# Report for ML project

Matini Gabriele, 1934803

July 14, 2024

## 1 Introduction: the environment

The Atari SpaceInvaders environment was used to train and test the RL agents. The action space is a discrete 6 moves actions space, with actions like fire, move left, move right, shoot left, shoot right and noop (no operation). The observation space is constituted by 128 bytes (integers from 0 to 255) that describe the state of the game through the RAM bytes. The environment has a frame skip in the game of 4 frames for every action by default. The reason frames are skipped, is twofold: we first have more learning efficiency, because we have less states to consider (less steps to do), and the changes that a certain action does to the environment are more clear and smoothed out to provide a stabler learning environment. Moreover, a human player would not be able to react at every frame with a different command, so we move closer to a normal reaction time. The setup is deterministic by itself, however some stochasticity is given by the "repeat" special action probability: there is a 0.25 probability that the previous action is repeated. This was introduced by the environment to avoid the agent simply memorizing an optimal sequence without paying attention to the environment observations. If we lose the environment sets a "truncated" boolean to true, and we will have to reset it. The game gives the agent a reward equal to how many points it makes, thus the only way to get a reward is to kill the space invaders. If the agent manages to kill all the enemies, a new batch of faster enemies will spawn. Thus the objective is solely to score as many points as possible. Notice that no negative reward is given for losing lives, of which we have up to 3.

## 2 Q-Table Agent

A highlight of the learning process of the Q-Table was the memory usage of the state: the environment returns a 240 bytes numpy array of unsigned 8bit integers. Thus, the state occupies a lot of space in the dictionary. Furthermore, converting the array to a tuple (hashable type) creates a massive overhead, going from 240bytes to 1064bytes, just to represent a 128bytes RAM state. This is why I opted instead for a bytes array representation of the state occupying only 161 bytes. The way it was learned of these facts was by detecting a memory problem through Google Colab's memory leaks, and using in locale a python memory profiling library to discover that the observations where occupying all the RAM's space. Using sys.getsizeof(), it was then printed the size of the single representations, finding out that the problem was using a tuple representation. In this way we are able to construct Q-Tables of 9 or 10 Gigabytes.

### 2.1 The learning process

First of all, the $\epsilon$ greedy policy was defined: this policy let the agent pick a random action for the state with probability of $\epsilon$, and the action with the maximum reward associated in that state in the Q-Table with a probability of 1 - $\epsilon$:

```python
class EpsGreedyPolicy:

    ...

    def __call__(self, obs, eps, Q) -> int:
        greedy = random.random() > eps
        if greedy:
            if np.all(Q[obs] == Q[obs][0]):
                return np.random.choice(len(Q[obs]))
            else:
                if self.show:
                    print("Chosen action ", np.argmax(Q[obs]))
                return np.argmax(Q[obs])
        else:
            return env.action_space.sample()
```

Notice how, if the state is new but the policy is greedy, all the actions associated to it will have reward 0. We decide to choose a random action in this case, to favour exploration in training and randomness (otherwise we would just pick the first action all the time). Next, we actually define the Q-Learning Agent and its training:

```python
class QTable:

    def __init__(self, eps=1.0, showAction=False):
        self.Q = defaultdict(lambda: np.zeros(env.action_space.n))
        self.eps = eps
        self.policy = EpsGreedyPolicy(show=showAction)

    def qlearn(
        self,
        env,
        alpha0,
        gamma,
```

```
13          max_steps
14      ):
15          rewards = []
16          obs, info = env.reset()
17          obs = bytes(obs)
18          done = False
19          tot_rew = 0
20          for step in tqdm(range(max_steps)):
21
22              if done:
23                  rewards.append(tot_rew)
24                  tot_rew = 0
25                  obs, info = env.reset()
26                  obs = bytes(obs)
27                  done = False
28              self.eps = (max_steps - step) / max_steps
29              action = self.policy(obs, self.eps, self.Q, env)
30              results = env.step(action)
31
32              if len(results) == 5:
33                  obs2, reward, terminated, truncated, info = results
34              else:
35                  obs2, reward, terminated, info = results
36                  truncated = False
37
38              obs2 = bytes(obs2)
39              done = terminated or truncated
40              max_next_q = np.max(self.Q[obs2])
41              tot_rew += reward
42
43              self.Q[obs][action] += alpha0 * (reward + gamma * max_next_q - self.Q[obs][action])
44              self.Q[obs][action] = round(self.Q[obs][action], 5)
45
46              obs = obs2
47          return rewards
```

The Q-Table is a dictionary $< state : rewards array >$, where each index of the rewards array defines the correspondent action. Wanting to control how many total steps we take, it was decided to have the training last a fixed number of steps, and if the game ends sooner, we just reset the environment instead of restarting the episode. For each episode, $\epsilon$ descends in line with the number of steps, and at each iteration, the action taken with the state is updated in reward value according to the Q-Function update rule:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right] \quad (1)$$

Where we maximize the sum of future rewards, where $\gamma$ is a weight telling how much we care about the next state rewards, between 1 and 0 (we use 0.99). $\alpha$ is the learning rate and we set that to 0.1, and r is the reward we got by doing our action. So this is going to be our non-deterministic update function. Every episode is defined from its start to the point at which the game ends.

The way we train the agent is by making it do 25 milion steps in total and plotting, for each episode, how the average collected reward in the last 100 episodes evolves through each episode from

the 99th to the last. This is because the actual plotting of the reward through each episode is not stable at all, and thus not very informative. A plotting of the average can actually give us a better idea. We also try to change gamma and see how the performances are affected. The idea of this plot came from the fact that, testing the final Q-Table N times and getting an average didn't seem to get a stable number for the average reward, so, to get a better idea of the performances, a better approach is to actually plot the evolution of this average.

## 2.2 The testing

We test the agent with a rollout function as follows:

```python
def rollout(
        self,
        env: gym.Env,
        policy,
        gamma: float,
        n_episodes: int,
        max_steps: int,
        render=[],
        takeVideo=False
    ) -> float:
    sum_returns = 0.0
    obs, info = env.reset()
    discounting = 1
    for ep in tqdm(range(n_episodes)):
        done = False
        obs, info = env.reset()
        discounting = 1
        while not done:
            obs = tuple(obs)
            action = self.policy(obs, self.eps, self.Q)
            obs, reward, terminated, truncated, info= env.step(action)
            done = terminated or truncated
            sum_returns += reward * discounting
            if takeVideo:
                frame = env.render()
                render.append(frame)
            discounting *= gamma
    print("Info: ",info)
    return sum_returns / n_episodes
```

The rollout function supports the rendering of the environment through taking singular frames of every state in the game. The function supports by itself more games than one, in order to test N times the training approaches.

