

A Framework for Mining Fuzzy Association Rules from Composite Items

Muhammad Sulaiman Khan¹
 Dr Maybin Muyebe¹
 Dr Frans Coenen²

¹Liverpool Hope University
²Liverpool University

ALSIP May 20, 2008 - Osaka Japan

Outline of the Presentation

Organised as follows:

- Introduction
 - Classical Association Rule Mining (ARM)
 - Quantitative Association Rule Mining
 - Fuzzy Association Rule Mining (FARM)
- Background & Related Work
- Problem definition
- Methodology & Application
- CFARM Algorithm
- Experimental Results
- Conclusion & Further work

ALSIP May 20, 2008 - Osaka Japan

Introduction

- **Association Rule Mining**
 - Data Mining Technique
 - Determine customer buying Patterns from market basket data/Transactions.
 - Association rules are of the form

$$X \rightarrow Y$$

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

ALSIP May 20, 2008 - Osaka Japan

Association Rule Mining

• Support

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

• Confidence

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

• Downward Closure Property (DCP)

- Subsets of a frequent set are also frequent, e.g. if {A,B,C} is a frequent set then {A,B}, {A,C} and {B,C} will also be frequent.

ALSIP May 20, 2008 - Osaka Japan

Quantitative Association Rule Mining

• Quantitative ARM

- Applies on non-boolean and relational databases
- 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

Quantitative Association Rule Mining (Fuzziness)

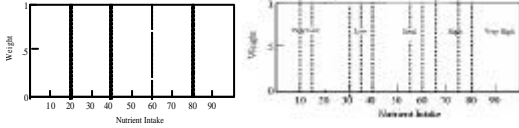
RecordID	Age	Income	Age	Intervals	Income	Intervals
100	31	25000	0-25	Young	0-20K	Low
200	16	10000	26-40	Middle	10K-30K	Medium
300	45	60000	41-60	Old	30K-45K	High
400	-	-	60+above	Very Old	45K+above	Very High

YID	Young	Middle	Old	Y-Old	Low	Medium	High	Y-High
100	0	1	0	0	0	1	0	0
200	1	0	0	0	1	0	0	0
300	0	0	1	0	0	0	0	1
400	-	-	-	-	-	-	-	-

ALSIP May 20, 2008 - Osaka Japan

Fuzzy Association Rule Mining (FARM)

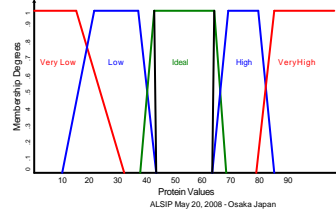
- Issues from partitioning quantitative attributes
 - Interval Partitions give boundary Problems to quantitative attributes using usual support measure (Supp)
 - To resolve, borders can overlap but this might over emphasise some intervals (Fuzziness is evident).
 - Example: Nutrition content of food – a protein or vitamin etc.



ALSIP May 20, 2008 - Osaka Japan

Fuzzy Association Rule Mining (FARM)

- Fuzzy sets used to resolve this by Providing a smooth change between boundaries using normalisation process (see later!).
- Firstly, fuzziness is defined by a membership mapping function $\mu(x) : A \rightarrow [0,1], x \in A$
- Trapezoidal membership function can be derived (example above)



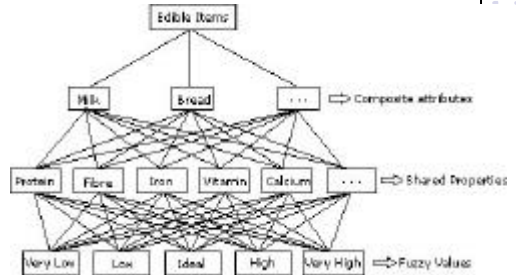
ALSIP May 20, 2008 - Osaka Japan

FARM using composite attributes

- FARM extended to composite Attributes
- Composite Attributes
 - Objects with different Properties
 - Properties can be quantitative and categorical
 - Attributes share the same Properties with other attributes
 - Quantitative Properties can be discretised into several ranges (fuzzy sets)
 - ARM usually on all data sets (edible, non-edible)

