

Data Mining in Learning Classifier Systems: Comparing XCS with GAssist

Jaume Bacardit and Martin V. Butz

Enginyeria i Arquitectura La Salle, Universitat Ramon Llull.
Passeig Bonanova 8, 08022, Barcelona.
jbacardit@salleURL.edu

Illinois Genetic Algorithms Laboratory (IlliGAL)
University of Illinois at Urbana-Champaign
61801 Urbana, IL, USA
butz@illigal.ge.uiuc.edu

Abstract

This paper compares performance of the Pittsburgh-style system GAssist with the Michigan-style system XCS on several datamining problems. Our analysis shows that both systems are suitable for datamining but have different advantages and disadvantages. The study does not only reveal important differences between the two systems but also suggests several structural properties of the underlying datasets.

1 Introduction

Successful data mining applications are important for modern-day learning classifier systems (LCSs). Additionally, the study and comparison of different types of data miners on various data sets may enable the identification of strengths and weaknesses of the respective data miners. Several types of problem difficulty can be distinguished in data mining including data volume, search space size and type, complexity of the concept, noise in the data, the handling of missing values, or the problem of over-fitting.

Successful datamining applications of learning classifier systems have been shown in the past (Bernadó, Llorà, & Garrell, 2001) investigating and comparing performance of the accuracy-based Michigan-style LCS XCS (Wilson, 1995) and the Pittsburgh-style LCS GALE (Llorà & Garrell, 2001). Both systems showed competent performance in comparison to six other machine learning systems.

Recently, new systems have appeared in the LCS field, like the Pitt-style LCS GAssist (Bacardit & Garrell, 2003a). Also, there are improved versions of already established systems, like the XCS with tournament selection (Butz, Sastry, & Goldberg, 2003). The objective of this paper is two-fold: (1) We provide further performance results of GAssist and XCS on several interesting datasets. (2) We compare and investigate performance of the two systems revealing problem dependencies, suitability of the respective approaches, as well as over-fitting or over-generalization tendencies.

2 Framework

Before we start with the datamining analysis, this section provides a short introduction to the LCSs under investigation.

2.1 GAssist

GAssist (Bacardit & Garrell, 2003a) is a Pittsburgh genetic-based machine learning system descendant of *GABIL* (DeJong, Spears, & Gordon, 1993). The system applies a near-standard *GA* that evolves individuals that represent complete problem solutions. An individual consists of an ordered, variable-length rule set. Bloat control is achieved by a combination of a fitness function based on the minimum description length (*MDL*) principle and a rule deletion operator (Bacardit & Garrell, 2003a).

The knowledge representation used for real-valued attributes is called *adaptive discretization intervals* rule representation (*ADI*) (Bacardit & Garrell, 2003b). This representation uses the semantics of the *GABIL* rules (conjunctive normal form predicates), but applies non-static intervals formed by joining several neighbor discretization intervals. These intervals can evolve through the learning process splitting or merging among them potentially using several discretizers at the same time.

The system also uses a windowing scheme called *ILAS* (incremental learning with alternating strata) (Bacardit & Garrell, 2003c). This scheme stratifies the training set into s subsets of equal size and approximately uniform class distribution. Each *GA* iteration uses a different strata to perform its fitness computation, using a round-robin policy. This method showed to introduce an additional implicit generalization pressure to GAssist. ¹

2.2 XCS

The XCS classifier system (Wilson, 1995; Wilson, 1999) evolves online a set of condition-action rules, that is, a *population of classifiers*. In difference to GAssist, in XCS the population as a whole represents the problem solution. XCS differs in two fundamental ways to other Michigan-style LCSs: (1) Rule fitness is derived from rule accuracy instead of rule reward prediction. (2) *GA* selection is applied in the subsets of currently active classifiers resulting in an implicit pressure towards more general rules.