## 2.3 The results

The training seems, to work, since the average reward per episode does go up with each iteration, even though slowly. We end up with a 9Gbytes Q-Table that is able to eliminate around 10 or 15 enemies per run.
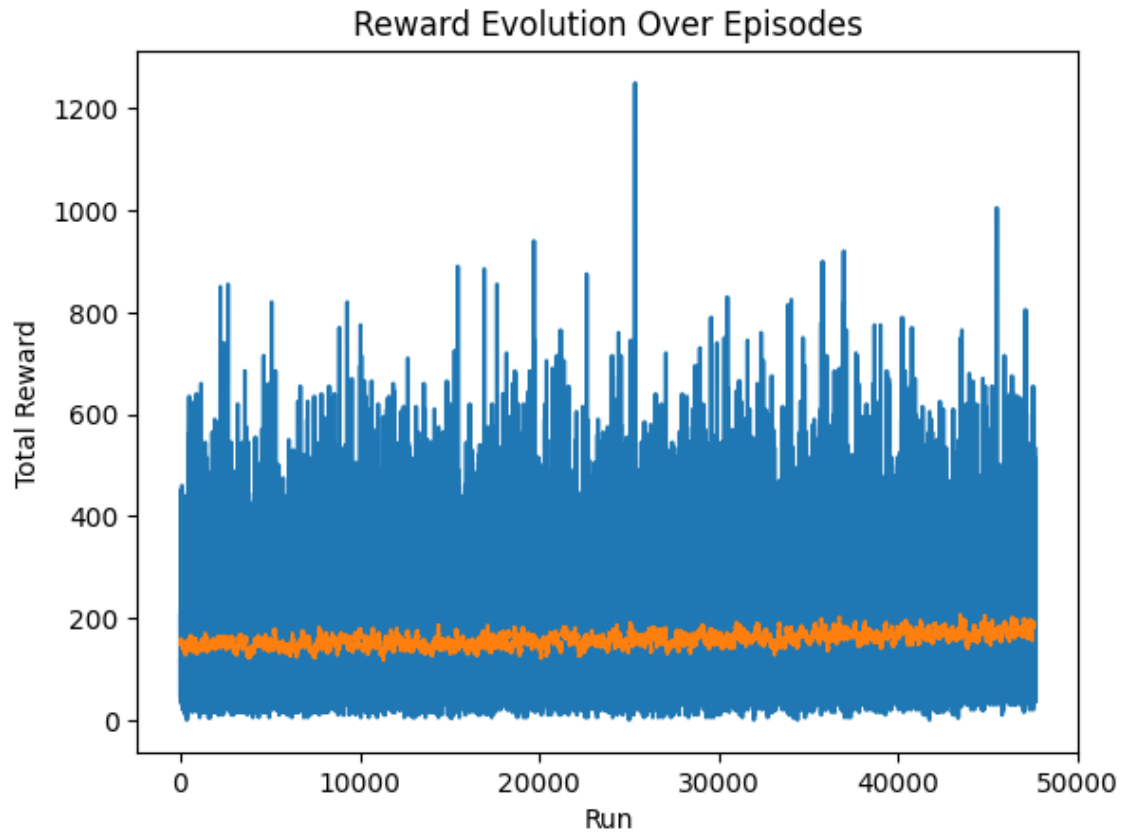
4

Figure 1: In orange: average reward evolution in the last 100 episodes, in blue we have the reward per episode, as one can see it's not that informative because of how unstable it is, that's why it was chosen to plot the average reward. The gamma used was 0.99, alpha 0.1.
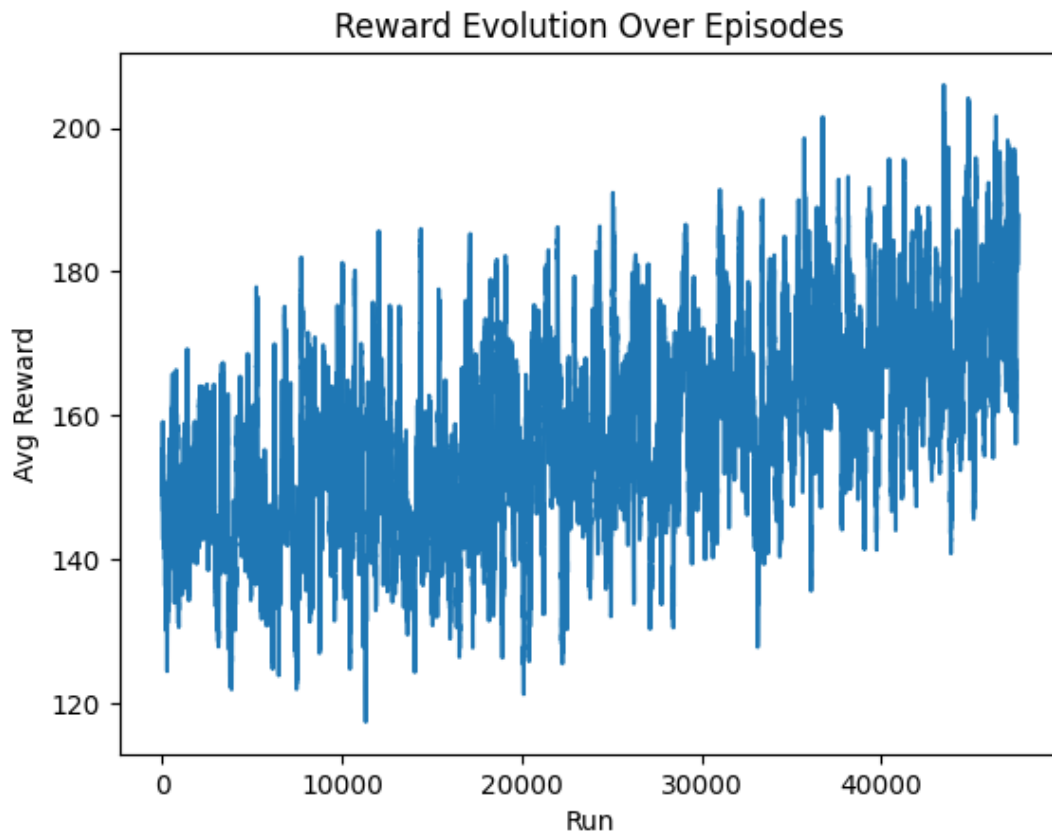
Figure 2: Here we have a zoomed-in plot of the average reward.

Finally, lowering gamma from 0.99 to 0.9 seems to slightly worsen the performances.

# 3 Deep Q-Learning Agent

A more sophisticated approach is the use of a neural network that replaces the Q-Table with respect to learning states and actions' Q-values from the environment. The states will still be updated with the cited Q-Function updated rule, but instead of updating a dictionary, we update the weights of a neural network, by computing a loss gradient between the estimated Q-values and the actual target ones and using the back-propagation algorithm to update the weights. The Q-Table method becomes infeasible in environments with large or continuous state spaces due to the sheer size of the table required, and our state space is very big. In contrast, the neural network can generalize across similar states, making it possible to handle much larger and more complex environments efficiently.

