

## 图书基本信息

书名：《智能Web算法》

13位ISBN编号：9787121139192

10位ISBN编号：7121139197

出版时间：2011-11

出版社：电子工业出版社

作者：Haralambos Marmanis,Dmitry Babenko

页数：374

译者：阿稳,陈钢

版权说明：本站所提供下载的PDF图书仅提供预览和简介以及在线试读，请支持正版图书。

更多资源请访问：[www.tushu000.com](http://www.tushu000.com)

# 《智能Web算法》

## 内容概要

本书涵盖了五类重要的智能算法：搜索、推荐、聚类、分类和分类器组合，并结合具体的案例讨论了它们在Web应用中的角色及要注意的问题。除了第1章的概要性介绍以及第7章对所有技术的整合应用外，第2~6章以代码示例的形式分别对这五类算法进行了介绍。

本书面向的是广大普通读者，特别是对算法感兴趣的工程师与学生，所以对于读者的知识背景并没有过多的要求。本书中的例子和思想应用广泛，所以对于希望从业务角度更好地理解有关技术的技术经理、产品经理和管理层来说，本书也有一定的价值。

## 作者简介

Haralambos (Babis) Marmanis 博士是一个把机器学习技术应用于工业界的先行者，也是供应链管理的世界级专家。Dmitry Babenko曾经为银行、保险、供应链管理与商务智能公司设计过应用与基础架构。本书拥有者可以通过 [www.manning.com/AlgorithmsoftheIntelligentWeb](http://www.manning.com/AlgorithmsoftheIntelligentWeb)在线获得作者的信息、样例代码与免费的电子版。

Dr. Haralambos (Babis) Marmanis is a pioneer in the adoption of machine learning techniques for industrial solutions, and also a world expert in supply management. He has about twenty years of experience in developing professional software. Currently, he is the director of R&D and chief architect, for expense management solutions, at Emptoris, Inc. Babis holds a Ph.D. in applied mathematics from Brown University, an M.S. degree in theoretical and applied mechanics from the University of Illinois at Urbana-Champaign, and B.S. and M.S. degrees in civil engineering from the Aristotle University of Thessaloniki in Greece. He was the recipient of the Sigma Xi award for innovative research in 2000, and he is the author of numerous publications in peer-reviewed international scientific journals, conferences, and technical periodicals.

Dmitry Babenko is the lead for the data warehouse infrastructure at Emptoris, Inc. He is a software engineer and architect with 13 years of experience in the IT industry. He has designed and built a wide variety of applications and infrastructure frameworks for banking, insurance, supply-chain management, and business intelligence companies. He received a M.S. degree in computer science from Belarussian State University of Informatics and Radioelectronics.

## 书籍目录

前言	XV
致谢	XIX
关于本书	XXI
1 什么是智能Web？	1
1.1 智能Web应用实例	3
1.2 智能应用的基本要素	4
1.3 什么应用会受益于智能？	5
1.3.1 社交网络	6
1.3.2 Mashup	7
1.3.3 门户网站	8
1.3.4 维基	9
1.3.5 文件分享网站	9
1.3.6 网络游戏	11
1.4 如何构建智能应用？	11
1.4.1 检查功能和数据	12
1.4.2 获取更多的数据	12
1.5 机器学习、数据挖掘及其他	16
1.6 智能应用中八个常见的误区	17
1.6.1 误区1：数据是可靠的	18
1.6.2 误区2：计算能马上完成	19
1.6.3 误区3：不用考虑数据规模	19
1.6.4 误区4：不考虑解决方案的可扩展性	19
1.6.5 误区5：随处使用同样的方法	19
1.6.6 误区6：总是能知道计算时间	

20	
1.6.7 误区7：复杂的模型更好	
20	
1.6.8 误区8：存在无偏见的模型	
20	
1.7 小结	
20	
1.8 参考资料	
21	
2 搜索	
22	
2.1 用Lucene实现搜索	
23	
2.1.1 理解Lucene代码	
24	
2.1.2 搜索的基本步骤	
31	
2.2 为什么搜索不仅仅是索引？	
33	
2.3 用链接分析改进搜索结果	
35	
2.3.1 PageRank简介	
35	
2.3.2 计算PageRank向量	
37	
2.3.3 alpha：网页间跳转的影响	
38	
2.3.4 理解幂方法	
40	
2.3.5 结合索引分值和PageRank分值	
45	
2.4 根据用户点击改进搜索结果	
47	
2.4.1 用户点击初探	
48	
2.4.2 朴素贝叶斯分类器的使用	
50	
2.4.3 整合Lucene索引、PageRank和用户点击	
54	
2.5 Word、PDF等无链接文档的排序	
58	
2.5.1 DocRank算法简介	
58	
2.5.2 DocRank的原理	
60	
2.6 大规模实现的有关问题	
65	
2.7 用户得到了想要的结果吗？精确度和查全率	
67	

