

RELATIONAL SUBGROUP DISCOVERY FOR GENE EXPRESSION DATA MINING

F. Železný*, J. Tolar**, N. Lavrač*** and O. Štěpánková*

* Gerstner Laboratory, Department of Cybernetics,
Czech Technical University, Technická 2,
166 27 Prague, Czech Republic

** Division of Hematology, Oncology, Blood and Marrow Transplant, Department of
Pediatrics, University of Minnesota Medical School
420 Delaware Street, 55455 Minneapolis, USA

*** Department of Knowledge Technologies, Jožef Stefan Institute
Jamova 39, 1000 Ljubljana, Slovenia
and

Nova Gorica Polytechnic, Vipavska 13, 5000 Nova Gorica, Slovenia

zelezny@fel.cvut.cz, tolar003@umn.edu, nada.lavrac@ijs.si, step@labe.felk.cvut.cz

Abstract: We propose a methodology for predictive classification from gene expression data, able to combine the robustness of high-dimensional statistical classification methods with the comprehensibility and interpretability of simple logic-based models. We first construct a robust classifier combining contributions of a large number of gene expression values, and then search for compact summarizations of subgroups among genes associated in the classifier with a given class. The subgroups are described by means of relational logic features extracted from publicly available gene annotations. The curse of dimensionality pertaining to the gene expression based classification problem due to the large number of attributes (genes) is turned into an advantage in the secondary subgroup discovery task, as here the original attributes become learning examples.

Introduction

Many tasks of automated knowledge discovery from gene expression microarray data by data mining algorithms aim at constructing classifiers able to diagnose a cancer type from a gene expression profile. See eg. the seminal papers [1, 2] for reference.

This task is characterized by the abundance of attributes (eg. simultaneously measured gene expression values) on one hand confronted with the shortage of the available samples (eg. patients/tissues subject to measurements) on the other hand. It is known from scientific discovery literature that such domains are prone to overfitting: overfitted classifiers are characterized by significantly decreased predictive accu-

racy on unseen samples compared to the training set accuracy.

Informally, in domains characterized by a small number of examples and a large number of attributes, overfitting occurs because some artefacts (“flukes”) of actually irrelevant attribute combinations can emerge simply by means of chance and appear significant with respect to the examples available to the machine learning algorithm.

To avoid the overfitting pitfall, state-of-art approaches construct complex classifiers that combine relatively weak contributions of up to thousands of genes (attributes) to classify a disease. Real-valued support vector machine (SVM) [3] models are currently specifically popular in the gene expression data mining domain. However, classifiers based on many real-valued attributes have an important drawback: they are not appropriate for expert interpretation. The complexity of such classifiers limits their transparency and consequently the biological insight they can provide. Although it is possible to extract the attributes with maximal contribution weight, the logical connections among the extracted attributes are then lost.

In our previous work [4], we have tested the feasibility of constructing simple yet robust logic-based classifiers, in the form of propositional rules amenable to direct expert interpretation, by an innovative algorithmic methodology of subgroup discovery. Such rules typically include two to five gene expression attributes and, in contrast to markers obtained from SVM schemes, these rules explicitly stress the importance of the correlation of the activity (or non-activity) of a narrow gene set. Following our work, further papers appeared [5, 6] proposing methods for constructing very simple, interpretable

models for gene expression data.

Despite the obvious attractiveness stemming from the interpretability of the classifiers yielded by the above mentioned approaches, they still fall short of the high dimensional (attribute-rich) classifiers in terms of predictive accuracy in some benchmark sub-tasks considered in [4, 2]. Preliminary experiments [7] suggest that certain accuracy gap manifests itself also in the predictive discovery tasks formulated in [8].

We thus seem to face an inevitable trade-off between interpretability on one hand and accuracy on the other hand pertaining to predictive gene expression data based models. Here we propose a methodology with the potential to overcome this challenge, combining the advantages of both the approaches of high-dimensional models and the simple logic models.

The fundamental idea is as follows. First, a high-dimensional classification model C is constructed from gene expression data D , wherein each sample is assigned a class label c out of a set of class labels \mathcal{C} . Depending on the particular inductive method employed, C may acquire different forms, but as a general rule it will associate a high number of genes ('predictors') to each of the target classes in D ; we denote this gene set by $G_C(c)$. In this paper we will confine ourselves to the straightforward way originally employed by [1], where $G_C(c)$ is simply a set of genes with expression highly correlated with the class c .

