

1 Markov Decision Process

1.1 Description

The Markov decision process is a process where the action of an agent only depends on the current state of the environment. There are a set of possible actions and a set of possible states. There is also an agent policy which chooses an action based on the given state. The agent receives information about the state, performs the action provided by its policy and may receive a reward for the action. The state is updated and executed another time. The MDP has fixed transition probabilities which depend as said only on the current state. The reward probability is also fixed.

The agent's policy will be adapted to gain a maximum Return meaning a maximum future reward. The reward is estimated by estimating the Return for the next step and adding the reward for the transition.

The best estimates for this transition are stored as V values depending on the state or Q depending on both the state and the action.

1.2 Example for a MDP

A robot navigates through a maze, the states are the positions and orientations of the robot. The actions are the movement commands. The reward is given based on the proximity to the goal position.

1.3 Example for a POMDP

A real life application where the robot from above is only able to measure a few of the values which are defining its state. Therefore it is only able to approximate a distribution of possible states and apply a more general form of the MDP.

2 Exploration vs. exploitation

In reinforcement learning exists a tradeoff between exploration and exploitation. Exploitation means optimizing the currently best working policy towards a local maximum. Exploration on the other hand means trying out new or worse strategies to explore unknown states and skip local maxima in the best case towards a global maximum.

This tradeoff can be implemented using a ϵ -greedy action selection, where a random action is chosen with a probability of $1 - \epsilon$ for exploration purposes. If no random action is selected the currently best fitting action is chosen.

An agent may be failing or take a very long time to converge if it explores too much. On the other hand, exploiting too much results in a local maximum

of the reward, meaning it maybe fails to choose an action that is suboptimal at the moment but leads to a much higher reward later on.

3 SARSA and Q-Learning

3.1 is temporal difference (TD) learning?

TD Learning is a class of model-free reinforcement learning methods. TD learning adjusts the estimated Return every few steps (e.g. every step) instead of waiting for a final Outcome. TD-learning can be used in combination with Neural Networks to optimize the adjusting of the estimates.

3.2 How does the SARSA algorithm work?

The SARSA Algorithm is a TD-Learning method. The Algorithm is used again and again with many initial states to learn the entire state space well. SARSA initializes a Round (Trial) by selecting an Action a for an initial State s und calculating the corresponding Q-Value. Next it repeats until the Round is over when e.g. the final State is reached or it got Out-Of-Bounds. Next it executes the selected Action a , reads the Reward r and the reached State s' and selects the next Action with the current Target Policy. It calculates the TD-Error for this iteration from the State s , the Action a , the Reward r obtained from this transition, the next State s' and the next Action a' , hence the acronym SARSA. The Q-Value for the State s and Action a is updated with the previously calculated TD-Error. It provides the latest Values for the next Iteration by updating the Variables $Q(s, a)$, s and a with the Values of $Q(s', a')$, s' and a' .

3.3 How does Q-learning work, and how is it different from SARSA?

Q-Learning works exactly the same as SARSA except that the TD-Error is calculated differently. SARSA is a On-Policy Algorithm as it uses the actually selected Action and State from the Policy for calculating the Q-Value. Q-Learning is an Off-Policy Algorithm instead and uses the best possible Q-Value for all next Actions and States instead of the ones selected by the Target Policy.

4 Implementing Q-Learning

The values in the Q-Table describe the estimated Return for taking a specific Action in a specific State/Observation. Looking at the output summary of `q.table` it seems that dropping passengers is penalized compared to moving

around and picking passengers up which is expected as dropping an passenger on a invalid field is penalized with a Reward of -10.

