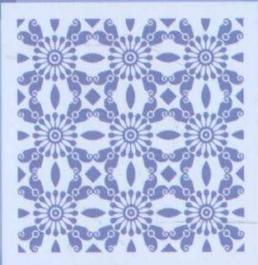


P

PRINCIPLES AND APPLICATIONS OF XML

XML原理与应用

夏天 编著

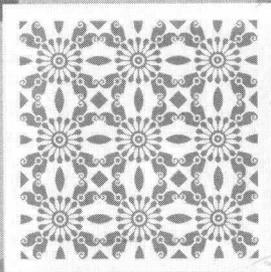


机械工业出版社
China Machine Press

高等院校计算机教材系列

P
RINCIPLES AND APPLICATIONS OF XML
XML原理与应用

夏天 编著



 机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

XML 原理与应用 / 夏天编著. —北京: 机械工业出版社, 2015.3
(高等院校计算机教材系列)

ISBN 978-7-111-49378-5

I. X… II. 夏… III. 可扩展语言—程序设计—高等学校—教材 IV. TP312

中国版本图书馆 CIP 数据核字 (2015) 第 031625 号

本书全面系统地介绍了 XML 的基本原理和关联技术, 注重 XML 技术体系的整体性和连续性, 便于读者快速把握各组成技术在 XML 中所起的主要作用。

全书共 10 章, 介绍了 XML 的基础知识、XML 的验证处理、XML 的呈现与转换、XML 的编程接口和应用等。各章最后均提供了习题, 便于读者巩固和拓展所学知识。

本书不仅可供信息管理与信息系统、计算机专业的本科生 / 专科生使用, 也可供信息技术领域的研究生和工程技术人员参考使用。

出版发行: 机械工业出版社 (北京市西城区百万庄大街 22 号, 邮政编码 100037)

责任编辑: 李 燕

印 刷: 三河市宏图印务有限公司

版 次: 2015 年 3 月第 1 版第 1 次印刷

开 本: 185mm × 260mm 1/16

印 张: 14.25

书 号: ISBN 978-7-111-49378-5

定 价: 35.00 元

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

客服热线: (010) 88378991 88361066

投稿热线: (010) 88379604

购书热线: (010) 68326294 88379649 68995259

读者信箱: hzjsj@hzbook.com

版权所有·侵权必究

封底无防伪标均为盗版

本书法律顾问: 北京大成律师事务所 韩光 / 邹晓东

前 言

XML (eXtensible Markup Language, 可扩展标记语言) 是由 W3C 于 1998 年 2 月发布的 Web 标准之一。经过十余年的发展与完善, XML 相关技术已经对当今社会的技术发展产生了深刻的影响, 从网页编写的标准化到 WebService, 从异构系统之间的数据交换到行业标准的制定, 从开放办公文档的表示到矢量图形的表示, 都能发现 XML 在其中的应用。不仅程序开发人员应该掌握 XML 技术的基本原理和组成体系, 系统架构人员、管理人员同样需要对 XML 有深入的理解。

本书从 XML 技术体系本身出发, 结合具体实例、历史背景和解决问题的思路, 逐一介绍 XML 最为重要的核心技术: DTD、Schema、CSS、XPath、XSLT、DOM, 同时结合技术的发展变化, 介绍 XML 在可缩放矢量图形方面的应用 SVG 和 XML 在数据表示领域的竞争技术 JSON。XML 作为信息技术的核心标准之一, 是发展新的 Web 应用、制定新的信息规范、描述新的行业应用必不可少的基础理论, 是 Web 标准体系的核心组成部分。

本书在写作过程中参考了大量国内外相关资料, 并结合作者多年的应用开发实践和教学经验, 重新设计了内容体系和实例, 所有章节均由作者独立完成, 保证了内容的原创性和连贯性。

在写作过程中, 作者力图加强对基本概念、基本理论和基本技术的介绍, 既强调 XML 技术体系的整体性, 也注重各个部分的相对独立性, 以适合教学使用。本书的特色如下:

- **问题驱动展开讲解。**对 XML 技术体系中的每一种特定技术, 首先介绍引入该技术的基本目的, 再进一步讲解其详细技术, 力求使读者能了解该技术的来龙去脉。
- **适用面广。**对 XML 技术的讲解尽量不涉及特定的操作系统、开发商或实现方式, 需要借助程序语言讲解的部分采用易于获取和测试的 JavaScript 语言, 适用于所有对 XML 技术感兴趣的人群。
- **实例丰富, 适合教学。**教材体系完整, 大量的实例和习题有助于读者对 XML 技术的理解和学习。
- **注重技术变化。**结合新流行的 Web 技术, 介绍与 XML 相关的 SVG 和 JSON 技术。

本书可作为高等院校电子政务、电子商务、信息管理与信息系统、计算机应用等专业本科生的教材, 也可作为相关标准制定人员、信息技术管理人员、应用系统设计和开发人员的参考书。

