

```
1 !pip install gym[all]
2 !pip install -U gym[atari,accept-rom-license]
3 !AutoROM --accept-license
```

```
1 from google.colab import drive
2 drive.mount('/content/gdrive')
```

Drive already mounted at /content/gdrive; to attempt to forcibly remount, call drive.mount("/content/gdrive", force_remount=True).

```
1 def save_weights_to_drive(model,DRIVE_PATH):
2     %cp $model $DRIVE_PATH
3
4 # save_weights_to_drive("/content/t.txt", DRIVE_PATH)
```

```
1 import gym
2 import cv2
3
4 import time
5 import json
6 import random
7 import numpy as np
8
9 import torch
10 import torch.nn as nn
11 import torch.optim as optim
12 import torch.nn.functional as F
13
14 from collections import deque
```

▼ hyper parameters

```
1 ENVIRONMENT = "PongDeterministic-v4"
2
3 DEVICE = torch.device("cuda" if torch.cuda.is_available() else "cpu")
4
5 SAVE_MODELS = True # Save models to file so you can test later
6 MODEL_PATH = "/content/pong-cnn-" # Models path for saving or loading
7 DRIVE_PATH= "/content/gdrive/MyDrive/Pong/"
8 SAVE_MODEL_INTERVAL = 10 # Save models at every X epoch
9 TRAIN_MODEL = True # Train model while playing (Make it False when testing a model)
10 #/content/pong-cnn-580.pkl
11 LOAD_MODEL_FROM_FILE = True # Load model from file
12 LOAD_FILE_EPISODE = 740 # Load Xth episode from file
13
14 BATCH_SIZE = 64 # Minibatch size that select randomly from mem for train nets
15 MAX_EPISODE = 100000 # Max episode
16 MAX_STEP = 100000 # Max step size for one episode
17
18 MAX_MEMORY_LEN = 50000 # Max memory len
19 MIN_MEMORY_LEN = 40000 # Min memory len before start train
20
21 GAMMA = 0.97 # Discount rate
22 ALPHA = 0.00025 # Learning rate
23 EPSILON_DECAY = 0.99 # Epsilon decay rate by step
24
25 RENDER_GAME_WINDOW = False # Opens a new window to render the game (Won't work on colab default)
```

the Double DQN algrtm that we are implementing <https://arxiv.org/abs/1511.06581>

Algorithm 1: Double DQN Algorithm.

```

input :  $\mathcal{D}$  – empty replay buffer;  $\theta$  – initial network parameters,  $\theta^-$  – copy of  $\theta$ 
input :  $N_r$  – replay buffer maximum size;  $N_b$  – training batch size;  $N^-$  – target network replacement freq.
for episode  $e \in \{1, 2, \dots, M\}$  do
    Initialize frame sequence  $\mathbf{x} \leftarrow ()$ 
    for  $t \in \{0, 1, \dots\}$  do
        Set state  $s \leftarrow \mathbf{x}$ , sample action  $a \sim \pi_{\theta}$ 
        Sample next frame  $x^t$  from environment  $\mathcal{E}$  given  $(s, a)$  and receive reward  $r$ , and append  $x^t$  to  $\mathbf{x}$ 
        if  $|\mathbf{x}| > N_r$  then delete oldest frame  $x_{t_{min}}$  from  $\mathbf{x}$  end
        Set  $s' \leftarrow \mathbf{x}$ , and add transition tuple  $(s, a, r, s')$  to  $\mathcal{D}$ ,
            replacing the oldest tuple if  $|\mathcal{D}| \geq N_r$ 
        Sample a minibatch of  $N_b$  tuples  $(s, a, r, s') \sim \text{Unif}(\mathcal{D})$ 
        Construct target values, one for each of the  $N_b$  tuples:
        Define  $a^{\max}(s'; \theta) = \arg \max_{a'} Q(s', a'; \theta)$ 
         $y_j = \begin{cases} r & \text{if } s' \text{ is terminal} \\ r + \gamma Q(s', a^{\max}(s'; \theta); \theta^-) & \text{otherwise.} \end{cases}$ 
        Do a gradient descent step with loss  $\|y_j - Q(s, a; \theta)\|^2$ 
        Replace target parameters  $\theta^- \leftarrow \theta$  every  $N^-$  steps
    end
end

