

# Course Outline

## I. Introduction

- Data Mining and KDD process
- Introduction to Data Mining
- Data Mining platforms

## II. Predictive DM Techniques

- Decision Tree learning
- Bayesian classifier
- Classification rule learning
- Classifier Evaluation

## III. Regression

## IV. Descriptive DM

- Predictive vs. descriptive induction
- Subgroup discovery
- Association rule learning
- Hierarchical clustering

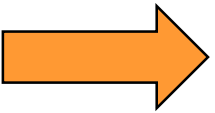
## V. Relational Data Mining

- RDM and Inductive Logic Programming
- Propositionalization
- Semantic data mining

## VI. Advanced Topics

# Part V:

## Relational Data Mining



What is RDM

- Propositionalization techniques
- Semantic Data Mining

# Relational Data Mining (Inductive Logic Programming) task

customer							
ID	Zip	Sex	SoSt	Income	Age	Club	Resp
...	...	...	...	...	...	...	...
3478	34677	m	si	60-70	32	me	nr
3479	43666	f	ma	80-90	45	nm	re
...	...	...	...	...	...	...	...

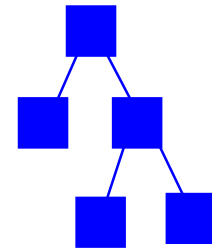
order				
Customer ID	Order ID	Store ID	Delivery Mode	Paymt Mode
...	...	...	...	...
3478	2140267	12	regular	cash
3478	3446778	12	express	check
3478	4728386	17	regular	check
3479	3233444	17	express	credit
3479	3475886	12	regular	credit
...	...	...	...	...

store			
Store ID	Size	Type	Location
...	...	...	...
12	small	franchise	city
17	large	indep	rural
...	...	...	...

knowledge discovery  
from data

Relational Data Mining



model, patterns, ...

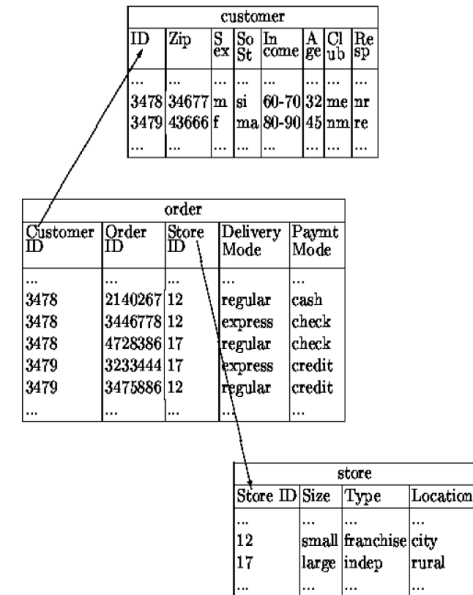
Relational representation of customers, orders and stores.

**Given:** a relational database, a set of tables. sets of logical facts, a graph, ...

**Find:** a classification model, a set of interesting patterns

# Relational data mining

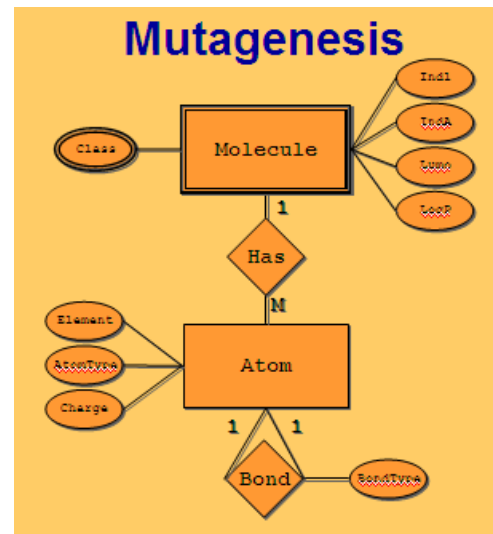
- ILP, relational learning, relational data mining
  - Learning from complex multi-relational data



Relational representation of customers, orders and stores.

# Relational data mining

- ILP, relational learning, relational data mining
  - Learning from complex multi-relational data
  - Learning from complex structured data: e.g., molecules and their biochemical properties



customer						
ID	Zip	Sex	Income	Age	Club	Resp
...	...	...	...	...	...	...
3478	34677	m	60-70	32	me	nr
3479	43666	f	ma	80-90	45	nr
...	...	...	...	...	...	...

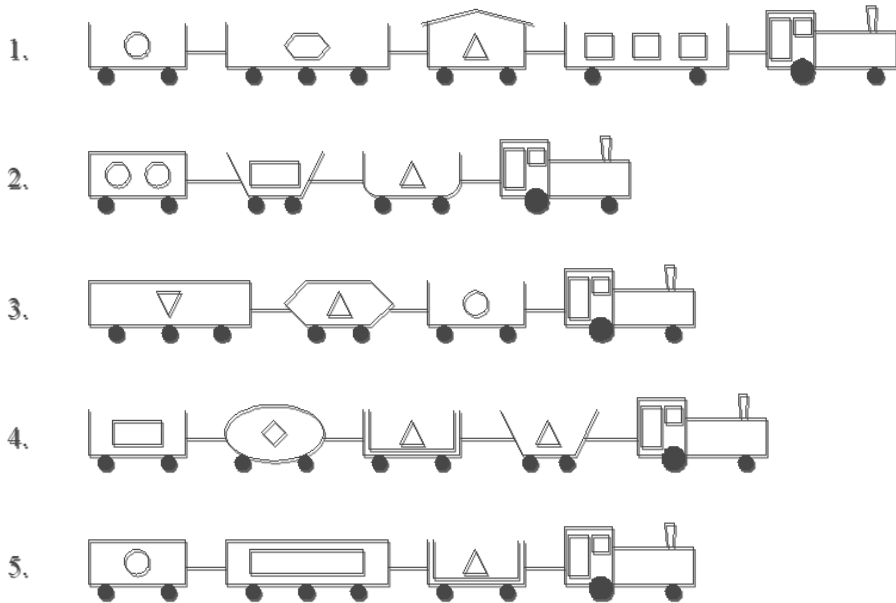
order				
Customer ID	Order ID	Store ID	Delivery Mode	Paymt Mode
...	...	...	...	...
3478	2140267	12	regular	cash
3478	3446778	12	express	check
3478	4728386	17	regular	check
3479	3233444	17	express	credit
3479	3475886	12	regular	credit
...	...	...	...	...

store			
Store ID	Size	Type	Location
...	...	...	...
12	small	franchise	city
17	large	indep	rural
...	...	...	...

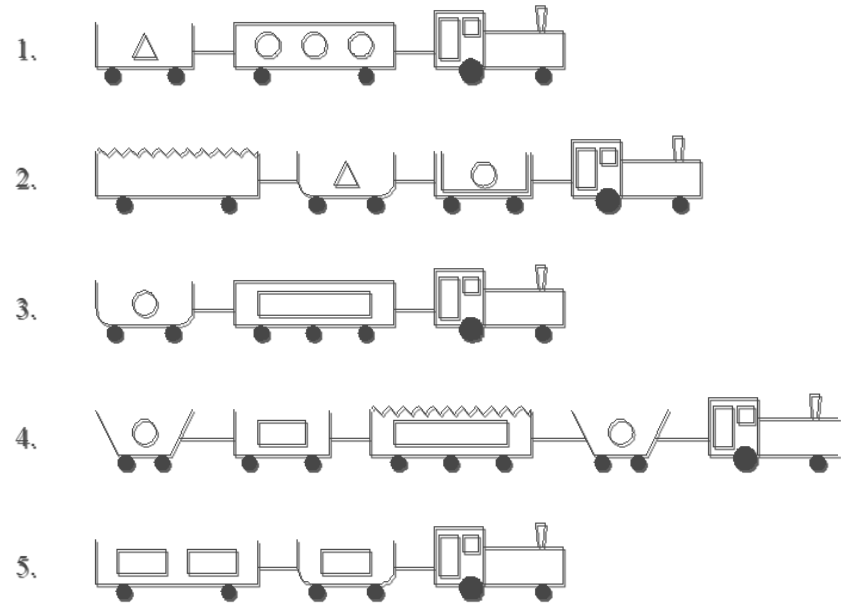
Relational representation of customers, orders and stores.

# Sample problem: East-West trains

## 1. TRAINS GOING EAST



## 2. TRAINS GOING WEST



# RDM knowledge representation (database)

## LOAD\_TABLE

LOAD	CAR	OBJECT	NUMBER
l1	c1	circle	1
l2	c2	hexagon	1
l3	c3	triangle	1
l4	c4	rectangle	3
...	...	...	...

## TRAIN\_TABLE

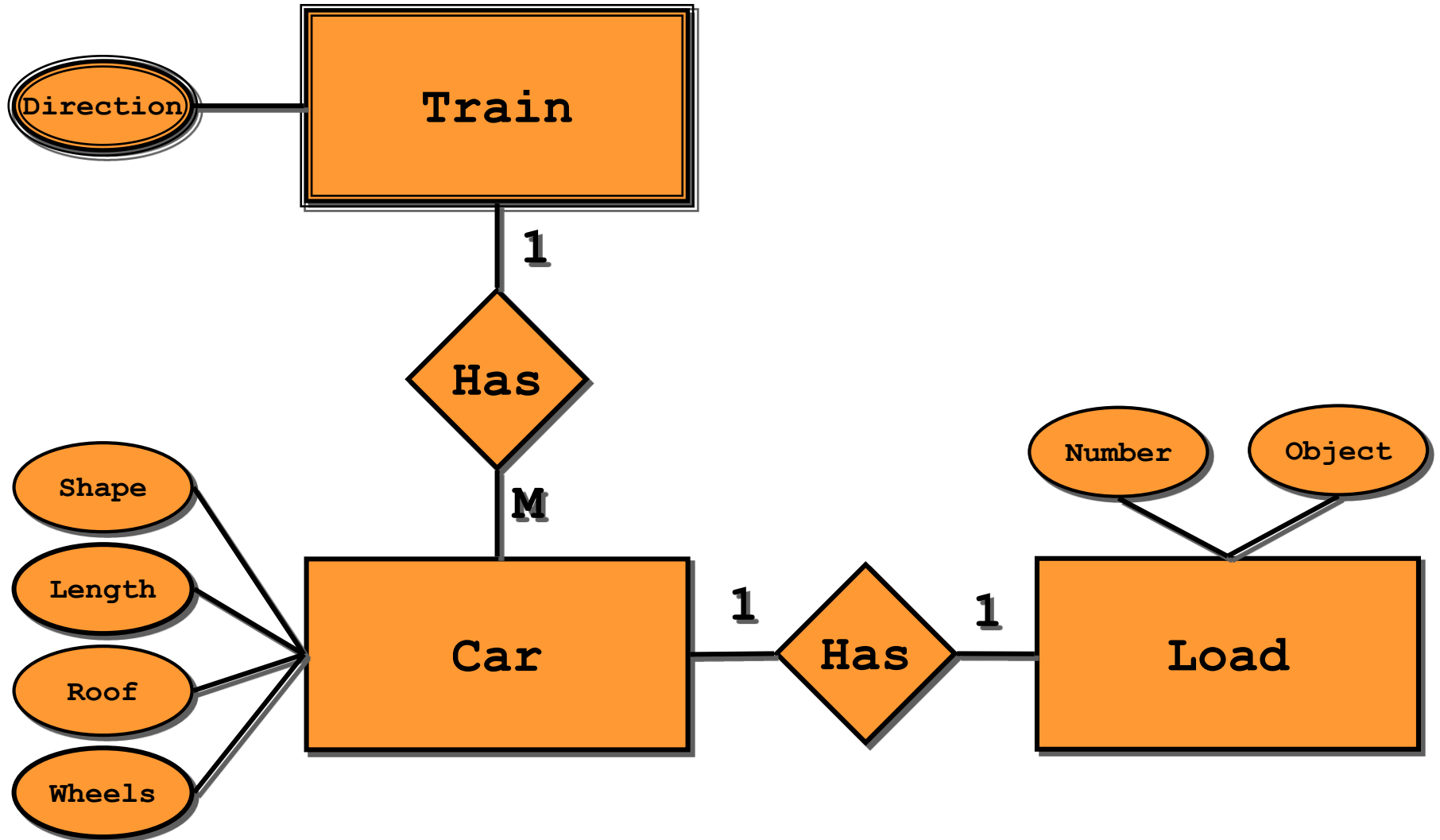
TRAIN	EASTBOUND
t1	TRUE
t2	TRUE
...	...
t6	FALSE
...	...

## CAR\_TABLE

CAR	TRAIN	SHAPE	LENGTH	ROOF	WHEELS
c1	t1	rectangle	short	none	2
c2	t1	rectangle	long	none	3
c3	t1	rectangle	short	peaked	2
c4	t1	rectangle	long	none	2
...	...	...	...	...	...



# ER diagram for East-West trains





# Relational data mining

- Relational data mining is characterized by using background knowledge (domain knowledge) in the data mining process
- Selected approaches:
  - Inductive logic programming - ILP (Muggleton, 1991; Lavrač & Džeroski 1994), ...
  - Relational learning (Quinlan, 1993)
  - Learning in DL (Lisi 2004), ...
  - Relational Data Mining (Džeroski & Lavrač, 2001),
  - Statistical relational learning (Domingos, De Raedt...)
  - Propositionalization approach to RDM (Lavrač et al.)

# Our early work:

## Semantic subgroup discovery

- Propositionalization approach: Using relational subgroup discovery in the SDM context
  - General purpose system **RSD** for **Relational Subgroup Discovery**, using a propositionalization approach to relational data mining
  - Applied to semantic data mining in a biomedical application by using the Gene Ontology as background knowledge in analyzing microarray data

(Železny and Lavrač, MLJ 2006)

# Part V: Relational Data Mining

- What is RDM
-  Propositionalization techniques
- Semantic Data Mining

# Relational Data Mining through Propositionalization

Step 1

Propositionalization

customer							
ID	Zip	Sex	St	Income	Age	Club	Resp
...	...	...	...	...	...	...	...
3478	34677	m	si	60-70	32	me	nr
3479	43666	f	ma	80-90	45	nm	re
...	...	...	...	...	...	...	...

order				
Customer ID	Order ID	Store ID	Delivery Mode	Paymt Mode
...	...	...	...	...
3478	2140267	12	regular	cash
3478	3446778	12	express	check
3478	4728386	17	regular	check
3479	3233444	17	express	credit
3479	3475886	12	regular	credit
...	...	...	...	...

store			
Store ID	Size	Type	Location
...	...	...	...
12	small	franchise	city
17	large	indep	rural
...	...	...	...

