

XCSF with Computed Continuous Action

Trung Hau Tran, Cédric Sanza, Yves Duthen

IRIT
Paul Sabatier University
Toulouse, France
hau@irit.fr

Dinh Thuc Nguyen

Faculty of Information Technology
University of Natural Sciences
HoChiMinh City, Vietnam

ABSTRACT

Wilson introduced XCSF as a successor to XCS. The major development of XCSF is the concept of a computed prediction. The efficiency of XCSF in dealing with numerical input and continuous payoff has been demonstrated. However, the possible actions must always be determined in advance. Yet domains such as robot control require numerical actions, so that neither XCS nor XCSF with their discrete actions can yield high performance. This paper studies computed action in XCSF, where the action is continuous with respect to the input state. In comparison with Wilson's architecture for continuous action, our XCSF version, called XCSFCA, proves to be more efficient.

Categories and Subject Descriptors

I.2.6 [Artificial Intelligence]: Learning—*Concept learning*; I.2.8 [Artificial Intelligence]: Problem Solving, Control Methods, and Search—*Heuristic methods*.

General Terms

Algorithms.

Keywords

XCSF, continuous actions, evolution strategy.

1. INTRODUCTION

Research on learning classifier systems has focused on the use of accuracy-based fitness after the introduction of XCS [11], in which (1) the classifier fitness is based on the accuracy of the prediction and (2) general and accurate classifiers are maximally obtained. The payoff is a discontinuous function of the state x . The predictions are estimated values representing the expected payoffs if the system matches x and activates an action. The system builds a complete map of predictions *Condition x Action* \Rightarrow *Prediction*. However, when the number of states and actions increases, it is clear that the system cannot construct such map of predictions. This is not only due to the memory allocated by the system, but also the time as well as the examples used for training

the system. The system should generalize from a subset of data to approximate a payoff function as the study of the generalization in reinforcement learning [7]. Wilson recently introduced a new learning classifier system XCSF [13][14][15] that differs from XCS in real value inputs for the condition part and the new concept of a computed classifier prediction. XCSF learns how to construct a function approximator which approximates a payoff function. It can be a discontinuous 0/1000 payoff function or a continuous payoff function of real inputs. When XCSF receives the input x , it builds a *system prediction* for each discrete action $\hat{P}(x,a)$ via the function approximator and chooses the action a^* with the highest prediction $\hat{P}(x,a^*)$ to maximize the payoff. Some XCSF versions adapting to different types of payoff functions are presented in [4][13][14][15].

In traditional learning classifier systems, the possible actions must always be determined in advance. An animat in *woods* problems [10][11] can take one of eight actions corresponding to eight directions (*go north, north-east, east, south-east...*). One can use three bits to encode these eight actions in the action part. However, when dealing with problems that require real value outputs, the systems show some limits. In the example of an animat which could move in any direction represented by an angle from $[0, 2\pi]$, the solution is to increase the number of bits encoding actions in order to represent real angles. This one is feasible but we might lose the reactive attribute in the learning classifier systems. Indeed, the population size of classifiers, the time of building a match set and the convergence time of the system increase. Another solution is to use a generalized classifier system (GCS) [16], in which the input x is linked with the action a in the computation of the prediction as well as in the process of building the match set M . Each classifier employs the form of $t(x,a):w \Rightarrow p(x,a)$ where $t(x,a)$ is the condition form allowing the classifier to be or not to be a member of M . The prediction weight vector w extends a new weight associated with the action a and the prediction $p(x,a)$ is computed linearly with x and a . The difficulty of GCS is to evolve a condition form $t(x,a)$ so that the action is continuous. Wilson stated that if $t(x,a)$ employs interval predicates, the best action a^* is discontinuous. An example of GCS with orientable elliptical form $t(x,a)$ is applied to the frog problem.

In this paper, we will investigate a new approach in which the action is computed directly. Thus, the action is continuous with respect to the input state. Basing on XCSF with a computed classifier prediction $\hat{P}(x,a)$, we extend a computed classifier action to XCSF. The classifier action is computed as a linear combination of the input and a vector of action weights. The

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

GECCO'07, July 7–11, 2007, London, England, United Kingdom.
Copyright 2007 ACM 978-1-59593-697-4/07/0007...\$5.00.

action weights are themselves adapted to the problem by using an evolution strategy. We apply our XCSF version, called XCSFCA, to the frog problem described in [15][16]. The results show that XCSFCA obtains better performances than GCS. In addition, the population size of classifiers is more compact.

The next section describes briefly XCSF. Section 3 presents XCSFCA, its application in the frog problem and the results. Section 4 concludes this paper.

2. REVIEW OF XCSF

XCSF is an extension to XCS for learning environments where the input is continuous and the payoff landscape is continuous with respect to the input. A new classifier condition form adapting to real value inputs is introduced that allows XCSF to evolve more appropriate conditions. The use of a computed prediction as a function of the input allows XCSF to evolve classifiers with accurate values of the prediction. Thus, the system performance improves, XCSF can yield more compact solutions, and the convergence time may increase.

XCSF runs like XCS. The system receives an input message x via the detectors, forms a match set M of classifiers whose condition part matches x , determines a winner action a_{win} between possible actions in M according to an action selection strategy, forms an action set A of the classifiers advocating a_{win} , and sends a_{win} to the effectors. The system receives a reward r at the next step to update the parameters of the classifiers in A . A new component updates the prediction weights in order to make the predictions more accurate. Covering operator is triggered according to a covering strategy. Either no existing classifier matches x as in [11] or M contains less than ϕ_{mna} discrete actions as in [3]. A niche genetic algorithm is invoked to generate potential classifiers adapting to the problem.

