

RESEARCH ON FAST ITEMSET MINING ALGORITHMS  
AND THEIR PERFORMANCE

# 高性能数据挖掘

——快速项集挖掘算法

及性能研究

屈俊峰 著



中国水利水电出版社

[www.waterpub.com.cn](http://www.waterpub.com.cn)

# 高性能数据挖掘——快速项集 挖掘算法及性能研究

屈俊峰 著



## 内 容 提 要

本书主要探讨数据挖掘中的项集挖掘问题，详细介绍了频繁项集、高可用项集、最大频繁项集、频繁闭项集的定义、挖掘算法、搜索空间剪枝技术、性能优化等方面的内容。本书的重点在于介绍如何提高挖掘速度、提升挖掘时的内存使用效率；本书的特色在于不仅对这些挖掘方法与技术在理论上进行描述，同时作者执行了严格的实验用以佐证结论。

本书可作为高年级本科生、数据挖掘方向的研究生、有兴趣青年学者的参考书。

### 图书在版编目 (C I P ) 数据

高性能数据挖掘：快速项集挖掘算法及性能研究 /  
屈俊峰著. -- 北京 : 中国水利水电出版社, 2018.8  
ISBN 978-7-5170-6691-0

I. ①高… II. ①屈… III. ①数据采集—研究 IV.  
①TP274

中国版本图书馆CIP数据核字(2018)第175240号

责任编辑：杨元泓

封面设计：李 佳

书 名	高性能数据挖掘——快速项集挖掘算法及性能研究 GAOXINGNENG SHUJU WAJUE——KUAISU XIANGJI WAJUE SUANFA JI XINGNENG YANJIU
作 者	屈俊峰 著
出版发行	中国水利水电出版社 (北京市海淀区玉渊潭南路1号D座 100038) 网址: www.waterpub.com.cn E-mail: mchannel@263.net (万水) sales@waterpub.com.cn 电话: (010) 68367658 (营销中心)、82562819 (万水) 全国各地新华书店和相关出版物销售网点
经 售	北京万水电子信息有限公司 三河市兴国印务有限公司
排 版	170mm×240mm 16开本 10.75印张 200千字
印 刷	2018年8月第1版 2018年8月第1次印刷
规 格	0001—2000册
版 次	48.00元
印 数	
定 价	

凡购买我社图书，如有缺页、倒页、脱页的，本社营销中心负责调换

版权所有·侵权必究

# 前　　言

目前，大部分的计算机应用系统都以数据库作为底层支撑。例如，工业控制系统通过传感网络不断收集工业过程的各种参数，将其存入数据库，这些参数是系统下一步的操作基础，也能为工程师提供重要的状态信息。在网站日志分析系统中，每个页面的访问情况或每个用户的访问历史被存储在数据库里，对这些信息进行分析可以发现用户感兴趣的内容，进而改进网站布局或向用户推介相关信息。在零售商业中，一个顾客一次购物的所有商品形成数据库中的一条记录，对这些记录进行分析可以提示短缺商品或更好地布局终端货架。

在大数据时代，如何从数据库中挖掘有用信息是一个吸引人且具有挑战性的任务。通常，一个数据库包含若干条记录，每条记录由若干项组成。例如，一个顾客的一次购物记录是其所购商品名称的集合。这里，一个商品名称即是一个项。由项组成的集合称为项集，项集是能够从数据库中挖掘出的重要信息。项集挖掘问题通常被定义为：给定一个数据库及特定标准下的一个条件，要求从库中挖掘出能满足这一条件的所有项集。项集挖掘任务通常非常耗时，因为搜索空间巨大。如果一个数据库中出现了 $n$ 个项，那么存在 $2^n$ 个项集在搜索空间中。

因此，如何高效挖掘数据库中的项集是一项有挑战性的任务。本书选定了四类项集挖掘问题。本书的第1章为引言，给出概要性的介绍。第2~6章介绍项集挖掘领域最经典的频繁项集挖掘问题。第7~9章介绍了最近的研究热点高可用项集挖掘问题。第10、11章分别介绍了最大频繁项集挖掘问题及频繁闭项集挖掘问题。对每一类问题，在介绍了问题定义之后，作者以挖掘性能提升为导向，分别从挖掘算法、搜索空间剪枝技术、性能优化等方面进行了详细阐述。本书的特色在于，大部分的方法与技术都配有实验验证，实验结论不仅能够佐证方法技术的有效性，而且也给了读者非常直观的认识。

本书图文并茂、以实用为主，力求能够成为高年级本科生、数据挖掘方向的研究生、有兴趣的青年学者在研究相关主题时的参考书。本书的编写得到了湖北文理学院教师科研能力培育基金项目（编号：2017kypy037）的支持。本书在编写内容与特点上均进行了一种新的尝试，缺点和错误在所难免，由于编者水平有限，希望广大读者给予批评和指正，对此作者深表谢意。