Relational representation of customers, orders and stores.

	f1	f2	f3	f4	f5	f6						f <sub>n</sub>
g1	1	0	0	1	1	1	0	0	1	0	1	1
g2	0	1	1	0	1	1	0	0	0	1	1	0
g3	0	1	1	1	0	0	1	1	0	0	0	1
g4	1	1	1	0	1	1	0	0	1	1	1	0
g5	1	1	1	0	0	1	0	1	1	0	1	0
g1	0	0	1	1	0	0	0	1	0	0	0	1
g2	1	1	0	0	1	1	0	1	0	1	1	1
g3	0	0	0	0	1	0	0	1	1	1	0	0
g4	1	0	1	1	1	0	1	0	0	1	0	1

# Relational Data Mining through Propositionalization

Step 1

Propositionalization

customer							
ID	Zip	Sex	State	Income	Age	Club	Response
...	...	...	...	...	...	...	...
3478	34677	m	si	60-70	32	me	nr
3479	43666	f	ma	80-90	45	nm	re
...	...	...	...	...	...	...	...

order				
Customer ID	Order ID	Store ID	Delivery Mode	Paymt Mode
...	...	...	...	...
3478	2140267	12	regular	cash
3478	3446778	12	express	check
3478	4728386	17	regular	check
3479	3233444	17	express	credit
3479	3475886	12	regular	credit
...	...	...	...	...

store			
Store ID	Size	Type	Location
...	...	...	...
12	small	franchise	city
17	large	indep	rural
...	...	...	...

Relational representation of customers, orders and stores.

	f1	f2	f3	f4	f5	f6	...	...	...	...	fn
g1	1	0	0	1	1	1	0	0	1	0	1
g2	0	1	1	0	1	1	0	0	0	1	1
g3	0	1	1	1	0	0	1	1	0	0	1
g4	1	1	1	0	1	1	0	0	1	1	1
g5	1	1	1	0	0	1	0	1	1	0	1
g1	0	0	1	1	0	0	0	1	0	0	1
g2	1	1	0	0	1	1	0	1	0	1	1
g3	0	0	0	0	1	0	0	1	1	1	0
g4	1	0	1	1	1	0	1	0	0	1	1

1. constructing relational features
2. constructing a propositional table

# Relational Data Mining through Propositionalization

customer							
ID	Zip	Sex	Status	Income	Age	Club	Response
...	...	...	...	...	...	...	...
3478	34677	m	si	60-70	32	me	nr
3479	43666	f	ma	80-90	45	nm	re
...	...	...	...	...	...	...	...

order				
Customer ID	Order ID	Store ID	Delivery Mode	Payment Mode
...	...	...	...	...
3478	2140267	12	regular	cash
3478	3446778	12	express	check
3478	4728386	17	regular	check
3479	3233444	17	express	credit
3479	3475886	12	regular	credit
...	...	...	...	...

store			
Store ID	Size	Type	Location
...	...	...	...
12	small	franchise	city
17	large	indep	rural
...	...	...	...

Relational representation of customers, orders and stores.

Step 1

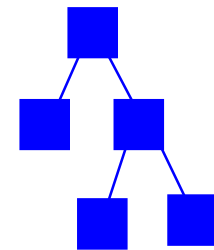
Propositionalization

	f1	f2	f3	f4	f5	f6	...	...	...	...	fn
g1	1	0	0	1	1	1	0	0	1	0	1
g2	0	1	1	0	1	1	0	0	0	1	1
g3	0	1	1	1	0	0	1	1	0	0	1
g4	1	1	1	0	1	1	0	0	1	1	1
g5	1	1	1	0	0	1	0	1	1	0	1
g1	0	0	1	1	0	0	0	1	0	0	1
g2	1	1	0	0	1	1	0	1	0	1	1
g3	0	0	0	0	1	0	0	1	1	1	0
g4	1	0	1	1	1	0	1	0	0	1	1

Step 2

Data Mining

	f1	f2	f3	f4	f5	f6	...	...	...	...	fn
g1	1	0	0	1	1	1	0	0	1	0	1
g2	0	1	1	0	1	1	0	0	0	1	1
g3	0	1	1	1	0	0	1	1	0	0	1
g4	1	1	1	0	1	1	0	0	1	1	1
g5	1	1	1	0	0	1	0	1	1	0	1
g1	0	0	1	1	0	0	0	1	0	0	1
g2	1	1	0	0	1	1	0	1	0	1	1
g3	0	0	0	0	1	0	0	1	1	1	0
g4	1	0	1	1	1	0	1	0	0	1	1



model, patterns, ...

# Relational Data Mining through Propositionalization

customer							
ID	Zip	Sex	Status	Income	Age	Club	Rep
...	...	...	...	...	...	...	...
3478	34677	m	si	60-70	32	me	nr
3479	43666	f	ma	80-90	45	nm	re
...	...	...	...	...	...	...	...

order				
Customer ID	Order ID	Store ID	Delivery Mode	Paymt Mode
...	...	...	...	...
3478	2140267	12	regular	cash
3478	3446778	12	express	check
3478	4728386	17	regular	check
3479	3233444	17	express	credit
3479	3475886	12	regular	credit
...	...	...	...	...

store			
Store ID	Size	Type	Location
...	...	...	...
12	small	franchise	city
17	large	indep	rural
...	...	...	...

Relational representation of customers, orders and stores.

Step 1

Propositionalization

	f1	f2	f3	f4	f5	f6					fn
g1	1	0	0	1	1	1	0	0	1	0	1
g2	0	1	1	0	1	1	0	0	0	1	1
g3	0	1	1	1	0	0	1	1	0	0	1
g4	1	1	1	0	1	1	0	0	1	1	1
g5	1	1	1	0	0	1	0	1	1	0	1
g1	0	0	1	1	0	0	0	1	0	0	1
g2	1	1	0	0	1	1	0	1	0	1	1
g3	0	0	0	0	1	0	0	1	1	1	0
g4	1	0	1	1	1	0	1	0	0	1	1

Step 2

Data Mining

	f1	f2	f3	f4	f5	f6					fn
g1	1	0	0	1	1	1	0	0	1	0	1
g2	0	1	1	0	1	1	0	0	0	1	1
g3	0	1	1	1	0	0	1	1	0	0	1
g4	1	1	1	0	1	1	0	0	1	1	1
g5	1	1	1	0	0	1	0	1	1	0	1
g1	0	0	1	1	0	0	0	1	0	0	1
g2	1	1	0	0	1	1	0	1	0	1	1
g3	0	0	0	0	1	0	0	1	1	1	0
g4	1	0	1	1	1	0	1	0	0	1	1

```
target(A) :-
    'Doctor'(A), 'Italy'(A).
```

```
target(A) :-
    'Public'(A), 'Gold'(A).
```

```
target(A) :-
    'Poland'(A), 'Deposit'(A), 'Gold'(A).
```

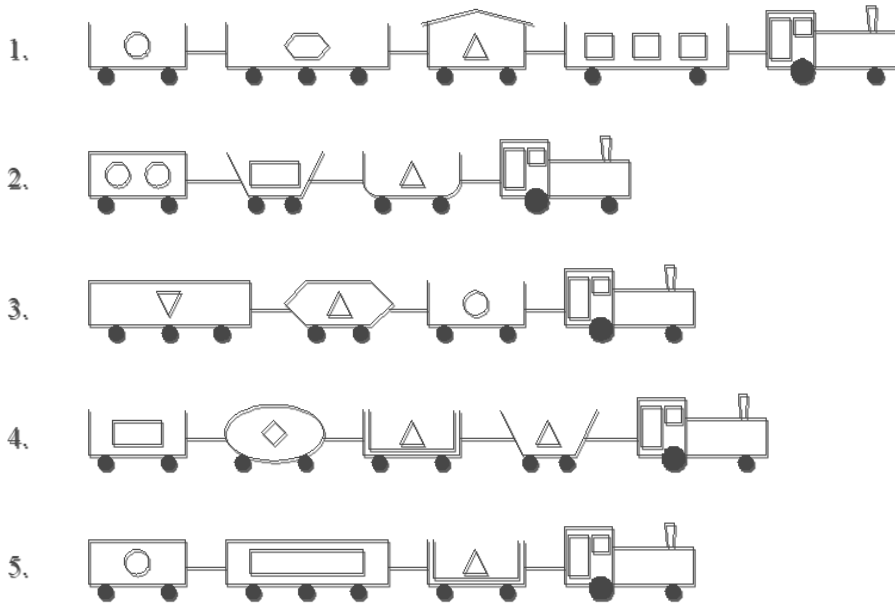
```
target(A) :-
    'Germany'(A), 'Insurance'(A).
```

```
target(A) :-
    'Service'(A), 'Germany'(A).
```

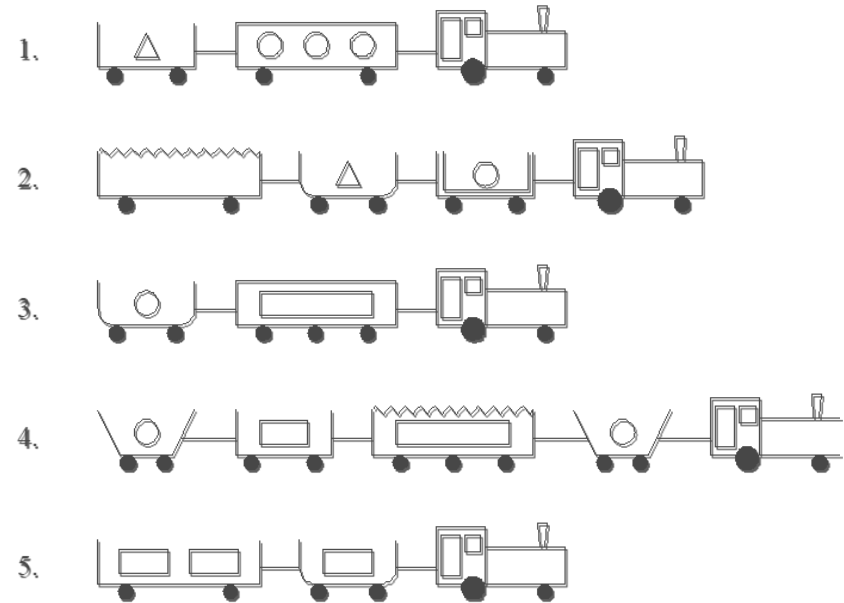
patterns (set of rules)

# Sample ILP problem: East-West trains

## 1. TRAINS GOING EAST



## 2. TRAINS GOING WEST





# Relational data representation



LOAD	CAR	OBJECT	NUMBER
l1	c1	circle	1
l2	c2	hexagon	1
l3	c3	triangle	1
l4	c4	rectangle	3
...	...	...	...

**TRAIN\_TABLE**

TRAIN	EASTBOUND
t1	TRUE
t2	TRUE
...	...
t6	FALSE
...	...

CAR	TRAIN	SHAPE	LENGTH	ROOF	WHEELS
c1	t1	rectangle	short	none	2
c2	t1	rectangle	long	none	3
c3	t1	rectangle	short	peaked	2
c4	t1	rectangle	long	none	2
...	...	...	...	...	...

# Propositionalization in a nutshell



## Propositionalization task

**Transform** a multi-relational  
(**multiple-table**)  
representation to a  
propositional representation  
(**single table**)

LOAD	CAR	OBJECT	NUMBER
l1	c1	circle	1
l2	c2	hexagon	1
l3	c3	triangle	1
l4	c4	rectangle	3
...	...	...	...

## TRAIN\_TABLE

TRAIN	EASTBOUND
t1	TRUE
t2	TRUE
...	...
t6	FALSE
...	...

CAR	TRAIN	SHAPE	LENGTH	ROOF	WHEELS
c1	t1	rectangle	short	none	2
c2	t1	rectangle	long	none	3
c3	t1	rectangle	short	peaked	2
c4	t1	rectangle	long	none	2
...	...	...	...	...	...

Proposed in ILP systems

LINUS (Lavrač et al. 1991, 1994),  
1BC (Flach and Lachiche 1999), ...

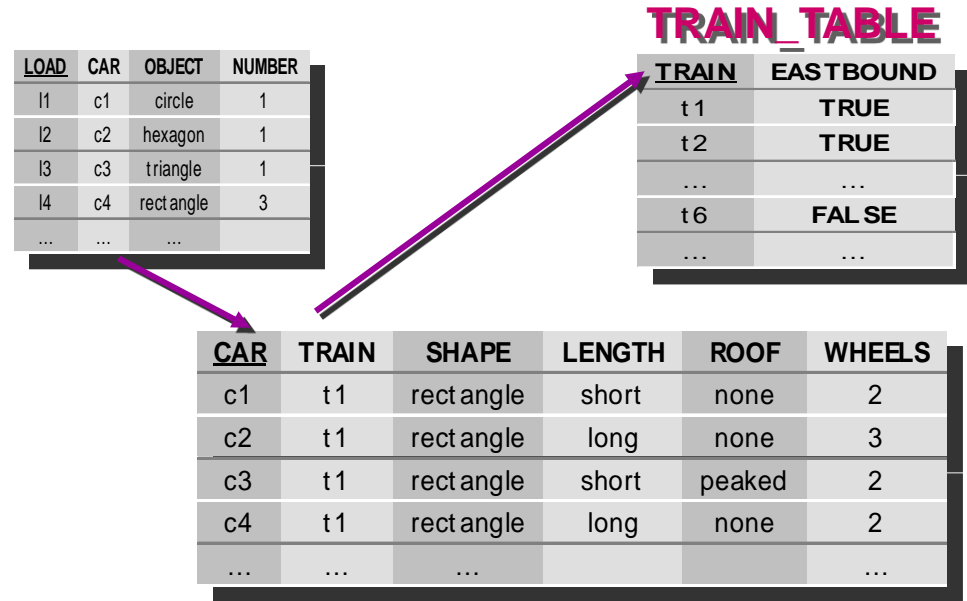
# Propositionalization in a nutshell

**Main propositionalization step:  
first-order feature construction**

$f_1(T) :- \text{hasCar}(T,C), \text{clength}(C, \text{short}).$

$f_2(T) :- \text{hasCar}(T,C), \text{hasLoad}(C,L),$   
 $\text{loadShape}(L, \text{circle})$

$f_3(T) :- \dots$



**Propositional learning:**

$t(T) \leftarrow f_1(T), f_4(T)$

**Relational interpretation:**

$\text{eastbound}(T) \leftarrow$

$\text{hasShortCar}(T), \text{hasClosedCar}(T).$

**PROPOSITIONAL TRAIN\_TABLE**

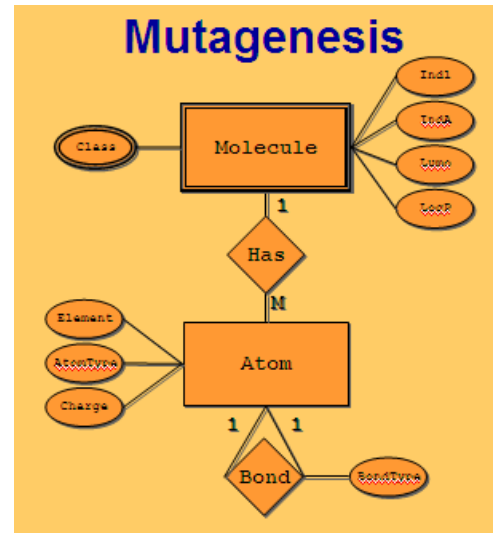
<u>train(T)</u>	f1(T)	f2(T)	f3(T)	f4(T)	f5(T)
t1	t	t	f	t	t
t2	t	t	t	t	t
t3	f	f	t	f	f
t4	t	f	t	f	f
...	...	...			...

