

# Fuzzy Weighted Association Rule Mining with Weighted Support and Confidence Framework

Muhammad Sulaiman Khan<sup>1</sup>  
Dr Maybin Mueyba<sup>1</sup>  
Dr Frans Coenen<sup>2</sup>

<sup>1</sup>Liverpool Hope University  
<sup>2</sup>Liverpool University

ALSIP May 20, 2008 - Osaka Japan

## Outline of the Presentation

Organised as follows:

- Introduction
  - Classical Association Rule Mining (ARM)
  - Fuzzy Association Rule Mining
  - Downward Closure Property (DCP)
  - Weighted Association Rule Mining
    - Weighted ARM Support & Confidence
    - Classical WARM Types
    - Classical WARM Issues
- Our Contribution
  - Binary Weighted Association Rule Mining (BWARM)
  - Fuzzy Weighted Association Rule Mining (FWARM)
- Problem Definition, Methodology & Application
- FWARM Algorithm
- Experimental Results
- Conclusion & Further work

ALSIP May 20, 2008 - Osaka Japan

## Association Rule Mining

- **Association Rule Mining (ARM)**
  - Data Mining Technique for rule extraction
  - Typically used to determine customer buying Patterns from **large** market basket data/Transactions.
  - Association rules are expressions of the form

$$X \rightarrow Y$$

where X and Y are item sets and  $X \cap Y = \emptyset$

- Classical ARM utilises binary/boolean data

ALSIP May 20, 2008 - Osaka Japan

## Classical ARM (Boolean)

- **Support**

$$\text{Supp}(X \rightarrow Y) = \text{Supp}(X \cup Y)$$

- **Confidence**

$$\text{Conf}(X \rightarrow Y) = \text{Supp}(X \cup Y) / \text{Supp}(X),$$

a conditional probability i.e. Given X, we can determine Y.

ALSIP May 20, 2008 - Osaka Japan

## Fuzzy ARM (Quantitative)

- **Fuzzy ARM**
  - Applies to non-boolean and relational databases with quantitative attributes.
  - Determine rules of the form:

if <X is A> then <Y is B>

where X and Y are attributes in a database and A and B are the discretised values of these attributes.

For example:

if <Age is Young> then <Salary is Low>

ALSIP May 20, 2008 - Osaka Japan

## Downward Closure Property (DCP)

- **DCP**
  - In a classical Apriori algorithm it is assumed that if the itemset is large, then all its subsets should also be large and is called Downward Closure Property (DCP).
  - e.g. if {A,B,C} is a frequent set then {A,B}, {A,C} and {B,C} will also be frequent.
- **Application**
  - DCP helps to generate large itemsets of increasing size by adding items to itemsets that are already large.
  - For example, if itemset {A, B} and {B, C} are not frequent, then {A, B, C} and {B, C, D} cannot be frequent so we don't consider generating the supersets that contain non-frequent subsets.
  - Thus, DCP plays important role in efficient ARM algorithm development.

ALSIP May 20, 2008 - Osaka Japan

# Weighted ARM

- Weighted association rule mining process deals with the significance/importance of individual items in a database
- Items are assigned weights (w) according to their significance as shown in table below.
- These weights are set according to items profit margins

Table1: Weighted Items Database

ID	Item	Profit	Weight	...
1	Scanner	10	0.1	...
2	Printer	30	0.3	...
3	Monitor	60	0.6	...
4	Computer	90	0.9	...

Table2: Transactions

TID	Item
1	1, 2, 4
2	2, 3
3	1, 2, 3, 4
4	2, 3, 4

# Weighted ARM

- Classical association rule mining (ARM) model assumes that all items have the same significance without taking into account their weight within a transaction or record.

For example rules:

**A:** [wine → salmon, 1%, 80%] may be more important than  
**B:** [bread → milk, 3%, 80%]

In classical ARM rule **B** is more important than rule **A** because rule **B** has higher support than rule **A**.

But in weighted ARM with weighted settings rule **A** may be more important than rule **B**, even though the former holds a lower support.

This is because those items in the first rule usually come with more profit per unit sale, but the standard ARM simply ignores this difference.

# Classical WARM Types

**Post Processing:** where all frequent sets are generated using any classical ARM algorithm and later aggregated weights of itemsets are multiplied with their supports in order to find the weighted support.

**Pre Processing:** In pre processing weighted support is calculated after each step (database scan) using the same formula used in post processing. The only obvious reason for pre processing seems to do early pruning and avoid un-necessary generation of frequent sets.

