

* * * * *
* * * * *
* *
* *
* *
* *
* * * * *

目 录

第一章 计算机情报检索概述	1—1
1.1 情报与社会的发展	1—1
1.2 情报检索与文献检索	1—2
1.3 文献情报检索系统的基本功能	1—3
1.4 文献情报检索系统的基本原理	1—4
1.4.1 文献与文献标识	1—6
1.4.2 文献——语词矩阵	1—6
1.4.3 三种基本的文献检索方式	1—7
1.5 联机情报检索	1—10
1.6 关于本课程的说明	1—11
第二章 基于倒排档的检索系统	2—3
2.1 倒排档检索技术发展简史	2—1
2.2 布尔逻辑	2—5
2.3 典型的文档结构	2—7
2.4 检索过程	2—11
2.5 检索式的逻辑运算	2—12
2.5.1 运算顺序的正确控制	2—13
2.5.2 集合的逻辑运算	2—17

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

第六章 聚类检索

6.1 问题的提出

在前面的章节中，我们讨论了基于倒排档的检索机制。倒排档检索方法有很多优点，但同时存在着以下一些缺点：

- ① 需要一个辅助检索用的倒排档；
- ② 基本上是一种完全匹配方式，不能表示命中文献对用户适用程度上的不同；
- ③ 不注意标引词的权值，对于同一篇文章的若干标引词往往是同等看待的，不注意它们重要程度的不同，也不注意词间关系；
- ④ 不能由一篇相关文献引导至另一篇相关文献；
- ⑤ 因为没有在机内存贮词间关系库，所以不易实现改善检索效果的自动化。

针对以上问题，国外近年来发展了一些聚类检索系统，它们在理论和技术上尚未十分成熟，因此又被称为实验性检索系统。这个名称也是相对传统的倒排档检索系统而言的。

我们将通过一些实际的检索系统来介绍这种检索机制。其中最主要的是 SMART 系统。

6.2 SMART 系统

SMART 系统与传统的基于倒排档的系统有以下几点主要的不同之处：

- ① 使用全自动的标引方式和检索提问构成方式；
- ② 集中与同一主题有关的文献，使之可以从一篇文献方便地寻至另一篇；

② 通过处理文献与文献的匹配度、文献与提问的匹配度来完成检索，并且将输出文献按匹配度的降序排列；

(1) 通过先检索文献，自动改善检索提问式。

6.2.1 文献的向量表示和匹配度计算

我们有文献——词矩阵

$$D = \begin{matrix} & T_1 & T_2 & \dots & T_M \\ D_1 & a_{11} & a_{12} & \dots & a_{1M} \\ D_2 & a_{21} & a_{22} & \dots & a_{2M} \\ \vdots & \dots & \dots & \dots & \dots \\ D_N & a_{N1} & a_{N2} & \dots & a_{NM} \end{matrix} \quad (6.1)$$

我们把其中的第1行当作一个M维向量，以代表文献 D_1 ，即有

$$D_1 = (a_{11}, a_{12}, \dots, a_{1n}) \quad (6.2)$$

其中 a_{1j} 称为向量 D_1 的分量，表示索引词 T_j 对文献 D_1 的重要程度，或权

文献与文献之间可以计算其匹配度。匹配度的计算方法有很多，

例如

①

$$\text{Similar}(D_i, D_j) = \sum_{t=1}^M a_{it} a_{jt} \quad (6.3)$$

②

$$\text{Similar}(D_i, D_j) = \frac{\sum_{t=1}^M a_{it} a_{jt}}{\sum_{t=1}^M a_{it}^2 \cdot \sum_{t=1}^M a_{jt}^2} \quad (6.4)$$

③

$$\text{Similar}(D_i, D_j) = \frac{\sum_{t=1}^M a_{it} a_{jt}}{\sum_{t=1}^M a_{it}^2 + \sum_{t=1}^M a_{jt}^2 - \sum_{t=1}^M a_{it} a_{jt}} \quad (6.5)$$

