

Redescription Mining

An Introduction

Esther Galbrun

2009-2013



2014



Redescriptions

An Example with World countries



Canada



Chile



China



France



United Kingdom



Mexico



Mozambique



Russia



United States

Redescriptions

An Example with World countries

— Countries outside the Americas with land area above 8 billion square kilometers



Canada



France



Mozambique



Chile



United Kingdom



Russia



China



Mexico



United States

Redescriptions

An Example with World countries

— Countries outside the Americas with land area above 8 billion square kilometers



Canada



Chile



China



France



United Kingdom



Mexico



Mozambique



Russia



United States

Redescriptions

An Example with World countries

— Permanent members of the UN Security Council with a history of state communism



Canada



Chile



China



France



United Kingdom



Mexico



Mozambique



Russia



United States

Redescriptions

An Example with World countries

— Permanent members of the UN Security Council with a history of state communism



Canada



Chile



China



France



United Kingdom



Mexico



Mozambique



Russia



United States

Redescriptions

An Example with World countries

- Countries outside the Americas with land area above 8 billion square kilometers
- Permanent members of the UN Security Council with a history of state communism



Canada



France



Mozambique



Chile



United Kingdom



Russia



China



Mexico



United States

Redescription Mining

- Countries outside the Americas with land area above 8 billion square kilometers
- Permanent members of the UN Security Council with a history of state communism

*Finding different ways
to characterize the same things ...*

Redescription Mining



Canada



Chile



China



France



United Kingdom



Mexico



Mozambique



Russia



United States








*...finding multiple things
that share common characterizations*

Redescriptions

An Example with World countries

Let's get a bit more specific...








Geographic attributes

-  South Hemisphere
-  Border to the Atlantic Ocean
-  Border to the Indian Ocean
-  Border to the Pacific Ocean
-  Continent
-  Land area
-  Highest elevation

Geographic attributes

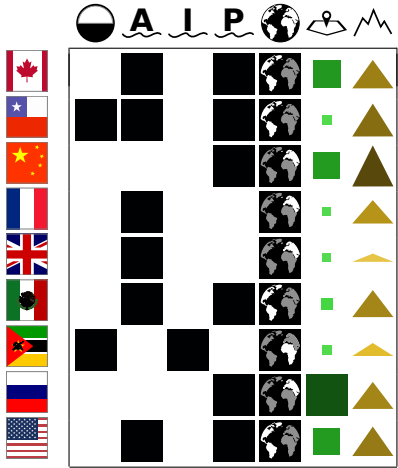


Canada

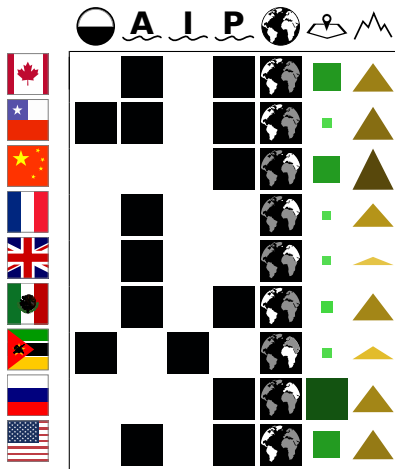
	South Hemisphere	No
	Border to the Atlantic Ocean	Yes
	Border to the Indian Ocean	No
	Border to the Pacific Ocean	Yes
	Continent	Americas
	Land area	9985M km ²
	Highest elevation	5959 m



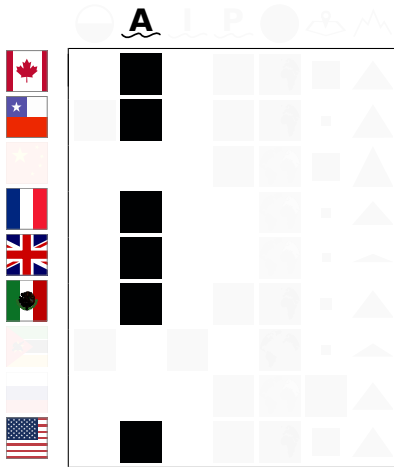
Geographic attributes



Geographic descriptions



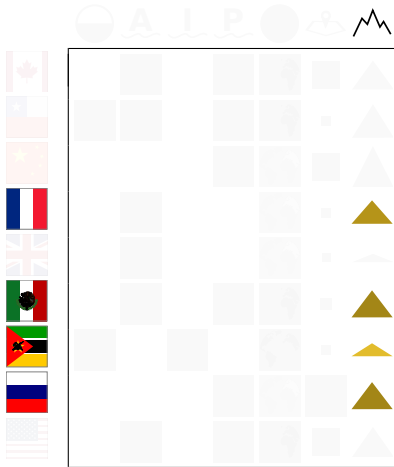
Geographic descriptions



A

*Countries bordering
the Atlantic Ocean*

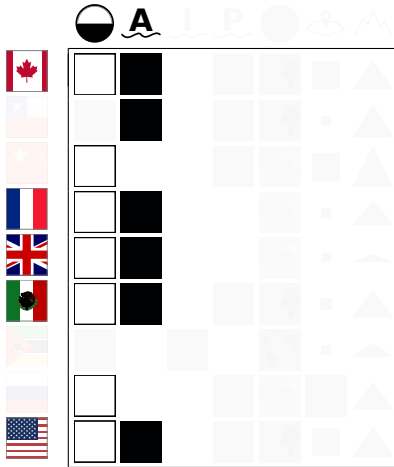
Geographic descriptions



$$\triangle \leq \text{mountain range} \leq \triangle$$

*Countries with highest elevation
between 2400 and 5600 meters*

Geographic descriptions



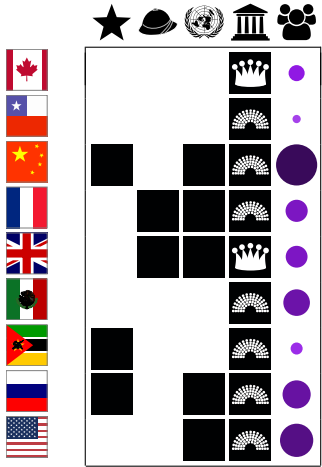
→  AND 

*Countries in the North hemisphere
bordering the Atlantic Ocean*

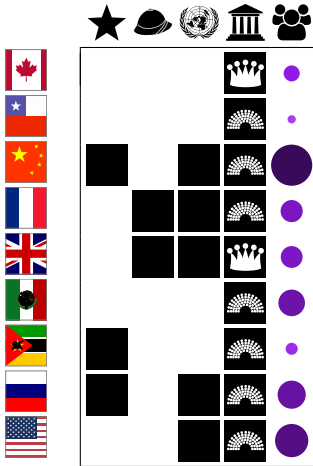
Geopolitical attributes

- ★ History of state communism
- 🌐 History of colonialism
- 🇺🇳 Permanent member of the UNSC
- 🏛️ Type of government
- 👥 Population

