

UNIVERSITÀ DI BOLOGNA



School of Engineering

Deep Learning

**Project 5: Flatland Challenge**

Professor: **Andrea Asperti**

Students:  
**Giovanni Montanari**  
**Lorenzo Sarti**  
**Alessandro Sitta**

Academic year 2019/2020



# Abstract

The Flatland Challenge is a competition to foster progress in multi-agent reinforcement learning for any re-scheduling problem (RSP). The challenge addresses a real-world problem faced by many transportation and logistics companies around the world (such as the Swiss Federal Railways, SBB). Different tasks related to RSP on a simplified 2D multi-agent railway simulation must be solved: the proposed solutions may shape the way modern traffic management systems (TMS) are implemented not only in railway but also in other areas of transportation and logistics. This will be the first of a series of challenges related to re-scheduling and complex transportation systems [1].

# Contents

<b>1</b>	<b>Problem Description</b>	<b>5</b>
<b>2</b>	<b>Proposed Solution</b>	<b>6</b>
<b>3</b>	<b>Network model</b>	<b>8</b>
<b>4</b>	<b>Single Agent</b>	<b>10</b>
4.1	Reducing distance reward . . . . .	12
4.2	Standstill penalty reward . . . . .	13
<b>5</b>	<b>Multi Agent</b>	<b>15</b>
5.1	Training 3 agents: first attempt . . . . .	16
5.2	Training 3 agents: deadlock penalty implementation . . . . .	17
5.3	Training 3 agents: stop on switch penalty implementation . . . . .	18
5.4	Training 5 agents: global10 + deadlock . . . . .	19
5.5	Training 5 agents: mixed approaches (Alternate solo descr.) . . . . .	19
5.6	Training 10 agents: global10 + deadlock . . . . .	19
<b>6</b>	<b>Final Result</b>	<b>20</b>
	<b>Conclusions</b>	<b>22</b>
	<b>References</b>	<b>23</b>

# Chapter 1

## Problem Description

The Swiss Federal Railways (SBB) operate the densest mixed railway traffic in the world. Due to the growing demand for mobility, SBB needs to increase the transportation capacity of the network by approximately 30% in the future. The "Flatland" Competition aims to address the vehicle rescheduling problem by providing a simplistic grid world environment[1].

Even if the challenge is open to any methodological approach, in this work the problem has been faced by means of reinforcement learning techniques: in a 2D grid environment representing the railway network, multiple agents with different goals must collaborate to maximize a global reward.

## Chapter 2

# Proposed Solution

The main problem of vehicle rescheduling is suitable for a Reinforcement Learning solution, where a large number of different environments are generated and trains (agents) are free to explore them and learn from the action they take, by means of reward functions, indicating positive and negative moves.

The considerable complexity of this task leads to start with a simpler scenario in which a single agent must reach his goal by following the shortest possible path, and then to consider a multi-agent scenario only after this problem is well-solved. The solution presented in this report is based on the *TreeObservation* proposed by the Flatland challenge with  $tree\_depth = 2$  and on a shortest path predictor with  $max\_depth = 20$ , as can be visualized in Figure 2.1 The followed approach is aimed at improving the performance as the reward functions vary, with particular attention to the scalability of the network when trained on more or less complex environments.

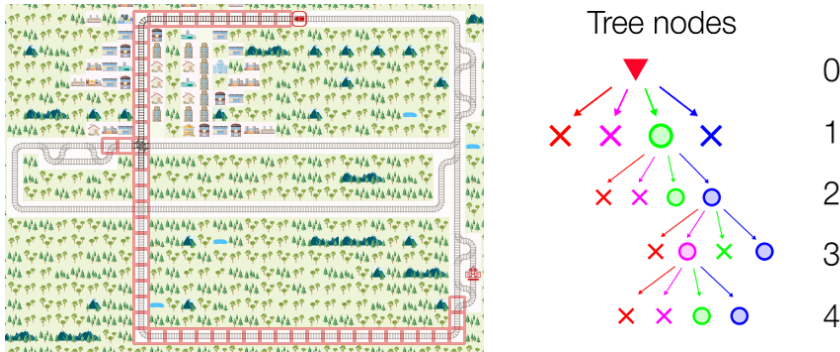


Figure 2.1: On the left, visual representation of the *TreeObservation* of an agent with  $tree\_depth = 2$  and  $max\_depth = 20$ . On the right, Tree nodes of a generic railway network.