```

CNN architecture

```

1 #CNN this will be the structure for both the online and target model.
2 class DuelCNN(nn.Module):
3     def __init__(self, h, w, output_size):
4         super(DuelCNN, self).__init__()
5         self.conv1 = nn.Conv2d(in_channels=4, out_channels=32, kernel_size=8, stride=4)
6         self.bn1 = nn.BatchNorm2d(32)
7         convw, convh = self.conv2d_size_calc(w, h, kernel_size=8, stride=4)
8         self.conv2 = nn.Conv2d(in_channels=32, out_channels=64, kernel_size=4, stride=2)
9         self.bn2 = nn.BatchNorm2d(64)
10        convw, convh = self.conv2d_size_calc(convw, convh, kernel_size=4, stride=2)
11        self.conv3 = nn.Conv2d(in_channels=64, out_channels=64, kernel_size=3, stride=1)
12        self.bn3 = nn.BatchNorm2d(64)
13        convw, convh = self.conv2d_size_calc(convw, convh, kernel_size=3, stride=1)
14
15        linear_input_size = convw * convh * 64 # Last conv layer's out sizes
16
17        # Action layer
18        self.Alinear1 = nn.Linear(in_features=linear_input_size, out_features=128)
19        self.Alrelu = nn.LeakyReLU() # Linear 1 activation funct
20        self.Alinear2 = nn.Linear(in_features=128, out_features=output_size)
21
22        # State Value layer
23        self.Vlinear1 = nn.Linear(in_features=linear_input_size, out_features=128)
24        self.Vlrelu = nn.LeakyReLU() # Linear 1 activation funct
25        self.Vlinear2 = nn.Linear(in_features=128, out_features=1) # Only 1 node
26
27        #calculate the Convelotional layers output size
28        def conv2d_size_calc(self, w, h, kernel_size=5, stride=2):
29            next_w = (w - (kernel_size - 1) - 1) // stride + 1
30            next_h = (h - (kernel_size - 1) - 1) // stride + 1
31            return next_w, next_h
32
33        def forward(self, x):
34            x = F.relu(self.bn1(self.conv1(x)))
35            x = F.relu(self.bn2(self.conv2(x)))
36            x = F.relu(self.bn3(self.conv3(x)))
37
38            x = x.view(x.size(0), -1) # Flatten every batch
39
40            Ax = self.Alrelu(self.Alinear1(x))
41            Ax = self.Alinear2(Ax) # No activation on last layer
42
43            Vx = self.Vlrelu(self.Vlinear1(x))
44            Vx = self.Vlinear2(Vx) # No activation on last layer
45
46            q = Vx + (Ax - Ax.mean())
47
48            return q

```

The agent class

```

1 class Agent:
2     def __init__(self, environment):
3         """
4         Hyperparameters definition for Agent
5         """