2.8 总结	69
2.9 To Do	70
2.10 参考资料	72
3 推荐系统	73
3.1 一个在线音乐商店：基本概念	74
3.1.1 距离与相似度的概念	75
3.1.2 走近相似度的计算	80
3.1.3 什么才是最好的相似度计算公式？	83
3.2 推荐引擎是怎么工作的	84
3.2.1 基于相似用户的推荐	85
3.2.2 基于相似条目的推荐	94
3.2.3 基于内容的推荐	98
3.3 推荐朋友、文章与新闻报道	104
3.3.1 MyDiggSpace.com简介	105
3.3.2 发现朋友	106
3.3.3 DiggDelphi的内部工作机制	108
3.4 像Netflix.com那样推荐电影	114
3.4.1 电影数据集的介绍及推荐器	114
3.4.2 数据标准化与相关系数	117
3.5 大规模的实现与评估	123
3.6 总结	124
3.7 To Do	125
3.8 参考资料	127
4 聚类：事物的分组	128
4.1 聚类的需求	

129	
4.1.1	网站中的用户组：案例研究
129	
4.1.2	用SQL order by子句分组
131	
4.1.3	用数组排序分组
132	
4.2	聚类算法概述
135	
4.2.1	基于分组结构的聚类算法分类
136	
4.2.2	基于数据类型和结构的聚类算法分类
137	
4.2.3	根据数据规模的聚类算法分类
137	
4.3	基于链接的算法
138	
4.3.1	树状图：基本的聚类数据结构
139	
4.3.2	基于链接的算法概况
141	
4.3.3	单链接算法
142	
4.3.4	平均链接算法
144	
4.3.5	最小生成树算法
147	
4.4	k-means算法
149	
4.4.1	初识k-means算法
150	
4.4.2	k-means的内部原理
151	
4.5	鲁棒的链接型聚类（ROCK）
153	
4.5.1	ROCK简介
154	
4.5.2	为什么ROCK这么强大？
154	
4.6	DBSCAN
159	
4.6.1	基于密度的算法简介
159	
4.6.2	DBSCAN的原理
162	
4.7	超大规模数据聚类
165	
4.7.1	计算复杂性
166	

4.7.2 高维度	167
4.8 总结	168
4.9 To Do	169
4.10 参考资料	171
5 分类：把事物放到它该在的地方	172
5.1 对分类的需求	173
5.2 分类器的概述	177
5.2.1 结构分类算法	178
5.2.2 统计分类算法	180
5.2.3 分类器的生命周期	181
5.3 邮件的自动归类与垃圾邮件过滤	182
5.3.1 朴素贝叶斯分类	184
5.3.2 基于规则的分类	197
5.4 用神经网络做欺诈检测	210
5.4.1 交易数据中关于欺诈检测的一个用例	210
5.4.2 神经网络概览	212
5.4.3 一个可用的神经网络欺诈检测器	214
5.4.4 神经网络欺诈检测器剖析	218
5.4.5 创建通用神经网络的基类	226
5.5 你的结果可信吗？	232
5.6 大数据集的分类	235
5.7 总结	237
5.8 To Do	239
5.9 参考资料	242
6 分类器组合	



244	
6.1	信贷价值：分类器组合案例研究
246	
6.1.1	数据的简要说明
247	
6.1.2	为真实问题生成人工数据
250	
6.2	用单分类器做信用评估
255	
6.2.1	朴素贝叶斯的基准线
255	
6.2.2	决策树基准线
258	
6.2.3	神经网络的基准线
260	
6.3	在同一个数据集中比较多个分类器
263	
6.3.1	McNemar检验
264	
6.3.2	差额比例检验
266	
6.3.3	Cochran Q检验与F检验
268	
6.4	bagging: bootstrap聚合 ( bootstrap aggregating )
270	
6.4.1	bagging实例
272	
6.4.2	bagging分类器底层细节
274	
6.4.3	分类器集成
276	
6.5	boosting：一种迭代提高的方法
279	
6.5.1	boosting分类器实例
280	
6.5.2	boosting分类器底层细节
282	
6.6	总结
286	
6.7	To Do
288	
6.8	参考资料
292	
7	智能技术大汇集：一个智能新闻门户
293	
7.1	功能概览
295	
7.2	获取并清洗内容
296	

