

Reinforcement Learning for the Soccer Dribbling Task

Arthur Carvalho and Renato Oliveira

Abstract—We propose a reinforcement learning solution to the *soccer dribbling task*, a scenario in which a soccer agent has to go from the beginning to the end of a region keeping possession of the ball, as an adversary attempts to gain possession. While the adversary uses a stationary policy, the dribbler learns the best action to take at each decision point. After defining meaningful variables to represent the state space, and high-level macro-actions to incorporate domain knowledge, we describe our application of the reinforcement learning algorithm *Sarsa* with CMAC for function approximation. Our experiments show that, after the training period, the dribbler is able to accomplish its task against a strong adversary around 58% of the time.

I. INTRODUCTION

Soccer dribbling consists of the ability of a soccer agent to go from the beginning to the end of a region keeping possession of the ball, while an adversary attempts to gain possession. In this work, we focus on the dribbler’s learning process, *i.e.*, the learning of an effective policy that determines a good action for the dribbler to take at each decision point.

We study the soccer dribbling task using the RoboCup soccer simulator [1]. Specific details of this simulator increase the complexity of the learning process. For example, besides the adversarial and real-time environment, agents’ perceptions and actions are noisy and asynchronous.

We model the soccer dribbling task as a *reinforcement learning* problem. Our solution to this problem combines the Sarsa algorithm with CMAC for function approximation. Despite the fact that the resulting learning algorithm is not guaranteed to converge to the optimal policy in all cases, many lines of evidence suggest that it converges to near-optimal policies (for example, see [2], [3], [4], [5]).

Besides this introductory section, the rest of this paper is organized as follows. In the next section, we describe the soccer dribbling task. In Section 3, we show how to map this task onto an episodic reinforcement learning framework. In Section 4 and 5, we present, respectively, the reinforcement learning algorithm and its results against a strong adversary. In Section 6, we review the literature related to our work. In Section 7, we conclude and present future research directions.

II. SOCCER DRIBBLING

Soccer dribbling is a crucial skill for an agent to become a successful soccer player. It consists of the ability of a soccer agent, henceforth called the *dribbler*, to go from the beginning to the end of a region keeping possession of the ball, while an adversary attempts to gain possession. We can see soccer dribbling as a subproblem of the complete soccer domain. The main simplification is that the players involved are only focused on specific goals, without worrying about team strategies or unrelated individual skills (*e.g.*, passing and

shooting). Nevertheless, a successful policy learned by the dribbler can be used in the complete soccer domain whenever a soccer agent faces a dribbling situation.

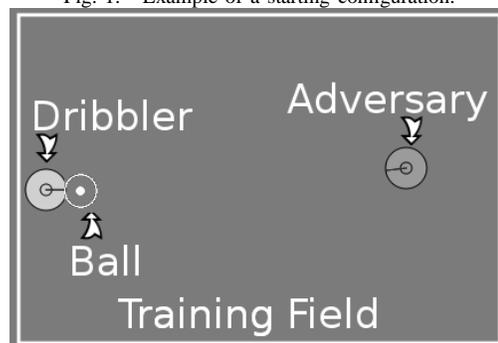
Since our focus is on the dribbler’s learning process, an omniscient coach agent is used to manage the play. At the beginning of each trial (*episode*), the coach resets the location of the ball and of the players within a *training field*. The dribbler is placed in the center-left region together with the ball. The adversary is placed in a random position with the constraint that it does not start with possession of the ball. An example of a starting configuration is shown in Figure 1.

Whenever the adversary gains possession for a set period of time or when the ball goes out of the training field by crossing either the left line or the top line or the bottom line, the coach declares the adversary as the winner of the episode. If the ball goes out of the training field by crossing the right line, then the winner is the first player to intercept the ball. After declaring the winner of an episode, the coach resets the location of the players and of the ball within the training field and starts a new episode. Thus, the dribbler’s goal is to reach the right line that delimits the training field with the ball. We call this task the *soccer dribbling task*.

We argue that the soccer dribbling task is an excellent benchmark for comparing different machine learning techniques since it involves a complex problem, and it has a well-defined objective, which is to maximize the number of episodes won by the dribbler. We study the soccer dribbling task using the RoboCup soccer simulator [1].

The RoboCup soccer simulator operates in discrete time steps, each representing 100 milliseconds of simulated time. Specific details of this simulator increase the complexity of the learning process. For example, random noise is injected into all perceptions and actions. Further, agents must sense and act asynchronously. Each soccer agent receives visual information about other objects every 150 milliseconds, *e.g.*, its distance

Fig. 1. Example of a starting configuration.



from other players in its current field of view. Each agent has also a body sensor, which detects its current “physical status” every 100 milliseconds, *e.g.*, that agent’s stamina and speed. Agents may execute a parameterized primitive action every 100 milliseconds, *e.g.*, *turn*(angle), *dash*(power), and *kick*(power, angle). Full details of the RoboCup soccer simulator are presented by Chen *et al.* [6].

Since possession is not well-defined in the RoboCup soccer simulator, we consider that an agent has possession of the ball whenever the ball is close enough to be kicked, *i.e.*, it is in a distance less than 1.085 meters from the agent.