Classifier condition. XCSF extends XCS to deal with continuous input. The components of a classifier condition are the interval predicates $int_i=(l_i, u_i)$ where l_i (“lower”) and u_i (“upper”) are real values. A classifier cl matches an input message $x=(x_1, \dots, x_n)$ if each element x_i belongs to the corresponding interval predicate $cl.int_i$ at the position i i.e. $l_i \leq x_i \leq u_i$. Fortunately, the system with the real input was studied in [12] and one can use again the methods manipulating interval predicates.

Computed classifier prediction. The new concept in XCSF is the computation of the classifier prediction instead of using a parameter estimating its value. A vector of weights $w=(w_0, w_1, \dots, w_n)$, called the prediction weight vector, is added to each classifier. The classifier prediction $cl.p$ is computed as a linear combination of the prediction weight vector $cl.w$ and the input message $x=(x_1, \dots, x_n)$ concatenated with a constant x_0 . Wilson’s studies [14] proposed that the value of x_0 should be the same order of the values of the elements of x .

$$\begin{aligned} cl.p(x) &= w \cdot x' \\ w &= (w_0, w_1, \dots, w_n) \\ x' &= (x_0, x_1, \dots, x_n) \end{aligned} \quad (1)$$

Covering operator. When a new covered classifier is created, each interval predicate $int_i=(l_i, u_i)$ is generated as $l_i = x_i - rand(r_0)$ and $u_i = x_i + rand(r_0)$, where $rand(r_0)$ is a value uniform randomly from $[0, r_0]$ and r_0 is a real constant, and each element w_i of the

weight vector is initialized to a value uniform randomly from $[-1, 1]$.

Discovery mechanism. A genetic algorithm works as in XCS. Crossover (with probability χ) permutes alleles of two parents between two crossover points. Since an allele is a real value, a new mutation operator is introduced. Mutation (with probability μ) modifies an allele by adding an amount $\pm rand(m_0)$ where m_0 is a real constant.

Domain control. Covering and mutation operators can generate intervals out of the condition range. If it happens, l_i and u_i of an interval predicate int_i are brought back to the extremes of the condition range and they are possibly permuted in order to respect the predicate constraint $l_i \leq u_i$.

Reinforcement. The prediction weight vector $cl.w$ of each classifier in A is updated using a *modified delta rule* [9]:

$$\Delta w_i = \left(\frac{\eta}{\|x'\|^2} \right) (P - cl.p(x)) x_i \quad (2)$$

$$cl.w_i \leftarrow cl.w_i + \Delta w_i \quad (3)$$

Where η is the correct rate. The classifier prediction error $cl.\varepsilon$ and the classifier fitness $cl.fit$ are then updated as follows:

$$cl.\varepsilon \leftarrow cl.\varepsilon + \beta (P - cl.p(x) - cl.\varepsilon) \quad (4)$$

$$cl.\kappa = \begin{cases} 1 & \text{if } cl.\varepsilon < \varepsilon_0 \\ \alpha \left(\frac{cl.\varepsilon}{\varepsilon_0} \right)^{-\nu} & \text{otherwise} \end{cases} \quad (5)$$

$$cl.\kappa' = \frac{cl.\kappa \times cl.num}{\sum_{j \in [A]} cl_j.\kappa \times cl_j.num} \quad (6)$$

$$cl.fit \leftarrow cl.fit + \beta (cl.\kappa' - cl.fit) \quad (7)$$

Where β is the learning rate. The constant ε_0 represents the threshold of tolerance for the prediction error $cl.\varepsilon$. If $cl.\varepsilon$ is below ε_0 , the error is accepted and the classifier cl is considered to be accurate ($cl.\kappa=1$). Otherwise its accuracy drops off and is controlled by an accuracy function with the parameters α and ν . The relative accuracy $cl.\kappa'$ is calculated with respect to the other classifiers cl_j of A . $cl.num$ is the numerosity of the classifier cl . Notice that $cl.w$ is updated first, so $cl.\varepsilon$ and $cl.fit$ take into account the change of $cl.w$.

3. EXTENSION TO XCSF

3.1 XCSF with Computed Continuous Action

As introduced in the first section, XCSF with discrete actions performs inefficiently when facing problems requiring continuous real values. We propose an extension to XCSF where actions are directly computed. The actions are then a continuous function of the input. The modified classifier system is called XCSFCA.

A vector of action weights $\zeta=(\zeta_0, \zeta_1, \dots, \zeta_n)$ is associated to each classifier action. The classifier action $cl.a$ is computed as a linear combination of ζ and the current input $x=(x_1, \dots, x_n)$ concatenated with a constant x_0 .

$$\begin{aligned} cl.a(x, \zeta) &= \zeta \cdot x' \\ x' &= (x_0, x_1, \dots, x_n) \end{aligned} \quad (8)$$

A separate process called an evolution strategy (ES) evolves the action weights to make the computed actions more accurate. The ES was developed by Schwefel [6] for optimizing real-vectors. Each individual in the ES consists of a pair of two real-vectors (ζ, σ) having the same length n_σ : a vector of action weights ζ and a vector of standard deviations σ used by mutation. The ES, which we use to evolve the action weights, is (1+1)-ES, in which one parent produces one offspring by mutation, the offspring is assigned a fitness value evaluating its quality with respect to the problem, the offspring competes with the parent based on the fitness, and the best individual becomes the parent of the next generation. Mutation on ζ is performed by adding a random value $N(0, \sigma_i)$ from a normal distribution with mean zero and standard deviation σ_i . The standard deviation vector is also itself adapted.

