

Frequent Pattern Mining

Charu C. Aggarwal • Jiawei Han
Editors

Frequent Pattern Mining

 Springer

Editors

Charu C. Aggarwal
IBM T. J. Watson Research Center
Yorktown Heights
New York
USA

Jiawei Han
University of Illinois at Urbana-Champaign
Urbana
Illinois
USA

ISBN 978-3-319-07820-5 ISBN 978-3-319-07821-2 (eBook)

DOI 10.1007/978-3-319-07821-2

Springer Cham Heidelberg New York Dordrecht London

Library of Congress Control Number: 2014944536

© Springer International Publishing Switzerland 2014

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

Preface

The field of data mining has four main “super-problems” corresponding to clustering, classification, outlier analysis, and frequent pattern mining. Compared to the other three problems, the frequent pattern mining model was formulated relatively recently. In spite of its shorter history, frequent pattern mining is considered the marquee problem of data mining. The reason for this is that interest in the data mining field increased rapidly soon after the seminal paper on association rule mining by Agrawal, Imielinski, and Swami. The earlier data mining conferences were often dominated by a large number of frequent pattern mining papers. This is one of the reasons that frequent pattern mining has a very special place in the data mining community. At this point, the field of frequent pattern mining is considered a mature one.

While the field has reached a relative level of maturity, very few books cover different aspects of frequent pattern mining. Most of the existing books are either too generic or do not cover frequent pattern mining in an exhaustive way. A need exists for an exhaustive book on the topic that can cover the different nuances in an exhaustive way.

This book provides comprehensive surveys in the field of frequent pattern mining. Each chapter is designed as a survey that covers the key aspects of the field of frequent pattern mining. The chapters are typically of the following types:

- *Algorithms*: In these cases, the key algorithms for frequent pattern mining are explored. These include join-based methods such as *Apriori*, and pattern-growth methods.
- *Variations*: Many variations of frequent pattern mining such as interesting patterns, negative patterns, constrained pattern mining, or compressed patterns are explored in these chapters.
- *Scalability*: The large sizes of data in recent years has led to the need for big data and streaming frameworks for frequent pattern mining. Frequent pattern mining algorithms need to be modified to work with these advanced scenarios.
- *Data Types*: Different data types lead to different challenges for frequent pattern mining algorithms. Frequent pattern mining algorithms need to be able to work with complex data types, such as temporal or graph data.

- *Applications*: In these chapters, different applications of frequent pattern mining are explored. These includes the application of frequent pattern mining methods to problems such as clustering and classification. Other more complex algorithms are also explored.

This book is, therefore, intended to provide an overview of the field of frequent pattern mining, as it currently stands. It is hoped that the book will serve as a useful guide for students, researchers, and practitioners.

Contents

| | | |
|----------|--|-----------|
| 1 | An Introduction to Frequent Pattern Mining | 1 |
| | Charu C. Aggarwal | |
| 1 | Introduction | 1 |
| 2 | Frequent Pattern Mining Algorithms | 3 |
| 2.1 | Frequent Pattern Mining with the Traditional Support Framework | 4 |
| 2.2 | Interesting and Negative Frequent Patterns | 6 |
| 2.3 | Constrained Frequent Pattern Mining | 7 |
| 2.4 | Compressed Representations of Frequent Patterns | 7 |
| 3 | Scalability Issues in Frequent Pattern Mining | 8 |
| 3.1 | Frequent Pattern Mining in Data Streams | 8 |
| 3.2 | Frequent Pattern Mining with Big Data | 9 |
| 4 | Frequent Pattern Mining with Advanced Data Types | 9 |
| 4.1 | Sequential Pattern Mining | 10 |
| 4.2 | Spatiotemporal Pattern Mining | 10 |
| 4.3 | Frequent Patterns in Graphs and Structured Data | 11 |
| 4.4 | Frequent Pattern Mining with Uncertain Data | 11 |
| 5 | Privacy Issues | 12 |
| 6 | Applications of Frequent Pattern Mining | 13 |
| 6.1 | Applications to Major Data Mining Problems | 13 |
| 6.2 | Generic Applications | 13 |
| 7 | Conclusions and Summary | 14 |
| | References | 14 |
| 2 | Frequent Pattern Mining Algorithms: A Survey | 19 |
| | Charu C. Aggarwal, Mansurul A. Bhuiyan and Mohammad Al Hasan | |
| 1 | Introduction | 19 |
| 1.1 | Definitions | 22 |
| 2 | Join-Based Algorithms | 23 |
| 2.1 | Apriori Method | 24 |
| 2.2 | DHP Algorithm | 27 |
| 2.3 | Special Tricks for 2-Itemset Counting | 28 |

| | | |
|----------|---|-----------|
| 2.4 | Pruning by Support Lower Bounding | 28 |
| 2.5 | Hypercube Decomposition | 29 |
| 3 | Tree-Based Algorithms | 29 |
| 3.1 | AIS Algorithm | 31 |
| 3.2 | TreeProjection Algorithms | 32 |
| 3.3 | Vertical Mining Algorithms | 36 |
| 4 | Recursive Suffix-Based Growth | 39 |
| 4.1 | The FP-Growth Approach | 41 |
| 4.2 | Variations | 45 |
| 5 | Maximal and Closed Frequent Itemsets | 47 |
| 5.1 | Definitions | 47 |
| 5.2 | Frequent Maximal Itemset Mining Algorithms | 48 |
| 5.3 | Frequent Closed Itemset Mining Algorithms | 55 |
| 6 | Other Optimizations and Variations | 57 |
| 6.1 | Row Enumeration Methods | 57 |
| 6.2 | Other Exploration Strategies | 58 |
| 7 | Reducing the Number of Passes | 58 |
| 7.1 | Combining Passes | 58 |
| 7.2 | Sampling Tricks | 59 |
| 7.3 | Online Association Rule Mining | 60 |
| 8 | Conclusions and Summary | 61 |
| | References | 61 |
| 3 | Pattern-Growth Methods | 65 |
| | Jiawei Han and Jian Pei | |
| 1 | Introduction | 66 |
| 2 | FP-Growth: Pattern Growth for Mining Frequent Itemsets | 68 |
| 3 | Pushing More Constraints in Pattern-Growth Mining | 72 |
| 4 | PrefixSpan: Mining Sequential Patterns by Pattern Growth | 74 |
| 5 | Further Development of Pattern Growth-Based Pattern Mining Methodology | 77 |
| 6 | Conclusions | 78 |
| | References | 79 |
| 4 | Mining Long Patterns | 83 |
| | Feida Zhu | |
| 1 | Introduction | 83 |
| 2 | Preliminaries | 84 |
| 3 | A Pattern Lattice Model | 86 |
| 4 | Pattern Enumeration Approach | 87 |
| 4.1 | Breadth-First Approach | 87 |
| 4.2 | Depth-First Approach | 88 |
| 5 | Row Enumeration Approach | 89 |
| 6 | Pattern Merge Approach | 92 |
| 6.1 | Piece-wise Pattern Merge | 93 |

| | | |
|----------|---|------------|
| 6.2 | Fusion-style Pattern Merge | 98 |
| 7 | Pattern Traversal Approach | 101 |
| 8 | Conclusion | 102 |
| | References | 103 |
| 5 | Interesting Patterns | 105 |
| | Jilles Vreeken and Nikolaj Tatti | |
| 1 | Introduction | 106 |
| 2 | Absolute Measures | 107 |
| 2.1 | Frequent Itemsets | 107 |
| 2.2 | Tiles | 112 |
| 2.3 | Low Entropy Sets | 114 |
| 3 | Advanced Methods | 114 |
| 4 | Static Background Models | 115 |
| 4.1 | Independence Model | 116 |
| 4.2 | Beyond Independence | 119 |
| 4.3 | Maximum Entropy Models | 120 |
| 4.4 | Randomization Approaches | 123 |
| 5 | Dynamic Background Models | 124 |
| 5.1 | The General Idea | 125 |
| 5.2 | Maximum Entropy Models | 125 |
| 5.3 | Tile-based Techniques | 126 |
| 5.4 | Swap Randomization | 128 |
| 6 | Pattern Sets | 128 |
| 6.1 | Itemsets | 129 |
| 6.2 | Tiles | 130 |
| 6.3 | Swap Randomization | 130 |
| 7 | Conclusions | 131 |
| | References | 132 |
| 6 | Negative Association Rules | 135 |
| | Luiza Antonie, Jundong Li and Osmar Zaiane | |
| 1 | Introduction | 135 |
| 2 | Negative Patterns and Negative Association Rules | 136 |
| 3 | Current Approaches | 138 |
| 4 | Associative Classification and Negative Association Rules | 143 |
| 5 | Conclusions | 143 |
| | References | 144 |
| 7 | Constraint-Based Pattern Mining | 147 |
| | Siegfried Nijssen and Albrecht Zimmermann | |
| 1 | Introduction | 147 |
| 2 | Problem Definition | 148 |
| 2.1 | Constraints | 149 |
| 3 | Level-Wise Algorithm | 152 |

| | | |
|----------|--|------------|
| 3.1 | Generic Algorithm | 153 |
| 4 | Depth-First Algorithm | 154 |
| 4.1 | Basic Algorithm | 154 |
| 4.2 | Constraint-based Itemset Mining | 155 |
| 4.3 | Generic Frameworks | 158 |
| 4.4 | Implementation Considerations | 159 |
| 5 | Languages | 159 |
| 6 | Conclusions | 162 |
| | References | 162 |
| 8 | Mining and Using Sets of Patterns through Compression | 165 |
| | Matthijs van Leeuwen and Jilles Vreeken | |
| 1 | Introduction | 165 |
| 2 | Foundations | 167 |
| 2.1 | Kolmogorov Complexity | 168 |
| 2.2 | MDL | 169 |
| 2.3 | MDL in Data Mining | 171 |
| 3 | Compression-based Pattern Models | 171 |
| 3.1 | Pattern Models for MDL | 172 |
| 3.2 | Code Tables | 173 |
| 3.3 | Instances of Compression-based Models | 179 |
| 4 | Algorithmic Approaches | 181 |
| 4.1 | Candidate Set Filtering | 181 |
| 4.2 | Direct Mining of Patterns that Compress | 184 |
| 5 | MDL for Data Mining | 185 |
| 5.1 | Classification | 186 |
| 5.2 | A Dissimilarity Measure for Datasets | 188 |
| 5.3 | Identifying and Characterizing Components | 189 |
| 5.4 | Other Data Mining Tasks | 191 |
| 5.5 | The Advantage of Pattern-based Models | 192 |
| 6 | Challenges Ahead | 193 |
| 6.1 | Toward Mining Structured Data | 193 |
| 6.2 | Generalization | 194 |
| 6.3 | Task- and/or User-specific Usefulness | 194 |
| 7 | Conclusions | 195 |
| | References | 196 |
| 9 | Frequent Pattern Mining in Data Streams | 199 |
| | Victor E. Lee, Ruoming Jin and Gagan Agrawal | |
| 1 | Introduction | 200 |
| 2 | Preliminaries | 201 |
| 2.1 | Frequent Pattern Mining: Definition | 201 |
| 2.2 | Data Windows | 202 |
| 2.3 | Frequent Item Mining | 203 |
| 3 | Frequent Itemset Mining Algorithms | 204 |

| | | |
|-----------|---|------------|
| 3.1 | Mining the Full Data Stream | 206 |
| 3.2 | Recently Frequent Itemsets | 209 |
| 3.3 | Closed and Maximal Itemsets | 214 |
| 3.4 | Mining Data Streams with Uncertain Data | 216 |
| 4 | Mining Patterns Other than Itemsets | 216 |
| 4.1 | Subsequences | 217 |
| 4.2 | Subtrees and Semistructured Data | 218 |
| 4.3 | Subgraphs | 219 |
| 5 | Concluding Remarks | 219 |
| | References | 220 |
| 10 | Big Data Frequent Pattern Mining | 225 |
| | David C. Anastasiu, Jeremy Iverson, Shaden Smith and George Karypis | |
| 1 | Introduction | 225 |
| 2 | Frequent Pattern Mining: Overview | 226 |
| 2.1 | Preliminaries | 226 |
| 2.2 | Basic Mining Methodologies | 228 |
| 3 | Paradigms for Big Data Computation | 232 |
| 3.1 | Principles of Parallel Algorithms | 232 |
| 3.2 | Shared Memory Systems | 233 |
| 3.3 | Distributed Memory Systems | 234 |
| 4 | Frequent Itemset Mining | 236 |
| 4.1 | Memory Scalability | 236 |
| 4.2 | Work Partitioning | 239 |
| 4.3 | Dynamic Load Balancing | 241 |
| 4.4 | Further Considerations | 242 |
| 5 | Frequent Sequence Mining | 242 |
| 5.1 | Serial Frequent Sequence Mining | 243 |
| 5.2 | Parallel Frequent Sequence Mining | 245 |
| 6 | Frequent Graph Mining | 250 |
| 6.1 | Serial Frequent Graph Mining | 250 |
| 6.2 | Parallel Frequent Graph Mining | 252 |
| 7 | Conclusion | 255 |
| | References | 256 |
| 11 | Sequential Pattern Mining | 261 |
| | Wei Shen, Jianyong Wang and Jiawei Han | |
| 1 | Introduction | 261 |
| 2 | Problem Definition | 263 |
| 3 | Apriori-based Approaches | 264 |
| 3.1 | Horizontal Data Format Algorithms | 264 |
| 3.2 | Vertical Data Format Algorithms | 268 |
| 4 | Pattern Growth Algorithms | 271 |
| 4.1 | FreeSpan | 271 |
| 4.2 | PrefixSpan | 272 |

| | | |
|-----------|---|------------|
| 5 | Extensions | 274 |
| 5.1 | Closed Sequential Pattern Mining | 274 |
| 5.2 | Multi-level, Multi-dimensional Sequential Pattern Mining . . . | 276 |
| 5.3 | Incremental Methods | 277 |
| 5.4 | Hybrid Methods | 278 |
| 5.5 | Approximate Methods | 279 |
| 5.6 | Top- k Closed Sequential Pattern Mining | 279 |
| 5.7 | Frequent Episode Mining | 280 |
| 6 | Conclusions and Summary | 281 |
| | References | 281 |
| 12 | Spatiotemporal Pattern Mining: Algorithms and Applications | 283 |
| | Zhenhui Li | |
| 1 | Introduction | 283 |
| 2 | Basic Concept | 284 |
| 2.1 | Spatiotemporal Data Collection | 284 |
| 2.2 | Data Preprocessing | 285 |
| 2.3 | Background Information | 286 |
| 3 | Individual Periodic Pattern | 286 |
| 3.1 | Automatic Discovery of Periodicity in Movements | 287 |
| 3.2 | Frequent Periodic Pattern Mining | 289 |
| 3.3 | Using Periodic Pattern for Location Prediction | 289 |
| 4 | Pairwise Movement Patterns | 290 |
| 4.1 | Similarity Measure | 290 |
| 4.2 | Generic Pattern | 292 |
| 4.3 | Behavioral Pattern | 294 |
| 4.4 | Semantic Patterns | 296 |
| 5 | Aggregate Patterns over Multiple Trajectories | 298 |
| 5.1 | Frequent Trajectory Pattern Mining | 298 |
| 5.2 | Detection of Moving Object Cluster | 300 |
| 5.3 | Trajectory Clustering | 302 |
| 6 | Summary | 304 |
| | References | 304 |
| 13 | Mining Graph Patterns | 307 |
| | Hong Cheng, Xifeng Yan and Jiawei Han | |
| 1 | Introduction | 307 |
| 2 | Frequent Subgraph Mining | 308 |
| 2.1 | Problem Definition | 308 |
| 2.2 | Apriori-Based Approach | 309 |
| 2.3 | Pattern-Growth Approach | 310 |
| 2.4 | Closed and Maximal Subgraphs | 311 |
| 2.5 | Mining Subgraphs in a Single Graph | 311 |
| 2.6 | The Computational Bottleneck | 313 |

| | | |
|-----------|--|------------|
| 3 | Mining Significant Graph Patterns | 314 |
| 3.1 | Problem Definition | 314 |
| 3.2 | gboost: A Branch-and-Bound Approach | 314 |
| 3.3 | gPLS: A Partial Least Squares Regression Approach | 317 |
| 3.4 | LEAP: A Structural Leap Search Approach | 319 |
| 3.5 | GraphSig: A Feature Representation Approach | 323 |
| 4 | Mining Representative Orthogonal Graphs | 326 |
| 4.1 | Problem Definition | 327 |
| 4.2 | Randomized Maximal Subgraph Mining | 327 |
| 4.3 | Orthogonal Representative Set Generation | 329 |
| 5 | Mining Dense Graph Patterns | 329 |
| 5.1 | Cliques and Quasi-Cliques | 330 |
| 5.2 | K-Core and K-Truss | 331 |
| 5.3 | Other Dense Subgraph Patterns | 332 |
| 6 | Mining Graph Patterns in Streams | 332 |
| 7 | Mining Graph Patterns in Uncertain Graphs | 334 |
| 8 | Conclusions | 336 |
| | References | 336 |
| 14 | Uncertain Frequent Pattern Mining | 339 |
| | Carson Kai-Sang Leung | |
| 1 | Introduction | 339 |
| 2 | The Probabilistic Model for Mining Expected Support-Based Frequent Patterns from Uncertain Data | 340 |
| 3 | Candidate Generate-and-Test Based Uncertain Frequent Pattern Mining | 343 |
| 4 | Hyperlinked Structure-Based Uncertain Frequent Pattern Mining .. | 344 |
| 5 | Tree-Based Uncertain Frequent Pattern Mining | 345 |
| 5.1 | UF-growth | 345 |
| 5.2 | UFP-growth | 346 |
| 5.3 | CUF-growth | 347 |
| 5.4 | PUF-growth | 349 |
| 6 | Constrained Uncertain Frequent Pattern Mining | 350 |
| 7 | Uncertain Frequent Pattern Mining from Big Data | 351 |
| 8 | Streaming Uncertain Frequent Pattern Mining | 353 |
| 8.1 | SUF-growth | 353 |
| 8.2 | UF-streaming for the Sliding Window Model | 354 |
| 8.3 | TUF-streaming for the Time-Fading Model | 355 |
| 8.4 | LUF-streaming for the Landmark Model | 356 |
| 8.5 | Hyperlinked Structure-Based Streaming Uncertain Frequent Pattern Mining | 356 |
| 9 | Vertical Uncertain Frequent Pattern Mining | 357 |
| 9.1 | U-Eclat: An Approximate Algorithm | 357 |
| 9.2 | UV-Eclat: An Exact Algorithm | 357 |
| 9.3 | U-VIPER: An Exact Algorithm | 358 |

| | | |
|-----------|---|------------|
| 10 | Discussion on Uncertain Frequent Pattern Mining | 360 |
| 11 | Extension: Probabilistic Frequent Pattern Mining | 361 |
| 11.1 | Mining Probabilistic Heavy Hitters | 361 |
| 11.2 | Mining Probabilistic Frequent Patterns | 362 |
| 12 | Conclusions | 364 |
| | References | 365 |
| 15 | Privacy Issues in Association Rule Mining | 369 |
| | Aris Gkoulalas-Divanis, Jayant Haritsa and Murat Kantarcioglu | |
| 1 | Introduction | 369 |
| 2 | Input Privacy | 370 |
| 2.1 | Problem Framework | 371 |
| 2.2 | Evolution of the Literature | 376 |
| 3 | Output Privacy | 379 |
| 3.1 | Terminology and Preliminaries | 380 |
| 3.2 | Taxonomy of ARH Algorithms | 381 |
| 3.3 | Heuristic and Exact ARH Algorithms | 382 |
| 3.4 | Metrics and Performance Analysis | 390 |
| 4 | Cryptographic Methods | 392 |
| 4.1 | Horizontally Partitioned Data | 394 |
| 4.2 | Vertically Partitioned Data | 396 |
| 5 | Conclusions | 398 |
| | References | 398 |
| 16 | Frequent Pattern Mining Algorithms for Data Clustering | 403 |
| | Arthur Zimek, Ira Assent and Jilles Vreeken | |
| 1 | Introduction | 403 |
| 2 | Generalizing Pattern Mining for Clustering | 406 |
| 2.1 | Generalized Monotonicity | 407 |
| 2.2 | Count Indexes | 410 |
| 2.3 | Pattern Explosion and Redundancy | 410 |
| 3 | Frequent Pattern Mining in Subspace Clustering | 412 |
| 3.1 | Subspace Cluster Search | 412 |
| 3.2 | Subspace Search | 414 |
| 3.3 | Redundancy in Subspace Clustering | 417 |
| 4 | Conclusions | 419 |
| | References | 419 |
| 17 | Supervised Pattern Mining and Applications to Classification | 425 |
| | Albrecht Zimmermann and Siegfried Nijssen | |
| 1 | Introduction | 425 |
| 2 | Supervised Pattern Mining | 427 |
| 2.1 | Explicit Class Labels | 428 |
| 2.2 | Classes as Data Subsets | 428 |
| 2.3 | Numerical Target Values | 431 |

| | | |
|-----------|--|------------|
| 3 | Supervised Pattern Set Mining | 432 |
| 3.1 | Local Evaluation, Local Modification | 434 |
| 3.2 | Global Evaluation, Global Modification | 435 |
| 3.3 | Local Evaluation, Global Modification | 436 |
| 3.4 | Data Instance-Based Selection | 437 |
| 4 | Classifier Construction | 437 |
| 4.1 | Direct Classification | 437 |
| 4.2 | Indirect Classification | 438 |
| 5 | Summary | 439 |
| | References | 440 |
| 18 | Applications of Frequent Pattern Mining | 443 |
| | Charu C. Aggarwal | |
| 1 | Introduction | 443 |
| 2 | Frequent Patterns for Customer Analysis | 445 |
| 3 | Frequent Patterns for Clustering | 446 |
| 4 | Frequent Patterns for Classification | 447 |
| 5 | Frequent Patterns for Outlier Analysis | 449 |
| 6 | Frequent Patterns for Indexing | 450 |
| 7 | Web Mining Applications | 451 |
| 7.1 | Web Log Mining | 451 |
| 7.2 | Web Linkage Mining | 452 |
| 8 | Frequent Patterns for Text Mining | 452 |
| 9 | Temporal Applications | 453 |
| 10 | Spatial and Spatiotemporal Applications | 455 |
| 11 | Software Bug Detection | 456 |
| 12 | Chemical and Biological Applications | 457 |
| 12.1 | Chemical Applications | 458 |
| 12.2 | Biological Applications | 458 |
| 13 | Resources for the Practitioner | 460 |
| 14 | Conclusions and Summary | 461 |
| | References | 461 |
| | Index | 469 |

Contributors

Charu C. Aggarwal IBM T. J. Watson Research Center, Yorktown Heights, NY, USA

Gagan Agrawal Ohio State University, Columbus, OH, USA

Luiza Antonie University of Guelph, Guelph, Canada

David C. Anastasiu Department of Computer Science and Engineering, University of Minnesota, Minneapolis, USA

Ira Assent Department of Computer Science, Aarhus University, Aarhus, Denmark

Mansurul A. Bhuiyan Indiana University–Purdue University, Indianapolis, IN, USA

Hong Cheng Department of Systems Engineering and Engineering Management, The Chinese University of Hong Kong, Hong Kong, China

Aris Gkoulalas-Divanis IBM Research-Ireland, Damastown Industrial Estate, Mulhuddart, Dublin, Ireland

Jiawei Han University of Illinois at Urbana-Champaign, Urbana, IL, USA

Department of Computer Science, University of Illinois at Urbana-Champaign, Champaign, USA

Jayant Haritsa Database Systems Lab, Indian Institute of Science (IISc), Bangalore, India

Mohammad Al Hasan Indiana University–Purdue University, Indianapolis, IN, USA

Jeremy Iverson Department of Computer Science and Engineering, University of Minnesota, Minneapolis, USA

Ruoming Jin Kent State University, Kent, OH, USA

Murat Kantarcioglu UTD Data Security and Privacy Lab, University of Texas at Dallas, Texas, USA

Victor E. Lee John Carroll University, University Heights, OH, USA

Matthijs van Leeuwen KU Leuven, Leuven, Belgium

Carson Kai-Sang Leung University of Manitoba, Winnipeg, MB, Canada

Jundong Li University of Alberta, Alberta, Canada

Zhenhui Li Pennsylvania State University, University Park, USA

George Karypis Department of Computer Science and Engineering, University of Minnesota, Minneapolis, USA

Siegfried Nijssen KU Leuven, Leuven, Belgium

Universiteit Leiden, Leiden, The Netherlands

Jian Pei Simon Fraser University, Burnaby, BC, Canada

Shaden Smith Department of Computer Science and Engineering, University of Minnesota, Minneapolis, USA

Wei Shen Tsinghua University, Beijing, China

Nikolaj Tatti HIIT, Department of Information and Computer Science, Aalto University, Helsinki, Finland

Jilles Vreeken Max-Planck Institute for Informatics and Saarland University, Saarbrücken, Germany

Jianyong Wang Tsinghua University, Beijing, China

Xifeng Yan Department of Computer Science, University of California at Santa Barbara, Santa Barbara, USA

Osmar Zaiane University of Alberta, Alberta, Canada

Feida Zhu Singapore Management University, Singapore, Singapore

Albrecht Zimmermann INSA Lyon, Villeurbanne CEDEX, France

Arthur Zimek Ludwig-Maximilians-Universität München, Munich, Germany

Chapter 1

An Introduction to Frequent Pattern Mining

Charu C. Aggarwal

Abstract The problem of frequent pattern mining has been widely studied in the literature because of its numerous applications to a variety of data mining problems such as clustering and classification. In addition, frequent pattern mining also has numerous applications in diverse domains such as spatiotemporal data, software bug detection, and biological data. The algorithmic aspects of frequent pattern mining have been explored very widely. This chapter provides an overview of these methods, as it relates to the organization of this book.