7.2.1 各就位、预备、开抓！	296
7.2.2 搜索预备知识回顾	298
7.2.3 一个抓取并处理好的新闻数据集	299
7.3 搜索新闻	301
7.4 分配新闻类别	304
7.4.1 顺序问题	304
7.4.2 使用NewsProcessor类进行分类	309
7.4.3 分类器	310
7.4.4 分类策略：超越底层的分类	313
7.5 用NewsProcessor类创建新闻分组	316
7.5.1 聚类全部文章	317
7.5.2 在一个新闻类别中聚类文章	321
7.6 基于用户评分的动态内容展示	325
7.7 总结	328
7.8 To Do	329
7.9 参考资料	333
附录A BeanShell简介	334
A.1 什么是BeanShell？	334
A.2 为什么使用BeanShell？	335
A.3 运行BeanShell	335
A.4 参考资料	336
附录B 网络采集	337
B.1 爬虫组件概况	337
B.1.1 采集的步骤	338
B.1.2 我们的简单爬虫	

338
B.1.3 开源Web爬虫
339
B.2 参考资料
340
附录C 数学知识回顾
341
C.1 向量和矩阵
341
C.2 距离的度量
342
C.3 高级矩阵方法
344
C.4 参考资料
344
附录D 自然语言处理
345
D.1 参考资料
347
附录E 神经网络
348
E.1 参考资料
349
索引
350

## 编辑推荐

算法是解决问题的一系列步骤。为实现有价值的Web应用（如推荐引擎、智能化搜索、内容组织系统等），本书提供了清晰的、精心组织过的算法模式。利用这些技术，你可以捕获用户原始而重要的信息，并把它们应用于实践中以获取相应的收益。用户数据中包含大量有价值的关联信息，它们往往无法通过人工观察而直观地获取，对于希望从这些数据中挖掘信息的Web开发者来说，玛若曼尼斯、巴宾寇编著的《智能Web算法》是一本很好的手册。作者作为一名Web开发者，拥有丰富的实践经验，加上多年来对机器学习领域技术的专研，使得本书对技术的解释清晰明了，读者可快速将其用于解决自己的问题。同时，本书提供的Java程序展示了如何搭建一个智能的应用，以及如何从用户的行为中进行学习，这是一笔现成的财富。

## 精彩短评

- 1、呵呵。这本书介绍了很多方面的算法。值得一看。
  - 2、书真心是好书。但是封面刮破半厘米的口子。
  - 3、书中的代码完全没必要。
  - 4、翻译的一塌糊涂
  - 5、去年年底看的，当时看了一半，是好书，以后慢慢啃。
  - 6、全是代码，不喜欢这样大段大段贴代码的书。
  - 7、书本质量不鏢
  - 8、偏重于实践
  - 9、作为web数据挖掘不错的入门读物
  - 10、对我来说开始阅读就是一个错误，坚持阅读下去是一个更大的错误。前150页左右尚在自己熟悉的范围之内，第一次真正观摩了rankpage和贝叶斯，认真阅读及理解书中代码，但是到了聚合，分类开始大量未知的概念涌现，一下子不知所措，跳过了大量实现细节，吸收了一下基础的概念算是此次阅读的目标吧。
- ps 附录中关于自然语言处理部分提供了不错的参考
- 11、个人觉得这本书写得实在不怎么样，译序中比较推崇作者大量使用代码而不使用公式的叙述方法，但在我看来一些零散的代码片断远不如一个公式来的简洁，易懂——其实这些公式并不涉及多么高深的数学知识，比如内积。我想，如果使用公式的话，这本书只需要三分之一的篇幅。

