

Data Mining: Concepts and Techniques

— edited by Manjunath —
— Chapter 6 —

©Jiawei Han and Micheline Kamber
Intelligent Database Systems Research Lab
School of Computing Science
Simon Fraser University, Canada
<http://www.cs.sfu.ca>

May 18, 2003

Data Mining: Concepts and Techniques

1

Chapter 6: Mining Association Rules in Large Databases

- Association rule mining
- Multilevel and Multidimensional association rules
- From association mining to correlation analysis
- Summary

May 18, 2003

Data Mining: Concepts and Techniques

2

What Is Association Mining?

- Association rule mining:
 - Finding frequent patterns, associations, correlations, or causal structures among sets of items or objects in transaction databases, relational databases, and other information repositories.
- Applications:
 - Basket data analysis, cross-marketing, catalog design, loss-leader analysis, clustering, classification, etc.
- Examples.
 - Rule form: "**Body** \rightarrow **Head** [support, confidence]".
 - buys(x, "diapers") \rightarrow buys(x, "beers") [0.5%, 60%]
 - major(x, "CS") \wedge takes(x, "DB") \rightarrow grade(x, "A") [1%, 75%]

May 18, 2003

Data Mining: Concepts and Techniques

3

Association Rule: Basic Concepts

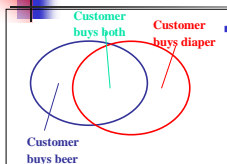
- Given: (1) database of transactions, (2) each transaction is a list of items (purchased by a customer in a visit)
- Find: all rules that correlate the presence of one set of items with that of another set of items
 - E.g., *98% of people who purchase tires and auto accessories also get automotive services done*
- Applications
 - * \Rightarrow *Maintenance Agreement* (What the store should do to boost Maintenance Agreement sales)
 - *Home Electronics* \Rightarrow * (What other products should the store stocks up?)
 - Attached mailing in direct marketing
 - Detecting "ping-pong"ing of patients, faulty "collisions"

May 18, 2003

Data Mining: Concepts and Techniques

4

Rule Measures: Support and Confidence



- Find all the rules $X \& Y \Rightarrow Z$ with minimum confidence and support
 - support, *s*, probability that a transaction contains $\{X \& Y \& Z\}$
 - confidence, *c*, conditional probability that a transaction having $\{X \& Y\}$ also contains *Z*
- Let minimum support 50%, and minimum confidence 50%, we have
- $A \Rightarrow C$ (50%, 66.6%)
 - $C \Rightarrow A$ (50%, 100%)

Transaction ID	Items Bought
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F

May 18, 2003

Data Mining: Concepts and Techniques

5

Association Rule Mining: A Road Map

- Boolean vs. quantitative associations (Based on the types of values handled)
 - buys(x, "SQLServer") \wedge buys(x, "DMBook") \rightarrow buys(x, "DBMiner") [0.2%, 60%]
 - age(x, "30..39") \wedge income(x, "42..48K") \rightarrow buys(x, "PC") [1%, 75%]
- Single dimension vs. multiple dimensional associations (see ex. Above)
- Single level vs. multiple-level analysis
 - What brands of beers are associated with what brands of diapers?
- Various extensions
 - Correlation, causality analysis
 - Association does not necessarily imply correlation or causality
 - Maxpatterns and closed itemsets
 - Constraints enforced
 - E.g., small sales (sum < 100) trigger big buys (sum > 1,000?)

May 18, 2003

Data Mining: Concepts and Techniques

6

Mining Association Rules—An Example

Transaction ID	Items Bought
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F

Min. support 50%
Min. confidence 50%

Frequent Itemset	Support
{A}	75%
{B}	50%
{C}	50%
{A,C}	50%

For rule $A \Rightarrow C$:

support = support({A C}) = 50%

confidence = support({A C})/support({A}) = 66.6%

The **Apriori** principle:

Any subset of a frequent itemset must be frequent

May 18, 2003

Data Mining: Concepts and Techniques

7

Mining Frequent Itemsets: the Key Step

- Find the **frequent itemsets**: the sets of items that have minimum support
 - A subset of a frequent itemset must also be a frequent itemset
 - i.e., if {AB} is a frequent itemset, both {A} and {B} should be a frequent itemset
 - Iteratively find frequent itemsets with cardinality from 1 to k (k -itemset)
- Use the frequent itemsets to generate association rules.