屈俊峰

湖北文理学院计算机工程学院

2018年5月

# 目 录

## 前言

<b>第1章 概述</b>	1
1.1 项集：数据挖掘研究领域的焦点之一	3
1.2 频繁项集挖掘问题的研究历史	5
1.3 高可用项集挖掘问题的研究历史	7
1.4 本书的主要内容	9
<b>第2章 频繁项集挖掘问题</b>	11
2.1 概述	12
2.1.1 问题形式化定义	12
2.1.2 搜索空间与方法	13
2.2 基础频繁项集挖掘算法介绍	14
2.2.1 经典的候选生成 Apriori 算法	15
2.2.2 以垂直视角处理数据库的 Eclat 算法	16
2.2.3 基于前缀树结构的 FP-growth 算法	17
2.3 性能测试的软硬件环境	19
2.3.1 数据库描述	19
2.3.2 参照算法介绍	20
2.3.3 其他软硬件设施	22
2.4 实验一：三种基础算法的性能测试	23
2.4.1 实验结果	23
2.4.2 性能评价	24
<b>第3章 BFP-growth：快速模式增长算法</b>	27
3.1 经典模式增长算法的性能分析	28
3.1.1 影响 FP-growth 性能的三个因素	28
3.1.2 ICDM 最佳算法：FPgrowth*	28
3.2 批量模式增长算法：BFP-growth	30
3.2.1 性能提升的途径	30
3.2.2 核心步骤：两次前缀树遍历	31
3.2.3 算法伪代码	34

3.3	BFP-growth 算法的性能分析.....	35
3.3.1	更少的遍历花费.....	35
3.3.2	FP-array 技术应该集成在 BFP-growth 中吗.....	36
3.3.3	无修饰的前缀树结构.....	37
3.4	实验二：BFP-growth 的性能测试及讨论.....	38
3.4.1	BFP-growth 及 FPgrowth*与基础算法的对比.....	38
3.4.2	实验结果讨论.....	38
3.5	小结.....	40
<b>第 4 章</b>	<b>基于结点集合结构的 NS 算法.....</b>	<b>41</b>
4.1	Eclat 及 FP-growth 算法的优缺点.....	42
4.2	结点集合结构（Node-set）.....	43
4.2.1	条件结点.....	44
4.2.2	结点拓扑序号.....	45
4.2.3	使用结点集合结构表示前缀树.....	46
4.3	NS 算法.....	47
4.3.1	映射前缀树到结点集合结构.....	47
4.3.2	从结点集合结构中挖掘频繁项集.....	48
4.3.3	一个例子.....	50
4.3.4	NS 算法的原子操作.....	51
4.4	实验三：NS 算法与其他快速挖掘算法的性能对比.....	51
4.4.1	实验结果.....	52
4.4.2	结果讨论：NS 算法的性能优势.....	53
4.5	小结.....	54
<b>第 5 章</b>	<b>用 Patricia*结构挖掘频繁项集.....</b>	<b>55</b>
5.1	研究动机.....	56
5.2	Patricia*结构.....	57
5.2.1	单孩子结点.....	58
5.2.2	构造 Patricia*结构.....	59
5.3	用 Patricia*结构挖掘频繁项集.....	60
5.3.1	先前的挖掘流程.....	60
5.3.2	改进的挖掘流程.....	61
5.3.3	PatriciaMine*算法.....	62
5.4	实验结果.....	63

5.4.1 结点数量统计 .....	64
5.4.2 性能对比 .....	65
5.5 小结 .....	66
<b>第 6 章 频繁项集挖掘算法的内存耗费 .....</b>	<b>68</b>
6.1 BFP-growth 算法内存使用情况分析 .....	69
6.2 NS 算法内存使用情况分析 .....	69
6.3 实验四：快速挖掘算法的内存耗费 .....	70
6.4 SP 算法 .....	71
6.4.1 研究动机 .....	71
6.4.2 基础知识 .....	72
6.4.3 挖掘频繁项集 .....	76
6.4.4 实验结果与结论 .....	79
<b>第 7 章 高可用项集挖掘问题 .....</b>	<b>80</b>
7.1 从频繁项集到高可用项集 .....	81
7.2 问题的形式化定义 .....	82
7.3 已有挖掘算法概述 .....	83
<b>第 8 章 非候选生成高可用项集挖掘算法 .....</b>	<b>87</b>
8.1 项集有用性列表结构 .....	88
8.1.1 初始有用性列表 .....	88
8.1.2 2-项集的有用性列表 .....	90
8.1.3 $k$ -项集有用性列表 ( $k \geq 3$ ) .....	91
8.2 HUI-Miner 算法 .....	92
8.2.1 剪枝策略 .....	93
8.2.2 算法伪代码 .....	94
8.3 HUI-Miner 算法的实现细节 .....	95
8.3.1 有用性列表表头 .....	95
8.3.2 重新标注 tid .....	95
8.3.3 交易权重有用性增加的顺序 .....	96
8.4 实验五：HUI-Miner 性能测试 .....	97
8.4.1 实验设置 .....	97
8.4.2 HUI-Miner 及对比算法的运行时间 .....	98
8.4.3 HUI-Miner 及对比算法的内存耗费 .....	99
8.4.4 项处理顺序对 HUI-Miner 性能的影响 .....	100