除去叙述风格外，我觉得书中所蕴含的知识也比较浅，各个算法都是点到为止。个人更推荐《大数据：互联网大规模数据挖掘与分布式处理》。

值得一提的是译者的水平确实很高，在我看过的译作中可算上乘！

- 12、很早买了，这次翻了翻。
  - 13、写的太差了。不过我第一次知道 PageRank 的细节是通过这本书。
  - 14、看看可以完全的颠覆传统的对数据库增删改操作的过程。
- 主要讲了以下四点web智能应用：搜索引擎，推荐系统，事物分组，分类器
- 15、难度较低，适合入门！
  - 16、书虽然看起来很厚，但是只有350来页，中间还有大量可有可无的代码，所以实际内容不是很多，涉及到的东西也基本点到为止，只适合做些web数据挖掘的了解了
  - 17、给老公买的工具书，还不错
  - 18、很浅显
  - 19、难得的一本web算法人们的好书
  - 20、本书带的代码很值得一看。。
  - 21、但是纸张也太对不起书的内容了吧
  - 22、对于我这种智商的人，入门还是挺好的。。。
  - 23、个人觉得一般，看完《集体智慧编程》再比较
  - 24、数据挖掘入门读物，通俗易懂，没深度
  - 25、书不错，里边主要是java实现，智能算法讲的也很不错，比较推荐。
  - 26、这本书比较不错，慢慢看，细细品。
  - 27、很一般
  - 28、很有用，常用算法都有介绍，学习中。
  - 29、主要的WEB应用算法都涉及到了
  - 30、其它的不说了，难得一见的好书
  - 31、尚未看，感觉挺多算法知识的，同志仍需努力
  - 32、纸张好差。。
  - 33、太过强调java代码而忽略了对算法本身的分析，中规中矩。
  - 34、搜索，推荐，聚类 and 分类算法入门。

- 35、纸张不错，拿起来很轻。
- 还只看了序言，译者对学习方法和很有心得
- 36、封面的设计很外国人文色彩
- 37、内容新颖，翻译还有待提高
- 38、我的专业需要这个，希望能有用
- 39、主要的WEB应用算法都涉及到了，值得细读
- 40、像个概论，代码太多反倒影响了阅读体验。有些东西还是要用数学公式的
- 41、买本看看，web编程不再是增删改。
- 42、对于智能web来说，算法的选择很重要
- 43、需要有一定的算法基础才能看哦
- 44、不错的书，只是还没有看
- 45、算法讲的听深入，不错
- 46、实用，有例子
- 47、里面的内容比较一般
- 48、刚收到书，大致翻了一下，应该很对胃口
- 49、搜索、推荐、聚类、分类等种种技术比较全面而先进，但缺乏理论分析，而通篇Java代码让人头疼。
- 50、读起来真心吃力。
- 51、全书只认真读了贝叶斯一章，自己实现了个分类器，里面居然有个公式错了。纸张很不好，像盗版的，翻译的也一般。比较好的是有代码。
- 52、通俗易懂，读完之后对于智能Web应用的开发，已经小有所成了
- 53、书真的挺好的，适合网络挖掘入门者学习与收藏。
- 54、翻了下，感觉挺好的，只是自己还木有时间深入的看
- 55、比较适合初学者看，尤其是对着作者写的源代码看，收获更大。
- 56、比较失望，堆砌了太多代码，而且很多不是算法代码，而是调用算法的代码.....
- 57、这本书看完更混乱了。。可能太旧了。
- 58、偏实践类书籍。
- 59、有点罗嗦
- 60、和我想的不太一样，貌似是我想差了，不过内容还不错。
- 61、智能技术的入门书，还有不少的网路资源，喜欢。
- 62、算是对数据挖掘有个了解吧。以后会复习。
- 63、书刚到，还没看，感觉还行！！
- 64、介绍的很清楚，值得看
- 65、很不错，可惜本人水平有限，看得不太懂。翻译的质量有待进一步的提高。
- 66、一句话，本书介绍的技术可操作性强。很适合对搜索引擎感兴趣的朋友。
- 67、有了人工智能的基础看这个可能比较轻松，没有相关经验可能看起来比较难。
- 68、这本书写得不错，值得一看。
- 69、给程序员的机器学习API说明书
- 70、这本书适合之前没有过相关经验，然后需要用最快的速度完成一个智能系统的人。这各系统只是模型，离实际应用还有相当的距离。
- 71、算法很重要，经典
- 72、结构清晰，层层递进，读起来很顺畅。
- 73、用beanshell简直上世纪
- 74、我只想说这是一本有用的书
- 75、我上午下的订单，晚上就拿到书了，很棒的物流。
- 书质量也不错，很好的初级启蒙式读本
- 76、搭配集体智慧编程阅读，了解计算机如何实现人工智能的入门书籍。两本书刻意回避数学公式，不过作为替代是通过代码呈现的，两者都不太感冒说不上好坏。吐槽下书中引用Aristotle的话，“We are what we repeatedly do. Excellence, then is not an act but a habit.” 诶不是Will Durant写的嘛。

