

Dipartimento di Informatica - Scienza e Ingegneria
Artificial Intelligence

Flatland Challenge

Project Presentation

Professor

Andrea Asperti

Students

Alessandro Lombardi
Fiorenzo Parascandolo

Academic Year 2019/2020

Contents

1	Flatland	1
1.1	The classes RailAgentStatus and EnvAgent	1
1.2	The class RailEnv	2
2	Multi Agent Reinforcement Learning (MARL)	5
2.1	Theoretical background	5
2.1.1	Introduction	5
2.1.2	Challenges	7
2.1.3	Algorithms	9

1 Flatland

The analyzed version is the 2.1.10

1.1 The classes RailAgentStatus and EnvAgent

RailAgentStatus extends Python IntEnum and assumes the following values:

- READY_TO_DEPART (0) the agent is not in the grid yet (position is None), the prediction is to stay at the starting position. If a MOVE_* action is performed during this state it becomes ACTIVE.
- ACTIVE (1) the agent is in the grid (position is not None) and hasn't reached the target yet, the prediction is the remaining path.
- DONE (2) the agent is still in the grid (position is not None) but has already reached the target, the prediction is to stay at the target forever.
- DONE_REMOVED (3) the agent has reached the target and it's removed from the grid.

Grid4TransitionsEnum extends Python IntEnum and assumes the following values: NORTH (0), EAST (1), SOUTH (2), WEST (3). Grid4TransitionsEnum is used to indicate absolute directions, related to the environment, like a compass. Possible usages are storing where the agent is facing or computing legal actions, for example including as observation a one hot encoding of the directions where the agent can move.

EnvAgent class models the agent and encapsulates in its internal state the following attributes:

- initial_position: Tuple[int, int], initial coordinate.
- initial_direction: Grid4TransitionsEnum, the initial agent facing direction.
- direction: Grid4TransitionsEnum, the current facing direction.
- target: Tuple[int, int], the final coordinate.
- moving: bool, True if the agent is in a moving state.
- speed_data: dictionary, TODO
- malfunction_data: dictionary, TODO
- status: RailAgentStatus, the current agent status.

The speed of an agent contains the keys 'position_fraction' used as a counter of the percentage of completion of an movement from a cell to another, 'speed' the value between 0 and 1 used to increment the 'position_fraction' and 'transition_action_on_cellexit'

which contains the action to perform on the next cell (if it completes the one in the current step).

The malfunction of an agent contains the keys 'malfunction' which contains how many steps are necessary to fix the agent, 'malfunction_rate', the mean rate (average number of events in an interval) of the Poisson distribution, 'next_malfunction' the number of steps the next malfunction will occur and 'nr_malfunctions' the number of previous malfunctions.

1.2 The class RailEnv

From the documentation

RailEnv is an environment inspired by a (simplified version of) a rail network, in which agents (trains) have to navigate to their target locations in the shortest time possible, while at the same time cooperating to avoid bottlenecks.

In the *step* function the number of steps is updated and if the overall task is still uncompleted, for each agent the associated reward is initially put to zero, a malfunction is tried to be induced and the specific step is performed. The info of the agent are prepared and finally the end malfunctions are repaired. Agents are handled in the order in which are passed.

Environment Actions

The available actions are:

- DO_NOTHING (0) Default action if None has been provided or the value is not within this list. If agent.moving is True then the agent will MOVE_FORWARD.
- MOVE_LEFT (1) If agent.moving is False then becomes True. If it's possible turn the agent left, changing its direction, otherwise if agent.moving is True tries the action MOVE_FORWARD.
- MOVE_FORWARD (2) If agent.moving is False then becomes True. It updates the direction of the agent and if the new cell is a dead-end the new direction is the opposite of the current.
- MOVE_RIGHT (3) If agent.moving is False then becomes True. If it's possible turn the agent right, changing its direction, otherwise if agent.moving is True tries the action MOVE_FORWARD.
- STOP_MOVING (4) If agent.moving is True then becomes False. A penalty will be added. Stop the agent in the current occupied cell.

```

1 if agent is in DONE or in DONE_REMOVED (1th case) then
2   | no reward is computed
3   | return
4 if agent is in READY_TO_DEPART (2th case) then
5   | if the provided action is a MOVE_* type and the initial cell is free the agent
6   |   become ACTIVE and is initialized
7   | reward is computed
8   | return
9 if agent is in malfunction (3th case) then
10  | reward is computed
11  | return
12 if agent is at the beginning of a cell then
13  | update agent.moving considering the observations above depending on the
14  |   different action types.
15  | if agent.moving then
16  |   | the wanted action validity is first checked and if it is valid (considering
17  |     also the possibility to backup from an invalid MOVE_RIGHT or
18  |     MOVE_LEFT to a valid MOVE_FORWARD) action is stored otherwise
19  |     agent.moving becomes False and penalties are added, in this process
20  |     agent.moving and agent.speed_data['transition_action_on_cellexit'] are
21  |     updated
22 if agent.moving (4th case) then
23  | Updates the percentage of completion then if it is completely arrived on the
24  |   next cell, before updating the position, the direction and clears the
25  |   completion percentage it checks whether the new cell is free. Until the cell
26  |   remains occupied in the future executions the agent will repeat this process.
27  | reward is computed
28 else
29  | reward is computed
30 end

```

Algorithm 1: The *_step_agent* algorithm

Some useful questions:

- An agent can stop during an action between two cells? Absolutely no, it's like any other action.
- Requested actions during a malfunction are ignored? Yes.
- Requested actions during a not completed movement are saved for after execution? No, because the condition at line 11 is not executed and the conditions in line 16 check whether the cell is free and possibly complete agent data. Actions are not allowed to change within the cell, each agent can only chose an action to be taken

when entering a cell. This action is then executed when a step to the next cell is valid.

- How is possible to understand if an agent is ready to perform an action? The entry `info_dict["action_required"]` returned by the function `step` of **RailEnv** contains True for the given agent. This doesn't mean that the action will be successfully executed due to the presence of malfunctions or blocking agents, in this case `info_dict["action_required"]` will remain True.
- When an agent reaches DONE is removed automatically the following step? TODO Yes.
- Does an agent automatically pass from READY_TO_DEPART to ACTIVE at the beginning? TODO No. A MOVE_* is necessary.
- Do collisions occur? No, agents check if the cell is free before moving. Deadlocks are possible when two agents are one in front the other without any chance to change path.

Malfunctions

The strategy depends on the passed *malfunction_generator_and_process_data*. TODO

- A malfunction can occur during the resolution of another? TODO
- An agent could have a malfunction during the completion of an action between two cells? TODO Yes.

Speed

The different speed profiles (speed is between 0 and 1) can be generated using the `schedule_generator`. Speed configurations can be build using **ScheduleGenerators**. TODO

Rewards

The rewards are based on the following values:

- `invalid_action_penalty` which is currently set to 0, penalty for requesting an invalid action
- `step_penalty` which is $-1 * \alpha$, penalty for a time step.
- `global_reward` which is $1 * \beta$, a sort of default penalty.
- `stop_penalty` which is currently set to 0, penalty for stopping a moving agent
- `start_penalty` which is currently set to 0, penalty for starting a stopped agent

The full step penalty is computed as the product between `step_penalty` and `agent.speed_data['speed']`. There are different rewards for different situations:

- single agents that are in `DONE` or in `DONE_REMOVED` have zero reward (1th *_step_agent* case).
- all agents that have finished in this episode (checked at the end of the *step*) or previously (checked at the beginning of the *step*), have reward equal to the `global_reward` (when in *step* all agents have reached their target)
- full step penalty is assigned when an agent is `READY_TO_DEPART` and in the current turn moves or stay there (2th *_step_agent* case), or when is in malfunction (3th *_step_agent* case).
- full step penalty plus the other penalties (`invalid_action_penalty`, `stop_penalty` and `start_penalty`) when the agent is finishing actions or start new ones (4th *_step_agent* case). Currently the other penalties are all set to zero.

So each train starts counting rewards since the beginning, not since it becomes `ACTIVE`. Currently it is possible to say that rewards are always full step excluding the end of the episode and the single agents that have finished which have reward equal to 0.

