

2.6	倒排档检索机制的加强	2-19
2.6.1	邻接	2-19
2.6.2	截词	2-21
2.6.3	范围检索	2-21
2.6.4	加权	2-21
2.7	商业性检索系统介绍	2-22
2.7.1	DIALOG 系统	2-23
2.7.2	STAIRS 系统	2-24
2.7.3	MEDLARS 系统	2-29

第三章	文献情报检索的数据结构和检索技术	3-1
3.1	情报检索中的数据结构	3-1
3.1.1	逻辑结构与物理结构	3-1
3.1.2	线性表	3-5
3.1.3	树	3-8
3.1.4	图	3-13
3.2	查找技术	3-14
3.2.1	顺序查找	3-15
3.2.2	基于索引的方法	3-17
3.2.2.1	二分法查找	3-18
3.2.2.2	分块查找法	3-20
3.2.2.3	索引顺序法	3-23
3.2.2.4	B-树	3-28
3.2.3	基于 Hash 的查找方法	3-29
3.2.3.1	碰撞问题及其解决	3-30

3.2.3.2	截词检索	3-34
3.2.3.3	Hash法与情报检索	3-36
第四章 检索效果及其改善		4-1
4.1	检索效果及其测量指标	4-1
4.2	影响检索效果的主要因素	4-5
4.2.1	情报提问对情报需求的表达程度	4-6
4.2.2	数据库的选择和比较	4-8
4.2.3	检索途径的选择	4-9
4.2.4	检索词的选择与调节	4-9
4.2.5	检索式的结构	4-11
4.3	提高检索效果的反馈调整方法	4-12
4.3.1	反馈调整在检索过程中的作用	4-12
4.3.2	调节检索策略的若干方法	4-15
第五章 自动标引		5-1
5.1	自动标引和人工标引	5-1
5.2	西文自动标引方案简介	5-3
5.2.1	词频统计原理	5-3
5.2.2	逆文献频率法	5-6
5.2.3	信号——噪音法	5-7
5.2.4	词辨别值法	5-10
5.2.5	词短语的构造	5-16
5.3	自动标引中的词表	5-18

张庆国

第六章 聚类检索	6-1
6.1 问题的提出	6-1
6.2 SMART 系统	6-1
6.2.1 文献的向量表示和匹配度计算	6-2
6.2.2 聚类文件的生成和 SMART 系统的文档结构	6-4
6.2.3 提问式的反馈调整	6-10
6.2.4 动态文献空间	6-14
6.2.5 聚类检索和分类检索的区别	6-15
6.3 倒排检索和聚类检索的结合	6-16
6.3.1 SIRE 系统	6-16
6.3.2 加权的布尔检索	6-20
第七章 检索效果的改善(续)	7-1
7.1 文献——语词矩阵的若干推论	7-1
7.1.1 词联接矩阵	7-1
7.1.2 词结合矩阵和改良型文献——语词矩阵	7-2
7.2 与词结合矩阵相关的权和提问向量	7-4
7.3 通过结合反馈进行的提问自动修正	7-8
7.4 检索策略的最优化	7-11
第八章 数据检索系统	8-1
8.1 概论	8-1
8.2 数据库管理系统的结构	8-4
8.2.1 信息项的结构	8-4
8.2.2 关系数据库模式	8-8

8.2.3	层次数据库模式	8-13
8.2.4	网络数据库模式	8-18
8.3	查询和查询语言	8-19
8.3.1	分步法	8-21
8.3.2	“菜单”方法	8-22
8.3.3	表查询法	8-23
8.3.4	例举查询	8-24
第九章 事实检索		9-1
9.1	事实检索和自然语言处理	9-2
9.2	自然语言处理的句法分析系统	9-3
9.2.1	自然语言的处理层次	9-3
9.2.2	短语结构语法	9-4
9.2.3	转换语法	9-9
9.2.4	扩充转换网络语法	9-12
9.3	知识的表示	9-18
9.4	目前水平上的事实检索系统	9-23
第十章 情报信息的存贮和输入输出		10-1
10.1	数据标识的代码化	10-1
10.2	数据库的存贮载体	10-1
10.2.1	磁带数据库	10-5
10.2.2	磁盘数据库	10-7
10.2.3	其他存贮设备	10-8
10.3	情报资料的输入手段	

