

C2: Mining Association Rules: Apriori and Its Related Issues

Ming-Syan Chen

September 25, 2018

Agenda for Classes 2018 (Might be revised as we progress)

- Class 1 – (9/11) Overview of data mining
- Class 2 – (9/18) [R \(I\)](#), [Wush Wu](#)
- Class 3 – (9/25) [Association, Apriori and its related issues](#)
- Class 4 – (10/2) Data stream mining, FP Tree, Vertical mining
- Class 5 – (10/9) Classification: decision tree, [GPGPU](#)
- Class 6 – (10/16) Description of Data, Project announcement
- Class 7 – (10/23) [R \(II\)](#), [Wush Wu](#)

Tentative Class Agenda (cont'd)

- Class 8 – (10/30) Data exploration, more on decision trees, rule-based classifiers
- Class 9 – (11/6) [Scikit learn, LibSVM, Preparation for HW3 and HW4](#)
- Class 10 – (11/13) KNN, Bays, Neural network, Concept of SVM
- Class 11 – (11/20) Abstract presentation, SVM, Clustering, K-means, PAM
- Class 12 – (11/27) More on clustering; Sequential pattern mining;
- Class 13 – (12/4) Web mining, PageRank, etc.

3

Tentative Class Agenda (cont'd)

- Class 14– (12/11) [When Database and Data Mining Meet, Prof. Mingling Lo](#)
- Class 15 – (12/18) Project presentation I
- Class 16 – (12/25) Project presentation II
(Final Exam according to Univ. Schedule)
(Project due 1/24/2019)
- Happy New Year!

4

Procedure of Data Mining

- Obtain and look over the data
- Decide your goal (usually a stretched and reachable one)
- Data cleaning/cleansing
- Choose data granularity, feature selection
- Apply mining methods
- Decide what to output and in what form
- Interpret your results (may have iterative refinements); convince your receiver/boss

5

Mining Capabilities

- Association
- Classification
- Clustering
- Sequential Pattern
- and more

Mining Association Rules

- Transaction data analysis: Mining association rules
 - Given: (1) a database of transactions
(2) each tx has a list of items purchased
- Find all asso. rules: the presence of one set of items implies the presence of another set of items
 - people who purchased hammers also purchased nails

7

Two Parameters

- Confidence (how true)
 - the rule $X \& Y \Rightarrow Z$ has 90% conf. means 90% of customers who bought X and Y also bought Z
- Support (how useful is the rule)
 - useful rules should have some minimum tx support

8

Mining Association Rules in Transaction DBs

- Measurement of rule strength in a transaction DB.

$$A \rightarrow B \text{ [support, confidence]}$$

$$\text{support} = \text{Prob}(A \cup B) = \frac{\text{\#_of_trans_that_contain_both } A \text{ and } B}{\text{total_}\#_of_trans}$$

$$\text{confidence} = \text{Prob}(B|A) = \frac{\text{\#_of_trans_that_contain_both } A \text{ and } B}{\text{\#_of_trans_containing } A}$$

- We are often interested in only strong associations, i.e.

$$\text{support} \geq \text{min_sup} \quad \text{and} \quad \text{confidence} \geq \text{min_conf.}$$

- Examples.

$$\text{milk} \rightarrow \text{bread} [5\%, 60\%].$$

$$\text{tire} \wedge \text{auto_accessories} \rightarrow \text{auto_services} [2\%, 80\%].$$

9

Two Steps for Mining Asso.