$$\sigma'_i = \sigma_i \exp(\tau' N(0,1) + \tau N_i(0,1)) \quad (9)$$

$$\zeta'_i = \zeta_i + N(0, \sigma'_i) \quad (10)$$

$$\tau' = \left(\sqrt{2n_\sigma}\right)^{-1} \quad (11)$$

$$\tau = \left(\sqrt{2\sqrt{n_\sigma}}\right)^{-1} \quad (12)$$

Where (ζ', σ') is an individual mutant, n_σ is the length of the vector σ , and the factors τ' and τ are set by default. The normal distribution with mean zero and standard deviation one $N_i(0,1)$ is called anew for each value of $i=0, \dots, (n_\sigma-1)$. Each element ζ_i is initialized to a value uniform randomly from $[-1,1]$ and each element σ_i is initialized to a constant value $\sigma^{(0)}$.

XCSFCA works like XCSF with a computed classifier prediction $\hat{P}(x, a)$, except for the following cases: (1) the process of building the match set M , covering operator, action selection in exploration and crossover operator are modified, and (2) a new component updating action weights is introduced.

Match set M . A classifier will become a member of M if its condition part matches the input as usual and its computed action $cl.a$ belongs to the range of the action a_{range} . The ranges of action weights are not given. The action weights are initially generated randomly and then updated by mutation. So, a classifier can compute *bad* actions out of a_{range} in response to the input. We modified the process of building M so that M consists of matching classifiers advocating *good* actions, which belong to a_{range} , and thus they will have a chance to correct their action weights.

Covering operator. Given the input x , a vector of action weights $cl.\zeta$ is repeatedly generated until an action $cl.a$ computed by Eq.(8) belongs to the action range allowed. This ensures that the match set M has at least one *good* action in response to the input. Interval predicates $cl.int_i$ are created as in XCSF.

Action selection in exploration. A pure random exploration is usually used for exploring actions of the environment. We propose to replace it by a deterministic strategy, i.e. the action with the highest prediction is selected. The classifiers advocating the best action will survive while an ES will allow the system to explore actions.

Crossover operator. The recombination of alleles of two parents does not generate new values of alleles. The system cannot explore values near the parents' values, and thus appropriate conditions are not generated. To solve this problem, we use a *mix* crossover [5]. The allele value v_i^o of an offspring is the result of a *mix* function of two allele values v_i^1 and v_i^2 of parents at the same position i .

$$v_i^o = (1 - \gamma_{GA})v_i^1 + \gamma_{GA}v_i^2 \quad (13)$$

$$\gamma_{GA} = (1 + 2\alpha_{GA})u - \alpha_{GA} \quad (14)$$

Where α_{GA} is a real constant from $[0,1]$ and u is a random value from $[0,1]$. When $\alpha_{GA}=0$, the *mix* crossover is a kind of crossover that interpolates two parents' values. Notice that XCSFCA still has the XCSF mutation operator.

Action weights update component. An (1+1)-ES is applied on each member of action set A to evolve actions during explore problems. Indeed, it will mutate the current action of each classifier of A to produce one mutant action and evaluate it. Accurate actions will be kept to the next time-step. Consider a classifier cl in A with its action weights $cl.\zeta$, its standard deviations $cl.\sigma$, its error $cl.\varepsilon$, its experience $cl.exp$, its numerosity $cl.num$ and its fitness $cl.fit$. The mutant classifier cl' is a copy of the parent classifier cl . It inherits the parameters as $cl.\varepsilon$, $cl.fit$, $cl.exp$ and $cl.num$ of the parent classifier cl except for its action represented by $(cl'.\zeta, cl'.\sigma)$ generated by mutation on $(cl.\zeta, cl.\sigma)$. This corresponds to step 1 (Algorithm 1). The evaluation of the mutant classifier's fitness works as follows. The mutant action represented by $(cl'.\zeta, cl'.\sigma)$ is calculated. The mutant classifier cl' computes its prediction and is tested in the context of the current payoff P . Then its prediction error $cl'.\varepsilon$ is estimated and its accuracy value $cl'.\kappa$ is calculated. Finally, its fitness $cl'.fit$ is updated by its relative accuracy value. These correspond to step 2 (Algorithm 1). If the mutant is accurate and better than the parent in the fitness sense, it will replace its parent (step 3 in Algorithm 1). If replacement doesn't occur, the mutant is discarded and a new mutant will be generated in next problems. This update process (Algorithm 1) is called after the reinforcement component.

Algorithm 1. Update action weights by (1+1)-ES.

