

INFORMATION RECOMMENDATION SYSTEM

信息推荐系统

曾子明 著



科学出版社

013058036

G202
229

信息推荐系统

曾子明 著



教育部人文社会科学青年项目(项目编号:08JC870011)
国家自然科学基金青年项目(项目编号:71103136)

G202
229

科学出版社

北京



北航

C1669396

013028038

内 容 简 介

信息推荐系统是解决互联网海量信息资源出现“信息过载”问题的非常有潜力的方法。本书根据国内外信息推荐系统的发展和在作者最新科研成果的基础上,较为系统地介绍信息推荐系统的原理、技术和应用,为用户提供个性化的信息推荐服务。本书首先介绍信息推荐系统的基础知识,在此基础上,探讨信息推荐系统在电子商务领域的应用,包括基于领域本体的商品信息推荐系统、基于 Web 挖掘的商品信息推荐系统和基于案例推理的商品信息推荐系统。此外,本书还对信息推荐系统的研究热点进行探讨,包括基于社会化标签的信息推荐系统和基于情境感知的信息推荐系统。

本书适合从事信息管理和应用、信息系统设计与开发、企业信息系统等相关领域的广大工程技术人员和管理人员参考,同时也可作为高等院校信息管理系统、计算机应用、电子商务等专业的高年级本科生、研究生的教学参考书。

图书在版编目(CIP)数据

信息推荐系统/曾子明著. —北京:科学出版社,2013

ISBN 978-7-03-037401-1

I. ①信… II. ①曾… III. ①信息系统 IV. ①G202

中国版本图书馆 CIP 数据核字(2013)第 092996 号

责任编辑:李 娜 童安齐 / 责任校对:王万红

责任印制:吕春珉 / 封面设计:耕者设计工作室

科 学 出 版 社 出 版

北京东黄城根北街16号

邮政编码:100717

<http://www.sciencep.com>

双 青 印 刷 厂 印 刷

科学出版社发行 各地新华书店经销

*

2013年5月第 一 版 开本:B5(720×1000)

2013年5月第一次印刷 印张:14 1/4

字数:271 000

定价:60.00 元

(如有印装质量问题,我社负责调换<双青>)

销售部电话 010-62136131 编辑部电话 010-62137026(BA08)

版权所有,侵权必究

举报电话:010-64030229; 010-64034315; 13501151303

前 言

我们生活在一个越来越依靠信息的时代,并正在向数字化时代迈进。数字化时代来临时,各种信息资源的电子化传递都将成为数字化经济的标志。信息社会化、社会信息化,信息生产与消费促进了信息产业和信息技术的飞速发展。当前,互联网已经成为人们获取信息的重要来源,是人们获取信息、改变生活方式、赢得商机的重要媒介。然而,互联网规模和信息资源的迅猛增长带来了信息过载的问题,人们面临“信息丰富、但有用信息获取困难”的窘境,从互联网中有效地获取信息日益困难。目前,搜索引擎是最普遍的辅助人们获取信息的工具,但它只能满足主流需求,没有考虑用户的个性化信息需求,仍然无法很好地解决信息过载的问题。以“信息推送”为服务模式的信息推荐系统,是当前解决信息过载问题的主要手段,它能够在分析预测用户需求的基础上主动推送用户可能需要但又无法获取的有用信息,并能够以用户为中心,通过研究用户行为、兴趣和环境等,为用户推荐更具针对性的信息,即实现信息的“按需定制服务”。

信息推荐系统作为一种人机交互系统,主要应用信息检索、信息过滤、数据挖掘、人工智能等多种技术和方法为用户提供“信息推送”服务,帮助用户在互联网海量信息中筛选符合其个性化需求的信息资源,为用户带来全新的信息服务体验。随着电子商务、Web 2.0 和社交网络的流行和发展,信息推荐系统作为信息服务科学的一个重要研究领域,已得到国内外学者、研究机构和企业界的广泛关注。因此,系统地探讨信息推荐系统的基本原理、技术以及研究热点,无疑将从理论和实践上推动信息推荐系统的进一步发展。

