

Search Engine Optimization Bible, 2nd Edition

搜索引擎优化宝典

(第2版)

(美) Jerri L. Ledford 著
马 煦 译

- 通过SEO提升网站在搜索引擎中的排名
- 针对移动网络和社会化媒体优化网站
- 锁定和获取目标客户



Bible

网站成功必备书籍



清华大学出版社

G354.4/54

2010

搜索引擎优化宝典

(第2版)

(美) Jerri L. Ledford 著

马 煜 译

清华大学出版社

北京

Jerri L. Ledford

Search Engine Optimization Bible, 2nd Edition

EISBN: 978-0-470-45264-6

Copyright © 2009 by Wiley Publishing, Inc.

All Rights Reserved. This translation published under license.

本书中文简体字版由 Wiley Publishing, Inc. 授权清华大学出版社出版。未经出版者书面许可，不得以任何方式复制或抄袭本书内容。

北京市版权局著作权合同登记号 图字： 01-2010-1180

本书封面贴有 Wiley 公司防伪标签，无标签者不得销售。

版权所有，侵权必究。侵权举报电话：010-62782989 13701121933

图书在版编目(CIP)数据

搜索引擎优化宝典(第 2 版)/(美)莱特福特(Ledford, J. L.) 著；马煜 译.—北京：清华大学出版社，2010.6

书名原文：Search Engine Optimization Bible, 2nd Edition

ISBN 978-7-302-22346-7

I. 搜… II. ①莱…②马… III. 互联网络—情报检索 IV. G354.4

中国版本图书馆 CIP 数据核字(2010)第 059273 号

责任编辑：王 军 赵利通

装帧设计：孔祥丰

责任校对：成凤进

责任印制：何 芊

出版发行：清华大学出版社 地 址：北京清华大学学研大厦 A 座

<http://www.tup.com.cn> 邮 编：100084

社 总 机：010-62770175 邮 购：010-62786544

投稿与读者服务：010-62776969, c-service@tup.tsinghua.edu.cn

质 量 反 馈：010-62772015, zhiliang@tup.tsinghua.edu.cn

印 刷 者：北京鑫丰华彩印有限公司

装 订 者：三河市金元印装有限公司

经 销：全国新华书店

开 本：185×260 **印 张：**26 **字 数：**633 千字

版 次：2010 年 6 月第 1 版 **印 次：**2010 年 6 月第 1 次印刷

印 数：1~4000

定 价：49.00 元

前　　言

欢迎阅读本书！本书为《搜索引擎优化宝典》的第2版。与所有的宝典系列图书一样，本书包含实用的教程和切实可用的信息，以及可供读者进一步学习的参考资料和背景信息。本书是搜索引擎优化的一本完全手册。读完本书，就可以优化网站或博客，获得最佳的搜索引擎排名。

搜索引擎优化的具体含义是一个仁者见仁、智者见智的问题。严格地说，SEO就是网页内外可以用来提高网站在搜索引擎中排名的各种设计策略。这通常意味着要利用各种设计元素和内容来调整网站。在大部分情况下，这也意味着无需任何额外的经费开销。

SEM——即搜索引擎营销(Search Engine Marketing)——并不仅仅是指SEO。准确地说，SEM还包括PPC，即竞价排名广告(pay-per-click advertising)。搜索引擎营销就是要利用各种方法，确保网站能在搜索引擎的搜索结果中获得尽可能高的排名。这意味着不仅需要对网站设计做出必要的修改，还需要利用其他各种策略，比如付费广告项目或是网站的内容策略。

本书将这些内容都归在一起。SEO和SEM的最终目标是将更多的人吸引到网站，而提高网站在搜索结果中的排名无疑能实现这个目的，此外也可以利用网络上日益流行的一种现象——社会化媒体(social media)。社会化媒体就是在网络上分享信息，并四处传播。这是一种与有同样兴趣的人分享信息和爱好的复杂方法。通过社会化媒体提高网站的访问量称为社会化媒体营销(Social Media Marketing, SMM)。

本书将详细介绍社会化媒体，此外还添加了移动网站营销的内容，因为移动网络正在迅速地发展。

抛开上面的这些文字游戏，所有的这些营销工作都有一个共同点：获得目标受众(target audience)。现在，所有人都将这些营销方法归于SEM，当然SEO纯粹论者不在此列。所有的这些方法都是针对目标受众对网站进行优化。随着社会化媒体和移动网络的逐渐盛行，社会化媒体自身的内容不仅会出现在搜索结果中，同时还会影响到网站在搜索引擎中的排名。

