



## Discovering Complex Matchings across Web Query Interfaces: A Correlation Mining Approach

Bin He, Kevin Chen-Chuan Chang, Jiawei Han

Presented by: Maryam Karimzadehgan  
mkarimzadehgan@cs.uwaterloo.ca

---

---

---

---

---

---

---

---

## Outline

- Motivation of integrating the deep web
- Mining Algorithm
- H-Measure
- Data Preparation Step
- Experiments
- Conclusion and Discussion

2

---

---

---

---

---

---

---

---

## Motivation: Matching $\rightarrow$ Mining

- Group attributes  
Larger concept e.g. { adults, seniors, children, infants} denotes the number of passengers.
- Synonym relationship  
Different sources may use different attributes for the same concept. e.g. {from} = {depart}, {to} = {destination}

Group attribute + Synonym attribute  $\rightarrow$  Complex Matching  
 $m:n \rightarrow m:n$  matching  
e.g. {adults, seniors, children, infants} = {passengers} 4:1 matching

3

---

---

---

---

---

---

---

---

# Correlation Mining

- Match more than two attribute groups

{adults, seniors, children, infants} = {passengers} = {number of tickets}  
4:1:1

*n-ary complex matching → aggregation of several binary m:n matchings*

- DCM (Dual Correlation Mining) framework for mining *n*-ary complex matchings

4

---

---

---

---

---

---

---

---

# Outline

- Introduction
- Motivation of integrating the deep web
- Mining Algorithm
- H-Measure
- Data Preparation Step
- Experiments
- Discussion

5

---

---

---

---

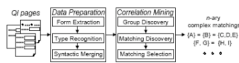
---

---

---

---

# DCM Framework (Dual Correlation Mining)



### Data Preparation:

schema transaction → mining by extracting the attribute entities from query interfaces

### Correlation Mining:

finds complex matchings considering both positive and negative correlations

6

---

---

---

---

---

---

---

---

## Mining Algorithm

- **Group discovery (individual attribute)**
  - Mining positively correlated attribute --> potential attribute group
  - Potential group is not suitable for matching --> synonym relationship (negative correlation)
- **Matching discovery (attribute groups)**
  - Given the potential group --> Mining negatively correlated attribute --> potential *n*-ary complex matching
- **Matching selection**
  - A potential matching may not be considered as correct due to existence of conflicts among matchings.
  - Select the consistent matchings from the mining result

7

---

---

---

---

---

---

---

---

---

---

## Complex Matching Discovery

- Negative correlation should exist between two groups. (e.g. {destination} = {to} = {arrival city})

$C_{\min}$  is the minimal value of the pairwise evaluation

- $C_{\min}$  satisfies the "Apriori" feature. e.g.  $C_{\min}(A, m) \leq C_{\min}(A', m)$   
 $\min\{1, 3, 5\} \leq \min\{3, 5\}$
- $C_{\min}$  can mine both positive and negative correlations.

- A set of attributes  $\{A_1, \dots, A_n\}$  is positively correlated if

$$C_{\min} \geq \text{Threshold}$$

- A set of attributes groups  $\{G_1, \dots, G_n\}$  is negatively correlated attribute groups if

