

# Structural Regression Trees <sup>1</sup>

Stefan Kramer

Austrian Research Institute for Artificial Intelligence

Schottengasse 3, A-1010 Vienna, Austria

E-mail: `stefan@ai.univie.ac.at`

**Keywords:** Inductive Learning, Inductive Logic Programming,  
Regression

## Abstract

In many real-world domains the task of machine learning algorithms is to learn a theory predicting numerical values. In particular several standard test domains used in Inductive Logic Programming (ILP) are concerned with predicting numerical values from examples and relational and mostly non-determinate background knowledge. However, so far no ILP algorithm except one can predict numbers and cope with non-determinate background knowledge. (The only exception is a covering algorithm called FORS.)

In this paper we present Structural Regression Trees (SRT), a new algorithm which can be applied to the above class of problems by integrating the statistical method of regression trees into ILP. SRT constructs a tree containing a literal (an atomic formula or its negation) or a conjunction of literals in each node, and assigns a numerical value to each leaf. SRT provides more comprehensible results than purely statistical methods, and can be applied to a class of problems most other ILP systems cannot handle. Experiments in several real-world domains demonstrate that the approach is competitive with existing methods, indicating that the advantages are not at the expense of predictive accuracy.

---

<sup>1</sup>This research is sponsored by the Austrian *Fonds zur Förderung der Wissenschaftlichen Forschung (FWF)* under grant number P10489-MAT. Financial support for the Austrian Research Institute for Artificial Intelligence is provided by the Austrian Federal Ministry of Science, Research, and Arts.

# 1 Introduction

Many real-world machine learning domains involve the prediction of a numerical value. In particular several test domains used in Inductive Logic Programming (ILP) (including the Mesh data sets [6] and the problem of learning quantitative structure-activity relations (QSAR) [12] [13]) are concerned with the prediction of numerical values from examples and relational background knowledge. This kind of learning problem is called Relational Regression in [7], and can be formulated in the “normal” ILP framework (i.e., it is not part of the non-monotonic ILP framework which includes the closed-world assumption). Nevertheless, Relational Regression differs from other ILP learning tasks in that there are no negative examples. So far, two methods have been applied to this problem: FORS [15] builds a first-order theory by a covering algorithm. The other approach, a combination of DINUS [16] and RETIS [14], transforms the learning problem into a propositional language, and subsequently applies a regression tree algorithm to the transformed problem [10]. The transformation, however, does not work for non-determinate background knowledge, which is a strict limitation of the approach.

In this paper we present Structural Regression Trees (SRT), a new algorithm for predicting numerical values from examples and relational background knowledge. SRT differs from FORS in its use of divide-and-conquer instead of separate-and-conquer (i.e., it works like a decision tree algorithm and not like a covering algorithm). In contrast to DINUS/RETIS, SRT solves the problem in its original representation, and does not require transforming the problem. Moreover, SRT can utilize non-determinate background knowledge.

To simplify the presentation, we first review work in statistics and machine learning that is related to our approach. In the third section we will describe the method, including the solution for the problem of non-determinate literals. Furthermore, we present a new method for detecting outliers by analogy. Subsequently, we discuss results of experiments in several real-world domains. Finally, we draw our conclusions, and sketch possible directions of further research.

## 2 Related Work

The classical statistical model for the prediction of numerical values is linear least-squares regression. Refinements and extensions like non-linear models are also well-known and used in many real-world applications. However, regression models have several limitations: First of all, regression models are often hard to understand. Secondly, classical statistical methods assume that all features are equally relevant for all parts of the instance space. Thirdly, regression models do not allow for easy utilization of domain knowledge. The only way to include knowledge is to “engineer” features, and to map these symbolic features to real-

valued features.

In order to solve some of these problems, regression tree methods (CART [3], RETIS [14], M5 [21]) have been developed. Regression trees are supposed to be more comprehensible than traditional regression models. Furthermore, regression trees by definition partition the instance space, so the features may be of different importance for different parts of the space. The basic idea of regression trees according to CART is to minimize the least squared error for the next split of a node in the tree, and to predict the average of the dependent variable of all covered training instances for unseen instances in a leaf. RETIS and M5 differ in that they do not assign single values to the leaves, but linear regression models.