Geopolitical attributes

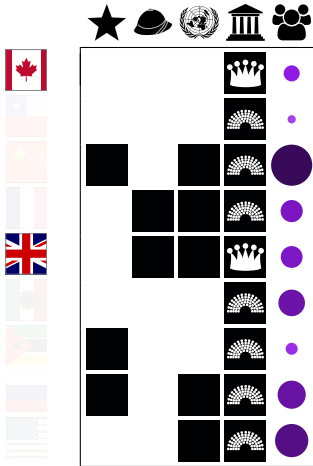


Geopolitical descriptions



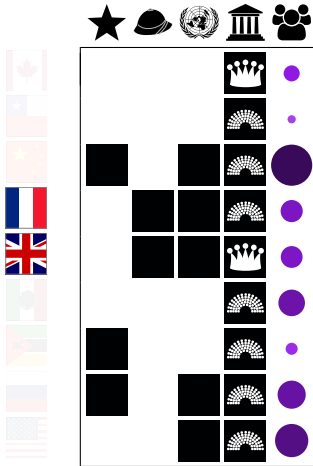


Geopolitical descriptions



 = 
Monarchies

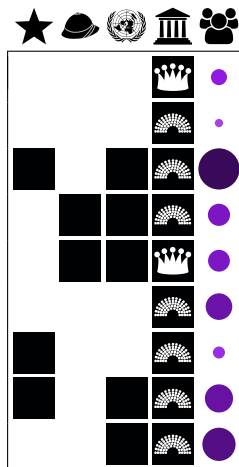
Geopolitical descriptions



 AND 

*Countries with a history of
colonialism members of UNSC*

Two views on the objects



Redescriptions

- Countries outside the Americas with land area above 8 billion square kilometers
- Permanent members of the UN Security Council with a history of state communism



Canada



France



Mozambique



Chile



United Kingdom



Russia



China



Mexico

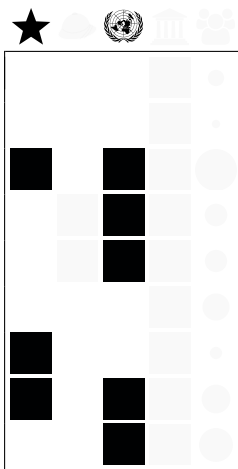
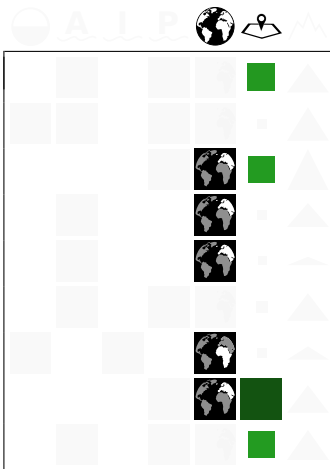


United States

Redescriptions

 \neq  AND  \leq 

 AND 

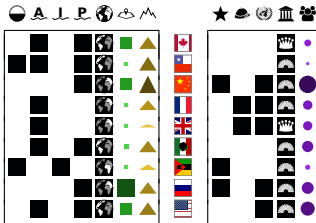


Definitions

Redescription Given two datasets with identity between the rows, a **redescription** is a pair of queries (q_L, q_R) over the columns characterizing approximately the same sets of rows.

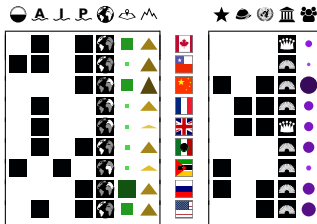
Redescription Mining Given such a pair of datasets and a set of constraints, find the best redescrptions satisfying the constraints.

Definitions



Dataset Two data matrices

Definitions



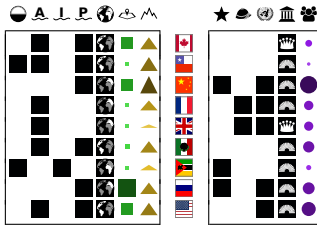
Dataset Two data matrices

Queries Logical formulae

$$\text{Globe} \neq \text{Globe} \text{ AND } \text{Green} \leq \text{Map}$$

$$\text{Star} \text{ AND } \text{UN}$$

Definitions



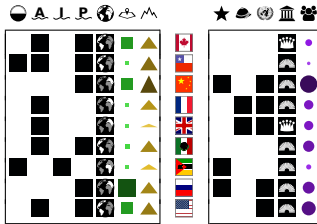
Dataset Two data matrices

Queries Logical formulae

Accuracy Jaccard coefficient

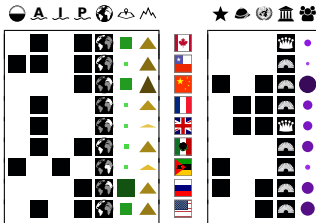
$$J(q_L, q_R) = \frac{|\text{supp}(q_L) \cap \text{supp}(q_R)|}{|\text{supp}(q_L) \cup \text{supp}(q_R)|}$$

Definitions



- Dataset** Two data matrices
- Queries** Logical formulae
- Accuracy** Jaccard coefficient
- Constraints** Support, accuracy, length of the query, p -value, ...

