

# Making Sense of (Multi-)Relational Data

Part IV: Exploration by Descriptive Modelling – Fully Relational  
Local Approaches

Jefrey Lijffijt

Eirini Spyropoulou

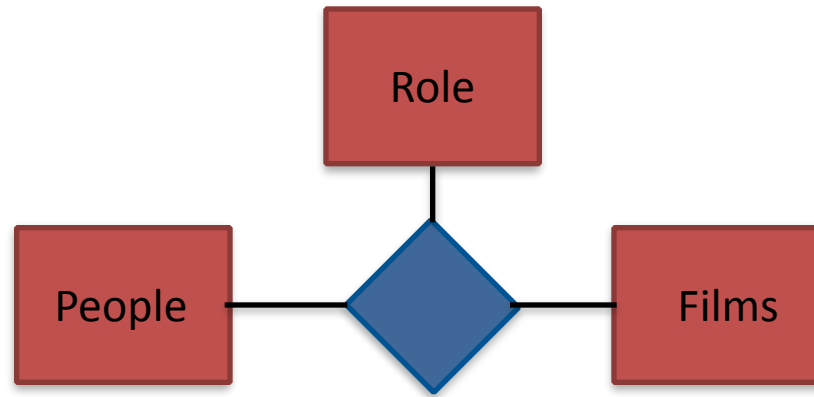
Tijl De Bie

## Fully relational local algorithms

- N-set mining
- RMiner & variants
- Constraint programming for closed relational sets
- Uncovering the plot

# N-set mining

# N-set mining (Cerf et al., Trans. Knowl. Discov. Data, 2009)



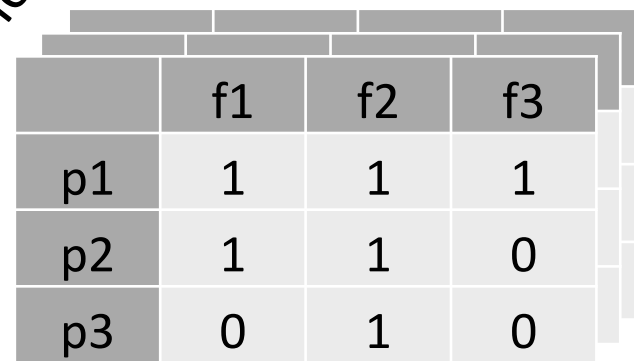
# Pattern Syntax, N-sets

roles

films

$\mathcal{R}$

	f1	f2	f3
p1	1	1	1
p2	1	1	0
p3	0	1	0



## N-sets

$\mathcal{R}$

		films											
		f1	f2	f3									
people	actor	p1	1	1	1	p1	1	1	1	p1	0	1	1
		p2	1	1	0	p2	1	1	0	p2	0	0	0
		p3	1	0	0	p3	0	1	0	p3	0	0	1
		actor			director				producer				

Sets of entities that are complete with respect to the relationship and maximal

---

## N-sets

$\mathcal{R}$

films

	f1	f2	f3
people p1	1	1	1
p2	1	1	0
p3	1	0	0

**actor**

	f1	f2	f3
p1	1	1	1
p2	1	1	0
p3	0	1	0

**director**

	f1	f2	f3
p1	0	1	1
p2	0	0	0
p3	0	0	1

**producer**

N-set: {p1, f2, f3, actor, director, producer}

---

## N-sets

$\mathcal{R}$

films

people

	f1	f2	f3
p1	1	1	1
p2	1	1	0
p3	1	0	0

actor

	f1	f2	f3
p1	1	1	1
p2	1	1	0
p3	0	1	0

director

	f1	f2	f3
p1	0	1	1
p2	0	0	0
p3	0	0	1

producer

N-set:  $\{p1, f2, f3, \text{actor}, \text{director}, \text{producer}\}$

$\{p1, f2, \text{actor}\} \in \mathcal{R}$      $\{p1, f2, \text{director}\} \in \mathcal{R}$   
 $\{p1, f3, \text{actor}\} \in \mathcal{R}$      $\{p1, f3, \text{director}\} \in \mathcal{R}$



## Interestingness of N-sets

- No interestingness measure defined.
- Constraints on the number of entities per entity type help to focus on a smaller pattern set.