其中式(6.4)称为余弦函数，是常用的一种计算匹配度值的公式，它的值域在0和1之间。

在一个有M个标引词的检索系统中，可以认为存在一个M维的文献空间，每篇文献可用这个文献空间中的一个向量来表示，即 $D_i = (a_{i1}, a_{i2}, \dots, a_{iM})$ 其中各分量是该向量在各相应标轴上的坐标。

例如，下图是一个三维的文献空间，在这个文献空间中，三个文献向量

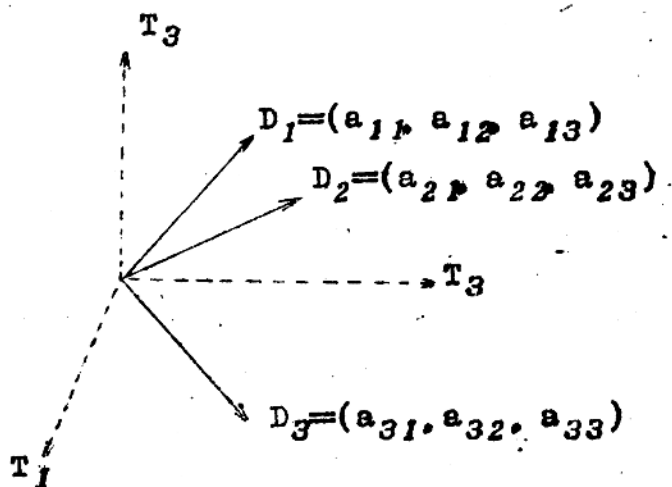


图 6.1 文献空间和文献向量

当文献被表示为文献空间中的向量之后，文献之间的相似程度（匹配度）还可以用其向量之间夹角的反比函数来度量。如果两篇文献完全相同，则其文献向量在文献空间中重合，夹角为0，而匹配度值为最大。

在SMART系统中，提问Q也被表成向量形式

$$Q = (q_1, q_2, \dots, q_M) \quad (6.5)$$

其中分量 q_j 表示检索词 T_j 对该提问的重要程度，或权。

需要指出的是，一般的布尔逻辑提问式可以转换成向量形式的提问式。例如

$$\begin{aligned} Q &= (A + B) * (C * D) \\ &= (A + B) * C * D \\ &= (A * C * D) + (B * C * D) \end{aligned}$$

于是这个提问式Q可用两个提问向量表示：

$$\begin{aligned} Q_1 &= (A, C, D) \\ Q_2 &= (B, C, D) \end{aligned}$$

如果系统中只有A、B、C、D四个标引词，那么提问式 Q_1 和 Q_2 也可表成

$$\begin{aligned} Q_1 &= (1, 0, 1, 1) \\ Q_2 &= (0, 1, 1, 1) \end{aligned}$$

提问向量与文献向量之间也可以计算匹配度，公式仍如式(6.3)、(6.4)、(6.5)所示。如需计算 D_i 与Q的匹配度，只要把公式中的 a_{jt} 换成 q_t 即可。常用的匹配度计算公式是余弦函数的式(6.4)。

用定量的方法计算文献与提问的匹配度之后，我们不再要求命中文献恰好包含全部检索词，我们这里只要求命中文献达到一定的与提问向量的匹配度，或者要求输出的文献满足一定的数量要求。例如，给定匹配度阈值为0.5，则所有与提问向量的匹配度值达到或超过0.5的文献都是命中的；或者，我们要求输出n篇命中文献，那么将系统中的文献分别与提问计算匹配度并将它们按匹配度值降序排列后，前面的n篇便是所需的。

6.2.2 聚类文件的生成和SMART系统的文档结构

计算了所有文献对之间的匹配度之后，我们可以构造一个聚类文件。这个聚类过程是一个从下到上的生成过程。

1. 给定一个阈值，由这个阈值生成若干初始类目。在每一个类目中，任意两篇文献的匹配度都不小于这个阈值。设生成的初始类目为 C_1, C_2, \dots, C_z 。