## 3.1 DQN Training and the used NN architecture

We start with a 4 layers network (2 inner layers) of 32 nodes. The agent is trained at each single step, provided the replay buffer is large enough, and each episode lasts until the agent loses at the game. This time we don't have any problems with memory space, since the NN is a smaller solution to approximate the environment, with respect to the Q-Table. ReLu was used as the activation function, a batch of 32 elements was used as input, and Adam as the optimizer with a learning rate of 0.00025, and an MSE loss function. A 0.99 $\gamma$ value was used, as well as a decay of $\epsilon = 0.999$, since the agent needs to train for 1000 episodes. $\epsilon$ must decay relatively slowly, as to expose the agent to as many states and situations as possible. The agent is trained in a twofold manner: first, only one main network is used to do everything and then a second NN is used to estimate the target Q-values for each batch. In the single-NN mode, the main NN is trained by sampling a batch of 32 experiences, computing the Q-values for these states, and then it also has to compute the target Q-values for the next states, meaning that it computes the maximum Q-values for the next states and then we use the Q-Function update rule that makes use of those maximum Q-values to update the target Q-values. We then calculate the loss between the initial estimated Q-values and the target Q-values and back-propagate the gradient results to update the weights of our deep Q-Network. In this way, we can use an NN as a shortened version of the Q-Table. This approach by itself suffers from instability, due to the fact that the main NN has to estimate the target Q-values too, and so whatever noise or error is propagated much more throughout the training. As we will see in the results section, this translates into a more or less consistent overestimation of some action with respect to others. A better approach is to leave the target Q-values estimation to another NN that in our case is the old (last episode) version of the main NN, so that whatever noise, and consequent overestimation may get smoothed out by the target NN, which gets updated once per episode, instead of once per step. As we will see, the target model will play a fundamental role in lowering the noise. To this purpose, it's also of paramount importance to expose the agent to as many different actions as possible, and that's another reason why $\epsilon$ was chosen to decay very slowly.

```python
class DQN(nn.Module):
    def __init__(self, input_shape, num_actions):
        super(DQN, self).__init__()
        self.flatten = nn.Flatten()
        self.fc1 = nn.Linear(input_shape[0] * input_shape[1], 32)
        self.fc2 = nn.Linear(32, 32)
        self.fc3 = nn.Linear(32, 32)
        self.fc4 = nn.Linear(32, num_actions)

    def forward(self, x):
        x = self.flatten(x)
        x = torch.relu(self.fc1(x))
        x = torch.relu(self.fc2(x))
        x = torch.relu(self.fc3(x))
```

```python
15          x = self.fc4(x)
16          return x
17
18
19  class ReplayBuffer:
20      def __init__(self, max_size):
21          self.buffer = deque(maxlen=max_size)
22
23      def add(self, experience):
24          self.buffer.append(experience)
25
26      def sample(self, batch_size):
27          return random.sample(self.buffer, batch_size)
28
29      def size(self):
30          return len(self.buffer)
31
32
33  class DQNAgent:
34      def __init__(self, state_shape, num_actions, buffer_size, batch_size,
35      gamma, alpha, epsilon, epsilon_decay, min_epsilon, use_target_model=False):
36          self.state_shape = state_shape
37          self.num_actions = num_actions
38          self.memory = ReplayBuffer(buffer_size)
39          self.batch_size = batch_size
40          self.gamma = gamma
41          self.alpha = alpha
42          self.use_target_model = use_target_model
43          self.epsilon = epsilon
44          self.epsilon_decay = epsilon_decay
45          self.min_epsilon = min_epsilon
46          self.model = DQN(state_shape, num_actions)
47          self.target_model = DQN(state_shape, num_actions)
48          self.update_target_model()
49          self.optimizer = optim.Adam(self.model.parameters(), lr=0.00025)
50          self.loss_fn = nn.MSELoss()
51
52      def update_target_model(self):
53          self.target_model.load_state_dict(self.model.state_dict())
54
55
56      def act(self, state):
57          if np.random.rand() < self.epsilon:
58              return np.random.randint(self.num_actions)
59          state = torch.FloatTensor(state).unsqueeze(0)
60          with torch.no_grad():
61              q_values = self.model(state)
```

```python
                return torch.argmax(q_values).item()

        def remember(self, state, action, reward, next_state, done):
            self.memory.add((state, action, reward, next_state, done))

        def train(self):

            batch = self.memory.sample(self.batch_size)
            states, actions, rewards, next_states, dones = zip(*batch)

            states = torch.FloatTensor(states)
            next_states = torch.FloatTensor(next_states)
            rewards = torch.FloatTensor(rewards)
            dones = torch.FloatTensor(dones)

            q_values = self.model(states)
            q_values = q_values.gather(1, torch.LongTensor(actions).unsqueeze(1)).squeeze(1)

            if self.use_target_model:
              max_next_q_values = self.target_model(next_states).max(1)[0]
              target_q_values = q_values * (1 - self.alpha) + self.alpha * (rewards + self.gamma *
              max_next_q_values * (1 - dones))
            else:
              max_next_q_values = self.model(next_states).max(1)[0]
              target_q_values = q_values * (1 - self.alpha) + self.alpha * (rewards + self.gamma *
              max_next_q_values * (1 - dones))

            loss = self.loss_fn(q_values, target_q_values.detach())

            self.optimizer.zero_grad()
            loss.backward()
            self.optimizer.step()
            return loss.item()

    def train_dqn(env, agent, episodes):
        rewards = []
        actions = []
        total_reward = 0
        avg_losses = []
        for episode in range(episodes):
            losses = []
            state, info = env.reset()
            state = np.reshape(state, [1, *state.shape])
            done = False
            while not done:
                action = agent.act(state)
                actions.append(action)
```

```
109            next_state, reward, terminated, truncated, info = env.step(action)
110            done = terminated or truncated
111            next_state = np.reshape(next_state, [1, *next_state.shape])
112            agent.remember(state, action, reward, next_state, done)
113            state = next_state
114            total_reward += reward
115            if agent.memory.size() > agent.batch_size:
116                loss = agent.train()
117                losses.append(loss)
118
119        avg_losses.append(np.mean(losses))
120        if agent.use_target_model:
121            agent.update_target_model()
122        print(f"Episode: {episode + 1}/{episodes}, Reward: {total_reward}, Epsilon:
123        {agent.epsilon} MemSize: {agent.memory.size()}")
124        rewards.append(total_reward)
125        total_reward = 0
126        if agent.epsilon > agent.min_epsilon:
127            agent.epsilon *= agent.epsilon_decay
128
129    return (rewards, avg_losses, actions)
```