本书附带的电子教案和实例源代码均可从华章网站 (www.hzbook.com) 免费下载, 与本书有关的问题和建议也可以直接通过 <https://github.com/iamxiatian/xml/> 网站在线反馈。同时, 由于 XML 的相关标准和应用还在不断发展, 我们将紧跟变化, 对本书内容不断优化和完善, 而每一位读者的反馈都将是我们的最大动力。

作者

2015 年 1 月于中国人民大学

教学建议

教学章节	考核要求	课 时
第1章 XML 引论	<ul style="list-style-type: none"> 了解标记语言及其发展历史; 理解 XML 的目标与特点; 掌握 XML 的整体技术体系; 了解 XML 的常见应用场景; 掌握 XML 的编辑器、解析器和浏览器的特点, 学会基本的 XML 编辑器的安装。 	4
第2章 XML 的基本语法	<ul style="list-style-type: none"> 掌握 XML 的文档组成结构; 掌握 XML 声明、处理指令、注释、元素、属性、命名空间等基本概念; 能够编写格式良好的 XML 文档。 	4 ~ 6
第3章 文档类型定义 DTD	<ul style="list-style-type: none"> 了解 DTD 的作用; 掌握 DTD 的语法规则和使用方法; 能够利用 XML 编辑工具实现 DTD 的编辑和验证。 	6
第4章 XML Schema	<ul style="list-style-type: none"> 理解 XML Schema 的含义及用途; 掌握 XML Schema 的元素、属性的作用及使用方式; 掌握 XML Schema 的数据类型; 理解 XML Schema 的命名空间概念。 	6 ~ 8
第5章 利用 CSS 格式化 XML	<ul style="list-style-type: none"> 掌握 CSS 与 XML、HTML 的关联方式; 掌握 CSS 的常见选择器, 理解 CSS 的属性继承与覆盖原则; 掌握 CSS 的颜色、字体和文本属性以及 CSS 的盒状模型。 	4 ~ 6
第6章 XML 路径语言 XPath	<ul style="list-style-type: none"> 了解 XPath 的工作原理; 掌握 XPath 的定位路径表达式; 熟悉 XPath 的常用函数和数据类型; 掌握 XPath 2.0 的部分新特性。 	6
第7章 可扩展样式语言转换 XSLT	<ul style="list-style-type: none"> 了解 XSLT 的发展历史、XSLT 与 XSL 的关系; 掌握 XSLT 的转换处理过程; 掌握 XSLT 的基本语法; 能够使用 Saxon 进行 XSLT 测试。 	6 ~ 8
第8章 JavaScript (选读)	<ul style="list-style-type: none"> 了解 JavaScript 的起源; 掌握 JavaScript 的测试方法; 熟悉 JavaScript 的基本语法, 能够借助于文档进行基本的脚本编写; 了解浏览器对象模型 BOM; 理解 JavaScript 的定时器操作。 	6 ~ 8

(续)

教学章节	考核要求	课时
第9章 文档对象模型 DOM	<ul style="list-style-type: none"> ● 了解 DOM 的特点及规范级别; ● 掌握常见的 DOM 对象; ● 熟悉浏览器操纵 HTML 文档和 XML 文档的基本方法。 	6 ~ 8
第10章 XML 的应用与挑战	<ul style="list-style-type: none"> ● 了解 SVG 的特点; ● 能够基于 SVG 进行简单的图形表示; ● 了解 d3.js 的功能特点; ● 掌握 JSON 的数据结构和类型; ● 能够通过 JavaScript 解析 JSON 字符串。 	6 ~ 8

说明:

- 1) 建议总课时为 54~68。
- 2) 建议课堂教学全部在多媒体机房内完成, 实现“讲-练”结合。
- 3) 本书配有实验指导书, 可通过访问 <http://www.github.com/iamxiatian/xml/> 获取, 任课教师在讲授完每一章的基本理论后, 可以安排学生练习实验指导书中的内容, 并根据学生的练习情况对重点内容做进一步讲解。
- 4) 各章最后都附有习题, 任课教师可根据情况给学生留一些基本和中等难度的习题作为课外作业。

目 录

前言	
教学建议	
第 1 章 XML 引论	1
1.1 XML 的起源	1
1.1.1 标记简介	1
1.1.2 过程标记	2
1.1.3 通用编码	3
1.1.4 SGML	4
1.1.5 HTML	5
1.1.6 XML	7
1.1.7 SGML、HTML 与 XML 的 关系	10
1.2 XML 的设计目标与特点	10
1.2.1 XML 的设计目标	10
1.2.2 XML 的主要特点	11
1.3 XML 的技术体系	14
1.3.1 DTD 与 XML Schema	14
1.3.2 CSS	15
1.3.3 XSLT	15
1.3.4 XML DOM 与 SAX	15
1.3.5 XPath、XLink、XPointer	16
1.4 XML 的应用与发展	16
1.4.1 行业标记语言设计	17
1.4.2 电子文件的长期保存	17
1.4.3 电子数据交换	18
1.4.4 Web 应用	18
1.5 XML 的相关工具	19
1.5.1 XML 编辑工具	19
1.5.2 XML 浏览工具	20
1.5.3 XML 验证工具	20
1.5.4 XML 解析器	21
1.6 小结	21
1.7 习题	22
第 2 章 XML 的基本语法	24
2.1 XML 文档结构	24
2.1.1 文档声明	25
2.1.2 处理指令	26
2.1.3 注释	27
2.2 XML 的元素	29
2.2.1 元素和标记	29
2.2.2 元素的内容	30
2.2.3 元素的嵌套	34
2.3 XML 的属性	34
2.3.1 属性的语法形式	35
2.3.2 属性的使用场景	35
2.3.3 属性的命名规则	37
2.3.4 属性值	37
2.4 XML 的命名空间	38
2.4.1 命名空间的引入	38
2.4.2 命名空间的使用	39
2.4.3 默认命名空间	40
2.4.4 命名空间的作用域	40
2.5 XML 文档的规范级别	41
2.5.1 格式良好的 XML 文档	41
2.5.2 有效的 XML 文档	42
2.5.3 规范化的 XML 文档	42
2.6 小结	43

