

## **Mining Oblique Data with XCS**

**Stewart W. Wilson**

IlliGAL Report No. 2000028  
July, 2000

Illinois Genetic Algorithms Laboratory  
University of Illinois at Urbana-Champaign  
117 Transportation Building  
104 S. Mathews Avenue Urbana, IL 61801  
Office: (217) 333-2346  
Fax: (217) 244-5705

# Mining Oblique Data with XCS

**Stewart W. Wilson**

Department of General Engineering  
The University of Illinois, Urbana-Champaign IL 61801, USA  
Prediction Dynamics, Concord, MA 01742 USA  
wilson@prediction-dynamics.com

## Abstract

The classifier system XCS was investigated for data mining applications where the dataset discrimination surface (DS) is generally oblique to the attribute axes. Despite the classifiers' hyper-rectangular predicates, XCS reached 100% performance on synthetic problems with diagonal DS's and, in a train/test experiment, competitive performance on the Wisconsin Breast Cancer dataset. Final classifiers in an extended WBC learning run were interpretable to suggest dependencies on one or a few attributes. For data mining of numeric datasets with partially oblique discrimination surfaces, XCS shows promise from both performance and pattern discovery viewpoints.

## 1 Introduction

Data mining has been described as “the process of discovering patterns in data” (Witten and Frank 2000). The discovered patterns are often represented by rules that, given a data instance, can be used to predict an outcome or consequence of interest. Similarly, XCS (Wilson 1995), a learning classifier system, evolves rules (classifiers) through which the system gradually improves its ability to obtain environmental reward. In effect, XCS *mines its environment* for patterns that, expressed in classifiers, allow it to make increasingly remunerative decisions. XCS is potentially applicable to data mining problems because in many environments it evolves accurate, maximally general rules in which the patterns are easily seen.

Some of the relevant XCS work has involved learning Boolean functions (Wilson 1995; Kovacs 1997; Wilson 1998). From a data-mining point of view, a Boolean function is a very strongly patterned dataset containing all two-valued attribute strings of a given length. The ability to learn Boolean functions indicates XCS's ability to find highly non-linear patterns among attributes. Saxon and Barry (2000) tested XCS on The Monk's Problems (Thrun et al. 1991), a widely used non-Boolean test set for comparing learning algorithms. The domain has six attributes taking up to four nominal values (e.g., the `head_shape` attribute takes values `round`, `square`, and `octagon`). The Monk's 1 problem is a relatively simple concept, (`head_shape = body_shape`) or (`jacket_color = red`). The second problem is a more complicated concept not simply described by a disjunction of conjunctions of attribute-value pairs. Monk's 3 is again a simple concept but contains noise in the form of approximately 5% misclassified examples in the training set. Comparing the results with published results for other systems led Saxon and Barry to conclude that on The Monk's Problems XCS “performed at least as well as traditional Machine Learning techniques”, and that XCS was “also able to produce and maintain an easily identifiable accurately general

set of classifiers representing the concepts within the data sets”. Wilson (2000a) showed that XCS could learn a problem in which there was a Boolean relation among real-valued attributes, showing that XCS could be adapted to take continuous inputs.

Promising as these results are, in all the problems addressed the concept or function sought is essentially “logical”. I.e., it can be expressed by a combination of **and**, **or**, and **not** applied to attribute values, simple recodings of them, or, in the case of Wilson (2000a), ranges of the values. XCS appears well suited to “logical” data-mining problems, including realistically large ones (as suggested by preliminary complexity results in Wilson (1998)). However, in many data-mining problems the function sought is *not* expressible logically. This is because the problem’s discrimination surface (DS) is *oblique*—neither parallel nor perpendicular to—the attribute axes. Datasets with oblique DS’s form a large class and their instances typically have numeric attributes.

XCS is good at logical problems because the classifier conditions define hyper-rectangles, which can evolve to fit the shape of a logical DS exactly. However, if the DS is oblique, hyper-rectangles can only approximate the shape. Many classifiers may be required, and the shape may be difficult to see in the classifiers. One approach to oblique problems is to modify the syntax of classifier conditions so that oblique DS’s can be approximated. The most general method is to use *S-classifiers* (Lanzi 1999; Wilson 1994), where the conditions are Lisp S-expression predicates and can be based on arithmetic function primitives. Our aim in this paper, however, is to apply basic XCS (modified for numeric inputs) to oblique data, using the well-known Wisconsin Breast Cancer dataset (Blake and Merz 1998) as test bed. The WBC dataset offers an opportunity to test XCS’s performance and the interpretability of its evolved rules on a realistic, oblique, data-mining problem.

The next section gives a brief introduction to XCS. Section 3 presents modifications needed for the WBC problem. We warm up on some synthetic data in Section 4. Section 5 presents results on the WBC dataset, including a stratified cross-validation (train/test) experiment, and an experiment that reveals patterns. Sections 6 and 7 contain discussion and conclusion.

## 2 XCS in Brief

The following, drawn from Wilson (2000a), is an abbreviated description of XCS. (For more detail, see Wilson (1995), Kovacs (1997), and Wilson (1998); Butz and Wilson (2000) gives an algorithmic description.) XCS is designed for both single- and multiple-step tasks, but the discussion here applies only to XCS for independent single-step tasks in which an input is presented, the system makes a decision, and the environment provides some reward.

Structurally, each classifier  $C_j$  in XCS’s population [P] has a *condition*, an *action*, and a set of associated parameters. The condition is a string from  $\{0,1,\#\}$ ; the action is an integer. The three principal parameters are: (1) *payoff prediction*  $p_j$ , which estimates the payoff the system will receive if  $C_j$  matches and its action is chosen by the system; (2) *prediction error*  $\epsilon_j$ , which estimates the error in  $p_j$  with respect to actual payoffs received; and (3) *fitness*  $F_j$ , computed as later explained. It is convenient to divide the description of a single operating cycle or time-step into the traditional performance, update (reinforcement), and discovery components.