Other than the neural netowork, the main architecture is also composed by a replay buffer, which is as a deque data structure, to provide a mechanism to store and randomly sample up to 100000 experiences ($< state, action, reward, next - state, done >$ tuples), used to train the NN. Additionally, a second architecture has been contemplated in the code, providing a second NN called "target neural network" or "target model", which can be activated or not during training, which will be updated using a "hard" update technique, which consists of copying the weights of the main NN into the target NN at each new episode. Functions to save and load the agent have also been provided too. Finally, the training function is expected to run for N episodes and collect important data, such as the evolution of the loss and the reward, and how many times the actions were chosen during training, so as to detect any biases towards certain actions. In the end, the collected data will be the average loss per episode, the reward per episode and the average reward evolution calculated on the last 100 episodes (to have a better understanding if the reward is going up or not), and finally, for each action, how many times that action was chosen during training.

```
1   def test_dqn(env, agent, episodes, render):
2       agent.epsilon = 0.0  # Disable exploration
3       total_reward = 0
4       for episode in range(episodes):
5           state, info = env.reset()
6           state = np.reshape(state, [1, *state.shape])
7           done = False
8           while not done:
9               action = agent.act(state)
10              next_state, reward, terminated, truncated, info = env.step(action)
11              done = terminated or truncated
12              next_state = np.reshape(next_state, [1, *next_state.shape])
13              state = next_state
```

```
14              total_reward += reward
15              print("Chosen action is ", action)
16              frame = env.render()
17              render.append(frame)
18              if done:
19                  print(f"Episode: {episode + 1}/{episodes}, Reward: {total_reward}")
20                  break
21      return total_reward / episodes
```

The test function is designed to not only collect the reward, but also the frames to then be able to construct a video of the agent playing in the environment. It was chosen to have the function print the actions chosen by the agent in order to understand better its behavior. $\epsilon$ is also set to 0 in order to fully exploit what the agent has learned.