May 18, 2003

Data Mining: Concepts and Techniques

8

The Apriori Algorithm

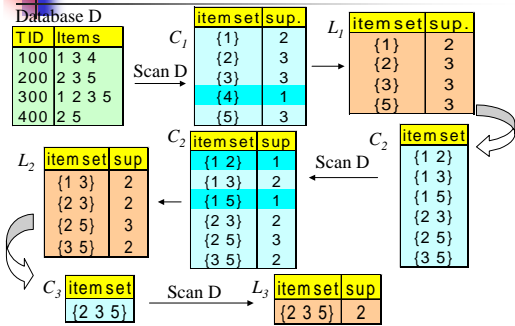
- Join Step:** C_k is generated by joining L_{k-1} with itself
- Prune Step:** Any $(k-1)$ -itemset that is not frequent cannot be a subset of a frequent k -itemset
- Pseudo-code:**
 - C_k : Candidate itemset of size k
 - L_k : frequent itemset of size k
 - $L_1 = \{\text{frequent items}\}$
 - for $(k = 1; L_k \neq \emptyset; k++)$ do begin
 - C_{k+1} = candidates generated from L_k ;
 - for each transaction t in database do
 - increment the count of all candidates in C_{k+1} that are contained in t
 - L_{k+1} = candidates in C_{k+1} with min_support
 - end
 - return $\cup_k L_k$;

May 18, 2003

Data Mining: Concepts and Techniques

9

The Apriori Algorithm — Example



May 18, 2003

Data Mining: Concepts and Techniques

10

How to Generate Candidates?

- Suppose the items in L_{k-1} are listed in an order
- Step 1: self-joining L_{k-1}
 - insert into C_k
 - select $p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$
 - from $L_{k-1} p, L_{k-1} q$
 - where $p.item_1 = q.item_1, \dots, p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$
- Step 2: pruning
 - forall **itemsets** c in C_k do
 - forall **(k-1)-subsets** s of c do
 - if (s is not in L_{k-1}) then delete c from C_k

May 18, 2003

Data Mining: Concepts and Techniques

11

Example of Generating Candidates

- $L_3 = \{abc, abd, acd, ace, bcd\}$
- Self-joining: $L_3 * L_3$
 - $abcd$ from abc and abd
 - $acde$ from acd and ace
- Pruning:
 - $acde$ is removed because ade is not in L_3
- $C_4 = \{abcd\}$

May 18, 2003

Data Mining: Concepts and Techniques

12

Is Apriori Fast Enough? — Performance Bottlenecks

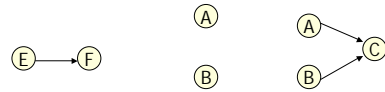
- The core of the Apriori algorithm:
 - Use frequent $(k - 1)$ -itemsets to generate **candidate** frequent k -itemsets
 - Use database scan and pattern matching to collect counts for the candidate itemsets
- The bottleneck of Apriori: **candidate generation**
 - Huge candidate sets:
 - 10^4 frequent 1-itemset will generate 10^7 candidate 2-itemsets
 - To discover a frequent pattern of size 100, e.g., $\{a_1, a_2, \dots, a_{100}\}$, one needs to generate $2^{100} \approx 10^{30}$ candidates.
 - Multiple scans of database:
 - Needs $(n + 1)$ scans, n is the length of the longest pattern

May 18, 2003

Data Mining: Concepts and Techniques

13

Generalizations: Finding Episodes from Sequences



Given a set of E of *event types*, an *event sequence* s is a sequence of pairs (e, t) , where $e \in E$ and t is an integer, the occurrence time of the event type e . An **episode** α is a partial order of event types. Episodes can be viewed as graphs.

May 18, 2003

Data Mining: Concepts and Techniques

14

Frequent Episodes

- **Frequency of an episode** : Given an event sequence S and a window width W , the frequency of an episode α is the fraction of slices of width W taken from S such that the slice contains events of types occurring in α in the order described by α .
- **Task**: given an event sequence s , a set E of episodes, a window width win , and a frequency threshold min_fr , find the collection $FE(s, win, min_fr)$ of all episodes from the set that occur at least in a fraction of min_fr of all the windows of width win on the sequence s .