8.4.5 可扩展性.....	101
8.4.6 实验结果讨论 .....	102
8.5 小结.....	103
<b>第 9 章 快速识别高可用项集.....</b>	<b>105</b>
9.1 先前算法的性能瓶颈 .....	106
9.2 基本识别算法（BIA） .....	107
9.3 基于候选树的快速识别算法（FIA） .....	110
9.3.1 候选树结构.....	110
9.3.2 快速识别算法.....	111
9.4 算法分析： BIA 与 FIA.....	114
9.5 实验六： BIA 与 FIA 的性能对比 .....	115
9.5.1 高可用项集识别时间 .....	116
9.5.2 候选项集生成时间 .....	117
9.5.3 内存耗费 .....	117
9.5.4 实验结果分析.....	117
9.6 实验七： FIA-UP-Growth+和 HUI-Miner 的性能对比 .....	118
9.6.1 运行时间&内存耗费 .....	118
9.6.2 实验结果分析.....	119
9.7 小结.....	120
<b>第 10 章 最大频繁项集挖掘 .....</b>	<b>122</b>
10.1 介绍 .....	122
10.2 基本概念 .....	124
10.3 MAFIA 算法 .....	125
10.3.1 深度优先遍历 .....	125
10.3.2 搜索空间剪枝 .....	126
10.3.3 有效的 MFI 超集检查 .....	130
10.4 挖掘非最大频繁项集 .....	131
10.4.1 挖掘所有的频繁项集 .....	132
10.4.2 挖掘所有的频繁闭项集 .....	132
10.5 实施细节 .....	133
10.6 结论 .....	134
<b>第 11 章 频繁闭项集挖掘 .....</b>	<b>135</b>
11.1 介绍 .....	135

11.2	频繁项集挖掘 .....	137
11.2.1	基本定义 .....	137
11.2.2	先前的解决方案 .....	138
11.3	项集—记录标识符集合搜索树与等价类 .....	139
11.4	CHARM 算法设计与实现 .....	141
11.4.1	快速的闭项集子集合检查 .....	144
11.4.2	使用差异集合快速进行频繁计数 .....	145
11.4.3	其他优化及正确性 .....	147
11.5	实验结果 .....	148
11.6	结论 .....	149
	参考文献 .....	150

# 1

## 概述



### 本章导读

交易数据库是最常见的一类数据库，这种数据库的每一个条目是一个项的集合。当数据库变得越来越大时，人们希望从数据中发现一些规律。对于交易数据库，这些规律可以表示成“项集”。站在不同的角度，人们需要的项集是不一样的。本章将详细描述什么是项集，给出频繁项集及高可用项集的介绍、挖掘方法与历史。在本章的最后，我们给出了本书的组织结构。



### 本章要点

- 项集：数据挖掘研究领域的焦点之一
- 频繁项集挖掘问题的研究历史
- 高可用项集挖掘问题的研究历史
- 本书的组织结构

计算机的广泛使用促使现实生活中各领域中的大量数据被数字化，数据库技术的飞速发展则推动了这些数字化数据的有效管理和快捷使用。随着时间的推移和采集渠道的增多，汇入各种数据库中的数据越来越多，于是，一个现实的问题产生了：如何从这些体积日益膨胀的数据库中找出重要的或用户感兴趣的信息？

例如，随着条形码的广泛使用，超市中商品的销售信息得以大量地进入计算机。通常一个顾客的一次购物行为信息能够形成超市交易数据库中的一条记录。这种记录着用户交易行为的数据库不仅仅是一张数字和字符组成的表，其中还蕴含着丰富的信息。再如，店员可能想知道哪些商品的组合被频繁地同时售出，基于这样的信息，他能更好地组织货架布局促进商品销售；销售经理可能想知道哪些商品可以带来较大的利润，基于这样的信息，他能更好地安排营销计划。然而，这些信息并非显而易见而是隐藏在库中大量数据的背后。这就需要合适的计算机技术协助人们从数据的海洋中发现潜在的有用信息。

