Relational and Semantic Data Mining

Nada Lavrač Anže Vavpetič

Jožef Stefan Institute Ljubljana, Slovenia

LPNMR-2015, Lexington, September 2015



Jožef Stefan Institute, Ljubljana, Slovenia

Jožef Stefan Institute

 leading national research organization in natural sciences and technology (~900 staff, researchers and students, ~500 PhDs)





Jožef Stefan Institute, Ljubljana, Slovenia

Jožef Stefan Institute

 leading national research organization in natural sciences and technology (~900 staff, researchers and students, ~500 PhDs)



Department of Knowledge Technologies

~ 40 researchers (~ 10 PhD students) + NAO robot





Talk outline

- Intro. to Machine Learning and Data Mining
- Inductive Logic Programming (ILP) and Relational Data Mining (RDM)
 - Propositionalization approach to RDM
- Semantic data mining
- Summary and future work



Machine Learning and Data Mining

data

	1000				
Person	Age	Spect. presc.	Astigm.	Tear prod.	Lenses
01	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
06-013					
O14	35	hypermetrope	no	normal	SOFT
O15	43	hypermetrope	yes	reduced	NONE
O16	39	hypermetrope	yes	normal	NONE
017	54	myope	no	reduced	NONE
O18	62	myope	no	normal	NONE
019-023					
024	56	hypermetrope	yes	normal	NONE

knowledge discovery from data Machine Learning

Data Mining

model, patterns, ...

Given: class labeled data Find: classification model or set of interesting patterns in the data



Example: Contact lens data

DATA

			(and the		
Person	Age	Spect. presc.	Astigm.	Tear prod.	Lenses
01	17	myope	no	reduced	NONE
02	23	myope	no	normal	SOFT
03	22	myope	yes	reduced	NONE
O4	27	myope	yes	normal	HARD
O5	19	hypermetrope	no	reduced	NONE
06-013		/			
014	35	hypermetrope	no	normal	SOFT
O15	43	hypermetrope	yes	reduced	NONE
O16	39	hypermetrope	yes	normal	NONE
017	54	myope	no	reduced	NONE
018	62	myope	no	normal	NONE
019-023				··· ··	
024	56	hypermetrope	yes	normal	NONE



Learning models from contact lens data

Person	Age	Spect, presc.	Astiam.	Tear prod.	Lenses	
01	17	myope	no	reduced	NONE	and being and been and been and
O2	23	myope	no	normal	SOFT	
O3	22	myope	yes	reduced	NONE	
O4	27	myope	yes	normal	HARD	Data Mining
O5	19	hypermetrope	no	reduced	NONE	Data Mining
06-013						
O14	35	hypermetrope	no	normal	SOFT	
O15	43	hypermetrope	yes	reduced	NONE	MEDIA
O16	39	hypermetrope	yes	normal	NONE	reduced pormal
017	54	myope	no	reduced	NONE	
O18	62	myope	no	normal	NONE	NONE
019-023						no yes
O24	56	hypermetrope	yes	normal	NONE	SOFT spect. pre.

HARD

NONE



Learning models from contact lens data

Person	Age	Spect. presc.	Astigm.	Tear prod.	Lenses	and a second sec
01	17	myope	no	reduced	NONE	and the second se
02	23	myope	no	normal	SOFT	
O3	22	myope	yes	reduced	NONE	
04	27	myope	yes	normal	HARD	Data Mining
O5	19	hypermetrope	no	reduced	NONE	Data winning
06-013						
O14	35	hypermetrope	no	normal	SOFT	
O15	43	hypermetrope	yes	reduced	NONE	tear pred
O16	39	hypermetrope	yes	normal	NONE	reduced normal
017	54	myope	no	reduced	NONE	
O18	62	myope	no	normal	NONE	NONE
019-023						no yes
O24	56	hypermetrope	yes	normal	NONE	SOFT Spect. pre.
line and	-				(m)	myope hypermet
	lanca		_ toor	nroducti	ion-red	

- lenses=NONE ← tear production=normal ^ astigmatism=yes ^ spect. pre.=hypermetrope
- lenses=SOFT ← tear production=normal ^ astigmatism=no
- lenses=HARD ← tear production=normal ^ astigmatism=yes ^

spect. pre.=myope

Finding patterns in Contact lens data

DATA

· · · · ·						
Person	Age	Spect. presc.	Astigm.	Tear prod.	Lenses	Section and Section 1
01	17	myope	no	reduced	NONE	
02	23	myope	no	normal	SOFT	
O3	22	myope	yes	reduced	NONE	Data Mining
O4	27	myope	yes	normal	HARD	Data Mining
O5	19	hypermetrope	no	reduced	NONE	
06-013		/				
014	35	hypermetrope	no	normal	SOFT	
O15	43	hypermetrope	yes	reduced	NONE	
O16	39	hypermetrope	yes	normal	NONE	
017	54	myope	no	reduced	NONE	
O18	62	myope	no	normal	NONE	
019-023						
O24	56	hypermetrope	yes	normal	NONE	