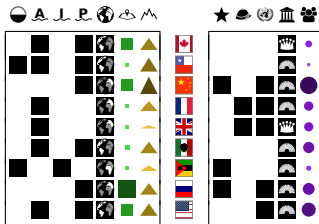
Special Cases



$\boxed{?} \text{ AND } \boxed{?} \Rightarrow ? \text{ AND } ? \text{ AND } ?$
 $\boxed{?} \text{ AND } \boxed{?} \Leftarrow ? \text{ AND } ? \text{ AND } ?$

Only conjunctive queries:
bi-directional
association rules

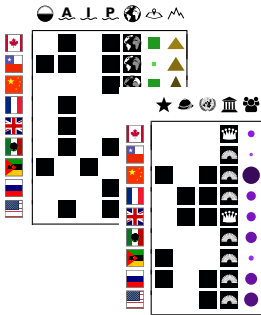
Special Cases



$$\boxed{?} \boxed{?} \boxed{?} \Rightarrow (\text{helmet} \text{ AND } \text{UN}) \text{ OR } \text{star}$$

One query given: classification task

Special Cases



(\neq AND \leq ,
 AND)
{ , }

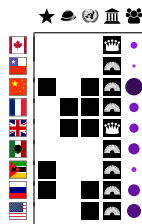
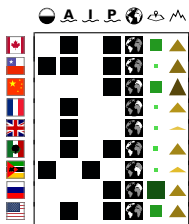
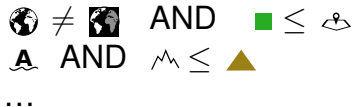
(AND \leq ,
 = OR \leq \leq)
{ , , , , }

...

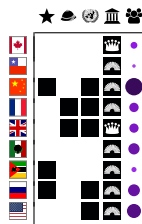
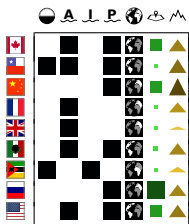
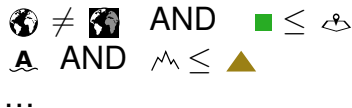
Exploration Strategies

How do we find redescriptions?

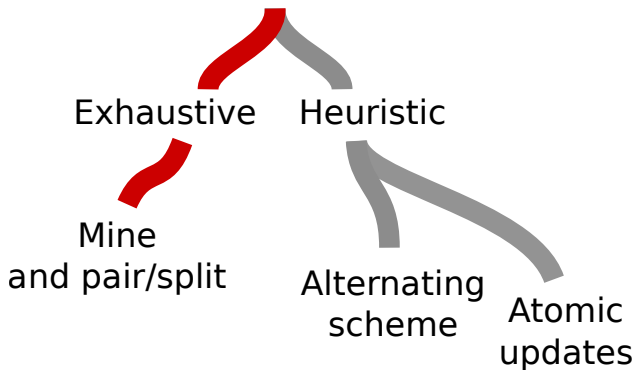
Exploration Strategies



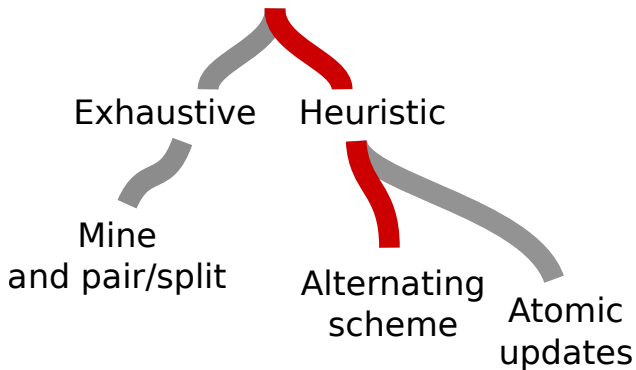
Exploration Strategies



Exploration Strategies

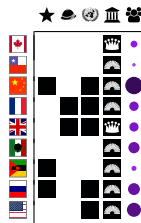
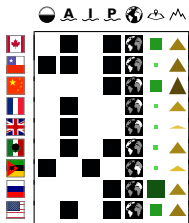


Exploration Strategies



Exploration Strategies

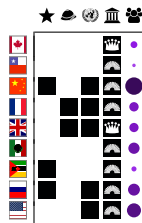
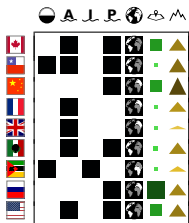
$$(\neg \underbrace{\mathbf{A}} \text{ AND } \neg \underbrace{\mathbf{I}}) \text{ OR } (\underbrace{\mathbf{A}} \text{ AND } \underbrace{\mathbf{P}})$$



Exploration Strategies

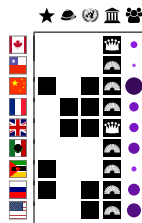
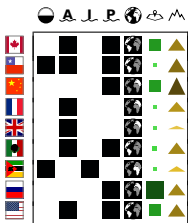
$$(\neg \underline{\mathbf{A}} \text{ AND } \neg \underline{\mathbf{I}}) \text{ OR } (\underline{\mathbf{A}} \text{ AND } \underline{\mathbf{P}})$$

$$\star \text{ OR } \neg \text{🌐}$$



Exploration Strategies

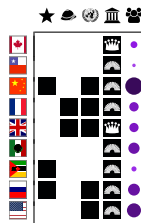
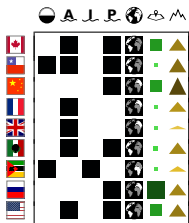
★ OR ↗ 🌐



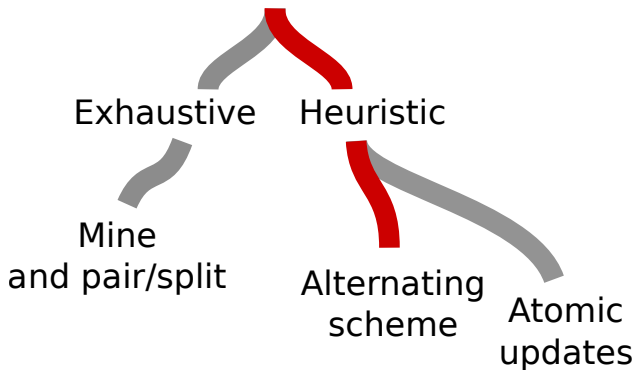
Exploration Strategies

$$\left(\text{🌐} = \text{🌐} \text{ AND } \text{♂} \right) \text{ OR } \text{🌐} \neq \text{🌐}$$

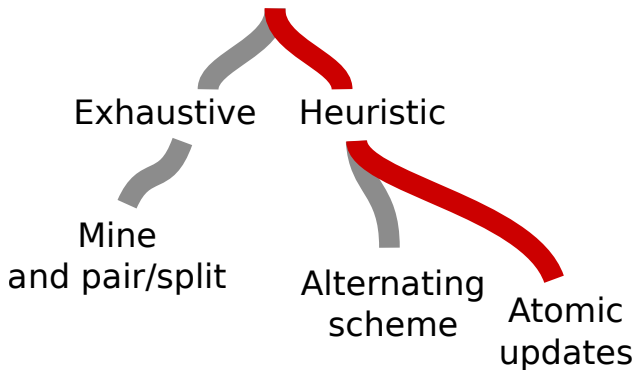
$$\star \text{ OR } \neg \text{🌐}$$



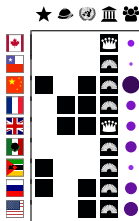
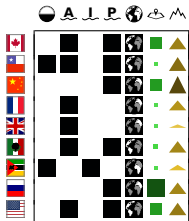
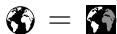
Exploration Strategies