诸如上面的问题广泛地存在于现实生活中的各个领域，再比如说如何从病例数据库中提取分类规则用于辅助病人的诊疗，如何在天文数据库中寻找宇宙演化的规律然后作出科学预测，等等。这类问题的提出开创了数据库应用的一个新的研究领域，即数据挖掘与知识发现（Data Mining and Knowledge Discovery, DMKD），有些地方也称其为基于数据库的知识发现（Knowledge Discovery from Databases, KDD）。

在该领域中，有一个非常有意思的现象：一个问题被提出之后，解决方案往往层出不穷。比如非常著名的频繁项集挖掘（Frequent Itemset Mining）问题，自 Rakesh Agrawal 等在 1993 年 ACM Special Interest Group on Management of Data (ACMSIGMOD) 会议上正式提出后<sup>[15]</sup>，现在至少有 10 个较著名的算法能够解决此问题。即便是到最近，解决此问题的新算法还在继续被提出<sup>[28, 101, 109, 119]</sup>。

造成这一现象的本质原因是算法性能提升的无极限性。对于相同的问题，人们总是在不停地探寻着最佳算法（通常以算法的运行速度来衡量），然而现实却是“没有最佳，只有更佳”。以下客观因素加剧了这一现象：

- (1) 数据库规模的不断膨胀。
- (2) 硬件体系结构的改进，比如说并行和多核。

总而言之，对于该领域的许多问题，人们不仅关心是否有解决方案，更为关心的是解决方案是否足够地快，同时耗费的其他资源（主要是内存资源）又在可以接受的范围之内。

本书研究数据挖掘领域中的两类项集挖掘问题，一是频繁项集挖掘问题；二是高可用项集挖掘（High Utility Itemset Mining）问题。如上所述，我们不仅致力于探索解决这两个问题的新算法，并且非常重视相关算法（包括先前的算法和我们提出算法）的性能问题。本章下面几节我们将依次阐述为什么要关注项集挖掘，项集挖掘问题目前的研究现状，我们在项集挖掘方面的一些工作，以及本书的组织结构。

## 1.1 项集：数据挖掘研究领域的焦点之一

从实际中应用最广泛的关系数据库的视角观察，一个数据库能够被看作是一张二维表。这张表是由一条条记录组成，每条记录是由若干个项构成，每一个项又涉及到零个或多个属性值。项集表示若干项的集合。项集的特征可以从许多方面来刻画，例如在数据库中出现的次数是否足够多、价值是否很大，是否满足一定的约束条件，等等。基于不同的特性，从数据库中可以导出不同的项集的集合。对于项集某一（些）方面特性感兴趣的用户或应用系统，关心如何从数据库中挖掘出所有的满足这一（些）方面特性要求的全部项集。

最基本同时也是最重要的一种项集特征是项集在数据库中的出现次数，称为项集的支持度<sup>[15]</sup>。在许多应用中，用户希望在一个数据库中找出所有支持度超过其指定阈值的项集。此任务即为频繁项集挖掘。频繁项集挖掘在数据挖掘研究领域中扮演着一个非常重要的基础性角色。

首先，频繁项集挖掘问题是数据挖掘领域最早提出的问题之一（1993年），也是该领域中最受研究者关注的问题之一。根据谷歌学术搜索的统计，正式定义该问题的论文<sup>[15]</sup>目前已经被引用了上万次。其次，由频繁项集挖掘问题派生出来的数据挖掘问题非常多。派生出的问题可以分为三类，一类是重定义问题内部约束条件而产生的新问题；另一类是扩展问题外部约束条件而产生的新问题；第三类是同时重定义问题内、外部约束条件而产生的新问题。

在第一类问题中，比较著名的有以下几个：最大频繁项集挖掘问题要求挖掘出所有的最大频繁项集<sup>[25, 37, 39, 49, 63, 134]</sup>，一个最大频繁项集是频繁的且没有任何频繁项集是它的超集；频繁闭项集挖掘问题要求挖掘出所有的频繁闭项集<sup>[80, 81, 84, 115, 131, 132]</sup>，一个频繁闭项集是频繁的且它的任一频繁超集都和此闭项集有不同的支持度；top-k 频繁项集挖掘问题要求挖掘出支持度最高的 k 个频繁项集<sup>[45, 86, 95]</sup>，等等。从上面的问题描述中我们能够看出，从最大频繁项集集合或频繁闭项集集合中可以导出完整的频繁项集集合，不同的是从前者中导出的频繁项集没有支持度信息而从后者导出的频繁项集包含着完整的支持度信息。对于 top-k 频繁项集的挖掘则不需要指定最小支持度阈值。

