



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	3
3	LTL_f/LDL_f Non-Markovian Rewards	5
3.1	Theoretical Background	5
3.2	Examples	5
4	OpenAI Gym	6
4.1	Framework	6
4.2	Examples	6
5	Atari Breakout	7
5.1	PyGame Breakout	7
5.2	Arcade Learning Environment	7
5.3	Implementation	7
5.4	Experiments	11
6	Conclusion	12

1 Introduction

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

2 Reinforcement Learning

Introduction to RL.

2.1 Q-Learning

Q-Learning algorithm.

```
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

SARSA algorithm.

```
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]
```

```
12         delta = reward + self.gamma * Qa - self.Q[state][  
13             ↪ action]  
13     for (s, a) in set(self.eligibility.traces.keys()):  
14         self.Q[s][a] += self.alpha.get(s,a) * delta *  
15             ↪ self.eligibility.get(s, a)  
15         self.eligibility.update(s, a)  
16  
17     return action2
```

3 LTL_f/LDL_f Non-Markovian Rewards

3.1 Theoretical Background

Introduction to the research paper.

3.2 Examples

How it can be used to train a RL model.

4 OpenAI Gym

Intro.

4.1 Framework

Introduction to the framework.

4.2 Examples

Examples like the one on the website (+CODE).

5 Atari Breakout

Intro.

5.1 PyGame Breakout

Original implementation of the paper (non-ATARI).

5.2 Arcade Learning Environment

ATARI Breakout (from ALE) and differences from the other one.

5.3 Implementation

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

```
1 class BreakoutNRobotFeatureExtractor(  
2     ↳ BreakoutRobotFeatureExtractor):  
3  
4     def __init__(self, obs_space):  
5         robot_feature_space = Tuple((  
6             Discrete(287),  
7             Discrete(157),  
8         ))  
9  
10        self.prev_ballX = 0  
11        self.prev_ballY = 0  
12        self.prev_paddleX = 0  
13        self.still_image = True  
14  
15        super().__init__(obs_space, robot_feature_space)  
16  
17        def _extract(self, input, **kwargs):  
18            self.still_image = not self.still_image  
19            if self.still_image:  
20                return (self.prev_ballX-self.prev_paddleX+143,  
21                    ↳ self.prev_ballY)  
22            # Extract position of the paddle:  
23            paddle_img = input[189:193,8:152,:]  
24            gray = cv2.cvtColor(paddle_img, cv2.COLOR_RGB2GRAY)  
25            thresh = cv2.threshold(gray, 60, 255, cv2.  
26                ↳ THRESH_BINARY)[1]  
27            cnts = cv2.findContours(thresh.copy(), cv2.  
28                ↳ RETR_EXTERNAL, cv2.CHAIN_APPROX_SIMPLE)  
29            cnts = cnts[0] if imutils.is_cv2() else cnts[1]  
30            min_distance = np.inf  
31            paddleX = self.prev_paddleX  
32            for c in cnts:  
33                M = cv2.moments(c)  
34                if M["m00"] == 0:
```



```

31         continue
32         pX = int(M["m10"] / M["m00"])
33         if abs(self.prev_paddleX - pX) < min_distance:
34             min_distance = abs(self.prev_paddleX - pX)
35             paddleX = pX
36
37     # Extract position of the ball:
38     ballX = self.prev_ballX
39     ballY = self.prev_ballY
40     ballspace_img = input[32:189,8:152,:]
41     lower = np.array([200, 72, 72], dtype=np.uint8)
42     upper = np.array([200, 72, 72], dtype=np.uint8)
43     mask = cv2.inRange(ballspace_img, lower, upper)
44     cnts = cv2.findContours(mask.copy(), cv2.
        ↳ RETR_EXTERNAL, cv2.CHAIN_APPROX_SIMPLE)
45     cnts = cnts[0] if imutils.is_cv2() else cnts[1]
46     for c in cnts:
47         M = cv2.moments(c)
48         # Avoid to compute position of the ball if M["
        ↳ m00"] is zero:
49         if M["m00"] == 0:
50             continue
51         # Calculate the centroid
52         cX = int(M["m10"] / M["m00"])
53         cY = int(M["m01"] / M["m00"])
54         # Check that the centroid is actually part of
        ↳ the ball:
55         left_black = False
56         right_black = False
57         if cX > 3:
58             if ballspace_img[cY][cX-3][0] != 200 or \
59                 ballspace_img[cY][cX-3][1] != 72 or \
60                 ballspace_img[cY][cX-3][2] != 72:
61                 left_black = True
62         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,
        ↳ self.prev_ballY)

```

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

```

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.