III. THE SOCCER DRIBBLING TASK AS A REINFORCEMENT LEARNING PROBLEM

In the soccer dribbling task, an episode begins when the dribbler may take the first action. When an episode ends (*e.g.*, when the adversary gains possession for a set period of time), the coach starts a new one, thereby giving rise to a series of episodes. Thus, the interaction between the dribbler and the environment naturally breaks down into a sequence of distinct episodes. This point, together with the fact that the RoboCup soccer simulator operates in discrete time steps, allows the soccer dribbling task to be mapped onto a discrete-time, episodic reinforcement-learning framework.

Roughly speaking, reinforcement learning is concerned with how an agent must take actions in an environment so as to maximize the expected long-term reward [7]. Like in a trial-and-error search, the learner must discover which action is the most rewarding one in a given state of the world. Thus, solving a reinforcement learning problem means finding a function (*policy*) that maps *states* to *actions* so that it maximizes a *reward* over the long run. As a way of incorporating domain knowledge, the actions available to the dribbler are the following high-level *macro-actions*, which are built on the simulator’s *primitive actions*¹:

- *HoldBall*(\cdot): The dribbler holds the ball close to its body, keeping it in a position that is difficult for the adversary to gain possession;
- *Dribble*(Θ, k): The dribbler turns its body towards the global angle Θ , kicks the ball k meters ahead of it, and moves to intercept the ball.

The global angle Θ is in the range $[0, 360]$. In detail, the center of the training field has been chosen as the origin of the system, where the zero-angle points towards the middle of the right line that delimits the training field, and it increases in the clockwise direction. Those macro-actions are based on high-level skills used by the UvA Trilearn 2003 team [8]. The first one maps directly onto the primitive action *kick*. Consequently, it usually takes a single time step to be performed. The second one, however, requires an extended sequence of the primitive actions *turn*, *kick*, and *dash*. To handle this situation, we treat the soccer dribbling task as a *semi-Markov decision process* (SMDP) [9].

¹Henceforth, we use the terms *action* and *macro-action* interchangeably, while always distinguishing *primitive actions*.

Formally, an SMDP is a 5-tuple $\langle S, A, P, r, F \rangle$, where S is a countable set of states, A is a countable set of actions, $P(s'|s, a)$, for $s', s \in S$, and $a \in A$, is a probability distribution providing the transition model between states, $r(s, a) \in \mathfrak{R}$ is a reward associated with the transition (s, a) , and $F(\tau|s, a)$ is a probability distribution indicating the sojourn time in a given state $s \in S$, *i.e.*, the time before transition provided that action a was taken in state s .

Let $a_i \in A$ be the i th macro-action selected by the dribbler. Thus, several simulator’s time steps may elapse between a_i and a_{i+1} . Let $s_{i+1} \in S$ and $r_{i+1} \in \mathfrak{R}$ be, respectively, the state and the reward following the macro-action a_i . From the dribbler’s point of view, an episode consists of a sequence of SMDP steps, *i.e.*, a sequence of states, macro-actions, and rewards: $s_0, a_0, r_1, s_1, \dots, s_i, a_i, r_{i+1}, s_{i+1}, \dots, a_{n-1}, r_n, s_n$, where a_i is chosen based exclusively on the state s_i , and s_n is a terminal state in which either the adversary or the dribbler is declared the winner of the episode by the coach. In the former case, the dribbler receives the reward $r_n = -1$, while in the latter case its reward is $r_n = 1$. The intermediate rewards are always equal to zero, *i.e.*, $r_1 = r_2 = \dots = r_{n-1} = 0$. Thus, our objective is to find a policy that maximizes the dribbler’s reward, *i.e.*, the number of episodes in which it is the winner.

A. Dribbler

The dribbler must take a decision at each SMDP step by selecting an available macro-action. Besides the macro-action *HoldBall*, the set of actions available to the dribbler contains four instances of the macro-action *Dribble*: *Dribble*($30^\circ, 5$), *Dribble*($330^\circ, 5$), *Dribble*($0^\circ, 5$), and *Dribble*($0^\circ, 10$). Thus, besides hiding the ball from the adversary, the dribbler can kick the ball forward (strongly and weakly), diagonally upward, and diagonally downward. If at some time step the dribbler has not possession of the ball and the current state is not a terminal state, then it usually means that the dribbler chose an instance of the macro-action *Dribble* before and it is currently moving to intercept the ball.

We turn now to the state representation used by the dribbler. It consists of a set of state variables which are based on information related to the ball, the adversary, and the dribbler itself. Let $ang(x)$ be the global angle of the object x , and $ang(x, y)$ and $dist(x, y)$ be, respectively, the relative angle and the distance between the objects x and y . Further, let w and h be, respectively, the width and the height of the training field. Finally, let $posY(x)$ be a function indicating whether the object x is close to (less than 1 meter away from) the top line or the bottom line that delimits the training field. In the former case, $posY(x) = 1$, whereas in the latter case $posY(x) = -1$, and otherwise $posY(x) = 0$. Table 1 shows the state variables together with their ranges.

