



农业信息 智能获取技术

岳峻 傅泽田 高文 ◎ 著



NLIC 2970735425



科学出版社

农业信息智能获取技术

岳 峻 傅 泽 田 高 文 著



NLIC 2970735425

科学出版社

北京

内 容 简 介

本书论述了农业领域信息的专题语义获取技术,共包含四部分内容。第一部分(第1~6章)从技术角度描述如何构建农业信息专用搜索引擎;第二部分(第7、8章)介绍本体理论与知识获取;在此基础上,第三部分(第9~13章)以蔬菜供应链为例,从技术角度提出了实现知识管理和知识语义获取的系统框架;第四部分(第14~18章)以鱼病诊断为例,提出了能够指导基于案例的鱼病诊断的智能案例获取和推理方法。

本书可供信息主管、知识主管、企业及政府管理人员、知识管理系统软件设计开发人员参考阅读,也可作为高等院校知识管理、企业管理、计算机软件等专业研究生与本科生的参考教材。

图书在版编目(CIP)数据

农业信息智能获取技术/岳峻,傅泽田,高文著.—北京:科学出版社,2011.5
ISBN 978-7-03-030860-3

I. ①农… II. ①岳… ②傅… ③高… III. ①农业科学—信息获取
IV. ①G252.7

中国版本图书馆 CIP 数据核字(2011)第 078261 号

责任编辑:孙 芳 / 责任校对:郭瑞芝

责任印制:赵 博 / 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>

丽 源 印 刷 厂 印 刷

科学出版社发行 各地新华书店经销

*

2011 年 5 月第 一 版 开本:B5(720×1000)

2011 年 5 月第一次印刷 印张:18 3/4

印数:1—2 000 字数:360 000

定 价:58.00 元

(如有印装质量问题,我社负责调换)

前　　言

有效获取和传播农业知识对我国现代农业的发展和新农村建设具有重要的意义。本书着重论述了利用现代信息技术进行农业信息获取的方法。全书共四部分。第一部分重点介绍了农业信息垂直搜索引擎的开发理论和方法。农业信息搜索引擎主要针对当前我国农业发展中所面临的小生产与大市场的矛盾,试图通过利用信息搜索促进农产品电子商务发展来缓解这一矛盾。针对我国农业信息搜索效率低、搜索成本高这一问题,本部分介绍了一种行之有效的服务于农业领域的农业信息垂直搜索引擎的开发理论和方法。与传统的农业信息搜索方法相比,该搜索引擎针对农业领域的特点进行了以下创新:①在分析网页时效性的基础上,提出了一种适用于农业专业领域网页数据采集模型。②在研究农业网页交易信息数据特点的基础上,提出了一种针对具有时空属性的信息表示与抽取模型。③在考虑垂直搜索引擎的索引对象具有领域性问题的前提下,提出了一种面向中文网页特点的高性能双字节倒排索引模型,保证了搜索的时间效率。④在借鉴用户个性化服务技术的基础上,针对当前搜索引擎检索质量存在的一些问题,引入检索结果自动分类与查询自动纠错技术来提高检索精度。第二至四部分别介绍了利用本体理论和智能推理理论来实现农业专业领域信息智能获取的方法。第二部分首先介绍了本体及知识获取的基本理论。第三部分和第四部分分别介绍了知识本体和智能推理理论在构建蔬菜供应链知识获取系统及鱼病诊断案例知识获取中的应用。第三部分从技术角度提出了一套能够实现蔬菜供应链知识获取系统的框架,主要包括蔬菜供应链本体模型、本体模型的形式化表示方法、领域概念获取推理方法、文本获取映射方法、蔬菜供应链知识获取系统的设计与开发。针对鱼病诊断案例的特点,第四部分提出了一套能够指导鱼病诊断案例知识获取系统的框架,主要包括诊断案例知识面向对象表示和向量空间模型、鱼病诊断本体论、鱼病诊断本体的构建与学习、诊断案例知识自动获取模型。这两部分的研究是当前语义网技术及智能推理技术在农业领域的具体应用。农业知识智能获取技术研究是当前信息技术在农业信息领域知识智能获取的探索,对开发专业领域的语义知识获取系统具有一定的借鉴意义。

