

Session 10

Review Practice Assignments

jmanero@faculty.ie.edu

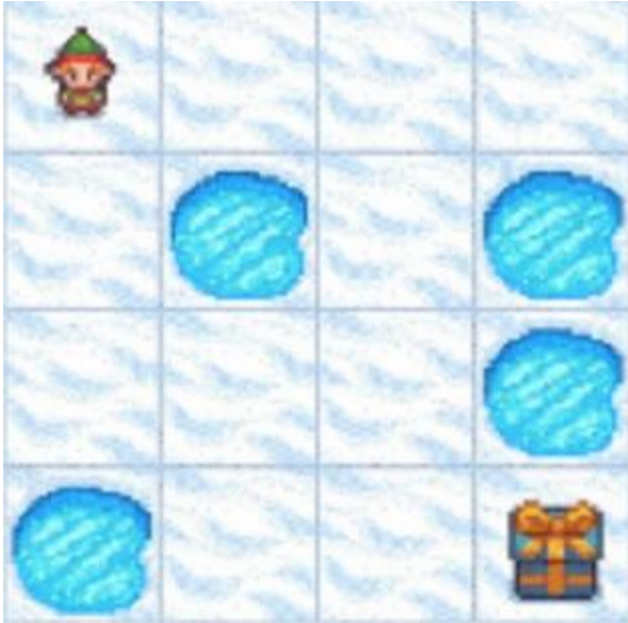
Let's see where we are

- Assignment 1 - Jump start
- Assignment 2 - Dynamic Programming
- Assignment 3 - Model Free models
- Assignment 4 - DQN / DDQN
- Group Practice
- Mountain Cart Challenge

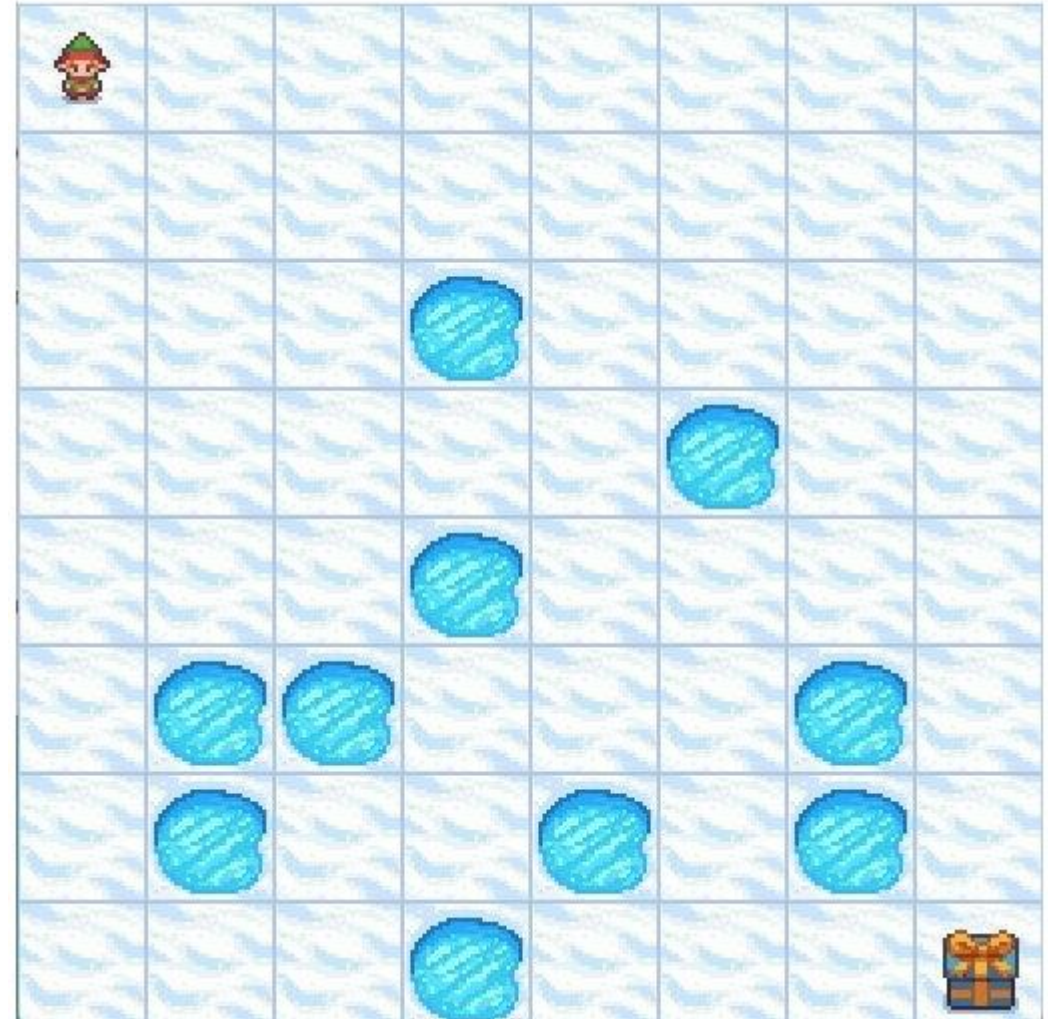
- All methods have different learning strategies
 - MC
 - TD
 - SARSA
 - Q-Learning