The first three variables help the dribbler to locate itself and the adversary inside the training field. Together, the last two variables can be seen as a point describing the position of the adversary in a polar coordinate system, where the ball is the pole. Thus, these variables are used by the dribbler to locate the adversary with respect to the ball. It is interesting to

TABLE I
DESCRIPTION OF THE STATE REPRESENTATION.

State Variable	Range
$posY(\text{dribbler})$	$\{-1, 0, 1\}$
$ang(\text{dribbler})$	$[0, 360]$
$ang(\text{dribbler}, \text{adversary})$	$[0, 360]$
$ang(\text{ball}, \text{adversary})$	$[0, 360]$
$dist(\text{ball}, \text{adversary})$	$[0, \sqrt{w^2 + h^2}]$

note that a more informative state representation can be used by adding more state variables, *e.g.*, the current speed of the ball and the dribbler’s stamina. However, large domains can be impractical due to the “curse of dimensionality”, *i.e.*, the general tendency of the state space to grow exponentially in the number of state variables [10]. Consequently, we focus on a state representation that is as concise as possible.

B. Adversary

The adversary uses a fixed, pre-specified policy. Thus, we can see it as part of the environment in which the dribbler is interacting with. When the adversary has possession of the ball, it tries to maintain possession for another time step by invoking the macro-action *HoldBall*. If it maintains possession for two consecutive time steps, then it is the winner of the episode. When the adversary does not have the ball, it uses an iterative scheme to compute a near-optimal interception point based on the ball’s position and velocity. Thereafter, the adversary moves to that point as fast as possible. This procedure is the same used by the dribbler when it is moving to intercept the ball after invoking the macro-action *Dribble*. More details about this iterative scheme can be found in the description of the UvA Trilearn 2003 team [8].

IV. THE REINFORCEMENT LEARNING ALGORITHM

Our solution to the soccer dribbling task combines the reinforcement learning algorithm Sarsa with CMAC for function approximation. In what follows, we briefly introduce both of them before presenting the final learning algorithm.

A. Sarsa

The Sarsa algorithm works by estimating the action-value function $Q^\pi(s, a)$, for the current policy π and for all state-action pairs (s, a) [7]. The Q -function assigns to each state-action pair the expected return from it. Given a quintuple of events, $(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})$, that make up the transition from the state-action pair (s_t, a_t) to the next one, (s_{t+1}, a_{t+1}) , the Q -value of the first state-action pair is updated according to the following equation:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \delta_t, \quad (1)$$

where δ_t is the traditional *temporal-difference error*,

$$\delta_t = r_{t+1} + \lambda Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t), \quad (2)$$

α is the learning rate parameter, and λ is a discount rate governing the weight placed on future, as opposed to immediate, rewards. Sarsa is an on-policy learning method, meaning that

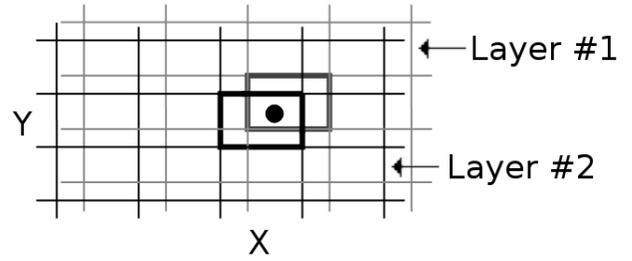


Fig. 2. Example of two layers overlaid over a two-dimensional state space. Any input vector (state) activates two receptive fields, one from each layer. For example, the state represented by the black dot activates the highlighted receptive fields.

it continually estimates Q^π , for the current policy π , and at the same time changes π towards greediness with respect to Q^π . A typical policy derived from the Q -function is an ϵ -greedy policy. Given the state s_t , this policy selects a random action with probability ϵ and, otherwise, it selects the action with the highest estimated value, *i.e.*, $a = \operatorname{argmax}_a Q(s_t, a)$.

B. CMAC

In tasks with a small number of state-action pairs, we can represent the action-value function Q^π as a table with one entry for each state-action pair. However, this is not the case of the soccer dribbling task. For illustration’s sake, suppose that all variables in Table 1 are discrete. If we consider the 5 actions available to the dribbler and a 20m x 20m training field, we end up with more than 1.9×10^{10} state-action pairs. This would not only require an unreasonable amount of memory, but also an enormous amount of data to fill up the table accurately. Thus, we need to *generalize* from previously experienced states to ones that have never been seen. For dealing with this task, we use a technique commonly known as *function approximation*.

By using a function approximation, the action-value function Q^π is now represented as a parameterized functional form [7]. Now, whenever we make a change in one parameter value, we also change the estimated value of many state-action pairs, thus obtaining generalization. In this work, we use the Cerebellar Model Arithmetic Computer (CMAC) for function approximation [11], [12].

CMAC works by partitioning the state space into multi-dimensional *receptive fields*, each of which is associated with a *weight*. In this work, receptive fields are hyper-rectangles in the state space. Nearby states share receptive fields. Thus, generalization occurs between them. Multiple partitions of the state space (*layers*) are usually used, which implies that any input vector falls within the range of multiple *excited receptive fields*, one from each layer.