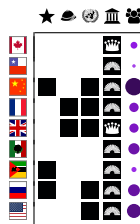
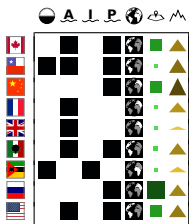
Exploration Strategies



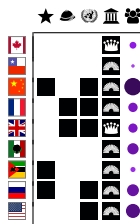
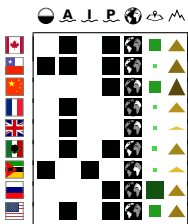
Exploration Strategies



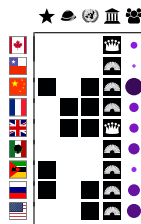
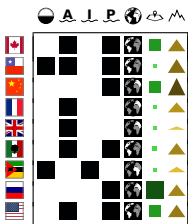
Exploration Strategies



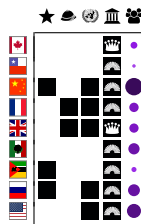
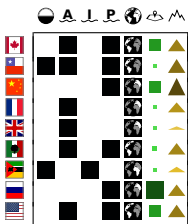
Exploration Strategies



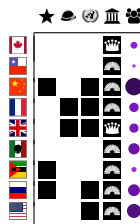
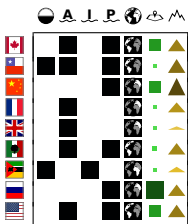
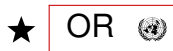
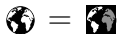
Exploration Strategies



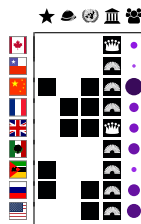
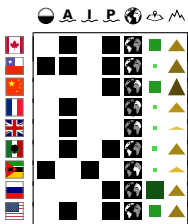
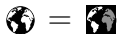
Exploration Strategies



Exploration Strategies



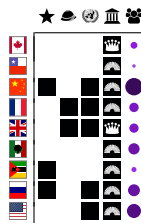
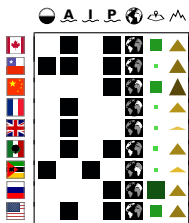
Exploration Strategies



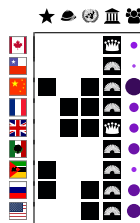
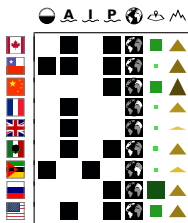
Exploration Strategies

$$\text{🌐} = \text{🌐} \text{ OR } \text{🟩} \leq \text{📍}$$

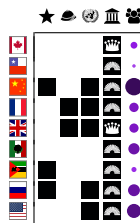
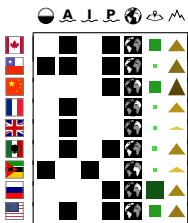
★



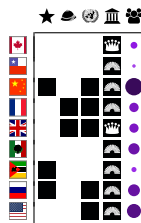
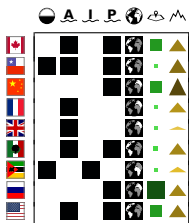
Exploration Strategies



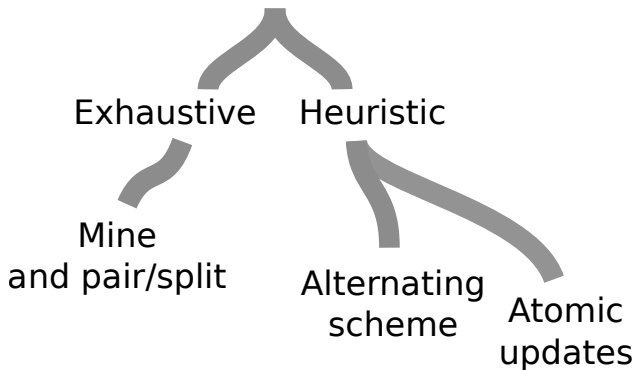
Exploration Strategies



Exploration Strategies



Exploration Strategies

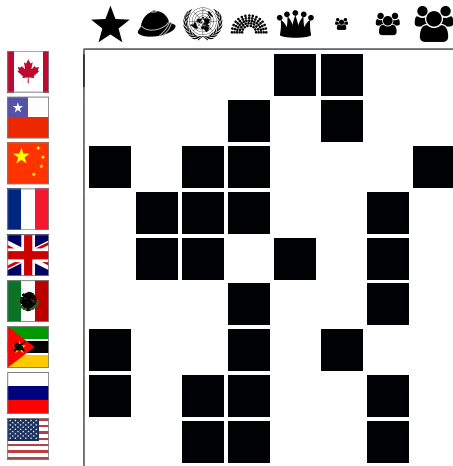


Related work

Redescription mining for Boolean data

Related work

Geopolitical Boolean attributes



Related work

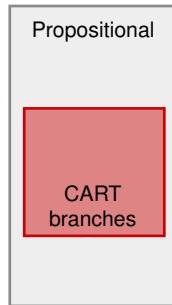
Turning CARTwheels: An Alternating Algorithm for Mining Redescriptions.

N. Ramakrishnan, D. Kumar, B. Mishra, M. Potts,
and R. F. Helm.
In *KDD*, 2004.

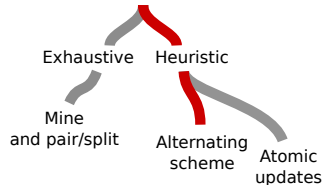
Redescription Mining: Algorithms and Applications in Bioinformatics.

