



SAPIENZA
UNIVERSITÀ DI ROMA

DEPARTMENT OF COMPUTER, CONTROL AND
MANAGEMENT ENGINEERING

Atari Breakout with LTL_f/LDL_f Goals

ELECTIVE IN ARTIFICIAL INTELLIGENCE:
REASONING ROBOTS

Professor:

Giuseppe De Giacomo

Students:

Ivan Bergonzani

Michele Cipriano

Armando Nania

Contents

1	Introduction	2
2	Reinforcement Learning	3
2.1	Q-Learning	3
2.2	SARSA	4
3	LTL_f/LDL_f Non-Markovian Rewards	6
3.1	Theoretical Background	6
3.2	Examples	7
4	OpenAI Gym	9
4.1	Framework	9
4.2	Examples	9
5	Atari Breakout	11
5.1	PyGame Breakout	11
5.2	Arcade Learning Environment	12
5.3	Implementation	12
5.3.1	Robot Features Extractor	12
5.3.2	Goal Features Extractor	14
5.3.3	Temporal Goals	15
6	Experiments	17
6.1	PyGame and Atari Results	17
6.2	Modified PyGame	18
6.3	The Atari Breakout Difficulties	19
6.3.1	3.1	19
7	Conclusion	23

1 Introduction

Introduction to the whole project, structure of the report and summary of the work.

2 Reinforcement Learning

Reinforcement learning [?] is an area of machine learning which aims at studying how to develop agents that can interact with their environment maximizing a cumulative reward. The environment can be formally defined as a Markov Decision Process (MDP), which is a tuple $\langle S, A, \delta, R \rangle$ where S is a finite set of states that can represent the environment, A is a finite set of actions that can be performed by an agent in the environment, δ is a probability function modeling the transition from a state to another when performing a certain action and R is a reward function which models the reward received by the environment when performing a certain action which makes the agent move from a state to another.

An interesting property of the MDP is that it satisfies the Markov property, hence, the future states that will be reached by the agent do not depend on the past interaction of the environment, but just on the current state. This makes it possible to define the transition and the reward function depending only on the current state (and of course the action and the future state of interest).

This section considers two common reinforcement learning algorithms, namely Q-Learning and SARSA, which have been used in our experiments in order to train an agent interacting with an Atari Breakout environment (section 6).

2.1 Q-Learning

Q-Learning is a temporal difference (TD) algorithm that directly approximates the optimal action-value function. This method guarantees to find an optimal behaviour under the assumption that the all state-action pairs are updated infinitely many times. It is defined [?] by the following equation:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right] \quad (1)$$

Let's briefly discuss the implementation used in our project by studying the Python implementation (Algorithm 1). The algorithm is defined by the class `QLearning` that extends the abstract class `TDBrain`. The constructor of the class (lines 2-4) simply calls its parent constructor that will initialize the parameters of the object, hence, the observation space and the action space (`gym` objects), the strategy used by the policy function, which is ε -greedy by default and the hyperparameters γ , α and λ of the upper class. The abstract method inherited from `TDBrain` is `update_Q`, which should be implemented in order to define how to update the action-state table. The method (lines 6-21) simply follows Eq. 1.

TODO: eligibility.

Algorithm 1: Q-Learning algorithm Python implementation.

```
1 class QLearning(TDBrain):
2     def __init__(self, observation_space:Discrete,
3         ↪ action_space, policy:Policy=EGreedy(),
4         ↪ gamma=0.99, alpha=None, lambda_=0):
5         super().__init__(observation_space, action_space,
6         ↪ policy, gamma, alpha, lambda_)
7
8     def update_Q(self, obs:AgentObservation):
9         state, action, reward, state2 = obs.unpack()
10
11         action2 = self.choose_action(state2)
12         Qa = np.max(self.Q[state2])
13         actions_star = np.argwhere(self.Q[state2] == Qa).
14         ↪ flatten().tolist()
15
16         delta = reward + self.gamma * Qa - self.Q[state][
17         ↪ action]
18         for (s, a) in set(self.eligibility.traces.keys()):
19             self.Q[s][a] += self.alpha.get(s,a) * delta *
20             ↪ self.eligibility.get(s, a)
21             if action2 in actions_star:
22                 self.eligibility.update(s, a)
23             else:
24                 self.eligibility.to_zero(s, a)
25
26         return action2
```

2.2 SARSA

A similar TD algorithm is the SARSA algorithm, which name comes from the fact that at each timestep a quintuple $\langle S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1} \rangle$ is considered. As before, SARSA converges to an optimal action-value function under the assumption that all state-action pairs are updated infinitely many times. It is defined [?] by the following equation:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)] \quad (2)$$

Let's briefly discuss the implementation used in our project by studying the implementation (Algorithm 2), as done before with the Q-Learning algorithm. Again, the algorithm is defined by the class **Sarsa** that extends the abstract class **TDBrain**. The constructor of the class (lines 2-4) calls its parent constructor initializing the parameters of the upper class exactly in the same way as the class **QLearning**. The class implements the inherited method **update_Q** by following Eq. 2 (lines 6-17). TODO: eligibility.