Keywords Frequent pattern mining · Association rules

1 Introduction

The problem of frequent pattern mining is that of finding relationships among the items in a database. The problem can be stated as follows.

Given a database \mathcal{D} with transactions $T_1 \dots T_N$, determine all patterns P that are present in at least a fraction s of the transactions.

The fraction s is referred to as the *minimum support*. The parameter s can be expressed either as an absolute number, or as a fraction of the total number of transactions in the database. Each transaction T_i can be considered a sparse binary vector, or as a set of discrete values representing the identifiers of the binary attributes that are instantiated to the value of 1. The problem was originally proposed in the context of market basket data in order to find frequent groups of items that are bought together [10]. Thus, in this scenario, each attribute corresponds to an item in a superstore, and the binary value represents whether or not it is present in the transaction. Because the problem was originally proposed, it has been applied to numerous other applications in the context of data mining, Web log mining, sequential pattern mining, and software bug analysis.

In the original model of frequent pattern mining [10], the problem of finding *association rules* has also been proposed which is closely related to that of frequent

C. C. Aggarwal (✉)
IBM T. J. Watson Research Center, Yorktown Heights, NY 10598, USA
e-mail: charu@us.ibm.com

patterns. In general association rules can be considered a “second-stage” output, which are derived from frequent patterns. Consider the sets of items U and V . The rule $U \Rightarrow V$ is considered an *association rule* at minimum support s and minimum confidence c , when the following two conditions hold true:

1. The set $U \cup V$ is a frequent pattern.
2. The ratio of the support of $U \cup V$ to that of U is at least c .

The minimum confidence c is always a fraction less than 1 because the support of the set $U \cup V$ is always less than that of U . Because the first step of finding frequent patterns is usually the computationally more challenging one, most of the research in this area is focussed on the former. Nevertheless, some computational and modeling issues also arise during the second step, especially when the frequent pattern mining problem is used in the context of other data mining problems such as classification. Therefore, this book will also discuss various aspects of association rule mining along with that of frequent pattern mining.

A related problem is that of *sequential pattern mining* in which an order is present in the transactions [5]. Temporal order is quite natural in many scenarios such as customer buying behavior, because the items are bought at specific time stamps, and often follow a natural temporal order. In these cases, the problem is redefined to that of sequential pattern mining, in which it is desirable to determine relevant and frequent *sequences* of items.

Some examples of important applications are as follows;

- *Customer Transaction Analysis*: In this case, the transactions represent sets of items that co-occur in customer buying behavior. In this case, it is desirable to determine frequent patterns of buying behavior, because they can be used for making decision about shelf stocking or recommendations.
- *Other Data Mining Problems*: Frequent pattern mining can be used to enable other major data mining problems such as classification, clustering and outlier analysis [11, 52, 73]. This is because the use of frequent patterns is so fundamental in the analytical process for a host of data mining problems.
- *Web Mining*: In this case, the Web logs may be processed in order to determine important patterns in the browsing behavior [24, 63]. This information can be used for Web site design, recommendations, or even outlier analysis.
- *Software Bug Analysis*: Executions of software programs can be represented as graphs with typical patterns. Logical errors in these bugs often show up as specific kinds of patterns that can be mined for further analysis [41, 51].
- *Chemical and Biological Analysis*: Chemical and biological data are often represented as graphs and sequences. A number of methods have been proposed in the literature for using the frequent patterns in such graphs for a wide variety of applications in different scenarios [8, 29, 41, 42, 69–75].

Since the publication of the original article on frequent pattern mining [10], numerous techniques have been proposed both for frequent and sequential pattern mining [5, 4, 13, 33, 62]. Furthermore, many variants of frequent pattern mining, such as

sequential pattern mining, constrained pattern mining, and graph mining have been proposed in the literature.

Frequent pattern mining is a rather broad area of research, and it relates to a wide variety of topics at least from an application specific-perspective. Broadly speaking, the research in the area falls in one of four different categories:

- **Technique-centered:** This area relates to the determination of more *efficient* algorithms for frequent pattern mining. A wide variety of algorithms have been proposed in this context that use different enumeration tree exploration strategies, and different data representation methods. In addition, numerous variations such as the determination of compressed patterns of great interest to researchers in data mining.
- **Scalability issues:** The scalability issues in frequent pattern mining are very significant. When the data arrives in the form of a stream, multi-pass methods can no longer be used. When the data is distributed or very large, then parallel or big-data frameworks must be used. These scenarios necessitate different types of algorithms.
- **Advanced data types:** Numerous variations of frequent pattern mining have been proposed for advanced data types. These variations have been utilized in a wide variety of tasks. In addition, different data domains such as graph data, tree structured data, and streaming data often require specialized algorithms for frequent pattern mining. Issues of interestingness of the patterns are also quite relevant in this context [6].
- **Applications:** Frequent pattern mining have numerous applications to other major data mining problems, Web applications, software bug analysis, and chemical and biological applications. A significant amount of research has been devoted to applications because these are particularly important in the context of frequent pattern mining.

This book will cover all these different areas comprehensively, so as to provide a comprehensive overview of this broader area.

This chapter is organized as follows. The next section discusses algorithms for the frequent pattern mining problem, and its basic variations. Section 3 discusses scalability issues for frequent pattern mining. Frequent pattern mining methods are advanced data types are discussed in Sect. 4. Privacy issues of frequent pattern mining are addressed in Sect. 5. The applications are discussed in Sect. 6. Section 7 gives the conclusions and summary.

2 Frequent Pattern Mining Algorithms

Most of the algorithms for frequent pattern mining have been designed with the traditional support-confidence framework, or for specialized frameworks that generate

more interesting kinds of patterns. These specialized framework may use different types of interestingness measures, model negative rules, or use constraint-based frameworks to determine more relevant patterns.

2.1 Frequent Pattern Mining with the Traditional Support Framework

The support framework is designed to determine patterns for which the raw frequency is greater than a minimum threshold. Although this is a simplistic way of defining frequent patterns, this model has an algorithmically convenient property, which is referred to as the *level-wise* property. The level-wise property of frequent pattern mining is algorithmically crucial because it enables the design of a bottom-up approach to exploring the space of frequent patterns. In other words, a $(k + 1)$ -pattern may not be frequent when any of its subsets is not frequent. This is a crucial observation that is used by virtually all the efficient frequent pattern mining algorithms.

Since the problem of frequent pattern mining was first proposed, numerous algorithms have been proposed in order to make the solutions to the problem more efficient. This area of research is so popular that an annual workshop *FIMI* was devoted to implementations of frequent pattern mining for a few years. This site [77] is now organized as a repository, where many efficient implementations of frequent pattern mining are available. The techniques for frequent pattern mining started with *Apriori*-like join-based methods. In these algorithms, candidate itemsets are generated in increasing order of itemset size. The generation in increasing order of itemset size is referred to as *level-wise exploration*. These itemsets are then tested against the underlying transaction database and the frequent ones satisfying the minimum support constraint are retained for further exploration. Eventually, it was realized that these *Apriori*-like methods could be more systematically explored as *enumeration trees*. This structure will be explained in detail in Chap. 2, and provides a methodology to perform systematic and non-redundant frequent pattern exploration. The enumeration tree provides a more flexible framework for frequent itemset mining because the tree can be explored in a variety of different strategies such as depth-first, breadth-first, or other hybrid strategies [13]. One property of the breadth-first strategy is that level-wise pruning can be used, which is not possible with other strategies. Nevertheless, strategies such as depth-first search have other advantages, especially for maximal pattern mining. This observation for the case of maximal pattern mining was first stated in [12]. This is because long patterns are discovered early, and they can be used for downward closure-based pruning of large parts of the enumeration tree that are already known to be frequent. It should be pointed out, that for the case where *all* frequent patterns are mined, the order of exploration of an enumeration tree does not affect the number of candidates that are explored because the size of the enumeration tree is fixed.

Join-based algorithms are always level-wise, and can be viewed as equivalent to breadth-first enumeration tree exploration. The algorithm proposed in the first frequent pattern mining paper [10] was an enumeration-tree based algorithm, whereas the second algorithm proposed was referred to as *Apriori*, and was a join-based algorithm [4]. Both algorithms are level-wise algorithms. Subsequently, many algorithms have been proposed in order to improve the implementations based on the enumeration tree paradigm with the use of techniques such as lookahead [17], depth-first search [12, 13, 33] and vertical exploration [62]. Some of these methods such as *TreeProjection*, *DepthProject* and *FP-growth* [33] use a projection strategy in which smaller transaction databases are explored at lower levels of the tree.

One of the challenges of frequent pattern mining is that a large number of redundant patterns are often mined. For example, the subset of a frequent pattern is also guaranteed to be frequent and by mining a maximal itemset, one is assured that the other frequent patterns can also be generated from this smaller set. Therefore, one possibility is to mine for only *maximal* itemsets [17]. However, the mining of maximal itemsets loses information about the exact value of support of the subsets of maximal patterns. Therefore, a further refinement would be to find *closed* frequent itemsets [58, 74]. Closed frequent itemsets are defined as frequent patterns, no superset of which have the same frequency as that itemset. By mining closed frequent itemsets, it is possible to significantly reduce the number of patterns found, without losing any information about the support level. Closed patterns can be viewed as the maximal patterns from each group of *equi-support* patterns (i.e., patterns with the same support). All maximal patterns are, therefore, closed.

The depth-first method has been shown to have a number of advantages in maximal pattern mining [12], because of the greater effectiveness of the pruning-based lookaheads in the depth-first strategy. Different techniques for frequent pattern mining will be discussed in Chaps. 2 and 3. The former chapter will generally focus on frequent pattern mining algorithms, whereas the latter chapter will focus on pattern-growth algorithms. An additional chapter with greater detail has been devoted to pattern-growth methods, because of it is considered a state-of-the-art technique in frequent pattern mining. The efficiency in frequent pattern mining algorithms can be gained in several ways:

1. Reducing the size of the candidate search space, with the use of pruning methods, such as maximality pruning. The notion of *closure* can also be used to prune large parts of the search space. However, these methods often do not exhaustively return the full set of frequent patterns. Many of these methods returned condensed representations such as maximal patterns or closed patterns.
2. Improving the efficiency of *counting*, with the use of database projection. Methods such as *TreeProjection* speed up the rate at which each pattern is counted, by reducing the size of the database with respect to which patterns are compared.
3. Using more efficient data structures, such as vertical lists, or an FP-Tree for more compressed database representation. In frequent pattern mining, both memory and computational speeds can be improved by judicious choice of data structures.

A particular scenario of interest is one in which the patterns to be mined are very long. In such cases, the number of subsets of frequent patterns can be extremely large. Therefore, a number of techniques need to be designed in order to mine very long patterns. In such cases, a variety of methods are used to explore the long patterns early, so that their subsets can be pruned effectively. The scenario of long pattern generation is discussed in detail in Chap. 4, though it is also discussed to some extent in the earlier Chaps. 2 and 3.

2.2 Interesting and Negative Frequent Patterns

A major challenge in frequent pattern mining is that the rules found may often not be very interesting, when quantifications such as support and confidence are used. This is because such quantifications do not normalize for the original frequency of the underlying items. For example, an item that occurs very rarely in the underlying database would naturally also occur in itemsets with lower frequency. Therefore, the *absolute* frequency often does not tell us much about the likelihood of items to *co-occur* together, because of the biases associated with the frequencies of the individual items. Therefore, numerous methods have been proposed in the literature for finding interesting frequent patterns that normalize for the underlying item frequencies [6, 26]. Methods for finding interesting frequent patterns are discussed in Chap. 5. The issue of interestingness is also related to compressed representations of patterns such as closed or maximal itemsets. These issues are also discussed in the chapter.

In negative associative rule mining, we attempt to determine rules such as $Bread \Rightarrow \neg Butter$, where the symbol \neg indicates negation. Therefore, in this case $\neg Butter$ becomes a pseudo-item denoting a “negative item.” One possibility is to add negative items to the data, and perform the mining in the same way as one would determine rules in the support-confidence framework. However, this is not a feasible solution. This is because traditional support frameworks are not designed for cases where an item is presented in the data 98 % of the time. This is the case for “negative items.” For example, most transactions may not contain the item *Butter*, and therefore even positively correlated items may appear as negative rules. For example, the rule $Bread \Rightarrow \neg Butter$ may have confidence greater than 50 %, even though *Bread* is clearly correlated in a positive way with *Butter*. This is because, the item $\neg Butter$ may have an even higher support of 98 %.

The issue of finding negative patterns is closely related to that of finding interesting patterns in the data [6] because one is looking for patterns that satisfy the support requirement in an interesting way. This relationship between the two problems tends to be under-emphasized in the literature, and the problem of negative pattern mining is often treated independently from interesting pattern mining. Some frameworks, such as collective strength, are designed to address both issues simultaneously. Methods for negative pattern mining are addressed in Chap. 6. The relationship between interesting pattern mining and negative pattern mining will be discussed in the same chapter.

2.3 *Constrained Frequent Pattern Mining*

Off-the-shelf frequent pattern mining algorithms discover a large number of patterns which are not useful when it is desired to determine patterns on the basis of more refined criteria. Frequent pattern mining methods are often particularly useful in the context of constrained applications, in which rules satisfying particular criteria are discovered. For example, one may desire specific items to be present in the rule. One solution is to first mine all the itemsets, and then enable online mining from this set of base patterns [3]. However, pushing constraints directly into the mining process has several advantages. This is because when constraints are pushed directly into the mining process, the mining can be performed at much lower support levels than can be performed by using a two-phase approach. This is especially the case when a large number of intermediate candidates can be pruned by the constraint-based pattern mining algorithm.

A variety of arbitrary constraints may also be present in the patterns. The major problem with such methods is that the constraints may result in the violation of the downward closure property. Because most frequent pattern mining algorithms depend crucially on this property, its violation is a serious issue. Nevertheless, many constraints have specialized properties because of which specialized algorithms can be developed. Methods for constrained frequent pattern mining method have been discussed in [55, 57, 60]. Constrained methods have also been developed for the sequential pattern mining problem [31, 61]. In real applications, the output of the vanilla frequent pattern mining problem may be too large, and it is only by pushing constraints into the pattern mining process, that useful application-specific patterns can be found. Constrained frequent pattern mining methods are closely related to the problem of pattern-based classification, because the latter problem requires us to discover discriminative patterns from the underlying data. Methods for constrained frequent pattern mining will be discussed in Chap. 2.

2.4 *Compressed Representations of Frequent Patterns*

A major problem in frequent pattern mining algorithms is that the volume of the mining patterns is often extremely large. This scenario creates numerous challenges for using these patterns in a meaningful way. Furthermore, different kinds of redundancy are present in the mined patterns. For example, maximal patterns imply the presence of all their subsets in the data. There is some information loss in terms of the exact support values of these subsets. Therefore, if it is not needed to preserve the values of the support across the patterns, then the determination of concise representations can be very useful.

A particularly interesting form of concise representation is that of *closed patterns* [56]. An itemset X is set to be closed if none of its supersets have the same support as X . Therefore, by determining all the closed frequent patterns, one can derive not only the exhaustive set of frequent itemsets, but also their supports. Note that

support values are lost by maximal pattern mining. In other words, the set of maximal patterns cannot be used to derive the support values of missing subsets. However, the support values of closed frequent itemsets can be used to derive the support values of missing subsets. Many interesting methods [58, 67, 74] have been designed for identifying frequent closed patterns. The general principle of determining frequent closed patterns has been generalized to that of determining δ -freesets [18]. This issue is closely related to that of mining all *non-derivable frequent itemsets* [20]. A survey on this topic may be found in [21]. These different forms of compression are discussed in Chaps. 2 and 5.

Finally, a formal way of viewing compression is from the perspective of information-theoretic models. Information-theoretic models are designed for compressing different kinds of data, and can therefore be used to compress itemsets as well. This basic principle has been used for methods such as *Krimp* [66]. The problem of determining compressed representations of frequent itemsets is discussed in Chap. 8. This chapter focusses mostly on the information-theoretic issues of frequent itemset compression.

3 Scalability Issues in Frequent Pattern Mining

In the modern era, the ability to collect large amounts of data has increased significantly because of advances in hardware and software platforms. The amount of data is often so large that specialized methods are required for the mining process. The streaming and big-data architectures are slightly different and pose different challenges for the mining process. The following discussion will address each of these challenges.

3.1 Frequent Pattern Mining in Data Streams

In recent years, data stream have become very popular because of the advances in hardware and software technology that can collect and transmit data continuously over time. In such cases, the major constraint on data mining algorithms is to execute the algorithms in a single pass. This can be significantly challenging because frequent and sequential pattern mining methods are generally designed as level-wise methods. There are two variants of frequent pattern mining for data streams:

- *Frequent Items or Heavy Hitters*: In this case, frequent 1-itemsets need to be determined from a data stream in a single pass. Such an approach is generally needed when the total number of distinct items is too large to be held in main memory. Typically, sketch-based methods are used in order to create a compress data structure in order to maintain *approximate* counts of the items [23, 27].
- *Frequent itemsets*: In this case, it is not assumed that the number of distinct items are too large. Therefore, the main challenge in this case is computational, because

the typical frequent pattern mining methods are multi-pass methods. Multiple passes are clearly not possible in the context of data streams [22, 39].

The streaming scenario also presents numerous challenges in the context of data of advanced types. For example, graph streams are often encountered in the context of network data. In such cases, methods need to be designed for determining dense groups of nodes in real time [16]. Methods for mining frequent items and itemsets in data streams are discussed in Chap. 9.

3.2 Frequent Pattern Mining with Big Data

The big data scenario poses numerous challenges for the problem of frequent pattern mining. A major problem arises when the data is large enough to be stored in a distributed way. Therefore, significant costs are incurred in shuffling around data or intermediate results of the mining process across the distributed nodes. These costs are also referred to as data transfer costs. When data sets are very large, then the algorithms need to be designed to take into account both the disk access constraint and the data transfer costs. In addition, many distributed frameworks such as *MapReduce* [28] require specialized algorithms for frequent pattern mining. The focus of big-data framework is somewhat different from streams, in that it is closely related to the issue of shuffling large amounts of data around for the mining process. Interestingly, it is sometimes easier to process the algorithms in a single pass in streaming fashion, than when they have already been stored in distributed frameworks where access costs become a major issue. Algorithms for frequent pattern mining with big data are discussed in detail in Chap. 10. This chapter discusses both the parallel algorithms and the big-data algorithms that are based on the *MapReduce* framework.

4 Frequent Pattern Mining with Advanced Data Types

although the frequent pattern mining problem is naturally defined on sets, it can be extended to various advanced data types. The most natural extension of frequent pattern mining algorithms is to the case of temporal data. This was one of the earliest proposed extensions and is referred to as *sequential pattern mining*. Subsequently, the problem has been generalized to other advanced data types, such as spatiotemporal data, graphs, and uncertain data. Many of the developed algorithms are basic variations of the frequent pattern mining problem. In general, the basic frequent pattern mining algorithms need to be modified carefully to address the variations required by the advanced data types.

4.1 Sequential Pattern Mining

The problem of sequential pattern mining is closely related to that of frequent pattern mining. The major difference in this case is that record contain *baskets of items* arranged sequential. For example, each record R_i may be of the following form:

$$R_i = \langle \{Bread\}, \{Butter, Cake\}, \{Chicken, Yogurt\} \rangle$$

In this case, each entity within $\{\}$ is a basket of items that are bought together and, therefore, do not have a temporal ordering. This basket of items is collectively referred to as an *event*. The length of a pattern is equal to the sum of the lengths of the complex items in it. For example, R_i is a 5-pattern, even though it has 3 events. The different complex entities (or events) do have a temporal ordering. In the aforementioned example, it is clear that $\{Bread\}$ has been bought earlier than $\{Butter, Cake\}$. The problem of sequential pattern mining is that of finding sequences of events that are present in at least a fraction s of the underlying records [5]. For example, the sequence $\langle \{Bread\}, \{Butter\}, \{Chicken\} \rangle$ is present in the afore-mentioned record, but not the sequence $\langle \{Bread\}, \{Cake\}, \{Butter\} \rangle$. The pattern may also contain complex events. For example, the pattern $\langle \{Bread\}, \{Chicken, Yogurt\} \rangle$ is present in R_i . The problem of sequential pattern mining is closely related to that of frequent pattern mining except that it is somewhat more complex to account for both the presence of complex baskets of items in the database, and the temporal ordering of the individual baskets. An extension of a sequential pattern may either be a set-wise extension of a complex item, or a temporal extension with an entirely new event. This affects the nature of the extensions of items in the transactions. Numerous modifications of known frequent pattern mining methods such as *Apriori* and its variants, *TreeProjection* and its variants [32], and the *FP-growth* method and its variants, can be used in order to solve the sequential pattern mining problem [5, 35, 36]. The enumeration tree concept can also be generalized to sequential pattern mining [32]. Therefore, in principle, all enumeration tree algorithms can be generalized to sequential pattern mining. This is a powerful ability because, as we will see in Chap. 2 all frequent pattern mining algorithms are, implicitly or explicitly, enumeration-tree algorithms. Sequential pattern mining methods will be discussed in detail in Chap. 11.

4.2 Spatiotemporal Pattern Mining

The advent of GPS-enabled mobile phones and wearable sensors has enabled the collection of large amounts of spatiotemporal data. Such data may include trajectory data, location-tagged images, or other content. In some cases, the spatiotemporal data exists in the form of RFID data [37]. The mining of patterns from such spatiotemporal data provides numerous insights in a wide variety of applications, such as traffic control and social sensing [2]. Frequent patterns are also used for trajectory

clustering classification and outlier analysis [38, 45–48]. Many trajectory analysis problems can be approximately transformed to sequential pattern mining with the use of appropriate transformations. Algorithms for spatiotemporal pattern mining are discussed in Chap. 12.

