



# Java 语言与 XML 处理教程

SAX, DOM, JDOM, JAXP 与 TrAX 指南

Processing XML with Java

A Guide to SAX, DOM, JDOM, JAXP and TrAX

[美] Elliotte Rusty Harold 著

刘文红 赵伟明 等译



电子工业出版社  
Publishing House of Electronics Industry  
<http://www.phei.com.cn>

Java 技术丛书

# Java 语言与 XML 处理教程

SAX, DOM, JDOM, JAXP 与 TrAX 指南

Processing XML with Java

A Guide to SAX, DOM, JDOM, JAXP and TrAX

[ 美 ] Elliotte Rusty Harold 著

刘文红 赵伟明 等译

电子工业出版社

Publishing House of Electronics Industry

北京 · BEIJING

## 内 容 简 介

本书是为要在系统中集成 XML 的 Java 编程人员编写的，是介绍如何使用 Java 编程语言处理 XML 文档的实用而且综合的指南与教程。书中简要概述了 XML 基础，包括 XML 语法、DTD、模式、有效性，样式单和 XML 协议 XML-RPC、SOAP 与 RSS。本书的核心内容是深入介绍了 Java 编程人员用 Java 生成与操纵 XML 文档时所用的关键 XML API，包括 SAX、DOM（文档对象模型）和 JDOM。此外，还介绍了这些核心 API 的许多重要补充，包括 XPath、XSLT、TrAX 与 JAXP。

本书详尽介绍了实用和面向任务的方法，是所有需要使用 XML 的 Java 编程人员的宝贵参考资料。

Simplified Chinese edition Copyright © 2003 by PEARSON EDUCATION ASIA LIMITED and Publishing House of Electronics Industry.

Processing XML with Java: A Guide to SAX, DOM, JDOM, JAXP and TrAX, ISBN: 0201771861 by Elliotte Rusty Harold. Copyright © 2003. All Rights Reserved.

Published by arrangement with the original publisher, Pearson Education, Inc., publishing as Addison-Wesley.

This edition is authorized for sale only in the People's Republic of China (excluding the Special Administrative Region of Hong Kong and Macau).

本书中文简体字翻译版由电子工业出版社和 Pearson Education 培生教育出版亚洲有限公司合作出版。未经出版者预先书面许可，不得以任何方式复制或抄袭本书的任何部分。

本书封面贴有 Pearson Education 培生教育出版集团激光防伪标签，无标签者不得销售。

版权贸易合同登记号：图字：01-2003-1045

## 图书在版编目 (CIP) 数据

Java 语言与 XML 处理教程：SAX, DOM, JDOM, JAXP 与 TrAX 指南 / (美) 哈罗德 (Harold, E. R.) 著；刘文红等译。—北京：电子工业出版社，2003.11  
(Java 技术丛书)

书名原文：Processing XML with Java: A Guide to SAX, DOM, JDOM, JAXP and TrAX  
ISBN 7-5053-9277-8

I.J... II.①哈... ②刘... III.①Java 语言 - 程序设计 ②可扩充语言，XML- 程序设计 IV. TP312

中国版本图书馆 CIP 数据核字 (2003) 第 098435 号

责任编辑：窦昊                   特约编辑：赵宏英

印 刷 者：北京兴华印刷厂

出版发行：电子工业出版社 <http://www.phei.com.cn>

北京市海淀区万寿路 173 信箱 邮编：100036

经 销：各地新华书店

开 本：787 × 1092 1/16      印张：42.25 字数：1082 千字

版 次：2003 年 11 月第 1 版      2003 年 11 月第 1 次印刷

定 价：68.00 元

凡购买电子工业出版社的图书，如有缺损问题，请向购买书店调换；若书店售缺，请与本社发行部联系。

联系电话：(010) 68279077。质量投诉请发邮件至 [zlts@phei.com.cn](mailto:zlts@phei.com.cn)，盗版侵权举报请发邮件至 [dbqq@phei.com.cn](mailto:dbqq@phei.com.cn)。

## 前　　言

一天晚上，五位新加入Elephants公司的开发人员坐在一起熟悉公司的订单处理系统。他们钻进主服务器中一大堆 XML 文档的目录中。“这是什么？”小组长激动地问。他们过去谁也没有听说过 XML，因此准备了解手头的文件，看看这个奇怪而奇妙的新技术到底是怎么回事。