1. Mutate one parent action represented by $(cl.\zeta, cl.\sigma)$ to produce one mutant action $(cl'.\zeta, cl'.\sigma)$ by Eq.(9) and Eq.(10)

2. Evaluate the mutant action $(cl'.\zeta, cl'.\sigma)$

$cl'.a(x, cl'.\zeta)$ // computed action of the mutant cl'

$good = \begin{cases} 1 & \text{if } cl'.a \in \text{action range allowed} \\ 0 & \text{otherwise} \end{cases}$

// the mutant cl' computes its prediction $cl'.p$

// and updates its error $cl'.\varepsilon$ and its fitness $cl'.fit$

$cl'.\varepsilon \leftarrow cl'.\varepsilon + \begin{cases} \frac{(P - cl'.p(x, cl'.a)) - cl'.\varepsilon}{cl'.exp} & \text{if } cl'.exp < \frac{1}{\beta} \\ \beta(P - cl'.p(x, cl'.a)) - cl'.\varepsilon & \text{otherwise} \end{cases}$

$$cl'.\kappa = \begin{cases} 1 & \text{if } cl'.\varepsilon < \varepsilon_0 \\ \alpha \left(\frac{cl'.\varepsilon}{\varepsilon_0} \right)^{-\nu} & \text{otherwise} \end{cases} \quad // \text{ accuracy}$$

$$cl'.fit = \delta_{ES} cl.fit \quad // \text{ reduce inherited fitness by } \delta_{ES}$$

$$cl'.fit \leftarrow cl'.fit + \beta \left(\frac{cl'.\kappa \times cl'.num}{cl.\kappa \times cl.num + cl'.\kappa \times cl'.num} - cl'.fit \right)$$

3. Replace the parent classifier if the mutant classifier is accurate and dominates the parent classifier
if (*good* = 1 and *cl'.κ* = 1 and *cl'.fit* > *cl.fit*)
replace cl by cl'

Where β , ε_0 , α , ν are parameters used in XCS and the factor δ_{ES} is used to reduce the mutant's fitness.

Computed classifier prediction. Notice that because the classifier prediction $cl.p$ is computed linearly with x and a , a weight associated with the action a is added to the prediction weight vector $cl.w$.

$$cl.p(x, a) = w \cdot x''$$

$$w = (w_0, w_1, \dots, w_n, w_{n+1}) \quad (15)$$

$$x'' = (x_0, x_1, \dots, x_n, a)$$

In summary, XCSFCA has a couple of computations: one for the action based on the input $a(x)$, the other one for the prediction based on the input and the result of the first computation $p(x, a(x))$.

3.2 Frog Problem

In this section, we summarize the frog problem introduced in [15]. A frog-like system learns to jump to catch a fly that is at a distance d far from the frog. The frog receives the signal input x related to d ($0 \leq d \leq 1$) and jumps a certain distance. A jump can be over or under the fly. Then the frog receives a payoff related to the remaining distance. Consider a jump as an action a that the system should choose. The frog problem is how to choose the best action a^* to maximize the payoff. Mathematically, the payoff function is given as follows:

$$P(x, a) = \begin{cases} x+a & \text{if } x+a \leq 1 \\ 2-(x+a) & \text{otherwise} \end{cases} \quad (16)$$

$$x(d) = 1-d \quad (17)$$

The payoff function is continuous and nonlinear composed of two linear planes. To facilitate the reading, we called the frog problem with the payoff function in Eq.(16) *frog1*. *frog2* is a modification of *frog1*. The payoff function is still continuous and nonlinear but composed of two nonlinear forms (Eq.(18)).

$$P(x, a) = \begin{cases} xe^a & \text{if } a \leq -\ln x \\ x^{-1}e^{-a} & \text{otherwise} \end{cases} \quad (18)$$

$$x(d) = e^{-d} \quad (19)$$

In both of *frog1* and *frog2* problems, the best value of a^* should be continuous with respect to x in order to be able to maximize the payoff.

3.3 Simulation Settings

In order to investigate the impact of computed action and facilitate the comparison of XCSFCA with XCSF in [16], the minimal change to parameter settings is made to the latter and the same measures measuring the system performance are used. The parameter settings used in our results were: $N=2000$, $\beta=0.2$, $\eta=0.2$, $\theta_{GA}=48$, $\chi=0.8$, $\mu=0.04$, $r_0=0.1$, $m_0=0.1$, $x_0=1$, $\varepsilon_0=0.01$, $\alpha=0.1$, $\nu=5$, $\theta_{del}=50$, $\delta=0.1$, $\delta_{ES}=0.1$, $\alpha_{GA}=0$, and initial standard deviation $\sigma^{(0)}=3$. $\sigma^{(0)}=3$ is chosen as the default value proposed by [2]. GA-subsumption and action set subsumption were not activated. Since the input x is one-dimensional, the vector length of action weights and standard deviations are $n_\sigma=2$. Explore and exploit problems are alternated with probability 0.5. One run is stopped after 100,000 explore problems. The payoff received after jumping is plotted. The system error measures the difference between the prediction of the expected payoff and the payoff received. The payoff and system error curves are plotted by using a 50-point running average from exploit problems, averaged over ten runs. In each problem, the fly was randomly placed at a distance d ($0 \leq d \leq 1$). The action range allowed is $[0,1]$. At the end of one run, x was scanned from 0 to 1 and from e^{-1} to 1 in *frog1* and *frog2* respectively, increased by 0.001, the best action a^* is plotted. The best action curve is plotted, averaged over ten runs. This curve exhibits the continuous actions related to the input.

3.4 Results

Since d is generated randomly from $[0,1]$, the sensory input x received by Eq.(17) gives the value from $[0,1]$. However, values of x in *frog2* received by Eq.(19) have the range $[e^{-1},1]$, they are scaled into $[0,1]$ when computing the prediction and the action so that the input and the action use the same range. In covering, given x , $cl.\zeta$ of a covered classifier is repeatedly generated until an action $cl.a$ calculated by Eq.(8) belongs to the action range allowed $[0,1]$.

frog1's results.