The second step aims at improving the interpretability of C . Informally, we do this by identifying subgroups of genes in $G_C(c)$ (for each $c \in \mathcal{C}$) which can be summarized in a compact way. Put differently, for each $c_i \in \mathcal{C}$ we search for compact descriptions of gene categories which correlate strongly with c_i and weakly with all $c_j \in \mathcal{C}, j \neq i$.

The *subgroup discovery* procedure just outlined is approached as another supervised machine learning process. This is, in a way, orthogonal to the primary discovery task in that the original attributes (genes) now become learning examples, each of which has a class label $c \in \mathcal{C}$. To apply a discovery algorithm, information about relevant features of the new examples is required. No such features (ie. 'attributes' of the original attributes – genes) are usually present in the gene expression microarray data sets themselves. However, this information can be extracted from a public database of gene annotations (in this paper, we use the Entrez Gene database maintained at the US National Center for Biotechnology Information).

In traditional machine learning, examples are expected to be described by a tuple of values corresponding to some predefined, fixed set of attributes. Note that a gene annotation does not straightforwardly correspond to a fixed attribute set, as it has an inherently *relational* character. For example, a

gene may be related to a variable number of cell processes, play role in variable numbers of regulatory pathways etc. This imposes 1-to-many relations hard to elegantly capture within an attribute set of fixed size. Furthermore, a useful piece information about a gene g may for instance be expressed by the feature

g interacts with another gene whose functions include protein binding.

Going even further, the feature may not include only a single interaction relation but rather consider entire chains of interactions. The difficulties of representing such features through attribute-value tuples is evident.

In summary, we are approaching the task of subgroup discovery from a relational data domain. For this purpose we employ the methodology of relational subgroup discovery we proposed in [9, 10] and implemented in the RSD algorithm [11]. Using RSD, we are able to discover knowledge such as

The expression of genes coding for proteins located in the integral to membrane cell component, whose functions include receptor activity, has a high correlation with the BCR class of acute lymphoblastic leukemia (ALL) and a low correlation with the other classes of ALL.

The RSD algorithm proceeds in two steps. First, it constructs a set of relational features in the form of first-order logic atom conjunctions. The entire set of features is then viewed as an attribute set, where an attribute has the value *true* for a gene (example) if the gene has the feature corresponding to the attribute. As a result, by means of relational feature construction we achieve the conversion of relational data into attribute-value descriptions.¹ In the second step, interesting gene subgroups are searched, such that each subgroup is represented as a conjunction of selected features. The subgroup discovery algorithm employed in this second step is an adaptation of the popular propositional rule learning algorithm CN2 [13].

Relational Feature Construction

The feature construction component of RSD aims at generating a set of relational features in the form for relational logic atom conjunctions. For example, the feature exemplified informally in the previous section has the relational logic form

$$\text{interaction}(g,G),$$
$$\text{function}(G,\text{protein_binding})$$

¹This process is referred to as *propositionalization* [12]

Here, upper cases denote existentially quantified variables and g is the key term that binds a feature to a specific example (here a gene).

The user specifies a grammar declaration which constraints the resulting set of constructed features. RSD accepts feature language declarations similar to those used in the inductive logic programming system Progol [14]. A declaration lists the predicates that can appear in a feature, and to each argument of a predicate a *type* and a *mode* are assigned. In a correct feature, if two arguments have different types, they may not hold the same variable. A mode is either *input* or *output*; every variable in an input argument of a literal must appear in an output argument of some preceding literal in the same feature. [15] further dictate the opposite constraint: every output variable of a literal must appear as an input variable of some subsequent literal. Furthermore, the maximum length of a feature (number of contained literals) is declared, along with optional constraints such as the maximum *variable depth* [14], maximum number of occurrences of a given predicate symbol in a feature, etc.