May 18, 2003

Data Mining: Concepts and Techniques

15

subepisode

- An episode $\beta \preceq \alpha$ is defined as a **subepisode** of an episode α if all the nodes in β occur also in α and if all the relationships between the nodes in β are also present in α .
- I.e., β is an induced subgraph of α
- We write $\beta \preceq \alpha$ if β is a subepisode of α , and $\beta \prec \alpha$ if $\beta \preceq \alpha$ and $\beta \neq \alpha$.

May 18, 2003

Data Mining: Concepts and Techniques

16

$FE(s, win, min_fr)$

- $C_1 = \{a \text{ in } E \mid |a| = 1\}$;
- $l := 1$;
- While C_l is not a null set do
 - /* database pass */
 - Compute $F_l := \{a \text{ in } C_l \mid fr(a, s, win) \geq min_fr\}$;
 - $l := l + 1$
 - /* candidate generation */
 - Compute $C_l := \{a \text{ in } E \mid |a| = l, \text{ and for all } b \text{ in } E \text{ such that } b \preceq a \text{ and } |b| < l \text{ we have } b \text{ in } F_{|b|}\}$;
- End
- For all l do output F_l .

May 18, 2003

Data Mining: Concepts and Techniques

17

- Algorithm performs a levelwise search in the class of episodes following the subepisode relation.
- Search starts from the most general episodes—episodes with only one event.
- On each level the algorithm first computes a collection of candidate episodes, and then checks their frequencies from the event sequence.
- Algorithm does at most $k + 1$ passes through the data, where k is the number of edges and vertices in the largest frequent episode.

May 18, 2003

Data Mining: Concepts and Techniques

18

Chapter 6: Mining Association Rules in Large Databases

- Association rule mining
- Multilevel and Multidimensional association rules**
- From association mining to correlation analysis
- Summary

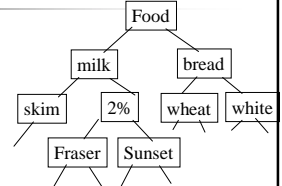
May 18, 2003

Data Mining: Concepts and Techniques

19

Multiple-Level Association Rules

- Items often form hierarchy.
- Items at the lower level are expected to have lower support.
- Rules regarding itemsets at appropriate levels could be quite useful.
- Transaction database can be encoded based on dimensions and levels
- We can explore shared multi-level mining



TID	Items
T1	{111, 121, 211, 221}
T2	{111, 211, 222, 323}
T3	{112, 122, 221, 411}
T4	{111, 121}
T5	{111, 122, 211, 221, 413}

May 18, 2003

Data Mining: Concepts and Techniques

20

Mining Multi-Level Associations

- A top_down, progressive deepening approach:
 - First find high-level strong rules:
 - milk \rightarrow bread [20%, 60%].
 - Then find their lower-level "weaker" rules:
 - 2% milk \rightarrow wheat bread [6%, 50%].
- Variations at mining multiple-level association rules.
 - Level-crossed association rules:
 - 2% milk \rightarrow *Wonder* wheat bread
 - Association rules with multiple, alternative hierarchies:
 - 2% milk \rightarrow *Wonder* bread

May 18, 2003

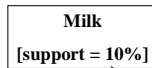
Data Mining: Concepts and Techniques

21

Uniform Support

Multi-level mining with uniform support

Level 1
min_sup = 5%



Level 2
min_sup = 5%

May 18, 2003

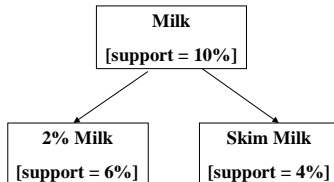
Data Mining: Concepts and Techniques

22

Reduced Support

Multi-level mining with reduced support

Level 1
min_sup = 5%



Level 2
min_sup = 3%

May 18, 2003

Data Mining: Concepts and Techniques

23

Multi-Dimensional Association: Concepts

- Single-dimensional rules:
 - $\text{buys}(X, \text{"milk"}) \Rightarrow \text{buys}(X, \text{"bread"})$
- Multi-dimensional rules: μ 2 dimensions or predicates
 - Inter-dimension association rules (*no repeated predicates*)
 - $\text{age}(X, \text{"19-25"}) \wedge \text{occupation}(X, \text{"student"}) \Rightarrow \text{buys}(X, \text{"coke"})$
 - hybrid-dimension association rules (*repeated predicates*)
 - $\text{age}(X, \text{"19-25"}) \wedge \text{buys}(X, \text{"popcorn"}) \Rightarrow \text{buys}(X, \text{"coke"})$
- Categorical Attributes
 - finite number of possible values, no ordering among values
- Quantitative Attributes
 - numeric, implicit ordering among values