10. 3. 1	键到纸介质方式	10—8
10. 3. 2	键到磁介质方式	10—9
10. 3. 3	联机终端输入方式	10—10
10. 3. 4	全自动字符识别方式	10—11
10. 3. 4. 1	光学字符识别法	10—11
10. 3. 4. 2	光学标记阅读装置	10—14
10. 4	情报资料的输出手段	10—15
结 语		10—17

第五章 自动标引

5.1 自动标引和人工标引

情报检索的最终目的在于找出满足用户需要的文献资料，但是在检索系统中，文献却是以其主题的标识化形式——标引词来表示的，检索过程也是通过情报提问中的检索词与检索系统中的标引词的匹配而实现的，标引词（检索词）是联接文献与提问的中介，是联系情报与其用户的中介，由此可见标引工作在情报检索工作中的重要地位。

对于标引工作，存在着标引质量和标引效率两方面的要求。

标引质量，所选中的标引词能否全面、准确地反映文献的主题，决定着该篇文献能否正确地为用户通过该篇文献的标引词集合而了解原文献，决定着检索系统能否正确地将检索提问与其匹配，从而决定着该篇文献能否为需要它的用户检出。

标引效率，由于标引工作在质量上的极高要求，由于标引工作在自身上的难度，标引工作的效率是比较低的，然而它又是检索系统中不可缺少的重要部分，是文献进入检索系统的必经环节，因此其效率问题是必须解决的。

对于现代化的、自动化的情报检索系统来说，手工方式的标引工作显然是不适宜的。长期以来，人们一直致力于自动标引的研究工作，也取得了一些成就，但就目前情况看来，自动的标引工作方式仍未占据标引工作的统治地位，在有些检索系统中，同时存在着自动和手工两种标引手段，在有些检索系统中，手工标引是主要的标引方式。

手工标引与目前水平上的自动标引（以下简称自动标引）相比较，各自特点如下：

1. 手工标引效率较低，自动标引效率较高。

标引是一个智力活动过程。对于人工标引来说，它需要标引者具有相当的普通知识、专业知识、文献知识和标引知识，认真领会文献内容，经过归纳、分析、综合等思维活动，选择出最恰当的若干个标引词，这种工作难度使得即使是较高水平的标引人员（甚至是标引专家）也需要耗费一定的精力和时间才能完成对一篇文献的成功标引；自动标引则不同。当把标引工作需要的方法、规律、要求以及待标引的文献输入计算机后，计算机会自动、快速地按照人们规定的方法、规律和要求对文献进行标引。

2. 手工标引质量较高，自动标引质量较低。

如前所述，标引是一个智力活动过程，因此，由具有较高水平的标引员标引出的文献，其可信程度是比较高的，能在一定程度上满足人们对标引工作的要求；自动标引则不同，目前的自动标引技术中并没有什么“智能”因素，它是通过文献及其内容成份（句子、短语、词组、词片、字等）的外部形式实现的，这些方法固然有一些理论根据和在实际工作中成功的经验，但它毕竟是沿着人工标引不同的道路发展的。目前水平上的标引质量还不太令人满意，而且，我们还应注意，标引的目的在于协调标引者和检索者，目前对提问的“标引”工作（选择检索词的工作）还是人工进行的，与自动的文献标引方式相对照，二者在形式上的结果是否完全一致，恐还值得研究。

3. 手工标引的一致性较差，自动标引的一致性较好。

所谓标引的一致性，是指在不同的时间和环境下，对同一篇文

文献标引的若干个结果的一致程度。手工标引除受标引工作的一般规律制约外，还受许多偶然因素的影响，这使得手工标引工作的“再现性”较差。同一篇文章就由不同的人同时标引，或由一个人在不同的时间和场合标引，结果往往会不一致；自动标引则不同，自动标引是机械的，在输入关于标引工作的一般性方法、规律、要求之后，除非有意识地输入新的指令性信息，可在任何时候再现出同样的标引结果。

4 手工标引多数是赋词标引，自动标引基本上是抽词标引

赋词标引是指标引用词是从词表中抽出赋给该篇文献的，换言之，赋词是规范词。抽词标引是指标引用词是从文献中抽取出来的，只要它被标引者认为恰当表达了文献主题，则不考虑词表对词汇的限制。换言之，抽词是自由词。自动标引工作在抽词标引领域里进展较快。