For the single agent case, only a penalty for each time step elapsed has been considered first, but clearly that's not representing the optimal solution. To improve the results, an important reward associated to the achievement of the destination has been added, as well as a reduced time-step penalty when the train gets closer to the target. Moreover, a standstill penalty has been evaluated, together with a lower decay rate for the epsilon variable, with the intention of encouraging a bit more the exploration of the map.

The model corresponding to the single agent best performance has been selected as the starting point for the multi-agent solution and evaluated. Some modifications on the reward function have been applied, mainly with the aim of strongly penalizing deadlocks configuration, getting to an improvement in terms of trains trajectories rescheduling. After some tests for the comparisons between different reward functions, involving penalties for stopping on switches, the optimal one for the multi-agent setup has been chosen. Initially, a simple 3-agents scenario with a small map was considered, then the complexity of the task, which depends on the number of trains, available railways and environment size, has been gradually increased, up to 10 trains. In order to find the optimal tradeoff between number of deadlocks and number of agents reaching their station, many training processes have been performed, including "multi-level" training approach.

the proposed solution has been tested in both low and high complexity environments, initially considering only one possible speed and without malfunctions. The tests were then extended to the case with four different velocities and a non-zero malfunction rate.

## Chapter 3

# Network model

The Network architecture that has been chosen is the Dueling Double DQN (Deep Q-Network), which is at the state of the art of Value-based Reinforcement Learning techniques [2] [3]. The usage of this network has been deployed to solve a variety of games and problems, such as Doom Deathmatch [4].

In Flatland's framework, that has been considered, the Network does not include the three convolutional layers as the input of the neural network is not an image but a flattened vector, built starting from the observation (see parte di lollo). The implemented Network is shown in Figure 3.1.

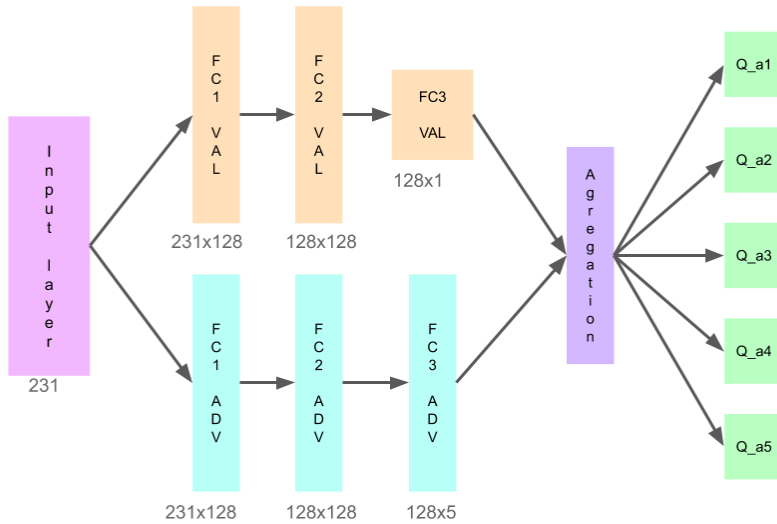


Figure 3.1: Dueling Double DQN

This choice reduces the complexity of the model as the useful information is already extracted and listed in the state. In fact the state,



considering the TreeObservation with a maximum depth of 2 [5], results in a vector of 231 elements encoding information on the environment.

## Chapter 4

# Single Agent

The first problem that has been considered is the one with a single agent involved. The idea was to verify that the choice of the network and observations was good enough to achieve the goal and see if satisfactory results can be achieved in this simplified framework.



Figure 4.1: Single agent in a  $35 \times 35$  environment

The first explorative training has been performed considering:

- $x\_dim = 35$
- $y\_dim = 35$
- $n\_agents = 1$
- $max\_num\_cities = 3$
- $max\_rails\_between\_cities = 2$
- $max\_rails\_in\_city = 3$

where  $x\_dim$  and  $y\_dim$  define the dimensions of the map and the last three values are parameters used by *sparse\_rail\_generator*, a railway generator provided by the organizers of the Challenge that mimics the real structure of a railway [5]. The velocity for the train is set to 1, this value is used by the *sparse\_schedule\_generator* to set properly the agent's goal, [5]. The observation that has been considered is the *TreeObservationForRailEnv*, [5], with  $max\_depth = 2$ .