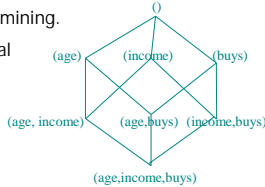
May 18, 2003

Data Mining: Concepts and Techniques

24

Static Discretization of Quantitative Attributes

- Discretized prior to mining using concept hierarchy.
- Numeric values are replaced by ranges.
- In relational database, finding all frequent k-predicate sets will require k or k+1 table scans.
- Data cube is well suited for mining.
- The cells of an n-dimensional cuboid correspond to the predicate sets.
- Mining from data cubes can be much faster.



May 18, 2003

Data Mining: Concepts and Techniques

25

Chapter 6: Mining Association Rules in Large Databases

- Association rule mining
- Multilevel and Multidimensional association rules
- From association mining to correlation analysis
- Summary

May 18, 2003

Data Mining: Concepts and Techniques

26

Interestingness Measurements

- Objective measures
 - Two popular measurements:
 - support; and
 - confidence
- Subjective measures (Silberschatz & Tuzhilin, KDD95)
 - A rule (pattern) is interesting if
 - it is *unexpected* (surprising to the user); and/or
 - actionable* (the user can do something with it)

May 18, 2003

Data Mining: Concepts and Techniques

27

Criticism to Support and Confidence

- Example 1: (Aggarwal & Yu, PODS98)
 - Among 5000 students
 - 3000 play basketball
 - 3750 eat cereal
 - 2000 both play basket ball and eat cereal
 - $play\ basketball \Rightarrow eat\ cereal$ [40%, 66.7%] is misleading because the overall percentage of students eating cereal is 75% which is higher than 66.7%.
 - $play\ basketball \Rightarrow not\ eat\ cereal$ [20%, 33.3%] is far more accurate, although with lower support and confidence

	basketball	not basketball	sum(row)
cereal	2000	1750	3750
not cereal	1000	250	1250
sum(col.)	3000	2000	5000

May 18, 2003

Data Mining: Concepts and Techniques

28

Criticism to Support and Confidence (Cont.)

- Example 2:
 - X and Y: positively correlated,
 - X and Z, negatively related
 - support and confidence of $X \Rightarrow Z$ dominates

X	1	1	1	1	0	0	0	0
Y	1	1	0	0	0	0	0	0
Z	0	1	1	1	1	1	1	1

- We need a measure of dependent or correlated events

$$corr_{A,B} = \frac{P(A \cup B)}{P(A)P(B)}$$

Rule	Support	Confidence
$X \Rightarrow Y$	25%	50%
$X \Rightarrow Z$	37.50%	75%

- $P(B|A)/P(B)$ is also called the *lift* of rule $A \Rightarrow B$

May 18, 2003

Data Mining: Concepts and Techniques

29

Other Interestingness Measures: Interest

- Interest (correlation, lift) $\frac{P(A \wedge B)}{P(A)P(B)}$
 - taking both P(A) and P(B) in consideration
 - $P(A \wedge B) = P(B) \cdot P(A)$, if A and B are independent events
 - A and B negatively correlated, if the value is less than 1; otherwise A and B positively correlated

X	1	1	1	0	0	0	0
Y	1	1	0	0	0	0	0
Z	0	1	1	1	1	1	1

Itemset	Support	Interest
X,Y	25%	2
X,Z	37.50%	0.9
Y,Z	12.50%	0.57

May 18, 2003

Data Mining: Concepts and Techniques

30

Chapter 6: Mining Association Rules in Large Databases

- Association rule mining
- Multilevel and Multidimensional association rules
- From association mining to correlation analysis
- **Summary**

May 18, 2003

Data Mining: Concepts and Techniques

31

Summary

- Association rule mining
 - probably the most significant contribution from the database community in KDD
 - A large number of papers have been published
- Many interesting issues have been explored
- An interesting research direction
 - Association analysis in other types of data: spatial data, multimedia data, time series data, etc.