Algorithm 2: SARSA algorithm Python implementation.

```
1 class Sarsa(TDBrain):
2     def __init__(self, observation_space:Discrete,
3         ↪ action_space, policy:Policy=EGreedy(),
4         ↪ gamma=0.99, alpha=None, lambda_=0.0):
5         super().__init__(observation_space, action_space,
6         ↪ policy, gamma, alpha, lambda_)
7
8     def update_Q(self, obs:AgentObservation):
9         state, action, reward, state2 = obs.unpack()
10
11         action2 = self.choose_action(state2)
12         Qa = self.Q[state2][action2]
13
14         delta = reward + self.gamma * Qa - self.Q[state][
15         ↪ action]
16
17         for (s, a) in set(self.eligibility.traces.keys()):
18             self.Q[s][a] += self.alpha.get(s,a) * delta *
19             ↪ self.eligibility.get(s, a)
20             self.eligibility.update(s, a)
21
22     return action2
```

3 LTL_f/LDL_f Non-Markovian Rewards

Recently, non-Markovian reward decision processes (NMRDPs) has attracted interest in the scientific community because of the possibility of specifying them as MDPs with LTL_f/LDL_f non-Markovian rewards [?]. In particular, it is possible to model the problem with two separate representations of the world, one for the agent (low-level) and one for the goal (expressed in terms of high-level fluents).

This section presents the approach used in [?], where an efficient method has been developed in order to work with NMRDPs. The theory behind the main idea is quickly described and an example on a theoretical Breakout environment is discussed in order to be used in the following sections easily.

3.1 Theoretical Background

Before describing the problem, let's give a formal definition of NMRDP. A non-Markovian reward decision process is a tuple $M = \langle S, A, \delta, \bar{R} \rangle$, with S finite set of states that can represent the environment, A is a finite set of actions that can be performed by an agent in the environment, δ is a probability function modeling the transition from a state to another when performing a certain action and $\bar{R} : (S \times A)^* \rightarrow \mathbb{R}$ is a function from finite state-action sequences (traces) to real-values that represents the reward given by the environment when performing a certain state-action sequence. Specifying a non-Markovian reward function explicitly is difficult even when considering a finite number of traces. Luckily, the LTL_f/LDL_f formalism allows to specify \bar{R} implicitly using a set of pairs $\{(\phi_i, r_i)\}_{i=1}^m$ with ϕ_i boolean proposition over the components of the state vector and r_i such that, given a current trace $\pi = \langle s, a_1, \dots, s_{n-1}, a_n \rangle$, the agent receives at s_n a reward r_i if ϕ_i is satisfied by π , hence:

$$\bar{R}(\pi) = \begin{cases} r_i & \text{if } \pi \models \phi_i \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

Since the NMRDP rewards are based on traces, instead of state-action pairs, typical learning algorithms like Q-learning or SARSA cannot be used. Nevertheless, it has been shown [?] that for any NMRDP $M = \langle S, A, \delta, \{(\phi_i, r_i)\}_{i=1}^m \rangle$ there exists an MDP $M' = \langle S', A', \delta' R' \rangle$ that is equivalent to M . The idea behind the proof consists in starting from an initial decision process $M_{ag}^{goal} = \langle S, A, R, \mathcal{L}, \delta_{ag}^g, \{(\phi_i, r_i)\}_{i=1}^m \rangle$ with LTL_f/LDL_f goals (with \mathcal{L} set of of configuration of the high-level features needed for expressing ϕ_i), transform it into a

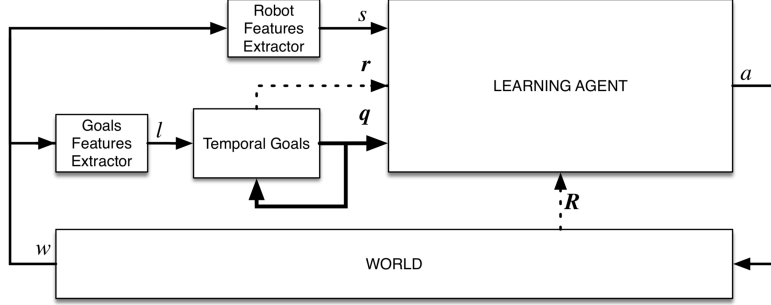


Figure 1: Pipeline describing how the agent is interacting with the world and how the robot features extractor and the goal features extractor are used in order to handle non-Markovian rewards.