The default rewards consisting in:

- +1 when the goal is reached
- -1 given at each step to the agent

have been considered. The result of the training, considering 7000 iterations and a diminishing epsilon Figure 4.3 which is important for the exploration/exploitation trade-off [4], is shown in Figure 4.2.

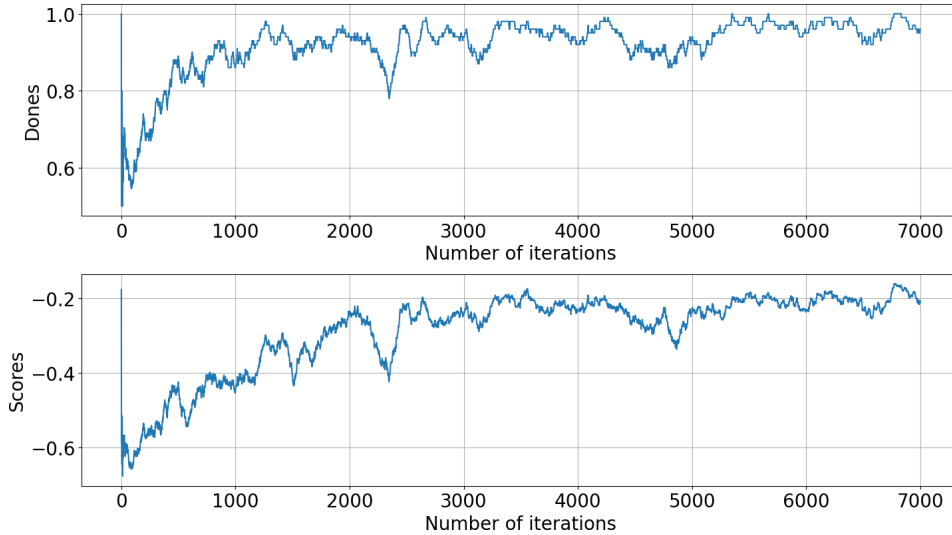


Figure 4.2: Dones and Scores metrics for the single agent first trial

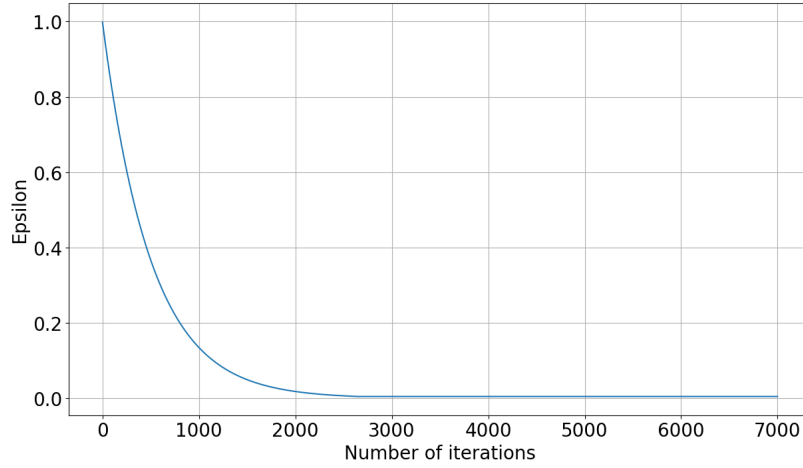


Figure 4.3: Decaying epsilon

The result seemed promising as the 100% of dones has been achieved at some iterations. Moreover, even with a local observation it has been possible to achieve good results, which confirms that at least in simple frameworks a global observation is not the only way to get results. At that point, since there was still room for improvements also in this simplified framework, some modifications have been tested and will be presented in the next sections.

## 4.1 Reducing distance reward

The first intuition was that, even if the state contains information about the agent's goal, it can be interesting to modify the step penalty in order to penalize less the agent in case it gets closer to the target. The same idea was common also to the Netcetera's group which reached the 5th position in the 2019 Challenge [6]. The advantage in introducing this penalty should be to speed up the learning, improve dones and scores and force the agent to move. The choice for the rewards has been:

- +10 when the goal is reached
- -1 given generally at each step to the agent
- -0.5 given instead of the general one, in case the agent gets closer to the target

The choice of increasing the target reward, together with the choice

of the reducing distance rewards, contributes to the result shown in Figure 4.4.