```

```

6      # State size for Pong environment is (210, 160, 3).
7      self.state_size_h = environment.observation_space.shape[0]
8      self.state_size_w = environment.observation_space.shape[1]
9      self.state_size_c = environment.observation_space.shape[2]
10
11     # actions size for Pong environment is 6
12     self.action_size = environment.action_space.n
13
14     # Image pre process params
15     self.target_h = 80 # Height after process
16     self.target_w = 64 # Widht after process
17
18     self.crop_dim = [20, self.state_size_h, 0, self.state_size_w] # Cut 20 px from top to get rid of the score table
19
20     # Trust rate to our experiences
21     self.gamma = GAMMA # Discount factor for future predictions
22     self.alpha = ALPHA # Learning Rate
23
24     # After many experinces epsilon will be 0.05
25     # So we will do less Explore more Exploit
26     self.epsilon = 1 # Explore or Exploit
27     self.epsilon_decay = EPSILON_DECAY # Adaptive Epsilon Decay Rate
28     self.epsilon_minimum = 0.05 # Minimum for Explore
29
30     # Deque to stor experience replay .
31     self.memory = deque(maxlen=MAX_MEMORY_LEN)
32
33     # initialize the two model for DDQN algorithm online model, target model
34     self.online_model = DuelCNN(h=self.target_h, w=self.target_w, output_size=self.action_size).to(DEVICE)
35     self.target_model = DuelCNN(h=self.target_h, w=self.target_w, output_size=self.action_size).to(DEVICE)
36     self.target_model.load_state_dict(self.online_model.state_dict())
37     #we put target model in evaluation mode because we don't want it to train
38     self.target_model.eval()
39
40     # Adam used as optimizer
41     self.optimizer = optim.Adam(self.online_model.parameters(), lr=self.alpha)
42
43 #Process image crop resize, grayscale and normalize the images
44 def preProcess(self, image):
45     frame = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY) # To grayscale
46     frame = frame[self.crop_dim[0]:self.crop_dim[1], self.crop_dim[2]:self.crop_dim[3]] # Cut 20 px from top
47     frame = cv2.resize(frame, (self.target_w, self.target_h)) # Resize
48     frame = frame.reshape(self.target_w, self.target_h) / 255 # Normalize
49
50     return frame
51
52 #epsilon greedy algorithm to explor and exploit
53 def act(self, state):
54     act_protocol = 'Explore' if random.uniform(0, 1) <= self.epsilon else 'Exploit'
55
56     if act_protocol == 'Explore':
57         action = random.randrange(self.action_size)
58     else:
59         with torch.no_grad():
60             state = torch.tensor(state, dtype=torch.float, device=DEVICE).unsqueeze(0)
61             q_values = self.online_model.forward(state) # (1, action_size)
62             action = torch.argmax(q_values).item() # Returns the indices of the maximum value of all elements
63
64     return action
65
66 #experience replay to train the model
67 def train(self):
68     if len(agent.memory) < MIN_MEMORY_LEN:
69         loss, max_q = [0, 0]
70         return loss, max_q
71     # sample a minibatch from the memory
72     state, action, reward, next_state, done = zip(*random.sample(self.memory, BATCH_SIZE))
73
74     # Concat batches in one array
75     # (np.arr, np.arr) ==> np.BIGarr
76     state = np.concatenate(state)
77     next_state = np.concatenate(next_state)
78
79     # Convert them to tensors
80     state = torch.tensor(state, dtype=torch.float, device=DEVICE)
81     next_state = torch.tensor(next_state, dtype=torch.float, device=DEVICE)
82     action = torch.tensor(action, dtype=torch.long, device=DEVICE)
83     reward = torch.tensor(reward, dtype=torch.float, device=DEVICE)
84     done = torch.tensor(done, dtype=torch.float, device=DEVICE)
85
86     # Make predictions
87     state_q_values = self.online_model(state)
88     next_states_q_values = self.online_model(next_state)

```

```

88     next_states_target_q_values = self.target_model(next_state)
89
90     # Find selected action's q_value
91     selected_q_value = state_q_values.gather(1, action.unsqueeze(1)).squeeze(1)
92     # Get indice of the max value of next_states_q_values
93     # Use that indice to get a q_value from next_states_target_q_values
94     # We use greedy for policy So it called off-policy
95     next_states_target_q_value = next_states_target_q_values.gather(1, next_states_q_values.max(1)[1].unsqueeze(1)).squeeze(1)
96     # Use Bellman function to find expected q value
97     expected_q_value = reward + self.gamma * next_states_target_q_value * (1 - done)
98
99     # Calc loss with expected_q_value and q_value
100    loss = (selected_q_value - expected_q_value.detach()).pow(2).mean()
101
102    self.optimizer.zero_grad()
103    loss.backward()
104    self.optimizer.step()
105
106    return loss, torch.max(state_q_values).item()
107
108    #add new experince to the replay memory
109    def storeResults(self, state, action, reward, nextState, done):
110        self.memory.append([state[None, :], action, reward, nextState[None, :], done])
111
112    #decay epsilon at every step to allow our model to exploit more as it trains
113    def adaptiveEpsilon(self):
114        if self.epsilon > self.epsilon_minimum:
115            self.epsilon *= self.epsilon decav
116
117    1 weightsPath = MODEL_PATH + str(LOAD_FILE_EPISODE) + '.pkl'
118    2 epsilonPath = MODEL_PATH + str(LOAD_FILE_EPISODE) + '.json'