第一位开发人员擅长于 Oracle 数据库优化，因此打印了一堆 FileMaker 数据库产生的 FMPXMLRESULT 文档，开始进行分析，这个数据库中存放了所有订单。“这就是 XML，我明白了，也没什么新东西。不难看出，XML 文档不过就是表格”。

“什么，就是表格？”第二位开发人员说，他非常熟悉面向对象理论，手头有一堆 XMI 文档对系统的 UML 框架进行了编码。“即使 Visual Basic 编程人员也能看出，XML 文档不是表格，表格关系中不允许复制，除非这是些奇怪的变种。这些 XML 文档明摆着就是类和对象。事实上，这是一目了然的。XML 文档是对象，DTD 是类”。

“对象？真是奇怪的对象”，第三位开发人员说，他是一位具有一定知名度的 Web 设计人员，把订单处理系统的 XHTML 用户文档装入 Mozilla 中。“我没有看到任何类型。如果你认为这是对象，我拒绝安装你的软件。但从这么多样式单不难看出，XML 就是更新的 HTML”。

“HTML？别开玩笑”，第四位说，他是来自 MIT 的计算机科学教授，曾经致力于用 XSLT 样式单和 Schematron 模式检验所有其他文档。“看看层次结构的整洁嵌套，每个标志与相应标志恰好匹配，HTML 什么时候有这么漂亮。这是个 S 表达式，没什么新东西，1882 年就已经出现这个思想了”。

“S 表达式？”技术作家问，他手头有一大堆用 DocBook 编写的项目文档。“对你们这些博学之士，也许是这样吧。但是我觉得这就像是 FrameMaker MIF 文件。不过熟悉 GUI 好像有点费事”。

他们就这样相互争吵，直到深夜，谁也不服谁，每个人都引经据典，却又都不屑看看别人的证据。事实上，他们可能现在还在争论，你甚至还能听到他们的争吵声。当然，他们的错误在于执意把 XML 归结到他们已经熟悉的技术，而不是把 XML 当做全新的技术。XML 可以存储数据，但不是数据库。XML 可以将对象序列化，但 XML 文档不是对象。XML 可以编写 Web 页面，但 XML 不是 HTML。XML 可以编写功能性（和其他）编程语言，但 XML 本身不是编程语言。XML 可以写书，但 XML 不是桌面出版软件。

XML 是个新东西，在计算世界中没有先例。XML 继承了过去的思想，一些固执的人总想通过数据库、对象、功能性或 S 表达式之类的有色眼镜来看 XML，但 XML 并不是这些东西，而是个新东西，是一个在计算世界中没有先例的新概念。要理解 XML，首先要把它看成一个新概念，而不是跨进昨天的旧框框中。

世界上已经有许多工具、API 和应用程序，使 XML 看上去像是开发人员熟悉的某个旧东西——像某种数据库、某种对象或远程过程调用。这些 API 在限制性的可预测环境中也许有用，

但不适合处理更一般格式的 XML。这些 API 在有限领域中能够顺利工作，但超出人为定义的边界时就行不通了。XML 是可扩展的，但许多为 XML 设计的工具并没有这么强的可扩展性。

本书要介绍如何处理完全一般性的 XML，而不是局限于某一个方面。XML 就是 XML，而不是任何别的东西。我们介绍如何设计程序，处理真正意义上的 XML：有效与无效、混合与非混合、类型与非类型，以及这些特性的不同组合。为了达到这个目标，本书主要考虑不会隐藏 XML 本质的 API。特别地，有三大类 Java API 正确地建模 XML，而不是建模 XML 文档的特定类或 XML 的某个狭窄子集。这些 Java API 分别是：

- SAX（XML 的简单 API，Simple API for XML）
- DOM（文档对象模型，Document Object Model）
- JDOM，Java 本地 API

这些 API 是本书的核心。此外，书中还将介绍基本 API 的几个前提与补充项目，包括：

- XML 语法
- DTD、模式与有效性
- XPath
- XSLT 与 TrAX API
- JAXP 是 SAX，DOM，TrAX 和几个工厂类的结合

由于书中要用几个 XML 应用程序例子演示这些 API，因此还要详细介绍 XML-RPC，SOAP 与 RSS。但是，本书介绍的技术不局限于这三类应用程序。

## 适用读者

本书是为要在系统中集成 XML 的高级 Java 开发人员编写的。Java 是处理 XML 文档的理想语言，特别是其强大的 Unicode 支持使其成为许多早期实现版本的首选语言。因此，用 Java 编写的 XML 工具比用其他任何语言编写的 XML 工具都多，用 Java 编写的开放源代码 XML 工具也比用其他任何语言编写的都多。大多数开发者都用 Java 处理 XML。