NMRDP in order to further transform it into a MDP where it is possible to execute learning algorithms such as Q-learning. All the details are out of the scope of the project and are discussed in [?]. The set of states S is used to express low-level features of the agent.

The new MDP is extended by using automata needed to track the satisfiability of a LTL_f/LDL_f formulas, as explained in [?]. Each formula is, hence, associated with an automaton that accepts exactly the traces satisfying the formula itself.

Once the MDP has been determined it is possible to train a model that will be able to manage non-Markovian rewards as well. Fig. 1 shows a pipeline of the agent interacting with the world. In particular, a robot feature extractor will analyze the state of the world w in order to pass a representation s to the agent, a goal feature extractor will extract higher-level features of the world l that can be used to determine if a temporal goal has been achieved or not. In positive case a reward r is given to the agent. The evaluation q of the formula that specifies the temporal goals is given to the learning agent as well, in order to construct an extended state of the MDP as specified by the theory above. Each time the learning agent performs an action a , the environment changes and it gives to the agent a reward R .

3.2 Examples

A typical example is breaking the bricks of a theoretical Breakout game in a certain order. Let's consider a 3×3 Breakout. A good idea is to use LTL_f/LDL_f goals to help the agent break to bricks from top to bottom in order to make it create a hole on the left (or the right) and let the ball do all the job without making any movements when it reaches the top part of the environment. The

idea is, hence, to help the agent learn the most effective strategy, which should also decrease the training time. Let's consider a formula that checks whether a row has been broken or not:

$$\begin{aligned} & \langle (\neg\varphi_0 \wedge \neg\varphi_1 \wedge \neg\varphi_2)^*; (\varphi_0 \wedge \neg\varphi_1 \wedge \neg\varphi_2); (\varphi_0 \wedge \neg\varphi_1 \wedge \neg\varphi_2)^*; \\ & (\varphi_0 \wedge \varphi_1 \wedge \neg\varphi_2); (\varphi_0 \wedge \varphi_1 \wedge \neg\varphi_2)^*; (\varphi_0 \wedge \varphi_1 \wedge \varphi_2) \rangle tt \end{aligned} \quad (4)$$

which, given φ_0 first row from the top broken, φ_1 second row from the top broken and φ_2 third row from the top broken, can be interpreted as:

- $(\neg\varphi_0 \wedge \neg\varphi_1 \wedge \neg\varphi_2)^*$: initially, and for an indefinite amount of time, none of the rows have been broken;
- $(\varphi_0 \wedge \neg\varphi_1 \wedge \neg\varphi_2)$: then, the first row from the top has been broken but the remaining two are intact;
- $(\varphi_0 \wedge \neg\varphi_1 \wedge \neg\varphi_2)^*$: then, for an indefinite amount of time, the first row remains broken and the remaining two remain intact;
- $(\varphi_0 \wedge \varphi_1 \wedge \neg\varphi_2)$: then, the first two rows from the top have been broken and the remaining one at the bottom is intact;
- $(\varphi_0 \wedge \varphi_1 \wedge \neg\varphi_2)^*$: then, for an indefinite amount of time, the first two rows from the top remain broken and the last one remains intact;
- $(\varphi_0 \wedge \varphi_1 \wedge \varphi_2)$: in the end, all the rows have been broken, hence, all the bricks have been broken.

In general, the formula $\langle \rho \rangle \phi$ states that, from the current step in the trace, there exists an execution satisfying ρ such that in its last step ϕ is satisfied. In the example above, at the end of the execution all the rows have been broken.

4 OpenAI Gym

OpenAI `gym` [?] is a toolkit for developing and comparing reinforcement learning algorithms, without making assumptions about the structure of the agent interacting with the environment, in order to keep development flexible to updates on both sides.

4.1 Framework

The framework of `gym` allows to interact easily with an environment, giving the developers the tools they need to perform actions and to observe the state of the environment itself. In this way it is possible to focus more on the development of the agent without spending time on the structure of the world.