Sophisticated post-pruning methods have been developed for CART, since the method grows the tree until every leaf is “pure”, i.e. the leaf covers exactly one instance. The regression tree resulting from the growing phase is usually bigger than a classification tree, since it takes more nodes to achieve pure leaves.

Manago’s KATE [17] learns decision trees from examples represented in a frame-based language that is equivalent to first-order predicate calculus. KATE makes extensive use of a given hierarchy and heuristics to generate the branch tests. To our knowledge, KATE was the first system to induce first-order theories in a divide-and-conquer fashion.

Watanabe and Rendell [29] also investigated the use of divide-and-conquer for learning first-order theories. Although their so-called structural decision trees are used for the prediction of categorical classes and not continuous classes, it is the closest work found in the literature.

## 3 Description of the Method

### 3.1 Overview

SRT is an algorithm which learns a theory for the prediction of numerical values from examples and relational (and even non-determinate) background knowledge. The algorithm constructs a tree containing a literal (an atomic formula or its negation) or a conjunction of literals in each node, and assigns a numerical value to each leaf.

More precisely, SRT generates a series of increasingly complex trees, and subsequently returns one of the generated trees according to a preference criterion. SRT’s preference criterion is based on *minimum description length* (MDL) principle [24]. In such a way, we try to avoid overfitting the data in the presence of noise.

For the construction of a single tree, SRT uses the same method as used for the usual top-down induction of decision trees [22]. The algorithm recursively builds a binary tree, selecting a literal or a conjunction of literals (as defined by user-defined schemata [26]) in each node of the tree until a stopping criterion is

fulfilled. With each selected literal or conjunction, the examples covered by a node are further partitioned, depending on the success or failure of the literal(s) on the example.

The selection of the literal or conjunction is performed as follows: Let  $I$  be the set of training instances covered by a leaf  $l$  in a partial tree, and  $c$  be the conjunction of all literals in the path from the root of the tree to  $l$ . (For a definition of all used terms see table 1.) Then every possible test  $t$  is evaluated according to the resulting partitioning of the training instances  $I$  in  $l$ . The instances  $I$  are partitioned into the instances  $I_1 \subseteq I$  for which proving  $c \wedge t$  succeeds, and into the instances  $I_2 \subseteq I$  for which proving  $c \wedge t$  fails. For every possible test  $t$  we calculate the sum of the squared differences between the actual values  $y_{i,j}$  of the training instances and the average  $\bar{y}_i$  of  $I_i$ . From all possible tests, SRT selects  $t^* \in T$  which minimizes this sum of squared differences (see equation 1). When the stopping criterion is fulfilled, the average  $\bar{y}_i$  is assigned to the leaf as the predicted value for unseen cases that reach the leaf.

$$\text{Sum Squared Errors} = \sum_{i=1}^2 \sum_{j=1}^{n_i} (y_{i,j} - \bar{y}_i)^2 \quad (1)$$

---

$I$	set of instances covered by a leaf $l$ in a partial tree
$c$	the conjunction of all literals in the path from the root of the tree to $l$
$I_1$	subset of $I$ for which proving $c \wedge t$ ( $t \in T$ , the set of possible tests) succeeds
$I_2$	subset of $I$ for which proving $c \wedge t$ fails ( $I = I_1 \cup I_2$ , $I_1 \cap I_2 = \emptyset$ )
$n_1$	number of instances in $I_1$ ( $n_1 =  I_1 $ )
$n_2$	number of instances in $I_2$ ( $n_2 =  I_2 $ )
$y_{1,j}$	value of the dependent variable of a training instance $j$ in $I_1$
$y_{2,j}$	value of the dependent variable of a training instance $j$ in $I_2$
$\bar{y}_1$	average of all instances in $I_1$
$\bar{y}_2$	average of all instances in $I_2$

---

Table 1: Definition of terms

From another point of view, each path starting from the root can be seen as a clause. Every time the tree is extended by a further literal or conjunction, two further clauses are generated: One of them is the current clause (i.e. the path in the current partial tree) extended by the respective literal or conjunction of literals. The other clause is the current clause extended by the negation of the literal(s). Table 2 shows a simple example of a structural regression tree in clausal form. The three clauses predict the biological activity of a compound from its structure and its characteristics. Depending on the conditions in the clauses, the

theory assigns either 8.273 or 6.844 or 6.176 to every unseen instance. In the following, we will use this latter, clausal view on the process.