Layers are identical in organization, but each one is offset relative to the others so that each layer cuts the state space in a different way. By overlapping multiple layers, it is possible to achieve quick generalization while maintaining the ability to learn fine distinctions. Figure 2 shows an example of two grid-like layers overlaid over a two-dimensional space.

The receptive fields excited by a given state s make up the

feature set \mathbb{F}_s , with each action a indexing their weights in a different way. In other words, each macro-action is associated with a particular CMAC. Clearly, the number of receptive fields inside each feature set is equal to the number of layers. The CMAC’s response to a feature set \mathbb{F}_s is equal to the sum of the weights of the receptive fields in \mathbb{F}_s . Formally, let $\theta_a(i)$ be the weight of the receptive field i indexed by the action a . Thus, the CMAC’s response to \mathbb{F}_s is equal to $\sum_{i \in \mathbb{F}_s} \theta_a(i)$, which represents the Q -value $Q(s, a)$.

CMAC is trained by using the traditional *delta rule* (also known as the least mean square). In detail, after selecting an action a , the weight of an excited receptive field i indexed by a , $\theta_a(i)$, is updated according to the following equation:

$$\theta_a(i) \leftarrow \theta_a(i) + \alpha\delta, \quad (3)$$

where δ is the temporal-difference error. A major issue when using CMAC is that the total number of receptive fields required to span the entire state space can be very large. Consequently, an unreasonable amount of memory may be needed. A technique commonly used to address this issue is called *pseudo-random hashing* [7]. It produces receptive fields consisting of noncontiguous, disjoint regions randomly spread throughout the state space, so that only information about receptive fields that have been excited during previous training is actually stored.

C. Linear, Gradient-Descent Sarsa

Our solution to the soccer dribbling task combines the Sarsa algorithm with CMAC for function approximation. We use an ϵ -greedy policy for action selection. Sutton and Barto [7] provide a complete description of this algorithm under the name of *linear, gradient-descent Sarsa*. Our implementation follows the solution proposed by Stone *et al.* [13]. It consists of three routines: *RLstartEpisode*, to be run by the dribbler at the beginning of each episode; *RLstep*, run on each SMDP step; and *RLendEpisode*, to be run when an episode ends. In what follows, we present each routine in detail.

1) *RLstartEpisode*: Given an initial state s_0 , this routine starts by iterating over all available actions. In line 2, it finds the receptive fields excited by s_0 , which compose the feature set \mathbb{F}_{s_0} . Next, in line 3, the estimated value of each macro-action a in s_0 is calculated as the sum of the weights of the excited receptive fields. In line 5, this routine selects a macro-action by following an ϵ -greedy policy and sends it to the RoboCup soccer simulator. Finally, the chosen action and the initial state s_0 are stored, respectively, in the variables *LastAction* and *LastState*.

Algorithm 1 RLstartEpisode

```

1: for each action  $a$  do
2:    $\mathbb{F}_{s_0} \leftarrow$  receptive fields excited by  $s_0$ 
3:    $Q_a \leftarrow \sum_{i \in \mathbb{F}_{s_0}} \theta_a(i)$ 
4: end for
5:  $LastAction \leftarrow \begin{cases} \operatorname{argmax}_a Q_a & \text{w/ prob. } 1 - \epsilon \\ \text{random action} & \text{w/ prob. } \epsilon \end{cases}$ 
6:  $LastState \leftarrow s_0$ 

```

2) *RLstep*: This routine is run on each SMDP step, whenever the dribbler has to choose a macro-action. Given the current state s , it starts by calculating part of the temporal-difference error (Equation 2), namely the difference between the intermediate reward r and the expected return of the previous SMDP step, $Q_{LastAction}$. In lines 2 to 5, this routine finds the receptive fields excited by s and uses their weights to compute the estimated value of each action a in s . In line 6, the next action to be taken by the dribbler is selected according to an ϵ -greedy policy. In line 7, this routine finishes to compute the temporal-difference error by adding the discount rate λ times the expected return of the current SMDP step, $Q_{CurrentAction}$. Next, in lines 8 to 10, this routine adjusts the weights of the receptive fields excited in the previous SMDP step by the learning factor α times the temporal-difference error δ (see Equation 3). Since the weights have changed, we must recalculate the expected return of the current SMDP step, $Q_{CurrentAction}$ (line 11). Finally, the chosen action and the current state are stored, respectively, in the variables *LastAction* and *LastState*.