`gym` makes it possible to interact with multiple kinds of environments. Among these, the authors of the framework developed the support for Arcade Learning Environment [?], which includes all the classic Atari games, including Breakout, which has been used in this project.

4.2 Examples

Let's consider a simple example to understand how `gym` works and how the framework can be used to interact with an environment. The description will follow Algorithm 3.

Algorithm 3: Example of a random interaction with the `gym` environment `BreakoutNoFrameskip-v4`, used also in our experiments of subsection 6.

```
1 import gym
2
3 env = gym.make("BreakoutNoFrameskip-v4")
4 env.reset()
5
6 for _ in range(1000):
7     env.render()
8     action = env.action_space.sample() # takes random actions
9     observation, reward, done, info = env.step(action)
10    if done == True:
11        env.reset()
12
13 env.close()
```

Initially (line 1) the framework is imported. Then (line 3-4) an environment is created specifying its name and initializing it. The program makes a random agent interact randomly with the environment for 1000 timesteps (lines 6-11) before closing the environment. Line 7 renders the current observation of the

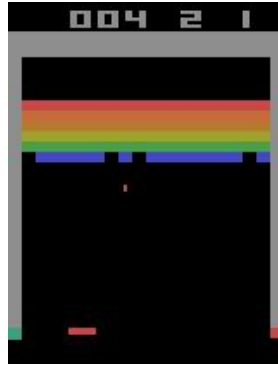


Figure 2: Observation of a frame of the environment `BreakoutNoFrameskip-v4`.

environment on screen, line 8-9 performs a random action between those available in this Brekout version, note that the method `step` return an `observation` (shown in Fig. 2), which is an array of pixels that represent the current state of the environment, a `reward`, which is a value return by the game after performing the specified action `action`, a boolean value `done`, which is `True` when the game is over, `False` otherwise, and `info` which contains extra information about the game. Lines 10-11 handles the case when the game is over, resetting the environment.

5 Atari Breakout

This section contains the main part of the project, describing in detail our starting point, how the program has been developed, the results achieved and the comparison with the original implementation, which uses a non-Atari version of the game Breakout [?] with the same number of bricks used in the Atari implementation of `gym`, already introduced in section 5.2.

Initially, the non-Atari version of Breakout (built on PyGame) is introduced, its implementation is discussed and the results on the training with a brick matrix of dimension 6×18 are presented. Then the `gym` environment `BreakoutNoFrameskip-v4` is presented and compared with the PyGame breakout. A detailed description of the implementation of the project is described, discussing in detail how the features have been extracted from the environment (both the robot and the goal features), how temporal goals have been used to evaluate the state of the bricks and how everything is connected together in order to make it work with `gym`. In the end, the main experiments performed on `BreakoutNoFrameskip-v4` are presented.

5.1 PyGame Breakout

As introduced above, in [?] a PyGame version of the Breakout game has been used in order to test the algorithms introduced in the paper. The implementation of this Breakout easily allows to determine the state of the environment, saving the status of the bricks (either in the scene or broken), the position of the ball, the direction of the ball and the position of the paddle. This makes it possible to reduce a lot the computational time of the implementation of the agent since the data it receives are already preprocessed in order to have a complete overview of the environment, allowing to focus on higher-level reasoning tasks.

Originally, the paper focused on a Breakout environment with a brick matrix of dimension 4×5 . Since the Atari version of Breakout is dealing with a brick matrix of dimension 6×18 , a new test has been performed to make the two environments comparable, in order to better understand the potentiality of LTL_f/LDL_f , the use of non-Markovian rewards and how to approach complex `gym` environments in the future.

The method managed to correctly break all the bricks in around one hour of training. TODO: add three figures with initial state, middle, final. TODO: add plot of rewards (shall this go to Experiments subsec?).

5.2 Arcade Learning Environment

This work aims at comparing the `gym` environment `BreakoutNoFrameskip-v4`, already introduced in section with the non-Atari Breakout used in [?]. In last years, the reinforcement learning community has grown a lot thanks to the introduction of `gym` and deep reinforcement learning algorithms [?] that managed to easily solve complex games that are considered difficult also for humans, often achieving better results than expert human gamers. More and more algorithms are introduced every year, exploiting GPU resources and managing to solve harder games like Montezuma Revenge [?]. The popularity of deep reinforcement learning begun with the introduction of Arcade Learning Environment (ALE) that includes most famous arcade Atari games [?].