本书主要介绍下列内容：

- 从 Java 应用程序中保存 XML 文档
- 读取其他程序产生的 XML 文档
- 搜索、查询与更新 XML 文档
- 将遗留平面数据转换成层次式 XML 格式
- 与发送和接收 XML 数据的网络服务器通信
- 用 DTD、模式和业务规则验证文档
- 组合功能性 XSLT 转换与传统的命令式 Java 代码

本书适用于用 XML 进行任何工作的 Java 开发人员，全面介绍了基础课题和高级课题，是 Java 语言与 XML 处理教程的综合性课程，使不熟悉 XML 的开发人员可以通过学习本书设计出

复杂的 XML 应用程序和分析复杂的 XML 文档。本书的例子中包括各种可能的应用，包括文件格式、数据交换、文档转换、数据库集成，等等。

## 预备知识

本书不是 Java 与 XML 的入门书籍，要求读者对 Java 与 XML 有一定的了解。在 Java 方面，书中随意使用这个语言的高级特性及其类库，不做太多解释。读者至少要熟悉下列内容：

- 面向对象编程，包括继承与多态。
- 包与 CLASSPATH，要习惯于没有 main( )方法的类和不在默认包中的类。
- I/O 包括流、阅读程序和写入程序，要了解 System.out 是 Java 程序中很常用的工具。
- Java Collections API 包括散列表、映射、集合、迭代器和列表。

此外，本书某些地方还用到 SQL 与 JDBC。但这些部分在书中相对独立，如果读者不熟悉 SQL，则通常不需要这些材料。

## 所需工具

XML 是独立于体系结构、平台、操作系统、GUI 和语言的（比 Java 更加独立），在 Mac OS, Windows, Linux, OS/2 和各种 UNIX 系统中都能顺利工作，可以用 Python, C++, Haskell, ECMAScript, C#, Perl, Visual Basic, Ruby 和 Java 等语言访问。在 PowerPC, X86 和其他体系结构之间进行切换时，不需要考虑字节顺序问题。几乎本书的任何内容在任何能运行 Java 的平台上都同样适用。

本书的大部分材料相对独立于特定的 Java 版本。Java 1.4 把 SAX, DOM 和另外几个重要类捆绑到核心 JDK 中，但它们在旧版 JVM 中也很容易安装，因为 Apache XML Project 和其他厂家提供了开放源代码库。通常，书中在测试例子时使用 Java 1.3 和 1.4，因此所用的有些类与方法可能在旧版中没有。但是在大多数情况下它们很容易向后移植。除 TrAX 以外的所有基本 XML API 都适用于 Java 1.1 以上版本，TrAX 适用于 Java 1.2 以上版本。

## 本书用法

本书编写成高级教程，也可以作为综合性的参考手册。第 1 章介绍使用 XML 所需的最基础材料，但不作为综合性引言，而是作为复习材料，要求读者先阅读一些更基础的书籍。第 2 章介绍 RSS, XML-RPC, SOAP 和本书例子中使用的 XML 应用程序。后面两章介绍从程序中产生 XML（这个课题常常被介绍得过于复杂）。第 3 章介绍从代码中直接生成 XML，第 4 章介绍将其他格式的遗留数据转换成 XML。本书其余内容介绍几个主要的 API，用于处理 XML：

- 基于事件的 SAX API
- 基于树的 DOM API
- 基于树的 JDOM API

- 搜索 XML 文档的 XPath API
- 用于 XSLT 处理的 TrAX API

本书最后的附录提供主要 API 的速查手册。

如果读者的 XML 经验不是太丰富，建议至少先依次读完本书前两章。然后如果喜欢特定的 API，可以直接转入相应部分：

- 第 6 章到第 8 章介绍 SAX。
- 第 9 章到第 13 章介绍 DOM。
- 第 14 章到第 15 章介绍 JDOM。

熟悉一个或几个 API 之后，可以学习第 16 章与第 17 章介绍的 XPath 与 XSLT。但是，这些章节和 API 的内容要求至少对这三大 API 之一有一定了解。

## 联机版本

可以从笔者的 Cafe con Leche 站点取得本书的联机版本（原英文版），网址为 <http://www.cafeconleche.org/books/xmljava/>，本书的每一个字都在上面，毫无保留。我希望你喜欢印刷版本，买上一本，比你自己用激光打印机打印 1120 页的原书成本要低得多，但这完全凭你自愿。我的目标是尽量推广本书内容，使其为更多人服务。