## 2.1 Performance

Upon presentation of an input—a binary string—XCS forms a *match set* [M] of classifiers in [P] whose conditions are satisfied by the input. The condition is satisfied if and only if each of its non-# positions is the same as the corresponding bit of the input string.

[If no classifiers match, the input is “covered” by creating matching classifiers with each of the possible actions, and placing them in [M]. Each new classifier’s condition is a copy of the input except that #’s (“don’t cares”) are inserted with probability  $P_{\#}$  per position. Parameters are given initial values. If [P] has reached its allowable maximum  $N$ , deletion of classifiers occurs to make room for the new ones (deletion is generally not necessary; initial populations are empty and covering normally occurs only at the very beginning of a run)].

Then, for each action  $a_k$  represented in [M], the system computes a fitness-weighted average  $P_k$  of the predictions  $p_j$  of each classifier in [M] having that action:  $P_k = \sum_j F_j p_j / \sum_j F_j$ .  $P_k$  is termed the *system prediction* for action  $k$ .

Next, XCS chooses an action from those represented in [M] and sends it to the environment. According to the *action-selection* regime in force, the action may be picked randomly or otherwise probabilistically based on the  $P_k$ , or it may be picked deterministically—i.e., the action with the highest  $P_k$  is chosen. Finally, an *action set* [A] is formed consisting of the subset of [M] having the chosen action.

## 2.2 Update

In this component, the parameters of the classifiers in [A] are re-estimated according to the *reward*  $R$  returned by the environment as a consequence of the system’s taking action  $a_k$ . First, the predictions are updated:  $p_j \leftarrow p_j + \beta(R - p_j)$ . Next, the errors:  $\epsilon_j \leftarrow \epsilon_j + \beta(|R - p_j|)$ . Third, for each  $C_j$ , an *accuracy*  $\kappa_j$  is computed:  $\kappa_j = 0.1(\epsilon_j/\epsilon_0)^{-n}$ , for  $\epsilon_j > \epsilon_0$ , else 1.0. Then, from the  $\kappa_j$ , each classifier’s *relative accuracy*  $\kappa'_j$  is computed:  $\kappa'_j = \kappa_j / \sum_j \kappa_j$ . Finally the fitnesses  $F_j$  are updated according to:  $F_j \leftarrow F_j + \beta(\kappa'_j - F_j)$ .

## 2.3 Discovery

On some time-steps, XCS executes a *genetic algorithm* (GA) within [A]. Two classifiers are chosen probabilistically based on their fitnesses and copied. The copies are crossed (two-point crossover) with probability  $\chi$ , and then mutated with probability  $\mu$  per allele. The resulting offspring are inserted into [P]; if the population size is already at its maximum value,  $N$ , two classifiers are deleted. The probability of deletion of a classifier is determined by Kovacs’s (1999) method and is designed to preferentially remove low-fitness classifiers that have participated in a threshold number of action sets—that is, have had sufficient time for their parameters to be accurately estimated.

Whether or not to execute the GA on a given time-step is determined as follows. The system keeps a count of the number of time-steps since the beginning of a run. Every time a GA occurs, the classifiers in that [A] are “time-stamped” with the current count. Whenever an [A] is formed, the time-stamp values of its members are averaged and subtracted from the current count; if the difference exceeds a threshold  $\theta_{GA}$ , a GA takes place.

A *macroclassifier* technique is used to speed processing and provide a more perspicuous view of population contents. Whenever a new classifier is generated by the GA (or covering), [P] is scanned to see if there already exists a classifier with the same condition and action.

If so, the *numerosity* parameter of the existing classifier is incremented by one, and the new classifier is discarded. If not, the new classifier is inserted into [P]. The resulting population consists entirely of structurally unique classifiers, each with numerosity  $\geq 1$ . If a classifier is chosen for deletion, its numerosity is decremented by 1, unless the result would be 0, in which case the classifier is removed from [P]. All operations in a population of macroclassifiers are carried out as though the population consisted of conventional classifiers; that is, the numerosity is taken into account. In a macroclassifier population, the sum of numerosities equals  $N$ , the traditional population size. [P]’s actual size in macroclassifiers,  $M$ , is of interest as a measure of the population’s space complexity.

## 2.4 Overall Picture

XCS’s overall objective is to maximize the rate of reward it obtains from the environment. If the environment (or trainer, teacher, etc.) gives high reward for actions (decisions) that it regards as “correct” or “best”, then XCS’s achieving a high reward rate will be equivalent to learning or solving the problem presented by the environment. This is the framework of *reinforcement learning* (Sutton and Barto 1998). Data mining problems can be cast in terms of reinforcement learning if rewards given the system correlate with the correctness of its decisions.

XCS is designed to evolve classifiers that accurately predict environmental reward in all the situations (inputs) the system might encounter. Having these it can choose the actions that maximize reward. XCS evolves accurate classifiers by a process of trial and error in which, given an input and a match set, a particular action is tried, the reward (if any) is received, and the associated classifiers are updated as previously described. Classifiers that predict accurately tend to be reproduced, while inaccurate classifiers are eventually deleted.

XCS not only evolves accurate classifiers, it evolves ones that are *maximally general*. That is, if several *different* input situations result in the same reward when a particular action is made by the system, then XCS will tend to evolve a *single* (macro)classifier having that action that matches all of those situations. As a result, the system achieves compactness of representation, and the evolved classifiers display the structure of the environment or problem in a way that is often readily interpretable by users. This important tendency of XCS occurs because if there are two equally accurate classifiers but one matches a subset of the states matched by the other, the more general classifier will win out since it occurs in more action sets and will have more reproductive opportunities (Wilson 1995; Kovacs 1997; Wilson 1998; Wilson 2000b).