- Determining “large itemsets”
 - the main factor for overall performance
- Generating rules

10

Two approaches for Large Itemset Counting

- Apriori-Based
 - R. Agrawal and R. Srikant
- FP-Tree-Based
 - J. Han and J. Pei, etc (SIGMOD 2000)

11

Methods for Mining Association Rules

- Apriori (Agrawal & Srikant'94).
- Itemset generation
 - derivation of large 1-itemsets L_1 : At the first iteration, scan all the transactions and count the number of occurrences for each item.
 - **level-wise derivation**: At the k -th iteration, the candidate set C_k are those whose every $(k - 1)$ -item subset is in L_{k-1} . Scan DB and count the # of occurrences for each candidate itemset.
 - the cardinality of C_2 is huge
 - the exe time for the first 2 iterations is the dominating factor to overall performance

12

Support=2 tx's (i.e., 50%)

Database D

TID	Items
100	A C D
200	B C E
300	A B C E
400	B E

Scan
D
→

C₁

Itemset	Sup.
{A}	2
{B}	3
{C}	3
{D}	1
{E}	3

L₁

Itemset	Sup.
{A}	2
{B}	3
{C}	3
{E}	3

C₂

Itemset
{A B}
{A C}
{A E}
{B C}
{B E}
{C E}

Scan
D
→

C₂

Itemset	Sup.
{A B}	1
{A C}	2
{A E}	1
{B C}	2
{B E}	3
{C E}	2

L₂

Itemset	Sup.
{A C}	2
{B C}	2
{B E}	3
{C E}	2

C₃

Itemset
{B C E}

Scan
D
→

C₃

Itemset	Sup.
{B C E}	2

L₃

Itemset	Sup.
{B C E}	2

BE=>C conf:66%

13

Two Steps for Mining Asso. (cont'd)

- for each large itemset m do
 - for each subset p of m do
 - if $(\text{sup}(m)/\text{sup}(m-p)) \geq \text{minconf}$ then
 - output the rule $(m-p) \Rightarrow p$
 - with $\text{conf} = \text{sup}(m)/\text{sup}(m-p)$ and
 - $\text{support} = \text{sup}(m)$
- $m = \{a, c, d, e, f, g\}$ 2000 tx's
 $p = \{a, d\}$ 5000 tx's
 $\{a, d\} \Rightarrow \{c, e, f, g\}$ conference: 40%, support: 2000 tx's

14

Properties of Apriori

- Downward closure for large (also called frequent) itemset generation
- The bottleneck is usually in C2
- Database scan is expensive
- The setting of “support” and “confidence”
- Using “top-k” itemsets instead of support
 - How to do itemset generation

15

Follow-ups of Apriori

- Data Stream mining
 - W.-G. Teng, M.-S. Chen and P. S. Yu, “A Regression-Based Temporal Pattern Mining Scheme for Data Streams,” *Proc. of the 29th Intern'l Conf. on Very Large Data Bases (VLDB-2003)*, September 9-12, 2003.
- Upper bound on the number of large itemsets
 - F. Geerts and B. Goethals and J. V. D. Bussche, “**Tight upper bounds on the number of candidate patterns**”, TODS 2005)
- “closed large itemset”
- Spawned many works to improve its efficiency and also to explore its variations

16

Closed Itemsets and Maximal Itemsets

- An itemset X is called closed if there does not exist an itemset Y , s.t. Y contains X and $s(Y)=s(X)$
- A large itemset X is called maximal if there does not exist a large itemset Y , s.t., Y contains X
- Q1: if a large itemset X is closed, is X always maximal?
- Q2: if a large itemset X is maximal large itemset, is X always closed?

17

Closed Itemsets and Maximal Itemsets (cont'd)

- DB1 is said to be “equivalent” to DB2 if DB1 can be obtained from tx permutation from DB1
- Q: Prove or disapprove with an example.
 - Q1. If DB1 and DB2 are equivalent, do they lead to the same set of closed itemsets?
 - Q2. If DB1 and DB2 lead to the same set of closed itemsets, are DB1 and DB2 always equivalent to each other?