文献资料自动标引的研究工作目前还远未完成，但它无疑是标引工作必然的发展方向，也是自动化情报检索中的一个重要领域。我们在这里选择介绍一些国内外的自动标引技术方案。

5.2 西文自动标引方案简介

5.2.1 词频统计原理

西文的自动标引方法基本上都是基于词频统计的，为此，我们先介绍词频统计的思想。

人们在长期的文献工作和标引工作中发现，在一篇文献中，某个词在该篇文献中的重要程度往往与该词在该篇文献中的出现频率有关，一个词越重要，它在文献中出现的相对频率就越高，因此，

有可能通过统计词汇在文献中的出现频率，判断它们对于该篇文献的重要程度。从而在其中选取某些频率适当的词作为该篇文献的主题标识。

进一步地，人们还发现了词与词频之间的一些数量上的关系定律。将一篇文献中的各个不同词按其出现频率降序排序，那么各词的频率与该词在这个序列中的序号的乘积近似地等于一个常数。用公式表示为

$$F_i \cdot r_i = C \quad (5-1)$$

其中 F_i 是某词在一篇文献中的出现频率， r_i 是该词在上述序列中的序号， C 是一常数。这个统计规律被称为 Zipf 定律。

例如，有人对出现的 1000000 个词汇进行了统计，其中出现频率最高的十个词的频率及其序号情况如下表所示

序号 (r)	词	词频 (F)	$r \cdot (F/1000000)$
1	the	69971	0.070
2	of	36411	0.073
3	and	28852	0.086
4	to	26149	0.104
5	a	23237	0.116
6	in	21341	0.128
7	that	10595	0.074
8	is	10099	0.081
9	was	9816	0.088
10	he	9543	0.095

表 5-1 Zipf 定律的一个实例

并不是频率越高的词越有主题标识意义和检索意义。因为这些频率最高的词往往是一些功能词（如上表中的词）和专指度不高的其他词（如某一学科中的通用和常用词汇），这些词显然不适合甚至不可能用作标引词。与此相反的一个情况是，频率太低的词也不适合作标引词，因为它们在文献中只是偶而出现或很少出现的，我们可以认为它们不代表文献的基本主题和主要内容。

用来作为标引词的那些词汇往往是中频词。

我们可以设想一个指标：标识能力。这个指标表示词作为标引词的合适程度。从高频词出发，频率最高的词作为标引词的能力是很弱的。随着频率的降低，相应的词逐渐能够表示该篇文献的主题和内容。在中频词的某一点上，标识能力达到最大值。但以后则随着频率的继续下降而下降。图示如下。

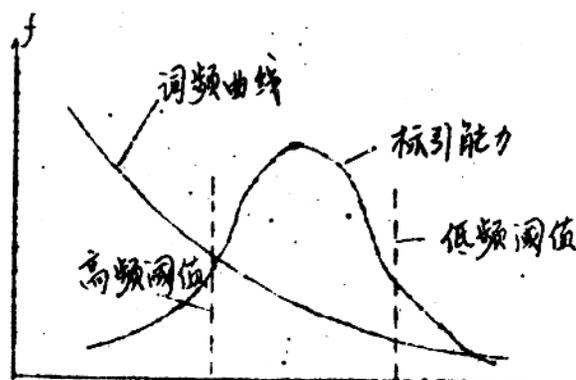


图 5.1 Zipf 定律和标引能力

选择对于一篇文献适合的标引词是如此。选择对于一个检索系统适合的标引词集合同样如此。设该检索系统共有 N 篇文献。选定标引词集合的步骤如下。

1. 对于每一篇文献。计算其中每一个不同词的出现频率。即

主题词 k 在文献 i 中的频率 F_{ik} 。

2. 对于每一个标引词，计算其在全部 N 篇文献中的出现频率 TF_k 。

$$TF_k = \sum_{i=1}^N F_{ik}$$

3. 按照 TF_k 的降序，将所有主题词排序。指定一个适当的高频阈值。弃去全部在这个高频阈值之上的词。这些被弃去的词由于在文献集中出现过频，其出现与否并不过分影响检索性能。