Due to the variable properties of the investigated datasets including real values, nominals, and binary features, we use a hybrid XCS/XCSR approach that can handle any feature combination as done before in (Bernadó, Llorà, & Garrell, 2001). Additionally, we apply tournament selection which proved to result in more robust fitness pressure toward accurate rules (Butz, Sastry, & Goldberg, 2003). ²

¹GAssist’s parameters were set as follows: Crossover probability 0.6; tournament selection; tournament size 3; population size 400; probability of mutating an individual 0.6; initial number of rules per individual 20; probability of “1” in initialization 0.75; Rule Deletion Operator: Iteration of activation: 5; minimum number of rules: number of classes of domain +3; MDL-based fitness function: Iteration of activation 25; initial theory length ratio: 0.075; weight relax factor: 0.9. ADI knowledge representation: split and merge probability: 0.05; reinitialize probability at initial iteration: 0.02; reinitialize probability at final iteration: 0; merge restriction probability: 0.5; maximum number of intervals: 5; set of uniform discretizers used: 4, 5, 6, 7, 8, 10, 15, 20 and 25 bins; iterations: maximum of 1500. Results are averaged over 150 experiments.

²XCS’s parameters are set as follows: $N = 6400$, $r_0 = 4(100)$, $P_{\#} = 0.6$, $\beta = 0.2$, $\chi = 1.0$ applying uniform crossover, $\mu = 0.04$, $m_0 = 0.2$, $\theta_{GA} = 48$, $\tau = 0.4$, $\varepsilon_0 = 1$, $\delta = 0.1$, $\theta_{del} = 50$, *GA* Subsumption is applied with $\theta_{sub} = 50$. Experiments are run applying either 100,000 learning steps (averaging over 150 experiments) or 500,000

3 Experiments

3.1 Tests setup

In Table 1 we show the datasets we have selected from the University of California at Irvine (UCI) repository (Blake, Keogh, & Merz, 1998). The selected datasets are:

- Annealing Data (*ann*)
- 1985 Auto Imports Database (*aut*)
- Balance Scale Weight & Distance (*bal*)
- Contraceptive Method Choice (*cmc*)
- Horse Colic (*col*)
- German Credit (*cr-g*)
- Glass Identification (*gls*)
- Cleveland Heart Disease (*h-c*)
- Hungarian Heart Disease (*h-h*)
- Johns Hopkins University Ionosphere database (*ion*)
- Sonar, Mines vs. Rocks database (*son*)
- Wisconsin Breast Cancer database (*wbcd*)
- Wisconsin Diagnostic Breast Cancer (*wdbc*)

The selection of datasets gives a representative overview over the phenomena we were able to detect using GAssist and XCS.

The test design for GAssist has two goals: Comparing the effect of using both different number of iterations and different degrees of generalization pressure. The latter goal is achieved by using the *ILAS* windowing scheme. However, our goal here is not run-time reduction, but maximizing as much as possible the generalization pressure introduced by *ILAS*. Thus, we will increase the number of iterations when using windowing proportional to the number of strates used. This means having constant number of learning steps (using the Michigan-LCS meaning of the term). We will also test another stratified setup using a number of iterations that makes it equivalent in run-time as the non-windowed setting.

3.2 Results

Results of GAssist and XCS are shown in Table 2. The comparison is not meant to determine which system is better in general but rather to show in which problem types which system appears to have advantages. Our comparison starts with a general data observation and then investigates separate datasets with respect to specific phenomena.

A look at the overall performance shows that XCS and GAssist show comparative performance results indicating the general difficulty of the respective datasets. XCS tends to learn the training

learning steps (averaging over 20 experiments).

Table 1: The dataset properties indicate complexity, size, and data distributions in the respective datasets. #Inst. = Number of Instances, #Attr. = Number of attributes, #Real = Number of real-valued attributes, #Nom. = Number of nominal attributes, #Cla. = Number of classes, Dev.cla. = Deviation of class distribution, Maj.cla. = Percentage of instances belonging to the majority class, Min.cla. = Percentage of instances belonging to the minority class, MV Inst. = Percentage of instance with missing values, MV Attr. = Number of attributes with missing values, MV values = Percentage of values ($\#instances \cdot \#attr$) with missing values