在本书的撰写过程中,北京语言大学的胡亮老师和中国农业大学的李振波副

教授提供了不少有价值的研究素材,在此表示衷心的感谢!同时,本书出版获得了国家自然科学基金项目(60875039)的支持,并得到了鲁东大学信息科学与工程学院的大力支持和帮助,在此深表感谢!

由于作者水平有限,书中难免存在不妥之处,请同行和读者批评指正。

目 录

前言

第一部分 农业信息垂直搜索引擎

第 1 章 国内外农业信息搜索引擎现状	3
1.1 国内外农业相关的信息搜索引擎	4
1.2 相关农业信息搜索引擎的对比	7
1.3 面向主题的专用搜索引擎系统核心技术研究	8
1.4 服务于电子商务的搜索引擎在专业领域的应用	8
1.5 价格搜寻理论	9
1.6 农业信息搜索引擎开发目标与技术路线	11
1.7 本章小结	13
第 2 章 垂直搜索引擎的基本原理与技术	14
2.1 垂直搜索引擎系统架构特点	16
2.2 垂直搜索引擎开发关键技术	18
2.2.1 主题型网页数据采集技术	18
2.2.2 专业领域信息抽取技术	23
2.2.3 大规模文件索引技术	24
2.2.4 检索个性化服务技术	25
2.3 检索质量评估标准	27
2.4 本章小结	28
第 3 章 农业信息主题网页采集技术	29
3.1 农业信息主题网页的特点分析	29
3.1.1 农业交易信息来源	29
3.1.2 农业交易信息分类	30
3.1.3 农产品电子交易信息搜寻成本	30
3.2 数据采集与更新模型	32
3.2.1 网页时效性问题	32
3.2.2 数据更新频率	33
3.2.3 队列排序	35

3.2.4 区域负责机制	37
3.3 性能测试与评估	39
3.3.1 测试环境	39
3.3.2 实验结果	40
3.4 本章小结	43
第4章 时空属性信息过滤与抽取技术	44
4.1 农业信息数据特点分析	44
4.1.1 农产品交易信息网页	44
4.1.2 交易数据的时间与空间属性	46
4.2 特定结构化信息过滤与抽取模型	48
4.2.1 网页信息表示	48
4.2.2 包装器定义	49
4.2.3 K-EA 算法设计	50
4.3 性能测试与评估	52
4.3.1 评价指标	52
4.3.2 试验结果	53
4.4 本章小结	54
第5章 大规模文件索引技术	55
5.1 全文索引结构	55
5.1.1 位图	56
5.1.2 署名文件	56
5.1.3 倒排文件	56
5.1.4 后缀数组	57
5.2 双字节倒排中文索引模型	58
5.2.1 双字节倒排	61
5.2.2 虚拟内存硬盘缓存	63
5.3 性能测试与评估	66
5.3.1 评价指标	66
5.3.2 实验结果	66
5.4 本章小结	67
第6章 面向垂直搜索引擎的个性化检索服务技术	69
6.1 检索结果自动分类	69
6.1.1 农产品概念与分类问题	69
6.1.2 分类算法选择	70
6.1.3 K-近邻算法应用与改进	73

6.1.4 性能测试与评估	75
6.2 查询自动纠错.....	77
6.2.1 拼写错误问题	77
6.2.2 纠错原理与算法设计	77
6.2.3 性能测试与评估	79
6.3 本章小结.....	79