联机版本受到版权法的保护，阅读时不需要注册，也不必下载特殊的电子键而在换一台电脑时失效。我希望更多人阅读本书，使用其中的知识。作为礼节，希望读者不要在自己的服务器上重新发表这个联机版本，否则很难让你的副本保持更新。如果你想在手提电脑中保存几页，以便在飞机上阅读，我不会介意，但不要把自己的副本发给别人，而要让你的同事和朋友到我的 Web 站点或者买一本书。

## 联系作者

我喜欢听到读者的意见、建议和与本书有关的问题。由于本书完全联机发表，联机版本可以比传统图书更快改版，因此纠正与勘误特别有用，我可以随时修正。发来纠正与勘误之前，请先看看联机版本是否已经解决了所提的问题。

请将所有意见、建议和问题发到 [elharo@metalab.unc.edu](mailto:elharo@metalab.unc.edu)。由于收到了很多 E-mail，因此不能保证有信必复，但我会尽力而为。最好能在主题中明确表示你是本书读者，否则你的消息很容易被当做垃圾邮件，我可能不加浏览就扔进垃圾箱中。另外，请确保你的邮件使用正确的答复地址，至少在发信后的几个月内保证这个地址有效。最令人失望的是，花好几个小时写了一封信之后，却因为地址不对而被退回。欢迎来信，盼望你的意见、建议和问题。

# 目 录

## 第一部分 XML

|  |           |
|--|-----------|
| <b>第 1 章 XML 与数据 .....</b>               | <b>2</b>  |
| 1.1 XML 的优点 .....                        | 2         |
| 1.2 XML 语法 .....                         | 8         |
| 1.3 有效性 .....                            | 20        |
| 1.4 样式单 .....                            | 27        |
| 1.5 小结 .....                             | 34        |
| <br>                                     |           |
| <b>第 2 章 XML 协议：XML-RPC 与 SOAP .....</b> | <b>36</b> |
| 2.1 XML 消息格式 .....                       | 36        |
| 2.2 HTTP 传输协议 .....                      | 40        |
| 2.3 RSS .....                            | 46        |
| 2.4 自定义请求 .....                          | 48        |
| 2.5 XML-RPC .....                        | 52        |
| 2.6 SOAP .....                           | 61        |
| 2.7 自定义协议 .....                          | 75        |
| 2.8 小结 .....                             | 75        |
| <br>                                     |           |
| <b>第 3 章 使用 XML 与 Java .....</b>         | <b>77</b> |
| 3.1 Fibonacci 数 .....                    | 77        |
| 3.2 编写 XML .....                         | 79        |
| 3.3 输出流、写入器与编码方式 .....                   | 84        |
| 3.4 简单 XML-RPC 客户程序 .....                | 88        |
| 3.5 简单 SOAP 客户程序 .....                   | 90        |
| 3.6 小服务 .....                            | 92        |
| 3.7 小结 .....                             | 95        |
| <br>                                     |           |
| <b>第 4 章 将平面文件转换成 XML .....</b>          | <b>96</b> |
| 4.1 预算 .....                             | 96        |
| 4.2 模型 .....                             | 97        |
| 4.3 输入 .....                             | 98        |
| 4.4 确定输出格式 .....                         | 100       |
| 4.5 从平面数据建立层次结构 .....                    | 111       |
| 4.6 不用 Java 的方法 .....                    | 123       |

|            |                              |            |
|------------|------------------------------|------------|
| 4.7        | 关系型数据库 .....                 | 129        |
| 4.8        | 小结 .....                     | 134        |
| <b>第5章</b> | <b>读取 XML .....</b>          | <b>136</b> |
| 5.1        | InputStreams 与 Readers ..... | 136        |
| 5.2        | XML 分析器 .....                | 139        |
| 5.3        | SAX .....                    | 147        |
| 5.4        | DOM .....                    | 151        |
| 5.5        | JAXP .....                   | 154        |
| 5.6        | JDOM .....                   | 156        |
| 5.7        | dom4j .....                  | 159        |
| 5.8        | ElectricXML .....            | 160        |
| 5.9        | XMLPULL .....                | 162        |
| 5.10       | 小结 .....                     | 164        |