## 3.2   The results

The agent does always score some points, however from the videos and the test function prints, it's clear that both with and without the use of the target model, the agent seems to think that the "right", "right-fire" and "fire" actions will award it the most points in any state. In particular, with the use of the target model the agent will consistently overestimate "right-fire" and "fire", while without target model the agent will use both "right" and "right-fire", but fundamentally the agent is experiencing the same problem. This translates to the agent moving to the right side of the screen while shooting at the start, and ultimately keep shooting until they get hit by the enemies. The agent might start taking other actions only when some enemies have been defeated, but it will still overwhelmingly choose "right-fire" and "fire". This stems from a documented problem of noise in some Atari environments, which makes the agent overestimate some actions at the expense of every other action. However, the agent does seem to learn something during training though, as documented by these graphs plotted from the data acquired during the training:
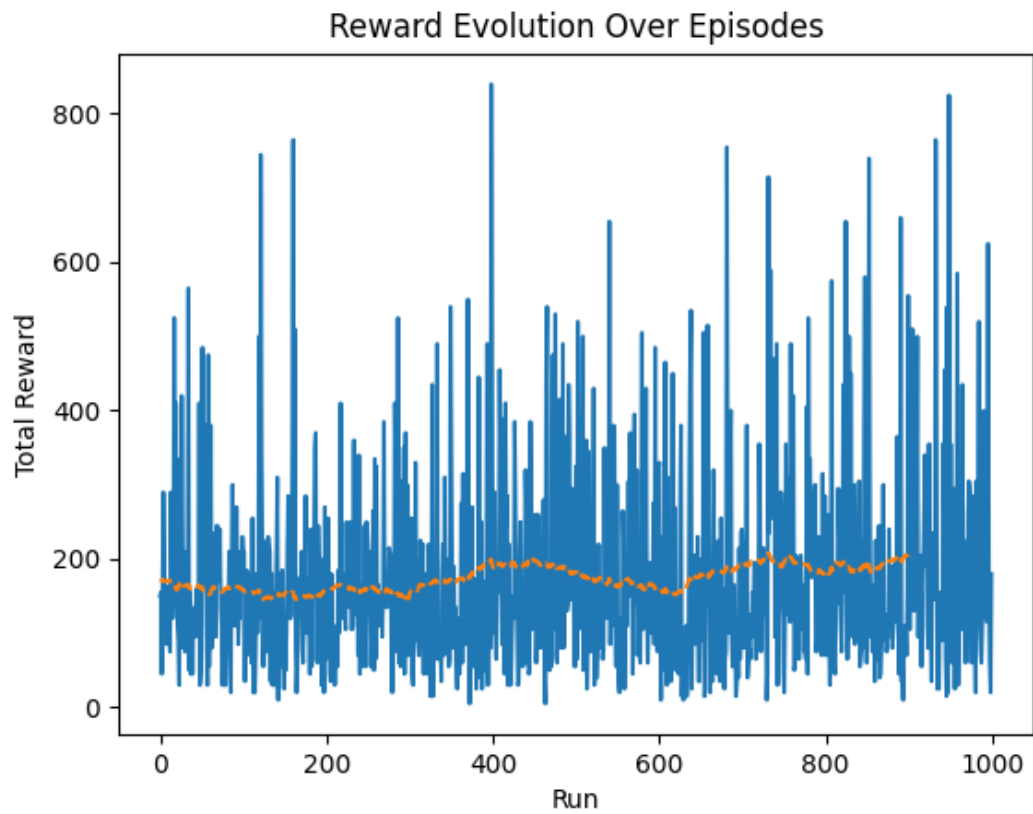
Figure 3: Single-NN DQN reward graph. The orange line represents the average reward of the last 100 runs compared to that run.
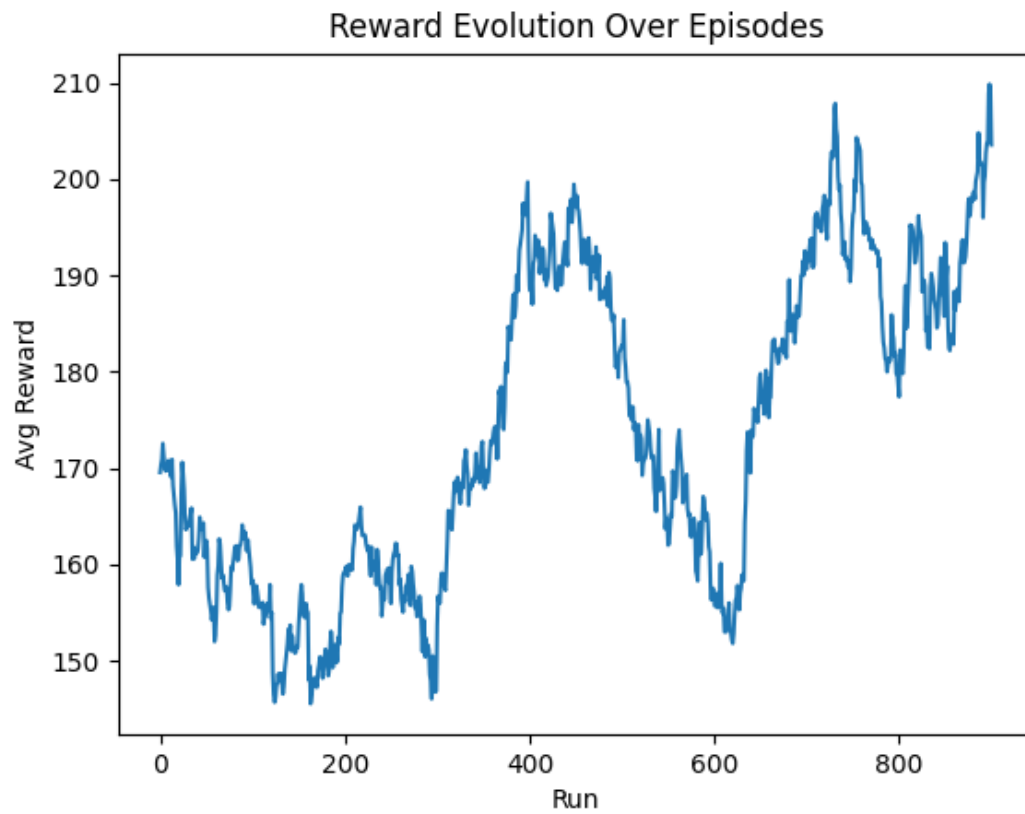
Figure 4: The average single-NN DQN reward plotted alone for a better understanding of its evolution. We can see the reward has almost, but not quite converged yet, hinting at a more unstable learning.
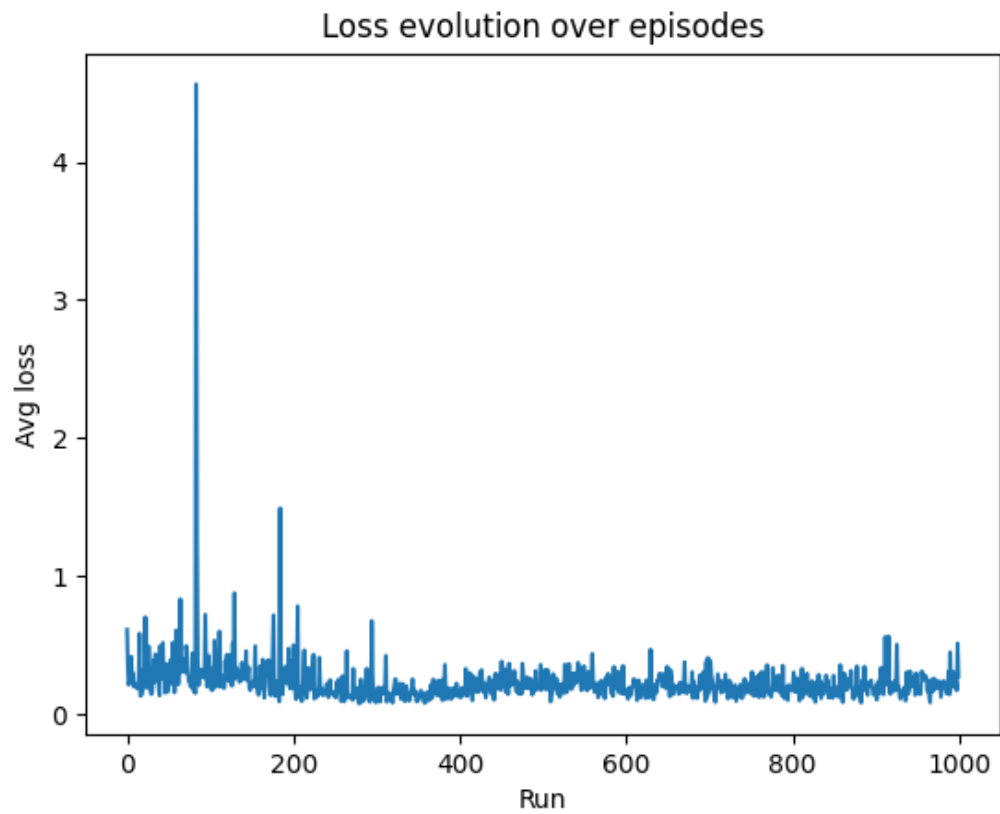
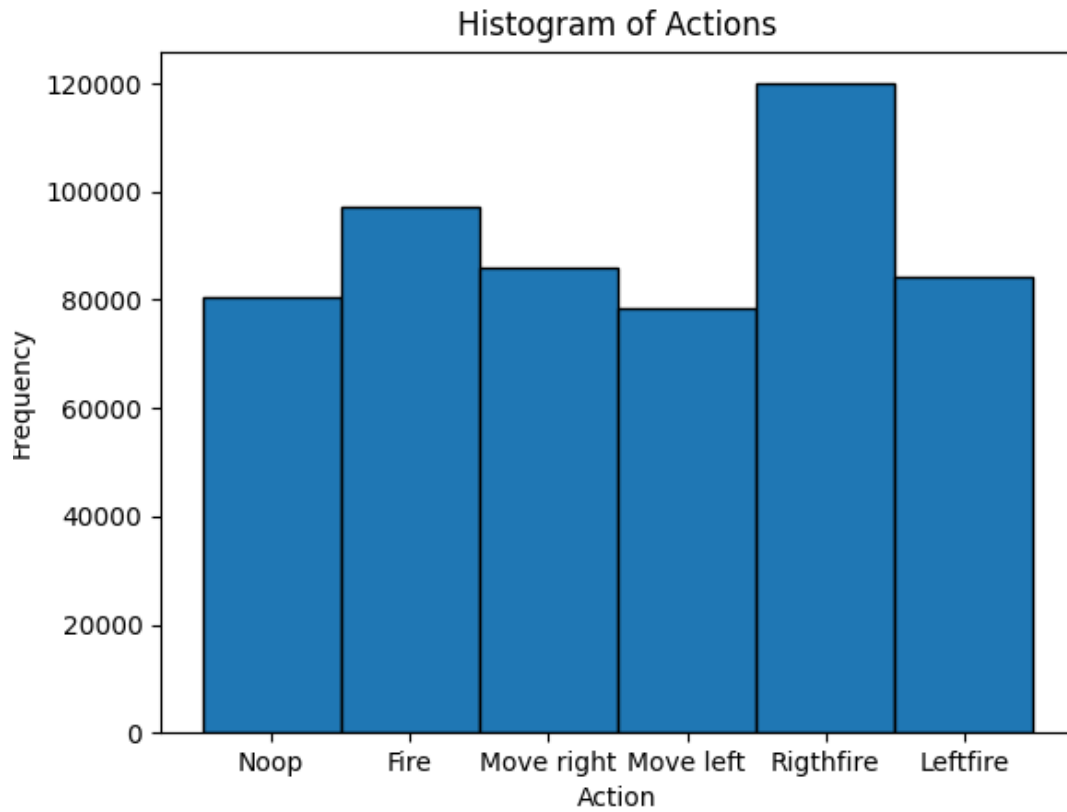Figure 5: Avg loss evolution for single-NN DQN. A peak of more than 4 can be seen, which might signify instability.