第二部分 本体论和知识获取

第 7 章 本体理论	83
7.1 本体的概念与内涵.....	83
7.2 本体的构建.....	86
7.3 本体表示语言.....	89
7.4 领域本体构建研究.....	93
7.5 本体自动获取相关理论.....	94
7.5.1 本体获取	94
7.5.2 本体获取分类	95
7.5.3 本体自动获取技术	96
7.6 本体学习	97
7.6.1 本体学习系统	98
7.6.2 本体学习基本原理与架构	99
7.6.3 本体学习系统结构	103
7.6.4 本体学习基本方法	104
7.7 本章小结	106
第 8 章 知识获取	107
8.1 知识获取方法	107
8.2 知识搜索推理方法	108
8.3 本章小结	110

第三部分 基于本体论的蔬菜供应链知识获取系统

第 9 章 蔬菜供应链	113
9.1 蔬菜供应链发展现状	113
9.1.1 发展现状	115
9.1.2 现状分析	117
9.2 蔬菜供应链知识获取系统构建框架	118

9.2.1 构建目标	118
9.2.2 技术路线	119
第 10 章 蔬菜供应链本体构建及形式化表示	121
10.1 蔬菜供应链及蔬菜供应链知识本体模型.....	121
10.1.1 我国蔬菜领域供应链模式	121
10.1.2 蔬菜供应链本体模型	123
10.1.3 蔬菜供应链知识本体模型	125
10.1.4 蔬菜供应链知识用户本体模型	126
10.1.5 知识、知识用户与知识背景本体间的关系	127
10.2 领域本体的形式化表示.....	128
10.2.1 RDF(S)形式化表示	128
10.2.2 Voronoi 图的形式化表示	129
10.3 本章小结.....	133
第 11 章 领域概念的获取推理方法	134
11.1 基于本体形式化表示的领域概念获取.....	134
11.1.1 RDF(S)下的定性推理	134
11.1.2 Voronoi 图下的定量推理	135
11.1.3 Voronoi 实验测评	136
11.2 基于模糊推理的领域概念获取.....	136
11.2.1 模糊推理	137
11.2.2 NSM 推理方法	143
11.2.3 改进的 NSM 推理算法	145
11.2.4 实验测评	146
11.3 基于 WordNet 的领域概念获取	147
11.3.1 同位关系与上下位关系	147
11.3.2 WordNet	150
11.3.3 局部线性嵌入领域概念提取算法	152
11.3.4 实验测评	154
11.4 本章小结.....	157
第 12 章 基于统计策略的文本搜索算法	158
12.1 统计语言建模.....	158
12.2 查询似然检索模型.....	158
12.2.1 投掷骰子的问题	159
12.2.2 基于查询似然的检索模型	160
12.2.3 数据平滑技术	161

12.3 查询似然检索模型在蔬菜供应链知识获取中的应用.....	161
12.4 本章小结.....	163
第 13 章 蔬菜供应链知识获取系统设计与实现	164
13.1 系统总体框架.....	164
13.2 系统开发工具与开发环境.....	166
13.2.1 Java 和 JDK	166
13.2.2 Eclipse	166
13.2.3 Tomcat	167
13.2.4 Protégé	167
13.2.5 Jena	167
13.3 系统模块设计.....	168
13.3.1 关键词检索	168
13.3.2 语义扩展检索	169
13.3.3 基于本体的语义检索	170
13.4 实验与结果分析.....	177
13.4.1 系统实现	177
13.4.2 结果分析	179
13.5 本章小结.....	181