4.3 Frequent Patterns in Graphs and Structured Data

Many kinds of chemical and biological data, XML data, software program traces, and Web browsing behaviors can be represented as structured graphs. In these cases, frequent pattern mining is very useful for making inferences in such data. This is because frequent structural patterns provide important insights about the graphs. For example, specific chemical structures result in particular properties, specific program structures result in software bugs, and so on. Such patterns can even be used for clustering and classification of graphs! [14, 73].

A variety of methods for structural frequent pattern mining are discussed in [41, 69–71, 72]. A major problem in the context of graphs is the problem of *isomorphism*, because of which there are multiple ways to match two graphs. An *Apriori*-like algorithm can be developed for graph pattern mining. However, because of the complexity of graphs and also because of issues related to isomorphism, the algorithms are more complex. For example, in an *Apriori*-like algorithm, pairs of graphs can be joined in multiple ways. Pairs of graphs can be joined when they have $(k - 1)$ nodes in common, or they have $(k - 1)$ edges in common. Furthermore, either kind of join between a pair of graphs can have multiple results. The counting process is also more challenging because of isomorphism. Pattern mining in graphs becomes especially challenging when the graphs are large, and the isomorphism problem becomes significant. Another particularly difficult case is the streaming scenario [16] where one has to determine dense patterns in the graphs stream. Typically, these problems cannot be solved exactly, and approximations are required.

Frequent pattern mining in graphs has numerous applications. In some cases, these methods can be used in order to perform classification and clustering of structured data [14, 73]. Graph patterns are used for chemical and biological data analysis, and software bug detection in computer programs. Methods for finding frequent patterns in graphs are discussed in Chap. 13. The applications of graph pattern mining are discussed in Chap. 18.

4.4 Frequent Pattern Mining with Uncertain Data

Uncertain or probabilistic data has become increasingly common over the last few years, as methods have been designed in order to collect data with very low quality. The attribute values in such data sets are *probabilistic*, which implies that the values are represented as probability distributions. Numerous algorithms have been

proposed in the literature for uncertain frequent pattern mining [15], and a computational evaluation of the different techniques is provided in [64]. Many algorithms such as *FP-growth* are harder to generalize to uncertain data [15] because of the difficulty in storing probability information with the FP-Tree. Nevertheless, as the work in [15] shows, other related methods such as *H-mine* [59] can be generalized easily to the case of uncertain data. Uncertain frequent pattern mining methods have also been extended to the case of graph data [76]. A variant of uncertain graph pattern mining discovers highly *reliable* subgraphs [40]. Highly reliable subgraphs are subgraphs that are hard to disconnect in spite of the uncertainty associated with the edges. A discussion of the different methods for frequent pattern mining with uncertain data is provided in Chap. 14.

5 Privacy Issues

Privacy has increasingly become a topic of concern in recent years because of the wide availability of personal data about individuals [7]. This has often led to reluctance to share data, share it in a constrained way, or share downgraded versions of the data. The additional constraints and downgrading translate to challenges in discovering frequent patterns. In the context of frequent pattern and association rule mining, the primary challenges are as follows:

1. When privacy-preservation methods such as randomization are used, it becomes a challenge to discover associations from the underlying data. This is because a significant amount of noise has been added to the data, and it is often difficult to discover the association rules in the presence of this noise. Therefore, one class of association rule mining methods [30] proposes effective methods to perturb the data, so that meaningful patterns may be discovered while retaining privacy of the perturbed data.
2. In some cases, the output of a privacy-preserving data mining algorithm can lead to violation of privacy. This is because association rules can reveal sensitive information about individuals when they relate sensitive attributes to other kinds of attributes. Therefore, one class of methods focusses on the problem of *association rule hiding* [65].
3. In many cases, the data to be mined is stored in a distributed way by competitors who may wish to determine global insights without, at the same time, revealing their local insights. This problem is referred to as that of distributed privacy preservation [25]. The data may be either horizontally partitioned across rows (different records) or vertically partitioned (across attributes). Each of these forms of partitioning require different methods for distributed mining.

Methods for privacy-preserving association rule mining are addressed in Chap. 15.

6 Applications of Frequent Pattern Mining

Frequent pattern mining has applications of two types. The first type of application is to other major data mining problems such as clustering, outlier detection, and classification. Frequent patterns are often used to determine relevant clusters from the underlying data. In addition, rule-based classifiers are often constructed with the use of frequent pattern mining methods. Frequent pattern mining is also used in generic applications, such as Web log analytics, software bug analysis, chemical, and biological data.

6.1 Applications to Major Data Mining Problems

Frequent pattern mining methods can also be applied to other major data mining problems such as clustering [9, 19], classification and outlier analysis. For example, frequent pattern mining methods are often used for subspace clustering [11], by discretizing the quantitative attributes, and then finding patterns from these discrete values. Each such pattern, therefore, corresponds to a rectangular region in a subspace of the data. These rectangular regions can then be integrated together in order to create a more comprehensive subspace representation.

Frequent pattern mining is also applied to problems such as classification, in which rules are generated by using patterns on the left hand side of the rule, and the class variable on the right hand side of the rule [52]. The main goal here is to find *discriminative* patterns for the purpose of classification, rather than simply patterns that satisfy the support requirements. Such methods have also been extended to structured XML data [73] by finding discriminative graph-structured patterns. In addition, sequential pattern mining methods can be applied to other temporal mining methods such as event detection [43, 44, 53, 54] and sequence classification [68]. Frequent pattern mining has also been applied to the problem of outlier analysis [1], by determining deviations from the expected patterns in the underlying data. Methods for clustering based on frequent pattern mining are discussed in Chap. 16, while rule-based classification are discussed in Chap. 17. It should be pointed out that constrained frequent pattern mining is closely related to the problem of classification with frequent patterns, and therefore both are discussed in the same chapter.

6.2 Generic Applications

Frequent pattern mining has applications to a variety of problems such as clustering, classification and event detection. In addition, specific application areas such as Web mining and software bug detection can also benefit from frequent pattern mining methods. In the context of Web mining, numerous methods have been proposed for finding useful patterns from Web logs in order to make recommendations [63]. Such

techniques can also be used to determine outliers from Web log sequences [1]. Frequent patterns are also used for trajectory classification and outlier analysis [49–48]. Frequent pattern mining methods can also be used in order to determine relevant rules and patterns in spatial data, as they related to spatial and non-spatial properties of objects. For example, an association rule could be created from the relationships of land temperatures of “nearby” geographical locations. In the context of spatiotemporal data, the relationships between the motions of different objects could be used to create spatiotemporal frequent patterns. Frequent pattern mining methods have been used for finding patterns in biological and chemical data [42, 29, 75]. In addition, because software programs can be represented as graphs, frequent pattern mining methods can be used in order to find logical bugs from program execution traces [51]. Numerous applications of frequent pattern mining are discussed in Chap. 18.

7 Conclusions and Summary

Frequent pattern mining is one of four major problems in the data mining domain. This chapter provides an overview of the major topics in frequent pattern mining. The earliest work in this area was focussed on determining the efficient algorithms for frequent pattern mining, and variants such as long pattern mining, interesting pattern mining, constraint-based pattern mining, and compression. In recent years scalability has become an issue because of the massive amounts of data that continue to be created in various applications. In addition, because of advances in data collection technology, advanced data types such as temporal data, spatiotemporal data, graph data, and uncertain data have become more common. Such data types have numerous applications to other data mining problems such as clustering and classification. In addition, such data types are used quite often in various temporal applications, such as the Web log analytics.

References

1. C. Aggarwal. Outlier Analysis, *Springer*, 2013.
2. C. Aggarwal. Social Sensing, *Managing and Mining Sensor Data*, Springer, 2013.
3. C. C. Aggarwal, and P. S. Yu. Online generation of Association Rules, *ICDE Conference*, 1998.
4. R. Agrawal, and R. Srikant. Fast Algorithms for Mining Association Rules in Large Databases, *VLDB Conference*, pp. 487–499, 1994.
5. R. Agrawal, and R. Srikant. Mining Sequential Patterns, *ICDE Conference*, 1995.
6. C. C. Aggarwal, and P. S. Yu. A New Framework for Itemset Generation, *ACM PODS Conference*, 1998.
7. C. Aggarwal and P. Yu. Privacy-preserving data mining: Models and Algorithms, *Springer*, 2008.
8. C. C. Aggarwal, and H. Wang. Managing and Mining Graph Data, *Springer*, 2010.
9. C. C. Aggarwal, and C. K. Reddy. Data Clustering: Algorithms and Applications, *CRC Press*, 2013.

10. R. Agrawal, T. Imielinski, and A. Swami. Database Mining: A Performance Perspective. *IEEE Transactions on Knowledge and Data Engineering*, 5(6), pp. 914–925, 1993.
11. R. Agrawal, J. Gehrke, D. Gunopulos, P. Raghavan. Automatic Subspace Clustering of High Dimensional Data for Data Mining Applications, *ACM SIGMOD Conference*, 1998.
12. R. Agarwal, C. C. Aggarwal, and V. V. V. Prasad. Depth-first Generation of Long Patterns, *ACM KDD Conference*, 2000: Also appears as IBM Research Report, RC, 21538, 1999.
13. R. Agarwal, C. C. Aggarwal, and V. V. V. Prasad. A Tree Projection Algorithm for Generation of Frequent Itemsets, *Journal of Parallel and Distributed Computing*, 61(3), pp. 350–371, 2001. Also appears as IBM Research Report, RC 21341, 1999.
14. C. C. Aggarwal, N. Ta, J. Wang, J. Feng, M. Zaki. Xproj: A framework for projected structural clustering of XML documents, *ACM KDD Conference*, 2007.
15. C. C. Aggarwal, Y. Li, J. Wang, J. Feng. Frequent Pattern Mining with Uncertain Data, *ACM KDD Conference*, 2009.
16. C. Aggarwal, Y. Li, P. Yu, and R. Jin. On dense pattern mining in graph streams, *VLDB Conference*, 2010.
17. R. J. Bayardo Jr. Efficiently mining long patterns from databases. *ACM SIGMOD Conference*, 1998.
18. J.-F. Boulicaut, A. Bykowski, and C. Rigotti. Free-sets: A Condensed Representation of Boolean data for the Approximation of Frequency Queries. *Data Mining and Knowledge Discovery*, 7(1), pp. 5–22, 2003.
19. G. Buehrer, and K. Chellapilla. A Scalable Pattern Mining Approach to Web Graph Compression with Communities. *WSDM Conference*, 2009.
20. T. Calders, and B. Goethals. Mining all non-derivable frequent itemsets, *Principles of Knowledge Discovery and Data Mining*, 2006.
21. T. Calders, C. Rigotti, and J. F. Boulicaut. A survey on condensed representations for frequent sets. In *Constraint-based mining and inductive databases*, pp. 64–80, Springer, 2006.
22. J. H. Chang, W. S. Lee. Finding Recent Frequent Itemsets Adaptively over Online Data Streams. *ACM KDD Conference*, 2003.
23. M. Charikar, K. Chen, and M. Farach-Colton. Finding Frequent Items in Data Streams, *Automata, Languages and Programming*, pp. 693–703, 2002.
24. M. S. Chen, J. S. Park, and P. S. Yu. Efficient data mining for path traversal patterns, *IEEE Transactions on Knowledge and Data Engineering*, 10(2), pp. 209–221, 1998.
25. C. Clifton, M. Kantarcioglu, J. Vaidya, X. Lin, and M. Zhu. Tools for privacy preserving distributed data mining. *ACM SIGKDD Explorations Newsletter*, 4(2), pp. 28–34, 2002.
26. E. Cohen. M. Datar, S. Fujiwara, A. Gionis, P. Indyk, R. Motwani, J. Ullman, and C. Yang. Finding Interesting Associations without Support Pruning, *IEEE TKDE*, 13(1), pp. 64–78, 2001.
27. G. Cormode, S. Muthukrishnan. What’s hot and what’s not: tracking most frequent items dynamically, *ACM TODS*, 30(1), pp. 249–278, 2005.
28. J. Dean and S. Ghemawat. *MapReduce*: Simplified Data Processing on Large Clusters. *OSDI*, pp. 137–150, 2004.
29. M. Deshpande, M. Kuramochi, N. Wale, and G. Karypis. Frequent substructure-based approaches for classifying chemical compounds. *IEEE TKDE*, 17(8), pp. 1036–1050, 2005.
30. A. Evmimievski, R. Srikant, R. Agrawal, and J. Gehrke. Privacy preserving mining of association rules. *Information Systems*, 29(4), pp. 343–364, 2004.
31. M. Garofalakis, R. Rastogi, and K. Shim.: Sequential Pattern Mining with Regular Expression Constraints, *VLDB Conference*, 1999.
32. V. Guralnik, and G. Karypis. Parallel tree-projection-based sequence mining algorithms. *Parallel Computing*, 30(4): pp. 443–472, April 2004.
33. J. Han, J. Pei, and Y. Yin. Mining Frequent Patterns without Candidate Generation, *ACM SIGMOD Conference*, 2000.
34. J. Han, H. Cheng, D. Xin, and X. Yan. Frequent Pattern Mining: Current Status and Future Directions, *Data Mining and Knowledge Discovery*, 15(1), pp. 55–86, 2007.

35. J. Han, J. Pei, B. Mortazavi-Asl, Q. Chen, U. Dayal, and M. C. Hsu. FreeSpan: frequent pattern-projected sequential pattern mining. *ACM KDD Conference*, 2000.
36. J. Han, J. Pei, H. Pinto, B. Mortazavi-Asl, Q. Chen, U. Dayal, and M. C. Hsu. PrefixSpan: Mining sequential patterns efficiently by prefix-projected pattern growth. *ICDE Conference*, 2001.
37. J. Han, J.-G. Lee, H. Gonzalez, X. Li. Mining Massive RFID, Trajectory, and Traffic Data Sets (Tutorial). *ACM KDD Conference*, 2008. Video of Tutorial Lecture at: http://videlectures.net/kdd08_han_mmrfd/
38. H. Jeung, M. L. Yiu, X. Zhou, C. Jensen, H. Shen, Discovery of Convoys in Trajectory Databases, *VLDB Conference*, 2008.
39. R. Jin, G. Agrawal. Frequent Pattern Mining in Data Streams, *Data Streams: Models and Algorithms*, pp. 61–84, Springer, 2007.
40. R. Jin, L. Liu, and C. Aggarwal. Discovering highly reliable subgraphs in uncertain graphs. *ACM KDD Conference*, 2011.
41. G. Kuramuchi and G. Karypis. Frequent Subgraph Discovery, *ICDM Conference*, 2001.
42. A. R. Leach and V. J. Gillet. *An Introduction to Chemoinformatics*. Springer, 2003.
43. W. Lee, S. Stolfo, and P. Chan. Learning Patterns from Unix Execution Traces for Intrusion Detection, *AAAI workshop on AI methods in Fraud and Risk Management*, 1997.
44. W. Lee, S. Stolfo, and K. Mok. A Data Mining Framework for Building Intrusion Detection Models, *IEEE Symposium on Security and Privacy*, 1999.
45. J.-G. Lee, J. Han, K.-Y. Whang, Trajectory Clustering: A Partition-and-Group Framework, *ACM SIGMOD Conference*, 2007.
46. J.-G. Lee, J. Han, X. Li. Trajectory Outlier Detection: A Partition-and-Detect Framework, *ICDE Conference*, 2008.
47. J.-G. Lee, J. Han, X. Li, H. Gonzalez. TraClass: trajectory classification using hierarchical region-based and trajectory-based clustering. *PVLDB*, 1(1): pp. 1081–1094, 2008.
48. X. Li, J. Han, and S. Kim. Motion-alert: Automatic Anomaly Detection in Massive Moving Objects, *IEEE Conference in Intelligence and Security Informatics*, 2006.
49. X. Li, J. Han, S. Kim and H. Gonzalez. ROAM: Rule- and Motif-based Anomaly Detection in Massive Moving Object Data Sets, *SDM Conference*, 2007.
50. Z. Li, B. Ding, J. Han, R. Kays. Swarm: Mining Relaxed Temporal Object Moving Clusters, *VLDB Conference*, 2010.
51. C. Liu, X. Yan, H. Lu, J. Han, and P. S. Yu. Mining Behavior Graphs for “backtrace” of non-crashing bugs, *SDM Conference*, 2005.
52. B. Liu, W. Hsu, Y. Ma. Integrating Classification and Association Rule Mining, *ACM KDD Conference*, 1998.
53. S. Ma, and J. Hellerstein. Mining Partially Periodic Event Patterns with Unknown Periods, *IEEE International Conference on Data Engineering*, 2001.
54. H. Mannila, H. Toivonen, and A. I. Verkamo. Discovering Frequent Episodes in Sequences, *ACM KDD Conference*, 1995.
55. R. Ng, L. V. S. Lakshmanan, J. Han, and A. Pang. Exploratory mining and pruning optimizations of constrained associations rules. *ACM SIGMOD Conference*, 1998.
56. N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Discovering frequent closed itemsets for association rules. *International Conference on Database Theory*, pp. 398–416, 1999.
57. J. Pei, and J. Han. Can we push more constraints into frequent pattern mining? *ACM KDD Conference*, 2000.
58. J. Pei, J. Han, R. Mao. CLOSET: An Efficient Algorithms for Mining Frequent Closed Itemsets, *DMKD Workshop*, 2000.
59. J. Pei, J. Han, H. Lu, S. Nishio, S. Tang, and D. Yang. H-mine: Hyper-structure mining of frequent patterns in large databases. In *Data Mining, ICDM Conference*, 2001.
60. J. Pei, J. Han, and L. V. S. Lakshmanan. Mining Frequent Patterns with Convertible Constraints in Large Databases, *ICDE Conference*, 2001.

61. J. Pei, J. Han, and W. Wang. Constraint-based Sequential Pattern Mining: The Pattern-Growth Methods, *Journal of Intelligent Information Systems*, 28(2), pp. 133–160, 2007.
62. P. Shenoy, J. Haritsa, S. Sudarshan, G. Bhalotia, M. Bawa, D. Shah. Turbo-charging Vertical Mining of Large Databases. *ACM SIGMOD Conference*, pp. 22–33, 2000.
63. J. Srivastava, R. Cooley, M. Deshpande, and P. N. Tan. Web usage mining: Discovery and applications of usage patterns from Web data. *ACM SIGKDD Explorations Newsletter*, 1(2), pp. 12–23, 2000.
64. Y. Tong, L. Chen, Y. Cheng, P. Yu. Mining Frequent Itemsets over Uncertain Databases. *PVLDB*, 5(11), pp. 1650–1661, 2012.
65. V. S. Verykios, A. K. Elmagarmid, E. Bertino, Y. Saygin, and E. Dasseni. Association rule hiding. *IEEE Transactions on Knowledge and Data Engineering*, pp. 434–447, 16(4), pp. 434–447, 2004.
66. J. Vreeken, M. van Leeuwen, and A. Siebes. Krimp: Mining itemsets that compress. *Data Mining and Knowledge Discovery*, 23(1), pp. 169–214, 2011.
67. J. Wang, J. Han, and J. Pei. CLOSET+: Searching for the Best strategies for mining frequent closed itemsets. *ACM KDD Conference*, 2003.
68. Z. Xing, J. Pei, and E. Keogh. A Brief Survey on Sequence Classification, *ACM SIGKDD Explorations*, 12(1), 2010.
69. X. Yan, P. S. Yu, and J. Han. Graph indexing: A frequent structure-based approach. *ACM SIGMOD Conference*, 2004.
70. X. Yan, P. S. Yu, and J. Han. Substructure similarity search in graph databases. *ACM SIGMOD Conference*, 2005.
71. X. Yan, F. Zhu, J. Han, and P. S. Yu. Searching substructures with superimposed distance, *ICDE Conference*, 2006.
72. M. Zaki. Efficiently mining frequent trees in a forest: Algorithms and applications. *IEEE Transactions on Knowledge and Data Engineering*, 17(8), pp. 1021–1035, 2005.
73. M. Zaki, C. Aggarwal. XRules: An Effective Classifier for XML Data, *ACM KDD Conference*, 2003.
74. M. Zaki, C. J. Hsiao. CHARM: An Efficient Algorithm for Closed Frequent Itemset Mining, *SDM Conference*, 2002.
75. S. Zhang, T. Wang. Discovering Frequent Agreement Subtrees from Phylogenetic Data. *IEEE Transactions on Knowledge and Data Engineering*, 20(1), pp. 68–82, 2008.
76. Z. Zou, J. Li, H. Gao, and S. Zhang. Mining Frequent Subgraph Patterns from Uncertain Graph Data, *IEEE Transactions on Knowledge and Data Engineering*, 22(9), pp. 1203–1218, 2010.
77. <http://fimi.ua.ac.be/>

Chapter 2

Frequent Pattern Mining Algorithms: A Survey

Charu C. Aggarwal, Mansurul A. Bhuiyan and Mohammad Al Hasan

Abstract This chapter will provide a detailed survey of frequent pattern mining algorithms. A wide variety of algorithms will be covered starting from *Apriori*. Many algorithms such as *Eclat*, *TreeProjection*, and *FP-growth* will be discussed. In addition a discussion of several maximal and closed frequent pattern mining algorithms will be provided. Thus, this chapter will provide one of most detailed surveys of frequent pattern mining algorithms available in the literature.

Keywords Frequent pattern mining algorithms · *Apriori* · *TreeProjection* · *FP-growth*

1 Introduction

In data mining, frequent pattern mining (FPM) is one of the most intensively investigated problems in terms of computational and algorithmic development. Over the last two decades, numerous algorithms have been proposed to solve frequent pattern mining or some of its variants, and the interest in this problem still persists [45, 75]. Different frameworks have been defined for frequent pattern mining. The most common one is the support-based framework, in which itemsets with frequency above a given threshold are found. However, such itemsets may sometimes not represent interesting positive *correlations* between items because they do not normalize for the absolute frequencies of the items. Consequently, alternative measures for interestingness have been defined in the literature [7, 11, 16, 63]. This chapter will focus on the support-based framework because the algorithms based on the interestingness

C. C. Aggarwal (✉)

IBM T. J. Watson Research Center, Yorktown Heights, NY 10598, USA
e-mail: charu@us.ibm.com

M. A. Bhuiyan · M. A. Hasan

Indiana University–Purdue University, Indianapolis, IN, USA
e-mail: mbhuiyan@cs.iupui.edu

M. A. Hasan

e-mail: alhasan@cs.iupui.edu

Algorithm *Baseline Mining*(Database: \mathcal{T} , Minimum Support: s)

```

begin
   $\mathcal{FP} = \{\}$ ;
  Insert length-one frequent pattern in  $\mathcal{FP}$ 
  until all frequent patterns in  $\mathcal{FP}$  are explored do
    begin
      Generate a candidate pattern  $P$  from one (or more) frequent
        pattern(s) in  $\mathcal{FP}$ 
      if  $\text{support}(P, \mathcal{T}) \geq s$ 
        Add  $P$  to frequent pattern set  $\mathcal{FP}$ ;
      end
    end
  end

```

Fig. 2.1 A generic frequent pattern mining algorithm

framework are provided in a different chapter. Surveys on frequent pattern mining may be found in [26, 33].

One of the main reasons for the high level of interest in frequent pattern mining algorithms is due to the computational challenge of the task. Even for a moderate sized dataset, the search space of FPM is enormous, which is exponential to the length of the transactions in the dataset. This naturally creates challenges for itemset generation, when the support levels are low. In fact, in most practical scenarios, the support levels at which one can mine the corresponding itemsets are limited (bounded below) by the memory and computational constraints. Therefore, it is critical to be able to perform the analysis in a space- and time-efficient way. During the first few years of research in this area, the primary focus of work was to find FPM algorithms with better computational efficiency.