RSD generates an exhaustive set of features satisfying the language declarations as well as a the *connectivity* requirement, which stipulates that no feature may be decomposable into a conjunction of two or more features. For example, the following expression does not form an admissible feature

```
interaction(g,G1),
function(G1,protein_binding),
interaction(g,G2), component(G2,membrane)
```

since it can be decomposed into two separate features. We do not construct such decomposable expressions, as these are redundant for the purpose of subsequent search for rules with conjunctive antecedents. Furthermore the concept of undecomposability allows for powerful search space pruning [9, 10]. Notice also that the expression above may be extended into an admissible undecomposable feature if a further logic atom is added:

```
interaction(g,G1),
function(G1,protein_binding),
interaction(g,G2), component(G2,membrane),
interaction(G1,G2)
```

The construction of features is implemented as depth-first, general-to-specific search where refinement corresponds to adding a literal to the currently examined expression. During the search, each search node found to be a correct feature is listed in the output.

A remark is in turn concerning the way constants (such as `protein_binding`) are employed in features. Rather than making the user responsible for declaring all possible constants that may occur in features,

RSD extracts them automatically from the learning data. The user marks the types of variables which should be replaced by constants. For each constant-free feature, a number of different features are then generated, each corresponding to a possible replacement of the combination of the indicated variables with constants. RSD then only proceeds with those combinations of constants, which make the feature true for at least a pre-specified number of examples.

Finally, to evaluate the truth value of each feature for each example for generating the attribute-value representation of the relational data, the first-order logic resolution procedure is used, provided by a Prolog language engine.

Subgroup Discovery

A subgroup discovery task is defined as follows: *Given a population of individuals and a property of individuals we are interested in, find population subgroups that are statistically ‘most interesting’, e.g., are as large as possible and have the most unusual statistical (distributional) characteristics with respect to the property of interest.*

Notice an important aspect of the above definition: there is a predefined property of interest, meaning that a subgroup discovery task aims at characterizing population subgroups of a given *target* class. This property indicates that standard classification rule learning algorithms could be used for solving the task. However, while the goal of classification rule learning is to generate models (sets of rules), inducing class descriptions in terms of properties occurring in the descriptions of training examples, in contrast, subgroup discovery aims at discovering individual patterns of interest (individual rules describing the target class).

Rule learning typically involves two main procedures: the search procedure that performs search to find a single rule (described in this section) and the control procedure (the covering algorithm) that repeatedly executes the search in order to induce a set of rules.

Inducing a single subgroup rule

Our algorithm is based on an adaptation of the standard propositional rule learner CN2 [16, 13]. Its search procedure used in learning a single rule performs beam search, starting from the empty conjunct, successively adding conditions (relational features). In CN2, classification accuracy of a rule is used as a heuristic function in the beam search. The accuracy² of an induced rule of the form $H \leftarrow B$ (where H is the rule head - the target class, and B is the rule body formed of a conjunction of relational

²In some contexts, this quantity is called *precision*.

features) is equal to the conditional probability of head H , given that body B is satisfied: $p(H|B)$.

The accuracy heuristic $Acc(H \leftarrow B) = p(H|B)$ can be replaced by the *weighted relative accuracy* heuristic. Weighted relative accuracy is a reformulation of one of the heuristics used in MIDOS [17] aimed at balancing the size of a group with its distributional unusualness [18].

The weighted relative accuracy heuristic is defined as follows:

$$WRAcc(H \leftarrow B) = p(B) \cdot (p(H|B) - p(H)). \quad (1)$$

Weighted relative accuracy consists of two components: generality $p(B)$, and relative accuracy $p(H|B) - p(H)$. The second term, relative accuracy, is the accuracy gain relative to fixed rule $H \leftarrow true$. The latter rule predicts all instances to satisfy H ; a rule is only interesting if it improves upon this 'default' accuracy. Another way of viewing relative accuracy is that it measures the utility of connecting rule body B with rule head H . Note that it is easy to obtain high relative accuracy with very specific rules, i.e., rules with low generality $p(B)$. To this end, generality is used as a 'weight' which trades off generality of the rule (rule coverage $p(B)$) and relative accuracy ($p(H|B) - p(H)$).