In Fig. 1, the system performance is greater than 99% after an averaged number of 30,000 explore problems. The system error drops to smaller than 1%. The population size of classifiers N is to about 37% of N . Fig. 2 shows the best action curve which is nearly coincident with the diagonal $1-x$. The system receives the input x (not the distance d). If the frog-like system wants to catch the fly at d , it must jump a distance $a=d$. Since x is given by Eq.(17), the best action, which the system needs to obtain, is the diagonal $a=1-x$. The averaged number of classifiers building the best action curve of Fig. 2a is 60.8 classifiers. Fig. 3 shows the error of *best action* a^* of Fig. 2b, which measures the difference between a^* and the diagonal $1-x$. In Fig. 5, the vectors of action weights (ζ_0, ζ_1) are close to $(-1,1)$, which implies that the best action curve of Fig. 2b is nearly coincident with the diagonal. The standard deviations σ_0 and σ_1 approach zero values and thus the values of ζ_0 and ζ_1 will be little updated by Eq.(10). For example, in the last classifier in Fig. 5, given $x = 0.999$, its action $cl.a$ and its prediction $cl.p$ are computed as follows: $cl.a = \zeta_0 x_0 + \zeta_1 x = 1.0006 - 1.0006 * 0.999 = 0.001$ and $cl.p = w_0 x_0 + w_1 x + w_2 a =$

$670.17 + 329.84 \cdot 0.999 + 329.38 \cdot 0.001 = 1000.0$. The results showed here are better than the ones of GCS [16] (Fig. 4 and Fig. 6).

The convergence speed is not fast because action weights are updated by mutation by small changes from one generation to another in the evolution strategy. In contrast to the *delta rule* method, the weights are updated by adding small changes in the direction that would reduce the error. Because of only one payoff used for evaluating one individual mutant action, the system should have strict replacement conditions to ensure that replacement will occur if the mutant is significantly better than its parent. This means a poor mutant could replace the effective classifier by mistake and then cause the system performance to be sometimes volatile. So, a good insurance is to make a reduction of the mutant's fitness by the factor δ_{ES} before updating it and use the accurate mutant in the replacement.

XCSFCA evolves thanks to the loop "reinforcement – ES – GA". This means a GA generates offspring that inherit the parameters of the parents, a reinforcement procedure corrects their parameters in the context of the current payoff, and their vectors of action weights will be updated through an ES procedure in order to produce actions more accurate. The GA will be triggered to generate new offspring. The loop of "reinforcement – ES – GA" continues.

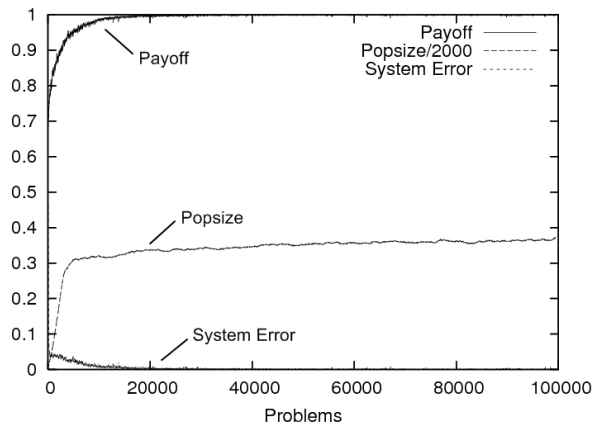
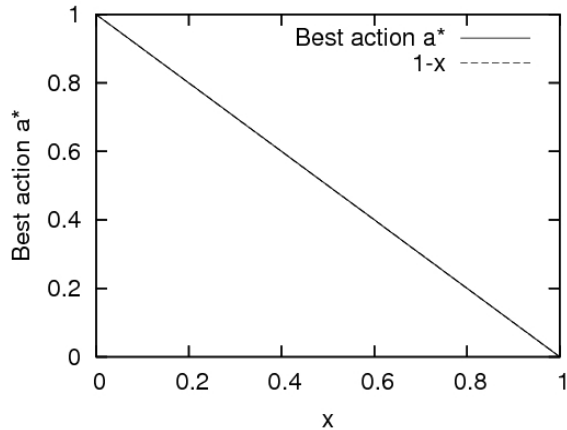
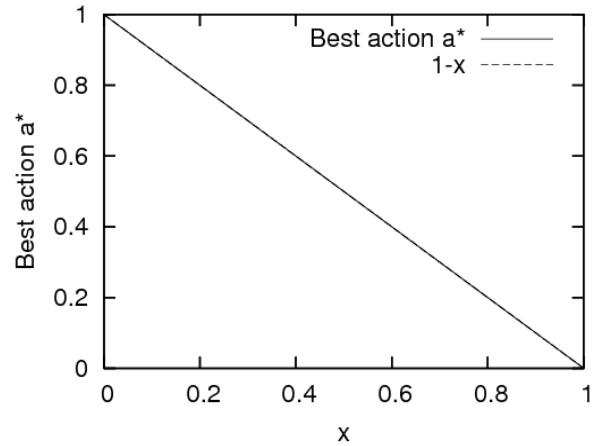


Figure 1: Results of *frog1* averaged over ten runs.



(2a) averaged over ten runs



(2b) from one run

Figure 2: *Best action a** of *frog1* is plotted by scanning the values of x from 0 to 1, increased by 0.001. The diagonal $1-x$ and *best action a** are nearly coincident. *Best action a** is slightly broken at some inputs x .