Figure 6: Histogram of actions to see if some action where used more than others: the "right-fire" action was used significantly more often, hinting as some over-evaluation of this action during the exploit phase. Indeed the agent will take this action significantly more than others in the testing phases.

The plots of the same data for the Target-Model-DQN is available in the source folder, and has not been published here since the graphs are very similar, and the overestimation problem is still very present. However, as anticipated, the target model will now be used to solve the overestimation problem. In fact, in the actions' histogram of the target-model DQN we can still see clearly that action "right-fire" is the most chosen action by the agent. However, we do have one difference, that being the target-model DQN loss graph:
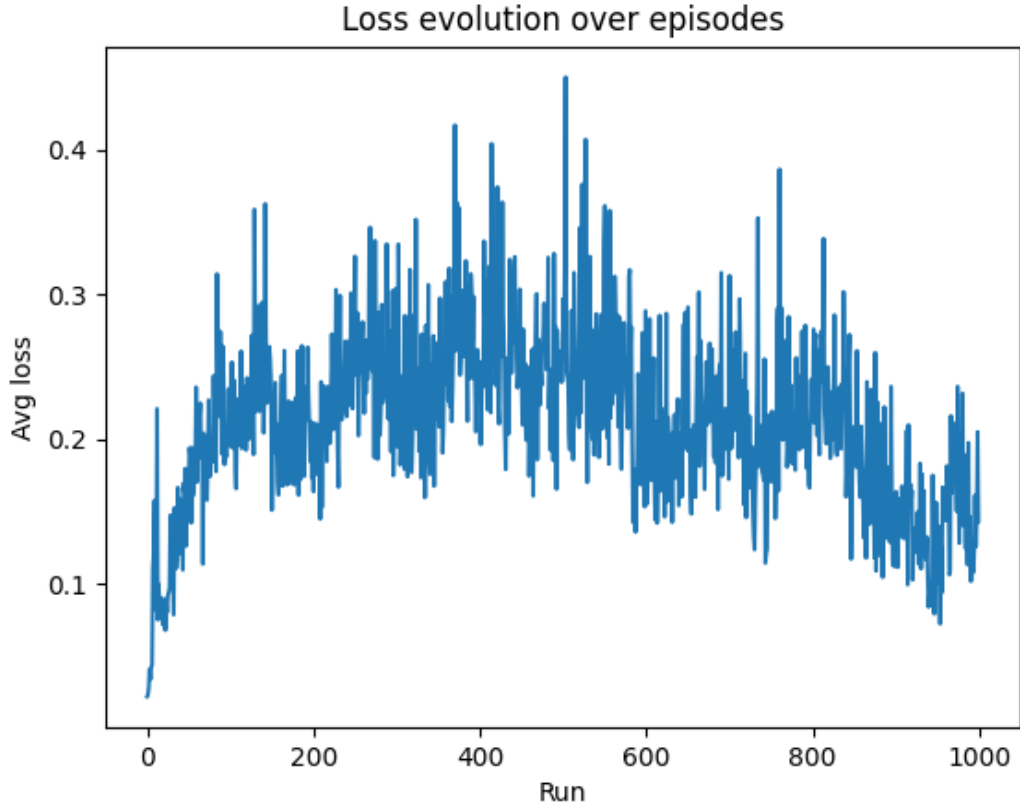


Figure 7: Average loss graph of the Target-Model DQN

As we can see, there are no peaks of loss that surpass 1, which, contrary to the single-NN DQN might be a sign of a stabler learning.

# 4  Double DQN: addressing the problem of overestimation

In order to solve the overestimation of some actions, it is required to smooth out the noise even more. It turns out that in the equation for the target values, the max function is the problem: it consistently overestimates the values of some actions. In traditional DQN, the target values for training the Q-network are computed using the maximum Q-value of the next state. This is achieved by using the max function, which selects the action with the highest estimated Q-value for the next state. However, this approach tends to overestimate the Q-values of some actions. This overestimation bias occurs because the same Q-network is used both to select and evaluate the action, which can lead to an overemphasis on certain actions, particularly in noisy and stochastic environments. This means that noise gets amplified too and in our case this was a major cause of action overestimation, and it might take a long time for the agent to figure out better policies. Double DQN relies on decoupling the action selection from the action evaluation, so that these steps are performed by two different neural networks. In particular, the main NN will select the action and the target NN will evaluate it. This, in practical terms, means using the argmax function when the main NN chooses the action, and having the target network evaluate said action for the next states:

**Traditional DQN Update Rule**

$$y_t = (1 - \alpha)Q(s_t, a_t) + \alpha(r_t + \gamma \max_{a'} Q(s_{t+1}, a'; \theta)) \tag{2}$$

**Double DQN Update Rule**

$$y_t = (1 - \alpha)Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, \arg\max_{a'} Q(s_{t+1}, a'; \theta); \theta^-)) \tag{3}$$

In the first equation, the target value $y_t$ is computed normally, by utilizing the Q-value for state $s_t$ and chosen action $a_t$, the reward $r_t$, and the max Q-value for the next state's ($s_{t+1}$) action to update our estimation. $\theta$ represents the parameters of the main NN. In the second equation, we input into the target model with parameters $\theta^-$ the output action from the main NN (that is why we use the *argmax* function) estimated for the next states, and thus use the target NN to estimate the Q-value for such action.

**Notice: the Double DQN approach is not the same as just using a target model. In fact, while the target model would just compute the maximum Q-value for the next action using the traditional update rule, the Double DQN approach purposefully decouples the action selection from the action evaluation for the next state. In contrast, the max function implicitly makes the same NN perform these two steps, may it be the main NN or the target model when the target model is used.**

```python
class DDQNAgent:

    ...same code of normal DQN...

    def train(self):

        batch = self.memory.sample(self.batch_size)
        states, actions, rewards, next_states, dones = zip(*batch)

        states = torch.FloatTensor(states)
        next_states = torch.FloatTensor(next_states)
        rewards = torch.FloatTensor(rewards)
        dones = torch.FloatTensor(dones)

```

```
15    q_values = self.model(states)
16    q_values = q_values.gather(1, torch.LongTensor(actions).unsqueeze(1)).squeeze(1)
17
18    # Double DQN Implementation
19    #Decoupling of action selection and action evaluation: the main network now will estimate
20    #the best index and feed it to the target model for evaluation
21    next_action_indices = self.model(next_states).max(1)[1] # Get the indices of the actions
22    #with the max Q-value from the primary model, max(*)[1] gives index of max value
23    max_next_q_values = self.target_model(next_states).gather(1, next_action_indices.unsqueeze(1)
24    .squeeze(1)
25    target_q_values = (1 - self.alpha ) * q_values + self.alpha * (rewards + self.gamma *
26    max_next_q_values * (1 - dones))
27
28    loss = self.loss_fn(q_values, target_q_values.detach())
29
30    self.optimizer.zero_grad()
31    loss.backward()
32    self.optimizer.step()
33    return loss.item()
```

The only real modification is in the training function, specifically in the target Q-values calculation.