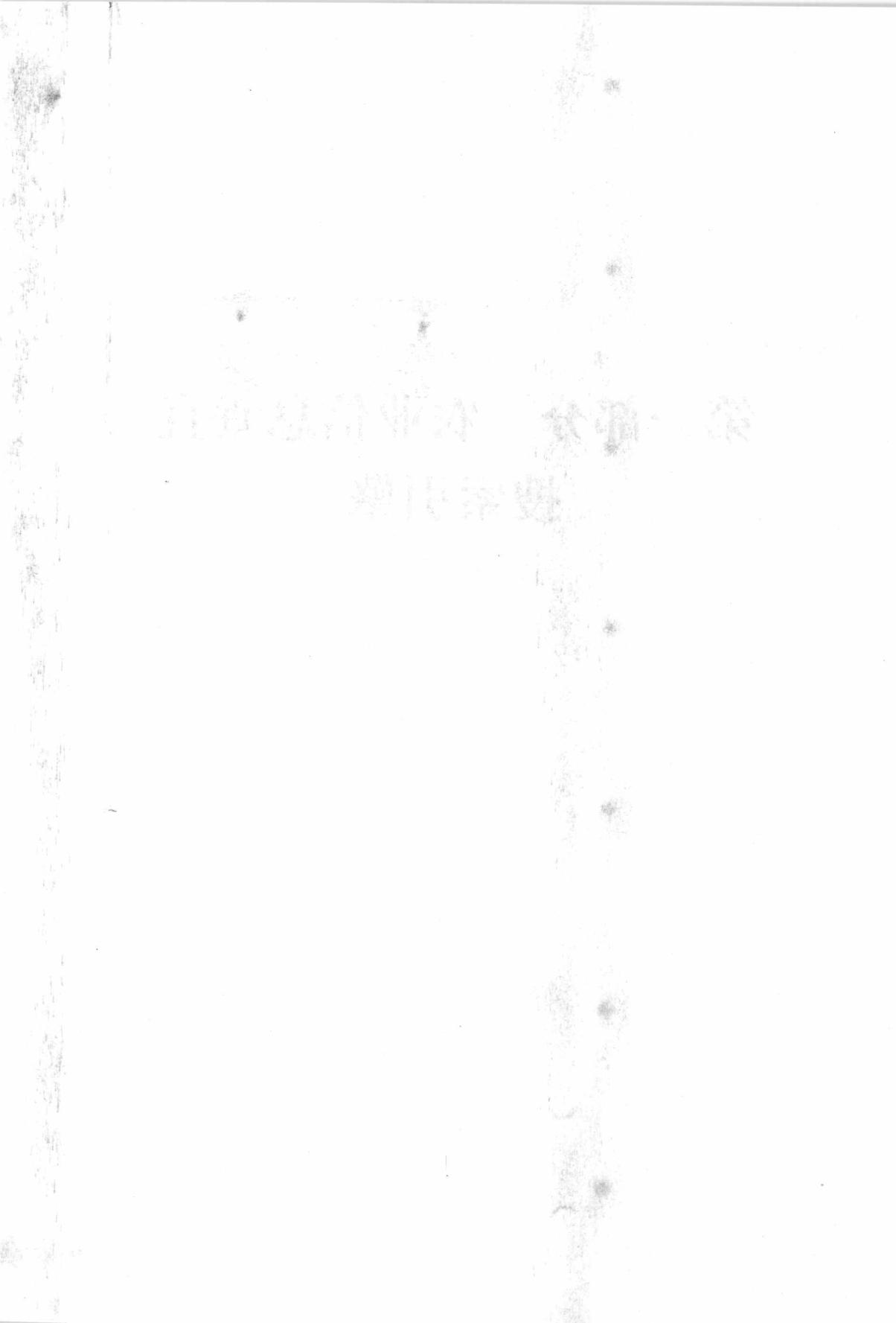
第四部分 基于本体的鱼病诊断案例知识获取

第 14 章 鱼病知识诊断与案例知识获取	185
14.1 鱼病知识诊断.....	185
14.2 CBR	186
14.3 机器学习方法与知识科学技术.....	187
14.4 鱼病诊断知识获取框架.....	190
第 15 章 基于本体的诊断案例知识表示	192
15.1 CBR 方法	192
15.1.1 CBR 系统	193
15.1.2 案例诊断系统中的案例知识获取	193
15.1.3 案例知识存在形式及源案例	194
15.2 诊断案例知识获取.....	195
15.2.1 诊断案例知识表示	195
15.2.2 诊断案例从非结构化到结构化的映射	196
15.3 诊断案例知识面向对象表示.....	197