---

<code>activity(Drug,8.273)</code>	<code>:-</code>	<code>struct(Drug,Group1,Group2,Group3),</code> <code>(pi_doner(Group3,X), X&lt;2).</code>
<code>activity(Drug,6.844)</code>	<code>:-</code>	<code>struct(Drug,Group1,Group2,Group3),</code> <code>\+ (pi_doner(Group3,X), X&lt;2),</code> <code>h_doner(Group1,0).</code>
<code>activity(Drug,6.176)</code>	<code>:-</code>	<code>struct(Drug,Group1,Group2,Group3),</code> <code>\+ (pi_doner(Group3,X), X&lt;2),</code> <code>\+ h_doner(Group1,0).</code>

---

Table 2: Example of a structural regression tree in clausal form

The simplest possible stopping criterion is used to decide if we should further grow a tree: SRT stops extending a clause, when no literal(s) can produce two further clauses that both cover more than a required minimum number of training instances. In the following this parameter will be called the *minimum coverage* of all clauses in the theory. Apart from its use as a stopping criterion, the minimum coverage parameter has the following benefits: We have direct control over the complexity of the trees being built. The smaller the value of the parameter, the more complex the tree will be, since we allow for more specific clauses in the tree. In such a way we can generate a series of increasingly complex trees, and return the one which optimizes a preference function. Furthermore, this solution prevents splits that are more asymmetric than the parameter allows: To a certain degree, the minimum coverage controls how balanced the resulting trees will be.

SRT generates a series of increasingly complex trees by varying the minimum coverage parameter. The algorithm starts with a high minimum coverage, and decreases it from iteration to iteration. Fortunately, many iterations can be skipped, since nothing would change for certain values of the minimum coverage parameter: From those literals and conjunctions that produce two clauses with an admissible coverage we select the one which yields the lowest sum of squared errors. There could be literals or conjunctions yielding an even lower sum of squared errors, but with a coverage that is too low. The maximum coverage of these literals or conjunctions is the next value of the parameter, for which the tree would be different from the current tree. So we choose this value as the next minimum coverage.

Finally, SRT returns the one tree from this series that obtains the best compression of the data. The compression measure is based on the minimum description length (MDL) principle [24], and will be discussed in the next section.

## 3.2 Tree Selection by MDL

The MDL principle tries to measure both the simplicity and the accuracy of a particular theory in a common currency, namely in terms of the number of bits needed for encoding theory and data. [4] defines the message length of a theory (called *model* in his article) as:

$$\begin{aligned} \text{Total message length} = & \\ & \text{Message length to describe the model} + \\ & \text{Message length to describe the data,} \\ & \text{given the model.} \end{aligned}$$

This way a more complex theory will need more bits to be encoded, but might save bits when encoding more data correctly.

The message length of the model consists of the encoding of the literals and the encoding of the predicted values in the leaves. The message length of the data, given the model, is simply the encoding of the errors.

The predicted values and the errors are real numbers, so we have to devise a suitable coding scheme for reals. In our coding scheme, we turn them into integers by multiplication and rounding. The factor is the minimum integer that still allows to discern the values in the training data after rounding. Subsequently, the integers are encoded by the *universal prior of integers* (UPI) [25] — in this way the coding length of the numbers roughly corresponds to their magnitude.

The encoding of the tree is simply the encoding of the choices made (for the respective literals) as the tree is built. For a single node, we encode the choice from all possible literals, so that the encoding considers predicates as well as all possible variabilizations of the predicates.