本书是教育部人文社会科学青年项目(项目编号:08JC870011)和国家自然科学基金青年项目(项目编号:71103136)的成果之一。本书较系统和全面地论述信息推荐系统的相关原理、技术和应用,全书内容较新颖,反映了信息系统和信息服务领域的发展动态以及作者多年来的研究成果。全书共9章:

第1章是概论部分,在介绍两种服务模式的基础上探讨基于“推送”模式的信息推荐系统,包括它的应用领域和研究热点。

第2章首先介绍信息推荐的基本技术,并从信息系统的角度探讨信息推荐系统的设计方法与开发的基本原则。

第3章介绍信息内容过滤推荐的模型和方法,包括信息内容过滤推荐的相关技术、信息内容过滤推荐的系统模型和相关算法、信息内容过滤推荐的用户反馈机制等。

第4章首先介绍两种基本的信息协同过滤推荐方法,在此基础上介绍基于模型的协同过滤推荐方法,并提出一种移动环境下基于隐式评分的博客信息推荐方法。

第5章在探讨基于领域本体的商品信息组织方法基础上,从不同研究视角提出相应的基于领域本体的商品信息推荐模型。

第6章针对顾客经常购买的商品,提出一种基于Web挖掘技术的商品信息推荐系统。

第7章针对专业知识较强、顾客购买频率较低的商品,提出一种基于案例推理的商品信息推荐系统,为顾客购物提供咨询服务和决策支持。

第8章首先介绍社会化标签系统以及目前基于社会化标签的信息推荐相关技术;在此基础上,提出一种基于社会化标签的信息推荐模型。

第9章针对泛在环境下用户个性化信息需求具有情境敏感性,研究基于情境感知的信息资源推荐服务的理论、模型与方法,从不同研究视角提出相应基于情境感知的信息推荐方法。

本书在撰写、成稿的过程中,参考了国内外许多专家、学者的论著,他们的成果为本书提供了丰富的素材和理论支撑,并在每章的参考文献中进行了标注,如果有不慎遗漏的,在此表示歉意。

感谢我的父母和妻子,他们在我的写作过程中给予了极大的支持,并为写作创造了良好的条件。

信息推荐系统是一个新的信息系统研究领域,发展迅速,需要进一步深入研究的问题很多,希望本书的出版能起到抛砖引玉的作用。尽管作者在项目研究和本书撰写过程中付出了艰辛的努力,但由于该领域研究内容新,一些理论方法和技术还在发展之中,书中难免存在不足之处,欢迎读者不吝赐教。

目 录

前言

第 1 章 信息推荐系统概论	1
1.1 网络信息资源及获取服务模式	1
1.1.1 网络信息资源	1
1.1.2 信息获取服务模式	3
1.2 基于“信息推送”模式的信息推荐系统	7
1.2.1 信息推荐系统的概念与通用模型	7
1.2.2 信息推荐系统与个性化信息服务	10
1.2.3 信息推荐系统的研究内容	13
1.2.4 信息推荐系统的分类	14
1.2.5 信息推荐系统的发展现状和实例	15
1.3 信息推荐系统的应用领域和研究热点	17
1.3.1 信息推荐系统的应用领域	17
1.3.2 信息推荐系统的研究热点	20
参考文献	22
第 2 章 信息推荐技术和系统设计	24
2.1 信息推荐系统的相关技术	24
2.1.1 信息检索和信息过滤	24
2.1.2 数据挖掘技术	26
2.1.3 信息推荐算法概述	28
2.2 信息推荐的系统分析与设计	34
2.2.1 信息推荐的系统分析	34
2.2.2 信息推荐的系统设计	36
2.3 信息推荐的系统开发方法	41
2.3.1 原型法的基本思想	41
2.3.2 基于原型法的信息推荐系统开发过程	42
参考文献	44
第 3 章 信息内容过滤推荐系统	46
3.1 引言	46