PATTERNs: rules describing subgroups of instances



Lenses = NONE ← Tear prod. = reduced

Why learn symbolic models and patterns

Use learned models to classify and explain classifications of new instances





Task reformulation: Binary class values

Person	Age	Spect. presc.	Astigm.	Tear prod.	Lenses
01	17	myope	no	reduced	NO
02	23	myope	no	normal	YES
O3	22	myope	yes	reduced	NO
O4	27	myope	yes	normal	YES
O5	19	hypermetrope	no	reduced	NO
06-013					
014	35	hypermetrope	no	normal	YES
O15	43	hypermetrope	yes	reduced	NO
O16	39	hypermetrope	yes	normal	NO
017	54	myope	no	reduced	NO
O18	62	myope	no	normal	NO
019-023					
O24	56	hypermetrope	yes	normal	NO

Binary classes (positive vs. negative examples of Target class) - simplifies single concept learning

- is used in "one class vs. all" multi-class learning methods



Classification rules learned from contact lens data with binary classes

lenses=YES \leftarrow tear production=normal \land astigmatism=no lenses=YES ← tear production=normal ^ astigmatism=yes ^ spect. pre.=myope lenses=NO ← tear production=reduced lenses=NO ← tear production=normal ^ astigmatism=yes ^ spect. pre.=hypermetrope

lenses=NO ←



Covering algorithm used in rule learning





Covering algorithm used in rule learning





Covering algorithm used in rule learning





Heuristics used in classification rule learning

Rule: Class ← Conditions



Precision (rel. freq. of correctly covered pos. examples)

Negative examples



Positive examples



e.g., trading off coverage and precision gain: Coverage x (Precision – Default) $q(Rule) = p(Conds) \times (p(Class|Conds) - p(Class))$



Coverage (rel. freq. of covered examples)Precision (rel. freq. of correctly covered pos. examples)

Subgroup discovery example: High CHD Risk Group Detection

 Input: Patient records described by anamnestic, laboratory and ECG attributes
 Task: Find and characterize population subgroups with high CHD risk (large, distributionally unusual subgroups)

From best induced descriptions, five were selected by the expert as most actionable for CHD risk screening (by GPs): high-CHD-risk ← male ^ pos. fam. history ^ age > 46 high-CHD-risk ← female ^ bodymassIndex > 25 ^ age > 63 high-CHD-risk ← ... high-CHD-risk ← ... high-CHD-risk ← ...



(Gamberger & Lavrač, JAIR 2002)

SD algorithms in the Orange DM Platform

Orange data mining toolkit

- classification and subgroup discovery algorithms
- data mining workflows
- visualization
- developed at FRI, Ljubljana



SD Algorithms in Orange

- SD (Gamberger & Lavrač, JAIR 2002
- APRIORI-SD (Kavšek & Lavrač, AAI 2006
- CN2-SD (Lavrač et al., JMLR 2004): Adapting CN2

classification rule learner to Subgroup Discovery

Other Data Mining Platforms

WEKA, KNIME, RapidMiner, Orange4WS, ...



– include numerous data mining algorithms
– enable data and model visualization
– enable complex workflow construction

Talk outline

- Intro. to Machine Learning and Data Mining
 Inductive Logic Programming (ILP) and Relational Data Mining (RDM)
 - Propositionalization approach to RDM
- Semantic data mining
- Summary and future work



ILP and Relational Data Mining



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 patterns

ILP and Relational Data Mining





ILP and Relational Data Mining

ILP, relational learning, relational data mining

- Learning from complex relational databases
- Learning from complex structured data, e.g. molecules and their biochemical properties





Relational representation of customers, orders and stores.



