

Subgroup Discovery in Defect Prediction

Rachel Harrison, Oxford Brookes University

Daniel Rodríguez, Univ of Alcalá

José Riquelme, Univ of Seville

Roberto Ruiz, Pablo de Olavide University



Universidad
de Alcalá



Outline

Supervised Description

Subgroup Discovery

Preliminary Experimental Work

- Datasets
- Algorithms (SD and CN₂-SD)
- Results

Conclusions and future work

Descriptive Models

Typically, ML algorithms have been divided into:

- Predictive (Classification, Regression, temporal series)
- Descriptive (Clustering, Association, summarisation)

Recently, ***supervised descriptive rule discovery*** is being introduced in the literature.

- The aim is to understand the underlying phenomena, not to classify new instances, i.e., to find information about a specific value in the class attribute.
- The information should be useful to the domain expert and easily interpretable.
- Types of supervised descriptive techniques include:
 - Contrast Set Mining (CSM)
 - Emerging Pattern Mining (EPM)
 - Subgroup discovery (SD)

SD – Definition

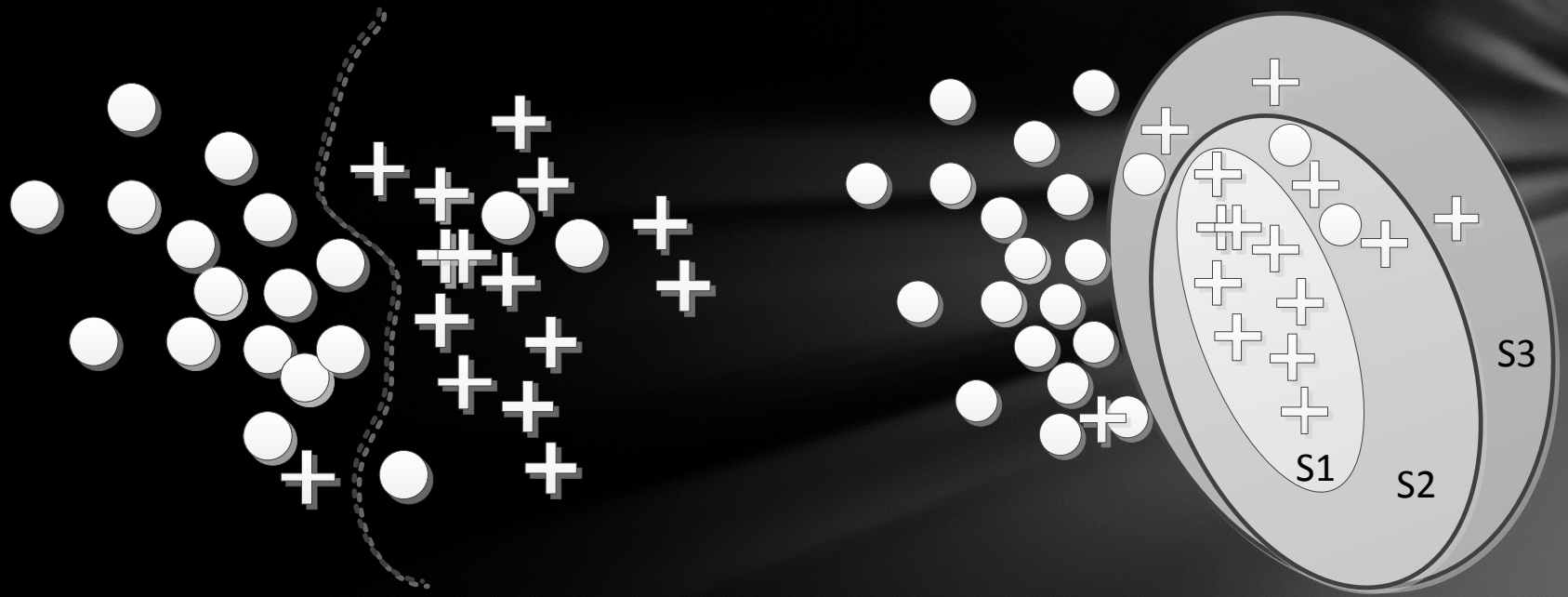
SD algorithms aims to find subgroups of data that are statistically different given a property of interest. [Klösger, 96; Wrobel, 97]

- SD lies between predictive (finding rules given historical data and a property of interest) and descriptive tasks (discovering interesting patterns in data).
- SD algorithms generally extract rules subsets of the data of previously specified the concept, for example defective modules from a software metrics repository.
- Rules have also the "**Condition** → **Class**" where the condition is the conjunction of a set of selected variables (pairs attribute–value) among all variables.
 - Advantages of rules include that are well known representation easily understandable by the domain experts
- So far, SD has majoritarily been applied to the medical domain.

