.2.3 结合朴素贝叶斯文本分类的 EM 迭代 .....	309
14.2.6 变量三次筛选 .....	266	16.2.4 数据获取与数据预处理 .....	313
14.2.7 字符变量压缩 .....	269		



16.2.5 文本特征化	318
16.2.6 模型开发：产生文本分类器	321
16.2.7 模型验证：测试分类器效果	336
16.2.8 模型评估：计算混淆矩阵	342
16.2.9 模型应用：对用户查询关键字进行分类	343
16.2.10 小结与展望	343
参考文献	345

# 第1章 SAS 系统简介

## 1.1 系统简介

SAS 系统是一个面向对象的、跨平台的、模块化和组合化的应用软件系统，具有完备的数据存取、数据管理、数据分析和数据展现功能，特别是在数据处理和统计分析领域，SAS 系统至今仍然被誉为全球标准软件系统。

### 1.1.1 SAS 系统与商务智能系统

随着商务智能的日益普及，传统的 SAS 系统也已经进化到与商务智能和数据挖掘保持同步，基于商务智能的商务解决方案已成为 SAS 产品的重要组成部分。传统的 SAS 系统已经演化成服务端的一个模块基础，SAS BI 通过包含 SAS 基础在内的所有服务模块，并以 Web 等作为中间媒介，最终提供给客户包括联机分析系统 OLAP、客户端分析软件 EG 和客户端报表工具 Portal 等在一揽子商务解决方案，如图 1-1 所示。

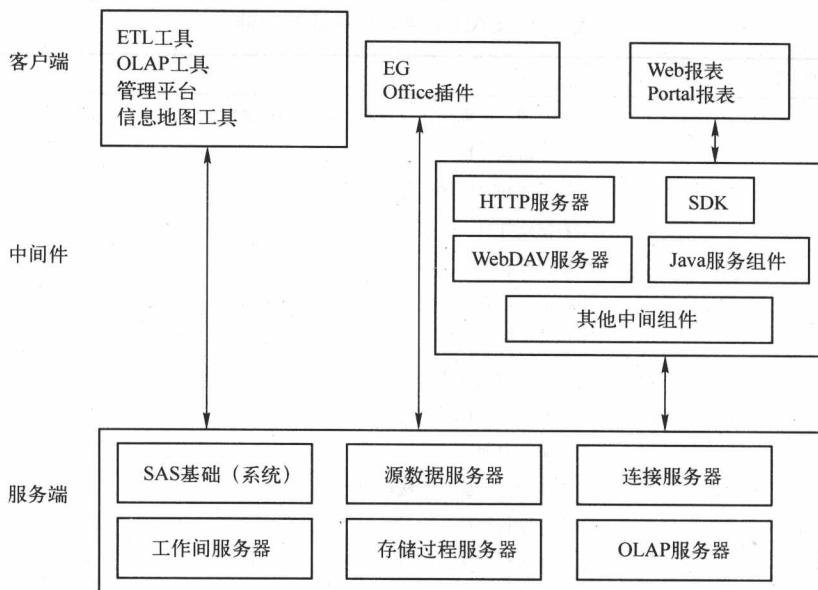


图 1-1 SAS 智能商务平台

本书内容主要集中在 SAS 基础（系统）模块，并未涉及其他任何模块。

### 1.1.2 SAS 系统与其他数据库的数据交换

SAS 系统提供了多种与现行数据库（如主流数据库 TERADATA、DB2、Oracle）的接



口，以及其他应用类软件（如 SPSS 等），如图 1-2 所示。所以在数据获取上除了自身产生的数据外，还可以通过与现行数据库的接口实现从其他数据库获取数据。

在商业实践中，从外部数据库获取数据的频率要远远大于 SAS 本身产生数据的频率，主要原因是商业数据往往都是海量的，一般情况下，需要处理的观测数都会在百万级以上。而 SAS 本身虽然具有建立逻辑数据库的机制，但是与真正意义上的数据库还存在一定差距。比如，对数据用户登录机制、查询优化机制、索引机制、数据存储机制、安全机制等都没有主流数据库（如 DB2、TERADATA 等）高效。因此，提供与这些主流数据库的接口是 SAS 系统处理大型数据的必然选择。

