

Motivation

Many recent publications have explored the use of reinforcement learning for Atari games, achieving comparable performance for many games, and superhuman performance for a select few including Pong Mnih et al. [2015], He et al. [2016], Guo et al. [2014], Silver et al. [2016], Mnih et al. [2015]. All of these algorithms utilized a convolutional neural network approach based on the pixel input of the game image. In this work we would like to focus on Pong, and see how the learning performance will be affected if the state space is instead composed of the positions and velocities of the game objects. Our hypothesis is that the richer state space of the game image will lead to higher performance. However, our revised state space is of interest because of its ability to learn policies based on the game objects, as humans do.

Problem Definition

The state, action, and reward definitions are below. The state space is composed of the y positions and y velocities of paddles one and two, and the x and y position and velocity of the ball. The action space is a discrete space of dimension three. The learning agent trains the control for paddle two, and plays against the default AI of paddle one, which is based on the position of the ball. The reward is one when paddle two hits the ball.

$$S = \{ P_{1,y_{pos}} \quad P_{1,y_{vel}} \quad P_{2,y_{pos}} \quad P_{2,y_{vel}} \quad B_{x_{pos}} \quad B_{y_{pos}} \quad B_{x_{vel}} \quad B_{y_{vel}} \} \quad R = \begin{cases} 1 & B_{x_{pos}} = P_{x_{pos}}, B_{y_{pos}} \in P_{1,y_{range}} \\ 0 & \text{else} \end{cases}$$

As stated above, the action space is discrete of dimension three, and controls the speed of paddle two. The paddle either stays still, goes down, or goes up, if A is zero, one, or two respectively. Paddle one either stays still if it is at the same position as the ball, goes down if the ball is below it, or goes up if the ball is above it.

$$A \in \{0, 1, 2\} \quad P_{1,y_{vel}} = \begin{cases} 0 & A = 0 \\ -P_{1,y_{vel_{max}}} & A = 1 \\ P_{1,y_{vel_{max}}} & A = 2 \end{cases} \quad P_{2,y_{vel}} = \begin{cases} 0 & B_{x_{pos}} = P_{2,y_{pos}} \\ -P_{2,y_{vel_{max}}} & B_{x_{pos}} < P_{2,y_{pos}} \\ P_{2,y_{vel_{max}}} & B_{x_{pos}} > P_{2,y_{pos}} \end{cases}$$

Proposed Experiments

In this work, we will test two hypothesis. The first is that the reduced dimension of the modified state space decreases the performance. The second is that the speed of the second paddle affects the performance, as faster velocities in the opponent paddle will lead to reduced reward for suboptimal policies. These are tested in the following experiments

1. Experiment 1: We will benchmark the performance of our custom environment against Pong-V0 in OpenAI gym, using the DQN algorithm from Keras RL.
2. Experiment 2: We will record the performance of our algorithm for multiple speeds of the second paddle, and test our algorithms robustness to this parameter variation.

1 Justification for RL Approach

1.1 Classic Pong Formulation

In the above Pong Formulation, where the positions and velocities of the game object are used, the cardinality of the state space is given by the following equation.

$$|S_{\text{classic}}| = \left[\text{Max} (P_{1,\text{pos}}) \right] \cdot \left[\text{Max} (P_{2,\text{pos}}) \right] \cdot \left[\text{Max} (B_{x,\text{pos}}) \right] \cdot \left[\text{Max} (B_{y,\text{pos}}) \right] \\ \cdot \left[2 \cdot \text{Max} (P_{1,\text{vel}}) \right] \cdot \left[2 \cdot \text{Max} (P_{2,\text{vel}}) \right] \cdot \left[2 \cdot \text{Max} (B_{x,\text{vel}}) \right] \cdot \left[2 \cdot \text{Max} (B_{y,\text{vel}}) \right]$$

We now substitute the parameters for our experiment, in which the screen is 400×600 pixels and the maximum velocity is 20 pixels per second for all objects.

$$|S_{\text{classic}}| = [400Px.] \cdot [400Px.] \cdot [600Px.] \cdot [400Px.] \cdot \left[2 \cdot 20 \frac{Px.}{\text{sec}} \right] \cdot \left[2 \cdot 20 \frac{Px.}{\text{sec}} \right] \cdot \left[2 \cdot 20 \frac{Px.}{\text{sec}} \right] \cdot \left[2 \cdot 20 \frac{Px.}{\text{sec}} \right]$$

The cardinality of the transition matrix is the square of the cardinality of the state matrix.

$$|P_{\text{classic}}| = |S_{\text{classic}}|^2$$

The enormous cardinality of the transition matrix necessitates the use of Reinforcement Learning for this problem.

1.2 Convolutional Pong Formulation

In the Convolutional Pong Formulation used in previous work, where the pixel input is used, we can see the cardinality of the state space is much larger. The problem is framed in an identical fashion, except for the state and reward definitions which are given below.

$$S' = \left\{ R_{x \in X, y \in Y} \quad G_{x \in X, y \in Y} \quad B_{x \in X, y \in Y} \right\} \quad R' = \left\{ P_{2,\text{score}} - P_{1,\text{score}} \right\}$$

The cardinality of the state space becomes the following

$$|S_{\text{convolutional}}| = \left[\text{Max} (X_{\text{pos}}) \right] \cdot \left[\text{Max} (Y_{\text{pos}}) \right] \cdot \left[\text{Max} (R_{\text{val}}) \right] \cdot \left[\text{Max} (G_{\text{val}}) \right] \cdot \left[\text{Max} (B_{\text{val}}) \right]$$

We now substitute the parameters for this experiment, in which the screen is 400×600 pixels and the RGB space is composed of 255 voxels.

$$|S_{\text{convolutional}}| = [400Px.] \cdot [600Px.] \cdot [255Voxels] \cdot [255Voxels] \cdot [255Voxels]$$

The cardinality of the transition matrix is the square of the cardinality of the state matrix.