SD vs. Classification

	<i>Classification</i>	<i>Subgroup Discovery</i>
Induction	Predictive	Descriptive
Output	Set of classification rules (dependent rules)	Individual Rules to describe subgroups (independent rules)
Purpose	To learn a model for classification or prediction	To find interesting and interpretable patterns with respect to a specific attribute

SD vs. Classification



Following [Herrera et al, 2011]

SD Algorithms

SD algorithms could be classified as:

- Exhaustive (e.g.: SD-map, Apriori-SD)
- Heuristic (e.g.: SD, CN2-SD)
 - Fuzzy genetic algorithms (SDIGA, MESDIF, EDER-SD)

Or from their origin, evolved from different communities:

- Extension of classification algorithms (SD, CN2-SD, etc.)
- Extension of association algorithms (Apriori-SD, SD₄TS, SD-Map, etc.)

Comprehensive survey by [Herrera et al. 2011]

Quality Measures in SD

Measures of Complexity

- Number of rules: It measures the number of induced rules.
- Number of conditions: It measures the number of conditions in the antecedent of the rule.

Measures of Generality

- Coverage: $Cov(R) = \frac{n(Cond)}{N}$

where N is the number of samples and $n(Cond)$ is the no. of instances that satisfy the antecedent of the rule.

- Support: $Sup(R) = \frac{n(Cond \cdot Class)}{N}$

where $n(Cond \cdot Class)$ is the no. of instances that satisfy both the condition and the class

Quality Measures in SD

Measures of precision

- Confidence: $Conf(R) = \frac{n(Cond \cdot Class)}{n(Cond)}$
- Precision Q_c : $Q_c = \frac{n(Class \cdot Cond) - c \cdot n(\neg Class \cdot Cond)}{n(Cond)}$
- Precision Q_g : $Q_g = \frac{n(Class \cdot Cond)}{n(\neg Class \cdot Cond) + g}$

Measures of interest

- Significance: $Sig(R) = 2 \sum_{k=1}^n n(cond \cdot Class_k) \cdot \log \frac{n(Cond \cdot Class_k)}{n(Class_k) \cdot p(Cond)}$

Other Measures

Sensitivity: $Sens(R) = TPr = \frac{TP}{Pos} = \frac{n(Class \cdot Cond)}{n(Class)}$

False alarm: $FA(R) = FPr = \frac{FP}{Neg} = \frac{n(\neg \overline{Class}Class \cdot Cond)}{n(\neg Class)}$

Specificity: $Spec(R) = \frac{TN}{TN + FP} = \frac{TN}{Neg} = \frac{n(\neg Class \cdot \neg Cond)}{n(\neg Class)}$

Unusualness: $WRAcc(R) = \frac{n(Cond)}{N} = \left(\frac{n(Class \cdot Cond)}{n(Cond)} - \frac{n(Class)}{N} \right)$

Experimental Work – Datasets

NASA Datasets

- Originally available from:
 - <http://mdp.ivv.nasa.gov/>
- From PROMISE, using the ARFF format (Weka – data mining toolkit):
 - <http://promisedata.org/>
 - Boetticher, T. Menzies, T. Ostrand, Promise Repository of Empirical Software Engineering Data, 2007.

Bug prediction dataset

- <http://bug.inf.usi.ch/>
- D'Ambros, M., Lanza, M., Robbes, Romain, Empirical Software Engineering (EMSE), In press, 2011

Datasets Characteristics

Some of these datasets are highly unbalanced, with duplicates and contradictory instances, and irrelevant attributes for defect prediction.