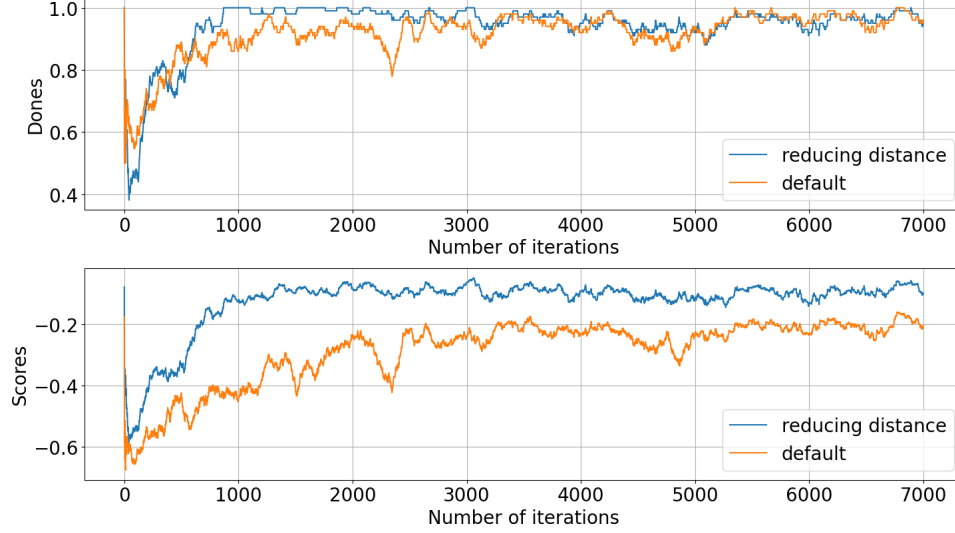


Figure 4.4: Dones and Scores for the single agent reducing distance

The 100% of dones is reached sooner and it is maintained until the performance degrades over 4000 iterations and the result becomes similar to the default one. It is worth mentioning that the scores improves significantly but this is partially due to the fact that the reducing distance penalty is less than -1.

The choice of the epsilon has been equal to the one shown in Figure 4.3, thus at less than 3000 iterations its value becomes 0.

In conclusion, in this single agent scenario the introduction of the reducing distance penalty does not degrades the performances in training but it improves them both in the dones plot and in the scores one.

## 4.2 Standstill penalty reward

Another intuition, coming from the observation of the behaviour of the agent trained using default parameters, is that there ins no reason why the train should be stopped at any time, since it will only be a cost and it wouldn't contribute in reaching the target.

Therefore it has been introduced a negative reward of -1 in case the agent does not perform any movement but it stays in standstill state. The result and comparison with the other results is shown in Figure 4.5.