ILP for Logic Programming

Given:

- A set of observations (ground facts)
 - positive examples E⁺
 - negative examples E⁻
- background knowledge B (definite clauses)
- hypothesis language (definite clauses) L_H
- covers relation (logical entailment)

Find:

A hypothesis (a theory) $H \in L_H$, such that (given *B*) *H* covers all positive and no negative examples

- In logic, find H such that
 - ∀e ∈ E^+ : B ∪ H |= e (*H* is complete)
 - $\forall e \in E^-$: $B \cup H \neq e$ (*H* is consistent)





Inductive Logic Programming Example

E⁺ = {sort([2,1,3],[1,2,3])}
E⁻ = {sort([2,1],[1]),sort([3,1,2],[2,1,3])}

B: definitions of permutation/2 and sorted/1

Predictive ILP: Learning a theory H

 $sort(X,Y) \leftarrow permutation(X,Y), sorted(Y).$

Descriptive ILP: Finding individual patterns

sorted(Y) \leftarrow sort(X,Y). permutation(X,Y) \leftarrow sort(X,Y). sorted(X) \leftarrow sort(X,X).



ILP for relational learning

Given:

- A set of observations (ground facts)
 - positive examples E⁺
 - negative examples E⁻
- background knowledge B (definite clauses)
- hypothesis language (definite clauses) L_H
- covers relation (theta-subsumption)
- quality criterion, e.g., predictive accuracy A(H)

Find:

A hypothesis (a set of clauses) $H \in L_H$, such that (given *B*) *H* is optimal w.r.t. given quality criterion

(relaxing the request for finding a hypothesis $H \in L_H$, such that (given B) H covers all positive and no negative examples)





Relational Learning Example

E * = {daughter(mary, ann), daughter(eve, tom) }
E * = {daughter(tom, ann), daughter(eve, ann) }

B = Facts: {mother(ann, mary), mother(ann, tom), father(tom, eve), father(tom, ian), female(ann), female(mary), female(eve), male(pat), male(tom) } **Rules:** {parent(X,Y) \leftarrow mother(X,Y), parent(X,Y) \leftarrow father(X,Y) } ann tom mar

ian

eve



Relational Learning Example

- Finding patterns: Induce individual rules (individual general clauses) ← daughter(X,Y), mother(X,Y). female(X) ← daughter(X,Y). mother(X,Y); father(X,Y) ← parent(X,Y).



ILP as search of program clauses

- ILP systems structure the hypothesis space based on syntactic generality relation (θ-subsumption)
 - Clause $c_1 \theta$ -subsumes $c_2 (c_1 \ge \theta c_2)$ iff $\exists \theta : c_1 \theta \subseteq c_2$
 - Hypothesis $H_1 \ge \theta H_2$ iff $\forall c_2 \in H_2$ exists $c_1 \in H_1$ such that $c_1 \ge \theta c_2$

Example

- $c1 = daughter(X,Y) \leftarrow parent(Y,X)$
- $c2 = daughter(mary,ann) \leftarrow female(mary), parent(ann,mary),$
- c1 θ -subsumes c_2 under θ = {X/mary,Y/ann}

Learning strategies

- Top-down search of refinement graphs (FOIL)
- Bottom-up search (building least general generalizations, inverting resolution (CIGOL), inverting entailment (PROGOL))
- Mixed strategy (Aleph)



Generality ordering of clauses

Training examples		Background knowledge	
daughter(mary,ann).	\oplus	parent(ann,mary).	female(ann).
daughter(eve,tom).	\oplus	parent(ann,tom).	female(mary).
daughter(tom,ann).	θ	parent(tom,eve).	female(eve).
daughter(eve,ann).	θ	parent(tom,ian).	

daughter(X,Y) \leftarrow daughter(X,Y) \leftarrow daughter(X,Y) \leftarrow daughter(X,Y) \leftarrow X=Y parent(Y,X) parent(X,Z) daughter(X,Y) \leftarrow female(X) Part of the refinement daughter(X,Y) \leftarrow daughter(X,Y) \leftarrow graph for the family KNOJemale (X) female(X) relations problem. female(Y) parent(Y.X)

Top-down search of refinement graphs

Training examples		Background knowledge	
daughter(mary,ann).	\oplus	parent(ann,mary).	female(ann.).
daughter(eve,tom).	\oplus	parent(ann,tom).	female(mary).
daughter(tom,ann).	θ	parent(tom,eve).	female(eve).
daughter(eve,ann).	θ	parent(tom,ian).	