Dataset Properties											
Name	#Inst.	#Attr.	#Real	#Nom.	#Cla.	Dev.cla.	Maj.cla.	Min.cla.	MV Inst.	MV Attr.	MV values
ann	898	38	6	32	5	28.28%	76.17%	0.89%	—	—	—
aut	205	25	15	10	6	10.25%	32.68%	1.46%	22.44%	7	1.11%
bal	625	4	4	—	3	18.03%	46.08%	7.84%	—	—	—
cmc	1473	9	2	7	3	8.26%	42.70%	22.61%	—	—	—
col	368	22	7	15	2	13.04%	63.04%	36.96%	98.10%	21	22.77%
cr-g	1000	20	8	12	2	20.00%	70.00%	30.00%	—	—	—
glc	214	9	9	—	6	12.69%	35.51%	4.21%	—	—	—
h-c1	303	13	6	7	2	4.46%	54.46%	45.54%	2.31%	2	0.17%
h-h	294	13	6	7	2	13.95%	63.95%	36.05%	99.66%	9	19.00%
ion	351	34	34	—	2	14.10%	64.10%	35.90%	—	—	—
son	208	60	60	—	2	3.37%	53.37%	46.63%	—	—	—
wbcd	699	9	9	—	2	15.52%	65.52%	34.48%	2.29%	1	0.23%
wdbc	569	30	30	—	2	12.74%	62.74%	37.26%	—	—	—

data much more precise which however is not necessarily advantageous for performance on the test data (using stratified ten-fold cross-validation). The solution representation differs (as expected) very significantly between GAssist and XCS: The number of rules in the best individual in GAssist is much smaller than the number of rules in XCS. However, it should be noted that GAssist maintains 400 individuals and thus the overall number of rules is actually similar to the number of rules in XCS. While we did not make explicit speed comparisons it appears that XCS runs take longer than GAssist’s. Again, this is expectable since XCS is an online learner that learns from each problem instance separately and iteratively. Thus, the number of necessary learning iterations are higher.

Taking a closer look at the particular datasets we see that in the anneal (ann) dataset, performance of both systems reaches a similar level if XCS is run long enough. As also indicated by XCS’s smaller population size in longer runs, generalization appears important and requires sufficient learning time. Generalization is even more important in the autos (aut) problem indicated by XCS’s poor performance when starting specific and its improved test performance and smaller population size in longer runs as well as in GAssist’s slight performance improvement and rule number decrease when using three strata. Additionally, the higher population size of XCS compared to the anneal problem indicates a general higher complexity of the problem. Balance-scale (bal) is a typical problem which can be over-fitted easily: XCS’s performance is worse when starting more specific and when performing longer runs. Note that the population size of XCS actually increases when starting general and running more iterations—a clear indication of over-fitting. GAssist’s performance points in the same direction in that generalization can slightly improve performance but longer runs are not helpful. The cmc problem appears to be a tough problem in general. XCS over-fits the data more than GAssist showing higher train performance but worse test performance. In the colic (col) as well as in the heart-h (h-h) problem, performance of XCS is significantly worse compared to GAssist. The major reason for this appears to be the missing value policy. While in GAssist a missing value is replaced by the majority value for the nominal case or by the average value in the real-valued case, XCS assumes a match in the missing value case. The latter strategy

Table 2: Performance results of GAssist and XCS show train and test performance using 10-folded cross-validation. Additionally, we show the number of rules in the best individual of GAssist and the number of (macro-)classifiers in XCS (at the end of a run). The different GAssist runs distinguish a different application of strata as well as number of iterations (609,1827, and 1447, respectively). In XCS, we compare long and short learning runs as well as a general and specific initialization of the population.