第二类问题是频繁项集挖掘问题在不同软硬件环境下的扩展，例如：从不确定数据中挖掘频繁项集<sup>[13, 20, 26, 107]</sup>，从数据流中挖掘频繁项集<sup>[31, 46, 69]</sup>，从敏感或隐私数据中挖掘频繁项集<sup>[21, 117]</sup>，在多核 CPU 上进行频繁项集挖掘<sup>[32, 71]</sup>，以及在内存受限条件下进行频繁项集挖掘<sup>[94, 101]</sup>。

第三类问题是前两类问题的叠加，例如在数据流上挖掘 top-k 的频繁项集<sup>[53]</sup>。解决这些派生问题的方法都直接受原始问题解决方法的影响，即只要有新的高效

的频繁项集挖掘算法被提出，那么就有可能通过改进新算法来解决这些派生问题。因此，频繁项集挖掘问题是以上所有问题的基础和核心。

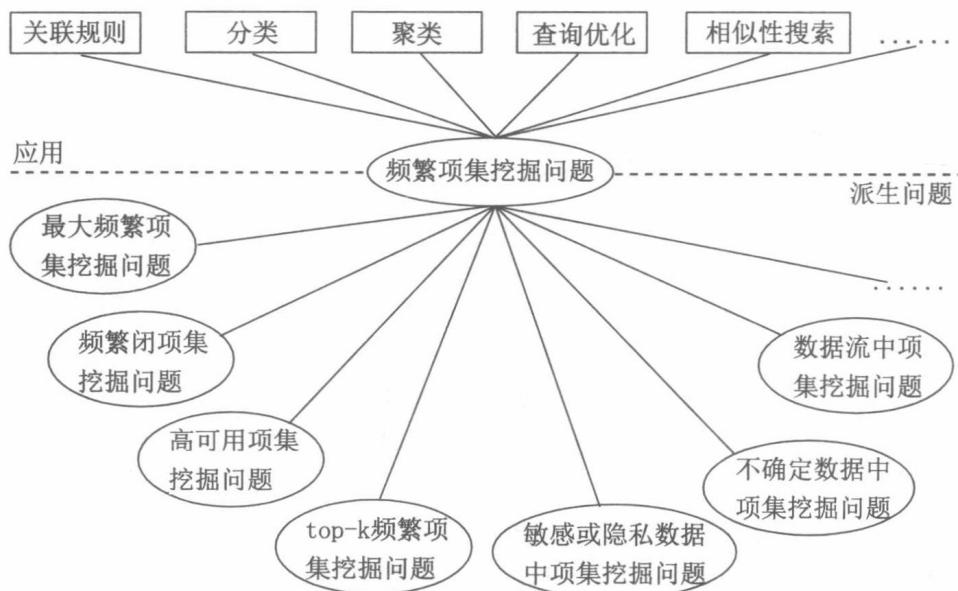


图 1-1 项集挖掘相关问题及应用

再者，挖掘出来的频繁项集有着广泛的应用。例如，频繁项集最开始是作为关联规则挖掘的中间结果而开始得到研究者的关注<sup>[15, 16]</sup>，后来发现也可以基于频繁项集做分类及聚类处理<sup>[14, 19, 30, 54, 116, 126]</sup>。频繁项集也能用在演绎数据库和查询处理的优化上<sup>[75, 124]</sup>。例如，对查询处理的一个非常简单的优化可以是先找出用户频繁执行的一些查询操作，然后在查询结果尺寸不大的情况下，对结果做直接保存。这样做的好处有两点，一是可以实现查询提示，即用户在输入部分关键字后，系统对剩下的可能关键字作出提示，既方便用户输入又简化系统处理流程，目前像著名的搜索引擎谷歌及百度都实现了这样的功能；二是当这样的查询提交到系统之后，系统可以快速直接返回查询结果。最近的一些文献表明频繁项集还能在下面一些应用中发挥作用：复杂结构的索引及相似性搜索<sup>[121-123]</sup>，时空数据挖掘<sup>[27, 55, 120, 135]</sup>，多媒体数据挖掘<sup>[127]</sup>，Web 数据挖掘<sup>[29, 34, 51, 85, 90, 106]</sup>，甚至软件缺陷的定位<sup>[60, 61, 62, 65]</sup>。

图 1-1 总结了本节的内容，可以看出，项集挖掘的研究不仅具有重要的理论意义，而且也有广泛的应用价值。如何挖掘具有某些特征的项集以及如何高效地挖掘出这些项集毫无疑问已经成为数据挖掘研究领域的焦点之一。本书主要研究频繁项集及高可用项集的挖掘算法及其性能问题，下面我们将详细介绍这两个问题的研究历史。