# Part V: Relational Data Mining

- What is RDM
- Propositionalization techniques
- Semantic Data Mining

# Semantic data mining

- **ILP, relational learning, relational data mining**
  - Learning from complex multi-relational data
  - Learning from complex structured data: e.g., molecules and their biochemical properties
  - Learning by using domain knowledge in the form of ontologies = **semantic data mining**



customer							
ID	Zip	Sex	St	In come	Age	Cl ub	Re sp
...	...	...	...	...	...	...	...
3478	34677	m	si	60-70	32	me	nr
3479	43666	f	ma	80-90	45	nm	re
...	...	...	...	...	...	...	...

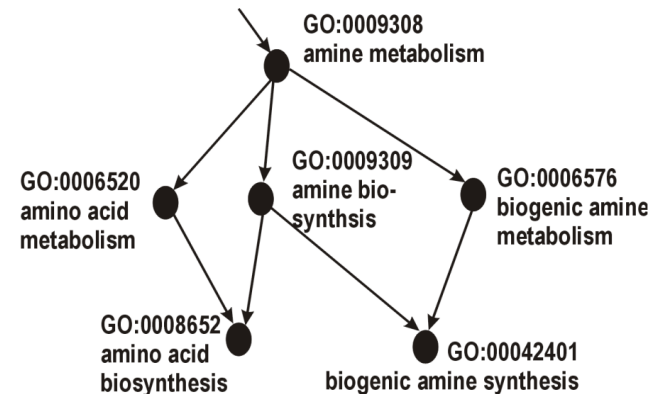
  

order				
Customer ID	Order ID	Store ID	Delivery Mode	Paymt Mode
...	...	...	...	...
3478	2140267	12	regular	cash
3478	3446778	12	express	check
3478	4728386	17	regular	check
3479	3233444	17	express	credit
3479	3475886	12	regular	credit
...	...	...	...	...

store			
Store ID	Size	Type	Location
...	...	...	...
12	small	franchise	city
17	large	indep	rural
...	...	...	...

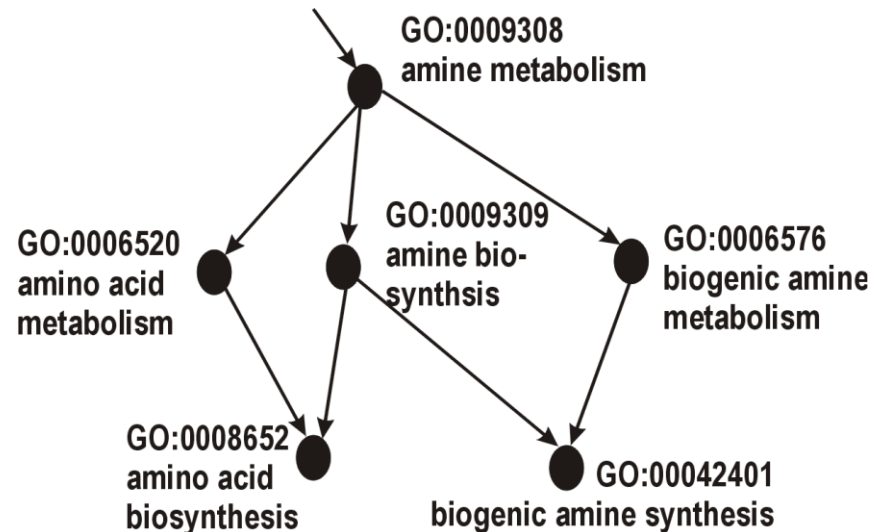
Relational representation of customers, orders and stores.



# Using domain ontologies in Semantic Data Mining

Using domain ontologies as background knowledge, e.g., using the Gene Ontology (GO)

- GO is a database of terms, describing gene sets in terms of their
  - functions (12,093)
  - processes (1,812)
  - components (7,459)
- Genes are annotated to GO terms
- Terms are connected (is\_a, part\_of)
- Levels represent terms generality

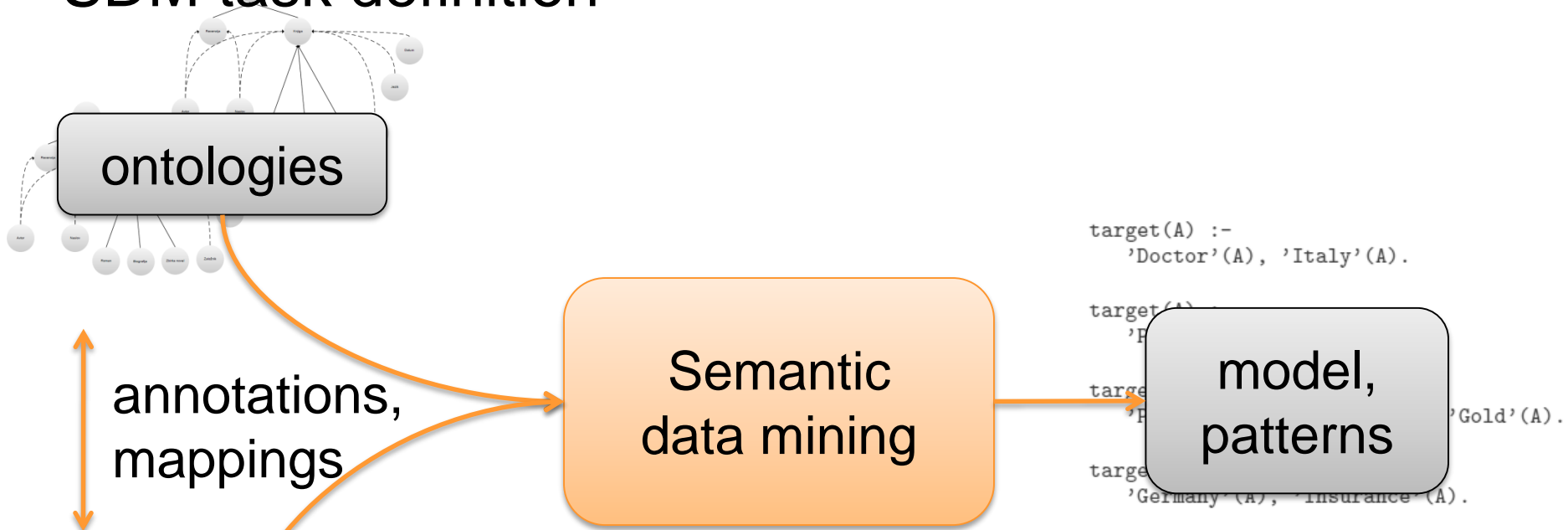


# What is Semantic Data Mining

- Ontology-driven (semantic) data mining is an emerging research topic
- Semantic Data Mining (SDM) - a new term denoting:
  - the new challenge of mining semantically annotated resources, with ontologies used as background knowledge to data mining
  - approaches with which semantic data are mined

# What is Semantic Data Mining

## SDM task definition



ID	occupation	location	account	loan	deposit	tax_deduct	invest	disposable	
1	Doctor	Milan	Cluster	No	Full	Irregular	Family	YES	
2	Doctor	Krakow	Gold	Car	ShortTerm	No	No	YES	
3	Military	Munich	Gold	No	No	No	Regular	YES	
4	Doctor	Catanzaro	Cluster	Car	LongTerm	Irregular	Senior	YES	
5	Energy	Prague	Gold	Apartment	LongTerm	No	No	YES	
6	Doctor	Rome	Gold	Apartment	ShortTerm	No	Regular	YES	
7	Finance	Berlin	Gold	No	ShortTerm	Gold	No	YES	
8	Health-care	Frankfurt	Cluster	Car	No	Child	Share	Family	YES
9	Military	Warsaw	Gold	No	ShortTerm	No	Regular	YES	
10	Education	Lublin	Gold	Apartment	ShortTerm	No	Family	YES	
11	Health-care	Kielce	Cluster	Apartment	No	GoldShare	No	YES	
12	Retail	Munich	Cluster	Car	LongTerm	IrregularShare	Regular	YES	
13	Education								
14	Doctor								
15	Police								
16	Retail								
17	Finance								
18	Doctor								
19	Manufacturing								
20	Doctor								
21	Admission								
22	Unemployed								
23	Military								
24	Manufacturing								
25	Police								
26	Police								
27	Police								
28	Transport								
29	Transport								
30	Police	Warsaw	Gold	Car	ShortTerm	IrregularShare	Regular	No	
		Catanzaro	Cluster	Car	No	No	No	No	

data

Semantic  
data mining

model,  
patterns

### Given:

- transaction data table, relational database, text documents, Web pages, ...
- one or more domain ontologies

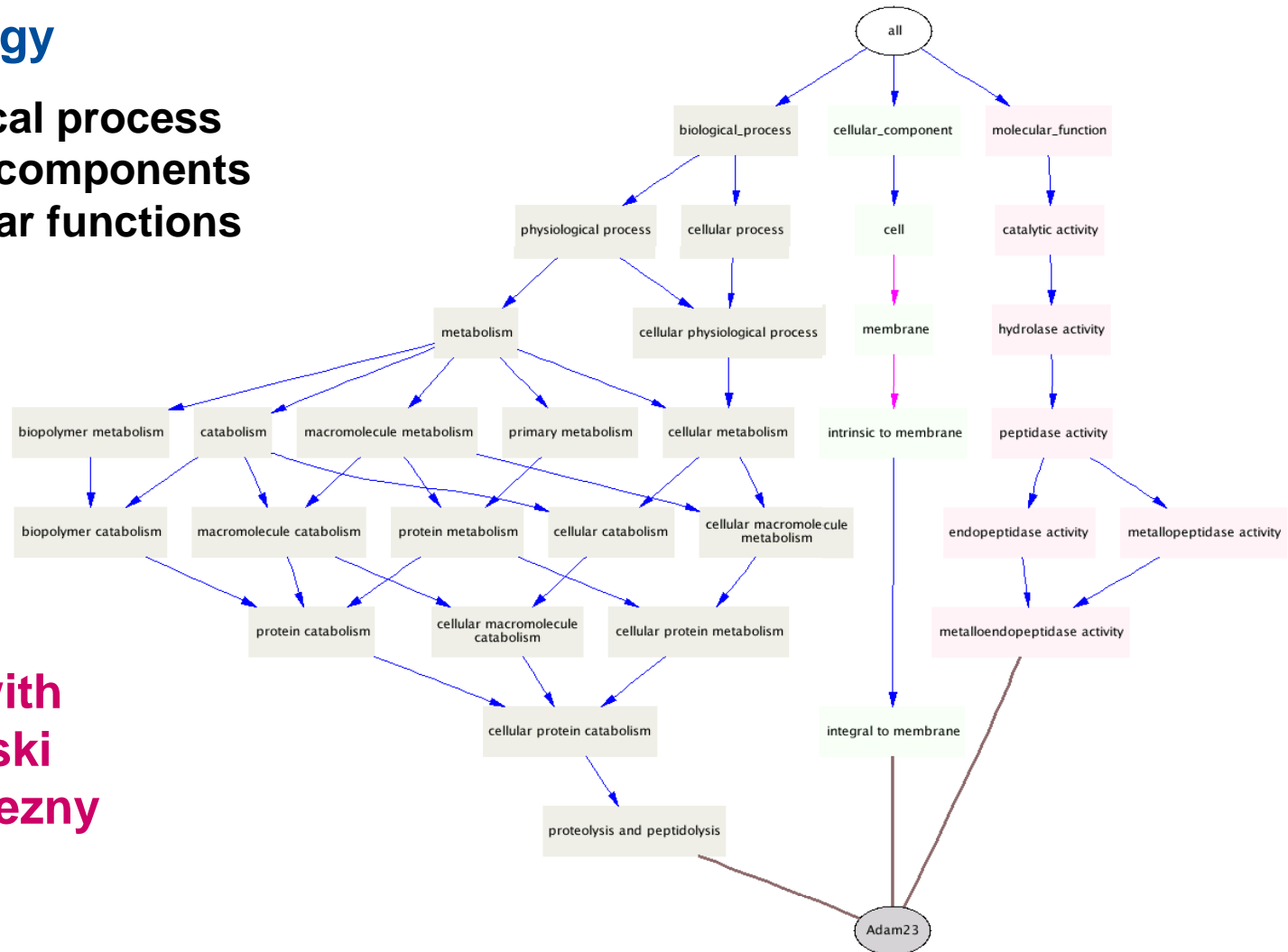
**Find:** a classification model, a set of patterns



# Using domain ontologies (e.g. Gene Ontology) as background knowledge for Data Mining

## Gene Ontology

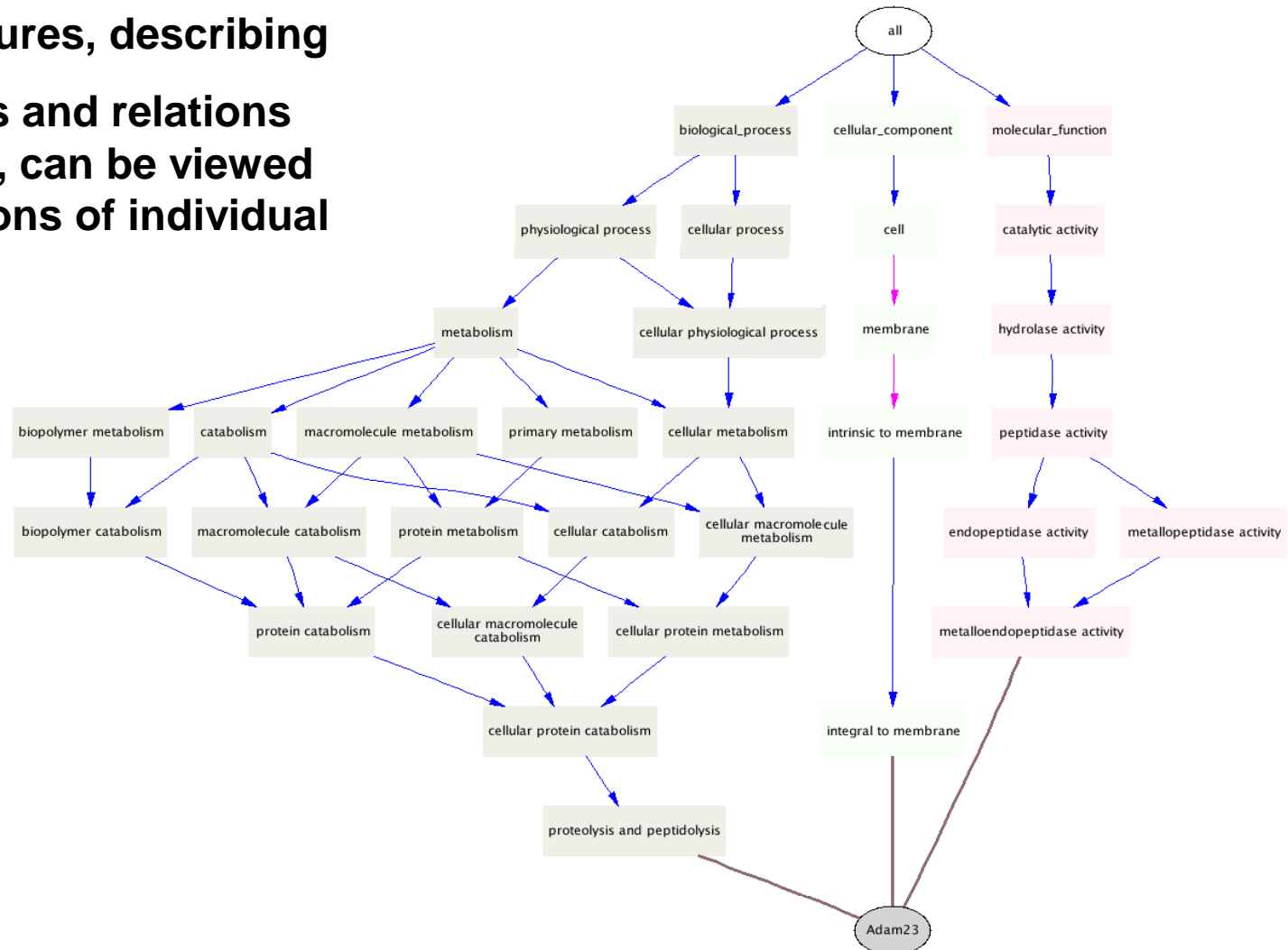
12093 biological process  
1812 cellular components  
7459 molecular functions