Data	Res.	GAssist			XCS (500,000)		XCS (100,000)	
		1 strata	3 str.(steps)	3 str.(time)	$r_0 = 100$	$r_0 = 4$	$r_0 = 100$	$r_0 = 4$
ann	Train	97.44±2.23	97.80±3.27	97.89±2.51	99.56±0.46	99.95±0.18	94.26±1.97	98.88±0.61
	Test	97.03±2.55	97.40±3.45	97.47±2.80	98.38±1.57	98.56±1.49	91.22±2.70	91.73±2.93
	#rules	6.9±0.9	6.3±0.7	6.3±0.5	2507.1±232.0	3210.8±146.0	4440.4± 86.8	5425.9± 51.4
aut	Train	85.54±2.93	84.66±3.16	82.82±3.73	99.76±0.23	99.64±0.39	99.27±0.67	99.38±0.56
	Test	67.54±9.82	68.79±9.66	67.50±9.46	71.54±9.47	68.83±12.05	64.69±9.64	13.40±6.91
	#rules	12.8±2.7	7.8±1.1	7.8±1.0	3403.2± 98.5	4679.1±216.7	4281.1± 87.3	5426.2± 36.9
bal	Train	87.67±0.49	85.97±0.69	85.93±0.73	98.37±0.72	98.59±0.64	90.63±2.18	97.96±0.86
	Test	78.98±4.22	78.80±3.76	79.17±4.38	81.37±3.58	80.97±3.82	84.57±3.27	81.96±3.49
	#rules	13.1±2.0	9.6±1.6	9.8±1.6	2060.8± 73.2	2013.8± 59.8	1611.2±168.5	2465.2± 65.9
cmc	Train	59.77±0.96	59.55±1.13	59.75±1.09	70.46±1.86	77.56±2.01	57.04±1.84	71.48±2.23
	Test	54.78±4.18	54.60±4.00	54.90±4.11	53.63±4.02	52.89±4.71	50.12±4.67	53.59±3.56
	#rules	7.7±1.4	9.3±3.0	9.1±2.9	3261.3± 88.1	3210.1± 84.3	3957.8± 91.4	3929.2± 64.7
col	Train	99.72±0.34	99.56±0.48	99.54±0.50	94.58±1.18	95.52±1.32	91.67±1.64	94.96±1.12
	Test	93.00±4.67	93.77±4.57	94.06±4.31	84.35±5.03	83.68±5.80	84.46±5.83	84.82±5.55
	#rules	7.4±1.6	7.0±1.4	7.0±1.4	3102.1±156.1	3685.1± 84.2	3612.4±168.8	4099.5± 96.3
cr-g	Train	81.95±0.76	83.72±0.94	84.32±0.83	98.24±1.19	99.61±0.34	89.69±3.15	94.42±1.39
	Test	72.30±3.61	72.03±4.21	72.20±3.78	70.15±3.63	72.30±4.16	71.39±3.85	72.45±3.11
	#rules	6.8±1.5	11.3±3.0	13.1±2.1	2015.6± 69.2	2622.9± 75.4	3217.0±105.8	4401.2±103.8
gls	Train	82.14±1.81	80.41±1.89	79.88±1.84	98.84±0.64	99.57±0.67	89.67±2.83	96.62±1.43
	Test	68.18±9.32	69.40±9.16	68.39±9.89	74.68±7.71	71.20±8.69	70.65±8.15	70.71±8.43
	#rules	8.8±1.4	6.6±0.8	6.6±0.8	1808.4± 86.5	2142.9± 78.4	3092.6±133.6	3137.1± 92.7
h-cl	Train	93.42±0.82	91.44±0.98	92.64±0.86	99.85±0.25	100.00±0.00	99.48±0.46	100.00±0.00
	Test	80.18±7.04	79.96±6.84	80.28±6.48	76.40±6.69	79.58±6.52	77.67±6.80	68.90±8.60
	#rules	9.3±1.5	6.9±1.1	7.4±1.2	2042.7± 69.2	2807.8± 89.6	2854.1± 99.6	2906.5± 68.2
h-h	Train	99.69±0.32	99.04±0.48	99.02±0.50	99.65±0.44	100.00±0.00	95.38±2.22	100.00±0.00
	Test	95.53±4.40	95.65±4.39	95.79±3.25	78.65±9.02	76.59±6.90	79.44±7.69	70.82±6.93
	#rules	6.1±0.7	6.3±0.5	6.0±0.2	2072.3±103.1	2686.2± 71.9	3090.8±135.8	2860.8± 67.5
ion	Train	98.24±0.46	96.77±0.63	96.78±0.59	99.94±0.19	99.68±0.41	99.71±0.32	99.78±0.34
	Test	92.51±4.93	92.70±4.74	92.97±4.83	89.33±4.81	57.39±6.39	90.73±5.29	57.07±6.81
	#rules	3.9±0.8	2.2±0.7	2.2±0.8	2934.6± 93.5	5613.1± 28.9	3479.2± 97.6	5685.4± 31.4
son	Train	97.00±0.96	96.62±1.19	96.25±1.17	100.00±0.00	100.00±0.00	99.89±0.30	100.00±0.00
	Test	74.35±8.89	76.81±9.00	77.47±9.19	78.35±7.42	82.61±8.28	77.27±8.08	81.58±7.88
	#rules	8.3±1.4	6.8±1.1	6.9±1.1	4958.9±119.8	4168.0±142.2	5148.2±107.0	4472.6± 89.8
wbcd	Train	99.05±0.27	97.82±0.50	97.85±0.47	99.84±0.24	100.00±0.00	97.68±0.89	99.94±0.13
	Test	95.15±2.93	96.05±2.59	96.04±2.37	96.06±2.83	96.22±2.23	96.19±2.18	96.45±1.90
	#rules	5.0±1.0	2.4±0.6	2.4±0.6	1562.3± 96.8	2131.1± 52.9	1107.9±143.5	3137.3± 81.8
wdbc	Train	98.60±0.45	97.59±0.68	97.57±0.78	99.98±0.09	100.00±0.00	99.84±0.22	99.85±0.24
	Test	94.06±3.01	94.18±2.91	94.10±2.84	96.13±2.48	96.67±2.20	95.85±2.62	92.90±3.29
	#rules	6.0±1.3	3.8±0.7	3.9±0.9	4104.1±111.5	5050.8± 50.9	4484.5± 85.5	5551.2± 87.7