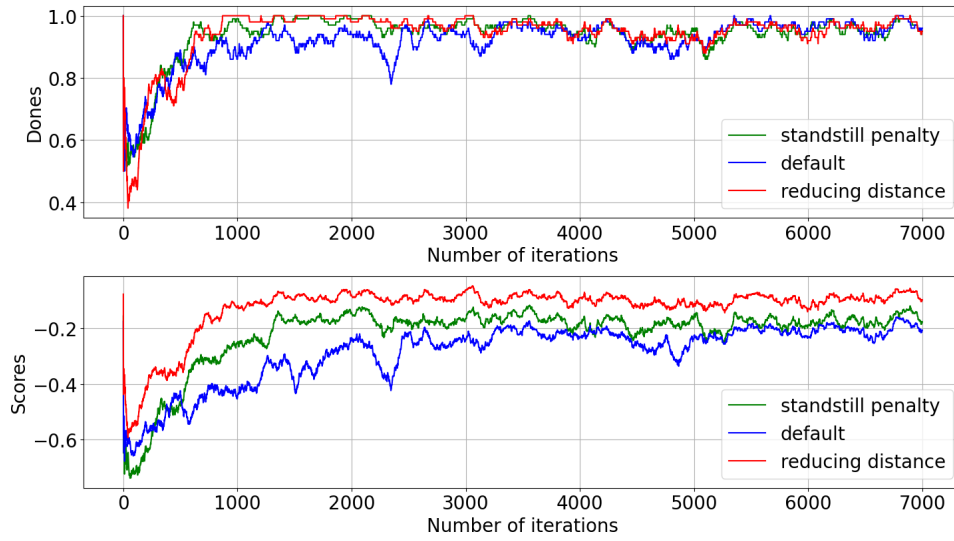


Figure 4.5: Dones and Scores for the single agent with standstill penalty

The result is not so different with respect to the reducing distance one, this is reasonable as the lower penalty should teach the agent to move at each step in the target direction and not to stay standstill. However, it seems that this new reward brings less information with respect to the reducing distance one, having also worst performances both in the dones and scores metrics. Therefore, the reducing distance rewards should be preferred.

The results on the single agent seems satisfactory as almost 100% of dones and a low value of scores has been achieved, moreover different possibilities has been explored and the chosen network model and observations gave positive feedback. Now, it is time to introduce higher complexities and move to the multi agent scenario.

## Chapter 5

# Multi Agent

After reaching acceptable results with single agents, the training environment has been revisited to allow a 3 agents training, with the following parameters:

- $x\_dim = 40$
- $y\_dim = 40$
- $n\_agents = 3$
- $max\_num\_cities = 4$
- $max\_rails\_between\_cities = 2$
- $max\_rails\_in\_city = 3$

First, it’s important to point out that the observations used in all of the following trainings also include a *predictor*, more specifically the *ShortestPathPredictor* with a depth of 20. This object, included in the Flatland repository, is useful only in a multi-agent setup. Indeed it predicts the agents shortest paths to destination to allow the observer to detect possible future conflicts. This gives some useful information for the agents interactions.

The first try was to start a training following the improvements learned with the single agents. Indeed two trainings (4000*iterationseach*) were compared, the first one setting the *reward* for reaching the goal to 10 and the latter also exploiting the *reduced step penalty* when agents' direction is towards the destination. The comparison, in 5.2 shows a counter-intuitive result with respect to what obtained with a single agent; indeed giving a reduced step penalty worsens the performance



of 3 agents.