3.2	内容过滤推荐系统的相关技术	48
3.2.1	信息检索模型	48
3.2.2	文本特征抽取	50
3.3	内容过滤推荐系统的模型和算法	52
3.3.1	基于内容过滤的信息推荐模型	52
3.3.2	基于向量空间模型匹配的信息推荐算法	54
3.3.3	基于朴素贝叶斯分类的信息推荐算法	56
3.4	内容过滤推荐系统的用户反馈	58
	小结	60
	参考文献	61
第4章	信息协同过滤推荐系统	63
4.1	引言	63
4.2	基于内存的信息协同过滤推荐	66
4.2.1	基于用户的信息协同过滤	66
4.2.2	基于项目的信息协同过滤	70
4.3	基于模型的信息协同过滤推荐	72
4.3.1	基于降维技术的协同过滤推荐	73
4.3.2	基于聚类的协同过滤推荐	74
4.3.3	基于贝叶斯的协同过滤推荐	77
4.4	移动环境下基于隐式评分的博客推荐	78
4.4.1	问题的提出	78
4.4.2	相关工作	79
4.4.3	隐式评分的计算方法	80
4.4.4	基于隐式评分的协同过滤推荐算法	81
4.4.5	实验及结果分析	83
	小结	85
	参考文献	85
第5章	基于领域本体的商品信息推荐系统	87
5.1	基于领域本体的商品信息组织方法	87
5.1.1	问题的提出	87
5.1.2	商务信息资源特点的研究	88
5.1.3	商务信息的本体建模	90
5.1.4	商务信息语义互操作及其本体映射方法	94
5.2	基于领域本体的商品信息内容过滤推荐模型	102
5.2.1	商品推荐中的信息语义标记	102

5.2.2 基于内容过滤的语义信息推荐	107
5.3 基于领域本体和多属性决策方法的商品信息推荐模型	108
5.3.1 商品信息推荐模型	108
5.3.2 实验及结果分析	114
5.4 基于领域本体的商品信息协同过滤推荐模型	116
5.4.1 语义信息协同过滤推荐模型	116
5.4.2 实验与结果分析	119
小结	121
参考文献	121
第6章 基于Web挖掘的商品信息推荐系统	123
6.1 问题的提出	123
6.2 点击流相关理论和技术	125
6.2.1 点击流简述	125
6.2.2 基于点击流的商品信息个性化推荐服务	126
6.3 Web挖掘技术	127
6.3.1 Web挖掘简述	127
6.3.2 Web挖掘与商品信息推荐系统	131
6.4 基于Web挖掘的商品信息推荐模型	132
6.4.1 商品信息推荐系统的体系结构	132
6.4.2 商品分类树	133
6.4.3 基于点击流的顾客偏好分析	134
6.4.4 基于点击流的商品关联规则挖掘	136
6.4.5 商品信息推荐算法	138
6.5 商品信息推荐的实验及结果分析	139
小结	141
参考文献	142
第7章 基于案例推理的商品信息推荐系统	143
7.1 问题的提出	143
7.2 智能Agent	144
7.2.1 Agent技术概述	144
7.2.2 Agent的抽象结构	145
7.3 案例推理的决策支持	146
7.3.1 案例推理技术	146
7.3.2 基于案例推理的决策支持流程	149
7.3.3 基于案例推理的智能信息推荐	150

7.4 基于案例推理的商品信息推荐模型	152
7.4.1 基于 CBR 的系统解决方案	153
7.4.2 基于 CBR 的商品信息推荐系统结构	154
7.4.3 实例分析	164
7.5 基于 JADE 平台的推荐系统集成与 Web 应用	165
7.5.1 Agent 的系统集成	166
7.5.2 Web 应用设计	167
小结	167
参考文献	168
第 8 章 基于社会化标签的信息推荐系统	170
8.1 社会化标签系统与信息推荐	170
8.1.1 社会化标签系统概述	170
8.1.2 社会化标签系统的特点	174
8.1.3 社会化标签系统的实例	175
8.1.4 社会化标签推荐——信息推荐研究的新视角	176
8.2 基于社会化标签的相关信息推荐技术	178
8.2.1 基于协同过滤的标签推荐	178
8.2.2 基于内容过滤的标签推荐	178
8.2.3 基于图的标签推荐	179
8.3 基于社区标签云的信息推荐模型	180
8.3.1 基于社会化标签的聚类	180
8.3.2 基于社区标签云的个性化推荐	183
8.3.3 实例分析	185
小结	186
参考文献	187
第 9 章 基于情境感知的信息推荐系统	189
9.1 情境感知信息推荐的提出	189
9.1.1 情境感知推荐——个性化信息服务新模式	189
9.1.2 情境感知推荐的研究现状	191
9.2 融合多种情境的信息多维推荐服务模型	193
9.2.1 情境信息识别获取与语义描述方法研究	193
9.2.2 信息资源多维推荐服务模型	194
9.2.3 基于情境感知的信息资源推荐算法	195
9.2.4 信息多维推荐服务的系统体系结构	198
9.3 基于情境感知的个性化信息协同过滤推荐	200