appears mediocre in the investigated data mining experiments explaining XCS’s poor performance in these settings.

Performance in the credit-g problem (cr-g) indicates that over-fitting is unlikely but in order to reach higher performance more specific initialization is helpful. Again, XCS reaches a much higher train performance but test performance is hardly influenced.

XCS’s behavior in the glass problem (gls) is similar to that of credit-g. However, generalization is more important as also indicated by the performance improvement in GAssist when using three strata. Similar to the autos problem, XCS outperforms GAssist in the glass problem indicating higher problem complexity which might partially stem from the large number of classes in the problem.

XCS’s performance in heart-c1 (h-c1) is actually very similar to the performance in heart-h (h-h) suggesting that besides the problem of missing values in heart-h, XCS tends to strongly over-fit the training data. GAssist does not suffer from this problem in these datasets.

Another interesting observation was made in the ionosphere problem (ion) in which the automatic default rule detection mechanism in GAssist is actually able to discover that the minority class results in a better problem performance. XCS tends to over-fit as indicated by the poor performance and large population size when starting too specific.

On the other hand, in the sonar problem (son) a start from the specific side is actually beneficial for XCS suggesting small special-case niches which can be separated only if the population is initialized more specific. The more generalized representation of GAssist is not advantageous in this dataset.

In the Wisconsin breast-cancer dataset (wdbc) performance of both systems is similar and the problem appears to be generally easy as indicated by the small number of rules in both systems.

Finally, wdbc is another problem in which the complexity of the problem makes it hard for GAssist to reach XCS’s performance level. XCS needs a large number of classifiers to solve the problem but is able to evolve the appropriate set. Slight generalizations are possible. GAssist on the other hand learns a very general—but slightly over-general solution.

4 Summary and Conclusions

In sum, both LCS systems showed that they are suitable for data-mining applications developing very different problem solutions that nonetheless perform similarly well on the test sets. Additionally, the comparison showed that regardless of offline (GAssist) or online (XCS) learning, LCSs are suitable data-miners.