2.7 习题	43	4.2.4 包含与导入	73
第 3 章 文档类型定义 DTD	45	4.3 XML Schema 的元素	74
3.1 DTD 的作用	45	4.3.1 schema 根元素	75
3.2 DTD 的关联方式	46	4.3.2 element 元素	75
3.2.1 内部 DTD 关联方式	46	4.3.3 element 元素的默认值和 固定值	75
3.2.2 外部 DTD 关联方式	47	4.3.4 元素的引用和替代	76
3.2.3 公用 DTD 关联方式	47	4.4 XML Schema 的属性	77
3.2.4 内外结合关联方式	48	4.4.1 属性声明	77
3.3 DTD 的元素	49	4.4.2 指派属性类型	78
3.3.1 元素类型声明	49	4.4.3 属性的默认值和固定值	79
3.3.2 空元素	50	4.5 XML Schema 的数据类型	79
3.3.3 文本类型元素	50	4.5.1 简单数据类型: SimpleType	80
3.3.4 元素内容模型与混合内容元素	51	4.5.2 复杂数据类型: ComplexType	84
3.4 DTD 的属性	53	4.6 XML Schema 与命名空间	89
3.4.1 属性声明	53	4.6.1 targetNamespace	89
3.4.2 属性类型	54	4.6.2 elementFormDefault 与 attributeFormDefault	90
3.4.3 属性的默认形态	57	4.6.3 form 属性	91
3.4.4 特殊属性	57	4.7 XML Schema 的注释与注解	92
3.5 DTD 的实体	59	4.7.1 注释	92
3.5.1 实体类型与实体引用	59	4.7.2 注解	92
3.5.2 内部可解析通用实体	60	4.8 小结	94
3.5.3 外部可解析通用实体	61	4.9 习题	94
3.5.4 外部非解析通用实体	62	第 5 章 利用 CSS 格式化 XML	96
3.5.5 内部参数实体	63	5.1 CSS 概述	96
3.5.6 外部参数实体	65	5.1.1 CSS 的基本概念	96
3.6 DTD NOTATION	66	5.1.2 CSS 的发展历史	96
3.7 DTD 的包含与忽略	66	5.2 关联 CSS 的方法	97
3.8 小结	67	5.2.1 CSS 与传统网页的关联方式	97
3.9 习题	67	5.2.2 CSS 与 XML 的关联方式	99
第 4 章 XML Schema	69	5.3 CSS 的语法基础	100
4.1 XML Schema 概述	69	5.3.1 CSS 的基本语法	100
4.2 XML Schema 快速入门	70	5.3.2 CSS 的选择器	100
4.2.1 快速入门实例	70	5.3.3 CSS 的继承与覆盖	102
4.2.2 Schema 文档结构	72		
4.2.3 引用方式	73		