Practices Review

Frozen Lake

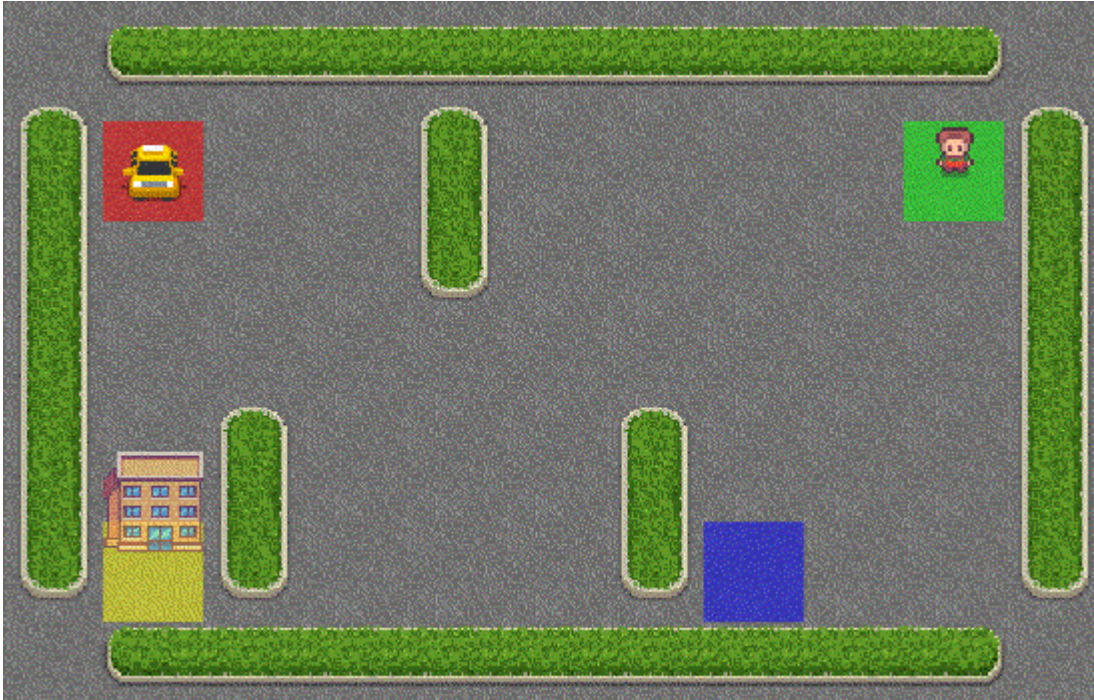


- 2 Dimensional
- Fixed start - end-point
- Stochastic and non-stochastic
- 8x8 large state space



Practices Review

Taxi



- Large observation space
- Grid 5x5
- State is encoded
 - Taxi row (5 possible positions)
 - Taxi column (5 possible positions)
 - Passenger location (5 possible values: at 4 locations or in the taxi)
 - Destination (4 possible values)
- 4 pick-up/drop-off places (red green yellow and blue)
- Taxi starts in any place
- Final cumulative reward around 8 if done properly

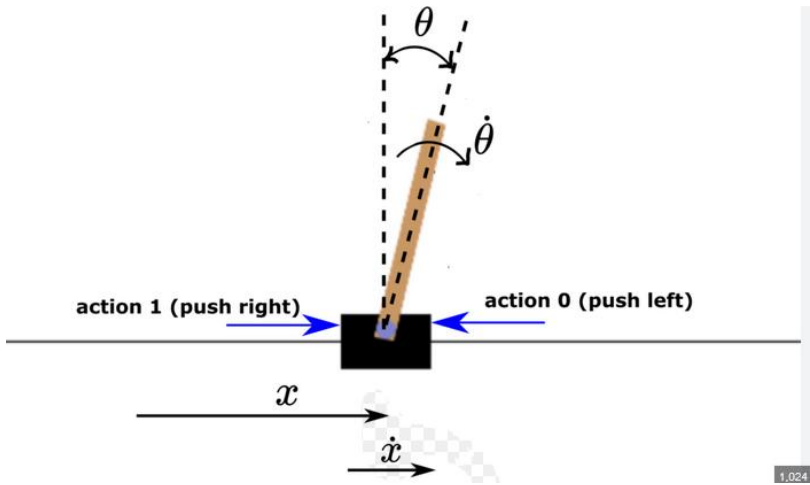
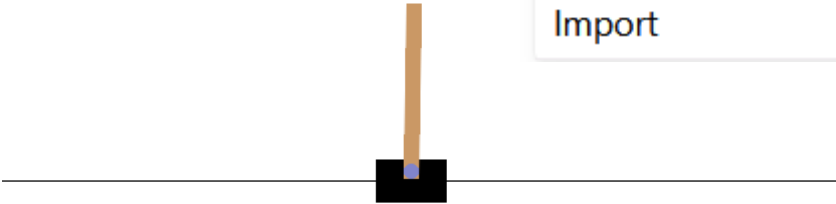
- The policy is difficult to visualize!!!