The results allowed us to infer problem properties as well as problem difficulties. We saw that the current policy of handling missing values in XCS can affect performance negatively. Also, while GAssist has the tendency to ignore additional problem complexity, XCS tends to over-fit the training data more often (dependent on the nature of the data). Additionally, GAssist has slight problems with handling many output classes as well as a huge search space suggesting the addition of special covering operators that could ensure that each individual in GAssist differentiates at least all classes in the problem at hand. On the other hand, XCS’s generalization tendency needs to be revisited in the data-mining domain. Especially in smaller datasets, XCS clearly tends to over-fit the data. Due to the small size of the datasets, the natural generalization pressure due to the niche reproduction mechanism hardly applies and pressure towards syntactic generality appears to become more important.

The insights gained from our study prepare the systems for a more general problem application suggesting initial testing with each learning approach for suitability and appropriate initialization.

XCS may need to be improved in terms of generalization to avoid over-fitting. GAssist may be endowed with further covering mechanism to ensure that all problem classes are covered by each individual and that it is able to detect additional small but significant problem subspaces.

Acknowledgments

The authors acknowledge the support provided by the Spanish Research Agency (CICYT) under grant numbers TIC2002-04160-C02-02 and TIC 2002-04036-C05-03, the support provided by the Department of Universities, Research and Information Society (DURSI) of the Autonomous Government of Catalonia under grants 2002SGR 00155 and 2001FI 00514.

The work was sponsored by the Air Force Office of Scientific Research, Air Force Materiel Command, USAF, under grant F49620-03-1-0129. The US Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation thereon. Additional funding from the German research foundation (DFG) under grant DFG HO1301/4-3 is acknowledged. Additional support from the Computational Science and Engineering graduate option program (CSE) at the University of Illinois at Urbana-Champaign is acknowledged.

The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the Air Force Office of Scientific Research or the U.S. Government.

References

- Bacardit, J., & Garrell, J. M. (2003a). Bloat control and generalization pressure using the minimum description length principle for a pittsburgh approach learning classifier system. In *Proceedings of the 6th International Workshop on Learning Classifier Systems* (in press), LNAI, Springer.
- Bacardit, J., & Garrell, J. M. (2003b). Evolving multiple discretizations with adaptive intervals for a pittsburgh rule-based learning classifier system. In *Proceedings of the Genetic and Evolutionary Computation Conference - GECCO2003* pp. 1818–1831. LNCS 2724, Springer.
- Bacardit, J., & Garrell, J. M. (2003c). Incremental learning for pittsburgh approach classifier systems. In *Proceedings of the “Segundo Congreso Español de Metaheurísticas, Algoritmos Evolutivos y Bioinspirados.”* pp. 303–311.
- Bernadó, E., Llorà, X., & Garrell, J. M. (2001). XCS and GALE: a comparative study of two learning classifier systems with six other learning algorithms on classification tasks. In *Fourth International Workshop on Learning Classifier Systems - IW LCS-2001* pp. 337–341.
- Blake, C., Keogh, E., & Merz, C. (1998). UCI repository of machine learning databases. (www.ics.uci.edu/mllearn/MLRepository.html).
- Butz, M. V., Sastry, K., & Goldberg, D. E. (2003). Tournament selection in XCS. *Proceedings of the Fifth Genetic and Evolutionary Computation Conference (GECCO-2003)*, 1857–1869.
- DeJong, K. A., Spears, W. M., & Gordon, D. F. (1993). Using genetic algorithms for concept learning. *Machine Learning*, 13(2/3), 161–188.
- Llorà, X., & Garrell, J. M. (2001). Knowledge-independent data mining with fine-grained parallel evolutionary algorithms. In *Proceedings of the Third Genetic and Evolutionary Computation Conference* pp. 461–468. Morgan Kaufmann.
- Wilson, S. W. (1995). Classifier fitness based on accuracy. *Evolutionary Computation*, 3(2), 149–175.
- Wilson, S. W. (1999). Get real! XCS with continuous-valued inputs. In Booker, L., Forrest, S., Mitchell, M., & Riolo, R. L. (Eds.), *Festschrift in Honor of John H. Holland* pp. 111–121. Center for the Study of Complex Systems.