Joint work with  
Igor Trajkovski  
and Filip Zelezny

# Using domain ontologies (e.g. Gene Ontology) as background knowledge for Data Mining

First-order features, describing gene properties and relations between genes, can be viewed as generalisations of individual genes



# Semantic subgroup discovery with RSD

1. Take ontology terms represented as logical facts in Prolog, e.g.

```
component (gene2532, 'GO:0016020') .  
function (gene2534, 'GO:0030554') .  
process (gene2534, 'GO:0007243') .  
interaction (gene2534, gene4803) .
```

2. Automatically generate generalized relational features:

```
f(2,A):-component(A,'GO:0016020') .  
f(7,A):-function(A,'GO:0030554') .  
f(11,A):-process(A,'GO:0007243') .  
f(224,A):- interaction(A,B), function(B,'GO:0016787'),  
            component(B,'GO:0043231') .
```

3. Propositionalization: Determine truth values of features
4. Learn rules by a subgroup discovery algorithm CN2-SD

## Step 2: RSD feature construction

Construction of first order features, with support  $> \textit{min\_support}$

f(7,A):-function(A,'GO:0046872').  
f(8,A):-function(A,'GO:0004871').  
f(11,A):-process(A,'GO:0007165').  
f(14,A):-process(A,'GO:0044267').  
f(15,A):-process(A,'GO:0050874').  
f(20,A):-function(A,'GO:0004871'), process(A,'GO:0050874').  
f(26,A):-component(A,'GO:0016021').  
f(29,A):- function(A,'GO:0046872'), component(A,'GO:0016020')  
f(122,A):-interaction(A,B),function(B,'GO:0004872').  
f(223,A):-interaction(A,B),function(B,'GO:0004871'),  
process(B,'GO:0009613').  
f(224,A):-interaction(A,B),function(B,'GO:0016787'),  
component(B,'GO:0043231').

existential



# Step 3: RSD Propositionalization

diffexp g1 (gene64499)

diffexp g2 (gene2534)

diffexp g3 (gene5199)

diffexp g4 (gene1052)

diffexp g5 (gene6036)

....

random g1 (gene7443)

random g2 (gene9221)

random g3 (gene2339)

random g4 (gene9657)

random g5 (gene19679)

....

	f1	f2	f3	f4	f5	f6	...				...	fn
g1	1	0	0	1	1	1	0	0	1	0	1	1
g2	0	1	1	0	1	1	0	0	0	1	1	0
g3	0	1	1	1	0	0	1	1	0	0	0	1
g4	1	1	1	0	1	1	0	0	1	1	1	0
g5	1	1	1	0	0	1	0	1	1	0	1	0
g1	0	0	1	1	0	0	0	1	0	0	0	1
g2	1	1	0	0	1	1	0	1	0	1	1	1
g3	0	0	0	0	1	0	0	1	1	1	0	0
g4	1	0	1	1	1	0	1	0	0	1	0	1

# Step 4: RSD rule construction with CN2-SD

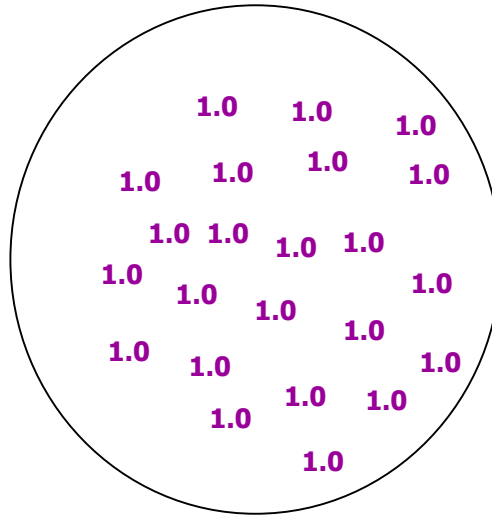
	f1	f2	f3	f4	f5	f6	...				...	fn
g1	1	0	0	1	1	1	0	0	1	0	1	1
g2	0	1	1	0	1	1	0	0	0	1	1	0
g3	0	1	1	1	0	0	1	1	0	0	0	1
g4	1	1	1	0	1	1	0	0	1	1	1	0
g5	1	1	1	0	0	1	0	1	1	0	1	0
g1	0	0	1	1	0	0	0	1	0	0	0	1
g2	1	1	0	0	1	1	0	1	0	1	1	1
g3	0	0	0	0	1	0	0	1	1	1	0	0
g4	1	0	1	1	1	0	1	0	0	1	0	1

Over-  
expressed  
IF  
f2 and f3  
[4,0]

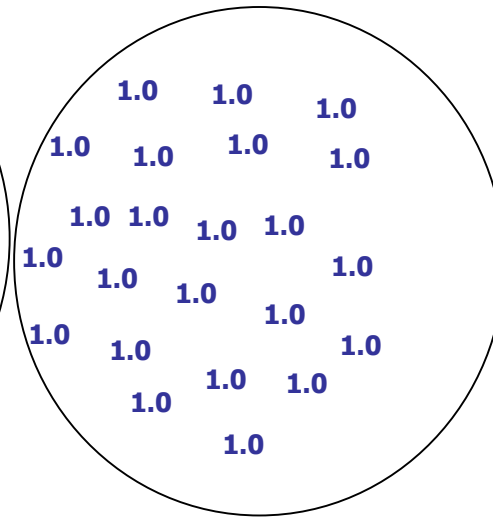
diffexp(A) :- interaction(A,B) & function(B,'GO:0004871')

# Subgroup Discovery

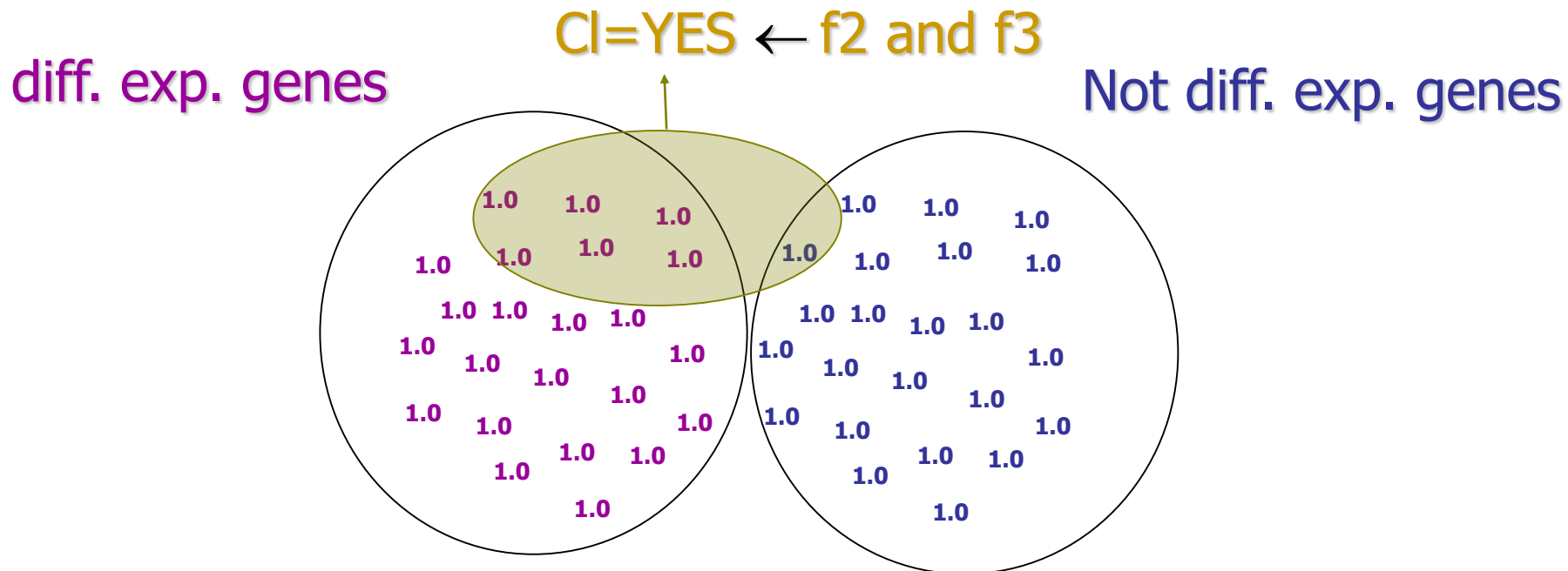
diff. exp. genes



Not diff. exp. genes



# Subgroup Discovery



In RSD (using propositional learner CN2-SD):

Quality of the rules = Coverage x Precision

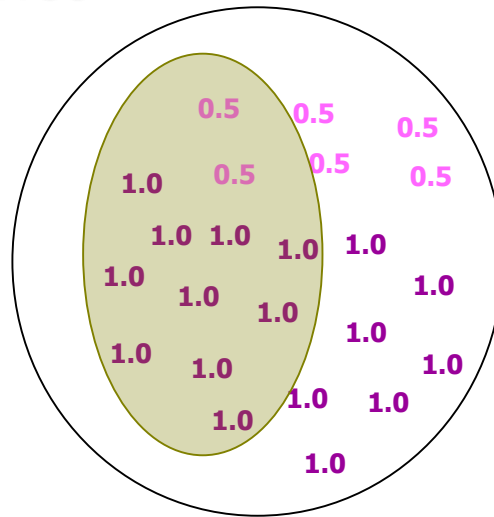
\*Coverage = sum of the covered weights

\*Precision = purity of the covered genes

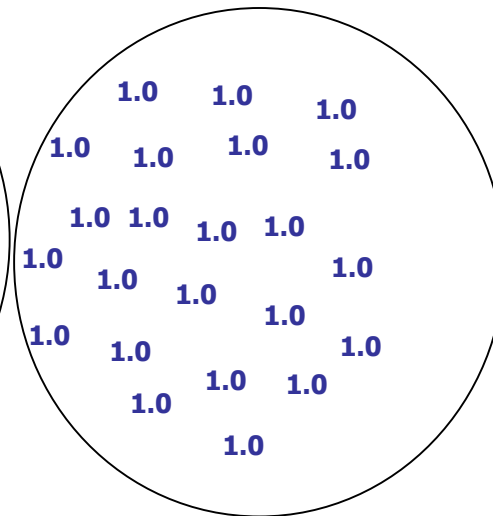


# Subgroup Discovery

diff. exp. genes



Not diff. exp. genes



RSD naturally uses gene weights in its procedure for repetitive subgroup generation, via its heuristic rule evaluation: weighted relative accuracy

# RSD Lessons learned

Efficient propositionalization can be applied to individual-centered, multi-instance learning problems:

- one free global variable (denoting an individual, e.g. molecule M)
- one or more structural predicates: (e.g. `has_atom(M,A)`), each introducing a new existential local variable (e.g. atom A), using either the global variable (M) or a local variable introduced by other structural predicates (A)
- one or more utility predicates defining properties of individuals or their parts, assigning values to variables

`feature121(M):- hasAtom(M,A), atomType(A,21)`

`feature235(M):- lumo(M,Lu), lessThr(Lu,-1.21)`

`mutagenic(M):- feature121(M), feature235(M)`

# Relational Data Mining in Orange4WS

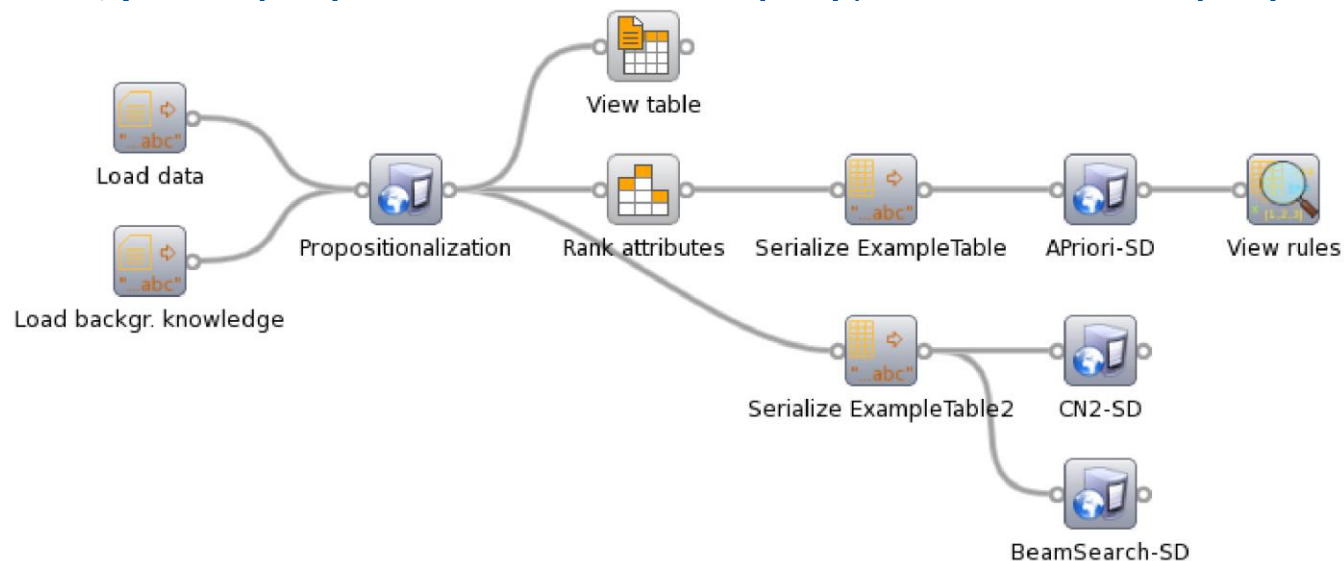
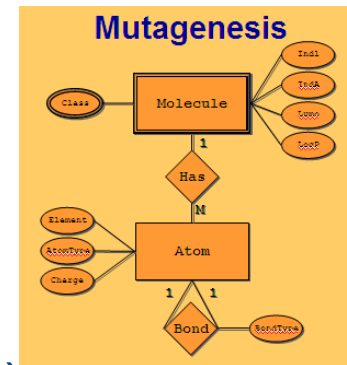
- service for propositionalization through efficient first-order feature construction (Železny and Lavrač, MLJ 2006)