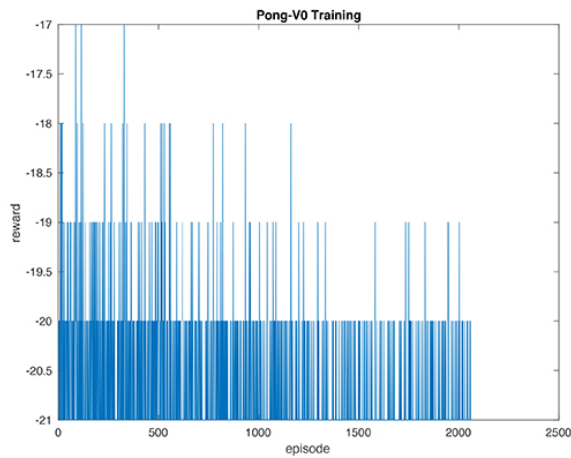
$$|P_{\text{convolutional}}| = |S_{\text{convolutional}}|^2$$

The enormous cardinality of the transition matrix necessitates the use of Reinforcement Learning for this problem as well.

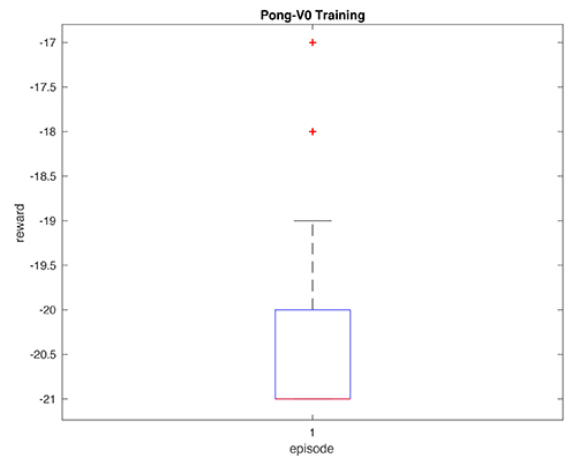
Results

Experiment 1

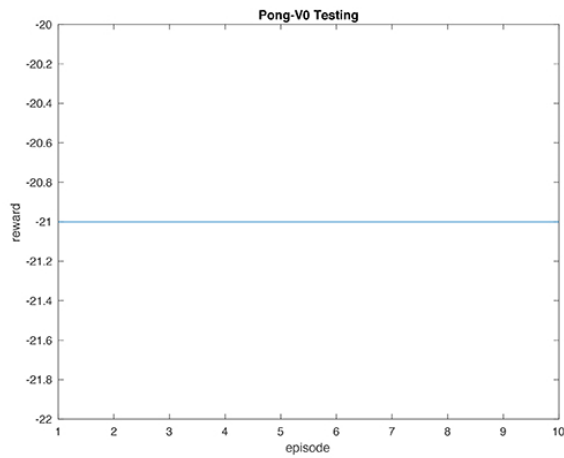
The baseline, the Pong-V0 environment, from Open AI Gym, is trained with the DQG Atari algorithm from Keras RL (Appendix A). Figure 1 shows plots and box plots for both the training and testing phases. The simulation was run for 170000 steps (pixel updates) for twenty-four hours using the CPU version of Tensor Flow on a 2016 MacBook Pro with a 2.9GHz Intel i7 processor, before being tested for ten episodes. It can be observed from Figure 1 that this simulation, which was trained with a reward of the difference in score between the two paddles and a state input of the game image, fails completely even after the twenty-four hour training period.