D. Kumar.
PhD Thesis, Virginia Tech, 2007.

Query Languages



Exploration Strategies

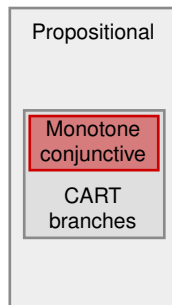


Related work

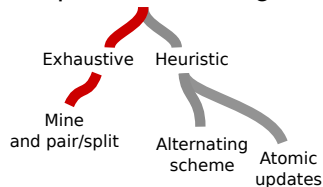
*Redescription Mining:
Structure Theory and Algorithms.*

L. Parida and N. Ramakrishnan.
In *AAAI*, 2005.

Query Languages



Exploration Strategies

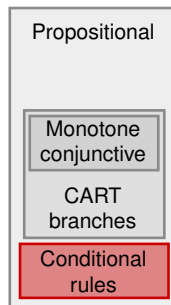


Related work

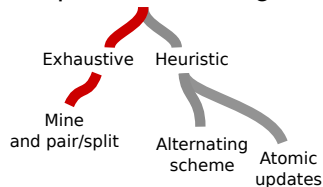
Reasoning About Sets Using Redescription Mining.

M. J. Zaki and N. Ramakrishnan.
In *KDD*, 2005.

Query Languages



Exploration Strategies

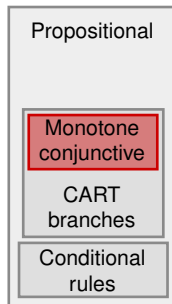


Related work

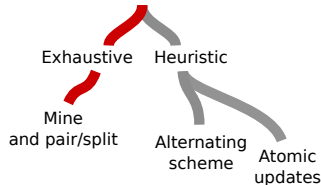
*Finding Subgroups
Having Several Descriptions:
Algorithms for Redescription Mining.*

A. Gallo, P. Miettinen, and H. Mannila.
In *SDM*, 2008.

Query Languages



Exploration Strategies

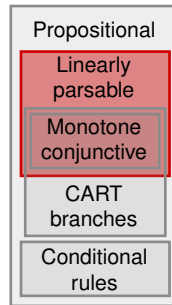


Related work

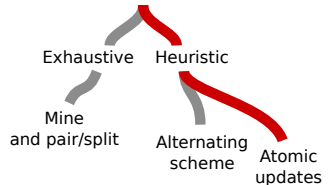
*Finding Subgroups
Having Several Descriptions:
Algorithms for Redescription Mining.*

A. Gallo, P. Miettinen, and H. Mannila.
In *SDM*, 2008.

Query Languages



Exploration Strategies

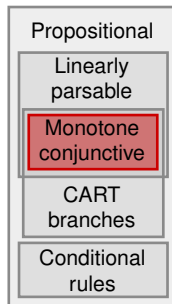


Own work

Selecting a good set of redescrptions

Own work

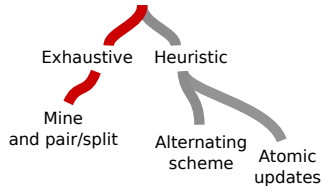
Query Languages



“MDL for Redescription Mining”

with Matthijs van Leeuwen,
Under review.

Exploration Strategies

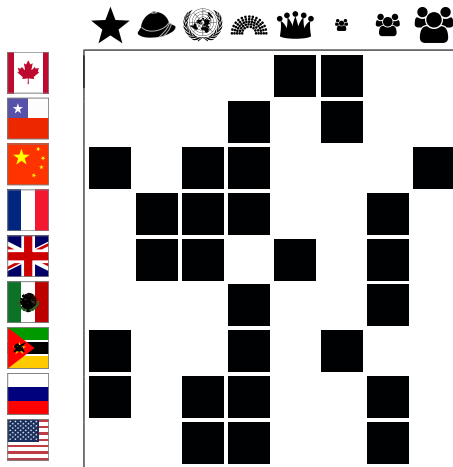


Own work

Extending redescription mining to non-Boolean data

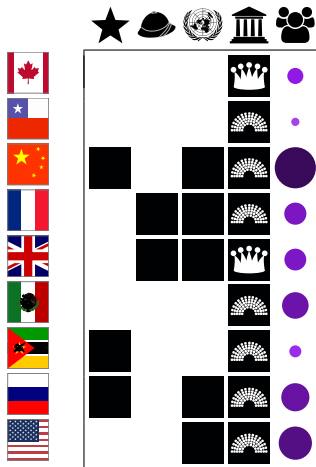
Own work

Geopolitical Boolean attributes



Own work

Geopolitical attributes

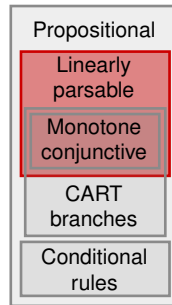


Related work

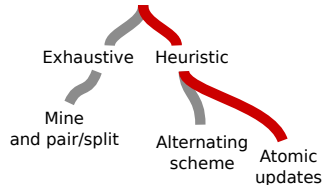
*Finding Subgroups
Having Several Descriptions:
Algorithms for Redescription Mining.*

A. Gallo, P. Miettinen, and H. Mannila.
In *SDM*, 2008.

Query Languages

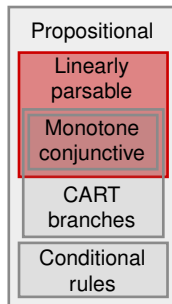


Exploration Strategies



Own work

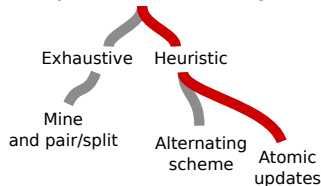
Query Languages



*From Black and White to Full Color:
Extending Redescription Mining
Outside the Boolean World*

with Pauli Miettinen,
In *Statistical Analysis and Data Mining*, 2012.