9.3.1 基于情境感知的协同过滤推荐·····	200
9.3.2 实验及结果分析·····	204
9.4 基于情境感知的移动数字图书馆信息推荐·····	205
9.4.1 情境感知的移动阅读推荐——数字图书馆个性化服务新模式·····	205
9.4.2 基于情境熵的情境感知度·····	206
9.4.3 基于情境感知的协同过滤推荐·····	211
9.4.4 实验及结果分析·····	213
小结·····	215
参考文献·····	215

第 1 章 信息推荐系统概论

随着互联网规模和数字信息资源的不断增长,信息数量呈几何级数激增,信息服务领域面临“信息丰富、但有用信息获取困难”的窘境。一方面,网络中的海量信息资源出现信息过载现象(information overload);另一方面,用户被这些信息所包围着,却无法从中有效获取自己所需的信息资源。针对信息过载的问题,网络信息获取以现代信息技术为手段,向用户提供所需的信息资源,其服务模式包括信息拉取和信息推送两种类型的服务^[1]。信息拉取包括门户网站、信息检索和搜索引擎等,但这些工具只能满足主流需求,没有个性化的考虑,仍然无法很好地解决信息过载的问题^[2]。推荐系统(recommendation system)作为一种“信息推送”模式的重要方法,是当前解决信息过载问题的主要手段,它能够在分析预测用户需求的基础上主动推送其可能需要但又无法获取的有用信息,并能够以用户为中心,通过研究用户行为、兴趣和环境等,为用户推荐更具针对性的信息,即实现信息的“按需定制服务”。本章首先介绍网络信息获取及其提供的两种主要服务模式,在此基础上重点介绍基于“推送”模式的信息推荐系统,最后介绍信息推荐系统的应用领域和研究热点。

1.1 网络信息资源及获取服务模式

我们生活在一个越来越依靠信息的时代,并在向数字化时代迈进。数字化时代来临时,各种信息资源的电子化传递都将成为数字化经济的标志^[3]。信息社会化、社会信息化,信息生产与消费促进了信息产业和信息技术的飞速发展。当前,互联网已经成为人们获取信息的重要来源,是人们获取信息、改变生活方式、赢得商机的重要媒介。然而,互联网规模和信息资源的迅猛增长使得人们从互联网中有效获取信息日益困难。因此,针对用户的信息需求,如何高效地为其提供高质量的信息获取服务就成为当前互联网可持续发展的关键。

1.1.1 网络信息资源

网络信息资源是指以数字化形式记录的,以多媒体形式表达的,存储在网络计算磁介质、光介质以及各类通信介质上的,并通过计算机网络通信方式进行传递的信息内容的集合^[4]。随着 Internet 的普及和发展,网络信息资源对日常生活及商务活动起到越来越重要的作用。人们可以足不出户,方便地接触到大量的信息,信

息资源不足的问题已不复存在。但是,人们也深刻感受到,目前最大的问题不是信息资源的缺乏和不足,而是网络信息资源严重膨胀,呈现“信息过载”的现象。另外,用户被这些信息所包围着,却无法从中有效获取自己所需的信息资源,即用户面对的信息资源远远超出其处理的能力。现代信息服务领域已经进入一个“信息丰富、但有用信息获取困难”的两难窘境。因此,针对用户的信息需求,如何高效地为其提供高质量的信息获取服务已成为当前网络信息资源健康、可持续发展过程中亟待解决的重要问题。

目前,Internet 上的信息资源主要有以下几个特点^[5]:

(1) 存储结构的无序性。网络信息资源分散存储在 Internet 网络中不同的 Web 服务器中,采用不同的系统平台、数据管理平台、存取方式及人机接口,并缺乏集中统一的管理控制和质量控制,使得网络信息资源处于无序的分布式存储状态。

(2) 形式多样性。网络信息资源形式广泛,涉及人们生活的各个领域。传统信息资源主要是以文字或数字形式表现出来的信息。网络信息资源则以文本、图像、音频、视频等多种形式存在,涉及科学研究、商务、教育、娱乐、体育和艺术等不同领域。例如, Yahoo! 网站,其主页就是按字母排列的 15 个大类、464 个二级类目,每个二级类目又分为若干个三级类目。Sohu 网站与 Yahoo! 网站相似,它将所提供的信息资源分为 18 个大类,各个大类又细分为二、三、四级类目,各级子类目总数达 5 万多个,其主题几乎涉及各个行业和领域。

(3) 网络信息资源的实时性。网络信息发布者根据各自的需要和迎合用户的需求,经常更新信息,甚至变换发布信息的网页形式和内容。在增加网络信息的同时,也增加了其动态性,这些都给信息获取增加了难度。CNNIC 于 2012 年 1 月发布的《第 29 次中国互联网络发展状况统计报告》指出,2011 年中国网页数量已达到 866 亿个,年增长率为 44.3%。网络信息更新容易、传播迅捷,其实时性也是任何传统载体所无法比拟的。

(4) 信息冗余度高。由于大量的网络信息是免费的,而且相关信息网站没有统一的规划,所以各个网络站点之间存在着大量的冗余信息,这就造成了信息检索过程中对相同信息的重复检索,这既增加了网络信息资源建设的成本和开销,同时也增加了用户获取资源的困难,浪费了用户的时间。

(5) 信息资源的异构性。Internet 上包含海量的 Web 站点,这些 Web 站点的信息格式各异,有结构化的数据(如关系型的 Web 数据库)、半结构的数据(如 HTML 网页等)和非结构化的数据(如图片、广告视频等多媒体信息)。由于 Internet 开放式网络架构和网际互联协议的成功,如今人们可以轻易地在 Internet 上许可的范围内访问任何信息资源。但由于网络信息资源在数据结构定义、语法格式、语义描述、不同信息系统中的异构性以及访问语言上的异构性,大量的信息

资源并不能得到有效开发和利用,即存在“信息过载”现象。

Internet 的发展带来了信息资源的极大丰富,成为人们获取信息的重要来源,但是网络信息资源的无序多样、异构,以及冗余度高等特点也给用户从 Internet 网络中获取信息带来了困难。人们通常借助于各种网络信息获取工具,如门户网站、搜索引擎、专业数据索引、RSS 等从网络获取信息。网络信息获取服务是指在互联网上,针对用户的信息需求,以现代信息技术为手段,向用户提供所需的信息资源及服务,其服务模式包括信息拉取(information pull)和信息推送(information push)服务^[6]。

1.1.2 信息获取服务模式

信息拉取和信息推送是人们从网上获取信息的两种主要手段。信息的拉取和推送的区别在于用户获取信息的模式不同。信息拉取是用户根据自身的需求,通过信息服务系统搜索或浏览信息,经过不断地筛选和重定位,找到所需信息资源的过程;信息推送是指信息服务系统无须用户表明自身的信息需求,就能根据用户的历史访问记录以及所处环境等,主动在网上搜索信息,将符合用户需求的信息以合适的方式主动推送给用户的过程。可以看出,从信息系统服务的视角上,信息拉取是系统根据用户的请求,被动地地提供信息服务;信息推送则是系统主动地提供信息服务。两种服务模式之间的关系如图 1-1 所示。

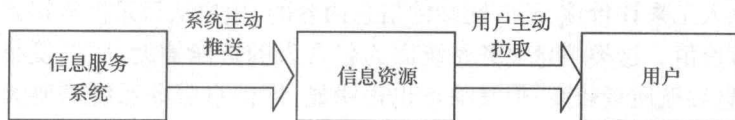


图 1-1 信息推送和信息拉取示意图

1. 信息拉取服务

信息拉取是用户获取网络信息资源的传统方式,它是指用户有目的地在网络上主动查询信息,其一般过程为:用户从浏览器给 Web 服务器发出请求,由 Web 服务器处理用户的请求,并将处理结果(即满足用户请求的信息或无法满足的信息)返回给用户^[7]。信息拉取服务的方式主要有两种:基于关键词检索的信息拉取服务和基于分类目录的信息拉取服务^[8]。