$f_{121}(M):- \text{hasAtom}(M,A), \text{atomType}(A,21)$

$f_{235}(M):- \text{lumo}(M,Lu), \text{lessThr}(Lu,1.21)$

- subgroup discovery using CN2-SD

$\text{mutagenic}(M) \leftarrow \text{feature}_{121}(M), \text{feature}_{235}(M)$

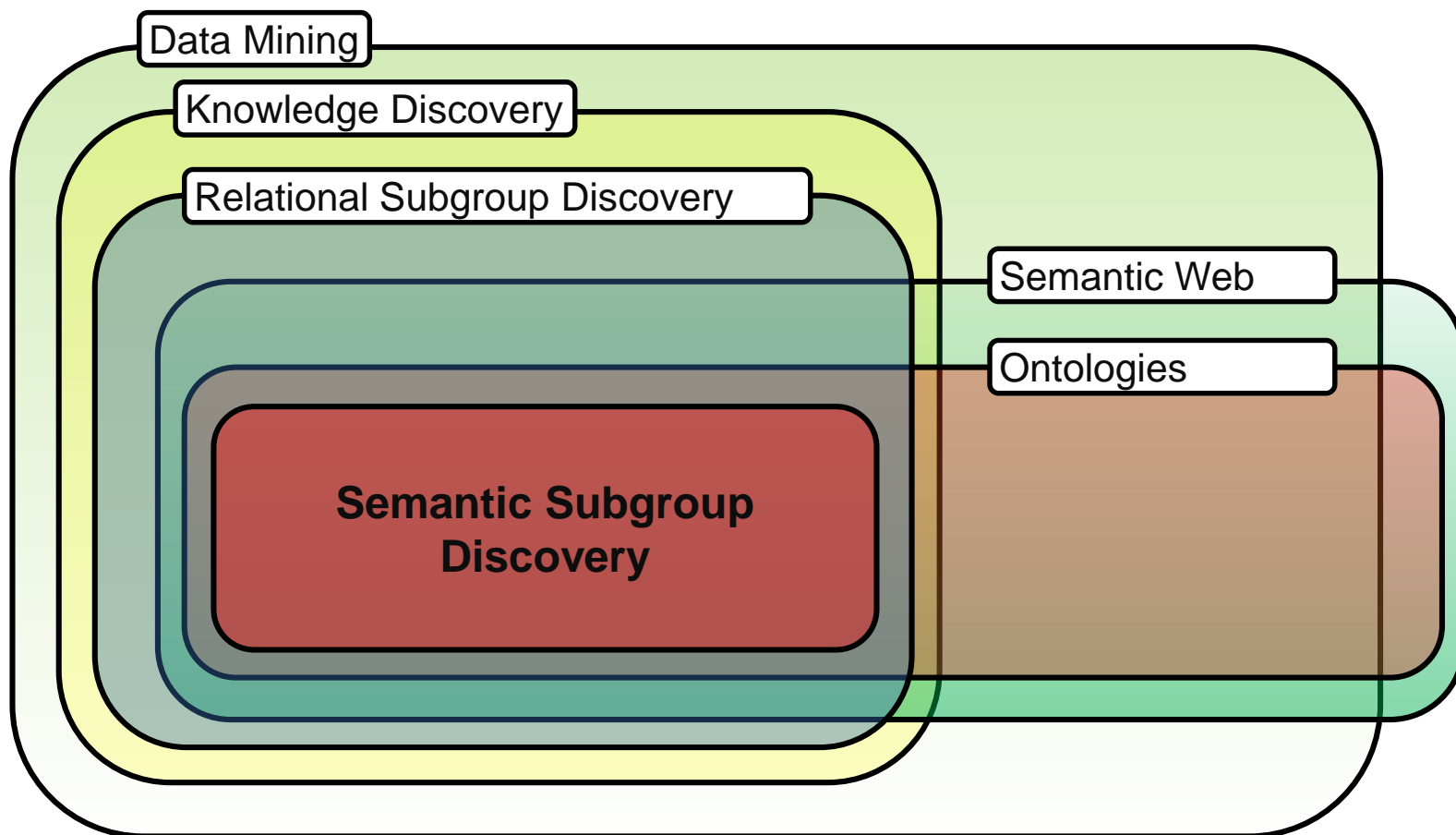


# Semantic Data Mining in Orange4WS

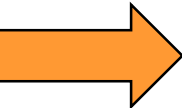
- A special purpose Semantic Data Mining algorithm SEGS
  - discovers interesting gene group descriptions as conjunctions of ontology concepts from GO, KEGG and Entrez
  - integrates public gene annotation data through relational features
  - SEGS algorithm (Trajkovski, Železny, Lavrač and Tolar, JBI 2008) is available in Orange4WS
- Recent developments:
  - Special purpose SDM algorithms: RSD, SDM-SEGS, SDM-Aleph, Hedwig
  - Implemented in web based DM platform ClowdFlows

# Semantic Data Mining

- Semantic subgroup discovery (Vavpetič et al., 2012)



# Advanced Topics

- 
- Outlier detection
    - Text mining: An introduction
    - Document clustering and outlier detection
    - Wordification approach to relational data mining



# Noise and outliers

- Errors in the data – noise

- Animals of white color



- Exceptions or Outliers

- Herd of sheep





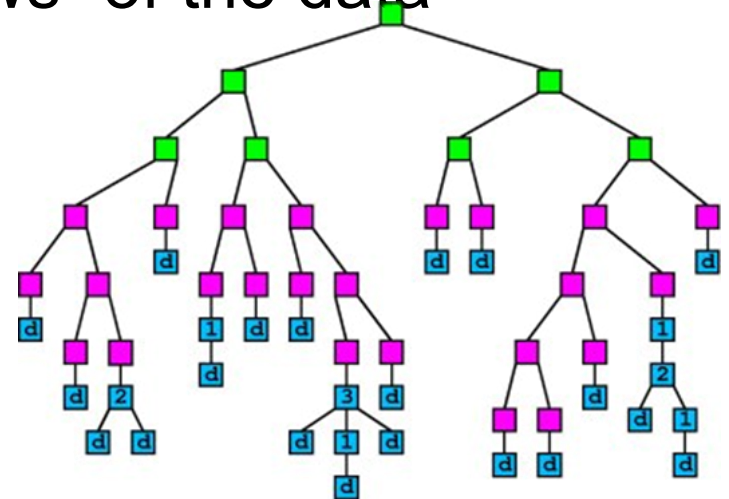
# Noise and outliers

- Data in nature
  - follows certain patterns
  - adheres to the laws of physics
  - is not random



- Build models to Identify the “laws” of the data
  - Patterns and rules =
  - = “laws” of the data

- Errors and outliers
  - Do NOT obey the laws (models)



# Noise and outlier detection

- **Noise** in data negatively affect data mining results. (Zhu et al., 2004)
- False medical diagnosis (**classification noise**) can have serious consequences (Gamberger et al. 2003)
- **Outlier** detection proved to be effective in detection of network intrusion and bank fraud. (Aggarwal and Yu, 2001)

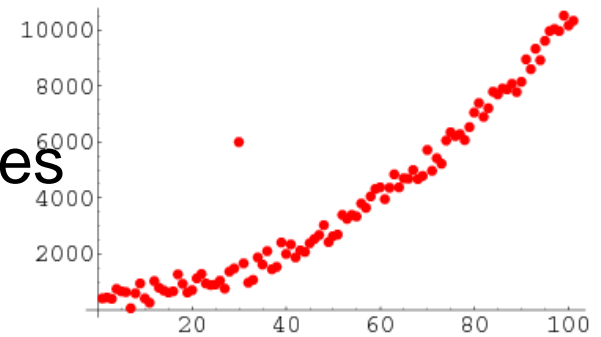
# Detecting noise and outliers

- Errors and exceptions are:

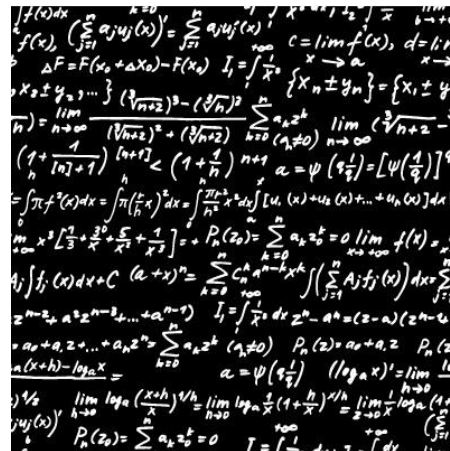
- Inconsistencies with common patterns



- Great deviations from expected values



- Hard to describe

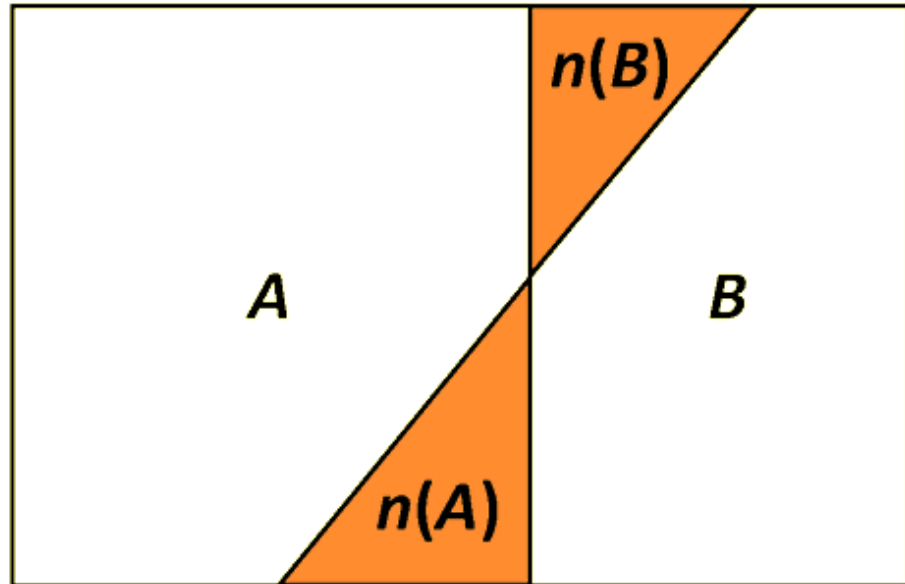
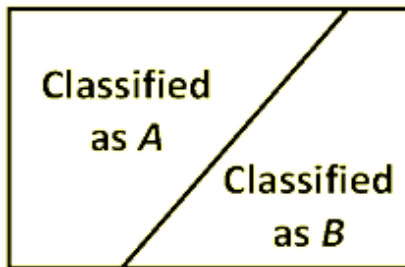
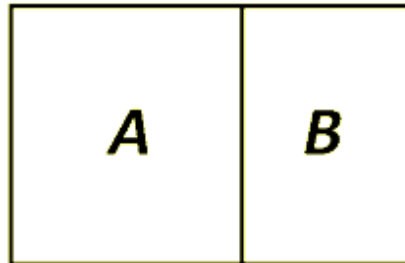


# Classification noise filtering

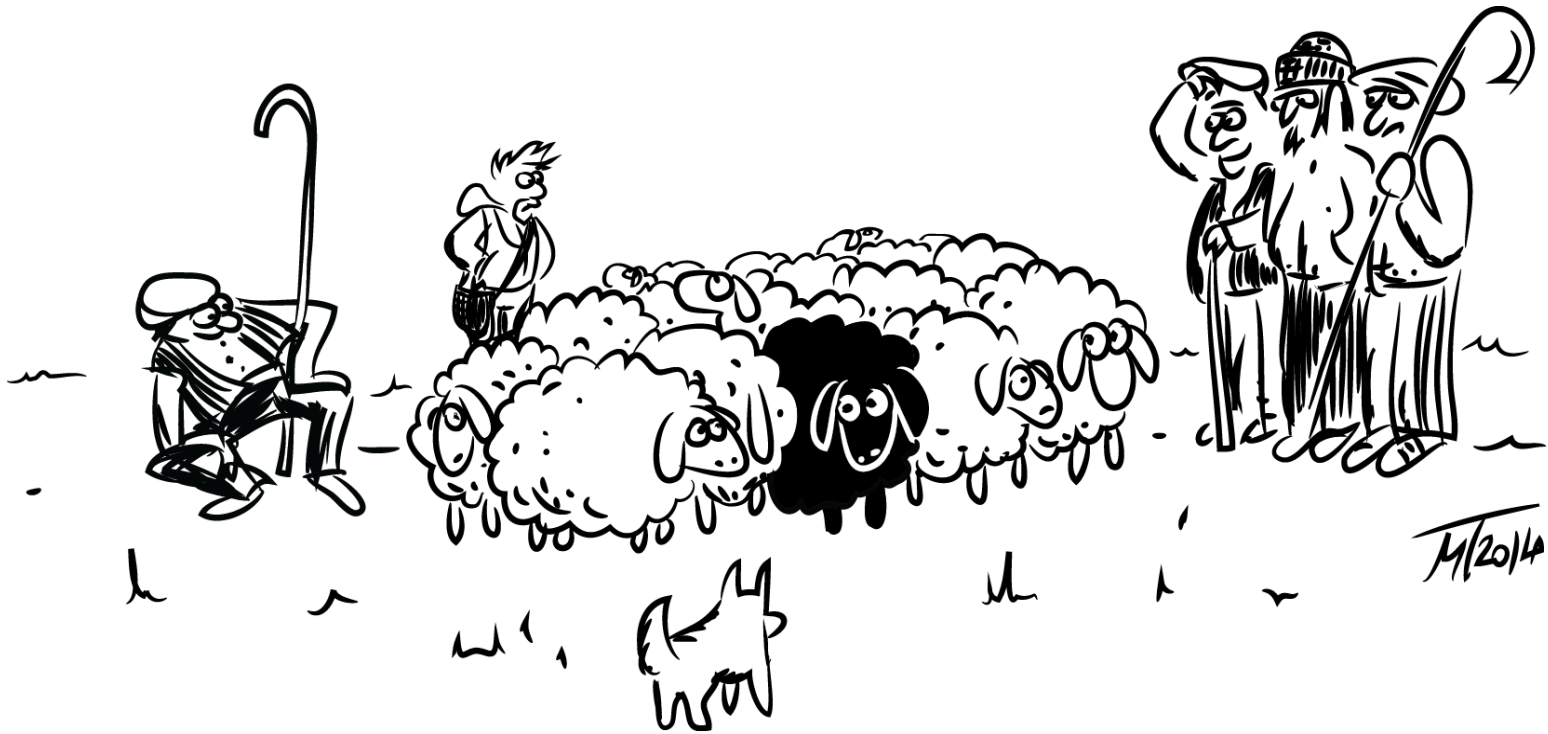
- Model the data
- What can't be modeled is considered noise

# Classification noise filtering

- Model the data, using any learning algorithm
- What can't be modeled is considered noise