18

Redundant Rules

- For the same support and confidence, if we have a rule $\{a,d\} \Rightarrow \{c,e,f,g\}$, do we have
 - $\{a,d\} \Rightarrow \{c,e,f\}$
 - $\{a\} \Rightarrow \{c,e,f,g\}$
 - $\{a,d,c\} \Rightarrow \{e,f,g\}$
 - $\{a\} \Rightarrow \{d,c,e,f,g\}$

19

Scan Reduction

- Use candidate sets to generate candidate sets
e.g., Instead of $C_i \rightarrow Li \rightarrow C_{i+1} \text{ (dbscan)} \rightarrow Li+1$
We use $C_i \rightarrow C_{i+1}' \rightarrow C_{i+2}' \text{ (dbscan)} \rightarrow Li+1, Li+2..$
- Save the runs of database scans
- May get back to use large itemsets to generate candidate sets if so necessary

20

Improvement for Apriori

- DHP

- J. Park, M.-S. Chen, and P. Yu. “*An effective hash based algorithm for mining association rules.*” *Proceedings of ACM SIGMOD*, May 1995. A complete version in Using A Hash-Based Method with Transaction Trimming for Mining Association Rules,” IEEE Trans. on Knowledge and Data Eng., vol. 9, no. 5, pp. 813-825, Sept./Oct. 1997.

- Hash table scheme

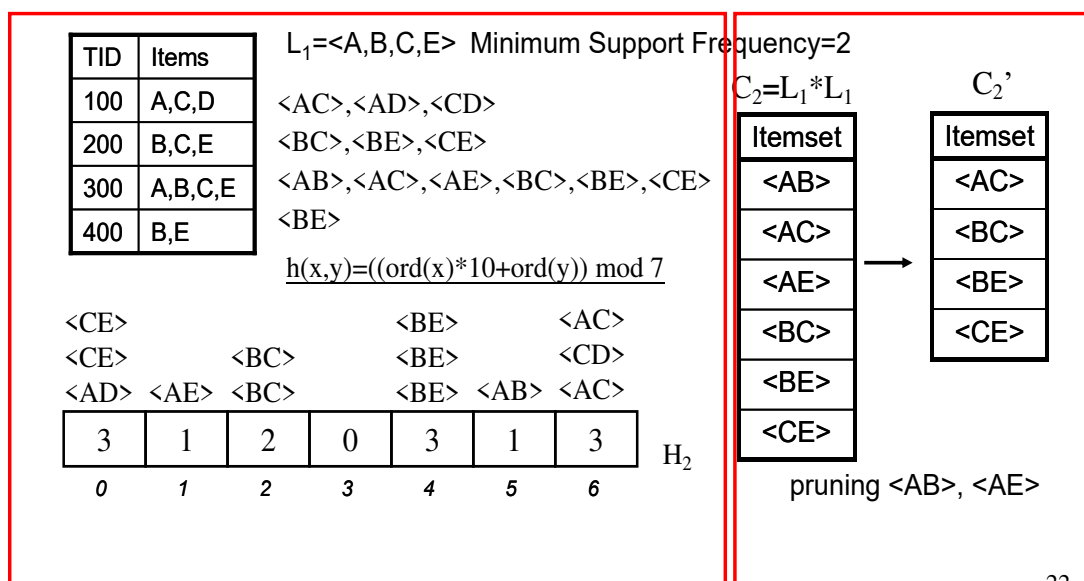
- Eliminate infrequent candidate itemsets in the early phase