We chose MDL instead of cross-validation, since it is computationally less expensive, and it can be used for pruning in search [19]. However, we are planning to compare both methods for model selection in the future.

## 3.3 Non-Determinate Background Knowledge

Literals are non-determinate if they introduce new variables that can be bound in several alternative ways. Non-determinate literals often introduce additional parts of a structure like adjacent nodes in a graph. (Other examples are "part-of"-predicates.) Clearly, non-determinate literals do not immediately reduce the error when they are added to a clause under construction. Thus, any greedy algorithm without look-ahead would ignore non-determinate literals. The problem is how to introduce non-determinate literals in a controlled manner.

In SRT, the user has to specify which literal(s) may be used to extend a clause. Firstly, the user can define conjunctions of literals that are used for a limited look-ahead. (These user-defined schemata are similar to relational cliches [26]). Furthermore, the user can constrain the set of possible literals depending

on the body of the clause so far. The conditions on the body are arbitrary Prolog clauses, and therefore the user has even more possibilities to define a language than by Antecedent Description Grammars (ADGs) [5]. To further reduce the number of possibilities, the set of literals and conjunctions is also constrained by modes, types of variables, and variable symmetries.

### 3.4 Outlier Detection by Analogy

Test instances that are outliers strongly deteriorate the average performance of learned regression models. Usually we cannot detect if test instances are outliers, because only little information is available for this task. If relational background knowledge is available, however, a lot of information can be utilized to detect, by “analogy”, if test instances are outliers. Intuitively, when a new prediction is made, we check if the test instance is similar to the training instances which are covered by the clause that fires. If we have not seen anything like the test instance before, we should consider a different prediction than the one suggested by the clause which succeeds on the instance. In this case, we interpret the regression tree as defining a hierarchy of clusters. SRT chooses the cluster which is most similar to the test instance, and predicts the average of this cluster for the test instance.

To implement this kind of reasoning by analogy, we first have to define similarity of “relational structures” (such as labeled graphs).<sup>2</sup> Our simple approximation of similarity is based on *properties* of such structures. In this context, we say that an instance  $i$  has a property  $P$  iff  $P$  is a literal or a conjunction (permitted by specified schemata) that immediately succeeds on  $i$  (i.e., it succeeds without the introduction of intermediate variables). The similarity is defined as follows: Let  $p_{instance}(i)$  denote the set of properties of an instance  $i$ . Let  $p_{in\_common}(I)$  be the set of properties all instances in a set  $I$  have in common. Then the similarity between a test instance  $i$  and a set (cluster) of training instances  $I$  is

$$similarity(i, I) = \frac{|p_{instance}(i) \cap p_{in\_common}(I)|}{|p_{in\_common}(I)|}$$

The similarity is simply defined as the number of properties that the test instance and the covered training instances have in common, divided by the number of properties that the training instances have in common.

SRT uses a parameter for the *minimum similarity* to determine if the similarity between a test instance and the training instances covered by the clause that fires is large enough.

This way of detecting and handling outliers adds an instance-based aspect to SRT. However, it is just an additional possibility, and can be turned off by means of the minimum similarity parameter.

---

<sup>2</sup>[1] defined a similarity measure for first-order logic, but it measures the similarity of two tuples in a relation, not of two “relational structures”.

## 4 Experimental Results

A common step in pharmaceutical development is forming a quantitative structure-activity relationship (QSAR) that relates the structure of a compound to its biological activity. Two QSAR domains, namely the inhibition of *Escherichia coli* dihydrofolate reductase (DHFR) by pyrimidines [12] and by triazines [13] have been used to test SRT.

The pyrimidine dataset consists of 2198 background facts and 55 instances (compounds), which are partitioned into 5 cross-validation sets. For the triazines, the background knowledge are 2933 facts, and 186 instances (compounds) are used to perform 6-fold cross validation. Hirst et al. made comprehensive comparisons of several methods in these domains, but they concluded there is no statistically significant difference between these methods.