The main characteristics of `gym` (or ALE) environments is that the world can be observed only from the pixels of the screen, putting the algorithms at the same level of the human, that can only observe the display while playing. This makes the game a lot more complex since a more abstract reasoning strategy is needed in order to solve the game. This hypothesis should make `gym` games a lot harder than the non-Atari version of Breakout that has been used to test algorithms that work with non-Markovian rewards.

5.3 Implementation

Intro to implementation.

5.3.1 Robot Features Extractor

`RobotFeatureExtractor` (OpenCV). Extracts features of the robot (robot and ball positions).

Algorithm 4: Robot feature extractor Python implementation.

```
1 class BreakoutNRobotFeatureExtractor(  
2     ↳ BreakoutRobotFeatureExtractor):  
3     def __init__(self, obs_space):  
4         robot_feature_space = Tuple((  
5             Discrete(287),  
6             Discrete(157),  
7         ))  
8  
9         self.prev_ballX = 0  
10        self.prev_ballY = 0  
11        self.prev_paddleX = 0  
12        self.still_image = True  
13  
14        super().__init__(obs_space, robot_feature_space)
```

```

15
16     def _extract(self, input, **kwargs):
17         self.still_image = not self.still_image
18         if self.still_image:
19             return (self.prev_ballX-self.prev_paddleX+143,
20                     ↪ self.prev_ballY)
21         # Extract position of the paddle:
22         paddle_img = input[189:193,8:152,:]
23         gray = cv2.cvtColor(paddle_img, cv2.COLOR_RGB2GRAY)
24         thresh = cv2.threshold(gray, 60, 255, cv2.
25             ↪ THRESH_BINARY)[1]
26         cnts = cv2.findContours(thresh.copy(), cv2.
27             ↪ RETR_EXTERNAL, cv2.CHAIN_APPROX_SIMPLE)
28         cnts = cnts[0] if imutils.is_cv2() else cnts[1]
29         min_distance = np.inf
30         paddleX = self.prev_paddleX
31         for c in cnts:
32             M = cv2.moments(c)
33             if M["m00"] == 0:
34                 continue
35             pX = int(M["m10"] / M["m00"])
36             if abs(self.prev_paddleX - pX) < min_distance:
37                 min_distance = abs(self.prev_paddleX - pX)
38                 paddleX = pX
39
40         # Extract position of the ball:
41         ballX = self.prev_ballX
42         ballY = self.prev_ballY
43         ballspace_img = input[32:189,8:152,:]
44         lower = np.array([200, 72, 72], dtype=np.uint8)
45         upper = np.array([200, 72, 72], dtype=np.uint8)
46         mask = cv2.inRange(ballspace_img, lower, upper)
47         cnts = cv2.findContours(mask.copy(), cv2.
48             ↪ RETR_EXTERNAL, cv2.CHAIN_APPROX_SIMPLE)
49         cnts = cnts[0] if imutils.is_cv2() else cnts[1]
50         for c in cnts:
51             M = cv2.moments(c)
52             # Avoid to compute position of the ball if M["
53             ↪ m00"] is zero:
54             if M["m00"] == 0:
55                 continue
56             # Calculate the centroid
57             cX = int(M["m10"] / M["m00"])
58             cY = int(M["m01"] / M["m00"])
59             # Check that the centroid is actually part of
60             ↪ the ball:
61             left_black = False
62             right_black = False
63             if cX > 3:
64                 if ballspace_img[cY][cX-3][0] != 200 or \
65                     ballspace_img[cY][cX-3][1] != 72 or \
66                     ballspace_img[cY][cX-3][2] != 72:
67                     left_black = True
68             else:

```

```

63         if ballspace_img[cY][cX+3][0] != 200 or \
64            ballspace_img[cY][cX+3][1] != 72 or \
65            ballspace_img[cY][cX+3][2] != 72:
66             right_black = True
67         if left_black or right_black:
68             ballX = cX
69             ballY = cY
70
71         self.prev_ballX = ballX
72         self.prev_ballY = ballY
73         self.prev_paddleX = paddleX
74
75         return (self.prev_ballX - self.prev_paddleX + 143,
76                 ↪ self.prev_ballY)

```

5.3.2 Goal Features Extractor

GoalFeatureExtractor (OpenCV). Extracts 6x18 table representation of the bricks in order to evaluate a formula.