## 2.5 Subsumption Deletion

The generalization mechanism just mentioned, while powerful, may sometimes not go as far as one would wish. Winning classifiers will be those that match more inputs—environmental *states*—than equally accurate rivals. Usually, a classifier whose condition is formally more general (has more #’s) will match more states. However, in certain environments (termed *sparse*), not all states may actually be present, with the result that some formally more general classifiers, while still perfectly accurate, cannot “drive out” their rivals. The resulting population will be larger than is strictly necessary to solve the problem, and its knowledge will not be as perspicuous. To deal with sparseness and evolve classifiers that are as formally

general as possible without sacrificing accuracy, an optional procedure called “subsumption deletion” was introduced (Wilson 1998; Wilson 2000b; Butz and Wilson 2000).

Two forms of subsumption deletion are used. In *GA subsumption*, a new offspring is examined to see if its condition is logically subsumed by the condition of either of its parents. If so, and if the subsuming parent is both highly accurate and sufficiently experienced (its parameters have been updated a threshold number of times), the offspring is not added to the population but the parent’s numerosity is incremented. *Action-set subsumption* is different from and independent of GA subsumption. Each action set is searched for the most general classifier that is both accurate and sufficiently experienced. Then all other classifiers in the set are tested against the general one to see if it subsumes them. Any that are subsumed are deleted from the population.

### 3 XCSI: XCS Modified for Integer Inputs

The attributes of WBC dataset instances are integer-valued, from 1 to 10 inclusive. To handle this format, XCS was modified in its input interface, the mutation operator, covering, and subsumption. For convenience, we call the modified system “XCSI”. The changes are analogous to those made for real-valued inputs in XCSR (Wilson 2000a). The classifier condition is changed from a string from  $\{0,1,\#\}$  to a concatenation of “interval predicates”,  $int_i = (l_i, u_i)$ , where  $l_i$  (“lower”) and  $u_i$  (“upper”) are integers. A classifier matches an input  $x$  with attributes  $x_i$  if and only if  $l_i \leq x_i \leq u_i$  for all  $x_i$ . A consequence of using interval predicates is that the number of numerical values or alleles in the condition is twice the number of components in  $x$ .

Crossover (two-point) operates in direct analogy to crossover in XCS. A crossover point can occur between any two alleles, i.e., within an interval predicate or between predicates, and also at the ends of the condition (the action is not involved in crossover). Mutation, however, is different. Preliminary experiments were done, and the best method appears to be to mutate an allele by adding an amount  $\pm rand(m_0)$ , where  $m_0$  is a fixed integer,  $rand$  picks a value uniform randomly from  $[0, m_0)$ , and the sign is chosen uniform randomly. If a new value of  $l_i$  is less than the minimum possible input value, in this case 1, the new value is set to 1. If the new value is greater than  $u_i$ , it is set equal to  $u_i$ . A corresponding rule holds for mutations of  $u_i$ .

The condition of a “covering” classifier has components  $l_0, u_0, \dots, l_n, u_n$ , where each  $l_i = x_i - rand_1(r_0)$ , but limited to the minimum possible input value, and each  $u_i = x_i + rand_1(r_0)$ , limited to the maximum possible input value;  $rand_1$  picks a random integer from  $(0, r_0)$ , with  $r_0$  a fixed integer.

Subsumption of one classifier by another occurs if every interval predicate in the first classifier’s condition subsumes the corresponding predicate in the second classifier’s condition. An interval predicate subsumes another one if its  $l_i$  is less than or equal to that of the other and its  $u_i$  is greater than or equal to that of the other. For purposes of action-set subsumption, a classifier is more general than another if its *generality* is greater. Generality is defined as the sum of the widths  $u_i - l_i + 1$  of the interval predicates, all divided by the maximum possible value of this sum.

## 4 Experiments on Synthetic Oblique Data

To get a feel for its behavior, we began by applying XCSI to some synthetic oblique data. In the first experiment the problem had just two dimensions, in order to give a clear picture of the resulting classifiers. In the other experiment, the problem had the same dimensionality as the WBC dataset—nine. In both cases the DS (discrimination surface) was diagonal to all the axes, i.e., maximally oblique.

### 4.1 Experiment 1—Two Dimensions

A dataset, “Random-Data2”, was constructed as follows. Random vectors  $(x_1, x_2)$  were created, with each  $x_i$  a random integer from  $(1, 10)$ . The correct outcome  $o$  for each vector was chosen according to

$$o = \text{if } x_1 + x_2 \geq 11 \text{ then } 1 \text{ else } 0.$$

Each vector and its outcome constituted an instance of Random-Data2. The dataset had 683 such instances (comparable to the size of the WBC dataset), and the attribute range from 1 through 10 corresponded to the WBC’s attribute range. Note that the DS for this problem is indeed diagonal to the  $x_1$  and  $x_2$  axes.

In the experiment, XCSI operated in a *pure explore* mode in which on each time-step the system was given a randomly selected instance from Random-Data2, it carried out the performance, update, and discovery procedures as explained in Sections 2 and 3, but its action selection rule was to choose its output action (its decision as to the correct outcome) randomly. If the action was correct, XCSI received a reward of 1000; if incorrect, 0. Both forms of subsumption deletion were enabled.

To determine whether the system was actually learning something, the explore steps alternated with *test* steps in which XCSI deliberately chose as its decision the action with the maximum system prediction. During test steps the update, discovery, and subsumption components were disabled.

The following parameters of XCSI were the same for all experiments in this paper:  $\beta = 0.2$ ,  $\epsilon_0 = 1$ ,  $n = 5$ ,  $\theta_{GA} = 48$ ,  $\chi = 0.8$ ,  $\mu = 0.04$ , and  $m_0 = 2$ . Also the same were the experience threshold for deletion and for GA subsumption,  $\theta_{del} = 50$  (*experience* is the number of times a classifier has been in an updated action set) and a separate experience threshold for action-set subsumption,  $\theta_{as} = 1000$ . The space of parameter settings is obviously large, and a number of experiments preceded those reported here. Principally, it was found that  $\epsilon_0$  should be quite small (to push toward high accuracy) and  $\theta_{as}$  should be quite large (to ensure that a classifier’s accuracy is reliable before allowing it to delete another classifier in action-set subsumption). Parameters for this experiment that were not the same through all experiments were  $N = 800$  and  $r_0 = 4$ .