	<i># inst</i>	<i>Non-def</i>	<i>Def</i>	<i>% Def</i>	<i>Lang</i>
CM ₁	498	449	49	9.83	C
KC ₁	2,109	1,783	326	15.45	C++
KC ₂	522	415	107	20.49	C++
KC ₃	458	415	43	9.39	Java
MC ₂	161	109	52	32.29	C++
MW ₁	434	403	31	7.14	C++
PC ₁	1,109	1,032	77	6.94	C
Eclipse JDT Core	997	791	206	20.66	Java
Eclipse PDE-UI	1,497	1,288	209	13.96	Java
Equinox	324	195	129	39.81	Java
Lucene	691	627	64	9.26	Java
Mylyn	1,862	1,617	245	13.15	Java

Metrics Used from the Datasets

For the NASA datasets:

	<i>Metric</i>	<i>Definition</i>
McCabe	loc	McCabe's Lines of code
	v(g)	Cyclomatic complexity
	ev(g)	Essential complexity
	iv(g)	Design complexity
Halstead	uniqOp	Unique operators, n_1
	uniqOpnd	Unique operands, n_2
	totalOp	Total operators, N_1
	totalOpnd	Total operands N_2
Branch	branchCount	No. branches of the flow graph
Class	defective?	Reported defects? (true/false)

For the OO datasets:

	<i>Metric</i>	<i>Definition</i>
C&K	wmc	Weighted Method Count
	dit	Depth of Inheritance Tree
	cbo	Coupling Between Objects
	noc	No. of Children
	lcom	Lack of Cohesion in Methods
	rfc	Response For Class
Class	defective?	Reported defects?

Algorithms

The algorithms used:

- The **Subgroup Discovery** algorithm (SD) [Gamberger, 02] is a covering rule induction algorithm that using beam search aims to find rules that maximise:

$$q_g = \frac{TP}{FP+g}$$

where TP and FP are the number of true and false positives respectively and g is a generalisation parameter that allow us to control the specificity of a rule, i.e., balance between the complexity of a rule and its accuracy.

- The **CN2-SD** [Lavrač, 04] algorithm is an adaptation of the CN2 classification rule algorithm [Clark, 89]. It induces subgroups in the form of rules using as a quality measure the relation between true positives and false positives. The original algorithm consists of a search procedure using beam search within a control procedure and the control procedure that iteratively performs the search.
 - The CN2-SD algorithm uses *Weighted Relative Accuracy* (explained next) as a covering measure of the quality of the induced rules.

Tool:

- Orange: <http://orange.biolab.si/>

Examples Rules – KC2 Dataset

	#	<i>pd</i>	<i>pf</i>	<i>TP</i>	<i>FP</i>	<i>Rules</i>
SD	0	.24	0	26	0	$ev(g) > 4 \wedge totalOpnd > 117$
	1	.28	.01	30	5	$iv(G) > 8 \wedge uniqOpnd > 34 \wedge ev(g) > 4$
	2	.27	.01	29	5	$loc > 100 \wedge uniqOpnd > 34 \wedge ev(g) > 4$
	3	.27	.01	29	5	$loc > 100 \wedge iv(G) > 8 \wedge ev(g) > 4$
	4	.27	.01	29	5	$loc > 100 \wedge iv(G) > 8 \wedge totalOpnd > 117$
	5	.24	.01	26	5	$iv(G) > 8 \wedge uniqOp > 11 \wedge totalOp > 80$
	6	.24	.01	26	5	$iv(G) > 8 \wedge uniqOpnd > 34$
	7	.23	.01	25	5	$totalOpnd > 117$
	8	.31	.01	34	5	$loc > 100 \wedge iv(G) > 8$
	9	.29	.01	32	5	$ev(g) > 4 \wedge iv(G) > 8$
	10	.29	.01	32	5	$ev(g) > 4 \wedge uniqOpnd > 34$
	11	.28	.01	30	5	$loc > 100 \wedge ev(g) > 4$
	12	.28	.01	30	5	$iv(G) > 8 \wedge uniqOp > 11$
	13	.35	.01	38	5	$ev(g) > 4 \wedge totalOp > 80 \wedge v(g) > 6 \wedge uniqOp > 11$
	14	.27	.01	29	5	$iv(G) > 8 \wedge totalOp > 80$
	15	.27	.01	29	5	$ev(g) > 4 \wedge totalOp > 80 \wedge uniqOp > 11$
	16	.26	.01	28	5	$ev(g) > 4 \wedge totalOp > 80 \wedge v(g) > 6$
	17	.26	.01	28	5	$loc > 100 \wedge uniqOpnd > 34$
	18	.31	.01	34	5	$ev(g) > 4 \wedge totalOp > 80$
	19	.31	.01	34	5	$iv(G) > 8$
CN2-SD	0	.35	.01	38	5	$uniqOpnd > 34 \wedge ev(g) > 4$
	1	.4	.02	43	9	$totalOp > 80 \wedge ev(g) > 4$
	2	.78	.21	84	88	$uniqOP > 11$