```

```

1 # create the environment
2 environment = gym.make(ENVIRONMENT, render_mode='rgb_array')
3 #create an instance of agent class
4 agent = Agent(environment)
5
6 #load pretrained weights if we already trined
7 if LOAD_MODEL_FROM_FILE:
8     agent.online_model.load_state_dict(torch.load(MODEL_PATH+str(LOAD_FILE_EPISODE)+".pkl"))
9
10    with open(MODEL_PATH+str(LOAD_FILE_EPISODE)+'.json') as outfile:
11        param = json.load(outfile)
12        agent.epsilon = param.get('epsilon')
13
14    startEpisode = LOAD_FILE_EPISODE + 1
15    #start with eps=1 and random weights because we didn't train yet
16    else:
17        startEpisode = 1
18
19    #data structure deque to store the last 100 episode rewards
20    last_100_ep_reward = deque(maxlen=100)
21    total_step = 1
22    for episode in range(startEpisode, MAX_EPISODE):
23
24        startTime = time.time() # store the time
25        # Reset env at the beginning of each episode
26        state = environment.reset()
27        #get the observations from the environment
28        state= environment.render()
29
30        #process the observations to get suitable images that can be fed to the CNNs
31        state = agent.preProcess(state)
32
33
34        #stack observations to create the state "4 consecutive observations makes a state"
35        state = np.stack((state, state, state, state))
36
37        total_max_q_val = 0 # Total max q vals
38        total_reward = 0 # Total reward for each episode
39        total_loss = 0 # Total loss for each episode
40        for step in range(MAX_STEP):
41
42            #if we want to show how is the model playing we can render the opervations
43            #this doesn't work probably yet since we are working with colab and we dont have a monitor to show the observations
44            if RENDER_GAME_WINDOW:
45                environment.render()
46
47            # use epsilon greedy to select an action to be perforemed on the environemt
48            action = agent.act(state)
49            #get the next observation, reward, and done to show of we reached a terminal state.
50            next_state, reward, done,_, info = environment.step(action)