5.4 CSS 重要属性	104	6.6.1 节点集函数	133
5.4.1 颜色属性	104	6.6.2 布尔函数	134
5.4.2 字体属性	106	6.6.3 数值函数	134
5.4.3 文本属性	109	6.6.4 字符串函数	135
5.4.4 盒状模型相关属性	110	6.7 XPath 2.0 的新特性	135
5.4.5 可视格式化模型相关属性	113	6.7.1 支持 XML Schema 的 数据类型	136
5.5 小结	114	6.7.2 更为丰富的处理函数	136
5.6 习题	114	6.7.3 支持序列	136
第 6 章 XML 路径语言 XPath	115	6.7.4 支持逻辑判断	137
6.1 XPath 概述	115	6.7.5 更多的节点测试	137
6.1.1 XPath 及其作用	115	6.7.6 调用自定义函数	137
6.1.2 XPath 的工作原理	116	6.8 小结	137
6.1.3 XPath 的表达式与操作符	118	6.9 习题	138
6.1.4 如何测试 XPath	118	第 7 章 可扩展样式语言转换 XSLT	139
6.2 XPath 节点与节点集	120	7.1 XSLT 概述	139
6.2.1 节点的基本属性	120	7.1.1 XSLT 与 XSL	139
6.2.2 节点类型	121	7.1.2 XSLT 的作用	140
6.2.3 节点集	122	7.1.3 XSLT 的工作流程	141
6.3 XPath 定位路径表达式	123	7.1.4 XSLT 的应用模式	141
6.3.1 XPath 定位步骤	123	7.1.5 XSLT 与 CSS 的区别	141
6.3.2 XPath 轴	124	7.2 如何测试 XSLT	142
6.3.3 节点测试	126	7.2.1 通过浏览器测试 XSLT	142
6.3.4 谓词	127	7.2.2 通过 XML 专业工具 测试 XSLT	142
6.3.5 定位路径缩写	128	7.2.3 通过 XSLT 处理器 测试 XSLT	144
6.4 XPath 基本表达式	128	7.3 XSLT 快速入门	145
6.4.1 布尔表达式	128	7.3.1 stylesheet 元素	147
6.4.2 等式表达式	128	7.3.2 template 元素	147
6.4.3 关系表达式	128	7.3.3 apply-templates 元素	147
6.4.4 数值表达式	129	7.3.4 value-of 元素	148
6.5 XPath 的数据类型	129	7.3.5 attribute 元素	148
6.5.1 字符串类型	129	7.4 XSLT 的输出格式控制	149
6.5.2 数值类型	130	7.5 XSLT 的逻辑处理元素	149
6.5.3 布尔类型	132		
6.5.4 节点集类型	133		
6.6 XPath 1.0 的常用函数	133		

7.5.1 条件处理元素	149	8.6.3 基本类型和引用类型	174
7.5.2 循环元素 for-each	151	8.6.4 原型与继承	175
7.5.3 排序元素 sort	152	8.6.5 类方法	176
7.6 XSLT 的模式	153	8.7 浏览器对象模型 BOM	176
7.7 XSLT 的命名模板	155	8.7.1 window 对象	176
7.8 XSLT 的函数	156	8.7.2 document 对象	177
7.9 XSLT 2.0 的新特性	157	8.7.3 navigator 对象	178
7.10 小结	157	8.7.4 location 对象	178
7.11 习题	158	8.7.5 screen 对象	179
第 8 章 JavaScript (选读)	160	8.7.6 history 对象	179
8.1 JavaScript 概述	160	8.8 定时器	179
8.1.1 JavaScript 的历史	160	8.8.1 一次性定时器的设置与取消	179
8.1.2 jQuery 概述	161	8.8.2 重复定时器的设置与取消	180
8.2 JavaScript 的测试方法	162	8.9 小结	180
8.2.1 JavaScript 与网页的关联 测试方法	162	8.10 习题	181
8.2.2 在页面加载之后运行 JavaScript	163	第 9 章 文档对象模型 DOM	182
8.2.3 利用浏览器内置的 JavaScript 控制台	163	9.1 DOM 概述	182
8.3 JavaScript 的变量和常量	164	9.1.1 DOM 的定义及作用	182
8.3.1 数据类型	164	9.1.2 DOM 的发展历史与规范 级别	183
8.3.2 变量的声明和赋值	165	9.2 DOM 的基本对象	184
8.3.3 变量的作用域	166	9.3 利用 Mongoose 搭建 DOM 测试环境	185
8.3.4 常量	166	9.4 利用 DOM 操纵 HTML	186
8.4 JavaScript 的基本语句	167	9.4.1 HTML DOM 及元素定位 方法	186
8.4.1 注释语句	167	9.4.2 改变元素节点内容	188
8.4.2 条件语句	167	9.4.3 改变属性节点内容	188
8.4.3 循环语句	169	9.4.4 节点的创建与删除	189
8.5 函数和数组	170	9.4.5 HTML DOM 示例	189
8.5.1 函数	170	9.5 利用 DOM 操纵 XML	190
8.5.2 数组	172	9.5.1 加载 XML 文档	191
8.6 对象	173	9.5.2 节点访问方法	192
8.6.1 创建对象	173	9.5.3 节点定位属性	192
8.6.2 属性和方法	174	9.5.4 节点常用属性	193

9.5.5 节点常用方法	195	10.2.5 基于 SVG 的 d3.js 图形 绘制库	206
9.5.6 XML DOM 示例	196	10.3 数据传输的挑战者——JSON	210
9.6 小结	198	10.3.1 JSON 的数据结构	210
9.7 习题	198	10.3.2 JSON 的值类型	211
第 10 章 XML 的应用与挑战	199	10.3.3 JSON 与 XML 的对比	212
10.1 概述	199	10.3.4 利用 JavaScript 解析 JSON	213
10.2 新流行应用——SVG	200	10.4 小结	214
10.2.1 SVG 的基本形状	201	10.5 习题	214
10.2.2 SVG 的样式设置	203	参考文献	216
10.2.3 SVG 的层与重叠	204		
10.2.4 SVG 的透明度	204		

第 1 章 XML 引论

本章主要讨论与 XML 相关的三个重要问题：①为什么会出现 XML，它的目标是什么？② XML 的技术体系是如何构成的？③ XML 有哪些典型应用场景？在此，我们只是对 XML 的体系结构做一个整体描述，其详细内容将在后续章节中依次展开讲解。