Algorithm 2 RLstep

```

1:  $\delta \leftarrow r - Q_{LastAction}$ 
2: for each action  $a$  do
3:    $\mathbb{F}_s \leftarrow$  receptive fields excited by  $s$ 
4:    $Q_a \leftarrow \sum_{i \in \mathbb{F}_s} \theta_a(i)$ 
5: end for
6:  $CurrentAction \leftarrow \begin{cases} \operatorname{argmax}_a Q_a & \text{w/ prob. } 1 - \epsilon \\ \text{random action} & \text{w/ prob. } \epsilon \end{cases}$ 
7:  $\delta \leftarrow \delta + \lambda Q_{CurrentAction}$ 
8: for each  $i \in \mathbb{F}_{LastState}$  do
9:    $\theta_{LastAction}(i) \leftarrow \theta_{LastAction}(i) + \alpha\delta$ 
10: end for
11:  $Q_{CurrentAction} \leftarrow \sum_{i \in \mathbb{F}_s} \theta_{CurrentAction}(i)$ 
12:  $LastAction \leftarrow CurrentAction$ 
13:  $LastState \leftarrow s$ 

```

3) *RLendEpisode*: This routine is run when an episode ends. Initially, it calculates the appropriate reward based on who won the episode. Next, it calculates the temporal-difference error in the action-value estimates (line 6). There is no need to add the expected return of the current SMDP step ($Q_{CurrentAction}$) since this value is defined to be 0 for terminal states. Lastly, this routine adjusts the weights of the receptive fields excited in the previous SMDP step.

Algorithm 3 RLendEpisode

```

1: if the dribbler is the winner then
2:    $r \leftarrow 1$ 
3: else
4:    $r \leftarrow -1$ 
5: end if
6:  $\delta \leftarrow r - Q_{LastAction}$ 
7: for each  $i \in \mathbb{F}_{LastState}$  do
8:    $\theta_{LastAction}(i) \leftarrow \theta_{LastAction}(i) + \alpha\delta$ 
9: end for

```

V. EMPIRICAL RESULTS

In this section, we report our experimental results with the soccer dribbler task. In all experiments, we used the standard RoboCup soccer simulator (version 14.0.3, protocol 9.3) and a 20m x 20m training region. In that simulator, agents typically have limited and noisy visual sensors. For example, each player can see objects within a 90° view cone, and the precision of an object’s sensed location degrades with distance. To simplify the learning process, we removed those restrictions. Both the dribbler and the adversary were given 360° of noiseless vision to ensure that they would always have complete and accurate knowledge of the environment.

Related to parameters of the reinforcement learning algorithm², we set $\epsilon = 0.01$, $\alpha = 0.125$, and $\lambda = 1$. By no means do we argue that these values are optimal. They were set based on results of brief, informal experiments.

The weights of first-time excited receptive fields were set to 0. The bounds of the receptive fields were set according to the generalization that we desired: angles were given widths of about 20 degrees, and distances were given widths of approximately 3 meters. We used 32 layers. Each dimension of every layer was offset from the others by $1/32$ of the desired width in that dimension. We used the CMAC implementation proposed by Miller and Glanz [14], which uses pseudo-random hashing. To retain previously trained information in the presence of subsequent novel data, we did not allow hash collisions.

To create episodes as realistic as possible, agents were not allowed to recover their staminas by themselves. This task was done by the coach after five consecutive episodes. This enabled agents to start episodes with different stamina values. We ran this experiment 5 independent times, each one lasting 50,000 episodes, and taking, on average, approximately 74 hours. Figure 3 shows the histogram of the average number of episodes won by the dribbler during the training process. Bins of 500 episodes were used.

Throughout the training process, the dribbler won, on average, 23,607 episodes ($\approx 47\%$). From Figure 3, we can see that it greatly improves its average performance as the number of episodes increases. At the end of the training process, it is winning slightly less than 53% of the time.

Qualitatively, the dribbler seems to learn two major rules. In the first one, when the adversary is at a considerable distance, the dribbler keeps kicking the ball to the opposite side in which the adversary is located until the angle between them is in the range $[90, 270]$, *i.e.*, when the adversary is behind the dribbler. After that, the dribbler starts to kick the ball forward. An illustration of this rule can be seen in Figure 4.

The second rule seems to occur when the adversary is relatively close to and in front of the dribbler. Since there is no way for the dribbler to move forward or diagonally without putting the possession at risk, it then holds the ball until the

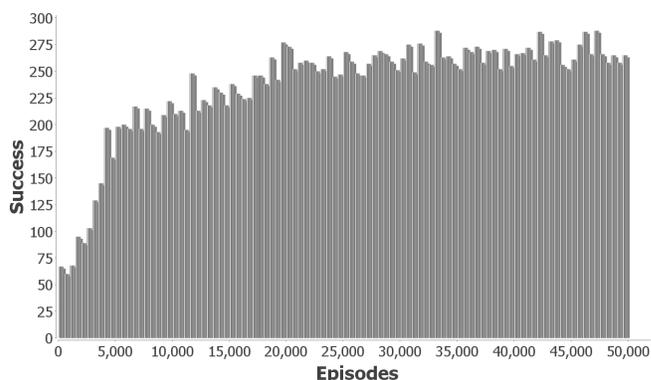


Fig. 3. Histogram of the average number of episodes won by the dribbler (success) during the training process (50,000 episodes). Bins of 500 episodes were used, and 5 independent simulations were performed.

angle between it and the adversary is in the range $[90, 270]$. Thereafter, it starts to advance by kicking the ball forward. An illustration of this rule can be seen in Figure 5.