For the experiments we set *minimum similarity* = 0.75. Table 3 shows the results of the methods compared in [12] and in [13], and the results of SRT. The table summarizes the test set performances in both domains as measured by the Spearman rank correlation coefficients. The Spearman rank correlation coefficient is a measure of how much the order of the test instances according to the target variable correlates with the order predicted by the induced theory. The only reason why Hirst et al. use the Spearman rank correlation coefficient instead of, say, the average error is to compare GOLEM [18] (which cannot predict numbers) with other methods.<sup>3</sup>

For the pyrimidines, SRT performs better than other methods, but the improvement is not statistically significant. Hirst et al. emphasize that a difference in Spearman rank correlation coefficient of about 0.2 would have been required for a data set of this size. The comparatively good performance of SRT is mostly due to the detection of two outliers that cannot be recognized by other methods. These two outliers were the only ones identified in these two domains. For the triazine dataset, SRT performs quite well, but again the differences are not statistically significant.<sup>4</sup>

Since the Spearman rank correlation coefficient does not measure the quantitative error of a prediction, we included several other measures as proposed by Quinlan [23]. Clearly, these measures have disadvantages, too, but they represent interesting aspects of how well a theory works for a given test set. Unfortunately, we do not yet have a full comparison with other methods that are capable of predicting numbers. Tables 4 and 7 contain the cross-validation test set performances of SRT in four test domains not only in terms of the Spearman rank

---

<sup>3</sup>Despite this disadvantage of GOLEM, Hirst et al. state that GOLEM is the only method that provides understandable rules about drug-receptor interactions. SRT can be seen as a step towards integrating both capabilities.

<sup>4</sup>Note that for the triazines the differences between the standard deviations are very high. Since the hypothesis of equal standard deviations must be rejected even with Bonferroni adjustment, we cannot perform an analysis of variance with the results.



Method	Pyr. Mean ( $\sigma$ )	Triaz. Mean ( $\sigma$ )
Lin. Regr. on Hansch parameters and squares	0.693 (0.170)	0.272 (0.220)
Lin. Regr. on attributes and squares	0.654 (0.104)	0.446 (0.181)
Neur. Netw. on Hansch parameters and squares	0.677 (0.118)	0.377 (0.190)
Neur. Netw. on attributes and squares	0.660 (0.225)	0.481 (0.145)
GOLEM	0.692 (0.077)	0.431 (0.166)
SRT	0.806 (0.110)	0.457 (0.089)

Table 3: Summary of all methods in the biomolecular domains of the inhibition of dihydrofolate reductase by pyrimidines and by triazines: Performances as measured by the Spearman rank correlation coefficients

Measure of Accuracy	Pyr. Mean ( $\sigma$ )	Triaz. Mean ( $\sigma$ )	Mutagen. Mean ( $\sigma$ )
Spearman rank corr. coeff.	0.806 (0.110)	0.457 (0.089)	0.683 (0.124)
Average error $ E $	0.435 (0.088)	0.514 (0.084)	1.103 (0.121)
Correlation $r$	0.818 (0.091)	0.457 (0.104)	0.736 (0.089)
Relative Error $RE$	0.218 (0.170)	0.381 (0.132)	0.170 (0.055)

Table 4: Performances of SRT in three domains in terms of several accuracy measures

correlation coefficient, but also in terms of several other accuracy measures.

Furthermore, we performed experiments in the domain of finite element mesh design (for details see [6]), where the background knowledge is non-determinate. Table 5 shows the results of SRT for the mesh dataset together with the results of FOSSIL [11] and results of other methods that were directly taken from [15]. SRT performs better than FOIL [20] and mFOIL [8], but worse than the other methods. However, statistical analysis shows that only the differences between FOIL and the other algorithms are significant.<sup>5</sup>

Struct.	FOIL	mFOIL	GOLEM	MILP	FOSSIL	FORS	SRT
A	17	23	21	21	23	22	23
B	6	12	12	12	13	12	11
C	7	9	10	11	6	8	9
D	0	6	16	16	16	16	16
E	6	12	21	30	32	29	9
$\Sigma$	36	62	80	90	90	87	68
%	12.9	22.3	28.8	32.4	32.4	31.3	24.4

Table 5: Summary of numbers and percentages of correctly classified edges in the domain of finite element mesh design

<sup>5</sup>For FOIL/GOLEM, FOIL/MILP and FOIL/FORS, the results of both t-test and paired t-test are significant. For FOIL/mFOIL and FOIL/SRT, only the paired t-test shows that the difference is significant. (For FOIL/mFOIL the paired t-test is highly significant, even with Bonferroni adjustment.) For FOIL/FOSSIL, only the t-test shows significance.