Practices Review

CartPole

Action Space	Discrete(2)
Observation Shape	(4,)
Observation High	[4.8 inf 0.42 inf]
Observation Low	[-4.8 -inf -0.42 -inf]
Import	<code>gym.make("CartPole-v1")</code>

Num	Observation	Min	Max
0	Cart Position	-4.8	4.8
1	Cart Velocity	-Inf	Inf
2	Pole Angle	~ -0.418 rad (-24°)	~ 0.418 rad (24°)
3	Pole Angular Velocity	-Inf	Inf



- Requires binning to use the continuous values in the observation for Q-Learning
- It may take time to converge
- Reward +1 for every step

Review Methods and environments

Methods and Environments

Monte Carlo

Feature	Experience	Conclusion
Speed to converge		
Use in several environments		
What went well and what went bad		
Overall rating		

Methods and Environments

SARSA

Feature	Experience	Conclusion
Speed to converge		
Use in several environments		
What went well and what went bad		
Overall rating		

Methods and Environments

Q-Learning

Feature	Experience	Conclusion
Speed to converge		
Use in several environments		
What went well and what went bad		
Overall rating		

Remember

- SARSA shares similarities with Q-learning
- Both algorithms aim to learn the optimal action-value function and use similar update rules. However, the key difference lies in their learning strategies.
- While Q-learning is an off-policy algorithm that updates the Q-values based on the maximum expected future reward, regardless of the action taken, SARSA updates the Q-values based on the actual action taken by the current policy. This difference can lead to different learning dynamics and performance characteristics in certain problems.

Methods and Environments

What should we observe

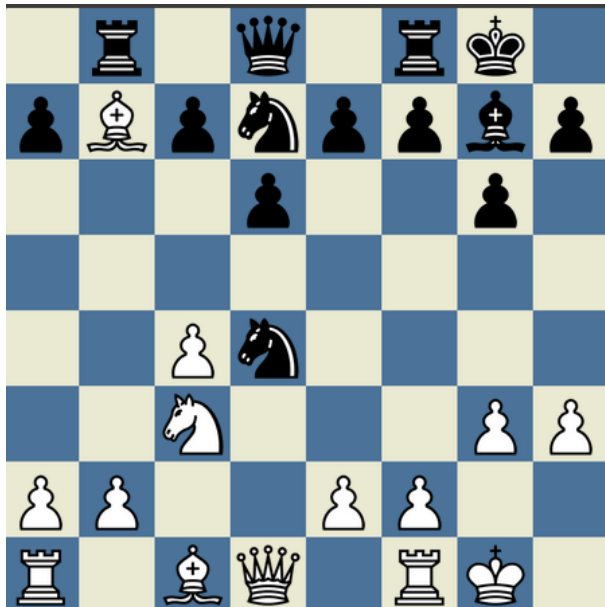
Aspect	Q-Learning	SARSA
Policy Type	Off-policy	On-policy
Next Q-value Estimation	Uses $\max Q(s', a')$ (greedy action)	Uses $Q(s', a')$ of the actual action taken
Exploration vs. Exploitation	Encourages more exploitation (greedy updates)	Encourages more exploration (policy-following updates)
Risk Sensitivity	More aggressive, assumes optimal future actions	More conservative, considers the actual exploration strategy
Convergence Speed	Typically, faster but less stable in stochastic environments	More stable but slower

Rewards

Reward Shaping

is the use of small intermediate 'fake' rewards given to the learning agent that help it converge more quickly.

- Each environment has its reward already shaped
- We may try 'reward shaping' to improve learning
- For instance, in Chess how do we evaluate the reward?



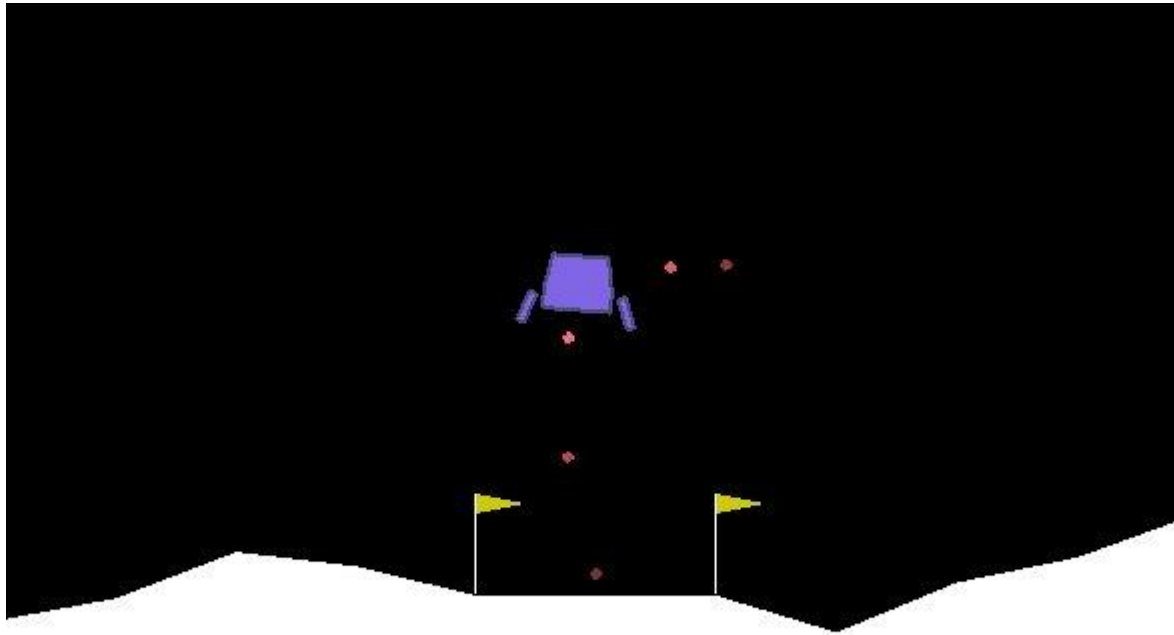
- King's safety.
- Material on the board.
- Pieces activity.
- Pawn structure.