## 1.2 频繁项集挖掘问题的研究历史

1993年，频繁项集挖掘作为关联规则挖掘的一个子问题由Rakesh Agrawal在ACMSIGMOD会议上首先提出来<sup>[15]</sup>。当发现这个子题是关联规则挖掘中最困难的步骤时，人们将这个子问题独立出来专门展开研究。解决这个问题的第一个算法是AIS<sup>[15]</sup>。

早期的Apriori算法已经被许多教科书收录并奉为经典<sup>[16, 42, 64, 76]</sup>。Apriori有两个经常被诟病的问题：生成了大量的候选项集及多遍的数据库扫描。针对这两个问题一些Apriori算法的改进版本相继被提出<sup>[100, 108]</sup>。例如，Partition算法首先把一个数据库水平切分成可以导入内存的一个个子库，接着依次挖掘每个子库上的频繁项集，最后通过一次数据库遍历来验证从子库中生成的频繁项集在整个数据库上是否仍是频繁的<sup>[100]</sup>。另一些Apriori的优化则采用哈希表来实现项集的快速定位用以加速计数过程<sup>[79]</sup>。

在1997年的ACM Special Interest Group on Knowledge Discovery and Data Mining(ACM SIGKDD)会议上由Mohammed J. Zaki提出的Eclat算法非常新颖地采用了垂直的数据库视角来挖掘频繁项集<sup>[133]</sup>，之后Mohammed J.Zaki连续提出了扩展性非常好的Eclat版本和集成了diffset技术的Eclat算法<sup>[129, 130]</sup>，其中后者称为dEclat算法。dEclat显著地减少了Eclat算法核心结构Tid-list的长度，从而获得了明显的性能提升。

三年之后，Jiawei Han等在2000年的ACM SIGMOD会议上提出了著名的FP-growth算法<sup>[43, 44]</sup>，在这个算法中，他们首次将一个数据库压缩成一种前缀树结构，并在此结构上采用递归构造条件数据库的方式来挖掘频繁项集。FP-growth算法在数据结构上与之前的算法大相径庭，此后，许多基于FP-growth的新算法被陆续开发出来。例如，Jiawei Han针对FP-growth算法在稀疏数据库上性能较差的问题，2001年提出了能更有效地处理这种数据库的H-struct结构，以及基于此结构的H-mine算法<sup>[82, 83]</sup>。

Junqiang Liu等人在2002年提出了OP算法<sup>[70]</sup>，此算法是FP-growth算法和H-mine算法的混合体，能根据数据库的特性自动地在FP-growth算法和H-mine算法之间切换。在这一时期其他较有影响的算法包括，采用位图表示数据库的VIPER算法<sup>[11, 103]</sup>，采用树投影策略的TreeProject算法<sup>[12]</sup>，以及采用直接计数的DCI算法<sup>[77, 78, 89]</sup>。

由于频繁项集的重要理论意义与广泛应用价值逐步地凸现出来（见上节），在认识到实际解决此问题需要大量的计算时间后，为了对大量出现的算法作性能上的对比，找出快速的算法，著名国际会议International Conference on Data Mining

(ICDM) 在 2003 年及 2004 年连续举办了两次 Frequent Itemset Mining Implementation (FIMI) 工作组。会议的组织者鼓励全世界所有对此问题感兴趣的研究者提交自己的算法和相关论文，然后由组委会独立地对所有提交的算法进行性能上的评测。这次会议不仅接受对已有算法的高性能实施，而且也接受新提出的算法。于是，大量性能令人激动的算法出现在这两次工作组上。例如，Ferenc Bodon 和 Walter A. Kosters 等分别提出了两种 Apriori 算法的快速实施方案<sup>[22,52]</sup>，Christian Borgelt 提出了基于前缀树的 Apriori 及 Eclat 算法<sup>[23]</sup>，Lars Schmidt-Thieme 详细地研究了 Eclat 算法的特性，并提出了自己的 Eclat 版本<sup>[102]</sup>。一系列 FP-growth 算法的改进和优化也在次会议上被提出。例如，Liu Guimei 详细地研究了 FP-growth 的特点及流程，识别出了其中几个制约其性能的关键因素，提出了 AFOPT 算法<sup>[66,67,68]</sup>；Osmar R. Zaiane 提出 COFI-tree 结构可以减少候选生成的花费<sup>[128]</sup>；Andrea Pietracaprina 提出了采用 Patricia Tries 结构来快速地实现 FP-growth 流程的 PatriciaMine 算法<sup>[50, 87]</sup>；Balazs Racz 采用了线性投影的方式提出的 nonordfp 算法能够快速地构造前缀树<sup>[96]</sup>。