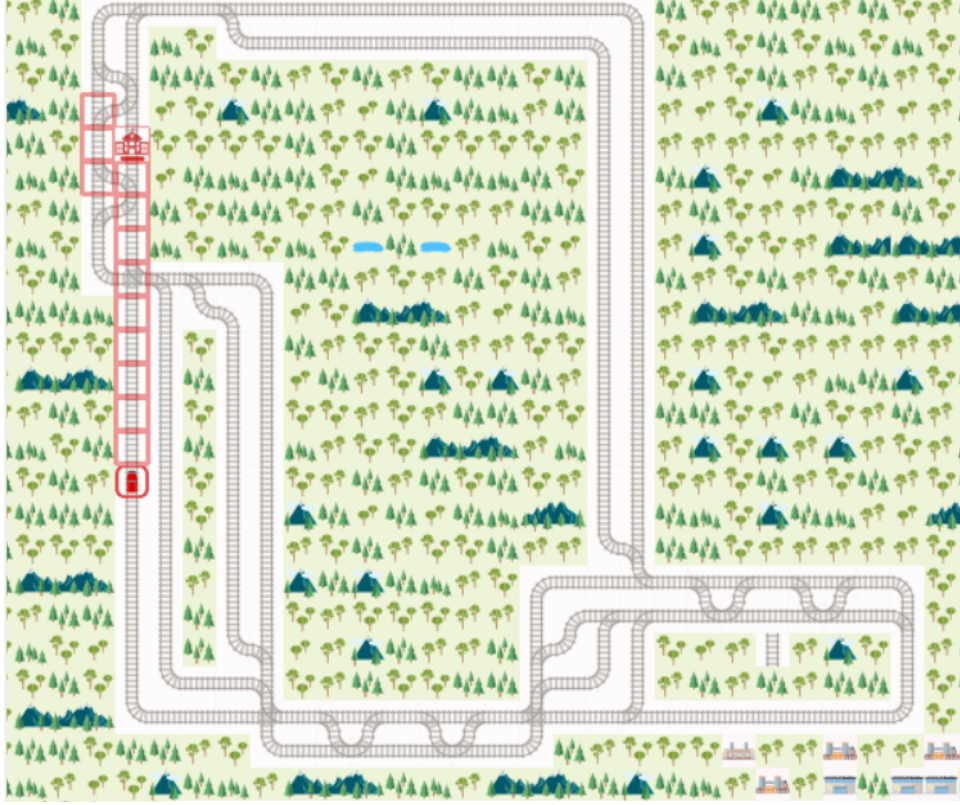


Figure 5.2: 3 agents comparison global10vsRedDist

Even the best between these two trainings didn't give good results, indeed the tests (see Results) reached a *som* of Dones. The main problem was in the interactions of the trains, which didn't seem to avoid each other, reaching deadlocks in several situations.

## 5.2 Training 3 agents: deadlock penalty implementation

In order to solve this issue, we decided to feed the agents with some information to understand whether they are in a deadlock. This info comes as a serious *penalty*, set to  $-10$  and implemented into *rail.env.py* from the Flatland repository. As expected, results were very encouraging, reaching the following learning rate compared to the previous case (5.3). Agents learned to avoid each other by choosing alternative rails, and more surprisingly to sometimes stop before a switch to let others

pass.

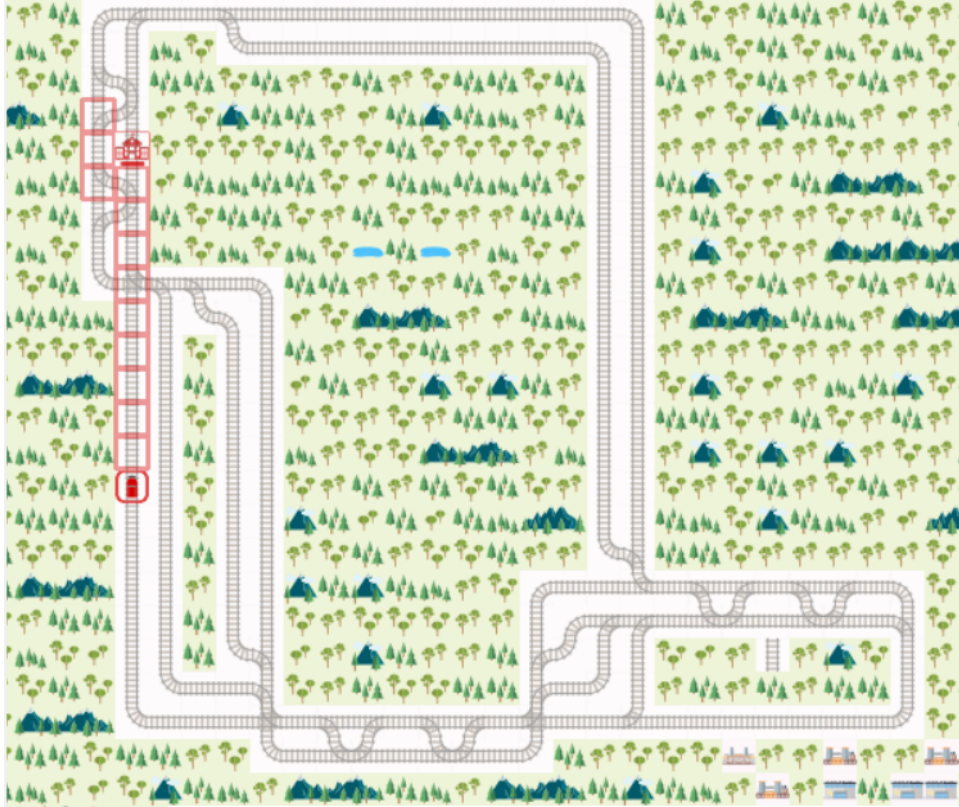


Figure 5.3: 3 agents training: deadlock penalty improvement

In the test, we reached the following scores, shown in Figure 5.4.

### 5.3 Training 3 agents: stop on switch penalty implementation

FORSE: altrimenti mettere in multi5.