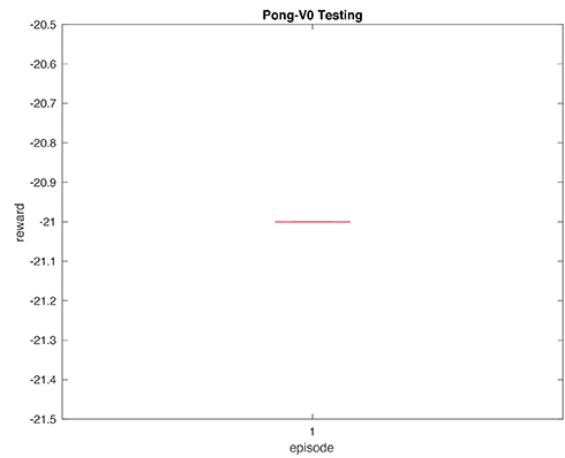
(a) Training Plot



(b) Training Box Plot



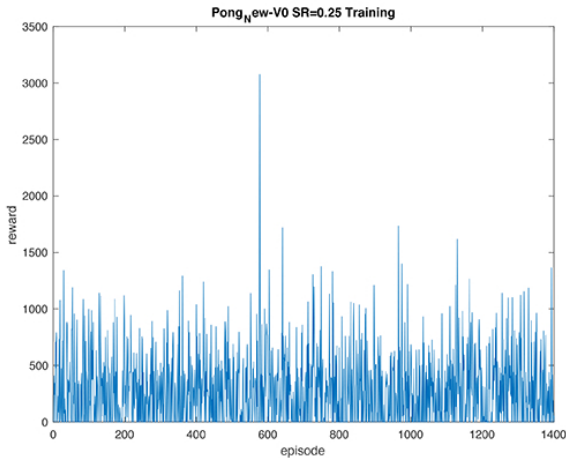
(c) Test Plot



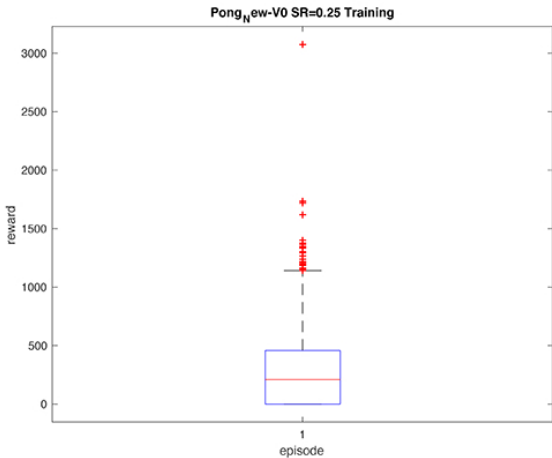
(d) Test Box Plot

Figure 1: Pong-V0

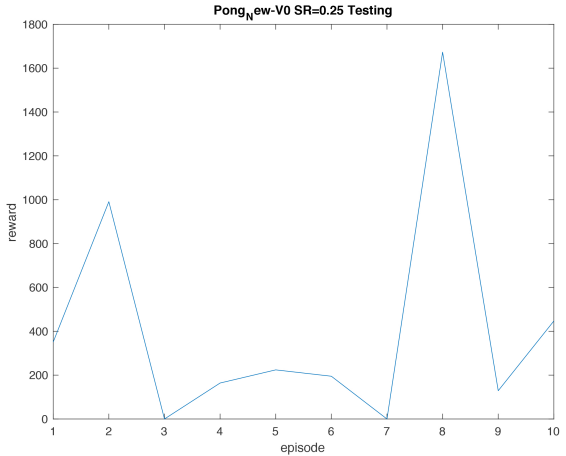
Our environment, Pong New - V0 (Appendix B), which uses the classic control framework from Open AI Gym, was trained with the Classic Control DQN Algorithm from Keras RL (Appendix C). We decreased the simulation time by a factor of ten, running the simulation for 175000 steps over a period of three hours on the same hardware system before testing for ten episodes. This revised learning algorithm, which learns a defensive strategy where the reward is given for hitting the ball, succeeds eighty percent of the time after only training for three hours, with the positions and velocities of the game objects as state input. During this experiment, the ratio of the maximum speed of the paddle we train to the one we compete against is four to one.