### 1.1.3 SAS 语言与 SAS 系统

SAS 语言是 SAS 系统的基础，是用户和系统直接对话的语言，其特点是用户不必告诉 SAS “怎么做”，只要告诉它“做什么”就行了。

SAS 语言结构主要由 DATA 和 PROC 两个基本步骤任意组合而成。其中，DATA 步完成对数据的获取、加工和处理；PROC 步用于数据分析和输出报告。

SAS 系统包含很多软件模块，见表 1-1。

表 1-1 SAS 主要软件模块及功能

软件模块	功能
BASE	数据管理、基础统计、报表生成和图形显示
STAT	统计分析软件包
ETS	计量经济和时间序列分析软件包
OR	运筹学软件包
QC	质量控制软件包
IML	矩阵语言运算软件包
GRAPH	图形软件包
FSP	数据处理的交互式菜单系统
AF	应用开发软件包
ASSIST	菜单驱动界面
ACCESS	与其他数据库接口的接口集
EIS	应用开发软件包
INSIGHT	可视化数据探索软件包
SHARE	进行数据库并发式控制的软件包
CONNECT	分布处理不同操作平台上的 SAS 系统

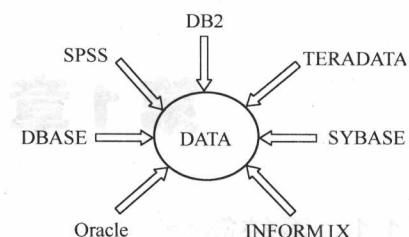


图 1-2 数据库交换系统

SAS 商务智能所有的解决方案及其产品都是以上面所列软件模块作为基础和核心的。而其中的 BASE 又是 SAS 基础（系统）的核心，主要功能是数据管理和数据处理，并有报表输出和基本统计量输出功能。STAT 是 SAS 的统计分析软件算法包，包括方差分析、回归分析、多变量分析、非参数分析等共 50 多个 PROC 步，每个 PROC 步提供多种不同的算法选

择，从而组成一个完整的统计分析方法集。历经 40 多年，至今该软件仍然是全球统计分析领域的领头羊。

BASE 模块和 STAT 模块是本书重点介绍的内容。

### 1.1.4 SAS 9 浏览窗口简介

SAS 9 是一个集成了几乎所有 SAS 软件模块的集成软件。在 Windows 操作系统环境下，如果安装了 SAS 软件，就可以在开始程序菜单里面看到快捷方式。单击快捷方式后，打开的主界面包括浏览器窗口、日志窗口、编程窗口以及通用的菜单栏和工具条，如图 1-3 所示（SAS 9.2）。

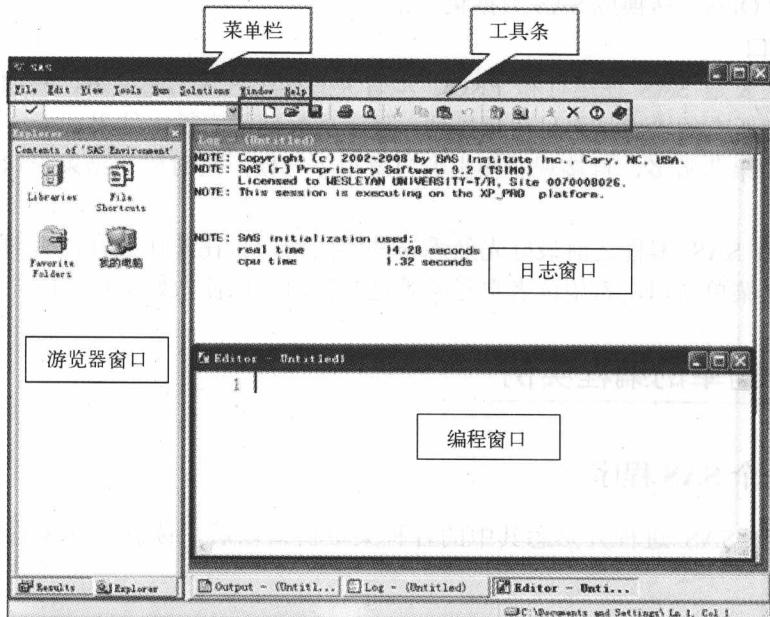


图 1-3 SAS 操作界面

#### 1. 菜单栏

菜单栏由多个子菜单组成，不同的子菜单实现不同的功能。所以，激活不同的子菜单，窗口菜单的内容也会不同。与菜单栏功能相同的是弹出菜单，只要在当前子窗口右击鼠标即可。

#### 2. 工具条

工具条包含了若干按钮，每一个按钮代表一个常用的命令，如保存程序、运行程序等。用户可以在菜单栏 Tools 子菜单中的 Customize 中设定工具条，非常简单。