After the training process, we randomly generated 10,000 initial configurations to test our solution. This time, the dribbler always selected the macro-action with the highest estimated value, *i.e.*, we set $\epsilon = 0$. Further, the weights of the receptive fields were not updated, *i.e.*, we set $\alpha = 0$. We used the receptive fields’ weights resulting from the simulation where the dribbler obtained the highest success rate. The result of this experiment was even better. The dribbler won 5,795 episodes, thus obtaining a success rate of approximately 58%.

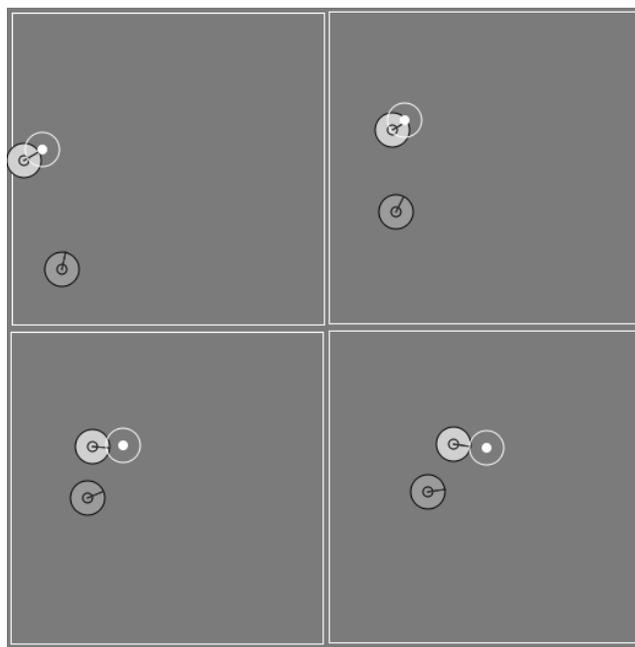


Fig. 4. Example of the first major rule learned by the dribbler. (Top Left) The adversary is at a considerable distance from the dribbler. (Top Right) The dribbler starts to kick the ball to the opposite side in which the adversary is located. (Bottom Left) The angle between the adversary and the dribbler is in the range $[90, 270]$. Consequently, the dribbler starts to kick the ball forward. (Bottom Right) The dribbler keeps kicking the ball forward.

²The implementation of the learning algorithm can be found at: <http://sites.google.com/site/soccerdribbling/>

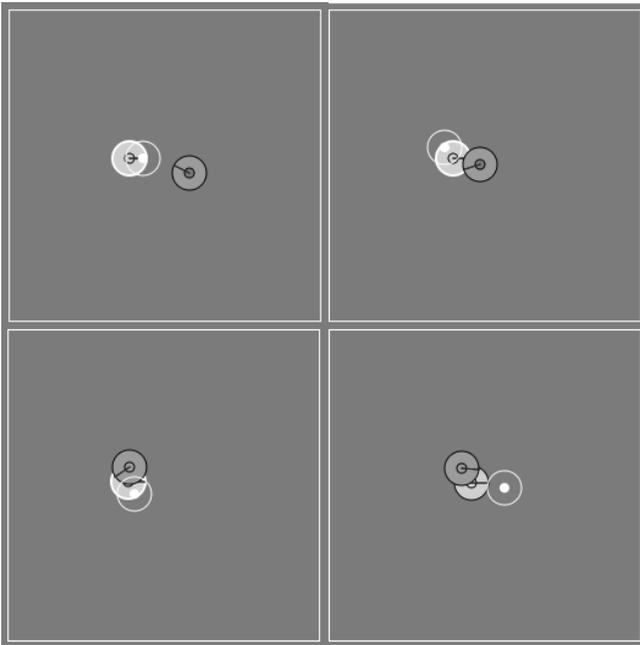


Fig. 5. Example of the second major rule learned by the dribbler. (Top Left) The adversary is close to and in front of the dribbler. (Top Right) The dribbler holds the ball so as not to lose possession. (Bottom Left) The dribbler keeps holding the ball. (Bottom Right) The angle between the adversary and the dribbler is the range $[90, 270]$. Consequently, the dribbler starts to advance by kicking the ball forward.

A. One-Dimensional CMACs

For comparison’s sake, we repeated the above experiment using the original solution proposed by Stone *et al.* [13]. It consists of the same learning algorithm presented in Section 3, but using one-dimensional CMACs. In detail, each layer is an interval along a state variable. In this way, the feature set \mathbb{F}_s is now composed by $32 \times 5 = 160$ excited receptive fields, *i.e.*, 32 excited receptive fields for each state variable.

One of the main advantages of using one-dimensional CMACs is that it is possible to circumvent the curse of dimensionality. In detail, the state space does not grow exponentially in the number of state variables because dependence between variables is not taken into account.

Figure 6 shows the histogram of the average number of episodes won by the dribbler during the training process. Each simulation took, on average, approximately 43 hours. Throughout the training process, the dribbler won, on average, 16,278 episodes ($\approx 33\%$). From Figure 6, we can see that the learning algorithm converges much faster when using one-dimensional CMACs. However, its average performance is considerably worse. At the end of the training process, the dribbler is winning, on average, less than 30% of the time.