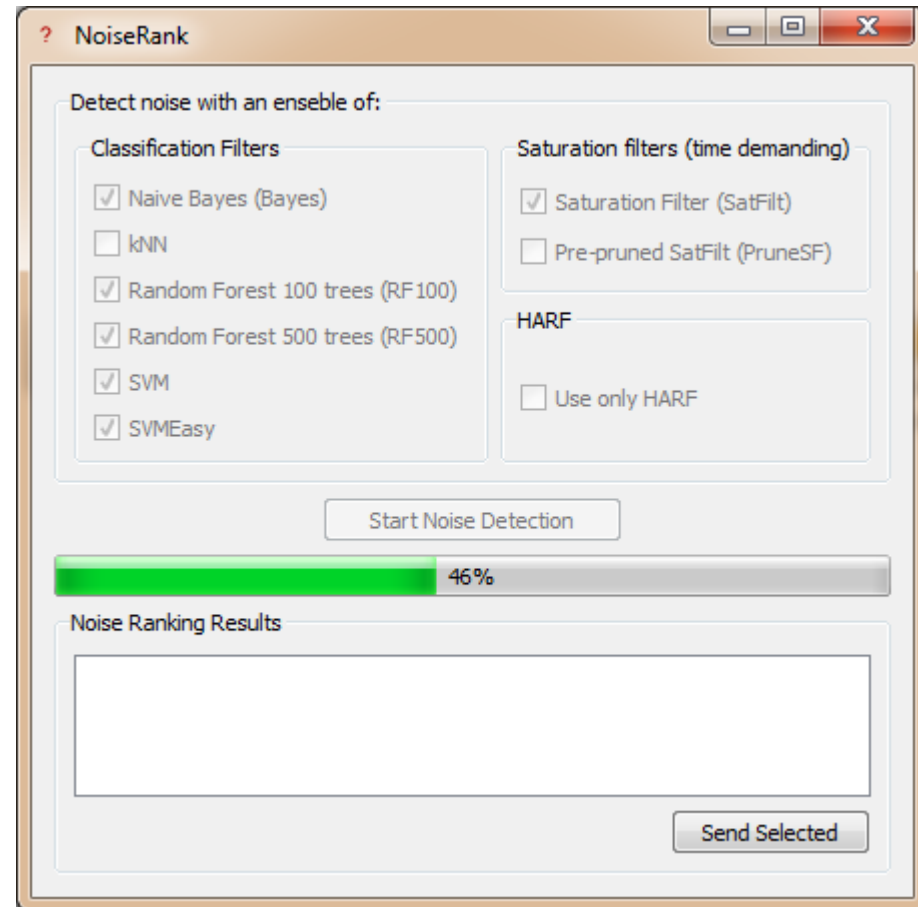
# Ensembles of classifiers



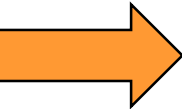
- Combine predictions of various models
- To overcome weaknesses or bias of individual models
- Averaging, Majority voting, Consensus voting, Ranking, etc.

# NoiseRank: Ensemble-based noise and outlier detection

- Misclassified document detection by an ensemble of diverse classifiers (e.g., Naive Bayes, Random Forest, SVM, ... classifiers)
- Ranking of misclassified documents by “voting” of classifiers



# Advanced Topics

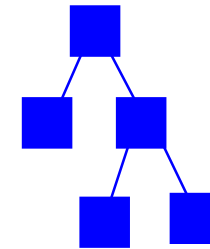
- Outlier detection
-  Text mining: An introduction
- Document clustering and outlier detection
- Wordification approach to relational data mining



# Background: Data mining

Person	Age	Spect. presc.	Astigm.	Tear prod.	Lenses
O1	17	myope	no	reduced	NONE
O2	23	myope	no	normal	SOFT
O3	22	myope	yes	reduced	NONE
O4	27	myope	yes	normal	HARD
O5	19	hypermetrope	no	reduced	NONE
O6-O13	...	...	...	...	...
O14	35	hypermetrope	no	normal	SOFT
O15	43	hypermetrope	yes	reduced	NONE
O16	39	hypermetrope	yes	normal	NONE
O17	54	myope	no	reduced	NONE
O18	62	myope	no	normal	NONE
O19-O23	...	...	...	...	...
O24	56	hypermetrope	yes	normal	NONE

knowledge discovery  
from data



model, patterns, clusters,

...

data

**Given:** transaction data table, a set of text documents, ...

**Find:** a classification model, a set of interesting patterns

# Data mining: Task reformulation

Person	Young	Myope	Astigm.	Reduced tea	Lenses
O1	1	1	0	1	NO
O2	1	1	0	0	YES
O3	1	1	1	1	NO
O4	1	1	1	0	YES
O5	1	0	0	1	NO
O6-O13	...	...	...	...	...
O14	0	0	0	0	YES
O15	0	0	1	1	NO
O16	0	0	1	0	NO
O17	0	1	0	1	NO
O18	0	1	0	0	NO
O19-O23	...	...	...	...	...
O24	0	0	1	0	NO

Binary features and class values

# Text mining: Words/terms as binary features

Document	Word1	Word2	...	WordN	Class
d1	1	1	0	1	NO
d2	1	1	0	0	YES
d3	1	1	1	1	NO
d4	1	1	1	0	YES
d5	1	0	0	1	NO
d6-d13	...	...	...	...	...
d14	0	0	0	0	YES
d15	0	0	1	1	NO
d16	0	0	1	0	NO
d17	0	1	0	1	NO
d18	0	1	0	0	NO
d19-d23	...	...	...	...	...
d24	0	0	1	0	NO

Instances = documents

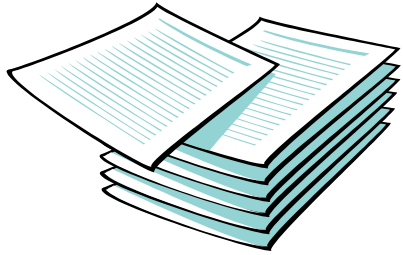
Words and terms = Binary features

# Text Mining from unlabeled data

Document	Word1	Word2	...	WordN	Class
d1	1	1	0	1	NO
d2	1	1	0	0	YES
d3	1	1	1	1	NO
d4	1	1	1	0	YES
d5	1	0	0	1	NO
d6-d13	...	...	...	...	...
d14	0	0	0	0	YES
d15	0	0	1	1	NO
d16	0	0	1	0	NO
d17	0	1	0	1	NO
d18	0	1	0	0	NO
d19-d23	...	...	...	...	...
d24	0	0	1	0	NO

**Unlabeled data** - clustering: grouping of similar instances  
- association rule learning

# Text mining



Step 1

BoW vector construction

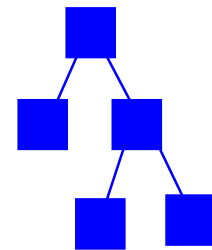
1. BoW features construction
2. Table of BoW vectors construction

Document	Word1	Word2	...	WordN	Class
d1	1	1	0	1	NO
d2	1	1	0	0	YES
d3	1	1	1	1	NO
d4	1	1	1	0	YES
d5	1	0	0	1	NO
d6-d13	...	...	...	...	...
d14	0	0	0	0	YES
d15	0	0	1	1	NO
d16	0	0	1	0	NO
d17	0	1	0	1	NO
d18	0	1	0	0	NO
d19-d23	...	...	...	...	...
d24	0	0	1	0	NO

Document	Word1	Word2	...	WordN	Class
d1	1	1	0	1	NO
d2	1	1	0	0	YES
d3	1	1	1	1	NO
d4	1	1	1	0	YES
d5	1	0	0	1	NO
d6-d13	...	...	...	...	...
d14	0	0	0	0	YES
d15	0	0	1	1	NO
d16	0	0	1	0	NO
d17	0	1	0	1	NO
d18	0	1	0	0	NO
d19-d23	...	...	...	...	...
d24	0	0	1	0	NO

Step 2

Data Mining



model, patterns, clusters,

...

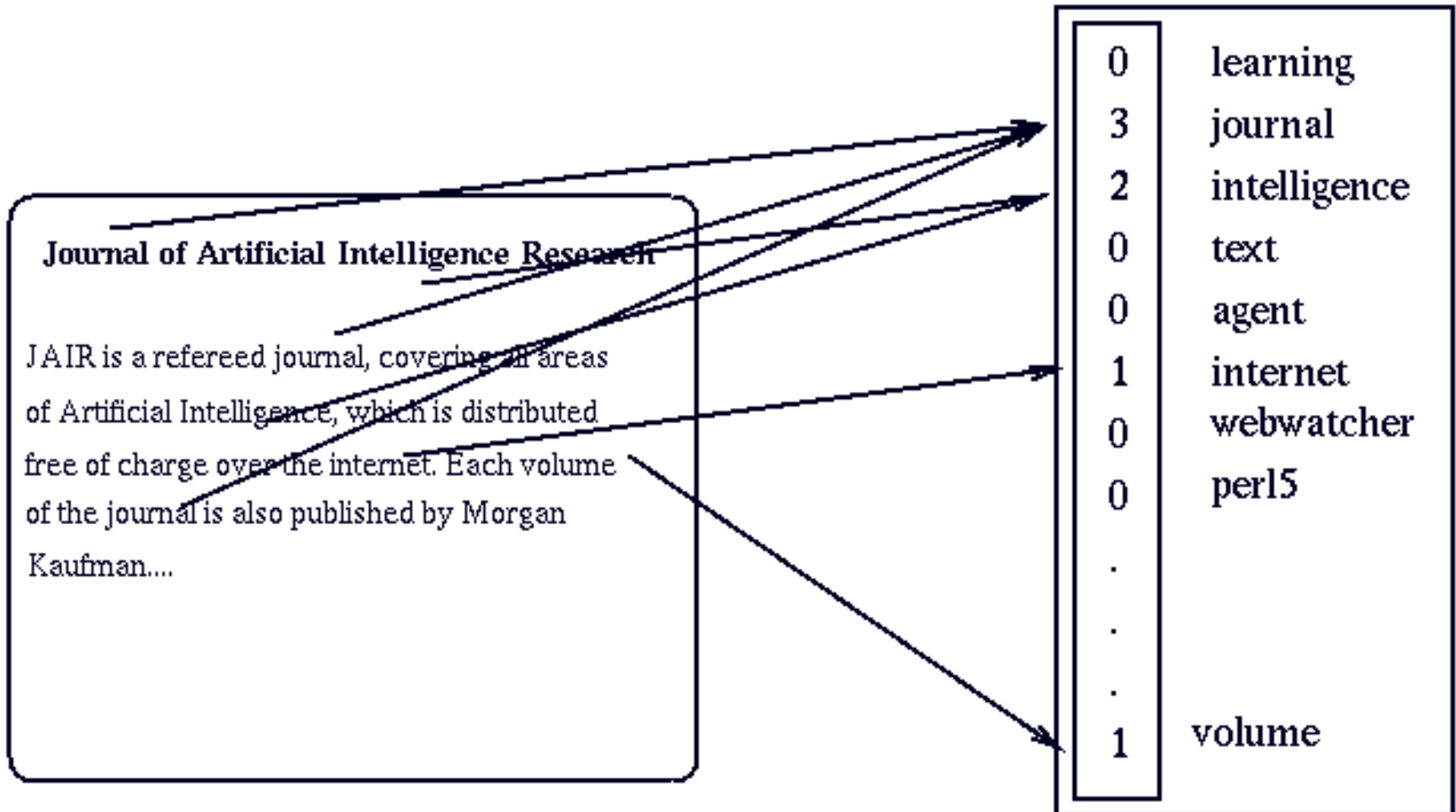
# Text Mining

- Feature construction
  - StopWords elimination
  - Stemming or lemmatization
  - Term construction by frequent N-Grams construction
  - Terms obtained from thesaurus (e.g., WordNet)
- BoW vector construction
- Mining of BoW vector table
  - Feature selection, Document similarity computation
  - Text mining: Categorization, Clustering, Summarization, ...

# Stemming and Lemmatization

- Different forms of the same word usually problematic for text data analysis
  - because they have different spelling and similar meaning (e.g. learns, learned, learning,...)
  - usually treated as completely unrelated words
- Stemming is a process of transforming a word into its stem
  - cutting off a suffix (eg., smejala -> smej)
- Lemmatization is a process of transforming a word into its normalized form
  - replacing the word, most often replacing a suffix (eg., smejala -> smejati)

# Bag-of-Words document representation





# Word weighting

- In bag-of-words representation each word is represented as a separate variable having numeric weight.
- The most popular weighting schema is normalized word frequency TFIDF:

$$tfidf(w) = tf \cdot \log\left(\frac{N}{df(w)}\right)$$

- $Tf(w)$  – term frequency (number of word occurrences in a document)
- $Df(w)$  – document frequency (number of documents containing the word)
- $N$  – number of all documents
- $Tfidf(w)$  – relative importance of the word in the document

The word is more important if it appears several times in a target document

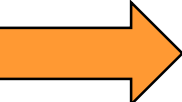
The word is more important if it appears in less documents

# Cosine similarity between document vectors

- Each document  $D$  is represented as a vector of TF-IDF weights
- Similarity between two vectors is estimated by the similarity between their vector representations (cosine of the angle between the two vectors):

$$\text{Similarity} (D_1, D_2) = \frac{\sum_i x_{1i} x_{2i}}{\sqrt{\sum_j x_j^2} \sqrt{\sum_k x_k^2}}$$

# Advanced Topics

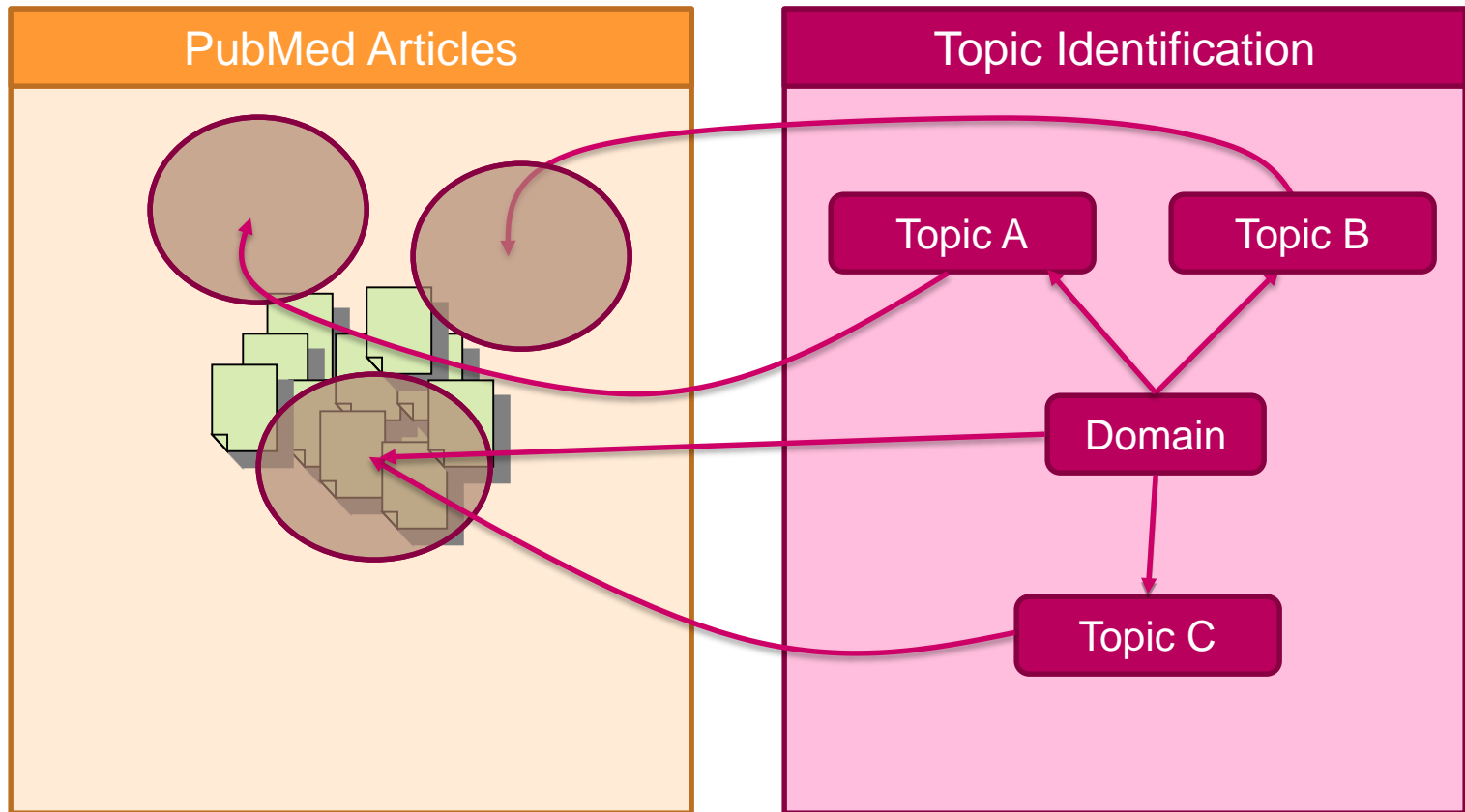
- Outlier detection
- Text mining: An introduction
-  Document clustering and outlier detection
- Wordification approach to relational data mining

# Document clustering

- Clustering is a process of finding natural groups in data in a unsupervised way (no class labels pre-assigned to documents)
- Document similarity is used
- Most popular clustering methods:
  - K-Means clustering
  - Agglomerative hierarchical clustering
  - EM (Gaussian Mixture)
  - ...