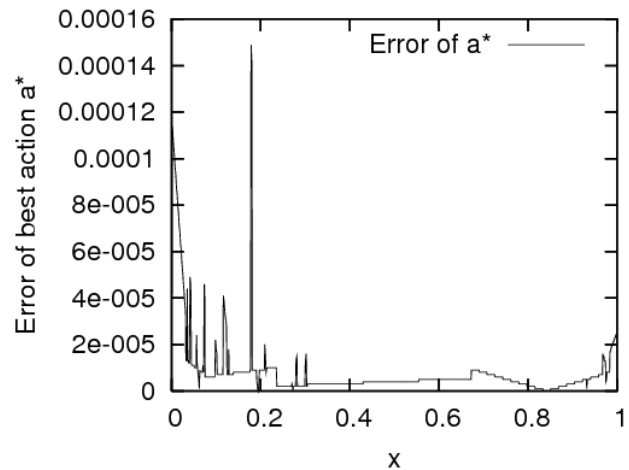


Figure 3: Error of *best action a** of Fig. 2b. The vertical axis is scaled to $[0..0.00016]$ to make visible.

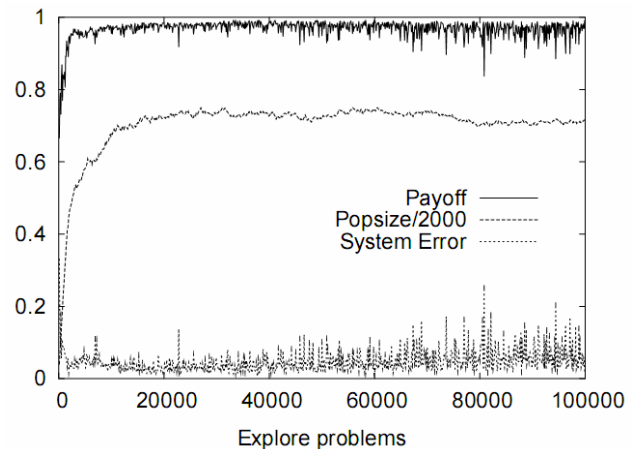


Figure 4: Results of *frog1* from one run in [16].

l_0	u_0	w_0	w_1	w_2	a	ζ_0	ζ_1	σ_0	σ_1	p	ε	fit	num
0.000	0.893	669.33	327.70	330.69	0.108	1.0000	-0.9995	0.0002	0.0009	997.5060	0.0000	1.0000	10.00
0.098	0.654	669.47	327.70	330.83	0.346	1.0002	-1.0022	0.0000	0.0001	997.8604	0.0000	0.9966	2.00
0.278	0.951	669.83	327.77	331.11	0.048	1.0005	-1.0018	0.0000	0.0000	997.3538	0.0000	0.5847	2.00
0.235	0.991	670.31	329.68	329.69	0.009	1.0000	-1.0000	0.0001	0.0016	999.9920	0.0000	1.0000	50.00
0.209	0.577	669.46	327.04	331.48	0.422	1.0009	-1.0040	0.0000	0.0000	997.8859	0.0000	0.9972	9.00
0.070	0.582	669.42	327.70	330.78	0.419	1.0001	-1.0010	0.0000	0.0009	998.2709	0.0000	0.4624	2.00
0.232	0.933	670.18	329.85	329.38	0.068	1.0005	-1.0006	0.0000	0.0000	999.9981	0.0000	0.3147	6.00
0.189	0.650	669.64	327.73	330.97	0.350	1.0003	-1.0014	0.0000	0.0000	998.2962	0.0000	0.6973	5.00
0.056	0.566	669.41	327.70	330.77	0.434	1.0001	-1.0017	0.0000	0.0003	998.1680	0.0000	0.9975	3.00
0.056	0.618	669.41	327.70	330.77	0.381	1.0001	-1.0021	0.0001	0.0004	997.8939	0.0000	0.9892	21.00
0.070	0.586	669.43	327.70	330.79	0.413	1.0001	-1.0015	0.0000	0.0001	998.1676	0.0000	0.7966	8.00
0.115	0.738	669.51	327.71	330.86	0.262	1.0003	-1.0020	0.0000	0.0001	997.6579	0.0000	0.5435	4.00
0.270	0.762	669.81	327.76	331.10	0.237	1.0005	-1.0017	0.0000	0.0001	998.0886	0.0000	0.7509	9.00
0.299	0.970	669.87	327.79	331.14	0.030	1.0005	-1.0017	0.0000	0.0000	997.3851	0.0000	0.4274	3.00
0.035	0.952	669.36	327.69	330.72	0.045	1.0001	-1.0030	0.0001	0.0006	996.3081	0.0000	0.2863	2.00
0.056	0.582	669.41	327.70	330.77	0.418	1.0001	-1.0015	0.0000	0.0001	998.1341	0.0000	0.5594	4.00
0.127	0.705	669.53	327.71	330.88	0.295	1.0002	-1.0021	0.0000	0.0001	997.7748	0.0000	0.8111	8.00
0.191	0.552	669.65	327.73	330.98	0.448	1.0003	-1.0014	0.0000	0.0000	998.6774	0.0000	0.7074	10.00
0.151	0.991	670.31	329.68	329.68	0.009	1.0000	-1.0000	0.0001	0.0006	999.9912	0.0000	0.4019	17.00
0.072	0.991	670.31	329.67	329.69	0.009	1.0000	-1.0000	0.0001	0.0006	999.9897	0.0000	0.7817	27.00
0.229	0.976	670.17	329.85	329.38	0.024	1.0006	-1.0006	0.0000	0.0000	999.9995	0.0000	0.3594	4.00
0.040	0.539	669.37	327.70	330.73	0.461	1.0001	-1.0015	0.0000	0.0002	998.2476	0.0000	0.8092	9.00
0.178	0.573	669.35	327.02	331.39	0.427	1.0009	-1.0041	0.0000	0.0000	997.7532	0.0001	0.6971	6.00
0.308	0.994	670.17	329.85	329.38	0.007	1.0006	-1.0006	0.0000	0.0000	1000.0013	0.0000	0.5493	5.00
0.032	0.991	670.30	329.67	329.68	0.009	1.0000	-1.0002	0.0000	0.0005	999.9342	0.0000	0.8500	12.00
0.229	0.975	670.17	329.84	329.38	0.026	1.0006	-1.0006	0.0000	0.0000	1000.0002	0.0000	0.2828	1.00
0.235	0.988	670.31	329.68	329.69	0.012	1.0000	-0.9999	0.0001	0.0004	999.9987	0.0000	0.6848	9.00
0.000	0.943	669.34	327.69	330.70	0.055	1.0001	-1.0027	0.0001	0.0015	996.4049	0.0001	0.2902	3.00
0.343	1.000	670.17	329.84	329.38	0.001	1.0006	-1.0006	0.0000	0.0000	1000.0005	0.0000	0.0501	1.00