**Comment:** postprocessing may be inefficient (generates all sets) pre processing may be more efficient, may also miss out interesting rules (under DCP in weighted settings, this is solved) i.e. Lemma given in paper: If an itemset is not frequent, then its superset can not be frequent

# Classical WARM Issues

- Many potential itemsets are not considered (considering their weights) during the mining step because both approaches (pre-, post) use classical ARM algorithms which normally follows **DCP** and the only items considered for weighted support are the already generated frequent sets.
- Itemsets are first generated using their occurrences in the data base and later their weights are considered for weighted support. This approach leads to loose many potential itemsets which could be important if their weights are considered as well as occurrence.
- After calculating weighted support, the weighted frequent sets do not hold **DCP**.
- Exhaustive search is not possible with high number of items due to computational limitation  $2^n - 1$ , n is the number of items

# Our Contribution: Weighted Support & Confidence framework

- We address the Weighted Association Rule Mining issues present in previous approaches and have proposed a **Weighted Support & Confidence framework** for databases with boolean and quantitative/fuzzy attributes.
- In our framework we consider items occurrences and weights together instead of just their occurrences (pruning process) in the database to calculate their support and confidence.
- Thus, our proposed framework reflects not only number of records supporting the itemsets, but also their degree of significance in the dataset.
- Frequent itemset generated using our approach holds valid **DCP**.

# Problem Definition 1

Let the input data **D** have transactions  $T = \{t_1, t_2, \dots, t_n\}$  with a set of items  $I = \{i_1, i_2, \dots, i_m\}$  and a set of non-negative, real number weights  $W = \{w_1, w_2, \dots, w_m\}$  associated with each item. Each  $t^j$  transaction is some subset of **I** and a weight **w** is attached to each item  $t_i^j$  (the " $i^j$ " item in the " $t^j$ " transaction).

Thus each item  $i_j$  will have associated with it a weight corresponding to the set **W**, i.e. a pair  $(i_j, w)$  is called a weighted item where  $i_j \in I$ . Weight for the " $i^j$ " item in the " $t^j$ " transaction is given by  $t_i^j [w]$ .

**Item Weight  $IW$**  is a real value given to each item  $i_j$  ranging [0..1] with some degree of importance, a weight  $i_j [w]$ .

**Itemset Transaction Weight  $ITW$**  is the aggregated weights of all the items in the itemset present in a single transaction. Itemset transaction weight for an itemset **X** can be calculated as:

$$\text{vote for } t_j \text{ satisfying } X = \prod_{k=1}^{|X|} (v_{\{t_i^j \in X\}} t_i^j [w])$$

## Binary Weighted Support & Confidence framework

**Binary Weighted Support** is the aggregated sum of itemset transaction weight of all the transactions in which itemset is present, divided by the total number of transactions. It is denoted as:

$$WS(X) = \frac{\sum_{i=1}^n \prod_{k=1}^{|X|} (\forall [i|w] \in X) t_i [i_k [w]]}{n}$$

**Binary Weighted Confidence** is the ratio of sum of votes satisfying both  $X \cup Y$  to the sum of votes satisfying  $X$  (with  $Z = (X \cup Y)$ ). It is formulated as:

$$WC(X \rightarrow Y) = \frac{\sum_{i=1}^n \prod_{k=1}^{|Z|} (\forall [z|w] \in Z) t_i [z_k [w]]}{\prod_{k=1}^{|X|} (\forall [i|w] \in X) t_i [x_k [w]]}$$

ALSP May 20, 2008 - Osaka Japan

## Fuzzy Weighted Support & Confidence framework

**Fuzzy Weighted Support FWS** is the aggregated sum of **FTW** of all the transactions in which itemset is present, divided by the **TOTAL** number of transactions. It is calculated as:

$$FWS(X) = \frac{\sum_{i=1}^n \prod_{k=1}^{|X|} (\forall [i|w] \in X) t'_i [i_k [w]]}{n}$$

**Fuzzy Weighted Confidence FWC** is the ratio of sum of votes satisfying both  $X \cup Y$  to the sum of votes satisfying  $X$  (with  $Z = (X \cup Y)$ ). It is formulated as:

$$FWC(X \rightarrow Y) = \frac{\sum_{i=1}^n \prod_{k=1}^{|Z|} (\forall [z|w] \in Z) t'_i [z_k [w]]}{\prod_{k=1}^{|X|} (\forall [i|w] \in X) t'_i [x_k [w]]}$$

ALSP May 20, 2008 - Osaka Japan

## Lemma

If an itemset is not frequent then its superset cannot be frequent and is always true

$WS(subset) \geq WS(superset)$

**Proof**

Let  $X$  be a subset of  $Y$ . Then  $WS(X) \geq WS(Y)$ . This is because the numerator of  $WS(X)$  is the sum of weights of all transactions containing  $X$ , and the numerator of  $WS(Y)$  is the sum of weights of all transactions containing  $Y$ . Since  $X$  is a subset of  $Y$ , any transaction containing  $Y$  also contains  $X$ . Therefore, the numerator of  $WS(Y)$  is a subset of the numerator of  $WS(X)$ . The denominator of both  $WS(X)$  and  $WS(Y)$  is the total number of transactions,  $n$ . Hence,  $WS(X) \geq WS(Y)$ .