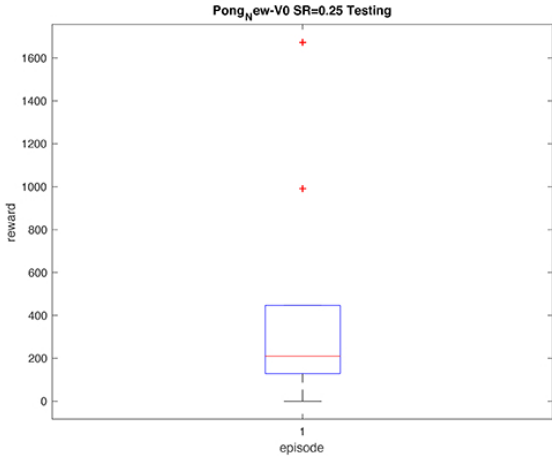
(a) Training Plot



(b) Training Box Plot



(c) Test Plot

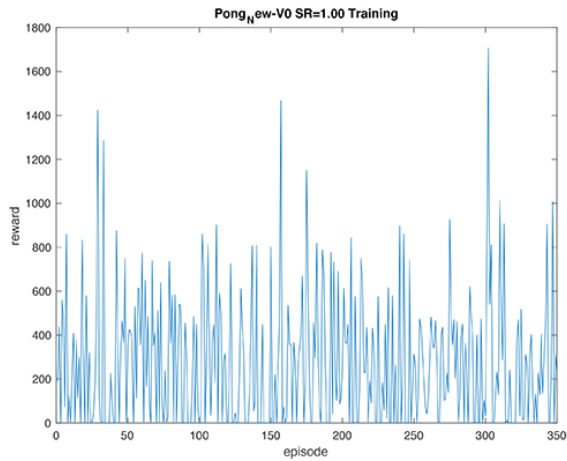


(d) Test Box Plot

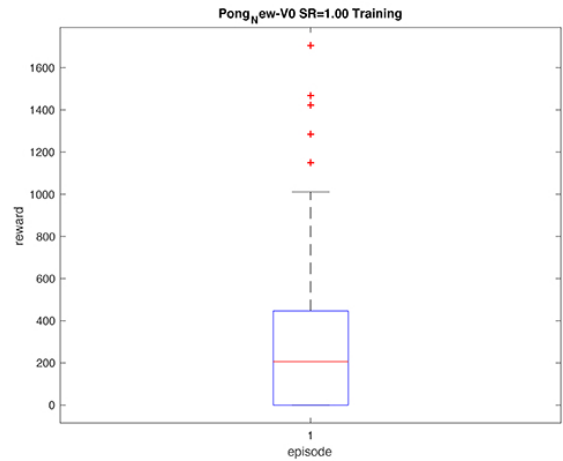
Figure 2: Pong New – V0 SR = 0.25

Experiment 2

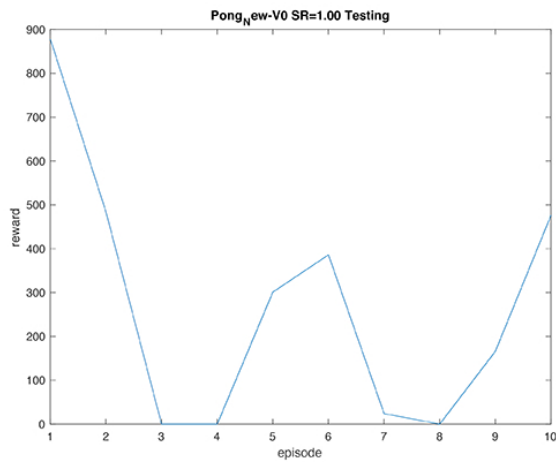
We now repeat the previous experiment on our Pong New - V0 environment, decreasing the ratio of the maximum speed of the paddle we train to the one we compete against to one to one. We observe seventy percent success rate for this experiment, showing the robustness of our environment to this parameter variation.



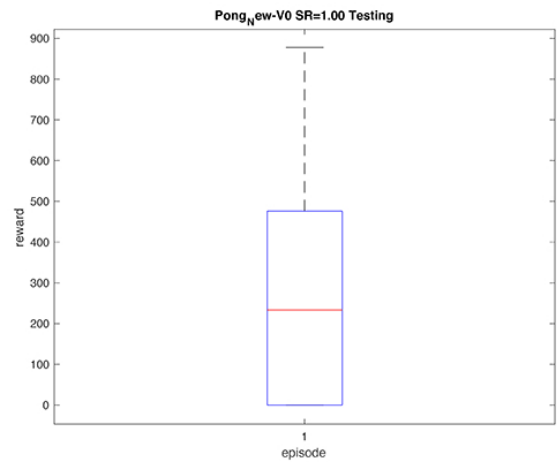
(a) Training Plot



(b) Training Box Plot



(c) Test Plot



(d) Test Box Plot

Figure 3: Pong New – V0 SR = 1.00

2 Utilized Environments

For the training environment, we utilized Open AI Gym Brockman et al. [2016]. For the DQN network, we utilized **Keras RL**, which implements the DQN algorithm utilized in Mnih et al. [2015] as part of its package of modern Deep RL algorithms built on the **Tensor Flow** training environment. Finally, for the Pong New - V0 environment, we created our new environment based on an existing python implementation of the classic **Pong Game** using **Pygames**, and our own previous implementation of **Pong in Java**.

3 Discussion

It is sensible that our approach trains faster, as the dimension of the state space and the one step return on our reward expedite training. The baseline Pong - V0, with its high dimensional image input, and a score based reward which will not give high return until a winning policy is found, cannot be expected to train as efficiently. Since all winning policies lie in the convex hull of all defensive policies, it is sensible to perform two-shot learning for the Pong environment, where the policy trained with our environment is used as a baseline for imitation learning in the original Pong - V0 environment. The state spaces of the two environments can be related to each other by either storing the game image data during training of our environment, or by modifying the Pong - V0 environment to first perform object classification via supervised learning as in Kulkarni et al. [2016], and then learn a policy based on the classified position and velocities of the game objects.

References

- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- Xiaoxiao Guo, Satinder Singh, Honglak Lee, Richard L Lewis, and Xiaoshi Wang. Deep learning for real-time atari game play using offline monte-carlo tree search planning. In *Advances in neural information processing systems*, pages 3338–3346, 2014.
- Frank S He, Yang Liu, Alexander G Schwing, and Jian Peng. Learning to play in a day: Faster deep reinforcement learning by optimality tightening. *arXiv preprint arXiv:1611.01606*, 2016.
- Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In *Advances in Neural Information Processing Systems*, pages 3675–3683, 2016.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.