We also applied SRT to the biological problem of learning to predict the mutagenic activity of a chemical, i.e., if it is harmful to DNA. (For details see [28] and [27]). This domain involves non-determinate background knowledge, too. In Table 6 we compiled results from [15] and [27], and filled the result of SRT. In the table, 'S' refers to structural background knowledge, 'NS' refers to non-structural features, 'PS' refers to predefined structural features that can be utilized by propositional algorithms, and 'MDL' refers to MDL pre-pruning. (Note that the results of FORS are the best that can be found for its three parameters.) In the experiments we used the 188 instances (compounds) for a 10-fold cross-validation. The accuracy concerns the problem to predict *if* a chemical is active or not. Since SRT learns a theory that predicts the *activity* (a number) instead, we had to evaluate it in a different way to compare the results. Summing up, the experiments showed that also in this domain SRT is competitive, although the differences between SRT and the rest are not statistically significant.

Finally, we applied SRT to a domain where we are trying to predict the half-rate of surface water aerobic aqueous biodegradation in hours [9]. To simplify the learning task, we discretized this quantity and mapped it to  $\{1, 2, 3, 4\}$ . The background knowledge is non-determinate, and except for the molecular weight there are no "global" features available. The dataset contains 62 chemicals, and we performed 6-fold cross-validation in our tests. The results of SRT can be found in table 7. SRT is the first algorithm to be tested on the data, and the results appear to indicate that there are too few instances to find good generalizations. Interestingly, SRT with outlier detection (see column 'Biod.1') improves the initial result of SRT without it in this domain (column 'Biod.2'). Note that neither a propositional algorithm (such as CART) nor an algorithm that cannot handle non-determinate background knowledge (such as FOIL, GOLEM and DINUS/RETIS) can be applied to this problem.

To sum up the experiments, SRT turned out to be quantitatively competitive, but its main advantages are that it yields comprehensible and explicit rules predicting numbers, even when given non-determinate background knowledge.

## 5 Conclusion and Further Research

In this paper we presented Structural Regression Trees (SRT), a new algorithm which can be applied to learning problems concerned with the prediction of numbers from examples and relational (and non-determinate) background knowledge. SRT can be viewed as integrating the statistical method of regression trees [3] into ILP. SRT can be applied to a class of problems no ILP systems except FORS can handle, and provides more comprehensible results than purely statistical methods. The main difference between SRT and FORS is that it is a tree-based and not a covering algorithm. Therefore, all the advantages and disadvantages known

Algorithm	Accuracy
Lin.Regr. + NS	0.85 (0.03)
Lin.Regr. + NS + PS	0.89 (0.02)
Neural Network + NS	0.86 (0.03)
Neural Network + NS + PS	0.89 (0.02)
CART + NS + PS	0.88 (0.02)
CART + NS	0.82 (0.03)
FOIL + NS + S	0.81 (0.03)
Progol + NS + S	0.88 (0.02)
FORS + NS + S	0.89 (0.06)
FORS + NS + S + MDL	0.84 (0.11)
SRT + NS + S	0.85 (0.08)

Table 6: Summary of accuracy of several systems in the mutagenicity domain

Measure of Accuracy	Biod.1 Mean ( $\sigma$ )	Biod.2 Mean ( $\sigma$ )
Spearman rank correlation coefficient	0.463 (0.213)	0.402 (0.232)
Average error $ E $	0.744 (0.190)	0.771 (0.210)
Correlation $r$	0.382 (0.247)	0.364 (0.223)
Relative Error $RE$	0.363 (0.139)	0.377 (0.141)

Table 7: Performances of SRT with and without outlier detection in the biodegradability domain

from other algorithms of these types apply.<sup>6</sup> Experiments in several real-world domains demonstrate that the approach is competitive with existing methods, indicating that its advantages (the applicability to relational regression given non-determinate background knowledge and the comprehensibility of the rules) are not at the expense of predictive accuracy.