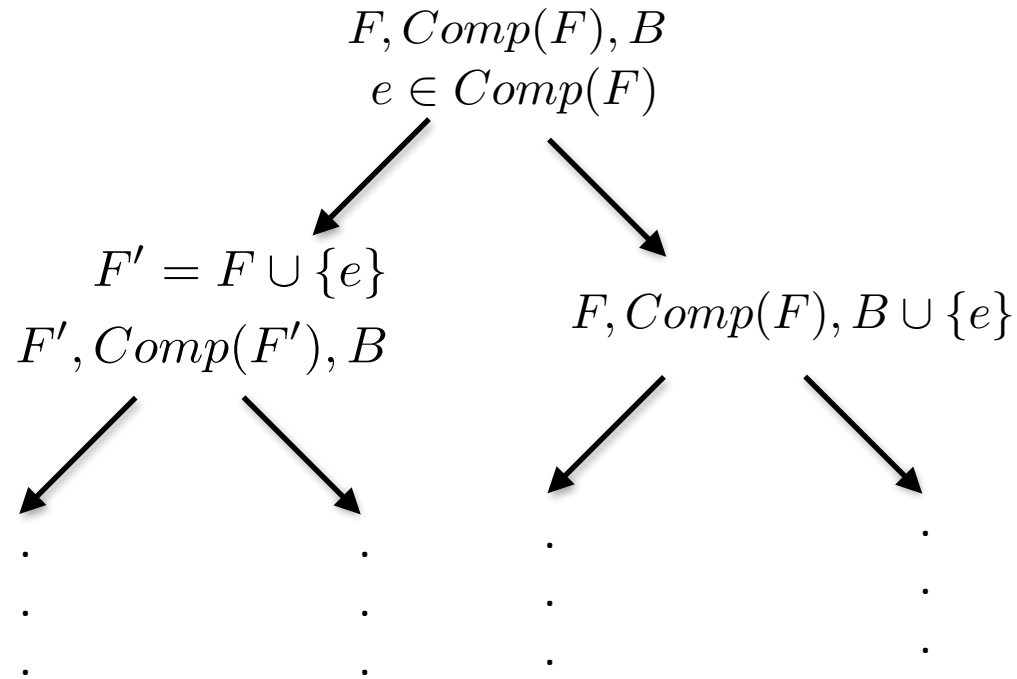
## N-set Mining Algorithm - DataPeeler

- Exhaustive Search
- Divide and conquer enumeration strategy
- Use completeness to reduce the search space

# DataPeeler

$F$  : current solution

$Comp(F) : \forall e \in E, F \cup \{e\}$  is complete

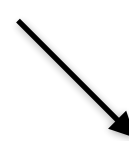


# DataPeeler

$F$  : current solution

$Comp(F) : \forall e \in E, F \cup \{e\}$  is complete

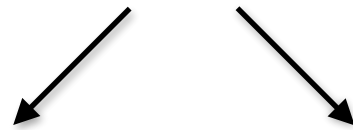
$F, Comp(F), B$   
 $e \in Comp(F)$



Add all elements of  $Comp(F')$   
that do not violate  
completeness in  
 $Comp(F')$

$F' = F \cup \{e\}$   
 $F', Comp(F'), B$

$F, Comp(F), B \cup \{e\}$

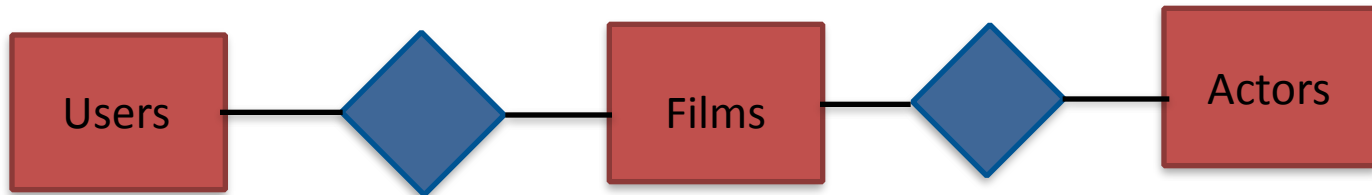


⋮  
⋮  
⋮

⋮  
⋮  
⋮

# RMiner & variants

## RMiner - E.S. et al., DMKD, 2014



## RMiner - Pattern Syntax (MCCSs)

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

users

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

films

# MCCSs (Maximal Complete Connected Subsets)

films

	f1	f2	f3
users			
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

actors

	a1	a2	a3
films			
f1	1	0	1
f2	1	0	1
f3	0	1	0



# MCCSs (Maximal Complete Connected Subsets)

Completeness

$$\mathcal{R} = R_{users,films} \cup R_{films,actors}$$

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

users

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

films

# MCCSs (Maximal Complete Connected Subsets)

Complete Subset: {u1, f1, f2, a1}

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

users

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

films

# MCCSs (Maximal Complete Connected Subsets)

Complete Subset: {u1, f1, f2, a1}

$\{u1, f1\} \in \mathcal{R}$   
 $\{u1, f2\} \in \mathcal{R}$   
 $\{f1, a1\} \in \mathcal{R}$   
 $\{f2, a1\} \in \mathcal{R}$