1) 基于关键词检索的信息拉取

这种信息拉取服务是搜索出符合检索条件的网页信息,信息服务系统服务器端通常采用基于机器人的技术,即使用一个被称为机器人(robot,也叫做 spider、Web crawler)的机器人程序自动访问 Web 站点。机器人程序以某种策略自动地在 Internet 中搜索和发现信息,由索引器为搜索到的信息建立索引,并根据用户的

查询输入检索索引库,并将查询结果返回给用户。这种拉取服务方式是基于网页的全文检索服务,其优点是信息量大、更新及时、不需要人工干预;缺点是返回信息过多,没有个性化的考虑,用户仍需要从检索出的信息中手工选择自己所需要的信息。另外,基于关键词检索的信息拉取仍会出现两个问题:一个是用户有时很难简单地用关键词来准确地表达所需要检索的内容,用户需求表达困难将导致检索困难;另一个是检索词的语义问题,同一概念可以用不同的检索形式来表达,如“计算机”和“电脑”,同时同一检索词在不同的上下文环境中可能语义不同,如“苹果”可以指一种水果,也可以指 Apple 公司旗下的“苹果手机”或“苹果电脑”。

2) 基于分类目录的信息拉取服务

这种信息拉取服务使用网站分类技术,即把网站进行树状的归类,对每个站点进行简略描述,形成分类目录。分类目录按网站的内容进行系统的分类整理,最终形成一个按类别编排的网站目录。在每一类中,排列着属于这一类别的网站的别名、网址链接、内容摘要以及子分类目录。同时,为了信息资源分类的科学准确,需要有相关各学科的专业人员对信息资源进行分类和维护。这些专业人员在访问了某个 Web 站点后撰写一段对该站点的描述,并根据站点的内容和性质将其进行归类,即把该 Web 站点的统一资源定位符(uniform resource location, URL)和描述放在这个类别中。当用户通过目录的某一节点进行检索时,用户安装目录的层次结构逐步细化直到找出满足用户需求的信息资源。由于基于分类目录的信息拉取服务是依靠人工来评价该 Web 网站的信息内容的,用户从目录搜索得到的结果往往更具参考价值。这类信息服务系统因为包含人的智能辅助,所以发布的信息较为准确,信息导航质量较高;但其缺点也很明显,即信息服务系统需要大量的人工介入,信息资源维护量大,因此信息更新不及时,维护成本较高。

搜索引擎是实现网络信息拉取服务的主要工具,它可以在面对拥有海量信息的网络环境时,辅助用户快速、高效地寻找有用信息。根据工作原理的不同,可以把搜索引擎分为两个基本类别:全文(fulltext)搜索引擎和分类目录(directory)搜索引擎。Google、百度都是比较典型的全文搜索引擎。分类目录搜索引擎则是通过人工的方式收集整理网站资料形成数据库,比如雅虎中国、搜狐、新浪、网易等均采用了分类目录和索引技术搜索信息资源。以 Google 为例,它是目前全球规模最大的搜索引擎,它支持用户以关键词的方式查询,提供包括网页、图像、视频、地图、学术、博客、电子商务等在内的信息搜索功能。同时,Google 也提供了许多智能化的搜索功能,如跨语言信息搜索功能,允许用户使用自己的本国语言搜索外文网站,从而解决了互联网上横亘在人们面前的语言障碍。百度搜索引擎是目前世界上最大的中文搜索引擎,它拥有超过千亿个的中文网页数据库,同样也支持用户以关键词的方式搜索信息资源,提供包括网页、图片、视频、MP3、地图、新闻等在内的信息搜索。目前,百度在国内各地分布的服务器,能直接从离用户所在地最近的