如下图，每一个×表示一篇文献。它们的位置关系表示这些文献向量在文献空间中的分布。两个×相距越近，表示这两篇文献的相似程度越高（匹配度值高）。每个圆圈表示一个初始类目。在一个类目（圆圈）中，任意两个×之间的距离都不大于该圆的直径。显然，这里采用的是类似于“词团”结构的“文献团”结构。在一个文献团结构的类中，任一篇文献与该类中其他任一篇文献的匹配度都不小于指定阈值。

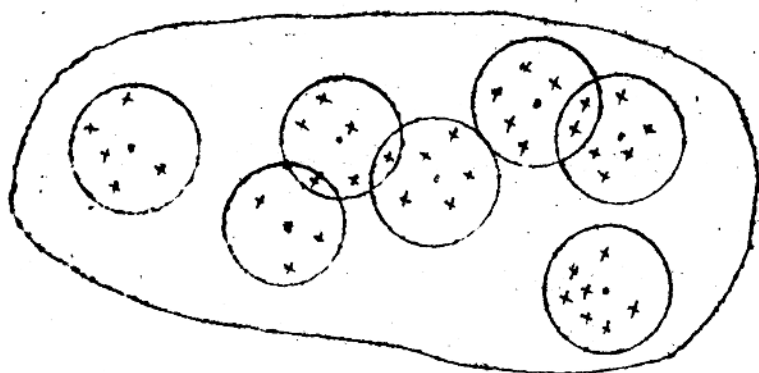


图 6.2 文献的聚类

可以想象，这若干个类目很可能不是互相独立的，即有的文献会同时出现在不止一个类目中。

在上图中，我们为了直观起见，仅给出了二维的文献空间。在一个 M 维的文献空间中，一类是一个 M 维的单位球域。其中任意两点之间的距离不超过二单位。

2. 对于每一个初始类目 C_j , 给出一个类目向量 $C_j = (T_{1j}^C, T_{2j}^C, \dots, T_{Mj}^C)$, 该类目向量各分量的值是各文献向量的相应分量 T_{ik} 的平均值, 即

$$T_{jk}^C = \frac{1}{m} \sum_{i=1}^m T_{ik} \quad (6.7)$$

其中 m 是该类目 C_j 中的文献篇数。

3. 计算这些初始类目之间的两两匹配度, 以作为“类对”的匹配度。

4. 给定一个阈值, 由这个阈值再生成若干高一级的类目, 每一个类目中包含若干上一步生成的初始类目, 而这些初始类目之间的两两匹配度都不小于指定的阈值。

5. 重复上述步骤 2、3、4, 直至归并成为一个最高级类目, 如下图所示:

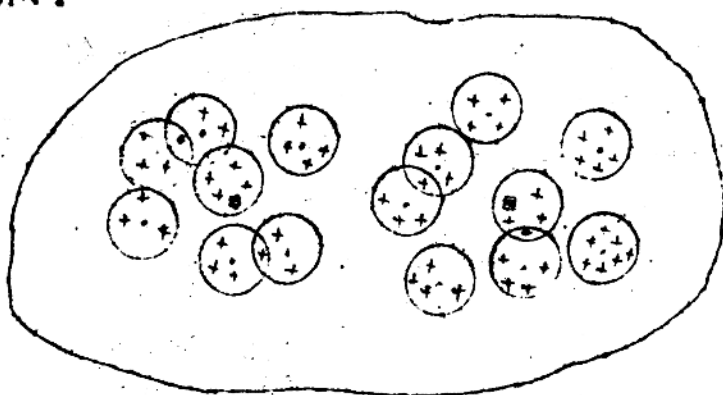


图 6.3 文献聚类的层次

我们将这样生成的各级类目按从上到下的顺序分别称为一、二、三、...级类目, 其中最后一级已是实际的文献。但由于文献和类目都被表示为向量, 所以在处理手段上并无区别。

