



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

### 第三章 文献情报检索的数据结构 和检索技术

文献情报检索系统的主要处理对象是大量的数据，能否合理地组织数据，能否在计算机上更有效地存贮和处理数据，直接影响系统运行的效率。本章首先介绍几种情报检索中常用的数据结构，然后讨论计算机软件技术中与情报检索关系最为密切的一个方面——查找技术。

#### 3.1 情报检索中的数据结构

##### 3.1.1 逻辑结构与物理结构

###### 1. 逻辑结构

我们已经涉及到了检索文档和文档中记录的关系，在这里，我们可以把记录看成为“基本元素”，把文档看成是由这些基本元素构成的整体。在这个文档中，记录与记录之间一般存在某种关系，这种关系是逻辑上的，由记录内容所决定的，于是，这些记录和它们之间的逻辑关系构成一个逻辑结构。

我们在其他方面也遇见过一个整体中的基本元素以及这些基本元素之间的关系，例如，把主题词和文献看成是一个集合中的元素，则在主题词和文献之间存在着逻辑关系——标引关系，在主题词与主题词之间存在着用、代、属、分、参等关系，在一个检索提问式中，各检索词之间也存在关系——逻辑组配关系，等等。

对数据间关系的描述称为数据的逻辑结构，形式地可以用一个

## 二元组

$$B = (K, R)$$

来表示, 其中  $K$  是基本元素的有穷集合,  $R$  是  $K$  上的关系的有穷集合。

在对数据结构的讨论中, 通常把基本元素称为“结点”, 注意结点的概念是相对的, 例如在文档和记录的关系中, 记录是结点, 而在记录和其属性的关系中, 属性是结点, 记录是由属性组成的整体。

设  $B = (K, R)$  是一个逻辑结构,  $r \in R, k, k' \in K$ , 如果  $\langle k, k' \rangle \in r$ , 则称  $k'$  是  $k$  的后继,  $k$  是  $k'$  的前驱, 称  $k$  和  $k'$  是相邻的结点 (都是相对  $r$  而言)。如果不存在一个  $k'$  使  $\langle k, k' \rangle \in r$ , 则称  $k$  为关于  $r$  的终端结点。如果不存在一个  $k'$  使  $\langle k', k \rangle \in r$ , 则称  $k$  为关于  $r$  的开始结点。如果  $k$  不是终端结点, 也不是开始结点, 则称  $k$  为内部结点。

我们只讨论包含一个关系的  $R$ , 即  $R = \{r\}$ 。

数据的逻辑结构常常简称为数据结构。数据结构分为线性结构和非线性结构。在线性结构里只有一个终端结点和一个开始结点, 并且所有的结点都最多只有一个前驱和一个后继。线性表是典型的线性结构。非线性结构中最重要的是树结构, 树结构中只有一个称为根的结点没有前驱, 其他结点有且仅有一个前驱。最一般的情况是对结点的前驱和后继都不作限制。这种结构称作图。对于线性表、树和图, 我们将分别在 3.1.2、3.1.3、3.1.4 中简单介绍。

### 2 物理结构

数据的逻辑结构是从逻辑的观点来观察数据, 它与数据的存贮无关, 是独立于计算机的。数据的物理结构是逻辑结构在计算机存



贮器里的实现。它又称为存贮结构。数据的物理结构是依赖于具体的计算机的。

计算机的存贮器是由有限多个存贮单元组成的，每个存贮单元有唯一的整数地址。各存贮单元的地址是连续编码的，每个存贮单元  $z$  都有唯一的后继单元  $\text{suc}(z)$ 。如果  $z'$  是  $z$  的后继单元，则称  $z$  和  $z'$  是相邻的。一片相邻的存贮单元的整体叫作一个存贮区域，记作  $M$ 。

设有逻辑结构  $B = (K, R)$ 。要把  $B$  存贮在计算机中，首先必须建立一个从  $K$  的结点到  $M$  的单元的映象  $S: K \rightarrow M$ ，即对于每一个  $k \in K$ ，都有唯一的  $z \in M$ ，使  $S(k) = z$ ，同时这个映象还必须明显地或隐含地体现关系  $R$ 。