$$C_{\min} \geq \text{Threshold}$$

8

---

---

---

---

---

---

---

---

---

---

## Algorithm N-ary (Discovering Complex Matching)

```

Algorithm: N-ARYSCHEMAMATCHING
Input: InputSchemas  $S_T = \{S_1, \dots, S_n\}$ ,
Measures  $m, m_n$ , Thresholds  $T_p, T_n$ 
Output: Potential n-ary complex matchings
begin
1 /* group discovery */
2  $G \leftarrow$  APRIORICORRMINING( $S_T, m_p, T_p$ )
3 /* adding groups into  $S_T^*$  */
4 for each  $S_i \in S_T^*$ 
5   for each  $G_k \in G$ 
6     if  $S_i \cap G_k \neq \emptyset$  then  $S_i \leftarrow S_i \cup \{G_k\}$ 
7 /* matching discovery */
8  $M \leftarrow$  APRIORICORRMINING( $S_T, m_n, T_n$ )
9 return  $M$ 
end
    
```

```

Algorithm: APRIORICORRMINING
Input: InputSchemas  $S_T = \{S_1, \dots, S_n\}$ ,
Measures  $m$ , Thresholds  $T$ 
Output: Correlated items
begin
1  $X \leftarrow \emptyset$ 
2  $V \leftarrow \bigcup_{i=1}^n S_i, S_i \in S_T$ 
3 for all  $A_p, A_q \in V, p \neq q$ 
4   if  $m(A_p, A_q) \geq T$  then  $X \leftarrow X \cup \{\{A_p, A_q\}\}$ 
5  $l \leftarrow 2$ 
6 /*  $X_l$ : correlated items with length =  $l$  */
7  $X_l \leftarrow X$ 
8 while  $X_l \neq \emptyset$ 
9   construct  $X_{l+1}$  from  $X_l$  using apriori feature
10   $X \leftarrow X \cup X_{l+1}$ 
11   $X_l \leftarrow X_{l+1}$ 
12 return  $X$ 
end
    
```

9

---

---

---

---

---

---

---

---

---

---

## Complex Matching Selection

- False semantic matching due to coincidental correlations
  - {author} = {first name, last name} M1 **correct one**
  - {subject} = {first name, last name} M2 **wrong one**

The more negatively correlated <sup>False ones should be removed</sup>  $\rightarrow$  higher confidence to be real synonyms

- Strategy for *ranking* the discovered matchings:
  - Score function to evaluate the discovered matchings under measure

The goal of qualifying  $\rightarrow$  correlation passes some threshold  
 The goal of ranking  $\rightarrow$  compare the strength of correlations

$C_{\max}$  The maximal measure value among pairs of groups in a matching

10

---

---

---

---

---

---

---

---

---

---

---

---

## Complex Matching Selection

- Strategy for tie breaking  $\rightarrow$  richer semantic information
  - Take "top K" approach
  - If  $C_{\max}$  value of two groups are the same  $\rightarrow$  compare their second highest to break the tie.
  - If two matching are tie after "top-k" comparison  $\rightarrow$  the one with richer semantic information.

- Rule for ranking matches

- If  $s(M_j, m_n) > s(M_k, m_n)$ ,
- If  $s(M_j, m_n) = s(M_k, m_n)$  and  $M_j \geq M_k$
- rank them arbitrary

11

---

---

---

---

---

---

---

---

---

---

---

---

## Complex Matching Selection

```

Algorithm: MATCHINGSELECTION
Input: Potential complex matchings  $M = \{M_1, \dots, M_n\}$ .
Measure  $m_n$ .
Output: Selected complex matchings
begin:
1  $R \leftarrow \emptyset$  /* selected n-ary complex matchings */
2 while  $M \neq \emptyset$ 
3   /* select the matching ranked the highest */
4    $M_i \leftarrow \text{GETMATCHINGRANKFIRST}(M, m_n)$ 
5    $R \leftarrow R \cup \{M_i\}$ 
6   for each  $M_j \in M$ 
7     /* remove the conflicting part */
8      $M_j \leftarrow M_j - M_i$ 
9     /* delete  $M_j$  if it contains no matching */
10    if  $|M_j| < 2$  then  $M \leftarrow M - \{M_j\}$ 
11 return R
end
    
```

```

Algorithm: GETMATCHINGRANKFIRST
Input: Potential complex matchings  $M = \{M_1, \dots, M_n\}$ .
Measure  $m_n$ .
Output: The matching with the highest ranking
begin:
1  $M_i \leftarrow M_1$ 
2 for each  $M_j \in M, 2 \leq j \leq n$ 
3   if  $s(M_j, m_n) > s(M_i, m_n)$  then
4      $M_i \leftarrow M_j$ 
5   if  $s(M_i, m_n) = s(M_j, m_n)$  and  $M_j \geq M_i$  then
6      $M_i \leftarrow M_j$ 
7 return  $M_i$ 
end
    
```

- Greedy selection strategy by choosing the highest ranked matching,  $M_i$ , in each iteration.