服务器上,把所搜索的信息返回用户,使用户享受方便、快捷的信息搜索服务。

除以 Google、百度等为代表的传统搜索引擎外,元搜索引擎(meta search engine)也是一种比较常用的信息搜索工具。元搜索引擎通过调用、控制和优化其他多个独立搜索引擎的搜索结果,并以统一的格式在同一界面集中显示。例如,国内开发的搜魅网(someta),推出一种全新的信息聚合服务,它集成 Google、百度、雅虎、搜狗等多家主流搜索引擎的结果,即将其他搜索引擎的返回结果利用自动聚类的方法聚合在一个独立的搜索界面上,为用户提供网页、图片、资讯等信息搜索服务。

从以上分析我们可以看出,信息拉取服务在当前网络环境下仍然居于主导地位,在网络数字化信息服务中依然会发挥重要的作用。在信息拉取过程中,用户是通过明确表达自身的需求后(如输入搜索关键词)才从信息服务系统中得到现有的信息服务。但是,信息拉取方式存在如下缺点:

(1) 信息拉取所获取的信息结果仍然是大量的,用户仍需对搜索结果进行人工过滤,以获取自己所需信息,这将耗费用户较多的时间和精力。

(2) 信息拉取从本质上是以“信息资源”为中心,而不是以“用户”为中心的信息获取行为,即需要用户主动从信息服务系统中获取信息资源。因此,用户不能随时进行搜索,无法得到及时的信息更新。

(3) 用户表达请求不准确,如采用基于关键词的信息搜索,无法用合适的关键词精准地表达自己的需求,因而导致信息服务的误差。

(4) 信息拉取方式没有考虑用户个性化需求,搜索的信息资源并不符合用户的个性化需求。

因此,从以上几点可以看出,现有的信息拉取及其实现工具(如信息资源门户网站、专业数据索引、搜索引擎)从本质上仅是帮助用户进行网络信息资源过滤的手段。这些工具只提供公共用户的一般需求,并没有针对目标用户的个性化需求提供定制的信息服务,因此仍然不能很好地解决信息过载的问题。

2. 信息推送服务

随着互联网上数字信息资源的迅速增长,基于“信息拉取”方式的搜索引擎返回的结果少则几百条多则上千条,甚至更多。用户通常需要不断手工构造复杂的查询条件以减少无关的返回结果。为减轻用户的负担,提高信息获取服务的质量,研究人员在研究信息服务系统新的搜索算法的同时,更关注用户的个性化需求和行为。因此,新的信息获取服务——信息推送服务应运而生。

1) 信息推送的定义

关于信息推送的研究目前已被国内学者所关注,关于信息推送的概念也有不同的解释和定义。其中具有代表性的有:①信息推送服务是利用推送技术(push

technology)自动搜索网络上用户感兴趣的信息,并主动推送到用户面前的服务,也可以称为基于“推”模式的网络信息服务。从技术上看,Push模式的网络信息服务是具有一定智能性、可以自动提供信息服务的一组计算软件,或者将其描述为网络环境下的一个高度专业化、智能化的网络专题信息服务系统^[9]。②信息推送服务利用推送技术主动把用户感兴趣的信息推送到用户端,与传统的信息拉取技术(pull technology)相比,减少了用户盲目的网上搜索时间,提高了信息检索效率^[10]。③信息推送服务就是通过一定的技术和协议,从网上的信息源或信息提供商获取信息,通过固定频道向用户发送信息的新型信息传播系统^[11]。④信息推送服务相对于传统的信息拉取服务而言,它是在“推”技术作用下信息找用户,而不是用户找信息^[12,13]。

我们可以将以上关于信息推送的定义进行综合,给出一个较为全面的定义:信息推送服务是信息服务系统通过识别和获取用户在信息检索过程中的行为和个性化需求特征,记录、学习并推导出用户的潜在需求和偏好,并及时动态追踪用户需求的更新情况,主动实时地把用户所需的信息资源推送给用户。

2) 信息推送的方式

根据信息推送采取媒介和方式的不同,可以将信息推送的方式分为以下几种^[10]:①频道推送服务。信息服务系统将互联网上一些内容相关的文档、网页以及多媒体信息等组合起来,通过特定的频道推送给用户查看。②页面推送服务。页面是互联网的基本组成单元,页面推送就是把一页面形式组成的信息内容推送给用户。③电子邮件推送服务。电子邮件推送主要利用了电子邮件的群发功能,将用户预定的或可能感兴趣的信息内容推送给相关用户群组。④专用式推送服务。专用式推送采取专门的信息收发软件进行推送,由信息员把信息直接推送给用户。⑤移动通信推送。将用户感兴趣的信息通过移动设备进行传送。