15.4 案例知识及其语义.....	198
15.4.1 知识与语义	198
15.4.2 诊断案例知识语义定义及其语义层次模型.....	198
15.5 本体与语义.....	199
15.5.1 语义与本体的关系	199
15.5.2 本体在知识系统中的作用.....	201
15.5.3 VSM 及其语义化改进	201
15.5.4 非结构化、半结构化和结构化诊断案例知识的语义特征向量空间表示	204
15.6 诊断案例知识相似性度量.....	206
15.6.1 案例知识相似性关系的种类	207
15.6.2 传统案例相似性度量方法	208
15.6.3 基于面向对象模型的案例相似度计算	208
15.6.4 基于 VSM 的案例知识相似度比较	212
15.7 本章小结.....	214
第 16 章 鱼病诊断知识本体论	215
16.1 鱼病诊断知识本体模型.....	215
16.1.1 一般本体模型	215
16.1.2 鱼病诊断知识本体元数据定义	215
16.1.3 鱼病诊断本体元关系定义	216
16.1.4 鱼病诊断知识本体模型	217
16.1.5 鱼病诊断知识本体建模思想	218
16.2 鱼病诊断知识核心本体构建.....	219
16.2.1 核心诊断本体建模步骤与方法	219
16.2.2 基于 OWL 的鱼病诊断本体形式化模型	224
16.3 本章小结.....	226
第 17 章 诊断本体概念学习	227
17.1 基于关系模式和种子概念的鱼病诊断知识本体学习系统.....	227
17.1.1 基于关系模式的本体概念学习规则	229
17.1.2 基于种子概念面向文本的本体学习系统	233
17.1.3 实验分析和验证	240
17.2 本章小结.....	243
第 18 章 基于向量中心距离和 K-近邻算法的案例知识自动获取	244
18.1 复杂案例知识获取系统框架.....	244
18.1.1 诊断案例知识的特点	244

18.1.2 诊断案例知识获取系统框架	244
18.2 非结构化、半结构化诊断案例预处理及语义特征向量提取	245
18.2.1 诊断案例知识源文本化	246
18.2.2 非结构化诊断案例知识语义特征向量提取	246
18.2.3 案例特征向量约减——特征抽取技术	247
18.3 结构化案例知识的语义 VSM 构建	249
18.4 诊断案例知识库结构与案例知识组织	249
18.4.1 诊断案例知识结构	249
18.4.2 诊断案例知识库的组织	250
18.5 基于语义特征向量模型的诊断案例检索策略	251
18.5.1 基于语义特征向量模型的诊断案例知识检索思想	251
18.5.2 基于中心向量距离的非结构化、半结构化新案例知识学习算法	252
18.6 基于语义向量模型的非结构化诊断案例多类分类	253
18.6.1 分类模型	253
18.6.2 案例相似度计算	254
18.6.3 K-近邻算法文本分类器	256
18.6.4 实验结果与分析	258
18.7 本章小结	258
参考文献	259
附录	271
附录 A 互联网在中国蔬菜供应链中应用情况调查问卷	271
附录 B 中华人民共和国国家标准物流术语	276

第一部分 农业信息垂直 搜索引擎



第1章 国内外农业信息搜索引擎现状

以计算机为代表的信息技术在农业领域的商业化应用改变了传统市场的竞争格局,带来了新的机遇与挑战。尽快提高网络经济下的农业市场效率,增强国家农业竞争力,占领世界经济制高点,已经成为各国政府关注的焦点。但与此同时,我国农产品交易市场体系功能落后,传统农产品流通领域中存在的信息不对称问题导致了中介环节多、交易成本高及流通效率低等问题,极大阻碍了国内农产品交易市场的发展。