Appendix

A: DQN Algorithm for Pong - V0: Atari DQN From Keras-RL

```
1 from __future__ import division
2 import argparse
3 from PIL import Image
4 import numpy as np
5 import gym
6 from keras.models import Sequential
7 from keras.layers import Dense, Activation, Flatten, Convolution2D,
   Permute
8 from keras.optimizers import Adam
9 import keras.backend as K
10 from rl.agents.dqn import DQNAgent
11 from rl.policy import LinearAnnealedPolicy, BoltzmannQPolicy,
   EpsGreedyQPolicy
12 from rl.memory import SequentialMemory
13 from rl.core import Processor
14 from rl.callbacks import FileLogger, ModelIntervalCheckpoint
15 import scipy
16 INPUT_SHAPE = (84, 84)
17 WINDOW_LENGTH = 4
18 class AtariProcessor(Processor):
19     def process_observation(self, observation):
20         assert observation.ndim == 3 # (height, width, channel)
21         img = Image.fromarray(observation)
22         img = img.resize(INPUT_SHAPE).convert('L') # resize and
   convert to grayscale
23         processed_observation = np.array(img)
24         assert processed_observation.shape == INPUT_SHAPE
25         return processed_observation.astype('uint8') # saves storage
   in experience memory
26     def process_state_batch(self, batch):
27         # We could perform this processing step in '
   process_observation'. In this case, however,
28         # we would need to store a 'float32' array instead, which is 4
   x more memory intensive than
29         # an 'uint8' array. This matters if we store 1M observations.
30         processed_batch = batch.astype('float32') / 255.
31         return processed_batch
32     def process_reward(self, reward):
33         return np.clip(reward, -1., 1.)
34 parser = argparse.ArgumentParser()
35 parser.add_argument('--mode', choices=['train', 'test'], default='
   train')
36 parser.add_argument('--env-name', type=str, default='Pong-v0')
```

```

37 parser.add_argument('--weights', type=str, default=None)
38 args = parser.parse_args()
39 # Get the environment and extract the number of actions.
40 env = gym.make(args.env_name)
41 env = env.unwrapped
42 np.random.seed(123)
43 env.seed(123)
44 nb_actions = env.action_space.n
45 def _step(a):
46     reward = 0.0
47     action = env._action_set[a]
48     lives_before = env.ale.lives()
49     for _ in range(4):
50         reward += env.ale.act(action)
51     ob = env._get_obs()
52     done = env.ale.game_over() or (args.mode == 'train' and
53         lives_before != env.ale.lives())
54     return ob, reward, done, {}
55 env._step = _step
56 input_shape = (WINDOW_LENGTH,) + INPUT_SHAPE
57 model = Sequential()
58 if K.image_dim_ordering() == 'tf':
59     model.add(Permute((2, 3, 1), input_shape=input_shape))
60 elif K.image_dim_ordering() == 'th':
61     model.add(Permute((1, 2, 3), input_shape=input_shape))
62 else:
63     raise RuntimeError('Unknown image_dim_ordering.')
64 model.add(Convolution2D(32, 8, 8, subsample=(4, 4)))
65 model.add(Activation('relu'))
66 model.add(Convolution2D(64, 4, 4, subsample=(2, 2)))
67 model.add(Activation('relu'))
68 model.add(Convolution2D(64, 3, 3, subsample=(1, 1)))
69 model.add(Activation('relu'))
70 model.add(Flatten())
71 model.add(Dense(512))
72 model.add(Activation('relu'))
73 model.add(Dense(nb_actions))
74 model.add(Activation('linear'))
75 print(model.summary())
76 memory = SequentialMemory(limit=1000000, window_length=WINDOW_LENGTH)
77 processor = AtariProcessor()
78 policy = LinearAnnealedPolicy(EpsGreedyQPolicy(), attr='eps',
79     value_max=1., value_min=.1, value_test=.05,
80     nb_steps=1000000)
81 dqn = DQNAgent(model=model, nb_actions=nb_actions, policy=policy,
82     memory=memory,

```



```

80         processor=processor, nb_steps_warmup=50000, gamma=.99,
           target_model_update=10000,
81         train_interval=4, delta_clip=1.)
82 dqn.compile(Adam(lr=.00025), metrics=['mae'])
83 if args.mode == 'train':
84     weights_filename = 'dqn_{}_weights.h5f'.format(args.env_name)
85     checkpoint_weights_filename = 'dqn_' + args.env_name + '_weights_{
           step}.h5f'
86     log_filename = 'dqn_{}_log.json'.format(args.env_name)
87     callbacks = [ModelIntervalCheckpoint(checkpoint_weights_filename,
           interval=250000)]
88     callbacks += [FileLogger(log_filename, interval=100)]
89     history_0 = dqn.fit(env, callbacks=callbacks, nb_steps=1750000,
           log_interval=10000)
90     dqn.save_weights(weights_filename, overwrite=True)
91     history_1 = dqn.test(env, nb_episodes=10, visualize=False)
92 elif args.mode == 'test':
93     weights_filename = 'dqn_{}_weights.h5f'.format(args.env_name)
94     if args.weights:
95         weights_filename = args.weights
96         dqn.load_weights(weights_filename)
97         history_1 = dqn.test(env, nb_episodes=10, visualize=False)
98 scipy.io.savemat('history_0.mat', history_0.history, appendmat=True,
           format='5', long_field_names=False, do_compression=False, oned_as='
           row')
99 scipy.io.savemat('history_1.mat', history_1.history, appendmat=True,
           format='5', long_field_names=False, do_compression=False, oned_as='
           row')