After the training process, we tested this solution using the same 10,000 initial configurations previously generated. Again, we set $\epsilon = \alpha = 0$, and used the receptive fields’ weights resulting from the simulation where the dribbler obtained the highest success rate. The result of this experiment was slightly better. The dribbler won 3,701 episodes, thus

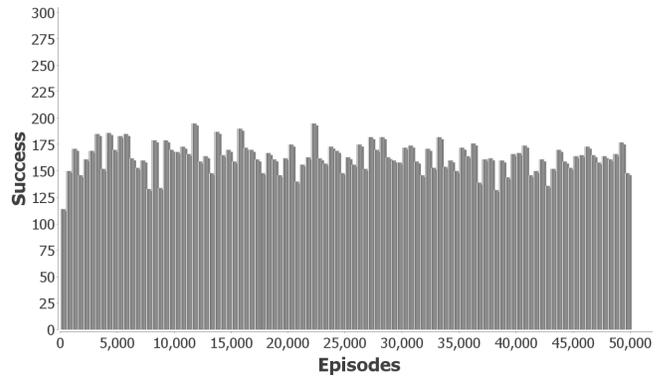


Fig. 6. Histogram of the average number of episodes won by the dribbler (success) during the training process (50,000 episodes) when using one-dimensional CMACs. Bins of 500 episodes were used, and 5 independent simulations were performed.

obtaining a success rate of approximately 37%.

Qualitatively, the dribbler seems to learn a rule similar to the one shown in Figure 4. The major difference is that it always kicks the ball to the opposite side in which the adversary is located, it does not matter its distance from the adversary’s location. Consequently, it is highly unlikely that the dribbler succeeds when the adversary is close to it.

We conjecture that one of the main reasons for such a poor performance of the reinforcement learning algorithm when using one-dimensional CMACs is that it does not take into account dependence between variables, *i.e.*, they are treated individually. Hence, such approach may throw away valuable information. For example, the variables $ang(\text{ball}, \text{adversary})$ and $dist(\text{ball}, \text{adversary})$ together describe the position of adversary with respect to the ball. However, they do not make as much sense when considered individually.

VI. RELATED WORK

Reinforcement learning has long been applied to the robot soccer domain. For example, Andou [15] uses “observational reinforcement learning” to refine a function that is used by the soccer agents for deciding their positions on the field. Riedmiller *et al.* [16] use reinforcement learning to learn low-level soccer skills, such as kicking and ball-interception. Nakashima *et al.* [17] propose a reinforcement learning method called “fuzzy Q-learning”, where an agent determines its action based on the inference result of a fuzzy rule-based system. The authors apply the proposed method to the scenario where a soccer agent learns to intercept a passed ball.

Arguably, the most successful application is due to Stone *et al.* [13]. They propose the “keepaway task”, which consists of two teams, the keepers and the takers, where the former tries to keep control of the ball for as long as possible, while the latter tries to gain possession. Our solution to the soccer dribbling task follows closely the solution proposed by those authors to learn the keepers’ behavior. Iscen and Eroglu [18] use similar solution to learn a policy for the takers.

Gabel *et al.* [19] propose a task which is the opposite of the soccer dribbling task, where a defensive player must interfere

and disturb the opponent that has possession of the ball. Their solution to that task uses a reinforcement learning algorithm with a multilayer neural network for function approximation.

Kalyanakrishnan *et al.* [20] present the “half-field offense task”, a scenario in which an offense team attempts to outplay a defense team in order to shoot goals. Those authors pose that task as a reinforcement learning problem, and propose a new learning algorithm for dealing with it.

More closely related to our work are reinforcement learning-based solutions to the task of conducting the ball (*e.g.*, [21]), which can be seen as a simplification of the dribbling task since it usually does not include adversaries.

VII. CONCLUSION

We proposed a reinforcement learning solution to the soccer dribbling task, a scenario in which an agent has to go from the beginning to the end of a region keeping possession of the ball, while an adversary attempts to gain possession. Our solution combined the Sarsa algorithm with CMAC for function approximation. Empirical results showed that, after the training period, the dribbler was able to accomplish its task against a strong adversary around 58% of the time.

Although we restricted ourselves to the soccer domain, dribbling, as defined in this paper, is also common in other sports, *e.g.*, hockey, basketball, and football. Thus, the proposed solution can be of value to dribbling tasks of other sports games. Furthermore, we believe that the soccer dribbling task is an excellent benchmark for comparing different machine learning techniques because it involves a complex problem, and it has a well-defined objective.

There are several exciting directions for extending this work. From a practical perspective, we intend to analyze the scalability of our solution, *i.e.*, to study how it performs with training fields of distinct sizes and against different adversaries. Further, we are considering schemes to extend our solution to the original partially observable environment, where the available information is incomplete and noisy.

As stated before, a more informative state representation could be obtained by using more state variables. The major problem of adding extra variables to our solution is that CMAC’s complexity increases exponentially with its dimensionality. Due to this fact, we are considering other solutions which use function approximations whose complexity is unaffected by dimensionality *per se*, *e.g.*, the Kanerva coding (for example, see Kostiadis and Hu’s work [22]).