#### 3. 日志窗口

日志窗口可显示程序运行过程中的所有信息。例如，图 1-3 中日志窗口的文字就是在 SAS 启动后看到的日志信息。日志窗口对程序的调试起着非常重要的作用，因为程序员在开发程序的过程中，很少有一次就能开发所有的代码并运行成功的，一般都需要调试若干次才获得通过，这时所有调试出现的信息包括程序出错的信息都出现在日志窗口中。通过日志窗口，一般就能知道程序出错的原因，并找出相应的 Bug。

#### 4. 浏览器窗口

单击图 1-3 中左下角的 Explorer 按钮，就可以显示对应的浏览器窗口，包含四个对象，其中最常用的是 Library 逻辑库对象。

#### 5. 编程窗口

这是最重要的人机交互界面。一般来说，用户所有的开发代码都是在这个窗口完成和调试的。

#### 6. 输出窗口

输出窗口主要显示 SAS 在运行和 PROC 步有关的程序时输出的结果，如包含一些统计量的输出，一些报表的输出等。后面将专门提到，输出窗口的输出结果还可以通过 SAS 的文件输出系统（ODS）转换成 SAS 数据集等格式文件。

#### 7. 结果窗口

结果窗口主要是 SAS 在运行和 PROC 步有关的程序时输出的结果，和输出窗口不同的是，结果窗口把所有的输出结果呈树状排列，便于更好地检索输出结果。

如果输出结果非常多，直接观看输出窗口就很不方便，这时候在结果窗口中检索输出会更加便捷。

读者在熟悉 SAS 编程之前最好先熟悉这些界面窗口。在此作者并没有对这些窗口做详细介绍，特别是菜单窗口，希望读者自行熟悉这些窗口，以便做好编程前的充分准备工作。

## 1.2 一个简单的编程实例

### 1.2.1 编写一个 SAS 程序

在打开一个 SAS 进程并熟悉其中的各种菜单窗口以后，读者可在编程窗口输入如下代码：

```
data test;
input x y z;
cards;
1 2 3
4 5 6
;
run;
```

这段程序仅包含七行代码：

第 1 行：产生数据集的 DATA 步语句。后面的 test 是要产生的 SAS 数据集。

第 2 行：输入变量的 INPUT 语句。该语句表明 test 数据集包含三个变量 x、y 和 z。

第 3 行：输入数据的 CARDS 语句。其后接具体的输入数据。

第 4, 5 行：数据块。是要输入的具体数据。

第 6 行：分号。分号是 SAS 不同语句之间间隔的标志。

第 7 行：RUN 语句。表明到此为止就可以提交上面六行语句。

## 1.2.2 提交一个 SAS 程序

有四种方式可以运行该语句：

第1种：单击工具条里面的小人图标。

第2种：单击菜单栏 RUN 下面的 SUBMIT。

第3种：按热键〈F3〉。

第4种：保存上面的代码到某个盘符目录下，然后右键单击文件名，可以通过批处理方式提交，此时选择 Batch Submit With SAS 9.2；当然也可以交互式提交，此时选择 Submit With SAS 9.2。批处理提交和直接提交的区别在于：前者是全部提交所有的 SAS 代码，而后者是逐行提交，如果没有遇到 RUN 语句，进程会一直停留在运行状态，并等待用户的下一个命令。所以，一般情况下，除非在所有代码都调试成功后才用批处理命令，否则交互式命令是常用的选择。

按〈F3〉键运行上面的代码后，会发现日志窗口有如下显示：

---

```

1  data test;
2  input x y z;
3  cards;
NOTE: The data set work.test has 2 observations and 3 variables.
NOTE: data statement used (Total process time):
      real time          0.01 seconds
      cpu time           0.01 seconds
5  ;
6  run;