Figure 1 shows results for a typical run of XCSI on Random-Data2. The quantities are all moving averages over the past 50 problems (instances seen); values are plotted every 100 problems. Individual plots appear in the top-to-bottom order shown in the legend.

Approximately 100,000 explore problems (instances) were required for the system to reach 100% performance. However, 98% performance (one error in 50 problems) was reached quickly, so most of the time was needed to get the last few instances—usually just one or two—right. “Generality” is a fitness-weighted average of classifier generality. It climbs slowly and seems to level off at around 0.5. “Popsiz/800” is the number of macroclassifiers divided

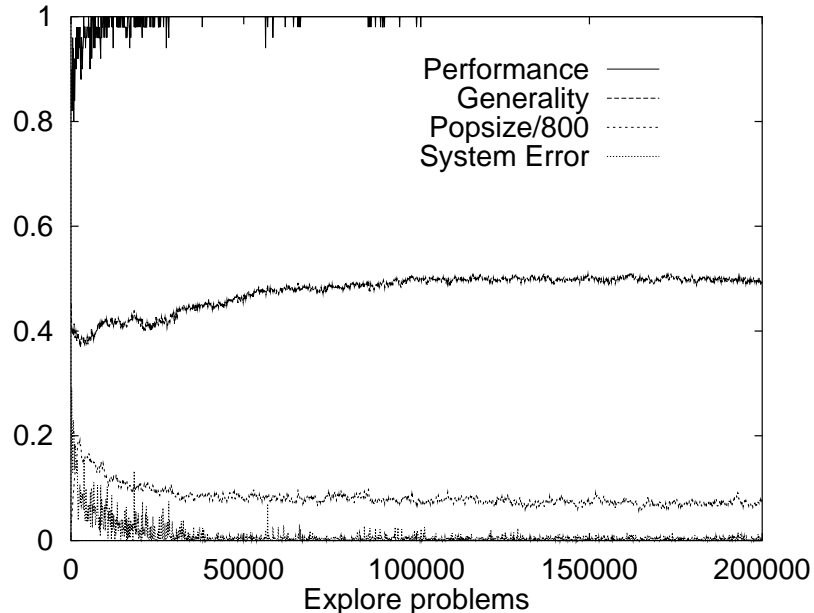


Figure 1: Performance (fraction correct), generality, population size (/800), and system error vs. number of explore problems in Experiment 1.

by  $N$ , the total numerosity. It falls quickly to approximately 0.08. The actual number of classifiers at the end was 63. Finally, “System Error” is the difference (divided by 1000) between the system prediction for the chosen action and the actual reward received. It falls quickly to a very small value.

Figure 2 depicts 29 representative classifiers from the 63 in the final population of Experiment 1, ordered by descending prediction. A graphic notation is used to represent the conditions. In each interval predicate (shown between “|”s), the range of input values accepted is indicated by a sequence of 0’s. Thus the first classifier accepts (matches) input vectors for which  $7 \leq x_1 \leq 10$  and  $4 \leq x_2 \leq 10$ .

Notice that the first classifier is always correct in its prediction: all accepted inputs satisfy  $x_1 + x_2 \geq 11$ , the classifier advocates action 1, and the reward received is indeed 1000. Furthermore, the classifier is maximally general: neither of its predicates can be expanded without causing an error. For the second classifier, all accepted inputs satisfy  $x_1 + x_2 < 11$ , it advocates action 0, and it is completely accurate predicting a reward of 1000. It, too, is maximally general. In fact all classifiers in the first group are accurate (error 0.000) and maximally general except no. 12 which is merely accurate.

The second group represents the 13 classifiers that were not accurate. Note that this is reflected in their very small fitnesses (all fitness values are shown multiplied by 1000). Because system predictions are a fitness-weighted sum (Section 2.1), low-fitness classifiers have minimal effect on the system, and are soon deleted. They represent GA offspring that turned out not to be accurate.

The classifier conditions in the first and third groups can be thought of as “boxes”—hyperrectangles in 2-d—fitting into the space on either side of the line defined by the dataset’s DS. The fit is awkward: apart from nos. 12 and 55, each box is as big as it can be without causing an error, but they overlap in such a way that none subsumes any other and all survive. In effect, the boxes are trying to approximate the diagonally divided input space

	CONDITION	ACT	PRED	ERR	FITN	NUM	GEN	EXPER
0.	.....0000 ...0000000	1	1000.	.000	282.	15	.55	29435
1.	000000.... 0000.....	0	1000.	.000	359.	26	.50	22546
2.	0000..... 000000....	0	1000.	.000	153.	13	.50	22780
3.	00000.... 00000....	0	1000.	.000	358.	19	.50	23150
4.	0000000.. 000.....	0	1000.	.000	295.	19	.50	18197
5.	....000000 ....00000	1	1000.	.000	150.	12	.55	29394
6.	000000000. 0.....	0	1000.	.000	435.	28	.50	7924
7.	.....000 ..00000000	1	1000.	.000	250.	23	.55	26022
8.	00000000.. 00.....	0	1000.	.000	236.	16	.50	12647
9.	.....00 ..000000000	1	1000.	.000	229.	29	.55	17336
10.	.....0 000000000	1	1000.	.000	429.	22	.55	9401
11.	..00000000 .....000	1	1000.	.000	213.	10	.55	19615
12.	0..... .....000.	0	1000.	.000	515.	29	.20	1658
13.	00..... 00000000..	0	1000.	.000	501.	24	.50	12346
25.	0000000.. 0000.....	0	991.	.032	0.	1	.55	46
26.	0000..... 00000000....	0	990.	.037	0.	1	.55	41
27.	000000.... 00000....	0	966.	.081	5.	3	.55	80
28.	00..... 000000000..	0	930.	.135	0.	1	.55	51
29.	000..... 00000000..	0	800.	.161	0.	2	.55	39
30.	...0000000 ...0000000	1	754.	.308	3.	1	.70	16
47.	00000.... 00000....	1	0.	.000	202.	15	.50	21652
48.	.....00 ..000000000	0	0.	.000	232.	20	.55	17411
49.	.....0 000000000	0	0.	.000	284.	23	.55	10208
50.	....000000 ....00000	0	0.	.000	407.	32	.55	27947
51.	..00000000 .....000	0	0.	.000	345.	8	.55	21179
52.	000..... 0000000....	1	0.	.000	251.	27	.50	18306
53.	0000000.. 000.....	1	0.	.000	188.	20	.50	16057
54.	....00000 ....000000	0	0.	.000	234.	20	.55	25685
55.	0..... .....000.	1	0.	.000	556.	28	.20	1735

Figure 2: Some classifiers from Experiment 1. Condition predicates shown graphically (see text). Also shown: ACTION, PREDiction, ERRor, FITNess, NUMerosity, GENerality, and EXPERience of each classifier.