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

users

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

films

# RMiner - Pattern Syntax (MCCSs)

Connectedness

films

users

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

actors

films

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

## RMiner - Pattern Syntax (MCCSs)

Complete Connected Subset (CCS): {u1, f1, f2, a1}

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

## RMiner - Pattern Syntax (MCCSs)

Complete **not connected** Subset: {u1, a2}

films

	f1	f2	f3
users			
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

actors

	a1	a2	a3
films			
f1	1	0	1
f2	1	0	1
f3	0	1	0

## RMiner - Pattern Syntax (MCCSs)

Complete Connected Subset (CCS): {u1, f1, f2, a1}

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

# RMiner - Pattern Syntax (MCCSs)

Maximality

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0



## RMiner - Pattern Syntax (MCCSs)

Maximal Complete Connected Subset (MCCS): {u1, u2, f1, f2, a1, a3}

films

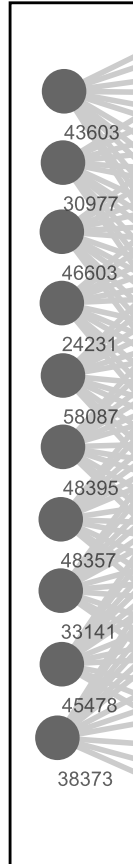
	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

# films

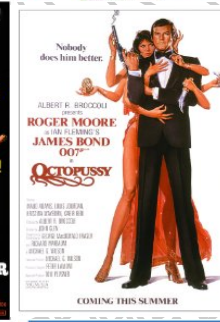
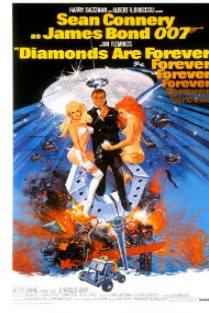
## users



## genres



## actors



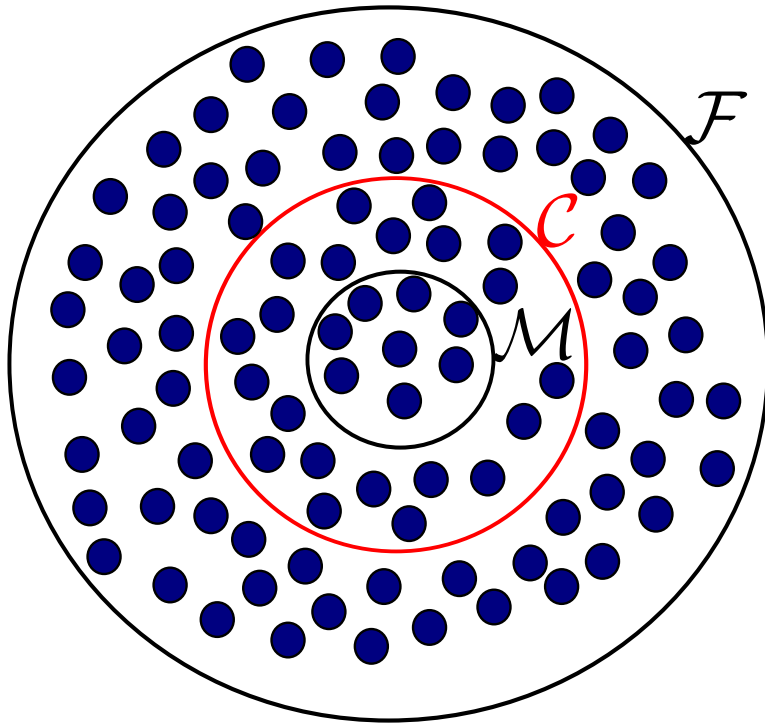
## RMiner - Interestingness

$$\text{Interestingness}(\gamma) = \frac{\text{Self Information}(\gamma)}{\text{Description Length}(\gamma)}.$$

## Interestingness - Background model

- Maximum Entropy distribution on the data
- Constraints on the expected number of “1s” in every row and every column to be equal to the actual number
- The probability between entities that exist in many relationship instances is going to be high

# RMiner - Algorithm



$\mathcal{F}$  : the set of CCSs

$\mathcal{M}$  : the set of MCCSs

$\mathcal{C}$  : the set of closed CCSs

$$\mathcal{M} \subseteq \mathcal{C} \subseteq \mathcal{F}$$

## Closed CCSs

- Fixpoints of a closure operator

Mapping  $\rho: \mathcal{F} \rightarrow \mathcal{F}$

---

Extensivity:  $F \subseteq \rho(F)$  for all  $F \in \mathcal{F}$

Monotonicity:  $\rho(F) \subseteq \rho(F')$  for all  $F, F' \in \mathcal{F}$  with  $F \subseteq F'$