### Example

$M_1$  {author} = {last name, first name}, 0.95  
 $M_2$  {author} = {last name}, 0.92  
 $M_3$  {subject} = {category}, 0.92  
 $M_4$  {author} = {first name}, 0.90  
 $M_5$  {subject} = {last name, first name}, 0.88  
 $M_6$  {subject} = {last name}, 0.88 and  
 $M_7$  {subject} = {first name}, 0.86.

12

---

---

---

---

---

---

---

---

---

---

---

---

# Outline

- Introduction
- Motivation of integrating the deep web
- Mining Algorithm
- H-Measure
- Data Preparation Step
- Experiments
- Discussion

13

---

---

---

---

---

---

---

---

---

---

# Measures for Association Patterns

#	Measure	Formula
1	Support	$\sum_{i,j} P(A_i B_j)$
2	Goussinman-Kravtsov's (A)	$\frac{\sqrt{P(A)P(B)}(\sqrt{P(A B)} + \sqrt{P(B A)}) - \min_x P(A_x) - \max_x P(B_x)}{2 - \min_x P(A_x) - \max_x P(B_x)}$
3	Odds ratio (o)	$\frac{P(A, B)P(C, D)}{P(A, C)P(B, D)}$
4	Yule's Q	$\frac{P(A, B)P(C, D) - P(A, C)P(B, D)}{P(A, B)P(C, D) + P(A, C)P(B, D)}$
5	Yule's Y	$\frac{\sqrt{P(A, B)P(C, D)} - \sqrt{P(A, C)P(B, D)}}{\sqrt{P(A, B)P(C, D)} + \sqrt{P(A, C)P(B, D)}}$
6	Kaplan's (k)	$\frac{P(A, B)P(C, D) - P(A, C)P(B, D)}{P(A, B)P(C, D) + P(A, C)P(B, D)}$
7	Mutual Information (M)	$\sum_{i,j} P(A_i B_j) \log \frac{P(A_i B_j)}{P(A_i)P(B_j)}$
8	J-Measure (J)	$\min \left( \sum_{i,j} P(A_i B_j) \log \frac{P(A_i B_j)}{P(A_i)P(B_j)}, \sum_{i,j} P(A_i B_j) \log \frac{P(A_i B_j)}{P(A_j)P(B_i)} \right)$
9	Gini Index (G)	$\frac{P(A, B) \log \left( \frac{P(A, B)}{P(A)P(B)} \right) + P(A, \bar{B}) \log \left( \frac{P(A, \bar{B})}{P(A)P(\bar{B})} \right) + P(\bar{A}, B) \log \left( \frac{P(\bar{A}, B)}{P(\bar{A})P(B)} \right) + P(\bar{A}, \bar{B}) \log \left( \frac{P(\bar{A}, \bar{B})}{P(\bar{A})P(\bar{B})} \right)}{P(A, B) + P(A, \bar{B}) + P(\bar{A}, B) + P(\bar{A}, \bar{B})}$
10	Support (s)	$P(A, B)$
11	Confidence (c)	$\frac{P(A, B)}{P(A)}$
12	Lift (L)	$\frac{P(A, B)}{P(A)P(B)}$
13	Conviction (V)	$\frac{P(A)P(B)}{P(A, \bar{B})}$
14	Interest (I)	$\frac{P(A, B)}{P(A)P(B)}$
15	Gini (G)	$\frac{P(A, B)}{P(A)P(B)}$
16	Plattensky-Shapirko's (PS)	$\frac{P(A, B)}{P(A)P(B)}$
17	Certainty factor (F)	$\frac{P(A, B) - P(A)P(B)}{P(A)}$
18	Added Value (AV)	$\frac{P(A, B) - P(A)P(B)}{P(A)}$
19	Collective strength (S)	$\frac{P(A, B) - P(A)P(B)}{P(A)}$
20	Accuracy (C)	$\frac{P(A, B) - P(A)P(B)}{P(A)}$
21	Khloegen (K)	$\sqrt{P(A, B) \max(P(B A) - P(B), P(A B) - P(A))}$

14

---

---

---

---

---

---

---

---

---

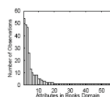
---

# Correlation Measure

- Contingency table.