SRT generates a series of increasingly complex trees, but currently every iteration starts from scratch. We are planning to extend the algorithm such that parts of the tree of one iteration can be reused in the next iteration.

We also plan to compare our way of coverage-based prepruning and tree selection by MDL with more traditional pruning methods a la CART [3].

Besides, we addressed the problem of non-determinate literals. We adopted and generalized solutions for this problem, but they involve the tiresome task of writing a new specification of admissible literals and conjunctions for each domain. We therefore think that a more generic solution would make the application of the method easier.

One of the current limitations of the approach is that only constants are

---

<sup>6</sup>For a good discussion of tree-based vs. covering algorithms in ILP we have to refer to [2]. For a comparison of tree induction and rule induction in propositional regression see [30].

assigned to the leaves, not linear models as in [14] and [21]. Since it could help to build more accurate models, one of the next steps will be to assign linear regression models to the leaves.

## Acknowledgements

I would like to thank Johannes Fürnkranz, Bernhard Pfahringer and Gerhard Widmer for valuable discussions. I also wish to thank Sašo Džeroski for providing the biodegradability dataset.

## References

- [1] G. Bisson, ‘Learning in FOL with a similarity measure’, in *Proc. Tenth National Conference on Artificial Intelligence (AAAI-92)*, (1992).
- [2] H. Boström, ‘Covering vs. Divide-and-Conquer for Top-Down Induction of logic programs’, in *Proc. Fourteenth International Joint Conference on Artificial Intelligence (IJCAI-95)*, pp. 1194–1200, San Mateo, CA, (1995). Morgan Kaufmann.
- [3] L. Breiman, J.H. Friedman, R.A. Olshen, and C.J Stone, *Classification and Regression Trees*, The Wadsworth Statistics/Probability Series, Wadsworth International Group, Belmont, CA, 1984.
- [4] P. Cheeseman, ‘On finding the most probable model’, in *Computational Models of Discovery and Theory Formation*, eds., J. Shrager and P. Langley, Morgan Kaufmann, Los Altos, CA, (1990).
- [5] W.W. Cohen, ‘Grammatically biased learning: Learning logic programs using an explicit antecedent description language’, *Artificial Intelligence*, **68**(2), (1994).
- [6] B. Dolsak, I. Bratko, and A. Jezernik, ‘Finite element mesh design: An engineering domain for ILP application’, in *Proceedings of the Fourth International Workshop on Inductive Logic Programming (ILP-94)*, GMD-Studien Nr. 237, pp. 305–320, (1994).
- [7] S. Džeroski, *Numerical Constraints and Learnability in Inductive Logic Programming*, Ph.D. dissertation, University of Ljubljana, Ljubljana, Slovenija, 1995.
- [8] S. Džeroski and I. Bratko, ‘Handling noise in Inductive Logic Programming’, in *Proceedings of the International Workshop on Inductive Logic Programming*, Tokyo, Japan, (1992).