4. 指定一个适当的低频阈值，同样弃去全部在这个低频阈值之下的词。

5. 剩下的中频词被用来标引这 N 篇文献。

下面给出几种标引方法。每一种标引方法的核心是一个赋权函数。按照这个赋权函数为一篇文献中的各“准标引词”赋权之后，其中权值较高的若干个词便作为该篇文献的标引词。

5.2.2 逆文献频率法

这种方法的理论基础是：假定词的重要程度与该词在该篇文献中的出现频率成正比，而与用这个词标引过的文献篇数成反比（因为，在越少的文献中出现这个词，则这个词对这些文献的标识能力就越强。或者说，区别这些文献与其他文献的能力就越强）。

先看后一个指标，记这个文献篇数为 D_k ，即在整个检索系统的 N 篇文献中，有 D_k 篇文献中出现了主题词 k 。我们用“逆文献频率因子”来表示 D_k 对标引词 k 的权值的影响因子。对逆文献频率因子的一个可能的计算方法是

$$\log_2 \frac{N}{D_k} + 1 = \log_2 N - \log_2 D_k + 1 \quad (5-2)$$

其中 N 是整个系统中的文献篇数。

例如，在一个有 1000 篇文献的集合中，Alpha 出现在 100 篇文献中，Beta 出现在 500 篇文献中，Gamma 出现在 900 篇文献中。于是与这三个词相应的三个逆文献频率因子的值分别为 4.322、2.000 和 1.132。我们看到，有某词出现的文献篇数越少，则该词的逆文献频率因子的值越高，从而对于这个词来说，可以更容易地把这些文献检索出来。

令 W_{1k} 为文献 1 中词 k 的权。我们可以规定以下公式

$$W_{1k} = F_{1k} \cdot (\log_2 N - \log_2 D_k + 1) \quad (5-3)$$

5.2.3 信号——噪音法

利用信息论的方法也可以给标引词赋权。我们知道，在信息论中，一个信号（或一个词）所含内容的信息量可以通过该信号（或该词）在一次通信（或一个文本）中出现概率的反比函数来表示。该词的出现概率越高，则该词的信息量越小。一个词的信息量定义为

$$I = -\log_2 p \quad (5-4)$$

其中 p 是该词的出现频率。

例如，词 Alpha 在 10000 个词中只出现了一次，则其出现概率为 0.0001，而其信息量为

$$\begin{aligned} I &= -\log_2(0.0001) \\ &= -(-13.278) \end{aligned}$$

仿照香农的信息量公式，我们可以定义一个标识词对一个有 N 篇文献的文献集合的“噪音”，

$$Noise_k = \sum_{i=1}^N \frac{F_{ik}}{TF_k} \log_2 \frac{TF_k}{F_{ik}} \quad (5-6)$$

这个值与该词在文献集合中的“聚集度”成反比变化。即对于绝对平均的分布，该词在每篇文献中出现的频率相同，这时的“噪音”值为最大。例如，设词 k 在每篇文献中都出现且仅出现一次 ($F_{ik} = 1, i = 1, 2, \dots, N$)，则

$$\begin{aligned} Noise_k &= \sum_{i=1}^N \frac{1}{N} \log_2 \frac{N}{1} \\ &= \log_2 N \end{aligned}$$

再看另一个极端的情况，即绝对聚集的情况，词 k 只出现在一篇文献 i 中（这时显然有 $F_{ik} = TF_k$ ），则噪音值

$$\begin{aligned} Noise_k &= \frac{TF_k}{TF_k} \log_2 \frac{TF_k}{TF_k} \\ &= 1 \log_2 1 \\ &= 0 \end{aligned}$$

显然，在噪音和标引词的专指度之间存在着联系，因为宽泛的、不专指的词总是趋向于在文献集合中较为分散，因而也有较高的噪音。所以，噪音的反比函数可用来作为标引词权值的因子。一种称为标引词 k 的信号函数可定义如下：

$$Signal_k = \log_2 TF_k - Noise_k \quad (5-7)$$

对于噪音的最大值 ($F_{ik} = 1, i = 1, 2, \dots, N$)，

“信号”值为0，因为那时有

$$\log_2 TF_k = \text{Noise}_k = \log_2 N$$

而另一方面，当标引词k仅出现在一篇文献中时，噪音值 $\text{Noise}_k = 0$ ，而信号值为最大

$$\begin{aligned} \text{Signal}_k &= \log_2 TF_k - \text{Noise}_k \\ &= \log_2 TF_k \\ &= \log_2 F_{1k} \end{aligned}$$