	$A_1$	$\neg A_1$	
$A_2$	$f_{11}$	$f_{10}$	$f_{1+}$
$\neg A_2$	$f_{01}$	$f_{00}$	$f_{0+}$
	$f_{+1}$	$f_{+0}$	$f_{++}$

- Design of a correlation measure is empirical
- No good correlation measure
- Not uniform attribute frequency



15

---

---

---

---

---

---

---

---

---

---



## H-Measure

$$m_n(A_p, A_q) = H(A_p, A_q) = \frac{f_{01}f_{10}}{f_{+1}f_{1+}}$$

- H-measure is as the negative correlation  $m_n$
- H value close to 0  $\rightarrow$  highly degree of positive correlation
- H value to 1  $\rightarrow$  high degree of negative correlation
- H-measure avoids the sparseness problem by ignoring  $f_{00}$

19

---

---

---

---

---

---

---

---

## Frequent attribute problem

- The ability to differentiate highly frequent attributes from really correlated ones.
- Generating false positives in group discovery

(c1) Example of frequent attribute problem with measure Jaccard.  $A_p$  and  $A_q$  are independent but a higher Jaccard (0.82)

	$A_p$	$\neg A_p$	
$A_q$	81	9	90
$\neg A_q$	9	1	10
	90	10	100

(c2) Example of frequent attribute problem with measure Jaccard.  $A_p$  and  $A_q$  are positively correlated but a lower Jaccard (0.8)

	$A_p$	$\neg A_p$	
$A_q$	8	1	9
$\neg A_q$	1	90	91
	9	91	100

20

---

---

---

---

---

---

---

---

## Outline

- Introduction
- Motivation of integrating the deep web
- Mining Algorithm
- H-Measure
- Data Preparation Step
- Experiments
- Discussion

21

---

---

---

---

---

---

---

---

## Data Preparation

- **Form extraction:** extracting attributes entities (names and domains) from query interfaces
  - Standard Normalization: Stemming attribute names and domains
  - Normalizing irregular nouns and verbs e.g. children → child
  - Delete common stop words
- **Type recognition**
  - Problem of homonyms (the same name with different meanings) → distinguish both names and types
  - Type identification is not declared in web interfaces → type recognizer to recognize types from domain values of attributes entities
- **Syntactic Merging**
  - Name-bases merging: merging two attribute entities if they are similar in names
  - Domain-based merging: merging two attribute entities if they are similar in domain values.

22

---

---

---

---

---

---

---

---

## Outline

- Introduction
- Motivation of integrating the deep web
- Mining Algorithm
- H-Measure
- Data Preparation Step
- Experiments
- Discussion

23

---

---

---

---

---

---

---

---

## Experiments

- Two datasets:
  - TEL-8 dataset
    - Raw web pages over 447 deep web sources in 8 domains.
    - Each domain has about 20-70 sources
  - BAMM dataset
    - Manually extracted attribute names over 211 sources in 4 domains ( around 50 sources per domain)