我们用  $\omega_k$  表示结点  $k$  的值，用  $\text{LOC}(k)$  表示结点  $k$  对应的存贮单元的地址，用  $\text{CON}(z)$  表示地址为  $z$  的存贮单元里存贮的内容，于是有

$$\omega_k = \text{CON}(\text{LOC}(k))$$

有四种基本的存贮映象方法，即有四种基本的物理结构，它们是：

### ① 顺序的方法

这种方法主要用于线性的数据结构，它把相邻的结点存贮在相邻的存贮单元里，结点之间的关系由存贮器的单元的邻接关系来体现。如果结点  $k$  所占存贮的第一个单元为  $z$ ，而最后一个单元为  $z_1$ ，那么  $z_1$  的后继单元  $z' = \text{suc}(z_1)$  就是  $k$  的后继  $k'$  的第一个存贮单元。

### ② 链接的方法

这种方法是给结点附加上指针字段，即是将对应于结点的存贮

单元分为两部分，一部分存放结点本身的信息，称数据项，另一部分存放此结点和后继结点所对应的存储单元的地址，称指针项。指针项可以包括一个或多个指针，以指向结点的—个或多个后继，或记录其他信息。

### ③ 索引的方法

在线性的结构里，结点可以排成一个序列， $k_1, k_2, \dots, k_n$ ，每个结点 $k_i$ 在序列里都有对应的位置数 $i$ ，这个位置数就是结点的索引。索引的存储结构就是用结点的索引号 $i$ （或关键字 $A_i$ ）和结点的存储地址来建立附加的索引表，使得索引表里第 $i$ 项的值就是第 $i$ 个结点（即结点 $k_i$ ）的存储地址。

### ④ 散列的方法（Hash法）

主要思想是根据结点的值来确定它的存储地址。在结点 $k$ 的各个字段里取一个或几个字段的值 $\omega_i k$ 作为关键码，结点 $k_j$ 对应的存储地址由函数 $H$ 确定

$$LOC(k_j) = H(\omega_i k)$$

这种方法又称KAT法（Key-to-Address Translation）用散列法进行存储映象，K的结点 $k$ 随机地，分散地存储在M的存储单元里。

在讨论数据的物理结构即存储结构时，有两点需要注意。

首先，数据的逻辑结构与物理结构是无关的，一种逻辑结构可以由多种存储结构实现。例如，一棵逻辑上的树可以由链接的方法或散列的方法实现，在经过适当的处理后，它也可以由顺序的方法或索引的方法实现。一种存储结构也可以实现多种逻辑结构，例如，散列的方法可分别实现线性表、树和图。

其次，数据的存储结构与数据的查找紧密相关。在设计存储结构

构时，我们给出了一个映象  $S$ ：  $S(k) = z$ ，于是把结点  $k$  存贮到  $z$  中，在查找数据时，我们同样利用这个映象。设我们要查找结点  $k$ ，  $S(k) = z$ ，于是到单元  $z$  中取出结点  $k$ 。

查找技术在 3.2 中讨论。

### 3.1.2 线性表

线性表的逻辑结构是  $B = (K, R)$ ，其中  $K$  是包含  $n$  个结点的集合， $R$  仅包含一个关系  $R$

$$R = \{ \langle k_{i-1}, k_i \rangle \mid 2 \leq i \leq n \}$$

即是说， $K$  里所有的结点都可以按关系  $R$  排成一个序列  $k_1, k_2, \dots, k_n$ 。在这个序列里，每个结点都有相对位置  $i$ ， $i$  叫作  $k_i$  的索引。表中的结点又称作表目，或元素。

线性表的每个表目可以包括一个或多个字段。能够唯一地标识表目的一个字段或一组字段叫作关键码。例如，按文献号顺序排列的主文档是一个表，每一篇文献记录是一个表目，每一个表目中有文献号、篇名、作者、出版年等字段，其中“文献号”字段可以唯一地标识一篇文献。文献号就可以作为关键码。非关键码字段的信息又称为属性。

用不同的方法（顺序、链接、索引、散列）存贮一个线性表，可以有不同的用途，下面介绍情报检索中常用的一些。

#### 1. 顺序表和目录表

用顺序的方法来存贮一个线性表，即把表目一个接一个地放在顺序的一片单元里，这样的表称为顺序表。

例如，情报检索中的主题词，如果考虑其字顺关系，则它们是一个线性结构。在目前的大多数检索系统中，也是按字顺把词典存

贮在存储器上的，因此，词典是一个顺序表。

再例如，在SDI中，新到文献也是排成一个顺序表的。

顺序表的逻辑结构和物理结构都很简单，并且存贮密度大，节省空间。由于结点与结点之间在物理上的关系反映了原有的逻辑关系，所以可以在顺序表上再加索引后成为索引顺序文件，这是目前词典常采取的一种形式。