古往今来，要成功就必须做到与众不同，跟着别人亦步亦趋是不可能获得成功的。在SEO领域中也是如此，独树一帜往往能为网站赢得更多的流量。大家的目标是一样的，但实现目标的方法则可以千变万化。用不同的方法来做同样的事情，正是本书的初衷。

接下来会展示搜索引擎优化的最佳实践，以及介绍的策略所依赖的理论知识。我们根据这些策略构建了数千个网站，实际结果表明这些策略都是很有效的。要根据不同的实际情况，因地制宜地应用这些策略。我们要吸收前人的经验，同时也要有自己的创新。正因

为如此，博客、社会化书签和社会化社区才会变得如此流行。不同的人从不同的角度出发，对营销的看法也会不同。

本书的第2版添加了几章新内容，包括长尾搜索、如何创建社区提高网页在搜索引擎中的排名、把货币化网站作为一个SEO策略，以及一些有用的SEO插件信息。另外还更新并添加了关于购买过程各阶段的参考资料和相关信息，以帮助读者理解网站访问者进入网站时会查看什么信息，这样就能更好地了解这些访问者。

读者可以根据接下来所介绍的方法提高网页在搜索引擎中的排名，增加网站的流量，但最重要的是，要将网站的目标受众吸引到网站来。客户永远都是最重要的。将吸引客户作为目标，时刻关注客户，所付出的努力才会得到回报。

本书读者对象

搜索引擎优化的工作并不适合所有人，这需要耗费大量的时间和努力的工作。但它并不需要什么专业知识。任何有时间和兴趣的人都能够学会大部分的SEO策略。这也就解释了为什么现在会有如此多的所谓“SEO顾问”。

任何人都可以成为SEO顾问。不存在任何官方的SEO顾问认证，也没有任何相关的行业标准。从一定程度上说，这是个好消息，因为这意味着任何人都可以成为自己的SEO顾问。入门的最好方法就是认真地学习本书。

当然，并不是每个人都希望成为SEO顾问。您阅读本书的目的可能只是想了解SEO的大概流程，以确定自己的SEO顾问或正考虑雇佣的SEO公司是不是在尽职尽责地提高您网站的排名。本书同样也适合这样的读者。

这两种类型的读者都能从本书中获益匪浅：想成为自己的SEO顾问的人，以及只是想了解SEO的人。如果您已经是SEO专家，那么可能已经非常熟悉本书中所介绍的内容。尽管如此，可能仍有一部分新信息值得注意，所以如果您想温故而知新的话，就请继续学习吧。

对于那些SEO的新手，本书讲解了所有有助于提高搜索引擎排名和吸引优质客户的基本SEO策略及相应的理论知识。

本书的组织方式

搜索引擎优化是一个很复杂的过程，但是其中的各个部分又是相对独立的。本书分成4个部分，每个部分分别讲解了SEO流程的一个部分。

每个部分中的每一章都介绍了各个SEO步骤的一些具体细节。而每章中的各个小节则能帮助您更细致地学习这些细节。本书还有4个单独的附录，为各种不同的策略和措施提供指导和支持。

第Ⅰ部分：理解 SEO。本部分假定要完成某个 SEO 任务，或许是为网站创建 SEO，或许是熟悉 SEO，以了解 SEO 专家的工作是否有效。本部分介绍如下内容：

搜索引擎和搜索引擎优化的基本原理(第 1 章)

长尾搜索以及它对 SEO 的影响(第 2 章)

如何制定 SEO 方案(第 3 章)

第Ⅱ部分：SEO 策略。本部分学习网站或博客使用的各种 SEO 策略，包括常见的策略，例如建立 SEO 友好的网站，还有比较先进的策略，例如把社区用作一种 SEO 工具。本部分介绍如下内容：

建立 SEO 友好的网站(第 4 章)

使用有效的关键词(第 5 章)

利用竞价排名(第 6 章)

最大限度地利用竞价排名广告(第 7 章)

使用关键词提高转化率(第 8 章)

正确定位竞价排名广告(第 9 章)

管理关键词(第 10 章)

使用 3 种主要的竞价排名程序(第 11 章)

有效地为网站添加标签(第 12 章)

创建吸引人的内容(第 13 章)

把社区用作一种 SEO 工具(第 14 章)

建立有效的链接策略(第 15 章)