2 Multi Agent Reinforcement Learning (MARL)

This section provides an overview of some useful works and theory behind MARL that we consider useful to suggest approaches or solutions to the Flatland Challenge.

2.1 Theoretical background

2.1.1 Introduction

"Specifically, MARL addresses the sequential decision-making problem of multiple autonomous agents that operate in a common environment, each of which aims to optimize its own long-term return by interacting with the environment and other agents."[1]

MARL algorithms can be divided into three groups [6]:

- **Fully cooperative**, where agents collaborate to optimize a common long-term return.
- **Fully competitive**, where the return of agents usually sum up to zero.
- **Mix of the two**, where both cooperative and competitive agents are involved.

We consider the Flatland Challenge a fully cooperative environment since each agent apparently compete to reach faster its destination and gain rails portions but the problem must consider the common time minimization goal. There exist two closely related theoretical frameworks for MARL [6]:

- **Markov/Stochastic Game.** All agents share the same state and differently from classical single agent’s **Markov Decision Process (MDP)**, the optimal performance of each agent is controlled not only by its own policy, but also the choices of all other players of the game. Usually agents share a common reward function but it is possible to have different functions like in the team-average reward setting.
- **Decentralized POMDP (Dec-POMDP).** As a MDP can be extended to a **Partially Observable Markov Decision Process (POMDP)** when information that can be accessed at a given state is incomplete, similarly a Markov Game can be extended to a Dec-POMDP. In multi-agent scenario, an agent may not only depend on the information it has autonomously gathered, it will also be influenced by the choices of other agents, which are partially observable. In Dec-POMDP each agent has its own local observation of the system state, that without other agents’ observations, leads to the impossibility to maintain a global belief state. To overcome this problem, agents can exploit levels of coordination among them to obtain the full observability of a state by combining the individual observations from each member [2]. Dec-POMDP approaches are usually considered more difficult to solve than others, especially when the number of agents is greater than two.
- **Extensive-Form Game** inspired from computational game theory, it handles imperfect information. It is usually used in mixed or competitive environments.

In the Flatland environment a state is surely represented by the positions and orientations of each agent in the map, since the railroad is static and agents influence it only by moving and interrupting paths. Flatland environment allows to personalize how agents perceive the world implementing custom observations.

The involvement of Deep Learning to tackle the problem of MARL defines a new specific subject called **Multi-agent Deep Reinforcement Learning (MADRL)**. [3] presents four categories of recent MADRL works:

- **Analysis of emergent behaviors**, in general, they do not propose learning algorithms, their main focus is to analyze and evaluate DRL algorithms in a multi-agent environment.
- **Learning communication**, they study communications techniques to share information.
- **Learning cooperation**, they directly explore approaches based on actions and observations to build multi-agent systems.

- **Agents modeling agents**, they study how agents reason about others to fulfill a task.

2.1.2 Challenges

MARL frameworks inevitably adds many challenging problems on the single-agent scenario.

"Learning in multiagent settings is fundamentally more difficult than the single-agent case due to the presence of multiagent pathologies, e.g., the moving target problem (non-stationarity) [2, 5, 10], curse of dimensionality [2, 5], multiagent credit assignment [31, 32], global exploration [8], and relative overgeneralization [33, 34, 35]."[3]

Below follow some problems that may affect the Flatland Challenge.

Non-Stationarity

"In Markov games, the state transition function T and the reward function of each agent r_i depend on the actions of all agents. During the training of multiple agents, the policy of each agent changes through time. As a result, each agents' perceived transition and reward functions change as well. Single-agent RL procedures which commonly assume stationarity of these functions might not quickly adapt to such changes."[5]

In a single-agent environment, an agent is concerning only the outcome of its own actions. In a multi-agent scenario, an agent observes not only the outcomes of its own action but also the behavior of other agents. Agents may interact with each other and learn concurrently leading to a continuous reshape of the environment and to non-stationarity [6]. For example the classical DQN does not provide working solutions, some derivations have been proposed to deal with this problem such as **Deep Repeated Update Q-network (DRUQN)** [2], **Deep Loosely Coupled Q-network (DLCQN)** [2] and **multi-agent concurrent DQN**. Other techniques to adapt classical experience replay to multi-agent environment have been proposed such as **Hysteretic-DQN (HDQN)** and **Lenient-DQN (LDQN)** [4].