```

```

51
52     #process the new observations to create the next state
53     next_state = agent.preProcess(next_state)
54
55     #add the new observation to the already defined state to create the new state
56     next_state = np.stack((next_state, state[0], state[1], state[2]))
57
58     # Store the transition in memory
59     agent.storeResults(state, action, reward, next_state, done) # Store to mem
60
61     # Update state
62     state = next_state
63
64     if TRAIN_MODEL:
65         # Perform one step of the optimization (on the target network)
66         #using the experience replay algorithm
67         loss, max_q_val = agent.train()
68     else:
69         loss, max_q_val = [0, 0]
70
71     total_loss += loss
72     total_max_q_val += max_q_val
73     total_reward += reward
74     total_step += 1
75     if total_step % 1000 == 0:
76         # decay epsilon each 1000 steps
77         agent.adaptiveEpsilon()
78
79     # we compleated an episode
80     if done:
81         #store the finish time
82         currentTime = time.time()
83         # get the episode duration
84         time_passed = currentTime - startTime
85         # Get current dateTIme as HH:MM:SS
86         current_time_format = time.strftime("%H:%M:%S", time.gmtime())
87         # Create epsilon dict to save model as file
88         epsilonDict = {'epsilon': agent.epsilon}
89
90         # Save model to the over come Ram craches
91         if SAVE_MODELS and episode % SAVE_MODEL_INTERVAL == 0:
92             weightsPath = MODEL_PATH + str(episode) + '.pkl'
93             epsilonPath = MODEL_PATH + str(episode) + '.json'
94
95             torch.save(agent.online_model.state_dict(), weightsPath)
96             with open(epsilonPath, 'w') as outfile:
97                 json.dump(epsilonDict, outfile)
98
99         if TRAIN_MODEL:
100             # Update target model at the end of the episode target_model = online_model
101             agent.target_model.load_state_dict(agent.online_model.state_dict())
102
103         last_100_ep_reward.append(total_reward)
104         avg_max_q_val = total_max_q_val / step
105
106         #create output file to show the results later
107         outStr = "Episode:{0} Time:{0} Reward:{:.2f} Loss:{:.2f} Last_100_Avg_Rew:{:.3f} Avg_Max_Q:{:.3f} Epsilon:{:.2f} Duration:{:0.1f}
108                 episode, current_time_format, total_reward, total_loss, np.mean(last_100_ep_reward), avg_max_q_val, agent.epsilon, time_passed
109             )
110
111         print(outStr)
112
113         #save the output file, model and epsilon value to the drive
114         if SAVE_MODELS:
115             outputPath = MODEL_PATH + "out" + '.txt' # Save outStr to file
116             with open(outputPath, 'a') as outfile:
117                 outfile.write(outStr+"\n")
118             save_weights_to_drive(weightsPath,DRIVE_PATH)
119             save_weights_to_drive(epsilonPath,DRIVE_PATH)
120             save_weights_to_drive(outputPath,DRIVE_PATH)
121
122
123         break
124

```

```

1 # after 3 hours alomest reashed the avwrage reward of 4.4
2 #after 346 episdpdes the average reward is 10.86 and it took 3:30 h to train
3 # after 5:46 trained 545 episodes and got reward of 14.26

```

▼ Results and processing output:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
```

```
1 #process the output file and create dictionary info to store the importabt information
2 info={}
3 with open("/content/pong-cnn-out.txt", "r") as input_file:
4     for line in input_file:
5         a,b,c,d,e,f,g,h,i,j=(item.strip() for item in line.split(' ',9))
6         a,a_ = a.split(":")
7         c,c_ = c.split(":")
8         d,d_ = d.split(":")
9         e,e_ = e.split(":")
10        f,f_ = f.split(":")
11        g,g_ = g.split(":")
12        h,h_ = h.split(":")
13        i,i_ = i.split(":")
14        j,j_ = j.split(":")
15        tempdict= dict(zip((c,d,e,f,g,h,i,j),(c_,d_,e_,f_,g_,h_,i_,j_)))
16        info[a_]= tempdict
17
```

```
1 #show the stored info for each episode
2 episode_='1'
3 info[episode_]

```

```
{'Reward': '-21.00',
 'Loss': '0.00',
 'Last_100_Avg_Rew': '-21.000',
 'Avg_Max_Q': '0.000',
 'Epsilon': '1.00',
 'Duration': '0.98',
 'Step': '883',
 'CStep': '885'}
```

```
1 #some values got missed up because of mulipule RAM craches so solving this problem
2 #stor the info into lists to plot them later and fix the problems
3 Reward=[]
4 Loss=[]
5 Last_100_Avg_Rew=[]
6 Avg_Max_Q=[]
7 Epsilon=[]
8 Duration=[]
9 Step=[]
10 for i in range(len(info)):
11     Reward.append(float(info[str(i+1)]["Reward"]))
12     Loss.append(float(info[str(i+1)]["Loss"]))
13     Last_100_Avg_Rew.append(float(info[str(i+1)]["Last_100_Avg_Rew"]))
14     Avg_Max_Q.append(float(info[str(i+1)]["Avg_Max_Q"]))
15     Epsilon.append(float(info[str(i+1)]["Epsilon"]))
16     Step.append(float(info[str(i+1)]["Step"]))
17     Duration.append(float(info[str(i+1)]["Duration"]))
18
```