在理解了搜索策略的基本知识之后，就可以根据这些策略来吸引访问者和搜索引擎的注意力。

第Ⅲ部分：搜索策略的优化。本部分包含的 6 章将帮助我们的 SEO 工作取得成效。本部分介绍如下内容：

把网站添加到分类目录中(第 16 章)

确定付费收录服务是否适合自己的网站(第 17 章)

利用搜索引擎爬虫(第 18 章)

避免 SEO 作弊(第 19 章)

利用社会化媒体(第 20 章)

为移动网络用户优化网站(第 21 章)

确定网站货币化是否是 SEO 方案的正确策略(第 22 章)

使用 SEO 插件监控策略是否成功(第 23 章)

自动优化(第 24 章)

第Ⅳ部分：SEO 的维护。搜索引擎优化并不是一劳永逸的工作，需要持续不断地维护。本部分介绍如下内容：

维护 SEO(第 25 章)

分析 SEO 的成效(第 26 章)

除了这些章节之外，本书还包含 4 个附录。它们提供了关于 SEO 流程的其他一些信息和资源：

针对三大搜索引擎(Google、MSN 和 Yahoo!)的优化技巧(附录 A)

业内专家访谈(附录 B)

SEO 软件、工具和资源(附录 C)

SEO 工作表(附录 D)

约定和标识

本书中会出现一些图标，这些图标都蕴含着重要的信息。要特别注意这些特殊的记号。为了防止这些信息被忽略，它们都独立于一般的文字。本书包括以下 4 种图标：

注意 —— 注意中含有非必要但有助于理解 SEO 流程或理论的信息。

提示 —— 提示可以是完成某项工作的简便方法，或是一小段有助于理解相关策略和措施的文字。通常提示中的内容都有助于简化 SEO 工作。

警告 —— 要特别注意这些警告信息。认真阅读警告中的信息和建议可以避免很多麻烦。警告中会告诉您哪些事情不能做，以及哪些事情要特别小心。

引用 —— 表示交叉引用，可以从本书的其他地方找到与当前内容有关的知识。

所有的这些标识都是为了能帮助读者更轻松地学习 SEO。在遇到这些标识时，最好能花点时间认真阅读上下文。这些额外的信息应该能帮助您更好地理解 SEO。

如何阅读本书

实际上，在真正读完本书之前就可以将本书中介绍的一些策略逐步地应用到网站中。试试吧，但要记得将本书放在手边，随时备查。还要记得回到书中将剩余的章节阅读完。

还要记住实施 SEO 是一个不间断的过程。可以立即开始，但是要不断地努力，甚至在达到目标之后还要不断地努力。网站流量的大幅增加就是对辛勤努力的最好回报。比单纯的流量增加更好的事情是转化率(conversion rate)的不断提高。换句话说，就是会有更多的访问者将表现出您所期望的行为。

要实现这些目标并不容易，但是只要开始采取行动，随着时间的推移，就一定会看到显著的效果！

祝您好运！

目 录

第 I 部分 理解 SEO

第 1 章 搜索引擎基础	3
1.1 什么是搜索引擎	4
1.2 搜索引擎的基本结构	5
1.2.1 查询界面	5
1.2.2 搜索引擎结果页面	7
1.2.3 爬虫、蜘蛛和机器人	7
1.2.4 数据库	8
1.2.5 搜索算法	9
1.2.6 检索和排序	12
1.3 搜索的特征	14
1.4 搜索引擎的分类	14
1.4.1 主流搜索引擎	14
1.4.2 二级搜索引擎	15
1.4.3 专用搜索引擎	16
1.5 让搜索引擎为自己服务	16
1.6 控制搜索引擎	17
1.6.1 SEO 是一项艰苦的工作	18
1.6.2 安排 SEO 工作	18
第 2 章 长尾搜索理论	19
2.1 什么是长尾搜索	20
2.1.1 长尾理论的应用	20
2.1.2 长尾关键词的特征	23
2.2 长尾和尖头	24
2.3 自下而上地工作	24
2.4 组合使用所有要素	25
第 3 章 制定 SEO 方案	27
3.1 为什么需要 SEO	28
3.2 设定 SEO 目标	29

3.3 制定 SEO 方案	30
3.3.1 挑剔的细节	30
3.3.2 确定网页的优先次序	31
3.3.3 网站评估	31
3.3.4 完成方案	32
3.3.5 监督方案	32
3.4 理解自然 SEO	33
3.5 实现自然 SEO	34
3.5.1 网站内容	34
3.5.2 Google Analytics	35
3.5.3 内外部链接	36
3.5.4 用户体验	36
3.5.5 网站交互性	37