Idempotence:  $\rho(\rho(F)) = \rho(F)$  for all  $F \in \mathcal{F}$

## RMiner - Algorithm

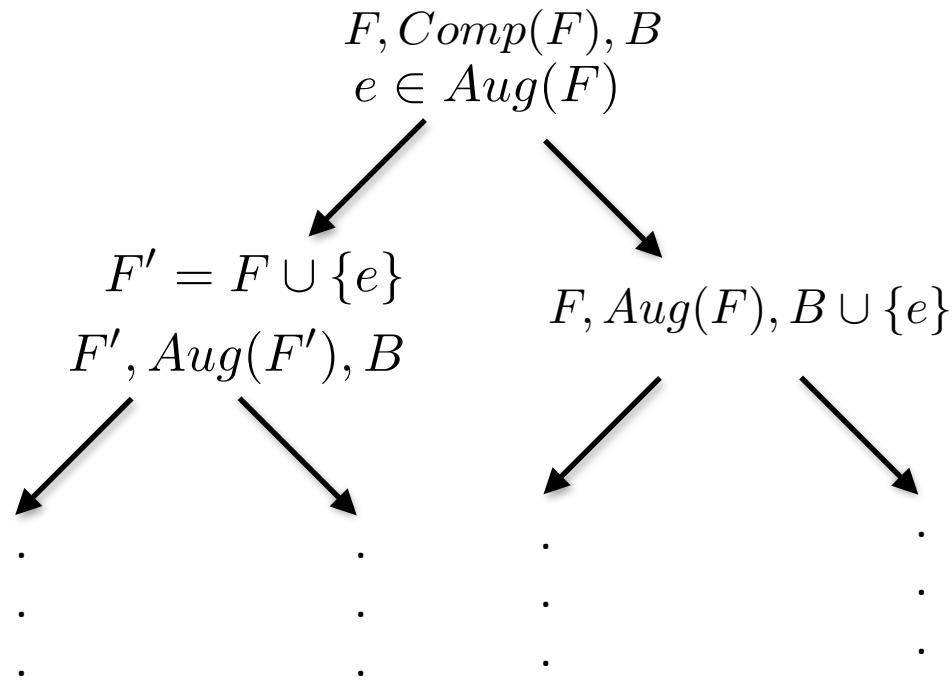
- Exhaustive search
- Based on divide and conquer algorithmic framework of Boley et al.