Selected ILP algorithms are available online in the ClowdFlows platform

 Example: ILP system Aleph in ClowdFlows available at http://clowdflows.org/workflow/2224/





Talk outline

- Intro. to Machine Learning and Data Mining
- Inductive Logic Programming (ILP) and Relational Data Mining (RDM)
- Propositionalization approach to RDM
- Semantic data mining
- Summary and future work



Relational Data Mining through Propositionalization




Relational Data Mining through Propositionalization

customer ID Zip So In ID Zip ex St come 3478 34677 msi 60-70 3479 43666 f ma. 80-90 Order Store Delivery Pr Mode M	A ge 32 45	Cl ub 	Re sp										
ID Zip Sx So In Zip Sx St Scome order Obstomer Order Store Delivery Presson	A ge 32 45	Cl ub 	Re sp										
<td> 32 45</td> <td> me</td> <td></td>	 32 45	 me											
3478 34677 m si 60-70 3479 43666 f ma 80-90 order Order Store Delivery P D D D Mode M	32 45	lme.											
3479 43666 f mal 80-90 order Order B B B Obstomer Order Store Delivery Pr D D Mode M	45	v	nr										
order Oustomer Order D D Mode		nm	re										
order Customer Order Store Delivery Pr DD Mode M													
Customer Order Store Delivery Particular Mode M	order												
ID ID ID Mode M	der Store Delivery Paymt												
i	lőd	le											
\		\neg											
3478 2140267 12 regular ca	sh												
3478 3446778 12 \express ch	lec]	k											
3478 4728386 17 regular ch	iecl	k											
3479 3233444 17 express cr	edi	it											
3479 3475886 12 regular cr	edi	it											
\													

. \			
7	5	store	
Store ID	Size	Туре	Location
12	small	franchise	city
17	large	indep	rural

Relational representation of customers, orders and stores.

	f1	f2	f3	f4	f5	f6		1		1		
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	0	
g4	1	1	1	0	1	10 ² 0	0	0	1	1	1	Γ
g5	1	1	1	0	0 /	0010	0	1	1	0	1	Γ
g1	0	٥	1	1	0	0	0	1	0	0	0	
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	0	Γ

Propositionalization

Step 1

									distant form			
	f1	f2	f3	f4	f5	f 6		11		1		\mathbf{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	10 ² 0	0	0	1	1	1	0
g5	1	1	1	0	0 4	0010	0	1	1	0	1	0
g1	0	٥	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 2 Data Mining

model, patterns, ...

Sample ILP problem: East-West trains





Relational representation of East-West trains

							RAIN	IABLE			
LOAD	CAR	OBJECT	NUMBE	R			TRAIN E	EAS TBOUND			
1	c1	circle	1				t 1	TRUE			
12	c2	hexagon	1	-			t2	TRUE			
13	c3	triangle	1								
14	c4	rectangle	3				t 6	FAL SE			
						_					
						1000					
				/			_				
		<u>c</u>	CAR	TRAIN	SHAPE	LENGTH	ROOF	WHEELS			
			c1	t1	rectangle	short	none	2			
			c2	t1	rectangle	long	none	3			
			c3	t 1	rectangle	short	neaked	2			
			00 04	14	rectangle	long	nono	2			
		_	64	L I	rectangle	long	none	2			
		_									
		e	Directio		Train						
		4		/							
					1						
					1						
					Has						
					nas		_				

1

Car

1

Has

Load



Length

Roof

Wheels

Propositionalization approach

Propositionalization through first-order feature construction (1BC, RSD, ...): f1(T):-hasCar(T,C),clength(C,short). f2(T):-hasCar(T,C), hasLoad(C,L), loadShape(L,circle)







f3(T) :-

Propositionalization approach

Standard propositionalization through first-order feature construction (1BC, RSD, ...): f1(T):-hasCar(T,C),clength(C,short). f2(T):-hasCar(T,C), hasLoad(C,L), loadShape(L,circle) f3(T):-....

Propositional learning: $t(T) \leftarrow f1(T), f4(T)$

Relational interpretation: eastbound(T) \leftarrow hasShortCar(T),hasClosedCar(T).



PROPOSITIONALIZED TRAIN_TABLE

train(T)	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

Propositionalization algorithms are available online in the ClowdFlows platform

- ClowdFlows browsed-based DM platform for data mining in the cloud and workflow sharing on the web (Kranjc et al. 2012)
- Example workflow: Propositionalization with RSD available in ClowdFlows at http://clowdflows.org/workflow/611/



Propositionalization in ClowdFlows

 Example workflow: Comparison of propositionalization algorithms (RSD, ReIF, ...), available in ClowdFlows at http://clowdflows.org/workflow/4018/



TECHNOLOGIES

Talk outline

- Intro, to Machine Learning and Data Mining
- Inductive Logic Programming (ILP) and Relational Data Mining (RDM)
- Propositionalization approach to RDM
- Semantic data mining
- Summary and future work