1.1 XML 的起源

XML (eXtensible Markup Language)，即可扩展标记语言，自从 1998 年 2 月成为 W3C 的推荐标准之后，便得到了人们的广泛关注。各个软件厂商纷纷宣布或者已经在自己的产品中增加对 XML 的支持，例如，微软公司在其新一代的办公软件中，采用 XML 格式代替原先的二进制格式对文档进行描述；Oracle 公司则将 XML 数据模型完全集成到其数据库中，并提供浏览和查询 XML 的新的标准访问方法，通过 Oracle XML DB，人们可以获得相关数据库技术的所有优势和 XML 的优势；原 Sun 公司更是在其开发语言 Java 当中内置了对 XML 的支持；与此同时，W3C 也在积极地致力于完善 XML 的标准体系。

近几年来，众多与 XML 相关的新名词、新技术如潮水般不断涌现，WebService、XHTML、手机 WAP、RSS、开放式电子公文格式、XML 数据交互接口，无一不与 XML 密切相关。即便是不参加具体开发的管理人员，甚至是与软件过程无直接关联的人员，也开始认识到 XML 的重要性，开始想方设法获取对 XML 的一个完整的认识。那么，到底什么是 XML，它从何方来，要往何方去？要弄清这个问题，我们就不得不回顾标记语言以及它的发展历史。

1.1.1 标记简介

标记语言 (Markup Language) 起源于传统印刷。文字在印刷之前必须先进行排版，在电子出版业出现以前，人们需要以手工或打字方式复制手稿副本，并对副本以人工方式进行标记，通过这些标记说明告诉排版印刷人员如何处理版面，如印刷选用的字体、行间距、对齐方式等格式信息，最终完成印刷前的制版工作。这些人工增加的注释就被称为“标记”(Mark-up)。

类似的，在计算机中，人们也需要对输入的文字指定其外观样式。例如，用户可以对写好的文本选择字体，设置字符间距、段落等格式信息，这些信息被称为电子标记（电子标记是“Markup”，以区别于传统的手写标记“Mark-up”），它们以特殊代码形式存储在文本中。计算机处理程序在读取文本的同时还要读出以电子形式标注的格式信息，最终告知计算机安装规定的格式展现文本，例如，大家熟悉的 Word、WPS、OpenOffice 等文本编辑器都是借助于标注代码来定义格式与外观的。与直接手工标注格式不同，在计算机中，用户通常通过在菜单中选择命令的方式来向文本中添加格式指令，进行标注。

通俗地说，标记语言就是一种用来给文本添加“格式标注”以指明文档中文本编排格式的语言，一般由定义文档格式的一些规定代码和控制标记组成。

1.1.2 过程标记

RTF (Rich Text Format, 富文本格式) 就是这样一种简单的过程标记语言, 它是由微软公司于 1987 年开始开发的一种跨平台文档格式, 大多数的文字处理软件都能读取和保存 RTF 文档。作为微软公司的标准文件, 在早期人们需要向微软付费才能购买一本薄薄的 RTF 标准文件。后来, 随着 RTF 的普及, 微软公司把标准文件公开, 放在网上供开发者下载, 无需支付任何费用, 目前 RTF 的最新版本是 RTF1.9.1^①。下面, 我们就通过一个示例来了解一下 RTF 文档的过程标记。

【例 1-1】创建、查看 RTF 文档的标记。

1) 打开 Word 软件, 输入“XML and RTF!”, 选择字体为 Times New Roman, 字号为 28, 并把“XML”的颜色设置为红色, 如图 1-1 所示。

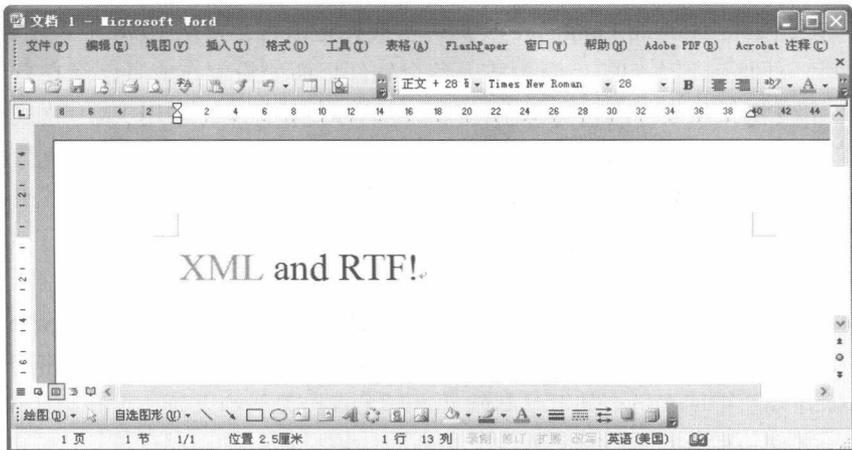


图 1-1 在 Word 中输入一段文字

2) 单击“保存”按钮, 在弹出的“另存为”对话框中输入文件名“1-1.rtf”, 注意必须把保存类型选择为“RTF 格式 (*.rtf)”, 如图 1-2 所示, 然后单击“保存”按钮。



图 1-2 保存为 RTF 文档

^① <http://www.microsoft.com/downloads/details.aspx?FamilyId=DD422B8D-FF06-4207-B476-6B5396A18A2B&displaylang=en>.