May 18, 2003

Data Mining: Concepts and Techniques

32

References

- R. Agarwal, C. Aggarwal, and V. V. V. Prasad. A tree projection algorithm for generation of frequent itemsets. In *Journal of Parallel and Distributed Computing (Special Issue on High Performance Data Mining)*, 2000.
- R. Agrawal, T. Imielinski, and A. Swami. Mining association rules between sets of items in large databases. *SIGMOD'93*, 207-216, Washington, D.C.
- R. Agrawal and R. Srikant. Fast algorithms for mining association rules. *VLDB'94* 487-499, Santiago, Chile.
- R. Agrawal and R. Srikant. Mining sequential patterns. *ICDE'95*, 3-14, Taipei, Taiwan.
- R. J. Bayardo. Efficiently mining long patterns from databases. *SIGMOD'98*, 85-93, Seattle, Washington.
- S. Brin, R. Motwani, and C. Silverstein. Beyond market basket: Generalizing association rules to correlations. *SIGMOD'97*, 265-276, Tucson, Arizona.
- S. Brin, R. Motwani, J. D. Ullman, and S. Tsur. Dynamic itemset counting and implication rules for market basket analysis. *SIGMOD'97*, 255-264, Tucson, Arizona, May 1997.
- K. Beyer and R. Ramakrishnan. Bottom-up computation of sparse and iceberg cubes. *SIGMOD'99*, 359-370, Philadelphia, PA, June 1999.
- D.W. Cheung, J. Han, V. Ng, and C.Y. Wong. Maintenance of discovered association rules in large databases: An incremental updating technique. *ICDE'96*, 106-114, New Orleans, LA.
- M. Fang, N. Shivakumar, H. Garcia-Molina, R. Motwani, and J. D. Ullman. Computing iceberg queries efficiently. *VLDB'98*, 299-310, New York, NY, Aug. 1998.

May 18, 2003

Data Mining: Concepts and Techniques

33

References (2)

- G. Grahne, L. Lakshmanan, and X. Wang. Efficient mining of constrained correlated sets. *ICDE'00*, 512-521, San Diego, CA, Feb. 2000.
- Y. Fu and J. Han. Meta-rule-guided mining of association rules in relational databases. *KDD'95*, 39-46, Singapore, Dec. 1995.
- T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama. Data mining using two-dimensional optimized association rules: Scheme, algorithms, and visualization. *SIGMOD'96*, 13-23, Montreal, Canada.
- E.-H. Han, G. Karypis, and V. Kumar. Scalable parallel data mining for association rules. *SIGMOD'97*, 277-288, Tucson, Arizona.
- J. Han, G. Dong, and Y. Yin. Efficient mining of partial periodic patterns in time series database. *ICDE'99*, Sydney, Australia.
- J. Han and Y. Fu. Discovery of multiple-level association rules from large databases. *VLDB'95*, 420-431, Zurich, Switzerland.
- J. Han, J. Pei, and Y. Yin. Mining frequent patterns without candidate generation. *SIGMOD'00*, 1-12, Dallas, TX, May 2000.
- T. Imielinski and H. Mannila. A database perspective on knowledge discovery. *Communications of ACM*, 39:58-64, 1996.
- M. Kamber, J. Han, and J. Y. Chiang. Meta-rule-guided mining of multi-dimensional association rules using data cubes. *KDD'97*, 207-210, Newport Beach, California.
- M. Klemettinen, H. Mannila, P. Ronkainen, H. Toivonen, and A.I. Verkamo. Finding interesting rules from large sets of discovered association rules. *CIKM'94*, 401-408, Gaithersburg, Maryland.

May 18, 2003

Data Mining: Concepts and Techniques

34

References (3)

- F. Korn, A. Labrinidis, Y. Kotidis, and C. Faloutsos. Ratio rules: A new paradigm for fast, quantifiable data mining. *VLDB'98*, 582-593, New York, NY.
- B. Lent, A. Swami, and J. Widom. Clustering association rules. *ICDE'97*, 220-231, Birmingham, England.
- H. Lu, J. Han, and L. Feng. Stock movement and n-dimensional inter-transaction association rules. *SIGMOD Workshop on Research Issues on Data Mining and Knowledge Discovery (DMKD'98)*, 12:1-12:7, Seattle, Washington.
- H. Mannila, H. Toivonen, and A. I. Verkamo. Efficient algorithms for discovering association rules. *KDD'94*, 181-192, Seattle, WA, July 1994.
- H. Mannila, H. Toivonen, and A. I. Verkamo. Discovery of frequent episodes in event sequences. *Data Mining and Knowledge Discovery*, 1:259-289, 1997.
- R. Meo, G. Psaila, and S. Ceri. A new SQL-like operator for mining association rules. *VLDB'96*, 122-133, Bombay, India.
- R.J. Miller and Y. Yang. Association rules over interval data. *SIGMOD'97*, 452-461, Tucson, Arizona.
- R. Ng, L. V. S. Lakshmanan, J. Han, and A. Pang. Exploratory mining and pruning optimizations of constrained associations rules. *SIGMOD'98*, 13-24, Seattle, Washington.
- N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Discovering frequent closed itemsets for association rules. *ICDT'99*, 398-416, Jerusalem, Israel, Jan. 1999.