Relational and semantic data mining

ILP, relational learning, relational data mining

- Learning from complex relational databases
- 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









Using domain ontologies

Using domain ontologies as background knowledge

- •E.g., 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





Using domain ontologies

- Using background knowledge in data mining has been a topic of extensive research
 - Hierarchical attribute values, hierarchy/taxonomy of attributes, since 1986
 - ILP, relational data mining, propositionalization, since 1991
 - Ontologies (Tim Berners-Lee), since 1989
 - accepted formalism for consensual knowledge representation for Semantic Web applications, a basic for the Semantic Web
 - Description logic, OWL, Protégé ontology editor
 - Using ontologies in data mining, since 2004



Semantic Data Mining

Ontology-driven (semantic) data mining is an emerging research topic
Semantic Data Mining (SDM) - a term denoting:

the new challenge of mining semantically annotated resources, with ontologies used as background knowledge in mining experimental data
approaches with which semantic data are mined





Find: a classification model, a set of patterns

Our early work: Semantic subgroup discovery

- The 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 (Železny and Lavrac, MLJ 2006)
 Applied to semantic data mining in a biomedical application by using the Gene Ontology as background knowledge in analyzing microarray data



RSD: Propositionalization approach to RDM and SDM

			_	сı	isto	mer			
		ID /	Zip	S ex	So St	\lim_{com}	ie ge	Cl ub	\mathbf{R}
	/								
		3478	34677	m	si	60-7	0 32	me	n1
	/	3479	43666	f	\mathbf{ma}	80-9	90 45	nm	re
/									
Custome ID	r Or ID	der	$\stackrel{\text{Store}}{\mathbb{D}}$	D M	eliv Iode	ery	Payı Mod	nt e	
			\	<u>†.</u>				-	
3478	214	10267	12	re	gula	ar	cash		
3478	344	6778	12	\ e:	rpre	ss	chec	k	
3478	472	8386	17	þ	gul	ar	chec	k	
3479	323	3444	17	- Re:	rpre	ss	cred	it	
3479	347	5886	12	r	gul	ar	cred	it	
				1	1				

/			
1	5	store	
Store ID	Size	Туре	Location
2	small	franchise	city
7	large	indep	rural
	Ŭ	1	

Relational representation of customers, orders and stores.

Propositionalization

Step 1

- constructing relational features
- 2. constructing a propositional table

	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	roto	0	0	1	1	1	0
g5	1	1	1	0	0 4	010	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



RSD: Propositionalization approach to RDM and SDM

					ietos	mor				
		ID 1	Zip	Sex	So St	In come	e ge	Cl ub	Re sp	
	/									
		3478	34677	m	si	60-70	0 32	me	nr	
	/	3479	43666	f	\mathbf{ma}	80-90) 45	nm	re	
						•••				
-/-			order							
Customer ID	B	der	$\frac{\text{Store}}{D}$	E M	eliv Iode	ery I	Payı Mod	nt e		
			\	1						
3478	214	10267	12	I	egula	ar (ash			
3478	344	6778	12	\le:	xpre	ss (hec	k		
3478	472	28386	17	h	gul	ar (hec	k		
3479	323	3444	17	- Re:	xpre	ss 🛛	red	it		
3479	347	5886	12	r	gul	ar (red	it		
					ł					
					1			s	tore	
					Sto	re ID	Siz	e (Туре	Lo
					12		sm	all	franchise	cit

Relational representation of customers, orders and stores

17

large indep

Location

TUTA

	f1	f2	f3	f4	f5	f 6		1		1	
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	0
g4	1	1	1	0	1	10 1 0	0	0	1	1	1
g5	1	1	1	0	0 /	010	0	1	1	0	1
g1	0	0	1	1	0	0	0	1	0	0	0
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	0

Propositionalization

Step 1

- 1. constructing relational features
- 2. constructing a propositional table

Step 2

Data Mining

	f1	f2	f3	f4	f5	f 6						fn
g 1	1	0	0	1	1	1	0	0	1	0	1	1
g 2	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	ro l o	0	0	1	1	1	0
g5	1	1	1	0	0 /	010	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

model, patterns, ...

Using GO as background knowledge in DNA microarray data analysis with relational subgroup discovery system RSD



Adam23

ožef Stefan Institute

Ontology terms (can be viewed as generalisations of individual genes) are described by firstorder features, presenting gene properties and relations between genes.

Jožef Stefan Institute



physiological process

all

cellular_component

cell

molecular_function

catalytic activity

biological process

cellular process

Application of RSD to microarray data analysis using GO as background knowledge (Zelezny et al. 2006, Trajkovski et al. 2008)

- 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).