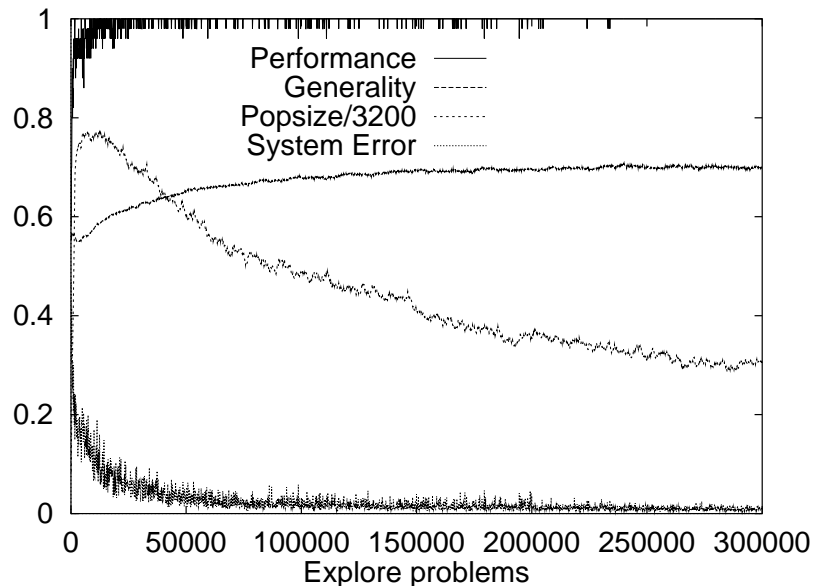


Figure 3: Performance (fraction correct), generality, population size (/3200), and system error vs. number of explore problems in Experiment 2

but can’t do it elegantly because they are rectangular. Despite this frustration, the system did its job: of the 63 classifiers in the set, all with substantial fitness (and thus able to affect the system prediction) were accurate, and all of these except nos. 12 and 55 were maximally general.

Still, this experiment’s input space had only two dimensions. Despite the system’s high performance and relatively small population (space complexity) in this case, what can be expected for more realistic data mining problems such as the WBC with its nine dimensions? Will XCSI’s “boxes” approximation be adequate, and what insight will it give into the problem’s structure? As a first step in answering these questions, we tested XCSI on a nine-dimensional problem with a “diagonal” DS.

## 4.2 Experiment 2—Nine Dimensions

A second random dataset, “Random-Data9”, was constructed like Random-Data2, except that there were nine dimensions and the expression determining the outcomes was

$$o = \text{if } x_1 + \dots + x_9 \geq 50 \text{ then } 1 \text{ else } 0.$$

The dataset contained 683 examples, 361 of which had outcome 1.

Experiment 2 was conducted like Experiment 1 except that  $N = 3200$  and  $r_0 = 6$ . Larger values for these parameters are generally needed as the input space becomes larger. In general, they are chosen empirically such that the total number of classifiers does not exceed approximately  $0.75N$  and any initial covering ceases quickly.

Results are presented in Figure 3, showing a typical run. In comparison with Experiment 1, performance took somewhat longer to reach 100%, system error fell more slowly, and the maximum and final population sizes were substantially larger. On the other hand, the input space is exponentially larger— $10^9$  vs.  $10^2$ —whereas the results differ by a factor of 10 or

less. Why the 9-d problem is not harder is not clear at present. It may have to do with the input space being only sparsely occupied by dataset instances, in comparison with the much fuller coverage in the 2-d case. Sparse coverage makes the DS irregular and thus less strictly diagonal; this may aid the classifiers in approximating it.

The 974 classifiers of the final population were examined, but examples are not shown here. As in Experiment 1, roughly 80% of the classifiers were accurate, with the remaining inaccurate ones having very low fitness. Because of the dimensionality and sparseness of the dataset, and unlike Experiment 1, it was not readily possible to judge if the accurate classifiers were maximally general. The accurate classifiers showed no clear pattern in their predicates, as would be expected from the obliqueness of the formula for creating the dataset, together with the sparseness. Nevertheless, as seen in Figure 3, the dataset was learned. It appears that XCSI can deal with nine-dimensional oblique data, at least from the point of view of performance.

## 5 Experiments on The WBC Dataset

We used the “original” Wisconsin Breast Cancer Database which was donated to the UCI Repository by Prof. Olvi Mangasarian and contains 699 cases (instances) collected over time by Dr. William H. Wolberg. For concreteness, the nine attributes of each instance are: Clump Thickness, Uniformity of Cell Size, Uniformity of Cell Shape, Marginal Adhesion, Single Epithelial Cell Size, Bare Nuclei, Bland Chromatin, Normal Nucleoli, and Mitoses. Each attribute has a value between 1 and 10 inclusive. Sixteen instances contain an attribute whose value is unknown. The outcome distribution is 458 Benign (65.5%), 241 Malignant (34.5%). The following is a small sample of the raw data.

```

1070935,3,1,1,1,1,1,2,1,1,2
1071760,2,1,1,1,2,1,3,1,1,2
1072179,10,7,7,3,8,5,7,4,3,4
1074610,2,1,1,2,2,1,3,1,1,2
1075123,3,1,2,1,2,1,2,1,1,2
1079304,2,1,1,1,2,1,2,1,1,2
1080185,10,10,10,8,6,1,8,9,1,4
1081791,6,2,1,1,1,1,7,1,1,2
1084584,5,4,4,9,2,10,5,6,1,4
1091262,2,5,3,3,6,7,7,5,1,4
1096800,6,6,6,9,6,?,7,8,1,2