# Document clustering with OntoGen

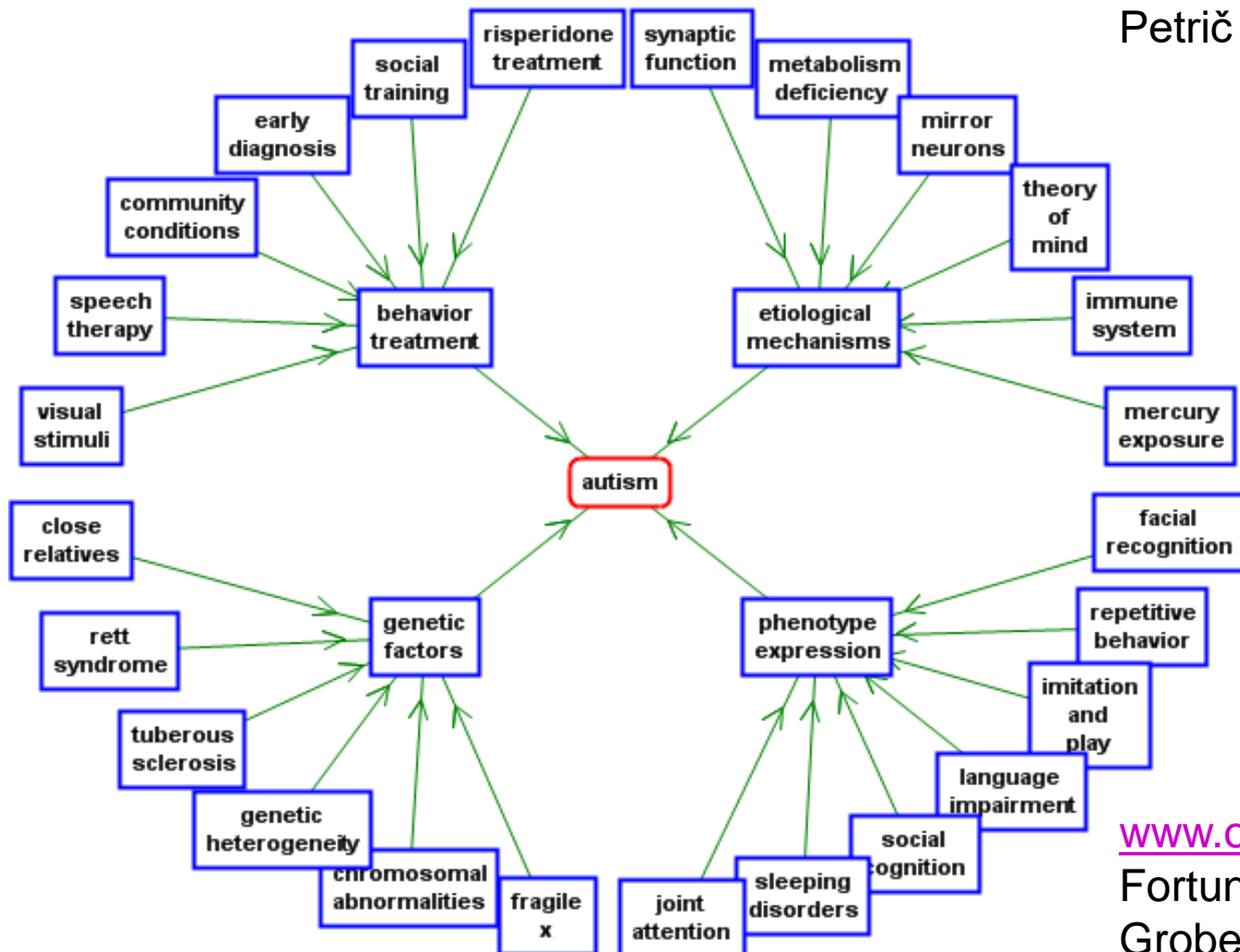
[ontogen.ijs.si](http://ontogen.ijs.si)



Slide adapted from D. Mladenić, JSI

# Using OntoGen for clustering PubMed articles on autism

Work by  
Petrič et al. 2009



[www.ontogen.si](http://www.ontogen.si)  
Fortuna, Mladenić,  
Grobelnik 2006

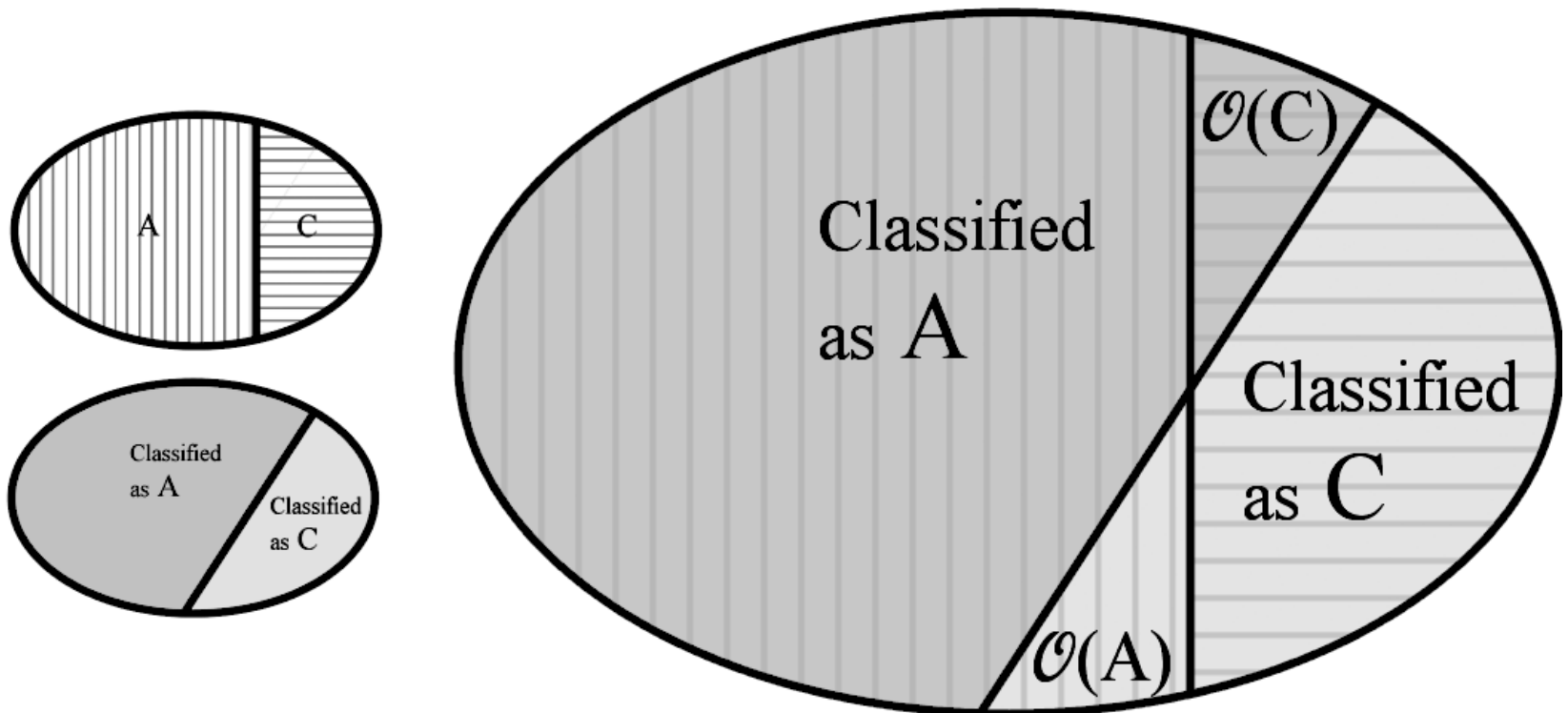
# K-Means clustering in OntoGen

OntoGen uses k-Means clustering for semi-automated topic ontology construction

- Given:
  - set of documents (eg., word-vectors with TFIDF),
  - distance measure (eg., cosine similarity)
  - K - number of groups
- For each group initialize its centroid with a random document
- While not converging
  - each document is assigned to the nearest group (represented by its centroid)
  - for each group calculate new centroid (group mass point, average document in the group)

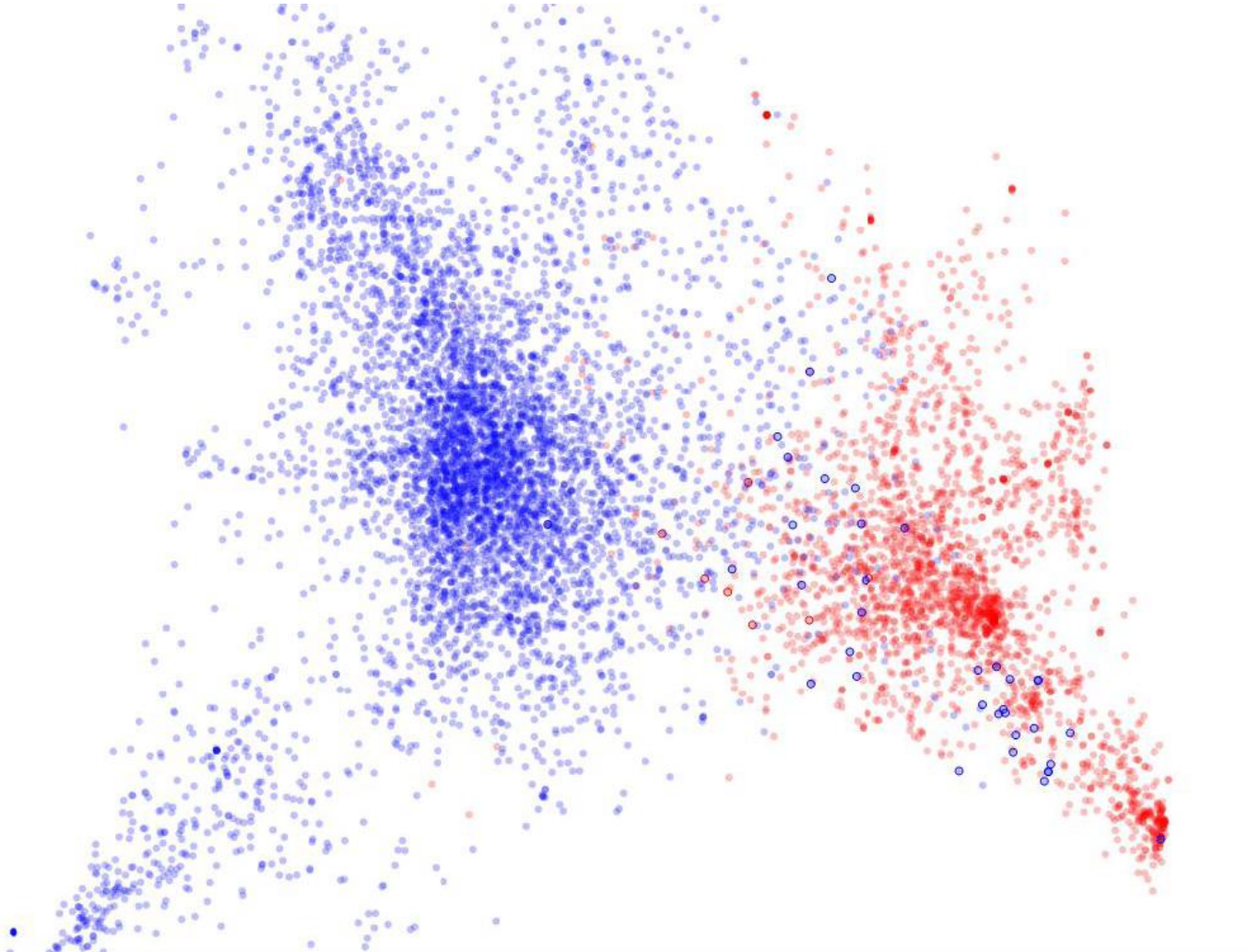
# Detecting outlier documents

- By classification noise detection on a domain pair dataset, assuming two separate document corpora A and C