3. Propositionalization: Determine truth values of features

A learn rules by a subgroup discovery algorithm CN2-SD

Step 2: Construction of first order features with supp. > min_supp.

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



RSD propositionalization step

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			U.			fn
1	0	0	1	1	1	0	0	1	0	1	1
0	1	1	0	1	1	0	0	0	1	1	0
0	1	1	1	0	0	1	1	0	0	0	1
1	1	1	0	1	1	0	0	1	1	1	0
1	1	1	0	0	1	0	1	1	0	1	0
0	0	1	1	0	0	0	1	0	0	0	1
1	1	0	0	1	1	0	1	0	1	1	1
0	0	0	0	1	0	0	1	1	1	0	0
1	0	1	1	1	0	1	0	0	1	0	1
	f1 1 0 1 1 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1	f1f210010111110011001100110010	f1f2f3100011011111111001110110110100101	f1f2f3f4100101100111111011101110001111000011110011001011	f1f2f3f4f51001101101011010110111001110011101011011001111011110111	f1f2f3f4f5f6100111011011011001110110110011110011001100110011001100100110101100101100	f1f2f3f4f5f610011100110110011001111011011101010110010110010110010110010100101101101101110	f1f2f3f4f5f61001110001101100011011001110111110110011101101110110011011101101100110110110101011010101101101101	f1f2f3f4f5f6100111001011011000011011001011011011110110111101101111001101110110101101101100110110101101011101101011101101001101101001101101000	f1f2f3f4f5f6II0100111001001101100110110011001110110011011101100111100110011110011001111011001101101101111101101111101101111101101011101101011101101011	f1 f2 f3 f4 f5 f6 I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I I

RSD: Rule construction with CN2-SD

	f1	f2	£3	f4	f5	f6		-	1			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

diffexp (Gene) \leftarrow f2 \land f3 [4,0]



RSD subgroup discovery in two steps

Step 1: Construct relational logic features of genes such as interaction(A, G) & function(G, protein_binding) (gene A interacts with another gene whose functions include protein binding) and propositional table construction with features as attributes

 Step 2: Using these features to discover and describe subgroups of genes that are differentially expressed (e.g., belong to class DIFF.EXP. of top 300 most differentially expressed genes) in contrast with RANDOM genes (randomly selected genes with low differential expression).

Sample subgroup description: diffexp(A) :- interaction(A,B) AND function(B,'GO:0004871') AND process(B,'GO:0009613')

RSD implementation in Orange4WS

RSD implemented as a workflow in Orange4WS

- propositionalization
- subgroup discovery algorithms: SD, Apriori-SD, CN2-SD
- applied also to standard ILP problems, e.g. mutagenicity prediction

mutagenic(M) ← feature121(M), feature235(M)



- More recent work: semantic subgroup discovery with SEGS
- Gene set enrichment: moving from single gene to gene set analysis
 - A gene set is enriched if the genes in the set are statistically significantly differentially expressed compared to the rest of the genes.
 - Observation: e.g., an 20% increase in all genes members of a biological pathway may alter the execution of this pathway ... and its impact on other processes ... significantly more then a 10fold increase in a single gene.



- Gene set enrichment methods:
 - Single GO terms:
 - Gene Set Enrichment Analysis (GSEA)
 - Parametric Analysis of Gene Set Enrichment (PAGE)

 Conjunctions of GO terms: SEGS - special purpose semantic subgroup discovery system for Searching for Enriched Gene Sets



• The SEGS approach:

– fuse information from GO, KEGG and ENTREZ

- Gene Ontology (GO): standardized biological terms used to annotate gene products: Molecular Functions, Biological Processes, Cellular Components
- Kyoto Encyclopedia of Genes and Genomes (KEGG): manually drawn pathway maps representing the knowledge on the molecular interaction and reaction networks
- ENTREZ: gene annotations with GO and KO terms and genegene interaction data
- generate gene set candidates by performing top-down search of rules, formed as conjunctions of ontology terms as conjunctions of GO, KEGG and ENTREZ terms
- combine Fisher, GSEA and PAGE enrichment tests to select most interesting groups of differentially expressed genes

 SEGS workflow is implemented in the Orange4WS data mining environment



 SEGS is also implemented also as a Web applications

(Trajkovski et al., IEEE TSMC 2008, Trajkovski et al., JBI 2008)