在 2003 年的工作组中，经过组委会独立的性能测试后，由 Gosta Grahne 和 Jianfei Zhu 提出的 FPgrowth\* 算法脱颖而出<sup>[38,40]</sup>，FPgrowth\* 是集成 FP-array 技术的 FP-growth 算法的一个高效变种。在 2004 年的工作组中，由 Takeaki Uno, Masashi Kiyomi 和 Hiroki Arimura 提出的 LCMv2 算法拔得头筹<sup>[112, 113]</sup>，LCMv2 算法是集多种优化策略于一体的杂交算法。

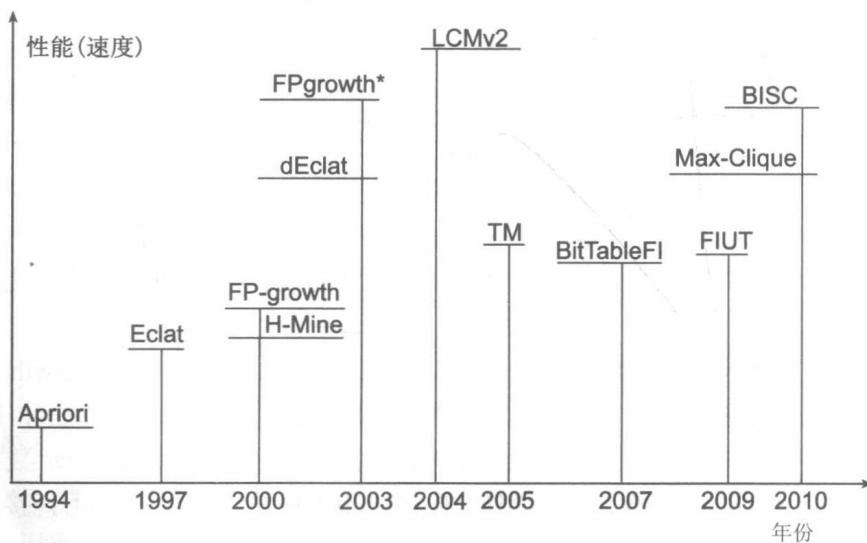


图 1-2 频繁项集挖掘算法发展历史

2003、2004 年 ICDM 的 FIMI 工作组把频繁项集挖掘算法的研究推到了一个顶峰。由于 FPgrowth\* 算法及 LCMv2 算法优异的性能表现（体现在极快的运行速

度上), 此后, 在性能上能够超越他们的算法可谓凤毛麟角。

2004 年之后提出的一些较有影响的算法如下。2005 年 Mingjun Song 提出了 TM 算法<sup>[104, 105]</sup>, TM 算法采用了交易树结构, 这是一种前缀树和 Tid-list 的混合结构。同样是 2005 年由 Amol Ghoting 等提出了 Tiling 算法<sup>[35, 36]</sup>, 此算法采用了 Cache 敏感的策略来提高性能。2007 年, Dong Jie 等提出了 BitTableFI 算法<sup>[33]</sup>, 这个算法是先前基于位图算法的扩展。

2009 年, Yuh-Jiuan Tsay 等提出了 FIUT 算法<sup>[109]</sup>, 采用从最大频繁项集往下分解的方式, 逐层生成长度较小的频繁项集。2010 年 Xie Yan 等提出了 Max-Clique 算法<sup>[119]</sup>, 采用图论中的相关理论来挖掘频繁项集。同年, Chen Jinlin 也提出了性能优异的杂交算法 BISC<sup>[28]</sup>。2011 年 Benjamin Schlegel 等提出了 CFP-growth 算法能够极大地减少 FP-growth 的内存消耗<sup>[101]</sup>。

然而, 这些算法在性能上, 主要指运行时间上, 很难胜过 FPgrowth\* 及 LCMv2 算法。例如, TM 算法在文献中实测的性能只是接近 FPgrowth\*; BitTableFI 和 FIUT 算法仅仅是优于基本的 FP-growth 算法 (不是 FPgrowth\*); Tiling 算法只是在特定机器上性能表现优异, 而 Max-Clique 算法则是在特定的数据库上能够和 LCMv2 匹敌。这些算法中最有希望的是 BISC 算法, 按照文献中的测试它的性能非常好, 但问题是对于杂交算法, 很难确定到底是其中的什么策略导致了较好的性能, 况且此算法的过程过于复杂不易实现。因此, 可以认为 FPgrowth\* 及 LCMv2 算法是当前最快的算法或属于最快算法的集合。图 1-2 给出了频繁项集挖掘算法发展历史。