我们在这里采用的类目生成方法需要计算每一对文献之间的匹配度, 计算量较大。另一种方法是先任意给出一些类目, 然后计算

每一篇文献与每一个类目的匹配度，调整文献对类目的归属，逐渐逼近最佳的类目体系。这种方法的具体介绍从略。

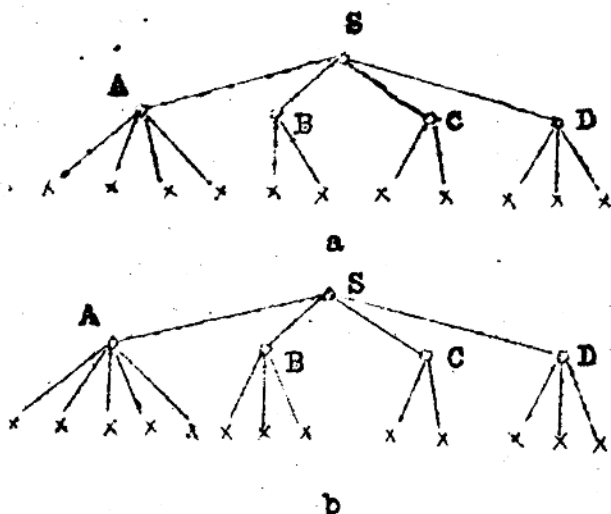
类目的生成并不是一次性的。当新的文献加入这个分类体系后，会给原类目体系带来以下两个影响：

1. 由于初始类目的向量是由该类目中各文献向量所决定的，所以加入一个新的文献向量后，会改变初始类目的向量表示；又由于上一级类目的向量表示是由下一级类目的向量表示所决定的，所以后者的变化会影响前者。这种影响将一直波及到根。波及范围是分类树中从一个树叶（新加入的文献）到根的一条路径。