Several classes of algorithms have been developed for frequent pattern mining, many of which are closely related to one another. In fact, the execution tree of all the algorithms is mostly different in terms of the order in which the patterns are explored, and whether the counting work done for different candidates is independent of one another. To explain this point, we introduce a primitive “baseline” algorithm that forms the heart of most frequent pattern mining algorithms.

Figure 2.1 presents the pseudocode for a very simple “baseline” frequent pattern mining algorithm. The algorithm takes the transaction database \mathcal{T} and a user-defined support value s as input. It first populates all length-one frequent patterns in a frequent pattern data-store, \mathcal{FP} . Then it generates a candidate pattern and computes its support in the database. If the support of the candidate pattern is equal or higher than the minimum support threshold the pattern is stored in \mathcal{FP} . The process continues until all the frequent patterns from the database are found.

In the aforementioned algorithm, candidate patterns are generated from the previously generated frequent patterns. Then, the transaction database is used to determine which of the candidates are truly frequent patterns. The key issues of computational efficiency arise in terms of generating the candidate patterns in an orderly and carefully designed fashion, pruning irrelevant and duplicate candidates, and using well chosen tricks to minimize the work in counting the candidates. Clearly, the

effectiveness of these different strategies depend on each other. For example, the effectiveness of a pruning strategy may be dependent on the order of exploration of the candidates (level-wise vs. depth first), and the effectiveness of counting is also dependent on the order of exploration because the work done for counting at the higher levels (shorter itemsets) can be reused at the lower levels (longer itemsets) with certain strategies, such as those explored in *TreeProjection* and *FP-growth*. Surprising as it might seem, virtually all frequent pattern mining algorithms can be considered complex variations of this simple baseline pseudocode. The major challenge of all of these methods is that the number of frequent patterns and candidate patterns can sometimes be large. This is a fundamental problem of frequent pattern mining although it is possible to speed up the counting of the different candidate patterns with the use of various tricks such as database projections. An analysis on the number of candidate patterns may be found in [25].

The candidate generation process of the earliest algorithms used joins. The original *Apriori* algorithm belongs to this category [1]. Although *Apriori* is presented as a join-based algorithm, it can be shown that the algorithm is a breadth first exploration of a structured arrangement of the itemsets, known as a *lexicographic tree* or *enumeration tree*. Therefore, later classes of algorithms explicitly discuss tree-based enumeration [4, 5]. The algorithms assume a lexicographic tree (or enumeration tree) of candidate patterns and explore the tree using breadth-first or depth-first strategies. The use of the enumeration tree forms the basis for understanding search space decomposition, as in the case of the *TreeProjection* algorithm [5]. The enumeration tree concept is very useful because it provides an understanding of how the search space of candidate patterns may be explored in a systematic and non-redundant way. Frequent pattern mining algorithms typically need to evaluate the support of frequent portions of the enumeration tree, and also rule out an additional layer of infrequent extensions of the frequent nodes in the enumeration tree. This makes the candidate space of all frequent pattern mining algorithms virtually invariant unless one is interested in particular types of patterns such as maximal patterns.

The enumeration tree is defined on the prefixes of frequent itemsets, and will be introduced later in this chapter. Later algorithms such as *FP-growth* perform suffix-based recursive exploration of the search space. In other words, the frequent patterns with a particular pattern as a suffix are explored at one time. This is because *FP-growth* uses the opposite item ordering convention as most enumeration tree algorithms though the recursive exploration order of *FP-growth* is similar to an enumeration tree.

Note that all classes of algorithms, implicitly or explicitly, explore the search space of patterns defined by an enumeration tree of frequent patterns with different strategies such as joins, prefix-based depth-first exploration, or suffix-based depth-first exploration. However, there are significant differences in terms of the order in which the search space is explored, the pruning methods used, and how the counting is performed. In particular, certain projection-based methods help in reusing the counting work for k -itemsets for $(k + 1)$ -itemsets with the use of the notion of projected databases. Many algorithms such as *TreeProjection* and *FP-growth* are able to achieve this goal.

Table 2.1 Toy transaction database and frequent items of each transaction for a minimum support of 3

| tid | Items | Sorted frequent items |
|-----|-------------|-----------------------|
| 2 | a,b,c,d,f,h | a,b,c,d,f |
| 3 | a,f,g | a,f |
| 4 | b,e,f,g | b,f,e |
| 5 | a,b,c,d,e,h | a,b,c,d,e |

This chapter is organized as follows. The remainder of this chapter discusses notations and definitions relevant to frequent pattern mining. Section 2 discusses join-based algorithms. Section 3 discusses tree-based algorithms. All the algorithms discussed in Sects. 2 and 3 extend prefixes of itemsets to generated frequent patterns. A number of methods that extend suffixes of frequent patterns are discussed in Sect. 4. Variants of frequent pattern mining, such as closed and maximal frequent pattern mining, are discussed in Sect. 5. Other optimized variations of frequent pattern mining algorithms are discussed in Sect. 6. Methods for reducing the number of passes, with the use of sampling and aggregation are proposed in Sect. 7. Finally, Sect. 8 concludes chapter with an overall summary.

1.1 Definitions

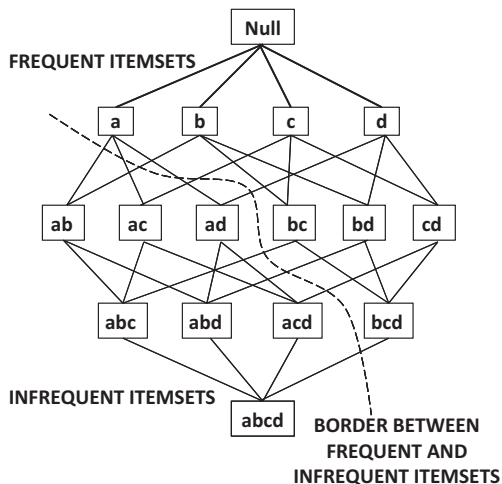
In this section, we define several key concepts of frequent pattern mining (FPM) that we will use in the remaining part of the chapter.

Let, $\mathcal{T} = \{T_1, T_2, \dots, T_n\}$ be a transaction database, where each $T_i \in \mathcal{T}, \forall i = \{1 \dots n\}$ consists of a set of items, say $T_i = \{x_1, x_2, x_3, \dots, x_l\}$. A set $P \subseteq T_i$ is called an itemset. The size of an itemset is defined by the number of items it contains. We will refer an itemset as *l-itemset* (or *l-pattern*), if its size is *l*. The number of transactions containing *P* is referred to as the *support* of *P*. A pattern *P* is defined to be frequent if its support is at least equal to the minimum threshold.

Table 2.1 depicts a toy database with 5 transactions (T_1, T_2, T_3, T_4 and T_5). The second column shows the items in each transaction. In the third column, we show the set of items that are frequent in the corresponding transaction for a minimum support value of 3. For example, the item *h* in transaction with *tid* value of 2 is an infrequent item with a support value of 2. Therefore, it is not listed in the third column of the corresponding row. Similarly, the pattern $\{a, b\}$ (or, *ab* in abbreviated form) is frequent because it has a support value of 3.

The frequent patterns are often used to generate *association rules*. Consider the rule $X \Rightarrow Y$, where *X* and *Y* are sets of items. The confidence of the rule $X \Rightarrow Y$ is equal to the ratio of the support of $X \cup Y$ to that of the support of *X*. In other words, it can be viewed as the conditional probability that *Y* occurs, given that *X* has occurred. The support of the rule is equal to the support of $X \cup Y$. Association rule-generation is a two-phase process. The first phase determines all the frequent patterns at a given minimum support level. The second phase extracts all the rules from these patterns. The second phase is fairly trivial and with limited sophistication. Therefore, most of the algorithmic work in frequent pattern mining focusses on the

Fig. 2.2 The lattice of itemsets



first phase. This chapter will also focus on the first phase of frequent pattern mining, which is generally considered more important and non-trivial.

Frequent patterns satisfy a *downward closure property*, according to which every subset of a frequent pattern is also frequent. This is because if a pattern P is a subset of a transaction, then every pattern $P' \subseteq P$ will also be a subset of T . Therefore, the support of P' can be no less than that of P . The space of exploration of frequent patterns can be arranged as a lattice, in which every node is one of the 2^d possible itemsets, and an edge represents an immediate subset relationship between these itemsets. An example of a lattice of possible itemsets for a universe of items corresponding to $\{a, b, c, d\}$ is illustrated in Fig. 2.2. The lattice represents the search of frequent patterns, and all frequent pattern mining algorithms must, in one way or another, traverse this lattice to identify the frequent nodes of this lattice. The lattice is separated into a frequent and an infrequent part with the use of a *border*. An example of a border is illustrated in Fig. 2.2. This border must satisfy the downward closure property.

The lattice can be traversed with a variety of strategies such as breadth-first or depth-first methods. Furthermore, *candidate nodes* of the lattice may be generated in many ways, such as using joins, or using lexicographic tree-based extensions. Many of these methods are conceptually equivalent to one another. The following discussion will provide an overview of the different strategies that are commonly used.

2 Join-Based Algorithms

Join-based algorithms generate $(k + 1)$ -candidates from frequent k -patterns with the use of joins. These candidates are then validated against the transaction database. The *Apriori* method uses joins to create candidates from frequent patterns, and is one of the earliest algorithms for frequent pattern mining.

2.1 Apriori Method

The most basic join-based algorithm is the *Apriori* method [1]. The *Apriori* approach uses a *level-wise* approach in which all frequent itemsets of length k are generated before those of length $(k + 1)$. The main observation which is used for the *Apriori* algorithm is that every subset of a frequent pattern is also frequent. Therefore, *candidates* for frequent patterns of length $(k + 1)$ can be generated from *known* frequent patterns of length k with the use of joins. A join is defined by pairs of frequent k -patterns that have at least $(k - 1)$ items in common. Specifically, consider a frequent pattern $\{i_1, i_2, i_3, i_4\}$ that is frequent, but has not yet been discovered because only itemsets of length 3 have been discovered so far. In this case, because the patterns $\{i_1, i_2, i_3\}$ and $\{i_1, i_2, i_4\}$ are frequent, they will be present in the set \mathcal{F}_3 of all frequent patterns with length $k = 3$. Note that this particular pair also has $k - 1 = 2$ items in common. By performing a join on this pair, it is possible to create the *candidate* pattern $\{i_1, i_2, i_3, i_4\}$. This pattern is referred to as a *candidate* because it might *possibly* be frequent, and one must either rule it in or rule it out by support counting. Therefore, this candidate is then *validated* against the transaction database by counting its support. Clearly, the design of an efficient support counting method plays a critical role in the overall efficiency of the process. Furthermore, it is important to note that the same candidate can be produced by joining multiple frequent patterns. For example, one might join $\{i_1, i_2, i_3\}$ and $\{i_2, i_3, i_4\}$ to achieve the same result. Therefore, in order to avoid duplication in candidate generation, two itemsets are joined only whether first $(k - 1)$ items are the same, based on a lexicographic ordering imposed on the items. This provides all the $(k + 1)$ -candidates in a non-redundant way.

It should be pointed out that some candidates can be pruned out in an efficient way, without validating them against the transaction database. For any $(k + 1)$ -candidates, it is checked whether *all* its k subsets are frequent. Although it is already known that two of its subsets contributing to the join are frequent, it is not known whether its remaining subsets are frequent. If all its subsets are not frequent, then the candidate can be pruned from consideration because of the downward closure property. This is known as the *Apriori* pruning trick. For example, in the previous case, if the itemset $\{i_1, i_3, i_4\}$ does not exist in the set of frequent 3-itemsets which have already been found, then the candidate itemset $\{i_1, i_2, i_3, i_4\}$ can be pruned from consideration with no further computational effort. This greatly speeds up the overall algorithm. The generation of 1-itemsets and 2-itemsets is usually performed in a specialized way with more efficient techniques.

Therefore, the basic *Apriori* algorithm can be described recursively in level-wise fashion. the overall algorithm comprises of three steps that are repeated over and over again, for different values of k , where k is the length of the pattern generated in the current iteration. The four steps are those of (i) generation of candidate patterns \mathcal{C}_{k+1} by using joins on the patterns in \mathcal{F}_k , (ii) the pruning of candidates from \mathcal{C}_{k+1} , for which all subsets to not lie in \mathcal{F}_k , and (iii) the validation of the patterns in \mathcal{C}_{k+1} against the transaction database \mathcal{T} , to determine the subset of \mathcal{C}_{k+1} which is truly frequent. The algorithm is terminated, when the set of frequent k -patterns \mathcal{F}_k in a given iteration is empty. The pseudo-code of the overall procedure is presented in Fig. 2.3.

Fig. 2.3 The *Apriori* algorithm

Algorithm *Apriori*(Database: \mathcal{T} , Support: s)
begin
 Generate frequent 1-patterns and 2-patterns
 using specialized counting methods and
 denote by \mathcal{F}_1 and \mathcal{F}_2 ;
 $k := 2$;
while \mathcal{F}_k is not empty **do**
begin
 Generate \mathcal{C}_{k+1} by using joins on \mathcal{F}_k ;
 Prune \mathcal{C}_{k+1} with *Apriori* subset pruning trick;
 Generate \mathcal{F}_{k+1} by counting candidates in
 \mathcal{C}_{k+1} with respect to \mathcal{T} at support s ;
 $k := k + 1$;
end
return $\cup_{i=1}^k \mathcal{F}_i$;
end

The computationally intensive procedure in this case is the counting of the candidates in \mathcal{C}_{k+1} with respect to the transaction database \mathcal{T} . Therefore, a number of optimizations and data structures have been proposed in [1] (and also the subsequent literature) to speed up the counting process. The data structure proposed in [1] is that of constructing a *hash-tree* to maintain the candidate patterns. A leaf node of the hash-tree contains a list of itemsets, whereas an interior node contains a hash-table. An itemset is mapped to a leaf node of the tree by defining a path from the root to the leaf node with the use of the hash function. At a node of level i , a hash function is applied to the i th item to decide which branch to follow. The itemsets in the leaf node are stored in sorted order. The tree is constructed recursively in top-down fashion, and a minimum threshold is imposed on the number of candidates in the leaf node.

To perform the counting, all possible k -itemsets which are subsets of a transaction are discovered in a *single* exploration of the hash-tree. To achieve this goal *all possible* paths in the hash tree that could correspond to subsets of the transaction, are followed in recursive fashion, to determine which leaf nodes are relevant to that transaction. After the leaf nodes have been discovered, the itemsets at these leaf nodes that are subsets of that transaction are isolated and their count is incremented. The actual selection of the relevant leaf nodes is performed by recursive traversal as follows. At the root node, all branches are followed such that *any* of the items in the transaction hash to one of branches. At a given interior node, if the i th item of the transaction was last hashed, then all items *following it* in the transaction are hashed to determine the possible children to follow. Thus, by following all these paths, the relevant leaf nodes in the tree are determined. The candidates in the leaf node are stored in sorted order, and can be compared efficiently to the hashed sequence of items in the transaction to determine whether they are relevant. This provides a count of the itemsets relevant to the transaction. This process is repeated for each transaction to determine the final support count for each itemset. It should be pointed out that the reason for using a hash function at the intermediate nodes is to reduce the branching factor of the hash tree. However, if desired, a trie can be used explicitly, in which the degree of a

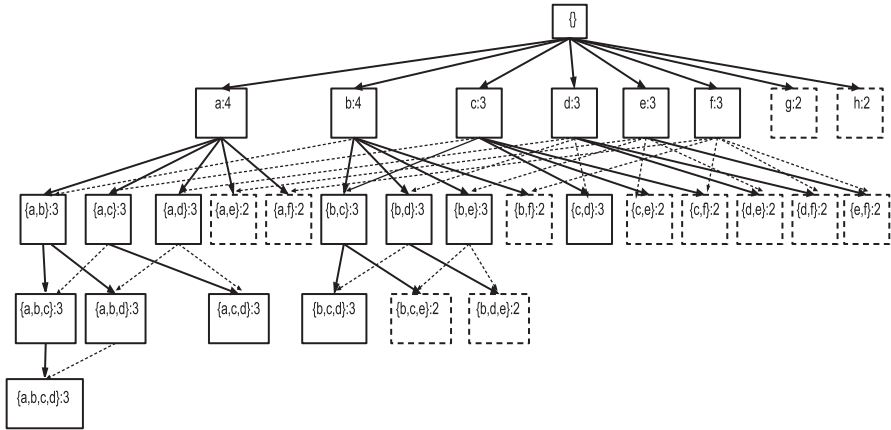


Fig. 2.4 Execution tree of *Apriori* algorithm

node is potentially of the order of the total number of items. An example of such an implementation is provided in [12], and it seems to work quite well. An algorithm that shares some similarities to the *Apriori* method, was independently proposed in [44], and subsequently a combined work was published in [3].

Figure 2.4 illustrates the execution tree of the join-based *Apriori* algorithm over the toy transaction database mentioned in Table 2.1 for minimum support value 3. As mentioned in the pseudocode of *Apriori*, a candidate k -patterns are generated by joining two frequent itemset of size $(k - 1)$. For example, at level 3, the pattern $\{a, b, c\}$ is generated by joining $\{a, b\}$ and $\{a, c\}$. After generating the candidate patterns, the support of the patterns is computed by scanning every transaction in the database and determining the frequent ones. In Fig. 2.4, a candidate patterns is shown in a box along with its support value. A frequent candidate is shown in a solid box, and an infrequent candidate is shown in a dotted box. An edge represents the join relationship between a candidate pattern of size k and a frequent pattern of size $(k - 1)$ such that the latter is used to generate the earlier. The figure also illustrates the fact that a pair of frequent patterns are used to generate a candidate pattern, whereas no candidates are generated from an infrequent pattern.

2.1.1 Apriori Optimizations

Numerous optimizations were proposed for the *Apriori* algorithm [1] that are referred to as *AprioriTid* and *AprioriHybrid* respectively. In the *AprioriTid* algorithm, each transaction is replaced by a shorter transaction or null transaction) during the k th phase. Let the set of $k + 1$ -candidates in C_{k+1} that are contained in transaction T be denoted by $\mathcal{R}(T, C_{k+1})$. This set $\mathcal{R}(T, C_{k+1})$ is added to a newly created transaction database \mathcal{T}'_k . If the set $\mathcal{R}(T, C_{k+1})$ is null, then clearly, a number of different tradeoffs exist with the use of such an approach.

- Because each newly created transaction in \mathcal{T}'_k is much shorter, this makes subsequent support counting more efficient.
- In some cases, no candidate may be a subset of the transaction. Such a transaction can be dropped from the database because it does not contribute to the counting of support values.
- In other cases, more than one candidate may be a subset of the transaction, which will actually increase the overhead of the algorithm. Clearly, this is not a desirable scenario.

Thus, the first two factors improve the efficiency of the new representation, whereas the last factor worsens it. Typically, the impact of the last factor is greater in the early iterations, whereas the impact of the first two factors is greater in the later iterations. Therefore, to maximize the overall efficiency, a natural approach would be to *not* use this optimization in the early iterations, and apply it only in the later iterations. This variation is referred to as the *AprioriHybrid* algorithm [1]. Another optimization proposed in [9] is that the support of many patterns can be inferred from those of key patterns in the data. This is used to significantly enhance the efficiency of the approach.

Numerous other techniques have been proposed that use different techniques to optimize the original implementation of the *Apriori* algorithm. As an example, the method in [1] and [44] share a number of similarities but are somewhat different at the implementation level. A work that combines the ideas from these different pieces of work is presented in [3].

2.2 DHP Algorithm

The DHP algorithm, also known as the *Direct Hashing and Pruning* method [50], was proposed soon after the *Apriori* method. It proposes two main optimizations to speed up the algorithm. The first optimization is to prune the candidate itemsets in each iteration, and the second optimization is to trim the transactions to make the support-counting process more efficient.

To prune the itemsets, the algorithm tracks partial information about candidate $(k+1)$ -itemsets, while explicitly counting the support of candidate k -itemsets. During the counting of candidate k -itemsets, all $(k+1)$ subsets of the transaction are found and hashed into a table that maintains the counts of the number of subsets hashed into each entry. During the phase of counting $(k+1)$ -itemsets, the counts in the hash table are retrieved for each itemset. Clearly, these counts are overestimates because of possible collisions in the hash table. Those itemsets for which the counts are below the user-specified support level are then pruned from consideration.

A second optimization proposed in DHP is that of transaction trimming. A key observation here is that if an item does not appear in at least k frequent itemsets in \mathcal{F}_k , then no frequent itemset in \mathcal{F}_{k+1} will contain that item. This follows from the fact that there should be at least k (immediate) subsets of each frequent pattern in \mathcal{F}_{k+1} .

containing a particular item that also occur in \mathcal{F}_k and also contain that item. This implies that if an item does not appear in at least k frequent itemsets in \mathcal{F}_k , then that item is no longer relevant to further support counting for finding frequent patterns. Therefore, that item can be trimmed from the transaction. This reduces the width of the transaction, and increases the efficiency of processing. The overhead from the data structures is significant, and most of the advantages are obtained for patterns of smaller length such as 2-itemsets. It was pointed out in later work [46, 47, 60] that the use of triangular arrays for support counting of 2-itemsets in the context of the *Apriori* method is even more efficient than such an approach.

2.3 *Special Tricks for 2-Itemset Counting*

A number of special tricks can be used to improve the effectiveness of 2-itemset counting. The case of 2-itemset counting is special and is often similar for the case of join-based and tree-based algorithms. As mentioned above, one approach is to use a triangular array that maintains the counts of the k -patterns explicitly. For each transaction, a nested loop can be used to explore all pairs of items in the transaction and increment the corresponding counts in the triangular array. A number of caching tricks can be used [5] to improve data locality access during the counting process. However, if the number of possible items are very large, this will still be a very significant overhead because it is needed to maintain an entry for each pair of items. This is also very wasteful, if many of the 1-items are not frequent, or some of the 2-item counts are zero. Therefore, a possible approach would be to first prune out all the 1-items which are not frequent. It is simply not necessary to count the support of a 2-itemset unless both of its constituent items are frequent. A hash table can then be used to maintain the frequency counts of the corresponding 2-itemsets. As before, the transactions are explored in a double nested loops, and all pairs of items are hashed into the table, with the caveat, that each of the individual items must be frequent. The set of itemsets which satisfy the support requirements are reported.

2.4 *Pruning by Support Lower Bounding*

Most of the pruning tricks discussed earlier prune itemsets when they are guaranteed *not* meet the required support threshold. It is also possible to skip the counting process for an itemset if the itemset is guaranteed to meet the support threshold. Of course, the caveat here is that the exact support of that itemset will not be available, beyond the knowledge that it meets the minimum threshold. This is sufficient in the case of many applications.

Consider two k -itemsets A and B that have $k - 1$ items $A \cap B$ in common. Then, the union of the items in A and B , denoted by $A \cup B$ will have exactly $k + 1$ items. Then, if $\text{sup}(\cdot)$ represent the support of an itemset, then the support of $A \cup B$ can

be lower bounded as follows:

$$\text{sup}(A \cup B) \geq \text{sup}(A) + \text{sup}(B) - \text{sup}(A \cap B) \quad (2.1)$$

This condition follows directly from set-theoretic considerations. Thus, the support of $(k+1)$ -candidates can be lower bounded in terms of the (already computed) support values of itemsets of length k or less. If the computed value on the right-hand side is greater than the required minimum support, then the counting of the candidate does not need to be performed explicitly, and therefore considerable savings can be achieved. An example of a method which uses this kind of pruning is the *Apriori_LB* method [10].

Another interesting rule is that if the support of an itemset X is the same as that of $X \cup Y$, then for any superset $X' \supseteq X$, it is the case that the support of the itemset X' is the same as that of $X' \cup Y$. This rule can be shown directly as a corollary of the equation above. This is very useful in a variety of frequent pattern mining algorithms. For example, once the support of $X \cup \{i\}$ has been shown to be the same as that of X , then, for any superset X' of X , it is no longer necessary to explicitly compute the support of $X' \cup \{i\}$, after the support of X' has already been computed. Such optimizations have been shown to be quite effective in the context of many frequent pattern mining algorithms [13, 51, 17]. As discussed later, this trick is not exclusive to join-based algorithms, and is often used effectively in tree-based algorithms such as *MaxMiner*, and *MAFIA*.