Finally, we note that when modeling the soccer dribbling task as a reinforcement learning problem, we do not directly use intermediate rewards (they are all set to zero). However, they may make the learning process more efficient (for example, see [23]). Thus, we intend to investigate the influence of intermediate rewards on the final solution in future work.

ACKNOWLEDGMENTS

We would like to thank W. Thomas Miller, Filson H. Glanz, and others from the Department of Electrical and Computer Engineering at the University of New Hampshire for making their CMAC code available.

REFERENCES

- [1] I. Noda, H. Matsubara, K. Hiraki, and I. Frank, “Soccer server: A tool for research on multiagent systems,” *Applied Artificial Intelligence*, vol. 12, no. 2, pp. 233–250, 1998.
- [2] G. J. Gordon, “Reinforcement learning with function approximation converges to a region,” in *Advances in Neural Information Processing Systems 13*, 2001, pp. 1040–1046.
- [3] R. S. Sutton, “Generalization in reinforcement learning: Successful examples using sparse coarse coding,” in *Advances in Neural Information Processing Systems 8*, 1996, pp. 1038–1044.
- [4] J. N. Tsitsiklis and B. V. Roy, “An analysis of temporal-difference learning with function approximation,” *IEEE Transactions on Automatic Control*, vol. 42, no. 5, pp. 674–690, 1997.
- [5] T. J. Perkins and D. Precup, “A convergent form of approximate policy iteration,” in *Advances in Neural Information Processing Systems 15*, 2003, pp. 1595–1602.
- [6] M. Chen, E. Foroughi, F. Heintz, S. Kapetanakis, K. Kostiadis, J. Kummeneje, I. Noda, O. Obst, P. Riley, T. Steffens, Y. Wang, and X. Yin, *Users manual: RoboCup soccer server manual for soccer server version 7.07 and later*, 2003, available at <http://sourceforge.net/projects/sserver/>.
- [7] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. The MIT Press, 1998.
- [8] R. de Boer and J. R. Kok, “The incremental development of a synthetic multi-agent system: the uva trilearn 2001 robotic soccer simulation team,” Master’s thesis, University of Amsterdam, The Netherlands, 2002.
- [9] M. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley, 2005.
- [10] R. Bellman, *Dynamic Programming*. Dover Publications, 2003.
- [11] J. S. Albus, “A Theory of Cerebellar Function,” *Mathematical Biosciences*, vol. 10, no. 1-2, pp. 25–61, 1971.
- [12] —, *Brain, Behavior, and Robotics*. Byte Books, 1981.
- [13] P. Stone, R. S. Sutton, and G. Kuhlmann, “Reinforcement Learning for RoboCup-Soccer Keepaway,” *Adaptive Behavior*, vol. 13, no. 3, pp. 165–188, 2005.
- [14] W. T. Miller and F. H. Glanz, *UNH_CMVAC Version 2.1: The University of New Hampshire implementation of the Cerebellar model arithmetic computer - CMAC*, 1994, available at <http://www.ece.unh.edu/robots/cmavac.htm>.
- [15] T. Andou, “Refinement of Soccer Agents’ Positions Using Reinforcement Learning,” in *RoboCup-97: Robot Soccer World Cup I*, 1998, pp. 373–388.
- [16] M. A. Riedmiller, A. Merke, D. Meier, A. Hoffman, A. Sinner, O. Thate, and R. Ehrmann, “Karlsruhe Brainstormers - A Reinforcement Learning Approach to Robotic Soccer,” in *RoboCup 2000: Robot Soccer World Cup IV*, 2001, pp. 367–372.
- [17] T. Nakashima, M. Udo, and H. Ishibuchi, “A fuzzy reinforcement learning for a ball interception problem,” in *RoboCup 2003: Robot Soccer World Cup VII*, 2004, pp. 559–567.
- [18] A. Iscen and U. Erogul, “A new perspective to the keepaway soccer: the takers,” in *Proceedings of the 7th international joint conference on Autonomous agents and multiagent systems*, 2008, pp. 1341–1344.
- [19] T. Gabel, M. Riedmiller, and F. Trost, “A Case Study on Improving Defense Behavior in Soccer Simulation 2D: The NeuroHassle Approach,” in *RoboCup-2008: Robot Soccer World Cup XII*, 2009, pp. 61–72.
- [20] S. Kalyanakrishnan, Y. Liu, and P. Stone, “Half field offense in RoboCup soccer: A multiagent reinforcement learning case study,” in *RoboCup-2006: Robot Soccer World Cup X*, 2007, pp. 72–85.
- [21] M. Riedmiller, R. Hafner, S. Lange, and M. Lauer, “Learning to dribble on a real robot by success and failure,” in *IEEE International Conference on Robotics and Automation*, 2008, pp. 2207–2208.
- [22] K. Kostiadis and H. Hu, “KaBaGe-RL: Kanerva-based generalisation and reinforcement learning for possession football,” in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2001, pp. 292–297.
- [23] A. Y. Ng, D. Harada, and S. Russell, “Policy invariance under reward transformations: Theory and application to reward shaping,” in *Proceedings of the 16th International Conference on Machine Learning*, 1999, pp. 278–287.