Algorithm 5: Goal feature extractor Python implementation.

```

1 class BreakoutGoalFeatureExtractor(FeatureExtractor):
2     def __init__(self, obs_space, bricks_rows=6,
3         ↪ bricks_cols=18):
4         self.bricks_rows = bricks_rows
5         self.bricks_cols = bricks_cols
6         output_space = Box(low=0, high=1, shape=(
7             ↪ bricks_cols, bricks_rows), dtype=np.uint8)
8         super().__init__(obs_space, output_space)
9
10    def _extract(self, input, **kwargs):
11        bricks_features = np.ones((self.bricks_cols, self.
12            ↪ bricks_rows))
13        for row, col in itertools.product(range(self.
14            ↪ bricks_rows), range(self.bricks_cols)):
15            # Pixel of the observation to check:
16            px_upper_left = int( 8 + 8 * col)
17            py_upper_left = int(57 + 6 * row)
18            px_upper_right = int(15 + 8 * col)
19            py_upper_right = int(57 + 6 * row)
20
21            # Checking max because the input has 3 channels
22            ↪ :
23            if max(input[py_upper_left][px_upper_left]) ==
24                ↪ 0 or \
25                max(input[py_upper_right][px_upper_right])
26                ↪ == 0:
27                bricks_features[col][row] = 0
28
29        return bricks_features

```

*Ext used to improve implementation.

5.3.3 Temporal Goals

LTL_f/LDL_f implementation (with Marco Favorito libraries).

Algorithm 6: LTL_f/LDL_f formulas Python implementation.

```

1 def get_breakout_lines_formula(lines_symbols):
2     # Generate the formula string
3     # E.g. for 3 line symbols:
4     # "<(!10 & !11 & !12)*;(10 & !11 & !12);(10 & !11 & !12
        ↪ )*(10 & 11 & !12); (10 & 11 & !12)*; 10 & 11 &
        ↪ 12>tt"
5     pos = list(map(str, lines_symbols))
6     neg = list(map(lambda x: "!" + str(x), lines_symbols))
7
8     s = "(%s)*" % " & ".join(neg)
9     for idx in range(len(lines_symbols)-1):
10        step = " & ".join(pos[:idx + 1]) + " & " + " & ".
            ↪ join(neg[idx + 1:])
11        s += ";({0});({0})*".format(step)
12    s += ";(%s)" % " & ".join(pos)
13    s = "<%s>tt" % s
14
15    return s
16
17 class BreakoutCompleteLinesTemporalEvaluator(
18     ↪ TemporalEvaluator):
19     """Breakout temporal evaluator for delete columns from
20     ↪ left to right"""
21
22     def __init__(self, input_space, bricks_cols=3,
23     ↪ bricks_rows=3, lines_num=3, gamma=0.99,
24     ↪ on_the_fly=False):
25         assert lines_num == bricks_cols or lines_num ==
26         ↪ bricks_rows
27         self.line_symbols = [Symbol("l%s" % i) for i in
28         ↪ range(lines_num)]
29         lines = self.line_symbols
30
31         parser = LDLfParser()
32
33         string_formula = get_breakout_lines_formula(lines)
34         print(string_formula)
35         f = parser(string_formula)
36         reward = 10000
37
38         super().__init__(BreakoutGoalFeatureExtractor(
39         ↪ input_space, bricks_cols=bricks_cols,
40         ↪ bricks_rows=bricks_rows),
41         ↪ set(lines),
42         ↪ f,
43         ↪ reward,
44         ↪ gamma=gamma,
45         ↪ on_the_fly=on_the_fly)

```


6 Experiments

The results are clear: the method is not efficient for the problem to be treated. Although the algorithm applied to the PyGame version of Breakout gives excellent results, the same algorithm applied to the Atari version does not give acceptable results, to the point that it was not even possible to test the temporal goals on the latter. In this section we will analyze the results obtained, also showing a some (not all) of the tests that were carried out in an attempt to find out what was the cause of the intractability of Atari's environment. In particular, we will first show the results on PyGame and Atari environments, comparing them. Then we will explain the changes made to the PyGame version during the subsequent test phases. Finally we will draw conclusions from these experiments, analyzing in detail the problems encountered, and we will show some possible changes to be used to address these problems

6.1 PyGame and Atari Results

As explained above, the results are very clear: the method applied to the environment of PyGame gives excellent results already from 5000 training periods, while the results obtained by applying the same method to the Atari environment are lower, even after 1,000,000 epochs. In particular, the best results obtained by PyGame are the following:

TODO: risultati su pygame

The results obtained from the use of the Atari environment are:

TODO: risultati su atari

As you can see, the results are much lower in the second case. In an attempt to improve the Atari results, changes in the code have been tested: we have reduced the size of the state space, either the ball space or the paddle space, or its movement space; we have tried to completely remove the information concerning the movement of the paddle; we tested a one-shot version of Atari, in order to make it similar to PyGame and finally we have tried combinations of the above ones. These are the results of these tests:

TODO: risultati test

Given the number of tests to be performed, we performed them with a low number of training periods, 5000, for obvious reasons of time. So it is necessary to take this data with a grain of salt. But it is still possible to see a trend: the reduction of the state space does not significantly improve the result you get, while giving Atari environment the opportunity to play all his 5 lives worsens

the score obtained at the end of the first life, but increases the total score. This is a consequence of the fact that there are no penalties in the reward associated with losing a life: there are only positive rewards given by the destruction of a brick or the achievement of a temporal goal. This makes possible scenarios in which, after losing a life, the paddle succeeds in destroying a brick in the next life, which leads to a reward for the sequence of state-action pairs that led to that result. All this makes the training of AI very random, because the life in which it will take the reward is random, with the risk of losing lives uselessly. We will draw further conclusions later in this section.

6.2 Modified PyGame

In the previous point we tried to identify strategies for managing the state space and lives to improve the score, but despite this the difference between the two environments is still large. It is true that PyGame Breakout has a very small state space, but nevertheless decreasing the state space in Atari Breakout to make it similar to PyGame does not produce any improvements. So instead of making Atari similar to PyGame, let's try to do the opposite. There are indeed considerable structural differences between the two environments: the size of paddle, ball, bricks, window and placements of objects in the field. Moreover, the movement of the ball is much slower in PyGame than in Atari: the ball, to reach the first bricks from the paddle, takes 35 frames in PyGame, and only 15 in Atari. In 1500 frames of PyGame it can only destroy 10 bricks, in 1500 frames of Atari it destroys 29 bricks, even 47 with 2000 frames! Remind that a frame corresponds to a state in which you choose an action to move to the next state. In order to better compare the two environments, in an attempt to understand the problems inherent to Atari, we gradually approximated PyGame to Atari, to make the two environments identical.

TODO: screen pygame vs atari

To do this we changed the PyGame code at some points. In doing so it was highlighted a bug that was already present in PyGame, but that did not show up in its original version:

TODO: codice correzione bug

We quickly solved it and compared the two environments that we obtained. Despite all these efforts, the difference between the results obtained by applying the method on one environment rather than another is still large: the modified PyGame still manages to produce good results, even with just 5000 epochs of training. It has been noted, however, that PyGame uses an internal reduction of

the actually used states: when the ball is lower in the screen, the ball and paddle positions are calculated with greater resolution. We tried to use this resolution reduction algorithm in Atari, but the results were once again disappointing. We even found a worsening effect of the score.

TODO: codice riduzione risoluzione

After all these futile attempts we have finally managed to identify the problems of Atari, which are summarized in the next point.

6.3 The Atari Breakout Difficulties

In this point we will analyze the intrinsic problems in Atari that prevent the achievement of good results with the method used. One has already been addressed in point 1, which concerned the use of an environment with more lives. Here we will describe the problems of the state space, of the local maximum, of non-determinism and of the collision of states.

6.3.1 The Dimension Space Problem

The problem of the size of the state space was clear from the beginning. Before testing it, PyGame had only been used with a 3x3 brick matrix. Atari, on the other hand, has a fixed-size matrix equal to 6x18, which makes both the automaton defining the temporal goals and the size of the playing field larger, in proportion to the paddle size. As explained up to here, there have been many attempts at reducing the state space in Atari, to the extent that it made that dimension identical to that of PyGame, and all attempts proved to be unsuccessful. They did not produce significant improvements in the score, and some times even showed a worsening. In fact, it is possible to notice from the number of states visited at each epoch that the speed with which the AI explores new states falls after some threshold, despite maintaining a low score. The algorithm therefore has difficulty in exploring new states, which can have two explanations: there are too many states, and the algorithm get lost, or the problem is not about the number of states, is about the inability of the algorithm to learn, probably due to a collision of states (see below).