3) 信息推送技术的应用

信息推送技术最早于1996年由美国PointCast公司提出,它因而成为第一个在Internet上使用推送技术发布信息的公司^[14]。该公司通过与一些媒体公司合作,利用信息推送软件通过Internet网络向读者发送预先打包好的新闻、经济、体育和其他信息,如CNN、纽约时报、生活时尚等信息会在预先的频道中循环播出。信息推送技术最成功的应用是在一些特定的领域,针对特定的用户群体,如通过E-mail或短消息方式向特定用户提供新闻、天气、广告等。

RSS是一种起源于网景(netscape)的信息推送技术。由于版本的不同,RSS全称既可以是Really Simple Syndication(真正简单聚合),又可以是Rich Site Summary(丰富站点摘要)或RDF Site Summary(RDF站点摘要)。虽然三种规范定义的结构不同,但是它们所包含的核心信息和技术实质却基本相同^[15]。从本质上来讲,RSS是一种数据规范或结构,该规范规定网站在发布新信息的时候要遵

循的标准格式,以 XML 文件形式呈现某网站内容更新的摘要信息,是一种用于共享新闻标题、摘要等内容的 XML 文件。作为互联网上信息推送方式的实现,RSS 能够将新内容在服务器中出现的第一时间推送到用户端阅读器中,极大地提高了信息的时效性和价值。网上信息发布者,无论是企业还是个人,都可以通过 RSS 服务平台向所有用户“推送”出他们所需要的信息内容。另外,RSS 能够实现信息的“聚合”,即能将互联网上很多不同源信息以 feeds 订阅的方式集中到同一点^[16]。因为 RSS 是一种被广泛采用的内容包装定义格式,所以任何内容源都可以采用这种方式来发布信息。而在用户端,RSS 阅读器软件按照用户的喜好,有选择地将用户感兴趣的内容来源聚合到软件界面中,从而为用户提供多来源信息的“一站式”服务。基于 RSS 的信息推送服务使得大量经过筛选的高质量信息能够及时满足用户的需求,同时信息流的方式也不再是用户单一方向的“拉”,还包括反方向的“推”,从而提高了信息服务与信息利用的效率和效益^[17]。目前,基于 RSS 的信息推送服务被广泛地应用于网络在线新闻、电子报刊、电子学习和数字图书馆等领域。

信息推送技术还被广泛地应用于博客、论坛和电子商务中,用来为不同类型的用户推送广告、新闻、朋友和商务等信息。此外,信息推送技术也被应用于企业情报搜集、信息资讯等服务中^[6]。例如,企业竞争情报系统(enterprise competitive intelligence system),它将反映企业自身、竞争对手和企业外部环境的时间状态和变化的数据、信息、情报进行收集、存储、处理和分析,并以适当的方式推送给企业有关战略管理人员。

从以上关于信息推送的定义、推送方式和推送技术的应用中可以看出,信息推送服务与信息拉取服务的本质区别在于:信息推送服务是信息服务系统根据用户的需求为其主动推送所需的信息资源,服务是以“用户”为中心;传统的信息拉取服务则相反,信息拉取是用户根据自身需求向信息服务系统提出服务请求,通过信息服务系统在网络上寻找相应的信息资源,因此服务是以“信息资源”为中心。因此,在当前互联网上数字信息资源呈几何级数增长,信息服务领域面临“信息丰富、但有用信息获取困难”的窘境背景下,信息推送服务为信息服务学科提供了一种崭新的信息服务方式和解决问题的思路。

1.2 基于“信息推送”模式的信息推荐系统

1.2.1 信息推荐系统的概念与通用模型

随着互联网技术的应用普及和电子商务的迅猛发展,充斥在网络中的信息资源数量呈现指数增长的态势。海量的信息同时呈现在用户面前,使得用户感觉无