第 II 部分 SEO 策略

第 4 章 为 SEO 构建网站	41
4.1 建站之前的准备工作	42
4.1.1 明确目标	42
4.1.2 页面元素	43
4.2 网站优化	47
4.2.1 主机服务提供商很重要吗	47
4.2.2 选取域名的技巧	48
4.2.3 可用性	50
4.3 SEO 友好网页的构成	52
4.3.1 网站的入口页面和出口页面	52
4.3.2 使用醒目的标题	54
4.3.3 优质的内容	55
4.3.4 利用图片提升网站排名	56
4.4 容易出问题的页面和解决方法	57

4.4.1 痛苦的门户网站	57	6.3.2 竞价排名的使用	89
4.4.2 烦人的框架	58	6.4 竞价排名的分类	90
4.4.3 可爱又可恨的 cookie	59	6.4.1 关键词竞价排名	90
4.5 编程语言和 SEO	59	6.4.2 商品竞价排名	90
4.5.1 JavaScript	59	6.4.3 服务竞价排名	91
4.5.2 Flash	60	6.5 关键词的考察和选择	92
4.5.3 动态 ASP	60	6.5.1 关键词建议工具	93
4.5.4 PHP	61	6.5.2 不断地测试关键词	96
4.6 其他注意事项	61	6.6 选择高效的关键词	99
4.6.1 域隐藏	61	6.6.1 创建第一份关键词列表	99
4.6.2 内容重复	62	6.6.2 禁用搜索词和毒药词	100
4.6.3 隐藏页面	62	6.6.3 预测搜索量	102
4.6.4 404 错误页面	63	6.6.4 最终确定关键词列表	103
4.7 验证 HTML	63	6.7 撰写广告描述	105
4.8 建站后的注意事项	64	6.8 监视和分析结果	106
4.8.1 防止网站内容被窃取	64		
4.8.2 网站更新和改版	65		
第 5 章 网站的关键词	67	第 7 章 竞价排名策略的优化	107
5.1 关键词的重要性	68	7.1 关键词的摆放	107
5.2 什么是“启发式方法”	69	7.2 alt 标签及其他标签	108
5.2.1 模式、相近性和衍生	69	7.2.1 图片链接中的 alt 标签	109
5.2.2 启发式方法和网站可用性	71	7.2.2 title 标签	110
5.3 自然语言和布尔搜索	73	7.2.3 description 元标签	113
5.3.1 开始时使用布尔搜索	74	7.2.4 锚链文本	114
5.3.2 搜索语言自然地成熟起来	76	7.2.5 标题标签的内容	118
5.4 选择合适的关键词	78	7.2.6 正文	119
5.5 什么样的关键词密度才合适	80	7.3 URL 和文件名	120
5.6 自然关键词的使用	82		
5.7 避免关键词堆砌	83		
5.8 更多关于关键词优化的内容	84		
第 6 章 竞价排名与 SEO	85	第 8 章 增加关键词的成功率	123
6.1 竞价排名对 SEO 的影响	86	8.1 访问量和转化率哪个 更重要？	124
6.2 在竞价排名之前	87	8.1.1 设定目标	124
6.3 竞价排名的工作方式	88	8.1.2 实现转化	125
6.3.1 判断访问者的价值	88	8.2 竞价排名广告的文本	125
		8.2.1 类别词和商品词	126
		8.2.2 撰写广告词	127
		8.3 创建优秀的着陆页面	130
		8.4 理解和使用 A/B 测试	132