But... How to define a quantitative measure?

Reward Shaping

Lunar Lander

- Sparse rewards (big reward far away)
- No truncation of the episode / infinite fuel
- May end-up learning to hover

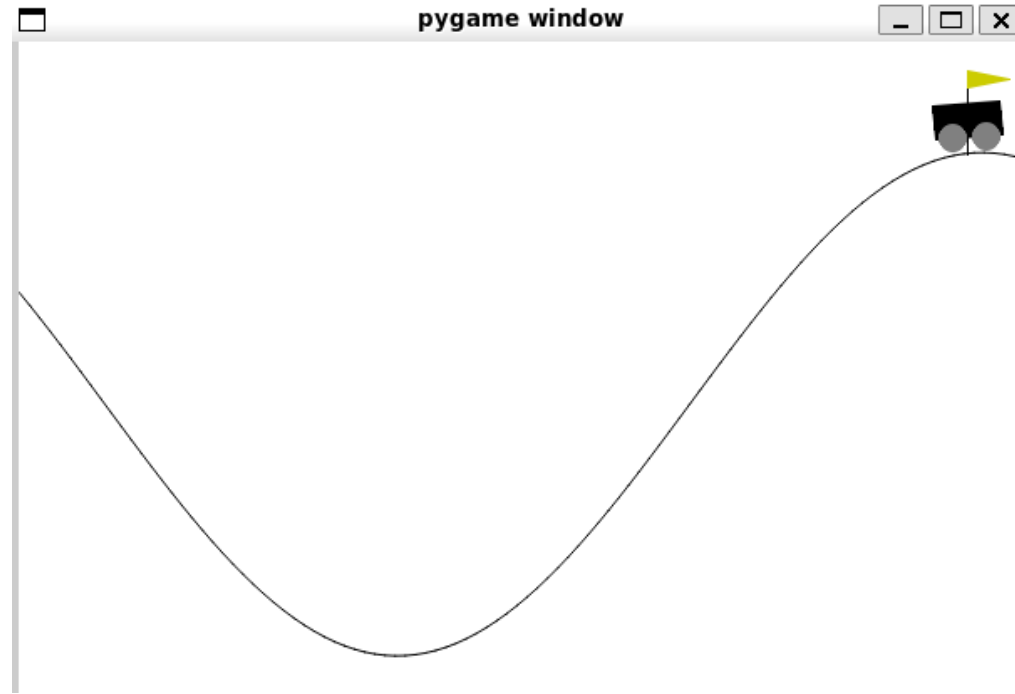


Solutions

- Reward shaping
- Include fuel consumption
- Increase gravity

Reward Shaping

Shaping Rewards – MOUNTAIN-CART



This environment is part of the Classic Control environments which contains general information about the environment.

Action Space	<code>Discrete(3)</code>
Observation Space	<code>Box([-1.2 -0.07], [0.6 0.07], (2,), float32)</code>
import	<code>gymnasium.make("MountainCar-v0")</code>

Reward Shaping

For the Assignment MOUNTAIN CAR or CARTPOLE Discretization

```
1 #Support Funtions
2
3
4 # Function to create bins
5 def create_bins(interval, num):
6     return np.linspace(interval[0], interval[1], num + 1)
7
8 # Updated intervals and bin sizes for discretization
9 intervals = [(-2.4, 2.4), (-3.0, 3.0), (-0.5, 0.5), (-2.0, 2.0)]
10 nbins = [12, 12, 24, 24] # Increased bins for finer state representation
11 bins = [create_bins(intervals[i], nbins[i]) for i in range(4)]
12
13 # Function to discretize state variables into bins
14 def discretize_bins(x):
15     return tuple(np.clip(np.digitize(x[i], bins[i]) - 1, 0, nbins[i] - 1) for i in range(4))
```

Reward shaping

Reward Shaping in Mountain Car

```
27
28 ##### SHAPING REWARDS #####
29
30     shaped_reward = reward
31     if done and step < 200:
32         # If episode is ended the we have won the game. So, give some large positive reward
33         shaped_reward = 250 + shaped_reward
34
35         # Velocity is important, we give positive reward for velocity (is the sign correct?)
36         velocity = next_obs[1]
37         distance = next_obs[0]
38         shaped_reward = shaped_reward + 1* abs(distance)
39     #     shaped_reward = shaped_reward + 10 * abs(velocity)
40
41 ##### END SHAPING REWARDS #####
42
```

2 Problems discretization and Shaping rewards in the Challenge

Discretization

- How many bins are ok?
- Hint: use this discretization strategy for the CARTPOLE Assignment !