There are different ways to tackle the non-stationary problem, as illustrated by [5].

- **Centralized Critic Architecture** based on an actor-critic algorithm. The critics' training is centralized and has access to the observations and actions of all agents, while the actors' training is decentralized. An example is the **Multi-Agent Deep Deterministic Policy Gradient (MADDPG)** algorithm. In MADDPG each agent uses a centralized critic and a decentralized actor. Since the training of each agent depends on the observations and actions of all the other agents, each agent perceives the environment as stationary.

- **Decentralized Learning** Techniques using self-play.
- Opponent Modelling.
- **Meta-Learning**.
- Communication. Either accessing hidden layers as **CommNet** or feeding other agents' neural networks as **Reinforced Inter-Agent Learning**.

Partial observability

As mentioned, in the single-agent scenario this type of problem is usually modelled with a POMDP. **Deep Recurrent Q-Networks (DRQN)** proposed using recurrent neural networks, in particular, Long Short-Term Memory (LSTMs) cells in DQN, to introduce a memory capability. An extension of DRQN for multi-agent environments is **Deep Distributed Recurrent Q-network (DDRQN)** [4]. Another technique to deal with partial observability is **Deep Recurrent Policy Inference Q-network (DRPIQN)** learned by adapting network's attention to policy features and their own Q-values at various stages of the training process. Methods to address partial observability in Markov Games usually involve communication (**Reinforced Inter-Agent Learning (RIAL)** and **Differentiable Inter-Agent Learning (DIAL)**) or parameter sharing (**PS-DQN**, **PS-DDPG**, **PS-A3C** and **PS-TRPO**). Alternatively to the Markov Game, the Dec-POMDP can be used to model this type of scenario in a more classical **centralized learning for decentralized execution** fashion or in a more modern **decentralized learning for decentralized execution**.

Scalability

To handle non-stationarity, each individual agent may need to account for the joint action space, whose dimension increases exponentially with the number of agents, this is also referred to as the combinatorial nature of MARL [6]. Many methods have been proposed to tackle this problem, one of them is the extension of **Curriculum Learning** for a multi-agent scenario.

Information Structures and Training Schemes

In the single-agent case is easier to understand what information is visible to the agent. In Markov games is sufficient to observe the current state, while on extensive-form games agents may need to recall the history of past decisions. In addition agents struggle to fully access information like rewards and policies of other agents, increasing the non-stationarity viewed by individual agents [4] [6]. These considerations led to the development of different training schemes such as centralized learning for decentralized execution, which originated from the works on the Dec-POMDP setting and has been

widely adopted in recent MADRL works, and **fully decentralized**. The former has become a standard as simulators are usually involved in MARL training, there are different types of communications between agents and the central controller such as centralized learning, concurrent learning and parameter sharing[4]. Usually decentralized settings may allow some sort of communication to address the non-convergence issue typical of the independent learning, this strategy is also referred as **decentralized setting with networked agents**.

Multi-agent credit assignment

2.1.3 Algorithms

References

- [1] Robert Babuska. A comprehensive survey of multiagent reinforcement learning. 38:156–172, 03 2008.
- [2] A. O. Castaneda. Deep reinforcement learning variants of multi-agent learning algorithms, 2016.
- [3] Pablo Hernandez-Leal, Bilal Kartal, and Matthew E. Taylor. A survey and critique of multiagent deep reinforcement learning. *Autonomous Agents and Multi-Agent Systems*, 33(6):750–797, Oct 2019.
- [4] Thanh Thi Nguyen, Ngoc Duy Nguyen, and Saeid Nahavandi. Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications. *IEEE Transactions on Cybernetics*, page 1–14, 2020.
- [5] Georgios Papoudakis, Filippos Christianos, Arrasy Rahman, and Stefano V. Albrecht. Dealing with non-stationarity in multi-agent deep reinforcement learning, 2019.
- [6] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms, 2019.