May 18, 2003

Data Mining: Concepts and Techniques

35

References (4)

- J.S. Park, M.S. Chen, and P.S. Yu. An effective hash-based algorithm for mining association rules. *SIGMOD'95*, 175-186, San Jose, CA, May 1995.
- J. Pei, J. Han, and R. Mao. CLOSET: An Efficient Algorithm for Mining Frequent Closed Itemsets. *DMKD'00*, Dallas, TX, 11-20, May 2000.
- J. Pei and J. Han. Can We Push More Constraints Into Frequent Pattern Mining? *KDD'00*. Boston, MA, Aug. 2000.
- G. Platetsky-Shapiro. Discovery, analysis, and presentation of strong rules. In G. Platetsky-Shapiro and W. J. Frawley, editors, *Knowledge Discovery in Databases*, 229-238. AAAI/MIT Press, 1991.
- B. Ozden, S. Ramaswamy, and A. Silberschatz. Cyclic association rules. *ICDE'98*, 412-421, Orlando, FL.
- J.S. Park, M.S. Chen, and P.S. Yu. An effective hash-based algorithm for mining association rules. *SIGMOD'95*, 175-186, San Jose, CA.
- S. Ramaswamy, S. Mahajan, and A. Silberschatz. On the discovery of interesting patterns in association rules. *VLDB'98*, 368-379, New York, NY.
- S. Sarawagi, S. Thomas, and R. Agrawal. Integrating association rule mining with relational database systems: Alternatives and implications. *SIGMOD'98*, 343-354, Seattle, WA.
- A. Savasere, E. Omiecinski, and S. Navathe. An efficient algorithm for mining association rules in large databases. *VLDB'95*, 432-443, Zurich, Switzerland.
- A. Savasere, E. Omiecinski, and S. Navathe. Mining for strong negative associations in a large database of customer transactions. *ICDE'98*, 494-502, Orlando, FL, Feb. 1998.

May 18, 2003

Data Mining: Concepts and Techniques

36



References (5)

- C. Silverstein, S. Brin, R. Motwani, and J. Ullman. Scalable techniques for mining causal structures. VLDB'98, 594-605, New York, NY.
- R. Srikant and R. Agrawal. Mining generalized association rules. VLDB'95, 407-419, Zurich, Switzerland, Sept. 1995.
- R. Srikant and R. Agrawal. Mining quantitative association rules in large relational tables. SIGMOD'96, 1-12, Montreal, Canada.
- R. Srikant, Q. Vu, and R. Agrawal. Mining association rules with item constraints. KDD'97, 67-73, Newport Beach, California.
- H. Toivonen. Sampling large databases for association rules. VLDB'96, 134-145, Bombay, India, Sept. 1996.
- D. Tsur, J. D. Ullman, S. Abitboul, C. Clifton, R. Motwani, and S. Nestorov. Query flocks: A generalization of association-rule mining. SIGMOD'98, 1-12, Seattle, Washington.
- K. Yoda, T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama. Computing optimized rectilinear regions for association rules. KDD'97, 96-103, Newport Beach, CA, Aug. 1997.
- M. J. Zaki, S. Parthasarathy, M. Ogihara, and W. Li. Parallel algorithm for discovery of association rules. Data Mining and Knowledge Discovery, 1:343-374, 1997.
- M. Zaki. Generating Non-Redundant Association Rules. KDD'00. Boston, MA, Aug. 2000.
- O. R. Zaiane, J. Han, and H. Zhu. Mining Recurrent Items in Multimedia with Progressive Resolution Refinement. ICDE'00, 461-470, San Diego, CA, Feb. 2000.