# RMiner - Algorithm

$F$  : current solution

$Comp(F) : \forall e \in E, F \cup \{e\}$  is complete

$Aug(F) : \forall e \in E, F \cup \{e\}$  is complete and connected



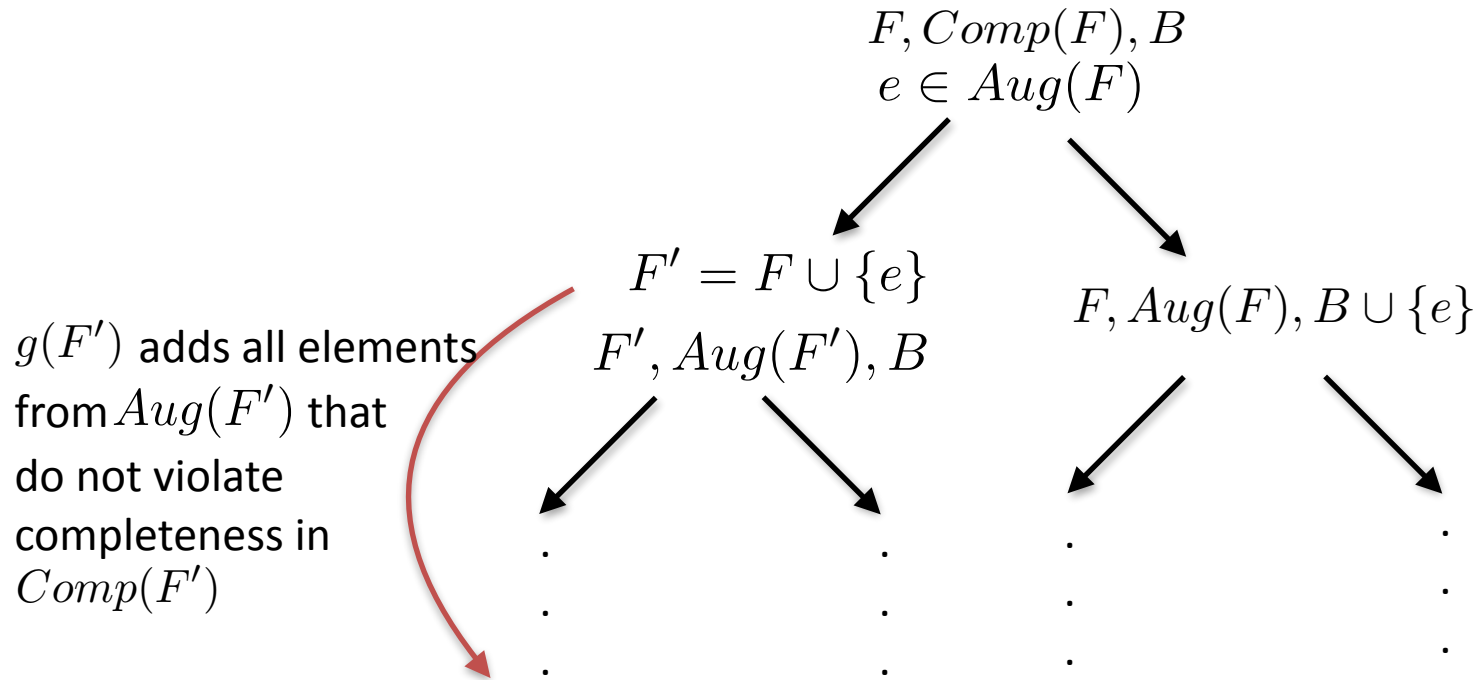


# RMiner - Algorithm

$F$  : current solution

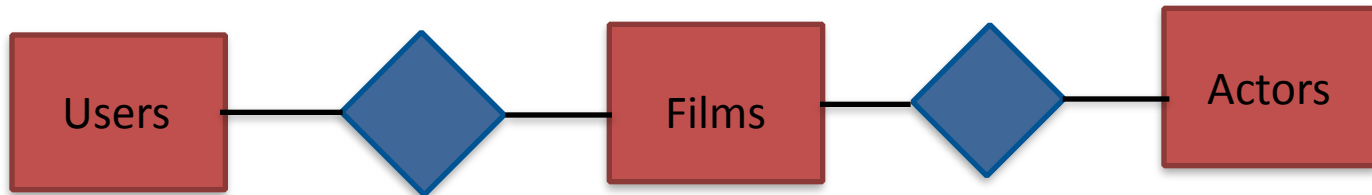
$Comp(F) : \forall e \in E, F \cup \{e\}$  is complete

$Aug(F) : \forall e \in E, F \cup \{e\}$  is complete and connected



# A-RMiner

## A-RMiner (E.S. et al., DSAA, 2014)



## A-RMiner Pattern Syntax (a-CCSs)

- Unions of MCCSs that are maximal extensions of a closed CCS.

## A-RMiner Pattern Syntax (a-CCSs)

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

users

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

films

## A-RMiner Pattern Syntax (a-CCSs)

$M1 = \{u1, u2, f1, f2, a1, a3\}$

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

## A-RMiner Pattern Syntax (a-CCSs)

$M2 = \{u3, u4, f2, f3\}$

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

## A-RMiner Pattern Syntax (a-CCSs)

a-CCS:  $M1 \cup M2 = \{u1, u2, u3, u4, f1, f2, f3, a1, a3\}$

films

	<b>f1</b>	<b>f2</b>	<b>f3</b>
<b>u1</b>	1	1	0
<b>u2</b>	1	1	0
<b>u3</b>	0	1	1
<b>u4</b>	0	1	1

actors

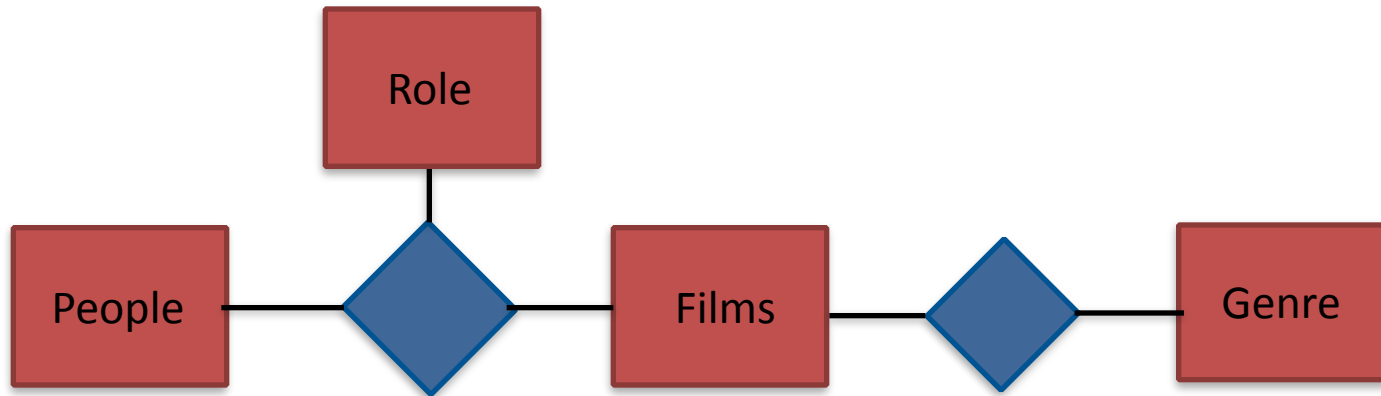
	<b>a1</b>	<b>a2</b>	<b>a3</b>
<b>f1</b>	1	0	1
<b>f2</b>	1	0	1
<b>f3</b>	0	1	0