ALSIP May 20, 2008 - Osaka Japan

FARM using composite attributes



ALSIP May 20, 2008 - Osaka Japan

Background & Related Work

- From literature the only works that mention "Composite ARM" was Presented in 1997 & 2006.
 - Ye, X.; Keane, J.A. "Mining association rules with composite items", IEEE International Conference on Systems, Man, and Cybernetics, PP 1367 – 1372, Oct 1997.
 - Ke Wang, James N. K. Liu, Wei-min Ma, "Mining the Most Reliable Association Rules with Composite Items," *icdmw*, PP. 749-754, Sixth IEEE International Conference on Data Mining - Workshops (ICDMW'06), 2006.
- In these works, Composite items expressed as combination of several attributes
 - Composite means put B and C to make a new item {BC}
 - If itemset {A,B} and {B,C} are not frequent (large)
 - B→A and C→A will not be generated
 - Rule {BC}→{A} may be generated.

ALSIP May 20, 2008 - Osaka Japan

Problem definition

- Given a Dataset D consist of set of transaction $t = \{t_1, t_2, t_3, \dots, t_n\}$, a set of composite items $I = \{i_1, i_2, i_3, \dots, i_m\}$ and a set of properties $P = \{p_1, p_2, p_3, \dots, p_m\}$.
- Each transaction t_i is subset of I , and each item $t_{[i,j]}$ is a subset of P .
- Thus each item i_j will have associated with it a set of numeric values corresponding to the set P , i.e. $t_{[i,j]} = \{v_1, v_2, v_3, \dots, v_m\}$.
- The "kth" property value for the "jth" item in "ith" transaction is given by $t_{[i,j]}[v_k]$.

ALSIP May 20, 2008 - Osaka Japan

Problem Definition

- Example

TID	Record
1	{<a,{2,4,6}>, <b,{4,5,3}>}
2	{<c,{1,2,5}>, <d,{4,1,3}>}
3	{<a,{2,4,6}>, <c,{1,2,5}>, <d,{4,1,3}>}
4	{<b,{4,5,3}>, <d,{4,1,3}>}

$D = \{t_1, t_2, t_3, t_4\}$
 $I = \{a, b, c, d\}$
 $P = \{x, y, z\}$

Problem Definition

- Property Dataset

- D is initially transformed into a Property dataset D^p .
- D^p consists of Property transactions $T^p = \{t_1^p, t_2^p, t_3^p, \dots, t_n^p\}$.
- For each transaction t_i^p , is subset of $P = \{p_1, p_2, p_3, \dots, p_n\}$.
- The value for each Property attribute $t_i^p[P_j]$ is obtained by aggregating the numeric values for all p_j in t_i . Thus

$$t_i^p[P_j] = \frac{\sum_{k=1}^{|t_i|} t_i^p[i_j][v_k]}{|t_i|}$$

Problem Definition

- Fuzzy Dataset

- D^p is further transformed into a fuzzy dataset D^f .
- A fuzzy dataset D^f consists of fuzzy transactions $T^f = \{t_1^f, t_2^f, t_3^f, \dots, t_n^f\}$ and fuzzy property attributes P^f .
- Each P^f has a number of fuzzy sets associated with it, identified by a set of linguistic labels $L = \{l_1, l_2, l_3, \dots, l_k\}$ e.g. {small, medium, large}.
- Each property attribute $t_i^f[P_j]$ is associated (to some degree) with several fuzzy sets, with a membership degree in the range [0,1].
- Membership degree indicates the correspondence between the value of a given $t_i^f[P_j]$ and the set of fuzzy linguistic labels.

Problem Definition

- Composite Item Value Table

- A composite item value table is table that allows us to get property values for specific items.

- Properties Table

- A properties table is a table that maps all possible values for each property attribute $t_i^f[P_j]$ onto fuzzy/overlapped ranges, e.g. [0-10], [7-19], [15-30]

Problem definition

- Fuzzy Normalisation Process

- The process of finding the contribution to the fuzzy support value, m_i^f , for individual property attributes $t_i^f[p_j][l_k]$ such that a partition of unity is guaranteed.