```
1 #define window length to calculate moving/running average
2 N=100
```

```
1 reward_moving_averages=[]
2 window_average=0
3 #after 500 the mean average 100 ep got missed up because of ram issues so we have to recalculate them
4 i=0
5 for i in range(500):
6     reward_moving_averages.append(Last_100_Avg_Rew[i])
7
8
9 i=500
10 while i < len(Reward):
11
12     # Store elements from i to i+window_size
13     # in list to get the current window
14     window = Reward[i-N : i ]
15
16     # Calculate the average of current window
17     window_average = round(sum(window) / N, 2)
18
19     # Store the average of current
```

```

20 # window in moving average list
21 reward_moving_averages.append(window_average)
22
23 # Shift window to right by one position
24 i += 1
25
26 #reward_moving_averages
27 print(len(reward_moving_averages))

```

717

```

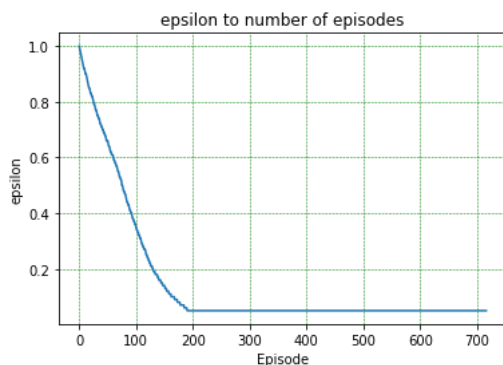
1 #calucate moving average for episodes durations
2 durations_moving_averages=[]
3 window_average=0
4 i=0
5 while i < len(Duration) - N + 1:
6
7     # Store elements from i to i+window_size
8     # in list to get the current window
9     window = Duration[i : i + N]
10
11     # Calculate the average of current window
12     window_average = round(sum(window) / N, 2)
13
14     # Store the average of current
15     # window in moving average list
16     durations_moving_averages.append(window_average)
17
18     # Shift window to right by one position
19     i += 1
20
21 #durations_moving_averages

```

```

1 plt.plot([Epsilon[i] for i in range(len(info))])
2 plt.grid(color = 'green', linestyle = '--', linewidth = 0.5)
3 plt.xlabel('Episode')
4 plt.ylabel('epsilon')
5 plt.title('epsilon to number of episodes')
6 plt.show()

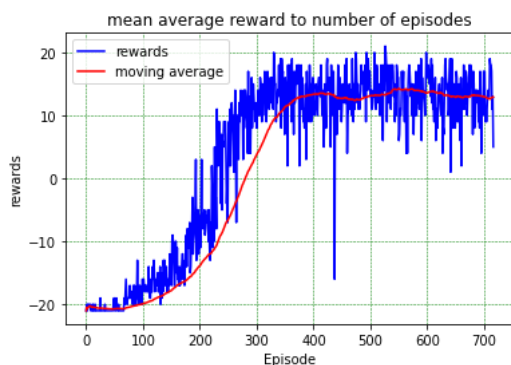
```



```

1 plt.plot([Reward[i] for i in range(len(info))], "b", label= "rewards")
2 plt.plot([reward_moving_averages[i] for i in range(len(info))], "r",linewidth=1.5, label="moving average")
3 plt.grid(color = 'green', linestyle = '--', linewidth = 0.5)
4 plt.xlabel('Episode')
5 plt.ylabel('rewards')
6 plt.title('mean average reward to number of episodes')
7 plt.legend()
8 plt.show()