Shaped Rewards

- -1 each timestep
- This reward does not allow good learning
- There is truncation and termination (careful)
- Can we define a shaped reward?

Objective

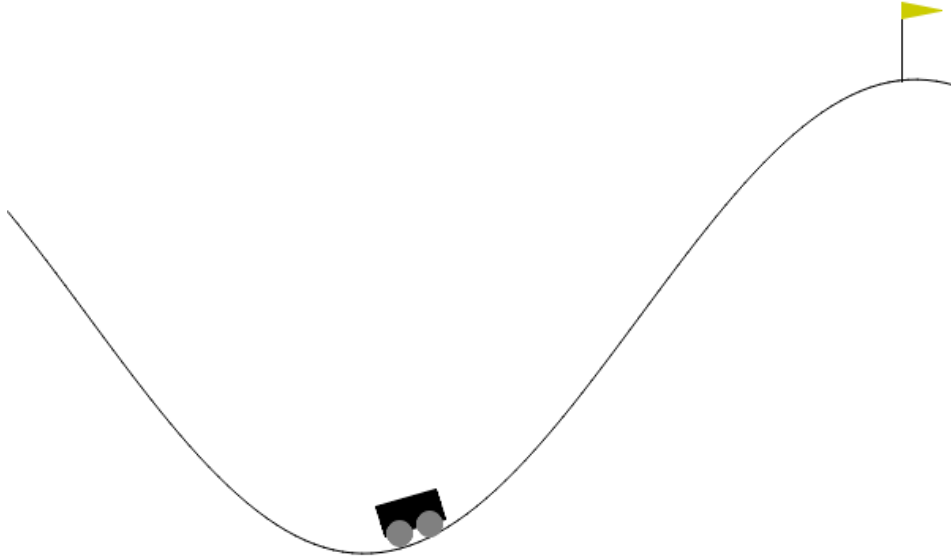
- **By changing BINS and Shaping reward obtain the BEST POSSIBLE LEARNING (FAST-SHORT)**
- Use ChatGPT, internet, your ideas, whatever
- Don't modify the code or the method. Just focus on shaped rewards

https://gymnasium.farama.org/environments/classic_control/mountain_car/

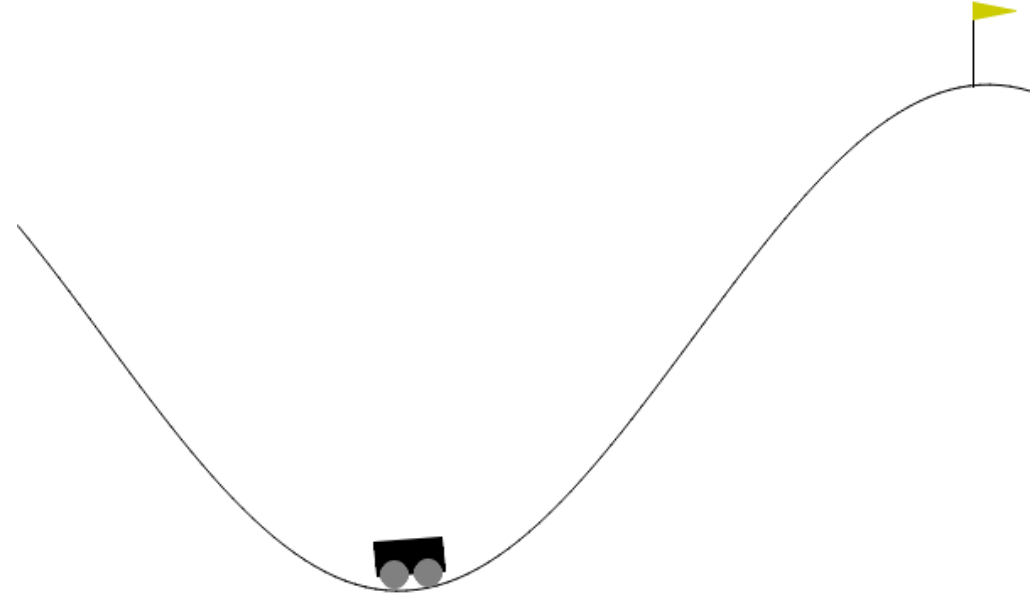
https://github.com/castorgit/RL_course/blob/main/021_Q_learning_MOUNTAIN_CAR.ipynb