8.5 组合使用所有要素	133	11.2.1 Dashboard	169
第 9 章 理解和使用行为定位	135	11.2.2 Campaigns	170
9.1 什么是行为定位	136	11.2.3 Reports	171
9.2 行为定位的优势	136	11.2.4 Administration	172
9.3 如何利用行为定位	137	11.3 Microsoft adCenter	172
9.3.1 及时满足顾客的需要	137	11.3.1 Campaign	173
9.3.2 适时的重要性	138	11.3.2 Accounts & Billing	173
9.4 行为定位的其他技巧	139	11.3.3 Research	174
9.4.1 多个用户需要多种		11.3.4 Reports	175
定位方法	139	第 12 章 网站的标签	177
9.4.2 行为定位与隐私	140	12.1 网站标签为什么很重要	178
9.5 网站定向	140	12.2 标签的工作原理	178
9.6 使用网站定向广告	141	12.3 其他 HTML 标签	180
第 10 章 关键词和竞价排名的管理	143	12.3.1 nofollow	180
10.1 关键词预算	143	12.3.2 strong 和 emphasis	181
10.1.1 转化的价值	144	12.3.3 noframes	182
10.1.2 基于转化的预算	145	12.3.4 表格的 summary 标签	183
10.2 关键词出价的管理	146	12.3.5 acronym 和 abbreviation	
10.2.1 手动出价管理	147	标签	184
10.2.2 自动化出价管理	148	12.3.6 虚拟包含	184
10.3 关键词和转化的跟踪	151	12.4 重定向页面	186
10.4 降低竞价排名的成本	154	第 13 章 内容为王	191
10.4.1 竞价排名的管理	154	13.1 网站内容对 SEO 的影响	192
10.4.2 否定关键词	156	13.2 高质量内容的基本元素	194
10.4.3 时间定向	156	13.3 使用重复内容	196
10.5 提高点进率	158	13.4 远离搜索引擎作弊	198
10.6 竞价排名的投资回报率	160	13.4.1 门页	199
第 11 章 关键词工具与相关服务	161	13.4.2 隐藏文本和小文本	199
11.1 Google AdWords	163	13.4.3 反复提交网站	200
11.1.1 Campaign Management	164	13.4.4 网页劫持	200
11.1.2 Reports	166	13.4.5 页面偷换	200
11.1.3 Analytics	167	13.4.6 隐藏	201
11.1.4 My Account	167	13.4.7 隐藏链接	201
11.1.5 Print Ads	168	13.5 多语言网站	201
11.2 Yahoo! Search Marketing	168	13.6 内容管理系统	202

13.6.1	什么时候应该使用 CMS	202	16.2	地理定位 SEO 策略	240
13.6.2	如何选择合适的 CMS	203	16.3	提交工具的使用	241
13.6.3	CMS 对 SEO 的影响	203	第 17 章	付费收录服务	243
13.7	理解和使用“病毒”内容	204	17.1	什么时候应该使用付费收录服务	245
第 14 章	利用社区改进 SEO	205	17.2	付费服务的商业模型	245
14.1	社区的价值	206	17.3	付费服务的管理	246
14.1.1	社区的统计	206	17.4	选用正确的专业服务	247
14.1.2	用户的期望	207	17.5	合同中的注意事项	247
14.2	利用社区改进 SEO	208	17.6	合作失败的原因	248
14.2.1	建立对话平台	208	第 18 章	机器人、蜘蛛和爬虫	249
14.2.2	提高关键词的效用	209	18.1	什么是机器人、蜘蛛和爬虫	249
14.3	选择正确类型的社区	211	18.2	什么是机器人排除标准	251
14.4	正确地关注和维护社区	212	18.3	robots 元标签	253
14.4.1	准备阶段	212	18.4	使用 XML 网站地图使网页被收录	254
14.4.2	社区运转起来后	214	18.4.1	创建 XML 网站地图	255
14.4.3	维护社区	214	18.4.2	提交网站地图	258
第 15 章	网站中的链接	217	第 19 章	SEO 作弊揭秘	259
15.1	链接对 SEO 的影响	218	19.1	什么是 SEO 作弊行为	260
15.2	链接的原理	221	19.2	为什么 SEO 作弊不是个好主意	263
15.2.1	争取导入链接	222	19.3	避免 SEO 作弊	263
15.2.2	创建导出链接	224	第 20 章	社会化媒体优化	265
15.2.3	交叉链接的使用	226	20.1	什么是社会化媒体优化	268
15.2.4	毫无意义的链接场	228	20.2	社会化媒体的价值	269
15.3	创建链接的基础知识	229	20.3	社会化媒体策略	270
15.4	内部链接的使用	229	20.4	社会化媒体优化的评测	272
15.5	判断链接的效果	230	第 21 章	移动 SEO	275
第 III 部分 搜索策略的优化					
第 16 章	将网站添加到分类目录	235	21.1	移动用户的体验	276
16.1	什么是分类目录	236	21.1.1	移动网络	276
16.1.1	提交网站到分类目录	237	21.1.2	移动设备	276
16.1.2	主要的在线分类目录	239			
16.1.3	付费和免费	239			