Example Rules – JDT Core Dataset

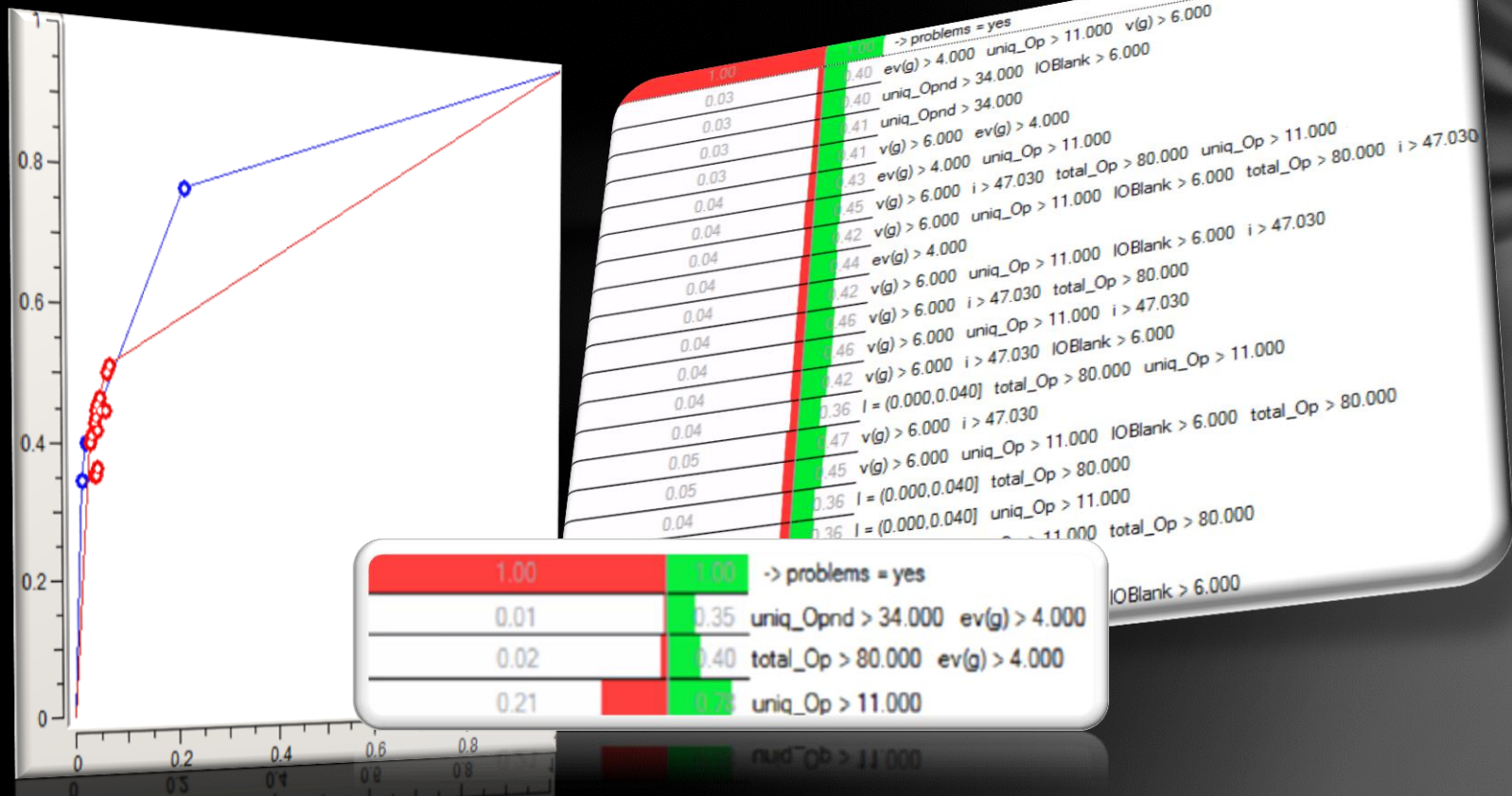
	#	pd	pf	TP	FP	Rules
SD	0	.27	.02	56	16	lcom > 171 \wedge rfc > 88 \wedge cbo > 16 \wedge wmc > 141
	1	.3	.02	62	16	rfc > 88 \wedge wmc > 141 cbo > 16
	2	.3	.02	62	16	cbo > 16 \wedge wmc > 141
	3	.29	.02	60	16	lcom > 171 \wedge rfc > 88 \wedge wmc > 141
	4	.29	.02	60	16	lcom > 171 \wedge wmc > 141
	5	.33	.03	68	24	rfc > 88 \wedge wmc > 141
	6	.32	.03	66	24	rfc > 88 \wedge wmc > 141 \wedge dit \leq 5
	7	.33	.03	68	24	wmc > 141
	8	.32	.03	66	24	dit \leq 5 \wedge wmc > 141
	9	.18	.02	38	16	wmc > 141 \wedge noc = 0 \wedge dit \leq 5
	10	.19	.02	40	16	wmc > 141 \wedge noc = 0
	11	.18	.02	38	16	cbo > 16 \wedge rfc > 88 \wedge noc > 0 dit \leq 5
	12	.42	.04	87	32	cbo > 16 \wedge rfc > 88 \wedge dit \leq 5
	13	.3	.03	62	24	lcom > 171 \wedge rfc > 88 \wedge cbo > 16 \wedge dit \leq 5
	14	.2	.02	42	16	cbo > 16 \wedge rfc > 88 \wedge noc > 0
	15	.24	.03	50	24	cbo > 16 \wedge rfc > 88 \wedge noc = 0 \wedge dit \leq 5
	16	.45	.05	93	40	cbo > 16 \wedge rfc > 88
	17	.32	.03	66	24	lcom > 171 \wedge rfc > 88 \wedge cbo > 16
	18	.25	.03	52	24	cbo > 16 \wedge rfc > 88 \wedge noc = 0
CN2-SD	19	.33	.05	68	40	cbo > 16 \wedge lcom > 171
	0	.45	.05	93	40	rfc > 88 \wedge cbo > 16
	1	.55	.09	114	72	rfc > 88