In the computation of Acc and $WRAcc$ all probabilities are estimated by relative frequencies³ as follows:

$$Acc(H \leftarrow B) = p(H|B) = \frac{p(HB)}{p(B)} = \frac{n(HB)}{n(B)} \quad (2)$$

$$WRAcc(H \leftarrow B) = \frac{n(B)}{N} \left(\frac{n(HB)}{n(B)} - \frac{n(H)}{N} \right) \quad (3)$$

where N is the number of all the examples, $n(B)$ is the number of examples covered by rule $H \leftarrow B$, $n(H)$ is the number of examples of class H , and $n(HB)$ is the number of examples of class H correctly classified by the rule (true positives).

Inducing a set of subgroup rules

In CN2, for a given class in the rule head, the rule with the best value of the heuristic function found in the beam search is kept. The algorithm then removes all examples of the target class satisfying the rule's conditions (i.e. *covered* by the rule) and invokes a new rule learning iteration on the remaining training set. All negative examples (i.e., examples that belong to other classes) remain in the training set.

In this classical covering algorithm, only the first few induced rules may be of interest as subgroup descriptors with sufficient coverage, since subsequently

³Alternatively, the Laplace [19] and the m -estimate [20] could also be used.

induced rules are induced from biased example subsets, i.e., subsets including only positive examples not covered by previously induced rules. This bias constrains the population of individuals in a way that is unnatural for the subgroup discovery process, which is aimed at discovering interesting properties of subgroups of the entire population.

In contrast, RSD uses the *weighted covering algorithm*, which allows for discovering interesting subgroup properties in the entire population. The weighted covering algorithm modifies the classical covering algorithm in such a way that covered positive examples are not deleted from the set of examples which is used to construct the next rule. Instead, in each run of the covering loop, the algorithm stores with each example a count that indicates how many times (with how many induced rules) the example has been covered so far.

Initial weights of all positive examples e_j equal 1. In the first iteration of the weighted covering algorithm all target class examples have the same weight, while in the following iterations the contributions of examples are inverse proportional to their coverage by previously constructed rules; weights of covered positive examples thus decrease according to the formula $\frac{1}{i+1}$, where i is the number of constructed rules that cover example e_j . In this way the target class examples whose weights have not been decreased will have a greater chance to be covered in the following iterations of the weighted covering algorithm.⁴

The combination of the weighted covering algorithm with the weighted relative accuracy thus implies the use of the following *modified WRAcc* heuristic:

$$WRAcc(H \leftarrow B) = \frac{n'(B)}{N'} \left(\frac{n'(HB)}{n'(B)} - \frac{n(H)}{N} \right) \quad (4)$$

where N is the number of examples, N' is the sum of the weights of all examples, $n(H)$ is the number of examples of class H , $n'(B)$ is the sum of the weights of all covered examples, and $n'(HB)$ is the sum of the weights of all correctly covered examples.

An experiment

Materials and methods

We follow here the predictive classification problem defined in [8] aiming at distinguishing among six classes of pediatric acute lymphoblastic leukemia from gene expression profiles obtained by the Affymetrix HG-U133A microarray

⁴Whereas this approach is referred to as *additive* in [21], another option is the *multiplicative* approach, where for a given parameter $\gamma < 1$, weights of positive examples covered by i rules decrease according to γ^i . Both approaches have been implemented in RSD, but additive weights lead to better results.

chip. The data contains 132 class-labelled samples of expression vectors and it can be obtained along with a detailed description from <http://www.stjuderesearch.org/data/ALL3/>.

For each class c we first extracted a set of genes $G(c)$ whose expression is highly correlated with c . More precisely, for each gene g and class c we evaluated the function

$$f(g,c) = \frac{P(c,g)}{P(c) \cdot P(g)} \quad (5)$$

where $P(c)$ and $P(g)$ denote the probability (estimated via relative frequencies) that a randomly drawn sample falls into class c , and that g is expressed in that sample, respectively. $P(c,g)$ is the joint probability of both events. If the two events are statistically independent, $f(g,c)$ equals 1. If class c and expression of g are more likely to occur together, $f(g,c)$ will be greater than 1. We set a fixed threshold on $f(g,c)$ so that $\forall g,c : g \in G(c)$ whenever $f(g,c) \geq 5$. As a result we obtained on average 257 correlating genes for every class.