21.1.3 移动用户使用网络的方式 277	25.1.4 内容更新 308
21.2 移动网站的设计 278	25.2 内容管理系统的使用 309
21.3 移动 SEO 280	25.3 SEO 的问题和解决方法 310
21.4 移动性的快速发展 281	25.3.1 网站被屏蔽啦 310
第 22 章 把访问量货币化作为 SEO 策略 283	25.3.2 内容剽窃 311
22.1 广告放置服务 284	25.3.3 点击欺诈 311
22.2 货币化服务概述 285	第 26 章 成效分析 313
22.3 SEO 的货币化策略 286	26.1 SEO 成效的分析 313
22.3.1 选择合适的货币化策略 286	26.1.1 SEO 期望的管理 314
22.3.2 在网站上添加货币化方法 287	26.1.2 自我定位 315
22.3.3 放置广告 287	26.2 网站统计分析 315
22.3.4 监控货币化的成功与否 288	26.2.1 基准统计 315
第 23 章 SEO 的插件 291	26.2.2 值得推荐的网站 315
23.1 理解插件 292	26.2.3 推荐关键词 (付费和自然) 316
23.2 选择正确的插件 292	26.2.4 访问的持续时间 316
23.2.1 Google 工具栏 293	26.2.5 访问的深度 316
23.2.2 Alexa 工具栏 294	26.2.6 重复访问 317
23.2.3 SEOQuake 296	26.2.7 其他统计 317
23.2.4 SEO for Firefox 296	26.3 竞争性分析 317
23.2.5 用于 Chrome 浏览器的 SEO 工具 299	26.4 转化分析 318
第 24 章 自动优化 301	26.5 服务器日志分析 319
24.1 应该自动化吗 301	
24.2 自动化工具介绍 303	
第 IV 部分 SEO 的维护	
第 25 章 SEO 维护基础 307	
25.1 还没有结束 307	附录 A 针对主流搜索引擎的网站优化 323
25.1.1 排名监控 308	附录 B 业内访谈 329
25.1.2 网站分析 308	附录 C SEO 软件、工具和相关资源 365
25.1.3 关键词和链接监控 308	附录 D 工作表 377
	术语表 391