```

B: Pong New - V0 Environment

```
1 import logging
2 import math
3 import gym
4 from gym import spaces
5 from gym.utils import seeding
6 import numpy as np
7 from os import path
8 import random
9 import pygame, sys
10 from pygame.locals import *
11 WHITE = (255,255,255)
12 RED = (255,0,0)
13 GREEN = (0,255,0)
14 BLACK = (0,0,0)
15 MAX_BALL_VEL = 20
16 WIDTH = 600
17 HEIGHT = 400
18 BALL_RADIUS = 20
19 PAD_WIDTH = 8
20 PAD_HEIGHT = 80
21 HALF_PAD_WIDTH = PAD_WIDTH // 2
22 HALF_PAD_HEIGHT = PAD_HEIGHT // 2
23 paddle1_pos = [HALF_PAD_WIDTH - 1, HEIGHT//2]
24 paddle2_pos = [WIDTH +1 - HALF_PAD_WIDTH, HEIGHT//2]
25 paddle1_vel = 0
26 paddle2_vel = 0
27 ball_pos = [WIDTH//2, HEIGHT//2]
28 ball_vel = [0, 0]
29 r_score_threshold = 100
30 l_score_threshold = 100
31 max_paddle1_vel = 20
32 max_paddle2_vel = 20
33 min_paddle1_vel = -20
34 min_paddle2_vel = -20
35 max_paddle1_pos = HEIGHT - HALF_PAD_HEIGHT
36 max_paddle2_pos = HEIGHT - HALF_PAD_HEIGHT
37 max_ball_pos = [WIDTH - BALL_RADIUS, HEIGHT - BALL_RADIUS]
38 max_ball_vel = [MAX_BALL_VEL, MAX_BALL_VEL]
39 min_paddle1_pos = HALF_PAD_HEIGHT
40 min_paddle2_pos = HALF_PAD_HEIGHT
41 min_ball_pos = [BALL_RADIUS, BALL_RADIUS]
42 min_ball_vel = [-MAX_BALL_VEL, -MAX_BALL_VEL]
43 logger = logging.getLogger(__name__)
44 global l_score, r_score, reward, reward_curr
45 global paddle1_pos, paddle2_pos, ball_pos, ball_vel
```

```

46 class PongEnv(gym.Env):
47     metadata = {
48         'render.modes' : ['human', 'rgb_array'],
49         'video.frames_per_second' : 30
50     }
51     def __init__(self):
52         global l_score, r_score, reward, reward_curr
53         global high, low
54         global paddle1_pos, paddle2_pos, ball_pos, ball_vel,
55             paddle1_vel, paddle2_vel
56         global ball_pos, ball_vel
57         def ball_init():
58             global ball_vel
59             horz      = 0
60             vert      = 0
61             while (horz == 0) or (vert == 0):
62                 horz      = random.randrange(-MAX_BALL_VEL, MAX_BALL_VEL
63                 )
64                 vert      = random.randrange(-MAX_BALL_VEL, MAX_BALL_VEL
65                 )
66                 if random.randrange(0,2) is not 0:
67                     horz = - horz
68                 ball_vel = [horz, -vert]
69         r_score = 0
70         l_score = 0
71         reward  = 0
72         high = np.array([max_paddle1_pos, max_paddle1_vel,
73             max_paddle2_pos, max_paddle2_vel, max_ball_pos[0],
74             max_ball_pos[1], max_ball_vel[0], max_ball_vel[1]])
75         low = np.array([min_paddle1_pos, min_paddle1_vel,
76             min_paddle2_pos, min_paddle2_vel, min_ball_pos[0],
77             min_ball_pos[1], min_ball_vel[0], min_ball_vel[1]])
78         self.action_space = spaces.Discrete(3)
79         self.observation_space = spaces.Box(low, high)
80         self._seed()
81         self.viewer = None
82         self.state = None
83         self.steps_beyond_done = None
84         self.steps_to_done = 0
85         self.state = self.np_random.uniform(low, high, size=(8,))
86         state = self.state
87         (paddle1_pos[1], paddle1_vel, paddle2_pos[1], paddle2_vel,
88             ball_pos[0], ball_pos[1], ball_vel[0], ball_vel[1]) = state
89         ball_pos = [WIDTH//2, HEIGHT//2]
90         ball_init()
91         self.state = (paddle1_pos[1], paddle1_vel, paddle2_pos[1],
92             paddle2_vel, ball_pos[0], ball_pos[1], ball_vel[0],

```

```

    ball_vel[1])
84 def _seed(self, seed=None):
85     self.np_random, seed = seeding.np_random(seed)
86     return [seed]
87 def _step(self, action):
88     self.steps_to_done += 1;
89     global paddle1_pos, paddle2_pos, paddle1_vel, paddle2_vel,
        l_score, r_score
90     global l_score, r_score, reward, reward_curr
91     global ball_pos, ball_vel
92     def ball_init():
93         global ball_vel
94         horz      = 0
95         vert      = 0
96         while (horz == 0) or (vert == 0):
97             horz      = random.randrange(-MAX_BALL_VEL, MAX_BALL_VEL
        )
98             vert      = random.randrange(-MAX_BALL_VEL, MAX_BALL_VEL
        )
99             if random.randrange(0,2) is not 0:
100                 horz = - horz
101                 ball_vel = [horz, -vert]
102     assert self.action_space.contains(action), "%r (%s) invalid"
        %(action, type(action))
103     state = self.state
104     (paddle1_pos[1], paddle1_vel, paddle2_pos[1], paddle2_vel,
        ball_pos[0], ball_pos[1], ball_vel[0], ball_vel[1]) = state
105     #update paddle velocity
106     if ball_pos[1] < paddle1_pos[1]:
107         paddle1_vel = -max_paddle2_vel
108     elif ball_pos[1] > paddle1_pos[1]:
109         paddle1_vel = max_paddle2_vel
110     else:
111         paddle1_vel = 0
112     if action == 0:
113         paddle2_vel = -max_paddle1_vel
114     elif action == 1:
115         paddle2_vel = max_paddle1_vel
116     elif action == 2:
117         paddle2_vel = 0
118     if paddle1_pos[1] > HALF_PAD_HEIGHT and paddle1_pos[1] <
        HEIGHT - HALF_PAD_HEIGHT:
119         paddle1_pos[1] += paddle1_vel
120     elif paddle1_pos[1] < HALF_PAD_HEIGHT and paddle1_vel > 0:
121         paddle1_pos[1] += paddle1_vel
122     elif paddle1_pos[1] > HEIGHT - HALF_PAD_HEIGHT and paddle1_vel
        < 0:

```

```

123     paddle1_pos[1] += paddle1_vel
124     elif paddle1_pos[1] < HALF_PAD_HEIGHT and paddle1_vel < 0:
125         paddle1_pos[1] -= paddle1_vel
126     elif paddle1_pos[1] > HEIGHT - HALF_PAD_HEIGHT and paddle1_vel
127         > 0:
128         paddle1_pos[1] -= paddle1_vel
129     if paddle2_pos[1] > HALF_PAD_HEIGHT and paddle2_pos[1] <
130         HEIGHT - HALF_PAD_HEIGHT:
131         paddle2_pos[1] += paddle2_vel
132     elif paddle2_pos[1] < HALF_PAD_HEIGHT and paddle2_vel > 0:
133         paddle2_pos[1] += paddle2_vel
134     elif paddle2_pos[1] > HEIGHT - HALF_PAD_HEIGHT and paddle2_vel
135         < 0:
136         paddle2_pos[1] += paddle2_vel
137     elif paddle2_pos[1] < HALF_PAD_HEIGHT and paddle2_vel < 0:
138         paddle2_pos[1] -= paddle2_vel
139     elif paddle2_pos[1] > HEIGHT - HALF_PAD_HEIGHT and paddle2_vel
140         > 0:
141         paddle2_pos[1] -= paddle2_vel
142     ball_pos[0] += int(ball_vel[0])
143     ball_pos[1] += int(ball_vel[1])
144     if int(ball_pos[1]) <= BALL_RADIUS:
145         ball_vel[1] = - ball_vel[1]
146     if int(ball_pos[1]) >= HEIGHT + 1 - BALL_RADIUS:
147         ball_vel[1] = - ball_vel[1]
148     if int(ball_pos[0])<=BALL_RADIUS+PAD_WIDTH and int(ball_pos
149         [1]) in range(int(paddle1_pos[1]-HALF_PAD_HEIGHT),int(
150         paddle1_pos[1]+HALF_PAD_HEIGHT),1):
151         ball_vel[0] = -ball_vel[0]
152         ball_vel[0] *= 1.1
153         ball_vel[1] *= 1.1
154     elif int(ball_pos[0]) <= BALL_RADIUS + PAD_WIDTH:
155         #reward += 2
156         r_score += 1
157         ball_pos = [WIDTH//2, HEIGHT//2]
158         ball_init()
159
160     if int(ball_pos[0])>=WIDTH+1-BALL_RADIUS-PAD_WIDTH and int(
161         ball_pos[1]) in range(int(paddle2_pos[1]-HALF_PAD_HEIGHT),
162         int(paddle2_pos[1]+HALF_PAD_HEIGHT),1):
163         reward += 1
164         ball_vel[0] = -ball_vel[0]
165         ball_vel[0] *= 1.1
166         ball_vel[1] *= 1.1
167     elif int(ball_pos[0]) >= WIDTH + 1 - BALL_RADIUS - PAD_WIDTH:
168         l_score += 1
169         ball_pos = [WIDTH//2, HEIGHT//2]

```

```

162     ball_init()
163     self.state = (paddle1_pos[1], paddle1_vel, paddle2_pos[1],
164                 paddle2_vel, ball_pos[0], ball_pos[1], ball_vel[0],
165                 ball_vel[1])
166 done = (l_score > l_score_threshold) or (r_score >
167       r_score_threshold)
168 done = bool(done)
169 done = bool(done)
170 if not done:
171     self.steps_beyond_done = self.steps_beyond_done
172     reward = reward
173 elif self.steps_beyond_done is None:
174     self.steps_beyond_done = 0
175     reward = reward
176 else:
177     if self.steps_beyond_done == 0:
178         logger.warning("You are calling 'step()' even though
179             this environment has already returned done = True.
180             You should always call 'reset()' once you receive '
181             done = True' — any further steps are undefined
182             behavior.")
183     self.steps_beyond_done += 1
184     self.steps_to_done = 0
185     reward = 0
186 if r_score > r_score_threshold:
187     r_score = 0
188     l_score = 0
189     reward_curr = 0
190 if l_score > l_score_threshold:
191     r_score = 0
192     l_score = 0
193     reward_curr = 0
194 return np.array(self.state), reward, done, {}
195 def _reset(self):
196     self.steps_to_done = 0
197     self.steps_beyond_done = None
198     global l_score, r_score, reward, reward_curr
199     global paddle1_pos, paddle2_pos, ball_pos, ball_vel
200     def ball_init():
201         global ball_vel
202         horz = 0
203         vert = 0
204         while (horz == 0) or (vert == 0):
205             horz = random.randrange(-MAX_BALL_VEL, MAX_BALL_VEL
206                 )
207             vert = random.randrange(-MAX_BALL_VEL, MAX_BALL_VEL
208                 )