## A-RMiner - Algorithm

- Same algorithmic framework as RMiner
- Finds approximate patterns at no extra computational cost

## N-RMiner (E.S. et al., Discovery Science 2013)



## N-RMiner Pattern Syntax (NMCCSs)

roles

films

	f1	f2	f3
p1	1	1	1
p2	1	1	0
p3	0	1	0

## NMCCSs

films

	f1	f2	f3
actor	1	1	1
actor	1	1	0
actor	1	0	0

	f1	f2	f3
director	1	1	1
director	1	1	0
director	0	1	0

	f1	f2	f3
producer	0	1	1
producer	0	0	0
producer	0	0	1

## NMCCSs

films

	f1	f2	f3	
actor	p1	1	1	1
p2	1	1	0	
p3	1	0	0	

	f1	f2	f3	
director	p1	1	1	1
p2	1	1	0	
p3	0	1	0	

	f1	f2	f3	
producer	p1	0	1	1
p2	0	0	0	
p3	0	0	1	

**Completeness:** Similar to N-sets with but allowing a subset of the entity types to be in the pattern.

---

## NMCCSs

films

	f1	f2	f3
<b>p1</b>	1	1	1
<b>p2</b>	1	1	0
<b>p3</b>	1	0	0

**actor**

	f1	f2	f3
<b>p1</b>	1	1	1
<b>p2</b>	1	1	0
<b>p3</b>	0	1	0

**director**

	f1	f2	f3
<b>p1</b>	0	1	1
<b>p2</b>	0	0	0
<b>p3</b>	0	0	1

**producer**

N-MCCS: {f1, f2, p1, p2, actor, director}

## NMCCSs

films

people

	f1	f2	f3
p1	1	1	1
p2	1	1	0
p3	1	0	0

**actor**

	f1	f2	f3
p1	1	1	1
p2	1	1	0
p3	0	1	0

**director**

	f1	f2	f3
p1	0	1	1
p2	0	0	0
p3	0	0	1

**producer**

N-MCCS: {p1, f2, f3, actor, director, producer}

---

## N-RMiner - Algorithm

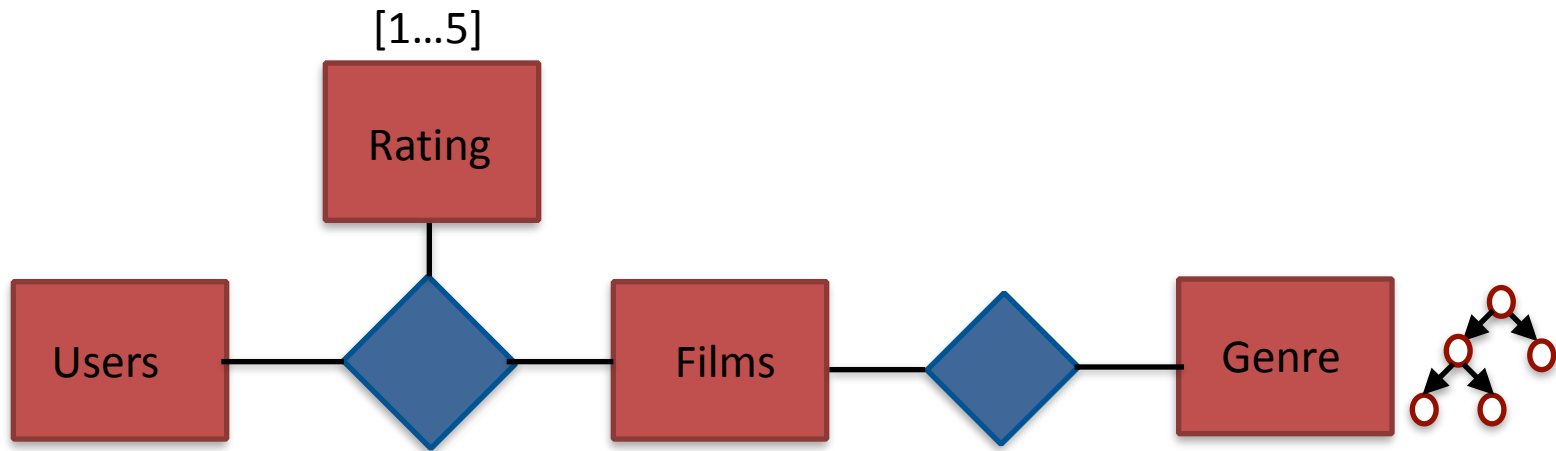
- Same enumeration algorithm
- Different in the way completeness is checked