其中 F_{1k} 是标引词k在这唯一的一篇文献中的出现频率。

在原理上，可以将一篇文献中的各标引词按其信号值的降序排列。这种函数有利于将少量文献（信号值高的标引词出现的文献）从其他文献中区别出来。于是，在考虑 F_{1k} 的值后，我们可以定义一个复合函数作为标引词k在文献1中的权值

$$W_{1k} = F_{1k} \cdot \text{Signal}_k \quad (5-8)$$

5.2.4 词辨别值法

标引工作有着双重目的，一个是标识文献使之在检索时能够实际被检索出来，另一个是区别各篇文献，即要求赋予文献的标引词具有辨别文献的功能，这后一个任务是非常重要的，因为它直接决定检索的查准率，上两节我们介绍的算法中已经考虑了这个因素，本节我们再介绍一种主要考虑这个因素的标引方法——词辨别值法。

考虑一个文献集合，令 D_i 和 D_j 为不同的两篇文献，每篇文献各自用一个标引词集合标引，我们可以用一个指标“匹配度”（ $\text{Similar}(D_i, D_j)$ ）来度量文献 D_i 和 D_j 的相似程度。在

典型的匹配度函数中，往往把两篇没有什么共同点的文献的匹配度规定为 0，把两篇完全一致的文献的匹配度规定为 1，而大多数文献对的匹配度值则在 0 与 1 之间。

能满足上述要求的匹配度计算公式有很多，例如余弦公式。设两文献——语词矩阵中 D_i 和 D_j 的两行分别是 $(a_{i1}, a_{i2}, \dots, a_{iM})$ 和 $(a_{j1}, a_{j2}, \dots, a_{jM})$ ，则令

$$\text{Similar}(D_i, D_j) = \frac{\sum_{k=1}^M (a_{ik} \cdot a_{jk})}{\sqrt{\sum_{k=1}^M a_{ik}^2 \cdot \sum_{k=1}^M a_{jk}^2}}$$

设除 $i = j$ 之外的所有文献对 D_i 和 D_j 的匹配度都已计算出，我们可以得到一个平均匹配度 $\bar{\text{Similar}}$ ，其值为

$$\bar{\text{Similar}} = C \cdot \sum_{i=1}^N \sum_{j=1}^N \text{Similar}(D_i, D_j) \quad (5-9)$$

其中 C 是一常数，其值可以任意指定（如指定为

$$\frac{1}{N(N-1)}$$

），这个表达式给出了

文献在文献空间中“聚集”的程度。当所有的文献都相同时，所有的 $\text{Similar}(D_i, D_j) = 1$ ，这时的平均匹配度达到最大值

$$\begin{aligned} \bar{\text{Similar}} &= C \cdot \sum_{i=1}^N \sum_{j=1}^N 1 \\ &= C \cdot N(N-1) \end{aligned}$$

空间密度可以以更方便的方法计算。这种方法是：构造一个人为的“平均文献” \bar{D} ，作为一个中心，在这个中心中，标引词被假定具有平均的频率，即，标引词 k 的平均频率被定义为

$$\bar{F}_k = \frac{1}{N} \sum_{i=1}^N F_{ik} \quad (5-10)$$

于是平均匹配度可以通过以下方式计算，

$$\bar{\text{Similar}} = C \cdot \sum_{i=1}^N \text{Similar}(\bar{D}, D_i) \quad (5-11)$$

考虑一个文献集合，其平均匹配度为 $\bar{\text{Similar}}$ ，现在把这个集合中的所有标引词 k 移除，并记移除掉标引词 k 后这个文献集合的平均匹配度为 $(\bar{\text{Similar}})_k$ 。首先设 k 是一个广泛的高频词，并且有均匀的分布，则把它除去后，会降低平均的文献匹配度，即

$$(\bar{\text{Similar}})_k < \bar{\text{Similar}}$$

其次，设 k 是一个专指度很高的低频词，在原文献集合中，它只在少量文献中有较高的权值，而在其他文献中却权值很小或为零，那么把它移除后，会增加文献集合的平均匹配度，即

$$(\bar{\text{Similar}})_k > \bar{\text{Similar}}$$

于是，对每个标引词 k ，可定义词辨别值 DV_k

$$DV_k = (\bar{\text{Similar}})_k - \bar{\text{Similar}} \quad (5-12)$$

在计算了每个标引词 k 的 DV_k 值之后，可将这些词按 DV_k 值的降序排列，排序后我们将会发现，序号越小的词越专指，序号越