Figure 5: Activated classifiers from one run of *frog1* of Fig. 2b. Values of the action and prediction columns are the values computed from the last input. $int_0=(l_0, u_0)$ is the interval predicate; w_0, w_1 and w_2 are the prediction weights for the constant x_0 , the input x and the action a , respectively; ζ_0 and ζ_1 are the action weights for x_0 and x , respectively; σ_0 and σ_1 are the standard deviations used by mutation on ζ_0 and ζ_1 , respectively; a, p, ε, fit and num indicate action, prediction, prediction error, fitness and numerosity, respectively.

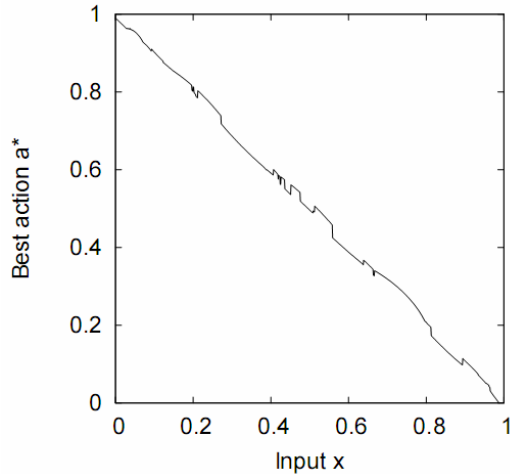


Figure 6: Best action a^* of *frog1* from one run in [16].

frog2's results.

frog2 is more difficult than *frog1* because the form of its payoff function is composed of two nonlinear tent-like forms on each side. But the system performance in Fig. 7 is still greater than 99%. The system error drops to about 1%. The population size is to about 37% of N . The best jump (action), which the system needs to obtain, is $a=-lnx$. The averaged number of classifiers building the best action curve of Fig. 8a is 73.6 classifiers. Fig. 9 shows the error of best action a^* of Fig. 8b. Fig. 10 shows activated classifiers with their main parameters that participated in the creation of best action a^* of Fig. 8b.

4. CONCLUSION

This paper presented an extension to XCSF, in which the action is computed directly as a linear combination of the input state and a vector of action weights. Thus, it is continuous when the input changes. The action weights are themselves adapted to the problem by using a separate process called an evolution strategy. The pure random exploration has been replaced by a deterministic strategy while the exploration of actions was guaranteed by the evolution strategy. XCSFCA has a couple of computations: one for the action and the other one for the prediction. We tested the performance of XCSFCA in the frog problems and the results proved to be more efficient than Wilson's architecture for continuous action.

XCSFCA needs more work. We proposed Algorithm 1 based on the evolution strategy to update action weights, although other methods can perform this update more efficiently. Thus, their integration into XCSFCA might improve the system performance further.

We will try to apply this new system to the prey – predator problem [8] to evolve weights controlling motor schemas [1]. This will be a promising solution to improve the learning capability of the learning predator further.

5. ACKNOWLEDGEMENTS

The authors would like to thank Stewart W. Wilson and the other anonymous reviewers for constructive comments on the submitted paper.