Exploration Strategies

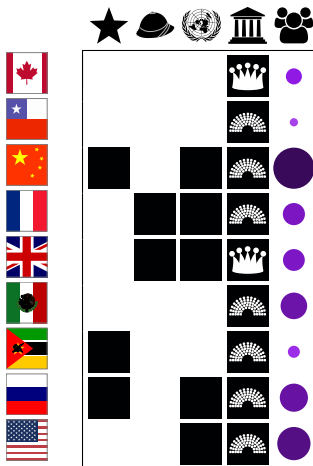


Own work

Extending redescription mining to relational data

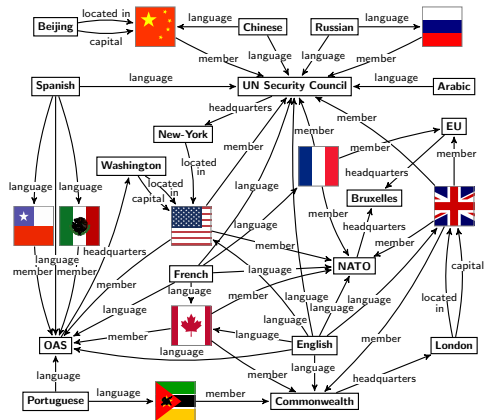
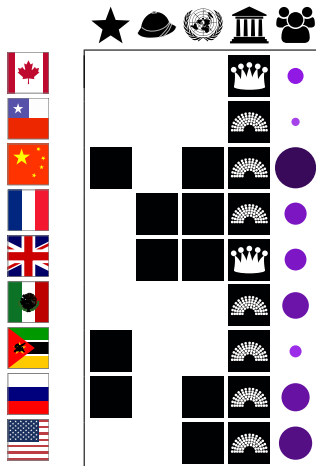
Own work

Geopolitical attributes



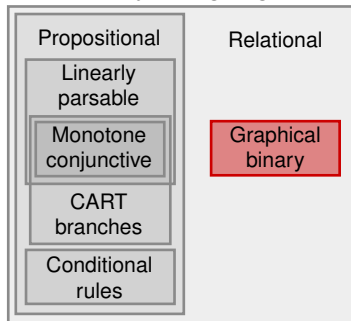
Own work

Geopolitical attributes and relations



Own work

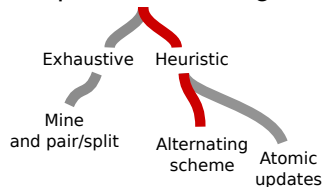
Query Languages



Finding Relational Redescriptions

with Angelika Kimmig,
In *Machine Learning*, 2013.

Exploration Strategies



Example Redescriptions

- Computer science bibliography

Researchers with multiple publications in SoCG and CCCG conferences often collaborate with Profs M. Overmars or E. D. Demaine.

Example Redescriptions

- Computer science bibliography
- Political candidates profiles

Candidates to the 2011 Finnish parliamentary election below age sixty accord little importance to the question of pension indices.

Example Redescriptions

- Computer science bibliography
- Political candidates profiles
- Bioclimatic niches

Scandinavia and Baltia, which are characterized by their specific cold climate, are the habitat of the European Elk.

Bioclimatic Niche Finding

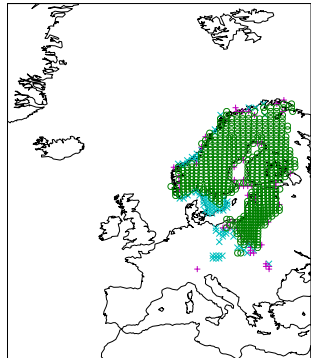
- Dataset:** Spatial land areas of Europe (2575 entities)
- Presence/absence of mammals
(194 species)
 - Climatic data
(48 temperature and rainfall variables)

Question: Find a query over climatic variables that describes the area inhabited by (a group of) mammal species (and vice versa)

Bioclimatic Niche Finding

European Elk

$$([-9.80 \leq t_{\text{Feb}}^{\text{max}} \leq 0.40] \wedge [12.20 \leq t_{\text{Jul}}^{\text{max}} \leq 24.60] \wedge [56.852 \leq p_{\text{Aug}}^{\text{avg}} \leq 136.46]) \vee [183.27 \leq p_{\text{Sep}}^{\text{avg}} \leq 238.78]$$

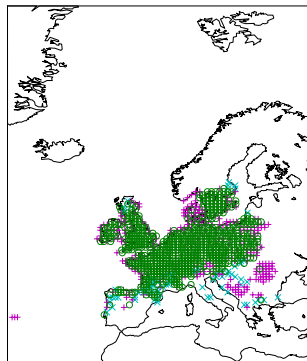


$J = 0.814$ $\text{supp} = 582$

Bioclimatic Niche Finding

Wood Mouse \wedge Natterer's Bat \wedge Eurasian Pygmy Shrew

$$([3.20 \leq t_{\text{Mar}}^{\text{max}} \leq 14.50] \wedge [17.30 \leq t_{\text{Aug}}^{\text{max}} \leq 25.20] \wedge [14.90 \leq t_{\text{Sep}}^{\text{max}} \leq 22.80]) \vee [19.60 \leq t_{\text{Jul}}^{\text{avg}} \leq 19.956]$$



$J = 0.623$ $\text{supp} = 681$

Example Redescriptions

- Computer science bibliography
- Political candidates profiles
- Bioclimatic niches

Scandinavia and Baltia, which are characterized by their specific cold climate, are the habitat of the European Elk.

Example Redescriptions

- Computer science bibliography
- Political candidates profiles
- Bioclimatic niches
- Ethnology

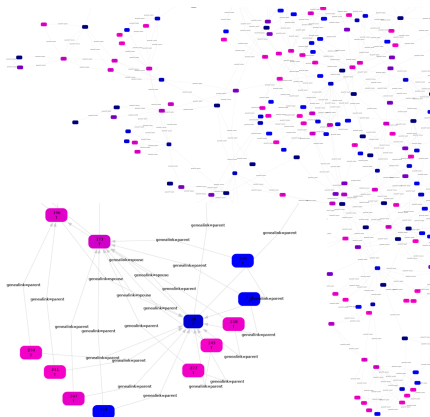
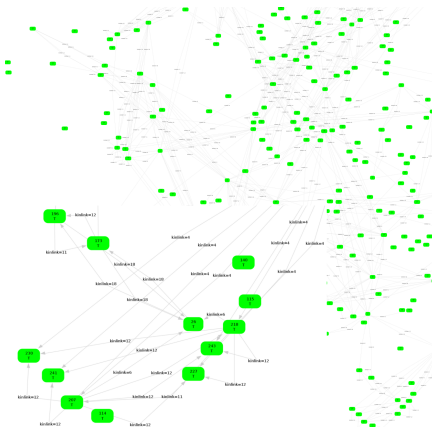
Among Alyawarra, “Aleriya” refers to the son of a male speaker or to the child of the speaker’s brother.