第Ⅰ部分

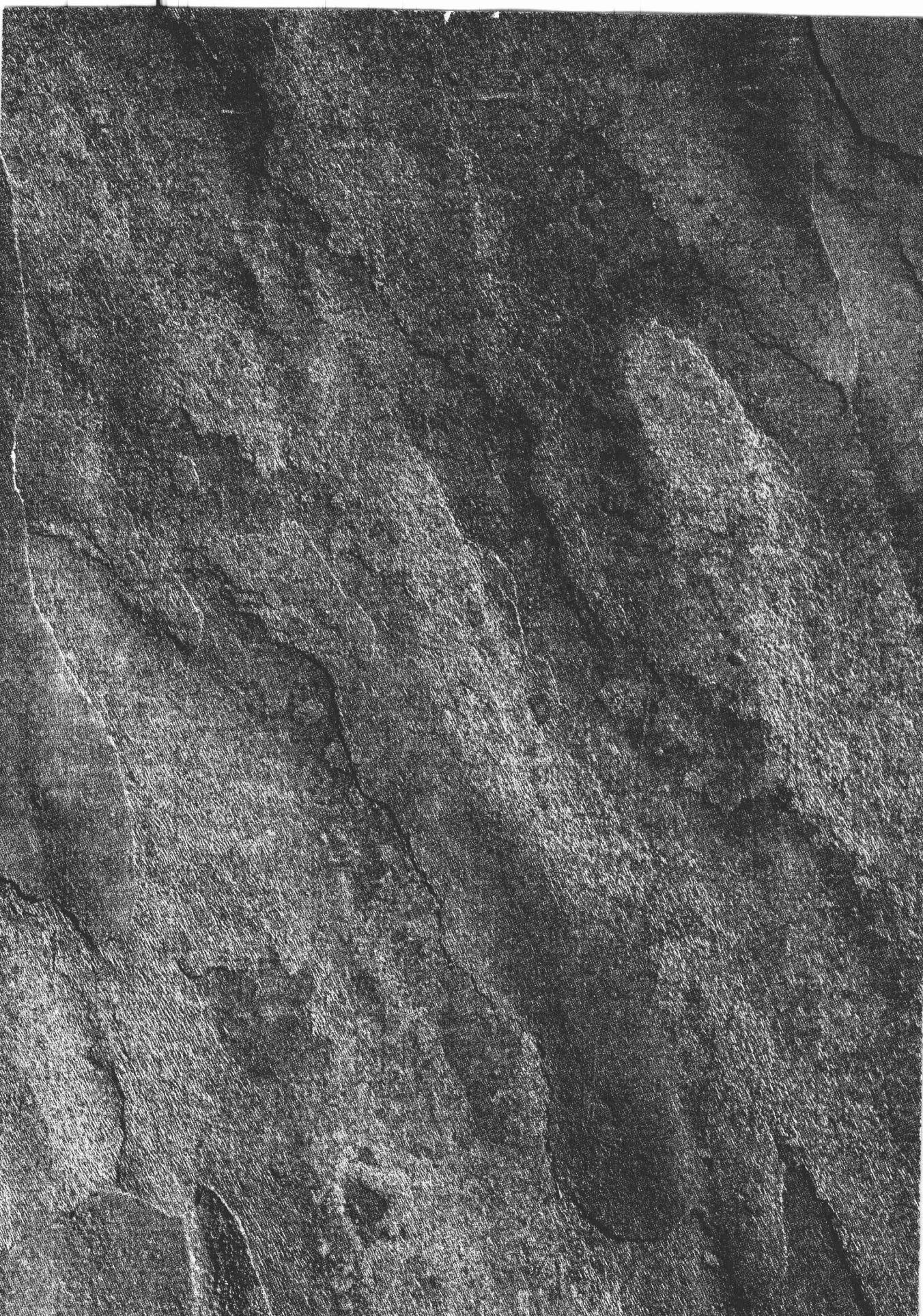
理 解 SEO

搜索引擎优化(Search Engine Optimization, SEO)是一个很宽泛的概念，很难简单地用几句话就说明其全部的含义。搜索引擎优化涉及到很多内容，包括搜索引擎的工作原理(以及不同搜索引擎之间的差异)、网页的设计等。要面面俱到地学习各个方面的知识，在时间上是不现实的。不过，搜索引擎优化并不是一项不可能完成的任务。但如果完全不知道它是什么以及它的原理，那就不可能实现 SEO。

第Ⅰ部分介绍了搜索引擎优化的基础知识。这部分内容对什么是搜索引擎以及搜索引擎的原理做了介绍，还解释了什么是长尾搜索以及 SEO 方案。将这些内容结合在一起，就能知道如何正确地实施 SEO 策略来提高网站的流量。

本部分包括

- 第 1 章 搜索引擎基础
- 第 2 章 长尾搜索理论
- 第 3 章 制定 SEO 方案



第1章

搜索引擎基础

内容提要：

- 什么是搜索引擎
- 搜索引擎的基本结构
- 搜索的特征
- 搜索引擎的分类
- 利用搜索引擎
- 控制搜索引擎

如果要在互联网上寻找信息——例如事件、统计数字、描述、商品甚至电话号码，您会怎么做？大部分情况下都会使用某个主流搜索引擎，输入需要查找的单词或短语，然后逐个点击搜索结果，不是吗？不一会儿，需要查找的信息就出现了，非常奇妙！当然，搜索引擎并不一定是最优的选择。

早期的互联网并不是现在这样的。实际上，当时的互联网并不像现在这样一个由相互连接的网站构成的网络，也没有成为如此庞大的商业助推器。当时所谓的互联网只是一些用户可以下载(或上传)文件的FTP(File Transfer Protocol, 文件传输协议)站点。

要在这些站点中寻找某个文件，用户只能逐个地浏览每个文件。当然，也有简便的方法。如果您认识某个知道您所需文件确切地址的人，就可以直接获得这个文件。这里的前提是明确地知道所寻找的是什么文件。

这使得在互联网上查找文件成了一件困难耗时、且极度考验耐性的事情。在坐落于蒙特利尔的 McGill 大学中，一个学生决定要简化这个工作。1990 年，这位名叫 Alan Emtage 的学生创建了互联网上的第一个搜索工具。他的杰作是一份互联网上各种文件的索引，名叫 Archie。

提到 Archie，如果您想到的还是 1941 年的那个漫画人物，那就有些落伍了(至少现在

是这样)。之所以叫 Archie，是因为原来的名称 Archives 太长了。不只是 Archie，实际上那一系列漫画书中的其他角色(Veronica 和 Jughead)随后也进入了搜索领域，这里就不详述了。

Archie 不同于现在所使用的搜索引擎，它并不是真正的搜索引擎。但在那个时候，它是很多互联网用户的挚爱。这个程序在指定的网络中下载匿名 FTP 站点的文件列表，然后将这些列表存储到可以进行搜索的网站的数据库中。

跟现代搜索引擎不同，Archie 没有自然语言处理能力(natural language capabilities)，但在那时这已经是一件伟大的工作了。Archie 索引计算机文件，方便了文件的查找。