## 1.3 高可用项集挖掘问题的研究历史

目前, 大部分的计算机应用系统都以数据库作为底层支撑。例如, 工业控制系统通过传感网络不断收集工业过程的各种参数, 将其存入数据库, 这些参数是系统下一步的操作基础也能为工程师提供重要的状态信息。在网站日志分析系统中, 每个页面的访问情况或每个用户的访问历史被存储在数据库里, 对这些信息进行分析可以发现用户感兴趣的内容, 进而改进网站布局或向用户推介相关信息。在零售商业中, 一个顾客一次购物的所有商品形成数据库中的一条记录, 对这些记录进行分析可以提示短缺商品或更好地布局终端货架。

在大数据时代, 如何从数据库中挖掘有用信息是一个吸引人且具有挑战性的任务。通常, 一个数据库包含若干条记录, 每条记录由若干项组成。例如, 一个顾客的一次购物记录是其所购商品名称的集合。这里, 一个商品名称即是一个项。由项组成的集合称为项集, 项集是能够从数据库中挖掘出的重要信息。项集挖掘问题通常被定义为: 给定一个数据库及特定标准下的一个条件, 要求从库中挖掘

出能满足这一条件的所有项集。项集挖掘任务通常非常耗时，因为搜索空间巨大。如果一个数据库中出现了  $n$  个项，那么存在  $2^n$  个项集在搜索空间中。最常见的一种项集重要性衡量标准是支持度。如果一个项集中的每个项都出现在一条记录中，那么称这条记录包含这个项集或者这个项集出现在这条记录中。给定一个数据库，一个项集的支持度是库中包含此项集记录的数量。按照支持度标准，项集的支持度越大就越重要。经典的频繁项集挖掘问题，即是给定一个数据库和一个最小支持度阈值，要求找出支持度大于等于阈值的所有项集，这些项集称为频繁项集。此问题是数据挖掘领域最基础的问题之一，我们已在上节中作了讨论。

在频繁项集挖掘问题中，数据库中出现的所有项被赋予了同等价值，支持度反映了项集的重要性。然而，在一些场景下，项集的重要性并不能简单地由支持度衡量。例如，需要从零售数据库中挖掘出利润最大的商品组合。通常，项集{牙膏、香皂}的支持度要高于项集{手机、内存卡}的支持度，但是后者产生的利润或许要远高于前者。于是，研究者提出了项集重要性衡量的效用标准。在此标准下，每一个项被赋予了一个权重值（外部效用），一条记录中的每一个项关联着一个计数（内部效用）。一个项集的效用由包含此项集的所有记录中相关项的内、外部效用确定。同支持度标准一样，多数情况下人们感兴趣的是效用高的项集，故提出了高效用项集挖掘问题：给定一个数据库及效用阈值，要求挖掘出效用大于等于此阈值的所有项集（称为高效用项集）。

高效用项集有各种重要的应用，例如，在上面的例子中，从交易数据库中挖掘出高效用项集，即利润大的商品组合，有助于策划更加赢利的商业活动。

高可用项集挖掘问题是频繁项集挖掘问题的扩展，但因问题定义的不同，频繁项集挖掘算法不能直接用于高可用项集的挖掘。高可用项集挖掘问题在 2005 年由 YaoHong 在 Society for Industrial and Applied Mathematics on Data Mining (SIAMDM) 国际会议上正式提出来之前<sup>[125]</sup>，已经有一些算法，可以挖掘份额频繁项集，例如 ZP 算法<sup>[18]</sup>。份额频繁项集是高可用项集的一种特殊形式，适用于挖掘份额频繁项集的算法经改造后也可用于高可用项集的挖掘。

后来出现的高可用频繁项集挖掘算法按照基本思路的不同可以分为两大类。一类是基于 Apriori 的算法，他们逐层生成候选项集再扫描数据库计算项集的可用性，例如：Brock Barber 等在 2003 年提出的 ZP 及 ZSP 算法<sup>[18]</sup>；Yu-Chiang Li 等在 2005、2008 年连续提出的 FSH 算法<sup>[58]</sup>、ShFSH 算法<sup>[57]</sup>、DCG 算法<sup>[56]</sup>、FUM 算法<sup>[59]</sup>、以及 DCG+算法<sup>[59]</sup>。另一类是基于前缀树生成候选项集然后再进行测试的算法，例如 Chowdhury Farhan Ahmed 等在 2009 年提出了 IHUPTWU 算法<sup>[17]</sup>，Vincent S. Tseng 等 2010 提出的 UP-Growth 算法<sup>[111]</sup>及 2012 年提出的 UPGrowth+ 算法<sup>[110]</sup>。