To access the annotation data for every gene considered, it was necessary to obtain unique gene identifiers from the microarray probe identifiers available in the original data. We achieved this by combining the *biobase*, *annotate* and *hgu133a* packages for the R system for statistical computing. The three packages are available from <http://www.bioconductor.org/> and R is available from <http://www.r-project.org/>.

Knowing the gene identifiers, the annotations can be accessed through hypertext queries to the Entrez Gene database, which is available at <http://www.ncbi.nlm.nih.gov/>. We developed a program script in the Python language, which automatically queries the server for the gene annotations, parses them and produces their structured, relational logic representations. This script is available on request to the first author.

Examples of subgroups discovered

Here we present two examples of gene subgroups discovered by RSD, converted to natural language from the original relational logic descriptions. We also attach a biologist's (2nd author) interpretation. The two examples are concerned with subgroups with unusually high frequency of the BCR (TEL, respectively) type of ALL.

BCR class: *Genes coding for proteins located in the integral to membrane cell component, whose functions include receptor activity.*

Comment: BCR/abl is a classic example of a leukemia driven by spurious expression of a fusion protein expressed as a continuously active kinase protein on the membrane of leukemic cells.

TEL class: *Genes coding for proteins located in the nucleus whose functions include protein binding and whose related processes include transcription.*

Comment: By contrast to BCR, the TEL leukemia is driven by expression of a protein, which is a transcription factor active in the nucleus.

As a result, our finding related to the location, function and processes associated to the subgroups, represent the most salient features of these respective types of acute lymphoblastic leukemia.

Discussion

In this paper we have proposed a methodology for predictive classification from gene expression data, able to combine the robustness of high-dimensional statistical classification methods with the comprehensibility and interpretability of simple logic-based models. Our methodology proposes to first construct a robust classifier combining contributions of a large number of gene expression values, and then finding compact, relational descriptions of subgroups among genes employed in the classifier.

It is noteworthy that the 'post-processing' step is also a machine learning task, in which the curse of dimensionality (the number of attributes – gene expressions measured) usually ascribed to the type of classification problem considered, actually turns into an advantage. The high number of *attributes*, incurring the risk of overfitting, turns into a high number of *examples*, which on the contrary works *against* overfitting in the subsequent subgroup discovery task. Furthermore, the dimensionality of the secondary attributes (relational features of genes extracted from gene annotations) can be conveniently controlled via suitable constraints of the language grammar used for the automatic construction of the gene features.

The limited experimental evaluation conducted within the present study obviously serves only as a first 'sanity check' and more elaborate assessment has yet to be conducted in a more extensive set of gene expression based classification problems. This is further emphasized by the fact that the division of ALL into the respective classes [8], which we accepted here, has not yet been fully established.

However, we have already confirmed the technical feasibility of the proposed approach as well as the fact that the subgroups discovered as early as in the first experimental run appear meaningful and amenable to interpretation to a biologist. We thus have high hopes on discovering novel, yet reliable knowledge from the relational combination of gene expression data with public gene annotation databases in future applications of our methodology.

Acknowledgements

F.Ž. is grateful to the Grant Agency of the Czech Academy of Sciences (CAS) for the support through the project KJB201210501 Logic Based Machine Learning for Genomic Data Analysis. J. T. is supported by Children Cancer Research Fund and University of Minnesota Cancer Center. O.Š. is supported by CAS through the project 1ET101210513 Relational Machine Learning for Biomedical Data Analysis.