24

---

---

---

---

---

---

---

---



## Experimental Result (result on TEL-8 Dataset)

Step	Rule of	Result	$C_{min}$	$C_{max}$
group discovery	G	$C_1 = \{last\ name\ (unknown),\ first\ name\ (any)\}$	0.94	0.94
		$C_2 = \{title\ (any),\ keyword\ (any)\}$	0.94	0.94
		$C_3 = \{last\ name\ (any),\ title\ (any)\}$	0.94	0.94
		$C_4 = \{first\ name\ (any),\ catalog\ (any)\}$	0.90	0.90
matching discovery	M	$C_5 = \{first\ name\ (any),\ keyword\ (any)\}$	0.87	0.87
		$M_1 = \{author\ (any)\} = \{last\ name\ (any),\ first\ name\ (any)\}$	0.87	0.87
		$M_2 = \{author\ (any)\} = \{last\ name\ (any)\}$	0.87	0.87
		$M_3 = \{tempet\ (string)\} = \{category\ (string)\}$	0.83	0.83
		$M_4 = \{author\ (any)\} = \{year\ number\ (any),\ catalog\ (any)\}$	0.82	0.82
matching selection	R	$M_5 = \{author\ (any)\} = \{first\ name\ (any)\}$	0.82	0.82
		$M_6 = \{category\ (string)\} = \{publisher\ (string)\}$	0.76	0.76
		$M_7 = \{author\ (any)\} = \{last\ name\ (any),\ first\ name\ (any)\}$	0.87	0.87
Domain	Final Output After Matching Selection	$R_1 = \{to\ (string)\} = \{arrival\ city\ (string)\}$		Correct?
		$R_2 = \{destination\ (string)\} = \{to\ (string)\} = \{arrival\ city\ (string)\}$		Y
		$R_3 = \{departure\ date\ (datetime)\} = \{depart\ (datetime)\}$		Y
		$R_4 = \{passenger\ (integer)\} = \{adult\ (integer),\ child\ (integer),\ infant\ (integer)\}$		N
		$R_5 = \{from\ (string),\ to\ (string)\} = \{departure\ city\ (string),\ arrival\ city\ (string)\}$		Y
		$R_6 = \{from\ (string)\} = \{depart\ (string)\}$		Y
		$R_7 = \{return\ date\ (datetime)\} = \{return\ (datetime)\}$		Y
Movies	$R_8 = \{artist\ (any)\} = \{actor\ (any)\} = \{star\ (any)\}$		Y	
	$R_9 = \{genre\ (string)\} = \{category\ (string)\}$		Y	
	$R_{10} = \{cast\ and\ crew\ (any)\} = \{actor\ (any),\ director\ (any)\}$		Y	

25

## Experimental Result (result on TEL-8 Dataset)

- Attributes above a frequency threshold T

Domain	$P_T$ (20%)	$R_T$ (20%)	$P_T$ (10%)	$R_T$ (10%)
Books	1	1	1	1
Airfares	1	1	1	0.71
Movies	1	1	1	1
MusicRecords	1	1	0.76	1
Hotels	0.86	1	0.86	0.87
CarRentals	0.72	1	0.72	0.80
Jobs	1	0.86	0.76	0.87
Automobiles	1	1	0.93	1

- Evaluating the H-measure

Domain	$P_T(H)$ (10%)	$R_T(H)$ (10%)	$P_T(C)$ (10%)	$R_T(C)$ (10%)
Books	1	1	0.80	1
Airfares	1	0.71	0.79	0.61
Movies	1	1	0.93	1
MusicRecords	0.76	1	0.76	1
Hotels	0.86	0.87	0.44	0.95
CarRentals	0.72	0.80	0.68	0.62
Jobs	0.76	0.87	0.64	0.87
Automobiles	0.93	1	0.78	1

26

## Conclusion

- Complex matching n-ary complex matching
- DCM framework
- Mining Algorithm
- H-measure

27

## Discussion Points

- Their approach works for the same domain, How it can cross the domain boundary is an open problem.
- Choosing the threshold values is empirically, they did it by testing the algorithm with various threshold, how it can be designed more systematically?
- H-Measure were derived based on the observation of the data, it should be systematic. Does it work the same if the datasets change?
- Compare their H-measure with other measures not only one.
- In removing the conflicted groups, how do they remove the conflicted one if both groups have the same score value?

28

---

---

---

---

---

---

---

---

*Questions ?*

29

---

---

---

---

---

---

---

---