# Outlier detection for cross-domain knowledge discovery

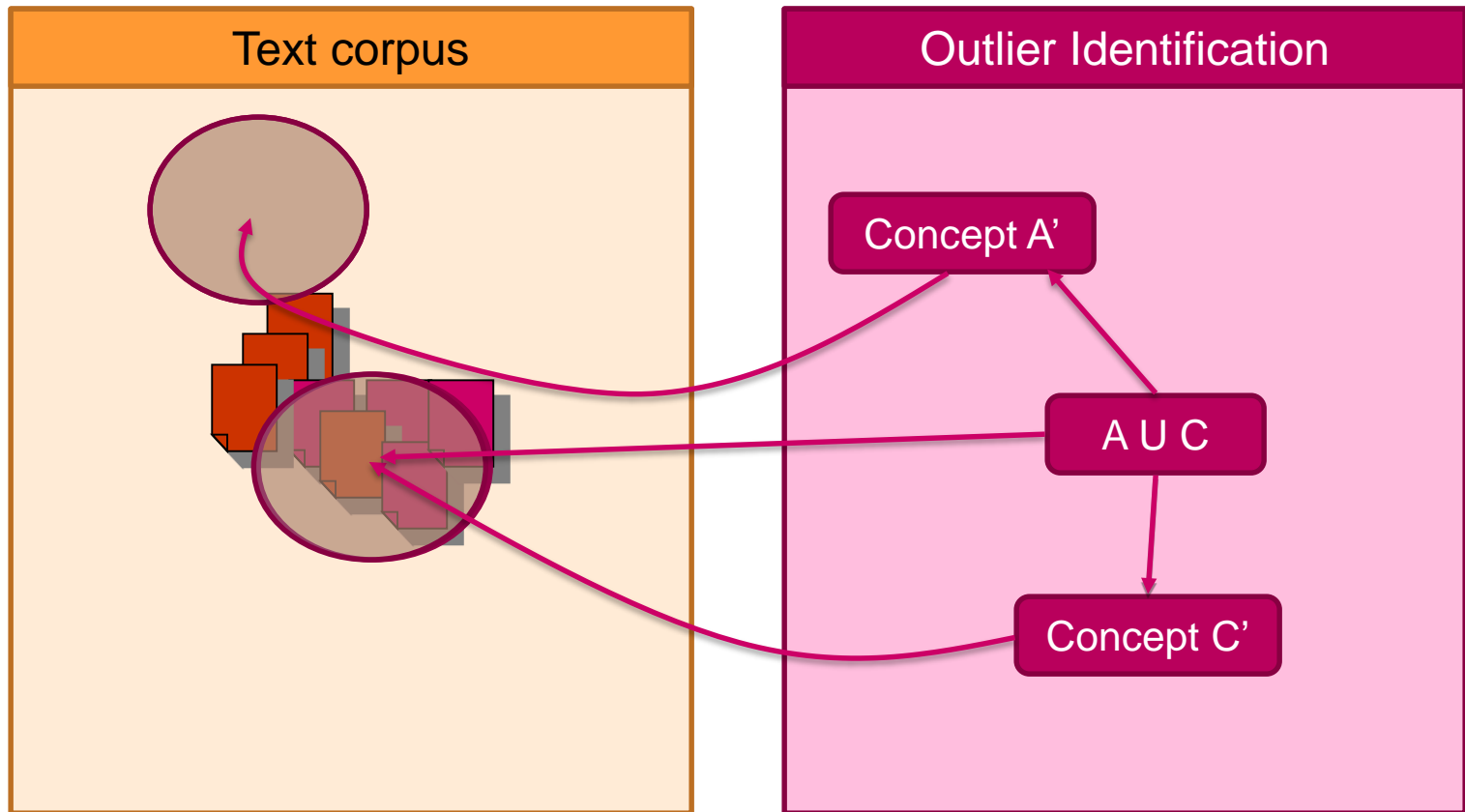


*2-dimensional projection of documents (about autism (red) and calcineurin (blue)). Outlier documents are bolded for the user to easily spot them.*

***Our research has shown that most domain bridging terms appear in outlier documents.***

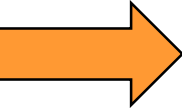
(Lavrač, Sluban, Grčar, Juršič 2010)

# Using OntoGen for outlier document identification



Slide adapted from D. Mladenić, JSI

# Advanced Topics

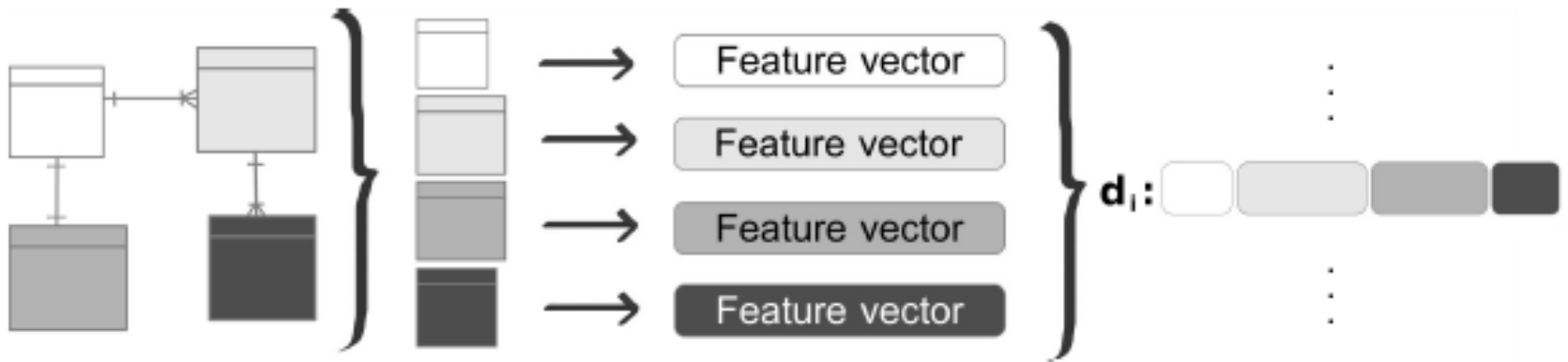
- Outlier detection
  - Text mining: An introduction
  - Document clustering and outlier detection
-  Wordification approach to relational data mining

# Propositionalization through Wordification: Motivation

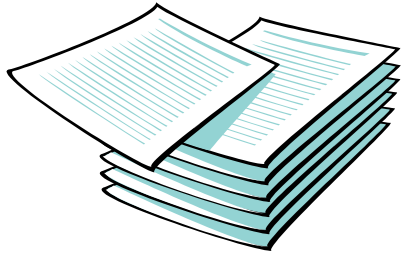
- Develop a RDM technique inspired by **text mining**
- Using a large number of simple, easy to understand features (**words**)
- Improved **scalability**, handling large datasets
- Used as a preprocessing step to propositional learners

# Wordification Methodology

- Transform a relational database to a document corpus
  - For each individual (row) in the main table, concatenate words generated for the main table with words generated for the other tables, linked through external keys



# Text mining



Step 1

BoW vector construction

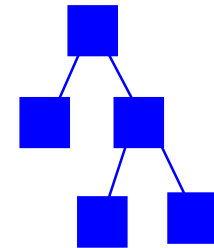
1. BoW features construction
2. Table of BoW vectors construction

Document	Word1	Word2	...	WordN	Class
d1	1	1	0	1	NO
d2	1	1	0	0	YES
d3	1	1	1	1	NO
d4	1	1	1	0	YES
d5	1	0	0	1	NO
d6-d13	...	...	...	...	...
d14	0	0	0	0	YES
d15	0	0	1	1	NO
d16	0	0	1	0	NO
d17	0	1	0	1	NO
d18	0	1	0	0	NO
d19-d23	...	...	...	...	...
d24	0	0	1	0	NO

Document	Word1	Word2	...	WordN	Class
d1	1	1	0	1	NO
d2	1	1	0	0	YES
d3	1	1	1	1	NO
d4	1	1	1	0	YES
d5	1	0	0	1	NO
d6-d13	...	...	...	...	...
d14	0	0	0	0	YES
d15	0	0	1	1	NO
d16	0	0	1	0	NO
d17	0	1	0	1	NO
d18	0	1	0	0	NO
d19-d23	...	...	...	...	...
d24	0	0	1	0	NO

Step 2

Data Mining



model, patterns, clusters,

...

# Wordification Methodology

- One individual of the main data table in the relational database ~ one text document
- Features (attribute values) ~ the words of this document
- Individual words (called **word-items** or **witems**) are constructed as combinations of:

*[table name]\_[attribute name]\_[value]*

- **n-grams** are constructed to model feature dependencies:

*[witem<sub>1</sub>]\_[witem<sub>2</sub>]\_ ... \_[witem<sub>n</sub>]*

# Wordification Methodology

- Transform a relational database to a document corpus
- Construct BoW vectors with TF-IDF weights on words  
(optional: Perform feature selection)
- Apply text mining or propositional learning on BoW table



# Wordification

TRAIN		CAR				
trainID	eastbound	carID	shape	roof	wheels	train
t1	east	c11	rectangle	none	2	t1
...	...	c12	rectangle	peaked	3	t1
...	...	...	...	...	...	...
t5	west	c51	rectangle	none	2	t5
...	...	c52	hexagon	flat	2	t5
...	...	...	...	...	...	...

**t1:** [car\_roof\_none, car\_shape\_rectangle, car\_wheels\_2, car\_roof\_none\_\_car\_shape\_rectangle, car\_roof\_none\_\_car\_wheels\_2, car\_shape\_rectangle\_\_car\_wheels\_2, car\_roof\_peaked, car\_shape\_rectangle, car\_wheels\_3, car\_roof\_peaked\_\_car\_shape\_rectangle, car\_roof\_peaked\_\_car\_wheels\_3, car\_shape\_rectangle\_\_car\_wheels\_3], **east**



# TF-IDF weights

- No explicit use of existential variables in features, TF-IDF instead
- The weight of a word indicates how relevant is the feature for the given individual
- The TF-IDF weights can then be used either for filtering words with low importance or for using them directly by a propositional learner (e.g. J48)

# Experiments

- Cross-validation experiments on 8 relational datasets: Trains (in two variants), Carcinogenesis, Mutagenensis with 42 and 188 examples, IMDB, and Financial.
- Results (using J48 for propositional learning)

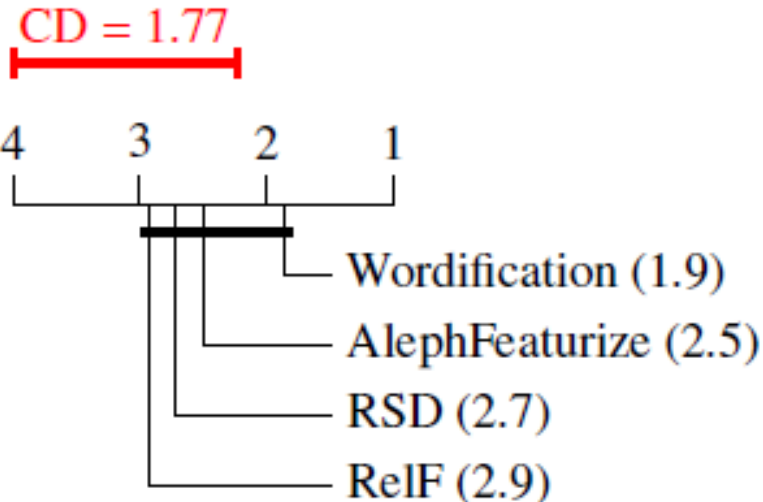
# Experiments

- Cross-validation experiments on 8 relational datasets: Trains (in two variants), Carcinogenesis, Mutagenensis with 42 and 188 examples, IMDB, and Financial.
- Results (using J48 for propositional learning)
  - first applying Friedman test to rank the algorithms,
  - then post-hoc test Nemenyi test to compare multiple algorithms to each other

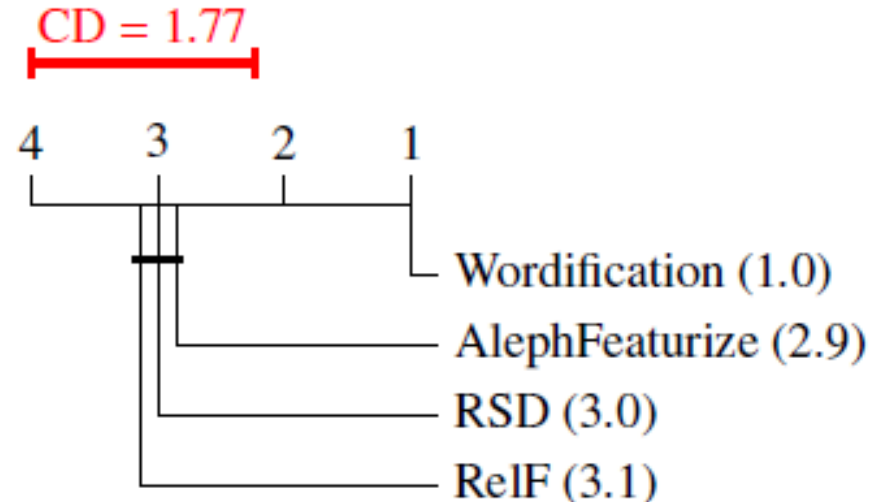
# Experiments

- Cross-validation experiments on 8 relational datasets: Trains (in two variants), Carcinogenesis, Mutagenensis with 42 and 188 examples, IMDB, and Financial.
- Results (using J48 for propositional learning)

MEASURE = CA



MEASURE = RUN-TIME



# Experiments

Domain	Algorithm	J48-Accuracy[%]	J48-AUC	Run-time[s]
Trains without position	Wordification	55.00	0.51	<b>0.11</b>
	RelF	65.00	0.65	1.04
	RSD	65.00	0.68	0.53
	AlephFeaturize	<b>75.00</b>	<b>0.82</b>	0.40
Trains	Wordification	<b>95.00</b>	<b>0.91</b>	<b>0.12</b>
	RelF	65.00	0.62	1.06
	RSD	50.00	0.53	0.47
	AlephFeaturize	85.00	0.74	0.38
Mutagenesis42	Wordification	<b>97.62</b>	<b>0.93</b>	<b>0.39</b>
	RelF	80.95	0.59	2.11
	RSD	<b>97.62</b>	<b>0.93</b>	2.63
	AlephFeaturize	<b>97.62</b>	<b>0.93</b>	2.07
Mutagenesis188	Wordification	<b>95.74</b>	0.90	<b>1.65</b>
	RelF	75.53	0.79	7.76
	RSD	94.15	<b>0.91</b>	10.10
	AlephFeaturize	87.23	0.88	19.27
IMDB	Wordification	<b>84.34</b>	<b>0.79</b>	<b>1.23</b>
	RelF	79.52	0.73	32.49
	RSD	73.49	0.47	4.33
	AlephFeaturize	73.49	0.47	4.96
Carcinogenesis	Wordification	<b>61.09</b>	<b>0.62</b>	<b>1.79</b>
	RelF	54.71	0.53	16.44
	RSD	58.05	0.56	9.29
	AlephFeaturize	55.32	0.49	104.70
Financial	Wordification	86.75	0.48	<b>4.65</b>
	RelF	<b>97.00</b>	<b>0.91</b>	260.93
	RSD	86.75	0.48	533.68
	AlephFeaturize	86.75	0.48	525.86

# Use Case: IMDB

- **IMDB subset:** Top 250 and bottom 100 movies
- Movies, actors, movie genres, directors, director genres
- Wordification methodology applied
- Association rules learned on BoW vector table



# Use Case: IMDB

goodMovie ← director\_genre\_drama, movie\_genre\_thriller,  
director\_name\_AlfredHitchcock. (Support: 5.38% Confidence: 100.00%)

movie\_genre\_drama ← goodMovie, actor\_name\_RobertDeNiro.  
(Support: 3.59% Confidence: 100.00%)

director\_name\_AlfredHitchcock ← actor\_name\_AlfredHitchcock.  
(Support: 4.79% Confidence: 100.00%)

director\_name\_StevenSpielberg ← goodMovie, movie\_genre\_adventure,  
actor\_name\_TedGrossman.  
(Support: 1.79% Confidence: 100.00%)

# Summary

- Wordification methodology
- Allows for solving non-standard RDM tasks, including RDM clustering, **word cloud visualization**, **association rule learning**, topic ontology construction, outlier detection, ...