Team Members	Metrics (avg over N tests)	No malfunctions - 1 speed (1.0)						Malfunctions - 4 speeds (0.25 each)		
		rail_env - Medium size			rail_env - Big size test			rail_env - Medium size		
		3 agents	5 agents	6 agents	5 agents	7 agents	10 agents	3 agents	5 agents	6 agents
Team 1	average % of Dones average % of Deadlocks mean normalized return									
Team 2	average % of Dones average % of Deadlocks mean normalized return									
Team 3	average % of Dones average % of Deadlocks mean normalized return									

Figure 5.4: TABELLA TESTS: 3 agents

- 5.4 Training 5 agents: global10 + deadlock
- 5.5 Training 5 agents: mixed approaches (Alternate solo descr.)
- 5.6 Training 10 agents: global10 + deadlock

## Chapter 6

# Final Result

The results obtained by training the network with 10 agents on the "Big Size" environment were the best and are presented in this section. The associated rewards/penalties are:

- *step\_penalty* = -1
- *global\_reward* = 10
- *deadlock\_penalty* = -10

The trained network has been tested on different sized maps and with different numbers of trains, initially with no malfunctions and with all the agents going at the same speed (1.0). In particular, two set of maps were selected, by mean of the official Flatland environment, specifically *RailEnv* from *flatland.envs.rail\_env.py* and *sparse\_rail\_generator* from *flatland.envs.rail\_generators* [5]. These two test frameworks are described below, according to the parameters specified in the *rail\_env* file:

- *rail\_env* - Medium Size:
  - *x\_dim* = 16\*3
  - *y\_dim* = 9\*3
  - *n\_agents* = 3 - 5 - 6
  - *max\_num\_cities* = 5
  - *max\_rails\_between\_cities* = 2
  - *max\_rails\_in\_city* = 3
  - *rand\_seed* = 14

- *rail\_env* - Big Size:
  - $x_{dim} = 16*4$
  - $y_{dim} = 9*4$
  - $n_{agents} = 5 - 7 - 10$
  - $max\_num\_cities = 9$
  - $max\_rails\_between\_cities = 5$
  - $max\_rails\_in\_city = 5$
  - $rand\_seed = 14$

500 iterations are considered for each test session and the scores are collected in Figure 6.1.

Metrics (avg over N tests)	No malfunctions - single speed (1.0)						Malfunctions - 4 speeds (0.25 each)		
	rail_env - Medium Size			rail_env - Big Size			rail_env - Medium Size		
	3 agents	5 agents	7 agents	5 agents	7 agents	10 agents	3 agents	5 agents	7 agents
fraction of done-agents	0.908667	0.8436	0.784	0.8236	0.787714	0.7016	0.612667	0.5324	0.459714
fraction of deadlocks	0.015333	0.038	0.070286	0.0344	0.053143	0.1112	0.046	0.1048	0.116
mean normalized return	-0.266862	-0.323886	-0.381059	-0.392504	-0.428224	-0.491795	-0.316406	-0.362855	-0.400246

Figure 6.1: caption

# Conclusions

# References

- [1] *Flatland Challenge*. URL: <https://www.aicrowd.com/challenges/flatland-challenge>.
- [2] Hado V. Hasselt. “Double Q-learning”. In: *Advances in Neural Information Processing Systems 23*. Ed. by J. D. Lafferty et al. Curran Associates, Inc., 2010, pp. 2613–2621. URL: <http://papers.nips.cc/paper/3964-double-q-learning.pdf>.
- [3] Ziyu Wang, Nando de Freitas, and Marc Lanctot. “Dueling Network Architectures for Deep Reinforcement Learning”. In: *CoRR* abs/1511.06581 (2015). arXiv: 1511.06581. URL: <http://arxiv.org/abs/1511.06581>.
- [4] *Deep Reinforcement Learning Course*. URL: <https://www.freecodecamp.org/news/improvements-in-deep-q-learning-dueling-double-dqn-prioritized-experience-replay-and-fixed-58b130cc5682/>.
- [5] *Flatland Challenge Documentation*. URL: <http://flatland-rl-docs.s3-website.eu-central-1.amazonaws.com/>.
- [6] *Netcetera Blog*. URL: <https://blog.netcetera.com/leverage-reinforcement-learning-for-building-intelligent-and-adaptive-trains-that-can-successfully-9f4cdef80162>.