```

The first number is a label, the next nine are the attribute values, and the last is the outcome, 2 for Benign and 4 for Malignant. Clearly, Malignant is associated with larger attribute values. But the precise borderline (DS) is not obvious, nor is it clear whether certain attributes are more important than others. The obliqueness of the data is reflected in the fact that a summation of attribute values or “weight of evidence” appears broadly to be decisive. On an earlier (smaller) version of the data, Wolberg and Mangasarian (1990) used a statistical method with multiple hyperplanes to obtain accuracies as high as 95.9%. The UCI Repository states “highest reported accuracy” on the current dataset to be 94%.

XCSI will first be applied to the WBC dataset in a “train/test” experiment using a well-known procedure for evaluating learning algorithms for data mining. After that, we will present an experiment in which XCSI extracted some of the patterns contained in the data.

### 5.1 Experiment 3—WBC Train/Test

A practical learning algorithm for data mining must of course be able to categorize instances correctly that it hasn’t seen before. In the experiments of Section 4, XCSI learned a dataset perfectly, but was not tested on further examples. In this experiment, we applied XCSI to the WBC data in a stratified tenfold cross-validation procedure in which the system learned on part of the data and was tested on the remainder.

According to Witten and Frank (2000), “Tenfold cross-validation is the standard way of measuring the error rate of a learning scheme on a particular dataset; for reliable results, ten times tenfold cross-validation.” We followed the procedure as they explain it, adding stratification, but left the 10 times repeat for later work.

Briefly, the procedure is as follows. The dataset is divided into 10 parts called “folds”. The system is then tested on each fold, after being trained on the balance of the dataset. Then the results on the 10 test folds are averaged giving the final score. The folds are made as equal in size as possible, given the size of the actual dataset. “Stratification” means that in each fold the possible outcomes have approximately the same prevalences as in the whole dataset.

The folds were created from the dataset by random selection without replacement. To get stratification in a fold, it was actually made up of two parts, one containing instances randomly selected without replacement from the remaining 0-outcomes in the dataset, the other from the remaining 1-outcomes. Then the two fold parts were shuffled randomly together. The actual numbers of outcomes selected for each fold were given by a predetermined pair of numbers. For eight of the ten folds, the pair was (46, 24) for 0-outcomes and 1-outcomes, respectively. For the last two folds, the pairs were (45, 24) and (45, 25), thus accounting for all the instances as evenly as possible. (The WBC outcome notation of 2 and 4 was replaced by 0 and 1.)

Parameters for the experiment were the same as previously described, except  $N = 6400$  and  $r_0 = 4$ . If an attribute value was missing, it was assumed to match. Calling each fold an “Out” set and the corresponding balance of the dataset an “In” set, XCSI learned on each of In sets for 50,000 explore problems, then was tested on the associated Out set. Each learning run reached 100% performance by 40,000 problems, with no mistakes after that, except for three runs that reached 100% at approximately 49,500, 49,000, and 40,100 problems. Testing consisted of a simple sweep through the Out set with XCSI in test mode. The test results are shown in Table 1.

With a mean performance of about 95.5%, XCSI appears on this problem to be in the same league as the best of other approaches. Of course, the present experiments are exploratory, and further testing is called for on the WBC data—using different fold sets—and on other oblique data problems.

### 5.2 Experiment 4—Discovering Patterns

As XCSI grinds through thousands of data instances, its generalization and subsumption mechanisms evolve increasingly general, accurate classifiers. On Boolean functions or the

Correct	Incorrect	Not Matched	Fraction Correct
68	2	0	0.9714
69	1	0	0.9857
65	5	0	0.9286
66	4	0	0.9429
65	3	2	0.9286
64	3	3	0.9143
70	0	0	1.0000
69	1	0	0.9857
65	3	1	0.9420
67	2	1	0.9571
		MEAN $\Rightarrow$	0.9556

Table 1: Results of stratified tenfold cross-validation test on WBC dataset. Each line has results for one fold. Last line is mean of fold scores.

Monk’s Problems the result is a small set of classifiers in which the problem structure is evident. On oblique data, however, the discrimination surface must be approximated by hyper-rectangles defined by the classifier conditions. To the extent the data are truly oblique, the evolved classifiers may reveal little about the DS. In Experiment 1, for instance, no individual classifier in Figure 2 offers any hint as to the shape of the DS, even though each one having high fitness is accurate and maximally general. On the other hand, if the data while largely oblique contains some “logical” structure—i.e., structure that hyper-rectangles can more directly approximate—we should be able to see it in the evolved classifiers.

An oblique dataset like WBC can be tested for logical structure by continuing the evolution far past the point of 100% performance. Evolutionary search does not stop if the GA is running and further generalization of existing classifiers is possible. We can tell that the search continues as long as population generality is rising and population size is falling. In Experiment 4 we carried a learning run on the entire WBC dataset out to 2,000,000 time-steps (instances seen) and examined the resulting classifiers. Graphs of the experiment are shown in Figure 4. Notice that 100% performance is reached by roughly 50,000 problems, but generality and population size are still changing at 2,000,000 problems.

Figure 5 shows the 27 most experienced classifiers in the final population (their experiences ranged from 50,810 to 414,215). These are rules that make no mistakes in predicting the payoff for any WBC dataset instance that they match. Because of their large experiences, it is quite safe to assume that they are generalized up to the point where expanding any of the interval predicates, even slightly, would cause an error. In short they are never contradicted by the dataset, and they are maximally general.