3) 在记事本中打开刚刚建立的“1-1.rtf”文档，如图 1-3 所示是加上标记后的文本格式。



```

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)
<Rtf1\ansi\ansicpg936\uc2\deff0\stshfdbch13\stshfloch0\stshfhich0
\stshfbi0\deflang1039\deflangfe2052\fonttbl{\f0\froman\charset0
\Fprq2{\*\panose 02020603050405020304}Times New Roman;}
\F13\fnil\charset134\Fprq2{\*\panose 02010600030101010101}
\cb\ce\cc\ce5{\*\falt SimSun}}{\f51\fnil\charset134\Fprq2
{\*\panose 0201060003010101010101}\cb\ce\cc\ce5}}{\f158
\froman\charset238\Fprq2 Times New Roman GE;}
{\f159\froman\charset204\Fprq2 Times New Roman Cyr}}{\f161
\froman\charset161\Fprq2 Times New Roman Greek}}{\f162
\froman\charset162\Fprq2 Times New Roman Tur}}{\f163
\froman\charset177\Fprq2 Times New Roman (Hebrew);}
{\f164\froman\charset178\Fprq2 Times New Roman (Arabic)}}{\f165
\froman\charset186\Fprq2 Times New Roman Baltic}}{\f166
\froman\charset163\Fprq2 Times New Roman (Vietnamese)}}{\f290
\fnil\charset0\Fprq2 SimSun Western{\*\falt SimSun}};
{\f67\fnil\charset0\Fprq2 @\cb\ce\cc\ce5 Western}}

```

图 1-3 用记事本打开“1-1.rtf”

从图 1-3 中可以看出，其中的内容全部是文本字符，并且比开始输入的“XML and RTF!”要多很多，这些多出来的字符正是用于描述字符串“XML and RTF!”的格式的标记信息。RTF 格式能够得到大多数文字编辑软件的支持，与它的优点是分不开的。首先，RTF 是一种得到广泛支持的文件格式，绝大多数文字处理程序都能够创建和读取 RTF 文档，如 Word、WPS 等；其次，RTF 可以存储复杂的格式信息，如图片、表格、超链接等。然而，RTF 格式的缺点也比较明显，如：

1) 通过 RTF 格式保存的文件通常都比较大，RTF 中对相同文本段落会有不同版本的格式标注，另外 RTF 本质上是一种文本格式，对于内嵌的图片，也只能用文本的方式对图片进行等价表示，于是导致文件体积变大。

2) RTF 格式不能保存一些特殊的信息，例如，对图片文件只能通过转换成位图的方式进行保存。我们在把一个复杂的 Word 文档转换为 RTF 时，经常会发生一些信息丢失的情况，正是这一原因造成的。

3) RTF 格式不够灵活，对格式规则的任何变更意味着需要对文档进行手工修改。还有，标记的多少是和系统有关的，这就增加了移植的困难性。

1.1.3 通用编码

1967 年，William W. Tunnicliffe 在加拿大政府印刷局的一次会议上提议出版商应该采用通用编码 (Generic Markup) 对文本进行编码，标识文本中不同部分的功能，例如哪段文字是标题，哪段文字属于引用等，而不是仅为最终的打印样式进行编码控制。相对于设置文本出版格式的传统标记，通用编码更侧重于文本不同部分的功能标记，力图实现文本的文字信息和呈现格式的分离。根据这一思想，以 Tunnicliffe 领导的小组在 1970 年开发出了出版业的通用编码标准“GenCode”，成为标准通用标记语言 SGML 的基础。

由著名计算机科学家 Donald E. Knuth 在 20 世纪 70 年代发明的排版系统 TeX，是另一种较有代表性的通用编码。TeX 引入了宏的概念，可以生成高质量的 DVI 格式文件并打印输出，还可以利用 dvips、dvi2pdf、pdflatex 等程序生成 PDF 或 PS 格式文件，也可以利用 latexhtml 生成 HTML 网页文件。同时，TeX 能够处理非常复杂的数学公式，排版出精美的打印效果，也可以借助于宏程序包制作幻灯片、定义模板等，因此，它在数学、物理学和计算

机科学等学术界十分流行。

【例 1-2】TeX 示例清单 Test.tex。

```
1 \noindent Xiatian\par
2 \noindent School of Information Resource Management\par
3 \noindent Renmin University of China\par
4 \smallskip
5
6 This is a book about XML. If you have some problems, contact us please.\par
7 \bye
```

在上述清单中，\noindent、\par、\bye 等都是些具有意义的标识符，例如 \par 代表的是段落标记，可见，TeX 采用了中性标记来标识文档的逻辑结构。