Reward Shaping

Difference prioritizing velocity or distance



Reward Shaping – Velocity ++

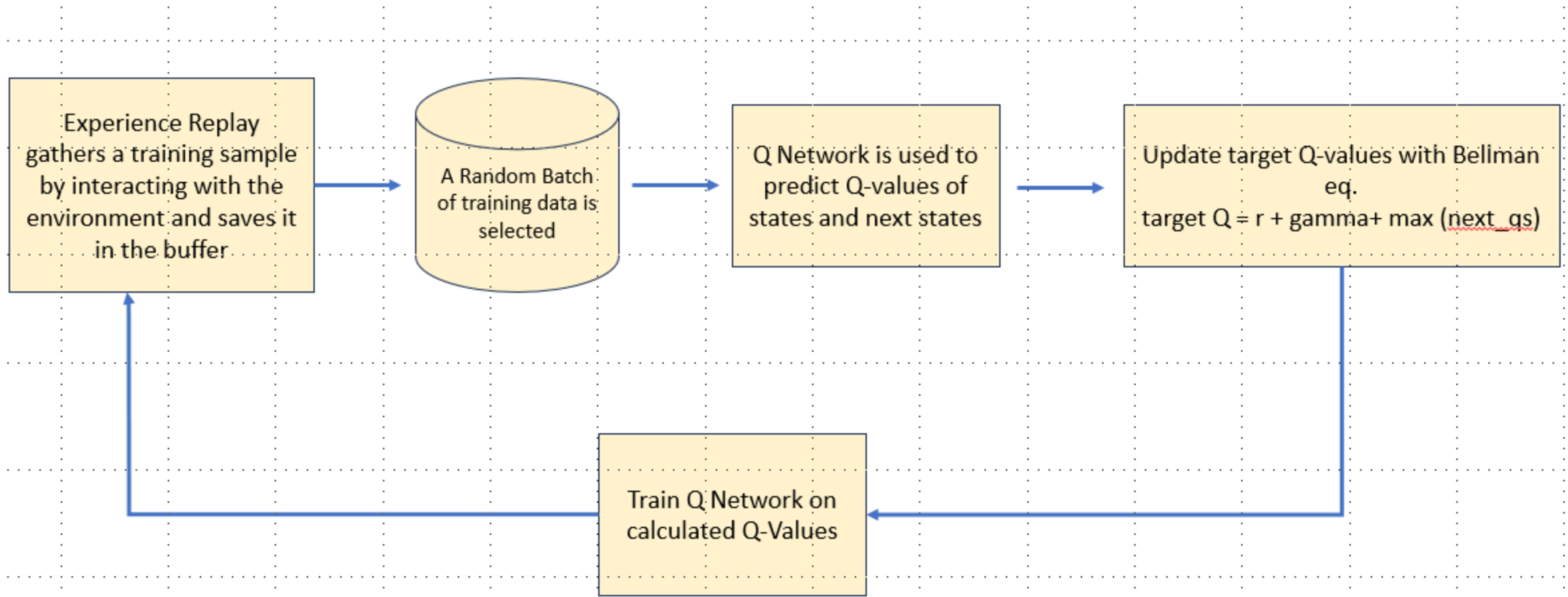


Reward Shaping – distance ++

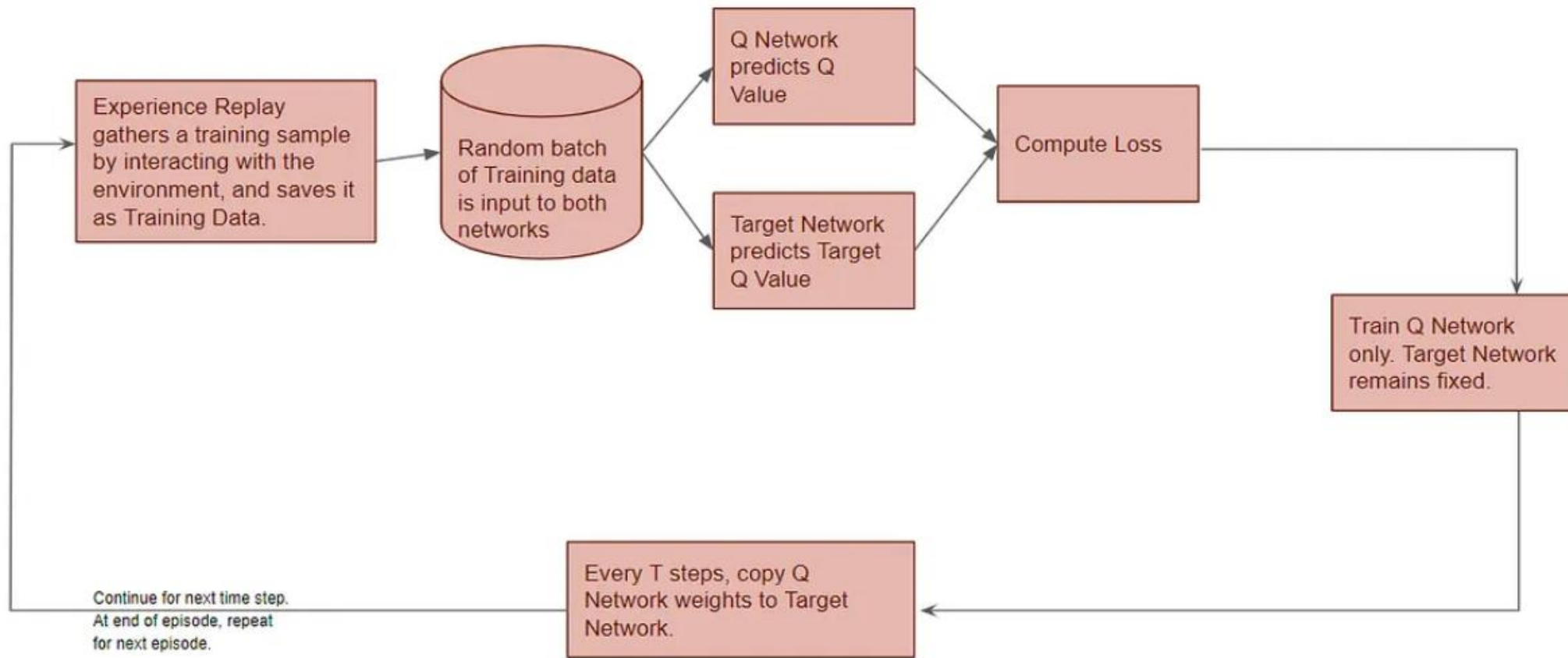
DQN

DQN

General Workflow



Using 2 networks instead of 1



DQN

Hyperparameters

```
1  # Parameters to fine tune
2  # Try your own parameters
3  # Remember epsilon and gamma are very important
4
5
6  MAX_EPISODES = 300
7  ROLLING_WINDOW = 20
8  MEMORY_SIZE = 2000
9  MAX_STEPS = 500
10
11  gamma = 0.99                # discount rate
12  epsilon = 1.0               # exploration rate
13  epsilon_min = 0.01
14  epsilon_decay = 0.99
15  learning_rate = 0.001
16  batch_size = 64
17  solved_threshold = 195
18
19  verb = 0                    # to see traces (verbosity)
```

DQN

Neural Network

Neural Network

```
1  ###
2  # [CREATE YOUR NEURAL NETWORK TRY 16/32 or 24/24]
3  ###
4
5  def build_model(state_size, action_size):
6      inputs = Input(shape=(state_size,), name="state_input")
7      x = Dense(1, activation='relu', name="dense_1")(inputs)
8      x = Dense(1, activation='relu', name="dense_2")(x)
9      outputs = Dense(action_size, activation='linear', name="output_layer")(x)
10
11     model = Model(inputs=inputs, outputs=outputs, name="Q_Network")
12     model.compile(loss='mse', optimizer=Adam(learning_rate=learning_rate))
13     return model
```

Action greedy selection

```
def select_action_greedy(state, DQN):  
    """  
    Selects the if agent takes a random action (explore) or a predicted action (exploit)  
    """  
    if np.random.rand() <= epsilon:  
        return random.randrange(action_size)  
    act_values = DQN.predict(state, verbose=verb)  
    return np.argmax(act_values[0]) # returns action selected with greedy strategy
```

Sample experiences

```
# Sample experiences from the replay buffer
def sample_experiences(batch_size):
    """
    Samples a batch_size of experiences from the Replay buffer.
    You MUST transform the data into numpy arrays as this accelerates the response time sensibly
    """
    indices = np.random.choice(len(replay_buffer), batch_size, replace=False)
    batch = [replay_buffer[i] for i in indices]
    states, actions, rewards, next_states, dones = zip(*batch)
    return (states, actions, rewards, next_states, dones)
###
# YOU MUST VECTIORIZE THE RETURN. TRANSFORM THE OUTPUTS IN np arrays otherwise it will be very slow
###
```