- 77、高级程序员必要书籍
- 78、数据挖掘算法了解更多了，之前做lbs的数据挖掘时收获很大。。
- 79、好书一本
- 80、这本书内容上比较容易懂，不想一般的翻译过来的书那样晦涩和恶心，但是在阅读本书之前建议先读一下《Lucene in Action》或者其他有关Lucene的书籍。
- 81、很好，比较深入，细节也不错
- 82、很不错，阿稳翻译的，看看！
- 83、人工智能的书很少
- 84、只有思想，没有算法，看过了，也没留下多少东西。
- 85、内容丰富，浅显易懂，文笔还行啊~
- 86、图书馆借的书，挺不错的，仔细看过里面的推荐系统的部分。
- 87、可能过于注重实现了，理论上很不系统。
- 88、非常不错入门书。重要的是理解书中提到的思想。
- 89、对于初学者来说，看完此书，会对搜索，智能推荐有明显的认识和提高。书中，还有核心算法的部分实现。对于想继续钻研此领域的人来说，是本好书
- 90、还可以
- 91、这本书还不错，但是代码占篇幅较多，可以作为一本不错的概述型书。
- 92、参考~
- 93、算法很全。
- 94、研究这方面内容，看上去不错，希望有所收获~~
- 95、整体脉络还是不错的，有大把的代码例子，其实这些例子仅仅也就是个示例而已，方便理解原理，但有的地方关键代码居然展示不全，感觉这些代码的作用有限，和线上的使用相距甚远，又何必占用那么多篇幅，画个好点的图效果更好，有些地方也是一笔带过。。。。
- 96、找了很久，终于找到了，很好
- 97、推荐算法部分的内容太单薄，离我的要求太远。
- 98、感觉一般般，唉，我有点吹毛求疵了
- 99、代码较多，图较少，对于没有相关背景的人来说有点难以理解；只能用来普及一下概念，真正想弄明白还是得实践。PS: 分类部分需要实践一下...
- 100、虽然基础，但是学到很多
- 101、入门还是挺不错的

## 精彩书评

- 1、有朋友对构建本书中的代码运行环境有疑问，特别准备了一点介绍，为了格式上的方便，请访问这里：<http://gossipcoder.com/?p=842>
- 2、可以作为智能算法学习的起点，覆盖了搜索、推荐、聚类、分类等领域，有大量实用的示例代码，提供了很多扩展阅读的资源，以此为线索可以帮助我们循序渐进的深入智能算法的领域。不足之处：书中代码的部分常常没有事先说明思路，直接先上代码，而代码中琐碎无关的部分，以及排版格式影响了阅读效果。个人认为书中通过伪码+文字说明算法的精髓就可以了，细节的部分自己去看代码，这样会更好
- 3、花了半个多月的时间断断续续地看完了这本书，说说感受。1. 先说这本书的适用人群，在译者序里说是学生和需要梳理的工作者，但是在我看来，我觉得最佳的订位，应该是之前没有过相关经验，然后需要用最快的速度完成一个智能系统的人。因为本书把所有的知识简单化，当然随之的也是把所有的知识最浅化，如果完全采用本书的知识，搭建的应该是一个最初步的智能Web模型，看清，是模型。2. 这是我第二次看到用Java语言，而非Python or 伪码写的算法书，即便所有语言中，我的C#是最熟的，也不得不说如C#/Java语言真的不适合写这类算法书籍，这类语言过度地纠缠于代码结构，而本书更是过度地强调代码的组织，类的层次，而忽略了算法本身，全书像是一个开源项目的讲解，而且没有重点，这是我认为全书，也是相比于《集体智慧编程》最大的败笔。当然，如果您想快速搭建网站，上面的代码也许都是可用的，而不需要重写，这也许是一个差异化竞争了吧。



## 版权说明

本站所提供下载的PDF图书仅提供预览和简介，请支持正版图书。

更多资源请访问:[www.tushu000.com](http://www.tushu000.com)