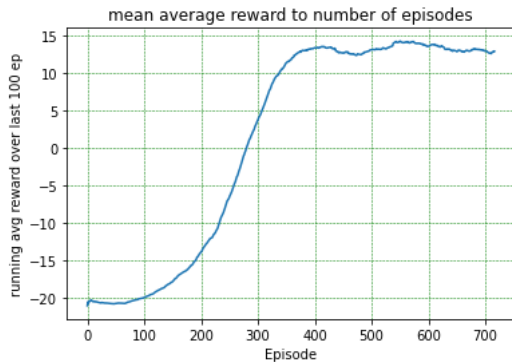
```



```
1 #maximum reward moving average
2 np.max(reward_moving_averages)
```

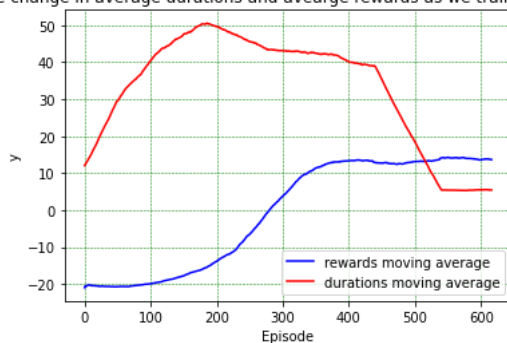
14.27

```
1 plt.plot([reward_moving_averages[i] for i in range(len(info))])
2 plt.grid(color = 'green', linestyle = '--', linewidth = 0.5)
3 plt.xlabel('Episode')
4 plt.ylabel('running avg reward over last 100 ep')
5 plt.title('mean average reward to number of episodes')
6 plt.show()
```

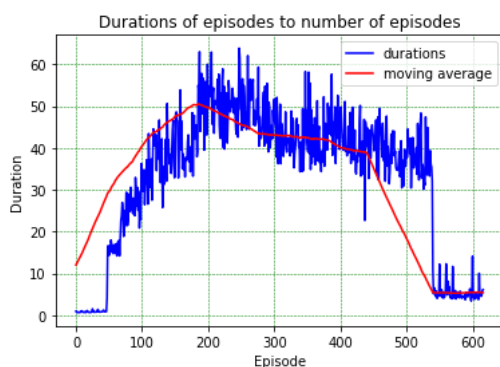


```
1 plt.plot([reward_moving_averages[i] for i in range(len(info)-100)], 'b', label= "rewards moving average")
2 plt.plot([durations_moving_averages[i] for i in range(len(info)-100)], 'r', label= "durations moving average")
3 plt.grid(color = 'green', linestyle = '--', linewidth = 0.5)
4 plt.xlabel('Episode')
5 plt.ylabel('y')
6 plt.title('The change in average durations and avearge rewards as we train the model')
7 plt.legend()
8 plt.show()
```

↳ The change in average durations and avearge rewards as we train the model



```
1 plt.plot([Duration[i] for i in range(len(info)-100)], 'b', label= "durations")
2 plt.plot([durations_moving_averages[i] for i in range(len(info)-100)], 'r', label= "moving average")
3 plt.grid(color = 'green', linestyle = '--', linewidth = 0.5)
4 plt.xlabel('Episode')
5 plt.ylabel('Duration')
6 plt.title('Durations of episodes to number of episodes')
7 plt.legend()
8 plt.show()
```



```
1 #print model surmmery
2 from torchsummary import summary
```



```
3 print("online model and target model architecture")
4 summary(agent.online_model, (4, 64, 80))
```

online model and target model architecture

Layer (type)	Output Shape	Param #
Conv2d-1	[-1, 32, 15, 19]	8,224
BatchNorm2d-2	[-1, 32, 15, 19]	64
Conv2d-3	[-1, 64, 6, 8]	32,832
BatchNorm2d-4	[-1, 64, 6, 8]	128
Conv2d-5	[-1, 64, 4, 6]	36,928
BatchNorm2d-6	[-1, 64, 4, 6]	128
Linear-7	[-1, 128]	196,736
LeakyReLU-8	[-1, 128]	0
Linear-9	[-1, 6]	774
Linear-10	[-1, 128]	196,736
LeakyReLU-11	[-1, 128]	0
Linear-12	[-1, 1]	129
Total params: 472,679		
Trainable params: 472,679		
Non-trainable params: 0		
Input size (MB): 0.08		
Forward/backward pass size (MB): 0.21		
Params size (MB): 1.80		
Estimated Total Size (MB): 2.09		

1