Elicit Kinship Terminology

Dataset: Ethnographic information about Australian *Alyawarra* tribe

- Kinship terminology
- Genealogic, age and sex informations

Elicit Kinship Terminology



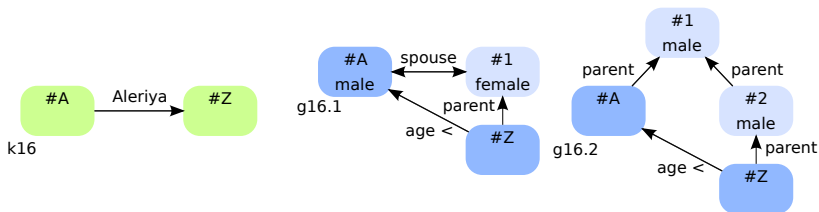
Elicit Kinship Terminology

Dataset: Ethnographic information about Australian *Alyawarra* tribe

- Kinship terminology
- Genealogic, age and sex informations

Question: Elicit the meaning of kinship terms

Elicit Kinship Terminology



Aleriya is used to refer to one's father or one's brother's child

The SIREN interface

The SIREN interface



Visualizing and interactively mining redescrptions

with Pauli Miettinen,

In *Instant Interactive Data Mining Workshop at ECML/PKDD*, 2012.

Demo at KDD 2012 and SIGMOD 2014.

The SIREN interface

SIREN :: tools

Entities	LHS Variables	RHS Variables	Redescriptions	Expansions	History	Log
	id	query LHS	query RHS	J	p-value	E _v t track
1	<input checked="" type="checkbox"/> R7	Wood mouse	$(([2.9 \leq t3+] \vee [9.7 \leq t7+ \leq 13.2]) \wedge [-3.26 \leq t0.836 \wedge [-16.85 \leq t1- \leq 2.6875]$		0.0	1712 0:8;1:14,18
2	<input checked="" type="checkbox"/> R12	Roe Deer	$(([-16.85 \leq t1- \leq 2.6875]$			
3	<input checked="" type="checkbox"/> R17	Red Squirrel	$(([t5- \leq 10.1] \wedge [12.1 \leq t6-$			
4	<input checked="" type="checkbox"/> R16	Eurasian Pygmy Shrew	$[-20.4 \leq t2- \leq 5.8] \wedge [t6-$			
5	<input checked="" type="checkbox"/> R10	Bank Vole	$[14.9 \leq t6+ \leq 27.4] \wedge [34$			
6	<input checked="" type="checkbox"/> R15	European Hare	$(([-6.7722 \leq t3- \leq 10.82]$			
7	<input checked="" type="checkbox"/> R9	Stoat	$[-11.3 \leq t4- \leq 6.8] \wedge [11.$			
8	<input checked="" type="checkbox"/> R13	Common Shrew	$[14.0 \leq t7+ \leq 28.7] \wedge [-10$			
9	<input checked="" type="checkbox"/> R21	Field Vole	$[11.5 \leq t6+ \leq 24.5] \wedge [12$			
10	<input checked="" type="checkbox"/> R30	European Pine Marten	$(([-6.5 \leq t10- \leq 7.5] \wedge [13$			
11	<input checked="" type="checkbox"/> R37	European Hedgehog	$(([4.5 \leq t4+ \leq 22.4] \wedge [4$			
12	<input checked="" type="checkbox"/> R22	Eurasian Water Shrew	$(([14.0 \leq t7+ \leq 26.9] \wedge [4$			
13	<input checked="" type="checkbox"/> R36	European Polecat	$(([18.5 \leq t8+ \leq 28.3] \wedge [4$			
14	<input checked="" type="checkbox"/> R19	Wild boar	$(([-6.1 \leq t11- \leq 10.0] \wedge [1$			
15	<input checked="" type="checkbox"/> R52	European Otter	$(([4.9 \leq t7- \leq 14.0] \wedge [-1.6$			
16	<input checked="" type="checkbox"/> R44	Common Pipistrelle	$(([t7- \leq 17.3] \wedge [8.0 \leq t4+$			
17	<input checked="" type="checkbox"/> R31	European Water Vole	$[-10.8 \leq t1+ \leq 7.1] \wedge [17.$			
18	<input checked="" type="checkbox"/> R27	Beech Marten	$(([-10.5 \leq t1- \leq 5.7] \wedge [20$			
19	<input checked="" type="checkbox"/> R18	European Mole	$(([14.4 \leq t9+ \leq 24.7] \wedge [3$			
20	<input checked="" type="checkbox"/> R41	European Rabbit	$(([0.7375 \leq t3-] \vee [127.$			
21	<input checked="" type="checkbox"/> R33	Yellow-necked Mouse	$[18.8 \leq t8+ \leq 28.1] \vee [1$			
22	<input checked="" type="checkbox"/> R20	House mouse	$(([3.5 \leq t1+ \leq 4.4 \leq t2+$			
23	<input checked="" type="checkbox"/> R51	Red Deer	$(([t10- \leq 7.4] \wedge [9.1 \leq t10$			
24	<input checked="" type="checkbox"/> R32	Brown long-eared bat	$(([13.7 \leq t9+ \leq 22.7] \vee [8$			
25	<input checked="" type="checkbox"/> R23	Common Vole	$(([20.5 \leq t8+ \leq 28.5] \wedge [6$			
26	<input checked="" type="checkbox"/> R24	Common Vole	$(([9.5 \leq t4+ \leq 9.5] \vee [20.$			
27	<input checked="" type="checkbox"/> R40	Daubenton's Bat	$(([14.0 \leq t5+ \leq 20.6] \wedge [1$			
28	<input checked="" type="checkbox"/> R42	Muskrat	$(([-5.5 \leq t3+ \leq -2.0] \vee [19$			
29	<input checked="" type="checkbox"/> R26	American Mink	$(([9.4 \leq t8+ \leq 22.0] \wedge [3$			
30	<input checked="" type="checkbox"/> R29	House mouse	$(([-9.6 \leq t1+ \leq 3.1] \wedge [-1.$			

Stoat

$[-11.3 \leq t4- \leq 6.8] \wedge [11.6 \leq t8+ \leq 25.3]$

J = 0.786 |E_v| = 116 |E_v| = 1367 |E_v| = 257
p-value = 0.000 |E_v| = 1740 |E_v| = 835 |E_v| = 2575