SEGS Descriptive Microarray Dat	a Analysis - Mozilla Firefox	1	Mo	zilla Firefox						
Eile Edit View History Bookmarks	Tools Help	Eik	e	Edit View History Book	marks	<u>T</u> ools <u>H</u> elp				
🔇 💽 - C 🗙 🏠 🗋	http://kt.ijs.si/software/SEGS/index.php?show=tool 🛛 🏠 🔹 biomine project 🔎		5	🕑 - C 🗙 🏠		http://kt.ijs.si/softwa	are/SEGS/work_dir/ph	prRlvFW.0.all.ht 🧹	3 • G•	piomine project 🎾
🔎 Most Visited 📄 Petra's Home Page		2	Mo	st Visited 📄 Petra's Home Pa	age					
💠 mands 📄 SEGS 🗵 🛛 🎼 Micr	roarray 😒 Gene set en 🔹 Biomine proj 🥨 Tulip Softwa 🤯 Tulip Softwa 🚽 🔹	4	man	ds 📄 htttml 🔯 📗	Microa	array 😒 Gene :	set en 🔹 🕸 Biomine	e proj 🥨 Tulip S	5oftwa 🥨 1	ulip Softwa 🔶
SEGS	×	P E	Pro En:	oject: [] riched geneset	s for	· class A				
Main page Publications	Project Name: (optional)	fe	эш	ad by Combining p	-valu	es				
Web tool Downloads GO & KEGG	Annotation data: Molecular Functions Biological Processes	1	#	Description	Set size	#DE_Genes	Fisher p-value (unadjusted p-value)	GSEA p-value (Enricment score)	PAGE p-value (Z-score)	Agregate p-value
Gene annotations Gene interactions Gene expression data	Cellular Components KEGG Orthology Gene interactions Constraints: Number of DE genes: 200		1	Func(monovalent inorganic cation transporter activity), Proc(monovalent inorganic cation	<u>26</u>	<u>10</u>	0.000 (9.20e-07)	0.010 (0.362)	0.020 (3.767)	0.010
Igot Hajkovski Nada Lavrac	Minimal set size: 20 (min=20) Output: Maximal p-value: 0.05 Combine p-values: Fisher 1.0 GSEA 1.0 PAGE 1.0 Report top 100 most enriched gene sets.		2	Func(monovalent inorganic cation transporter activity), Proc(monovalent inorganic cation transport), Comp(integral to	<u>24</u>	<u>9</u>	0.010 (4.23e-06)	0.010 (0.352)	0.020 (3.671)	0.013
DEPARTMENT OF KNOWLEDGE TECHNOLOGIES Jodef Stefan Institute	Summarize descriptions Upload: input file: Browse SEND	-	3	Func(monovalent inorganic cation transporter activity), Proc(transport), Comp(integral to membrane),	<u>26</u>	<u>9</u>	0.010 (9.10e-06)	0.040 (0.323)	0.020 (3.801)	0.023
× Find: garr	Jext 👚 Previous 🖌 Highlight all 🦵 Match case	×	Fi	nd: garr	↓ <u>N</u> e:	kt 👚 Previous 🖌	Highlight <u>a</u> ll 🧮 Ma	t <u>c</u> h case		
Done		Do	ne							



Biomine graph exploration (Toivonnen et al.)

- SEGS can be combined with other biomedical resources, such as BioMine
- **BioMine graph** contains information from public databases, including annotated sequences, proteins, orthology groups, genes and gene expressions, gene and protein interactions, PubMed articles, and different ontologies.
 - nodes (~1 mio) correspond to different concepts (such as gene, protein, domain, phenotype, biological process, tissue)
 - semantically labeled edges (~7 mio) connect related concepts
- BioMine query engine answers queries to potentially discover new links between entities by sophisticated graph exploration algorithms