相对于 RTF 这类过程标记语言，通用编码具有如下优点：

1) 通用编码具有更好的可移植性，灵活性高。为改变文档的外观，人们可以不修改文档本身，而只对宏的定义进行修改即可实现，编辑一个宏，其变化会自动应用于整篇文档。特别是我们不需要对标记进行重新编码，避免了这个既费时间又容易出错的过程。

2) 标记本身更倾向于描述文档的逻辑结构。例如，通过自定义 \title 宏，声明这是一篇文章的标题，\abstract 宏表示文档的摘要等。对标题和摘要的格式信息则交给宏的具体定义予以描述。用户在编写自定义的宏时，一般都会为标记起一个比较有意义的名字，如 title、abstract 等，清楚地体现了结构对于格式的优越性。

1.1.4 SGML

SGML (Standard Generalized Markup Language)，即标准通用标记语言，是一种定义电子文档结构和描述其内容的国际标准语言，早在 Web 发明之前 SGML 就已存在。

注意：SGML 中的 G、M 和 L 分别代表 SGML 的三位创始人：Goldfarb、Mosher 和 Lorie。

20 世纪 60 年代，在 IBM 计算机之间存在大量格式不同的电子文件，为了方便数据交换，IBM 公司的研究人员经过研究认为：要提高系统的移植性，必须采用一种通用的文档格式，这种文档的格式必须遵守特定的规则，并对文档的内容和样式实现分离。1969 年，IBM 研究人员开发出了通用标记语言 (Generalized Markup Language, GML)，GML 是一种自参考的语言，可以描述任何其他语言的语法和词汇，能够标识出任何数据集合的结构。在标记语言的概念基本达成共识的基础上，IBM 研究人员 Charles Goldfarb 带领的开发团队完善了 GML，并将其称为 SGML，SGML 成为了 IBM 内部格式化和维护合法化文件的手段。通过对 SGML 不断地拓展和修改，SGML 在 1986 年被国际标准化组织 (ISO) 采纳成为正式标准。

与通用编码不同，SGML 不仅能够对文档本身进行标注，还能够描述文档的层次结构信息。SGML 实际上是一个通用的文档结构描述符号化语言，用来定义文档模型的逻辑和物理类结构。一个 SGML 文件由三部分组成，即语法定义、文档类型定义 (Document Type Definition, DTD) 和文档实例。语法定义部分定义了文档类型定义和文档实例的语法结构；文档类型定义部分定义了文档实例的结构和组成结构的元素类型；文档实例是 SGML 的文档

主体部分。

在很多大型组织中 SGML 被用来创建、处理和发布大量的文本信息，得到了广泛应用。虽然 SGML 并没有给文档强加一种结构，但标准委员会、企业联盟及其他组织却开发了基于 SGML 的结构应用，它们为特定的应用程序定义了一些标准的 DTD。一些著名的应用包括：

1) HTML。用于描述 Web 文档的著名的超文本标记语言，尽管很少有 HTML 的开发人员知道 SGML，但 HTML 确实被定义为一种 SGML 的 DTD。

2) DocBook[⊖]。由结构化信息标准促进组织 OASIS 开发，用于技术文章和书籍的结构化表示。

1.1.5 HTML

毫无疑问，HTML 是 SGML 最流行的一种应用。这就是说，HTML 是遵守 SGML 规则的一组标记集合，这组标记适应于对互联网文档的结构描述。最初的 HTML 是于 1989 年由欧洲量子实验室的研究人员 Tim Berners Lee 在 SGML 的基础上开发的一个简化子集。HTML 继承了 SGML 许多重要的特点，比如结构化、独立于具体实现、可描述性，但是 HTML 也存在很多缺陷。HTML 只能使用固定的有限的标记，最初设计标记的重点在于如何表示网页的展示样式，随着网络的普及，人们开始对简单的 HTML 有所不满，因此，又开始了对 HTML 的改造过程。HTML 从最开始的 1.0 版本一直到最近发布的 HTML 5.0，扩展了许多命令，同时抛弃了一些用于单纯描述显示格式的标记，使得 HTML 得到了更为广泛的应用。

HTML 本质上是通过标记的方式，以纯文本形式对网页进行描述，而对标记的解释则由浏览器执行，如现在流行的 IE (Internet Explorer) 浏览器、FireFox 浏览器、Opera 浏览器等。浏览器读取网页源代码，即 HTML 标记文本，并通过渲染呈现给用户，就形成了用户最终看到的网页。

【例 1-3】HTML 编写方式及显示效果。

1) 在记事本中输入一段 HTML 代码，如图 1-4 所示，其中通过尖括号包括起来的的就是用于网页描述的 HTML 标记。

2) 选择“文件”菜单中的“保存”菜单项，在“另存为”对话框中，选择保存位置和保存文件名称，此处将文件命名为“1-2.html”，保存类型选择为“所有文件”，如图 1-5 所示。

注意：保存文件时需将“保存类型”选择为“所有文件”，且文件名必须包含“.html”。正常保存完毕，生成的新文件的图标类型为网页类型，一般情况下都表示刚才保存的是文本类型。

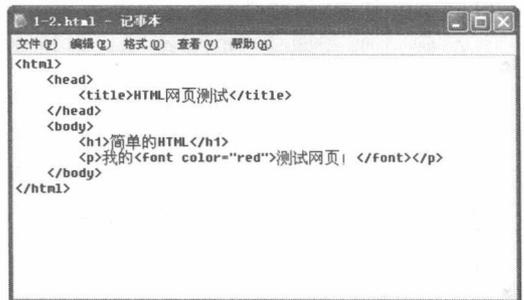


图 1-4 用记事本编辑 HTML 代码



图 1-5 保存 HTML 文档

⊖ <http://www.docbook.org>.