2. 类目中增加了新的文献或下位类后，可能会使该类目变得过大，因而影响检索。

对于前一个问题，解决办法是自加入文献的这个初始类目起，向上逐级重新计算其类目向量，直至树根。

对于后一个问题，解决办法是进行一些必要的类分裂工作，如下图所示：



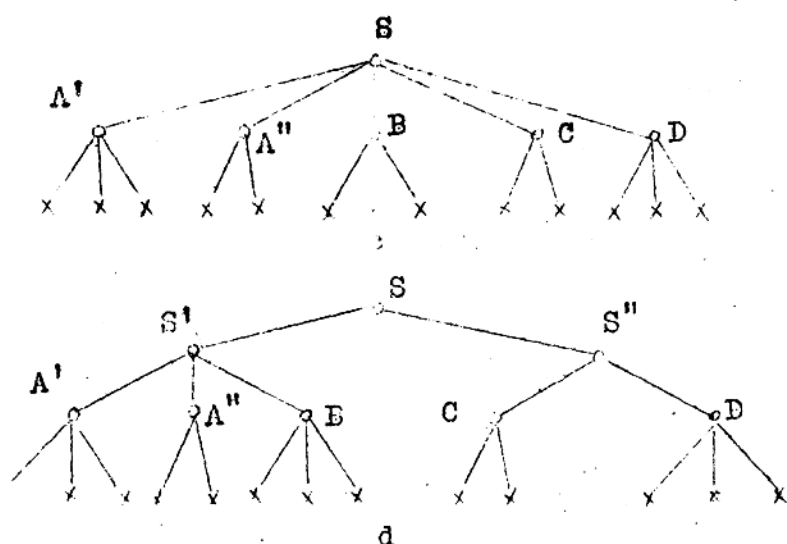


图 6.4 类的分裂

在 SMART 系统中，由于文献和提问都被表示成向量形式，所以增加新文献记录的工作机制与检索工作机制是基本相同的。下面给出它们的工作流程：

① 输入一个向量。输入向量可以是准备加入分类体系的文献向量，也可以是准备检索的提问向量。将所有二级类目排成一个待考察的类目队列。

② 在类目队列中选择一个结点，与输入向量计算匹配度。

③ 如果匹配度值小于指定阈值，则转②；否则转④。

④ 当前比较的结点是初始类目吗？如果是，则继续执行；否则转②。

⑤ 输入向量是文献向量还是提问向量？如是前者，则继续执行；否则转①。

⑥ 将输入的文献向量加入当前的初始类目。如加入后这一初始类目太大，则继续执行；否则转①。

- ⑦ 在准备分裂的结点路径中插入一个结点。
- ⑧ 新的文献记录已经加入所有应该加入的类了吗? 如是, 则转①, 否则转②。
- ⑨ 将当前类目的所有儿子加进待比较的类的队列, 转②。
- ⑩ 将当前初始类目的所有文献加进待比较的结点队列, 检索出那些与提问有足够大匹配度值的文献。若已检索出足够的文献记录, 则停止, 否则转②。
- ⑪ 出口, 转类分裂程序。

本节最后给出 SMART 系统的文档结构:

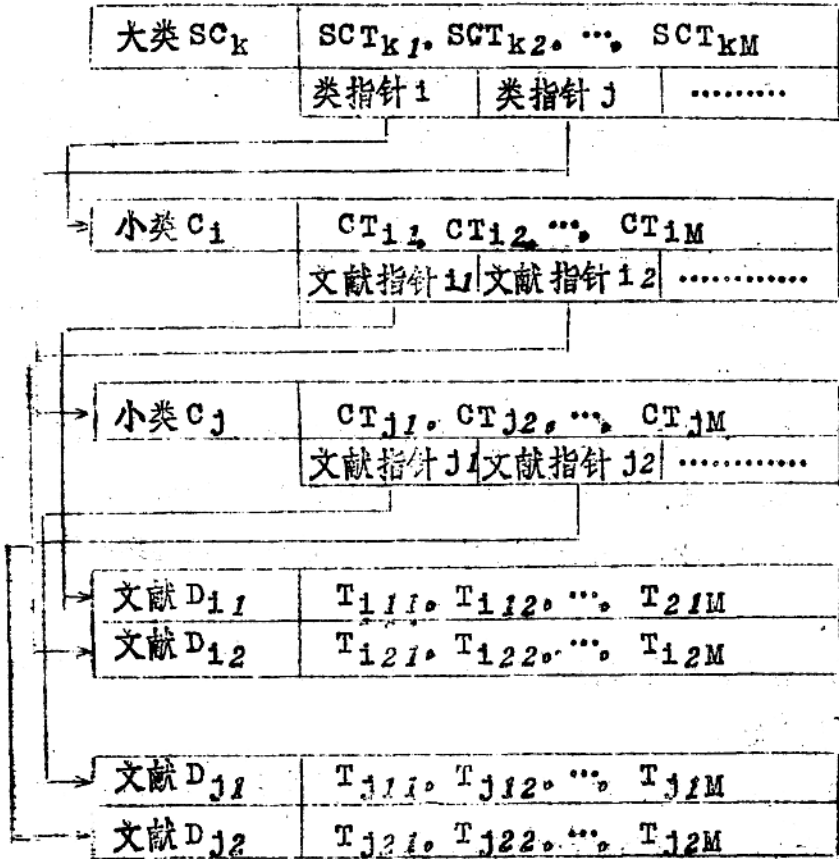


图 6.5 SMART 系统的层次聚类文档结构

6.2.3 提问式的反馈调整

提问式的反馈调整是 SMART 系统的一个重要特点。因为 SMART 系统的检索机制并不是“完全匹配”的，而是通过匹配度值来表示输出文献与提问的满足程度。因此，我们可以并且应该尽量提高这种满足程度。