到了 1991 年，明尼苏达大学一位名叫 Mark McCahill 的学生觉得，如果能在互联网上搜索文件，也就一定能在文件中搜索纯文本。为了填补这个空白，他创建了 Gopher 用于索引纯文本文档，并发展成了互联网的最早网站之一。

随着 Gopher 的创立，人们还需要能在 Gopher 创建的索引中查找相关信息的程序，于是 Archie 的小伙伴们登场了。人们开发了在 Gopher 索引系统中搜索文件的 Veronica(Very Easy Rodent-Oriented Net-wide Index to Computerized Archives)和 Jughead(Jonzy's Universal Gopher Hierarchy Excavation and Display)。

这些程序的基本原理都是一样的，用户使用关键词搜索文件的索引信息。至此，搜索技术开始逐步走向成熟。第一个现代意义上的搜索引擎是 Matthew Gray 于 1993 年开发的名为 Wandex 的搜索引擎。Wandex 是第一个同时具有网页索引和搜索功能的程序。这是第一个使用了网络爬虫技术的程序，也成为后来各种搜索爬虫的基础。从那以后，搜索引擎就开始蓬勃发展起来。从 1993 年到 1998 年，为人们所熟知的搜索引擎主要有：

- Excite——1993 年
- Yahoo! ——1994 年
- Web Crawler ——1994 年
- Lycos ——1994 年
- Infoseek——1995 年
- AltaVista——1995 年
- Inktomi——1996 年
- Ask Jeeves——1997 年
- Google——1997 年
- MSN Search——1998 年

今天，搜索引擎已经非常成熟，可以用日常的单词或短语来搜索各种文件和文档。看看现在搜索引擎的强大搜索能力，很难让人相信搜索引擎只有 15 年的短暂历史。

1.1 什么是搜索引擎

好的，现在已经知道了搜索引擎的基本概念。在搜索框中输入单词或短语，然后单击按钮，稍等片刻，就会看到成千上万的相关网页。接着要做的就是打开这些网页，寻找所

需要的内容。但是除了“搜索即可找到”这个泛泛的概念外，搜索引擎的准确定义是什么？

这有点复杂。在搜索引擎的后台，有一些用于搜集网页信息的程序。所收集的信息一般是能表明网站内容(包括网页本身、网页的 URL 地址、构成网页的代码以及进出网页的链接)的关键词或短语。接着将这些信息的索引存放到数据库中。

而在前端，是供用户输入搜索词(单词或短语)的用户界面。当用户单击“搜索”按钮时，算法就会在后台的数据库中查找信息，将与用户输入的搜索词相匹配的网页链接呈现给用户。

引 用

第 18 章会深入地介绍网络爬虫、网络蜘蛛以及网络机器人。

搜集网页信息的程序称为爬虫(crawler)、蜘蛛(spider)或机器人(robot)。爬虫会遍历网络中未屏蔽的 URL，并收集每个网页中的关键词和短语，然后将这些信息存放到搜索引擎的数据库中。想一下，互联网上的网站数量早已超过 1 亿个，而且还在以每月超过 150 万个新网站的速度增长。这就像是要用大脑将所见到的每一个单词都进行分类，需要的时候再将所有相关的信息调出来。

简单点说，这几乎是不可能完成的任务。

1.2 搜索引擎的基本结构

现在读者应该对搜索引擎的原理有了粗略的了解，但真正的搜索引擎远比您想象的复杂。实际上，搜索引擎是由多个部分组成的。然而，很难找到关于搜索引擎结构的资料，这些资料是搜索引擎公司严密保守的商业秘密，但这些资料对于搜索引擎优化(SEO)是非常重要的。

1.2.1 查询界面

查询界面(query interface)是人们最熟悉的部分。当人们提起“搜索引擎”时，想到的通常也是搜索引擎的查询界面。查询界面就是用户访问搜索引擎时输入搜索词的页面。

以前搜索引擎的界面就像图 1-1 中 Ask.com 的网页。其界面只是一个简单的网页，只有一个搜索框和一个启动搜索的按钮。

现在，网络上的很多搜索引擎的查询界面中都加入了越来越多的个性化内容，以增强其功能。例如图 1-2 中的 Yahoo!Search，用户可以根据自己的需求自行定制搜索页面，包括免费电子邮箱账户、天气信息、时政新闻、体育新闻等各种能吸引用户使用搜索引擎的元素。

另一种定制搜索引擎界面的方式类似于 Google 提供的功能。用户可以根据自己的需求和喜好在 Google 搜索引擎的主页上添加各种小程序。