6.3.2 The Local Maximum Problem

The artificial intelligence algorithm that must learn to maximize the reward must face a problem of local maximum. In fact, there are many states from which it is difficult to escape. Take the case in which the paddle hits the ball in the lower left corner. If the ball tends to go to the right as a result of the

action, there is the possibility that it will fall back into the opposite corner, lower right. If this is the case, the paddle will have to move multiple times, before it can get a new reward. In particular, we need a series of subsequent "go right" actions, which make it possible to continue the game. This is very unlikely, because, in the presence of new states, the choice of action to be taken is random. And there's a second problem: the difficulty introduced by a greater number of possible actions: while on PyGame there are only three actions ("go left", "go right", "stay still"), on Atari there is a fourth possible action, which is needed to spawn the ball at the start of the game, or when a life is lost. This action is automatically performed by the environment on the aforementioned occasions, but it has not been possible to prevent the execution of these actions during the normal course of the game. So the correct action to be taken ("go right") must be chosen randomly by a group of 4 actions (25% probability of choosing the correct action) instead of 3 (33%). And if, say, the number of times to perform the "go right" action is N , the probability of this happening is $25\%^N$ EPSILON, which falls below 1% already for $N = 4$. If we analyze the moments in which the paddle drops the ball in the training carried out, we note that most of them are just the cases described here, in which the paddle is in a corner of the screen and does not move, sign that no rewards have been discovered for any action, and the ball falls in the opposite corner.

TODO: sequenza immagini palla che cade

But we have found a trick to alleviate this problem. By default, Atari returns an observation every 4 frames, and for each of those frames Atari performs the same action we told it to execute. In practice a state-frame corresponds to 4 game-frame: at each state-frame the AI chooses an action that will be performed for 4 game-frames, after which the environment will return a new observation, that is a new state-frame. Increasing this number increases the number of times the same action is performed, increasing the paddle movement from one state-frame to another. This allows you to choose the same action for fewer times in order to reach the opposite corner of the screen. That is, it reduces N , thereby increasing the probability of this happening ($25\%^N$). On the other hand, this makes it more difficult to hit the ball. Tests were performed and optimal intermediate values were found for the frame skip, between 10 and 14, which lead to a good increase in the score.

6.3.3 The Non-Determinism Problem

During the frame skip tests we noticed the problem that is the main difference between the two environments: non-determinism. When you perform an

action in a state, it is not possible to know for sure what the next state will be like. This is not the case for PyGame. For Atari, this is because the amount of paddle movement is not fixed, but rather varies, partly randomly and partly according to the actions performed in previous states. It is even possible that, given a state and a movement action to be performed, the paddle will move into the next state in the opposite direction to that indicated, due to the inertia that the paddle has at that moment. In particular, without a frame skip, the paddle moves from -6 to +6 pixels. The skip of the frames, however, accentuates this problem, increasing the uncertainty of the movement from -6f to + 6f, with f the number of skipped frames from one state to another. To solve this problem at least partially, we inserted another variable in the state: move. It indicates the amount of movement performed with the last action, in order to have an estimate of the inertia of the paddle, and therefore know at least approximately where the paddle will go afterwards. The results of these tests are in point 1 of this section.

6.3.4 The State Collision Problem

The most serious problem, however, is the one introduced by non-determinism: the collision of states, that is when two different states are represented in the same way by the AI. The introduction of the variable move serves precisely this, that is, to provide the algorithm of a sort of memory that indicates more or less what the previous state was and, consequently, help to predict the future state based on its own actions. However, this is not a definitive solution. Suppose in fact that move indicates that the paddle has moved an average amount between no movement (0) and maximum movement (6f). We do not know if the paddle was accelerating or braking. We do not know if the previous action was "stand still" and the paddle moved randomly by a lot, or if it was "go right" and the paddle randomly moved by a little bit. Suppose now that the AI has learned that in a certain state, by performing a certain action, it will get a reward. What happens now if, during the training phase, the AI meets again the same state encountered some frames before, but in a different circumstance? If the epsilon policy does not tell us to choose randomly, we will obviously do the same action as before, but the result may not be the same, we could lose the ball. Even worse is the case in which the epsilon policy tells us to perform a random action, different from the optimal one, but leads to a reward. In this case we are going to overwrite the information of the optimal action on that state and moreover, since this reward is more recent, it will have a higher value. So it will happen that, the next restart of the game, when we meet the first occurrence of

the double state, the AI will perform the wrong action, losing not only the ball, but also all the rewards we would have obtained following that path of actions leading to the second occurrence of the state, the very same path that we had already explored at the cost of so many clock cycles. Furthermore, the greater the distance between the two occurrences of the conflicting state, the greater the loss of rewards. This explains the fluctuating trend of the score during the training phase, which is notably different from the PyGame case. So this is the main reason for the difficulty of applying the Q-Learning or SARSA algorithms on Atari Breakout.

TODO: grafici andamento reward per ogni iterazione, atari vs pygame

7 Conclusion

Why it does not work.

Summary + differences between the two environments.

Future works (neural networks and parallel computation).