需要补充说明的是，Windows 系统中默认会隐藏已知文件类型的扩展名，比如创建了一个“1-2.html.txt”文件，本质上是一个文本文件，但由于文本文件是 Windows 中默认可以识别的文件类型，如果不更改 Windows 的默认选项，则在资源管理器当中看到的文件名称是“1-2.html”，为避免出现混淆，请去掉 Windows 中的“隐藏已知文件类型的扩展名”选项。Windows XP、Windows 7、Windows 8 等不同 Windows 版本的取消方式相同，下面以 Windows XP 为例来具体说明取消方式。

如图 1-6 所示，打开资源管理器，进入任意一个目录，单击“工具”菜单中的“文件夹选项”命令，随即弹出如图 1-7 所示的窗口。

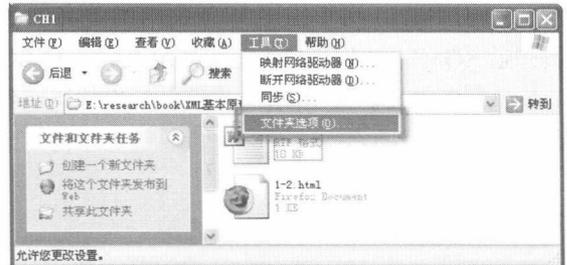


图 1-6 更改文件夹选项

选择“查看”选项卡，取消“高级设置”区域内的“隐藏已知文件类型的扩展名”复选框，然后单击“确定”按钮即可。

3) 双击刚刚创建的“1-2.html”文件，默认会由浏览器打开该文件，并解析 HTML 标记，以网页的形式呈现最终的结果，在 Firefox 浏览器中的显示效果如图 1-8 所示。如果个人计算机上没有安装 Firefox 浏览器，也可以通过 IE 浏览器浏览。

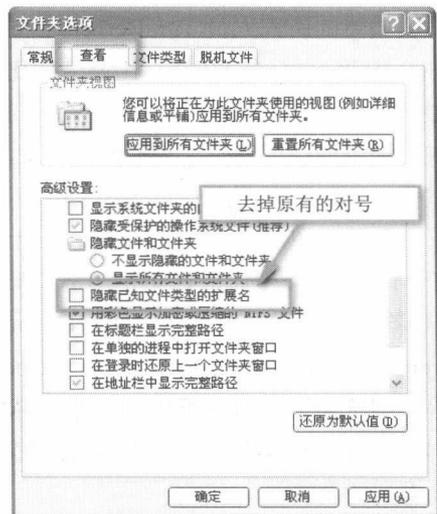


图 1-7 显示已知文件类型的扩展名

在浏览器中我们可以看到“测试网页”文字是红色的，对应于源代码中的标记即为“测试网页”，类似于这种标记就是一种过程标记。在 HTML 中，这类过程标记也可以通过通用标记予以实现，我们可以把“测试网页”更改为“测试网页”，然后再定义 highlight 样式如下：

```
.highlight {
    color:red;
}
```

这种描述样式在有关 CSS 的后续内容中 (1.3.2 节) 还会详细介绍。这个例子则表明现在的 HTML 是一种既支持过程标记又支持通用标记的混合标记语言。由于该示例文件没有明确指定网页的编码，读者在自己系统中打开该文件时，可能会遇到乱码问题，此时，只需要选择浏览器中的编码菜单，把当前文档的编码选择为 GBK 即可。

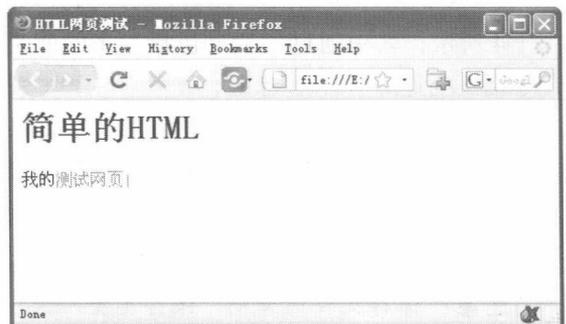


图 1-8 在 Firefox 浏览器中显示“1-2.html”文件

由于该示例文件没有明确指定网页的编码，读者在自己系统中打开该文件时，可能会遇到乱码问题，此时，只需要选择浏览器中的编码菜单，把当前文档的编码选择为 GBK 即可。