```

```

200         if random.randrange(0,2) is not 0:
201             horz = - horz
202             ball_vel = [horz, -vert]
203     r_score = 0
204     l_score = 0
205     reward = 0
206     self.state = self.np_random.uniform(low, high, size=(8,))
207     state = self.state
208     (paddle1_pos[1], paddle1_vel, paddle2_pos[1], paddle2_vel,
209      ball_pos[0], ball_pos[1], ball_vel[0], ball_vel[1]) = state
210     ball_pos = [WIDTH//2, HEIGHT//2]
211     ball_init()
212     self.state = (paddle1_pos[1], paddle1_vel, paddle2_pos[1],
213                  paddle2_vel, ball_pos[0], ball_pos[1], ball_vel[0],
214                  ball_vel[1])
215     return np.array(self.state)
216 def _render(self, mode='human', close=False):
217     if close:
218         if self.viewer is not None:
219             self.viewer.close()
220             self.viewer = None
221             pygame.quit()
222             sys.exit()
223         return
224     if self.viewer is None:
225         pygame.init()
226         fps = pygame.time.Clock()
227         window = pygame.display.set_mode((WIDTH, HEIGHT), 0, 32)
228         pygame.display.set_caption('Hello World')
229         window.fill(BLACK)
230         pygame.draw.line(window, WHITE, [WIDTH // 2, 0],[WIDTH //
231          2, HEIGHT], 1)
232         pygame.draw.line(window, WHITE, [PAD_WIDTH, 0],[PAD_WIDTH,
233          HEIGHT], 1)
234         pygame.draw.line(window, WHITE, [WIDTH - PAD_WIDTH, 0],[
235          WIDTH - PAD_WIDTH, HEIGHT], 1)
236         pygame.draw.circle(window, WHITE, [WIDTH//2, HEIGHT//2],
237          70, 1)
238         ball_pos_int = [int(ball_pos[0]), int(ball_pos[1])]
239         pygame.draw.circle(window, RED, ball_pos_int, 20, 0)
240         pygame.draw.polygon(window, GREEN, [[paddle1_pos[0] -
241          HALF_PAD_WIDTH, paddle1_pos[1] - HALF_PAD_HEIGHT], [
242          paddle1_pos[0] - HALF_PAD_WIDTH, paddle1_pos[1] +
243          HALF_PAD_HEIGHT], [paddle1_pos[0] + HALF_PAD_WIDTH,
244          paddle1_pos[1] + HALF_PAD_HEIGHT], [paddle1_pos[0] +
245          HALF_PAD_WIDTH, paddle1_pos[1] - HALF_PAD_HEIGHT]], 0)

```

```
234     pygame.draw.polygon(window, GREEN, [[paddle2_pos[0] -
        HALF_PAD_WIDTH, paddle2_pos[1] - HALF_PAD_HEIGHT], [
        paddle2_pos[0] - HALF_PAD_WIDTH, paddle2_pos[1] +
        HALF_PAD_HEIGHT], [paddle2_pos[0] + HALF_PAD_WIDTH,
        paddle2_pos[1] + HALF_PAD_HEIGHT], [paddle2_pos[0] +
        HALF_PAD_WIDTH, paddle2_pos[1] - HALF_PAD_HEIGHT]], 0)
235 myfont1 = pygame.font.SysFont("Comic Sans MS", 20)
236 label1 = myfont1.render("Score: "+str(l_score), 1,
        (255,255,0))
237 window.blit(label1, (50,20))
238 myfont2 = pygame.font.SysFont("Comic Sans MS", 20)
239 label2 = myfont2.render("Reward: "+str(reward), 1,
        (255,255,0))
240 window.blit(label2, (210, 20))
241 myfont2 = pygame.font.SysFont("Comic Sans MS", 20)
242 label2 = myfont2.render("Score: "+str(r_score), 1,
        (255,255,0))
243 window.blit(label2, (470, 20))
244 pygame.display.update()
245 fps.tick(200)
246 if self.state is None: return None
247 else: return None
```

C: DQN Algorithm for Pong New - V0: Classic Control DQN From Keras-RL

```
1 import numpy as np
2 import gym
3 import scipy
4 from keras.models import Sequential
5 from keras.layers import Dense, Activation, Flatten
6 from keras.optimizers import Adam
7 from rl.agents.dqn import DQNAgent
8 from rl.policy import EpsGreedyQPolicy
9 from rl.memory import SequentialMemory
10 ENV_NAME = 'pong_new-v0'
11 env = gym.make(ENV_NAME)
12 np.random.seed(123)
13 env.seed(123)
14 nb_actions = env.action_space.n
15 model = Sequential()
16 model.add(Flatten(input_shape=(1,) + env.observation_space.shape))
17 model.add(Dense(16))
18 model.add(Activation('relu'))
19 model.add(Dense(16))
20 model.add(Activation('relu'))
21 model.add(Dense(16))
22 model.add(Activation('relu'))
23 model.add(Dense(16))
24 model.add(Activation('relu'))
25 model.add(Dense(16))
26 model.add(Activation('relu'))
27 model.add(Dense(16))
28 model.add(Activation('relu'))
29 model.add(Dense(nb_actions))
30 model.add(Activation('linear'))
31 print(model.summary())
32 memory = SequentialMemory(limit=50000, window_length=1)
33 policy = EpsGreedyQPolicy()
34 dqn = DQNAgent(model=model, nb_actions=nb_actions, memory=memory,
35               nb_steps_warmup=50000, target_model_update=1e-2, policy=policy)
36 dqn.compile(Adam(lr=1e-3), metrics=['mae'])
37 hist_0 = dqn.fit(env, nb_steps=175000, visualize=False, verbose=2,
38               nb_max_episode_steps=10000)
39 dqn.save_weights('dqn_{}_weights.h5f'.format(ENV_NAME), overwrite=True)
40 hist_1 = dqn.test(env, nb_episodes=10, visualize=False)
41 scipy.io.savemat('hist_0.mat', history_0.history, appendmat=True, format='5',
42               long_field_names=False, do_compression=False, oned_as='row')
43 scipy.io.savemat('hist_1.mat', history_1.history, appendmat=True, format='5',
44               long_field_names=False, do_compression=False, oned_as='row')
```