## 第二部分 SAX

|            |                            |            |
|------------|----------------------------|------------|
| <b>第6章</b> | <b>SAX .....</b>           | <b>168</b> |
| 6.1        | SAX 简介 .....               | 168        |
| 6.2        | 分析 .....                   | 169        |
| 6.3        | 回调接口 .....                 | 171        |
| 6.4        | 接收文档 .....                 | 175        |
| 6.5        | 接收元素 .....                 | 177        |
| 6.6        | 处理属性 .....                 | 180        |
| 6.7        | 接收字符 .....                 | 184        |
| 6.8        | 接收处理指令 .....               | 187        |
| 6.9        | 接收名字空间映射 .....             | 189        |
| 6.10       | 可忽略空白符 .....               | 191        |
| 6.11       | 接收跳过的实体 .....              | 192        |
| 6.12       | 接收定位器 .....                | 193        |
| 6.13       | ContentHandler 所缺的信息 ..... | 196        |
| 6.14       | 小结 .....                   | 197        |
| <b>第7章</b> | <b>XMLReader 接口 .....</b>  | <b>198</b> |
| 7.1        | 建立分析器对象 .....              | 198        |
| 7.2        | 输入 .....                   | 200        |
| 7.3        | 异常与错误 .....                | 204        |
| 7.4        | 特性与属性 .....                | 211        |
| 7.5        | DTDHandler .....           | 231        |
| 7.6        | 小结 .....                   | 238        |
| <b>第8章</b> | <b>SAX 过滤器 .....</b>       | <b>240</b> |
| 8.1        | 过滤器体系结构 .....              | 240        |

|     |                       |     |
|-----|-----------------------|-----|
| 8.2 | XMLFilter 接口 .....    | 242 |
| 8.3 | 内容过滤器 .....           | 252 |
| 8.4 | XMLFilterImpl 类 ..... | 268 |
| 8.5 | 分析非 XML 文档 .....      | 271 |
| 8.6 | 多处理器适配器 .....         | 277 |
| 8.7 | 小结 .....              | 283 |

## 第三部分 DOM

|               |                                |            |
|---------------|--------------------------------|------------|
| <b>第 9 章</b>  | <b>文档对象模型 .....</b>            | <b>286</b> |
| 9.1           | DOM 的演变 .....                  | 286        |
| 9.2           | DOM 模块 .....                   | 287        |
| 9.3           | 应用程序特定的 DOM .....              | 290        |
| 9.4           | 树 .....                        | 290        |
| 9.5           | DOM Java 分析器 .....             | 298        |
| 9.6           | 用 DOM 分析器分析文档 .....            | 299        |
| 9.7           | Node 接口 .....                  | 308        |
| 9.8           | NodeList 接口 .....              | 318        |
| 9.9           | JAXP 序列化 .....                 | 319        |
| 9.10          | DOMException .....             | 321        |
| 9.11          | 选择 SAX 与 DOM .....             | 323        |
| 9.12          | 小结 .....                       | 324        |
| <b>第 10 章</b> | <b>用 DOM 建立 XML 文档 .....</b>   | <b>326</b> |
| 10.1          | DOMImplementation .....        | 326        |
| 10.2          | 定位 DOMImplementation .....     | 327        |
| 10.3          | Document 接口作为抽象工厂 .....        | 330        |
| 10.4          | Document 接口作为节点类型 .....        | 339        |
| 10.5          | 规范化 .....                      | 354        |
| 10.6          | 小结 .....                       | 356        |
| <b>第 11 章</b> | <b>DOM 核心 .....</b>            | <b>357</b> |
| 11.1          | Element 接口 .....               | 357        |
| 11.2          | NamedNodeMap 接口 .....          | 365        |
| 11.3          | CharacterData 接口 .....         | 370        |
| 11.4          | Text 接口 .....                  | 373        |
| 11.5          | CDATASection 接口 .....          | 376        |
| 11.6          | EntityReference 接口 .....       | 379        |
| 11.7          | Attr 接口 .....                  | 380        |
| 11.8          | ProcessingInstruction 接口 ..... | 382        |
| 11.9          | Comment 接口 .....               | 385        |
| 11.10         | DocumentType 接口 .....          | 388        |
| 11.11         | Entity 接口 .....                | 389        |

|                         |     |
|-------------------------|-----|
| 11.12 Notation 接口 ..... | 391 |
| 11.13 小结 .....          | 394 |

## 第 12 章 DOM 遍历模块 ..... 396

|                         |     |
|-------------------------|-----|
| 12.1 NodeIterator ..... | 396 |
| 12.2 NodeFilter .....   | 401 |
| 12.3 TreeWalker .....   | 404 |
| 12.4 小结 .....           | 408 |