The updated Code:

```
1 # This is the main loop.
2 num_train_episodes = 1000 # Modify to tune for convergence
3 render = True # Set to False to speed up training
4 returns = []
5 # This creates our Q-table: A two-dimensional array of all
   zeros.
6 q_table = np.zeros([num_observations, num_actions])
7 print('Shape of the Q-table:', q_table.shape)
8
9 def best_action(observation):
10     q_value, action = float("-inf"), -1
11
12     actions = q_table[observation]
13
14     for a, q in enumerate(actions):
15         if q > q_value:
16             q_value = q
17             action = a
18
19     # select best action with epsilon greed policy
20     best_probability = 1 - epsilon
21     other = epsilon / (num_actions - 1)
22
23     probabilities = []
24
25     for a in range(num_actions):
26         if a == action:
27             p = best_probability
28         else:
29             p = other
30         probabilities.append(p)
31
32     action = np.random.choice(list(range(num_actions)), p=
probabilities)
33
34     # return the action and its current q_value for the given
state
35     return action, actions[action]
36
37 for episode in range(1, num_train_episodes):
38     # By resetting the environment you also
39     # get the initial observation or "state".
40     observation = env.reset()
41
42     # Repeat until the environment is done.
43     done = False
44     episode_return = 0
45
```

```
46     # ME: select best action for initial state with epsilon-
greedy
47     action, q_value = best_action(observation)
48
49     while not done:
50         # Perform the chosen action in the environment
51         # to get reward values and the next observation.
52         next_observation, reward, done, _ = env.step(action)
53
54         # Update your Q-table based on the information you
gathered!
55         next_action, next_q_value = best_action(
next_observation)
56
57         # ME: calculate the td_error depending
58         # on the maximum q_value for the next state
59         best_next_q = np.max(q_table[next_observation])
60         td_error = reward + (gamma * best_next_q) - q_value
61         # ME: update the q_table for the current state and
action
62         q_table[observation, action] = q_value + (alpha *
td_error)
63
64         # Step the observation
65         observation = next_observation
66         # ME: update action
67         action = next_action
68         # ME: update q_value
69         q_value = next_q_value
70
71         # For performance reasons, only display every 200th
episode.
72         episode_return += reward
73         if render and episode % 200 == 0:
74             clear_output(wait=True)
75             env.render()
76             print("Episode: ", episode)
77             print("Reward: ", reward)
78             print("Return so far: ", episode_return)
79             # Wait a little bit. We don't want to go too fast.
80             sleep(0.1)
81     returns.append(episode_return)
```

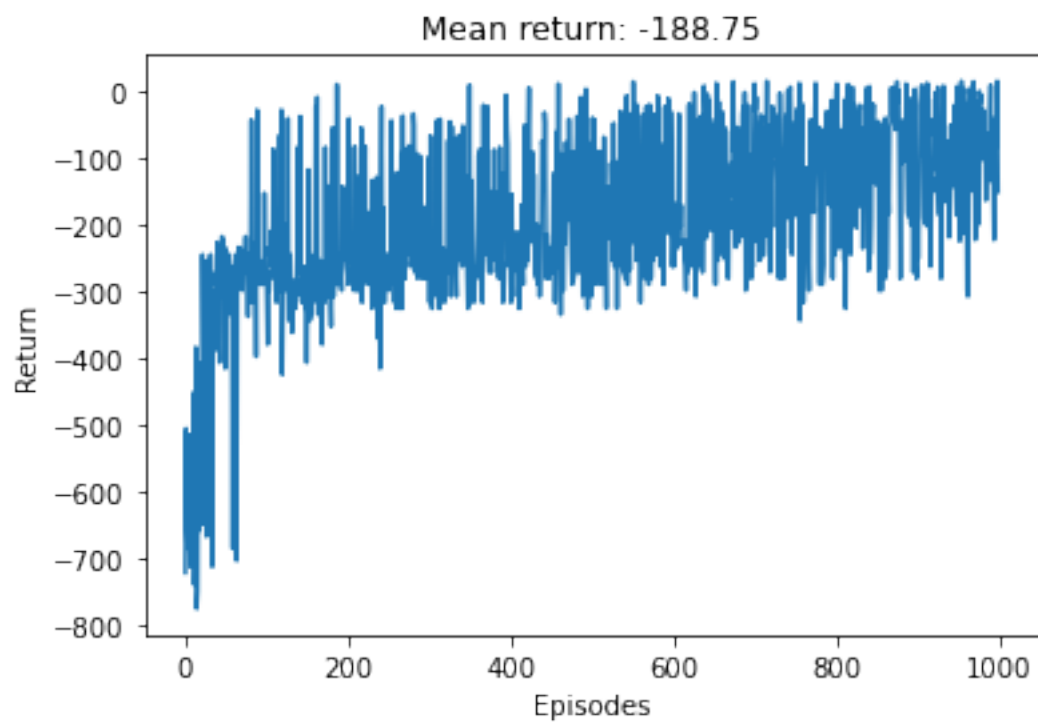


Figure 1: 1000 Episodes. In the first 100 Episodes the Return increases dramatically and slows down afterwards. It does not seem to converge yet.