LTL_f/LDL_f implementation (with Marco Favorito libraries).

```

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
5             ↪ )*(10 & 11 & !12); (10 & 11 & !12)*; 10 & 11 &
6             ↪ 12>tt"
7         pos = list(map(str, lines_symbols))
8         neg = list(map(lambda x: "!" + str(x), lines_symbols))
9
10        s = "(%s)*" % " & ".join(neg)
11        for idx in range(len(lines_symbols)-1):
12            step = " & ".join(pos[:idx + 1]) + " & " + " & ".
13                ↪ join(neg[idx + 1:])
14            s += "({0});({0})".format(step)
15        s += "({0})" % " & ".join(pos)
16        s = "<%s>tt" % s
17
18        return s
19
20    class BreakoutCompleteLinesTemporalEvaluator(
21        ↪ TemporalEvaluator):
22        """Breakout temporal evaluator for delete columns from

```

```

19         ↪ left to right"""
20     def __init__(self, input_space, bricks_cols=3,
21         ↪ bricks_rows=3, lines_num=3, gamma=0.99,
22         ↪ on_the_fly=False):
23         assert lines_num == bricks_cols or lines_num ==
24             ↪ bricks_rows
25         self.line_symbols = [Symbol("l%s" % i) for i in
26             ↪ range(lines_num)]
27         lines = self.line_symbols
28
29         parser = LDLfParser()
30
31         string_formula = get_breakout_lines_formula(lines)
32         print(string_formula)
33         f = parser(string_formula)
34         reward = 10000
35
36         super().__init__(BreakoutGoalFeatureExtractor(
37             ↪ input_space, bricks_cols=bricks_cols,
38             ↪ bricks_rows=bricks_rows),
39             set(lines),
40             f,
41             reward,
42             gamma=gamma,
43             on_the_fly=on_the_fly)
44
45     @abstractmethod
46     def fromFeaturesToPropositional(self, features, action,
47         ↪ *args, **kwargs):
48         """map the matrix bricks status to a propositional
49             ↪ formula
50         first dimension: columns
51         second dimension: row
52         """
53         matrix = features
54         lines_status = np.all(matrix == 0.0, axis=kwargs["
55             ↪ axis"])
56         result = set()
57         sorted_symbols = reversed(self.line_symbols) if
58             ↪ kwargs["is_reversed"] else self.line_symbols
59         for rs, sym in zip(lines_status, sorted_symbols):
60             if rs:
61                 result.add(sym)
62
63         return frozenset(result)
64
65 class BreakoutCompleteRowsTemporalEvaluator(
66     ↪ BreakoutCompleteLinesTemporalEvaluator):
67     """Temporal evaluator for complete rows in order"""
68
69     def __init__(self, input_space, bricks_cols=3,
70         ↪ bricks_rows=3, bottom_up=True, gamma=0.99,

```

```

60         ↪ on_the_fly=False):
        super().__init__(input_space, bricks_cols=
        ↪ bricks_cols, bricks_rows=bricks_rows,
        ↪ lines_num=bricks_rows, gamma=gamma,
        ↪ on_the_fly=on_the_fly)
61     self.bottom_up = bottom_up
62
63     def fromFeaturesToPropositional(self, features, action,
64     ↪ *args, **kwargs):
65         """complete rows from bottom-to-up or top-to-down,
        ↪ depending on self.bottom_up"""
        return super().fromFeaturesToPropositional(features
        ↪ , action, axis=0, is_reversed=self.bottom_up
        ↪ )

```

Atari wrappers (OpenAI).

5.4 Experiments

Results with 6x18 non-ATARI Breakout (+CODE).

Results with our experiments (+CODE).

6 Conclusion

Why it does not work.

Summary + differences between the two environments.

Future works (neural networks and parallel computation).

References