- Transaction items pruning

- Eliminate infrequent items from the database

21

Candidate Itemsets Pruning

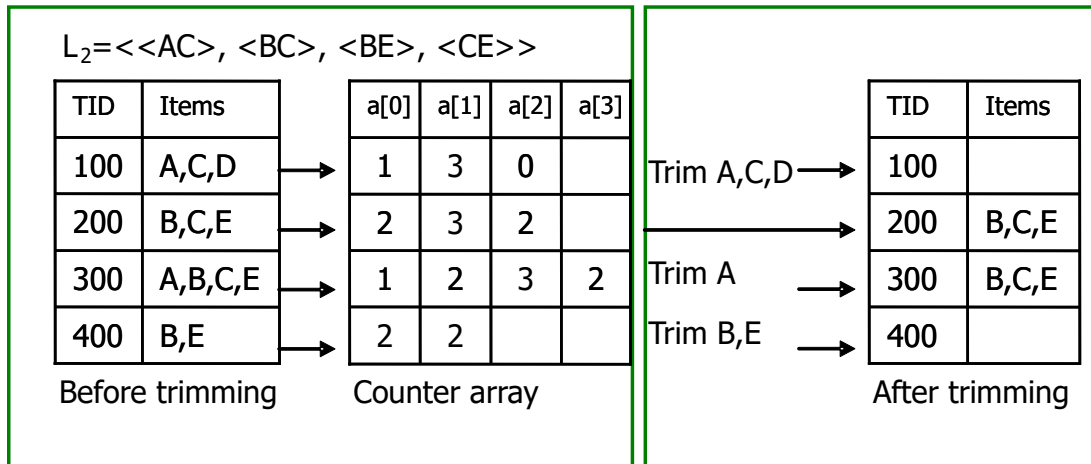


Hash table building

Candidate pruning

22

Transaction Items Pruning (from L2 to L3)



Trimming information collecting
(A appears in L2 once, C in L2 3 times, D not in L2)

Transaction trimming

23

A misleading “strong” association rule

- 10000 transactions
 - 6000 of them included computer games.
 - 7500 of them included video.
 - 4000 of them included computer games and video.
- Minimum support: 30%, minimum confidence: 60%

$buys(\text{computer games}) \Rightarrow buys(\text{videos})$
[support = 40%, confidence = 66%]
- However, $P(\{\text{video}\}) = 0.75$

From Association Analysis to Correlation Analysis

- The support and confidence measures are insufficient at filtering out uninteresting association rules.

$$A \Rightarrow B[\textit{support}, \textit{confidence}, \textit{correlation}]$$

Lift

- The **lift** between the occurrence of A and B can be measured by computing

The probability of a transaction contains the *union* of sets A and B.
↑
It doesn't mean $P(A \text{ or } B)$.

$$\textit{lift}(A, B) = \frac{P(A \cup B)}{P(A)P(B)} = \frac{P(B|A)}{P(B)} = \frac{\textit{conf}(A \Rightarrow B)}{\textit{sup}(B)}$$

< 1, negatively correlated

> 1, positively correlated

= 1, no correlation (A and B are independent)

- **Lift** assesses the degree to which the occurrence of one “lifts” the occurrence of the other.

Interestingness Measure: Correlations (Lift)

- *play basketball* \Rightarrow *eat cereal* [40%, 66.7%] is misleading
 - The overall % of students eating cereal is 75% > 66.7%.
- *play basketball* \Rightarrow *not eat cereal* [20%, 33.3%] is more accurate, although with lower support and confidence
- Measure of dependent/correlated events: [lift](#)

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

$$lift(B, C) = \frac{2000/5000}{3000/5000 * 3750/5000} = 0.89$$

$$lift(B, \neg C) = \frac{1000/5000}{3000/5000 * 1250/5000} = 1.33$$

	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

27

Generalized Association Rules

- Given the class hierarchy (taxonomy), one would like to choose proper data granularities for mining.
- Different confidence/support may be considered.

28

<pre> graph TD Clothes --> Outerwear Clothes --> Shirts Outerwear --> Jackets Outerwear --> SkiPants[Ski Pants] Footwear --> Shoes Footwear --> HikingBoots[Hiking Boots] </pre>		Freq. Itemset	Itemset support
Database		Jacket	2
		Outerwear	3
		Clothes	4
		Shoes	2
		Hiking Boots	2
		Footwear	4
		Outerwear, Hiking Boots	2
		Clothes, Hiking Boots	2
		Outerwear, Footwear	2
		Clothes, Footwear	2
Tx	Items bought		
100	Shirt		
200	Jacket, Hiking Boots	Outerwear → Hiking Boots	sup(30%) conf(60%) 33% 66%
300	Ski Pants, Hiking Boots	Outerwear → Footwear	33% 66%
400	Shoes	Hiking Boots → Outerwear	33% 100%
500	Shoes	Hiking Boots → Clothes	33% 100%
600	Jacket	However, Jacket → Hiking Boots	16% 50%
		Ski Pants → Hiking Boots	16% 100%