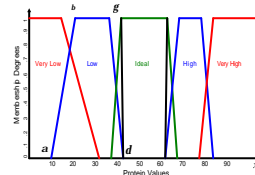
$$t_i^f[p_j][l_k] = \frac{m(t_i^f[p_j][l_k])}{\sum_{x=1}^{|L|} m(t_i^f[p_j][l_x])}$$

TID	VL	L	MD	H	VH	...
1	0.0	0.0	0.0	1.0	0.2	...
2	0.3	0.38	0.0	0.0	0.0	...
3

Problem Definition

- Membership Function (Trapezoidal)

- Membership degree to a particular fuzzy set (described by linguistic label), $t_i^f[p_j][l_k]$ is determined by a membership function $m(x): A \rightarrow [0,1], x \in A$



$$m(x, a, b, g, d) = \begin{cases} 0, & d < x < a \\ \frac{(x-a)}{(b-a)}, & a \leq x \leq b \\ \frac{(d-x)}{(d-g)}, & g \leq x \leq d \\ 1, & b < x < g \end{cases}$$

Problem definition

- Fuzzy Support

- Fuzzy support is calculated as

$$FS(A) = \frac{\text{Sum of votes satisfying } A}{\text{Number of records in } T'}$$

votes for t_i satisfying $A = \sum_{j=1}^m \prod_{l \in A} t'_{ij}[l_k]$

$$FS(A) = \frac{\sum_{i=1}^m \prod_{l \in A} t'_{ij}[l_k]}{n}$$

ALSIP May 20, 2008 - Osaka Japan

Problem definition

- Fuzzy Confidence

- Fuzzy confidence (FC) is calculated in the same manner that confidence is calculated in traditional ARM.
- Fuzzy confidence is calculated as:

$$FC(A \rightarrow B) = \frac{FS(A \cup B)}{FS(A)}$$

ALSIP May 20, 2008 - Osaka Japan

Problem definition

- Fuzzy Correlation Measure

- Fuzzy Confidence does not take into account $FS(B)$.
- The Fuzzy Correlation ($FCORR$) addresses this.
- Similar to statistical correlation but different in meaning.

$$FCORR(A \rightarrow B) = \frac{Cov(A, B)}{\sqrt{Var(A) \times Var(B)}}$$

$$Cov(A, B) = E[C] - E[A] \times E[B] \quad \text{where } C=A \cup B$$

$$Cov(A, B) = FS(C) - FS(A) \times FS(B)$$

ALSIP May 20, 2008 - Osaka Japan

Problem Definition

The variance of A and B can be obtained as follows:

$$Var(A) = E[A^2] - E[A]^2$$

$$Var(B) = E[B^2] - E[B]^2$$

where

$$E[X]^2 = FS(X)^2 = \left(\frac{\sum_{i=1}^m \prod_{l \in X} t'_{ij}[l_k]}{n} \right)^2$$

$$E[X^2] = FS(X^2) = \frac{\sum_{i=1}^m \prod_{l \in X} (t'_{ij}[l_k])^2}{n}$$

ALSIP May 20, 2008 - Osaka Japan

Proposed Methodology

- Data Transformation

- Transformation of raw dataset T into property dataset T^P .
- Transformation of property dataset T^P into a database containing fuzzy extensions T' .
- Normalization of fuzzy dataset.
- Candidate Generation i.e. search for all fuzzy frequent itemsets that have support higher than user specified threshold.
- Use frequent itemsets to generate all possible rules using fuzzy confidence or fuzzy correlation interestingness measures.

ALSIP May 20, 2008 - Osaka Japan

Proposed Methodology

TID	Items	TID	It	Ic	Ca	-
1	X,Z	1	53	41	77	-
2	Z	2	27	19	30	-
3	X,Y,Z	3	99	2	67	-
4	-	4	-	-	-	-

TID	VL	L	Ideal	H	VH	VL	L	Ideal	H	VH	-
1	0.0	0.0	0.9	0.1	0.0	0.0	0.4	0.6	0.2	0.8	-
2	0.4	0.6	0.0	0.0	0.0	0.5	0.5	0.0	0.0	0.0	-
3	0.0	0.0	0.0	0.7	0.3	0.0	0.0	0.3	0.7	0.0	-
4	-	-	-	-	-	-	-	-	-	-	-

ALSIP May 20, 2008 - Osaka Japan

CFARM Algorithm

Properties Converter(T, RDA, I_{data})

1. $\forall T$ in D
2. $\forall I$ in \mathcal{I}
3. $T' = T' \cup P$
4. $T' = \text{average Contribution}(T)$
5. $D' = \text{write}(T')$
6. end;