Cross-validation Results (10 CV)

[illegible]

Visualisation of SD

ROC and Rule visualisation for KC2 (SD & CN2-SD)



Conclusions

Rules obtained using SD are intuitive but needed to be analysed by an expert.

The metrics used for classifiers cannot be directly applied in SD and need to be adapted.

Current and future work

- Further validation and application in other software engineering domains, e.g., project management.
- SD is a search problem!
 - Development of new algorithms and metrics
 - EDER-SD (Evolutionary Decision Rules SD) in Weka
 - Unbalanced data (ROC, AUC metrics?), etc.
 - Feature Selection (as a pre-processing step, part of the algorithm?, which metrics really influence defects)
 - Discretisation
 - Different search strategies and fitness functions (and multi-objective!)
 - Use of global optimisation (set of metrics) vs. local metrics (individual metrics)

References

Kralj, P., Lavrac, N., Webb GI (2009) Supervised Descriptive Rule Discovery: A Unifying Survey of Contrast Set, Emerging Pattern and Subgroup Mining. *Journal of Machine Learning Research* 10: 377–403

Kloesgen, W. (1996), Explora: A Multipattern and Multistrategy Discovery Assistant. In: *Advances in Knowledge Discovery and Data Mining*, American Association for Artificial Intelligence, pp 249–271

Wrobel, S. (1997), An Algorithm for Multi-relational Discovery of Subgroups. *Proceedings of the 1st European Symposium on Principles of Data Mining and Knowledge Discovery*, Springer, LNAI, vol 1263, pp 78–87

Bay S., Pazzani, M. (2001) Detecting Group Differences: Mining Contrast Sets. *Data Mining and Knowledge Discovery* 5: 213–246

Dong, G., Li, J. (1999) Efficient Mining of Emerging Patterns: Discovering Trends and Differences. In: *Proceedings of the 5th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM Press, pp 43–52

Herrera, F., Carmona del Jesus, C.J., Gonzalez, P., and del Jesus, M.J., An overview on subgroup discovery: Foundations and applications, *Knowledge and Information Systems*, 2011 – In Press.

Gamberger, D., Lavrac, N.: Expert-guided subgroup discovery: methodology and application. *Journal of Artificial Intelligence Research* 17 (2002) 501–527

Lavrac, N., Kavsek, B., Flach, P., Todorovski, L.: Subgroup discovery with CN2-SD. *The Journal of Machine Learning Research* 5 (2004) 153–188

Clark, P., Niblett, T. (1989) , The CN2 induction algorithm, *Machine Learning* 3 261–283