References

- [1] T. R. GOLUB, D. K. SLONIM, P. TAMAYO, C. HUARD, M. GAASENBEEK, J. P. MESIROV, H. COLLIER, M. L. LOH, J. R. DOWNING, M. A. CALIGIURI, C. D. BLOOMFIELD, and E. S. LANDER. Molecular classification of cancer: Class discovery and class prediction by gene expression monitoring. *Science*, 286:531–537, 1999.
- [2] S. RAMASWAMY, P. TAMAYO, R. RIFKIN, S. MUKHERJEE, C. H. YEANG, M. ANGELO, C. LADD, M. REICH, E. LATULIPPE, J. P. MESIROV, T. POGGIO, W. GERALD, M. LODA, E. S. LANDER, and T. R. GOLUB. Multiclass cancer diagnosis using tumor gene expression signatures. In *Proceedings of the National Academy of Science USA*, volume 98, pages 15149–54, 2001.
- [3] T. HASTIE and J. TIBSHIRANI, R. ANDD FRIEDMAN. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*.
- [4] D. GAMBERGER, N. LAVRAČ, F. ŽELEZNÝ, and J. TOLAR. Induction of comprehensible models for gene expression datasets by subgroup discovery methodology. *Jr Biomedical Informatics*, 37(5):269–284, 2004.
- [5] S. A. VINTERBO and L. KIM, E. Y. OHNO-MACHADO. Small, fuzzy and interpretable gene expression based classifiers. *Bioinformatics*, 21(9):531–537, 2005.
- [6] M. MRAMOR, G. LEBAN, J. DEMSAR, and B. ZUPAN. Conquering the curse of dimensionality in gene expression cancer diagnosis: Tough problem, simple models. In *Proceedings of the 10th Conference on Artificial Intelligence in Medicine*, pages 514–523, 2005.
- [7] D. GAMBERGER. (personal communication).
- [8] M. E. ROSS, X. ZHOU, G. SONG, S. A. SHURTLEFF, K. GIRTMAN, W. K. WILLIAMS, H. C. LIU, R. MAHFOUZ, S. C. RAIMONDI, N. LENNY, A. PATEL, and J. R. DOWNING. Classification of pediatric acute lymphoblastic leukemia by gene expression profiling. *Blood*, 102(8):2961–9, 2003.
- [9] N. LAVRAČ, F. ŽELEZNÝ, and P. FLACH. RSD: Relational subgroup discovery through first-order feature construction. In *Proceedings of the 12th International Conference on Inductive Logic Programming*, pages 149–165. Springer, 2002.
- [10] N. LAVRAČ and F. ŽELEZNÝ. Propositionalization-based relational subgroup discovery with RSD. *Machine Learning*. (to appear in the 2005 special issue on statistical relational learning).
- [11] F. ŽELEZNÝ. RSD user's manual. Available at: <http://labe.felk.cvut.cz/~zelezny/rsd/rsd.pdf>.
- [12] N. LAVRAČ and P. A. FLACH. An extended transformation approach to inductive logic programming. *ACM Transactions on Computational Logic*, 2(4):458–494, October 2001.
- [13] P. CLARK and T. NIBLETT. The cn2 induction algorithm. *Machine Learning*, pages 261–283, 1989.
- [14] S. MUGGLETON. Inverse entailment and Progol. *New Generation Computing, Special issue on Inductive Logic Programming*, 13(3-4):245–286, 1995.
- [15] P. FLACH and N. LACHICHE. 1BC: A first-order Bayesian classifier. In S. Džeroski and P. Flach, editors, *Proceedings of the 9th International Workshop on Inductive Logic Programming*, pages 92–103. Springer, 1999.
- [16] P. CLARK and T. NIBLETT. Induction in noisy domains. In *Progress in Machine Learning (Proceedings of the 2nd European Working Session on Learning)*, pages 11–30. Sigma Press, 1987.
- [17] S. WROBEL. An algorithm for multi-relational discovery of subgroups. In Jan Komorowski and Jan Zytkow, editors, *Proceedings of the First European Symposium on Principles of Data Mining and Knowledge Discovery (PKDD-97)*, pages 78–87. Springer, 1997.
- [18] W. KLOESGEN. Explora: A multipattern and multistrategy discovery assistant. In *Advances in Knowledge Discovery and Data Mining*, pages 249–271. AAAI Press, Menlo Park, CA, 1996.
- [19] P. CLARK and R. BOSWELL. Rule induction with CN2: Some recent improvements. In *Proceedings of the 5th European Working Session on Learning*, pages 151–163. Springer, 1991.
- [20] B. CESTNIK. Estimating probabilities: A crucial task in machine learning. In *Proceedings of the 9th European Conference on Artificial Intelligence*, pages 147–149. Pitman, 1990.
- [21] N. LAVRAČ, B. KAVŠEK, P. FLACH, and L. TODOROVSKI. Subgroup discovery with CN2-SD. *Journal of Machine Learning Research*, 5:153–188, 2004.