但是，顺序表的结点增删比较麻烦。增加一个结点到与其相应的位置后，要向后移动其后面的所有结点，删去一个结点后，为了填补空闲区域，要往前移动其后面的所有结点。由于这个原因，大型词典的词汇更新和删除工作往往是成批处理的。

在一个表中，表目的大小往往是不一致的，如倒排档中的主题词记录，主文档中的文献记录等，可变长的记录处理起来很不方便。为此，可以把表分成两部分，一部分叫目录表，它的表目包括两个字段，一个字段是关键码，另一个字段是地址。另一部分叫属性数据，它包括表里关键码以外的其他信息。属性数据与目录表分别存在不同的存贮区域里，目录表里结点的地址字段指向属性数据存贮区域里对应于这个关键码的信息。目录表的结点是固定长的。

例如，主文档中的文献记录是可变长的，检索不便，我们便为其建一个目录表。这个目录表就是在第二章中提到的主文档索引文档。

## 2 栈

栈是情报检索中常用的一种线性表。我们在2.5.1中曾给出了它的非形式化定义，在这个定义中包含了它的存贮结构特点。

栈可以用顺序表的方式实现。分配一块连续的存贮区域存放栈里的表目，并用一个变量T来指明当前的栈顶（图3.1）

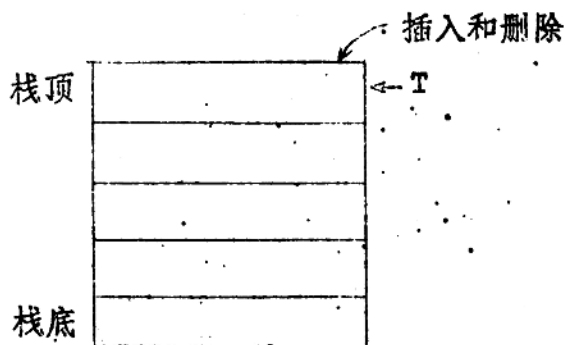


图 3.1 栈

栈的基本运算有：

- ① Push，插入（或叫推入）一个元素到栈顶。
- ② Pop，从栈顶删除一个元素，也就是托出栈顶元素。
- ③ Top，看一下当前的栈顶元素。

我们在 2.5.1 中曾给出一种检索逻辑式的运算方法，首先将检索式转换成逆波兰表示法，然后对逆波兰检索式进行运算。运算规则是，每个运算符的两个运算项目是逆波兰形式中紧挨在这个运算符之前的两个运算项目。我们下面给出一个算法，这个算法把上述的两个步骤并作一步完成，即，它不形成独立的逆波兰形式，而是在转换过程中就正确地选择运算项目及它们的运算符并对其进行运算。

我们用两个栈。一个是我们曾用过的运算符栈，在这里称它为 sop 栈，逆波兰表示区在这里也用一个栈，我们称它为 svar 栈，是表示运算项目无条件进栈的意思。三种逻辑运算符的优先数仍如 2.5.1 中所规定。在下面的算法中，C 是逻辑式中的字符，可能是是运算项目，也可能是运算符，OP 是 sop 栈中的元素，即运算符；

$V_1, V_2, R$  都是  $svar$  栈中的元素, 即运算项目, 其中  $R$  是运算的中间结果。

### 算法 3.1 检索逻辑式的运算

- ① 从左向右扫描逻辑表达式。若表达式已取完则转步骤⑦, 否则取来一个字符  $C$ 。
  - ② 若  $C$  是运算项目, 则  $Push(svar, C)$ , 转去执行步骤①。
  - ③ 若  $C$  是运算符且运算符栈为空, 则  $Push(sop, C)$ , 转步骤 1。
  - ④ 若  $C$  是运算符且运算符栈非空, 则  $Top(sop, op)$ , 比较栈顶运算符  $op$  和新取来的运算符  $C$  的优先数, 若  $C$  的优先数不小, 则  $Push(sop, C)$ , 转步骤 1。
  - ⑤  $Pop(sop, op), Pop(sop),$   
 $Pop(svar, V_1), Pop(svar),$   
 $Pop(svar, V_2), Pop(svar),$   
执行运算  $V_1 op V_2$ , 运算结果为  $R$   
 $Push(svar, R)$
  - ⑥ 转步骤③
  - ⑦ 反复执行步骤⑤, 直到  $sop$  为空则结束。  $svar$  的栈顶元素值就是逻辑式的计算结果。
- 除了顺序表、目录表和栈之外, 线性表还包括链表、队列等形式, 介绍从略。