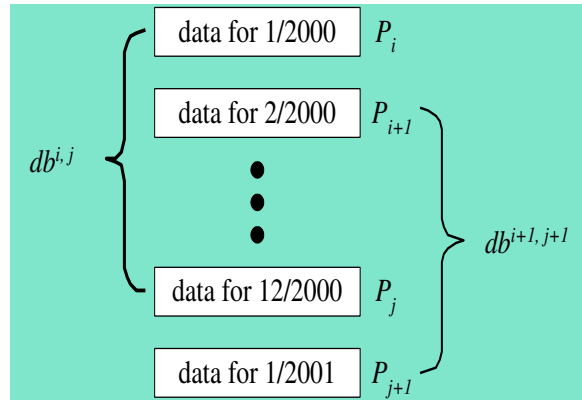
29

Incremental Mining

- Due to the increasing use of the record-based databases, recent important applications have called for the need of incremental mining
 - Such applications include Web log records, stock market data, grocery sales data, transactions in electronic commerce, and daily weather/traffic records, to name a few

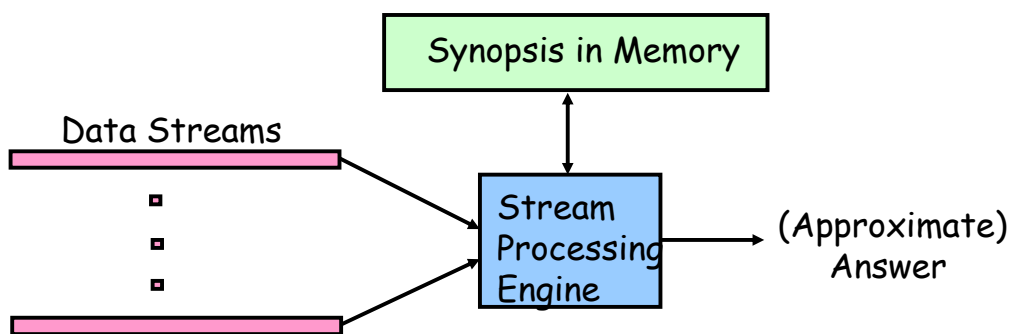
Incremental Mining

- To mine the transaction database for a fixed amount of most recent data (say, data in the last 12 months)
- One has to not only include new data (i.e., data in the new month) into, but also remove the old data (i.e., data in the most obsolete month) from the mining process.
- Google for “Incremental mining”



31

Data Streams: Computation Model



- Stream processing requirements
 - **Single pass:** Each record is examined at most once
 - **Bounded storage:** Limited Memory for storing synopsis
 - **Real-time:** Per record processing time must be low

32

Big Data (Recalled)

- Big data refers to a huge amount of data which is fast accumulated from various sources.
 - **Volume** - Scale from terabytes to zettabytes
 - **Variety** - Relational and non-relational data types from an ever-expanding variety of sources (e.g., IOT, Social Networks, Multimedia applications)
 - **Velocity** - Streaming data and fast movement of large volume data
- **Veracity**
- **Value**

33

Related Papers (just some)

- The easiest way is to use **Google**
- R. Agrawal and R. Srikant, "Fast Algorithms for Mining Association Rules in Large Databases," Proc. 20th Int'l Conf. Very Large Data Bases, pp. 478-499, Sept. 1994.
- M.-S. Chen, J. Han, and P. S. Yu. Data mining: An overview from a database perspective. IEEE Trans. Knowledge and Data Engineering, 8:866-883, 1996.
- J. Park, M.-S. Chen, and P. S. Yu. Using A Hash-Based Method with Transaction Trimming for Mining Association Rules," IEEE Trans. on Knowledge and Data Eng., vol. 9, no. 5, pp. 813-825, Sept./Oct. 1997.

34

About reading papers

- Not the same as reading textbook (important!)
- First, know who, why, where and when the paper was written
- What the problem is solved and how (in a high level)
- What is the primary contribution (method proposed, new finding, etc)
- How this paper was cited later?
- Anything that I should/could do with a **stretched** effort