## 第 13 章 DOM 输出 ..... 410

|                         |     |
|-------------------------|-----|
| 13.1 Xerces 序列化 .....   | 410 |
| 13.2 OutputFormat ..... | 411 |
| 13.3 DOM Level 3 .....  | 417 |
| 13.4 小结 .....           | 425 |

# 第四部分 JDOM

## 第 14 章 JDOM ..... 428

|                             |     |
|-----------------------------|-----|
| 14.1 JDOM 简介 .....          | 428 |
| 14.2 用 JDOM 生成 XML 元素 ..... | 430 |
| 14.3 用 JDOM 建立 XML 文档 ..... | 432 |
| 14.4 用 JDOM 编写 XML 文档 ..... | 432 |
| 14.5 文档类型声明 .....           | 436 |
| 14.6 名字空间 .....             | 438 |
| 14.7 用 JDOM 读取 XML 文档 ..... | 441 |
| 14.8 导航 JDOM 树 .....        | 444 |
| 14.9 与 DOM 程序通信 .....       | 448 |
| 14.10 与 SAX 程序通信 .....      | 449 |
| 14.11 Java 集成 .....         | 452 |
| 14.12 JDOM 的缺点 .....        | 454 |
| 14.13 小结 .....              | 455 |

## 第 15 章 JDOM 模型 ..... 456

|                                    |     |
|------------------------------------|-----|
| 15.1 Document 类 .....              | 456 |
| 15.2 Element 类 .....               | 458 |
| 15.3 Attribute 类 .....             | 478 |
| 15.4 Text 类 .....                  | 481 |
| 15.5 CDATA 类 .....                 | 484 |
| 15.6 ProcessingInstruction 类 ..... | 485 |
| 15.7 Comment 类 .....               | 487 |
| 15.8 名字空间 .....                    | 489 |
| 15.9 DocType 类 .....               | 491 |

|                         |     |
|-------------------------|-----|
| 15.10 EntityRef 类 ..... | 495 |
| 15.11 小结 .....          | 497 |

## 第五部分 XPath/XSLT

|                              |            |
|------------------------------|------------|
| <b>第 16 章 XPath .....</b>    | <b>500</b> |
| 16.1 查询 .....                | 500        |
| 16.2 XPath 数据模型 .....        | 502        |
| 16.3 定位路径 .....              | 504        |
| 16.4 表达式 .....               | 512        |
| 16.5 XPath 引擎 .....          | 517        |
| 16.6 DOM Level 3 XPath ..... | 524        |
| 16.7 Jaxen .....             | 530        |
| 16.8 小结 .....                | 534        |

|                           |            |
|---------------------------|------------|
| <b>第 17 章 XSLT .....</b>  | <b>535</b> |
| 17.1 XML 转换 .....         | 535        |
| 17.2 TrAX .....           | 546        |
| 17.3 用 Java 扩展 XSLT ..... | 565        |
| 17.4 小结 .....             | 575        |

## 第六部分 附录

|                                |            |
|--------------------------------|------------|
| <b>附录 A XML API 速查手册 .....</b> | <b>578</b> |
| <b>附录 B SOAP 1.1 模式 .....</b>  | <b>648</b> |
| <b>附录 C 推荐读物 .....</b>         | <b>662</b> |
| <b>译后记 .....</b>               | <b>664</b> |

# **第一部分**

## **XML**

- 第 1 章 XML 与数据**
- 第 2 章 XML 协议：XML-RPC 与 SOAP**
- 第 3 章 使用 XML 与 Java**
- 第 4 章 将平面文件转换成 XML**
- 第 5 章 读取 XML**

# 第1章 XML 与数据

XML 原先只是“Web 上的 SGML（标准通用标记语言）”，最初和 SGML 与 HTML 一样在同一种叙述性文档中使用：文章、图书、故事、诗歌、技术手册、Web 页面，等等。但是，让 XML 发明者大为吃惊的是，XML 的巨大成功不是在写作与发表领域，而是在更为枯燥乏味的数据格式世界中。通常，这不像文章与故事之类的叙述性数据，而是面向记录的数据，如数据库数据，用于对象序列化、财务记录、矢量图形、远程过程调用等任务。本章介绍这些数据传统格式中的一些缺陷，看看 XML 的特性是如何恰好适合这些任务的。