2.5 Hypercube Decomposition

One feasible way to reduce the computation cost of support counting is to find support of multiple frequent patterns at one time. LCM [66] devise a technique referred to as hypercube decomposition in this purpose. The multiple itemsets obtained at one time, comprise a hypercube in the itemset lattice. Suppose that P is a frequent pattern, $\text{tidset}(P)$ contains the transactions that P is part of, and $\text{tail}(P)$ denotes the latest item extension to the itemset P . $H(P)$ is the set of items e satisfying $e > \text{tail}(P)$ and $\text{tidset}(P) = \text{tidset}(P \cup e)$. The set $H(P)$ is referred to as the hypercube set. Then, for any $P' \subseteq H(P)$, $\text{tidset}(P \cup P') = \text{tidset}(P)$ is true, and $P \cup P'$ is frequent. The work in [66] uses this property in the candidate generation phase. For two itemsets P and $P \cup P'$, we say that P'' is between P and $P \cup P'$ if $P \subseteq P'' \subseteq P \cup P'$. In the phase with respect to P , we output all P'' between P and $P \cup H(P)$. This technique saves significant time in counting.

3 Tree-Based Algorithms

The tree-based algorithm is based on set-enumeration concepts. The candidates can be explored with the use of a subgraph of the lattice of itemsets (see Fig. 2.2), which is also referred to as the lexicographic tree or enumeration tree [5]. These terms will,

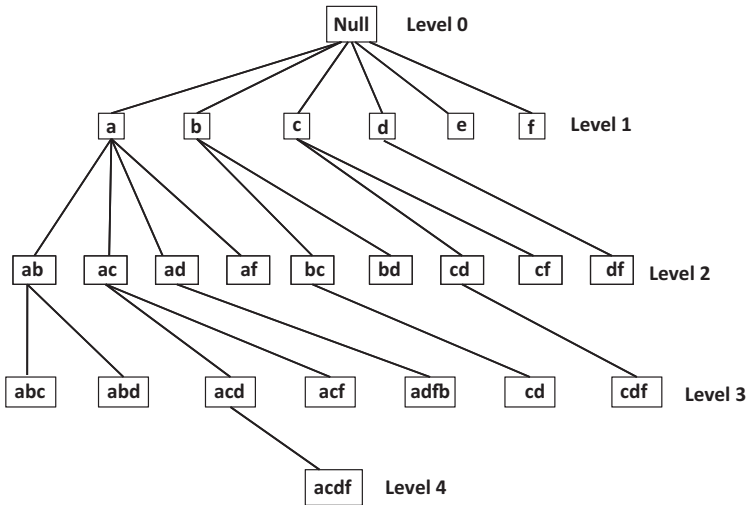


Fig. 2.5 The lexicographic tree (also known as enumeration tree)

therefore, be used interchangeably. Thus, the problem of frequent itemset generation is equivalent to that of constructing the enumeration tree. The tree can be grown in a wide variety of ways such as breadth-first or depth-first order. Because most of the discussion in this section will use this structure as a base for algorithmic development, this concept will be discussed in detail here. The main characteristic of tree-based algorithms is that the enumeration tree (or lexicographic tree) provides a certain order of exploration that can be extremely useful in many scenarios.

It is assumed that a lexicographic ordering exists among the items in the database. This lexicographic ordering is essential for efficient set enumeration without repetition. To indicate that an item i occurs lexicographically earlier than j , we will use the notation $i \leq_L j$. The lexicographic tree is an abstract representation of the large itemsets with respect to this ordering. The lexicographic tree is defined in the following way:

- A node exists in the tree corresponding to each large itemset. The root of the tree corresponds to the *null* itemset.
- Let $I = \{i_1, \dots, i_k\}$ be a large itemset, where i_1, i_2, \dots, i_k are listed in lexicographic order. The parent of the node I is the itemset $\{i_1, \dots, i_{k-1}\}$.

This definition of ancestral relationship naturally defines a tree structure on the nodes that is rooted at the *null* node. A frequent 1-extension of an itemset such that the last item is the contributor to the extension will be called a *frequent lexicographic tree extension*, or simply a tree extension. Thus, each edge in the lexicographic tree corresponds to an item which is the frequent lexicographic tree extension to a node. The frequent lexicographic extensions of node P are denoted by $E(P)$. An example of the lexicographic tree is illustrated in Fig. 2.5. In this example, the frequent lexicographic extensions of node a are b, c, d , and f .

Let Q be the immediate ancestor of the itemset P in the lexicographic tree. The set of *prospective branches* of a node P is defined to be those items in $E(Q)$ which occur lexicographically after the node P . These are the *possible* frequent lexicographic extensions of P . We denote this set by $F(P)$. Thus, we have the following relationship: $E(P) \subseteq F(P) \subset E(Q)$. The value of $E(P)$ in Fig. 2.5, when $P = ab$ is $\{c, d\}$. The value of $F(P)$ for $P = ab$ is $\{c, d, f\}$, and for $P = af$, $F(P)$ is empty.

It is important to point out that virtually all non-maximal and maximal algorithms, starting from *Apriori*, can be considered enumeration-tree methods. In fact, there are few frequent pattern mining algorithms which do not use the enumeration tree, or a subset thereof (in maximal pattern mining) for frequent itemset generation. However, the *order of exploration* of the different algorithms of the lexicographic tree is quite different. For example, *Apriori* uses a breadth-first strategy, whereas other algorithms discussed later in this chapter use a depth-first strategy. Some methods are explicit about the relationship about the candidate generation process with the enumeration tree, whereas others, such as *Apriori*, are not. For example, by examining Fig. 2.4, it is evident that *Apriori* candidates can be generated by joining two frequent siblings of a lexicographic tree. In fact, all candidates can be generated in an exhaustive and non-redundant way by joining frequent siblings. For example, the two itemsets $acdfh$ and $acdfg$ are siblings, because they are children of the node $acdf$. By joining them, one obtains the candidate pattern $acdfgh$. Thus, while the *Apriori* algorithm is a join-based algorithm, it can also be explained in terms of the enumeration tree.

Parts of the enumeration tree may be removed by some of the algorithms by pruning methods. For example, the *Apriori* algorithm uses a levelwise pruning trick. For *maximal* pattern mining the advantages gained from pruning tricks can be very significant. Therefore, the number of candidates in the execution tree of different algorithms is different only because of pruning optimization tricks. However, some methods are able to achieve better *counting* strategies by using the structure of the enumeration tree to avoid re-doing the counting work already done for k -candidates to $(k+1)$ -candidates. Therefore, explicitly introducing the enumeration tree is helpful because it allows a more flexible way to visualize candidate exploration strategies than join-based methods. The explicit introduction of the enumeration tree also helps in understanding whether the gains in different algorithms arise as a result of fewer number of candidates, or whether they arise as a result of better counting strategies.

3.1 AIS Algorithm

The original AIS algorithm [2] is a simple version of the lexicographic-tree algorithm, though it is not directly presented as such. In this approach, the tree is constructed in levelwise fashion and the corresponding itemsets at a given level are counted with the use of the transaction database. The algorithm does not use any specific optimizations to improve the efficiency of the counting process. As will be discussed later, a variety of methods can be used to further improve the efficiency of tree-based algorithms. Thus, this is a primitive approach that explores the entire search space with no optimization.

3.2 *TreeProjection Algorithms*

Two variants of an algorithm which use recursive projections of the transactions down the lexicographic tree structure are proposed in [5] and [4], respectively. The goal of using these recursive projections is to reuse the counting work down at a given level for lower levels of the tree. This reduces the counting work at the lower levels by orders of magnitude, as long as it is possible to successfully manage the memory requirements of the projected transactions. The main difference between the different versions of *TreeProjection* is the exploration strategy used. *TreeProjection* can be viewed as a generic framework that advocates the notion of database projection, in the context of several different strategies for constructing the enumeration tree, such as a breadth-first, depth-first, or a combination of the two. The depth-first version, described in detail in [4], also incorporates maximal pruning, though the disabling of the pruning options can also materialize all the patterns. The breadth-first and depth-first algorithms have different advantages. The former allows level-wise pruning which is not possible in depth-first methods though it is often not used in projection-based methods. The depth-first version allows better memory management. The depth-first approach works best when the itemsets are very long, and it is desirable to quickly discover maximal patterns, so that portions of the lexicographic tree can be pruned off quickly during exploration and it can also be used for discovering all patterns including non-maximal ones. When all patterns are required, including non-maximal ones, the primary difference between different strategies is not one of the size of the candidate space, but that of effective memory management of the projected transactions. This is because the size of the candidate space is defined by the size of the enumeration tree, which is fixed, and is agnostic to the strategy used for tree exploration. On the other hand, memory management of projected transactions is easier with the depth-first strategy because one only needs to maintain a small number of projected transaction sets along the depth of the tree. The notion of database projection is common to *TreeProjection* and *FP-growth*, and helps reduce the counting work by restricting the size of the database used for support counting. *TreeProjection* was developed independently from *FP-growth*. While the *FP-growth* paper provides a brief discussion of *TreeProjection*, this chapter will provide a more detailed discussion of the similarities and differences between the two methods. One major difference between the two methods is that the internal representation of the corresponding projected databases is different in the two cases.

The basic database projection approach is very similar in both cases of *TreeProjection* and *FP-growth*. An important observation is that if a transaction is not relevant for counting at a given node in the enumeration tree, then it will not be relevant for counting in any descendent of that node. Therefore, only those transactions are retained that contain all items in P for counting at the node P in the projected transactions. Note that this set strictly reduces as we move to lower levels of the tree, and the set of relevant transactions at the lower level of the enumeration tree is a subset of the set at a higher level. Furthermore, only the presence of items corresponding to the candidate extensions of a node are relevant for counting at any of the subtrees rooted

Fig. 2.6 Enumeration tree exploration

Algorithm *ExplorePrefix*(Database: \mathcal{T} , Minimum Support: s , Current Pattern Prefix: P)

```

begin
  Count support of 1-items in  $\mathcal{T}$ ;
  Remove infrequent items from  $\mathcal{T}$ ;
  for each frequent item  $i$  in  $\mathcal{T}$  do
    begin
      Append  $i$  to end of  $P$  and add to
        set of frequent patterns;
      Construct conditional database  $\mathcal{T}_i$ 
        with all transactions in  $\mathcal{T}$  containing item  $i$ ;
      Remove items lexicographically  $\leq i$  from  $\mathcal{T}_i$ ;
      ExplorePrefix( $\mathcal{T}_i$ ,  $s$ ,  $P \cup \{i\}$ );
    end
  end
end

```

at that node. Therefore, the database is also projected in terms of attributes, in which only items which are candidate extensions at a node are retained. The candidate set $F(P)$ of item extensions of node P is a very small subset of the universe of items at lower levels of the enumeration tree. In fact, even the items in the node P need not be retained explicitly in the transaction, because they are known to always be present in all the selected transactions based on the first condition. This projection process is performed recursively in top-down fashion down the enumeration tree for counting purposes, where lower level nodes inherit the projections from higher level nodes and add one additional item to the projection at each level. The idea of this inheritance-based approach is that the projected database remembers the counting work done at higher levels of the enumeration tree by (successively) removing irrelevant transactions and irrelevant items at each level of the projection. Such an approach works efficiently because it never repeats the counting work which has already been done at the higher levels. Thus, the primary savings in the strategy arise from avoiding repetitive and wasteful counting.

A bare-bones depth-first version of *TreeProjection*, that is similar to *DepthProject*, but without maximal pruning, is described in Fig. 2.6. A more detailed description with maximal pruning and other optimizations is provided later in this chapter. Because the algorithm is described recursively, the current prefix P (node of the lexicographic tree) being extended is one of the arguments to the algorithm. In the initial call, the value of P is *null* because one intends to determine all frequent descendants at the root of the lexicographic tree. This algorithm recursively extends frequent prefixes and maintains only the transaction database relevant to the prefix. The frequent prefixes are extended by determining the items i that are frequent in \mathcal{T} . Then the itemset $P \cup \{i\}$ is reported. The extension of the frequent prefix can be viewed as a recursive call at a node of the enumeration tree. Thus, at a given enumeration tree node, one now has a completely independent problem of extending the prefix with the projected database that is relevant to all descendants of that node. The conditional database \mathcal{T}_i refers to the subset of the original transaction database \mathcal{T} corresponding to transactions containing item i . Furthermore, the item i and any item occurring lexicographically earlier to it is not retained in the database because

these items are not relevant to counting the extensions of $P \cup \{i\}$. This independent problem is similar in structure to the original problem, and can be solved recursively. Although it is natural to use recursion for the depth-first versions of *TreeProjection*, the breadth-first versions are not defined recursively. Nevertheless, the breadth-first versions explore a pattern space of the same size as the depth-first versions, and are no different either in terms of the tree size or the counting work done over the entire algorithm. The major challenge in the breadth-first version is in maintaining the projected transactions along the breadth of the tree, which is storage-intensive. It is shown in [5], how many of these issues can be resolved with the use of a combination of exploration strategies for tree growth and counting. Furthermore, it is also shown in [5] how breadth-first and depth-first methods may be combined.

Note that this concept of database projection is common between *TreeProjection* and *FP-growth* although there are some differences in the internal representation of the projected databases. The aforementioned description is designed for discovering all patterns, and does not incorporate maximal pattern pruning. When generating *all* the itemsets, the main advantage of the depth-first strategy over the breadth-first strategy is that it is less memory intensive. This is because one does not have to *simultaneously* handle the large number of candidates along the breadth of the enumeration tree at any point in the course of algorithm execution when combined with counting data structures. The overall size of the candidate space is fixed, and defined by the size of the enumeration tree. Therefore, over the entire execution of the algorithm, there is no difference between the two strategies in terms of search space size, beyond memory optimization.

Projection-based algorithms, such as *TreeProjection*, can be implemented either recursively or non-recursively. Depth-first variations of projection strategies, such as *DepthProject* and *FP-growth*, are generally implemented recursively in which a particular prefix (or suffix) of frequent items is grown recursively (see Fig. 2.6). For recursive variations, the structure and size of the recursion tree is the same as the enumeration tree. Non-recursive variations of *TreeProjection* methods directly present the projection-based algorithms in terms of the enumeration tree by storing projected transactions at the nodes in the enumeration tree. Describing projection strategies directly in terms of the enumeration tree is helpful, because one can use the enumeration tree explicitly to optimize the projection. For example, one does not need to project at every node of the enumeration tree, but project only when the size of the database reduces by a particular factor with respect to the nearest ancestor node where the last projection was stored. Such optimizations can reduce the space-overhead of repeated elements in the projected databases at different levels of the enumeration (recursion) tree. It has been shown how to use this optimization in different variations of *TreeProjection*. Furthermore, breadth-first variations of the strategy are naturally defined non-recursively in terms of the enumeration tree. The recursive depth-first versions may be viewed either as divide-and-conquer strategies (because they recursively solve a set of smaller subproblems), or as projection-based counting reuse strategies. The notion of projection-based counting reuse clearly describes how computational savings are achieved in both versions of the algorithm.

When generating *maximal* patterns, the depth-first strategy has clear advantages in terms of pruning as well. We refer the reader to a detailed description of the *DepthProject* algorithm, described later in this chapter. This description describes how several specialized pruning techniques are enabled by the depth-first strategy for maximal pattern mining. The *TreeProjection* algorithm has also been generalized to sequential pattern mining [31]. There are many different types of data structures that may be used in projection-style algorithms. The choice of data structure is sensitive to the data set. Two common choices that are used with *TreeProjection* family of algorithms are as follows:

1. *Arrays*: In this case, the projected database is maintained as 2-dimensional array. One of the dimensions of the array is equal to the number of relevant transactions and the other dimension is equal to the number of relevant items in the projected database. Both dimensions of the projected database reduce from top level to lower levels of the enumeration tree with successive projection.
2. *BitStrings*: In this case, the projected database is maintained as a 0–1 bit string whose width is fixed to the total number of frequent 1-items, but the number of projected transactions reduces with successive projection. Such an approach loses the power of item-wise projection, but this is balanced by the fact that the bit-strings can be used more efficiently for counting operations.

Assume that each transaction T contains n bits, and can therefore be expressed in the form of $\lceil n/8 \rceil$ bytes. Each byte of the transaction contains the information about the presence or absence of eight items, and the integer value of the corresponding bitstring can take on any value from 0 to $2^8 - 1 = 255$. Correspondingly, for each byte of the (projected) transaction at a node, 256 counters are maintained and a value of 1 is added to the counter corresponding to the integer value of that transaction byte. This process is repeated for each transaction in the projected database at node P . Therefore, at the end of this process, one has $256 * \lceil d/8 \rceil$ counts for the d different items. At this point, a postprocessing phase is initiated in which the support of an item is determined by adding the counts of the $256/2 = 128$ counters which take on the value of 1 for that bit. Thus, the second phase requires $128 * d$ operations only, and is independent of database size. The first phase, (which is the bottleneck) is the improvement over the naive counting method because it performs only one operation for each *byte* in the transaction, which contains eight items. Thus, the method would be a factor of eight faster than the naive counting technique, which would need to scan the entire bitstring. Projection is also very efficient in the bitstring representation with simple AND operations.

The major problem with fixed width bitstrings is that they are not efficient representations at lower levels of the enumeration tree at which only a small number of items are relevant, and therefore most entries in these bitstrings are 0. One approach to speed this up is to perform the item-wise projection only at selected nodes in the tree, when the reduction in the number of items from the last ancestor at which the item-wise projection was performed is at particular multiplicative factor. At this point, a shorter bit string is used for representation for the descendants at that node,

Table 2.2 Vertical representation of transactions. Note that the support of itemset ab can be computed as the length of the intersection of the *tidlists* of a and b

| Item | tidlist |
|------|------------|
| a | 1, 2, 3, 5 |
| b | 1, 2, 4, 5 |
| c | 1, 2, 5 |
| d | 1, 2, 5 |
| e | 1, 4, 5 |
| f | 2, 3, 4 |
| g | 3, 4 |
| h | 2, 5 |

until the width of the bitstring is reduced even further by the same multiplicative factor. This ensures that the bit strings representations are not sparse and wasteful.

The key issue here is that different representations provide different tradeoffs in terms of memory management and efficiency. Later in this chapter, an approach called *FP-growth* will be discussed which uses the trie data structure to achieve compression of projected transactions for better memory management.

3.3 Vertical Mining Algorithms

The vertical pattern mining algorithms use a vertical representation of the transaction database to enable more efficient counting. The basic idea of the vertical representation is that one can express the transaction database as an inverted list. In other words, for each transaction identifiers, one can have a list of items that are contained in it. This is referred to as a *tidset* or *tidlist*. An example of a vertical representation of the transactions in Table 2.1 is illustrated in Table 2.2.

The key idea in vertical pattern mining algorithms is that the support of k -patterns can be computed by intersection of the underlying *tidlists*. There are two different ways in which this can be done.

- The support of a k -itemset can be computed as a k -way set intersection of the lists of the individual items.
- The support of a k -itemset can be computed as an intersection of the *tidlists* two $(k - 1)$ -itemsets that join to that k -itemset.

The latter approach is more efficient. The credit for both the notion of vertical *tidlists* and the advantages of recursive intersection of *tidlists* is shared by the *Monet* [56] and the *Partition* algorithms [57]. Not all vertical pattern mining algorithms use an enumeration tree concept to describe the algorithm. Many of the algorithms directly use joins to generate a $(k + 1)$ -candidate pattern from a frequent k -pattern, though even a join-based algorithm, such as *Apriori*, can be explained in terms of an enumeration tree. Many of the later variations of vertical methods use an enumeration tree concept to explore the lattice of itemsets more carefully and realize the full power of the vertical approach. The individual ensemble component of Savasere et al.'s [57] *Partition* algorithm is the progenitor of all vertical pattern mining algorithms today, and the original *Eclat* algorithm is a memory-optimized and candidate partitioned version of this Apriori-like algorithm.

3.3.1 Eclat

Eclat uses a breadth-first approach like Savasere et al.'s algorithm [57] on lattice partitions, after partitioning the candidate set into disjoint groups, using a candidate partitioning approach similar to earlier parallel versions of the *Apriori* algorithm. The *Eclat* [71] algorithm is best described with the concept of an enumeration tree because of the wide variation in the different strategies used by the algorithm. An important contribution of *Eclat* [71] is to recognize the earlier pioneering work of the *Monet* and *Partition* algorithms [56, 57] on recursive intersection of tid lists, and propose many efficient variants of this paradigm.

Different variations of *Eclat* explore the candidates in different strategies. The earliest description of *Eclat* may be found in [74]. A journal paper exploring different aspects of *Eclat* may be found in [71]. In the earliest versions of the work [74], a breadth-first strategy is used. The journal version in [71] also presents experimental results for only the breadth-first strategy, although the possibility of a depth-first strategy is mentioned in the paper. Therefore, the original *Eclat* algorithm should be considered a breadth-first algorithm. More recent depth-first versions of *Eclat*, such as *dEclat*, use recursive *tidlist* intersection with differencing [72], and realize the full benefit of the depth-first approach. The *Eclat* algorithm, as presented in [74], uses a levelwise strategy in which all $(k + 1)$ -candidates within a lattice partition are generated from frequent k -patterns in level-wise fashion, as in *Apriori*. The *tidlists* are used to perform support counting. The frequent patterns are determined from these *tidlists*. At this point, a new levelwise phase is initiated for frequent patterns of size $(k + 1)$.

Other variations and depth-first exploration strategies of *Eclat*, along with experimental results, are presented in later work such as *dEclat* [72]. The *dEclat* work in [72] presents some additional enhancements such as *diffsets* to improve counting. In this chapter, we present a simplified pseudo-code of this version of *Eclat*. The algorithm is presented in Fig. 2.8. The algorithm is structured as a recursive algorithm. A pattern set \mathcal{FP} is part of the input, and is set to the set of all frequent 1-items at the top level call. Therefore, it may be assumed that, at the top level, the set of frequent 1-items and *tidlists* have already been computed, though this computation is not shown in the pseudocode. In each recursive call of *Eclat*, a new set of candidates \mathcal{FP}_i is generated for every pattern (itemset) P_i , which extends the itemset by one unit. The support of a candidate is determined with the use of *tidlist* intersection. Finally, if P_i is frequent, it is added to a pattern set \mathcal{FP}_i for the next level.

Figure 2.7 illustrates the itemset generation tree with support computation by *tidlist* intersection for the sample database from Table 2.1. The corresponding *tidlists* in the tree are also illustrated. All infrequent itemsets in each level are denoted by dotted, and bordered rectangles. For example, an itemset ab is generated by joining b to a . The *tidlist* of (a) is $\{1, 2, 3, 5\}$, and the *tidlist* of b is $\{1, 2, 4, 5\}$. We can determine the support of ab by intersecting the two *tidlists* to obtain the *tidlist* $\{1, 2, 5\}$ of these candidates. Therefore, the support of ab is given by the length of this *tidlist*, which is 3.