High experience means a classifier matches a large fraction of the dataset. The classifiers in Figure 5 should offer some insight about the dataset’s structure and by extension, the domain from which the data were drawn. For example, No. 3 appears to say, “Regardless of the values of the remaining attributes, if attribute 1 is 7 or above and attribute 2 is 5 or above, the case is malignant.” In medical terms, this would be, “If clump thickness is 7 or above and uniformity of cell size is 5 or above, malignancy is indicated.”

The statement is an hypothesis derived from the data, certainly not a proven fact. In fuller terms, it would be: “To the extent the dataset accurately samples the universe of breast cancer cases, it is hypothesized that no case fitting the ranges of No. 3 will not

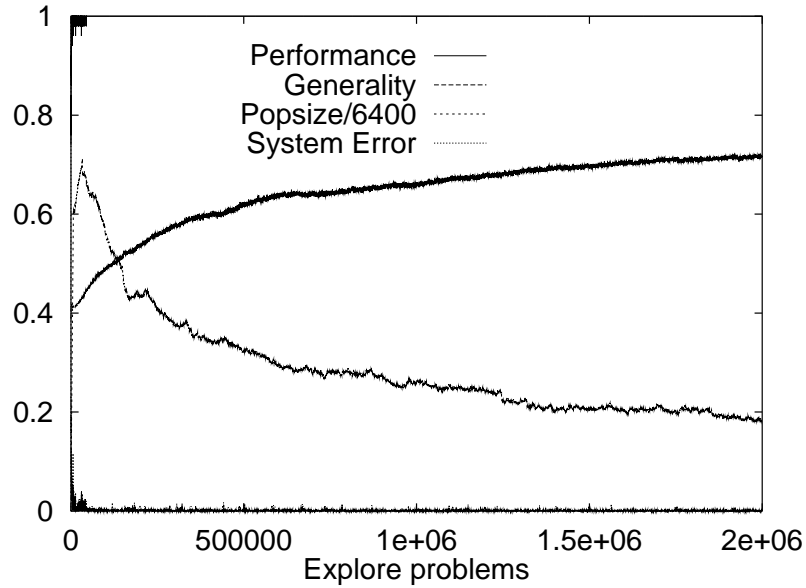


Figure 4: Performance (fraction correct), generality, population size (/6400), and system error vs. number of explore problems in Experiment 4.

0.	00000000... 0000000000 0000000000 00..... 00..... 0000000000 0000000000 0000000000 0000000000 0 1000
1.	0000000000 0..... 0000000000 0000000000 0000000000 0000..... 00..... 0000000000 0000000000 1 0
2.	0000000000... 0000000000 0000000000 0000000000 0000000000 0000..... 00..... 0000000000 0000000000 0 1000
3.	.....0000 ...000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 1 1000
4.	.....0000 ...000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0 0
5.	0000000000 0000000000 .....0000 .0000000000 0000000000 0000000000 0000000000 ..00000000 0000000000 0 0
6.	0000000000 ...000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 .0000000000 0 0
7.	0000000000 0000000000 0000..... 0000000000 0000000000 0000..... 000..... 0..... 0000000000 1 0
8.	0000000000 ...000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 .0000000000 1 1000
9.	.....00 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 1 1000
10.	.....00 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0 0
11.	0000000000 0000000000 .....000 .0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 1 1000
12.	0000000000 0000000000 0000..... 0000000000 0000000000 00..... 0000000000 00..... 0000000000 0 1000
13.	0000000000 0000000000 0000000000 .0000000000 0000000000 0000000000 0000000000 .....0 0000000000 1 1000
14.	0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 .....0 0000000000 0 0
15.	.....0000 0000000000 0000000000 0000000000 0000000000 .....000 0000000000 0000000000 0000000000 1 1000
16.	.....0000 0000000000 0000000000 0000000000 0000000000 .....000 0000000000 0000000000 0000000000 0 0
17.	0000000000 0000000000 ..00000000 .0000000000 0000000000 0000000000 0000000000 00000000... 0000000000 0 0
18.	0000000000 0000000000 0000..... 0000000000 0000000000 00..... 0000000000 00..... 0000000000 1 0
19.	0000000000 0000000000 0000000000 0000000000 0000000000 .....00000 .00000000 00000000... 0000000000 0 0
20.	00000000... 0000000000 0000000000 00..... 00..... 0000000000 0000000000 0000000000 0000000000 1 0
21.	0000000000 0000000000 0000000000 0000000000 0000000000 .....00000 .00000000 00000000... 0000000000 1 1000
22.	000000... 0000000000 ...00000 .0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 0 0
23.	000000... 0000000000 ...00000 .0000000000 0000000000 0000000000 0000000000 0000000000 0000000000 1 1000
24.	...00000 0000000000 0000000000 0000..... 0000000000 ...000000 0000000000 0000000000 0000000000 1 1000
25.	00000000... 0000000000 0000000000 0000000000 0000000000 0000..... 00..... 0000000000 0000000000 1 0
26.	0000000000 0000000000 0000000000 0000000000 0000000000 .....000 0000000000 0000000000 1 1000

Figure 5: The most experienced classifiers from the final population of Experiment 4.

be malignant. Furthermore, there are cases that fit No. 3 except for clump thickness less than 7 or uniformity of cell size less than 5 that are benign.” Thus the classifier makes an hypothesis about cases that fit its condition, and a weaker hypothesis about cases that nearly fit its condition.

Figure 5 includes some classifiers in which a single attribute is decisive. For example No. 9 says, “If clump thickness is 9 or above, then malignant.” No. 26 says, “If bland chromatin is 8 or greater, then malignant.” Some are very close to depending on a single attribute. No. 11 says, “If uniformity of cell shape is 8 or above and marginal adhesion is not 1, then malignant. Similarly for Nos. 6 and 8. Furthermore, rules dependent on just two attributes can be identified, for instance Nos. 1 (which indicates benign), 5, and 15. In all of these we have omitted the associated weaker hypothesis.