```

---

前3行是代码，第4行的 NOTE 说明有一个数据集 Work.test，在什么地方呢？有两种方式可以找到该数据集：

第1种：双击浏览窗口里面的 Libraries，进去以后会发现有四个逻辑库（关于什么是逻辑库，在第2章有详细说明，读者可暂时把它理解为类似于一个 Windows 系统下的目录），其中有个叫 Work 的逻辑库正是要寻找的对象，双击该逻辑库，就会发现里面藏有 test 文件集。双击该文件集，结果会显示该文件集有三个变量 x、y 和 z，对应值正是刚才输入的三个变量值。

第2种：单击工具条中的 SAS Explorer 图标，在浏览器窗口中会呈现一个树状的结构，如图 1-4 所示。

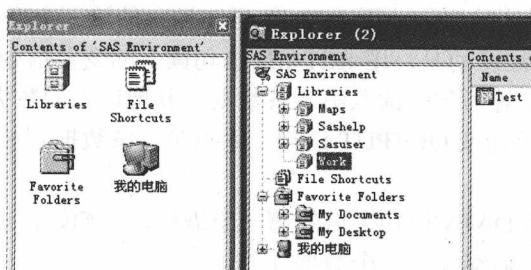


图 1-4 浏览器窗口

读者注意到，其实这个浏览器窗口就是图 1-3 中左边的浏览器窗口，只不过呈树状结构，看起来更方便。双击其中的 Work 逻辑库，在最右边也能看到 Test 数据集。

接下来的一个 NOTE 观测了该程序运行的时间。一般情况下，在项目实施过程中，会把这些时间观测输出，以便调试和优化现有的程序。

### 1.2.3 保存和打开一个 SAS 程序

刚才编写的这个程序在编程窗口中的名字仍然是 Untitled1，如果需要保存这个程序，有两种方式：

第 1 种：单击工具条的保存按钮，保存到指定目录。

第 2 种：单击菜单栏 File 菜单，选择 Save，保存到指定目录。

需要注意的是，SAS 程序文件的扩展名是 SAS，但可以用任何文本文件名作为扩展名，如 TXT，因为 SAS 程序文件实际上是一个文本文件。

如果要打开一个已经保存的 SAS 程序文件，只要选择 File 菜单下的 Open 选项，打开指定的文件即可。

## 1.3 DATA 步的数据指针和 PDV 流程

### 1.3.1 数据指针和 PDV 流程

SAS 语言是先编译后执行的语言，每一个 DATA 步开头标志着编译的开始，接下来的所有语句都被称为“程序”，SAS 总是把这些语句转换成计算机可识别的语言，然后通过逐行语句编译、转换，最后输出成 SAS 数据集。在理解一个 DATA 步所需要走过的正常流程之前，读者需要理解两个概念：

第一个概念是数据指针，为简化起见，读者可以把它理解成在当前的内存缓存区，输入数据所在的位置。

第二个概念是 PDV，它是 Program Data Vector 的缩写，也就是在 DATA 步中所有涉及的变量被变成当前向量的一部分。来自输入数据每一个原始数据或数据集的这些变量，和其他程序语句建立的新变量一起，都放在 PDV 里面。

下面以刚才产生的 Test 数据集为例讲解上面的流程。

第 1 步：编译程序，开辟内存空间。在 Test 数据集中，七行程序首先被编译，并在内存中开辟一块变量空间给三个变量 x、y 和 z。

第 2 步：初始化变量为缺失值。用 INPUT 语句输入的这三个变量在 DATA 步执行之前被置为缺失值。但是需要注意的是，用 RETAIN 语句读入的变量值则被保留。

第 3 步：执行 INPUT 等语句。读入第一条数据，并把读入的数据放入 PDV。

第 4 步：执行 RUN 语句或 OUTPUT 语句。输出第一条数据到被创建的 DATA 步后面的数据集中。

第 5 步：SAS 返回到 DATA 语句之后的第一个语句，并初始化 PDV 中的非 RETAIN 变量 x、y 和 z 为缺失值，然后继续下一个数据读入。

第 6 步：当 SAS 读完并处理完来自 INPUT 语句后面的 1;2;3 和 4;5;6 后，程序到了输入