Fuzzy Converter(T' , Properties Table)

1. $\forall T'$ in D'
2. $\forall P$ in T'
3. $\text{fuzzyattr} = \text{getFuzzyAttr}(P, \text{propTable})$
4. $T = T' \cup (\text{fuzzyattr})$
5. $D = \text{write}(T)$
6. end;

CFARM($\text{minsup}, \text{minconf}, \text{mincorr}, T'$)

1. $k=0; C_k = f = f_k = f$
2. do
3. $k=k+1$
4. $\beta[k]=I$
5. $C_k = \text{GenerateFirstCandidates}(T')$
6. else
7. $C_k = \text{GenerateCandidates}(F_{k-1})$
8. $\forall C_k$
9. $\text{count} = \text{CountSupport}(C_k)$
10. $C_k = \text{PruneCandidates}(C_k, \text{count}, \text{minsup})$
11. $C_k = \text{CalculateSignificance}(C_k, \text{minconf})$
12. $F_k = \text{GenerateFrequentItemsets}(C_k, \text{minconf})$
13. $F = F \cup F_k$
14. while($C_k.\text{count} > k$)
15. $\text{FuzzyCorr} = \text{Calculate}(F, \text{mincorr})$
16. Output (Rules($F, \text{mincorr}, \text{FuzzyCorr}$))

Example Application

Table with 10 columns and 10 rows showing transaction data for items like Bread, Butter, Milk, etc.

Three smaller tables showing intermediate results of the CFARM algorithm, such as candidate sets and frequent itemsets.

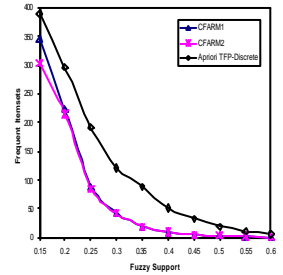
Example Application

Table showing the output of the CFARM algorithm, listing rules and their associated metrics like support, confidence, and correlation.

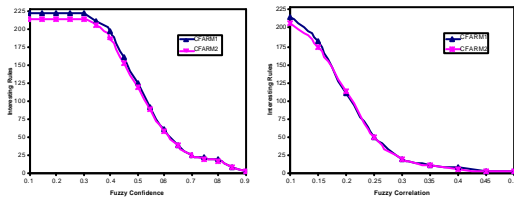
Experimental Results

Quality Measures

- We use all standard 27 nutrients present in edible items.
- We use T1014D100K dataset with 100K transactions.
- Fuzzy Support 30%
- Fuzzy Confidence 50%
- Fuzzy Correlation 25%



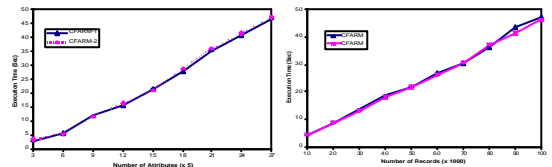
Experimental Results



Experimental Results

Performance Measures

We set a support threshold to 0.30, confidence 0.5 and correlation value to 0.25.



Experimental Results



- Some interesting fuzzy rules produced by our approach with 30% support, 50% confidence and 25% correlation are as follows:
 - IF *Protein* intake is *Ideal* THEN *Carbohydrate* intake is *low*.
 - IF *Protein* intake is *Low* THEN *Vitamin A* intake is *High*.
 - IF *Protein* intake is *High* AND *Vitamin A* intake is *Low* THEN *Fat* intake is *High*.
- Depending on expert analysis from a health practitioner, these rules are useful in analysing customer buying behavior concerning their nutrition.

Conclusion & Further Work



- We have presented a novel approach for extracting hidden information from composite items.
- We showed that with such items, common properties can be defined as quantitative itemsets themselves, which are transformed into fuzzy sets.
- We use fuzzy correlation measure as interestingness measure and showed more interesting rules.
- Overall, the approach presented here is effective and efficient for analysing databases with composite items.
- Further work will evaluate our approach on real and larger datasets and compare real performance with other common fuzzy ARM algorithms.
- There is potential to apply this to other applications with composite items or attributes even with varying fuzzy sets between attributes e.g. image analysis, inventory control database.
- We are expanding our work with the possibilities to extend it for Fuzzy Utility Association Rule Mining.