## 1.1 XML 的优点

如果你是位开发人员（至少我希望是，否则书中的内容可能让你没兴趣），毫无疑问，你在工作中已经编写了许多读取与写入文件的程序。每次编写新程序时，都要发明或学习新文件格式。几年来，我处理过的文件格式包括 RTF, Word .doc 文件，制表符分隔文件，FITS, PDF, PostScript，等等。毫无疑问，你一定还遇到过许多其他格式。

每次遇到新文件格式时，总是件麻烦事。如果这个文件格式有文档，通常也是文档不全，甚至可能错误连篇，通常没有指定字节顺序、行末规则等重要细节。号称能够读取与写入同一种文件格式的不同工具实际上会产生稍有不同的变体，实际应用时常常互不兼容。当你以为自己终于消灭了代码中最后一处缺陷时，却可能发现无法阅读别人的软件写成的文件。你发现自己对文件格式做了太多假设，又要回到草稿，重新开始。

因此，设计新文件格式时，开发人员通常采用最简单的格式，通常是制表符分隔文本或逗号分隔值。然而，即使如此简单的无修饰格式也常常遇到意想不到的问题。例如，把行中的两个制表符解释为空字符串、null 还是等于一个制表符？实际应用中，这几种变体都在使用。Java 语言的 StringTokenizer 类采用最后一种解释：连续两个制表符等同于一个制表符，但这只是实际数据文件中的最不常用方法，使许多 Java 编程人员大为吃惊，在 Java 程序中带来不少缺陷<sup>①</sup>。

### 1.1.1 一个思想试验

了解上述情况后，下面来做一个思想试验。假设任务是编写一个服务器方程序，通过 Internet 接受一个电子商务站点的订单。Web 服务器要将每个完成的订单发送到内部系统中，一次发送一个订单。你负责编写服务器上的代码，将订单发送到内部系统，并负责编写内部系统上的代码，接收和处理订单。这两个系统间唯一的连接是 TCP/IP 网络，即没有 JDBC 之类的高级 API 帮你在两个系统之间移动数据。你要发明一种数据格式，在一端产生，在另一端分析，

---

<sup>①</sup> 要理解这一点，注意 java.util.StringTokenizer 的设计目标是分析 Java 源代码，而不是读取制表符分隔的数据文件。但许多 Java 编程人员用它读取制表符分隔数据文件。

这种数据格式应当足够灵活以包含典型订单中的所有信息，包括客户名、所订产品、价格、制造商的库存单元（SKU）数、发货地址、税收和发货与处理费用。一种方法是将每个信息放在一行，如例子1.1所示。

---

**例子1.1 表示12 Birdsong Clocks, SKU 244订单的普通文本文档**

---

```
c32
Chez Fred
Birdsong Clock
244
12
USD
21.95
135 Airline Highway
Narragansett
RI
02882
USD
263.40
7.0
USD
18.44
USPS
USD
8.95
USD
290.79
```

另一种方法是使用更复杂更冗长的XML格式，如例子1.2所示。

---

**例子1.2 表示12 Birdsong Clocks, SKU 244订单的XML文档**

---

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<Order>
    <Customer id="c32">Chez Fred</Customer>
    <Product>
        <Name>Birdsong Clock</Name>
        <SKU>244</SKU>
        <Quantity>12</Quantity>
        <Price currency="USD">21.95</Price >
    </Product>
    <ShipTo>
        <Street>135 Airline Highway</Street >
        <City>Narragansett</City> <State>RI</State> <Zip>02882</Zip>
    </ShipTo>
    <Subtotal currency='USD'>263.40</Subtotal>
    <Tax rate="7.0"
        currency='USD'>18.44</Tax>
    <Shipping method="USPS" currency='USD'>8.95</Shipping>
    <Total currency='USD' >290.79</Total>
</Order>
```

编写发送与接收订单的代码时，是用例子1.1所示的简单分行符分隔文件还是用例子1.2所示的复杂标记XML格式？这两个文档包含相同信息。大多数没经验的开发人员喜欢第一种简

单格式，因为每个信息另起一行，不必增加多余的标识符。我的目的是要说明，与大多数开发人员的最初直觉相反，第二种格式更健壮、更可扩展和更容易使用。

### 1.1.2 健壮性

先看看健壮性。假设程序接收例 1.3 所示的订单。

#### 例子 1.3 表示 12 Birdsong Clocks, SKU 244 订单的文档

---