具体来讲,农业市场交易中的信息不对称问题主要体现在三个方面:①市场形成价格稳定性差。充分竞争的市场环境是有效价格形成的必要条件。我国农产品批发市场大多数规模较小,地域分散,产品和信息隔离,难以形成交易集中、市场透明度高、竞争充分的市场环境,导致农产品价格波动较大,区域价格差异明显。②市场交易方式原始。我国农产品批发市场仍主要采用“一对一”的对手交易方式,交易规模小,次数多,信息搜寻难度大,市场透明度低,形成的价格不能正确反映供求关系。③市场价格信息系统落后。许多批发市场的信息建设基本处于空白,信息的传递效率和共享度低,对信息的搜集、加工、处理、发布能力低下,容易造成信息的扭曲与失真。

加入WTO后,传统的农产品市场交易方式已经无法满足我国农业所面临的持续激烈的市场竞争环境及千差万别的顾客需求,使得中国农业发展所面临的小生产与大市场的矛盾更加突出,成为阻碍我国农业和农村经济健康发展、影响农民增收乃至农村稳定的重要因素(潘明等,2007)。怎样构建新型农产品市场交易模式、降低交易成本、提高交易效益等成为我国农业进一步发展亟待解决的重大现实问题(李晓等,2001)。农业电子商务作为一种运用电子化手段来进行农业商务活动的经济运行方式,能够将农产品交易活动网络化,减少传统方式中的大量中间环节,实现企业-企业、企业-消费者的直接交易,从而最大限度地降低经济活动中的交易成本,提高经济运转的效率和效益(廖咸真等,2005)。但是,传统的农业网上交易方式是商家在自己的网站上发布农产品价格信息,消费者从网站上找到自己感兴趣的农产品,直接与供应商联系购买,这样的模式简单易行,但非常依赖网站的访问量。之后又出现了集中式的农业电子交易中心,即由大量的商家登录自己的商品信息,由网站负责整个站点的宣传工作。尽管这种模式降低了成本,丰富了农产品信息,但消费者的选择范围仍然有限,特别是随着农业电子商务的普及与电子交易市场网站数量的增长,网页规模越来越大,导致用户要从海量信息中查找自

已需要的农产品的难度也越来越大,消费者在众多网站中寻找值得信赖的商家、获得最好的价值、最好的服务等问题也成为制约农业电子商务发展的因素,一般除了选择一些影响力较大的农业网站外,大多只能通过那些具有资金优势的网站的宣传来得到农产品信息。因此,在这种背景下,要求出现一种新的交易模式(郭云升,2007),能够从多个方面进行评比,包括网站访问量、顾客服务、商品价格、配送、客户反馈等,使得消费者在几秒钟内获得自己所需农产品的详细信息。

为了充分利用与扩展现有的农业信息资源,农业信息垂直搜索引擎综合服务平台将开发农业信息服务系统汇聚成虚拟数据库作为数据来源基础,从各农产品市场中采集原始的或经过分析处理的有用信息,进行筛选、加工、处理,再将各类信息进行整合、发布,并与其他农产品市场进行信息联网,使用户能从同一平台上获得即时、全面、有价值的信息。另外,平台还能提供各类信息增值服务,如信息的搜索、查询,同类产品销量、价格等的汇总、比较等,帮助用户减少信息搜寻成本,提高信息利用率,满足用户的多样化需求。

将所开发的农业信息垂直搜索引擎平台部署到网络上,通过系统应用为用户提供一个实际可用的农产品交易信息搜索比较平台,可以方便用户查找希望购买的农产品,降低消费者网络购物的时间成本,增强便利性;同时,提供丰富廉价的信息服务,打破了商家与消费者之间信息不对称的局面。平台包含的价格比较功能可以为商家提供最低成本的消费者来源渠道;系统还可以通过记录,分析大规模的消费者查找、比较和购买商品的行为过程,为农业企业提供精确细致的分析报告,借以提高企业的运营能力和销售技巧,从而进一步促进农业经济的全面发展,并为带动全国农产品交易的发展及农业信息化做出应有的贡献。