## 4.1   The results

The results are surprising, as now the DoubleDQN agent actually chooses a plethora of different actions each time. We can see also the average reward converging more clearly to around 190 points, indicating a stabler learning technique. This is confirmed by the loss as well. The videos of the agent playing the game are available in the source code.
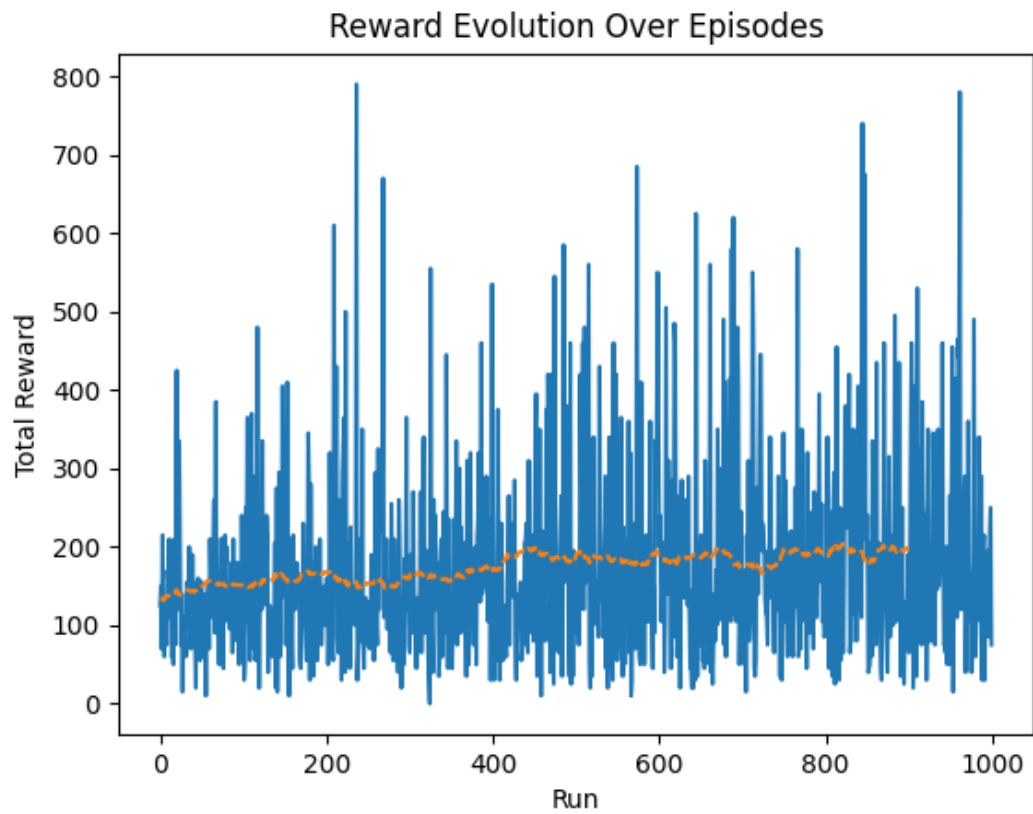
Figure 8: DoubleDQN agent reward evolution. We can see clearly more stability in the average reward evolution.
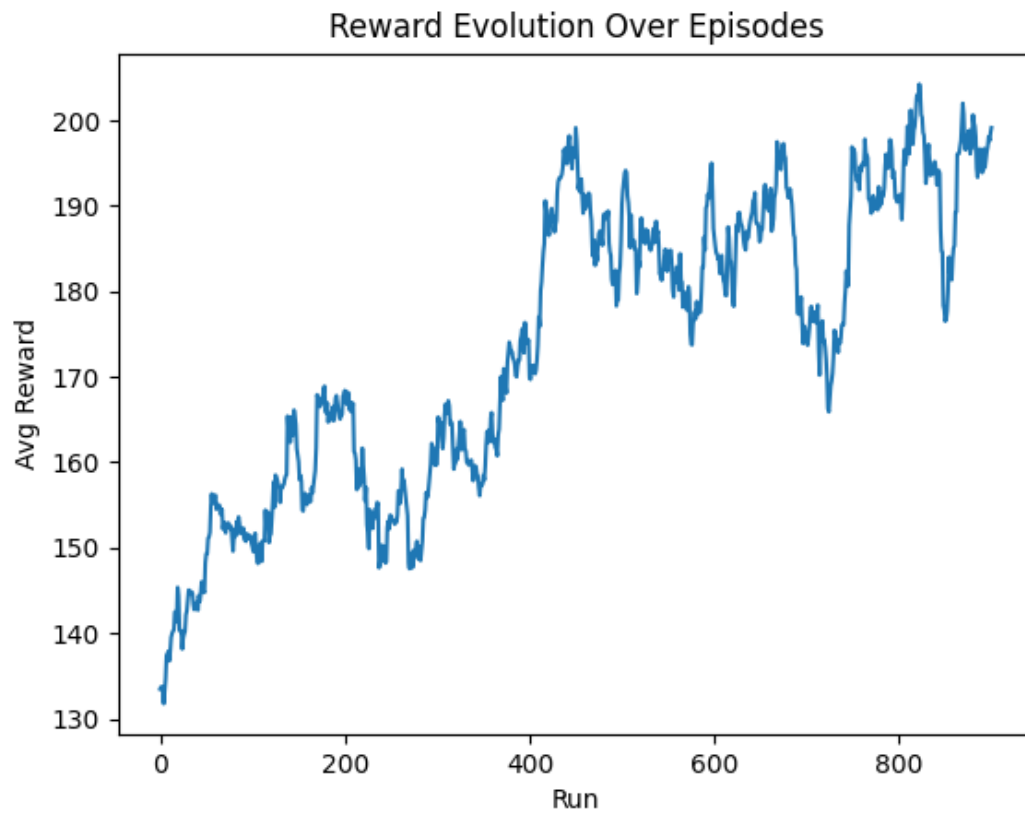
Figure 9: A zoom-in on the average reward evolution reveals a faster and more stable convergence with respect to the single-NN DQN.
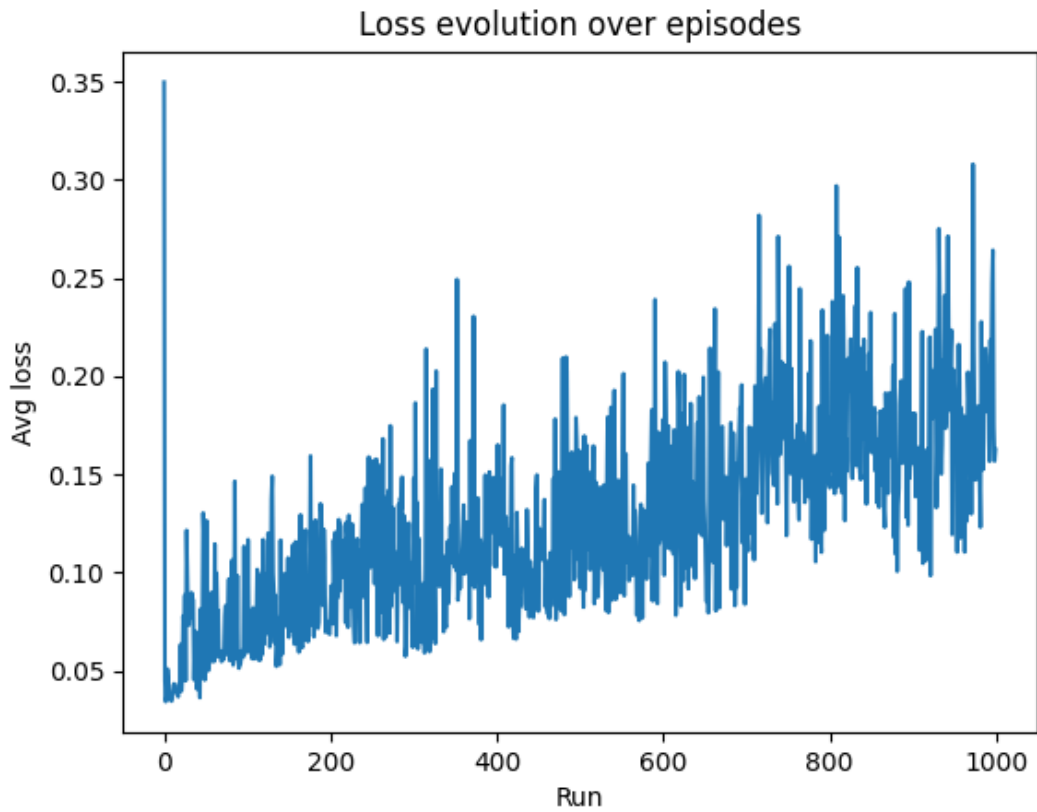
Figure 10: The loss is also stabler, as there are no peaks that surpass 1, in contrast to the single-NN DQN. Also, again the loss sits at around 0.20 at the end of the training, just like the Target-Model DQN did.
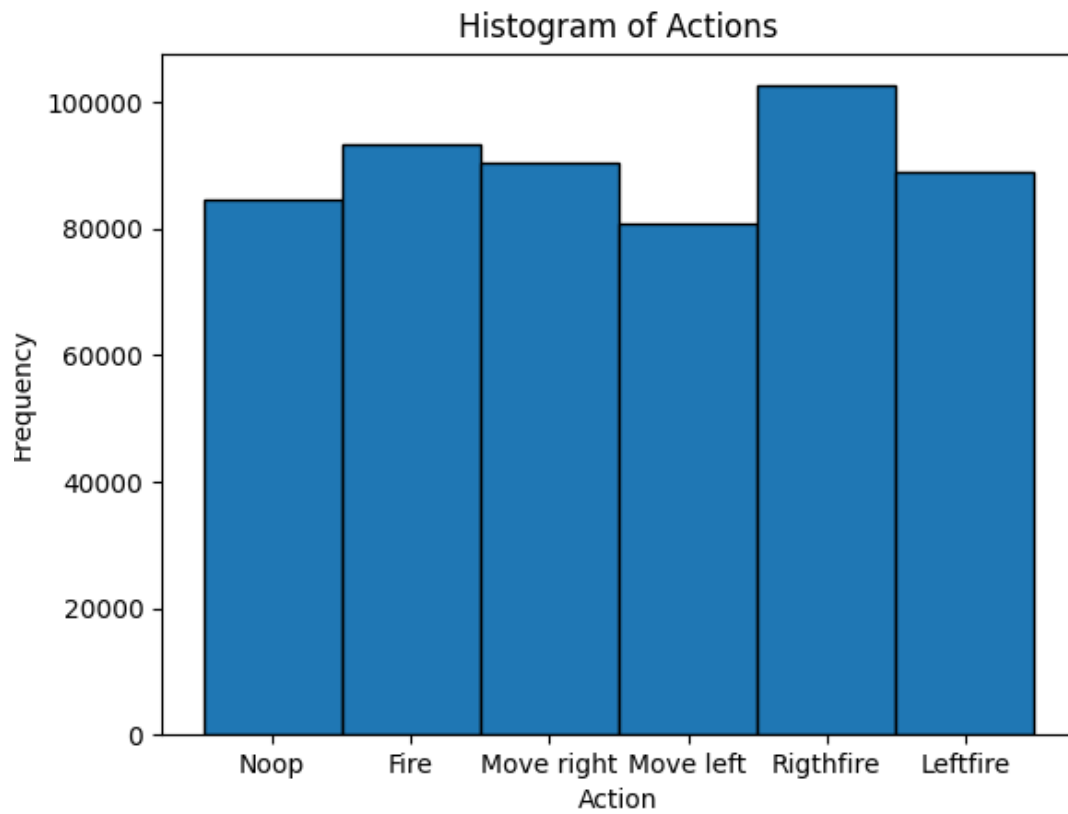
Figure 11: Even the histogram of actions shows a better trend with respect to the agent not choosing an action much more than all the other actions during training.