```
c32
Chez Fred
Birdsong Clock
12
244
USD
21.95
135 Airline Highway
Narragansett
RI
02882
USD
263.40
7.0
USD
18.44
USPS
USD
290.79
USD
8.95
```

它看上去和例子 1.1 相同，但如果认真比较，则可以发现 12 与 244 的位置互换了，原先 12 个报时钟的订单变成 244 个坐垫的订单。这个问题也许能在订单发货之前被发现，也许不能在订单发货之前被发现。更糟的是，发货费与总价也改变了，整个订单的成本变成 8.95 美元。这个问题也许能及时发现，也许不能及时发现。这类问题是实实在在的，许多电子商务站点因为价格错误而失去收入和信誉。

在 XML 版本中，根本不会有这个问题，因为每个数据项目都标明含义，可以随意打乱数量与 SKU、发货成本与总价的顺序，不会产生任何混乱，如例子 1.4 所示。传统系统中可能铸成大错的问题在 XML 中不成问题。

#### 例子 1.4 另一个表示 12 Birdsong Clocks, SKU 244 订单的文档

---

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<Order>
  <Customer id="c32">Chez Fred</Customer>
  <Product>
    <Name>Birdsong Clock</Name>
    <Quantity>12</Quantity>
    <SKU>244</SKU>
    <Price currency="USD">21.95</Price >
  </Product>
  <ShipTo>
    <Street>135 Airline Highway</Street >
    <City>Narragansett</City> <State>RI</State> <Zip>02882</Zip>
```

```
</ShipTo>
<Subtotal currency='USD'>263.40</Subtotal>
<Tax rate="7.0"
      currency='USD'>18.44</Tax>
<Total currency='USD'>290.79</Total>
<Shipping method="USPS" currency='USD'>8.95</Shipping>
</Order>
```

也许有人会说自己的系统中不会出现这类错误，因为他会检查每个值的意义。你会检查公司数据库中的SKU，保证其符合产品名和价格，然后再完成订单。也会检查每个方法调用的返回值，看看其是否为null，并捕获每个异常。你会编写大量测试，验证每个方法正常工作。会用源代码控制系统，因此总是可以撤销改变，代码要在进行所有回归测试并成功之后才能准入。每一行代码都文档齐全。事实上，所写的文档比实际代码还多。

即使你如此慎重，你能保证其他人发送与接收数据时也能同样慎重吗？倒不如采取更健壮的格式，使系统在遇到不可避免的错误时，不会造成更大的损害。

当然，XML对顽固的开发人员也大有帮助。定义限制时，如“每个订单要有一个发货地址”，“币值单位应为三字符代码USD, CAN或GBP”或“总价应为单价与数量的积加上税收和发货成本”，一种最简单的方法是使用声明式语言，指定所要的限制，而不是通过实际代码检查这些限制。对于XML格式数据，可以用声明式模式语言定义和测试这些限制。事实上，可以选择几种模式语言。最简单和最广泛支持的是经典的DTD（文档类型定义），可以检验所要的元素是否全部按所需的顺序出现，并且具有任何必要的属性。W3C XML模式语言更进了一步，可以限制特定元素与属性的内容，保证总价是大于1.00的十进制数。Schematron是最强大的模式语言，可以指定多元素限制，如“实际价格应小于或等于建议零售价”。

本章稍后和本书其余部分将详细介绍所有这些语言。目前只要知道，可以用简单格式在文档中列出所有限制，然后检查这些限制，而不需要大量多余的代码。可以先用验证程序检查文档之后再对其进行操作。验证是单独的、模块化和更易维护的过程。甚至可以不重新编译代码而改变限制或增加新限制。

### 1.1.3 扩展性

健壮性不是XML方法的惟一好处。XML还具有更好的扩展性。例如，假设突然需要在一些产品中增加折扣百分比。XML代码很容易改变，只要增加另一个元素：

```
<Product>
  <Name>Birdsong Clock</Name>
  <Quantity>12</Quantity>
  <SKU>244</SKU>
  <Price currency="USD">21.95</Price>
  <Discount>.10</Discount>
</Product>
```

而改变普通文本文件（或相应的二进制文件）则麻烦得多。尽管可以增加一行数据，但这样会使后面的所有内容顺序出错。可以将新信息放在文档末尾，但这样就不能让它靠近逻辑上归属的项目。假设并不是所有订单都有折扣，则没有折扣的产品是否使用空行？程序怎么知道空字符串要转换成0折扣，而不是转换成NaN（Not a Number）或抛出异常？这个问题不是无法解决，但会把简单问题复杂化。