## N-RMiner - Interestingness

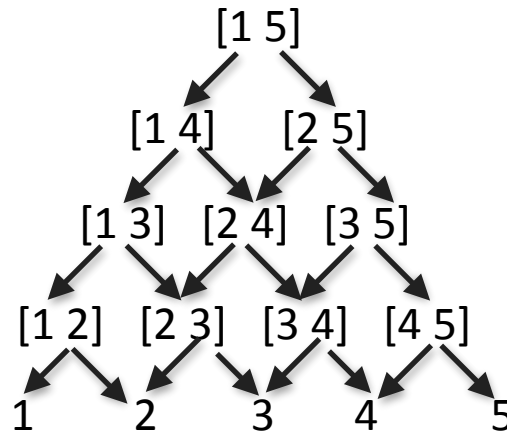
- Similar to RMiner

## P-N-RMiner (Lijffijt et al., DSAA, 2015)



## P-N-RMiner, Pattern Syntax

- Represents numerical attributes as partial orders



- It can handle both numerical and hierarchical attributes in the same way

## P-N-RMiner, Pattern Syntax (MCCPSs)

users

films

ratings

	f1	f2	f3
r1	1	0	0
r2	0	1	0
r3	0	0	0
r4	0	0	0
r5	0	0	1

genres

films

	comedy	history	drama
f1	0	1	1
f2	0	1	1
f3	1	0	0

## P-N-RMiner, Pattern Syntax (MCCPSs)

		films		
		f1	f2	f3
ratings	users			
	r1	1	0	0
	r2	0	1	0
	r3	0	0	0
	r4	0	0	0
r5	0	0	1	

		genres		
		comedy	history	drama
films	f1	0	1	1
	f2	0	1	1
	f3	1	0	0

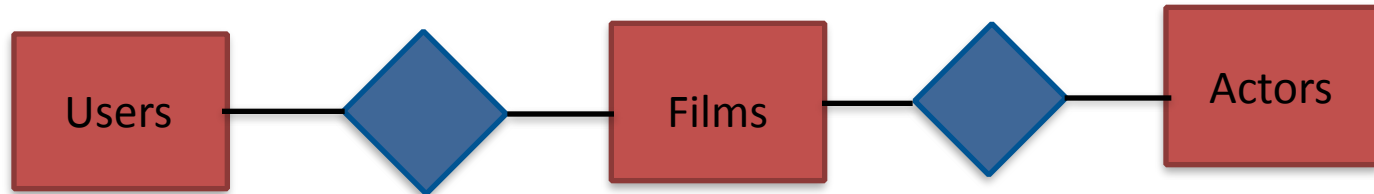
MCCPS: {user1, user2, [1 2], f1, f2, history, drama}

## P-N-RMiner, Algorithm

- Same enumeration framework as RMiner
- Uses the partial order to reduce the search space

# Constraint programming for closed relational sets

# Constraint programming for closed relational sets (Nijssen et al., ICDM Workshops 2011)





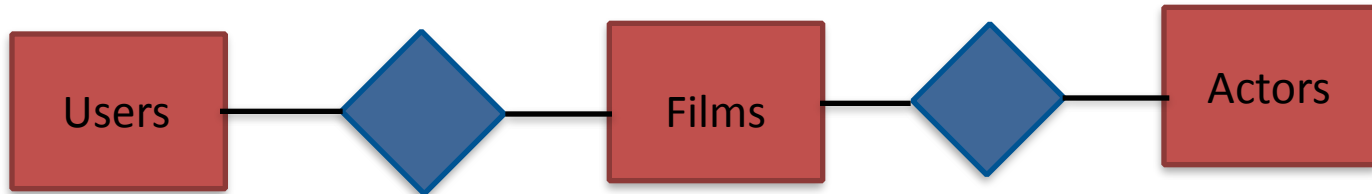
---

## Constraint programming for closed relational sets (Nijssen et al., ICDM Workshops 2011)

- Constraint programming approach
- Defines both pattern syntax and pattern quality constraints
- Pattern syntax is the same as that of MCCSs
- Pattern quality constraints are on the minimum number of entities per entity type

# Uncovering the plot

## Uncovering the plot (Wu et al. DMKD, 2014)



## Uncovering the plot (Wu et al. DMKD, 2014)

- Approximate multi-relational patterns similar to a-CCSs.
- Defined as chains of bi-clusters such that they overlap on one entity type.
- Works only for binary relationships.