Expand - opac. disabled +

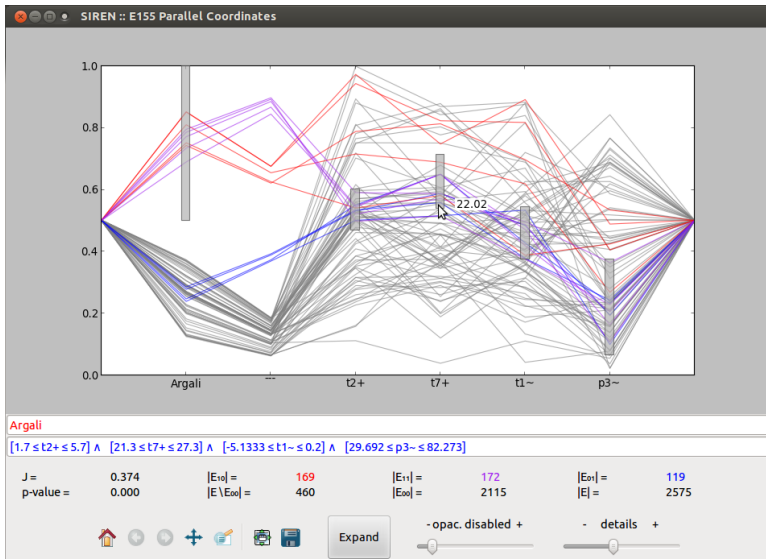
@13: Done...

The SIREN interface

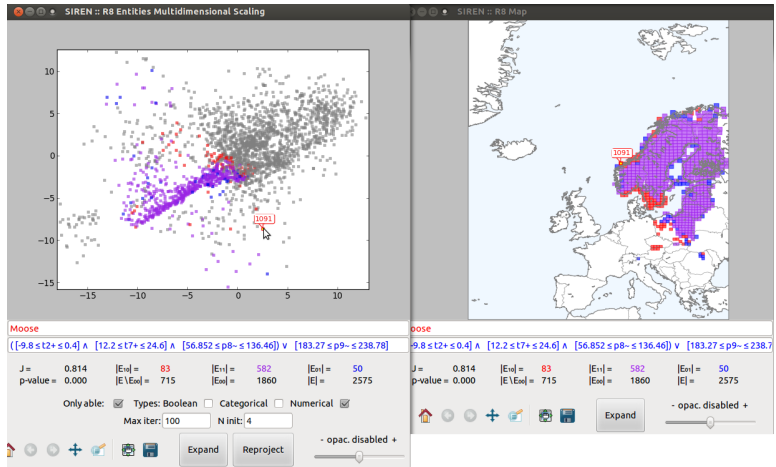
SIREN :: tools											
Entities	LHS Variables			RHS Variables			Redescriptions				
		id		LHS:Wisent ↑		LHS:Wild goat				LHS:European Hedgehog	
1	<input checked="" type="checkbox"/>	1464		True		False				False	
2	<input checked="" type="checkbox"/>	1494		True		False				False	
3	<input checked="" type="checkbox"/>	1891		True		False				False	
4	<input checked="" type="checkbox"/>	2000		True		False				False	
5	<input checked="" type="checkbox"/>	2005		True		False				False	
6	<input checked="" type="checkbox"/>	2016		True		False				False	
7	<input checked="" type="checkbox"/>	2032		True		False				False	
8	<input type="checkbox"/>	2033		True		False				False	
9	<input type="checkbox"/>	2037		True		False				False	
10	<input type="checkbox"/>	2049		True		False				False	
11	<input type="checkbox"/>	2050		True		False				False	
12	<input type="checkbox"/>	2264		True		False				False	
13	<input type="checkbox"/>	2270		True		False				False	
14	<input checked="" type="checkbox"/>	2311		True		False				False	
15	<input checked="" type="checkbox"/>	2357		True		False				False	
16	<input checked="" type="checkbox"/>	2361		True		False				False	
17	<input checked="" type="checkbox"/>	2362		True		False				False	
18	<input checked="" type="checkbox"/>	0		False		False				True	
19	<input checked="" type="checkbox"/>	1		False		False				True	
20	<input checked="" type="checkbox"/>	2		False		False				True	
21	<input checked="" type="checkbox"/>	3		False		False				True	
22	<input checked="" type="checkbox"/>	4		False		False				False	
23	<input checked="" type="checkbox"/>	5		False		False				True	
24	<input checked="" type="checkbox"/>	6		False		False				True	
25	<input checked="" type="checkbox"/>	7		False		False				True	

Loading done

The SIREN interface

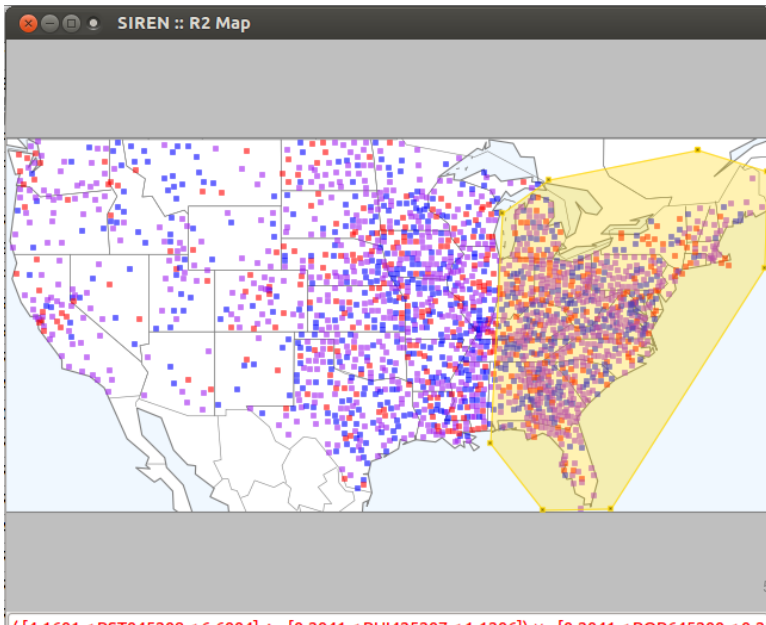


The SIREN interface



The SIREN interface

The SIREN interface





Conclusion

Redescription Mining is a versatile and powerful data-mining tool, applicable in various domains.

For more details:

- `galbrun@cs.helsinki.fi`

- `http://www.cs.helsinki.fi/u/galbrun/`

Conclusion

Redescription Mining is a versatile and powerful data-mining tool, applicable in various domains.

For more details:

- `galbrun@cs.helsinki.fi`

- `http://www.cs.helsinki.fi/u/galbrun/`

Thank you!