Further gains may be obtained with the use of the notion of *diffsets* [72]. This approach realizes the true power of vertical pattern mining. The basic idea, in *diffsets* is to maintain only the portion of the *tidlists* at a node, that correspond to the change in the inverted list from the parent node. Thus, the *tidlists* at a node can be reconstructed by examining the *tidlists* at the ancestors of a node in the tree. The major advantage

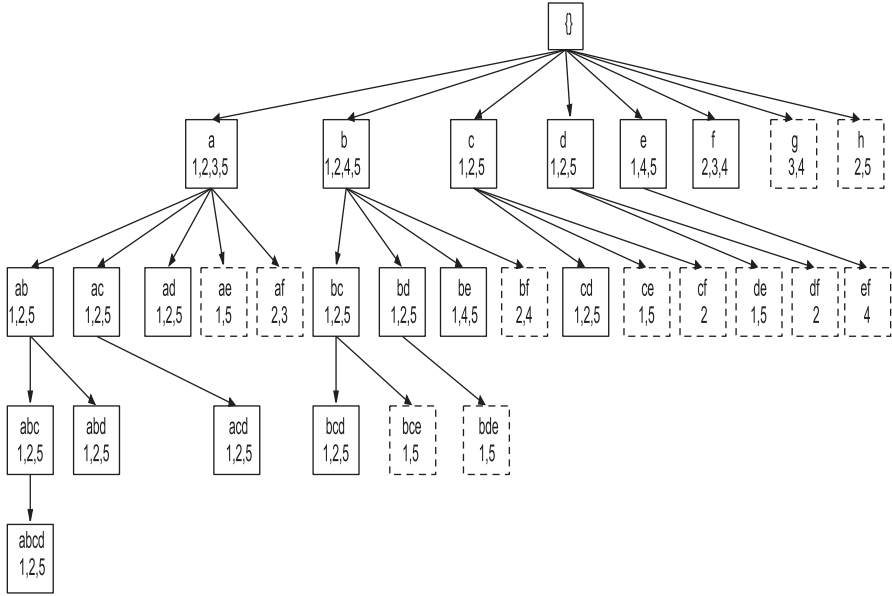


Fig. 2.7 Execution of *Eclat*

Fig. 2.8 The *Eclat* algorithm

```

Algorithm Eclat( $\mathcal{FP}$ , Support:  $s$ )
begin
  for each  $P_i \in \mathcal{FP}$  do
    begin
       $\mathcal{FP}_i = \{\}$ 
      for each  $P_j \in \mathcal{FP}$ , such that  $j > i$  do
        begin
           $P_{ij} = P_i \cup P_j$ 
           $\text{tidset}(P_{ij}) = \text{tidset}(P_i) \cap \text{tidset}(P_j)$ 
           $\text{support}(P_{ij}) = |\text{tidset}(P_{ij})|$ 
          if ( $\text{support}(P_{ij}) \geq s$ )
            Add  $P_{ij}$  to  $\mathcal{FP}_i$ ;
          end
          Eclat( $\mathcal{FP}_i, s$ )
        end
      end
    end
  end

```

of *diffsets* is that they save significant storage in requirements in terms of the size of the data structure required (Fig. 2.8).

Fig. 2.9 Suffix-based pattern exploration

Algorithm *SuffixGrowth*(Database of Frequent Items: \mathcal{T} , Minimum Support: s , Current Pattern Suffix: P)

```

begin
  for each item  $i$  in  $\mathcal{T}$  do
    begin
      Append  $i$  to beginning of  $P$  and add to
        set of frequent patterns;
      Let  $\mathcal{T}_i$  be transactions of  $\mathcal{T}$  containing  $i$ ;
      Remove any item lexicographically  $\geq i$  from  $\mathcal{T}_i$ ;
      Remove infrequent items in  $\mathcal{T}_i$ ;
      if  $\mathcal{T}_i \neq \phi$  SuffixGrowth( $\mathcal{T}_i$ ,  $s$ ,  $\{i\} \cup P$ );
    end
  end
end

```

3.3.2 VIPER

The *VIPER* algorithm [58] uses a vertical approach to mining frequent patterns. The basic idea in the *VIPER* algorithm is to represent the vertical database in the form of compressed bit vectors that are also referred to as *snakes*. These snakes are then used for efficient counting of the frequent patterns. The different compressed representation of the *tidlists* provide a number of optimization advantages that are leveraged by the algorithm. Intrinsically, *VIPER* is not very different from *Eclat* in terms of the basic counting approach. The major difference is in terms of the choice of the compressed bit vector representation, and the efficient handling of this representation. Details may be found in [58].

4 Recursive Suffix-Based Growth

In these algorithms recursive suffix-based exploration of the patterns is performed. Note that in most frequent pattern mining algorithms, the enumeration tree (execution tree) of patterns explores the patterns in the form of a lexicographic tree of itemsets built on the prefixes. Suffix-based methods use a different convention in which the suffixes of frequent patterns are extended. As in all projection-based methods, one only needs to use the transaction database containing itemset P in order to count itemsets that have the suffix P . Itemsets are extended from the suffix backwards. In each iteration, the conditional transaction database (or projected database) of transactions containing the current suffix P being explored is an input to the algorithm. Furthermore, it is assumed that the conditional database contains only frequent extensions of P . For the top-level call, the value of P is null, and the frequent items are determined using a single preprocessing pass that is not shown in the pseudo-code. Because each item is already known to be frequent, the frequent patterns $\{i\} \cup P$ can be immediately generated for each item $i \in \mathcal{T}$. The database is projected further to include only transactions containing i , and a recursive call is initiated with the pattern $\{i\} \cup P$. The projected database \mathcal{T}_i corresponding to transactions containing $\{i\} \cup P$ is determined. Infrequent items are removed from \mathcal{T}_i . Thus, the transactions are recursively projected to reflect the addition of an item in the suffix. Thus, this is a

smaller subproblem that can be solved recursively. The *FP-growth* approach uses the suffix-based pattern exploration, as illustrated in Fig. 2.9. In addition, the *FP-growth* approach uses an efficient data structure, known as the FP-Tree to represent the conditional transaction database \mathcal{T}_i with the use of compressed *prefixes*. The *FP-Tree* will be discussed in more detail in a later section. The suffix in the top level call to the algorithm is the null itemset.

Recursive suffix-based exploration of the pattern space is, in principle, no different from prefix-based exploration of the enumeration tree space with the ordering of the items reversed. In other words, by using a reverse ordering of items, suffix-based recursive pattern space exploration can be simulated with prefix-based enumeration tree exploration. Indeed, as discussed in the last section, prefix-based enumeration tree methods order items from the least frequent to the most frequent, whereas the suffix-based methods of this section order items from the most frequent to the least frequent, to account for this difference. Thus, suffix-based recursive growth has an execution tree that is identical in structure to a prefix-based enumeration tree. This is a difference only of convention, but it does not affect the pattern space that is explored.

It is instructive to compare the suffix-based exploration with the pseudocode of the prefix-based *TreeProjection* algorithm in Fig. 2.6. The two pseudocodes are structured differently because the initial pre-processing pass of removing frequent items is not assumed in the *TreeProjection* algorithm. Therefore, in each recursive call of the prefix-based *TreeProjection*, frequent itemsets must be counted before they are reported. In suffix-based exploration, this step is done as a preprocessing step (for the top-level call) and just before the recursive call for deeper calls. Therefore, each recursive call always starts with a database of frequent items. This is, of course, a difference in terms of how the recursive calls are structured but is not different in terms of the basic search strategy, or the amount of overall computational work required, because infrequent items need to be removed in either case. A few other key differences are evident:

- *TreeProjection* uses database projections on top of a prefix-based enumeration tree. Suffix-based recursive methods have a recursion tree whose structure is similar to an enumeration tree on the frequent suffixes instead of the prefixes. The removal of infrequent items from \mathcal{T}_i in *FP-growth* is similar to determining which branches of the enumeration tree to extend further.
- The use of suffix-based exploration is a difference only of convention from prefix-based exploration. For example, after reversing the item order, one might implement *FP-growth* by growing patterns on the prefixes, but constructing a compressed FP-Tree on¹ the suffixes. The resulting exploration order and execution in the two different implementations of *FP-growth* will be identical, but the latter can be more easily related to traditional enumeration tree methods.

¹ The resulting FP-Tree will be a suffix-based trie.

- Various database projection methods are different in terms of the specific data structures used for the projected database. The different variations of *TreeProjection* use arrays and bit strings to represent the projected database. The *FP-growth* method uses an FP-Tree. The FP-Tree will be discussed in the next section. Later variations of FP-Tree also use combinations of arrays and pointers to represent the projected database. Some variations, such as *OpportuneProject* [38], combine different data structures in an optimized way to obtain the best result.
- Suffix-based recursive growth is inherently defined as a depth-first strategy. On the other hand, as is evident from the discussion in [5], the specific choice of exploration strategy on the enumeration tree is orthogonal to the process of database projection. The overall size of the enumeration tree is the same, no matter how it is explored, unless maximal pattern pruning is used. Thus, *TreeProjection* explores a variety of strategies such as breadth-first and depth-first strategies, with no difference to the (overall) work required for counting. The major challenge with the breadth-first strategy is the simultaneous maintenance of projected transaction sets along the breadth of the tree. The issue of effective memory management of breadth-first strategies is discussed in [5], which shows how certain optimizations such as cache-blocking can improve the effectiveness in this case. Breadth-first strategies also allow certain kinds of pruning such as level-wise pruning.
- The major advantages of depth-first strategies arise in the context of maximal pattern mining. This is because a depth-first strategy discovers the maximal patterns very early, which can be used to prune the smaller non-maximal patterns. In this case, the size of the search space explored truly reduces because of a depth-first strategy. This issue is discussed in the section on maximal pattern mining. The advantages for maximal pattern mining were first proposed in the context of the *DepthProject* algorithm [4].

Next, we will describe the FP-Tree data structure that uses compressed representations of the transaction database for more efficient counting.

4.1 The FP-Growth Approach

The *FP-growth* approach combines suffix-based pattern exploration with a compressed representation of the projected database for more efficient counting. The prefix-based FP-Tree is a compressed representation of the database which is built by considering a fixed order among the items in an itemset [32]. This tree is used to represent the conditional transaction sets \mathcal{T} and \mathcal{T}_i of Fig. 2.9. An FP-Tree may be viewed as a prefix-based trie data structure of the transaction database of frequent items. Just as each node in a trie is labeled with a symbol, a node in the FP-Tree is labeled with an item. In addition, the node holds the support of the itemset defined by the items of the nodes that are on the path from the root to u . By consolidating the prefixes, one obtains compression. This is useful for effective memory management. On the other hand, the maintenance of counts and pointers with the prefixes is an

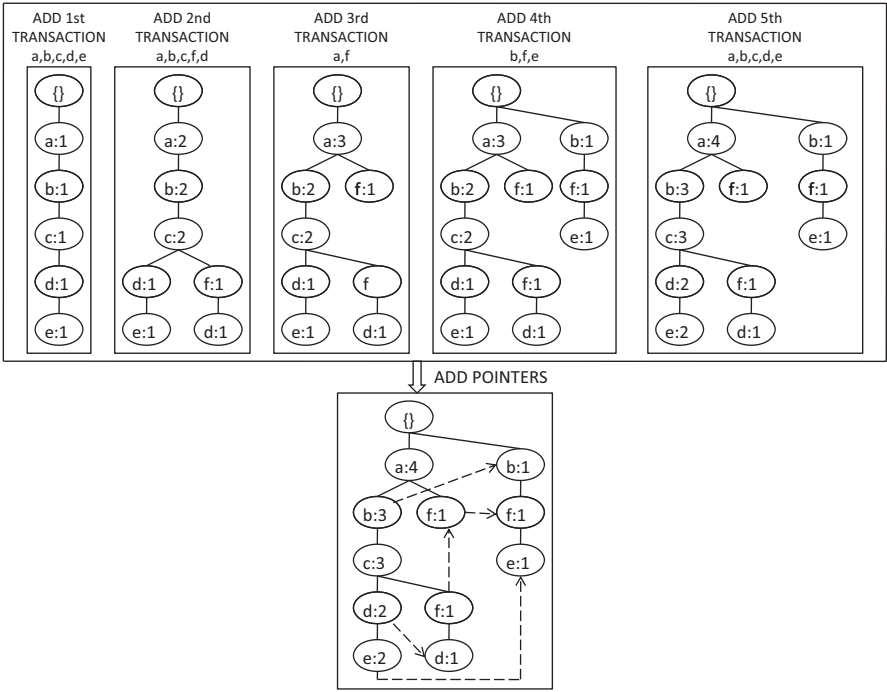


Fig. 2.10 FP-Tree construction

additional overhead. This results in a different set of trade-offs as compared to the array representation.

The initial FP-Tree is constructed as follows. We start with the empty FP-Tree *FPT*. Before constructing the FP-Tree, the database is scanned and infrequent items are removed. The frequent items are sorted in decreasing order of support. The initial construction of FP-Tree is straightforward, and similar to how one might insert a string in a trie. For every insertion, the counts of the relevant nodes that are affected by the insertion are incremented by 1. If there has been any sharing of prefix between the current transaction t being inserted, and a previously inserted transaction then t will be in the same path until the common prefix. Beyond this common prefix, new nodes are inserted in the tree for the remaining items in t , with support count initialized to 1. The above procedure ends when all transactions have been inserted.

To store the items in the final FP-Tree, a list structure called header table is maintained. A chain of pointers threads through the occurrence of the item in the FP-Tree. Thus, this chain of pointers need to be constructed in addition to the trie data structure. Each entry in this table stores the item label and pointers to the node representing the leftmost occurrence of the item in the FP-Tree (first item in the pointer chain). The reason for maintaining these pointers is that it is possible to determine the conditional FP-Tree for an item by chasing the pointers for that item. An example of the initial construction of the FP-Tree data structure from a

Fig. 2.11 The *FP-growth* algorithm

Algorithm *FP-growth*(FP-Tree on Frequent Items: FPT , Minimum Support; s , Current Itemset Suffix: P)

```

begin
  if  $FPT$  is a single path or empty
    for each combination  $C$  of nodes in path do
      report all patterns  $C \cup P$ ;
    else
      for each item  $i$  in  $FPT$  do
        begin
          Generate pattern  $P_i = \{i\} \cup P$ ;
          report pattern  $P_i$  as frequent;
          Use pointer-chasing to extract conditional
            prefix paths for item  $i$ ;
          Construct conditional FP-Tree  $FPT_i$  from conditional
            prefix paths after removing infrequent items;
          if ( $FPT_i \neq \phi$ ) FP-growth( $FPT_i$ ,  $P_i$ ,  $s$ )
        end
      end
    end
  end
end

```

database of five transactions is illustrated in Fig. 2.10. The ordering of the items is a, b, c, d, e, f . It is clear that a trie data structure is created, and the node counts are updated by the insertion of each transaction in the FP-Tree. Figure 2.10 also shows all the pointers between the different items. The sum of the counts on the items on this pointer path is the support of the item. This support is always larger than the minimum support because a full constructed FP-Tree (with pointers) contains only frequent items. The actual counting of the support of item-extensions and the removal of infrequent items must be done during conditional transaction database (and the relevant FP-Tree) creation. The pointer paths are not available during the FP-Tree creation process. For example, the item e has two nodes on this pointer path, corresponding to $e : 2$ and $e : 1$. By summing up these counts, a total count of three for the item e is obtained. It is not difficult to verify that three transactions contain the item e .

With this new compressed representation of the conditional transaction database of frequent items, one can directly extract the frequent patterns. The pseudo-code of the *FP-growth* algorithm is presented in Fig. 2.11. Although this pseudo-code looks much more complex to understand than the earlier pseudocode of Fig. 2.9, the main difference is that more details of the data structure (FP-Tree), used to represent the conditional transaction sets, have been added.

The algorithm accepts a FP-Tree FPT , current itemset suffix P and user defined minimum support s as input. The additional suffix P has been added to the parameter set P to facilitate the recursive description. At the top level call made by the user, the value of P is ϕ . Furthermore, the conditional FP-Tree is constructed on a database of frequent items rather than all the items. This property is maintained across different recursive calls.

For an FP-Tree FPT , the conditional FP-Trees are built for each item i in FPT (which is already known to be frequent). The conditional FP-Trees are constructed by chasing pointers for each item in the FP-Tree. This yields all the conditional prefix

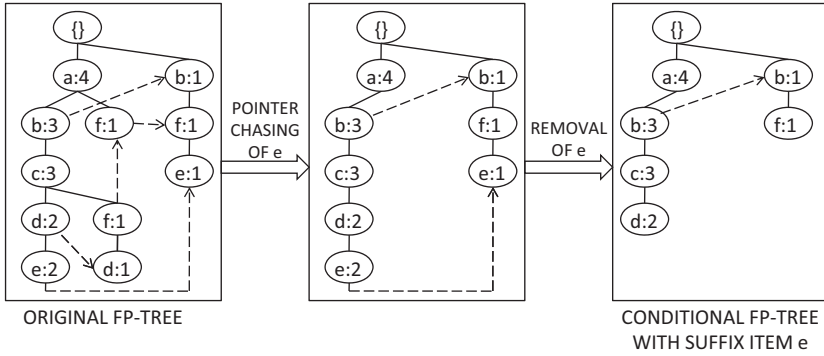


Fig. 2.12 Generating a conditional FP-Tree by pointer chasing

paths for the item i . The infrequent nodes from these paths are removed, and they are put together to create a conditional FP-Tree FPT_i . Because the infrequent items have already been removed from FPT_i the new conditional FP-Tree also contains only frequent items. Therefore, in the next level recursive call, any item from FPT_i can be appended to P_i to generate another pattern. The supports of those patterns can also be reconstructed via pointer chasing during the process of reporting the patterns. Thus, the current pattern suffix P is extended with the frequent item i appended to the front of P . This extended suffix is denoted by P_i . The pattern P_i also needs to be reported as frequent. The resulting conditional FP-Tree FPT_i is the compressed database representation of \mathcal{T}_i of Fig. 2.9 in the previous section. Thus, FPT_i is a smaller conditional tree that contains information relevant only to the extraction of various prefix paths relevant to different items that will extend the suffix P_i further in the backwards direction. Note that infrequent items are removed from FPT_i during this step, which requires the support counting of all items in FPT_i . Because the pointers have not yet been constructed for FPT_i , the support of each item-extension of $\{i\} \cup P$ corresponding to the items in FPT_i must be explicitly determined by locating each instance of an item in FPT_i . This is the primary computational bottleneck step. The removal of infrequent items from FPT_i may result in a different structure of the FP-Tree in the next step.

Finally, if the conditional FP-Tree FPT_i is not empty, the *FP-growth* method is called recursively with parameters corresponding to the conditional FP-Tree FPT_i , extended suffix P_i , and minimum support s . Note that successive projected transaction databases (and corresponding conditional FP-Trees) in the recursion will be smaller because of the recursive projection. The base case of the recursion occurs when the entire FP-Tree is a single path. This is likely to occur when the projected transaction database becomes small enough. In that case, *FP-growth* determines all combinations of nodes on this path, appends the suffix P to them, and reports them.

An example of how the conditional FP-Tree is created for a minimum support of 1 unit, is illustrated in Fig. 2.12. Note that if the minimum support were 2, then the right branch (nodes b and f) would not be included in the conditional FP-Tree. In this case, the pointers for item e are chased in the FP-Tree to create the conditional prefix paths of the relevant conditional transaction database. This represents all transactions

containing e . The counts on the prefix paths are re-adjusted because many branches are pruned. The removal of infrequent items and that of the item e might lead to a conditional FP-Tree that looks very different from the conditional prefix-paths. These kinds of conditional FP-trees need to be generated for each conditional frequent item, although only a single item has been shown for the sake of simplicity. Note that, in general, the pointers may need to be recreated every time a conditional FP-Tree is created.

4.2 Variations

As the database grows larger, the construction of the FP-Tree become challenging both from runtime and space complexity. There have been many works [8, 24, 27, 29, 30, 36, 39, 55, 59, 61, 62] to tackle these challenges. These variations of *FP-growth* method can be classified into two categories. Methods belonging to the first category design memory-based mining process using a memory-resident data structure that holds partitioned database. Methods belonging to the second category improve the efficiency of the FP-Tree representation. In this subsection, we will present these approaches briefly.

4.2.1 Memory-Resident Variations

In the following, a number of different memory-resident variations of the basic *FP-growth* idea will be described.

CT-PRO Algorithm In this work [62], the authors introduced a new FP-Tree like data structure called Compact FP-Tree (CFP-Tree) that holds the same information as FP-Tree but with 50 % less storage. They also designed a mining algorithm called CT-PRO which follows a non-recursive procedure unlike *FP-growth*. As discussed earlier, during the mining process, *FP-growth* constructs many conditional FP-Trees, which becomes an overhead as the patterns get longer or the support gets lower. To overcome this problem, the *CT-PRO* algorithm divides the database into several disjoint projections where each projection is represented as a CFP-Tree. Then a non-recursive mining process is executed over each projection independently. Significant modifications were made to the header Table 4.1 data structure. In the original FP-Tree, the nodes store the support and item label. However, in the CFP-Tree, item labels are mapped to an increasing sequence of integers that is actually the index of the header table. The header table of CFP-Tree stores the support of each item. To compress the original FP-Tree, all identical subtrees are removed by accumulating them and storing the relevant information in the leftmost branch. The header table contains a pointer to each node on the leftmost branch of the CFP-Tree, as these nodes are roots of subtrees starting with different items.

The mining process starts from the pointers of the least frequent items in the header table. This prunes a large number of nodes at an early stage and shrinks the tree

structure. By following the pointers to the same item, a projection of all transactions ending with the corresponding item is built. This projection is also represented as a CFP-Tree called local CFP-Tree. The local CFP-Tree is then traversed to extract the frequent patterns in the projection.

H-Mine Algorithm The authors in [54] proposed an efficient algorithm called *H-Mine*. It uses a memory efficient hyper-structure called H-Struct. The fundamental strategy of *H-Mine* is to partition the database and mine each partition in the memory. Finally, the results from different partitions are consolidated into global frequent patterns. An intelligent module of *H-Mine* is that it can identify whether the database is dense or sparse, and it is able to make dynamic choices between different data structures based on this identification. More details may be found in Chap. 3 on pattern-growth methods.

4.2.2 Improved Data Structure Variations

In this section, several variations of the basic algorithm by improving the underlying data structure will be described.

Using Arrays A significant part of the mining time in *FP-growth* is spent on traversing the tree. To reduce this time, the authors in [29] designed an array based implementation of *FP-growth*, named *FP-growth** which drastically reduces the traversal time of the mining algorithm. It uses the FP-Tree data structure in combination with an array-like data structure and it incorporates various optimization schemes. It should be pointed out that the *TreeProjection* family of algorithms also uses arrays, though the optimizations used are quite different.

When the input database is sparse, the array based technique performs well because the array saves the traversal time for all the items; moreover the initialization of the next level of FP-Trees is easier using an array. But in case of dense database, the tree base representation is more compact. To deal with the situation, *FP-growth** devises a mechanism to identify whether the database is sparse or not. To do so, *FP-growth** counts the number of nodes in each level of the tree. Based on experiments, they found that if the upper quarter of the tree contains less than 15% of the total number of nodes, then the database is most likely dense. Otherwise, it is sparse. If the database turns out to be sparse, *FP-growth** allocates an array for each FP-Tree in the next level of mining.

The nonordfp Approach This work [55] presented an improved implementation of the well known *FP-growth* algorithm using an efficient FP-Tree like data structure that allows faster allocation, traversal and optional projection. The tree nodes do not store their labels (item identifiers). There is no concept of header table. The data structure stores less administrative information in the tree node which allow the recursive step of mining without rebuilding the tree.