1.1 国内外农业相关的信息搜索引擎

1) 美国农业网络信息中心

美国农业网络信息中心(AgNIC)(<http://www.agnic.org>)是美国国家农业图书馆与一些大学、研究机构及政府机构自愿组成的联合体,其每个成员都负责农业科学中某一领域的信息工作,各成员单位间互相提供信息服务,每个成员在享受服务的同时,也有为其他成员提供服务的义务,服务方式主要是通过互联网相互提供电子形式的农业信息和检索服务。

(1) 数据库规模和范围。AgNIC 建有多个较为系统和完整的与农业相关的数据库,目前拥有 1100 多个数据库的记录。AgNIC 早先提供的也是数据库检索服务,此后升级为 Web 检索服务。

(2) 信息采集方式。AgNIC 的各个成员单位在其负责的专业领域下不断维护和增加数据资源,并存储在本地服务器上,同时,他们还要将这些新增加的数据

资源在 AgNIC 位于 Purdue 大学的中心服务器上反映出来,在其负责维护的 30 个大类下面的某个子类中添加相关信息,做好链接使别的单位可以访问。这种数据采集方式不仅使整个系统拥有的数据资源不断增加,同时将收集数据资源的工作由农业各个领域内处于领先地位的成员单位完成,因而能够确保数据的权威性和完整性。

(3) 检索功能。AgNIC 提供简单检索、高级检索和词表(thesaurus)检索三种界面。简单检索界面支持布尔逻辑检索,默认逻辑关系为逻辑与。高级检索界面比较复杂,有两个检索词输入项,中间有逻辑符号选择项(AND、OR、NOT)来确定两个检索词的逻辑关系,在下面还有三个选择项:前截断、中截断与完全匹配。高级检索还提供了检索字段选项,用户可以限定从某个字段中进行检索,共有题名、作者、摘要、关键词、主题词、全部字段等 6 个选择,这也为用户提供了更多的检索途径。由于 AgNIC 有自己的主题词表,用以描述它所收录的资源,用户输入自己的检索词后,系统会自动列出该词的相关词、上位词和下位词。例如,用户输入“crops”时,系统会提示“crops”的上位词有“plants”,下位词有“field crops”、“horticultural crops”与“specialty crops”,相关词有“crop acreages”、“crop losses”、“crop models”、“crop prices”等,这样就可以实现概念检索。

(4) 结果显示形式。AgNIC 检索的显示结果十分专业,在显示符合检索条件的记录总数后,再列出各条记录,每条记录首先给出题名,再给出简介、主题词、关键词,最后列出记录来源,即 URL 地址。

(5) 结果排列规则。AgNIC 检索结果的相关度排列规则是:首先看每条记录简介、主题词、关键词中出现的检索词的总数,然后根据出现频次的多少高低排序,AgNIC 检索认为,检索词出现频次高的记录相关性高,排列在前,检索词出现频次相同的记录再根据记录(网页)创建时间的先后排列。

2) Agriscape

Agriscape 网站(<http://www.Agriscape.com>)于 1999 年 4 月在美国普林斯顿建立,主要提供农业及相关产业的导航服务,其目标是发展成为农业信息、农业贸易和农业技术的信息中心。

(1) 数据库规模与范围。Agriscape 将所收录的内容划分为观光农业、公司、教育科研、图片、有机农业、拍卖、园艺、期刊、图书馆、政府、组织、导航目录等大类,各大类再进行细分,并提供会议、市场、新闻、天气等方面的服务。

(2) 信息采集方式。与其他搜索引擎直接用 robot 攀取网页并自动加入数据库不同,Agriscape 的信息完全是手工采集的。编辑人员对网站进行访问和获取,准确地描述网站的内容特征,并将其准确地归入不同类目之中。这种分析方法与 robot 从网页内容中提取关键词对网页进行描述的方法相比,确保了所收集网页信息的专指性,提高了搜索引擎的查准率,它还允许用户登记网站,在编辑人员审