Complex data mining methodology SegMine = SEGS + BioMine

SegMine overview

100mm1-P2 25.71 8.15 7.69 95.46 1.53 50.94 2.89 184.58 5.45 292.55 9.34 7.04 4.41 0.35 130 98	1donor2-P2 41.29 41.84 108.73 86.82 1.11 53.07 0.64 150.62 1.51 359.93 12.14 52.98 39.9 0.65 43.62	1donor3-P2 33.11 12.85 291.82 110.13 15.98 36.16 4.24 119.35 0.72 465.48 9.67 47.63 17.72 2.2 2.3 151.49	2donor1-P11 49,53 6.7 9,71 118.57 1.41 43.25 1.63 141.87 0.34 289,49 26.42 0.34 95,51	2donor2-P7 54.89 7.61 105.98 92.53 1.25 73.51 6.91 155.45 2.83 344.66 5.39 55.46 19.17 0.41 101.89	2donor3-P8 36.59 9.82 84.38 118.26 5.03 32.19 4.41 157.76 0.65 291.91 8.37 40.43 12.15 1.95 26.77	SEGS AND INTERACT: transcription coactivator activity RULE 1 := organetic organization AND INTERACT: transcription coactivator activity RULE 2 := cellular macromolecule metabolic process AND INTERACT: thomatin binding RULE 3 := cellular response to stimulus AND INTERACT: RNA binding Expert AND INTERACT: RNA binding	다 5551111 다 551111 다 다 다 다 다 다 다 다 다 다 다
mia (ex	raw croar cpres	data ray ssior	a froi expe i of (m a erime gene	ent es)	knowledge interpretation of gene expression dat from ontologies rules, clusters, genesets	a:
Ę	er an	xper alys	t is	3		Biomine public databases	

Podpečan et al., BMC Bioinformatics 2011

SEGS + BioMine outputsSEGS output:BioMine query output:



Project: []

Enriched genesets for class A

found by Combining p-values

#	Description	Set size	#DE_Genes	Fisher p-value (unadjusted p-value)	GSEA p-value (Enricment score)	PAGE p-value (Z-score)	Agregate p-value			
1	Func(monovalent inorganic cation transporter activity), Proc(monovalent inorganic cation transport),	<u>26</u>	<u>10</u>	0.000 (9.20e-07)	0.010 (0.362)	0.020 (3.767)	0.010			
2	Func(monovalent inorganic cation transporter activity), Proc(monovalent inorganic cation transport), Comp(integral to membrane),	<u>24</u>	<u>9</u>	0.010 (4.23e-06)	0.010 (0.352)	0.020 (3.671)	0.013			
3	Func(monovalent inorganic cation transporter activity), Proc(transport), Comp(integral to membrane),	<u>26</u>	<u>9</u>	0.010 (9.10e-06)	0.040 (0.323)	0.020 (3.801)	0.023			
X	× Find: garr 🕹 Next 👚 Previous 🖌 Highlight all 🗖 Match case									





SegMine methodology implemented in Orange4WS





http://segmine.ijs.si/

Biomedical applications of SegMine methodology





- Combination of advanced data processing and mining algorithms
- Enables semantic analysis of gene expression using background knowledge in the form of ontologies
- SegMine tool is actively used at the National Institute for Biology
- Successful application in human stem cell data analysis: new hypotheses, enabling better understanding of cell senescence mechanisms
- General purpose Semantic Data Mining algorithm g-SEGS is also available in Orange4WS



From SEGS to SDM-SEGS: Generalizing SEGS

 SDM-SEGS: a semantic data mining system generalizing SEGS



- Discovers subgroups both for ranked and labeled data
- Exploits input ontologies in OWL format
- Implemented as a web service in Orange4WS
 - Can also be used e.g. in Taverna

Recent work

 Semantic Subgroup Discovery workflows in Orange4WS and ClowdFlows (Vavpetič et al., 2012)




Recent biomedical applications

 Subgroup discovery and semantic explanation of subgroups on breast cancer data (Vavpetič et al., JIIS 2014)



 The workflow, implemented in ClowdFlows, is available for sharing at http://clowdflows.org/workflow/1283/

Future challenges for Semantic Data Mining

- Current SDM scenario: Mining empirical data with ontologies as background knowledge
 - abundant empirical data, but
 - scarce background knowledge
- Future SDM scenarios:
 - envisioning a growing amount of semantic data
 - abundance of ontologies and semantically anotated data collections
 - e.g. Linked Data

 –over 6 billion RDF triples
 –over 148 million links



Future work

- We may envision a paradigm shift from data mining to knowledge mining
- The envisioned future Semantic data mining scenario in mining the Semantic Web:
 - mining knowledge encoded in domain ontologies,
 - constrained by annotated (empirical) data collections.



Summary: RDM and SDM in Context



Acknowledgements

- Work on subgroup discovery was done jointly with D. Gamberger (RBI) and P. Kralj Novak (JSI)
- Work on relational data mining and semantic data mining was done jointy with I.Trajkovski (Skopje Uni.), F. Železny (CTU, Prague), J. Tolar (Univ. of Minnesota), I. Mozetič and A. Vavpetič (JSI), and colleagues from the National Institute of Biology
- Work on new data mining platforms and advanced workflows was done jointly with A. Vavpetič, V.
 Podpečan and J. Kranjc (JSI)

Acknowledgements