### 3.1.3 树

树形结构是结点之间具有分支关系的结构，它很类似于自然界中的树。在非线形结构中树形结构是最重要，应用最广泛的一种，尤其是在大量的数据处理中更为突出。

树的逻辑结构定义为  $B = (K, R)$ ，其中  $K$  是包含  $n$  个结点的有穷集合。 $R$  是  $K$  上关系的集合，它只包含一个元素  $N$ ，并满足

a. 有且仅有一个结点  $\omega \in K$ ， $\omega$  被称为树的根，它相对关系  $N$  没有前驱。

b. 除根结点  $\omega$  外， $K$  中每个结点相对关系  $N$  都有且仅有一个前驱。

c. 对于  $K$  里除  $\omega$  外的任何结点  $k$ ，都存在一个结点序列  $k_0, k_1, \dots, k_s$ ，使得

$$k_0 = \omega, k_s = k, \langle k_{i-1}, k_i \rangle \in N, 1 \leq i \leq s$$

这样的结点序列称为从根到结点  $k$  的一条路径。

树结构在客观世界和情报检索领域中大量存在。

例如，分类法是一种树形结构

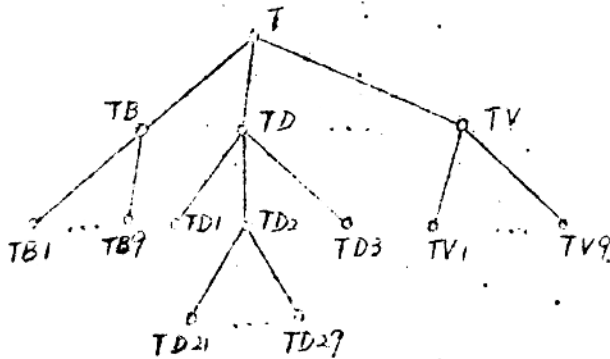


图 3.2 分类树

再例如，主题词若按其字母的等级顺序，则也是一棵树。

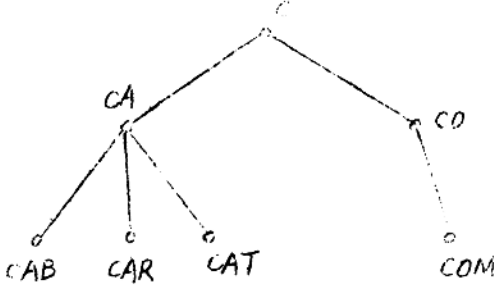


图 3.3 主题词树

再例如，逻辑表达式其实是树形结构而不是线形结构。表达式  $(A + B * (C + D)) + E * F$  可表示为

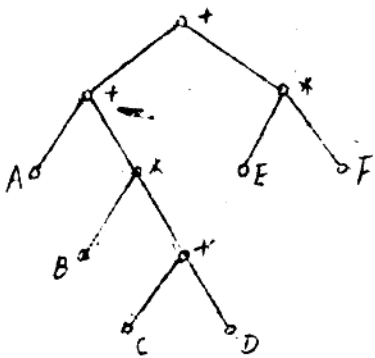


图 3.4 表达式的树形结构

在一棵树中，终端结点即没有后继的结点称为树叶。在实际应用中，树叶相对于其他结点往往有特殊意义。如在图 3.2 中，树叶是最低级的类目；在图 3.3 中，树叶是实际的主题词；在图 3.4 中，树叶是运算项目。

在图 3.4 中，我们给出了一个每个内部结点只有两个分支的树形结构，这称之为二叉树。二叉树与树不同，它的严格定义是

定义：二叉树是结点的有限集合，这个有限集合或者为空，或



者由一个根及两个不相交的分别称作这个根的左子树和右子树的二叉树组成。

这是个递归的定义。二叉树可以是空集合，因此根可以有空的左子树或右子树，或左右子树皆为空。图 3.5 是二叉树的五种基本形态。



图 3.5 二叉树的五种基本形态

对于二叉树，有两点需要特别注意：①二叉树的结点的子树个数只能是 0、1、2；②结点的子树要区分为左子树和右子树。

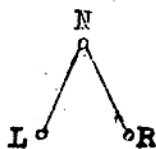
树和二叉树都可以用多种存贮方法实现。

#### 1. 顺序的方法

树和二叉树本身都是非线性的数据结构，要用顺序的方法实现存贮，首先要将它们线性化。线性化的主要方法是对树或二叉树进行周游。

#### 二叉树的周游方法

考虑到二叉树的基本组成部分是



所以，周游二叉树有三种主要的方式

1. 前序法 (NLR 次序)，其递归定义是