**Note:** WS is greater with fewer items because of product operator e.g. in  $WS(tx)$  against  $WS(ty)$ , where  $tx$  is a subset of  $ty$

ALSP May 20, 2008 - Osaka Japan

## Proof - DCP Holds

Itemset	Support	Weights				
		A	C	E	B	D
A	0.1	0.1	0.1	0.1	0.1	0.1
C	0.1	0.1	0.1	0.1	0.1	0.1
E	0.1	0.1	0.1	0.1	0.1	0.1
B	0.1	0.1	0.1	0.1	0.1	0.1
D	0.1	0.1	0.1	0.1	0.1	0.1

min. sup	min. conf	min. sup	min. conf	min. sup	min. conf	min. sup	min. conf
0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1

## FWARM Algorithm

- The proposed Fuzzy Weighted ARM (FWARM) algorithm is required to be efficient as processing of fuzzy weights
- The algorithm uses tree data structure [13] and similar in fashion to Apriori algorithm (with valid DCP).
- In Algorithm:
  - $C_k$  is the set of candidate itemsets of cardinality  $k$
  - $w$  is the set of weights associated to items  $I$ .
  - $F$  is the set of frequent item sets.
  - $R$  is the set of potential rules.
  - $R'$  is the final set of generated fuzzy ARs

```

Input:
I = item set
w = item weight
WS = weighted support
WC = weighted confidence
Output:
R' = Set of Weighted ARs
1.  $k = 1, C_k = I, F_k = I$ 
2.  $C_k = \text{Set of item sets}$ 
3.  $k \leq |I|$ 
4. Loop
5.  $\# C_k = |C_k|$ 
6. Add-weighted support values to  $C_k$ 
7.  $\forall i \in C_k$ 
   a.  $\text{weightedSupport} = \text{weighted support count}$ 
   b.  $\text{weightedConfidence} = \text{weighted confidence ratio}$ 
8.  $\text{weightedConfidence} > \text{min\_conf} \implies F_k \cup \{i\}$ 
9.  $k = k + 1$ 
10.  $C_{k+1} = \text{generateCandidate}(F_k)$ 
11. End Loop
12.  $\forall i \in F$ 
13. generate set of candidate rules  $\{r_1, \dots, r_n\}$ 
14.  $R = R \cup \{r_1, \dots, r_n\}$ 
15.  $\forall i \in R$ 
16.  $\text{weightedConfidence} = \text{weighted confidence ratio}$ 
17.  $\text{weightedConfidence} > \text{min\_conf} \implies R' \cup \{i\}$ 

```

ALSP May 20

FWARM Algorithm

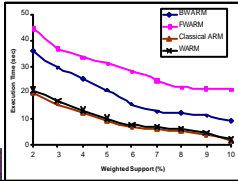
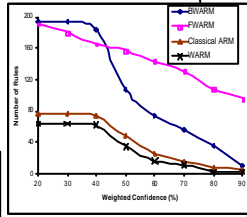
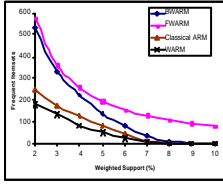
## Experimental Results

Experimental data is a transactional database containing 100K records and 1K (1000) items. Two sets of experiments were undertaken:

- First experiment show that the output behaviour of our proposed framework is quite similar to classical ARM because we use the Apriori approach in our algorithm but results are better than WARM. Experiments show
  - (i) the number of frequent sets generated (using weighted support measure) and
  - (ii) the number of rules generated (using weighted confidence measure)
- Experiment two shows comparison of execution times using different weighted supports and data sizes.

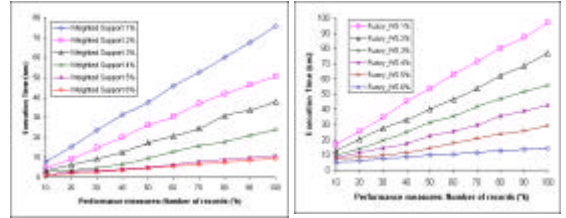
ALSP May 20, 2008 - Osaka Japan

# 1. Quality Measures



Osaka Japan

# 2. Performance Measures



Execution time, vary threshold 0.01 to 0.06, conf =0.5  
Algorithm scales linearly with FWS, similar to classical ARM

ALSIP May 20, 2008 - Osaka Japan

# Conclusion & Further Work

- We have presented novel approach for extracting hidden information from Weighted items, our proposed framework reflects not only number of records supporting the itemsets, but also their degree of significance in the dataset.
- The problem of invalidation of downward closure property (DCP) is solved using improved model of weighted support and confidence framework for binary and fuzzy association rule mining
- We showed the application of our method on different datasets.
- Overall, the approach presented here is effective and efficient for analysing databases with Weighted items (boolean or quantitative/fuzzy).

## Future Work

- There is potential to apply FWARM to other applications with Composite attributes even with varying fuzzy sets between attributes.
- Different measures for validating DCP, normalisation of values e.t.c.
- We intend to extend our work by integrating Weighted ARM with Utility ARM in transactional and relational databases containing weights and utilities of items. It is both a new approach and algorithmically challenging.
- Mechanisms for tuning membership degrees of itemsets without bias or human expert involvement e.g. genetic algorithms

ALSIP May 20, 2008 - Osaka Japan