## Uncovering the plot - Bi-cluster Chains

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

users

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

films

## Uncovering the plot - Bi-cluster Chains

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

users

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

films

# Uncovering the plot - Bi-cluster Chains

films

	f1	f2	f3
u1	1	1	0
u2	1	1	0
u3	0	1	1
u4	0	1	1

users

actors

	a1	a2	a3
f1	1	0	1
f2	1	0	1
f3	0	1	0

films

{u1, u2, f1, f2, a1}

## Uncovering the plot - Algorithm

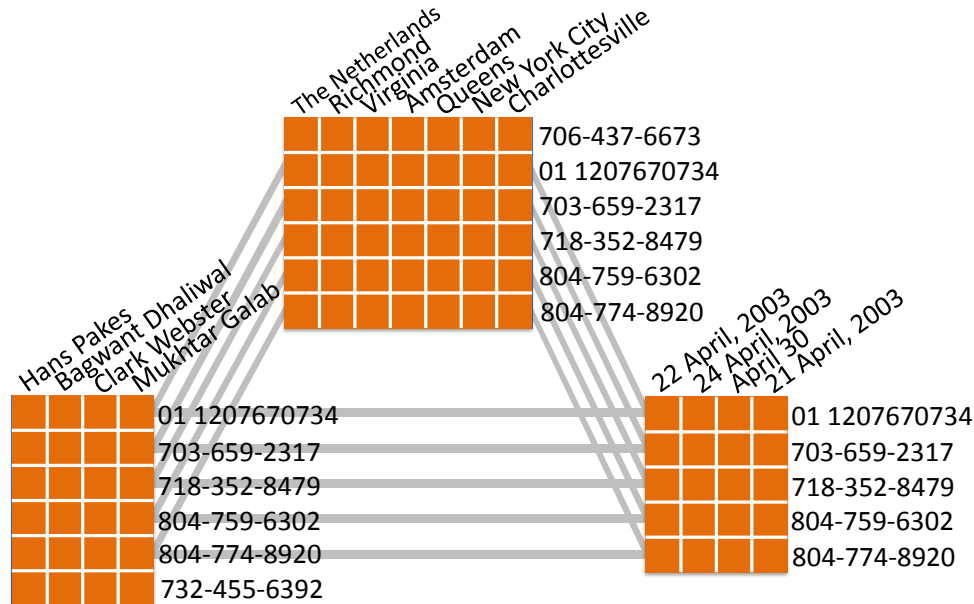
- Based on mining all bi-clusters on every relationship first.
- Then greedily combining them to form a chain.



## Uncovering the plot - Pattern interestingness

- Information content with respect to MaxEnt background model
- Different prior knowledge to RMiner and variants
  - Defined over the join of all relationships
  - Has the form of constraints over the area of a tile.

# Uncovering the plot - Pattern example



---

# References

1. L. Cerf, J. Besson, C. Robardet, and J.-F. Boulicaut. Closed patterns meet n-ary relations. *ACM Trans. Knowl. Discov. Data*, 3(1):3:1–3:36, 2009.
2. E. Spyropoulou, T. De Bie and M. Boley. Interesting Pattern Mining in Multi-Relational Data. *Data Mining & Knowledge Discovery*, 28(3), pp.808-849, 2014.
3. E. Spyropoulou and T. De Bie. Approximate Multi-Relational Patterns. *Proceedings of the International Conference on Data Science and Advanced Analytics (DSAA)*, 2014.
4. E. Spyropoulou, T. De Bie and M. Boley. Mining Interesting Patterns in Multi-Relational Data with N-ary Relationships. *Proceedings of the International Conference on Discovery Science*, pp. 217-232, 2013.
5. S. Nijssen, A. Jimenez and T Guns. Constraint-based pattern mining in multi-relational databases. *Proceedings of Data Mining Workshops (ICDMW)*, 2011.
6. Uncovering the Plot: Detecting Surprising Coalitions of Entities in Multi-Relational Schemas. *Data Mining and Knowledge Discovery vol.28(5)*, pp 1398-1428, Springer, 2014.
7. J. Lijffijt, E. Spyropoulou and T. De Bie. P-N-RMiner: A Generic Framework for Mining Interesting Structured Relational Patterns. *Proceedings of the International Conference on Data Science and Advanced Analytics (DSAA)*,