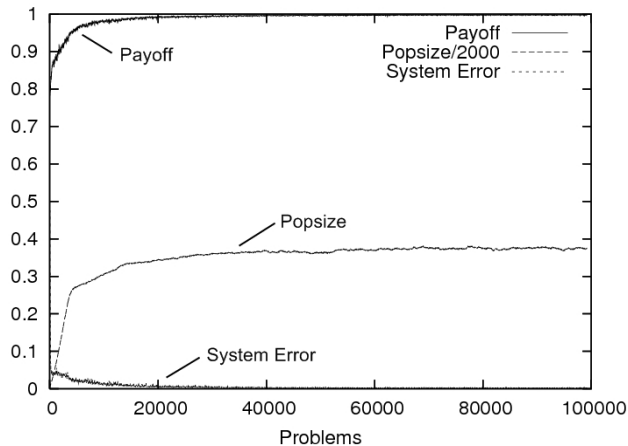
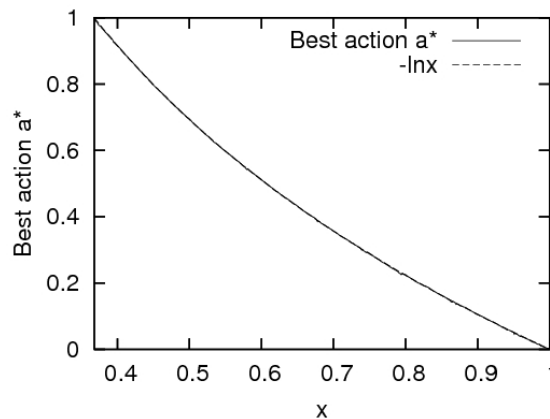


Figure 7: Results of *frog2* averaged over ten runs.

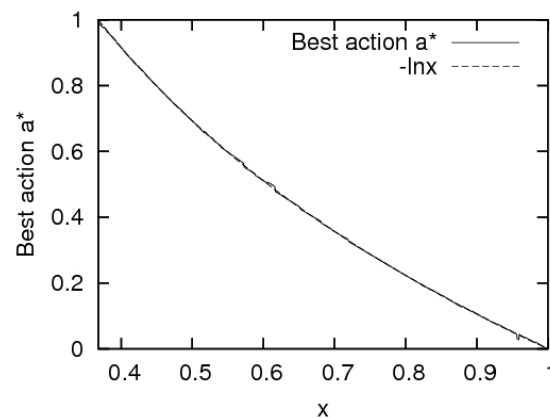
6. REFERENCES

- [1] Arkin, R. C., Behavior-Based Robot Navigation for Extended Domains, *Adaptive Behavior*, 1(2):201–225, 1992.
- [2] Bäck, T., *Evolutionary Algorithms in Theory and Practice*. Oxford University Press, New York, 1996.
- [3] Butz, M. V., and Wilson, S. W., An Algorithmic Description of XCS, *Soft Computing*, 6(3-4), pp. 144-153, 2002.
- [4] Lanzi, P. L., Loiacono, D., Wilson, S. W., and Goldberg, D. E., XCS with Computed Prediction for the Learning of Boolean Functions, *Proceedings of the IEEE Congress on Evolutionary Computation Conference (CEC2005)*, 2005.
- [5] Sanchez, S., *Mécanismes Evolutionnistes pour la Simulation Comportementale d'Acteurs Virtuels*, PhD thesis, pp. 93-94, University Toulouse I, IRIT (Institut de Recherche en Informatique de Toulouse), France, 2004.
- [6] Schwefel, H.-P., *Numerical Optimizaion of Computer Models*, Chichester: Wiley, 1981.
- [7] Sutton, R. S., Barto, A. G., Reinforcement Learning: An Introduction, MIT Press, 1998.
- [8] Tran, T. H., Sanza, C., and Duthen, Y., Evolving the Motor Schema Approach for Learning Cooperation, *CASA '2006 The nineteenth Conference on Computer Animation and Social Agents*, pp. 219-228, Genève, Suisse, 2006.
- [9] Widrow, B., and Hoff, M. E., Adaptive Switching Circuits, In James A. Anderson and Edward Rosenfeld, editors, *Neurocomputing: Foundations of Research*, pp. 126-134, The MIT Press, Cambridge, MA, 1988.
- [10] Wilson, S. W., ZCS: A Zeroth Level Classifier System, *Evolutionary Computation*, 2(1), pp. 1-18, 1994.
- [11] Wilson, S. W., Classifier Fitness Based on Accuracy, *Evolutionary Computation*, 3(2), pp. 149-175, 1995.
- [12] Wilson, S. W., Get Real! XCS with Continuous-Valued Inputs, From *Festschrift in Honor of John H. Holland*, L. Booker, S. Forrest, M. Mitchell, and R. Riolo (eds.), pp. 111-121, 1999.

- [13] Wilson, S. W., Function Approximation with a Classifier System, *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-2001)*, L. Spector et al, eds. Morgan Kaufmann, San Francisco, CA, pp. 974-981, 2001.



(8a) averaged over ten runs



(8b) from one run

Figure 8: *Best action a** of *frog2* is plotted by scanning the values of x from e^{-1} to 1, increased by 0.001. *Best action a** is close to the curve $-\ln x$.

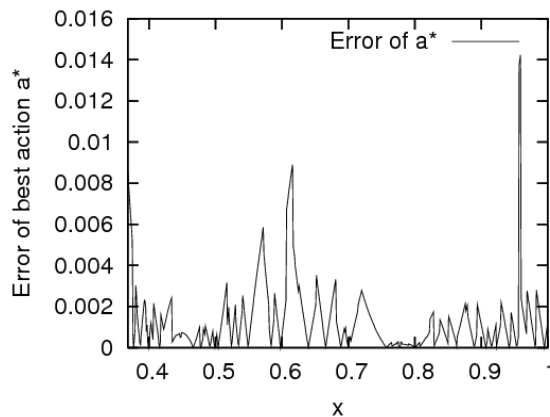


Figure 9: Error of *best action a** of Fig. 8b. The vertical axis is scaled to [0..0.0016] to make visible.