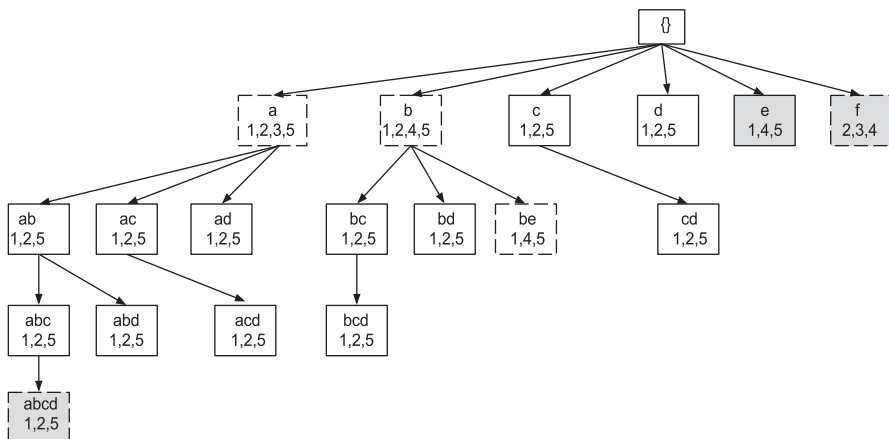


Fig. 2.13 Frequent, maximal and closed itemsets

5 Maximal and Closed Frequent Itemsets

One of the major challenges of frequent itemset mining is that, most of the itemsets mined are subset of the set of single length frequent items. Therefore, a significant amount of time is spent on counting redundant itemsets. One solution to this problem is to discover condensed representations of the frequent itemsets. It will be such representations that synopsizes the property of the set of itemsets completely or partially. The compact representation not only save computational and memory resource but also paved a much easier way towards knowledge discovery stage after mining. Another interesting observation by [53] was that, instead of mining the complete set of frequent itemsets and their associations, association mining only needs to find frequent closed itemsets and their corresponding rules. So, mining frequent closed itemset can fulfill the objectives of mining all frequent itemsets but with less redundancy and better efficiency and effectiveness in mining. In this section, we will discuss two types of condensed representation of itemset: maximal and closed frequent itemset.

5.1 Definitions

Maximal Frequent Itemset Suppose, \mathcal{T} is the transaction database, \mathcal{I} is the set of all items in the database and \mathcal{F} is the set of all frequent itemsets. A frequent itemset $P \in \mathcal{F}$ is called maximal if it has no frequent superset. let \mathcal{M} be the set of all frequent maximal itemsets, which is denoted by

$$\mathcal{M} = \{P \mid P \in \mathcal{F} \text{ and } \nexists Q \supset P, \text{ such that } Q \in \mathcal{F}\}$$

For the toy transaction database in Table 2.1 the frequent maximal itemsets at minimum support 3 are $abcd$, e , f , as illustrated in Fig. 2.13. All the rectangles filled with grey color represent maximal frequent patterns. As we can see in Fig. 2.4, that there are no frequent supersets of $abcd$, e or f .

Closed Frequent Itemset The closure operator γ induces an equivalence relation on the power set of items partitioning it into disjoint subsets called equivalence classes. The largest element with respect to the number of items in each equivalence class is called a closed itemset. A frequent itemset P is closed if $\gamma(P) = P$. From the closure property it can be said that both $\gamma(P)$ and P have the same tidset. In simpler terms, an itemset is closed if it does not have any frequent superset with the same support. A closed itemset C can be written as:

$$C = \{P \mid P \in \mathcal{F} \text{ and } \nexists Q \supset P, \text{ such that } support(Q) = support(P)\}$$

Because maximal itemsets have no frequent superset, they are vacuously closed frequent itemsets. Thus, all maximal patterns are closed. However, there is a key difference between mining maximal itemsets and closed itemsets. Mining maximal itemsets loses information about the support of the underlying itemsets. On the other hand, mining closed itemsets does not lose any information about the support. The support of the missing subsets can be derived from the closed frequent pattern database. One way of viewing closed frequent patterns is as the maximal patterns from each *equi-support* group of frequent patterns. Closed frequent itemsets are a condensed representation of frequent itemsets that is lossless.

For the toy transaction database of Table 2.1 the frequent closed patterns are a , b , $abcd$, be for minimum support value of 3, as illustrated in Fig. 2.13. All the rectangles with dotted border represent closed frequent patterns. The remaining nodes in the tree (not filled and dotted border) represent frequent itemsets.

5.2 Frequent Maximal Itemset Mining Algorithms

In this subsection, we will discuss some of maximal frequent itemset mining algorithms.

5.2.1 MaxMiner Algorithm

The *MaxMiner* algorithm was the first algorithm that used a variety of optimizations to improve the effectiveness of tree explorations [10]. This algorithm is generally focussed on determining maximal patterns rather than all patterns. The author of [10] observed that it is usually sufficient to only report maximal patterns, when frequent patterns are long. This is because of the combinatorial explosion in examining all subsets of patterns. Although the exploration of the tree is still done in breadth-first fashion, a number of optimizations are used to improve the efficiency of exploration:

- The concept of *lookaheads* is defined. Let $F(P)$ be the set of candidate items that might extend node P . Before counting, it is checked whether $P \cup F(P)$ is a subset of any of the frequent patterns found so far. If such is indeed the case, then it is known that the entire subtree rooted at P is frequent, and can be pruned from consideration (for maximal pattern mining). During counting the support of individual item extensions of P , the support of $P \cup F(P)$ is also determined. If the set $P \cup F(P)$ is frequent, then it is known that all itemsets in the entire subtree rooted at that node are frequent. Therefore, the tree does not need to be explored further, and can be pruned.
- The support lower bounding trick discussed earlier can be used to quickly determine patterns which are frequent without explicit counting. The counts of extensions of nodes can be determined without counting in many cases, where the count does not change by extending an item.

It has been shown in [10], that these simple optimizations can improve over the original *Apriori* algorithm by orders of magnitude.

5.2.2 DepthProject Algorithm

The *DepthProject* algorithm is based on the notion of the lexicographic tree, defined in [5]. Unlike *TreeProjection*, the approach aggressively explores the candidates in a depth-first strategy both to ensure better pruning and faster counting. As in *TreeProjection*, the database is recursively projected down the lexicographic tree to ensure more efficient counting. This kind of projection ensures that the counting information for k -candidates is reused for $(k + 1)$ -candidates, as in the case of *FP-growth*.

For the case of the *DepthProject* method [4], the lexicographic tree is explored in depth-first order to maximize the advantage of lookaheads in which entire subtrees can be pruned because it is known that all patterns in them are frequent. The overall pseudocode for the depth-first strategy is illustrated in Fig. 2.14. The pseudocodes for candidate generation and counting are not provided because they are similar to the previously discussed algorithms. However, one important distinction in counting is that projected databases are used for counting. This is similar to the *FP-growth* class of algorithms. Note that the recursive transaction projection is particularly effective with a depth-first strategy because a smaller number of projected databases need to be stored along a path in the tree, as compared to the breadth of the tree.

To reduce the overhead of counting long patterns, the notion of lookaheads are used. At any node P of the tree, let $F(P)$ be its possible (candidate) item extensions. Then, it is checked whether $P \cup F(P)$ is frequent in two ways:

1. Before counting the support of the individual extensions of P (i.e., $\{P \cup \{i\} : \forall i \in F(P)\}$), it is checked whether $P \cup F(P)$ occurs as subset of a frequent itemset that has already been discovered earlier during depth-first exploration. If such is the case, then the entire subtree rooted at P is pruned because it is known

to be frequent and it is not a maximal pattern. This type of pruning is particularly effective with a depth-first strategy.

2. During support counting of the item extensions, the support of $P \cup F(P)$ is also determined. If after support counting, $P \cup F(P)$ turns out to be frequent, then the entire subtree rooted at node P can be pruned. Note that the projected database at node P (as in *TreeProjection*) is used.

Although lookaheads are also used in the *MaxMiner* algorithm, it should be pointed out that the effectiveness of lookaheads is maximized with a depth-first strategy. This is true of the first of the two aforementioned strategies, in which it is checked whether $P \cup F(P)$ is a subset of an already existing frequent pattern. This is because a depth-first strategy tends to explore the itemsets in dictionary order. In dictionary order, maximal itemsets are usually explored much earlier than *most* of their subsets. For example, for a 10-itemset $abcdefghij$, only 9 of the 1024 subsets of the itemsets will be explored before exploring the itemset $abscdefghij$. These 9 itemsets are the immediate prefixes of the itemset. When, the longer itemsets are explored early they become available to prune shorter itemsets.

The following information is stored at each node during the process of construction of the lexicographic tree:

1. The itemset P at that node.
2. The set of lexicographic tree extensions at that node which are $E(P)$.
3. A pointer to the projected transaction set $\mathcal{T}(Q)$, where Q is some ancestor of P (including itself). The root of the tree points to the entire transaction database.
4. A bitvector containing the information about which transactions contain the itemset for node P as a subset. The length of this bitvector is equal to the total number of transactions in $\mathcal{T}(Q)$. The value of a bit for a transaction is equal to one, if the itemset P is a subset of the transaction. Otherwise it is equal to zero. Thus, the number of 1 bits is equal to the number of transactions in $\mathcal{T}(Q)$ which project to P . The bitvectors are used to make the process of support counting more efficient.

After all the projected transactions at a given node have been identified, then finding the subtree rooted at that node is a completely independent itemset generation problem with a *substantially reduced* transaction set. The number of transactions at a node is proportional to the support at that node.

The description in Fig. 2.14 shows how the depth first creation of the lexicographic tree is performed. The algorithm is described recursively, so that the call from each node is a completely independent itemset generation problem that finds all frequent itemsets that are descendants of a node. There are three parameters to the algorithm, a pointer to the database \mathcal{T} , the itemset node N , and the bitvector B . The bitvector B contains one bit for each transaction in $T \in \mathcal{T}$, and indicates whether or not the transaction T should be used in finding the frequent extensions of N . A bit for a transaction T is one, if the itemset at that node is a subset of the corresponding transaction. The first call to the algorithm is from the *null* node, the parameter \mathcal{T} is the entire transaction database. Because each transaction in the database is relevant to perform the counting, the bitvector B consists of all “one” values. One property

Fig. 2.14 The depth first strategy

Algorithm *DepthFirst*(*Itemset Node*: N ,
PointerToDatabase: \mathcal{T} , *Bitvector*: B)

```

begin
 $C = \text{GenerateCandidates}(N)$ ;
 $E = \text{Count}(N, \mathcal{T}, B, C)$ ;
 $\{ \text{Let } E = \{i_1, \dots, i_{|E|}\}, \text{ in lexicographic order } \}$ 
Store frequent itemsets  $N \cup \{i_r\}$  for  $r \in \{1, \dots, |E|\}$ ;
 $B' = \text{CreateBitvector}(N, B, \mathcal{T})$ ;
if (ProjectionCondition) then
  begin
 $\mathcal{T}' = \text{Project}(\mathcal{T}, E, N, B')$ ;
Modify  $B'$  to be a set of  $|\mathcal{T}'|$  ones;
  end;
  else  $\mathcal{T}' = \mathcal{T}$ ;
for  $r := 1$  to  $|E|$  do DepthFirst( $N \cup \{i_r\}, \mathcal{T}', B'$ );
end

```

Subroutine *Project*(*Database*: \mathcal{T} ,
FrequentExtensions: E , *Bitvector*: B)

```

begin
 $\mathcal{T}' = \text{Empty set of transactions}$ ;
for each transaction  $T \in \mathcal{T}$  do
  begin
if corresponding bit in  $B$  is 1 then add  $T \cap E$  to  $\mathcal{T}'$ ;
  end
return( $\mathcal{T}'$ );
end

```

Subroutine *CreateBitvector*(N, B, \mathcal{T})

```

begin
Initialize  $B' = B$ ;
Let  $n$  be the lexicographically largest item in  $N$ ;
for each transaction  $T \in \mathcal{T}$  do
  if  $n \notin T$  then set the corresponding bit in  $B'$  to 0;
return( $B'$ );
end

```

of the *DepthProject* algorithm is that the projection is performed only when the transaction database reduces by a certain size. This is the *ProjectionCondition* in Fig. 2.14.

Most of the nodes in the lexicographic tree correspond to the lower levels. Thus, the counting times at these levels account for most of the CPU times of the algorithm. For these levels, a strategy called bucketing can substantially improve the counting times. The idea is to change the counting technique at a node in the lexicographic tree, if $|E(P)|$ is less than a certain value. In this case, an upper bound on the number of distinct *projected* transactions is $2^{|E(P)|}$. Thus, for example, when $|E(P)|$ is nine, then there are only 512 distinct projected transactions at the node P . Clearly, this is because the projected database contains several repetitions of the same (projected)

Fig. 2.15 Aggregating bucket counts

Algorithm AggregateCounts(*Counts.bucket*[...])

```

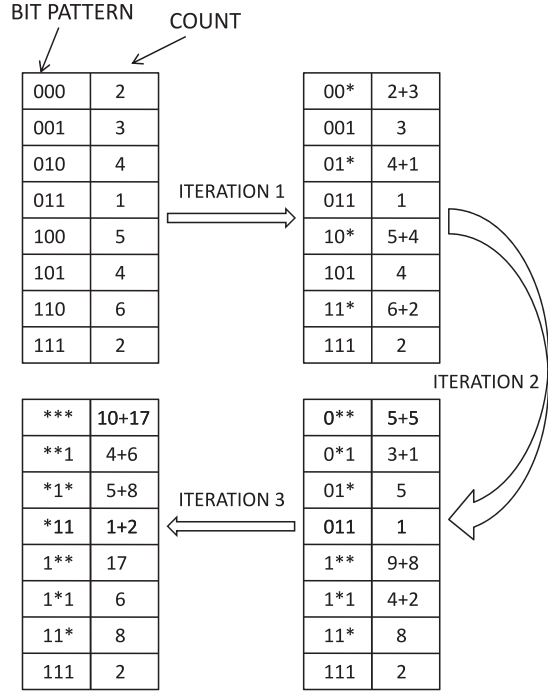
begin
  { We assume that there are  $2^{|E(P)|}$  buckets, one
    corresponding to each bitstring of length  $|E(P)|$  }
   $k = |E(P)|$ ;
  for  $i := 1$  to  $k$  do
    begin
      for  $j := 1$  to  $2^k$  do
        if the  $i$ th bit of bitstring representation
          of  $j$  is 0 then
            begin
               $bucket[j] = bucket[j] + bucket[j + 2^{i-1}]$ ;
            end
          end
        end
      end
    end
  end
end

```

transaction. The fact that the number of *distinct* transactions in the projected database is small can be exploited to yield substantially more efficient counting algorithms. The aim is to count the support for the entire subtree rooted at P with a quick pass through the data, and an additional postprocessing phase which is independent of database size. The process of performing bucket counting consists of two phases:

1. In the first phase, the counts of each distinct transaction present in the projected database are determined. This can be accomplished easily by maintaining $2^{|E(P)|}$ buckets or counters, scanning the transactions one by one, and adding counts to the buckets. The time for performing this set of operations is linear in the number of (projected) database transactions.
2. In the second phase, the counts of the $2^{|E(P)|}$ transaction are used to determine the aggregate support counts for each itemset. In general, the support count of an itemset may be obtained by adding the counts of all the supersets of that itemset to it. A skillful algorithm (from the efficiency perspective) for performing these operations is illustrated in Fig. 2.15.

Consider a string composed of 0, 1, and * that refers to an itemset in which the positions with 0 and 1 are fixed to those values (corresponding to presence or absence of items), while a position with a * is a “don’t care”. Thus, all itemsets can be expressed in terms of 1 and * because itemsets are traditionally defined with respect to presence of items. Consider for example, the case when $|E(P)| = 4$, and there are four items, numbered $\{1, 2, 3, 4\}$. An itemset containing items 2 and 4 is denoted by $*1*1$. We start off with the information on $2^4 = 16$ bitstrings which are composed of 0 and 1. These represent all possible distinct transactions. The algorithm aggregates the counts in $|E(P)|$ iterations. The count for a string with a “*” in a particular position may be obtained by adding the counts for the strings with a 0 and 1 in those positions. For example, the count for the string $*1*1$ may be expressed as the sum of the counts of the strings $01*1$ and $11*1$.

Fig. 2.16 Performing the second phase of bucketing

The procedure in Fig. 2.15 works by starting with the counts of the 0–1 strings, and then converts them to strings with 1 and *. The algorithm requires $|E(P)|$ iterations. In the i th iteration, it increases the counts of all those buckets with a 0 in the i th bit, so that the count now corresponds to a case when that bucket contains a * in that position. This can be achieved by adding the counts of the buckets with a 0 in the i th position to that of the bucket with a 1 in that position, with all other bits having the same value. For example, the count of the string $0*1*$ is obtained by adding the counts of the buckets $001*$ and $011*$. In Fig. 2.15, the process of adding the count of the bucket j to that of the bucket $j + 2^{i-1}$ achieves this.

The second phase of the bucketing operation requires $|E(P)|$ iterations, and each iteration requires $2^{|E(P)|}$ operations. Therefore, the total time required by the method is proportional to $2^{|E(P)|} \cdot |E(P)|$. When $|E(P)|$ is sufficiently small, the time required by the second phase of postprocessing is small compared to the first phase, whereas the first phase is essentially proportional to reading the database for the current projection.

We have illustrated the second phase of bucketing by an example in which $|E(P)| = 3$. The process illustrated in Fig. 2.16 illustrates how the second phase of bucketing is efficiently performed. The exact strings and the corresponding counts in each of the $|E(P)| = 3$ iterations are illustrated. In the first iteration, all those bits with 0 in the lowest order position have their counts added with the count of the bitstring with a 1 in that position. Thus, $2^{|E(P)|-1}$ pairwise addition operations

take place during this step. The same process is repeated two more times with the second and third order bits. At the end of three passes, each bucket contains the support count for the appropriate itemset, where the ‘0’ for the itemset is replaced by a “don’t care” which is represented by a ‘*’. Note that the number of transactions in this example is 27. This is represented by the entry for the bucket ***. Only two transactions contain all three items that is represented by the bucket 111.

The projection-based methods were shown to have an order of magnitude improvement over the *MaxMiner* algorithm. The depth-first approach has subsequently been used in the context of many tree-based algorithms. Other examples of such algorithms include those in [17, 18, 14]. Among these, the MAFIA algorithm [14] is discussed in some detail in the next subsection. An approach which varies on the projection methodology, and uses *opportunistic projection* is discussed in [38]. This algorithm opportunistically chooses between array-based and tree-based representations to represent projected transaction subsets. Such an approach has been shown to be more efficient than many state of the art methods such as the FP-Growth method. Other variations of tree-based algorithms have also been proposed [70] that use different strategies in tree exploration.

5.2.3 MAFIA Algorithm

The MAFIA algorithm proposed in [14] shares a number of similarities to the *Depth-Project* approach, though it uses a bitmap based approach for counting, rather than the use of a projected transaction database. In the bitmap-based approach, a sequence of bits is maintained for each itemset that corresponds to whether or not that transaction contains that particular item. Sparse representations (such as a list of transaction identifiers) may also be used, when the fraction of transactions containing the itemset is small. Note that such an approach may be considered a special case of database projection [5], in which vertical projection is used but horizontal projection is not. This has the advantage of requiring less memory, but it reuses a smaller fraction of the counting information from higher level nodes. A number of other pruning optimizations have also been proposed in this work that further improve the effectiveness of the algorithm. In particular, it has been pointed out that when the support of the extension of a node is the same as that of its parent, then that subtree can be pruned away, because of the counts of all the itemsets in the subtree can be derived from those of other itemsets in the data. This is the same as the support lower bounding trick discussed in Sect. 2.4, and also used in *MaxMiner* for pruning. Thus, the approach in [14] uses many of the same strategies used in *MaxMiner* and *TreeProjection*, but with in a different combination, and with some variations on specific implementation details.

5.2.4 GenMax

Like MAFIA, *GenMax* is a uses the vertical representation to speed up counting. Specifically the *tidlists* are used by *GenMax* to speed up the counting approach. In particular the more recent notion of *diffsets* [72] was used, and a depth-first exploration strategy was used. An approach known as successive focussing was used to further improve the efficiency of the algorithm. The details of the *GenMax* approach may be found in [28].

5.3 Frequent Closed Itemset Mining Algorithms

There are several frequent closed itemset mining algorithms [41, 42, 51–53, 64, 66–69, 73] exist to date. Most of the maximal and closed pattern mining algorithms are based on different variations of the non-maximal pattern mining algorithms. Typically pruning strategies are incorporated within the non-maximal pattern mining algorithms to yield more efficient algorithms.

5.3.1 Close

In this algorithm [52] authors apply *Apriori* based pattern generation over the closed itemset search space. The usage of closed itemset lattice (search space) significantly reduces the overall search space of the algorithm. *Close* operates in iterative manner. Each iteration consists of three phases. First, the closure function is applied for obtaining the candidate closed itemsets and their support. Next, the obtained set of candidate closed itemsets are tested against the minimum support constraint. If succeed, the candidates are marked as frequent closed itemset. Finally the same procedure is initiated to generate the next level of candidate closed itemsets. This process continues until all frequent closed itemsets have been generated.

5.3.2 CHARM

CHARM [73] is a frequent closed itemset mining algorithm, that takes advantage of the vertical representation of database as in the case of *Eclat* [71] for efficient closure checking operation. For pruning the search space *CHARM* uses the following three properties. Suppose for itemset P and Q , if $\text{tidset}(P) = \text{tidset}(Q)$, then it replaces every occurrence of P by $P \cup Q$ and prune the whole branch under Q . On the other hand if $\text{tidset}(P) \subset \text{tidset}(Q)$, it replaces every occurrence of P by $P \cup Q$, but does not prune the branch under Q . Finally if, $\text{tidset}(P) \not\subset \text{tidset}(Q)$, none of the aforementioned prunings can be applied. The initial call of *CHARM* accepts a set(I) of single length frequent item and minimum support as input. As a first step, it sorts I by the increasing the order of support of the items. For each item P ,

CHARM tries to extend it by another item Q from the same set and applies three conditions for pruning. If the newly create itemset by extension is frequent, *CHARM* performs closure-checking to identify whether the itemset is closed. *CHARM* also updates the set I accordingly. In other words, it replaces P with $P \cup Q$, if the corresponding pruning condition is met. If the set I is the not empty, then *CHARM* is called recursively.

5.3.3 CLOSET and CLOSET+

CLOSET [53] and *CLOSET+* [69] frequent closed itemset mining algorithms are inspired by the *FP-growth* method. The *CLOSET* algorithm makes use of the principles of the FP-Tree data structure to avoid the candidate generation step during the process of mining frequent closed itemsets. This work introduces a technique, referred to as single prefix path compression, that quickly assists the mining process. *CLOSET* also applies partition-based projection mechanisms for better scalability. The mining procedure of *CLOSET* follows the *FP-growth* algorithm. However, the algorithm is able to extract only the closed patterns by careful book-keeping. *CLOSET* treats items appearing in every transaction of the conditional database specially. For example, if Q is the set of items that appear in every transaction of the P conditional database then $P \cup Q$ creates a frequent closed itemset if it is not a proper subset of any frequent closed itemset with the equal support. *CLOSET* also prunes the search space. For example, if P and Q are frequent itemset with the equal support where Q is also a closed itemset and $P \subset Q$, then it does not mine the conditional database of P because the latter will not produce any frequent closed itemsets.

CLOSET+ is a follow-up work after *CLOSET* by the same group of authors. *CLOSET+* attempts to design the most optimized frequent closed itemset mining algorithm by finding the best trade-off between depth-first search versus breadth-first search, vertical formats versus horizontal formats, tree structure versus other data structures, top-down versus bottom-up traversal, and pseudo projection versus physical projection of the conditional database. *CLOSET+* keeps track of the unpromising prefix itemsets for generating potential closed frequent itemsets and prunes the search space by deleting them. *CLOSET+* also applies “item merging,” and “sub-itemset” based pruning. To save the memory of the closure checking operation, *CLOSET+* uses the combination of the 2-level hash-indexed tree based method and the pseudo-projection based upward checking method. Interested readers are encouraged to refer to [69] for more details.