为设计这种反馈调节机制，我们先看一种理想的然而也是不可能的情况：相关文献已知情况下的最优询问表达式。我们记这个最优询问式为 Q_{opt} 。

设文献数据库中共有 N 篇文献。对于某一用户提问，系统中有相关文献 R 篇，不相关文献 $N-R$ 篇。记这个相关文献群为 D_R ，不相关文献群为 D_{N-R} 。考虑一个检索词（标引词） T_k ，在相关文献群 D_R 中，标引词 T_k 的平均权值为

$$\frac{1}{R} \sum_{i \in D_R} T_{ik} \quad (6.8)$$

其中 $i \in D_R$ 表示这里的文献 D_i 属于 D_R 。

类似地，在不相关文献群 D_{N-R} 中，标引词 T_k 的平均权值为

$$\frac{1}{N-R} \sum_{i \in D_{N-R}} T_{ik} \quad (6.9)$$

可以理解，在最优询问式 Q_{opt} 中，检索词 T_k 的权值应规定为

$$(Q_{opt})_k = c \left(\frac{1}{R} \sum_{i \in D_R} T_{ik} - \frac{1}{N-R} \sum_{i \in D_{N-R}} T_{ik} \right) \quad (6.10)$$

其中 c 是一个常数。

实际上，无论是相关文献群 D_R 还是不相关文献群 D_{N-R} ，在检索前都是不知道的，因为否则也就不需要检索了。并且，一般也

不可能用几何的方法在文献空间中划分出两个互相独立的区域。其中一个恰好包含全部相关文献，另一个恰好包含全部不相关文献。聚类文件虽然在一定程度上聚集了学科相近，主题相关的文献，但也并非绝对如此。与同一检索课题相关的文献分布与不相关的文献分布仍然会有各种复杂的情况，如下图所示：

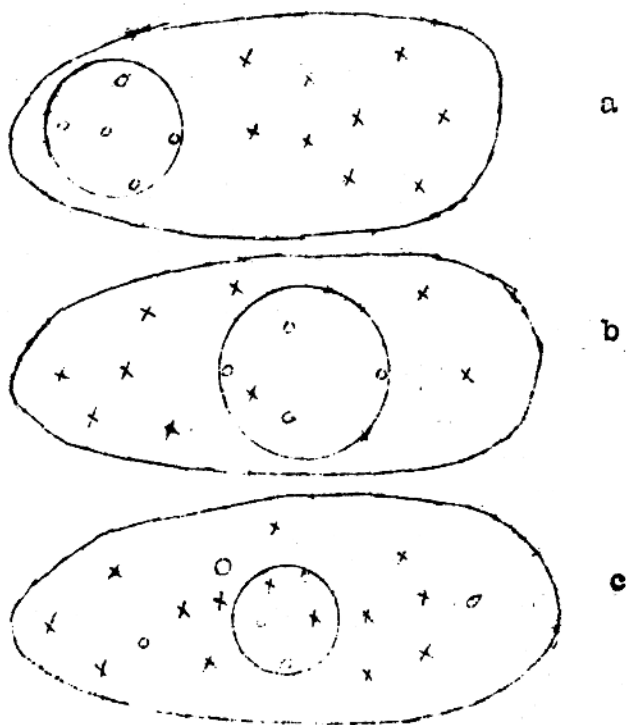


图 6.6 相关文献与不相关文献的分布

上图中， \circ 表示相关文献， \times 表示不相关文献。容易看出，a 的情况最好。我们可以画出一个圆，其中包含了全部相关文献，而没有包含不相关文献。c 的情况最差，我们无论在什么地方画圆和画多大的圆，都不可避免地要遗漏一些相关文献和包含进一些不相关文献。提问向量实际上是文献空间中的一个点（提问向量的端点）。输出时用户指定的阈值实际上是一个半径，这两个值共同确