- [9] S. Džeroski and B. Kompare, 1995. Personal Communication.
- [10] S. Džeroski, L. Todoroski, and T. Urbancic, ‘Handling real numbers in inductive logic programming: A step towards better behavioural clones’, in *Machine Learning: ECML-95*, eds., N. Lavrac and S. Wrobel, pp. 283–286, Berlin Heidelberg New York, (1995). Springer.
- [11] J. Fürnkranz, ‘FOSSIL: A robust relational learner’, in *Machine Learning: ECML-94*, eds., F. Bergadano and L. De Raedt, pp. 122–137, Berlin Heidelberg New York, (1994). Springer.
- [12] J.D. Hirst, R.D. King, and M.J.E. Sternberg, ‘Quantitative structure-activity relationships by neural networks and inductive logic programming. the inhibition of dihydrofolate reductase by pyrimidines’, *Journal of Computer-Aided Molecular Design*, **8**, 405–420, (1994).
- [13] J.D. Hirst, R.D. King, and M.J.E. Sternberg, ‘Quantitative structure-activity relationships by neural networks and inductive logic programming: The inhibition of dihydrofolate reductase by triazines’, *Journal of Computer-Aided Molecular Design*, **8**, 421–432, (1994).
- [14] A. Karalic, ‘Employing linear regression in regression tree leaves’, in *Proc. Tenth European Conference on Artificial Intelligence (ECAI-92)*, ed., B. Neumann, pp. 440–441, Chichester, UK, (1992). Wiley.
- [15] A. Karalic, *First Order Regression*, Ph.D. dissertation, University of Ljubljana, Ljubljana, Slovenija, 1995.
- [16] N. Lavrac and S. Džeroski, *Inductive Logic Programming*, Ellis Horwood, Chichester, UK, 1994.
- [17] M. Manago, ‘Knowledge-intensive induction’, in *Proceedings of the Sixth International Workshop on Machine Learning*, ed., A.M. Segre, pp. 151–155. Morgan Kaufman, (1989).
- [18] S. Muggleton and C. Feng, ‘Efficient induction of logic programs’, in *Inductive Logic Programming*, ed., S. Muggleton, 281–298, Academic Press, London, U.K., (1992).
- [19] B. Pfahringer and S. Kramer, ‘Compression-based evaluation of partial determinations’, in *Proceedings of the First International Conference on Knowledge Discovery and Data Mining*. AAAI Press, (1995).
- [20] J.R. Quinlan, ‘Learning logical definitions from relations’, *Machine Learning*, **5**, 239–266, (1990).

- [21] J.R. Quinlan, ‘Learning with continuous classes’, in *Proceedings AI’92*, ed., Sterling Adams, pp. 343–348, Singapore, (1992). World Scientific.
- [22] J.R. Quinlan, *C4.5: Programs for Machine Learning*, Morgan Kaufmann, San Mateo, CA, 1993.
- [23] J.R. Quinlan, ‘A case study in machine learning’, in *Proceedings ACSC-16 Sixteenth Australian Computer Science Conference*, (1993).
- [24] J. Rissanen, ‘Modeling by shortest data description’, *Automatica*, **14**, 465–471, (1978).
- [25] J. Rissanen, ‘Stochastic complexity and modeling’, *The Annals of Statistics*, **14**(3), 1080–1100, (1986).
- [26] G. Silverstein and M.J. Pazzani, ‘Relational cliches: Constraining constructive induction during relational learning’, in *Machine Learning: Proceedings of the Eighth International Workshop (ML91)*, eds., L.A. Birnbaum and G.C. Collins, pp. 203–207, San Mateo, CA, (1991). Morgan Kaufmann.
- [27] A. Srinivasan, S. Muggleton, and R.D. King, ‘Comparing the use of background knowledge by Inductive Logic Programming systems’, in *Proceedings of the 5th International Workshop on Inductive Logic Programming (ILP-95)*, pp. 199–230. Katholieke Universiteit Leuven, (1995).
- [28] A. Srinivasan, S. Muggleton, R.D. King, and M. Sternberg, ‘Mutagenesis: ILP experiments in a non-determinate biological domain’, in *Proceedings of the Fourth International Workshop on Inductive Logic Programming (ILP-94)*, GMD-Studien Nr. 237, pp. 217–232, (1994).
- [29] L. Watanabe and L. Rendell, ‘Learning structural decision trees from examples’, in *Proc. Twelfth International Joint Conference on Artificial Intelligence (IJCAI-91)*, pp. 770–776, San Mateo, CA, (1991). Morgan Kaufmann.
- [30] S.M. Weiss and N. Indurkha, ‘Rule-based machine learning methods for functional prediction’, *Journal of Artificial Intelligence Research*, **3**, 383–403, (1995).