DQN

Experience Replay

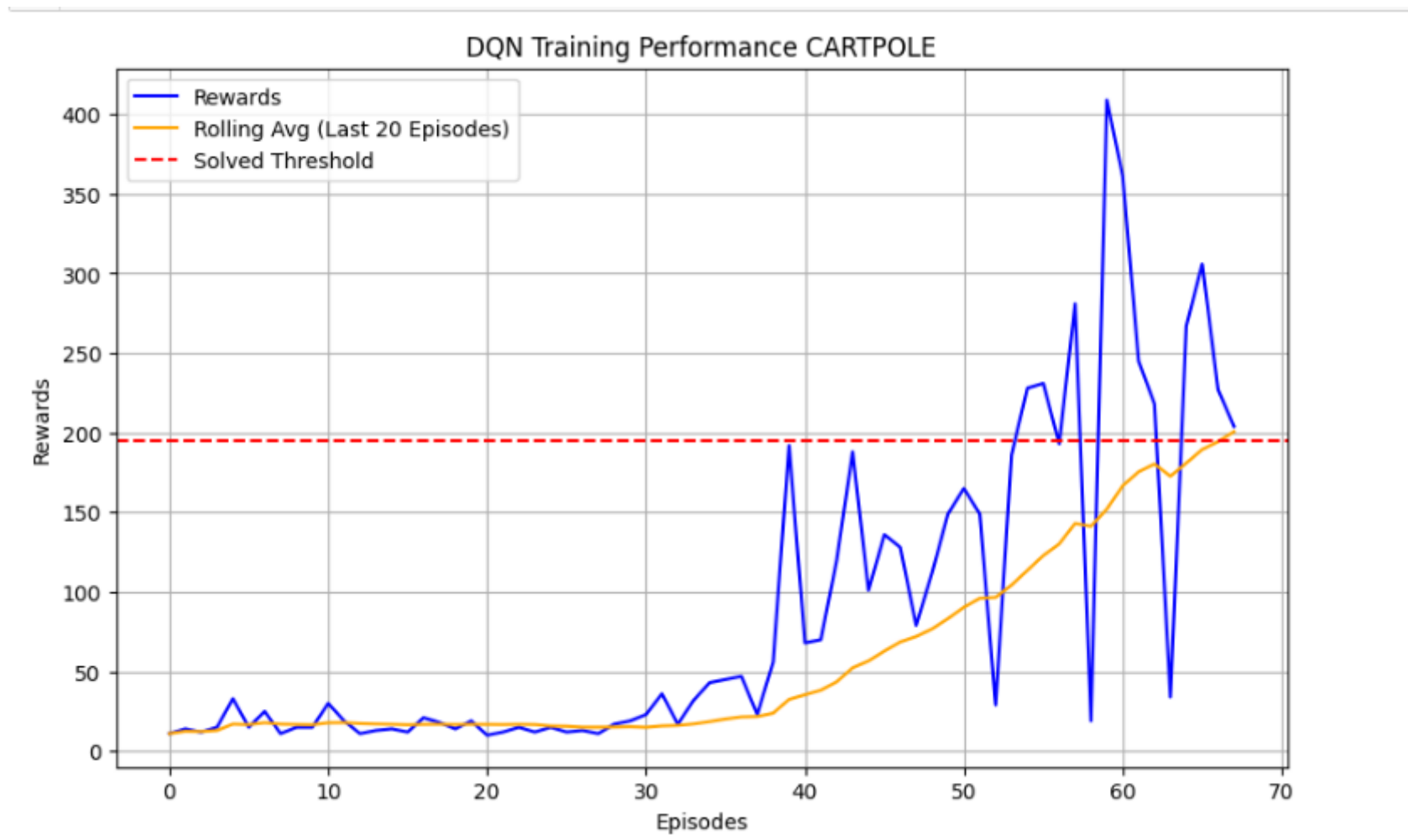
```
33 def experience_replay(batch_size, model, epsilon):
34     """
35     The critical function in the whole program
36     1. gets a minibatch from replay buffer
37     2. Predicts target_qs from states (we need the full batch predicted but we'll avoid dones)
38     3. predicts next_qs from next_states
39     4. Bellman equation on next_qs obtains new target_qs  $Q_{\text{target}}(s_t, a_t) = r_t + \gamma \cdot \max_{a'} Q(s_t, a')$ 
40
41     6. Train DQN input states, output target_qs
42
43     """
44     if len(replay_buffer) < batch_size:
45         return
46
47     states, actions, rewards, next_states, dones = sample_experiences(batch_size)
48
49     # Predict Q-values for current and next states using vectorized operations
50
51     #####
52     # [ Predict target_qs with model predict states, predict next_qs with model predict next states]
53     #####
54     target_qs =
55     next_qs =
56
57     # Update target Q-values using standard DQN logic
58     target_qs[np.arange(batch_size), actions] = rewards + gamma * np.max(next_qs, axis=1) * (1 - dones)
59
60     # Train the model on the Q-values
61     model.fit(states, target_qs, epochs=1, verbose=0)
62
```

DQN

Main Loop

In []:

```
1 rewards_per_episode= []
2
3 done = False
4 rolling_avg = 0
5 rolling_avg_rewards = []
6
7 start_time = time.time()
8
9 for e in range(MAX_EPISODES):           # Should be While True, however we limit number of eps
10     state, _ = env.reset()
11     state = np.reshape(state, [1, state_size])
12     total_reward = 0
13
14     for step in range(MAX_STEPS):
15
16         action = select_action_greedy(state, DQN)
17         next_state, reward, done, truncated, _ = env.step(action)
18
19         next_state = np.reshape(next_state, [1, state_size])
20         store(state, action, reward, next_state, done)
21         state = next_state
22         total_reward = total_reward + reward
23         if done:
24             break
25
26         if len(replay_buffer) > batch_size:
27             experience_replay(batch_size, DQN, epsilon)
28
29     epsilon = max(epsilon_min, epsilon * epsilon_decay)        # decay epsilon
30
31     rewards_per_episode.append(total_reward)
32     rolling_avg = np.mean(rewards_per_episode[-ROLLING_WINDOW:]) # append rewards
33     rolling_avg_rewards.append(rolling_avg)
34     print(f"Episode: {e+1:3}/{MAX_EPISODES}, Reward: {total_reward:+7.2f}, "
35           f"Epsilon: {epsilon:.2f}, Rolling Avg: {rolling_avg:6.2f}, Steps: {step:3} Terminated: {done} ")
36
37     # Check if environment is solved
38     if rolling_avg >= solved_threshold:
39         print(f"Environment solved in {e+1} episodes!")
40         # model.save("lunarlander_ddqn_model1.keras")
41     # break
42
43
44 end_time = time.time()
45 testing_duration = (end_time - start_time) / 60 # Convert to minutes
46 print(f"Testing completed in {testing_duration:.2f} minutes")
```



END

Session 10