## 6 Discussion

XCSI’s good performance on both the synthetic oblique problems and the WBC dataset suggests that the apparent mismatch between the obliqueness of the discrimination surface and XCSI’s hyper-rectangular predicates is not a fundamental limitation. The boxy predicates might have been expected to lead to *overfitting*, i.e., excessively detailed learning of the training data causing poor performance on new instances. But test performance on the WBC data was at a competitive level, suggesting that overfitting was not in fact a problem. Still, it is not clear why this pleasant outcome should have occurred. Perhaps there is something about the “boxes” approximation that actually tends to prevent overfitting.

The emergence of some pattern information in the WBC data suggests that the set is not totally oblique, but instead contains regions where there is a strong dependence on one or a few attributes. XCSI seems to zero in on such patterns. Hopefully, the patterns will make sense to clinicians or practitioners, aiding confidence in the system’s predictions, and providing suggestive domain hypotheses.

Beyond the patterns shown in Figure 5, the final population contained of course a great many others having a more intricate structure. Some of these classifiers are concerned with the region close to the DS, where discrimination is complex. The patterns of Figure 5, though they cover huge subspaces, may not in fact apply in “close cases” near the DS. There, however, the match set will often indicate the degree of uncertainty of the decision, because it will contain high fitness classifiers that vote strongly each way. A measure of match set agreement could assist practitioners in the close cases. Holmes (1997) used this sort of technique for evaluating risk of disease.

Next steps for research on XCS and oblique data should include further fold tests on WBC and tests on additional representative datasets. Classifier patterns such as those seen in Figure 5 should be shown to domain practitioners for plausibility. Analysis is needed of the boxes approximation, including the result that the learning complexity of XCSI on oblique data is apparently far below exponential. The conditions of the classifiers of evolved populations contain a great deal of information about the structure of the dataset, but only a fraction of it is directly interpretable as in Section 5.2. It would be desirable to find algorithms for extracting all the classifiers’ implications as rules of thumb and in other representations.

## 7 Conclusion

XCS, suitably modified for numeric inputs, appears to have considerable potential for mining oblique datasets. On tests so far, performance was high, and simple as well as more complex patterns describing the structure of the data emerged. Further investigation on a wide range of problems may show that systems based on XCS are a very good choice for high performance combined with clear pattern discovery.

## Acknowledgment

The author acknowledges helpful conversations with John H. Holmes and Pier Luca Lanzi, and support from the Automated Learning Group at the National Center for Supercomputer Applications in Urbana-Champaign.

## References

- Banzhaf, W., J. Daida, A. E. Eiben, M. H. Garzon, V. Honavar, M. Jakiela, and R. E. Smith (Eds.) (1999). *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-99)*. Morgan Kaufmann: San Francisco, CA.
- Blake, C. and C. Merz (1998). UCI repository of machine learning databases. <http://www.ics.uci.edu/~mlearn/MLRepository.html>.
- Butz, M. V. and S. W. Wilson (2000). An algorithmic description of XCS. IlliGAL Report No. 2000017, Illinois Genetic Algorithms Laboratory, University of Illinois at Urbana-Champaign.
- Holmes, J. H. (1997). Discovering Risk of Disease with a Learning Classifier System. In T. Bäck (Ed.), *Proceedings of the 7th International Conference on Genetic Algorithms (ICGA97)*. Morgan Kaufmann: San Francisco CA.
- Kovacs, T. (1997). XCS classifier system reliably evolves accurate, complete, and minimal representations for boolean functions. In Roy, Chawdhry, and Pant (Eds.), *Soft Computing in Engineering Design and Manufacturing*, pp. 59–68. Springer-Verlag, London.
- Kovacs, T. (1999). Deletion schemes for classifier systems. See Banzhaf, Daida, Eiben, Garzon, Honavar, Jakiela, and Smith (1999), pp. 329–336.
- Lanzi, P. L. (1999). Extending the representation of classifier conditions. Part II: From messy coding to S-expressions. See Banzhaf, Daida, Eiben, Garzon, Honavar, Jakiela, and Smith (1999), pp. 345–352.
- Lanzi, P. L., W. Stolzmann, and S. W. Wilson (Eds.) (2000). *Learning Classifier Systems: From Foundations to Applications*, Volume 1813 of *LNAI*. Berlin: Springer-Verlag.
- Saxon, S. and A. Barry (2000). XCS and the Monk’s Problems. See Lanzi, Stolzmann, and Wilson (2000), pp. 223–242.
- Sutton, R. S. and A. G. Barto (1998). *Reinforcement Learning: An Introduction*. Cambridge, MA: The MIT Press/Bradford Books.
- Thrun, S. B. et al. (1991). The monk’s problems: A performance comparison of different learning algorithms. Technical Report CMU-CS-91-197, Carnegie Mellon University, Pittsburgh, PA.

- Wilson, S. W. (1994). ZCS: A zeroth level classifier system. *Evolutionary Computation* 2(1), 1–18.
- Wilson, S. W. (1995). Classifier fitness based on accuracy. *Evolutionary Computation* 3(2), 149–175.
- Wilson, S. W. (1998). Generalization in the XCS classifier system. In J. R. Koza, W. Banzhaf, K. Chellapilla, K. Deb, M. Dorigo, D. B. Fogel, M. H. Garzon, D. E. Goldberg, H. Iba, and R. Riolo (Eds.), *Genetic Programming 1998: Proceedings of the Third Annual Conference*, pp. 665–674. Morgan Kaufmann: San Francisco, CA.
- Wilson, S. W. (2000a). Get Real! XCS with Continuous-Valued Inputs. See Lanzi, Stolzmann, and Wilson (2000), pp. 209–220.
- Wilson, S. W. (2000b). State of XCS Classifier System Research. See Lanzi, Stolzmann, and Wilson (2000), pp. 63–82.
- Witten, I. H. and E. Frank (2000). *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. San Francisco, CA: Morgan Kaufmann.
- Wolberg, W. H. and O. L. Mangasarian (1990). Multisurface method of pattern separation for medical diagnosis applied to breast cytology. *Proceedings of the National Academy of Sciences* 87, 9193–9196.