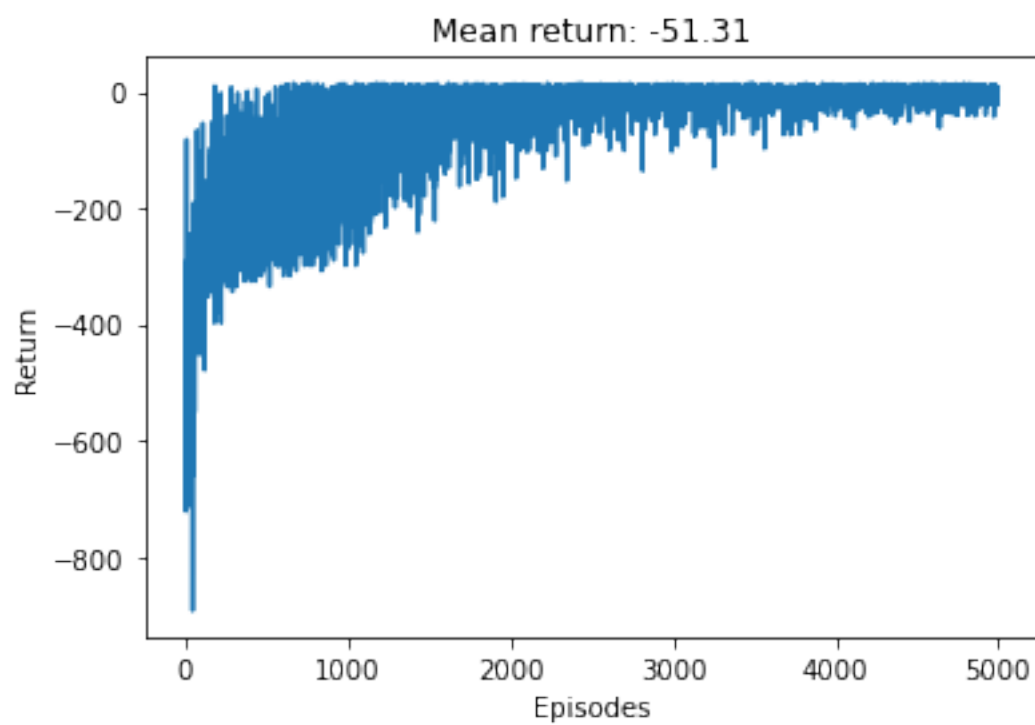


Figure 2: 5000 Episodes. It does not seem to converge around zero which is expected as the maximum Return is at most 20 for dropping the passenger (if picking up and dropping happens on the same field).

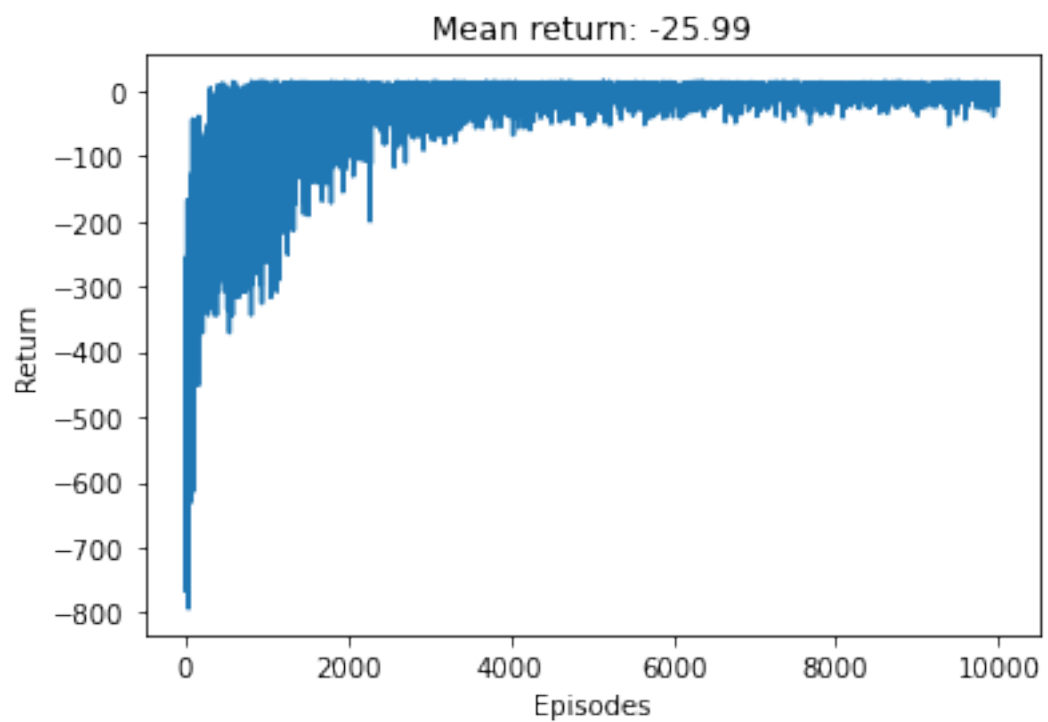


Figure 3: 10000 Episodes. No major improvements are seen compared to 5000 Episodes. The Agent seems to have converged.