## 4.2 Q-Table vs Traditional DQN vs Double DQN

It's clear from the videos that the Double DQN is a superior alternative, thanks to it solving the overestimation problem: the agent can now be seen actively pursuing the purple ship that awards more points than the other enemies, dodging enemy shots and exploiting the defenses to avoid being hit. In double DQN the agent also will sometimes try to remain on the left of the map, waiting for the enemies to be right on it so that it can eliminate the first column almost immediately, and at the same time making small movements to the right (under the shield) and back to the left when an enemy shot is approaching in order to dodge said shot. Another tactic adopted by the agent in both Traditional and Double DQN techniques, is that the agent will shoot a hole through its defenses and exploit that small gap to keep shooting the enemies whilst remaining in cover at the same time. Also, turning on the mode that reveals the agent's shots will reveal that the agent will start shooting more often if a column of enemies is directly onto it, which is understandable since in that way it will take out many enemies in little time. In any case, the overall performance of the agent does not change much: the agent still scores about 200 points before losing the game, so this is not much different from when the agent overestimated the "right-fire" action. However, comparing the results with the Q-Table it's apparent that the DQN is a much better alternative, since whilst the Q-Table reaches around 170 as average reward with a lot more training than the DQN, the DQN is able to reach 190 or 200 as average reward in comparatively little time: the DQN employs circa 45 minutes of training vs the 6 hours of training needed by the Q-Table. Also, the Q-Table is expensive to train and store memory-wise, while the DQN is a much more agile approach.