5.3.4 DCI_CLOSED

DCI_CLOSED [41, 42] uses a bitwise vertical representation of the input database. *DCI_CLOSED* can be executed independently on each partition of the database in any order and, thus, also in parallel. *DCI_CLOSED* is designed to improve memory-efficiency by avoiding the storage of duplicate closed itemsets. *DCI_CLOSED* designs a novel strategy for searching the lattice that can detect and discard duplicate closed patterns on the fly. Using the concept of order-preserving generators

of frequent closed itemsets, a new visitation scheme of the search space is introduced. Such a visitation scheme results a disjoint sub division of the search space. This also facilitates parallelism. *DCI_CLOSED* applies several optimization tricks to improve execution time, such as the bitwise intersection of tidsets to compute support and closure. Where possible, it reuses previously computed intersections to avoid redundant computations.

6 Other Optimizations and Variations

In this section, a number of other optimizations and variations of frequent pattern mining algorithms will be discussed. Many of these methods are discussed in detail in other chapters of this book, and therefore they will be discussed only briefly here.

6.1 Row Enumeration Methods

Not all frequent pattern mining algorithms follow the fundamental steps of baseline algorithm, there exists a number of special cases, for which specialized frequent pattern mining algorithms have been designed. An interesting case is that of micro-array data sets, in which the columns are very long but the number of rows are not very large. In such cases, a method called *row-enumeration* is used [22, 23, 40, 48, 49] instead of the usual column enumeration, in which combinations of rows are examined during the search process. There are two categories of row enumeration algorithm. One category algorithm perform bottom-up [22, 23, 48] search over the row enumeration tree whereas other category algorithms perform top-down[40] search strategy.

Row enumeration algorithms perform mining over the transpose of the transaction database. In transpose database, each transaction id become item and each item corresponds a transaction. Mining over the transposed database is basically the bottom up search for frequent patterns by enumeration of row sets. However, the bottom-up search strategy cannot take advantage of user-specified minimum support threshold to effectively prune the search space, and therefore leads to longer running time and large memory overhead. As a solution [40] introduce a top-down approach of mining using a novel row enumeration tree. Their approach can take full advantage of user-defined minimum support value and prune the search space efficiently hence lower down the execution time.

Note that, both of the search strategies are applied over the transposed transaction database. Most of developed algorithm using row enumeration technique concentrate on mining frequent closed itemset (explained in Sect. 5). The reason behind this motivation is that due to the nature of micro-array data there exists a large number of redundancy among the frequent patterns for a minimum support threshold and closed patterns are capable of summarizing the whole database. These strategies will be discussed in detail in Chap. 4, and therefore only a brief discussion is provided here.

6.2 Other Exploration Strategies

The advantage of tree-enumeration strategies is that they facilitate the exploration of candidates in the tree in an arbitrary order. A method known as *Pincer-Search* is proposed in [37] that combines top-down and bottom-up exploration in “pincer” fashion to avail of the advantages of both subset and superset pruning. Two primary observations are used in pincer search:

1. Any subset of a frequent itemset is frequent.
2. Any superset of an infrequent itemset is infrequent.

In pincer-search, top-down and bottom-up exploration are combined and irrelevant itemsets are pruned using both observations. More details of this approach are discussed in [37]. Note that, for sparse transaction data, superset pruning is likely to be inefficient. Other recent methods have been proposed for long pattern mining with methods such as “leap search.” These methods are discussed in the chapter on long pattern mining in this book.

7 Reducing the Number of Passes

A major challenge in frequent pattern mining is when the data is disk resident. In such cases, it is desirable to use level-wise methods to ensure that random accesses to disk are minimized. This is the reason that most of the available algorithms use level-wise methods, which ensure that the number of passes over the database are bounded by the size of the longest pattern. Even so, this can be significant, when many long patterns are present in the database. Therefore, a number of methods have been proposed in the literature to reduce the number of passes over the data. These methods could be used in the context of join-based algorithms, tree-based algorithms, or even other classes of frequent pattern mining methods. These correspond to combining the level-wise database passes, using sampling, and using a preprocess-once-query-many paradigm.

7.1 Combining Passes

The earliest work on combining passes was proposed in the original *Apriori* algorithm [1]. The key idea in combining passes is that it is possible to use joins to create candidates of higher order than $(k + 1)$ in a single pass. For example, $(k + 2)$ -candidates can be created from $(k + 1)$ -candidates before actual validation of the $(k + 1)$ -candidates over the data. Then, the candidates of size $(k + 1)$ and $(k + 2)$ can be validated together in a single pass over the data. Although such an approach reduces the number of passes over the data, it has the downside that the number of spurious $(k + 2)$ candidates will be far larger because the $(k + 1)$ candidates were not confirmed to be frequent before they were joined. Therefore, the saving of database

passes comes at an increased computational cost. Therefore, it was proposed in [1] that the approach should be used for later passes, when the number of candidates has already reduced significantly. This reduces the likelihood that the number of candidates blows up too much with this approach.

7.2 *Sampling Tricks*

A number of sampling tricks can be used to greatly improve the efficiency of the frequent pattern mining process. Most sampling methods require two passes over the data, the first of which is used for sampling. An interesting approach that uses two passes with the use of sampling is discussed in [65]. This method generates the approximately frequent patterns over the data, using a sample. False negatives can be reduced by lowering the minimum support level appropriately, so that bounds can be defined on the likelihood of false negatives. False positives can be removed with the use of a second pass over the data. The major downside of the approach is that the reduction in the minimum support level to reduce the number of false negatives can be significant. This also reduces the computational efficiency of the approach. The method however requires only two passes over the data, where the first pass is used to create the sample, and the second pass is used to remove the false positives.

An interesting approach proposed in [57] divides the disk resident database into smaller memory-resident partitions. For each partition, more efficiency algorithms can be used, because of the memory-resident nature of the partition. It should be pointed out that each frequent pattern over the entire database will appear as a frequent pattern in at least one transaction. Therefore, the union of the itemsets over the different transactions provides a superset of the true frequent patterns. A post-processing phase is then used to filter out the spurious itemsets, by counting this candidate set against the transaction database. As long as the partitions are reasonably large, the superset found approximates the true frequent patterns very well, and therefore the additional time spent in counting irrelevant candidates is relatively small. The main advantage of this approach is it requires only two passes over the database. Therefore, such an approach is particularly effective when the data is resident on disk.

The *Dynamic Itemset Counting (DIC)* algorithm [15] divides the database into intervals, and generates longer candidates when it is known that the subsets of these candidates are already frequent. These are then validated over the database. Such an approach can reduce the number of passes over the data, because it implicitly combines the process of candidate generation and counting.

7.3 Online Association Rule Mining

In many applications, a user may wish to query the transaction data to find the association rules or the frequent patterns. In such cases, even at high support levels, it is often impossible to create the frequent patterns in online time because of the multiple passes required over a potentially large database. One of the earliest algorithms for online association rule mining was proposed in [6]. In this approach, an augmented lexicographic tree is stored either on disk or in main-memory. The lexicographic tree is augmented with all the edges represented the subset relationships between itemsets, and is also referred to as the itemset *lattice*. For any given query, the itemset lattice may be traversed to determine the association rules. It has been shown in [6], that such an approach can also be used to determine the non-redundant association rules in the underlying data. A second method [40] uses a condensed frequent pattern tree (instead of a lattice) to pre-process and store the itemsets. This structure can be queried to provide online responses.

A very different approach for online association rule mining has been proposed in [34], in which the transaction database is processed in real time. In this case, an incremental approach is used to mine the transaction database. This is a *Continuous Association Rule Mining Algorithm*, which is referred to as *CARMA*. In this case, transactions are processed as they arrive, and candidate itemsets are generated on the fly, by examining the subsets of that transaction. Clearly, the downside is that such an approach is that it will create a lot more candidates than any of the offline algorithms which use levelwise methods to generate the candidates. This general characteristic is of course true of any algorithm which tries to reduce the number of passes with approximate candidate generation. One interesting characteristic of the *CARMA* algorithm is that it allows the user to change the minimum support level during execution. In that case, the algorithm is guaranteed to have generated the supersets of the true itemsets in the data. If desired, a second pass over the data can be used to remove the spurious frequent itemsets.

Many streaming methods have also been proposed that use only one pass over the transaction data [19–21, 35, 43]. It should be pointed out that it is often difficult to find even 1-itemsets exactly over a data stream because of the one-pass constraint [21], when the number of distinct items is larger than the main memory availability. This is often true of k -itemsets as well, especially at low support levels. Furthermore, if the patterns in the stream change over time, then the frequent k -itemsets will change significantly as well. These methods therefore have the challenge of finding the frequent itemsets efficiently, maintaining them, and handling issues involving evolution of the data stream. Given the numerous challenges of pattern mining in this scenario, most of these methods find the frequent items approximately. These issues will be discussed in detail in Chap. 9 on streaming pattern mining algorithms.

8 Conclusions and Summary

This chapter provides a survey of different frequent pattern mining algorithms. Most frequent pattern algorithms, implicitly or explicitly, explore the enumeration tree of itemsets. Algorithms such as *Apriori* explore the enumeration tree in breadth-first fashion with join-based candidate generation. Although the notion of an enumeration tree is not explicitly mentioned by the *Apriori* algorithm, the execution tree explores the candidates according to an enumeration tree constructed on the prefixes. Other algorithms such as *TreeProjection* and *FP-growth* use the hierarchical relationships between the projected databases for patterns of different lengths, and avoid re-doing the counting work done for the shorter patterns. Maximal and closed versions of frequent pattern mining algorithms are also able to achieve much better pruning performance. A number of efficiency-based optimizations of frequent pattern mining algorithms were also discussed in this chapter.

References

1. R. Agrawal, and R. Srikant. Fast Algorithms for Mining Association Rules in Large Databases, *VLDB Conference*, pp. 487–499, 1994.
2. R. Agrawal, T. Imielinski, and A. Swami. Mining association rules between sets of items in large databases. *ACM SIGMOD Conference*, 1993.
3. R. Agrawal, H. Mannila, R. Srikant, H. Toivonen, and A.I. Verkamo. Fast discovery of association rules, *Advances in Knowledge Discovery and Data Mining*, pp. 307–328, 1996.
4. R. Agarwal, C. C. Aggarwal, and V. V. V. Prasad. Depth-first Generation of Long Patterns, *ACM KDD Conference*, 2000. Also available as IBM Research Report, RC21538, July 1999.
5. R. Agarwal, C. C. Aggarwal, and V. V. V. Prasad. A Tree Projection Algorithm for Generation of Frequent Itemsets, *Journal of Parallel and Distributed Computing*, 61(3), pp. 350–371, 2001. Also available as IBM Research Report, RC21341, 1999.
6. C. C. Aggarwal, P. S. Yu. Online Generation of Association Rules, *ICDE Conference*, 1998.
7. C. C. Aggarwal, P. S. Yu. A New Framework for Itemset Generation, *ACM PODS Conference*, 1998.
8. E. Azkural and C. Aykanat. A Space Optimization for FP-Growth, *FIMI workshop*, 2004.
9. Y. Bastide, R. Taouil, N. Pasquier, G. Stumme, and L. Lakhal. Mining Frequent Patterns with Counting Inference. *ACM SIGKDD Explorations Newsletter*, 2(2), pp. 66–75, 2000.
10. R. J. Bayardo Jr. Efficiently mining long patterns from databases, *ACM SIGMOD Conference*, 1998.
11. J. Blanchard, F. Guillet, R. Gras, and H. Briand. Using Information-theoretic Measures to Assess Association Rule Interestingness. *ICDM Conference*, 2005.
12. C. Borgelt, R. Kruse. Induction of Association Rules: Apriori Implementation, *Conference on Computational Statistics*, 2002. <http://fuzzy.cs.uni-magdeburg.de/~borgelt/software.html>.
13. J.-F. Boulicaut, A. Bykowski, and C. Rigotti. Free-sets: A Condensed Representation of Boolean data for the Approximation of Frequency Queries. *Data Mining and Knowledge Discovery*, 7(1), pp. 5–22, 2003.
14. D. Burdick, M. Calimlim, and J. Gehrke. MAFIA: A Maximal Frequent Itemset Algorithm for Transactional Databases, *ICDE Conference*, 2000. Implementation URL: <http://himalaya-tools.sourceforge.net/Mafia/>.
15. S. Brin, R. Motwani, J.D. Ullman, and S. Tsur. Dynamic itemset counting and implication rules for market basket data. *ACM SIGMOD Conference*, 1997.

16. S. Brin, R. Motwani, and C. Silverstein. Beyond Market Baskets: Generalizing Association Rules to Correlations. *ACM SIGMOD Conference*, 1997.
17. T. Calders, and B. Goethals. Mining all non-derivable frequent itemsets *Principles of Data Mining and Knowledge Discovery*, pp. 1–42, 2002.
18. T. Calders, and B. Goethals. Depth-first Non-derivable Itemset Mining, *SDM Conference*, 2005.
19. T. Calders, N. Dexters, J. Gillis, and B. Goethals. Mining Frequent Itemsets in a Stream, *Informations Systems*, to appear, 2013.
20. J. H. Chang, and W. S. Lee. Finding Recent Frequent Itemsets Adaptively over Online Data Streams, *ACM KDD Conference*, 2003.
21. M. Charikar, K. Chen, and M. Farach-Colton. Finding Frequent Items in Data Streams. *Automata, Languages and Programming*, pp. 693–703, 2002.
22. G. Cong, A. K. H. Tung, X. Xu, F. Pan, and J. Yang. FARMER: Finding interesting rule groups in microarray datasets. *ACM SIGMOD Conference*, 2004.
23. G. Cong, K.-L. Tan, A. K. H. Tung, X. Xu. Mining Top-*k* covering Rule Groups for Gene Expression Data. *ACM SIGMOD Conference*, 2005.
24. M. El-Hajj and O. Zaiane. COFI-tree Mining: A New Approach to Pattern Growth with Reduced Candidacy Generation. *FIMI Workshop*, 2003.
25. F. Geerts, B. Goethals, J. Bussche. A Tight Upper Bound on the Number of Candidate Patterns, *ICDM Conference*, 2001.
26. B. Goethals. Survey on frequent pattern mining, *Technical report, University of Helsinki*, 2003.
27. R. P. Gopalan and Y. G. Sucahyo. High Performance Frequent Pattern Extraction using Compressed FP-Trees, *Proceedings of SIAM International Workshop on High Performance and Distributed Mining*, 2004.
28. K. Gouda, and M. Zaki. Genmax: An efficient algorithm for mining maximal frequent itemsets. *Data Mining and Knowledge Discovery*, 11(3), pp. 223–242, 2005.
29. G. Grahne, and J. Zhu. Efficiently Using Prefix-trees in Mining Frequent Itemsets, *IEEE ICDM Workshop on Frequent Itemset Mining*, 2004.
30. G. Grahne, and J. Zhu. Fast Algorithms for Frequent Itemset Mining Using FP-Trees. *IEEE Transactions on Knowledge and Data Engineering*. 17(10), pp. 1347–1362, 2005, vol. 17, no. 10, pp. 1347–1362, October, 2005.
31. V. Guralnik, and G. Karypis. Parallel tree-projection-based sequence mining algorithms. *Parallel Computing*, 30(4): pp. 443–472, April 2004.
32. J. Han, J. Pei, and Y. Yin. Mining Frequent Patterns without Candidate Generation, *ACM SIGMOD Conference*, 2000.
33. J. Han, H. Cheng, D. Xin, and X. Yan. Frequent Pattern Mining: Current Status and Future Directions, *Data Mining and Knowledge Discovery*, 15(1), pp. 55–86, 2007.
34. C. Hidber. Online Association Rule Mining, *ACM SIGMOD Conference*, 1999.
35. R. Jin, and G. Agrawal. An Algorithm for in-core Frequent Itemset Mining on Streaming Data, *ICDM Conference*, 2005.
36. Q. Lan, D. Zhang, and B. Wu. A New Algorithm For Frequent Itemsets Mining Based On Apriori And FP-Tree, *IEEE International Conference on Global Congress on Intelligent Systems*, pp. 360–364, 2009.
37. D.-I. Lin, and Z. Kedem. Pincer-search: A New Algorithm for Discovering the Maximum Frequent Set, *EDBT Conference*, 1998.
38. J. Liu, Y. Pan, K. Wang. Mining Frequent Item Sets by Opportunistic Projection, *ACM KDD Conference*, 2002.
39. G. Liu, H. Lu and J. X. Yu. AFOPT: An Efficient Implementation of Pattern Growth Approach, *FIMI Workshop*, 2003.
40. H. Liu, J. Han, D. Xin, and Z. Shao. Mining frequent patterns on very high dimensional data: a top- down row enumeration approach. *SDM Conference*, 2006.
41. C. Lucchesse, S. Orlando, and R. Perego. DCI-Closed: A fast and memory efficient algorithm to mine frequent closed itemsets. *FIMI Workshop*, 2004.

42. C. Lucchese, S. Orlando, and R. Perego. Fast and memory efficient mining of frequent closed itemsets. *IEEE TKDE Journal*, 18(1), pp. 21–36, January 2006.
43. G. Manku, R. Motwani. Approximate Frequency Counts over Data Streams. *VLDB Conference*, 2002.
44. H. Mannila, H. Toivonen, and A.I. Verkamo. Efficient algorithms for discovering association rules. *Proceedings of the AAAI Workshop on Knowledge Discovery in Databases*, pp. 181–192, 1994.
45. B. Negrevergne, T. Guns, A. Dries, and S. Nijssen. Dominance Programming for Itemset Mining. *IEEE ICDM Conference*, 2013.
46. S. Orlando, P. Palmerini, R. Perego. Enhancing the a-priori algorithm for frequent set counting. *Third International Conference on Data Warehousing and Knowledge Discovery*, 2001.
47. S. Orlando, P. Palmerini, R. Perego, and F. Silvestri. Adaptive and resource-aware mining of frequent sets. *ICDM Conference*, 2002.
48. F. Pan, G. Cong, A. K. H. Tung, J. Yang, and M. J. Zaki. Finding closed patterns in long biological datasets. *ACM KDD Conference*, 2003.
49. F Pan, A. K. H. Tung, G. Cong, X. Xu. COBBLER: Combining column and Row Enumeration for Closed Pattern Discovery. *SSDBM*, 2004.
50. J.-S. Park, M. S. Chen, and P. S. Yu. An Effective Hash-based Algorithm for Mining Association Rules, *ACM SIGMOD Conference*, 1995.
51. N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Discovering frequent closed itemsets for association rules. *ICDT Conference*, 1999.
52. N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Efficient mining of association rules using closed itemset lattices. *Journal of Information Systems*, 24(1), pp. 25–46, 1999.
53. J. Pei, J. Han, and R. Mao. CLOSET: An Efficient Algorithm for Mining Frequent Closed Itemsets, *DMKD Workshop*, 2000.
54. J. Pei, J. Han, H. Lu, S. Nishio, S. Tang, D. Yang. H-mine: Hyper-structure mining of frequent patterns in large databases, *ICDM Conference*, 2001.
55. B. Racz. nonordfp: An FP-Growth Variation without Rebuilding the FP-Tree, *FIMI Workshop*, 2004.
56. M. Holsheimer, M. Kersten, H. Mannila, and H. Toivonen. A Perspective on Databases and Data Mining, *ACM KDD Conference*, 1995.
57. A. Savasere, E. Omiecinski, and S. Navathe. An efficient algorithm for mining association rules in large databases. *VLDB Conference*, 1995.
58. P. Shenoy, J. Haritsa, S. Sudarshan, G. Bhalotia, M. Bawa, D. Shah. Turbo-charging Vertical Mining of Large Databases. *ACM SIGMOD Conference*, pp. 22–33, 2000.
59. Z. Shi, and Q. He. Efficiently Mining Frequent Itemsets with Compact FP-Tree, IFIP International Federation for Information Processing, V-163, pp. 397–406, 2005.
60. R. Srikant. Fast algorithms for mining association rules and sequential patterns. *PhD thesis, University of Wisconsin, Madison*, 1996.
61. Y. G. Sucahyo and R. P. Gopalan. CT-ITL: Efficient Frequent Item Set Mining Using a Compressed Prefix Tree with PatternGrowth, *Proceedings of the 14th Australasian Database Conference*, 2003.
62. Y. G. Sucahyo and R. P. Gopalan. CT-PRO: A Bottom Up Non Recursive Frequent Itemset Mining Algorithm Using Compressed FP-Tree Data Structures. *FIMI Workshop*, 2004.
63. P.-N. Tan, V. Kumar, and J. Srivastava. Selecting the Right Interestingness Measure for Association Patterns. *ACM KDD Conference*, 2002.
64. I. Taouil, N. Pasquier, Y. Bastide, and L. Lakhal. Mining Basis for Association Rules using Closed Sets, *ICDE Conference*, 2000.
65. H. Toivonen. Sampling large databases for association rules. *VLDB Conference*, 1996.
66. T. Uno, M. Kiyomi and H. Arimura. Efficient Mining Algorithms for Frequent/Closed/Maximal Itemsets, *FIMI Workshop*, 2004.
67. J. Wang, J. Han. BIDE: Efficient Mining of Frequent Closed Sequences. *ICDE Conference*, 2004.

68. J. Wang, J. Han, Y. Lu, and P. Tzvetkov. TFP: An efficient algorithm for mining top- k frequent closed itemsets. *IEEE Transactions on Knowledge and Data Engineering*, 17, pp. 652–664, 2002.
69. J. Wang, J. Han, and J. Pei. CLOSET+: Searching for the Best strategies for mining frequent closed itemsets. *ACM KDD Conference*, 2003.
70. G. I. Webb. Efficient Search for Association Rules, *ACM KDD Conference*, 2000.
71. M. J. Zaki. Scalable algorithms for association mining, *IEEE Transactions on Knowledge and Data Engineering*, 12(3), pp. 372–390, 2000.
72. M. Zaki, and K. Gouda. Fast vertical mining using diffsets. *ACM KDD Conference*, 2003.
73. M. J. Zaki and C. Hsiao. CHARM: An efficient algorithm for closed association rule mining. *SDM Conference*, 2002.
74. M. Zaki, S. Parthasarathy, M. Ogihara, and W. Li. New Algorithms for Fast Discovery of Association Rules. *KDD Conference*, pp. 283–286, 1997.
75. C. Zeng, J. F. Naughton, and JY Cai. On Differentially Private Frequent Itemset Mining. In *Proceedings of 39th International Conference on Very Large data Bases*, 2012.

Chapter 3

Pattern-Growth Methods

Jiawei Han and Jian Pei

Abstract Mining frequent patterns has been a focused topic in data mining research in recent years, with the development of numerous interesting algorithms for mining association, correlation, causality, sequential patterns, partial periodicity, constraint-based frequent pattern mining, associative classification, emerging patterns, etc. Many studies adopt an Apriori-like, candidate generation-and-test approach. However, based on our analysis, candidate generation and test may still be expensive, especially when encountering long and numerous patterns.

A new methodology, called *frequent pattern growth*, which mines frequent patterns without candidate generation, has been developed. The method adopts a divide-and-conquer philosophy to project and partition databases based on the currently discovered frequent patterns and grow such patterns to longer ones in the projected databases. Moreover, efficient data structures have been developed for effective database compression and fast in-memory traversal. Such a methodology may eliminate or substantially reduce the number of candidate sets to be generated and also reduce the size of the database to be iteratively examined, and, therefore, lead to high performance.

In this paper, we provide an overview of this approach and examine its methodology and implications for mining several kinds of frequent patterns, including association, frequent closed itemsets, max-patterns, sequential patterns, and constraint-based mining of frequent patterns. We show that *frequent pattern growth* is efficient at mining large data-bases and its further development may lead to scalable mining of many other kinds of patterns as well.

Keywords Scalable data mining methods and algorithms · Frequent patterns · Associations · Sequential patterns · Constraint-based mining

J. Han (✉)

University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA

e-mail: hanj@cs.uiuc.edu

J. Pei

Simon Fraser University, Burnaby, BC V5A 1S6, Canada

e-mail: jpei@cs.sfu.ca